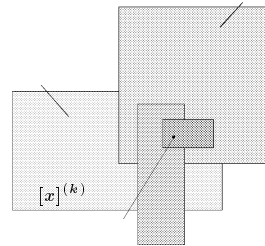


Institut für **A**ngewandte **M**athematik  
Universität **K**arlsruhe (TH)

# PASCAL-XSC BCD-Version 1.0

Benutzerhandbuch  
für das  
dezimale Laufzeitsystem

Frithjof Blomquist



$[x]^{(k+1)}$

**F**orschungsschwerpunkt  
**C**omputerarithmetik,  
**I**ntervallrechnung und  
**N**umerische Algorithmen mit  
**E**rgebnisverifikation

$x^*$



# PASCAL-XSC BCD-Version 1.0

## Benutzerhandbuch für das dezimale Laufzeitsystem

Frithjof Blomquist

Jeder Rechner liefert wegen unvermeidbarer Rundungsfehler z.B. für die Addition zweier Gleitkommazahlen eine i. a. fehlerbehaftete Summe, so daß man durch **naive** Gleitkommarechnungen Maschinenergebnisse erhält, über deren Fehler **gesicherte** Aussagen bei umfangreichen Algorithmen mit möglicherweise Hunderttausenden von Gleitkommaoperationen nicht möglich sind.

Zur Behebung dieser aus mathematischer Sicht sehr unbefriedigenden Situation stellt die **PASCAL-XSC** Umgebung Werkzeuge zur Verfügung, mit denen für viele mathematische Standardprobleme eine sehr genaue und vor allem **garantierte** Einschließung der jeweils exakten Lösung berechnet werden kann. Dazu müssen lediglich entsprechende Routinen, etwa zur Lösung linearer Gleichungssysteme, aus Hilfsmodulen aufgerufen werden, wobei ein Hintergrundwissen bzgl. der benutzten, selbstverifizierenden Algorithmen nicht erforderlich ist!

In diesem Handbuch findet der Anwender zur Entwicklung eigener, selbstverifizierender Algorithmen eine ausführliche Beschreibung der verfügbaren Zahlenformate mit 13 bzw. 21 Dezimalziffern und 47 **hochgenaue** Standardfunktionen für Punkt- und Intervallargumente. Mit einer Vielzahl von Funktionen, Prozeduren und Operatoren für reelle und komplexe Operanden läßt sich die **PASCAL-XSC** Umgebung mit dem Modulkonzept beliebig erweitern.

Die vorliegende **BCD**-Version ermöglicht das **rundungsfehlerfreie** Abspeichern dezimaler Eingabedaten, so daß die bei der **IEEE**-Version oft so lästigen Konversionsfehler entfallen!

© 1997 Frithjof Blomquist  
Adlerweg 6  
66346 Püttlingen 3

Text, Abbildungen, Tabellen und Programme wurden mit größter Sorgfalt erarbeitet. Der Autor kann jedoch für eventuell verbleibende fehlerhafte Angaben und deren Folgen weder eine juristische Verantwortung noch irgendeine Haftung übernehmen.

Das vorliegende Handbuch ist urheberrechtlich geschützt. Alle Rechte vorbehalten. Kein Teil dieses Buches darf ohne schriftliche Genehmigung des Autors in irgendeiner Form durch Fotokopie, Mikrofilm oder andere Verfahren reproduziert oder in eine für Maschinen, insbesondere Datenverarbeitungsanlagen, verwendbare Sprache übertragen werden. Auch die Rechte der Wiedergabe durch Vortrag, Funk und Fernsehen oder andere Medien sind vorbehalten.

Die in diesem Buch erwähnten Software- und Hardwarebezeichnungen sind in den meisten Fällen eingetragene Warenzeichen und unterliegen damit den gesetzlichen Bestimmungen.

# Vorwort

Zweifellos hat die Computer-Technik in den letzten Jahrzehnten gewaltige Fortschritte gemacht; so kann man ohne Übertreibung sagen, daß sich die Prozessorleistungen geradezu explosionsartig entwickelt haben. Betrachtet man insbesondere noch die zukünftigen Möglichkeiten der Parallelverarbeitung mit mehreren Prozessoren, so ist ein Ende dieser Leistungssteigerungen heute keinesfalls absehbar.

Nun könnte man vermuten, daß der Computer, ähnlich etwa wie in der Multi-Media-Technik, auch für den Bereich der Mathematik die gleichen gewaltigen Fortschritte gebracht hätte. Aus der Sicht der Mathematiker ist jedoch die Prozessorentwicklung, von der Geschwindigkeitssteigerung einmal abgesehen, in ihrem Anfangsstadium steckengeblieben. So werden bis heute lediglich die vier Grundoperationen in maximaler Genauigkeit und eine recht kleine Anzahl von Standardfunktionen bereitgestellt, deren Fehlerschranken, aus welchen Gründen auch immer, nicht zu erfahren sind. Mit diesen minimalen Prozessormöglichkeiten wurden bisher numerische Algorithmen entwickelt, die wegen unvermeidbarer Rundungsfehler bei den Grundoperationen für die mathematisch exakten Ergebnisse lediglich Näherungen geliefert haben. Für den Mathematiker sind solche Näherungen jedoch völlig unbrauchbar, wenn deren Güte nicht bekannt ist. Man benötigt daher entweder die Berechnung einer garantierten Fehlerschranke, oder man muß das anstehende numerische Problem mit einem selbstverifizierenden Algorithmus lösen, der automatisch eine möglichst enge Intervalleinschließung der exakten Lösung liefert. Beide Aufgaben lassen sich mit den XSC-Sprachen wie **PASCAL-XSC**, **FORTRAN-XSC** oder **C-XSC** realisieren, die am Institut für Angewandte Mathematik der Universität Karlsruhe unter der Leitung von Professor U. Kulisch entwickelt wurden.

Mit diesen Werkzeugen ist man heute in der Lage, nur aus numerischen Gleitkommaergebnissen wirkliche mathematische Schlüsse zu ziehen. Hat man z.B. für  $x \in [a, b]$  eine Einschließung  $Y = [y_1, y_2]$  des Wertebereichs der ersten Ableitung  $f'$  einer Funktion  $f : [a, b] \rightarrow \mathbb{R}$  berechnet, und gilt für das Gleitkommaergebnis  $y_1$  die Beziehung  $y_1 > 0$ , so ist mit dem Computer der Beweis geführt worden, daß die Funktion  $f$  in  $[a, b]$  streng monoton wächst. Beachten Sie bitte, daß sich diese Arbeitsweise von den sonst üblichen numerischen Methoden ganz wesentlich unterscheidet. Bisher hat man nämlich meist darauf vertraut, daß das Ergebnis einer Computerrechnung wenigstens auf den ersten Ziffern korrekt sein wird, wenn man instabile Teile des Algorithmus z.B. in doppelter Genauigkeit auswertet.

Mit diesen Methoden lassen sich natürlich keine mathematischen Beweise führen, zumal man auch bei doppeltgenauen Rechnungen immer wieder Problemstellungen findet, die nicht einmal annähernd korrekt gelöst werden können.

Die schon genannten **XSC**-Sprachen liefern zwar alle wesentlichen Hilfsmittel, die für eine mathematisch zuverlässige Numerik erforderlich sind, diese Hilfsmittel müssen jedoch alle softwaremäßig simuliert werden, da sie bis heute nicht direkt durch einen Prozessor realisiert worden sind. Durch die notwendigen Software-Simulationen besitzen die **XSC**-Systeme natürlich im Vergleich zur herkömmlichen Numeriksoftware einen großen Laufzeitnachteil. Aus Sicht des Mathematikers wäre es daher wünschenswert, zunächst wenigstens die vier Grundoperationen der Intervallrechnung direkt durch den Prozessor bereitzustellen, da die Intervallarithmetic insbesondere bei den selbstverifizierenden Algorithmen ein unentbehrliches Hilfsmittel ist.

In dem folgenden Beispiel wird aus mathematischer Sicht eine weitere ganz wesentliche Schwäche der heutigen Prozessorgeneration demonstriert. Mit einem Taschenrechner soll dazu in zehnstelliger Dezimalarithmetik die Differenz  $a \cdot b - c$  berechnet werden. Mit den Werten  $a = b = 1234567891$  und  $c = 1.524157877 \cdot 10^{18}$  erhält man an Stelle des exakten Wertes  $a \cdot b - c = 4.88187881 \cdot 10^8$  das völlig falsche Taschenrechnerergebnis 0, d.h. schon nach zwei Grundoperationen erhält man noch nicht einmal die erste Ziffer des korrekten Ergebnisses. Der Grund liegt darin, daß das exakte Produkt  $a \cdot b = 1524157877488187881$  vor der Subtraktion von  $c$  bei diesem Beispiel genau auf den Wert von  $c$  abgerundet wird, d.h. der Mantissenanteil, der das exakte Ergebnis repräsentiert, geht durch den Rundungsfehler verloren. Ganz entsprechende Beispiele lassen sich auch für PC's oder Workstations angeben, wobei jedoch die Zahlenangaben in dem für uns unlesbaren Binärsystem erfolgen müßten. Für den Mathematiker ist diese Situation natürlich eine einzige Katastrophe, da die im obigen Beispiel angesprochene Skalarproduktberechnung z.B. beim Lösen von linearen Gleichungssystemen ein zentraler Punkt ist, und auch Rechnungen im doppeltlangen Zahlenformat die kontrollierte Auswertung eines Skalarproduktes bei höheren Dimensionen nicht gewährleisten. Dies wird erst durch die **XSC**-Sprachen ermöglicht, die softwaremäßig einen Festkommaspeicher (Akkumulator) bereitstellen, in den man alle möglichen Teilprodukte in beliebiger Anzahl **rundungsfehlerfrei** addieren kann. Da die Simulation dieses Akkumulators in Software recht zeitaufwendig ist, wurde unter Federführung des Instituts für Angewandte Mathematik der Universität Karlsruhe ein Coprozessor entwickelt, mit dem Skalarprodukte beliebiger Dimension in Hardware exakt berechnet werden können. Dabei zeigte sich, daß die Laufzeit sogar kürzer ist als bei der Auswertung mit Hilfe der Grundoperationen  $*$ ,  $+$ , die zudem noch keine kontrollierten Ergebnisse liefern. Es wäre daher schon ein großer Fortschritt, wenn man auf einem Mathematik-Prozessor neben einer Intervallarithmetic auch ein oder zwei (für Parallelverarbeitung) Festkommaspeicher realisieren würde!

Aber selbst ein solcher Prozessor hätte für die Anforderungen der Mathematiker immer noch einen schwerwiegenden Nachteil, da die Arithmetiken aller heutigen Prozessoren in PC's und Workstations auf dem binären Zahlensystem basieren, während alle Mathematiker ihre Probleme im dezimalen Zahlensystem formulieren,

ganz einfach, weil man die Größe einer Zahl im Zehnersystem am einfachsten erfassen kann. Soll nun z.B. mit  $a = 9.000000000000001$  und  $b = 9$  die Differenz  $a - b$  berechnet werden, so müssen  $a, b$  nach der Tastatureingabe jeweils zur nächsten Binärzahl gerundet werden, was mit der Dezimalzahl  $a$  nicht fehlerfrei möglich ist. Die anschließende exakte Subtraktion liefert dann das völlig unbrauchbare Ergebnis  $1.776.. \cdot 10^{-15}$ , das nur durch den Konversionsfehler bei der Rundung von  $a$  ins Binärsystem begründet ist. Im Gegensatz dazu liefert nun jeder Taschenrechner mit einer dezimalen Arithmetik bei 16stelliger Mantisse die exakte Differenz  $1 \cdot 10^{-15}$ , und vermutlich könnte man nicht einen einzigen Taschenrechner verkaufen, der für obige Differenz den falschen Wert  $1.776.. \cdot 10^{-15}$  berechnet. Auf den Seiten 4 ... findet man weitere Beispiele, die den oft so lästigen Einfluß der Konversionsfehler eindrucksvoll demonstrieren. Der Leser sollte diese Beispiele sehr sorgfältig studieren, da die Erfahrung zeigt, daß selbst versierte Anwender numerischer Software die Folgen dieser Konversionsfehler nicht immer richtig einschätzen!

Es war somit ein naheliegender Schritt, eine **PASCAL-XSC** BCD-(Binary Coded Decimal)-Version zu entwickeln, die alle Sprachelemente der Binärversion umfaßt. Dem Anwender stehen zwei Datenformate (real, Longreal) mit 13 bzw. 21 Dezimalziffern zur Verfügung, wobei das vollständige Akkumulatorkonzept wie bei der Binärversion auch für das real-Format existiert. Im Longreal-Format sind alle Additionen in den Akkumulator möglich, während die Multiplikation bzgl. des Exponentenbereichs auf das real-Format eingeschränkt ist. Mit einer Vielzahl von Hilfsroutinen läßt sich das Akkumulatorkonzept vollständig auf das Longreal-Format übertragen. Damit wird z.B. auch die komplexe Intervallrechnung mit 21 Dezimalstellen ermöglicht. Für beide Datenformate existieren jeweils 47 hochgenaue Standardfunktionen mit Angabe der jeweiligen relativen Fehlerschranke. Dem Anwender steht ein weiteres Datenformat (rrG) mit 26 Dezimalstellen zur Verfügung, dessen Exponentenbereich nur durch die integer-Zahlen begrenzt ist. In diesem Format können daher Zwischenrechnungen praktisch ohne Über- oder Unterlauf durchgeführt werden, und zu den 47 schon existierenden Standardfunktionen lassen sich weitere Funktionen für reelle oder komplexe Punkt- bzw. Intervallargumente entwickeln. Durch Rundung der Funktionswerte ins Longreal-Format erhält man dann zusätzliche **hochgenaue** Anwenderfunktionen, d.h. das BCD-System läßt sich beliebig erweitern. Zusammen mit den vielen mitgelieferten Hilfsroutinen dieser Version besitzt der Anwender ein sehr mächtiges Werkzeug, mit dem sich aus garantierten numerischen Ergebnissen wirklich mathematische Aussagen gewinnen lassen. Wenn später einmal ein Prozessor mit einer dezimalen Intervallarithmetik und einigen Festkommaspeichern zur Verfügung steht, dann wird man die jetzt bei aufwendigen Anwendungen unbefriedigenden Laufzeiten noch wesentlich verbessern können.

Abschließend sei noch erwähnt, daß auch mir bekannt ist, daß eine Programmiersprache niemals abgeschlossen und vollkommen und stets noch verbesserungsfähig ist. Wohlwollende kritische Bemerkungen und Verbesserungsvorschläge werden daher gerne entgegengenommen.

Karlsruhe, im Juli 1997

Dr. Frithjof Blomquist

## Zur Entstehung dieses Buches

In einer ersten Version wurden die Standardfunktionen im 21-stelligen Longreal-Format mit einer relativen Fehlerschranke von ca.  $10^{-18}$  entwickelt. Für Punktar-gumente ergaben sich daher im Vergleich zu den exakten Funktionswerten Abwei-chungen bei den letzten zwei oder drei Mantissenziffern, und entsprechend waren die Funktionswerteinschließungen aufgebläht. Es schien mir dann jedoch ratsam zu sein, dem Anwender eines Mathematik-Paketes, das verifizierendes Rechnen ermöglicht, Standardfunktionen anzubieten, deren Maschinenergebnisse sich vom optimal ge-rundeten Funktionswert auf der letzten Mantissenstelle höchstens um **eine Einheit** unterscheidet. Zur Implementierung solcher **hochgenauen** Standardfunktionen im Longreal-Bereich wurden dann in einer zweiten Version die entsprechenden Algorithmen in einem doppeltlangen real-Format mit 26 BCD-Stellen realisiert. Insbeson-dere für die Intervallfunktionen waren zusätzliche Hilfsfunktionen und Programme notwendig, so daß die entsprechende Sammlung von Notizen und Anmerkungen immer unübersichtlicher wurde. Auch die Vielzahl der notwendigen Operatoren für die verschiedenen Datentypen `real`, `Longreal`, `interval`, `Linterval`, .... veranlaßte mich schließlich, die vorliegende Blattsammlung mit dem Textsystem  $\LaTeX$  in Buch-form zu bringen. Dabei war es nicht mein Ziel, eine neue Sprachbeschreibung für PASCAL-XSC zu entwickeln, da dies bereits durch die Autoren R. Klatte, U. Kulisch, M. Neaga, D. Ratz und Ch. Ullrich mit dem Buch: *PASCAL-XSC, Sprachbe-schreibung mit Beispielen* realisiert wurde. Im vorliegenden Buch werden vielmehr Abweichungen von der Binärversion und Besonderheiten der BCD-Version beschrie-ben und mit vielen Beispielen belegt. Dadurch erscheint die Gliederung nicht immer zwingend systematisch; ein umfangreiches Stichwortverzeichnis soll dem Leser des-halb eine zusätzliche Orientierungstütze sein. Mit Hilfe dieses Buches, zusammen mit der oben genannten Sprachbeschreibung, sollte der Leser in der Lage sein, verifi-zierende Algorithmen selbst zu realisieren und damit den Rechner als ein wirkliches Hilfsmittel des Mathematikers weiterzuentwickeln.

Ich möchte an dieser Stelle meinen Dank an U. Kulisch, W. Krämer und H. Berlejung aussprechen, die mit vielen Anregungen und Hilfen am Zustandekommen dieses Buches beteiligt waren.



# Inhaltsverzeichnis

<b>1</b>	<b>PASCAL-XSC im Internet</b>	<b>1</b>
<b>2</b>	<b>Vorteile eines BCD-Formats</b>	<b>3</b>
2.1	Binär-Format $\longleftrightarrow$ BCD-Format . . . . .	3
<b>3</b>	<b>BCD-Zahlenformate</b>	<b>11</b>
3.1	Der Datentyp REAL . . . . .	11
3.1.1	Mathematische Konstanten . . . . .	15
3.1.2	Mathematische Standardfunktionen . . . . .	15
3.1.3	Funktionen, Prozeduren und Operatoren . . . . .	17
3.2	Der Datentyp LONGREAL . . . . .	22
3.2.1	Standardoperatoren . . . . .	25
3.2.2	Mathematische Standardfunktionen . . . . .	26
3.2.3	Hilfsfunktionen und Operatoren . . . . .	28
3.2.4	Mathematische Konstanten . . . . .	33
3.2.5	Wertzuweisungen . . . . .	34
3.2.6	Hornerschema . . . . .	34
3.2.7	Suche eines Teilintervalls . . . . .	35
3.2.8	Einschließung eines exakten Funktionswertes . . . . .	36
3.3	Ein-/Ausgabeprozeduren . . . . .	37
3.4	Exakte Auswertung von Ausdrücken . . . . .	38
3.4.1	Operatoren für dotprecision-Ausdrücke . . . . .	42
3.5	<b>Der Datentyp rrG</b> . . . . .	46
3.5.1	Monadische Operatoren . . . . .	47
3.5.2	Addition, Subtraktion . . . . .	47
3.5.3	Multiplikation . . . . .	49
3.5.4	Division . . . . .	50
3.5.5	Vergleiche . . . . .	51
3.5.6	Transferfunktionen und Zuweisungsoperatoren . . . . .	53
3.5.7	Hilfsfunktionen . . . . .	53
3.5.8	Hornerschema . . . . .	55
3.5.9	Ein-/Ausgabeprozeduren . . . . .	56
3.5.10	Standardfunktionen . . . . .	58
3.5.11	Mathematische Konstanten . . . . .	64

3.6	<b>Der Datentyp <math>\text{IrG}</math></b>	66
3.6.1	Monadische Operatoren	67
3.6.2	Addition, Subtraktion	67
3.6.3	Multiplikation	68
3.6.4	Division	68
3.6.5	Vergleiche	69
3.6.6	Transferfunktionen und Zuweisungsoperatoren	69
3.6.7	Hilfsfunktionen	69
3.6.8	Hornerschema	70
3.6.9	Ein-/AusgabeprozEDUREN	71
3.6.10	Standardfunktionen	73
4	<b>Intervallrechnung</b>	<b>75</b>
4.1	Intervallrechnung mit dem Typ $\text{real}$	75
4.1.1	Der Datentyp $\text{interval}$	75
4.1.2	Operatoren	75
4.1.3	Transferfunktionen	77
4.1.4	Überladungen des Zuweisungsoperators	77
4.1.5	Standardfunktionen	77
4.1.6	Ein-/AusgabeprozEDUREN	79
4.2	Intervallrechnung mit dem Typ $\text{Longreal}$	81
4.2.1	Der Datentyp $\text{Linterval}$	81
4.2.2	Operatoren	81
4.2.3	Transferfunktionen	82
4.2.4	Standardfunktionen	83
4.2.5	Überladungen des Zuweisungsoperators	85
4.2.6	Ein-/AusgabeprozEDUREN	85
4.2.7	Exakte Auswertung von Intervallausdrücken	86
4.2.7.1	Operatoren	86
4.2.7.2	Transferfunktionen	89
4.2.7.3	Simulation von SUM-Ausdrücken	89
4.2.7.4	Beispiele	90
4.2.8	Wertebereich monotoner Funktionen	92
4.3	Intervallrechnung mit dem Typ $\text{rrG}$	93
4.3.1	Der Datentyp $\text{rrGinterval}$	93
4.3.2	Operatoren	93
4.3.3	Transferfunktionen	93
4.3.4	Standardfunktionen	94
4.3.5	Überladungen des Zuweisungsoperators	96
4.3.6	Ein-/AusgabeprozEDUREN	96
4.3.7	Wertebereich einer Funktion	97
4.3.7.1	Die Funktion ist monoton	97
4.3.7.2	Die Funktion besitzt ein relatives Maximum	98
4.3.7.3	Die Funktion besitzt ein relatives Minimum	100
4.3.7.4	Die Funktion besitzt ein relatives Maximum und Minimum	102

4.3.7.5	Die Funktion besitzt ein reallives Minimum und Maximum . . . . .	103
4.3.7.6	Die Funktion besitzt mehr als zwei Extrema . . . . .	105
<b>5</b>	<b>Komplexe Arithmetik</b>	<b>107</b>
5.1	Das Modul C_ARI . . . . .	107
5.1.1	Operatoren . . . . .	107
5.1.2	Transferfunktionen . . . . .	109
5.1.3	Überladungen des Zuweisungsoperators . . . . .	109
5.1.4	Standardfunktionen . . . . .	110
5.1.5	Ein-/AusgabeprozEDUREN . . . . .	111
5.2	Das Modul LC_ARI . . . . .	112
5.2.1	Der Typ Lcomplex . . . . .	112
5.2.1.1	Operatoren . . . . .	112
5.2.1.2	Transferfunktionen . . . . .	114
5.2.1.3	Überladungen des Zuweisungsoperators . . . . .	114
5.2.1.4	Standardfunktionen . . . . .	115
5.2.1.5	Ein-/AusgabeprozEDUREN . . . . .	116
5.2.2	Exakte Auswertung komplexer Ausdrücke . . . . .	116
5.2.2.1	Operatoren . . . . .	117
5.2.2.2	Transferfunktionen . . . . .	118
5.2.2.3	Simulation von SUM-Ausdrücken . . . . .	118
5.2.2.4	Beispiele . . . . .	121
5.2.3	Der Typ rrGcomplex . . . . .	123
5.2.3.1	Operatoren . . . . .	123
5.2.3.2	Transferfunktionen . . . . .	123
5.2.3.3	Standardfunktionen . . . . .	124
5.2.3.4	Überladungen des Zuweisungsoperators . . . . .	126
5.2.3.5	Ein-/AusgabeprozEDUREN . . . . .	126
<b>6</b>	<b>Komplexe Intervallarithmetik</b>	<b>127</b>
6.1	Das Modul CI_ARI . . . . .	127
6.1.1	Operatoren . . . . .	128
6.1.2	Transferfunktionen . . . . .	129
6.1.3	Überladungen des Zuweisungsoperators . . . . .	130
6.1.4	Standardfunktionen . . . . .	130
6.1.5	Ein-/Ausgabeanweisungen . . . . .	135
6.2	Das Modul LCI_ARI . . . . .	136
6.2.1	Der Typ Lcinterval . . . . .	136
6.2.1.1	Operatoren . . . . .	136
6.2.1.2	Transferfunktionen . . . . .	139
6.2.1.3	Überladungen des Zuweisungsoperators . . . . .	141
6.2.1.4	Standardfunktionen . . . . .	141
6.2.1.5	Ein-/Ausgabeanweisungen . . . . .	143
6.2.2	Exakte komplexe Intervallausdrücke . . . . .	144
6.2.2.1	Transferfunktionen . . . . .	144

6.2.2.2	Wertzuweisungen . . . . .	146
6.2.2.3	Operatoren . . . . .	146
6.2.2.4	Simulation von SUM-Ausdrücken . . . . .	149
6.2.2.5	Beispiele . . . . .	151
6.2.3	Komplexe Intervallrechnung mit dem Typ rrG . . . . .	153
6.2.3.1	Der Datentyp rrGcinterval . . . . .	153
6.2.3.2	Transferfunktionen . . . . .	153
6.2.3.3	Operatoren . . . . .	154
6.2.3.4	Überladungen des Zuweisungsoperators . . . . .	154
6.2.3.5	Standardfunktionen . . . . .	155
6.2.3.6	Wertebereich separierbarer Funktionen . . . . .	157
6.2.3.7	Ein-/Ausgabeprozeduren . . . . .	158
<b>7</b>	<b>Matrix/-Vektorarithmetik</b>	<b>159</b>
7.1	Übersicht . . . . .	159
7.2	Das Modul MV_ARI . . . . .	161
7.2.1	Datentypen . . . . .	161
7.2.2	Operatoren . . . . .	161
7.2.3	Überladungen des Zuweisungsoperators . . . . .	162
7.2.4	Standardfunktionen . . . . .	163
7.2.5	Ein-/Ausgabeprozeduren . . . . .	164
7.3	Das Modul LMV_ARI . . . . .	165
7.3.1	Datentypen . . . . .	165
7.3.2	Operatoren . . . . .	165
7.3.3	Überladungen des Zuweisungsoperators . . . . .	166
7.3.4	Standardfunktionen . . . . .	167
7.3.5	Ein-/Ausgabeprozeduren . . . . .	168
7.3.6	Exakte Auswertung von Ausdrücken . . . . .	168
7.4	Das Modul MVL_ARI . . . . .	172
7.4.1	Datentypen . . . . .	172
7.4.2	Operatoren . . . . .	172
7.4.3	Transferfunktionen . . . . .	174
7.4.4	Überladungen des Zuweisungsoperators . . . . .	175
7.4.5	Standardfunktionen . . . . .	175
7.4.6	Ein-/Ausgabeprozeduren . . . . .	177
7.5	Das Modul LMVL_ARI . . . . .	178
7.5.1	Datentypen . . . . .	178
7.5.2	Operatoren . . . . .	178
7.5.3	Transferfunktionen . . . . .	180
7.5.4	Überladungen des Zuweisungsoperators . . . . .	181
7.5.5	Standardfunktionen . . . . .	182
7.5.6	Ein-/Ausgabeprozeduren . . . . .	183
7.5.7	Exakte Auswertung von Ausdrücken . . . . .	183
7.6	Das Modul MVC_ARI . . . . .	186
7.6.1	Datentypen . . . . .	186
7.6.2	Operatoren . . . . .	186

7.6.3	Transferfunktionen . . . . .	188
7.6.4	Überladungen des Zuweisungsoperators . . . . .	188
7.6.5	Standardfunktionen . . . . .	189
7.6.6	Ein-/Ausgabeprozeduren . . . . .	190
7.7	Das Modul LMVC_ARI . . . . .	191
7.7.1	Datentypen . . . . .	191
7.7.2	Operatoren . . . . .	191
7.7.3	Transferfunktionen . . . . .	193
7.7.4	Überladungen des Zuweisungsoperators . . . . .	194
7.7.5	Standardfunktionen . . . . .	195
7.7.6	Ein-/Ausgabeprozeduren . . . . .	196
7.7.7	Exakte Auswertung von Ausdrücken . . . . .	197
7.8	Das Modul MVCL_ARI . . . . .	198
7.8.1	Datentypen . . . . .	198
7.8.2	Operatoren . . . . .	198
7.8.3	Transferfunktionen für komplexe Intervallvektoren . . . . .	202
7.8.4	Transferfunktionen für komplexe Intervallmatrizen . . . . .	203
7.8.5	Überladungen des Zuweisungsoperators . . . . .	204
7.8.6	Standardfunktionen . . . . .	205
7.8.7	Ein-/Ausgabeprozeduren . . . . .	206
7.9	Das Modul LMVCL_ARI . . . . .	207
7.9.1	Datentypen . . . . .	207
7.9.2	Operatoren . . . . .	207
7.9.3	Transferfunktionen für komplexe Intervallvektoren . . . . .	211
7.9.4	Transferfunktionen für komplexe Intervallmatrizen . . . . .	212
7.9.5	Überladungen des Zuweisungsoperators . . . . .	213
7.9.6	Standardfunktionen . . . . .	214
7.9.7	Ein-/Ausgabeprozeduren . . . . .	215
7.9.8	Exakte Auswertung von Ausdrücken . . . . .	216
<b>8</b>	<b>PASCAL-XSC Module</b>	<b>217</b>
8.1	Systemmodule . . . . .	217
8.1.1	Modulhierarchie . . . . .	217
8.1.2	Modulbeschreibungen . . . . .	219
8.1.2.1	Basismodul STDMOD . . . . .	219
8.1.2.2	TIMER . . . . .	219
8.1.2.3	X_REAL . . . . .	220
8.1.2.3.1	Konstanten . . . . .	221
8.1.2.3.2	Speicherformat für real/Longreal-Konstanten	221
8.1.2.3.3	Hexadezimale Ein-/Ausgabe . . . . .	222
8.1.2.3.4	Mantisse/Exponent . . . . .	223
8.1.2.3.5	Klassifizierung von real/Longreal-Werten . . . . .	223
8.1.2.3.6	IEEE Ausnahmebehandlungs-Routinen . . . . .	225
8.1.2.4	IOSTD, X_INTG, X_STRG . . . . .	225
8.1.2.5	LI_ARI . . . . .	225
8.1.2.6	I_ARI . . . . .	226

8.1.2.7	RRGI_ARI . . . . .	226
8.1.2.8	C_ARIAUX . . . . .	226
8.1.2.9	LC_HELP . . . . .	227
8.1.2.10	LC_ARI . . . . .	229
8.1.2.11	LCI_ARI . . . . .	229
8.1.2.12	C_ARI . . . . .	229
8.1.2.13	CI_ARI . . . . .	229
8.1.2.14	MV_ARI . . . . .	229
8.1.2.15	MVI_ARI . . . . .	229
8.1.2.16	MVC_ARI . . . . .	230
8.1.2.17	MVCI_ARI . . . . .	230
8.1.2.18	LMV_ARI . . . . .	230
8.1.2.19	LMVI_ARI . . . . .	230
8.1.2.20	LMVC_ARI . . . . .	230
8.1.2.21	LMVCI_ARI . . . . .	230
8.2	Hilfsmodule . . . . .	231
<b>Stichwortverzeichnis</b>		<b>234</b>

# Kapitel 1

## PASCAL-XSC im Internet

Da ein PASCAL-XSC-Compiler nach C übersetzt, kann ein PASCAL-XSC-System praktisch auf jedem Rechner installiert werden, der über einen C-Compiler verfügt.

Frei im Internet abrufbar sind derzeit PASCAL-XSC Compiler für DOS, OS/2 und LINUX. Die Software und weitere Informationen finden Sie auf der Homepage der Firma

Numerik Software GmbH

P.O. BOX 2232

D-76492 Baden-Baden

Germany

Tel.: +49 7221 949217

Fax: +49 7221 949218

email: 100145.2667@compuserve.com

www: <http://come.to/xsc-software>

## Auszug aus der Homepage der Firma Numerik Software GmbH

**The PASCAL-XSC package contains:**

The compiler from PASCAL-XSC to C, modules for interval arithmetic, complex arithmetic, complex interval arithmetic and the corresponding matrix/vector arithmetic modules and the PASCAL-XSC runtime library for the corresponding C compiler. Try to test PASCAL-XSC and download a DOS or OS/2 Public Domain version of the PASCAL-XSC compiler system running under EMX GNU-C emx0.9b or a Linux version running under GNU-C. You can also download a Toolbox for PASCAL-XSC or C-XSC including some numerical sample programs which are described in the book Numerical Toolbox for Verified Computing I resp. C++ Toolbox for Verified Computing.

**The C-XSC software package contains:**

Object modules and header files for interval arithmetic, complex arithmetic, complex interval arithmetic and the corresponding matrix/vector arithmetic modules, a staggered multiple precision arithmetic, application examples, the C-XSC runtime library for the corresponding C++ compiler and the sources from the Toolbox book (see CXSCTOOL95). Try to test the C-XSC class library and download a DOS Public Domain version of the C-XSC compiler system running under Borland C++ 4.0.

The corresponding language descriptions are not part of the software package and must be ordered separately. Schools and universities can have a discount of 40% for Unix based system packages. This rebate can be given only for educational use of the system. Beneficiaries are students, teachers and professors at schools, universities and comparable institutions. The prices are valid for one CPU of the corresponding target computer system. Licences can be granted for multiprocessor machines as well as for network systems. Campus licences are available on request. The systems can be adapted to further target machines on request. The license includes free updates for the first license year. We only charge for shipment costs.

Shipment for companies and public institutions inside Europe on account, outside Europe payment in advance, per C.O.D. or with credit cards.

Forwarding Expense: inside Europe: 30,00 DM, outside Europe: 50,00 DM.

Accepted credit cards: American Express, Master Card and Visa. For EC countries the VAT ID-number is necessary to allow VAT free invoices.



## Kapitel 2

# Vorteile eines BCD-Formats

### 2.1 Binär-Format $\longleftrightarrow$ BCD-Format

Die Vorteile eines BCD-Zahlenformats lassen sich am einfachsten durch die Nachteile z.B. eines internen Binärsystems beschreiben. Die folgende Abbildung 2.1 zeigt dazu das typische Schema der numerischen Datenverarbeitung mit einem binären Zahlenformat. Die ergebnisverifizierenden Algorithmen beziehen sich dabei auf die Binärversion von PASCAL-XSC:

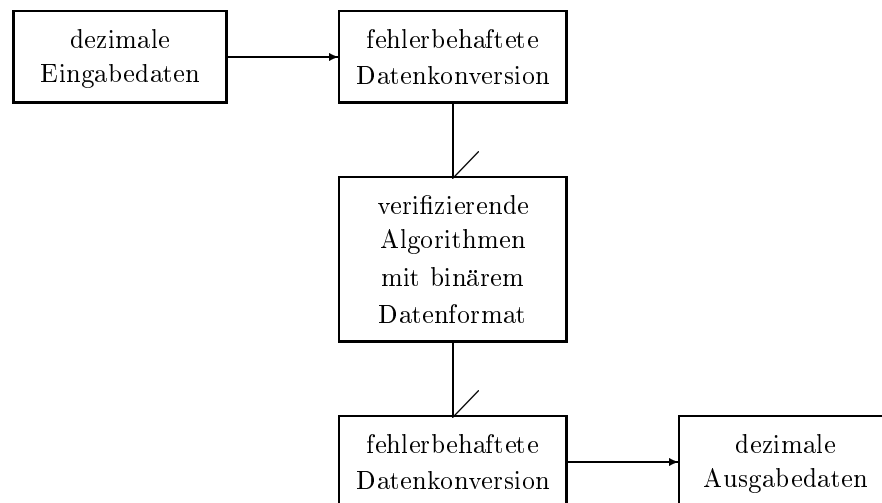


Abbildung 2.1: Datenverarbeitung mit binärem Zahlenformat

Daß die heutigen elektronischen Rechenanlagen intern ein duales Datenformat benutzen, hat rein historische Gründe, da jeder Rechner intern nur über **zwei** Grundzustände verfügt, die z.B. mit 0 oder 1 bezeichnet werden. Es war daher sehr

naheliegender, interne Zahlenformate mit der Basis  $2^n$  zu benutzen (Binär-, Oktal-, Hexadezimal-Systeme). Man braucht lediglich Konversionsroutinen, um die Eingabedaten aus dem uns gewohnten Zehnersystem z.B. ins Dual-System zu übertragen. Nach dem Programmablauf ist dann noch die Umwandlung der Ergebnisdaten in das von uns lesbare Zehnersystem notwendig. Zur Berechnung garantierter Ergebniseinschließungen muß man dabei in PASCAL-XSC natürlich entsprechend aufrunden, wenn z.B. eine Obergrenze ausgegeben werden soll!

Die Vorteile der internen Zahlensysteme mit der Basis  $2^n$  sind einmal die effiziente Speichernutzung und die kurzen Laufzeiten der Grundoperationen  $+$ ,  $-$ ,  $*$ ,  $/$ . Man könnte also der Auffassung sein, daß damit alle Grundvoraussetzungen erfüllt sind, einen Rechner in allen mathematischen Bereichen erfolgreich einzusetzen!

Leider ist aber die exakte Umwandlung der Dezimalzahlen in ein duales Zahlenformat endlicher Länge **nur in Ausnahmefällen** möglich, d.h. schon die Umwandlung von 0.1 ins Binärsystem erzeugt **zwangsläufig** einen Fehler, der zwar um so kleiner ist, je breiter die benutzte, interne Mantisse ist, der aber immer von Null verschieden sein wird!

Bei der Verarbeitung von Meßdaten, die grundsätzlich schon selbst mit Fehlern behaftet sind, haben die zusätzlichen Konversionsfehler keinen entscheidenden Einfluß auf die Ergebnisse. Die folgenden fünf mit der Binärversion von PASCAL-XSC durchgerechneten Beispiele zeigen jedoch, daß Konversionsfehler bei der Umwandlung von Dezimalzahlen ins Binär-System schon bei einfachen **mathematischen** Problemstellungen eine ganz entscheidende **negative** Rolle spielen können.

In den nachfolgenden Beispielen bedeutet  $\tilde{a}$  die zur vorgegebenen Dezimalzahl  $a$  nächstgelegene Maschinenzahl. Für  $a = 0.1$  gilt dann  $\tilde{a} \neq a$ , und für  $a = 8$  sind  $\tilde{a}$  und  $a$  identisch. Eingekreiste Operatorsymbole bedeuten Maschinenoperationen, die i.a. mit Rundungsfehlern behaftet sind, während nicht eingekreiste Operatoren stets exakte Ergebnisse erzeugen.

### 1. Beispiel

$a = 9.000000000000001$ ;  $b = 9$ ; exakter Wert:  $a - b = 10^{-15}$ ;  
 Das interne Rechnerergebnis  $\tilde{a} \ominus \tilde{b}$  der Binär-Version lautet jedoch  
 ins Zehnersystem zurückgerundet:  $1.776356839400250 \cdot 10^{-15}$ . Die  
 notwendige Zahlenkonversion zusammen mit der nachfolgenden Sub-  
 traktion erzeugt damit den sehr großen relativen Fehler von 0.776...!  
 Eine entsprechende intervallmäßige Auswertung liefert daher für den  
 exakten Wert  $10^{-15}$  die grobe Einschließung:

$$10^{-15} \in [0; 1.776356839400251 \cdot 10^{-15}].$$

### 2. Beispiel

$x_0 = 9.424777960769380 \approx 3 \cdot \pi$ ;  
 exakter Wert:  $\sin(x_0) = -2.846120698501614913474 \dots \cdot 10^{-16}$ ;  
 Das interne Rechnerergebnis  $\widetilde{\sin(\tilde{x}_0)}$  der Binär-Version lautet ins  
 Zehnersystem zurückgerundet:  $+3.673940397442059 \cdot 10^{-16}$ , d.h. der

Rechner liefert noch nicht einmal das richtige Vorzeichen des gesuchten Funktionswertes! Eine Intervallauswertung ergibt entsprechend die sehr grobe Einschließung:

$$\sin(x_0) \in [-1.408962799656045 \cdot 10^{-15}; +3.673940397442060 \cdot 10^{-16}].$$

Diese, durch Konversionsfehler bedingten Überschätzungen, findet man immer dann, wenn eine Funktion in der unmittelbaren Umgebung einer Nullstelle auszuwerten ist!

### 3. Beispiel

$a = b = 1.234567890123456$ ;  $c = 1.234567890123457$ ;  
 $d = 1.234567890123455$ ; exakter Wert:  $a \cdot b - c \cdot d = 10^{-30}$ ;  
 Das interne Skalarprodukt  $\#*(\tilde{a} \cdot \tilde{b} - \tilde{c} \cdot \tilde{d})$  lautet ins Zehnersystem zurückgerundet:  $2.741291394155932 \cdot 10^{-16}$ , d.h. im Vergleich zum exakten Wert  $10^{-30}$  erhalten wir mit dem binären Ergebnis eine Abweichung um 14 Größenordnungen, bedingt durch die anfänglichen Konversionen:  $a \rightarrow \tilde{a}$ ,  $b \rightarrow \tilde{b}$ , ...

Betrachtet man die große Abweichung vom exakten Ergebnis, so erinnert der obige Einsatz des langen Akkumulators deutlich an das Schießen mit Kanonen auf Spatzen, denn in der Tat liefert die Auswertung von  $\tilde{a} \odot \tilde{b} \ominus \tilde{c} \odot \tilde{d}$  den Wert  $2.22 \dots \cdot 10^{-16}$ , der sich vom Akkuergebnis nur unwesentlich unterscheidet!

Um eine **Einschließung** von  $a \cdot b - c \cdot d$  zu erhalten, müssen zunächst die vorgegebenen Dezimalzahlen  $a, b, c, d$  durch Binärintervalle  $ai, bi, ci, di$  kleinsten Durchmessers eingeschlossen werden. Die anschließende Auswertung des Ausdrucks  $\#\#(ai \cdot bi - ci \cdot di)$  liefert dann die für die Praxis viel zu grobe Einschließung:

$$a \cdot b - c \cdot d \in [-8.22 \dots \cdot 10^{-16}, +2.74 \dots \cdot 10^{-16}]$$

Das Beispiel zeigt weiter, welche Überschätzungen etwa beim Lösen fast singulärer, linearer Gleichungssysteme zu erwarten sind, wenn die vorgegebenen dezimalen Punktkoeffizienten zunächst in entsprechende Binärintervalle kleinsten Durchmessers eingeschlossen werden müssen. Wir betrachten dazu das

### 4. Beispiel

Zu lösen ist das lineare Gleichungssystem:  $\mathbf{A} \cdot \mathbf{x} = \mathbf{b}$  mit:

$$\mathbf{A} = \begin{pmatrix} c & a \\ b & d \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 0 \\ 1 \end{pmatrix},$$

wobei die Zahlen  $a, b, c, d$  die gleichen Werte haben wie im Beispiel 3. Die exakte Lösung lautet:

$$x_1 = \frac{a}{a \cdot b - c \cdot d} = +1.234567890123456 \cdot 10^{30}$$

$$x_2 = \frac{-c}{a \cdot b - c \cdot d} = -1.234567890123457 \cdot 10^{30}$$

Versucht man nun eine Lösung des Systems mit Hilfe der Prozedur LSS aus dem Modul ILSS der Binär-Version von PASCAL-XSC, so wird das Programm mit dem Errorcode 2 abgebrochen, da die interne Matrix  $\tilde{\mathbf{A}}$  für das System singulär erscheint. Die Koeffizienten  $a, b, c, d$  müssen nämlich vor Beginn der eigentlichen Rechnung in Binärintervalle eingeschlossen werden, deren Durchmesser alle von Null verschieden sind, so daß intern nicht mit  $\mathbf{A}$  sondern mit der Intervallmatrix  $\tilde{\mathbf{A}}$  gerechnet wird. Da  $\mathbf{A}$  wegen  $\text{Det}(\mathbf{A}) = -10^{-30}$  schon fast singulär ist, enthält  $\tilde{\mathbf{A}}$  wegen der Koeffizientenintervalle eine Punktmatrix, die tatsächlich singulär ist, so daß eine Lösung des Systems zwangsläufig unmöglich wird!

### 5. Beispiel

Mit Hilfe des langen Akkumulators und geeigneten Operatoren kann mit dem PASCAL-XSC-System ein Polynom **rundungsfehlerfrei** ausgewertet werden, wenn man **vorher** Argument und Polynomkoeffizienten in das jeweilige interne Zahlensystem des Rechners umgewandelt hat. Während die BCD-Version daher wegen fehlender Konversionsfehler **exakte** Polynomwerte liefern muß, sind Abweichungen der Ergebnisse bei der Binär-Version alleine auf die beschriebenen Konversionsfehler bei Argument und Polynomkoeffizienten zurückzuführen!

Mit dem ausmultiplizierten Polynom  $P(x) = \frac{1}{100} \cdot (3x - 1)^{13}$  und dem Argument  $x_1 = 0.3333333333333333$  erhält man die folgenden Ergebnisse:

$$\begin{array}{ll} \text{BCD-Version:} & P(x_1) = -1.000000000000000 \cdot 10^{-210} \quad \mathbf{exakt!} \\ \text{Binär-Version:} & \tilde{P}(x_1) = -1.002225955526144 \cdot 10^{-15} \end{array}$$

Beachten Sie bitte, daß in der Binärversion die Abweichung von **195 Größenordnungen** nur durch Konversionsfehler begründet ist, da die Polynomauswertung mit den binären Daten im langen Akkumulator rundungsfehlerfrei erfolgt! Weitere Informationen findet man dazu auf Seite 43.

Die obigen Beispiele zeigen recht eindrucksvoll den möglichen negativen Einfluß von Konversionsfehlern schon bei sehr einfachen numerischen Problemstellungen. Es soll an dieser Stelle aber ausdrücklich darauf hingewiesen werden, daß dieser Einfluß keine spezielle Eigenschaft etwa der Binärversion von PASCAL-XSC ist, denn man erhält z.B. mit TURBO-PASCAL und dem vergleichbaren Datentyp `double` für die Beispiele 1 bis 3 fast die gleichen Punktergebnisse wie mit dem PASCAL-XSC-System.

**TURBO-PASCAL-Ergebnisse mit dem Datentyp double:**

1. **Beispiel:**  $\tilde{a} \ominus \tilde{b} = 1.776 \dots \cdot 10^{-15};$
2. **Beispiel:**  $\widetilde{\sin}(x_0) = +3.6754 \dots \cdot 10^{-16};$
3. **Beispiel:**  $\tilde{a} \odot \tilde{b} \ominus \tilde{c} \odot \tilde{d} = +2.7408 \dots \cdot 10^{-16};$
4. **Beispiel:**  $\widetilde{x}_1 = -4.5043 \dots \cdot 10^{15}, \quad \widetilde{x}_2 = +4.5043 \dots \cdot 10^{15};$
5. **Beispiel:**  $\widetilde{P}(x_1) = -9.59645977591760 \cdot 10^{-17};$

Die TURBO-PASCAL-Lösungen der Beispiele 2 bis 5 sind im Vergleich zu den exakten Lösungen reine Phantasieergebnisse! Wertet man in Beispiel 4 die Ausdrücke für  $x_1, x_2$  zusätzlich mit dem genaueren Datentyp `extended` aus, so scheitert die TURBO-PASCAL-Berechnung völlig, da die Koeffizientendeterminante im Nenner an Stelle von  $-10^{-30}$  mit 0 ausgewertet wird! Im Beispiel 5 kann man mit TURBO-PASCAL natürlich kein sinnvolles Ergebnis erwarten, da neben den Konversionsfehlern noch zusätzliche Auslöschungseffekte beim Horner-Schema zu berücksichtigen sind, da ein `#` Konzept in dieser Sprachumgebung nicht zur Verfügung steht!

Vergleicht man die PASCAL-XSC-Ergebnisse mit den entsprechenden Ergebnissen von TURBO-PASCAL bzgl. der Beispiele 1 bis 5, so erkennt man noch zusätzlich den großen Vorteil des PASCAL-XSC-Systems:

Grobe Ergebniseinschließungen oder ein Programmabbruch mit Fehlermeldung des PASCAL-XSC-Systems weisen den Anwender darauf hin, daß der Algorithmus optimiert werden sollte, während die nur punktförmigen TURBO-PASCAL-Ergebnisse **keine Aussage** über die Zuverlässigkeit des Ergebnisses zulassen!

Zur **Vermeidung** der alleine durch Konversionsfehler bedingten Überschätzungen von Ergebnisintervallen gibt es grundsätzlich drei Möglichkeiten:

1. Man rechnet weiterhin mit der Binär-Version von PASCAL-XSC, jedoch mit Hilfe der Module `mp_ari`, `mpi_ari` in erhöhter Genauigkeit, wodurch die Laufzeiten jedoch drastisch anwachsen.
2. Man benutzt ein Algebra-System wie z.B. MATHEMATICA oder MAPLE. Die Nachteile sind hier jedoch die im Vergleich zu PASCAL-XSC schwerfälligen Programmier Techniken und vor allem die zu langen Laufzeiten.
3. Alle internen Rechnungen werden mit einer PASCAL-XSC-BCD-Version durchgeführt ( BCD = Binary Coded Decimal ), wodurch alle Datenkon-

versionen entfallen, vorausgesetzt, daß die Eingabedaten  $a, b, \dots$  nicht mehr Dezimalstellen enthalten, als das benutzte BCD-System zur Verfügung stellt.

Im Vergleich zu einer entsprechenden Binärversion ist ein BCD-System zwar merklich langsamer; wenn jedoch zur Vermeidung der beschriebenen Konversionsfehler mit den Modulen `mp_ari` bzw. `mpi_ari` in erhöhter Genauigkeit gerechnet wird oder wenn man ein Algebra-System benutzt, so hat man im Vergleich zum BCD-System doch merklich längere Laufzeiten und verliert mit einem Algebra-System zusätzlich die Mathematik bezogene einfache Programmierbarkeit eines PASCAL-XSC-Systems.

Zur Vermeidung der beschriebenen Konversionsfehler ist die BCD-Version von PASCAL-XSC daher ein geeignetes Hilfsmittel. Für die Beispiele 1,3,4 und 5 erhält man mit ihr die jeweils angegebenen **exakten Ergebnisse**, und im 2. Beispiel findet man für  $\sin(x_0)$  jetzt die **optimale Einschließung**:

$$\sin(x_0) \in [-2.84612069850161491348; -2.84612069850161491347] \cdot 10^{-16}$$

Die BCD-Version von PASCAL-XSC bietet für mathematische Anwendungen folgende Vorteile:

1. Verfügbarkeit des Systems für Personal Computer, Workstations, Großrechner und Supercomputer mittels einer Implementierung in C.
2. Vermeidung von Konversionsfehlern, wodurch viele mathematische Problemstellungen erst zufriedenstellend lösbar werden!
3. Mit ergebnisverifizierenden Algorithmen werden gesicherte, mathematische Aussagen möglich, z.B. über die Existenz und Eindeutigkeit von Lösungen linearer oder nichtlinearer Gleichungen oder Gleichungssysteme.
4. Algorithmen der Binär-Version lassen sich ohne größeren Aufwand auf die BCD-Version übertragen.
5. Das Modul- und Operator-Konzept zusammen mit den implementierten Überladungen von Wertzuweisungen, Operatoren, Prozeduren und Funktionen ermöglichen die Realisierung gut lesbarer PASCAL-XSC-Programme.
6. Mit dem zusätzlichen Datentyp Longreal lassen sich alle Rechnungen im Bedarfsfall auch in erhöhter Genauigkeit durchführen.
7. Für die Datentypen `real`, `Longreal`, `interval`, `Linterval`, `complex`, `Lcomplex`, `cinterval`, `Lcinterval` steht eine Vielzahl von hochgenauen Standardfunktionen mit garantierten Fehlerschranken zur Verfügung.
8. Spezielle Funktionen der mathematischen Physik werden zusätzlich in hochgenauer Darstellung durch Hilfsmodule bereitgestellt.
9. Zur Realisierung hochgenauer Standard- oder spezieller Funktionen steht ein doppeltlanges real-Format mit  $2 \cdot 13 = 26$  Mantissenstellen und ein nur durch die integer-Zahlen begrenzter Exponentenbereich zur Verfügung, zusammen

mit allen notwendigen Grundoperationen der Punkt- und Intervallrechnung für reelle und komplexe Zahlen. Zusätzlich existiert in diesem Zahlenformat eine Vielzahl von Standardfunktionen für reelle und komplexe Punkt- und Intervallrechnungen mit garantierten Fehlerschranken. Der Anwender kann Module mit Funktionen seiner eigenen Wahl erstellen und so das System beliebig erweitern.

10. Mit Hilfe der Toolbox-Bände I,II steht eine Vielzahl von Problemlösungen zur Verfügung, und an Hand der vielen Beispiele dieses Handbuches sollte der mit den Methoden des verifizierenden Rechnens vertraute Anwender in der Lage sein, Lösungsverfahren für eigene Aufgabenstellungen zu entwickeln.





# Kapitel 3

## BCD–Zahlenformate

### 3.1 Der Datentyp REAL

Gegeben sei die Dezimalzahl  $x = m \cdot 10^{ex}$ , dann gilt im Standard–Datenformat **REAL** für die normalisierten Zahlen:

$$m = 0 \quad \text{oder} \quad 0.1 \leq |m| \leq 0.9999999999999, \quad 13 \text{ BCD–Ziffern,}$$

wobei für den Zehnerexponenten  $ex$  folgender Bereich zur Verfügung steht:

$$-254 \leq ex \leq +256; \quad ex \in \{0, \pm 1, \pm 2, \dots\}$$

Bezeichnet man mit **NZ** das oben beschriebene normalisierte Zahlensystem, so gilt für alle positiven  $x \in \mathbf{NZ}$  die Beziehung:

$$1.000000000000 \cdot 10^{-255} \leq x \leq 9.999999999999 \cdot 10^{+255}, \quad \text{d.h.:$$

MIN\_REAL :=  $1.000000000000 \cdot 10^{-255}$  ist die kleinste positive Zahl in **NZ**  
MAXREAL :=  $9.999999999999 \cdot 10^{+255}$  ist die größte (positive) Zahl in **NZ**

Bedeutet für  $a, b \in \mathbf{NZ}$  z.B.  $a \oplus b$  das Maschinenergebnis der Addition, so gilt i.a.  $a \oplus b \neq a + b$ , d.h. die Grundoperationen  $+, -, *, /$  sind in **NZ nicht** abgeschlossen, so daß z.B.  $(a + b) \notin \mathbf{NZ}$  zur nächsten Rasterzahl zu runden ist. Für die dabei auftretenden relativen Fehler  $\varepsilon_G$  gelten die folgenden Aussagen:

- Unter der Voraussetzung

$$\text{MIN\_REAL} \leq |a \circ b| \leq \text{MAXREAL}, \quad \text{mit: } \circ \in \{+, -, *, /\}$$

gilt z.B. für den relativen Fehler  $\varepsilon_G := [a \oplus b - (a + b)] / (a + b)$  der Addition:

$$|\varepsilon_G| \leq 0.5 \cdot 10^{-12} = \varepsilon(G): \quad \text{Fehlerschranke der Grundoperationen.}$$

- Unter der Voraussetzung  $a, b \in \mathbf{NZ}$  und  $a \circ b = 0$  ist auch stets das Maschinenergebnis gleich Null. Das exakte Ergebnis stimmt also mit dem Maschinenergebnis überein, wobei aber der relative Fehler  $\varepsilon_G$  jetzt nicht definiert ist!
- Unter der Voraussetzung  $|a \circ b| > \text{MAXREAL}$ , mit:  $\circ \in \{+, -, *, /\}$  erfolgt eine Overflow-Meldung mit Programmabbruch.

Für den Datentyp **REAL** ist noch besonders zu beachten, daß es neben den normalisierten Zahlen auch noch die **denormalisierten** Zahlen gibt, die betragsmäßig kleiner als  $\text{MIN\_REAL} = 10^{-255}$  sind. Dazu betrachten wir die folgenden Beispiele:

$$\begin{aligned}
 1\text{E}-255 \oslash 3\text{E}+000 &= 3.333333333330\text{E}-256 \notin \mathbf{NZ} \\
 1\text{E}-255 \oslash 3\text{E}+001 &= 3.333333333300\text{E}-257 \notin \mathbf{NZ} \\
 1\text{E}-255 \oslash 3\text{E}+002 &= 3.333333333000\text{E}-258 \notin \mathbf{NZ} \\
 &\vdots \\
 1\text{E}-255 \oslash 3\text{E}+011 &= 3.000000000000\text{E}-267 \notin \mathbf{NZ} \\
 1\text{E}-255 \oslash 3\text{E}+012 &= 0.000000000000\text{E}+000 \in \mathbf{NZ}
 \end{aligned}$$

Man erkennt, daß mit kleiner werdendem Ergebnisexponenten die hinteren Mantissenstellen statt mit der Ziffer 3 nur noch mit Nullen besetzt werden, wodurch der relative Fehler  $\varepsilon_G$  jeweils um den Faktor 10 **anwächst!** Liegen die Ergebnisse  $x$  einer Operation daher im denormalisierten Bereich:

$$(3.1) \quad \text{MINREAL} = 10^{-267} \leq |x| < 10^{-255} = \text{MIN\_REAL},$$

wobei die Besetzung der Mantisse mit von Null verschiedenen Ziffern nach obigen Beispielen entsprechend eingeschränkt ist, so kann für diese Grundoperation die im **NZ**-System geltende relative Fehlerschranke  $\varepsilon(G) = 0.5 \cdot 10^{-12}$  nicht mehr garantiert werden!

Hat man also für einen Algorithmus eine Fehlerabschätzung auf der Basis der Fehlerschranke  $\varepsilon(G) = 0.5 \cdot 10^{-12}$  für die Grundoperationen vorgenommen, so dürfen keine Zwischen- oder Endergebnisse im denormalisierten Bereich (3.1) liegen, es sei denn, daß man sicher ist, daß die Fehlerschranken der Operationen mit Zahlen aus diesem kritischen Bereich nicht größer als  $\varepsilon(G)$  sind!

Dazu einige Beispiele:

- $1 \oplus 1\text{E}-267 = 1$ , d.h. der relative Fehler ist kleiner als  $\varepsilon(G) = 0.5 \cdot 10^{-12}$ .

Es ist zu beachten, daß im denormalisierten Bereich (3.1) und damit für alle Zahlen vom Typ REAL die Fehlerschranke für **Addition und Subtraktion** durch  $\varepsilon(G) = 0.5 \cdot 10^{-12}$  gegeben ist!

Im denormalisierten Bereich haben alle Rasterzahlen im Gegensatz zu den normalisierten Zahlen den konstanten Abstand  $10^{-267}$ , dessen Wert die kleinste positive Rasterzahl darstellt.

- $1\text{E}-255 \odot 3\text{E}+011 = 3.000000000000\text{E}-267$ , d.h. der relative Fehler 0.1 ist jetzt größer als  $\varepsilon(G)$ .
- $1.11\text{E}-255 \odot 3\text{E}-010 = 3.330000000000\text{E}-265$ , d.h. der relative Fehler ist 0 und damit kleiner als  $\varepsilon(G)$ .

Bei **Multiplikation** und **Division** kann der Betrag des zugehörigen relativen Fehlers die Fehlerschranke  $\varepsilon(G) = 0.5 \cdot 10^{-12}$  überschreiten, wenn das Ergebnis der Operation in den denormalisierten Bereich (3.1) fällt!

Es ist daher sinnvoll, für Multiplikation und Division Prozeduren und Funktionen bereitzustellen, die überprüfen, ob bei diesen Operationen die Fehlerschranken  $\varepsilon(G) = 0.5 \cdot 10^{-12}$  bzw.  $\varepsilon(G) = 1.0 \cdot 10^{-12}$  für die Rundung zur nächsten Gleitkommazahl bzw. für gerichtete Rundungen überschritten werden:

1. **function** MULNEXT\_TEST(x,y: real; var error: boolean): real;  
Rundet das exakte Produkt  $x \cdot y$  zur nächsten real-Zahl;  
error = TRUE  $\rightsquigarrow$  relativer Fehler  $> 0.5 \cdot 10^{-12}$ ;
2. **function** MULUP\_TEST(x,y: real; var error: boolean): real;  
Rundet das exakte Produkt  $x \cdot y$  zur nächstgrößeren real-Zahl;  
error = TRUE  $\rightsquigarrow$  relativer Fehler  $\geq 1 \cdot 10^{-12}$ ;
3. **function** MULDOWN\_TEST(x,y: real; var error: boolean): real;  
Rundet das exakte Produkt  $x \cdot y$  zur nächstkleineren real-Zahl;  
error = TRUE  $\rightsquigarrow$  relativer Fehler  $\geq 1 \cdot 10^{-12}$ ;
4. **function** PRODUCT(m1,m2: real; ex1,ex2: integer;  
var error: boolean): real;  
Zu gegebenen real-Mantissen m1,m2 und den zugehörigen Zehnerexponenten ex1,ex2 berechnet die Funktion PRODUCT das zur nächsten real-Zahl gerundete Produkt  $(m1 \cdot m2)^{ex1+ex2}$ ;  
Es muß gelten:  $0.1 \leq |m1|, |m2| < 1$  oder  $|m1|, |m2| = 0$ .  
error = TRUE  $\rightsquigarrow$  relativer Fehler  $> 0.5 \cdot 10^{-12}$ ;

5. **function** ABC(a,b,c: real; var error: boolean): real;  
 Zu gegebenen real-Zahlen a,b,c berechnet die Funktion ABC den Wert  $(a \odot b) \odot c$  ohne vorzeitigen Overflow oder Underflow.  
 error = TRUE  $\leadsto$  Im denormalisierten Bereich ist eine eigene Fehlerabschätzung notwendig!
6. **function** ABCD(a,b,c,d: real; var error: boolean): real;  
 Zu gegebenen real-Zahlen a,b,c,d berechnet die Funktion ABCD den Wert  $(a \odot b) \odot (c \odot d)$  ohne vorzeitigen Overflow oder Underflow.  
 error = TRUE  $\leadsto$  Im denormalisierten Bereich ist eine eigene Fehlerabschätzung notwendig!
7. **function** DIVNEXT\_TEST(x,y: real; var error: boolean): real;  
 Rundet den exakten Quotienten  $x/y$  zur nächsten real-Zahl;  
 error = TRUE  $\leadsto$  relativer Fehler  $> 0.5 \cdot 10^{-12}$ ;
8. **function** DIVUP\_TEST(x,y: real; var error: boolean): real;  
 Rundet den exakten Quotienten  $x/y$  zur nächstgrößeren real-Zahl;  
 error = TRUE  $\leadsto$  relativer Fehler  $\geq 1 \cdot 10^{-12}$ ;
9. **function** DIVDOWN\_TEST(x,y: real; var error: boolean): real;  
 Rundet den exakten Quotienten  $x/y$  zur nächstkleineren real-Zahl;  
 error = TRUE  $\leadsto$  relativer Fehler  $\geq 1 \cdot 10^{-12}$ ;
10. **function** QUOTIENT(m1,m2: real; ex1,ex2: integer;  
 var error: boolean): real;  
 Zu gegebenen real-Mantissen m1,m2 und den zugehörigen Zehnerexponenten ex1,ex2 berechnet die Funktion QUOTIENT den zur nächsten real-Zahl gerundeten Quotienten  $(m1/m2)^{ex1-ex2}$ ;  
 Es muß gelten:  $0.1 \leq |m1|, |m2| < 1$  oder  $|m1| = 0$ .  
 error = TRUE  $\leadsto$  relativer Fehler  $> 0.5 \cdot 10^{-12}$ ;
11. **function** ab\_DIV\_c(a,b,c: real; var error: boolean): real;  
 Zu gegebenen real-Zahlen a,b,c berechnet die Funktion ab\_DIV\_c den Wert  $(a \odot b) \odot c$  ohne vorzeitigen Overflow oder Underflow.  
 error = TRUE  $\leadsto$  Im denormalisierten Bereich ist eine eigene Fehlerabschätzung notwendig!
12. **function** a\_DIV\_bc(a,b,c: real; var error: boolean): real;  
 Zu gegebenen real-Zahlen a,b,c berechnet die Funktion a\_DIV\_bc den Wert  $a \odot (b \odot c)$  ohne vorzeitigen Overflow oder Underflow.  
 error = TRUE  $\leadsto$  Im denormalisierten Bereich ist eine eigene Fehlerabschätzung notwendig!
13. **function** ab\_DIV\_cd(a,b,c,d: real; var error: boolean): real;  
 Zu den vorgegebenen real-Zahlen a,b,c,d berechnet die Funktion ab\_DIV\_cd den Wert  $(a \odot b) \odot (c \odot d)$  ohne vorzeitigen Overflow oder Underflow.  
 error = TRUE  $\leadsto$  Im denormalisierten Bereich ist eine eigene Fehlerabschätzung notwendig!

### 3.1.1 Mathematische Konstanten

Dem Anwender stehen mehrere mathematische real-Konstanten zur Verfügung, die in Tabelle 3.1 zusammengestellt sind. Zusätzlich existiert im Longreal- und rrG-Format eine Vielzahl weiterer mathematischer Konstanten, die durch Aufruf entsprechender Funktionen benutzt werden können; vergl. dazu die Tabellen auf den Seiten 33,64;

Real – Konstanten			
Name	Wert	Name	Wert
MIN_EXP	−255	MIN_EXP_	−267
MIN_EXPL	−511	MIN_EXPL_	−531
MAX_EXP	+255	MAX_EXPL	+511
MANT_W	+13	MANT_WL	+21
TEST_26	−240	TEST_34	−232
PI_1	$3.141592653589 \cdot 10^{+00}$	PI_2	$7.932384626433 \cdot 10^{-13}$
PI_3	$8.327950288419 \cdot 10^{-26}$	PI_4	$7.169399375105 \cdot 10^{-39}$
PI_5	$8.209749445923 \cdot 10^{-52}$	PI_REST	$2.643383279503 \cdot 10^{-21}$
R_PI_R	$3.183098861838 \cdot 10^{-01}$	MIN_REAL	$1 \cdot 10^{-255}$

Tabelle 3.1: real-Konstanten.

#### Anmerkungen:

- Addiert man die Konstanten PI\_1 bis PI\_5 in den Akkumulator, so addiert man zu seinem Wert die  $5 \cdot 13 = 65$ -stellige, abgerundete Näherung für  $\pi$  mit der relativen Fehlerschranke  $2.4881 \cdot 10^{-66}$ .
- R\_PI\_R ist die zur nächsten real-Zahl aufgerundete Näherung für  $1/\pi$ .

### 3.1.2 Mathematische Standardfunktionen

Für das Argument  $x$  vom Typ **REAL** stehen mathematische Standardfunktionen zur Verfügung, die in Tabelle 3.2 zusammengestellt sind.

Alle mathematischen Funktionen aus Tabelle 3.2 sind mit Ausnahme von `power(x,n,rndmode)` **hochgenau**, d.h. das Maschinenergebnis unterscheidet sich vom optimal gerundeten Funktionsergebnis um maximal eine Einheit in der letzten Mantissenstelle. Diese Abweichung von 1 ulp tritt allerdings nur in **sehr seltenen Ausnahmefällen** auf!

Fallen Funktionswerte  $y := f(x)$  in den Bereich (3.1), d.h. gilt:

$$0 \leq |y| < 10^{-255} = \text{MIN\_REAL},$$

Funktion	Aufruf	Funktion	Aufruf
$ x $	<code>abs(x)</code>	$x^y$	<code>power(x, y)</code>
$x^2$	<code>sqr(x)</code>	$\sin(x)$	<code>sin(x)</code>
$x^2 - y^2$	<code>x2_y2(x, y)</code>	$\cos(x)$	<code>cos(x)</code>
$\sqrt{x}$	<code>sqrt(x), x \geq 0</code>	$\tan(x)$	<code>tan(x)</code>
$\sqrt{1+x} - 1$	<code>sqrt1pm1(x); x \geq -1</code>	$\cot(x)$	<code>cot(x)</code>
$\sqrt{1-x^2}$	<code>sqrt1mx2(x);  x  \leq 1</code>	$\sin(\pi x)$	<code>sinpi(x)</code>
$\sqrt{1+x^2}$	<code>sqrt1px2(x)</code>	$\cos(\pi x)$	<code>cospi(x)</code>
$\sqrt{x^2-1}$	<code>sqrtx2m1(x);  x  \geq 1</code>	$\tan(\pi x)$	<code>tanpi(x)</code>
$\sqrt{x^2+y^2}$	<code>sqrtx2y2(x, y)</code>	$\cot(\pi x)$	<code>cotpi(x)</code>
$e^x$	<code>exp(x)</code>	$\arcsin(x)$	<code>arcsin(x)</code>
$e^x - 1$	<code>expm1(x)</code>	$\arccos(x)$	<code>arccos(x)</code>
$e^{-x^2}$	<code>exp_x2(x)</code>	$\arctan(x)$	<code>arctan(x)</code>
$2^x$	<code>exp2(x)</code>	$\arctan(x/y)$	<code>arctan2(x, y)</code>
$10^x$	<code>exp10(x)</code>	$\operatorname{arccot}(x)$	<code>arccot(x)</code>
$\ln(x)$	<code>ln(x); x &gt; 0</code>	$\sinh(x)$	<code>sinh(x)</code>
$\ln(e \cdot x)$	<code>ln_ex(x)</code>	$\cosh(x)$	<code>cosh(x)</code>
$\ln(1+x)$	<code>ln1p(x); x &gt; -1</code>	$\tanh(x)$	<code>tanh(x)</code>
$0.5 \cdot \ln(x^2 + y^2)$	<code>ln_sqrtx2y2(x, y)</code>	$\operatorname{coth}(x)$	<code>coth(x)</code>
$0.5 \cdot \ln[(1+x)^2 + y^2]$	<code>ln_sqrt1px2y2(x, y)</code>	$\operatorname{arsinh}(x)$	<code>arsinh(x)</code>
$\log_a(x)$	<code>loga(a, x)</code>	$\operatorname{arcosh}(x)$	<code>arcosh(x)</code>
$\log_2(x)$	<code>log2(x); x &gt; 0</code>	$\operatorname{arcosh}(1+x)$	<code>arcosh1p(x)</code>
$\log_{10}(x)$	<code>log10(x); x &gt; 0</code>	$\operatorname{artanh}(x)$	<code>artanh(x)</code>
$\log_{10}(1+x)$	<code>log1p(x); x &gt; -1</code>	$\operatorname{arcoth}(x)$	<code>arcoth(x)</code>
$x^n$	<code>power(x, n, rndmode)</code>		

Tabelle 3.2: Hochgenaue Funktionen vom Typ real; a,x,y: real-Ausdrücke;

und werden hintere, von Null verschiedene Ziffern der Ergebnismantisse gleich Null gesetzt oder entsteht Underflow, so daß die Funktionsergebnisse dadurch nicht mehr hochgenau sind, so erfolgt die Fehlermeldung:

**rel. error exceeds error bound**

mit Programmabbruch.

**Anmerkungen zu den Standardfunktionen:**

1. Mit Ausnahme von `power(x,n,rndmode)` zur Berechnung von Unter- bzw. Oberschranken (`rndmode = -1` bzw. `+1`) für  $x^n, n = 0, 1, 2 \dots$  haben alle Funktionen aus Tabelle 3.2 die relative Fehlerschranke  $5.0000001 \cdot 10^{-12}$ , d.h. wenn etwa `exp(x)` das i.a. fehlerbehaftete Maschinenergebnis der Exponentialfunktion zum Argument  $x$ : real bedeutet, so gilt:

$$\exp(x) = e^x \cdot (1 + \varepsilon), \quad \text{mit } |\varepsilon| < 5.0000001 \cdot 10^{-12}$$

2. Sei  $x = 1.1111 \cdot 10^{-130}$ ; Während das Produkt  $x \odot x$  mit einem relativen Fehler  $> 0.5 \cdot 10^{-12}$  ohne Fehlermeldung in den Bereich (3.1) fällt, liefert  $y := \text{SQR}(x)$  eine Fehlermeldung, da die Quadratfunktion `SQR(x)` für obiges Argument  $x$  nicht mehr hochgenau ausgewertet werden kann.
3. Bei den trigonometrischen Funktionen `sin(x), cos(x), tan(x), cot(x)` ist das Argument  $x$  eingeschränkt auf den Bereich (Polstellen ausgenommen):

$$|x| < 10^8;$$

4. Bei den Funktionen `tan( $\pi \cdot x$ )` und `cot( $\pi \cdot x$ )` ist das Argument  $x$  eingeschränkt auf den Bereich (Polstellen ausgenommen):

$$|x| \leq 1.0737418235 \cdot 10^9;$$

5. Bei den Funktionen `sin( $\pi \cdot x$ )` und `cos( $\pi \cdot x$ )` ist das Argument  $x$  eingeschränkt auf den Bereich:

$$|x| \leq 2.147483647 \cdot 10^9;$$

6. Argumente außerhalb der drei obigen Bereiche erzeugen bei den entsprechenden Funktionen eine Fehlermeldung. Verletzungen des Definitionsbereiches, z.B. bei der  $\sqrt{x}$ -Funktion oder Argumente im Overflowbereich erzeugen ebenfalls Fehlermeldungen mit Programmabbruch.

**3.1.3 Funktionen, Prozeduren und Operatoren**

In der folgenden Auflistung werden alle diejenigen Funktionen, Prozeduren und Operatoren angegeben, die mit den Datentypen **integer,real** in Verbindung stehen.

Wie im Binär-System existieren auch in der PASCAL-XSC-BCD-Version die Vergleichsoperatoren `= <> < <= > >=` für Operanden vom Typ `real` und `integer`.

Für die Grundoperationen `+ - * /` mit Rundung der exakten Ergebnisse zur nächsten real-Zahl einschließlich gerichteter Rundungen stehen alle arithmetischen Operatoren für `real`- und `integer`-Operanden zur Verfügung und werden hier nicht weiter beschrieben.

1. **function** CEIL (b: real): real;  
 $y := \text{CEIL}(x)$  liefert die kleinste ganze real-Zahl  $y$ , die größer oder gleich  $x$  ist:  $x \leq y$ ;
2. **function** FLOOR (b: real): real;  
 $y := \text{FLOOR}(x)$  liefert die größte ganze real-Zahl  $y$ , die kleiner oder gleich  $x$  ist:  $y \leq x$ ;
3. **function** FRAC (b: real): real;  
 $y := \text{FRAC}(x)$  liefert den Nachkommawert von  $x$ ;  
 $y := \text{FRAC}(-12.012)$  liefert:  $y = -1.2\text{E}-002$ ;  
 $0 \leq |y| \leq 0.999999999999999$ ;
4. **function** ROUND (b: real): integer;  
 $n := \text{ROUND}(x)$  liefert die zu  $x$  nächstgelegene integer-Zahl  $n$ ;  
 $n := \text{ROUND}(3.5)$  liefert:  $y = 4$   
integer-Bereich:  $-2147483648 \leq n \leq +2147483647$ ;
5. **function** ROUND (b: real): real;  
 $y := \text{ROUND}(x)$  liefert die zu  $x$  nächstgelegene, ganze real-Zahl  $y$ ;  
 $y := \text{ROUND}(3.5)$  liefert:  $y = 4$   
 $y$ -Bereich:  $-2147483648 \leq y \leq +2147483647$ ;
6. **function** TRUNC (b: real): integer;  
 $n := \text{TRUNC}(x)$  liefert die integer-Zahl  $n$ , die durch Abschneiden der Nachkommastellen von  $x$  entsteht.  
integer-Bereich:  $-2147483648 \leq n \leq +2147483647$ ;
7. **function** TRUNC (b: real): real;  
 $y := \text{TRUNC}(x)$  liefert die ganze real-Zahl  $y$ , die durch Abschneiden der Nachkommastellen von  $x$  entsteht.  
 $y$ -Bereich:  $-2147483648 \leq y \leq +2147483647$ ;
8. **function** SIGN (b: real): integer;  
 $y := \text{SIGN}(x)$  liefert das Vorzeichen von  $x$ ;  $y \in \{-1, 0, +1\}$ ;
9. **function** MANT (b: real): real;  
 $m := \text{MANT}(x)$  liefert die vorzeichenbehaftete Mantisse von  $x$ ;  
es gilt:  $m = 0$  oder  $0.1 \leq |m| < 1$ ;
10. **function** EXPO (b: real): integer;  
 $ex := \text{EXPO}(x)$  liefert den Zehnerexponenten  $ex$  bzgl. der Darstellung:  $x = m \cdot 10^{ex}$ , wobei die Mantisse  $m$  folgende Bedingungen erfüllen muß:  $m = 0$  oder  $0.1 \leq |m| < 1$ ;  
Für  $x \neq 0$  gilt:  $-266 \leq ex \leq +256$ ;  
Im Fall  $x = 0$  gilt:  $ex = -2147483647$ ;



11. **function** EXPO\_0 (b: real): integer;  
 ex:=EXPO\_0(x) liefert den Zehnerexponenten ex bzgl. der Darstellung:  $x = m \cdot 10^{\text{ex}}$ , wobei die Mantisse m folgende Bedingungen erfüllen muß:  $m = 0$  oder  $0.1 \leq |m| < 1$ ;  
 Für  $x \neq 0$  gilt:  $-266 \leq \text{ex} \leq +256$ ;  
 Im Falle  $x = 0$  gilt:  $\text{ex} = 0$ , im Gegensatz zur EXPO-Funktion!
12. **function** EXPO\_ADD (a,b: integer): integer;  
 s := EXPO\_ADD(a,b) liefert die Summe  $s := a + b$ ;  
 integer-Unterlauf:  $s := -2147483647$ ;  
 integer-Überlauf: Fehlermeldung;
13. **procedure** INTEGER\_ADD (a,b: integer; var s,k: integer);  
 INTEGER\_ADD(a,b,s,k) liefert die Summe  $s := a + b$ ;  
 integer-Unterlauf:  $k := -1$ ;  
 integer-Überlauf:  $k := +1$ ;  
 $s := a + b$  wird exakt berechnet:  $k := 0$ ;
14. **function** COMP (m: real; ex: integer): real;  
 y:=COMP(m,ex) liefert die real-Zahl y, zusammengesetzt aus der Mantisse m und dem Zehnerexponenten ex bzgl. der Darstellung:  $x = m \cdot 10^{\text{ex}}$ , wobei die Mantisse m folgende Bedingungen erfüllen muß:  $m = 0$  oder  $0.1 \leq |m| < 1$ ;  
 Es muß gelten:  $-266 \leq \text{ex} \leq +256$ , und im Bereich  $-266 \leq \text{ex} \leq -255$  muß die Mantisse m, wie auf S. 12 beschrieben, zusätzlich eingeschränkt sein!  
**Beispiele:**  
 y:=COMP(0.1, -267)  $\longrightarrow$  y = 0; „Exponent too small“-Meldung; kein Programmabbruch!  
 y:=COMP(0.12...23, -255)  $\longrightarrow$  y = 1.23...20E - 256;  
 „Mantissa bits lost...“-Meldung; kein Programmabbruch!  
 y:=COMP(Long(0.1), +257)  $\longrightarrow$  „Overflow-Meldung“ mit Programmabbruch.
15. **function** PRED (b: real): real;  
 y:=PRED(x) liefert bzgl. x die nächstkleinere Maschinenzahl y.
16. **function** SUCC (b: real): real;  
 y:=SUCC(x) liefert bzgl. x die nächstgrößere Maschinenzahl y.
17. **function** RVAL (s: string): real;  
 y:=RVAL(s) wandelt den String s in eine real-Zahl y. Der String s muß eine Zeichenfolge enthalten, die eine real-Konstante darstellt. Dabei werden führende Leerzeichen überlesen. Ein eventuell verbleibender Reststring wird durch Runden zur nächsten real-Zahl berücksichtigt.

18. **function** RVAL (s: string; round: integer): real;  
 $y := \text{RVAL}(s, \text{round})$  wandelt den String  $s$  in eine real-Zahl  $y$ . Der String  $s$  muß eine Zeichenfolge enthalten, die eine real-Konstante darstellt. Dabei werden führende Leerzeichen überlesen. Ein eventuell verbleibender Reststring wird durch Runden gemäß  $\text{round}$  berücksichtigt:

$$\text{round} = \left\{ \begin{array}{ll} -1 & \text{nach unten gerundet} \\ 0 & \text{nächstliegend gerundet} \\ +1 & \text{nach oben gerundet} \end{array} \right\}$$

19. **function** IMAGE (s: real): string;  
 $s := \text{IMAGE}(y)$  wandelt die real-Zahl  $y$  in einen String  $s$  entsprechend der Standardausgabe von real-Werten.

**Beispiel:**

$s := \text{IMAGE}(-2.35, E - 7)$  liefert:  $s = -2.4E - 007$ ;

20. **function** IMAGE (s: real; width: integer): string;  
 $s := \text{IMAGE}(y, w)$  wandelt die real-Zahl  $y$  in einen String  $s$  aus mindestens  $w$  Zeichen entsprechend der Ausgabe von real-Werten, eventuell mit Auffüllen von Nullen am Mantissenende.

**Beispiele:**

$s := \text{IMAGE}(-2.35, E - 7, 0) \rightarrow s = -2.4E - 007$ ;

$s := \text{IMAGE}(-2.35, E - 7, 11) \rightarrow s = -2.350E - 007$ ;

21. **function** IMAGE (s: real; width, fracs: integer): string;  
 $s := \text{IMAGE}(y, w, f)$  wandelt die real-Zahl  $y$  in einen String  $s$  aus wenigstens  $w$  Zeichen (eventuell mit führenden Leerzeichen aufgefüllt) entsprechend der Ausgabe von real-Werten in Festkomma-darstellung mit  $f$  Nachkommastellen, wobei Nullen angehängt werden, wenn  $f$  unnötig groß ist.

**Beispiele:**

$s := \text{IMAGE}(-2.35, E - 3, 12, 7) \rightarrow s = \square\square - 0.0023500$ ;

$s := \text{IMAGE}(-2.35, E - 3, 7, 4) \rightarrow s = -0.0024$ ;

$s := \text{IMAGE}(-2.35, E - 3, 3, 4) \rightarrow s = -0.0024$ ;

22. **function** IMAGE (s: real; width, fracs, round: integer): string;  
 $s := \text{IMAGE}(y, w, f, \text{round})$  wandelt die real-Zahl  $y$  in einen String  $s$  aus wenigstens  $w$  Zeichen (eventuell mit führenden Leerzeichen aufgefüllt) entsprechend der Ausgabe von real-Werten in Festkomma-darstellung mit  $f$  Nachkommastellen, wobei Nullen angehängt werden, wenn  $f$  unnötig groß ist. Bei zu kleinem  $f$  wird entsprechend  $\text{round}$  gerundet, vgl. die entsprechende Funktion LVAL.

**Beispiele:**

$s := \text{IMAGE}(-2.35, E - 3, 7, 4, +0) \rightarrow s = -0.0024$ ;

$s := \text{IMAGE}(-2.35, E - 3, 3, 4, +1) \rightarrow s = -0.0023$ ;

$s := \text{IMAGE}(-2.35, E - 3, 7, 4, -1) \rightarrow s = -0.0024$ ;

23. **function** ENTIRE (x: real; var gerade: boolean): boolean;  
 bl:=ENTIRE(x, gerade) liefert an bl den Wert true, wenn die Real-Zahl x eine **ganze** Zahl ist, sonst erhält bl den Wert false.  
 Falls x ganzzahlig ist, erhält gerade den Wert true (false), wenn x zusätzlich gerade (ungerade) ist.  
 Ist x nicht ganzzahlig, so bleibt die Variable gerade bedeutungslos.
24. **procedure** MANT\_EXPO\_PRODUCT(F1,F2: real; var MP: real;  
 var exp: integer);  
 MANT\_EXPO\_PRODUCT liefert zu den real-Zahlen F1,F2 die zur nächsten real-Zahl gerundete Mantisse MP des exakten Produkts  $F1 \cdot F2$  und den zugehörigen Zehnerexponenten exp.  
 Im Fall  $F1 \cdot F2 = 0$  gilt  $MP = 0$  und  $exp = -maxint = -2147483647$ , und im Fall  $F1 \cdot F2 \neq 0$  gilt:  $0.1 \leq |MP| < 1$ .
25. **priority** M\_10HOCH = \*;  
**operator** M\_10HOCH (x: real; k: integer) rint: real;  
 $y := x$  M\_10HOCH  $k$  berechnet zur real-Zahl x das Produkt  $y = x \cdot 10^k$ ,  $k \in \{0, \pm 1, \pm 2, \dots\}$ ;  
 Im Fall  $10^{-267} \leq |x \cdot 10^k| < 10^{-255}$  wird ohne Fehlermeldung zur nächsten Rasterzahl gerundet, wenn das Ergebnis wegen einer zu langen Mantisse von x nicht mehr exakt dargestellt werden kann.  
 Im Fall  $0 < |x \cdot 10^k| < 10^{-267}$  wird das Ergebnis auf Null gesetzt.  
 Der Operand x darf auch vom Typ integer sein!

### 3.2 Der Datentyp LONGREAL

Durch das Modul `STDMOD` wird der Datentyp **LONGREAL** zur Verfügung gestellt. Gegeben sei die Dezimalzahl  $x = m \cdot 10^{ex}$ , dann gilt im **LONGREAL**-Datenformat für die normalisierten Zahlen:

$$m = 0 \quad \text{oder} \quad 0.1 \leq |m| \leq 0.99999999999999999999, \quad 21 \text{ BCD-Ziffern};$$

wobei für den Zehnerexponenten  $ex$  folgender Bereich zur Verfügung steht:

$$-510 \leq ex \leq +512; \quad ex \in \{0, \pm 1, \pm 2, \dots\}$$

Bezeichnet man mit **NZ** das oben beschriebene normalisierte Zahlensystem, so gilt für alle positiven  $x \in \mathbf{NZ}$  die Beziehung:

$$1.00000000000000000000 \cdot 10^{-511} \leq x \leq 9.99999999999999999999 \cdot 10^{+511}, \quad \text{d.h.:$$

MINLONGREAL :=  $1.0 \dots 0 \cdot 10^{-511}$  ist die kleinste positive Zahl in **NZ**  
 MAXLONGREAL :=  $9.9 \dots 9 \cdot 10^{+511}$  ist die größte (positive) Zahl in **NZ**

`idxminLongreal` (Longreal-Konstante) `idxmaxLongreal` (Longreal-Konstante)  
 Bedeutet für  $a, b \in \mathbf{NZ}$  z.B.  $a \oplus b$  das Maschinenergebnis der Addition, so gilt i.a.  $a \oplus b \neq a + b$ , d.h. die Grundoperationen  $+, -, *, /$  sind in **NZ** **nicht** abgeschlossen, so daß z.B.  $(a + b) \notin \mathbf{NZ}$  zur nächsten Rasterzahl zu runden ist. Für die dabei auftretenden relativen Fehler  $\varepsilon_G$  gelten die folgenden Aussagen:

- Unter der Voraussetzung

$$\text{MINLONGREAL} \leq |a \circ b| \leq \text{MAXLONGREL}, \quad \text{mit: } \circ \in \{+, -, *, /\}$$

gilt z.B. für den relativen Fehler  $\varepsilon_G := [a \oplus b - (a + b)] / (a + b)$  der Addition:

$$|\varepsilon_G| \leq 0.5 \cdot 10^{-20} = \varepsilon(G) : \text{ Fehlerschranke der Grundoperationen.}$$

- Unter der Voraussetzung  $a, b \in \mathbf{NZ}$  und  $a \circ b = 0$  ist auch stets das Maschinenergebnis gleich Null. Das exakte Ergebnis stimmt also mit dem Maschinenergebnis überein, wobei aber der relative Fehler  $\varepsilon_G$  jetzt nicht definiert ist!
- Unter der Voraussetzung

$$|a \circ b| > \text{MAXLONGREAL}, \quad \text{mit: } \circ \in \{+, -, *, /\}$$

erfolgt eine Overflow-Meldung mit Programmabbruch.

Für den Datentyp **LONGREAL** ist noch besonders zu beachten, daß es neben den normalisierten Zahlen auch noch die **denormalisierten** Zahlen gibt, die betragsmäßig kleiner als  $\text{MINLONGREAL} = 10^{-511}$  sind. Dazu betrachten wird die

folgenden Beispiele:

$$\begin{aligned}
 1\text{E}-511 \oslash 3\text{E}+000 &= 3.333333333333333330\text{E}-512 \notin \text{NZ} \\
 1\text{E}-511 \oslash 3\text{E}+001 &= 3.333333333333333300\text{E}-513 \notin \text{NZ} \\
 1\text{E}-511 \oslash 3\text{E}+002 &= 3.333333333333333000\text{E}-514 \notin \text{NZ} \\
 &\vdots \\
 1\text{E}-511 \oslash 3\text{E}+019 &= 3.000000000000000000\text{E}-531 \notin \text{NZ} \\
 1\text{E}-511 \oslash 3\text{E}+020 &= 0.000000000000000000\text{E}+000 \in \text{NZ}
 \end{aligned}$$

Man erkennt, daß mit kleiner werdendem Ergebnisexponenten die hinteren Mantissenstellen statt mit der Ziffer 3 nur noch mit Nullen besetzt werden, wodurch der relative Fehler  $\varepsilon_G$  jeweils um den Faktor 10 **anwächst!** Liegen die Ergebnisse  $x$  einer Operation daher im denormalisierten Bereich

$$(3.2) \quad \text{MIN\_LONGREAL} = 10^{-531} \leq |x| < 10^{-511} = \text{MINLONGREAL},$$

wobei die Besetzung der Mantisse mit von Null verschiedenen Ziffern nach obigen Beispielen entsprechend eingeschränkt ist, so kann für alle Grundoperationen die im **NZ**-System geltende Fehlerschranke  $\varepsilon(G) = 0.5 \cdot 10^{-20}$  nicht mehr garantiert werden!

Hat man also für einen Algorithmus eine Fehlerabschätzung auf der Basis der Fehlerschranke  $\varepsilon(G) = 0.5 \cdot 10^{-20}$  für die Grundoperationen vorgenommen, so dürfen keine Zwischen- oder Endergebnisse im denormalisierten Bereich (3.2) liegen, es sei denn, daß man sicher ist, daß die Fehlerschranken der Operationen mit Zahlen aus diesem kritischen Bereich nicht größer als  $\varepsilon(G)$  sind!

Dazu einige Beispiele:

- $1 \oplus 1\text{E}-531 = 1$ , d.h. der relative Fehler ist kleiner als  $\varepsilon(G) = 0.5 \cdot 10^{-20}$ .

Es ist zu beachten, daß im denormalisierten Bereich (3.2) und damit für alle Zahlen vom Typ **LONGREAL** die Fehlerschranke für **Addition und Subtraktion** durch  $\varepsilon(G) = 0.5 \cdot 10^{-20}$  gegeben ist! Im denormalisierten Bereich haben alle Rasterzahlen im Gegensatz zu den normalisierten Zahlen den konstanten Abstand  $10^{-531}$ , dessen Wert die kleinste positive Rasterzahl darstellt.

- $1\text{E}-511 \oslash 3\text{E}+019 = 3.000\dots000\text{E}-531$ , d.h. der relative Fehler 0.1 ist jetzt größer als  $\varepsilon(G)$ .
- $1.11\text{E}-511 \odot 3\text{E}-010 = 3.330000000000\text{E}-521$ , d.h. der relative Fehler ist 0 und damit kleiner als  $\varepsilon(G)$ .

Bei **Multiplikation** und **Division** kann der Betrag des zugehörigen relativen Fehlers die Fehlerschranke  $\varepsilon(G) = 0.5 \cdot 10^{-20}$  überschreiten, wenn das Ergebnis der Operation in den denormalisierten Bereich (3.2) fällt!

Es ist daher sinnvoll, für Multiplikation und Division Prozeduren und Funktionen bereitzustellen, die überprüfen, ob bei diesen Operationen die Fehlerschranken  $\varepsilon(G) = 0.5 \cdot 10^{-20}$  bzw.  $\varepsilon(G) = 1.0 \cdot 10^{-20}$  für die Rundung zur nächsten Gleitkommazahl bzw. für gerichtete Rundungen überschritten werden:

1. **function** MULNEXT\_TEST(x,y: Longreal; var error: boolean): Longreal;  
Rundet das exakte Produkt  $x \cdot y$  zur nächsten Longreal-Zahl;  
error = TRUE  $\rightsquigarrow$  relativer Fehler  $> 0.5 \cdot 10^{-20}$ ;
2. **function** MULUP\_TEST(x,y: Longreal; var error: boolean): Longreal;  
Rundet das exakte Produkt  $x \cdot y$  zur nächstgrößeren Longreal-Zahl;  
error = TRUE  $\rightsquigarrow$  relativer Fehler  $\geq 1 \cdot 10^{-20}$ ;
3. **function** MULDOWN\_TEST(x,y: Longreal; var error: boolean): Longreal;  
Rundet das exakte Produkt  $x \cdot y$  zur nächstkleineren Longreal-Zahl;  
error = TRUE  $\rightsquigarrow$  relativer Fehler  $\geq 1 \cdot 10^{-20}$ ;
4. **function** PRODUCT(m1,m2: Longreal; ex1,ex2: integer;  
var error: boolean): Longreal;  
Zu gegebenen Longreal-Mantissen m1,m2 und den zugehörigen Zehnerexponenten ex1,ex2 berechnet die Funktion PRODUCT das zur nächsten Longreal-Zahl gerundete Produkt  $(m1 \cdot m2)^{ex1+ex2}$ ;  
Es muß gelten:  $0.1 \leq |m1|, |m2| < 1$  oder  $|m1|, |m2| = 0$ .  
error = TRUE  $\rightsquigarrow$  relativer Fehler  $> 0.5 \cdot 10^{-20}$ ;
5. **function** ABC(a,b,c: Longreal; var error: boolean): Longreal;  
Zu gegebenen Longreal-Zahlen a,b,c berechnet die Funktion ABC den Wert  $(a \odot b) \odot c$  ohne vorzeitigen Overflow oder Underflow.  
error = TRUE  $\rightsquigarrow$  Im denormalisierten Bereich ist eine eigene Fehlerabschätzung notwendig!
6. **function** ABCD(a,b,c,d: Longreal; var error: boolean): Longreal;  
Zu gegebenen Longreal-Zahlen a,b,c,d berechnet die Funktion ABCD den Wert  $(a \odot b) \odot (c \odot d)$  ohne vorzeitigen Overflow oder Underflow.  
error = TRUE  $\rightsquigarrow$  Im denormalisierten Bereich ist eine eigene Fehlerabschätzung notwendig!
7. **function** DIVNEXT\_TEST(x,y: Longreal; var error: boolean): Longreal;  
Rundet den exakten Quotienten  $x/y$  zur nächsten Longreal-Zahl;  
error = TRUE  $\rightsquigarrow$  relativer Fehler  $> 0.5 \cdot 10^{-20}$ ;

8. **function** DIVUP\_TEST(x,y: Longreal; var error: boolean): Longreal;  
Rundet den exakten Quotienten  $x/y$  zur nächstgrößeren Longreal-Zahl;  
error = TRUE  $\rightsquigarrow$  relativer Fehler  $\geq 1 \cdot 10^{-20}$ ;
9. **function** DIVDOWN\_TEST(x,y: Longreal; var error: boolean): Longreal;  
Rundet den exakten Quotienten  $x/y$  zur nächstkleineren Longreal-Zahl;  
error = TRUE  $\rightsquigarrow$  relativer Fehler  $\geq 1 \cdot 10^{-20}$ ;
10. **function** QUOTIENT(m1,m2: Longreal; ex1,ex2: integer;  
var error: boolean): Longreal;  
Zu gegebenen Longreal-Mantissen m1,m2 und den zugehörigen Zehnerexponenten ex1,ex2 berechnet die Funktion QUOTIENT den zur nächsten Longreal-Zahl gerundeten Quotienten  $(m1/m2)^{ex1-ex2}$ ;  
Es muß gelten:  $0.1 \leq |m1|, |m2| < 1$  oder  $|m1| = 0$ .  
error = TRUE  $\rightsquigarrow$  relativer Fehler  $> 0.5 \cdot 10^{-20}$ ;
11. **function** ab\_DIV\_c(a,b,c: Longreal; var error: boolean): Longreal;  
Zu den vorgegebenen Longreal-Zahlen a,b,c berechnet die Funktion ab\_DIV\_c den Wert  $(a \odot b) \oslash c$  ohne vorzeitigen Overflow oder Underflow.  
error = TRUE  $\rightsquigarrow$  Im denormalisierten Bereich ist eine eigene Fehlerabschätzung notwendig!
12. **function** a\_DIV\_bc(a,b,c: Longreal; var error: boolean): Longreal;  
Zu den vorgegebenen Longreal-Zahlen a,b,c berechnet die Funktion a\_DIV\_bc den Wert  $a \oslash (b \odot c)$  ohne vorzeitigen Overflow oder Underflow.  
error = TRUE  $\rightsquigarrow$  Im denormalisierten Bereich ist eine eigene Fehlerabschätzung notwendig!
13. **function** ab\_DIV\_cd(a,b,c,d: Longreal; var error: boolean): Longreal;  
Zu den vorgegebenen Longreal-Zahlen a,b,c,d berechnet die Funktion ab\_DIV\_cd den Wert  $(a \odot b) \oslash (c \odot d)$  ohne vorzeitigen Overflow oder Underflow.  
error = TRUE  $\rightsquigarrow$  Im denormalisierten Bereich ist eine eigene Fehlerabschätzung notwendig!

### 3.2.1 Standardoperatoren

Für integer-, real- und Longreal-Operanden stehen alle in PASCAL-XSC üblichen Vergleichsoperatoren zur Verfügung:

=   <>   >   >=   <   <=

Der rechte oder linke Vergleichsoperand kann dabei wahlweise ein integer-, real- oder Longrealausdruck sein, d.h. möglich ist z.B.:

$3 + \sin(4.1) <= \text{LONG}(-4.33) - 2.1113\text{E}-13$

Für alle Kombinationen aus integer-, real- und Longreal-Operanden stehen die arithmetischen Operatoren

$+$ ,  $-$ ,  $*$ ,  $/$  mit den gerichteten Rundungen  $+>$ ,  $+<$ ,  $->$ ,  $-<$ ,  $*>$ ,  $*<$ ,  $/>$ ,  $/<$

zur Verfügung.

### 3.2.2 Mathematische Standardfunktionen

Für das Argument  $x$  vom Typ **LONGREAL** stehen hochgenaue mathematische Standardfunktionen zur Verfügung, die in nachfolgender Tabelle 3.3 zusammengestellt sind. Mit Ausnahme von `power(x,n,rndmode)` zur Berechnung von Unter- bzw. Obergrenzen (`rndmode = -1` bzw. `+1`) für  $x^n$ ,  $n = 0, 1, 2 \dots$  besitzen alle Funktionen die relative Fehlerschranke  $\varepsilon(f) = 5.0078 \cdot 10^{-21}$ .

#### Definition des relativen Fehlers:

Bedeutet  $\tilde{f}(x)$  das fehlerbehaftete Maschinenergebnis des exakten Funktionswertes  $f(x)$ , so gilt für den relativen Fehler  $\varepsilon_f$ :

$$\varepsilon_f := \frac{\tilde{f}(x) - f(x)}{f(x)}, \quad f(x) \neq 0; \quad |\varepsilon_f| \leq \varepsilon(f) : \text{Rel. Fehlerschranke für } f(x);$$

Im Fall  $f(x) = 0$  gilt für alle Standardfunktionen stets  $\tilde{f}(x) = 0$ , so daß obige Einschränkung  $f(x) \neq 0$  keine praktische Bedeutung hat!

Fallen Funktionswerte  $y := f(x)$  in den Bereich (3.2), d.h. gilt:

$$0 \leq |y| < 10^{-511} = \text{MINLONGREAL},$$

und werden hintere, von Null verschiedene Ziffern der Ergebnismantisse gleich Null gesetzt oder entsteht Underflow, so daß die angegebenen Fehlerschranken nicht mehr garantiert werden können, so erfolgt die Fehlermeldung:

**rel. error exceeds error bound**

mit Programmabbruch.



Funktion	Aufruf	Funktion	Aufruf
$ x $	<code>abs(x)</code>	$x^y$	<code>power(x, y)</code>
$x^2$	<code>sqr(x)</code>	$\sin(x)$	<code>sin(x)</code>
$x^2 - y^2$	<code>x2_y2(x, y)</code>	$\cos(x)$	<code>cos(x)</code>
$\sqrt{x}$	<code>sqrt(x), x ≥ 0</code>	$\tan(x)$	<code>tan(x)</code>
$\sqrt{1+x} - 1$	<code>sqrt1pm1(x); x ≥ -1</code>	$\cot(x)$	<code>cot(x)</code>
$\sqrt{1-x^2}$	<code>sqrt1mx2(x);  x  ≤ 1</code>	$\sin(\pi x)$	<code>sinpi(x)</code>
$\sqrt{1+x^2}$	<code>sqrt1px2(x)</code>	$\cos(\pi x)$	<code>cospi(x)</code>
$\sqrt{x^2-1}$	<code>sqrtx2m1(x);  x  ≥ 1</code>	$\tan(\pi x)$	<code>tanpi(x)</code>
$\sqrt{x^2+y^2}$	<code>sqrtx2y2(x, y)</code>	$\cot(\pi x)$	<code>cotpi(x)</code>
$e^x$	<code>exp(x)</code>	$\arcsin(x)$	<code>arcsin(x)</code>
$e^x - 1$	<code>expm1(x)</code>	$\arccos(x)$	<code>arccos(x)</code>
$e^{-x^2}$	<code>exp_x2(x)</code>	$\arctan(x)$	<code>arctan(x)</code>
$2^x$	<code>exp2(x)</code>	$\arctan(x/y)$	<code>arctan2(x, y)</code>
$10^x$	<code>exp10(x)</code>	$\operatorname{arccot}(x)$	<code>arccot(x)</code>
$\ln(x)$	<code>ln(x); x &gt; 0</code>	$\sinh(x)$	<code>sinh(x)</code>
$\ln(e \cdot x)$	<code>ln_ex(x)</code>	$\cosh(x)$	<code>cosh(x)</code>
$\ln(1+x)$	<code>ln1p(x); x &gt; -1</code>	$\tanh(x)$	<code>tanh(x)</code>
$0.5 \cdot \ln(x^2 + y^2)$	<code>ln_sqrtx2y2(x, y)</code>	$\operatorname{coth}(x)$	<code>coth(x)</code>
$0.5 \cdot \ln[(1+x)^2 + y^2]$	<code>ln_sqrt1px2y2(x, y)</code>	$\operatorname{arsinh}(x)$	<code>arsinh(x)</code>
$\log_a(x)$	<code>loga(a, x)</code>	$\operatorname{arcosh}(x)$	<code>arcosh(x)</code>
$\log_2(x)$	<code>log2(x); x &gt; 0</code>	$\operatorname{arcosh}(1+x)$	<code>arcosh1p(x)</code>
$\log_{10}(x)$	<code>log10(x); x &gt; 0</code>	$\operatorname{artanh}(x)$	<code>artanh(x)</code>
$\log_{10}(1+x)$	<code>log1p(x); x &gt; -1</code>	$\operatorname{arcoth}(x)$	<code>arcoth(x)</code>
$x^n$	<code>power(x, n, rndmode)</code>		

Tabelle 3.3: Hochgenaue Funktionen vom Typ Longreal; a,x,y: Longreal;

**Anmerkungen zu den Standardfunktionen:**

- Bei den trigonometrischen Funktionen  $\sin(x), \cos(x), \tan(x), \cot(x)$  ist das Argument  $x$  eingeschränkt auf den Bereich (Polstellen ausgenommen):

$$|x| < 10^{+8}$$

2. Bei den Funktionen  $\tan(\pi \cdot x)$  und  $\cot(\pi \cdot x)$  ist das Argument  $x$  eingeschränkt auf den Bereich (Polstellen ausgenommen):

$$|x| \leq 1.0737418235 \cdot 10^9;$$

3. Bei den Funktionen  $\sin(\pi \cdot x)$  und  $\cos(\pi \cdot x)$  ist das Argument  $x$  eingeschränkt auf den Bereich:

$$|x| \leq 2.147483647 \cdot 10^9;$$

4. Argumente außerhalb der drei obigen Bereiche erzeugen bei den entsprechenden Funktionen eine Fehlermeldung. Verletzungen des Definitionsbereiches, z.B. bei der  $\sqrt{x}$ -Funktion oder Argumente im Overflowbereich erzeugen ebenfalls Fehlermeldungen mit Programmabbruch.

### 3.2.3 Hilfsfunktionen und Operatoren

Die folgenden Funktionen stehen mit dem Datentyp **LONGREAL** in Verbindung:

1. **function CEIL** (b: Longreal): Longreal;  
 $y := \text{CEIL}(x)$  liefert die kleinste ganze Longreal-Zahl  $y$ , die größer oder gleich  $x$  ist:  $x \leq y$ ;
2. **function FLOOR** (b: Longreal): Longreal;  
 $y := \text{FLOOR}(x)$  liefert die größte ganze Longreal-Zahl  $y$ , die kleiner oder gleich  $x$  ist:  $y \leq x$ ;
3. **function FRAC** (b: Longreal): Longreal;  
 $y := \text{FRAC}(x)$  liefert den Nachkommawert von  $x$ ;  
 $y := \text{FRAC}(-12.012)$  liefert:  $y = -1.2\text{E}-002$ ;  
 $0 \leq |y| \leq 0.99999999999999999999$ ;
4. **function LONG** (b: real): Longreal;  
 $y := \text{LONG}(x)$  liefert die zu  $x$  wertgleiche Longreal-Zahl  $y$ ;
5. **function LONG** (b: integer): Longreal;  
 $y := \text{LONG}(n)$  liefert die zu  $n$  wertgleiche Longreal-Zahl  $y$ ;
6. **function ROUND** (b: Longreal): integer;  
 $n := \text{ROUND}(x)$  liefert die zu  $x$  nächstgelegene integer-Zahl  $n$ ;  
 $n := \text{ROUND}(\text{LONG}(3.5))$  liefert:  $y = 4$   
integer-Bereich:  $-2147483648 \leq n \leq +2147483647$ ;
7. **function ROUND** (b: Longreal): Longreal;  
 $y := \text{ROUND}(x)$  liefert die zu  $x$  nächstgelegene, ganze Longreal-Zahl  $y$ ;  
 $y := \text{ROUND}(\text{LONG}(3.5))$  liefert:  $y = 4$   
 $y$ -Bereich:  $-2147483648 \leq y \leq +2147483647$ ;

8. **function** SHORT (b: Longreal): real;  
 $y := \text{SHORT}(x)$  rundet  $x$  zur nächstgelegenen real-Zahl  $y$ ;
9. **function** SHORTDOWN (b: Longreal): real;  
 $y := \text{SHORTDOWN}(x)$  rundet  $x$  zur nächstkleineren real-Zahl  $y$ ;
10. **function** SHORTUP (b: Longreal): real;  
 $y := \text{SHORTUP}(x)$  rundet  $x$  zur nächstgrößeren real-Zahl  $y$ ;
11. **function** SHORT\_TEST (x: Longreal; var error: boolean): real;  
 $y := \text{SHORT\_TEST}(x, \text{error})$  rundet  $x$  zur nächsten real-Zahl  $y$ ;  
 $\text{error} = \text{TRUE} \iff \text{relativer Fehler} > 0.5 \cdot 10^{-12}$ ;
12. **function** SHORTUP\_TEST (x: Longreal; var error: boolean): real;  
 $y := \text{SHORTUP\_TEST}(x, \text{error})$  rundet  $x$  zur nächstgrößeren real-Zahl  $y$ ;  
 $\text{error} = \text{TRUE} \iff \text{relativer Fehler} \geq 1 \cdot 10^{-12}$ ;
13. **function** SHORTDOWN\_TEST (x: Longreal; var error: boolean): real;  
 $y := \text{SHORTDOWN\_TEST}(x, \text{error})$  rundet  $x$  zur nächstkleineren real-Zahl  $y$ ;  
 $\text{error} = \text{TRUE} \iff \text{relativer Fehler} \geq 1 \cdot 10^{-12}$ ;
14. **function** TRUNC (b: Longreal): integer;  
 $n := \text{TRUNC}(x)$  liefert die integer-Zahl  $n$ , die durch Abschneiden der Nachkommastellen von  $x$  entsteht.  
integer-Bereich:  $-2147483648 \leq n \leq +2147483647$ ;
15. **function** TRUNC (b: Longreal): Longreal;  
 $y := \text{TRUNC}(x)$  liefert die ganze Longreal-Zahl  $y$ , die durch Abschneiden der Nachkommastellen von  $x$  entsteht.  
 $y$ -Bereich:  $-2147483648 \leq y \leq +2147483647$ ;
16. **function** SIGN (b: Longreal): integer;  
 $y := \text{SIGN}(x)$  liefert das Vorzeichen von  $x$ ;  $y \in \{-1, 0, +1\}$ ;
17. **function** MANT (b: Longreal): Longreal;  
 $m := \text{MANT}(x)$  liefert die vorzeichenbehaftete Mantisse von  $x$ ;  
es gilt:  $m = 0$  oder  $0.1 \leq |m| < 1$ ;
18. **function** EXPO (b: Longreal): integer;  
 $\text{ex} := \text{EXPO}(x)$  liefert den Zehnerexponenten  $\text{ex}$  bzgl. der Darstellung:  $x = m \cdot 10^{\text{ex}}$ , wobei die Mantisse  $m$  folgende Bedingungen erfüllen muß:  $m = 0$  oder  $0.1 \leq |m| < 1$ ;  
Für  $x \neq 0$  gilt:  $-530 \leq \text{ex} \leq +512$ ;  
Im Fall  $x = 0$  gilt:  $\text{ex} = -2147483647$ ;

19. **function** EXPO\_0 (b: Longreal): integer;  
 ex:=EXPO\_0(x) liefert den Zehnerexponenten ex bzgl. der Darstellung:  $x = m \cdot 10^{\text{ex}}$ , wobei die Mantisse m folgende Bedingungen erfüllen muß:  $m = 0$  oder  $0.1 \leq |m| < 1$ ;  
 Es gilt:  $-530 \leq \text{ex} \leq +512$ ;  
 Im Fall  $x = 0$  gilt:  $\text{ex} = 0$ , im Gegensatz zur EXPO-Funktion!
20. **procedure** MANT\_EXPO\_PRODUCT(F1,F2: Longreal; var MP: Longreal; var exp: integer);  
 MANT\_EXPO\_PRODUCT liefert zu den Longreal-Zahlen F1,F2 die zur nächsten Longreal-Zahl gerundete Mantisse MP des exakten Produkts  $F1 \cdot F2$  und den zugehörigen Zehnerexponenten exp.  
 Im Fall  $F1 \cdot F2 = 0$  gilt  $\text{MP} = 0$  und  $\text{exp} = -\text{maxint} = -2147483647$ , und im Fall  $F1 \cdot F2 \neq 0$  gilt:  $0.1 \leq |\text{MP}| < 1$ .
21. **function** COMP (m: Longreal; ex: integer): Longreal;  
 y:=COMP(m,ex) liefert die Longreal-Zahl y, zusammengesetzt aus der Mantisse m und dem Zehnerexponenten ex bzgl. der Darstellung:  $x = m \cdot 10^{\text{ex}}$ , wobei die Mantisse m folgende Bedingungen erfüllen muß:  $m = 0$  oder  $0.1 \leq |m| < 1$ ;  
 Es muß gelten:  $-530 \leq \text{ex} \leq +512$ , und im Bereich  $-530 \leq \text{ex} \leq -511$  muß die Mantisse m, wie auf S. 23 beschrieben, zusätzlich eingeschränkt sein!  
**Beispiele:**  
 y:=COMP(Long(0.1), -531)  $\rightarrow y = 0$ ; „Exponent to small“-Meldung; kein Programmabbruch!  
 y:=COMP(Long(0.123...23) + Long(4.5678901 E - 14), -511)  $\rightarrow y = 1.23...900\text{E} - 512$ ; „Mantissa bits lost...“-Meldung; kein Programmabbruch!  
 y:=COMP(Long(0.1), +513)  $\rightarrow$  „Overflow-Meldung“ mit Programmabbruch.
22. **function** PRED (b: Longreal): Longreal;  
 y:=PRED(x) liefert bzgl. x die nächstkleinere MaschinenZahl  $y < x$
23. **function** SUCC (b: Longreal): Longreal;  
 y:=SUCC(x) liefert bzgl. x die nächstgrößere MaschinenZahl  $y > x$
24. **function** LVAL (s: string): Longreal;  
 y:=LVAL(s) wandelt den String s in eine Longreal-Zahl y. Der String s muß eine Zeichenfolge enthalten, die eine Longreal-Konstante darstellt. Dabei werden führende Leerzeichen überlesen. Ein eventuell verbleibender Reststring wird durch Runden zur nächsten Longreal-Zahl berücksichtigt.

25. **function** LVAL (s: string; round: integer): Longreal;  
 $y := \text{LVAL}(s, \text{round})$  wandelt den String  $s$  in eine Longreal-Zahl  $y$ . Der String  $s$  muß eine Zeichenfolge enthalten, die eine Longreal-Konstante darstellt. Dabei werden führende Leerzeichen überlesen. Ein eventuell verbleibender Reststring wird durch Runden gemäß  $\text{round}$  berücksichtigt:
- $$\text{round} = \left\{ \begin{array}{ll} -1 & \text{nach unten gerundet} \\ 0 & \text{nächstliegend gerundet} \\ +1 & \text{nach oben gerundet} \end{array} \right\}$$
26. **function** LVAL (s: string; round: integer; var rest: string): Longreal;  
 Wie obige Funktion; ein eventuell vorhandener Reststring wird an  $\text{rest}$  zurückgegeben.
27. **function** IMAGE (s: Longreal): string;  
 $s := \text{IMAGE}(y)$  wandelt die Longreal-Zahl  $y$  in einen String  $s$  entsprechend der Standardausgabe von Longreal-Werten.  
**Beispiel:**  
 $s := \text{IMAGE}(\text{LONG}(-2.35, E - 7))$  liefert:  $s = -2.4E - 007$ ;
28. **function** IMAGE (s: Longreal; width: integer): string;  
 $s := \text{IMAGE}(y, w)$  wandelt die Longreal-Zahl  $y$  in einen String  $s$  aus mindestens  $w$  Zeichen entsprechend der Ausgabe von Longreal-Werten, eventuell mit Auffüllen von Nullen am Mantissenende.  
**Beispiele:**  
 $s := \text{IMAGE}(\text{LONG}(-2.35, E - 7), 0) \rightarrow s = -2.4E - 007$ ;  
 $s := \text{IMAGE}(\text{LONG}(-2.35, E - 7), 11) \rightarrow s = -2.350E - 007$ ;
29. **function** IMAGE (s: Longreal; width, fracs: integer): string;  
 $s := \text{IMAGE}(y, w, f)$  wandelt die Longreal-Zahl  $y$  in einen String  $s$  aus wenigstens  $w$  Zeichen (eventuell mit führenden Leerzeichen aufgefüllt) entsprechend der Ausgabe von Longreal-Werten in Festkommadarstellung mit  $f$  Nachkommastellen, wobei Nullen angehängt werden, wenn  $f$  unnötig groß ist.  
**Beispiele:**  
 $s := \text{IMAGE}(\text{LONG}(-2.35, E - 3), 12, 7) \rightarrow s = \square\square - 0.0023500$ ;  
 $s := \text{IMAGE}(\text{LONG}(-2.35, E - 3), 7, 4) \rightarrow s = -0.0024$ ;  
 $s := \text{IMAGE}(\text{LONG}(-2.35, E - 3), 3, 4) \rightarrow s = -0.0024$ ;
30. **function** IMAGE (s: Longreal; width, fracs, round: integer): string;  
 $s := \text{IMAGE}(y, w, f, \text{round})$  wandelt die Longreal-Zahl  $y$  in einen String  $s$  aus wenigstens  $w$  Zeichen (eventuell mit führenden Leerzeichen aufgefüllt) entsprechend der Ausgabe von Longreal-Werten in Festkommadarstellung mit  $f$  Nachkommastellen, wobei Nullen angehängt werden, wenn  $f$  unnötig groß ist. Bei zu kleinem  $f$  wird entsprechend  $\text{round}$  gerundet, vgl die entsprechende Funktion LVAL.

**Beispiele:**

$s := \text{IMAGE}(\text{LONG}(-2.35, E - 3), 7, 4, 0) \longrightarrow s = -0.0024;$   
 $s := \text{IMAGE}(\text{LONG}(-2.35, E - 3), 3, 4, +1) \longrightarrow s = -0.0023;$   
 $s := \text{IMAGE}(\text{LONG}(-2.35, E - 3), 7, 4, -1) \longrightarrow s = -0.0024;$

31. **function** ENTIRE (x: Longreal; var gerade: boolean): boolean;  
 bl:=ENTIRE(x, gerade) liefert an bl den Wert true, wenn die Longreal-Zahl x eine **ganze** Zahl ist, sonst erhält bl den Wert false. Falls x ganzzahlig ist, erhält gerade den Wert true (false), wenn x zusätzlich gerade (ungerade) ist. Ist x nicht ganzzahlig, so bleibt die Variable gerade bedeutungslos.
32. **function** EM\_GANZ (x: Longreal): Longreal;  
 y:=EM\_GANZ(x) berechnet zur Longrealzahl x die zur Mantisse von x gehörige **ganze** Zahl y.

**Beispiele:**

$x = 12.3456789012345678901E - 55 \longrightarrow$   
 $y = 1.23456789012345678901E + 20;$   
 $x = -12.34 \longrightarrow y = -1234.0$

33. **procedure** ZERL (x: Longreal; var x1, x2: real);  
 ZERL (x, x1, x2) zerlegt die Longrealzahl x in zwei Summanden x1,x2 vom Typ real, mit:  $x = x1 + x2;$   
 Im Fall  $|x| > 1.9999999999999999E+256$  erfolgt eine Fehlermeldung mit Programmabbruch.  
 Im Fall  $|x| < 1.0E-267$  gilt:  $x1 = x2 = 0$  **ohne Fehlermeldung!!** Dies ist für viele Anwendungen nützlich!

**Beispiele:**

$x = 1.9999999999999999E+256 \longrightarrow$   
 $x1 = 9.9999999999999999E+255; \quad x2 = 9.9999999999999999E+255;$   
 $x = 0.9999E-267 \longrightarrow x1 = x2 = 0$  **ohne Fehlermeldung!!**  
 $x = 1.1E-267 \longrightarrow$  Fehlermeldung wegen zu breiter Mantisse von x, vergl. dazu die Beispiele auf Seite 12.  
 $x = 1.0E-267 \longrightarrow x1 = 1.0E-267, \quad x2 = 0;$

34. **priority** M\_10HOCH = \*;  
**operator** M\_10HOCH (x: Longreal; k: integer) lrint: Longreal;  
 $y := x \text{ M\_10HOCH } k$  berechnet zur Longreal-Zahl x das Produkt  $y = x \cdot 10^k, \quad k \in \{0, \pm 1, \pm 2, \dots\};$   
 Im Fall  $10^{-531} \leq |x \cdot 10^k| < 10^{-511}$  wird ohne Fehlermeldung zur nächsten Rasterzahl gerundet, wenn das Ergebnis wegen einer zu langen Mantisse von x nicht mehr exakt dargestellt werden kann.  
 Im Fall  $0 < |x \cdot 10^k| < 10^{-531}$  wird das Ergebnis auf Null gesetzt.

### 3.2.4 Mathematische Konstanten

Es existieren 29 mathematische Konstanten im Longreal-Format, die durch Aufruf entsprechender Funktionsnamen (ohne Argument) zur Verfügung stehen:

Konstanten im Longreal – Format			
Name	Wert	Funktionswert	Rundung
MAXLONGREAL		+9.9...9 · 10 <sup>+511</sup>	
MINLONGREAL		+1.0...0 · 10 <sup>-511</sup>	
MIN_LONGREAL		+1.0...0 · 10 <sup>-531</sup>	
PL_DOWN	$\pi$	+3.1415...23846	* PL_DOWN < $\pi$
PL_UP	$\pi$	+3.1415...23847	↑ PL_UP > $\pi$
PI2	$2\pi$	+6.2831...47693	* PI2 > $2\pi$
PL_HALF	$\pi/2$	+1.5707...61923	* PL_HALF < $\pi/2$
PL_HALF1	$\pi/2$	+1.5707...61924	↑ PL_HALF1 > $\pi/2$
PL_DIV4	$\pi/4$	+0.7853...09615	* PL_DIV4 < $\pi/4$
R_PI	$1/\pi$	+0.3183...71538	* R_PI > $1/\pi$
R_2PI	$1/(2\pi)$	+0.1591...35769	* R_2PI > $1/(2\pi)$
R_SQRTPI2	$2/\sqrt{\pi}$	+1.1283...57390	* R_SQRTPI2 > $2/\sqrt{\pi}$
SQRT_PI	$\sqrt{\pi}$	+1.7724...02730	* SQRT_PI > $\sqrt{\pi}$
R_SQRTPI	$1/\sqrt{\pi}$	+0.5641...86948	* R_SQRTPI < $1/\sqrt{\pi}$
SQRT_PL_D2	$\sqrt{\pi}/2$	+0.8862...13649	* SQRT_PL_D2 < $\sqrt{\pi}/2$
LN_PI	$\ln(\pi)$	+1.1447...17414	* LN_PI < $\ln(\pi)$
EUL	$e$	+2.7182...23536	* EUL < $e$
LN_10	$\ln(10)$	+2.3025...68402	* LN_10 > $\ln(10)$
R_LN_10	$1/\ln(10)$	+0.4342...27651	* R_LN_10 < $1/\ln(10)$
LN_2	$\ln(2)$	+0.6931...09417	* LN_2 < $\ln(2)$
R_LN_2	$1/\ln(2)$	+1.4426...40736	* R_LN_2 > $1/\ln(2)$
GAMMA_C	$C$	+0.5772...60607	* GAMMA_C > $C$
LN_GAM_C	$\ln(C)$	-0.5495...22338	* LN_GAM_C > $\ln(C)$
SQRT2	$\sqrt{2}$	+1.4142...04880	* SQRT2 < $\sqrt{2}$
SQRT3	$\sqrt{3}$	+1.7320...29353	* SQRT3 > $\sqrt{3}$
SQRT5	$\sqrt{5}$	+2.2360...69641	* SQRT5 > $\sqrt{5}$
SQRT10	$\sqrt{10}$	+3.1622...33200	* SQRT10 > $\sqrt{10}$
R_SQRT2	$1/\sqrt{2}$	+0.7071...24401	* R_SQRT2 > $1/\sqrt{2}$
R_SQRT2PI	$1/\sqrt{2\pi}$	+0.3989...77940	* R_SQRT2PI > $1/\sqrt{2\pi}$

\* : Rundung zur nächsten, ↑ : Rundung zur nächstgrößeren Longreal-Zahl.

### 3.2.5 Wertzuweisungen

Durch Überladen des Zuweisungsoperators `:=` sind folgende Wertzuweisungen möglich:

```
Longrealvariable := integer-Ausdruck;
Longrealvariable := real-Ausdruck;
Longrealvariable := Longreal-Ausdruck;
```

```
Realvariable := integer-Ausdruck;
Realvariable := real-Ausdruck;
Realvariable := Longreal-Ausdruck;
```

Wird ein Longreal-Ausdruck in einer real-Variablen gespeichert, so wird im Bedarfsfall zur nächsten real-Zahl gerundet.

Das folgende Beispiel zeigt die PASCAL-interne Realisierung der Wertzuweisung eines real-Ausdrucks an eine Longrealvariable; dabei ist die Benutzung von **var** vor dem linken Longreal-Operanden zwingend vorgeschrieben!

```
operator := (var x: Longreal; r: real);
begin x := LONG(r) end.
```

### 3.2.6 Hornerschema

Gegeben sei das Polynom:

$$P_{10}(x) = b[0] + \dots + b[5] \cdot x^5 + a[6] \cdot x^6 + \dots + a[10] \cdot x^{10},$$

das für das Longreal-Argument `x` nach dem Hornerschema auszuwerten ist. Zur Laufzeitverkürzung sollen die ersten 4 Schritte im groben real-Format durchgeführt werden. Im fünften, gemischten Hornerschritt wird die Multiplikation im real-Format und die Addition im Longrealformat ausgeführt. Die restlichen Horner-schritte werden dann vollständig im Longreal-Format berechnet.

Im folgenden PASCAL-XSC-Programm **HORNER** benutzen wir die im Sprachkern von **PASCAL-XSC** definierten Datentypen:

```
type rvector = dynamic array[*] of real;
type lvector = dynamic array[*] of Longreal;
```

Die Polynomauswertung nach dem Hornerschema erfolgt mit der Funktion:

```
function HORNER_L (a: rvector; b: lvector; x: Longreal): Longreal;
```

die vom Modul **STDMOD** exportiert wird.



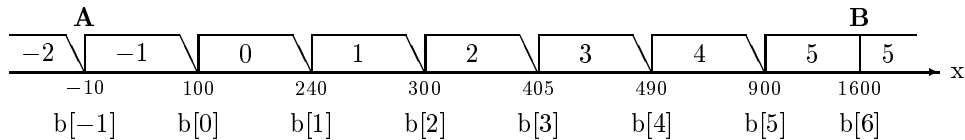
```

program HORNER(input,output);
var a : rvector[6..10];
      b : Lrvector[0..5];
      null : rvector[0..0];
      x, y, p10 : Longreal;
begin
  a[6] := 1.2; ... a[10] := 5.3E-6;
  b[0] := LONG(1) + LONG(1.2E-13); ...
  b[5] := LONG(4) + LONG(3.15E-13);
  x := 0.1;
  p10 := HORNER_L (a,b,x);
  y := HORNER_L (null,b,x);
end.

```

**Anmerkungen:**

1. Bei 4 groben Schritten benötigt man also  $4 + 1 = 5$  Komponenten  $a[k]$ .
2. Ist  $N$  der Polynomgrad und  $g \geq 1$  die Zahl der groben Schritte im real-Format, so muß gelten:  $1 \leq g \leq N - 2$  mit  $N \geq 3$ ;
3.  $p10$  ist eine Maschinennäherung für den exakten Wert  $P_{10}(0.1)$ ;
4.  $y := \text{HORNER\_L}(\text{null}, b, x)$ ; berechnet eine Maschinennäherung für das Polynom 5-ter Ordnung, dessen Longreal-Koeffizienten im  $b$ -Feld gespeichert sind, wenn **alle Hornerschnitte hochgenau** im Longrealformat ausgeführt werden.

**3.2.7 Suche eines Teilintervalls**Abbildung 3.1:  $[A, B]$  mit numerierten, halboffenen Teilintervallen.

Für die Auswertung einer reellwertigen Funktion über einem Intervall  $[A, B]$  ist es oft sinnvoll,  $[A, B]$  mittels eines  $b$ -Feldes vom Typ `rvector` (vergl. Seite 34) in Teilintervalle aufzuteilen und für jedes Teilintervall einen eigenen Algorithmus anzugeben. Dadurch entsteht die Aufgabe, zu gegebenem Longrealargument  $x \in [A, B]$  die entsprechend richtige Teilintervallnummer  $n$  zu bestimmen. In der Abb. 3.1 besitzt z.B. das halboffene Teilintervall  $[300; 405[$  die Nummer 2 entsprechend dem Indexwert 2 des linken Randpunktes  $b[2] = 300$ . Die gesuchte Teilintervallnummer  $n$  wird bestimmt mit Hilfe der Funktion:

function INT\_NR (b: rvector; x: Longreal): integer;

die vom Modul STDMOD exportiert wird. Das folgende Programm NUMMER zeigt eine Anwendung gemäß Abb. 3.1

```

program NUMMER(input,output);
var  b : rvector[-1..6];
      x : Longreal;
      n : integer;
begin
      b[-1] := -10;   ...   b[6] := 1600;
      x := -10.1;
      n := INT_NR (b,x);
end.

```

**Ergebnisse:**

n := INT_NR (b, LONG(1240.3))	liefert:	n = 5;
n := INT_NR (b, LONG(1601))	liefert:	n = 5;
n := INT_NR (b, LONG(1600))	liefert:	n = 5;
n := INT_NR (b, LONG(-10.1))	liefert:	n = -2;
n := INT_NR (b, LONG(-10))	liefert:	n = -1;
n := INT_NR (b, LONG(240))	liefert:	n = 1;

**Anmerkungen:**

1. Es muß gelten:  $A = b[-1]$  und  $B = b[6]$ ;
2. Durch geeignete Indexgrenzen des b-Feldes kann man sich einen entsprechenden Bereich für die Teilintervallnummern aussuchen!

### 3.2.8 Einschließung eines exakten Funktionswertes

Für einen reellwertigen Funktionswert  $f(x)$  sei eine Maschinennäherung  $\tilde{f}(x)$  vom Typ Longreal berechnet. Der relative Fehler  $\varepsilon_f$  ist definiert durch:

$$\varepsilon_f := \frac{\tilde{f}(x) - f(x)}{f(x)}; \quad f(x) \neq 0; \quad |\varepsilon_f| \leq \varepsilon(f);$$

Die Fehlerschranke  $\varepsilon(f)$  sei für alle  $x \in [a, b]$  bekannt. Gesucht sind eine garantierte Unterschranke US und eine Oberschranke OS für den exakten Funktionswert  $f(x)$ :

$$US \leq f(x) \leq OS; \quad US, OS = ?$$

Für  $\tilde{f}(x) > 0$  und  $\varepsilon(f) < 0.5 \cdot 10^{-3}$  gilt:

$$f(x) = \frac{\tilde{f}(x)}{1 + \varepsilon_f} \leq \frac{\tilde{f}(x)}{1 - \varepsilon(f)} < \tilde{f}(x) \cdot [1 + 1.001 \cdot \varepsilon(f)]$$

und mit  $\alpha := 1.001 \cdot \varepsilon(f)$  erhält man:  $f(x) < \tilde{f}(x) \cdot [1 + \alpha]$ . Setzt man  $\tilde{f}(x) = m \cdot 10^{ex}$ , so gilt:  $m := \text{MANT}(\tilde{f}(x)); \quad ex := \text{EXPO}(\tilde{f}(x))$ . Nach Seite 29 folgt:  $0.1 \leq m < 1; \implies m \cdot [1 + \alpha] = m + \alpha \cdot m < m + \alpha \leq m + \alpha \implies$

$$f(x) < 10^{ex} \cdot [m + \alpha] =: OS$$

Die Obergrenze OS kann daher **ohne zeitaufwendige Multiplikation** berechnet werden, wenn man  $[m + \alpha]$  und  $10^{e_x}$  mit dem schnellen Operator M\_10HOCH verknüpft!

Die Vermeidung der Multiplikation im Ausdruck  $\tilde{f}(x) \cdot [1 + \alpha]$  ergibt jedoch im ungünstigsten Fall eine Überschätzung der Obergrenze als wäre  $\varepsilon(f)$  verzehnfacht worden. Im Fall  $\tilde{f}(x) < 0$  können die Unter- und Obergrenzen US, OS ganz analog berechnet werden. Mit Hilfe der Funktion

function **BOUND** (y: Longreal; alpha: real): Longreal,

die vom Modul STDMOD exportiert wird, lassen sich US, OS bestimmen.

Sei  $y := \tilde{f}(x)$  und  $\alpha := 1.001 * \varepsilon(f)$ , dann liefert:

OS := BOUND (y, +alpha)      eine Obergrenze OS und  
 US := BOUND (y, -alpha)      eine Unterschranke US für  $f(x)$ .

#### Anmerkungen:

- Im Fall  $\tilde{f}(x) = 0$  gilt:  $US = -2 \cdot 10^{-531}$ ;  $OS = +2 \cdot 10^{-531}$ ; dadurch wird berücksichtigt, daß  $\tilde{f}(x) = 0$  durch Underflow entstanden sein kann!
- LONG sollte nicht mit  $\alpha = 0$  aufgerufen werden.

### 3.3 Ein-/AusgabeprozEDUREN

Mit Hilfe der bekannten Ein-/AusgabeprozEDUREN *read*, *readln*, *write*, *writeln* können Dateien (Filevariable) in bekannter Weise zur Ein- bzw. Ausgabe benutzt werden. Genaue Informationen findet man in: *PASCAL-XSC, Sprachbeschreibung mit Beispielen*.

Die Standardausgabe bei Verwendung der Prozedur *write* erfolgt für die Typen *integer*, *real*, *Longreal*, *char* und *boolean* entsprechend der Formatspezifikationstabelle 3.4.

Typ	<i>integer</i>	<i>real</i>	<i>Longreal</i>	<i>char</i>	<i>boolean</i>
Format	: 11	20 : 0 : 0	28 : 0 : 0	: 1	: 5

Tabelle 3.4: Standardausgabeformat für *write*.

#### Anmerkungen:

1. Mit *read*, *readln*, *write*, *writeln* ergeben sich **keine Konversionsfehler** wie sie z.B. bei einem Binär-System unvermeidlich sind!
2. Wird durch *read(x)* vom Textfile ein Wert gelesen, der betragsmäßig kleiner als  $10^{-267} (10^{-531})$  ist, so erhält die real-(Longreal)Variable den Wert 0.

Besitzt die Mantisse des Textfile-Wertes mehr Ziffern, als die interne Darstellung zuläßt, so wird zur nächstgelegenen Maschinenzahl gerundet.

3. Wird durch `read(x:r)` vom Textfile ein Wert gelesen, dessen Mantisse mehr Ziffern enthält, als die interne Darstellung zuläßt, so wird entsprechend `r` gerundet:

$$r = \left\{ \begin{array}{ll} -1 & \text{Rundung zur nächstkleineren Zahl} \\ 0 & \text{Rundung zur nächsten Rasterzahl} \\ +1 & \text{Rundung zur nächstgrößeren Zahl} \end{array} \right\}$$

`x` kann vom Typ `real` oder `Longreal` sein!

### 3.4 Exakte Auswertung von Ausdrücken

Um das Skalarprodukt

$$s := \sum_{k=1}^{10} a[k] * b[k]$$

mit nur einer einzigen Rundung berechnen zu können, gibt es in PASCAL-XSC das Konzept des **dotprecision-Ausdrucks** oder auch **Lattenkreuzausdrucks**. In der BCD-Version stehen solche **#-Ausdrücke** ebenfalls in vollem Umfang zur Verfügung; allerdings dürfen die `a[k]`, `b[k]` in #-Ausdrücken **nicht** vom Datentyp `Longreal` gewählt werden!

Mit Hilfe des Moduls `STDMOD` werden dem Anwender jedoch zusätzlich folgende Prozeduren und Operatoren zur Verfügung gestellt, mit denen unter gewissen Einschränkungen auch `Longreal`-Zahlen und exakte Produkte aus `Longreal`-Zahlen zu einer `dotprecision`-Variablen addiert (subtrahiert) werden können. Zusätzlich läßt sich der Akkumulator ins `real`- und `Longreal`-Format auslesen, wobei die Bereiche (3.1) und (3.2) entsprechend berücksichtigt werden können:

- operator `+` (d: `dotprecision`; a: `Longreal`) s : `dotprecision`;  
Zur `dotprecision`-Variablen `d` wird ein `Longreal`-Ausdruck addiert.
- operator `+` (a: `Longreal`; d: `dotprecision`) s : `dotprecision`;  
Zu einem `Longreal`-Ausdruck wird die `dotprecision`-Variable `d` addiert.
- operator `-` (d: `dotprecision`; a: `Longreal`) s : `dotprecision`;  
Von der `dotprecision`-Variablen `d` wird ein `Longreal`-Ausdruck subtrahiert.
- operator `-` (a: `Longreal`; d: `dotprecision`) s : `dotprecision`;  
Von einem `Longreal`-Ausdruck wird die `dotprecision`-Variable `d` subtrahiert.
- procedure **ADDACCU** (var d: `dotprecision`; a,b: `Longreal`);  
Zur `dotprecision`-Variablen `d` wird das **exakte** `Longreal`-Produkt `a · b` addiert.
- procedure **SUBACCU** (var d: `dotprecision`; a,b: `Longreal`);  
Von der `dotprecision`-Variablen `d` wird das **exakte** `Longreal`-Produkt `a · b` subtrahiert.

- procedure **ADDNACCU** (var d: dotprecision; a: Lrvector);  
Zur dotprecision-Variablen d werden alle Longrealzahlen des Feldes a vom Typ Lrvector (vergl. Seite 34) addiert.
- procedure **SUBNACCU** (var d: dotprecision; a: Lrvector);  
Von der dotprecision-Variablen d werden alle Longrealzahlen des Feldes a vom Typ Lrvector (vergl. Seite 34) subtrahiert.
- procedure **PADDNACCU** (var d: dotprecision; a,b: Lrvector);  
Zur dotprecision-Variablen d wird das aus den Feldern a,b gebildete **exakte** Skalarprodukt addiert.
- procedure **PSUBNACCU** (var d: dotprecision; a,b: Lrvector);  
Von der dotprecision-Variablen d wird das aus den Feldern a,b gebildete **exakte** Skalarprodukt subtrahiert.
- procedure **Akku\_times\_10power\_k** (var d: dotprecision; k: integer);  
Multipliziert den Akkuinhalt mit  $10^k$ ;  $k = 0, 1, 2, \dots$  ; negative k-Werte ändern den Akkuinhalt nicht.
- function **LONGNEXT** (var Akku: dotprecision): Longreal;  
 $y := \text{LONGNEXT}(\text{Akku})$  liest den Akku-Inhalt in die nächste Longreal-Zahl y;
- function **LONGDOWN** (var Akku: dotprecision): Longreal;  
 $y := \text{LONGDOWN}(\text{Akku})$  liest den Akku-Inhalt in die nächstkleinere Longreal-Zahl y;
- function **LONGUP** (var Akku: dotprecision): Longreal;  
 $y := \text{LONGUP}(\text{Akku})$  liest den Akkumulator-Inhalt in die nächstgrößere Longreal-Zahl y;
- function **LONGNEXT\_TEST** (var Akku: dotprecision;  
var error: boolean): Longreal;  
Liest den Akkumulator-Inhalt in die nächste Longreal-Zahl;  
 $\text{error} = \text{TRUE} \iff \text{relativer Fehler} > 0.5 \cdot 10^{-20}$ ; (durch Nullsetzen hinterer Ziffern der Ergebnismantisse im Bereich (3.2) ).
- function **LONGUP\_TEST** (var Akku: dotprecision;  
var error: boolean): Longreal;  
Liest den Akkumulator-Inhalt in die nächstgrößere Longreal-Zahl;  
 $\text{error} = \text{TRUE} \iff \text{relativer Fehler} \geq 1 \cdot 10^{-20}$ ; (durch Nullsetzen hinterer Ziffern der Ergebnismantisse im Bereich (3.2) ).
- function **LONGDOWN\_TEST** (var Akku: dotprecision;  
var error: boolean): Longreal;  
Liest den Akkumulator-Inhalt in die nächstkleinere Longreal-Zahl;  
 $\text{error} = \text{TRUE} \iff \text{relativer Fehler} \geq 1 \cdot 10^{-20}$ ; (durch Nullsetzen hinterer Ziffern der Ergebnismantisse im Bereich (3.2) ).

- function **REALNEXT\_TEST** (var Akku: dotprecision;  
var error: boolean): real;  
Liest den Akku-Inhalt in die nächste real-Zahl;  
error = TRUE  $\iff$  relativer Fehler  $> 0.5 \cdot 10^{-12}$ ; (durch Nullsetzen hinterer Ziffern der Ergebnismantisse im Bereich (3.1) ).
- function **REALUP\_TEST** (var Akku: dotprecision;  
var error: boolean): real;  
Liest den Akku-Inhalt in die nächstgrößere real-Zahl;  
error = TRUE  $\iff$  relativer Fehler  $\geq 1 \cdot 10^{-12}$ ; (durch Nullsetzen hinterer Ziffern der Ergebnismantisse im Bereich (3.1) ).
- function **REALDOWN\_TEST** (var Akku: dotprecision;  
var error: boolean): real;  
Liest den Akku-Inhalt in die nächstkleinere real-Zahl;  
error = TRUE  $\iff$  relativer Fehler  $\geq 1 \cdot 10^{-12}$ ; (durch Nullsetzen hinterer Ziffern der Ergebnismantisse im Bereich (3.1) ).
- function **SIGN** (Akku: dotprecision) : integer;  
Liefert das Vorzeichen des Akkumulators.

Da der lange Akkumulator ursprünglich nur für real-Zahlen konzipiert wurde, kann im Gegensatz zu einer Longreal-Zahl  $a$  ein exaktes Produkt  $a \cdot b$  aus den Longreal-Zahlen  $a, b$  nur unter bestimmten Einschränkungen in den Akku addiert werden.

Schreibt man z.B. die Longreal-Zahl  $a$  in der normalisierten Form

$$a = ma \cdot 10^{exa}, \quad \text{mit: } ma = 0 \quad \text{oder} \quad 0.1 \leq |ma| < 1,$$

so lautet die 1. Einschränkung:

$$exa + exb \leq 512$$

andernfalls erfolgt eine Fehlermeldung mit Programmabbruch.

Bedeutet z.B.  $Za$  die **Maximalzahl** der möglichen, **von Null verschiedenen** Mantissenziffern einer Longreal-Zahl  $a$ , so gelten für  $Za$  z.B. die Werte der nachfolgenden Tabelle :

$a$	$ma \cdot 10^{-530}$	$ma \cdot 10^{-529}$	$ma \cdot 10^{-511}$	$ma \cdot 10^{-510}$	$ma \cdot 10^{-509}$
$Za$	1	2	20	21	21

Tabelle 3.5:  $Za$  =Mantissenlänge einer Longreal-Zahl  $a$

und die zweite Einschränkung lautet:

$$exa + exb - Za - Zb \geq -552$$

andernfalls erfolgt eine Fehlermeldung mit Programmabbruch.  
Die zweite Einschränkung lässt sich auch wie folgt formulieren:

Für **beliebige**, zulässige Mantissen  $ma, mb$  darf das entsprechende exakte Produkt  $a \cdot b$  zweier Longrealzahlen in Festkommadarstellung höchstens bis zur 552. -ten Nachkommastelle reichen.

Da Verstöße gegen die beiden Einschränkungen durch Programmabbruch geahndet werden, sind Rechnerergebnisse mit den oben genannten Prozeduren und Operatoren absolut zuverlässig und ermöglichen z.B. den Aufbau einer hochgenauen, komplexen Longreal-Arithmetik.

Im Gegensatz zur Multiplikation sind **Addition** und **Subtraktion** von Longreal-Zahlen in den Akkumulator **ohne Einschränkungen** möglich!

#### Beispiele:

Vereinbarungen:                    var a,b : Lrvector[1..2];     d : dotprecision;  
Die Ergebnisse in Tabelle 3.7 entstanden durch die Anweisungen:

d := #(0);     PADDNACCU (d,a,b);

Nr.	a[1]	b[1]	a[2]	b[2]
1	1E-1	1E-531	0	0
2	1E-2	1E-531	0	0
3	1E-256	1E-256	0	0
4	1E-256	1E-257	0	0
5	1.234..8901E-256	a[1]	-1.234..8902E-256	1.234..8900E-256
6	1E+510	1	0	0
7	9.99..99E+510	9.99..99	0	0
8	1E+511	1	0	0
9	9.99..99E+510	9.99..99	1E+491	1
10	9.99..99E+510	9.99..99	2E+491	1

Tabelle 3.6: Komponenten vom Typ Longreal.

Nr.	$d := \#(0); \text{ paddnaccu}(d, a, b); \iff d = 0 + a[1] \cdot b[1] + a[2] \cdot b[2];$
1	1E-532 liefert Underflow beim Auslesen des Akkus.
2	Exponent range restricted ... ; Programmabbruch.
3	1E-512
4	Exponent range restricted ... ; Programmabbruch.
5	1E-552; kleinstmögliche, positive Zahl im Akkumulator.
6	1E+510
7	$9.99 \dots 998 \text{ } 00 \dots 001\text{E}+511 < 1\text{E}+512$
8	Invalid Operation... ; Programmabbruch.
9	$9.99 \dots 999 \text{ } 00 \dots 001\text{E}+511 < 1\text{E}+512$
10	$1.000 \dots (41\text{Ziffern } 0) \dots 0001\text{E}+512 > 1\text{E}+512$

Tabelle 3.7: Akkumulator-Inhalte.

### 3.4.1 Operatoren für dotprecision-Ausdrücke

Da ein #-Konzept für real- **und** Longreal-Ausdrücke nicht zur Verfügung steht, soll dieses mit Hilfe geeigneter Operatoren simuliert werden.

Um z.B. das **exakte** Produkt  $a \cdot b$  der Longreal-Zahlen  $a = b = 1.23456789012345678901 \cdot 10^{250}$  in die dotprecision-Variable Akku zu addieren, darf in der (erlaubten!) Anweisung

$$\text{Akku} := \text{Akku} + a * b;$$

der \* Operator **nicht** benutzt werden, da dieser das exakte Produkt  $a \cdot b$  **zuerst** ins Longrealformat rundet und **danach** das **gerundete** Longreal-Produkt in den Akku addiert! Will man daher das Problem ohne die Prozedur ADDACCU übersichtlicher mit einem **Operator** lösen, so muß man für den Multiplikationsoperator mit dem Ergebnistyp dotprecision einen **neuen** Namen wählen! In nachfolgender Tabelle 3.8 sind die entsprechenden neuen Operatornamen zusammengestellt:

Operation	Operator-Name	Ergebnistyp
Addition	<b>plus</b>	dotprecision
Subtraktion	<b>minus</b>	dotprecision
Multiplikation	<b>times</b>	dotprecision

Tabelle 3.8: Operatornamen mit Ergebnistyp dotprecision

In Tabelle 3.9 sind für die Operatoren **plus, minus, times, \*** alle erlaubten Operandentypen angegeben.



rechter Operand	integer	real	Longreal	dotprecision
linker Operand				
integer real	times plus, minus	times plus, minus	times plus, minus	* plus, minus
Longreal	times plus, minus	times plus, minus	times plus, minus	+, -, * plus, minus
dotprecision	* plus, minus	* plus, minus	+, -, * plus, minus	* plus, minus

Tabelle 3.9: Operatoren vom Ergebnistyp dotprecision

**1. Beispiel:**

Mit den Longreal-Zahlen

$a = b = 1.234567890123456$ ;  $c = 1.234567890123457$ ;  $d = 1.234567890123455$

liefern die Anweisungen:

```

yn := LONGNEXT ( a times b minus c times d );
yd := LONGDOWN ( a times b minus c times d );
yu := LONGUP ( a times b minus c times d );

```

jeweils den gleichen Longrealwert  $yn = yd = yu = 1 \cdot 10^{-30}$ , d.h. yn liefert sogar den **exakten** Wert von  $a \cdot b - c \cdot d$  (vergl. Beispiel 3. von Seite 5)

**2. Beispiel:**

Mit Hilfe der Operatoren `*` und `plus` zwischen *dotprecision*- und *Longreal*-Operanden kann ein Polynom

$$P(x) = \sum_{k=0}^n a_k \cdot x^k$$

mit der folgenden Funktion nach dem Horner Schema **rundungsfehlerfrei** berechnet werden:

```

function HORNER_dot (var a: Lrvector; xl: Longreal): dotprecision;
var k : integer;
    y : dotprecision;
begin
  y := a[ub(a)];
  for k := ub(a)-1 downto lb(a) do
    y := y * xl plus a[k];
  HORNER_dot := y
end;

```

Die Koeffizienten

$$\begin{array}{llll}
 a_0 = -0.01; & a_1 = +0.39; & a_2 = -7.02; & a_3 = +77.22; \\
 a_4 = -579.15; & a_5 = +3127.41; & a_6 = -12509.64; & a_7 = +37528.92; \\
 a_8 = -84440.07; & a_9 = +140733.45; & a_{10} = -168880.14; & a_{11} = +138174.66; \\
 a_{12} = -69087.33; & a_{13} = +15943.23; & & 
 \end{array}$$

realisieren das Polynom  $P(x) = \frac{1}{100} \cdot (3x - 1)^{13}$  mit der Nullstelle  $x_0 = \frac{1}{3}$ ;  
 Mit den Deklarationen:

```
var a: Lrvector[0..13];
    x1,y: Longreal; dot: dotprecision;
```

und den Anweisungen:

```
x1 := 0.3333333333333333; {x1 ≈ x0}
dot := HORNER_dot(a,x1);
y := LONGNEXT(dot);
```

erhält man mit der BCD-Version das **exakte Ergebnis**:

$$y = P(x1) = -1.00 \dots 00 \cdot 10^{-210};$$

Führt man ganz entsprechende Rechnungen mit der Binärversion von PASCAL-XSC durch, so erhält man:

$$\tilde{P}(\tilde{x1}) = -1.002225955526114 \cdot 10^{-15};$$

Die enorme Abweichung von **195 Größenordnungen** vom exakten BCD-Wert ist alleine bedingt durch den Umstand, daß **vor** der eigentlichen rundungsfehlerfreien Polynomauswertung neben dem Argument  $x1$  auch alle Polynomkoeffizienten mit Konversionsfehlern ins Binärsystem umgewandelt werden müssen.

In der Binärversion kann man bei diesem Beispiel die Konversionsfehler bei den  $a[k]$  vermeiden, indem man diese Polynomkoeffizienten mit dem Faktor 100 multipliziert und damit ganzzahlig macht:

$$Q(x) := 100 \cdot P(x) \equiv (3x - 1)^{13}$$

Der noch verbleibende Konversionsfehler beim Argument  $x1$  erzeugt dann mit der Binärversion im Vergleich zum exakten BCD-Wert immer noch einen Fehler von etwa vier Größenordnungen:

$$\begin{array}{ll}
 \text{Exakter BCD-Wert :} & Q(x1) = -1.0000000000000000 \cdot 10^{-208} \\
 \text{Binär-Ergebnis :} & \tilde{Q}(\tilde{x1}) = -4.752728915737900 \cdot 10^{-212}
 \end{array}$$

Das Beispiel zeigt, daß man trotz Vermeidung der Konversionsfehler bei den Polynomkoeffizienten (was zu vorzeitigem Overflow führen kann!) nur mit einer BCD-Version akzeptable Ergebnisse erhält, wenn man ein Polynom in der Nähe einer Nullstelle auswerten will.

**Vergleichsoperatoren:**

Bezüglich einer beliebigen Kombination von Operanden des Typs

**integer, real, Longreal, dotprecision**

existieren alle Vergleichsoperatoren:  $=$   $<>$   $>$   $>=$   $<$   $<=$  ;

**Wertzuweisungen:**

An eine Variable vom Typ dotprecision sind Wertzuweisungen von Ausdrücken folgender Typen möglich:

dotprecision-Variable := integer, real, Longreal, dotprecision-Ausdruck;

### 3.5 Der Datentyp rrG

Zur Implementierung **hochgenauer** Standardfunktionen im 21-stelligen Longreal-Format wurden die entsprechenden Algorithmen im **rrG**-Format realisiert, das wie folgt definiert

```
Type rrG = record
    r1,r2 : real;
    G      : integer;
end;
```

und durch das Standardmodul **STDMOD** zur Verfügung gestellt wird.

Durch `a : rrG` und

```
a.r1 := 1.234567890123; a.r2 := 4.567890123456E-13; a.G := +123456789;
```

wird eine BCD-Zahl mit mindestens 26 Mantissenstellen dargestellt, deren Exponentenbereich nur durch die integer-Zahlen begrenzt ist:

$$a = 1.2345678901234567890123456 \cdot 10^{+123456789} \\ -2147483648 \leq a.G \leq +2147483647$$

Für die korrekte Durchführung der arithmetischen Grundoperationen muß ein **rrG**-Operand `a` die beiden folgenden Bedingungen erfüllen:

$$a.r1 = 0 \implies a.r2 = 0 \quad (\text{Beding. I})$$

$$|a.r1| > 0 \implies \left| \frac{a.r2}{a.r1 + a.r2} \right| < 5.0001 \cdot 10^{-13} \quad (\text{Beding. II})$$

Die beiden obigen Bedingungen werden dabei automatisch realisiert, falls wie in allen praktischen Anwendungen `a.r1` und `a.r2` durch die folgenden Anweisungen aus dem Akkumulator ausgelesen werden:

```
a.r1 := #( Akku );          { a: rrG; Akku: dotprecision }
a.r2 := #( Akku - a.r1 );
```

Bedeutet `Akku` den numerischen Wert des Akkumulators, so gilt mit den obigen Anweisungen zusätzlich:

$$|a.r1| \geq 10^{-241}, \quad (a.r1 + a.r2) = \text{Akku} \cdot (1 + \varepsilon) \implies |\varepsilon| \leq 5 \cdot 10^{-26}$$

Das zweimalige Auslesen des Akkumulators erfolgt also mit dem relativen Höchstfehler  $\varepsilon(26) := 5 \cdot 10^{-26}$ , falls  $|a.r1| \geq 10^{-241}$ . Nach den beiden obigen **PASCAL-XSC**-Anweisungen erreicht man daher mit der Abfrage

```

if (SIGN(Akku) <> 0) and (EXP0(a.r1) < TEST_26)
  then Fehlermeldung

```

daß das zweimalige Auslesen des Akkumulators in die rrG-Komponenten a.r1 und a.r2 den relativen Fehler  $5 \cdot 10^{-26}$  nicht übersteigt; TEST\_26 = -240 wird vom Modul STDMOD als real-Konstante exportiert.

### 3.5.1 Monadische Operatoren

Für den Datentyp **rrG** stehen die monadischen Operatoren  $+$   $-$  in gewohnter Weise zur Verfügung.

### 3.5.2 Addition, Subtraktion

Mit Hilfe der folgenden Prozeduren und Operatoren lassen sich Variablen vom Typ **rrG** addieren bzw. subtrahieren. Falls beide Operanden den gleichen Zehnerexponenten besitzen, sind aus Laufzeitgründen die Prozeduren den Operatoren vorzuziehen. Es existieren auch die wichtigsten Operatoren mit gemischten Operanden.

- procedure **ADD\_11** (a,b: real; var r1,r2: real);  
 Es gilt:  $(a+b) \equiv (r1+r2)$  mit der relativen Fehlerschranke 0, falls kein Overflow auftritt.  $|r2/(r1 + r2)| < 5.0001 \cdot 10^{-13}$  wird erfüllt.
- procedure **ADD\_21** (a1,a2,b: real; var r1,r2: real);  
 $(r1 + r2) = (a1 + a2 + b) \cdot (1 + \varepsilon)$ ,  $|r1| \geq 10^{-241} \implies |\varepsilon| \leq 5 \cdot 10^{-26}$ ,  
 falls kein Overflow auftritt.  $|r2/(r1 + r2)| < 5.0001 \cdot 10^{-13}$  wird erfüllt.
- procedure **ADD\_22** (a1,a2,b1,b2: real; var r1,r2: real);  
 $(r1 + r2) = (a1 + a2 + b1 + b2) \cdot (1 + \varepsilon)$ ,  $|r1| \geq 10^{-241} \implies |\varepsilon| \leq 5 \cdot 10^{-26}$ ,  
 falls kein Overflow auftritt.  $|r2/(r1 + r2)| < 5.0001 \cdot 10^{-13}$  wird erfüllt.
- operator  $+$  (x,y: rrG) s: rrG;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $s = (x + y) \cdot (1 + \varepsilon)$   
 $\implies |\varepsilon| \leq 5.0001 \cdot 10^{-26}; |s.r2/(s.r1 + s.r2)| < 5.0001 \cdot 10^{-13}$
- operator  $+$  (x: real; y: rrG) s: rrG;  
 Vorauss.:  $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $s = (x + y) \cdot (1 + \varepsilon)$   
 $\implies |\varepsilon| \leq 5.0001 \cdot 10^{-26}; |s.r2/(s.r1 + s.r2)| < 5.0001 \cdot 10^{-13}$
- operator  $+$  (x: Longreal; y: rrG) s: rrG;  
 Vorauss.:  $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $s = (x + y) \cdot (1 + \varepsilon)$   
 $\implies |\varepsilon| \leq 5.0001 \cdot 10^{-26}; |s.r2/(s.r1 + s.r2)| < 5.0001 \cdot 10^{-13}$

- operator  $-$  (x,y: rrG) d: rrG;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $d = (x - y) \cdot (1 + \varepsilon)$   
 $\implies |\varepsilon| \leq 5.0001 \cdot 10^{-26}; \quad |d.r2/(d.r1 + d.r2)| < 5.0001 \cdot 10^{-13}$
- operator  $-$  (x: real; y: rrG) d: rrG;  
 Vorauss.:  $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $d = (x - y) \cdot (1 + \varepsilon)$   
 $\implies |\varepsilon| \leq 5.0001 \cdot 10^{-26}; \quad |d.r2/(d.r1 + d.r2)| < 5.0001 \cdot 10^{-13}$
- operator  $-$  (x: Longreal; y: rrG) d: rrG;  
 Vorauss.:  $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $d = (x - y) \cdot (1 + \varepsilon)$   
 $\implies |\varepsilon| \leq 5.0001 \cdot 10^{-26}; \quad |d.r2/(d.r1 + d.r2)| < 5.0001 \cdot 10^{-13}$
- operator  $-$  (y: rrG; x: real) d: rrG;  
 Vorauss.:  $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $d = (y - x) \cdot (1 + \varepsilon)$   
 $\implies |\varepsilon| \leq 5.0001 \cdot 10^{-26}; \quad |d.r2/(d.r1 + d.r2)| < 5.0001 \cdot 10^{-13}$
- operator  $-$  (y: rrG; x: Longreal) d: rrG;  
 Vorauss.:  $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $d = (y - x) \cdot (1 + \varepsilon)$   
 $\implies |\varepsilon| \leq 5.0001 \cdot 10^{-26}; \quad |d.r2/(d.r1 + d.r2)| < 5.0001 \cdot 10^{-13}$
- operator  $+>$  (x,y: rrG) s: rrG;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $\implies x + y < s; \quad |s.r2/(s.r1 + s.r2)| < 5.0001 \cdot 10^{-13}$
- operator  $->$  (x,y: rrG) d: rrG;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $\implies x - y < d; \quad |d.r2/(d.r1 + d.r2)| < 5.0001 \cdot 10^{-13}$
- operator  $+<$  (x,y: rrG) s: rrG;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $\implies x + y > s; \quad |s.r2/(s.r1 + s.r2)| < 5.0001 \cdot 10^{-13}$

- operator  $-<$  (x,y: rrG) d: rrG;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $\implies x - y > d; \quad |d.r2/(d.r1 + d.r2)| < 5.0001 \cdot 10^{-13}$
- procedure **ADD\_11G** (a,b: Longreal; var s: rrG);  
 $s = (a + b) \cdot (1 + \varepsilon) \implies |\varepsilon| \leq 5 \cdot 10^{-26}, |s.r2/(s.r1 + s.r2)| < 5.0001 \cdot 10^{-13}$   
 Für b kann auch ein real-Parameter eingesetzt werden.

### 3.5.3 Multiplikation

Mit Hilfe der folgenden Prozeduren und Operatoren lassen sich Variablen vom Typ rrG multiplizieren. Falls die Summe der Zehnerexponenten beider Operanden verschwindet, sind aus Laufzeitgründen die Prozeduren den Operatoren vorzuziehen. Es existieren auch die wichtigsten Operatoren mit gemischten Operanden.

- procedure **MUL\_11** (a,b: real; var r1,r2: real);  
 Es gilt:  $a \cdot b \equiv (r1 + r2)$  mit der relativen Fehlerschranke 0, falls  
 $|r1| \geq 10^{-241}, \quad |r2/(r1 + r2)| < 5.0001 \cdot 10^{-13}$  wird erfüllt.
- procedure **MUL\_12** (a,b1,b2: real; var r1,r2: real);  
 $(r1 + r2) = a \cdot (b1 + b2) \cdot (1 + \varepsilon), \quad |r1| \geq 10^{-241}$   
 $\implies |\varepsilon| \leq 5 \cdot 10^{-26}, \quad |r2/(r1 + r2)| < 5.0001 \cdot 10^{-13}$  wird erfüllt.
- procedure **MUL\_22** (a1,a2,b1,b2: real; var r1,r2: real);  
 Vorauss.:  $|a2/(a1 + a2)| < 5.0001 \cdot 10^{-13}$   
 $|b2/(b1 + b2)| < 5.0001 \cdot 10^{-13}$   
 $(r1 + r2) = (a1 + a2) \cdot (b1 + b2) \cdot (1 + \varepsilon), \quad |r1| \geq 10^{-241}$   
 $\implies |\varepsilon| \leq 3.0003 \cdot 10^{-25}; \quad |r2/(r1 + r2)| < 5.0001 \cdot 10^{-13}$
- procedure **MUL\_21** (a1,a2: real; xl: Longreal; var r1,r2: real);  
 Vorauss.:  $|a2/(a1 + a2)| < 5.0001 \cdot 10^{-13}, \quad |xl| \geq 10^{-247}$   
 $(r1 + r2) = (a1 + a2) \cdot xl \cdot (1 + \varepsilon), \quad |r1| \geq 10^{-241}$   
 $\implies |\varepsilon| \leq 3.0003 \cdot 10^{-25}; \quad |r2/(r1 + r2)| < 5.0001 \cdot 10^{-13}$
- procedure **MUL\_22S** (a1,a2,b1,b2: real; var r1,r2: real);  
 Vorauss.:  $|a2/(a1 + a2)| < 5.0001 \cdot 10^{-13}$   
 $|b2/(b1 + b2)| < 5.0001 \cdot 10^{-13}$   
 $(r1 + r2) = (a1 + a2) \cdot (b1 + b2) \cdot (1 + \varepsilon), \quad |r1| \geq 10^{-241}$   
 $\implies |\varepsilon| \leq 5 \cdot 10^{-26}; \quad |r2/(r1 + r2)| < 5.0001 \cdot 10^{-13}$
- operator  $*$  (x,y: rrG) p: rrG;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $p = (x \cdot y) \cdot (1 + \varepsilon)$   
 $\implies |\varepsilon| \leq 5.0001 \cdot 10^{-26}; \quad |p.r2/(p.r1 + p.r2)| < 5.0001 \cdot 10^{-13}$

- operator  $*$  (x: real; y: rrG) p: rrG;  
 Vorauss.:  $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $p = (x \cdot y) \cdot (1 + \varepsilon)$   
 $\implies |\varepsilon| \leq 5.0001 \cdot 10^{-26}; \quad |p.r2/(p.r1 + p.r2)| < 5.0001 \cdot 10^{-13}$
- operator  $*$  (x: Longreal; y: rrG) p: rrG;  
 Vorauss.:  $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $p = (x \cdot y) \cdot (1 + \varepsilon)$   
 $\implies |\varepsilon| \leq 5.0001 \cdot 10^{-26}; \quad |p.r2/(p.r1 + p.r2)| < 5.0001 \cdot 10^{-13}$
- operator  $*>$  (x,y: rrG) p: rrG;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $\implies x \cdot y < p; \quad |p.r2/(p.r1 + p.r2)| < 5.0001 \cdot 10^{-13}$
- operator  $*<$  (x,y: rrG) p: rrG;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $\implies x \cdot y > p; \quad |p.r2/(p.r1 + p.r2)| < 5.0001 \cdot 10^{-13}$

### 3.5.4 Division

Mit Hilfe der folgenden Prozeduren und Operatoren lassen sich Quotienten vom Typ rrG berechnen. Falls die Zehnerexponenten beider rrG-Operanden übereinstimmen, sind aus Laufzeitgründen die Prozeduren den Operatoren vorzuziehen. Es existieren auch die wichtigsten Operatoren mit gemischten Operanden.

- procedure **DIV\_22** (a1,a2,b1,b2: real; var q1,q2: real);  
 Vorauss.:  $|a2/(a1 + a2)| < 5.0001 \cdot 10^{-13}$   
 $|b2/(b1 + b2)| < 5.0001 \cdot 10^{-13}$   
 $(q1 + q2) = [(a1 + a2)/(b1 + b2)] \cdot (1 + \varepsilon), \quad |q1| \geq 10^{-241}$   
 $\implies |\varepsilon| \leq 2.2504 \cdot 10^{-24}; \quad |r2/(r1 + r2)| < 5.0001 \cdot 10^{-13}$
- procedure **REZIP** (x: rrG; var R: rrG);  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $R = (1/x) \cdot (1 + \varepsilon)$   
 $\implies |\varepsilon| \leq 2.2504 \cdot 10^{-24}; \quad |R.r2/(R.r1 + R.r2)| < 5.0001 \cdot 10^{-13}$
- operator  $/$  (x,y: rrG) q: rrG;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $q = (x/y) \cdot (1 + \varepsilon)$   
 $\implies |\varepsilon| \leq 2.2505 \cdot 10^{-24}; \quad |q.r2/(q.r1 + q.r2)| < 5.0001 \cdot 10^{-13}$



- operator / (x: real; y: rrG) q: rrG;  
 Vorauss.:  $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $q = (x/y) \cdot (1 + \varepsilon)$   
 $\implies |\varepsilon| \leq 2.2505 \cdot 10^{-24}; \quad |q.r2/(q.r1 + q.r2)| < 5.0001 \cdot 10^{-13}$
- operator / (x: Longreal; y: rrG) q: rrG;  
 Vorauss.:  $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $q = (x/y) \cdot (1 + \varepsilon)$   
 $\implies |\varepsilon| \leq 2.2505 \cdot 10^{-24}; \quad |q.r2/(q.r1 + q.r2)| < 5.0001 \cdot 10^{-13}$
- operator / (x: rrG; y: real) q: rrG;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $q = (x/y) \cdot (1 + \varepsilon)$   
 $\implies |\varepsilon| \leq 2.2505 \cdot 10^{-24}; \quad |q.r2/(q.r1 + q.r2)| < 5.0001 \cdot 10^{-13}$
- operator / (x: rrG; y: Longreal) q: rrG;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $q = (x/y) \cdot (1 + \varepsilon)$   
 $\implies |\varepsilon| \leq 2.2505 \cdot 10^{-24}; \quad |q.r2/(q.r1 + q.r2)| < 5.0001 \cdot 10^{-13}$
- operator /> (x,y: rrG) q: rrG;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $\implies x/y < q; \quad |q.r2/(q.r1 + q.r2)| < 5.0001 \cdot 10^{-13}$
- operator /< (x,y: rrG) q: rrG;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $\implies x/y > q; \quad |q.r2/(q.r1 + q.r2)| < 5.0001 \cdot 10^{-13}$

### 3.5.5 Vergleiche

Mit Hilfe der folgenden Funktionen und Operatoren lassen sich Operanden vom Typ real,Longreal,rrG vergleichen:

- function **greater\_21** (x1,x2,x: real): boolean;  
 Vorauss.:  $|x2/(x1 + x2)| < 5.0001 \cdot 10^{-13}, \quad x1 = 0 \implies x2 = 0;$   
 $\text{greater\_21} = \text{true} \iff (x1 + x2) > x;$
- function **less\_21** (x1,x2,x: real): boolean;  
 Vorauss.:  $|x2/(x1 + x2)| < 5.0001 \cdot 10^{-13}, \quad x1 = 0 \implies x2 = 0;$   
 $\text{less\_21} = \text{true} \iff (x1 + x2) < x;$

- operator = (x,y: rrG) eq: boolean;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $\implies \text{eq} = \text{true} \iff x = y;$
- operator > (x,y: rrG) gr: boolean;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $\implies \text{gr} = \text{true} \iff x > y;$
- operator < (x,y: rrG) le: boolean;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $\implies \text{le} = \text{true} \iff x < y;$
- operator  $\geq$  (x,y: rrG) ge: boolean;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $\implies \text{ge} = \text{true} \iff x \geq y;$
- operator  $\leq$  (x,y: rrG) le: boolean;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $\implies \text{le} = \text{true} \iff x \leq y;$
- operator <> (x,y: rrG) neq: boolean;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$   
 $\implies \text{neq} = \text{true} \iff x <> y;$
- operator = (x: rrG; a: real) eq: boolean;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $\implies \text{eq} = \text{true} \iff x = a;$
- operator <> (x: rrG; a: real) neq: boolean;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $\implies \text{neq} = \text{true} \iff x <> a;$
- operator > (x: rrG; a: real) gr: boolean;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $\implies \text{gr} = \text{true} \iff x > a;$
- operator < (x: rrG; a: real) le: boolean;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $\implies \text{le} = \text{true} \iff x < a;$

- operator  $\geq$  (x: **rrG**; a: real) ge: boolean;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $\implies \text{ge} = \text{true} \iff x \geq a;$
- operator  $\leq$  (x: **rrG**; a: real) le: boolean;  
 Vorauss.:  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$   
 $\implies \text{le} = \text{true} \iff x \leq a;$

Bei den letzten Operatoren kann a auch vom Typ **Longreal** gewählt werden, so daß Vergleiche auch zwischen **rrG**- und **Longreal**-Operanden möglich sind.

### 3.5.6 Transferfunktionen und Zuweisungsoperatoren

Die folgenden Funktionen und Operatoren ermöglichen die Verbindung zwischen den Datentypen **real**, **Longreal**, **rrG** :

- operator := (var a: **rrG**; r: real);
- operator := (var a: **rrG**; r: **Longreal**);
- function **LongNext** (x: **rrG**) : **Longreal**; Rundg. zur nächsten **Longreal**zahl
- function **LongUp** (x: **rrG**) : **Longreal**; zur nächstgrößeren **Longreal**zahl;
- function **LongDown** (x: **rrG**) : **Longreal**; zur nächstkleineren **Longreal**zahl;
- function **RealNext** (x: **rrG**) : real; Rundung zur nächsten **real**-Zahl;
- function **Trans\_rrG** (x1,x2: **Longreal**) : **rrG**;  
 Vorauss.:  $x1 = 0 \Rightarrow x2 = 0$ ;  $x1 <> 0 \Rightarrow |x2/(x1 + x2)| < 5.0001 \cdot 10^{-21}$ ;  
 $y := \text{Trans\_rrG}(x1,x2)$  liefert:  $y = (x1 + x2) \cdot (1 + \varepsilon) \implies |\varepsilon| \leq 5 \cdot 10^{-26}$ ;

### 3.5.7 Hilfsfunktionen

Bei den auftretenden Funktionsparametern vom Typ **rrG** wird vorausgesetzt, daß die Bedingungen I,II von Seite 46 erfüllt sind. Dies wird immer dann der Fall sein, wenn die **real**-Komponenten dieser Parameter durch das auf der gleichen Seite beschriebene zweimalige Auslesen des Akkumulators ins **real**-Format definiert sind. Besitzen Funktionen den Ergebnistyp **rrG**, so erfüllen auch die Komponenten aller Funktionswerte automatisch die Bedingungen I,II von Seite 46.

- function **Expo** (x: **rrG**) : integer;  
 Berechnet den Zehnerexponenten analog zu  $\text{Expo}(x: \text{real}): \text{integer}$ ;
- procedure **Mant\_26** (x1,x2: real; var r1,r2: real);  
 Berechnet die Mantisse von  $(x1 + x2)$  analog zu  $\text{Mant}(x: \text{real}): \text{real}$ , falls  $|x.r2/(x.r1 + x.r2)| < 5.0001 \cdot 10^{-13}$  erfüllt ist.  
 Es gilt also:  $|r1 + r2| = 0$  oder  $0.1 \leq |r1 + r2| < 1$ ;

- function **Scale** (x: rrG): rrG;  
y := Scale(x) liefert folgende Ergebnisse:

1.  $y = x$ , '=' im Sinne der Numerik;
2.  $y.r1 = 0 \implies y.r2 = 0, y.G = -\text{Maxint} = -2147483647$ ;
3.  $0.1 \leq |y.r1| < 1$ , falls  $y.r1 \neq 0$ ; (Skalierung);
4.  $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$ , falls  $y.r1 \neq 0$ ;

Falls das Ergebnis  $y = x$  in extremen Ausnahmefällen nicht erfüllt werden kann, erfolgt eine entsprechende Fehlermeldung.

- function **Normalisiere** (x: rrG): rrG;  
y := Normalisiere(x) liefert folgende Ergebnisse:

1.  $y.r1 = 0 \implies y.r2 = 0, y.G = -\text{Maxint} = -2147483647$ ;
2.  $0.1 \leq |y.r1| < 1$ , falls  $y.r1 \neq 0$ ; (Normalisierung);
3.  $|y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$ , falls  $y.r1 \neq 0$ ;
4.  $y = x \cdot (1 + \varepsilon) \implies |\varepsilon| < 1.0001 \cdot 10^{-266}$ ;
5.  $y = x + \delta \implies |\delta| < 10^{-267+x.G+\text{Expo}(x.r1)}$ ;

Bei extrem kleinen oder großen  $|x|$ -Werten erfolgt eine Fehlermeldung, wenn  $y.G$  im integer-Bereich nicht mehr darstellbar ist, was für die Praxis jedoch keine Einschränkung bedeutet.

- function **Bound** (a: rrG; eps: real): rrG;  
a sei eine Maschinennäherung für den exakten Wert  $f(x)$ , und es gelte:

1.  $a = f(x) \cdot (1 + \varepsilon)$ ,  $|\varepsilon| \leq \varepsilon(f)$ ,  $\text{epsup} := 1.00001 * \varepsilon(f)$ ;
2.  $\text{epsup} < 10^{-13}$ , was keine praktische Einschränkung bedeutet;

Mit:  $B1 := \text{Bound}(a, -\text{epsup})$ ;  $B2 := \text{Bound}(a, +\text{epsup})$  erhält man garantierte Schranken für  $f(x)$ :  $B1 < f(x) < B2$ ,  $B1, B2: \text{rrG}$ ;  
 $a = 0$  wird durch die Funktion **Bound** nicht auf- oder abgerundet. Bei extrem kleinen oder großen  $|a|$ -Werten erfolgt (durch die Funktion **Normalisiere**) eine Fehlermeldung, was für die Praxis jedoch keine Einschränkung bedeutet.

- function **Bound\_toL** (a: rrG; eps: real): Longreal;  
a sei eine Maschinennäherung für den exakten Wert  $f(x)$ , und es gelte:

$$a = f(x) \cdot (1 + \varepsilon), \quad |\varepsilon| \leq \varepsilon(f), \quad \text{epsup} := 1.00001 * \varepsilon(f);$$

Mit:  $L1 := \text{Bound\_toL}(a, -\text{epsup})$ ;  $L2 := \text{Bound\_toL}(a, +\text{epsup})$  erhält man garantierte Schranken für  $f(x)$ :  $L1 < f(x) < L2$ ,  $L1, L2: \text{Longreal}$ ;  
 $a = 0$  liefert:  $L1 = -10^{-531}$ ,  $L2 = +10^{-531}$ ;  
 $0 < a < +10^{-531}$  liefert:  $L1 = 0$ ,  $L2 = +2 \cdot 10^{-531}$ ;  
 $-10^{-531} < a < 0$  liefert:  $L1 = -2 \cdot 10^{-531}$ ,  $L2 = 0$ ;

- function **Prod\_Sum** (a,b,c,d: rRG): rRG;  
 $y := \text{Prod\_Sum}(a,b,c,d)$  liefert:  
 $y = (a \cdot b + c \cdot d) \cdot (1 + \varepsilon)$ ,  $|\varepsilon| < 1.1112 \cdot 10^{-25}$ ;
- function **Transform** (x: rRG): rRG;  
 Vorauss.: x.r1, x.r2, x.G seien beliebig;  
 Durch  $y := \text{Transform}(x)$  wird x so transformiert, daß gilt:
  1.  $y = x$ ;
  2.  $y.r1 = 0 \implies y.r2 = 0$ ;
  3.  $y.r1 \neq 0 \implies |y.r2/(y.r1 + y.r2)| < 5.0001 \cdot 10^{-13}$ ;

Falls  $y = x$  nicht erfüllt werden kann, erfolgt eine Fehlermeldung.

- procedure **Arg\_Red** (x1,x2,m: real; var r1,r2: real);  
 Zu gegebenem  $x = (x1 + x2)$ , mit  $\pi/2 < |x| < 10^{+8}$  ist  $(r1 + r2)$  eine Näherung für das reduzierte Argument  $xred := x - m \cdot \pi$ ;  
 $m := n/2$ , mit  $n \in \mathbb{Z}$ , muß dabei so vorgegeben sein, daß gilt:  
 $|xred| < \pi/2 + t$ ,  $0 \leq t \ll 1$ ;  
 $(r1 + r2) = xred \cdot (1 + \varepsilon)$ ,  $|\varepsilon| < 5.0002 \cdot 10^{-26}$ ;
- function **ENTIRE** (x: rRG; var gerade: boolean): boolean;  
 $bl := \text{ENTIRE}(x, \text{gerade})$  liefert an bl den Wert true, wenn die rRG-Zahl x eine **ganze** Zahl ist, sonst erhält bl den Wert false.  
 Falls x ganzzahlig ist, erhält gerade den Wert true (false), wenn x zusätzlich gerade (ungerade) ist.  
 Ist x nicht ganzzahlig, so bleibt die Variable gerade bedeutungslos.

### 3.5.8 HornerSchema

Gegeben sei das Polynom vom Grade  $N = 9$ :

$$P_9(x) = (b1[0] + b2[0]) + (b1[1] + b2[1]) \cdot x^1 + a[2] \cdot x^2 + \dots + a[9] \cdot x^9$$

das für das 26-stellige BCD-Argument  $x = (x1 + x2)$  nach dem HornerSchema auszuwerten ist. Die real-Zahlen  $x1, x2$  müssen folgende Bedingungen erfüllen:

$$x1 = 0 \implies x2 = 0; \quad x1 \neq 0 \implies |x2/(x1 + x2)| < 5.0001 \cdot 10^{-13}$$

Aus Laufzeitgründen werden die ersten  $gr = 7$  groben HornerSchritte im Longrealformat durchgeführt:

$$\text{var } a: \text{Lrvector}[2..9]; \quad \{ \text{LB}(a) := N - gr = 2; \quad \text{UB}(a) := N = 9; \}$$

Im achten, gemischten Schritt wird die Multiplikation im Longrealformat ausgeführt und die Addition im 26-stelligen BCD-Format. Alle restlichen HornerSchritte laufen im 26-stelligen BCD-Format wie folgt ab:

- Die Multiplikation wird im Akkumulator ausgeführt, wobei aus Laufzeitgründen das betragsmäßig kleinste Teilprodukt weggelassen wird. Da der Akku nach dieser Multiplikation nicht ins (real|real)-Format ausgelesen wird, gilt für die relative Fehlerschranke:  $\varepsilon(\text{Mul}) = 2.5002 \cdot 10^{-25}$ ;
- Auf das obige Multiplikationsergebnis wird im Akku der jeweilige Polynomkoeffizient rundungsfehlerfrei addiert. Diese Summe wird anschließend in das (real|real)-Format ausgelesen, weshalb die Fehlerschranke der Addition mit  $\varepsilon(\text{Add}) = 5 \cdot 10^{-26}$  anzusetzen ist. Dabei muß nach Seite 46 vorausgesetzt werden, daß für jedes Ergebnis *erg* eines Hornerschnittes gilt:  $|\text{erg}| \geq 10^{-241}$ . Diese Bedingung ist in allen praktischen Fällen erfüllt, wenn  $|x1 + x2|$  nicht zu klein gewählt wird und wenn das Polynom nicht in der Nähe einer seiner Nullstellen ausgewertet wird. Eine genaue Überprüfung kann mit dem **PASCAL-XSC** Programm **Horner\_T** vorgenommen werden.

```
var b1,b2: rvector[0..1];    { LB(b1) := 0;  UB(b1) := N - gr - 1 = 1; }
```

Die oben beschriebene Polynomauswertung erfolgt durch:

```
procedure HORNER (a: Lrvector; b1,b2: rvector; x1,x2: real; var r1,r2: real;
                 l_add: boolean);
```

Der Aufruf obiger Prozedur mit den deklarierten Eingangsparametern *a*, *b1*, *b2*, *x1*, *x2* liefert mit den Ausgangsgrößen *r1*, *r2* eine Näherung  $(r1 + r2) \approx P_9(x1 + x2)$  für den Polynomwert  $P_9(x1 + x2)$ . Es gilt:

- $r1 = 0 \implies r2 = 0$ ;  $r1 \ll 0 \implies |r2/(r1+r2)| < 5.0001 \cdot 10^{-13}$
- $l\_add = \text{true} \iff$  Im letzten Hornerschnitt wird  $(b1[0] + b2[0])$  addiert.

Die Berechnung von garantierten Oberschranken des absoluten oder relativen Polynomauswertefehlers, der beim Aufruf der Prozedur **HORNER** entsteht, ist mit folgenden Programmen möglich:

```
IFO_BCD   IEEE-Programm, das Polynominformationen sammelt;
AWF_BCD   IEEE-Programm, das den Auswertefehler abschätzt;
Horner_T   BCD-Programm überprüft AWF_BCD-Ergebnisse.
```

### 3.5.9 Ein-/AusgabeprozEDUREN

Es stehen die gewohnten Prozeduren

```
procedure read (var f: text; var a: rrG);
procedure write (var f: text; a: rrG);
```

mit optionalem Fileparameter, beliebig vielen Ein-/Ausgabeparametern, jedoch ohne Formatspezifikationen zur Verfügung. *a: rrG* muß in der Form

$\{a.r1, a.r2, a.G\}$  oder in der Form *a.r1*

eingegeben werden. Im zweiten Fall gilt:  $a.r2 = a.G = 0$ . Die Komponente  $a.G$  erhält jedoch stets den Wert  $-\text{maxint} = -2147483647$ , falls  $a.r1$  verschwindet.

Erfolgt die Eingabe in der ersten Form mit beliebigen real-Werten  $a.r1, a.r2$ , so wird auf  $a : \text{rrG}$  automatisch die Funktion  $\text{Transform}(a : \text{rrG}) : \text{rrG}$  angewandt, so daß ohne Änderung des numerischen Werts von  $a$  die beiden folgenden Bedingungen von Seite 46 erfüllt werden:

$$a.r1 = 0 \implies a.r2 = 0 \quad (\text{Beding. I})$$

$$|a.r1| > 0 \implies \left| \frac{a.r2}{a.r1 + a.r2} \right| < 5.0001 \cdot 10^{-13} \quad (\text{Beding. II})$$

Die Ausgabe einer rrG-Zahl erfolgt stets in der Form:

$$\{a.r1, a.r2, a.G\}$$

#### Beispiele:

Sei  $a$  vom Typ  $\text{rrG}$ , dann wird mit den Anweisungen `read(a); write(a)` und der Eingabe

$$\{1.234567890123E35, 5.934567890123E22, 1234567890\}$$

die rrG-Zahl  $a$  in der folgenden Form ausgegeben:

$$\{1.234567890124E+035, -4.065432109877E+022, 1234567890\}$$

Die Eingabe

$$\{4, -4, 1234567890\}$$

ergibt mit den gleichen obigen Anweisungen die Ausgabe des Wertes Null:

$$\{0.000000000000E+000, 0.000000000000E+000, -2147483647\}$$

Die Eingabe

$$\{9e+255, 9e+255, +2147483647\}$$

ergibt mit der obigen `read`-Anweisung eine Fehlermeldung, da der numerische Eingabewert wegen Overflow im  $\text{rrG}$ -Format nicht darstellbar ist.

Eine entsprechende Underflow-Fehlermeldung läßt sich mit den `read`- und `write`-Anweisungen jedoch nicht konstruieren!

## 3.5.10 Standardfunktionen

Standardfunktionen vom Typ rrG; x,y: rrG; a: Longreal			
Funktion	Aufruf	Argumentbereich	Rel. Fehler
$ x $	abs(x)	rrG-Bereich	0.0000
$x^2$	sqr(x)	***	$5.0001 \cdot 10^{-26}$
$x^2 - y^2$	x2_y2(x, y)	***	$1.5001 \cdot 10^{-25}$
$\sqrt{x}$	sqrt(x)	$x \geq 0$	$1.2253 \cdot 10^{-24}$
$\sqrt{1+x} - 1$	sqr1pm1(x)	$x \geq -1$ ; ***	$1.8403 \cdot 10^{-24}$
$\sqrt{1-x^2}$	sqr1mx2(x)	$ x  \leq 1$	$1.3004 \cdot 10^{-24}$
$\sqrt{1+x^2}$	sqr1px2(x)	rrG-Bereich	$1.2754 \cdot 10^{-24}$
$\sqrt{x^2-1}$	sqrtox2m1(x)	$ x  \geq 1$	$1.3004 \cdot 10^{-24}$
$\sqrt{x^2+y^2}$	sqrtox2y2(x, y)	$x^2 + y^2 \geq 0$ ; ***	$1.2754 \cdot 10^{-24}$
$e^x$	exp(x)	$ x  < 5 \cdot 10^{+3}$ ; *	$4.6213 \cdot 10^{-25}$
$e^x - 1$	expm1(x)	$x < 5 \cdot 10^{+3}$	$5.5478 \cdot 10^{-25}$
$e^{-x^2}$	exp_x2(x)	$ x  \leq 70.710516$ ; **	$4.6213 \cdot 10^{-25}$
$2^x$	exp2(x)	$ x  < 5 \cdot 10^{+3}$ ; *	$4.6833 \cdot 10^{-25}$
$10^x$	exp10(x)	$-10^8 < x < \text{maxint}$ ; *	$6.3644 \cdot 10^{-24}$
$\ln(x)$	ln(x)	$10^{-5000} \leq x < 10^{+5001}$	$2.7329 \cdot 10^{-24}$
$\ln(e \cdot x)$	ln_ex(x)	$10^{-5000} \leq x < 10^{+5001}$	$4.5584 \cdot 10^{-24}$
$\ln(1+x)$	ln1p(x)	$-1 < x < 10^{+5001}$	$2.7846 \cdot 10^{-24}$
$\ln(x^2 + y^2)$	lnx2y2(x, y)	$x^2 + y^2 > 0$ ; ***	$3.2816 \cdot 10^{-24}$
$0.5 \cdot \ln(x^2 + y^2)$	ln_sqrtx2y2(x, y)	$x^2 + y^2 > 0$ ; ***	$3.3317 \cdot 10^{-24}$
$\ln[(1+x)^2 + y^2]$	ln1px2y2(x, y)	$(1+x)^2 + y^2 > 0$ ; ***	$3.8301 \cdot 10^{-24}$
$0.5 \cdot \ln[(1+x)^2 + y^2]$	ln_sqrt1px2y2(x, y)	$(1+x)^2 + y^2 > 0$ ; ***	$3.8802 \cdot 10^{-24}$
$\log_a(x)$	loga(a, x)	$10^{-5000} \leq x < 10^{+5001}$	$7.7163 \cdot 10^{-24}$
$\log_2(x)$	log2(x)	$10^{-5000} \leq x < 10^{+5001}$	$3.0830 \cdot 10^{-24}$
$\log_{10}(x)$	log10(x)	$10^{-5000} \leq x < 10^{+5001}$	$3.0830 \cdot 10^{-24}$
$\log_{10}(1+x)$	log1p(x)	$-1 < x < 10^{+5001}$	$2.8371 \cdot 10^{-24}$
$x^n$	power(x,n,rnd)	$n \geq 0$ ; rnd = +1, -1	



Funktion	Aufruf	Argumentbereich	Rel. Fehler
$x^y$	power( $x, y$ )	$ y \cdot \ln(x)  < 5 \cdot 10^{+3}$ ; *	$5.6503 \cdot 10^{-25}$
$\sin(x)$	sin( $x$ )	$ x  < 10^{+8}$	$1.0196 \cdot 10^{-24}$
$\cos(x)$	cos( $x$ )	$ x  < 10^{+8}$	$1.0196 \cdot 10^{-24}$
$\tan(x)$	tan( $x$ )	$ x  < 10^{+8}$	$2.9035 \cdot 10^{-24}$
$\cot(x)$	cot( $x$ )	$ x  < 10^{+8}$	$2.9035 \cdot 10^{-24}$
$\sin(\pi \cdot x)$	sinpi( $x$ )	$ x  \leq 2147483647$	$1.0401 \cdot 10^{-24}$
$\cos(\pi \cdot x)$	cospi( $x$ )	$ x  \leq 2147483647$	$1.1556 \cdot 10^{-24}$
$\tan(\pi \cdot x)$	tanpi( $x$ )	$ x  < 1073741823.5$	$2.9035 \cdot 10^{-24}$
$\cot(\pi \cdot x)$	cotpi( $x$ )	$ x  < 1073741823.5$	$2.9035 \cdot 10^{-24}$
$\arcsin(x)$	arcsin( $x$ )	$ x  \leq 1$	$3.0671 \cdot 10^{-24}$
$\arccos(x)$	arccos( $x$ )	$ x  \leq 1$	$3.1225 \cdot 10^{-24}$
$\arctan(x)$	arctan( $x$ )	rrG-Bereich	$3.0170 \cdot 10^{-24}$
$\arctan(x/y)$	arctan2( $x, y$ )	***	$5.2675 \cdot 10^{-24}$
$\operatorname{arccot}(x)$	arccot( $x$ )	rrG-Bereich	$5.2675 \cdot 10^{-24}$
$\sinh(x)$	sinh( $x$ )	$ x  \leq 5 \cdot 10^{+3}$	$2.6320 \cdot 10^{-24}$
$\cosh(x)$	cosh( $x$ )	$ x  \leq 5 \cdot 10^{+3}$	$1.6875 \cdot 10^{-24}$
$\tanh(x)$	tanh( $x$ )	rrG-Bereich	$4.5980 \cdot 10^{-24}$
$\operatorname{coth}(x)$	coth( $x$ )	$x \neq 0$ ; ***	$6.8486 \cdot 10^{-24}$
$\operatorname{arsinh}(x)$	arsinh( $x$ )	$ x  < 10^{+5001}$	$4.5070 \cdot 10^{-24}$
$\operatorname{arcosh}(x)$	arcosh( $x$ )	$1 \leq x < 10^{+5001}$	$3.8908 \cdot 10^{-24}$
$\operatorname{arcosh}(1+x)$	arcosh1p( $x$ )	$0 \leq x < 10^{+5001}$	$6.7674 \cdot 10^{-24}$
$\operatorname{artanh}(x)$	artanh( $x$ )	$ x  < 1$	$4.5687 \cdot 10^{-24}$
$\operatorname{arcoth}(x)$	arcoth( $x$ )	$ x  > 1$	$4.6827 \cdot 10^{-24}$

- \* Wird der angegebene Definitionsbereich nach links überschritten, so wird der Funktionswert auf Null gesetzt.
- \*\* Wird der angegebene Definitionsbereich nach links und rechts überschritten, so wird der Funktionswert auf Null gesetzt.
- \*\*\* Falls der rrG-Bereich wegen Overflow- oder Underflow verlassen wird, erfolgt eine entsprechende Fehlermeldung.

Standardfunktionen vom Typ rrG; x,y,a: Longreal			
Funktion	Aufruf	Argumentbereich	Rel. Fehler
$ x $	abs_rrG( $x$ )	Longreal-Bereich	0.0000
$x^2$	sqr_rrG( $x$ )	Longreal-Bereich	$5.0000 \cdot 10^{-26}$
$x^2 - y^2$	x2_y2_rrG( $x, y$ )	Longreal-Bereich	$1.5001 \cdot 10^{-25}$
$\sqrt{x}$	sqr_rrG( $x$ )	$x \geq 0$	$1.2253 \cdot 10^{-24}$
$\sqrt{1+x} - 1$	sqr1pm1_rrG( $x$ )	$x \geq -1$	$1.8403 \cdot 10^{-24}$
$\sqrt{1-x^2}$	sqr1mx2_rrG( $x$ )	$ x  \leq 1$	$1.2504 \cdot 10^{-24}$
$\sqrt{1+x^2}$	sqr1px2_rrG( $x$ )	Longreal-Bereich	$1.2504 \cdot 10^{-24}$
$\sqrt{x^2 - 1}$	sqr x2m1_rrG( $x$ )	$ x  \geq 1$	$1.2504 \cdot 10^{-24}$
$\sqrt{x^2 + y^2}$	sqr x2y2_rrG( $x, y$ )	Longreal-Bereich	$1.2504 \cdot 10^{-24}$
$e^x$	exp_rrG( $x$ )	$ x  < 5 \cdot 10^{+3}$ ; *	$4.6213 \cdot 10^{-25}$
$e^x - 1$	expm1_rrG( $x$ )	$x < 5 \cdot 10^{+3}$	$5.5478 \cdot 10^{-25}$
$e^{-x^2}$	exp_x2_rrG( $x$ )	$ x  \leq 70.710516$ ; **	$4.6213 \cdot 10^{-25}$
$2^x$	exp2_rrG( $x$ )	$ x  < 5 \cdot 10^{+3}$ ; *	$4.6833 \cdot 10^{-25}$
$10^x$	exp10_rrG( $x$ )	$-10^8 < x < \text{maxint}$ ; *	$6.3644 \cdot 10^{-24}$
$\ln(x)$	ln_rrG( $x$ )	$x > 0$	$2.7329 \cdot 10^{-24}$
$\ln(e \cdot x)$	ln_ex_rrG( $x$ )	$x > 0$	$4.5584 \cdot 10^{-24}$
$\ln(1+x)$	ln1p_rrG( $x$ )	$x > -1$	$2.7846 \cdot 10^{-24}$
$\ln(x^2 + y^2)$	lnx2y2_rrG( $x, y$ )	$x^2 + y^2 > 0$	$5.1539 \cdot 10^{-24}$
$\ln(x^2 + y^2)/2$	ln_sqrtx2y2_rrG( $x, y$ )	$x^2 + y^2 > 0$	$5.2040 \cdot 10^{-24}$
$\ln[(1+x)^2 + y^2]$	ln1px2y2_rrG( $x, y$ )	$(1+x)^2 + y^2 > 0$	$2.8929 \cdot 10^{-24}$
$\ln[(1+x)^2 + y^2]/2$	ln_sqrt1px2y2_rrG( $x, y$ )	$(1+x)^2 + y^2 > 0$	$2.9430 \cdot 10^{-24}$
$\log_a(x)$	loga_rrG( $a, x$ )	$a, x > 0$ ; $a \neq 1$	$7.7163 \cdot 10^{-24}$
$\log_2(x)$	log2_rrG( $x$ )	$x > 0$	$3.0830 \cdot 10^{-24}$
$\log_{10}(x)$	log10_rrG( $x$ )	$x > 0$	$3.0830 \cdot 10^{-24}$
$\log_{10}(1+x)$	log1p_rrG( $x$ )	$x > -1$	$2.8371 \cdot 10^{-24}$
$x^n$	power_rrG( $x, n, \text{rnd}$ )	$n \geq 0$ ; $\text{rnd} = +1, -1$	
$x^y$	power_rrG( $x, y$ )	$ y \cdot \ln(x)  < 5 \cdot 10^{+3}$ ; *	$4.6212 \cdot 10^{-25}$

Funktion	Aufruf	Argumentbereich	Rel. Fehler
$\sin(x)$	<code>sin_rrG(x)</code>	$ x  < 10^{+8}$	$1.0196 \cdot 10^{-24}$
$\cos(x)$	<code>cos_rrG(x)</code>	$ x  < 10^{+8}$	$1.0196 \cdot 10^{-24}$
$\tan(x)$	<code>tan_rrG(x)</code>	$ x  < 10^{+8}$	$2.9035 \cdot 10^{-24}$
$\cot(x)$	<code>cot_rrG(x)</code>	$ x  < 10^{+8}$	$2.9035 \cdot 10^{-24}$
$\sin(\pi \cdot x)$	<code>sinpi_rrG(x)</code>	$ x  \leq 2147483647$	$1.0401 \cdot 10^{-24}$
$\cos(\pi \cdot x)$	<code>cospi_rrG(x)</code>	$ x  \leq 2147483647$	$1.1556 \cdot 10^{-24}$
$\tan(\pi \cdot x)$	<code>tanpi_rrG(x)</code>	$ x  < 1073741823.5$	$2.9035 \cdot 10^{-24}$
$\cot(\pi \cdot x)$	<code>cotpi_rrG(x)</code>	$ x  < 1073741823.5$	$2.9035 \cdot 10^{-24}$
$\arcsin(x)$	<code>arcsin_rrG(x)</code>	$ x  \leq 1$	$6.5180 \cdot 10^{-24}$
$\arccos(x)$	<code>arccos_rrG(x)</code>	$ x  \leq 1$	$6.5734 \cdot 10^{-24}$
$\arctan(x)$	<code>arctan_rrG(x)</code>	Longreal-Bereich	$3.0170 \cdot 10^{-24}$
$\arctan(x/y)$	<code>arctan2_rrG(x, y)</code>	Longreal-Bereich	$5.2675 \cdot 10^{-24}$
$\operatorname{arccot}(x)$	<code>arccot_rrG(x)</code>	Longreal-Bereich	$5.2675 \cdot 10^{-24}$
$\sinh(x)$	<code>sinh_rrG(x)</code>	$ x  \leq 5 \cdot 10^{+3}$	$2.6320 \cdot 10^{-24}$
$\cosh(x)$	<code>cosh_rrG(x)</code>	$ x  \leq 5 \cdot 10^{+3}$	$1.6875 \cdot 10^{-24}$
$\tanh(x)$	<code>tanh_rrG(x)</code>	Longreal-Bereich	$4.5980 \cdot 10^{-24}$
$\operatorname{coth}(x)$	<code>coth_rrG(x)</code>	$x \neq 0$	$6.8486 \cdot 10^{-24}$
$\operatorname{arsinh}(x)$	<code>arsinh_rrG(x)</code>	Longreal-Bereich	$4.5070 \cdot 10^{-24}$
$\operatorname{arcosh}(x)$	<code>arcosh_rrG(x)</code>	$x \geq 1$	$3.8513 \cdot 10^{-24}$
$\operatorname{arcosh}(1+x)$	<code>arcosh1p_rrG(x)</code>	$x \geq 0$	$6.7674 \cdot 10^{-24}$
$\operatorname{artanh}(x)$	<code>artanh_rrG(x)</code>	$ x  < 1$	$4.5687 \cdot 10^{-24}$
$\operatorname{arcoth}(x)$	<code>arcoth_rrG(x)</code>	$ x  > 1$	$4.6827 \cdot 10^{-24}$

\* Wird der angegebene Definitionsbereich nach links überschritten, so wird der Funktionswert auf Null gesetzt.

\*\* Wird der angegebene Definitionsbereich nach links und rechts überschritten, so wird der Funktionswert auf Null gesetzt.

Mit Hilfe obiger Funktionen vom Ergebnistyp **rrG** und Longreal-Argumenten  $x$  können auch verschachtelte Funktionen für den Ergebnistyp Longreal **hochgenau** implementiert werden.

Als Beispiel betrachten wir  $\ln(\cos(x))$  für  $0 \leq x < \pi/2$ ; Mit der naiven Realisierung:

```
function ln_cos(x: Longreal) : Longreal;
begin
  ln_cos := LongNext( ln( cos_rrG(x) ) );
end;
```

erhält man zum Argument  $x_0 = 10^{-14}$  den im Vergleich zum exakten Funktionswert  $-5.00 \dots 0083 \dots \cdot 10^{-29}$  völlig unbrauchbaren Maschinenwert  $0$ . Dieses schlechte Ergebnis ist dadurch begründet, daß die  $\ln$ -Funktion mit fehlerbehafteten  $\cos$ -Funktionswerten in der Nähe ihrer Nullstelle  $1$  ausgewertet wird, wobei  $\cos\_rrG(x)$  in einer ganzen Umgebung von Null den fehlerbehafteten Funktionswert  $1$  liefert. Benutzt man jedoch für  $x < 1$  die Identität

$$\ln(\cos(x)) \equiv \frac{1}{2} \cdot \ln[1 - \sin^2(x)]$$

so liefert die entsprechende Implementierung

```
function ln_cos_(x: Longreal): Longreal;
begin
  if x < 1 then
    ln_cos_ := LongNext( 0.5*ln1p( -sqr(sin_rrG(x)) ) ) else
    ln_cos_ := LongNext( ln( cos_rrG(x) ) );
end;
```

für das gleiche Argument  $x_0 = 10^{-14}$  den im Longreal-Format hochgenauen Funktionswert

$$\ln\_cos\_ (x_0) = -5.00000000000000000000 \cdot 10^{-29}$$

**Fehlerabschätzung** für  $0 \leq x < 1$ :

Mit den Bezeichnungen  $\tilde{y} := \text{sqr}(\text{sin\_rrG}(x))$ ,  $y := \sin(x)$  gilt zunächst:

$$\begin{aligned} \tilde{y} &= \sin^2(x) \cdot (1 + \varepsilon_{\text{sin}})^2 \cdot (1 + \varepsilon_{\text{sqr}}) \\ &= y \cdot (1 + \varepsilon_y); \quad |\varepsilon_y| < 2.0893 \cdot 10^{-24} = \varepsilon(y) \end{aligned}$$

$$\begin{aligned} 0.5 \odot \ln1p(-\tilde{y}) &= 0.5 \cdot \ln(1 - \tilde{y}) \cdot (1 + \varepsilon_{\ln1p}) \cdot (1 + \varepsilon_M) \\ |\varepsilon_{\ln1p}| &< 2.7846 \cdot 10^{-24} = \varepsilon(\ln1p); \quad |\varepsilon_M| < 5.0001 \cdot 10^{-26} = \varepsilon(M); \end{aligned}$$

Mit der Abkürzung:  $\alpha := \ln \left[ 1 - \frac{y \cdot \varepsilon_y}{1 - y} \right] / \ln(1 - y)$  gilt dann die Gleichung:

$\ln(1 - \tilde{y}) = \ln(1 - y) \cdot (1 + \alpha)$ , und einige Abschätzungen liefern:

$$|\alpha| < \varepsilon(y) / \left[ 1 - \frac{y}{1 - y} \cdot \varepsilon(y) \right]$$

$$x < 1 \rightsquigarrow y = \sin^2(x) < 0.7081 \rightsquigarrow |\alpha| < 2.0896 \cdot 10^{-24} \rightsquigarrow$$

$$\begin{aligned} 0.5 \odot \ln 1p(-\tilde{y}) &= 0.5 \cdot \ln(1-y) \cdot (1+\alpha) \cdot (1+\varepsilon_{\ln 1p}) \cdot (1+\varepsilon_M) \\ &= \ln(\cos(x)) \cdot (1+\varepsilon_1); \quad |\varepsilon_1| < 4.9243 \cdot 10^{-24} = \varepsilon(1) \end{aligned}$$

Für  $x < 1$  wird also  $\ln(\cos(x))$  im 26-stelligen rrG-Format mit einem relativen Höchstfehler von  $4.9243 \cdot 10^{-24} = \varepsilon(1)$  berechnet. Die Funktion **LongNext** rundet diese Werte schließlich zur nächsten Longreal-Zahl, was mit einem zusätzlichen Höchstfehler von  $5 \cdot 10^{-21}$  verbunden ist:

$$\begin{aligned} \ln\_cos\_ (x) &= \ln(\cos(x)) \cdot (1+\varepsilon_1) \cdot (1+\varepsilon_{21}) \\ &= \ln(\cos(x)) \cdot (1+\varepsilon_2); \quad |\varepsilon_2| < 5.0050 \cdot 10^{-21} = \varepsilon(2) \end{aligned}$$

Damit haben wir in  $x < 1$  auch eine relative Fehlerschranke für das Longreal-Format berechnet.

**Fehlerabschätzung** für  $1 \leq x < \pi/2$ :

Mit den Bezeichnungen  $\tilde{y} := \cos\_rrG(x)$ ,  $y := \cos(x)$  gilt zunächst:

$$\begin{aligned} \tilde{y} &= \cos(x) \cdot (1+\varepsilon_y); \quad |\varepsilon_y| < 1.0196 \cdot 10^{-24} = \varepsilon(y) \\ \tilde{\ln}(\tilde{y}) &= \ln(\tilde{y}) \cdot (1+\varepsilon_{\ln}); \quad |\varepsilon_{\ln}| < 2.7329 \cdot 10^{-24} \end{aligned}$$

Mit der Abkürzung:  $\beta := \ln(1+\varepsilon_y)/\ln(y)$  erhält man folgende Gleichung:

$\ln(\tilde{y}) = \ln[y(1+\varepsilon_y)] = \ln(y) \cdot (1+\beta)$ , wirklich einfache Abschätzungen liefern:

$$|\beta| < \frac{\varepsilon(y)}{-\ln(\cos(1))} < 1.6563 \cdot 10^{-24} \rightsquigarrow$$

$$\begin{aligned} \tilde{\ln}(\tilde{y}) &= \ln(y) \cdot (1+\beta) \cdot (1+\varepsilon_{\ln}) \\ &= \ln(\cos(x)) \cdot (1+\varepsilon_3); \quad |\varepsilon_3| < 4.3893 \cdot 10^{-24} = \varepsilon(3) \end{aligned}$$

Für  $1 \leq x < \pi/2$  wird damit  $\ln(\cos(x))$  im 26-stelligen rrG-Format mit einem relativen Höchstfehler von  $4.3893 \cdot 10^{-24} = \varepsilon(3)$  berechnet. Die Funktion **LongNext** rundet diese Werte schließlich zur nächsten Longreal-Zahl, was mit einem zusätzlichen Höchstfehler von  $5 \cdot 10^{-21}$  verbunden ist:

$$\begin{aligned} \ln\_cos\_ (x) &= \ln(\cos(x)) \cdot (1+\varepsilon_3) \cdot (1+\varepsilon_{21}) \\ &= \ln(\cos(x)) \cdot (1+\varepsilon_4); \quad |\varepsilon_4| < 5.0044 \cdot 10^{-21} = \varepsilon(4) \end{aligned}$$

Damit haben wir in  $1 \leq x < \pi/2$  auch eine relative Fehlerschranke für das Longreal-Format berechnet.

Beachten Sie bitte, daß im Bereich  $0 \leq x < \pi/2$  die Gesamtfehlerschranke  $\varepsilon(\ln\_cos\_ ) = 5.005 \cdot 10^{-21}$  ihre Gültigkeit verliert, wenn Funktionswerte in den denormalisierten Longreal-Bereich fallen, d.h. wenn die Argumente betragsmäßig zu klein gewählt werden! Der Leser möge die obige Funktion  $\ln\_cos\_$  so verbessern, daß das Auftreten denormalisierter Funktionswerte durch eine entsprechende Fehlermeldung verhindert wird.

### 3.5.11 Mathematische Konstanten

Es existieren 23 mathematische Konstanten im 26-stelligen rrG-Format, die durch Aufruf entsprechender Funktionsnamen (ohne Argument) zur Verfügung stehen:

Mathematische Konstanten im rrG – Format			
Funktionsname	Wert	Fehlerschranke	Vergleich
PI_rrG	$\pi$	$5.323 \cdot 10^{-27}$	PI_rrG > $\pi$
PI2_rrG	$2\pi$	$5.323 \cdot 10^{-27}$	PI2_rrG > $2\pi$
PI_HALF_rrG	$\pi/2$	$5.323 \cdot 10^{-27}$	PI_HALF_rrG > $\pi/2$
PI_DIV4_rrG	$\pi/4$	$5.323 \cdot 10^{-27}$	PI_DIV4_rrG > $\pi/4$
R_PI_rrG	$1/\pi$	$8.011 \cdot 10^{-28}$	R_PI_rrG > $1/\pi$
R_2PI_rrG	$1/(2\pi)$	$2.341 \cdot 10^{-27}$	R_2PI_rrG < $1/(2\pi)$
SQRT_PI_rrG	$\sqrt{\pi}$	$9.399 \cdot 10^{-27}$	SQRT_PI_rrG > $\sqrt{\pi}$
SQRT_PI_D2_rrG	$\sqrt{\pi}/2$	$1.886 \cdot 10^{-27}$	SQRT_PI_D2_rrG < $\sqrt{\pi}/2$
R_SQRTPI_rrG	$1/\sqrt{\pi}$	$2.767 \cdot 10^{-27}$	R_SQRTPI_rrG < $1/\sqrt{\pi}$
R_SQRTPI2_rrG	$2/\sqrt{\pi}$	$2.767 \cdot 10^{-27}$	R_SQRTPI2_rrG < $2/\sqrt{\pi}$
R_SQRT2PI_rrG	$1/\sqrt{2\pi}$	$1.645 \cdot 10^{-28}$	R_SQRT2PI_rrG > $1/\sqrt{2\pi}$
LN_PI_rrG	$\ln(\pi)$	$4.250 \cdot 10^{-26}$	LN_PI_rrG > $\ln(\pi)$
SQRT2_rrG	$\sqrt{2}$	$2.977 \cdot 10^{-27}$	SQRT2_rrG > $\sqrt{2}$
R_SQRT2_rrG	$1/\sqrt{2}$	$2.977 \cdot 10^{-27}$	R_SQRT2_rrG < $1/\sqrt{2}$
LN_2_rrG	$\ln(2)$	$2.104 \cdot 10^{-27}$	LN_2_rrG < $\ln(2)$
LN_10_rrG	$\ln(10)$	$2.035 \cdot 10^{-27}$	LN_10_rrG < $\ln(10)$
R_LN_10_rrG	$1/\ln(10)$	$2.495 \cdot 10^{-27}$	R_LN_10_rrG > $1/\ln(10)$
EUL_rrG	$e$	$4.977 \cdot 10^{-28}$	EUL_rrG < $e$
R_EUL_rrG	$1/e$	$4.389 \cdot 10^{-28}$	R_EUL_rrG < $1/e$
SQRTEUL_rrG	$\sqrt{e}$	$7.392 \cdot 10^{-27}$	SQRTEUL_rrG > $\sqrt{e}$
GAMMA_C_rrG	$C$	$1.428 \cdot 10^{-28}$	GAMMA_C_rrG < $C$
LN_GAM_C_rrG	$\ln(C)$	$2.179 \cdot 10^{-27}$	LN_GAM_C_rrG < $C$
CATALAN_rrG	<i>Catalan</i>	$5.385 \cdot 10^{-27}$	CATALAN_rrG < <i>Catalan</i>

$$C := \lim_{m \rightarrow \infty} \left( \sum_{k=1}^m \frac{1}{k} - \ln(m) \right); \quad Catalan := \sum_{k=0}^{\infty} \frac{(-1)^k}{(2k+1)^2};$$

Im folgenden sind die numerischen Werte der Konstanten aus vorheriger Tabelle zusammengestellt:

<b>x :=</b>	<b>Komponenten – Werte (x.G=0)</b>	
PI_rrG	x.r1 = +3.141592653590	x.r2 = -2.067615373566E-13
PI2_rrG	x.r1 = +6.283185307180	x.r2 = -4.135230747132E-13
PI_H_rrG	x.r1 = +1.570796326795	x.r2 = -1.033807686783E-13
PI_DIV4_rrG	x.r1 = +0.7853981633974	x.r2 = +4.830961566085E-14
R_PI_rrG	x.r1 = +0.3183098861838	x.r2 = -9.328462232473E-15
R_2PI_rrG	x.r1 = +0.1591549430919	x.r2 = -4.664231116237E-15
SQRT_PI_rrG	x.r1 = +1.772453850906	x.r2 = -4.839727018325E-13
SQRT_PI_D2_rrG	x.r1 = +0.8862269254528	x.r2 = -4.198635091626E-14
R_SQRTPI_rrG	x.r1 = +0.5641895835478	x.r2 = -4.371305192055E-14
R_SQRTPI2_rrG	x.r1 = +1.128379167096	x.r2 = -4.874261038411E-13
R_SQRT2PI_rrG	x.r1 = +0.3989422804014	x.r2 = +3.267793994606E-14
LN_PI_rrG	x.r1 = +1.144729885849	x.r2 = +4.001741434274E-13
SQRT2_rrG	x.r1 = +1.414213562373	x.r2 = +9.504880168872E-14
R_SQRT2_rrG	x.r1 = +0.7071067811865	x.r2 = +4.752440084436E-14
LN_2_rrG	x.r1 = +0.6931471805599	x.r2 = +4.530941723212E-14
LN_10_rrG	x.r1 = +2.302585092994	x.r2 = +4.568401799145E-14
R_LN_10_rrG	x.r1 = +0.4342944819033	x.r2 = -4.817234887108E-14
EUL_rrG	x.r1 = +2.718281828459	x.r2 = +4.523536028747E-14
R_EUL_rrG	x.r1 = +0.3678794411714	x.r2 = +4.232159552377E-14
SQRTEUL_rrG	x.r1 = +1.648721270700	x.r2 = +1.281468486508E-13
GAMMA_C_rrG	x.r1 = +0.5772156649015	x.r2 = +3.286060651209E-14
LN_GAM_C_rrG	x.r1 = -0.5495393129816	x.r2 = -4.482233766177E-14
CATALAN_rrG	x.r1 = +0.9159655941772	x.r2 = +1.901505460351E-14

### 3.6 Der Datentyp lrG

Bei der Implementierung **hochgenauer** Standardfunktionen im Longreal-Format kann es z.B. bei der power-Funktion notwendig werden, Zwischenrechnungen mit mehr als 26 Dezimalstellen durchzuführen. Analog zum rrG-Format bietet sich dann an, eine Longreal- und real-Zahl mit insgesamt 34 Dezimalstellen aneinanderzuhängen, was durch den Datentyp **lrG** realisiert wird

```
Type lrG = record
    l : Longreal; r : real;
    G : integer
end;
```

der durch das Standardmodul `STDMOD` zur Verfügung gestellt wird. Durch `a : lrG` und

```
a.l := 1.23456789012345678901; a.r := 2.345678901234E-21; a.G := +1234567;
```

wird eine BCD-Zahl mit mindestens 34 Mantissenstellen dargestellt, deren Exponentenbereich nur durch die integer-Zahlen begrenzt ist:

$$a = 1.234567890123456789012345678901234 \cdot 10^{+1234567} \\ -2147483648 \leq a.G \leq +2147483647$$

Für die korrekte Durchführung der arithmetischen Grundoperationen muß ein lrG-Operand `a` die beiden folgenden Bedingungen erfüllen:

$$a.l = 0 \implies a.r = 0 \quad (\text{Beding. I})$$

$$|a.l| > 0 \implies \left| \frac{a.r}{a.l + a.r} \right| < 5.0001 \cdot 10^{-21} \quad (\text{Beding. II})$$

Die beiden obigen Bedingungen werden dabei automatisch realisiert, falls wie in allen praktischen Anwendungen `a.l` und `a.r` durch die folgenden Anweisungen aus dem Akkumulator ausgelesen werden:

```
a.l := LongNext( Akku ); { a : lrG; Akku: dotprecision }
Akku := Akku - a.l;     { a.l : Longreal; }
a.r := #*( Akku );     { a.r : real; }
```

Bedeutet `Akku` den numerischen Wert des Akkumulators, so gilt mit den obigen Anweisungen zusätzlich:

$$|a.l| \geq 10^{-233}, \quad (a.l + a.r) = \text{Akku} \cdot (1 + \varepsilon) \implies |\varepsilon| \leq 5 \cdot 10^{-34}$$



Das zweimalige Auslesen des Akkumulators erfolgt also mit dem relativen Höchstfehler  $\varepsilon(34) := 5 \cdot 10^{-34}$ , falls  $|a.l| \geq 10^{-233}$ . Nach den drei obigen **PASCAL-XSC** Anweisungen erreicht man daher mit der Abfrage

```
if (SIGN(Akku) <> 0) and (EXP0(a.l) < TEST_34)
    then Fehlermeldung
```

daß das zweimalige Auslesen des Akkumulators in die lrG-Komponenten a.l und a.r den relativen Fehler  $5 \cdot 10^{-34}$  nicht übersteigt; TEST\_34 = -232 wird vom Modul STDMOD als real-Konstante exportiert.

Der Anwender findet für den Datentyp **lrG** nur einen Minimalsatz von Hilfsfunktionen und Operatoren, die zur Realisierung der benötigten Standardfunktionen in 34-stelliger Dezimalarithmetik wirklich erforderlich waren. Da Intervallrechnungen im lrG-Format nicht vorgesehen sind, fehlen z.B. bei den Grundoperatoren alle gerichteten Rundungen.

### 3.6.1 Monadische Operatoren

Für den Datentyp **lrG** stehen die monadischen Operatoren  $+ -$  in gewohnter Weise zur Verfügung.

### 3.6.2 Addition, Subtraktion

Mit Hilfe der folgenden Prozeduren und Operatoren lassen sich Variablen vom Typ lrG addieren bzw. subtrahieren. Falls beide Operanden den gleichen Zehnerexponenten besitzen, sind aus Laufzeitgründen die Prozeduren den Operatoren vorzuziehen.

- procedure **ADD\_21** (xl: Longreal; x,y: real; var rl: Longreal; var r: real);  
 $(rl + r) = (xl + x + y) \cdot (1 + \varepsilon)$ ,  $|rl| \geq 10^{-233} \implies |\varepsilon| \leq 5 \cdot 10^{-34}$ ,  
 falls kein Overflow auftritt.  $|r/(rl + r)| < 5.0001 \cdot 10^{-21}$  wird erfüllt.
- procedure **ADD\_22** (xl: Longreal; x: real; yl: Longreal; y: real;  
 var rl: Longreal; var r: real);  
 $(rl + r) = (xl + x + yl + y) \cdot (1 + \varepsilon)$ ,  $|rl| \geq 10^{-233} \implies |\varepsilon| \leq 5 \cdot 10^{-34}$ ,  
 falls kein Overflow auftritt.  $|r/(rl + r)| < 5.0001 \cdot 10^{-21}$  wird erfüllt.
- operator  $+$  (x,y: lrG) s: lrG;  
 Vorauss.:  $|x.r/(x.l + x.r)| < 5.0001 \cdot 10^{-21}$   
 $|y.r/(y.l + y.r)| < 5.0001 \cdot 10^{-21}$   
 $s = (x + y) \cdot (1 + \varepsilon)$   
 $\implies |\varepsilon| \leq 5.0001 \cdot 10^{-34}$ ;  $|s.r/(s.l + s.r)| < 5.0001 \cdot 10^{-21}$
- operator  $-$  (x,y: lrG) d: lrG;  
 Vorauss.:  $|x.r/(x.l + x.r)| < 5.0001 \cdot 10^{-21}$   
 $|y.r/(y.l + y.r)| < 5.0001 \cdot 10^{-21}$   
 $d = (x - y) \cdot (1 + \varepsilon)$   
 $\implies |\varepsilon| \leq 5.0001 \cdot 10^{-34}$ ;  $|d.r/(d.l + d.r)| < 5.0001 \cdot 10^{-21}$

### 3.6.3 Multiplikation

Mit Hilfe der folgenden Prozeduren und Operatoren lassen sich Variablen vom Typ `lrG` multiplizieren. Falls die Summe der Zehnerexponenten beider Operanden verschwindet, sind aus Laufzeitgründen die Prozeduren den Operatoren vorzuziehen. Es existieren auch die wichtigsten Operatoren mit gemischten Operanden.

- procedure **MUL\_12** (x: real; yl: Longreal; y: real;  
var rl: Longreal; var r: real);  

$$(rl + r) = x \cdot (yl + y) \cdot (1 + \varepsilon), \quad |rl| \geq 10^{-233}$$

$$\implies |\varepsilon| \leq 5 \cdot 10^{-34}, \quad |r/(rl + r)| < 5.0001 \cdot 10^{-21} \text{ wird erfüllt.}$$
- procedure **MUL\_12** (xl: Longreal; yl: Longreal; y: real;  
var rl: Longreal; var r: real);  

$$(rl + r) = xl \cdot (yl + y) \cdot (1 + \varepsilon), \quad |rl| \geq 10^{-233}$$

$$\implies |\varepsilon| \leq 5 \cdot 10^{-34}, \quad |r/(rl + r)| < 5.0001 \cdot 10^{-21} \text{ wird erfüllt.}$$
- procedure **MUL\_22** (xl: Longreal; x: real; yl: Longreal; y: real;  
var rl: Longreal; var r: real);  
 Vorauss.:  $|x/(xl + x)| < 5.0001 \cdot 10^{-21}$   
 $|y/(yl + y)| < 5.0001 \cdot 10^{-21}$   
 $(rl + r) = (xl + x) \cdot (yl + y) \cdot (1 + \varepsilon), \quad |rl| \geq 10^{-233}$   

$$\implies |\varepsilon| \leq 5.0001 \cdot 10^{-34}, \quad |r/(rl + r)| < 5.0001 \cdot 10^{-21}$$
- operator **\*** (x,y: lrG) p: lrG;  
 Vorauss.:  $|x.r/(x.l + x.r)| < 5.0001 \cdot 10^{-21}$   
 $|y.r/(y.l + y.r)| < 5.0001 \cdot 10^{-21}$   
 $p = (x \cdot y) \cdot (1 + \varepsilon)$   

$$\implies |\varepsilon| \leq 5.0001 \cdot 10^{-34}, \quad |p.r/(p.l + p.r)| < 5.0001 \cdot 10^{-21}$$

### 3.6.4 Division

Mit Hilfe der folgenden Prozeduren und Operatoren lassen sich Quotienten vom Typ `lrG` berechnen. Falls die Zehnerexponenten beider `lrG`-Operanden übereinstimmen, sind aus Laufzeitgründen die Prozeduren den Operatoren vorzuziehen.

- procedure **DIV\_22** (xl: Longreal; x: real; yl: Longreal; y: real;  
var rl: Longreal; var r: real);  
 Vorauss.:  $|x/(xl + x)| < 5.0001 \cdot 10^{-21}$   
 $|y/(yl + y)| < 5.0001 \cdot 10^{-21}$   
 $(rl + r) = (xl + x)/(yl + y) \cdot (1 + \varepsilon), \quad |rl| \geq 10^{-233}$   

$$\implies |\varepsilon| \leq 5.0001 \cdot 10^{-34}, \quad |r/(rl + r)| < 5.0001 \cdot 10^{-21}$$
- operator **/** (x,y: lrG) q: lrG;

$$\begin{aligned}
\text{Voraus.: } & |x.r/(x.l + x.r)| < 5.0001 \cdot 10^{-21} \\
& |y.r/(y.l + y.r)| < 5.0001 \cdot 10^{-21} \\
& q = (x/y) \cdot (1 + \varepsilon) \\
\Rightarrow & |\varepsilon| \leq 5.0002 \cdot 10^{-34}; \quad |q.r/(q.l + q.r)| < 5.0001 \cdot 10^{-21}
\end{aligned}$$

### 3.6.5 Vergleiche

Mit Hilfe der beiden folgenden Funktionen sind lediglich zwei Vergleiche vorgesehen:

- function **greater\_21** (xl: Longreal; x,a: real): boolean;

$$\text{Voraus.: } |x/(xl + x)| < 5.0001 \cdot 10^{-21}$$

$$\text{greater\_21} = \text{true} \iff (xl + x) > a$$

- function **less\_21** (xl: Longreal; x,a: real): boolean;

$$\text{Voraus.: } |x/(xl + x)| < 5.0001 \cdot 10^{-21}$$

$$\text{less\_21} = \text{true} \iff (xl + x) < a$$

### 3.6.6 Transferfunktionen und Zuweisungsoperatoren

Die folgenden Prozeduren und Operatoren ermöglichen die Verbindung zwischen den Datentypen **rrG**, **lrG**:

- procedure **Trans** (rl: Longreal; r: real; var x1,x2: real);

$$\text{Voraus.: } |r/(rl + r)| < 5.0001 \cdot 10^{-21}; \quad (x1 + x2) = (rl + r) \cdot (1 + \varepsilon)$$

$$\Rightarrow |\varepsilon| \leq 5 \cdot 10^{-26}; \quad |x2/(x1 + x2)| < 5.0001 \cdot 10^{-13}$$

- procedure **Trans** (x1,x2: real; var rl: Longreal; var r: real);

$$\text{Voraus.: } |x2/(x1 + x2)| < 5.0001 \cdot 10^{-13}; \quad (rl + r) = (x1 + x2) \cdot (1 + \varepsilon)$$

$$\Rightarrow |\varepsilon| = 0; \quad |r/(rl + r)| < 5.0001 \cdot 10^{-21}$$

- operator := (var a: rrG; b: lrG);

$$\text{Voraus.: } |b.r/(b.l + b.r)| < 5.0001 \cdot 10^{-21}; \quad a = b \cdot (1 + \varepsilon)$$

$$\Rightarrow |\varepsilon| \leq 5 \cdot 10^{-26}; \quad |a.r2/(a.r1 + a.r2)| < 5.0001 \cdot 10^{-13}$$

- operator := (var a: lrG; b: rrG);

$$\text{Voraus.: } |b.r2/(b.r1 + b.r2)| < 5.0001 \cdot 10^{-13}; \quad a = b \cdot (1 + \varepsilon)$$

$$\Rightarrow |\varepsilon| = 0; \quad |a.r/(a.l + a.r)| < 5.0001 \cdot 10^{-21}$$

### 3.6.7 Hilfsfunktionen

Bei den auftretenden Funktionsparametern vom Typ **lrG** wird vorausgesetzt, daß die Bedingungen I,II von Seite 66 erfüllt sind. Dies wird immer dann der Fall sein, wenn die Komponenten dieser Parameter durch das auf der gleichen Seite beschriebene zweimalige Auslesen des Akkumulators ins Longreal- bzw. real-Format definiert sind. Besitzen Funktionen den Ergebnistyp lrG, so erfüllen auch die Komponenten aller Funktionswerte automatisch die Bedingungen I,II von Seite 66.

- function **Expo** (x: lrG) : integer;  
Berechnet den Zehnerexponenten analog zu `Expo(x: real): integer`;
- procedure **Mant\_34** (xl: Longreal; x: real; var rl: Longreal; var r: real);  
Berechnet die Mantisse von  $(xl + x)$  analog zu `Mant(x: real): real`, falls  $|x/(xl + x)| < 5.0001 \cdot 10^{-21}$  erfüllt ist.  
Es gilt also:  $|rl + r| = 0$  oder  $0.1 \leq |rl + r| < 1$ ;
- function **Scale** (x: lrG): lrG;  
`y := Scale(x)` liefert folgende Ergebnisse:
  1.  $y = x$ , ' = ' im Sinne der Numerik;
  2.  $y.l = 0 \implies y.r = 0, y.G = -\text{Maxint} = -2147483647$ ;
  3.  $0.1 \leq |y.l| < 1$ , falls  $y.l \neq 0$ ; (Skalierung);
  4.  $|y.r/(y.l + y.r)| < 5.0001 \cdot 10^{-21}$ , falls  $y.l \neq 0$ ;

Falls das Ergebnis  $y = x$  in extremen Ausnahmefällen nicht erfüllt werden kann, erfolgt eine entsprechende Fehlermeldung.
- function **Transform** (x: lrG): lrG;  
Voraus.: `x.l, x.r, x.G` seien beliebig;  
Durch `y := Transform(x)` wird `x` so transformiert, daß gilt:
  1.  $y = x$ ;
  2.  $y.l = 0 \implies y.r = 0$ ;
  3.  $y.l \neq 0 \implies |y.r/(y.l + y.r)| < 5.0001 \cdot 10^{-21}$ ;

### 3.6.8 Horner Schema

Gegeben sei das Polynom vom Grade  $N = 9$ :

$$P_9(x) = (bl[0] + b[0]) + (bl[1] + b[1]) \cdot x^1 + a[2] \cdot x^2 + \dots + a[9] \cdot x^9$$

das für das 34-stellige BCD-Argument  $x = (xl + x)$  nach dem Horner Schema auszuwerten ist. Die Zahlen  $xl, x$  müssen folgende Bedingungen erfüllen:

$$xl = 0 \implies x = 0; \quad xl \neq 0 \implies |x/(xl + x)| < 5.0001 \cdot 10^{-21}$$

Aus Laufzeitgründen werden die ersten  $gr = 7$  groben Horner Schritte im Longrealformat durchgeführt:

$$\text{var a: Lrvector}[2..9]; \quad \{ \text{LB}(a) := N - gr = 2; \quad \text{UB}(a) := N = 9; \}$$

Im achten, gemischten Schritt wird die Multiplikation im Longrealformat ausgeführt und die Addition im 34-stelligen BCD-Format. Alle restlichen Horner Schritte laufen im 34-stelligen BCD-Format wie folgt ab:

- Die Multiplikation wird im Akkumulator ausgeführt, wobei aus Laufzeitgründen das betragsmäßig kleinste Teilprodukt weggelassen wird. Da der Akku nach dieser Multiplikation nicht ins (Longreal|real)-Format ausgelesen wird, gilt für die relative Fehlerschranke:  $\varepsilon(\text{Mul}) = 2.5002 \cdot 10^{-41}$ ;
- Auf das obige Multiplikationsergebnis wird im Akku der jeweilige Polynomkoeffizient rundungsfehlerfrei addiert. Diese Summe wird anschließend in das (Longreal|real)-Format ausgelesen, weshalb die Fehlerschranke der Addition mit  $\varepsilon(\text{Add}) = 5 \cdot 10^{-34}$  anzusetzen ist. Dabei muß nach Seite 66 vorausgesetzt werden, daß für jedes Ergebnis `erg` eines Hornerschrittes gilt:  $|\text{erg}| \geq 10^{-233}$ .  
Diese Bedingung ist in allen praktischen Fällen erfüllt, wenn  $|xl + x|$  nicht zu klein gewählt wird und wenn das Polynom nicht in der Nähe einer seiner Nullstellen ausgewertet wird. Eine genaue Überprüfung kann mit dem **PASCAL-XSC** Programm **HORN\_TS** vorgenommen werden.

```
var bl: Lrvector[0..1]; b: rvector[0..1];
{ LB(bl) := 0; UB(bl) := N - gr - 1 = 1; }
```

Die oben beschriebene Polynomauswertung erfolgt durch:

```
procedure HORNER (a,bl: Lrvector; b: rvector; xl: Longreal; x: real;
var rl: Longreal; var r: real);
```

Der Aufruf obiger Prozedur mit den deklarierten Eingangsparametern `a,bl,b,xl,x` liefert mit den Ausgangsgrößen `rl,r` eine Näherung  $(rl + r) \approx P_9(xl + x)$  für den Polynomwert  $P_9(xl + x)$ . Es gilt:

- $rl = 0 \implies r = 0$ ;  $rl \neq 0 \implies |r/(rl + r)| < 5.0001 \cdot 10^{-21}$

Die Berechnung von garantierten Oberschranken des absoluten oder relativen Polynomauswertefehlers, der beim Aufruf der Prozedur **HORNER** entsteht, ist mit folgenden Programmen möglich:

```
IFO_BCDS   IEEE-Programm, das Polynominformationen sammelt;
AWF_BCDS   IEEE-Programm, das den Auswertefehler abschätzt;
HORN_TS    BCD-Programm überprüft AWF_BCDS-Ergebnisse.
```

### 3.6.9 Ein-/Ausgabeprozeduren

Es stehen die gewohnten Prozeduren

```
procedure read (var f: text; var a: lrG);
procedure write (var f: text; a: lrG);
```

mit optionalem Fileparameter, beliebig vielen Ein-/Ausgabeparametern, jedoch ohne Formatspezifikationen zur Verfügung. `a: lrG` muß in der Form

```
{a.l,a.r,a.G}    oder in der Form    a.l
```

eingegeben werden. Im zweiten Fall gilt:  $a.r = a.G = 0$ . Die Komponente  $a.G$  erhält jedoch stets den Wert  $-\text{maxint} = -2147483647$ , falls  $a.l$  verschwindet.

Erfolgt die Eingabe in der ersten Form mit beliebigen Werten  $a.l, a.r$ , so wird auf  $a:lrG$  automatisch die Funktion  $\text{Transform}(a:lrG):lrG$  angewandt, so daß ohne Änderung des numerischen Werts von  $a$  die beiden folgenden Bedingungen von Seite 66 erfüllt werden:

$$a.l = 0 \implies a.r = 0 \quad (\text{Beding. I})$$

$$|a.l| > 0 \implies \left| \frac{a.r}{a.l + a.r} \right| < 5.0001 \cdot 10^{-21} \quad (\text{Beding. II})$$

Die Ausgabe einer  $lrG$ -Zahl erfolgt stets in der Form:

$$\{a.l, a.r, a.G\}$$

### Beispiele:

Sei  $a$  vom Typ  $lrG$ , dann wird mit den Anweisungen `read(a); write(a)` und der Eingabe

$$\{1.23456789012345678901E35, 5.934567890123E14, -1234567890\}$$

die  $lrG$ -Zahl  $a$  in der folgenden Form ausgegeben:

$$\{1.23456789012345678902E+035, -4.065432109877E+014, -1234567890\}$$

Die Eingabe

$$\{4e+500, -4e + 500, 1234567890\}$$

ergibt mit den gleichen obigen Anweisungen die Ausgabe des Wertes Null:

$$\{0.000000000000E+000, 0.000000000000E+000, -2147483647\}$$

Die Eingabe

$$\{9e+511, 9e+511, +2147483647\}$$

ergibt mit der obigen `read`-Anweisung eine Fehlermeldung, da der numerische Eingabewert wegen Overflow im  $lrG$ -Format nicht darstellbar ist.

Eine entsprechende Underflow-Fehlermeldung läßt sich mit den `read`- und `write`-Anweisungen jedoch nicht konstruieren!

## 3.6.10 Standardfunktionen

Standardfunktionen vom Typ lrG; x: lrG			
Funktion	Aufruf	Argumentbereich	Rel. Fehler
$\sqrt{x}$	sqrt(x)	$x \geq 0$	$6.6850 \cdot 10^{-31}$
$\sqrt{1-x^2}$	sqrt1mx2(x)	$ x  \leq 1$	$6.6926 \cdot 10^{-31}$
$\sqrt{1+x^2}$	sqrt1px2(x)	rrG-Bereich	$6.6901 \cdot 10^{-31}$
$\sqrt{x^2-1}$	sqrtx2m1(x)	$ x  \geq 1$	$6.6926 \cdot 10^{-31}$
$\ln(x)$	ln(x)	$10^{-5000} \leq x < 10^{+5001}$	$2.3299 \cdot 10^{-30}$
$\log_{10}(x)$	log10(x)	$10^{-5000} \leq x < 10^{+5001}$	$2.3310 \cdot 10^{-30}$

Standardfunktionen vom Typ lrG; x: Longreal			
Funktion	Aufruf	Argumentbereich	Rel. Fehler
$\sqrt{1-x^2}$	sqrt1mx2_lrG(x)	$ x  \leq 1$	$6.6876 \cdot 10^{-31}$
$\sqrt{1+x^2}$	sqrt1px2_lrG(x)	Longreal-Bereich	$6.6927 \cdot 10^{-31}$
$\sqrt{x^2-1}$	sqrtx2m1_lrG(x)	$ x  \geq 1$	$6.6876 \cdot 10^{-31}$

Die Komponenten  $x.l, x.r$  eines Funktionsarguments  $x$ : lrG müssen die beiden Bedingungen von Seite 66 erfüllen:

$$x.l = 0 \implies x.r = 0 \quad (\text{Beding. I})$$

$$|x.l| > 0 \implies \left| \frac{x.r}{x.l + x.r} \right| < 5.0001 \cdot 10^{-21} \quad (\text{Beding. II})$$

Ist  $y$ : lrG das Ergebnis einer der obigen Funktionen, so gelten die beiden obigen Bedingungen **automatisch** auch für die Komponenten  $y.l, y.r$ .





# Kapitel 4

## Intervallrechnung

Intervallarithmetiken existieren bzgl. der Grunddatentypen `real`, `Longreal`, `rrG` und `dotprecision`.

### 4.1 Intervallrechnung mit dem Typ `real`

Mit dem Modul `I_ARI` werden die für das Rechnen mit reellen Intervallen notwendigen Operatoren, Funktionen und Prozeduren bereitgestellt.

#### 4.1.1 Der Datentyp `interval`

Der Datentyp `interval` ist definiert durch:

```
type interval = record inf, sup : real end;
```

Mit der Deklaration: `var x : interval;` bedeutet `x` dann das reelle Intervall:

```
x = [x.INF, x.SUP]
```

#### 4.1.2 Operatoren

Sämtliche in diesem Modul vordefinierten arithmetischen und Verbands-Operatoren liefern den Ergebnistyp `interval`. Als arithmetische Operatoren stehen die monadischen Operatoren `+`, `-` und die vier Grundoperationen `+`, `-`, `*`, `/` mit der Rundung zum kleinsten einschließenden Intervall vom Typ `interval` zur Verfügung.

Ist bei den Grundoperatoren `+`, `-`, `*`, `/` ein Operand vom Typ `interval`, so sind für den zweiten Operanden die folgenden Typen möglich:

**interval, integer oder real.**

Die Vergleichsoperatoren `=`, `<>`, `<`, `<=`, `>`, `>=` sind mengentheoretisch zu interpretieren. Dabei bedeutet:

<code>=</code> gleich	<code>&lt;&gt;</code> ungleich	<code>&lt;</code> echte Teilmenge von
<code>&lt;=</code> Teilmenge von	<code>&gt;</code> echte Obermenge von	<code>&gt;=</code> Obermenge von

Der Operator **in** steht für die Relation „liegt in“ zwischen einem integer/real- und einem *interval*-Operanden oder für die Relation „echt enthalten in“ zwischen zwei *interval*-Operanden zur Verfügung.

Es gilt:

$$\begin{aligned} x \quad \mathbf{in} \quad y &\iff (x.\text{inf} > y.\text{inf}) \quad \mathbf{and} \quad (x.\text{sup} <> y.\text{sup}). \\ 3.2 \quad \mathbf{in} \quad y &\iff (3.2 \geq y.\text{inf}) \quad \mathbf{and} \quad (3.2 \leq y.\text{sup}). \end{aligned}$$

Der Operator **><** testet auf Disjunktheit zweier Intervalle. Dabei heißen zwei Intervalle  $x, y$  disjunkt, wenn gilt  $x \cap y = \emptyset$  (leere Menge).

Die Verbandsoperatoren **+\*** bzw. **\*\*** bezeichnen die Bildung der Intervallhülle bzw. des Durchschnitts, d.h. der Operator **+\*** liefert das kleinste, beide Operanden umfassende Intervall, und der Operator **\*\*** liefert das Schnittintervall. Ein leerer Schnitt führt zu einem Laufzeitfehler.

rechter Operand	integer real	interval	Bemerkungen
linker Operand			
monadisch		+, -	
integer real	+*	◇ <b>in</b> , =, <> +*	◇ ∈ {+, -, *, /}
interval	◇ =, <> +*	◇ <b>in</b> , ∨, >< +*, **	∨ ∈ {=, <>, <, <=, >, >=}

Tabelle 4.1: Die Operatoren des Moduls `L_ARI`

### Beispiele:

Seien  $a, b$  vom Typ *interval* mit:  $a = [-1, 3]$ ;  $b = [3, 4]$ , dann liefern die Operatoren **+**, **-**, **\***, **/**, **><**, **+\***, **\*\*** die folgenden Ergebnisse:

Ausdruck	Ergebnis	Ausdruck	Ergebnis
$a + b$	$[2, 7]$	$a +* b$	$[-1, 4]$
$a - b$	$[-5, 0]$	$a ** b$	$[3, 3]$
$a * b$	$[-4, 12]$	$a >< b$	<i>false</i>
$a / b$	$[-0.33333333333334, 1]$		

### 4.1.3 Transferfunktionen

Zur Wandlung zwischen den Typen *real* und *interval* werden folgende Transferfunktionen bereitgestellt:

Funktion	Ergebnistyp	Bedeutung
intval (r1,r2)	<i>interval</i>	Intervall mit $inf = r1$ und $sup = r2$ ; $r1 \leq r2$
intval (r)	<i>interval</i>	Punktintervall mit $inf = sup = r$
inf (x)	<i>real</i>	Untergrenze von x
sup (x)	<i>real</i>	Obergrenze von x

$r, r1, r2 = real$ -Ausdruck;  $x = interval$ -Ausdruck

### 4.1.4 Überladungen des Zuweisungsoperators

Die Wandlung von *real* nach *interval* wird auch in Form einer überladenen Zuweisung bereitgestellt:

Zuweisung	Bedeutung
$x := r$	$x := \text{intval}(r)$

$x = interval$ -Variable;  $r = real$ -Ausdruck

### 4.1.5 Standardfunktionen

Alle für *interval*-Argumente verfügbaren, hochgenauen Standardfunktionen sind in Tabelle 4.2 zusammengestellt.

Darüberhinaus sind Funktionen für die Berechnung des Mittelpunktes, des Vorzeichens und des Durchmessers von Intervallen verfügbar:

- Intervall-Vorzeichen:  $\text{sign}(x) = \begin{cases} -1 & x.\text{sup} < 0 \\ 0 & 0 \in x \\ +1 & x.\text{inf} > 0 \end{cases}$
- Intervall-Durchmesser:  $\text{diam}(x) := \text{sup}(x) - \text{inf}(x)$ ;
- Die Mittelpunktsfunktion ergibt sich aus der Anweisung:  
 $\text{mid}(x) := \# * (0.5 * \text{inf}(x) + 0.5 * \text{sup}(x))$ ;

Funktion	Aufruf	Funktion	Aufruf
$ x $	<code>abs(x)</code>	$\cos(x)$	<code>cos(x)</code>
$x^2$	<code>sqr(x)</code>	$\tan(x)$	<code>tan(x)</code>
$x^2 - y^2$	<code>x2_y2(x, y)</code>	$\cot(x)$	<code>cot(x)</code>
$\sqrt{x}$	<code>sqrt(x)</code>	$\sin(\pi \cdot x)$	<code>sinpi(x)</code>
$\sqrt{1+x} - 1$	<code>sqrt1pm1(x)</code>	$\cos(\pi \cdot x)$	<code>cospi(x)</code>
$\sqrt{x^2 + y^2}$	<code>sqrtx2y2(x, y)</code>	$\tan(\pi \cdot x)$	<code>tanpi(x)</code>
$\sqrt{1+x^2}$	<code>sqrt1px2(x)</code>	$\cot(\pi \cdot x)$	<code>cotpi(x)</code>
$\sqrt{1-x^2}$	<code>sqrt1mx2(x)</code>	$\arcsin(x)$	<code>arcsin(x)</code>
$\sqrt{x^2 - 1}$	<code>sqrtx2m1(x)</code>	$\arccos(x)$	<code>arccos(x)</code>
$e^x$	<code>exp(x)</code>	$\arctan(x)$	<code>arctan(x)</code>
$e^x - 1$	<code>expm1(x)</code>	$\arctan2(x, y)$	<code>arctan2(x, y)</code>
$e^{-x^2}$	<code>exp_x2(x)</code>	$\operatorname{arccot}(x)$	<code>arccot(x)</code>
$2^x$	<code>exp2(x)</code>	$\sinh(x)$	<code>sinh(x)</code>
$10^x$	<code>exp10(x)</code>	$\cosh(x)$	<code>cosh(x)</code>
$x^y$	<code>power(x, y)</code>	$\tanh(x)$	<code>tanh(x)</code>
$\ln(x)$	<code>ln(x)</code>	$\operatorname{coth}(x)$	<code>coth(x)</code>
$\ln(e \cdot x)$	<code>ln_ex(x)</code>	$\operatorname{arsinh}(x)$	<code>arsinh(x)</code>
$\ln(1+x)$	<code>ln1p(x)</code>	$\operatorname{arcosh}(x)$	<code>arcosh(x)</code>
$\ln(x^2 + y^2)/2$	<code>ln_sqrtx2y2(x, y)</code>	$\operatorname{arcosh}(1+x)$	<code>arcosh1p(x)</code>
$\ln[(1+x)^2 + y^2]/2$	<code>ln_sqrt1px2y2(x, y)</code>	$\operatorname{artanh}(x)$	<code>artanh(x)</code>
$\log_a(x)$	<code>loga(a, x)</code>	$\operatorname{arcoth}(x)$	<code>arcoth(x)</code>
$\log_2(x)$	<code>log2(x)</code>	$\operatorname{sign}(x)$	<code>sign(x)</code>
$\log_{10}(x)$	<code>log10(x)</code>	$\operatorname{diam}(x)$	<code>diam(x)</code>
$\log_{10}(1+x)$	<code>log1p(x)</code>	$\operatorname{mid}(x)$	<code>mid(x)</code>
$\sin(x)$	<code>sin(x)</code>	$\operatorname{blow}(x, \text{eps})$	<code>blow(x, eps)</code>

Tabelle 4.2: Mathematische Standardfunktionen des Moduls `LARI`;  $x, y$ : *interval*-Ausdrücke;  $\text{eps}, a$ : *real*-Ausdrücke;

- Mit Ausnahme von  $\operatorname{sign}(x)$ ,  $\operatorname{diam}(x)$ , und  $\operatorname{mid}(x)$  liefern alle anderen Funktionen aus Tabelle 4.2 ein Ergebnis vom Typ *interval* mit **hochgenau** berechneten Intervallgrenzen, d.h. nur in ganz seltenen Ausnahmefällen wird z.B. die untere Intervallgrenze um 1 ulp zu klein berechnet.

- Bei Einschließungsmethoden wird die Funktion `blow` für die sogenannte Epsilonaufblähung eines Intervalls benötigt. Mit einem *real*-Ausdruck `eps` und einem *interval*-Ausdruck `x` wird die `blow(x, eps)`-Funktion aufgebaut durch die Anweisungen:

$$y := (1 + \text{eps}) * x - \text{eps} * x;$$

$$\text{blow } (x, \text{eps}) := \text{intval } (\text{pred } (\text{inf}(y)), \text{succ } (\text{sup}(y)));$$

- $y := \text{abs}(x); \quad \{|r| \mid r \in x \cap \mathbb{R}\} = y$
- $y := \text{sqr}(x); \quad \{r^2 \mid r \in x \cap \mathbb{R}\} \subseteq y$
- $y := \text{sqrt}(x); \quad \{\sqrt{r} \mid r \in x \cap \mathbb{R}, \text{inf}(x) \geq 0\} \subseteq y$

#### Beispiele:

Seien  $a, b, c$  vom Typ *interval*, erzeugt durch die Anweisungen:

$$a := \text{intval}(-1, 3); \quad b := \text{intval}(1, 4) \quad c := \text{intval}(-5, -4);$$

dann liefern die Funktionen:

$$\text{abs, sqr, mid, sign, sqrt, diam, x2\_y2, sqrtx2y2, arctan2}$$

die folgenden Ergebnisse:

Ausdruck	Ergebnis	Ausdruck	Ergebnis
<code>abs (a)</code>	[0, 3]	<code>sqr (b)</code>	[1; 16]
<code>abs (b)</code>	[1, 4]	<code>mid (a)</code>	1
<code>sqr (a)</code>	[0, 9]	<code>diam (a)</code>	4
<code>sign (a)</code>	0	<code>sign (b)</code>	1
<code>sign (c)</code>	-1	<code>sqrt (b)</code>	[1, 2]
<code>x2_y2 (a,b)</code>	[-16, 8]	<code>sqrtx2y2 (a,b)</code>	[1, 5]
<code>arctan2 (b,a)</code>	[-1.57..95; 1.57..95]	<code>arctan2 (a,b)</code>	[-0.78..75; 1.24..99]

#### 4.1.6 Ein-/AusgabeprozEDUREN

Es stehen die bekannten Prozeduren

```
procedure read (var f: text; var a: interval);
procedure write (var f: text; a: interval);
```

mit optionalem Fileparameter, beliebig vielen Ein-/Ausgabeparametern, jedoch ohne Formatspezifikationen zur Verfügung.

Die Eingabe eines Intervalls  $x = [a, b]$  muß in der Form

$$[a, b] \quad \text{oder in der Form} \quad a$$

erfolgen. Der zweite Fall dient zur vereinfachten Eingabe eines Punktintervalls  $x = [a, a]$ .

Die Ausgabe eines Intervalls durch die write-Prozedur erfolgt stets in der Form

$$[a, b]$$

Das Ausgabeformat ist dabei durch folgende Regelung bestimmt:

Ist  $x$  eine Variable vom Typ *interval* und stimmen die Werte von  $\inf(x)$  und  $\sup(x)$  auf den ersten  $n \geq 2$  Ziffern überein, so gibt die Anweisung  $\text{write}(x)$  die Unterschranke  $a$  und die Oberschranke  $b$  jeweils auf  $n$  gerundete Mantissenstellen aus.

Für  $n = 0$  und  $n = 1$  werden  $a$  und  $b$  stets auf zwei Mantissenstellen gerundet.

Werden also durch  $\text{write}(x)$   $n > 2$  Mantissenstellen ausgegeben, so weiß der Anwender, daß  $\inf(x)$  und  $\sup(x)$  genau in den ersten  $n$  Ziffern übereinstimmen!

In folgender Tabelle sind für den Wert einer Variablen  $x$  vom Typ *interval* die jeweiligen Ergebnisse der Ausgabeprozedur  $\text{write}(x)$  zusammengestellt:

$x = [\inf(x), \sup(x)]$	$\text{write}(x)$
[2, 3]	[2.0E+000, 3.0E+000]
[1.1, 1.2]	[1.1E+000, 1.2E+000]
[1.11, 1.13]	[1.1E+000, 1.2E+000]
[2.000, 2.0001]	[2.000E+000, 2.001E+000]

Die *real*-Werte  $\inf(x)$  und  $\sup(x)$  können natürlich auch (**ohne** zusätzliche Rundungen durch die  $\text{write}(x)$ -Prozedur) wie folgt ausgegeben werden:

- $\text{write}(\inf(x), \sup(x));$
- $\text{write}(x.\text{inf}, x.\text{sup});$

$x := \text{intval}(1.11, 1.13)$  und  $\text{write}(x.\text{inf}, x.\text{sup})$  ergeben:

$$[1.110000000000E+000, 1.130000000000E+000]$$

im Gegensatz zur entsprechenden Ausgabe von  $\text{write}(x)$  in obiger Tabelle.

## 4.2 Intervallrechnung mit dem Typ Longreal

Mit dem Modul **LI\_ARI** werden die für das Rechnen mit reellen Intervallen notwendigen Operatoren, Funktionen und Prozeduren bereitgestellt.

### 4.2.1 Der Datentyp Linterval

Der Datentyp *Linterval* ist definiert durch:

```
type Linterval = record inf, sup : Longreal end;
```

Mit der Deklaration: **var** x : Linterval; bedeutet x dann das reelle Intervall:

```
x = [x.INF, x.SUP]
```

### 4.2.2 Operatoren

Wir betrachten zunächst alle im Modul LI\_ARI vordefinierten arithmetischen und Verbands-Operatoren mit dem Ergebnistyp *Linterval*. Als arithmetische Operatoren stehen die monadischen Operatoren  $+$ ,  $-$  und die vier Grundoperationen  $+$ ,  $-$ ,  $*$ ,  $/$  mit der Rundung zum jeweils kleinsten einschließenden Intervall vom Typ *Linterval* zur Verfügung.

Ist bei den Grundoperatoren  $+$ ,  $-$ ,  $*$ ,  $/$  ein Operand vom Typ *Linterval*, so sind für den zweiten Operanden die folgenden Typen möglich:

**integer, real, Longreal, interval** oder **Linterval**.

Die Vergleichsoperatoren  $=$   $<>$   $<$   $<=$   $>$   $>=$  sind mengentheoretisch zu interpretieren. Dabei bedeutet:

$=$ gleich	$<>$ ungleich	$<$ echte Teilmenge von
$<=$ Teilmenge von	$>$ echte Obermenge von	$>=$ Obermenge von

Der Operator **in** steht wieder für die Relation „liegt in“ zwischen einem integer/real/Longreal- und einem *Linterval*-Operanden oder für die Relation „echt enthalten in“ zwischen zwei *Linterval*-Operanden zur Verfügung. Es gilt:

$$\begin{aligned} x \text{ in } y &\iff (x.\text{inf} > y.\text{inf}) \text{ and } (x.\text{sup} <> y.\text{sup}). \\ 3.2 \text{ in } y &\iff (3.2 \geq y.\text{inf}) \text{ and } (3.2 \leq y.\text{sup}). \end{aligned}$$

Der Operator  $><$  testet auf Disjunktheit zweier Intervalle. Dabei heißen zwei Intervalle  $x, y$  disjunkt, wenn gilt  $x \cap y = \emptyset$  (leere Menge).

Die Verbandsoperatoren  $+*$  bzw.  $**$  bezeichnen die Bildung der Intervall-Hülle bzw. des Durchschnitts, d.h. der Operator  $+*$  liefert das kleinste, beide Operanden umfassende Intervall, und der Operator  $**$  liefert das Schnittintervall. Ein leerer Schnitt führt zu einem Laufzeitfehler.

<b>rechter Operand</b>	integer	Linterval	Bemerkungen
	real		
	Longreal	interval	
<b>linker Operand</b>	Ergebnistypen: <i>Linterval</i> oder <i>boolean</i>		
monadisch		+, -	
integer		◇	◇ ∈ {+, -, *, /}
real		<b>in</b> , =, <>	
Longreal		+*	
Linterval	◇	◇	V ∈ {=, <>, <, <=, >, >=}
interval	=, <>	<b>in</b> , V, ><	
	+*	+*, **	

**Beispiele:**

Seien  $a, b$  vom Typ *Linterval* mit:  $a = [-1, 3]$ ;  $b = [3, 4]$ , dann liefern die Operatoren  $+, -, *, /, ><, +*, **$  die folgenden Ergebnisse:

Ausdruck	Ergebnis	Ausdruck	Ergebnis
$a + b$	$[2, 7]$	$a +* b$	$[-1, 4]$
$a - b$	$[-5, 0]$	$a ** b$	$[3, 3]$
$a * b$	$[-4, 12]$	$a >< b$	<i>false</i>
$a / b$	$[-0.33333333333333333333333333333334, 1]$		

**4.2.3 Transferfunktionen**

In nachfolgender Tabelle sind alle die Funktionen angegeben, die den Zusammenhang zwischen dem Typ *Linterval* einerseits und den Typen *integer*, *real*, *Longreal* und *interval* andererseits herstellen:

Funktion	Ergebnistyp	Bedeutung
SHORT (Lx)	<i>interval</i>	$Lx \subseteq \text{SHORT}(Lx)$
LONG (x)	<i>Linterval</i>	$\text{Long}(x) \equiv x$ , aber Datentypen verschieden!
LINTVAL (a,b)	<i>Linterval</i>	Intervall mit $\text{inf} = a$ und $\text{sup} = b$ ; $a \leq b$
LINTVAL (a)	<i>Linterval</i>	Punktintervall mit $\text{inf} = \text{sup} = a$
INF (Lx)	<i>Longreal</i>	Untergrenze von Lx
SUP (Lx)	<i>Longreal</i>	Obergrenze von Lx

$a, b = \text{integer/real/Longreal}$ -Ausdruck  
 $Lx = \text{Linterval}$ -Ausdruck       $x = \text{interval}$ -Ausdruck;



### 4.2.4 Standardfunktionen

Dem Anwender stehen die in nachfolgender Tabelle 4.3 zusammengestellten hochgenauen, mathematischen Intervall-Standardfunktionen zur Verfügung. Argumente  $x, y$  vom Typ `Linterval` liefern Funktionswerte ebenfalls vom Typ `Linterval`.

Funktion	Aufruf	Funktion	Aufruf
$ x $	<code>abs(x)</code>	$\cos(x)$	<code>cos(x)</code>
$x^2$	<code>sqr(x)</code>	$\tan(x)$	<code>tan(x)</code>
$x^2 - y^2$	<code>x2_y2(x, y)</code>	$\cot(x)$	<code>cot(x)</code>
$\sqrt{x}$	<code>sqrt(x)</code>	$\sin(\pi \cdot x)$	<code>sinpi(x)</code>
$\sqrt{1+x} - 1$	<code>sqrt1pm1(x)</code>	$\cos(\pi \cdot x)$	<code>cospi(x)</code>
$\sqrt{x^2 + y^2}$	<code>sqrtx2y2(x, y)</code>	$\tan(\pi \cdot x)$	<code>tanpi(x)</code>
$\sqrt{1+x^2}$	<code>sqrt1px2(x)</code>	$\cot(\pi \cdot x)$	<code>cotpi(x)</code>
$\sqrt{1-x^2}$	<code>sqrt1mx2(x)</code>	$\arcsin(x)$	<code>arcsin(x)</code>
$\sqrt{x^2 - 1}$	<code>sqrtx2m1(x)</code>	$\arccos(x)$	<code>arccos(x)</code>
$e^x$	<code>exp(x)</code>	$\arctan(x)$	<code>arctan(x)</code>
$e^x - 1$	<code>expm1(x)</code>	$\arctan2(x, y)$	<code>arctan2(x, y)</code>
$e^{-x^2}$	<code>exp_x2(x)</code>	$\operatorname{arccot}(x)$	<code>arccot(x)</code>
$2^x$	<code>exp2(x)</code>	$\sinh(x)$	<code>sinh(x)</code>
$10^x$	<code>exp10(x)</code>	$\cosh(x)$	<code>cosh(x)</code>
$x^y$	<code>power(x, y)</code>	$\tanh(x)$	<code>tanh(x)</code>
$\ln(x)$	<code>ln(x)</code>	$\operatorname{coth}(x)$	<code>coth(x)</code>
$\ln(e \cdot x)$	<code>ln_ex(x)</code>	$\operatorname{arsinh}(x)$	<code>arsinh(x)</code>
$\ln(1+x)$	<code>ln1p(x)</code>	$\operatorname{arcosh}(x)$	<code>arcosh(x)</code>
$\ln(x^2 + y^2)/2$	<code>ln_sqrtx2y2(x, y)</code>	$\operatorname{arcosh}(1+x)$	<code>arcosh1p(x)</code>
$\ln[(1+x)^2 + y^2]/2$	<code>ln_sqrt1px2y2(x, y)</code>	$\operatorname{artanh}(x)$	<code>artanh(x)</code>
$\log_a(x)$	<code>loga(a, x)</code>	$\operatorname{arcoth}(x)$	<code>arcoth(x)</code>
$\log_2(x)$	<code>log2(x)</code>	$\operatorname{sign}(x)$	<code>sign(x)</code>
$\log_{10}(x)$	<code>log10(x)</code>	$\operatorname{diam}(x)$	<code>diam(x)</code>
$\log_{10}(1+x)$	<code>log1p(x)</code>	$\operatorname{mid}(x)$	<code>mid(x)</code>
$\sin(x)$	<code>sin(x)</code>	$\operatorname{blow}(x, \text{eps})$	<code>blow(x, eps)</code>

Tabelle 4.3: Hochgenaue Standardfunktionen des Moduls `LLARI`;  
 $x, y$ : `Linterval`;  $\text{eps}$ : `real`;  $a$ : `Longreal`;

Die Funktionen `mid`, `diam` sind für die Berechnung des Mittelpunktes und des Durchmessers von Intervallen vorgesehen.

Die Mittelpunktsfunktion liefert mit: `m := MID (x)` den Longreal-Mittelpunkt `m` eines Intervalls `x`. Kann `m` mit dem Akkumulator nicht maximalgenau berechnet werden, erfolgt eine entsprechende Fehlermeldung mit Programmabbruch.

Der Intervalldurchmesser ist definiert durch: `diam(x) := sup(x) -> inf(x)`;

Bei Einschließungsmethoden wird oft die Funktion `blow` für die sogenannte Epsilonaufblähung benötigt. Mit einem *real*-Ausdruck `eps` und einem *Linterval*-Ausdruck `x` wird die *blow*-Funktion aufgebaut durch die Anweisungen:

```
y := (Long(1) + eps) * x - eps * x;
blow (x,eps) := intval (pred (inf(y)), succ (sup(y)));
```

### Beispiele:

Seien `a, b` vom Typ *Linterval*, erzeugt durch die Anweisungen:

```
a := intval (-1, 3); b := intval (2);
```

dann liefern die Funktionen `abs`, `sqr`, `mid`, `diam` die folgenden Ergebnisse:

Ausdruck	Ergebnis	Ausdruck	Ergebnis
<code>abs (a)</code>	<code>[0, 3]</code>	<code>sqr (b)</code>	<code>[4; 4]</code>
<code>abs (b)</code>	<code>[2, 2]</code>	<code>mid (a)</code>	<code>1</code>
<code>sqr (a)</code>	<code>[0, 9]</code>	<code>diam (a)</code>	<code>4</code>

### Anmerkung:

Mit Hilfe der obigen Intervallfunktionen vom Typ *Linterval* können **hochgenaue** Intervallfunktionen vom Typ *interval* sehr einfach bereitgestellt werden, sofern dies nicht schon nach Tabelle 4.2 geschehen ist. Das folgende Beispiel zeigt, wie z.B. die Potenzfunktion `xy` für *interval*-Argumente `x, y` zu deklarieren ist:

```
program test (input,output);
use i_ari, li_ari;
function power (x,y: interval): interval;
begin power := SHORT ( power(Long(x),Long(y)) ); end;
begin { test } ... end.
```

### 4.2.5 Überladungen des Zuweisungsoperators

Die Wandlung von *integer/real/Longreal/interval* nach *Linterval* wird auch in Form einer überladenen Zuweisung bereitgestellt:

Zuweisung	Bedeutung
$Lx := r$	$Lx := \text{Lintval}(r)$
$Lx := x$	$Lx := \text{intval}(r1,r2)$

$Lx = \text{Linterval-Variable}; \quad x = \text{interval-Ausdruck};$   
 $r = \text{integer/real/Longreal-Ausdruck}; \quad r1,r2 = \text{integer/real-Ausdruck};$

### 4.2.6 Ein-/Ausgabeprozeduren

Es stehen die bekannten Prozeduren

```
procedure read (var f: text; var a: Linterval);
procedure write (var f: text; a: Linterval);
```

mit optionalem Fileparameter, beliebig vielen Ein-/Ausgabeparametern, jedoch ohne Formatspezifikationen zur Verfügung.

Die Eingabe eines Intervalls  $x = [a, b]$  muß in der Form

$[a, b]$  oder in der Form  $a$

erfolgen. Der zweite Fall dient zur vereinfachten Eingabe eines Punktintervalls  $x = [a, a]$ .

Die Ausgabe eines Intervalls durch die write-Prozedur erfolgt stets in der Form

$[a, b]$

Das Ausgabeformat ist dabei durch folgende Regelung bestimmt:

Ist  $Lx$  eine Variable vom Typ *Linterval* und stimmen die Werte von  $\text{inf}(Lx)$  und  $\text{sup}(Lx)$  auf den ersten  $n \geq 2$  Ziffern überein, so gibt die Anweisung `write(Lx)` die Unterschranke  $a$  und die Oberschranke  $b$  jeweils auf  $n$  gerundete Mantissenstellen aus.

Für  $n = 0$  und  $n = 1$  werden  $a$  und  $b$  stets auf zwei Mantissenstellen gerundet.

Werden also durch `write(Lx)`  $n > 2$  Mantissenfiguren ausgegeben, so weiß der Anwender, daß  $\text{inf}(Lx)$  und  $\text{sup}(Lx)$  genau in den ersten  $n$  Ziffern übereinstimmen!

In folgender Tabelle sind für den Wert einer Variablen  $Lx$  vom Typ *Linterval* die jeweiligen Ergebnisse der Ausgabeprozedur `write(Lx)` zusammengestellt:

$Lx = [\inf(Lx), \sup(Lx)]$	$\text{write}(Lx)$
[2, 3]	[2.0E+000, 3.0E+000]
[1.1, 1.2]	[1.1E+000, 1.2E+000]
[1.11, 1.13]	[1.1E+000, 1.2E+000]
[2.000, 2.0001]	[2.000E+000, 2.001E+000]

Die *Longreal*-Werte  $\inf(Lx)$  und  $\sup(Lx)$  können natürlich auch (**ohne** zusätzliche Rundungen durch die  $\text{write}(Lx)$ -Prozedur) wie folgt ausgegeben werden:

- $\text{write}(\inf(Lx), \sup(Lx));$
- $\text{write}(Lx.\text{inf}, Lx.\text{sup});$

$Lx := \text{Lintval}(1.11, 1.13)$  und  $\text{write}(Lx.\text{inf}, Lx.\text{sup})$  ergeben:

[1.11000000000000000000000000000000E+000, 1.13000000000000000000000000000000E+000]

im Gegensatz zur entsprechenden Ausgabe von  $\text{write}(Lx)$  in obiger Tabelle.

#### 4.2.7 Exakte Auswertung von Intervallausdrücken

In **PASCAL-XSC** können mit Hilfe der Operatoren  $+$ ,  $-$ ,  $*$  Intervall-Ausdrücke vom Typ **interval** mit Hilfe des langen Akkus rundungsfehlerfrei berechnet und durch die Anweisung

$i := \#\#(\text{exakter Intervallausdruck});$

mit nur einer einzigen Rundung pro Intervallgrenze in eine Variable  $i$  vom Typ **interval** abgespeichert werden. Dieses  $\#\#$ -Konzept steht auch für die BCD-Version in **vollem Umfang** zur Verfügung; es ist jedoch **nicht** anwendbar, wenn nur einer der Operanden im exakt auszuwertenden *Intervallausdruck* vom Typ *Longreal* oder *Linterval* ist!

Um das  $\#\#$ -Konzept für diese Datentypen wenigstens simulieren zu können, benötigt man den Datentyp **idotprecision**:

**type** idotprecision = **record** INF, SUP: dotprecision **end**;

##### 4.2.7.1 Operatoren

Wie auf Seite 42 beschrieben, muß auch jetzt z.B. für einen Additionsoperator mit *Linterval*-Operanden und einem Ergebnis vom Typ **idotprecision** ein **neuer** Operatorname (**plus\_i**) gewählt werden, da der normale  $+$  Operator mit den gleichen Operanden schon für den Ergebnistyp *Linterval* vergeben ist. Beachten Sie, daß Operatoren, Prozeduren und Funktionen vom Compiler nicht bzgl. des Ergebnistyps unterschieden werden! In nachfolgender Tabelle 4.4 sind die entsprechenden neuen Operatornamen zusammengestellt:

Operation	Operator-Name	Ergebnistyp
Addition	<b>plus_i</b>	idotprecision
Subtraktion	<b>minus_i</b>	idotprecision
Multiplikation	<b>times_i</b>	idotprecision

Tabelle 4.4: Operatornamen mit Ergebnistyp idotprecision

In Tabelle 4.5 sind für die Operatoren **plus\_i**, **minus\_i**, **times\_i**, \* die erlaubten Operanden-Typen angegeben. Mit einem *dotprecision*-Operanden existiert jedoch der \* Operator nur mit dem zweiten Operandentyp idotprecision; vergl. auch Tabelle 4.6 !

rechter Operand	integer/real Longreal	interval Linterval	idotprecision (dotprecision)
<b>linker Operand</b>			
integer/real Longreal	times_i plus_i, minus_i	times_i plus_i, minus_i	* plus_i, minus_i
interval Linterval	times_i plus_i, minus_i	times_i plus_i, minus_i	* plus_i, minus_i
idotprecision (dotprecision)	* plus_i, minus_i	* plus_i, minus_i	* plus_i, minus_i

Tabelle 4.5: Operatoren des Moduls LI-ARI mit Ergebnistyp idotprecision

**Wertzuweisungen** an eine Variable vom Typ idotprecision:

Durch Überladen des Zuweisungsoperators := sind folgende Wertzuweisungen möglich:

```
idotprecision-Variable := Linterval-Ausdruck;
idotprecision-Variable := interval-Ausdruck;
idotprecision-Variable := dotprecision-Ausdruck;
idotprecision-Variable := Longreal-Ausdruck;
idotprecision-Variable := real/integer-Ausdruck;
```

In nachfolgender Tabelle 4.6 können zu den Operatoren mit den Ergebnistypen:

**boolean, interval, Linterval, idotprecision**

alle möglichen Operanden abgelesen werden:

rechter Operand	integer real Longreal	Linterval interval	idotprecision	dotprecision
linker Operand				
monadisch		+, -		
integer real Longreal		$\diamond, \Delta$ <b>in</b> , =, <> +*	$\square$ <b>in</b> , =, <> +*	
Linterval interval	$\diamond, \Delta$ =, <> +*	$\diamond, \Delta$ <b>in</b> , $\vee$ , $><$ +*, **	$\square$ <b>in</b> , $\vee$ , $><$ +*, **	
idotprecision	$\square$ =, <> +*	$\square$ <b>in</b> , $\vee$ , $><$ +*, **	$\square$ <b>in</b> , $\vee$ , $><$ +*, **	$\square$
dotprecision			$\square$	

Tabelle 4.6: Operatoren vom Typ: *interval*, *Linterval*, *idotprecision* oder *boolean***Bedeutung der Operatoren:**

1. +, -; Dies sind monadische Operatoren mit den Ergebnistypen: *interval*, *Linterval*, *idotprecision*;
2.  $\diamond \in \{+, -, *, /\}$ , Grundoperatoren mit Ergebnistyp: *interval*, *Linterval*;
3.  $\square \in \{\text{plus}_i, \text{minus}_i, *\}$ ; Additions-, Subtraktions- und Multiplikationsoperatoren mit Ergebnistyp *idotprecision*;
4.  $\Delta \in \{\text{plus}_i, \text{minus}_i, \text{times}_i\}$ ; Additions-, Subtraktions- und Multiplikationsoperatoren mit Ergebnistyp *idotprecision*;
5. **in** -Operator testet:  $[a, b] \subset [A, B]$  oder  $x \in [A, B]$ ;
6.  $><$  Operator testet auf Disjunktheit zweier Intervalle;
7. +\* Operator liefert die konvexe Hülle zweier Intervalle oder eines Intervalls mit einer Zahl;
8. \*\* Operator liefert das Schnittintervall; ein leeres Schnittintervall erzeugt einen Laufzeitfehler;
9.  $\vee \in \{=, <>, <, <=, >, >=\}$ , Vergleichsoperatoren zwischen Intervallen;

### 4.2.7.2 Transferfunktionen

Zur Wandlung zwischen den Typen

*integer/real/Longreal/dotprecision/interval/Linterval*

einerseits und dem Typ *idotprecision* andererseits werden folgende Transferfunktionen bereitgestellt: *idxsup* (Funktion)

Funktion	Ergebnistyp	Bedeutung
DOTINTVAL(d1,d2)	<i>idotprecision</i>	Intervall mit $inf = d1$ und $sup = d2$ ;
DOTINTVAL(d)	<i>idotprecision</i>	Punktintervall mit $inf = sup = d$ ;
DOTINTVAL(r1,r2)	<i>idotprecision</i>	Intervall mit $inf = r1$ und $sup = r2$ ;
DOTINTVAL(r)	<i>idotprecision</i>	Punktintervall mit $inf = sup = r$ ;
INF (di)	<i>dotprecision</i>	Untergrenze von di
SUP (di)	<i>dotprecision</i>	Obergrenze von di
LONG (di)	<i>Linterval</i>	di-Einschl. durch Intervall: <i>Linterval</i>
SHORT (di)	<i>interval</i>	di-Einschl. durch Intervall: <i>interval</i>

Tabelle 4.7: Transferfunktionen mit dem Typ *idotprecision*

$r,r1,r2 = integer/real/Longreal$ -Ausdruck,  $r1 \leq r2$ , sonst Fehlermeldung;

$di = idotprecision$ -Ausdruck;

$d,d1,d2 = dotprecision$ -Variable,  $d1 \leq d2$ , sonst Fehlermeldung;

### 4.2.7.3 Simulation von SUM-Ausdrücken

Zur Simulation von **SUM**-Ausdrücken im **##**-Konzept stehen die folgenden Datentypen und Prozeduren zur Verfügung:

- **type** *rvector* = **dynamic** array [\*] of **real**;
  - **type** *Lrvector* = **dynamic** array [\*] of **Longreal**;
  - **type** *ivector* = **dynamic** array [\*] of **interval**;
  - **type** *Livector* = **dynamic** array [\*] of **Linterval**;
1. procedure **ADDNACCU** (var di: *idotprecision*; var a: *rvector*);  
**Exakte** Addition des real-Feldes a zum Intervall di;
  2. procedure **ADDNACCU** (var di: *idotprecision*; var a: *Lrvector*);  
**Exakte** Addition des Longreal-Feldes a zum Intervall di;
  3. procedure **ADDNACCU** (var di: *idotprecision*; var a: *ivector*);  
**Exakte** Addition des interval-Feldes a zum Intervall di;

4. procedure **ADDNACCU** (var di: idotprecision; var a: Livector);  
**Exakte** Addition des Linterval-Feldes a zum Intervall di;
5. procedure **SUBNACCU** (var di: idotprecision; var a: rvector);  
**Exakte** Subtraktion des real-Feldes a vom Intervall di;
6. procedure **SUBNACCU** (var di: idotprecision; var a: Lrvector);  
**Exakte** Subtraktion des Longreal-Feldes a vom Intervall di;
7. procedure **SUBNACCU** (var di: idotprecision; var a: ivector);  
**Exakte** Subtraktion des interval-Feldes a vom Intervall di;
8. procedure **SUBNACCU** (var di: idotprecision; var a: Livector);  
**Exakte** Subtraktion des Linterval-Feldes a vom Intervall di;
9. procedure **PADDNACCU** (var di: idotprecision; var a,b: rvector);  
**Exakte** Addition des real-Skalarproduktes zum Intervall di;
10. procedure **PADDNACCU** (var di: idotprecision; var a,b: Lrvector);  
**Exakte** Addition des Longreal-Skalarproduktes zum Intervall di;
11. procedure **PADDNACCU** (var di: idotprecision; var a,b: ivector);  
**Exakte** Addition des interval-Skalarproduktes zum Intervall di;
12. procedure **PADDNACCU** (var di: idotprecision; var a,b: Livector);  
**Exakte** Addition des Linterval-Skalarproduktes zum Intervall di;
13. procedure **PSUBNACCU** (var di: idotprecision; var a,b: rvector);  
**Exakte** Subtraktion des real-Skalarproduktes vom Intervall di;
14. procedure **PSUBNACCU** (var di: idotprecision; var a,b: Lrvector);  
**Exakte** Subtraktion des Longreal-Skalarproduktes vom Intervall di;
15. procedure **PSUBNACCU** (var di: idotprecision; var a,b: ivector);  
**Exakte** Subtraktion des interval-Skalarproduktes vom Intervall di;
16. procedure **PSUBNACCU** (var di: idotprecision; var a,b: Livector);  
**Exakte** Subtraktion des Linterval-Skalarproduktes vom Intervall di;

#### 4.2.7.4 Beispiele

##### 1. Beispiel:

Mit den Deklarationen:     var doti : idotprecision;  
                                  a,b : Livector [1..10];  
                                  Li : Linterval;

sind z.B. folgende Anweisungen möglich:

```
doti := 0;  PADDNACCU (doti,a,b);
Li := LONG( doti plus_i a[1] times_i b[10] minus_i a[2] times_i b[9] );
```

Das Intervall Li liefert eine Einschließung der exakten Summe mit **nur einer Rundung** bzgl. Intervallunter- und -Obergrenze!



**2. Beispiel:**

Mit Hilfe der Operatoren `*` und `plus_i` zwischen *idotprecision*- und *Linterval*-Operanden kann ein Polynom

$$P(x) = \sum_{k=0}^n a_k \cdot x^k$$

mit folgender Funktion **rundungsfehlerfrei** berechnet werden:

```
function HORNER_IDOT(var a:Livector; x:Linterval):idotprecision;
var k : integer;    y : idotprecision;
begin
  y := a[ub(a)];
  for k := ub(a)-1 downto lb(a) do
    y := y * x plus_i a[k];    { Intervall-Hornerschema }
  HORNER_IDOT := y
end;
```

Die Koeffizienten

$a_0 = -0.01;$	$a_1 = +0.39;$	$a_2 = -7.02;$	$a_3 = +77.22;$
$a_4 = -579.15;$	$a_5 = +3127.41;$	$a_6 = -12509.64;$	$a_7 = +37528.92;$
$a_8 = -84440.07;$	$a_9 = +140733.45;$	$a_{10} = -168880.14;$	$a_{11} = +138174.66;$
$a_{12} = -69087.33;$	$a_{13} = +15943.23;$		

realisieren das Polynom  $P(x) = \frac{1}{100} \cdot (3x - 1)^{13}$  mit der Nullstelle  $x_0 = \frac{1}{3}$ ;  
Mit den Deklarationen:

```
var a: Livector[0..13];    x,y: Linterval;    doti: idotprecision;
```

und den Anweisungen:

```
x := 0.3333333333333333;    {x ≈ x0; x ist Punktintervall!}
y := LONG( HORNER_IDOT (a,x) );
```

erhält man mit der BCD-Version das **exakte Intervallergebnis:**

$$y = P(x) \in [-1.00 \dots 00 \cdot 10^{-210}; -1.00 \dots 00 \cdot 10^{-210}]$$

Führt man ganz entsprechende Rechnungen mit der Binärversion von PASCAL-XSC durch, wobei die  $a_i$  und das Argument  $x$  vorher durch Binärintervalle einzuschließen sind, so erhält man:

$$P(x) \in [-4.41283 \dots 593 \cdot 10^{-15}; +2.40837 \dots 365 \cdot 10^{-15}]$$

Die enorme Abweichung von **195 Größenordnungen** von der exakten BCD-Einschließung ist alleine bedingt durch den Umstand, daß **vor** der eigentlichen rundungsfehlerfreien Polynomauswertung neben dem Argument  $x$  auch alle Polynomkoeffizienten in Binärintervalle eingeschlossen werden müssen, deren Durchmesser i.a. von Null verschieden sind, da Dezimalzahlen i.a. nicht exakt als Binärzahlen endlicher Länge darstellbar sind.

### 4.2.8 Wertebereich monotoner Funktionen

Zur Realisierung der hochgenauen Intervallfunktionen aus Tabelle 4.3 mit Hilfe monotoner Punktfunktionen vom Ergebnistyp **rrG** sind im Modul **LIARI** folgende Werkzeuge bereitgestellt:

- `function Lintval_f (`

<code>{ Wertebereichs-Einschließg. }</code>
<code>function f(x: Longreal): rrG; { Monotone Punktfunktion }</code>
<code>I: Linterval; { Monotonie-Intervall }</code>
<code>up: boolean; { Monotoniebeschreibung }</code>
<code>eps: real { Vergrößerte Fehlerschranke }</code>
<code>) : Linterval; { Typ des Wertebereichs }</code>

Die Funktion `f(x: Longreal): rrG` muß im Intervall `I` monoton sein. Mit der relativen Fehlerschranke  $\varepsilon(f)$  der Funktion  $f$  muß gelten:

$$\text{eps} := 1.00001 * > \varepsilon(f);$$

Ist  $f(x)$  in `I` monoton wachsend (fallend), so ist der Parameter `up true` (`false`) zu wählen. Die Anweisung

```
y := Lintval_f(f,I,up,eps);
```

liefert dann mit `y: Linterval` eine hochgenaue Einschließung aller Funktionswerte  $f(x)$ , mit  $x \in I$ .

- `function Lintval_frrG (`

<code>{ Wertebereichs-Einschließg. }</code>
<code>function f(x: rrG): rrG; { Monotone Punktfunktion }</code>
<code>I: Linterval; { Monotonie-Intervall }</code>
<code>up: boolean; { Monotoniebeschreibung }</code>
<code>eps: real { Vergrößerte Fehlerschranke }</code>
<code>) : Linterval; { Typ des Wertebereichs }</code>

Die Funktion `f(x: rrG): rrG` muß im Intervall `I` monoton sein. Mit der relativen Fehlerschranke  $\varepsilon(f)$  der Funktion  $f$  muß gelten:

$$\text{eps} := 1.00001 * > \varepsilon(f);$$

Ist  $f(x)$  in `I` monoton wachsend (fallend), so ist der Parameter `up true` (`false`) zu wählen. Die Anweisung

```
y := Lintval_frrG(f,I,up,eps);
```

liefert dann mit `y: Linterval` eine hochgenaue Einschließung aller Funktionswerte  $f(x)$ , mit  $x \in I$ .

### 4.3 Intervallrechnung mit dem Typ rrG

Mit dem Modul **RRGLARI** werden nur die Operatoren, Funktionen und Prozeduren bereitgestellt, die zur Implementierung von mathematischen Intervallfunktionen notwendig sind. Es ist nicht vorgesehen, die komplette Intervallrechnung wie etwa im Modul **LIARI** zu realisieren, so daß z.B. Vergleichsoperatoren für Intervalle oder Intervall-Grundoperationen mit gemischten Operandentypen nur in dem Umfang bereitgestellt werden, wie sie tatsächlich erforderlich waren.

#### 4.3.1 Der Datentyp rrGinterval

Der Datentyp *rrGinterval* ist definiert durch:

```
type rrGinterval = record inf, sup : rrG end;
```

Mit der Deklaration: **var** x : rrGinterval; bedeutet x dann das reelle Intervall:

$$x = [x.INF, x.SUP]$$

#### 4.3.2 Operatoren

Für Operanden vom Typ *rrGinterval* stehen die monadischen Operatoren  $+$   $-$  und die arithmetischen Operatoren  $+$   $-$   $*$   $/$  zur Verfügung. Bei der Multiplikation kann zusätzlich der erste Operand vom Typ *rrG,Longreal,real* und bei der Division der zweite Operand ebenfalls vom Typ *rrG,Longreal,real* gewählt werden.

Als Vergleichsoperatoren sind lediglich die Operatoren  $=$   $<$   $>$  definiert, wobei beide Operanden vom Typ *rrGinterval* sein müssen.

Beim Operator **in** ist der zweite Operand vom Typ *rrGinterval*, und der erste Operand kann vom Typ *real,rrGinterval* gewählt werden.

Der Operator  $><$  testet zwei *rrGinterval*-Operanden auf Disjunktheit, wobei zwei Intervalle disjunkt heißen, wenn sie kein gemeinsames Element enthalten.

Die Verbandsoperatoren  $+$   $*$   $**$  sind als Intervall-Hülle bzw. als Durchschnitt zweier *rrGinterval*-Operanden definiert.

#### 4.3.3 Transferfunktionen

In nachfolgender Tabelle sind alle die Funktionen angegeben, die den Zusammenhang zwischen dem Typ *rrGinterval* einerseits und den Typen *integer*, *real*, *Longreal* und *rrG* andererseits herstellen:

Funktion	Ergebnistyp	Bedeutung
RRGINTVAL (a,b)	rrGinterval	$[a, b]$ ; $a, b : \text{rrG}$
RRGINTVAL (a)	rrGinterval	$[a, a]$ ; $a : \text{integer/real/Longreal/rrG}$

## 4.3.4 Standardfunktionen

Es stehen die in nachfolgender Tabelle 4.8 angegebenen mathematischen Intervall-Standardfunktionen zur Verfügung. Argumente  $x, y$  vom Typ `rrGInterval` liefern, mit Ausnahme von `sign(x)`, Funktionswerte ebenfalls vom Typ `rrGInterval`.

$$\text{sign}(x) = \begin{cases} 0 & \text{falls } 0 \in x \\ +1 & \text{falls } x.\text{inf} > 0 \\ -1 & \text{falls } x.\text{sup} < 0 \end{cases}$$

Funktion	Aufruf	Funktion	Aufruf
$ x $	<code>abs(x)</code>	$\sin(x)$	<code>sin(x)</code>
$x^2$	<code>sqr(x)</code>	$\cos(x)$	<code>cos(x)</code>
$x^2 - y^2$	<code>x2_y2(x, y)</code>	$\tan(x)$	<code>tan(x)</code>
$\sqrt{x}$	<code>sqr(x)</code>	$\cot(x)$	<code>cot(x)</code>
$\sqrt{1+x} - 1$	<code>sqr1pm1(x)</code>	$\sin(\pi \cdot x)$	<code>sinpi(x)</code>
$\sqrt{x^2 + y^2}$	<code>sqr(x2y2(x, y))</code>	$\cos(\pi \cdot x)$	<code>cospi(x)</code>
$\sqrt{1+x^2}$	<code>sqr1px2(x)</code>	$\tan(\pi \cdot x)$	<code>tanpi(x)</code>
$\sqrt{1-x^2}$	<code>sqr1mx2(x)</code>	$\cot(\pi \cdot x)$	<code>cotpi(x)</code>
$\sqrt{x^2 - 1}$	<code>sqr(x2m1(x))</code>	$\arcsin(x)$	<code>arcsin(x)</code>
$e^x$	<code>exp(x)</code>	$\arccos(x)$	<code>arccos(x)</code>
$e^x - 1$	<code>expm1(x)</code>	$\arctan(x)$	<code>arctan(x)</code>
$e^{-x^2}$	<code>exp_x2(x)</code>	$\arctan2(x, y)$	<code>arctan2(x, y)</code>
$2^x$	<code>exp2(x)</code>	$\text{arccot}(x)$	<code>arccot(x)</code>
$10^x$	<code>exp10(x)</code>	$\sinh(x)$	<code>sinh(x)</code>
$x^y$	<code>power(x, y)</code>	$\cosh(x)$	<code>cosh(x)</code>
$\ln(x)$	<code>ln(x)</code>	$\tanh(x)$	<code>tanh(x)</code>
$\ln(e \cdot x)$	<code>ln_ex(x)</code>	$\coth(x)$	<code>coth(x)</code>
$\ln(1+x)$	<code>ln1p(x)</code>	$\text{arsinh}(x)$	<code>arsinh(x)</code>
$\ln(x^2 + y^2)/2$	<code>ln_sqrt(x2y2(x, y))</code>	$\text{arcosh}(x)$	<code>arcosh(x)</code>
$\ln[(1+x)^2 + y^2]/2$	<code>ln_sqrt1px2y2(x, y)</code>	$\text{arcosh}(1+x)$	<code>arcosh1p(x)</code>
$\log_a(x)$	<code>loga(a, x)</code>	$\text{artanh}(x)$	<code>artanh(x)</code>
$\log_2(x)$	<code>log2(x)</code>	$\text{arcoth}(x)$	<code>arcoth(x)</code>
$\log_{10}(x)$	<code>log10(x)</code>	$\text{sign}(x)$	<code>sign(x)</code>
$\log_{10}(1+x)$	<code>log1p(x)</code>	$\text{diam}(x)$	<code>diam(x)</code>

Tabelle 4.8: Mathematische Standardfunktionen des Moduls `RRGLARI`;  $x, y$ : `rrGInterval`;  $a$ : `Longreal`;

Um bei verschachtelten Intervall-Funktionen Zwischenrechnungen auch im genaueren rrGinterval-Format durchführen zu können, sind in nachfolgender Tabelle 4.9 alle mathematischen Intervall-Standardfunktionen angegeben, die für Argumente  $x, y$  vom Typ Linterval Funktionswerte vom Typ rrGinterval liefern:

Funktion	Aufruf	Funktion	Aufruf
$ x $	abs_rrG( $x$ )	$\sin(x)$	sin_rrG( $x$ )
$x^2$	sqr_rrG( $x$ )	$\cos(x)$	cos_rrG( $x$ )
$x^2 - y^2$	x2_y2_rrG( $x, y$ )	$\tan(x)$	tan_rrG( $x$ )
$\sqrt{x}$	sqrt_rrG( $x$ )	$\cot(x)$	cot_rrG( $x$ )
$\sqrt{1+x} - 1$	sqrt1pm1_rrG( $x$ )	$\sin(\pi \cdot x)$	sinpi_rrG( $x$ )
$\sqrt{x^2 + y^2}$	sqrtx2y2_rrG( $x, y$ )	$\cos(\pi \cdot x)$	cospi_rrG( $x$ )
$\sqrt{1+x^2}$	sqrt1px2_rrG( $x$ )	$\tan(\pi \cdot x)$	tanpi_rrG( $x$ )
$\sqrt{1-x^2}$	sqrt1mx2_rrG( $x$ )	$\cot(\pi \cdot x)$	cotpi_rrG( $x$ )
$\sqrt{x^2 - 1}$	sqrtx2m1_rrG( $x$ )	$\arcsin(x)$	arcsin_rrG( $x$ )
$e^x$	exp_rrG( $x$ )	$\arccos(x)$	arccos_rrG( $x$ )
$e^x - 1$	expm1_rrG( $x$ )	$\arctan(x)$	arctan_rrG( $x$ )
$e^{-x^2}$	exp_x2_rrG( $x$ )	$\arctan2(x, y)$	arctan2_rrG( $x, y$ )
$2^x$	exp2_rrG( $x$ )	$\operatorname{arccot}(x)$	arccot_rrG( $x$ )
$10^x$	exp10_rrG( $x$ )	$\sinh(x)$	sinh_rrG( $x$ )
$x^y$	power_rrG( $x, y$ )	$\cosh(x)$	cosh_rrG( $x$ )
$\ln(x)$	ln( $x$ )	$\tanh(x)$	tanh_rrG( $x$ )
$\ln(e \cdot x)$	ln_ex_rrG( $x$ )	$\coth(x)$	coth_rrG( $x$ )
$\ln(1+x)$	ln1p_rrG( $x$ )	$\operatorname{arsinh}(x)$	arsinh_rrG( $x$ )
$\ln(x^2 + y^2)/2$	ln_sqrtx2y2_rrG( $x, y$ )	$\operatorname{arcosh}(x)$	arcosh_rrG( $x$ )
$\ln[(1+x)^2 + y^2]/2$	ln_sqrt1px2y2_rrG( $x, y$ )	$\operatorname{arcosh}(1+x)$	arcosh1p_rrG( $x$ )
$\log_a(x)$	loga_rrG( $a, x$ )	$\operatorname{artanh}(x)$	artanh_rrG( $x$ )
$\log_2(x)$	log2_rrG( $x$ )	$\operatorname{arcoth}(x)$	arcoth_rrG( $x$ )
$\log_{10}(x)$	log10_rrG( $x$ )	$\operatorname{diam}(x)$	diam_rrG( $x$ )
$\log_{10}(1+x)$	log1p_rrG( $x$ )		

Tabelle 4.9: Mathematische Standardfunktionen des Moduls RRG.L.ARI;  
 $x, y$ : Linterval;  $a$ : Longreal;

### 4.3.5 Überladungen des Zuweisungsoperators

Die Wandlungen von *integer/real/Longreal/rrG/Linterval* nach *rrGinterval* bzw. von *rrGinterval* nach *Linterval* werden auch durch Überladungen des Zuweisungsoperators bereitgestellt:

Zuweisung	Bedeutung
$rrGi := rrGx$	$rrGi := RRGINTVAL(rrGx);$ $rrGi : rrGinterval;$ $rrGx : rrG$
$rrGi := Lx$	$rrGi := RRGINTVAL(Lx);$ $rrGi : rrGinterval;$ $Lx : Longreal$
$rrGi := x$	$rrGi := RRGINTVAL(x);$ $rrGi : rrGinterval;$ $x : real$
$rrGi := Li$	$rrGi \equiv Li;$ $rrGi : rrGinterval;$ $LI : Linterval$
$Li := rrGi$	$rrGi \subseteq Li;$ $rrGi : rrGinterval;$ $LI : Linterval$

### 4.3.6 Ein-/AusgabeprozEDUREN

Es stehen die bekannten Prozeduren

```
procedure read (var f: text; var a: rrGinterval);
procedure write (var f: text; a: rrGinterval);
```

mit optionalem Fileparameter, beliebig vielen Ein-/Ausgabeparametern, jedoch ohne Formatspezifikationen zur Verfügung.

Die Eingabe eines Intervalls  $x = [a, b]$  mit  $a, b : rrG$  muß in der Form

$$[a, b] \quad \text{oder in der Form} \quad a$$

erfolgen. Der zweite Fall dient zur vereinfachten Eingabe eines Punktintervalls  $x = [a, a]$ .

Die Eingabe eines Wertes  $a : rrG$  muß dabei in der Form

$$\{ a.r1, a.r2, a.G \} \quad \text{oder in der Form} \quad a.r1$$

erfolgen. Weitere Einzelheiten findet man auf Seite 56.

Die Ausgabe eines Intervalls durch die write-Prozedur erfolgt stets in der Form

$$[a, b]$$

wobei z.B.  $a : rrG$  in der Form

$$\{ a.r1, a.r2, a.G \}$$

ausgegeben wird.

### 4.3.7 Wertebereich einer Funktion

Gegeben sei eine über einem Intervall  $\mathbb{D} = [a, b]$  definierte, reellwertige Funktion  $f$

$$f : \mathbb{D} = [a, b] \longrightarrow \mathbb{R};$$

Der Wertebereich  $\mathbb{W}$  von  $f$  ist dann definiert durch:

$$\mathbb{W} := \{y \mid y = f(x) \wedge x \in [a, b]\}.$$

Wird  $f(x)$  für ein rrG-Argument  $x \in [a, b]$  mit einem geeigneten Algorithmus ausgewertet, so ist  $\tilde{f}(x) \neq f(x)$  das fehlerbehaftete Rechnerergebnis mit dem relativen Fehler  $\varepsilon_f$ , definiert durch:

$$\varepsilon_f := \frac{\tilde{f}(x) - f(x)}{f(x)}; \quad f(x) \neq 0; \quad |\varepsilon_f| \leq \varepsilon(f), \quad x \in [a, b].$$

Die Fehlerschranke  $\varepsilon(f)$  wird als bekannt angesehen. Die PASCAL-XSC-Funktion

```
function f (x: rrG): rrG;
```

liefert zum exakten Funktionswert  $f(x)$  die rrG-Näherung  $\tilde{f}(x)$ .

Gesucht ist eine Intervallfunktion vom Ergebnistyp *rrGinterval*, die zu einem gegebenem Intervall  $I \subset [a, b]$  vom Typ *rrGinterval* eine **Einschließung** des entsprechenden Wertebereichs

$$\mathbb{W}_I := \{y \mid y = f(x) \wedge x \in I \subset [a, b]\}.$$

liefert. Unter bestimmten Voraussetzungen lassen sich die gesuchten Intervallfunktionen mit vordefinierten Funktionen aus dem Modul **rrGi\_ari** realisieren:

#### 4.3.7.1 Die Funktion ist monoton

Mit Hilfe der Funktion **RRGINTVAL\_f\_rrG**(...) aus dem Modul **rrGi\_ari** erhält man die gewünschte Einschließung des Wertebereichs:

```
function RRGINTVAL_f_rrG ( function f (x: rrG): rrG; I: rrGinterval;
                           bl: boolean; eps: real): rrGinterval;
```

Zu einer in einem Intervall  $I = [x_1, x_2]$  monotonen Funktion  $f$  berechnet

```
w := RRGINTVAL_f_rrG (f, I, bl, eps)
```

eine garantierte Einschließung  $w$  der Funktionswerte  $f(x)$ , mit  $x \in I = [x_1, x_2]$ :

$$\mathbb{W}_I := \{y \mid y = f(x) \wedge x \in I\} \subseteq w;$$

Ist  $f$  in  $I$  monoton wachsend (fallend), so muß  $bl = true$  (*false*) gesetzt werden.

Ist  $\varepsilon(f)$  die relative Fehlerschranke von  $f$ , so muß  $eps$  wie folgt berechnet werden:

$$eps := 1.00001 * > \varepsilon(f);$$

Im nachfolgenden Beispiel wird gezeigt, wie die Exponentialfunktion

$$\text{function exp (x: rrGinterval): rrGinterval;}$$

aus dem Modul **rrGi\_ari** realisiert wurde:

```
function EXP(x: rrGinterval): rrGinterval;
var c : rrGinterval;
begin
  c := rrGintval_f_rrG(EXP,x,TRUE,4.6214E-25);
  if SIGN(c.SUP.r1) = 0 then
    begin { Oberschranke falls Underflow: }
      c.SUP.r1 := 1;  c.SUP.r2 := 0;
      c.SUP.G  := -2171;
    end;
  EXP := c;
end; { EXP }
```

EXP ist dabei der Name der Exponentialfunktion für rrG-Argumente aus der Tabelle von Seite 58. Mit der Fehlerschranke  $\varepsilon(f) = 4.6213 \text{ E} - 25$  ergibt sich dann eps zu:  $\text{eps} = 4.6214 \text{ E} - 25$ .

Mit den beiden folgenden Funktionen aus dem Modul **rrGi\_ari** lassen sich Wertebereiche von monotonen Funktionen ganz entsprechend berechnen, wenn das Monotonie-Intervall oder das Argument von  $f$  andere Datentypen besitzen:

$$\text{function RRGINTVAL\_f ( function f (x: Longreal): rrG; I: Linterval;}$$

$$\text{bl: boolean; eps: real): rrGinterval;}$$

$$\text{function RRGINTVAL\_frrG ( function f (x: rrG): rrG; I: Linterval;}$$

$$\text{bl: boolean; eps: real): rrGinterval;}$$

#### 4.3.7.2 Die Funktion besitzt ein relatives Maximum

$f : [a, b] \rightarrow \mathbb{R}$  sei eine stetige Funktion und habe in  $x_0 \in ]a, b[$  ein relatives Maximum. Für  $x < x_0$  sei  $f$  monoton wachsend, und für  $x > x_0$  sei  $f$  monoton fallend.

$f(x)$  wird auf dem Rechner durch folgende Funktion realisiert:

$$\text{function f (x: rrG): rrG;}$$

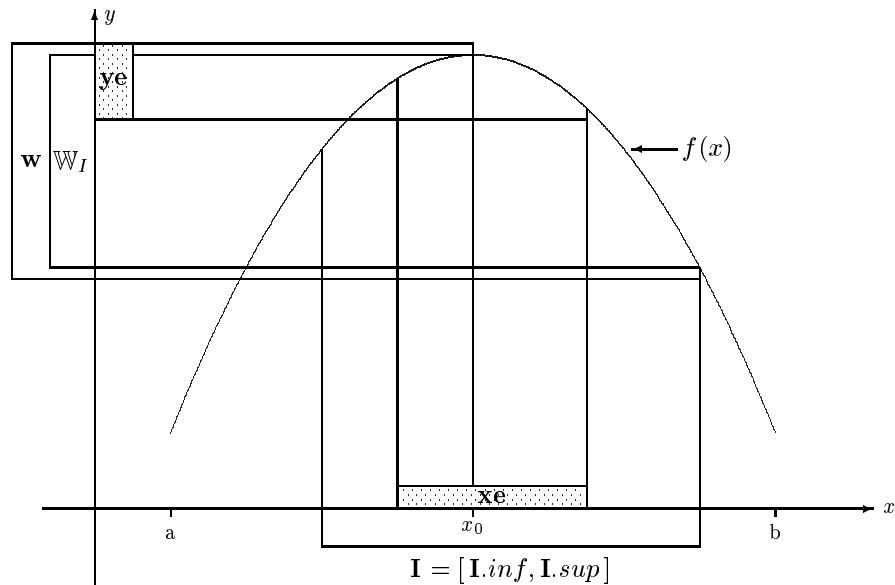
und liefert zum rrG-Argument  $x \in [a, b]$  den i.a. fehlerbehafteten Funktionswert  $\tilde{f}(x) \neq f(x)$ .

Die zugehörige relative Fehlerschranke  $\varepsilon(f)$  wird für alle  $x \in [a, b]$  als bekannt angenommen.

Zu einem beliebig vorgegebenen Intervall I vom Typ *rrGinterval*, mit:

$$I = [I.\text{inf}, I.\text{sup}] \subseteq [a, b]$$



Abbildung 4.1: Einschließung des Wertebereichs  $\mathbb{W}_I$  durch das Intervall  $w$ 

ist gesucht eine Einschließung  $w = [w.inf, w.sup]$  vom Typ *rrGinterval* für den Wertebereich:

$$\mathbb{W}_I := \{y \mid y = f(x) \wedge x \in I = [I.inf, I.sup]\} \subseteq w = [w.inf, w.sup];$$

In Abb. 4.1 ist  $xe$  ein Intervall vom Typ *rrGinterval*, welches die Maximumstelle  $x_0$  einschließt. Das Intervall  $ye$ , ebenfalls vom Typ *rrGinterval*, ist eine Einschließung aller Funktionswerte  $f(x)$ , mit  $x \in xe$ :

$$\{y \mid y = f(x) \wedge x \in xe\} \subseteq ye;$$

Ist die relative Fehlerschranke  $\varepsilon(f)$  bekannt, so lassen sich die Einschließungen  $xe, ye$  mit dem Programm **MIN\_26** bestimmen, wenn man dort  $-f(x)$  an Stelle von  $f(x)$  benutzt, da **MIN\_26** die genannten Einschließungen nur unter der Voraussetzung berechnet, daß die betrachtete Funktion ein relatives **Minimum** besitzt.

Mit Hilfe der Funktion

```
function INCL_MAX (function f (x: rRG): rRG; eps: real;
                  xe, ye, I: rrGinterval): rrGinterval;
```

die durch **rrGi\_ari** bereitgestellt wird, erhält man dann die Einschließung  $w$  durch den Funktionsaufruf:

```
w := INCL_MAX (f, eps, xe, ye, I);
```

wobei  $eps$  durch  $eps := 1.00001 * \varepsilon(f)$  zu berechnen ist.

### 4.3.7.3 Die Funktion besitzt ein relatives Minimum

$f : [a, b] \rightarrow \mathbb{R}$  sei eine stetige Funktion und habe in  $x_0 \in ]a, b[$  ein relatives Minimum. Für  $x < x_0$  sei  $f$  monoton fallend, und für  $x > x_0$  sei  $f$  monoton wachsend.

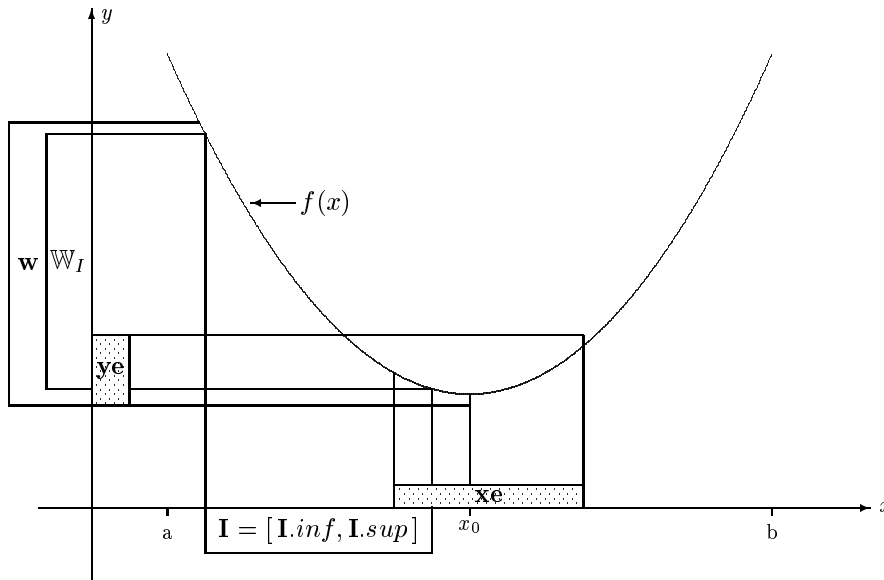


Abbildung 4.2: Einschließung des Wertebereichs  $\mathbb{W}_I$  durch das Intervall  $w$

$f(x)$  wird auf dem Rechner durch folgende Funktion realisiert:

```
function f (x: rrG): rrG;
```

und liefert zum  $rrG$ -Argument  $x \in [a, b]$  den i.a. fehlerbehafteten Funktionswert  $\tilde{f}(x) \neq f(x)$ .

Die zugehörige relative Fehlerschranke  $\varepsilon(f)$  wird für alle  $x \in [a, b]$  als bekannt angenommen.

Zu einem beliebig vorgegebenen Intervall  $I$  vom Typ  $rrGinterval$ , mit:

$$I = [I.inf, I.sup] \subseteq [a, b]$$

ist gesucht eine Einschließung  $w = [w.inf, w.sup]$  vom Typ  $rrGinterval$  für den Wertebereich:

$$\mathbb{W}_I := \{y \mid y = f(x) \wedge x \in I = [I.inf, I.sup]\} \subseteq w = [w.inf, w.sup];$$

In Abb. 4.2 ist  $xe$  ein Intervall vom Typ  $rrGinterval$ , welches die Minimumstelle  $x_0$  einschließt. Das Intervall  $ye$ , ebenfalls vom Typ  $rrGinterval$ , ist eine

Einschließung aller Funktionswerte  $f(x)$ , mit  $x \in xe$  :

$$\{y \mid y = f(x) \wedge x \in xe\} \subseteq ye;$$

Ist die relative Fehlerschranke  $\varepsilon(f)$  bekannt, so lassen sich die Einschließungen  $xe, ye$  mit dem Programm **MIN\_26** bestimmen.

Mit Hilfe der Funktion

```
function INCL_MAX (function f (x: rrG): rrG; eps: real;
                  xe, ye, I: rrGinterval): rrGinterval;
```

aus dem Modul **rrGi\_ari** erhält man dann die Einschließung  $w$  durch folgendes Programmbeispiel:

```
program test(input,output);
use rrGi_ari;
var xe, ye, I, w : rrGinterval;    eps: real;

function f(x: rrG): rrG;
begin { Definition von f;    x > 1 }
  f := exp(x) / ln(x);
end;    { Rel. Fehlerschr.: 5.4457E-24 }

function fI(xI: rrGinterval): rrGinterval;
{ fI liefert zu gegebenem Intervall I }
{ die Einschliessung w von W_I      }
function fn(x: rrG): rrG;
begin fn := -f(x) end;
begin fI := -INCL_MAX(fn, eps, xe, -ye, xI) end;

begin { Haupt }
  eps := 5.4458E-24;    { = 1.00001 *> 5.4457E-24 }
  { * Ergebnisse des Programms MIN_26 :      * }
  xe.INF.r1 := 1.763222834346;
  xe.INF.r2 := 3.793862611114E-13;  xe.INF.G := 0;
  xe.SUP.r1 := 1.763222834357;
  xe.SUP.r2 := 2.933226823867E-13;  xe.SUP.G := 0;

  ye.INF.r1 := 1.028170523089E+1;
  ye.INF.r2 := -2.754504467672E-12; ye.INF.G := 0;
  ye.SUP.r1 := 1.028170523089E+1;
  ye.SUP.r2 := -2.754504466670E-12; ye.SUP.G := 0;

  I.inf := 1.6; I.sup := 1.8; {Eingangsintervall }
  w := fI(I);    { w: Wertebereichseinschliessung }
end. { f(x) > 1.028170523089E+1 - 2.754504467672E-12 }
```

#### 4.3.7.4 Die Funktion besitzt ein relatives Maximum und Minimum

$f : [a, b] \rightarrow \mathbb{R}$  sei eine stetige Funktion und habe in  $x_0 \in ]a, b[$  ein relatives Maximum und in  $x_1 > x_0$  ein relatives Minimum. Vor, zwischen und nach den relativen Extrema sei  $f(x)$  jeweils monoton.

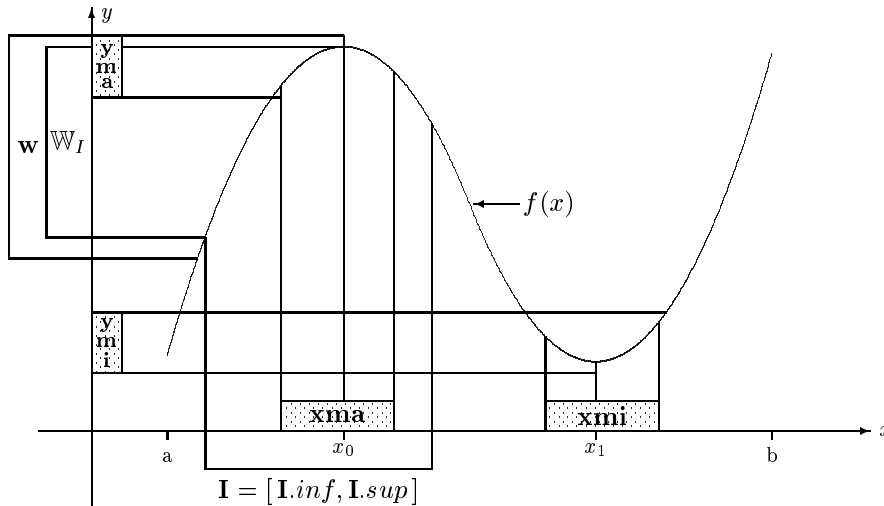


Abbildung 4.3: Einschließung des Wertebereichs  $\mathbb{W}_I$  durch das Intervall  $\mathbf{w}$

$f(x)$  wird auf dem Rechner durch folgende Funktion realisiert:

```
function f (x: rrG): rrG;
```

und liefert zum  $\text{rrG}$ -Argument  $x \in [a, b]$  den i.a. fehlerbehafteten Funktionswert  $\tilde{f}(x) \neq f(x)$ .

Die zugehörige relative Fehlerschranke  $\varepsilon(f)$  wird für alle  $x \in [a, b]$  als bekannt angenommen.

Zu einem beliebig vorgegebenen Intervall  $I$  vom Typ  $\text{rrGinterval}$ , mit:

$$I = [I.\text{inf}, I.\text{sup}] \subseteq [a, b]$$

ist gesucht eine Einschließung  $\mathbf{w} = [\mathbf{w}.\text{inf}, \mathbf{w}.\text{sup}]$  vom Typ  $\text{rrGinterval}$  für den Wertebereich:

$$\mathbb{W}_I := \{y \mid y = f(x) \wedge x \in I = [I.\text{inf}, I.\text{sup}]\} \subseteq \mathbf{w} = [\mathbf{w}.\text{inf}, \mathbf{w}.\text{sup}];$$

In Abb. 4.3 ist  $x_{m.a.}$  ein Intervall vom Typ  $\text{rrGinterval}$ , welches die Maximumstelle  $x_0$  einschließt. Das Intervall  $y_{m.a.}$ , ebenfalls vom Typ  $\text{rrGinterval}$ , ist eine Einschließung aller Funktionswerte  $f(x)$ , mit  $x \in x_{m.a.}$ :

$$\{y \mid y = f(x) \wedge x \in x_{m.a.}\} \subseteq y_{m.a.};$$

Entsprechend ist  $xmi$  ein Intervall vom Typ *rrGinterval*, welches die Minimumstelle  $x_1$  einschließt. Das Intervall  $y_m i$ , ebenfalls vom Typ *rrGinterval*, ist eine Einschließung aller Funktionswerte  $f(x)$ , mit  $x \in xmi$  :

$$\{y \mid y = f(x) \wedge x \in xmi\} \subseteq y_m i;$$

Ist die relative Fehlerschranke  $\varepsilon(f)$  bekannt, so lassen sich die Einschließungen  $xma, yma$  bzw.  $xmi, ymi$  mit dem Programm **MIN\_26** bestimmen.

Mit Hilfe der Funktion

```
function INCL_MAMI (function f (x: rrG): rrG; eps: real;
                   xma, yma, xmi, ymi, I: rrGinterval): rrGinterval;
```

aus dem Modul **rrGi\_ari** erhält man dann die Einschließung  $w$  durch folgenden Funktionsaufruf:

```
w := INCL_MAMI (f, eps, xma, yma, xmi, ymi, I);
```

wobei  $\text{eps}$  wie folgt zu berechnen ist:  $\text{eps} := 1.00001 * \varepsilon(f)$  ;

#### 4.3.7.5 Die Funktion besitzt ein relatives Minimum und Maximum

$f : [a, b] \rightarrow \mathbb{R}$  sei eine stetige Funktion und habe in  $x_0 \in ]a, b[$  ein relatives Minimum und in  $x_1 > x_0$  ein relatives Maximum.

Vor, zwischen und nach den relativen Extrema sei  $f$  jeweils monoton.

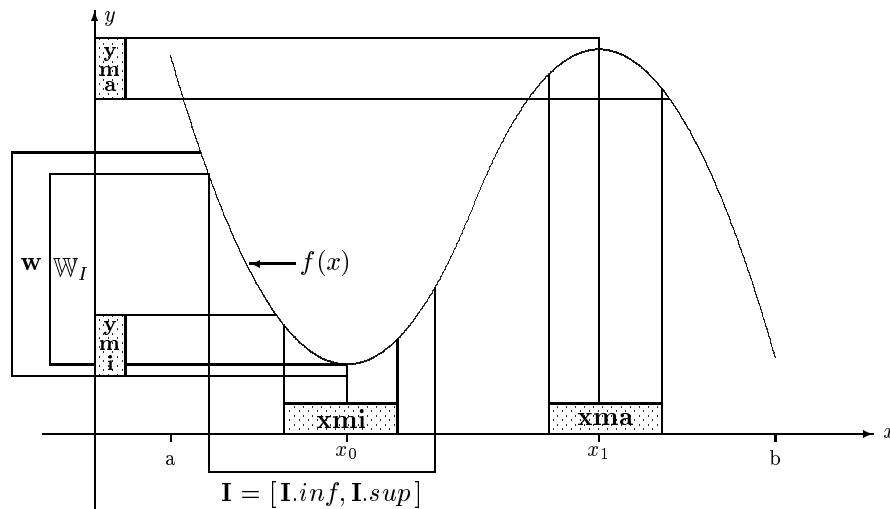


Abbildung 4.4: Einschließung des Wertebereichs  $W_I$  durch das Intervall  $w$

$f(x)$  wird auf dem Rechner durch folgende Funktion realisiert:

```
function f (x: rrG): rrG;
```

und liefert zum rrG-Argument  $x \in [a, b]$  den i.a. fehlerbehafteten Funktionswert  $\tilde{f}(x) \neq f(x)$ .

Die zugehörige relative Fehlerschranke  $\varepsilon(f)$  wird für alle  $x \in [a, b]$  als bekannt angenommen.

Zu einem beliebig vorgegebenen Intervall  $I$  vom Typ *rrGinterval*, mit:

$$I = [I.\text{inf}, I.\text{sup}] \subseteq [a, b]$$

ist gesucht eine Einschließung  $w = [w.\text{inf}, w.\text{sup}]$  vom Typ *rrGinterval* für den Wertebereich:

$$\mathbb{W}_I := \{y \mid y = f(x) \wedge x \in I = [I.\text{inf}, I.\text{sup}]\} \subseteq w = [w.\text{inf}, w.\text{sup}];$$

In Abb. 4.4 ist  $xmi$  ein Intervall vom Typ *rrGinterval*, welches die Minimumstelle  $x_0$  einschließt. Das Intervall  $y mi$ , ebenfalls vom Typ *rrGinterval*, ist eine Einschließung aller Funktionswerte  $f(x)$ , mit  $x \in xmi$ :

$$\{y \mid y = f(x) \wedge x \in xmi\} \subseteq y mi;$$

Entsprechend ist  $xma$  ein Intervall vom Typ *rrGinterval*, welches die Maximumstelle  $x_1$  einschließt. Das Intervall  $y ma$ , ebenfalls vom Typ *rrGinterval*, ist eine Einschließung aller Funktionswerte  $f(x)$ , mit  $x \in xma$ :

$$\{y \mid y = f(x) \wedge x \in xma\} \subseteq y ma;$$

Ist die relative Fehlerschranke  $\varepsilon(f)$  bekannt, so lassen sich die Einschließungen  $xmi, y mi$  bzw.  $xma, y ma$  mit dem Programm **MIN\_26** bestimmen.

Mit Hilfe der Funktion

```
function INCL_MAMI (function f (x: rrG): rrG; eps: real;
                   xma, yma, xmi, ymi, I: rrGinterval): rrGinterval;
```

aus dem Modul **rrGi\_ari** erhält man dann die Einschließung  $w$  durch folgendes Programmbeispiel:

```
program test (input,output);
use rrGi_ari;
var xmi, ymi, xma, yma, I, w : rrGinterval;
    eps : real;

function f (x: rrG): rrG;
begin
  f := ... { Definition von f }
end;
```

```

function fI (I: rrGinterval): rrGinterval;
  { fI liefert zu gegebenem Intervall I die Einschließung w von  $\mathbb{W}_I$  }
  function fn (x: rrG): rrG;
    begin fn := -f(x) end;
begin
  fI := -INCL_MAMI (fn, eps, xmi, -ymi, xma, -yma, I);
end;

begin { Haupt }
  eps := ... ; { eps := 1.00001 * >  $\varepsilon(f)$  }
  xmi := rrGintval (... , ...);
  ymi := rrGintval (... , ...);
  xma := rrGintval (... , ...);
  yma := rrGintval (... , ...);
  I := rrGintval (... , ...);
  w := fI(I); { w liefert eine Einschließung von  $\mathbb{W}_I$  }
end.

```

#### 4.3.7.6 Die Funktion besitzt mehr als zwei Extrema

$f : [a, b] \rightarrow \mathbb{R}$  sei eine reellwertige, stetige Funktion und habe in  $x_0 \in ]a, b[$  nach Abb. 4.5 z.B. drei relative Extrema.

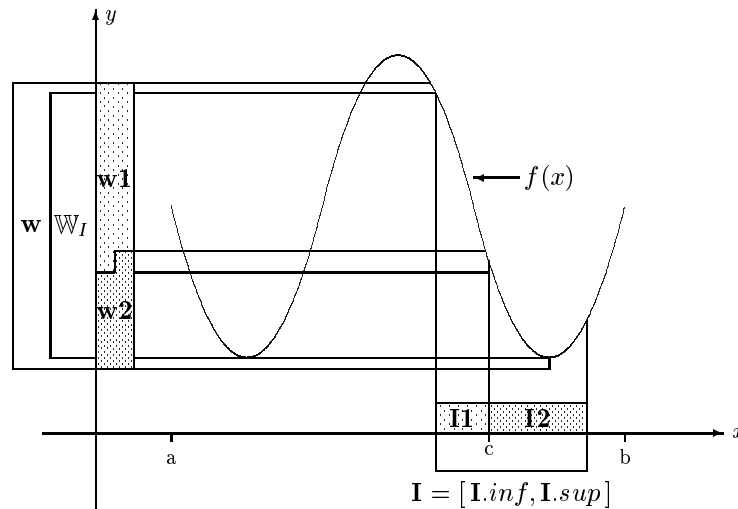


Abbildung 4.5: Funktion mit drei relativen Extrema

$f(x)$  wird auf dem Rechner durch folgende Funktion realisiert:

```

function f (x: rrG): rrG;

```

und liefert zum rrG–Argument  $x \in [a, b]$  den i.a. fehlerbehafteten Funktionswert  $\tilde{f}(x) \neq f(x)$ .

Die zugehörige relative Fehlerschranke  $\varepsilon(f)$  wird für alle  $x \in [a, b]$  als bekannt angenommen.

Zu einem beliebig vorgegebenen Intervall  $I$  vom Typ *rrGinterval*, mit:

$$I = [I.\text{inf}, I.\text{sup}] \subseteq [a, b]$$

ist gesucht eine Einschließung  $w = [w.\text{inf}, w.\text{sup}]$  vom Typ *rrGinterval* für den Wertebereich:

$$\mathbb{W}_I := \{y \mid y = f(x) \wedge x \in I = [I.\text{inf}, I.\text{sup}]\} \subseteq w;$$

Zur Lösung unterteilt man zunächst  $[a, b]$  durch  $c \in [a, b]$  in die beiden Teilintervalle  $[a, c]$  und  $[c, b]$ . Mit den Schnittintervallen

$$I_1 := I \cap [a, c] \quad \text{und} \quad I_2 := I \cap [c, b]$$

berechnet man dann entsprechend den Abbildungen 4.4 bzw. 4.2 die jeweiligen Einschließungen  $w_1$  und  $w_2$ . Die konvexe Hülle von  $w_1$  und  $w_2$  ist dann die gesuchte Einschließung  $w$ :

$$w := w_1 + * w_2 \supseteq I$$

Bei mehr als drei Extrema unterteilt man  $[a, b]$  so in Teilintervalle, daß jedes höchstens zwei Extrema enthält und verfährt dann analog zum obigen Beispiel.

Bei bekannter Fehlerschranke  $\varepsilon(f)$  lassen sich so für alle 'gutmütigen' Funktionen zu einem beliebigen Intervall–Argument  $I$  eine Einschließung  $w$  des entsprechenden Wertebereichs  $\mathbb{W}_I$  bestimmen, wenn man vorher alle Einschließungen der relativen Extrema mit dem Programm **MIN\_26** berechnet hat.



# Kapitel 5

## Komplexe Arithmetik

Mit dem Modul `C_ARI` werden die für das Rechnen mit komplexen Zahlen notwendigen Operatoren, Funktionen und Prozeduren bezüglich des Datentyps `real` bereitgestellt.

Das Modul `LC_ARI` ermöglicht zusätzlich komplexe Rechnungen auf Basis der Typen `real`, `Longreal`, `dotprecision` und `rrG`.

### 5.1 Das Modul C\_ARI

Für Rechnungen mit komplexen Zahlen wird der Datentyp `complex` entsprechend der Definition

```
type complex = record re,im : real end;
```

im Sprachkern von PASCAL-XSC bereitgestellt. Für  $z = x + i \cdot y$ ,  $i = \sqrt{-1}$  und  $z : \text{complex}$  gilt dann:

```
z.re = x    und    z.im = y;
```

#### 5.1.1 Operatoren

Sämtliche in diesem Modul vordefinierten arithmetischen Operatoren liefern den Ergebnistyp `complex`. Als arithmetische Operatoren stehen die monadischen Operatoren `+`, `-` und die vier Grundoperationen `+`, `-`, `*`, `/` mit den drei verschiedenen Rundungsarten zu Verfügung. Dabei werden die Rundungen jeweils auf Real- und Imaginärteil angewandt. Die Ergebnisse der Operationen `+`, `-`, `*` sind **maximalgenau**, d.h. das exakte Ergebnis wird stets zur nächsten Maschinenzahl gerundet. Das Ergebnis der komplexen Division ist jedoch nur **hochgenau**, d.h. zwischen dem exakten Quotienten und dem Maschinenergebnis liegt bzgl. Real- und Imaginärteil **keine** weitere Rasterzahl.

Das Beispiel  $c := (2.00\dots003 - i \cdot 10^{-13}) / (2 + i \cdot 10^{-13})$  liefert z.B. für den Realteil die nur hochgenaue Maschinennäherung  $\widetilde{\text{Re}}(c) = 1.00\dots002$  im Gegensatz zur maximalgenauen Näherung: `1.0000000000001`.

Bezeichnen wir mit  $\varepsilon(G)$  die relative Fehlerschranke der Operationen  $+$ ,  $-$ ,  $*$  und mit  $\varepsilon(D)$  die relative Fehlerschranke der Division und sind  $\hat{\varepsilon}(G), \hat{\varepsilon}(D)$  die entsprechenden Fehlerschranken für die gerichteten Rundungen, so gelten für den normalisierten bzw. denormalisierten real-Bereich die in folgender Tabelle angegebenen relativen Fehlerschranken:

Oper.	Normal.	Denormal.	Rel. Fehlerschranke
+/-	$\varepsilon(G); \hat{\varepsilon}(G)$	$\varepsilon(G); \hat{\varepsilon}(G)$	$\varepsilon(G) = 5 \cdot 10^{-13}; \hat{\varepsilon}(G) = 1 \cdot 10^{-12}$
*	$\varepsilon(G); \hat{\varepsilon}(G)$	$> \varepsilon(G); \hat{\varepsilon}(G)$	$\varepsilon(G) = 5 \cdot 10^{-13}; \hat{\varepsilon}(G) = 1 \cdot 10^{-12}$
/	$\varepsilon(D); \hat{\varepsilon}(D)$	$> \varepsilon(D); \hat{\varepsilon}(D)$	$\varepsilon(D) = 1 \cdot 10^{-12}; \hat{\varepsilon}(D) = 1.5 \cdot 10^{-12}$

Liegen die Ergebnisse von Multiplikation oder Division also im denormalisierten Bereich (3.1), so kann der relative Betragsfehler die Fehlerschranken aus dem normalisierten Bereich überschreiten, vergl. dazu Seite 12. Es ist daher sinnvoll, für die komplexe Multiplikation und Division Funktionen bereitzustellen, die für Real- und Imaginärteilergebnisse aus dem denormalisierten real-Bereich überprüfen, ob die Beträge der aufgetretenen relativen Fehler die Fehlerschranken aus dem normalisierten Bereich übertreffen:

1. **function** MULNEXT\_TEST(var A,B: complex; var error: boolean): complex;  
Berechnet  $A * B$ ; error = TRUE  $\leadsto$  |rel. Fehler|  $> 0.5 \cdot 10^{-12}$ ;
2. **function** MULUP\_TEST(var A,B: complex; var error: boolean): complex;  
Berechnet  $A * > B$ ; error = TRUE  $\leadsto$  |rel. Fehler|  $\geq 1 \cdot 10^{-12}$ ;
3. **function** MULDOWN\_TEST(var A,B: complex; var error: boolean): complex;  
Berechnet  $A * < B$ ; error = TRUE  $\leadsto$  |rel. Fehler|  $\geq 1 \cdot 10^{-12}$ ;
4. **function** DIVNEXT\_TEST(var A,B: complex; var error: boolean): complex;  
Berechnet  $A / B$ ; error = TRUE  $\leadsto$  |rel. Fehler|  $> 1 \cdot 10^{-12}$ ;
5. **function** DIVUP\_TEST(var A,B: complex; var error: boolean): complex;  
Berechnet  $A / > B$ ; error = TRUE  $\leadsto$  |rel. Fehler|  $> 1.5 \cdot 10^{-12}$ ;
6. **function** DIVDOWN\_TEST(var A,B: complex; var error: boolean): complex;  
Berechnet  $A / < B$ ; error = TRUE  $\leadsto$  |rel. Fehler|  $> 1.5 \cdot 10^{-12}$ ;

Die Vergleichsoperatoren  $= <> < <= > >=$  beziehen sich auf Real- und Imaginärteil. Für a,b vom Typ *complex* gilt demnach:

$$a <= b \iff (a.re <= b.re) \text{ and } (a.im <= b.im)$$

Entsprechende Vergleiche auch mit einem *integer*- oder *real*-Operanden sind zulässig.

rechter Operand	integer real	complex
linker Operand		
monadisch		+, -
integer real		◦ ∨
complex	◦ ∨	◦ ∨

Tabelle 5.1: Die Operatoren des Moduls C\_ARI

$$\circ \in \{+, +<, +>, -, -<, ->, *, *<, *>, /, /<, />\}$$

$$\vee \in \{=, <>, <, <=, >, >=\}$$

### 5.1.2 Transferfunktionen

Zur Wandlung zwischen den Typen *integer/real* und *complex* werden folgende Transferfunktionen bereitgestellt:

Funktion	Ergebnistyp	Bedeutung
compl (r1,r2)	<i>complex</i>	$z = r1 + i \cdot r2; \quad i = \sqrt{-1}$
compl (r)	<i>complex</i>	$z = r1 + i \cdot 0; \quad i = \sqrt{-1}$
re (z)	<i>real</i>	Realteil von z
im (z)	<i>real</i>	Imaginärteil von z

Tabelle 5.2: Transferfunktionen des Moduls C\_ARI  
 r, r1, r2 = *real*-Ausdruck;     z = *complex*-Ausdruck

### 5.1.3 Überladungen des Zuweisungsoperators

Die Wandlung von *real/integer* nach *complex* wird auch in Form einer überladenen Zuweisung bereitgestellt:

Zuweisung	Bedeutung
z := r	z := compl (r)

z = *complex*-**Variable**;     r = *real*-Ausdruck

### 5.1.4 Standardfunktionen

Es stehen die folgenden Standardfunktionen zur Verfügung:

Funktion	Bedeutung	Aufruf	Ergebnistyp
$z^2 = z \cdot z$	Quadrat	<code>sqr (z)</code>	complex
$\sqrt{z}$	Quadratwurzel (Realteil $\geq 0$ )	<code>sqrt (z)</code>	complex
$\sqrt{z}$	Quadratwurzel (Realteil $\geq 0$ )	<code>sqrt_u (z)</code>	complex
$\sqrt{z}$	Quadratwurzel	<code>sqrt_negim (z)</code>	complex
$e^z$	Exponentialfunktion	<code>exp (z)</code>	complex
$\ln(z)$	Natürlicher Logarithmus	<code>ln (z)</code>	complex
$\ln(1 + z)$	Natürlicher Logarithmus	<code>ln1p (z)</code>	complex
$\sin(z)$	Sinus	<code>sin (z)</code>	complex
$\cos(z)$	Kosinus	<code>cos (z)</code>	complex
$\tan(z)$	Tangens	<code>tan (z)</code>	complex
$\cot(z)$	Kotangens	<code>cot (z)</code>	complex
$\sinh(z)$	Hyperbolischer Sinus	<code>sinh (z)</code>	complex
$\cosh(z)$	Hyperbolischer Kosinus	<code>cosh (z)</code>	complex
$\tanh(z)$	Hyperbolischer Tangens	<code>tanh (z)</code>	complex
$\coth(z)$	Hyperbolischer Kotangens	<code>coth (z)</code>	complex
$\bar{z} = x - i \cdot y$	Konjugation von $z = x + i \cdot y$	<code>conj (z)</code>	complex
$\varphi$	Argument von $z = r \cdot e^{i\varphi}$	<code>arg (z)</code>	real
$\varphi$	Argument von $1 + z = r \cdot e^{i\varphi}$	<code>arg1p (z)</code>	real
$r = \sqrt{x^2 + y^2}$	Absolutbetrag von $z = r \cdot e^{i\varphi}$	<code>abs (z)</code>	real

Tabelle 5.3: Standardfunktionen des Moduls C\_ARI,  $z$ : complex

**Anmerkungen** zu den Funktionen: `sqrt(z)`, `sqrt_u(z)`, `sqrt_negim`, `ln(z)`, `ln1p(z)`:

- **sqrt(z)** Der Verzweigungsschnitt ist die negative, reelle Achse, wobei diese Achse zur oberen Halbebene zählt.
- **sqrt\_u(z)** Der Verzweigungsschnitt ist die negative, reelle Achse, wobei diese Achse zur unteren Halbebene zählt.
- **sqrt\_negim(z)** Der Verzweigungsschnitt ist die negative, imaginäre Achse, wobei diese Achse zur linken Halbebene zählt.
- **ln(z)** Der Verzweigungsschnitt ist die negative, reelle Achse, wobei diese Achse zur oberen Halbebene zählt.

- **ln1p(z)** Der Verzweigungsschnitt ist der negative Teil der reellen Achse von  $-\infty$  bis  $-1$ , wobei diese Achse zur oberen Halbebene zählt.

Für alle Funktionen aus Tabelle 5.3 ist der Betrag des relativen Fehlers kleiner als  $5.00000006 \cdot 10^{-13}$ , so daß alle Funktionen hochgenau sind, d.h. das Maschinenergebn unterscheidet sich vom optimal gerundeten Funktionsergebnis um maximal eine Einheit in der letzten Mantissenstelle. Diese Abweichung von 1 ulp tritt allerdings nur in **sehr seltenen Ausnahmefällen** auf!

Fallen Real- oder Imaginärteil in den denormalisierten Bereich (3.1), d.h.

$$0 \leq |\operatorname{Re}(f(z))|, |\operatorname{Im}(f(z))| < 10^{-255} = \text{MIN\_REAL},$$

und werden dadurch hintere, von Null verschiedene Ziffern der Ergebnismantisse gleich Null gesetzt oder entsteht Underflow, so daß die Funktionsergebnisse dadurch nicht mehr hochgenau sind, so erfolgt die Fehlermeldung:

**rel. error exceeds error bound**

mit anschließendem Programmabbruch.

Die noch unvollständige Liste der komplexen Standardfunktionen aus Tabelle 5.3 wird in folgenden Updates erweitert.

### 5.1.5 Ein-/Ausgabeprozeduren

Es stehen die gewohnten Prozeduren

```
procedure read (var f: text; var a: complex);
procedure write (var f: text; a: complex);
```

mit optionalem Fileparameter, beliebig vielen Ein-/Ausgabeparametern, jedoch ohne Formatspezifikationen zur Verfügung.

Die Eingabe einer komplexen Zahl  $z = x + i \cdot y$  muß in der Form

$(x, y)$       oder in der Form       $x$

erfolgen. Im zweiten Fall wird der Imaginärteil  $y$  automatisch auf 0 gesetzt.  $x$  und  $y$  sind dabei *real*-Konstanten. Die Ausgabe einer komplexen Zahl erfolgt stets in der Form

$(x, y)$

mit dem Standardformat für die *real*-Größen  $x$  und  $y$  aus Tabelle 3.4.

## 5.2 Das Modul LC\_ARI

### 5.2.1 Der Typ Lcomplex

Für allgemeine Gleitkommarechnungen mit komplexen Zahlen wird neben dem Typ *complex* zusätzlich der Datentyp *Lcomplex* entsprechend der Definition

```
type Lcomplex = record Re,Im : Longreal end;
```

bereitgestellt. Für  $z = x + i \cdot y$ ,  $i = \sqrt{-1}$  und  $z : Lcomplex$  gilt dann:  
 $z.Re = x$  und  $z.Im = y$ ;

#### 5.2.1.1 Operatoren

Wir betrachten zunächst alle in diesem Modul vordefinierten arithmetischen Operatoren mit dem Ergebnistyp *Lcomplex*. Als arithmetische Operatoren stehen die monadischen Operatoren  $+$ ,  $-$  und die vier Grundoperationen  $+$ ,  $-$ ,  $*$ ,  $/$  mit den drei verschiedenen Rundungsarten zu Verfügung. Dabei werden die Rundungen jeweils auf Real- **und** Imaginärteil angewandt.

rechter Operand	integer real	Longreal	complex	Lcomplex
linker Operand				
monadisch				$+, -$
integer real				$\circ$ $\vee$
Longreal			$\circ$ $\vee$	$\circ$ $\vee$
complex		$\circ$ $\vee$		$\circ$ $\vee$
Lcomplex	$\circ$ $\vee$	$\circ$ $\vee$	$\circ$ $\vee$	$\circ$ $\vee$

Tabelle 5.4: Operatoren des Moduls LC\_ARI mit Ergebnistyp *Lcomplex*, *boolean*

$\circ \in \{+, +<, +>, -, -<, ->, *, *<, *>, /, /<, />\}$

$\vee \in \{=, <>, <, <=, >, >=\}$

Die Ergebnisse der Operationen  $+$ ,  $-$ ,  $*$  sind **maximalgenau**, d.h. das exakte Ergebnis wird stets zur nächsten Maschinenzahl gerundet. Das Ergebnis der komplexen Division ist im Gegensatz dazu jedoch nur **hochgenau**, d.h. zwischen dem exakten Quotienten und dem Maschinenergebnis liegt bezüglich Real- und Imaginärteil **keine** weitere Rasterzahl.

Das Beispiel  $c := (2.00\dots003 - i \cdot 10^{-21}) / (2 + i \cdot 10^{-21})$  liefert z.B. für den Realteil die nur hochgenaue Maschinennäherung  $\widetilde{\text{Re}}(c) = 1.00\dots002$  im Gegensatz zur maximalgenauen Näherung:  $1.000000000000000000001$ .

Bezeichnen wir mit  $\varepsilon(G)$  die relative Fehlerschranke der Operationen  $+$ ,  $-$ ,  $*$  und mit  $\varepsilon(D)$  die relative Fehlerschranke der Division und sind  $\hat{\varepsilon}(G), \hat{\varepsilon}(D)$  die entsprechenden Fehlerschranken für die gerichteten Rundungen, so gelten für den normalisierten bzw. denormalisierten real-Bereich die in folgender Tabelle angegebenen relativen Fehlerschranken:

Oper.	Normal.	Denormal.	Rel. Fehlerschranke
+/-	$\varepsilon(G); \hat{\varepsilon}(G)$	$\varepsilon(G); \hat{\varepsilon}(G)$	$\varepsilon(G) = 5 \cdot 10^{-21}; \hat{\varepsilon}(G) = 1 \cdot 10^{-20}$
*	$\varepsilon(G); \hat{\varepsilon}(G)$	$> \varepsilon(G); \hat{\varepsilon}(G)$	$\varepsilon(G) = 5 \cdot 10^{-21}; \hat{\varepsilon}(G) = 1 \cdot 10^{-20}$
/	$\varepsilon(D); \hat{\varepsilon}(D)$	$> \varepsilon(D); \hat{\varepsilon}(D)$	$\varepsilon(D) = 1 \cdot 10^{-20}; \hat{\varepsilon}(D) = 1.5 \cdot 10^{-20}$

Liegen die Real- bzw. Imaginärteil-Ergebnisse von Multiplikation oder Division also im denormalisierten Bereich, so kann der relative Betragsfehler die Fehlerschranken aus dem normalisierten Bereich überschreiten, vergl. dazu Seite 23. Es ist daher sinnvoll, für die komplexe Multiplikation und Division Funktionen bereitzustellen, die für Real- und Imaginärteilergebnisse aus dem denormalisierten real-Bereich überprüfen, ob die Beträge der aufgetretenen relativen Fehler die Fehlerschranken aus dem normalisierten Bereich übertreffen:

1. **function** MULNEXT\_TEST (var A,B : Lcomplex;  
var error : boolean) : Lcomplex;  
Berechnet A \* B; error = TRUE  $\leadsto$  |rel. Fehler|  $> 0.5 \cdot 10^{-20}$ ;
2. **function** MULUP\_TEST (var A,B : Lcomplex;  
var error : boolean) : Lcomplex;  
Berechnet A \* $>$  B; error = TRUE  $\leadsto$  |rel. Fehler|  $\geq 1 \cdot 10^{-20}$ ;
3. **function** MULDOWN\_TEST (var A,B : Lcomplex;  
var error : boolean) : Lcomplex;  
Berechnet A \* $<$  B; error = TRUE  $\leadsto$  |rel. Fehler|  $\geq 1 \cdot 10^{-20}$ ;
4. **function** DIVNEXT\_TEST (var A,B : Lcomplex;  
var error : boolean) : Lcomplex;  
Berechnet A / B; error = TRUE  $\leadsto$  |rel. Fehler|  $> 1 \cdot 10^{-20}$ ;
5. **function** DIVUP\_TEST (var A,B : Lcomplex;  
var error : boolean) : Lcomplex;  
Berechnet A / $>$  B; error = TRUE  $\leadsto$  |rel. Fehler|  $> 1.5 \cdot 10^{-20}$ ;
6. **function** DIVDOWN\_TEST (var A,B : Lcomplex;  
var error : boolean) : Lcomplex;  
Berechnet A / $<$  B; error = TRUE  $\leadsto$  |rel. Fehler|  $> 1.5 \cdot 10^{-20}$ ;

Die Vergleichsoperatoren  $=, <>, <, <=, >, >=$  beziehen sich auf Real- und Imaginärteil. Für  $a, b$  vom Typ *Lcomplex* gilt demnach:

$$a <= b \iff (a.Re <= b.Re) \text{ \textbf{and} } (a.Im <= b.Im)$$

Entsprechende Vergleiche auch mit einem *integer/real*- oder *Longreal*-Operanden sind zulässig.

### 5.2.1.2 Transferfunktionen

Zur Wandlung zwischen den Typen *integer/real/Longreal/complex* und dem Typ *Lcomplex* werden folgende Transferfunktionen bereitgestellt:

Funktion	Ergebnistyp	Bedeutung
LCOMPL (x1,x2)	<i>Lcomplex</i>	$z = x1 + i \cdot x2; \quad i = \sqrt{-1}$
LCOMPL (x)	<i>Lcomplex</i>	$z = x + i \cdot 0; \quad i = \sqrt{-1}$
LONG (c)	<i>Lcomplex</i>	$z = \text{LONG}(c.Re) + i \cdot \text{LONG}(c.Im)$
RE (z)	<i>Longreal</i>	Realteil von z
IM (z)	<i>Longreal</i>	Imaginärteil von z
SHORT (z)	<i>complex</i>	Rundung zum nächsten c
SHORTUP (z)	<i>complex</i>	Rundung zum nächstgrößeren c
SHORTDOWN (z)	<i>complex</i>	Rundung zum nächstkleineren c
<b>short_test</b> ( z,error ), <b>shortup_test</b> ( z,error ), <b>shortdown_test</b> ( z,error ) Rundgn. wie oben; error = true $\Leftrightarrow$ rel. Fehler $> 5 \cdot 10^{-13}$ bzw. $\geq 10^{-12}$		

Tabelle 5.5: Transferfunktionen des Moduls LC\_ARI

$x, x1, x2 = \textit{integer/real/Longreal}$ -Ausdruck;  $c = \textit{complex}$ -Ausdruck;  
 $z = \textit{Lcomplex}$ -Ausdruck;  $\textit{error} = \textit{boolean}$ -Variable;

### 5.2.1.3 Überladungen des Zuweisungsoperators

Die Wandlung von *integer/real/Longreal/complex* nach *Lcomplex* wird auch in Form einer überladenen Zuweisung bereitgestellt:

Zuweisung	Bedeutung
$z := x$	$z := \text{LCOMPL}(x)$
$z := c$	$z := \text{LONG}(c)$

$z = \textit{Lcomplex-Variable}$ ;  $x = \textit{integer/real/Longreal}$ -Ausdruck;  
 $c = \textit{complex}$ -Ausdruck



## 5.2.1.4 Standardfunktionen

Es stehen die folgenden Standardfunktionen zur Verfügung:

Aufruf	Ergebnistyp	Funktion	$\varepsilon(\operatorname{Re}\{f(z)\})$	$\varepsilon(\operatorname{Im}\{f(z)\})$
sqr (z)	Lcomplex	$z^2 = z \cdot z$	$5.0002 \cdot 10^{-21}$	$5.0002 \cdot 10^{-21}$
sqrt (z)	Lcomplex	$\sqrt{z}$	$5.0045 \cdot 10^{-21}$	$5.0045 \cdot 10^{-21}$
sqrt_u (z)	Lcomplex	$\sqrt{z}$	$5.0045 \cdot 10^{-21}$	$5.0045 \cdot 10^{-21}$
sqrt_negim (z)	Lcomplex	$\sqrt{z}$	$5.0045 \cdot 10^{-21}$	$5.0045 \cdot 10^{-21}$
exp (z)	Lcomplex	$e^z$	$5.0016 \cdot 10^{-21}$	$5.0016 \cdot 10^{-21}$
ln (z)	Lcomplex	$\ln(z)$	$5.0053 \cdot 10^{-21}$	$5.0054 \cdot 10^{-21}$
ln1p (z)	Lcomplex	$\ln(1+z)$	$5.0030 \cdot 10^{-21}$	$5.0059 \cdot 10^{-21}$
sin (z)	Lcomplex	$\sin(z)$	$5.0028 \cdot 10^{-21}$	$5.0038 \cdot 10^{-21}$
cos (z)	Lcomplex	$\cos(z)$	$5.0028 \cdot 10^{-21}$	$5.0038 \cdot 10^{-21}$
tan (z)	Lcomplex	$\tan(z)$	$5.0098 \cdot 10^{-21}$	$5.0120 \cdot 10^{-21}$
cot (z)	Lcomplex	$\cot(z)$	$5.0098 \cdot 10^{-21}$	$5.0120 \cdot 10^{-21}$
sinh (z)	Lcomplex	$\sinh(z)$	$5.0038 \cdot 10^{-21}$	$5.0028 \cdot 10^{-21}$
cosh (z)	Lcomplex	$\cosh(z)$	$5.0028 \cdot 10^{-21}$	$5.0038 \cdot 10^{-21}$
tanh (z)	Lcomplex	$\tanh(z)$	$5.0120 \cdot 10^{-21}$	$5.0098 \cdot 10^{-21}$
coth (z)	Lcomplex	$\coth(z)$	$5.0120 \cdot 10^{-21}$	$5.0098 \cdot 10^{-21}$
conj (z)	Lcomplex	$\bar{z} = x - i \cdot y$	0.0	0.0
arg (z)	Longreal	$\varphi ; \quad z = r \cdot e^{i \cdot \varphi} ;$	$5.0054 \cdot 10^{-21}$	
arg1p (z)	Longreal	$\varphi ; \quad 1 + z = r \cdot e^{i \cdot \varphi} ;$	$5.0059 \cdot 10^{-21}$	
abs (z)	Longreal	$r = \sqrt{x^2 + y^2} ;$	$5.0013 \cdot 10^{-21}$	

Tabelle 5.6: Standardfunktionen des Moduls LC\_ARI,  $z$ : Lcomplex

**Anmerkungen** zu den Funktionen: sqrt(z), sqrt\_u(z), sqrt\_negim, ln(z), ln1p(z):

- **sqrt(z)** Der Verzweigungsschnitt ist die negative, reelle Achse, wobei diese Achse zur oberen Halbebene zählt.
- **sqrt\_u(z)** Der Verzweigungsschnitt ist die negative, reelle Achse, wobei diese Achse zur unteren Halbebene zählt.
- **sqrt\_negim(z)** Der Verzweigungsschnitt ist die negative, imaginäre Achse, wobei diese Achse zur linken Halbebene zählt.
- **ln(z)** Der Verzweigungsschnitt ist die negative, reelle Achse, wobei diese Achse zur oberen Halbebene zählt.
- **ln1p(z)** Der Verzweigungsschnitt ist der negative Teil der reellen Achse von  $-\infty$  bis  $-1$ , wobei diese Achse zur oberen Halbebene zählt.

Für alle Funktionen aus Tabelle 5.6 ist im normalisierten Bereich der Betrag des relativen Fehlers kleiner als die jeweils angegebene Fehlerschranke. Fallen jedoch Real- oder Imaginärteil in den denormalisierten Bereich (3.2), d.h.

$$0 \leq |\operatorname{Re}(f(z))|, |\operatorname{Im}(f(z))| < 10^{-511} = \text{MINLONGREAL},$$

und werden dadurch hintere, von Null verschiedene Ziffern der Ergebnismantisse gleich Null gesetzt oder entsteht Underflow mit dem relativen Fehler 1, so erfolgt die Fehlermeldung:

**rel. error exceeds error bound**

mit anschließendem Programmabbruch.

### 5.2.1.5 Ein-/Ausgabeprozeduren

Es stehen die gewohnten Prozeduren

```
procedure read (var f: text; var a: Lcomplex);
procedure write (var f: text; a: Lcomplex);
```

mit optionalem Fileparameter, beliebig vielen Ein-/Ausgabeparametern, jedoch ohne Formatspezifikationen zur Verfügung.

Die Eingabe einer komplexen Zahl  $z = x + i \cdot y$  muß in der Form

$(x, y)$       oder in der Form       $x$

erfolgen. Im zweiten Fall wird der Imaginärteil  $y$  automatisch auf 0 gesetzt.  $x$  und  $y$  sind dabei *real*-Konstanten. Die Ausgabe einer komplexen Zahl erfolgt stets in der Form

$(x, y)$

mit dem Standardformat für die *Longreal*-Größen  $x$  und  $y$  aus Tabelle 3.4.

### 5.2.2 Exakte Auswertung komplexer Ausdrücke

In **PASCAL-XSC** können mit Hilfe der Operatoren  $+$ ,  $-$ ,  $*$  Variablen, Konstanten und spezielle Funktionsaufrufe vom Typ *integer*, *real*, *complex*, *dotprecision* oder Produkte vom Typ *integer*, *real*, *complex* oder Skalarprodukte vom Typ *real* oder *complex* mit Hilfe des langen Akkus **rundungsfehlerfrei** berechnet und durch die Anweisungen

```
z := # *( exakter Ausdruck vom Typ complex);
z := # >( exakter Ausdruck vom Typ complex);
z := # <( exakter Ausdruck vom Typ complex);
```

mit nur einer einzigen Rundung bei Real- und Imaginärteil in eine Variable  $z$  vom Typ *complex* abgespeichert werden. Dieses #-Konzept steht auch für die BCD-Version in **vollem Umfang** zur Verfügung; es ist jedoch **nicht** anwendbar, wenn nur einer der Operanden im exakt auszuwertenden, komplexen Ausdruck vom Typ Longreal oder Lcomplex ist!

Um das #-Konzept für diese Datentypen wenigstens simulieren zu können, benötigt man den Datentyp **cdotprecision**:

```
type cdotprecision = record Re, Im: dotprecision end;
```

### 5.2.2.1 Operatoren

Wie auf Seite 42 beschrieben, muß auch jetzt z.B. für einen Additionsoperator mit Lcomplex-Operanden und einem Ergebnis vom Typ cdotprecision ein **neuer** Operatorname (*plus\_c*) gewählt werden, da der normale + Operator mit den gleichen Operanden schon für den Ergebnistyp Lcomplex vergeben ist. In nachfolgender Tabelle 5.7 sind die entsprechenden neuen Operatornamen zusammengestellt:

Operation	Operator-Name	Ergebnistyp
Addition	<b>plus_c</b>	cdotprecision
Subtraktion	<b>minus_c</b>	cdotprecision
Multiplikation	<b>times_c</b>	cdotprecision

Tabelle 5.7: Operatornamen mit Ergebnistyp cdotprecision

rechter Operand	integer/real Longreal	complex Lcomplex	cdotprecision (dotprecision)
linker Operand			
integer/real Longreal	times_c plus_c, minus_c	times_c plus_c, minus_c	* plus_c, minus_c
complex Lcomplex	times_c plus_c, minus_c	times_c plus_c, minus_c	* plus_c, minus_c
cdotprecision (dotprecision)	* plus_c, minus_c	* plus_c, minus_c	* plus_c, minus_c

Tabelle 5.8: Operatoren des Moduls LC\_ARI mit Ergebnistyp cdotprecision

In Tabelle 5.8 sind für die Operatoren *plus\_c*, *minus\_c*, *times\_c*, \* alle erlaubten Operandentypen angegeben. Mit einem *dotprecision*-Operanden existiert jedoch nur der \* Operator mit dem zweiten Operandentyp *cdotprecision*!

**Wertzuweisungen** an eine Variable vom Typ `cdotprecision`:

Durch Überladen des Zuweisungsoperators `:=` sind folgende Wertzuweisungen möglich:

```
cdotprecision-Variable := Lcomplex-Ausdruck;
cdotprecision-Variable := complex-Ausdruck;
cdotprecision-Variable := dotprecision-Ausdruck;
cdotprecision-Variable := Longreal-Ausdruck;
cdotprecision-Variable := real/integer-Ausdruck;
```

### 5.2.2.2 Transferfunktionen

Zur Wandlung zwischen den Typen

*integer/real/Longreal/dotprecision*

einerseits und dem Typ `cdotprecision` andererseits werden folgende Transferfunktionen bereitgestellt:

Funktion	Ergebnistyp	Bedeutung
DOTCOMPL (x1,x2)	<i>cdotprecision</i>	$dc = x1 + i \cdot x2; \quad i = \sqrt{-1}$
DOTCOMPL (x)	<i>cdotprecision</i>	$dc = x + i \cdot 0; \quad i = \sqrt{-1}$
DOTCOMPL (d1,d2)	<i>cdotprecision</i>	$dc = d1 + i \cdot d2; \quad i = \sqrt{-1}$
DOTCOMPL (d)	<i>cdotprecision</i>	$dc = d + i \cdot 0; \quad i = \sqrt{-1}$
Re (dc)	<i>dotprecision</i>	Realteil von dc
Im (dc)	<i>dotprecision</i>	Imaginärteil von dc

$x, x1, x2 = \text{integer/real/Longreal-Ausdruck}; \quad dc = \text{cdotprecision-Variablen};$   
 $d, d1, d2 = \text{dotprecision-Variablen};$

### 5.2.2.3 Simulation von SUM-Ausdrücken

Zur Auswertung von **SUM**-Ausdrücken im `#`-Konzept stehen die folgenden Datentypen und Prozeduren zur Verfügung:

- **type** `rvector` = **dynamic** array [\*] of **real**;
  - **type** `Lrvector` = **dynamic** array [\*] of **Longreal**;
  - **type** `cvector` = **dynamic** array [\*] of **complex**;
  - **type** `Lcvector` = **dynamic** array [\*] of **Lcomplex**;
1. procedure **ADDNACCU** (var dc: `cdotprecision`; var a: `rvector`);  
**Exakte** Addition des real-Feldes a zur `cdotprecision`-Variablen dc;
  2. procedure **ADDNACCU** (var dc: `cdotprecision`; var a: `Lrvector`);  
**Exakte** Addition des Longreal-Feldes a zur `cdotprecision`-Variablen dc;

3. procedure **ADDNACCU** (var dc: cdotprecision; var a: cvector);  
**Exakte** Addition des complex-Feldes a zur cdotprecision-Variablen dc;
4. procedure **ADDNACCU** (var dc: cdotprecision; var a: Lcvector);  
**Exakte** Addition des Lcomplex-Feldes a zur cdotprecision-Variablen dc;
5. procedure **SUBNACCU** (var dc: cdotprecision; var a: rvector);  
**Exakte** Subtraktion des real-Feldes a von der cdotprecision-Variablen dc;
6. procedure **SUBNACCU** (var dc: cdotprecision; var a: Lrvector);  
**Exakte** Subtraktion des Longreal-Feldes a von der cdotprecision-Variablen dc;
7. procedure **SUBNACCU** (var dc: cdotprecision; var a: cvector);  
**Exakte** Subtraktion des complex-Feldes a von der cdotprecision-Variablen dc;
8. procedure **SUBNACCU** (var dc: cdotprecision; var a: Lcvector);  
**Exakte** Subtraktion des Lcomplex-Feldes a von der cdotprecision-Variablen dc;
9. procedure **PADDNACCU** (var dc: cdotprecision; var a,b: rvector);  
**Exakte** Addition des real-Skalarproduktes zur cdotprecision-Variablen dc;
10. procedure **PADDNACCU** (var dc: cdotprecision; var a,b: Lrvector);  
**Exakte** Addition des Longreal-Skalarproduktes zur cdotprecision-Variablen dc;
11. procedure **PADDNACCU** (var dc: cdotprecision; var a,b: cvector);  
**Exakte** Addition des complex-Skalarproduktes zur cdotprecision-Variablen dc;
12. procedure **PADDNACCU** (var dc: cdotprecision; var a,b: Lcvector);  
**Exakte** Addition des Lcomplex-Skalarproduktes zur cdotprecision-Variablen dc;
13. procedure **PSUBNACCU** (var dc: cdotprecision; var a,b: rvector);  
**Exakte** Subtraktion des real-Skalarproduktes von der cdotprecision-Variablen dc;
14. procedure **PSUBNACCU** (var dc: cdotprecision; var a,b: Lrvector);  
**Exakte** Subtraktion des Longreal-Skalarproduktes von der cdotprecision-Variablen dc;
15. procedure **PSUBNACCU** (var dc: cdotprecision; var a,b: cvector);  
**Exakte** Subtraktion des complex-Skalarproduktes von der cdotprecision-Variablen dc;
16. procedure **PSUBNACCU** (var dc: cdotprecision; var a,b: Lcvector);  
**Exakte** Subtraktion des Lcomplex-Skalarproduktes von der cdotprecision-Variablen dc;

Mit Hilfe der folgenden Funktionen kann der Wert einer *cdotprecision*-Variablen durch die drei verschiedenen Rundungsarten in einen Wert vom Typ *Lcomplex* umgewandelt werden. Fallen dabei Real- oder Imaginärteil in den denormalisierten Zahlenbereich (3.2), so kann man mit den Funktionen

LONGNEXT\_TEST, LONGUP\_TEST, LONGDOWN\_TEST

überprüfen, ob der aktuelle, relative Fehler bei den entsprechenden Rundungen die jeweiligen Fehlerschranken für den normalisierten Bereich übersteigen; vergleiche dazu die Bemerkungen von Seite 23.

- function **LONGNEXT** (var dc: *cdotprecision*): *Lcomplex*;  
y:=LONGNEXT(dc) liest den dc-Wert in die nächste *Lcomplex*-Zahl y;
- function **LONGDOWN** (var dc: *cdotprecision*): *Lcomplex*;  
y:=LONGDOWN(dc) liest den dc-Wert in die nächstkleinere *Lcomplex*-Zahl y;
- function **LONGUP** (var dc: *cdotprecision*): *Lcomplex*;  
y:=LONGUP(dc); liest den dc-Wert in die nächstgrößere *Lcomplex*-Zahl y;
- function **LONGNEXT\_TEST** (var dc: *cdotprecision*;  
var error: boolean): *Lcomplex*;  
Liest den dc-Wert in die nächste *Lcomplex*-Zahl;  
error = TRUE  $\iff$  relativer Fehler  $> 0.5 \cdot 10^{-20}$ ; (durch Nullsetzen hinterer Ziffern der Ergebnismantisse im Bereich (3.2) ).
- function **LONGUP\_TEST** (var dc: *cdotprecision*;  
var error: boolean): *Lcomplex*;  
Liest den dc-Wert in die nächstgrößere *Lcomplex*-Zahl;  
error = TRUE  $\iff$  relativer Fehler  $\geq 1 \cdot 10^{-20}$ ; (durch Nullsetzen hinterer Ziffern der Ergebnismantisse im Bereich (3.2) ).
- function **LONGDOWN\_TEST** (var dc: *cdotprecision*;  
var error: boolean): *Lcomplex*;  
Liest den dc-Wert in die nächstkleinere *Lcomplex*-Zahl;  
error = TRUE  $\iff$  relativer Fehler  $\geq 1 \cdot 10^{-20}$ ; (durch Nullsetzen hinterer Ziffern der Ergebnismantisse im Bereich (3.2) ).

## 5.2.2.4 Beispiele

## 1. Beispiel:

```

program test (input,output);
var Lz      : Lcomplex;      a,b      : Lcvector[1..100];
    error   : boolean;      dc1,dc2  : cdotprecision;
    xl      : Longreal;     z       : complex;

begin
  ...
  dc1 := 0;      PADDNACCU(dc1,a,b);
  dc2 := xl times_c 3 minus_c Lz times_c z plus_c dc1;
  Lz := LONGNEXT_TEST (dc2,error);

```

**end.**

Obiges Programm berechnet den Ausdruck

$$dc2 = 3 \cdot xl - z \cdot Lz + \sum_{k=1}^{100} a[k] \cdot b[k]$$

**rundungsfehlerfrei** und schreibt seinen Wert mit **nur einer einzigen Rundung** in die nächstgelegene *Lcomplex*-Zahl der Variablen Lz. Falls beim Real- oder Imaginärteil durch das jeweilige Runden zur nächstgelegenen *Longreal*-Zahl der relative Betragsfehler größer wird als  $\varepsilon(G) = 0.5 \cdot 10^{-20}$ , erhält die Variable error den Wert TRUE.

## 2. Beispiel:

Mit Hilfe der Operatoren **\*** und **plus\_c** zwischen *cdotprecision*- und *Lcomplex*-Operanden kann ein komplexwertiges Polynom

$$P(z) = \sum_{k=0}^n a_k \cdot z^k$$

mit der folgenden Funktion nach dem Horner Schema **rundungsfehlerfrei** berechnet werden:

```

function HORNER_cdot (var a: Lcvector;
                    z: Lcomplex): cdotprecision;

var k : integer;
    y : cdotprecision;

begin
  y := a[ub(a)];
  for k := ub(a)-1 downto lb(a) do
    y := y * z plus_c a[k];
  HORNER := y
end;

```

Die Koeffizienten

$$\begin{aligned} a_0 &= 120 + 90i; & a_1 &= -54 + 22i; & a_2 &= 6 + 7i; \\ a_3 &= -5 - 3i; & a_4 &= 1 + 0i; \end{aligned}$$

realisieren ein komplexwertiges Polynom  $P(z)$  mit den Nullstellen:

$$z_{1,2} = \pm\sqrt{2} \cdot (1 + 2i), \quad z_3 = 5, \quad z_4 = 3i.$$

Mit den Deklarationen:

```
var a: Lcvector[0..4];
    z,y: Lcomplex;   dotz: cdotprecision;
```

und den Anweisungen:

```
z      := 5.00000000000000000001;   { z ≈ z3 }
dotz   := HORNER_cdot (a,z);
y      := LONGNEXT(dotz);
```

erhält man mit der BCD-Version das Ergebnis:

$$y = \tilde{P}(z) = 1.31000000000000000001 \cdot 10^{-18} - 1.3300 \dots 00 \cdot 10^{-18} \cdot i$$

mit **nur einer Rundung** bzgl. Real- und Imaginärteil.



### 5.2.3 Der Typ `rrGcomplex`

Zur Implementierung **hochgenauer**, komplexwertiger Standardfunktionen mit dem 21-stelligen Ergebnistyp `Lcomplex` wurden die entsprechenden Algorithmen im **rrGcomplex**-Format realisiert, das wie folgt definiert ist:

```
type rrGcomplex = record Re, Im : rrG end;
```

Es wurden nur die wichtigsten Operatoren, Funktionen und Prozeduren bereitgestellt, die zur Realisierung der Standardfunktionen notwendig waren.

#### 5.2.3.1 Operatoren

Neben den monadischen Operatoren `+` `-` existieren die in folgender Tabelle zusammengestellten, arithmetischen Grundoperatoren:

1. Operand		2. Operand	Ergebnistyp	Rel. Fehler
<code>rrGcomplex</code>	<code>+</code>	<code>rrGcomplex</code>	<code>rrGcomplex</code>	$5.0001 \cdot 10^{-26}$
<code>rrGcomplex</code>	<code>-</code>	<code>rrGcomplex</code>	<code>rrGcomplex</code>	$5.0001 \cdot 10^{-26}$
<code>rrGcomplex</code>	<code>*</code>	<code>rrGcomplex</code>	<code>rrGcomplex</code>	$1.1112 \cdot 10^{-25}$
<code>real</code>	<code>*</code>	<code>rrGcomplex</code>	<code>rrGcomplex</code>	$5.0001 \cdot 10^{-26}$
<code>Longreal</code>	<code>*</code>	<code>rrGcomplex</code>	<code>rrGcomplex</code>	$5.0001 \cdot 10^{-26}$
<code>rrG</code>	<code>*</code>	<code>rrGcomplex</code>	<code>rrGcomplex</code>	$5.0001 \cdot 10^{-26}$
<code>rrGcomplex</code>	<code>/</code>	<code>rrGcomplex</code>	<code>rrGcomplex</code>	$2.4617 \cdot 10^{-24}$
<code>rrGcomplex</code>	<code>/</code>	<code>real</code>	<code>rrGcomplex</code>	$2.2505 \cdot 10^{-24}$
<code>rrGcomplex</code>	<code>/</code>	<code>Longreal</code>	<code>rrGcomplex</code>	$2.2505 \cdot 10^{-24}$
<code>rrGcomplex</code>	<code>/</code>	<code>rrG</code>	<code>rrGcomplex</code>	$2.2505 \cdot 10^{-24}$

Als Vergleichsoperatoren existieren lediglich die Operatoren `=` `<>`.

#### 5.2.3.2 Transferfunktionen

Zur Wandlung zwischen den Typen `rrG` und `rrGcomplex` werden folgende Transferfunktionen bereitgestellt:

Funktion	Ergebnistyp	Bedeutung
<code>RRGCOMPL (x1,x2)</code>	<code>rrGcomplex</code>	$z = x1 + i \cdot x2; \quad i = \sqrt{-1}$
<code>RRGCOMPL (x)</code>	<code>rrGcomplex</code>	$z = x + i \cdot 0; \quad i = \sqrt{-1}$
<code>RE (z)</code>	<code>rrG</code>	$z.Re = \text{Realteil von } z$
<code>IM (z)</code>	<code>rrG</code>	$z.Im = \text{Imaginärteil von } z$

$x, x1, x2 = \text{rrG-Ausdruck}; \quad z = \text{rrGcomplex-Ausdruck};$

## 5.2.3.3 Standardfunktionen

Es stehen folgende Standardfunktionen mit Argumenttyp **rrGcomplex** und Ergebnistyp **rrGcomplex** zur Verfügung:

Aufruf	Ergebnistyp	Funktion	$\varepsilon(\operatorname{Re}\{f(z)\})$	$\varepsilon(\operatorname{Im}\{f(z)\})$
<code>sqr(z)</code>	rrGcomplex	$z^2 = z \cdot z$	$1.5001 \cdot 10^{-25}$	$1.5001 \cdot 10^{-25}$
<code>sqrt(z)</code>	rrGcomplex	$\sqrt{z}$	$4.4551 \cdot 10^{-24}$	$4.4551 \cdot 10^{-24}$
<code>sqrt_u(z)</code>	rrGcomplex	$\sqrt{z}$	$4.4551 \cdot 10^{-24}$	$4.4551 \cdot 10^{-24}$
<code>sqrt_negim(z)</code>	rrGcomplex	$\sqrt{z}$	$4.4551 \cdot 10^{-24}$	$4.4551 \cdot 10^{-24}$
<code>exp(z)</code>	rrGcomplex	$e^z$	$1.5318 \cdot 10^{-24}$	$1.5318 \cdot 10^{-24}$
<code>ln(z)</code>	rrGcomplex	$\ln(z)$	$3.3317 \cdot 10^{-24}$	$5.3282 \cdot 10^{-24}$
<code>ln1p(z)</code>	rrGcomplex	$\ln(1+z)$	$3.8802 \cdot 10^{-24}$	$5.8285 \cdot 10^{-24}$
<code>sin(z)</code>	rrGcomplex	$\sin(z)$	$2.7572 \cdot 10^{-24}$	$3.7017 \cdot 10^{-24}$
<code>cos(z)</code>	rrGcomplex	$\cos(z)$	$2.7572 \cdot 10^{-24}$	$3.7017 \cdot 10^{-24}$
<code>tan(z)</code>	rrGcomplex	$\tan(z)$	$9.7041 \cdot 10^{-24}$	$1.1985 \cdot 10^{-23}$
<code>cot(z)</code>	rrGcomplex	$\cot(z)$	$9.7041 \cdot 10^{-24}$	$1.1985 \cdot 10^{-23}$
<code>sinh(z)</code>	rrGcomplex	$\sinh(z)$	$3.7017 \cdot 10^{-24}$	$2.7572 \cdot 10^{-24}$
<code>cosh(z)</code>	rrGcomplex	$\cosh(z)$	$2.7572 \cdot 10^{-24}$	$3.7017 \cdot 10^{-24}$
<code>tanh(z)</code>	rrGcomplex	$\tanh(z)$	$1.1985 \cdot 10^{-23}$	$9.7041 \cdot 10^{-24}$
<code>coth(z)</code>	rrGcomplex	$\coth(z)$	$1.1985 \cdot 10^{-23}$	$9.7041 \cdot 10^{-24}$
<code>conj(z)</code>	rrGcomplex	$\bar{z} = x - i \cdot y$	0.0	0.0
<code>arg(z)</code>	rrG	$\varphi; z = r \cdot e^{i \cdot \varphi};$		$5.3282 \cdot 10^{-24}$
<code>arg1p(z)</code>	rrG	$\varphi; 1+z = r \cdot e^{i \cdot \varphi};$		$5.8285 \cdot 10^{-24}$
<code>abs(z)</code>	rrG	$r = \sqrt{x^2 + y^2};$		$1.2754 \cdot 10^{-24}$

Tabelle 5.9: Standardfunktionen des Moduls LC\_ARI,  $z$ : rrGcomplex

**Anmerkungen** zu den Funktionen: `sqrt(z)`, `sqrt_u(z)`, `sqrt_negim(z)`, `ln(z)`, `ln1p(z)`:

- **sqrt(z)** Der Verzweigungsschnitt ist die negative, reelle Achse, wobei diese Achse zur oberen Halbebene zählt.
- **sqrt\_u(z)** Der Verzweigungsschnitt ist die negative, reelle Achse, wobei diese Achse zur unteren Halbebene zählt.
- **sqrt\_negim(z)** Der Verzweigungsschnitt ist die negative, imaginäre Achse, wobei diese Achse zur linken Halbebene zählt.
- **ln(z)** Der Verzweigungsschnitt ist die negative, reelle Achse, wobei diese Achse zur oberen Halbebene zählt.

- **ln1p(z)** Der Verzweigungsschnitt ist der negative Teil der reellen Achse von  $-\infty$  bis  $-1$ , wobei diese Achse zur oberen Halbebene zählt.

Es stehen folgende Standardfunktionen mit Argumenttyp **Lcomplex** und Ergebnistyp **rrGcomplex** zur Verfügung:

Aufruf	Ergebnistyp	Funktion	$\varepsilon(\operatorname{Re}\{f(z)\})$	$\varepsilon(\operatorname{Im}\{f(z)\})$
sqr_rrG (z)	rrGcomplex	$z^2 = z \cdot z$	$1.5001 \cdot 10^{-25}$	$1.5001 \cdot 10^{-25}$
sqrt_rrG (z)	rrGcomplex	$\sqrt{z}$	$4.4551 \cdot 10^{-24}$	$4.4551 \cdot 10^{-24}$
sqrt_u_rrG (z)	rrGcomplex	$\sqrt{z}$	$4.4551 \cdot 10^{-24}$	$4.4551 \cdot 10^{-24}$
sqrt_negim_rrG(z)	rrGcomplex	$\sqrt{z}$	$4.4551 \cdot 10^{-24}$	$4.4551 \cdot 10^{-24}$
exp_rrG (z)	rrGcomplex	$e^z$	$1.5318 \cdot 10^{-24}$	$1.5318 \cdot 10^{-24}$
ln_rrG (z)	rrGcomplex	$\ln(z)$	$5.2040 \cdot 10^{-24}$	$5.3282 \cdot 10^{-24}$
ln1p_rrG (z)	rrGcomplex	$\ln(1+z)$	$2.9430 \cdot 10^{-24}$	$5.8285 \cdot 10^{-24}$
sin_rrG (z)	rrGcomplex	$\sin(z)$	$2.7572 \cdot 10^{-24}$	$3.7017 \cdot 10^{-24}$
cos_rrG (z)	rrGcomplex	$\cos(z)$	$2.7572 \cdot 10^{-24}$	$3.7017 \cdot 10^{-24}$
tan_rrG (z)	rrGcomplex	$\tan(z)$	$9.7041 \cdot 10^{-24}$	$1.1985 \cdot 10^{-23}$
cot_rrG (z)	rrGcomplex	$\cot(z)$	$9.7041 \cdot 10^{-24}$	$1.1985 \cdot 10^{-23}$
sinh_rrG (z)	rrGcomplex	$\sinh(z)$	$3.7017 \cdot 10^{-24}$	$2.7572 \cdot 10^{-24}$
cosh_rrG (z)	rrGcomplex	$\cosh(z)$	$2.7572 \cdot 10^{-24}$	$3.7017 \cdot 10^{-24}$
tanh_rrG (z)	rrGcomplex	$\tanh(z)$	$1.1985 \cdot 10^{-23}$	$9.7041 \cdot 10^{-24}$
coth_rrG (z)	rrGcomplex	$\coth(z)$	$1.1985 \cdot 10^{-23}$	$9.7041 \cdot 10^{-24}$
conj_rrG (z)	rrGcomplex	$\bar{z} = x - iy$	0.0	0.0
arg_rrG (z)	rrG	$\varphi; z = r \cdot e^{i\varphi};$	$5.3282 \cdot 10^{-24}$	
arg1p_rrG (z)	rrG	$\varphi; 1+z = r \cdot e^{i\varphi};$	$5.8285 \cdot 10^{-24}$	
abs_rrG (z)	rrG	$r = \sqrt{x^2 + y^2};$	$1.2504 \cdot 10^{-24}$	

Tabelle 5.10: Standardfunktionen des Moduls LC\_ARI,  $z$ : Lcomplex

**Anmerkungen** zu den Funktionen: sqrt\_rrG(z), sqrt\_u\_rrG(z), sqrt\_negim\_rrG, ln\_rrG(z), ln1p\_rrG(z):

- **sqrt\_rrG(z)** Der Verzweigungsschnitt ist die negative, reelle Achse, wobei diese Achse zur oberen Halbebene zählt.
- **sqrt\_u\_rrG(z)** Der Verzweigungsschnitt ist die negative, reelle Achse, wobei diese Achse zur unteren Halbebene zählt.
- **sqrt\_negim\_rrG(z)** Der Verzweigungsschnitt ist die negative, imaginäre Achse, wobei diese Achse zur linken Halbebene zählt.
- **ln\_rrG(z)** Der Verzweigungsschnitt ist die negative, reelle Achse, wobei diese Achse zur oberen Halbebene zählt.

- **ln1p\_rrG(z)** Der Verzweigungsschnitt ist der negative Teil der reellen Achse von  $-\infty$  bis  $-1$ , wobei diese Achse zur oberen Halbebene zählt.

#### 5.2.3.4 Überladungen des Zuweisungsoperators

Die Wandlung zwischen den Typen *Lcomplex* und *rrGcomplex* wird in Form einer überladenen Zuweisung bereitgestellt:

Zuweisung	Bedeutung
$Lc := rrGc$	Rundung von <i>rrGc</i> zur nächsten <i>Lcomplex</i> -Zahl
$rrGc := Lc$	$rrGc \equiv Lc$ , aber unterschiedliche Typen von <i>rrGc</i> , <i>Lc</i>

$rrGc = rrGcomplex\text{-Variable/Ausdruck}$ ;  $Lc = Lcomplex\text{-Variable/Ausdruck}$ ;

#### 5.2.3.5 Ein-/Ausgabeprozeduren

Es stehen die gewohnten Prozeduren

```
procedure read (var f: text; var z: rrGcomplex);
procedure write (var f: text; z: rrGcomplex);
```

mit optionalem Fileparameter, beliebig vielen Ein-/Ausgabeparametern, jedoch ohne Formatspezifikationen zur Verfügung.

Die Eingabe einer komplexen Zahl  $z = x + i \cdot y$  muß in der Form

$(x, y)$  oder in der Form  $x$

erfolgen. Im zweiten Fall wird der Imaginärteil  $y$  automatisch auf 0 gesetzt.  $x$  und  $y$  sind dabei *rrG*-Werte, deren Eingabe auf Seite 56 beschrieben wird.

Die Ausgabe einer *rrGcomplex*-Zahl erfolgt stets in der Form

$(x, y)$

wobei z.B. der Realteil  $x : rrG$  stets in der Form

$\{ x.r1, x.r2, x.G \}$

ausgegeben wird.

# Kapitel 6

## Komplexe Intervallarithmetik

Mit dem Modul `CI_ARI` werden die für das Rechnen mit komplexen Intervallen notwendigen Operatoren, Funktionen und Prozeduren bezüglich des Datentyps `REAL` bereitgestellt.

Das Modul `LCI_ARI` ermöglicht zusätzlich komplexe Intervallrechnungen auf Basis der Typen `REAL`, `LONGREAL`, `DOTPRECISION` und `RRG`.

### 6.1 Das Modul `CI_ARI`

Für Rechnungen mit komplexen Intervallen wird der Datentyp `cinterval` entsprechend der Definition

```
type cinterval = record re,im : interval end;
```

im Sprachkern von PASCAL-XSC bereitgestellt. Für  $z : cinterval$  sind damit Real- und Imaginärteil von  $z$  **Intervalle**, so daß

$$z = z.re + i \cdot z.im = [z.re.inf, z.re.sup] + i \cdot [z.im.inf, z.im.sup]; \quad i = \sqrt{-1}$$

in der komplexen Ebene als Rechteck-Intervall gedeutet werden kann:

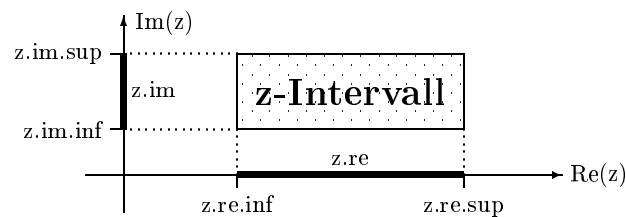


Abbildung 6.1: Komplexes Intervall

### 6.1.1 Operatoren

Sämtliche in diesem Modul vordefinierten arithmetischen und Verbands-Operatoren liefern den Ergebnistyp *cinterval*. Als arithmetische Operatoren stehen die monadischen Operatoren  $+$ ,  $-$  und die vier Grundoperationen  $+$ ,  $-$ ,  $*$ ,  $/$  mit der Rundung zum kleinsten, einschließenden komplexen Intervall vom Typ *cinterval* zur Verfügung.

Die Vergleichsoperatoren  $=$ ,  $<>$ ,  $<$ ,  $<=$ ,  $>$ ,  $>=$  sind mengentheoretisch zu interpretieren. Dabei bedeutet:

$=$  gleich                       $<>$  ungleich                       $<$  echte Teilmenge von  
 $<=$  Teilmenge von               $>$  echte Obermenge von               $>=$  Obermenge von

Für  $v, w$  vom Typ *cinterval* gilt:

$$v \leq w \iff (v.re \leq w.re) \text{ and } (v.im \leq w.im).$$

Die Operatoren auf der rechten Seite der Äquivalenzbeziehung sind dabei die Operatoren für Intervalle vom Typ *interval*.

rechter Operand	integer real	complex	interval	cinterval
<b>linker Operand</b>				
monadisch				$+, -$
integer real		$+$ *		$\diamond$ <b>in</b> , $=$ , $<>$ $+$ *
complex	$+$ *	$+$ *	$\diamond$ <b>in</b> , $=$ , $<>$ $+$ *	$\diamond$ <b>in</b> , $=$ , $<>$ $+$ *
interval		$\diamond$ $=$ , $<>$ $+$ *		$\diamond$ <b>in</b> , $\vee$ , $><$ $+$ *, $**$
cinterval	$\diamond$ $=$ , $<>$ $+$ *	$\diamond$ $=$ , $<>$ $+$ *	$\diamond$ $\vee$ , $><$ $+$ *, $**$	$\diamond$ <b>in</b> , $\vee$ , $><$ $+$ *, $**$

Tabelle 6.1: Die Operatoren des Moduls CLARI

$$\diamond \in \{+, -, *, /\} \quad \vee \in \{=, <>, <, <=, >, >=\}$$

Zusätzlich stehen die Operatoren **in** für die Relation „liegt in“ bzw. „enthalten im Innern“ sowie  $\gg$  für den Test auf Disjunktheit zweier komplexer Intervalle zur Verfügung. Dabei heißen zwei komplexe Intervalle  $v, w$  disjunkt, wenn gilt:  $v \cap w = \emptyset$  (leere Menge). Für zwei komplexe Intervalle  $x$  und  $y$  gilt:

$$x \text{ in } y \iff (x.re \text{ in } y.re) \text{ and } (x.im \text{ in } y.im).$$

Die Verbandsoperatoren  $+*$  bzw.  $**$  bezeichnen die Bildung der Intervall-Hülle bzw. des Durchschnitts, d.h. der Operator  $+*$  liefert das kleinste, beide Operanden umfassende komplexe Intervall, und der Operator  $**$  liefert das komplexe Schnittintervall. Ein leerer Schnitt führt zu einem Laufzeitfehler.

### 6.1.2 Transferfunktionen

Zur Wandlung zwischen den Typen *integer*, *real*, *complex*, *interval* und *cinterval* werden folgende Transferfunktionen bereitgestellt:

Funktion	Ergebnistyp	Bedeutung
compl (I1,I2)	<i>cinterval</i>	$I1 + i \cdot I2$
compl (R,I)	<i>cinterval</i>	$R + i \cdot I$
compl (I,R)	<i>cinterval</i>	$I + i \cdot R$
compl (I)	<i>cinterval</i>	$I + i \cdot 0$
intval(C1,C2)	<i>cinterval</i>	$[C1.re, C2.re] + i \cdot [C1.im, C2.im];$ $C1.re \leq C2.re$ <b>und</b> $C1.im \leq C2.im$
intval(R,C)	<i>cinterval</i>	$[R, C.re] + i \cdot [0, C.im];$ $R \leq C.re$ <b>und</b> $0 \leq C.im$
intval(C,R)	<i>cinterval</i>	$[C.re, R] + i \cdot [C.im, 0];$ $C.re \leq R$ <b>und</b> $C.im \leq 0$
intval(C)	<i>cinterval</i>	$[C.re, C.re] + i \cdot [C.im, C.im]$
re (CI)	<i>interval</i>	Realteil von CI
im (CI)	<i>interval</i>	Imaginärteil von CI
inf (CI)	<i>complex</i>	Komplexe Untergrenze $z$ von CI, mit: $z = ( CI.re.inf, CI.im.inf )$
sup (CI)	<i>complex</i>	Komplexe Obergrenze $z$ von CI, mit: $z = ( CI.re.sup, CI.im.sup )$

Tabelle 6.2: Transferfunktionen des Moduls CLARI

$R = \text{real-Ausdruck}$ ,  $I, I1, I2 = \text{interval-Ausdruck}$ ,  $i = \sqrt{-1}$ ,  
 $C, C1, C2 = \text{complex-Ausdruck}$ ,  $CI = \text{cinterval-Ausdruck}$

### 6.1.3 Überladungen des Zuweisungsoperators

Die Wandlungen von *integer/real/complex/interval* nach *cinterval* wird auch in Form überladener Zuweisungen bereitgestellt:

Zuweisung	Bedeutung
CI := R	CI := compl (intval(R))
CI := C	CI := intval (C)
CI := I	CI := compl (I)

Tabelle 6.3: Überladungen des Zuweisungsoperators

CI = *cinterval*-Variable,    I = *interval*-Ausdruck,  
 C = *complex*-Ausdruck,    R = *integer/real*-Ausdruck

### 6.1.4 Standardfunktionen

Es stehen die folgenden Standardfunktionen zur Verfügung:

Funktion	Ergb.-Typ	Bedeutung	
sqr (ci)	cinterval	$ci^2 = ci \cdot ci$	Quadrat
sqrt (ci)	cinterval	$\sqrt{ci}$	Quadratwurzel
exp (ci)	cinterval	$e^{ci}$	Exponentialfunktion
ln (ci)	cinterval	$\ln(ci)$	Natürlicher Logarithmus
ln1p (ci)	cinterval	$\ln(1 + ci)$	Natürlicher Logarithmus
power (w,z)	cinterval	$w^z$	Potenzfunktion, nicht immer hochgenau
xp1powy (w,z)	cinterval	$(1 + w)^z$	Potenzfunktion, nicht immer hochgenau
sin (ci)	cinterval	$\sin(ci)$	Sinus
cos (ci)	cinterval	$\cos(ci)$	Kosinus
sinh (ci)	cinterval	$\sinh(ci)$	Hyperbolischer Sinus
cosh (ci)	cinterval	$\cosh(ci)$	Hyperbolischer Kosinus
conj (ci)	cinterval	$\overline{ci} = x - i \cdot y$	Konjugation von $ci = x + i \cdot y$
arg (ci)	interval	$\varphi$	Argumentintervall für das <i>ci</i> -Rechteck



Funktion	Ergb.-Typ	Bedeutung	
arglp (ci)	interval	$\varphi$	Argumentintervall für das $(1 + ci)$ -Rechteck
abs (ci)	interval	$\sqrt{ci.re^2 + ci.im^2}$	Absolutbetrag von $ci$
mid (ci)	complex	m	Mittelpunkt von $ci$
diam (ci)	real	d	Durchmesser von $ci$
blow (ci,r)	cinterval		Epsilonaufblähung

Tabelle 6.4: Standardfunktionen des Moduls CLARI

$w, z, ci = cinterval$ -Ausdruck,  $r = real$ -Ausdruck,  $i = \sqrt{-1}$

#### Anmerkungen zu den Standardfunktionen:

1. Mit Ausnahme von  $diam(ci)$ ,  $blow(ci, r)$ ,  $power(w, z)$ ,  $xp1powy(w, z)$  liefern alle anderen Funktionen aus Tabelle 6.4 **hochgenaue** Ergebnisse, d.h. nur in ganz seltenen Ausnahmefällen wird z.B. die untere Intervallgrenze von  $e^{ci}$  um 1 ulp zu klein berechnet.
2.  $x := diam(CI)$  liefert zum komplexen Intervall  $CI: cinterval$  eine Obergrenze  $x: real$  für die Länge  $L$  der entsprechenden Rechteckdiagonalen mit:

$$0 \leq L \leq x$$

und einem relativen Höchstfehler:  $\varepsilon(DIAM) = 1.0000004 \cdot 10^{-12}$ ;

3. function **sqrt** (var  $CI: cinterval$ ):  $cinterval$ ;  
 $CI$  sei ein komplexes Intervallargument mit:

$$CI = [x_1, x_2] + i \cdot [y_1, y_2], \quad i = \sqrt{-1};$$

Zu vorgegebenem  $CI$  liefert die Anweisung  $EI := sqrt(CI)$  ein komplexes Einschließungsintervall  $EI$  mit:

$$\sqrt{z} \in EI \quad \text{für alle } z \in CI;$$

Bezüglich der Lage des Verzweigungsschnitts in der komplexen  $z$ -Ebene werden drei Fälle unterschieden:

- (a)  $y_1 < 0$ ,  $y_2 = 0$  und  $x_1 < 0$ ;  
Der Verzweigungsschnitt ist die negative, reelle Achse, wobei diese Achse der **unteren** Halbebene zugeordnet wird.
- (b)  $y_1 < 0 < y_2$  und  $x_1 < x_2 \leq 0$ ;  
Der Verzweigungsschnitt ist die negative, imaginäre Achse, wobei diese Achse der **linken** Halbebene zugeordnet wird.

- (c) In allen anderen Fällen ist der Verzweigungsschnitt wie üblich die negative, reelle Achse, wobei diese Achse der oberen Halbebene zugeordnet wird.

Durch die Wahl der Verzweigungsschnitte in den Fällen (a) und (b) wird eine unnötige Überschätzung der Einschließungsintervalle  $EI$  verhindert, die mit dem üblichen Verzweigungsschnitt im Fall (c) auftreten würde.

Allerdings geht in den Fällen (a) oder (b) die Inclusionsmonotonie verloren, denn es gilt z.B.:

$$A = \{[-2, -1] + i \cdot 0\} \subset \{[-2, -1] + i \cdot [-1, 0]\} = B, \quad \text{aber} \quad \text{sqrt}(A) \not\subset \text{sqrt}(B).$$

4. function **ARG** (var CI: cinterval): interval;  
CI sei ein komplexes Intervallargument mit:

$$CI = [x_1, x_2] + i \cdot [y_1, y_2], \quad i = \sqrt{-1};$$

Zu vorgegebenem  $CI$  liefert die Anweisung  $EI := \text{ARG}(CI)$  ein Einschließungsintervall  $EI$  mit:

$$\arg(z) \in EI \quad \text{für alle } z \in CI;$$

Gilt  $CI = 0$ , oder ist 0 **Innenpunkt** von  $CI$ , so erhält man das Intervall  $EI = [-\pi, +\pi]$ , wobei die Intervallgrenzen entsprechend gerundet werden.

In den restlichen Fällen müssen noch drei prinzipielle Lagen des komplexen Argumentintervalls unterschieden werden, in denen jeweils ein geeigneter Verzweigungsschnitt zu wählen ist:

- (a)  $y_1 < 0$ ,  $y_2 = 0$  und  $x_1 < 0$ ;  
Der Verzweigungsschnitt ist die negative, reelle Achse, wobei diese Achse der **unteren** Halbebene zugeordnet wird.
- (b)  $y_1 < 0 < y_2$  und  $x_1 < x_2 \leq 0$ ;  
Der Verzweigungsschnitt ist die negative, imaginäre Achse, wobei diese Achse der **linken** Halbebene zugeordnet wird.
- (c) In allen anderen Fällen ist der Verzweigungsschnitt wie üblich die negative, reelle Achse, wobei diese Achse der oberen Halbebene zugeordnet wird.

Durch die Wahl der Verzweigungsschnitte in den Fällen (a) und (b) wird eine unnötige Überschätzung der Einschließungsintervalle  $EI$  verhindert, die mit dem üblichen Verzweigungsschnitt im Fall (c) auftreten würde.

Allerdings geht in den Fällen (a) oder (b) die Inclusionsmonotonie verloren, denn es gilt z.B.:

$$A = \{[-2, -1] + i \cdot 0\} \subset \{[-2, -1] + i \cdot [-1, 0]\} = B, \quad \text{aber} \quad \text{ARG}(A) \not\subset \text{ARG}(B).$$

5. function **ARG1P** (var CI: cinterval): interval;  
 CI sei ein komplexes Intervallargument mit:

$$CI = [x_1, x_2] + i \cdot [y_1, y_2], \quad i = \sqrt{-1};$$

Zu vorgegebenem  $CI$  liefert die Anweisung  $EI := \text{ARG1P}(CI)$  ein Einschließungsintervall  $EI$  mit:

$$\arg(1 + z) \in EI \quad \text{für alle } z \in CI;$$

Gilt  $CI = -1$ , oder ist  $-1$  **Innenpunkt** von  $CI$ , so erhält man das Intervall  $EI = [-\pi, +\pi]$ , wobei die Intervallgrenzen entsprechend gerundet werden.

In den restlichen Fällen müssen noch drei prinzipielle Lagen des komplexen Argumentintervalls unterschieden werden, in denen jeweils ein geeigneter Verzweigungsschnitt zu wählen ist:

- (a)  $y_1 < 0$ ,  $y_2 = 0$  und  $x_1 < -1$ ;  
 Der Verzweigungsschnitt ist das Intervall  $(-\infty, -1]$  auf der negativen, reellen Achse, wobei dieses Intervall der **unteren** Halbebene zugeordnet wird.
- (b)  $y_1 < 0 < y_2$  und  $x_1 < x_2 \leq 0$ ;  
 Der Verzweigungsschnitt ist die Parallele zur negativen, imaginären Achse durch  $z = -1 + 0 \cdot i$ , wobei diese Parallele der **linken** Halbebene zugeordnet wird.
- (c) In allen anderen Fällen ist der Verzweigungsschnitt wieder das Intervall  $(-\infty, -1]$  auf der negativen, reellen Achse, wobei dieses Intervall der oberen Halbebene zugeordnet wird.

Durch die Wahl der Verzweigungsschnitte in den Fällen (a) und (b) wird eine unnötige Überschätzung der Einschließungsintervalle  $EI$  verhindert, die mit dem üblichen Verzweigungsschnitt im Fall (c) auftreten würde.

Allerdings geht in den Fällen (a) oder (b) die Inclusionsmonotonie verloren, denn es gilt z.B.:

$$A = \{-3, -2\} + i \cdot 0 \subset \{-3, -2\} + i \cdot [-1, 0] = B, \quad \text{aber}$$

$$\text{ARG1P}(A) \not\subset \text{ARG1P}(B).$$

6. function **LN** (var CI: cinterval): cinterval;  
 Zu vorgegebenem Argument  $CI$  liefert die Anweisung  $EI := \text{LN}(CI)$  ein Einschließungsintervall  $EI$  mit:

$$\ln(z) \in EI \quad \text{für alle } z \in CI;$$

Enthält das komplexe Argumentintervall  $CI$  den Ursprung, so erfolgt eine Fehlermeldung mit Programmabbruch. Wegen  $\text{Im}\{\text{LN}(CI)\} \equiv \text{ARG}(CI)$  gelten für die Funktionen  $\text{LN}(CI)$  und  $\text{ARG}(CI)$  bzgl. der Verzweigungsschnitte **identische** Aussagen!

7. function **LN1P** (var CI: cinterval): cinterval;  
 Zu vorgegebenem  $CI$  liefert die Anweisung  $EI := \text{LN1P}(CI)$  ein  
 Einschließungsintervall  $EI$  mit:

$$\ln(1 + z) \in EI \quad \text{für alle } z \in CI;$$

Enthält das komplexe Argumentintervall  $CI$  die Zahl  $-1$ , so erfolgt eine Fehlermeldung mit Programmabbruch.

Wegen  $\text{Im}\{\text{LN1P}(CI)\} \equiv \text{ARG1P}(CI)$  gelten für die Funktionen  $\text{LN1P}(CI)$  und  $\text{ARG1P}(CI)$  bezüglich der Verzweigungsschnitte **identische** Aussagen!

8. function **POWER** (var w,z: cinterval): cinterval;  
 Zu vorgegebenen komplexen Intervallargumenten  $w, z$  liefert die Anweisung  $EI := \text{POWER}(w, z)$  ein nicht immer hochgenau berechnetes  
 Einschließungsintervall  $EI$  mit:

$$w^z \in EI \quad \text{für alle komplexen Zahlen } u \in w \text{ und } v \in z;$$

Enthält das komplexe Argumentintervall  $w$  die Zahl  $0$ , so erfolgt eine Fehlermeldung mit Programmabbruch. Wegen

$$w^z \equiv e^{z \cdot \ln(w)}$$

gelten für die Funktionen  $\text{POWER}(w, z), \text{LN}(w), \text{ARG}(w)$  bezüglich der Verzweigungsschnitte **identische** Aussagen!

9. function **XP1POWY** (var w,z: cinterval): cinterval;  
 Zu vorgegebenen komplexen Intervallargumenten  $w, z$  liefert die Anweisung  $EI := \text{XP1POWY}(w, z)$  ein nicht immer hochgenau berechnetes  
 Einschließungsintervall  $EI$  mit:

$$(1 + w)^z \in EI \quad \text{für alle komplexen Zahlen } u \in w \text{ und } v \in z;$$

Enthält das komplexe Argumentintervall  $w$  die Zahl  $-1$ , so erfolgt eine Fehlermeldung mit Programmabbruch. Wegen

$$(1 + w)^z \equiv e^{z \cdot \ln(1+w)}$$

gelten für die Funktionen  $\text{XP1POWY}(w, z), \text{LN1P}(w), \text{ARG1P}(w)$  bezüglich der Verzweigungsschnitte **identische** Aussagen!

### 6.1.5 Ein-/Ausgabeweisungen

Es stehen die bekannten Prozeduren

```
procedure read (var f: text; var ci: cinterval);
procedure write (var f: text; ci: cinterval);
```

mit optionalem Fileparameter, beliebig vielen Ein-/Ausgabeparametern, jedoch ohne Formatspezifikationen zur Verfügung.

Die Eingabe eines komplexen Intervalls  $ci = [x, y] + i \cdot [v, w]$  kann mit dem Prozeduraufruf **read**(ci) in der Form

$([x, y], [v, w])$	allgemeines komplexes Intervall
oder in der Form	
$(x, [v, w])$	mit Punktintervall als Realteil, d.h. $x = y$
oder in der Form	
$([x, y], v)$	mit Punktintervall als Imaginärteil, d.h. $v = w$
oder in der Form	
$[x, y]$	rein reelles Intervall, d.h. $v = w = 0$
oder in der Form	
$(x, v)$	komplexes Punktintervall, d.h. $x = y$ und $v = w$
oder in der Form	
$x$	rein reelles Punktintervall, d.h. $x = y$ und $v = w = 0$

erfolgen.

Die Ausgabe eines komplexen Intervalls  $ci$  erfolgt durch den Prozeduraufruf **write**(ci) stets in der Form

$$([x, y], [v, w])$$

Das Ausgabeformat für die Real- und Imaginärteilintervalle ist dabei bestimmt durch das entsprechende Ausgabeformat für Intervalle, vergleiche Seite 80.

Die Ausgabe der Intervallgrenzen  $x, y$  und  $v, w$  kann natürlich auch erfolgen durch:

```
write(ci.re.inf,ci.re.sup)   und   write(ci.im.inf,ci.im.sup);
```

wobei entsprechende Formatspezifikationen für die *real*-Zahlen

$$x = ci.re.inf, \quad y = ci.re.sup, \quad v = ci.im.inf, \quad w = ci.im.sup$$

wieder möglich sind.

## 6.2 Das Modul LCI\_ARI

### 6.2.1 Der Typ Lcinterval

Für allgemeine Gleitkommarechnungen mit komplexen Intervallen wird neben dem Typ *cinterval* zusätzlich der Datentyp *Lcinterval* entsprechend der Definition

```
type Lcinterval = record re,im : Linterval end;
```

bereitgestellt. Für  $z : Lcinterval$  sind damit Real- und Imaginärteil von  $z$  **Intervalle**, so daß

$$z = z.re + i \cdot z.im = [z.re.inf, z.re.sup] + i \cdot [z.im.inf, z.im.sup]; \quad i = \sqrt{-1}$$

in der komplexen Ebene als Rechteck-Intervall gedeutet werden kann:

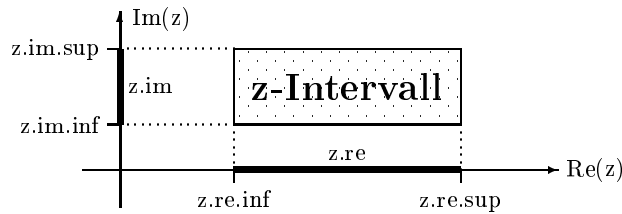


Abbildung 6.2: Komplexes Intervall

#### 6.2.1.1 Operatoren

Wir betrachten zunächst alle in diesem Modul vordefinierten arithmetischen und Verbands-Operatoren mit dem Ergebnistyp *Lcinterval*. Als arithmetische Operatoren stehen die monadischen Operatoren  $+$ ,  $-$  und die vier Grundoperationen  $+$ ,  $-$ ,  $*$ ,  $/$  mit der Rundung zum kleinsten, einschließenden komplexen Intervall vom Typ *Lcinterval* zur Verfügung.

Die Vergleichsoperatoren  $=$ ,  $<>$ ,  $<$ ,  $<=$ ,  $>$ ,  $>=$  sind mengentheoretisch zu interpretieren. Dabei bedeutet:

$=$ gleich	$<>$ ungleich	$<$ echte Teilmenge von
$<=$ Teilmenge von	$>$ echte Obermenge von	$>=$ Obermenge von

Für  $v, w$  vom Typ *Lcinterval* gilt:

$$v <= w \iff (v.re <= w.re) \text{ and } (v.im <= w.im).$$

Die Operatoren auf der rechten Seite der Äquivalenzbeziehung sind dabei die Operatoren für Intervalle vom Typ *Linterval*.

rechter Operand	integer real	Longreal	complex	Lcomplex
linker Operand	Ergebnistypen: <i>Lcinterval</i> oder <i>boolean</i>			
monadisch				
integer real				+*
Longreal			+*	+*
complex		+*		+*
Lcomplex	+*	+*	+*	+*
interval				◇ =, <> +*
Linterval			◇ =, <> +*	◇ =, <> +*
cinterval		◇ =, <> +*		◇ =, <> +*
Lcinterval	◇ =, <> +*	◇ =, <> +*	◇ =, <> +*	◇ =, <> +*

Tabelle 6.5: Die Operatoren des Moduls LCLARI (Teil 1)

$$\diamond \in \{+, -, *, /\} \quad \forall \in \{=, <>, <, <=, >, >=\}$$

Zusätzlich stehen die Operatoren **in** für die Relation „liegt in“ bzw. „enthalten im Innern“ sowie  $><$  für den Test auf Disjunktheit zweier komplexer Intervalle zur Verfügung. Dabei heißen zwei komplexe Intervalle  $v, w$  disjunkt, wenn gilt:  $v \cap w = \emptyset$  (leere Menge). Für zwei komplexe Intervalle  $x$  und  $y$  gilt:

$$x \text{ in } y \iff (x.re \text{ in } y.re) \text{ and } (x.im \text{ in } y.im).$$

Die Verbandsoperatoren  $+*$  bzw.  $**$  bezeichnen die Bildung der Intervall-Hülle bzw. des Durchschnitts, d.h. der Operator  $+*$  liefert das kleinste, beide Operanden umfassende komplexe Intervall, und der Operator  $**$  liefert das komplexe Schnittintervall. Ein leerer Schnitt führt zu einem Laufzeitfehler.

rechter Operand	interval	Linterval	cinterval	Lcinterval
linker Operand	Bezeichnungen wie in Tabelle 6.5			
monadisch				$+, -$
integer real				$\diamond$ $\mathbf{in}, =, \langle \rangle$ $+*$
Longreal			$\diamond$ $\mathbf{in}, =, \langle \rangle$ $+*$	$\diamond$ $\mathbf{in}, =, \langle \rangle$ $+*$
complex		$\diamond$ $\mathbf{in}, =, \langle \rangle$ $+*$		$\diamond$ $\mathbf{in}, =, \langle \rangle$ $+*$
Lcomplex	$\diamond$ $\mathbf{in}, =, \langle \rangle$ $+*$	$\diamond$ $\mathbf{in}, =, \langle \rangle$ $+*$	$\diamond$ $\mathbf{in}, =, \langle \rangle$ $+*$	$\diamond$ $\mathbf{in}, =, \langle \rangle$ $+*$
interval				$\diamond$ $\mathbf{in}, \vee, \rangle \langle$ $+*, **$
Linterval			$\diamond$ $\mathbf{in}, \vee, \rangle \langle$ $+*, **$	$\diamond$ $\mathbf{in}, \vee, \rangle \langle$ $+*, **$
cinterval		$\diamond$ $\vee, \rangle \langle$ $+*, **$		$\diamond$ $\mathbf{in}, \vee, \rangle \langle$ $+*, **$
Lcinterval	$\diamond$ $\vee, \rangle \langle$ $+*, **$	$\diamond$ $\vee, \rangle \langle$ $+*, **$	$\diamond$ $\mathbf{in}, \vee, \rangle \langle$ $+*, **$	$\diamond$ $\mathbf{in}, \vee, \rangle \langle$ $+*, **$

Tabelle 6.6: Die Operatoren des Moduls LCLARI (Teil 2)



## 6.2.1.2 Transferfunktionen

Zur Wandlung zwischen den Typen *integer*, *real*, *Longreal*, *complex*, *Lcomplex*, *interval*, *Linterval*, *cinterval* und *Lcinterval* werden folgende Transferfunktionen bereitgestellt:

Funktion	Ergebnistyp	Bedeutung
LCOMPL (I1,I2)	<i>Lcinterval</i>	$I1 + i \cdot I2$
LCOMPL (R,I)	<i>Lcinterval</i>	$R + i \cdot I$
LCOMPL (I,R)	<i>Lcinterval</i>	$I + i \cdot R$
LCOMPL (I)	<i>Lcinterval</i>	$I + i \cdot 0$
LCOMPL (I,LI)	<i>Lcinterval</i>	$I + i \cdot LI$
LCOMPL (LI,I)	<i>Lcinterval</i>	$LI + i \cdot I$
LCOMPL (LI1,LI2)	<i>Lcinterval</i>	$LI1 + i \cdot LI2$
LCOMPL (R,LI)	<i>Lcinterval</i>	$R + i \cdot LI$
LCOMPL (LI,R)	<i>Lcinterval</i>	$LI + i \cdot R$
LCOMPL (LR,I)	<i>Lcinterval</i>	$LR + i \cdot I$
LCOMPL (I,LR)	<i>Lcinterval</i>	$I + i \cdot LR$
LCOMPL (LR,LI)	<i>Lcinterval</i>	$LR + i \cdot LI$
LCOMPL (LI,LR)	<i>Lcinterval</i>	$LI + i \cdot LR$
LCOMPL (LI)	<i>Lcinterval</i>	$LI + i \cdot 0$
LINTVAL (C1,C2)	<i>Lcinterval</i>	$[C1.re, C2.re] + i \cdot [C1.im, C2.im];$ $C1.re \leq C2.re$ <b>und</b> $C1.im \leq C2.im$
LINTVAL (C,LC)	<i>Lcinterval</i>	$[C.re, LC.re] + i \cdot [C.im, LC.im];$ $C.re \leq LC.re$ <b>und</b> $C.im \leq LC.im$
LINTVAL (LC,C)	<i>Lcinterval</i>	$[LC.re, C.re] + i \cdot [LC.im, C.im];$ $LC.re \leq C.re$ <b>und</b> $LC.im \leq C.im$
LINTVAL (LC1,LC2)	<i>Lcinterval</i>	$[LC1.re, LC2.re] + i \cdot [LC1.im, LC2.im];$ $LC1.re \leq LC2.re$ <b>und</b> $LC1.im \leq LC2.im$
LINTVAL (R,C)	<i>Lcinterval</i>	$[R, C.re] + i \cdot [0, C.im];$ $R \leq C.re$ <b>und</b> $0 \leq C.im$

Tabelle 6.7: Transferfunktionen des Moduls LCLARI (Teil 1)

Funktion	Ergebnistyp	Bedeutung
LINTVAL (R,LC)	<i>Lcinterval</i>	$[R, LC.re] + i \cdot [0, LC.im];$ $R \leq LC.re$ <b>und</b> $0 \leq LC.im$
LINTVAL (LR,C)	<i>Lcinterval</i>	$[LR, C.re] + i \cdot [0, C.im];$ $LR \leq C.re$ <b>und</b> $0 \leq C.im$
LINTVAL (LR,LC)	<i>Lcinterval</i>	$[LR, LC.re] + i \cdot [0, LC.im];$ $LR \leq LC.re$ <b>und</b> $0 \leq LC.im$
LINTVAL (C,R)	<i>Lcinterval</i>	$[C.re, R] + i \cdot [C.im, 0];$ $C.re \leq R$ <b>und</b> $C.im \leq 0$
LINTVAL (LC,R)	<i>Lcinterval</i>	$[LC.re, R] + i \cdot [LC.im, 0];$ $LC.re \leq R$ <b>und</b> $LC.im \leq 0$
LINTVAL (C,LR)	<i>Lcinterval</i>	$[C.re, LR] + i \cdot [C.im, 0];$ $C.re \leq LR$ <b>und</b> $C.im \leq 0$
LINTVAL (LC,LR)	<i>Lcinterval</i>	$[LC.re, LR] + i \cdot [LC.im, 0];$ $LC.re \leq LR$ <b>und</b> $LC.im \leq 0$
LINTVAL (C)	<i>Lcinterval</i>	$[C.re, C.re] + i \cdot [C.im, C.im]$
LINTVAL (LC)	<i>Lcinterval</i>	$[LC.re, LC.re] + i \cdot [LC.im, LC.im]$
LINTVAL (CI)	<i>Lcinterval</i>	$LONG(CI.re) + i \cdot LONG(CI.im)$
LONG (CI)	<i>Lcinterval</i>	$LONG(CI.re) + i \cdot LONG(CI.im)$
SHORT (LCI)	<i>cinterval</i>	$SHORT(LCI.re) + i \cdot SHORT(LCI.im)$
RE (LCI)	<i>Linterval</i>	Realteil von <i>LCI</i>
IM (LCI)	<i>Linterval</i>	Imaginärteil von <i>LCI</i>
INF (LCI)	<i>Lcomplex</i>	Komplexe Untergrenze <i>z</i> von <i>LCI</i> , mit: $z = (LCI.re.inf, LCI.im.inf)$
SUP (LCI)	<i>Lcomplex</i>	Komplexe Obergrenze <i>z</i> von <i>LCI</i> , mit: $z = (LCI.re.sup, LCI.im.sup)$

Tabelle 6.8: Transferfunktionen des Moduls LCI\_ARI (Teil 2)

$LC, LC1, LC2 = Lcomplex$ -Ausdruck,  $LCI = Lcinterval$ -Ausdruck,  
 $C, C1, C2 = complex$ -Ausdruck,  $CI = cinterval$ -Ausdruck,  $i = \sqrt{-1}$ ,  
 $I, I1, I2 = interval$ -Ausdruck,  $R = real, Longreal$ -Ausdruck,  
 $LI, LI1, LI2 = Linterval$ -Ausdruck,  $LR = Longreal$ -Ausdruck

### 6.2.1.3 Überladungen des Zuweisungsoperators

Die Wandlungen der Typen *integer*, *real*, *Longreal*, *complex*, *Lcomplex*, *interval*, *Linterval* nach *Lcinterval* wird auch in Form überladener Zuweisungen bereitgestellt:

Zuweisung	Bedeutung
$LCI := R$	$LCI := LCOMPL (INTVAL(R))$
$LCI := LR$	$LCI := LCOMPL (INTVAL(LR))$
$LCI := CI$	$LCI := LINTVAL (CI)$
$LCI := C$	$LCI := LINTVAL (C)$
$LCI := LC$	$LCI := LINTVAL (LC)$
$LCI := I$	$LCI := LCOMPL (I)$
$LCI := LI$	$LCI := LCOMPL (LI)$

Tabelle 6.9: Überladungen des Zuweisungsoperators

$LCI$  = *Lcinterval*-Variable,       $R$  = *integer/real*-Ausdruck,  
 $LR$  = *Longreal*-Ausdruck,       $CI$  = *cinterval*-Ausdruck,  
 $C$  = *complex*-Ausdruck,       $LC$  = *Lcomplex*-Ausdruck,  
 $I$  = *interval*-Ausdruck,       $LI$  = *Linterval*-Ausdruck;

### 6.2.1.4 Standardfunktionen

Es stehen die folgenden Standardfunktionen zur Verfügung:

Funktion	Ergb.-Typ	Bedeutung
$\text{sqr} (LCI)$	Lcinterval	$LCI^2 = LCI \cdot LCI$ Quadrat
$\text{sqrt} (LCI)$	Lcinterval	$\sqrt{LCI}$ Quadratwurzel
$\text{exp} (LCI)$	Lcinterval	$e^{LCI}$ Exponentialfunktion
$\text{ln} (LCI)$	Lcinterval	$\ln(LCI)$ Natürlicher Logarithmus
$\text{ln1p} (LCI)$	Lcinterval	$\ln(1 + LCI)$ Natürlicher Logarithmus
$\text{power} (W,Z)$	Lcinterval	$W^Z$ Potenzfunktion, nicht immer hochgenau
$\text{xplpowy} (W,Z)$	Lcinterval	$(1 + W)^Z$ Potenzfunktion, nicht immer hochgenau

Funktion	Ergb.-Typ	Bedeutung	
sin (LCI)	Lcinterval	sin(LCI)	Sinus
cos (LCI)	Lcinterval	cos(LCI)	Kosinus
sinh (LCI)	Lcinterval	sinh(LCI)	Hyperbolischer Sinus
cosh (LCI)	Lcinterval	cosh(LCI)	Hyperbolischer Kosinus
conj (LCI)	Lcinterval	$\overline{LCI} = x - i \cdot y$	Konjugation von $LCI = x + i \cdot y$
arg (LCI)	Linterval	$\varphi$	Argumentintervall für das LCI-Rechteck
arg1p (LCI)	Linterval	$\varphi$	Argumentintervall für das (1 + LCI)-Rechteck
abs (LCI)	Linterval	$\sqrt{LCI.re^2 + LCI.im^2}$	Absolutbetrag von LCI
mid (LCI)	Lcomplex	m	Mittelpunkt von LCI
diam (LCI)	Longreal	d	Durchmesser von LCI
blow (LCI,R)	Lcinterval		Epsilonaufblähung

Tabelle 6.10: Standardfunktionen des Moduls LCI\_ARI

W,Z,LCI = *Lcinterval*-Ausdruck, R = *real*-Ausdruck,  $i = \sqrt{-1}$

#### Anmerkungen zu den Standardfunktionen:

1. Mit Ausnahme der Funktionen `power`, `xp1powy` sind alle Intervallfunktionen aus obiger Tabelle 6.10 **hochgenau**.
2. Für die Funktionen **sqrt**, **arg**, **arg1p**, **ln**, **ln1p**, **power**, **xp1powy** vom Ergebnistyp *Lcinterval* gelten bezüglich der Verzweigungsschnitte in Abhängigkeit von der Lage des komplexen Argumentintervalls die gleichen Aussagen wie für die entsprechenden *cinterval*-Funktionen von Seite 131.
3. `xl := diam (LCI)` liefert zum komplexen Intervall LCI: *Lcinterval* eine Obergrenze `xl`: *Longreal* für die Länge L der entsprechenden Rechteckdiagonalen mit:  $0 \leq L \leq xl$ ;
4. Die Anweisung `mc := MID(LCI)` liefert den maximalgenauen komplexen Mittelpunkt `mc`: *Lcomplex* des komplexen Intervalls LCI: *Lcinterval*.
5. Bei Einschließungsmethoden wird oft die Funktion **blow** für die sogenannte Epsilonaufblähung benötigt. Mit einem *real*-Ausdruck R wirkt die `blow(LCI,R)`-Funktion auf die Intervalle LCI.re und LCI.im wie die entsprechende `blow`-Funktion aus dem Modul LI\_ARI, vergl. Seite 84.

**6.2.1.5 Ein-/Ausgabeeweisungen**

Es stehen die bekannten Prozeduren

```
procedure read (var f: text; var LCI: Lcinterval);
procedure write (var f: text; LCI: Lcinterval);
```

mit optionalem Fileparameter, beliebig vielen Ein-/Ausgabeparametern, jedoch ohne Formatspezifikationen zur Verfügung.

Die Eingabe eines komplexen Intervalls  $LCI = [x, y] + i \cdot [v, w]$  kann mit dem Prozeduraufruf **read**(LCI) in der Form

$([x, y], [v, w])$	allgemeines komplexes Intervall
oder in der Form	
$(x, [v, w])$	mit Punktintervall als Realteil, d.h. $x = y$
oder in der Form	
$([x, y], v)$	mit Punktintervall als Imaginärteil, d.h. $v = w$
oder in der Form	
$[x, y]$	rein reelles Intervall, d.h. $v = w = 0$
oder in der Form	
$(x, v)$	komplexes Punktintervall, d.h. $x = y$ und $v = w$
oder in der Form	
$x$	rein reelles Punktintervall, d.h. $x = y$ und $v = w = 0$

erfolgen.

Die Ausgabe eines komplexen Intervalls  $LCI$  erfolgt durch den Prozeduraufruf **write**(LCI) stets in der Form

$$([x, y], [v, w])$$

Das Ausgabeformat für die Real- und Imaginärteilintervalle ist dabei bestimmt durch das entsprechende Ausgabeformat für Intervalle, vergleiche Seite 80.

Die Ausgabe der Intervallgrenzen  $x, y$  und  $v, w$  kann natürlich auch erfolgen durch:

```
write(LCI.re.inf,LCI.re.sup)   und   write(LCI.im.inf,LCI.im.sup);
```

wobei entsprechende Formatspezifikationen für die *Longreal*-Zahlen

$$x = LCI.re.inf, \quad y = LCI.re.sup, \quad v = LCI.im.inf, \quad w = LCI.im.sup$$

wieder möglich sind.

## 6.2.2 Exakte komplexe Intervallausdrücke

In **PASCAL-XSC** können mit Hilfe der Operatoren  $+$ ,  $-$ ,  $*$  komplexe Intervall-Ausdrücke vom Typ *cinterval* mit Hilfe des langen Akkumulators rundenfehlerfrei berechnet und durch die Anweisung

```
CI := ##( exakter, komplexer Intervallausdruck );
```

mit nur einer einzigen Rundung pro Intervallgrenze in eine Variable CI vom Typ *cinterval* abgespeichert werden. Dieses **##**-Konzept steht auch für die BCD-Version in **vollem Umfang** zur Verfügung; es ist jedoch **nicht** anwendbar, wenn nur einer der Operanden im exakt auszuwertenden, komplexen Intervallausdruck vom Typ *Longreal*, *Lcomplex*, *cdotprecision*, *Linterval*, *idotprecision*, *Lcinterval* oder *cidotprecision* ist!

Um das **##**-Konzept für diese Datentypen wenigstens simulieren zu können, benötigt man den Datentyp **cidotprecision**:

```
type cidotprecision = record re, im: idotprecision end;
```

### 6.2.2.1 Transferfunktionen

Zur Wandlung zwischen den Typen

*integer, real, Longreal, dotprecision, interval, Linterval, idotprecision, complex, Lcomplex, cdotprecision, cinterval, Lcinterval*

einerseits und dem Typ *cidotprecision* andererseits werden folgende Transferfunktionen bereitgestellt:

Funktion	Ergebnistyp	Bedeutung
DOTICOMPL(I1,I2)	<i>cidotprecision</i>	$I1 + i \cdot I2$ ;
DOTICOMPL(LI1,LI2)	<i>cidotprecision</i>	$LI1 + i \cdot LI2$ ;
DOTICOMPL(DI1,DI2)	<i>cidotprecision</i>	$DI1 + i \cdot DI2$ ;
DOTICOMPL(R,I)	<i>cidotprecision</i>	$[R, R] + i \cdot I$ ;
DOTICOMPL(L,LI)	<i>cidotprecision</i>	$[L, L] + i \cdot LI$ ;
DOTICOMPL(D,DI)	<i>cidotprecision</i>	$[D, D] + i \cdot DI$ ;
DOTICOMPL(I,R)	<i>cidotprecision</i>	$I + i \cdot [R, R]$ ;
DOTICOMPL(LI,L)	<i>cidotprecision</i>	$LI + i \cdot [L, L]$ ;
DOTICOMPL(DI,D)	<i>cidotprecision</i>	$DI + i \cdot [D, D]$ ;
DOTICOMPL(I)	<i>cidotprecision</i>	$I + i \cdot [0, 0]$ ;

Tabelle 6.11: Transferfunktionen mit dem Typ *cidotprecision* (Teil 1)

Funktion	Ergebnistyp	Bedeutung
DOTICOMPL(LI)	<i>cidotprecision</i>	$LI + i \cdot [0, 0];$
DOTICOMPL(DI)	<i>cidotprecision</i>	$DI + i \cdot [0, 0];$
DOTICOMPL(DC1,DC2)	<i>cidotprecision</i>	$[DC1.re, DC2.re]+$ $+i \cdot [DC1.im, DC2.im];$
DOTICOMPL(LC1,LC2)	<i>cidotprecision</i>	$[LC1.re, LC2.re]+$ $+i \cdot [LC1.im, LC2.im];$
DOTICOMPL(C1,C2)	<i>cidotprecision</i>	$[C1.re, C2.re]+$ $+i \cdot [C1.im, C2.im];$
DOTICOMPL(D,DC)	<i>cidotprecision</i>	$[D, DC.re] + i \cdot [0, DC.im];$
DOTICOMPL(L,LC)	<i>cidotprecision</i>	$[L, LC.re] + i \cdot [0, LC.im];$
DOTICOMPL(R,C)	<i>cidotprecision</i>	$[R, C.re] + i \cdot [0, C.im];$
DOTICOMPL(DC,D)	<i>cidotprecision</i>	$[DC.re, D] + i \cdot [DC.im, 0];$
DOTICOMPL(LC,L)	<i>cidotprecision</i>	$[LC.re, L] + i \cdot [LC.im, 0];$
DOTICOMPL(C,R)	<i>cidotprecision</i>	$[C.re, R] + i \cdot [C.im, 0];$
DOTICOMPL(DC)	<i>cidotprecision</i>	$[DC.re, DC.re]+$ $+i \cdot [DC.im, DC.im];$
DOTICOMPL(LC)	<i>cidotprecision</i>	$[LC.re, LC.re]+$ $+i \cdot [LC.im, LC.im];$
DOTICOMPL(C)	<i>cidotprecision</i>	$[C.re, C.re] + i \cdot [C.im, C.im];$
RE(DCI)	<i>idotprecision</i>	$[DCI.re.inf, DCI.re.sup];$
IM(DCI)	<i>idotprecision</i>	$[DCI.im.inf, DCI.im.sup];$
INF(DCI)	<i>cdotprecision</i>	$DCI.re.inf + i \cdot DCI.im.inf;$
SUP(DCI)	<i>cdotprecision</i>	$DCI.re.sup + i \cdot DCI.im.sup;$
LONG(DCI)	<i>Lcinterval</i>	$\text{Long}(DCI.re)+$ $+i \cdot \text{Long}(DCI.im)$
SHORT(DCI)	<i>cinterval</i>	$\text{SHORT}(DCI.re)+$ $+i \cdot \text{SHORT}(DCI.im)$

Tabelle 6.12: Transferfunktionen mit dem Typ *cidotprecision* (Teil 2)

$I, I1, I2 = \text{interval}$ -Ausdruck,  $LI, LI1, LI2 = \text{Linterval}$ -Ausdruck,  
 $DI, DI1, DI2 = \text{idotprecision}$ -Variable,  $R = \text{integer/real}$ -Ausdruck,  
 $L = \text{Longreal}$ -Ausdruck,  $D = \text{dotprecision}$ -Variable,  
 $C, C1, C2 = \text{complex}$ -Ausdruck,  $LC, LC1, LC2 = \text{Lcomplex}$ -Ausdruck,  
 $DC, DC1, DC2 = \text{cdotprecision}$ -Variable,  $DCI = \text{cidotprecision}$ -Variable.

### 6.2.2.2 Wertzuweisungen

Durch Überladen des Zuweisungsoperators `:=` sind folgende Wertzuweisungen an eine Variable vom Typ `cidotprecision` möglich:

```

cidotprecision-Variable := integer/real-Ausdruck;
cidotprecision-Variable := Longreal-Ausdruck;
cidotprecision-Variable := dotprecision-Ausdruck;
cidotprecision-Variable := complex-Ausdruck;
cidotprecision-Variable := Lcomplex-Ausdruck;
cidotprecision-Variable := cdotprecision-Ausdruck;
cidotprecision-Variable := interval-Ausdruck;
cidotprecision-Variable := Linterval-Ausdruck;
cidotprecision-Variable := idotprecision-Ausdruck;
cidotprecision-Variable := cinterval-Ausdruck;
cidotprecision-Variable := Lcinterval-Ausdruck;

```

### 6.2.2.3 Operatoren

Auch jetzt muß z.B. für einen Multiplikationsoperator mit `Lcinterval`-Operanden und einem Ergebnis vom Typ `cidotprecision` ein **neuer** Operatorname (`times_ci`) gewählt werden, da der normale `*` Operator mit gleichen Operanden schon für den Ergebnistyp `Lcinterval` vergeben ist.

Entsprechende Überlegungen bezüglich der Additions- und Subtraktionsoperatoren (z.B mit Namen `plus_ci`, `minus_ci`) führen zum Ergebnis, daß bei Berücksichtigung aller Operandenkombinationen neben den 81 `times_ci` Operatoren noch zusätzlich  $2 \cdot 13 \cdot 13 = 338$  Operatoren zu definieren sind! Die große Zahl von insgesamt 419 Operatoren läßt sich auf  $81 + 2 \cdot 13 = 107$  reduzieren, wenn man bei Addition und Subtraktion die folgende für die Praxis unerhebliche Einschränkung bezüglich der Operandenkombinationen macht. Dazu betrachten wir die Variable `S` vom Typ `cidotprecision` und die Variablen `a,b,c` mit den möglichen Typen:

```

cidotprecision, Lcinterval, cinterval, cdotprecision, Lcomplex, complex,
dotprecision, Longreal, real, integer, idotprecision, Linterval, interval;

```

Da für den Compiler die beiden folgenden Anweisungen äquivalent sind

$$S := a + b + c \iff S := (a + b) + c,$$

kann man mit der Variablen `NULL_DCI`: `cidotprecision` die Anweisung `S := a + b` auch in der Form schreiben:

$$\begin{aligned} \text{NULL\_DCI} &:= 0; \\ S &:= \text{NULL\_DCI} + a + b \iff S := ((\text{NULL\_DCI} + a) + b), \end{aligned}$$

und jetzt werden in der Anweisung für `S` **nur** Additionsoperatoren mit einem linken Operanden vom Typ `cidotprecision` benutzt! Für eine Subtraktion `S := a - b` gilt dann entsprechend:

$$\text{NULL\_DCI} := 0; \quad S := \text{NULL\_DCI} + a - b;$$



In den nachfolgenden Tabellen 6.13, 6.14 sind alle Operanden der Operatoren aus dem Modul LCLARI mit Ergebnistyp *cidotprecision* zusammengestellt:

rechter Operand	real Longreal	interval Linterval	complex Lcomplex	cinterval Lcinterval
<b>linker Operand</b>				
cidotprecision	* +, -	* +, -	* +, -	* +, -
real Longreal	times_ci	times_ci	times_ci	times_ci
interval Linterval	times_ci	times_ci	times_ci	times_ci
complex Lcomplex	times_ci	times_ci	times_ci	times_ci
cinterval Lcinterval	times_ci	times_ci	times_ci	times_ci

Tabelle 6.13: Operatoren aus LCLARI mit Ergebnistyp *cidotprecision* (Teil 1)

rechter Operand	dotprecision idotprecision	cdotprecision	cidotprecision
<b>linker Operand</b>			
cidotprecision	* +, -	* +, -	* +, -
real Longreal			*
interval Linterval			*
complex Lcomplex			*
cinterval Lcinterval			*
dotprecision idotprecision cdotprecision			*

Tabelle 6.14: Operatoren aus LCLARI mit Ergebnistyp *cidotprecision* (Teil 2)

Beachten Sie, daß jetzt neue Namen für Additions- und Subtraktionsoperatoren **nicht** notwendig sind, da die Einschränkung besteht, daß für diese Operatoren mit Ergebnistyp *cidotprecision* der linke Operand vom Typ *cidotprecision* sein muß. Durch diese Einschränkung ergibt sich noch der wichtige Satz:

In einem Ausdruck mit Ergebnistyp *cidotprecision* dürfen keine Klammern gesetzt werden.

Aus den Tabellen 6.13, 6.14 ergibt sich noch:

Ist bei einer Multiplikation mit Ergebnistyp *cidotprecision* wenigstens ein Operand vom Typ *cidotprecision*, so wird zur Vereinfachung das übliche Operatorsymbol \* benutzt und nicht der Name *times\_ci*.

**Beispiele:**

1. Sind *a, b* vom Typ *real* und *S* vom Typ *cidotprecision*, so ist die Anweisung

$$S := b + (a - b)$$

zwar zulässig, aber die Subtraktion in der Klammer und die anschließende Addition liefern i.a. fehlerbehaftete *real*-Ergebnisse, so daß in *S* i.a. **nicht** der exakte Differenzwert *a* gespeichert wird! Mit den Anweisungen

$$\text{NULL\_DCI} := 0; \quad S := \text{NULL\_DCI} + b + a - b$$

erhält man jedoch ohne Klammern den **exakten** Differenzwert *a*.

2. Sind *a, b, c* vom Typ *real* und *S* vom Typ *cidotprecision*, so ist die Anweisung

$$S := c * (a - b)$$

ebenfalls zulässig, liefert aber wie im 1. Beispiel i.a. **nicht** das exakte Produkt! Mit der Anweisung

$$S := c \text{ times\_ci } a - c \text{ times\_ci } b$$

erhält man jedoch ohne Klammern das **exakte** Produkt  $c \cdot (a - b)$ . Natürlich liefert die Anweisung

$$S := \#(c \cdot a - c \cdot b);$$

ebenfalls das exakte Produkt  $c \cdot (a - b)$  mit dem Ergebnistyp *cidotprecision*.

### 6.2.2.4 Simulation von SUM-Ausdrücken

Zur Simulation von **SUM**-Ausdrücken im **##**-Konzept stehen die folgenden Datentypen und Prozeduren zur Verfügung:

- **type** `civector` = **dynamic** array [\*] of **cinterval**;
  - **type** `Lcivector` = **dynamic** array [\*] of **Lcinterval**;
1. procedure **ADDNACCU** (var dci: cidotprecision; var a: Lcivector);  
**Exakte** Addition des *Lcinterval*-Feldes a zur *cidotprecision*-Variablen dci;
  2. procedure **ADDNACCU** (var dci: cidotprecision; var a: civector);  
**Exakte** Addition des *cinterval*-Feldes a zur *cidotprecision*-Variablen dci;
  3. procedure **ADDNACCU** (var dci: cidotprecision; var a: Livector);  
**Exakte** Addition des *Linterval*-Feldes a zur *cidotprecision*-Variablen dci;
  4. procedure **ADDNACCU** (var dci: cidotprecision; var a: ivector);  
**Exakte** Addition des *interval*-Feldes a zur *cidotprecision*-Variablen dci;
  5. procedure **ADDNACCU** (var dci: cidotprecision; var a: Lcvector);  
**Exakte** Addition des *Lcomplex*-Feldes a zur *cidotprecision*-Variablen dci;
  6. procedure **ADDNACCU** (var dci: cidotprecision; var a: cvector);  
**Exakte** Addition des *complex*-Feldes a zur *cidotprecision*-Variablen dci;
  7. procedure **ADDNACCU** (var dci: cidotprecision; var a: Lrvector);  
**Exakte** Addition des *Longreal*-Feldes a zur *cidotprecision*-Variablen dci;
  8. procedure **ADDNACCU** (var dci: cidotprecision; var a: rvector);  
**Exakte** Addition des *real*-Feldes a zur *cidotprecision*-Variablen dci;
  9. procedure **SUBNACCU** (var dci: cidotprecision; var a: Lcivector);  
**Exakte** Subtraktion des *Lcinterval*-Feldes a zur *cidotprecision*-Variablen dci;
  10. procedure **SUBNACCU** (var dci: cidotprecision; var a: civector);  
**Exakte** Subtraktion des *cinterval*-Feldes a zur *cidotprecision*-Variablen dci;
  11. procedure **SUBNACCU** (var dci: cidotprecision; var a: Livector);  
**Exakte** Subtraktion des *Linterval*-Feldes a zur *cidotprecision*-Variablen dci;
  12. procedure **SUBNACCU** (var dci: cidotprecision; var a: ivector);  
**Exakte** Subtraktion des *interval*-Feldes a zur *cidotprecision*-Variablen dci;
  13. procedure **SUBNACCU** (var dci: cidotprecision; var a: Lcvector);  
**Exakte** Subtraktion des *Lcomplex*-Feldes a zur *cidotprecision*-Variablen dci;
  14. procedure **SUBNACCU** (var dci: cidotprecision; var a: cvector);  
**Exakte** Subtraktion des *complex*-Feldes a zur *cidotprecision*-Variablen dci;

15. procedure **SUBNACCU** (var dci: cidotprecision; var a: Lrvector);  
**Exakte** Subtraktion des *Longreal*-Feldes a zur *cidotprecision*-Variablen dci;
16. procedure **SUBNACCU** (var dci: cidotprecision; var a: rvector);  
**Exakte** Subtraktion des *real*-Feldes a zur *cidotprecision*-Variablen dci;
17. procedure **PADDNACCU** (var dci: cidotprecision; var a,b: Lcivector);  
**Exakte** Addition des *Lcinterval*-Skalarproduktes zur Variablen dci vom Typ *cidotprecision*;
18. procedure **PADDNACCU** (var dci: cidotprecision; var a,b: civector);  
**Exakte** Addition des *cinterval*-Skalarproduktes zur Variablen dci vom Typ *cidotprecision*;
19. procedure **PADDNACCU** (var dci: cidotprecision; var a,b: Livector);  
**Exakte** Addition des *Linterval*-Skalarproduktes zur Variablen dci vom Typ *cidotprecision*;
20. procedure **PADDNACCU** (var dci: cidotprecision; var a,b: ivector);  
**Exakte** Addition des *interval*-Skalarproduktes zur Variablen dci vom Typ *cidotprecision*;
21. procedure **PADDNACCU** (var dci: cidotprecision; var a,b: Lvvector);  
**Exakte** Addition des *Lcomplex*-Skalarproduktes zur Variablen dci vom Typ *cidotprecision*;
22. procedure **PADDNACCU** (var dci: cidotprecision; var a,b: cvector);  
**Exakte** Addition des *complex*-Skalarproduktes zur Variablen dci vom Typ *cidotprecision*;
23. procedure **PADDNACCU** (var dci: cidotprecision; var a,b: Lrvector);  
**Exakte** Addition des *Longreal*-Skalarproduktes zur Variablen dci vom Typ *cidotprecision*;
24. procedure **PADDNACCU** (var dci: cidotprecision; var a,b: rvector);  
**Exakte** Addition des *real*-Skalarproduktes zur Variablen dci vom Typ *cidotprecision*;
25. procedure **PSUBNACCU** (var dci: cidotprecision; var a,b: Lcivector);  
**Exakte** Subtraktion des *Lcinterval*-Skalarproduktes zur Variablen dci vom Typ *cidotprecision*;
26. procedure **PSUBNACCU** (var dci: cidotprecision; var a,b: civector);  
**Exakte** Subtraktion des *cinterval*-Skalarproduktes zur Variablen dci vom Typ *cidotprecision*;
27. procedure **PSUBNACCU** (var dci: cidotprecision; var a,b: Livector);  
**Exakte** Subtraktion des *Linterval*-Skalarproduktes zur Variablen dci vom Typ *cidotprecision*;

28. procedure **PSUBNACCU** (var dci: cidotprecision; var a,b: ivector);  
**Exakte** Subtraktion des *interval*-Skalarproduktes zur Variablen dci vom Typ *cidotprecision*;
29. procedure **PSUBNACCU** (var dci: cidotprecision; var a,b: Lcvector);  
**Exakte** Subtraktion des *Lcomplex*-Skalarproduktes zur Variablen dci vom Typ *cidotprecision*;
30. procedure **PSUBNACCU** (var dci: cidotprecision; var a,b: cvector);  
**Exakte** Subtraktion des *complex*-Skalarproduktes zur Variablen dci vom Typ *cidotprecision*;
31. procedure **PSUBNACCU** (var dci: cidotprecision; var a,b: Lrvector);  
**Exakte** Subtraktion des *Longreal*-Skalarproduktes zur Variablen dci vom Typ *cidotprecision*;
32. procedure **PSUBNACCU** (var dci: cidotprecision; var a,b: rvector);  
**Exakte** Subtraktion des *real*-Skalarproduktes zur Variablen dci vom Typ *cidotprecision*;

### 6.2.2.5 Beispiele

#### 1. Beispiel:

Mit den Deklarationen:

```
var dotci,dci : cidotprecision;
    a,b       : Lcvector [1..10];
    Lci       : Lcinterval;
```

sind z.B. folgende Anweisungen möglich:

```
dci := 0;
PADDNACCU (dci,a,b);
dotci := dci + a[1] times_ci b[10] - a[2] times_ci b[9];
Lci := LONG (dotci);
```

Das komplexe Intervall `Lci : Lcinterval` liefert eine Einschließung der exakten, komplexen Intervallsumme `dotci` mit **nur einer Rundung** bezüglich der reellen bzw. imaginären Intervallunter- und -Obergrenze!

#### 2. Beispiel:

Mit Hilfe der Operatoren `*` und `+` zwischen Operanden vom Typ *cidotprecision* und *Lcinterval* kann ein Polynom vom Ergebnistyp *cidotprecision* mit Argument und Koeffizienten vom Typ *Lcinterval*

$$P(z) = \sum_{k=0}^n a_k \cdot z^k$$

mit der folgenden Funktion nach dem Horner Schema **rundungsfehlerfrei** berechnet werden:

```

function HORNER_CIDOT (var a: Lcivector;
                      z: Lcinterval); cidotprecision;

var k : integer;
    y : cidotprecision;
begin
    y := a[ub(a)];
    for k := ub(a)-1 downto lb(a) do
        y := y * z + a[k];
    HORNER_CIDOT := y
end;

```

Die komplexen Punktintervall-Koeffizienten

$$\begin{aligned}
 a_0 &= 120 + 90i; & a_1 &= -54 + 22i; & a_2 &= 6 + 7i; \\
 a_3 &= -5 - 3i; & a_4 &= 1 + 0i;
 \end{aligned}$$

realisieren ein komplexwertiges Polynom  $P(z)$  mit den Nullstellen:

$$z_{1,2} = \pm\sqrt{2} \cdot (1 + 2i), \quad z_3 = 5, \quad z_4 = 3i.$$

Mit den Deklarationen:

```

var a: Lcivector[0..4];
    z,y: Lcinterval;
    dci: cidotprecision;

```

und den Anweisungen:

```

z      := 5.00000000000000000001; { z ≈ z3 }
dci    := HORNER_CIDOT (a,z);
y      := LONG(dci);

```

erhält man mit der BCD-Version das Ergebnis:

$$y = P(z) \in ( [+1.3100 \dots 000 \cdot 10^{-18}, +1.3100 \dots 001 \cdot 10^{-18}], \\
 \quad \quad \quad [-1.3300 \dots 001 \cdot 10^{-18}, -1.3300 \dots 000 \cdot 10^{-18}] )$$

mit **nur einer Rundung** bzgl. Unter- und Obergrenze der Real- und Imaginärteilintervalle.

### 6.2.3 Komplexe Intervallrechnung mit dem Typ rrG

Mit dem Modul `LCLARI` werden nur diejenigen Operatoren, Funktionen und Prozeduren bereitgestellt, die zur Implementierung von komplexen Intervall-Standardfunktionen zum Basistyp `rrG` notwendig sind. Es ist daher nicht vorgesehen, die komplette Intervallrechnung, wie sie etwa zwischen den Typen `cinterval` und `Lcinterval` definiert wurde, ebenfalls zu realisieren, so daß z.B. Vergleichsoperatoren für Intervalle oder Intervall-Grundoperationen mit gemischten Operandentypen nur in dem Umfang bereitgestellt werden, wie sie tatsächlich erforderlich waren.

#### 6.2.3.1 Der Datentyp rrGcinterval

Der Datentyp `rrGcinterval` ist definiert durch:

```
type rrGcinterval = record re, im : rrGinterval end;
```

Für  $z : rrGcinterval$  sind Real- und Imaginärteil von  $z$  **Intervalle**, so daß

$$z = z.re + i \cdot z.im = [z.re.inf, z.re.sup] + i \cdot [z.im.inf, z.im.sup]; \quad i = \sqrt{-1}$$

in der komplexen Ebene als Rechteck-Intervall gedeutet werden kann:

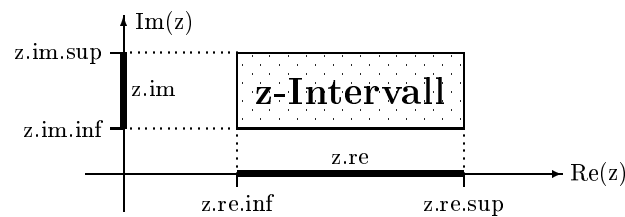


Abbildung 6.3: Komplexes Intervall

#### 6.2.3.2 Transferfunktionen

Zur Wandlung zwischen den Typen `rrGcinterval` und `rrGinterval` werden folgende Transferfunktionen bereitgestellt:

Funktion	Ergebnistyp	Bedeutung
RRGCINTVAL (x1,x2)	<code>rrGcinterval</code>	$z = x1 + i \cdot x2; \quad i = \sqrt{-1}$
RRGCINTVAL (x)	<code>rrGcinterval</code>	$z = x + i \cdot 0; \quad i = \sqrt{-1}$
RE (z)	<code>rrGinterval</code>	z.Re = Realteil von z
IM (z)	<code>rrGinterval</code>	z.Im = Imaginärteil von z

$x, x1, x2 = rrGinterval$ -Ausdruck;  $z = rrGcinterval$ -Ausdruck;

### 6.2.3.3 Operatoren

Neben den monadischen Operatoren  $+ -$  existieren die in folgender Tabelle zusammengestellten, arithmetischen Grundoperatoren:

1. Operand		2. Operand	Ergebnistyp
rrGcinterval	+	rrGcinterval	rrGcinterval
rrGcinterval	-	rrGcinterval	rrGcinterval
rrGcinterval	*	rrGcinterval	rrGcinterval
real	*	rrGcinterval	rrGcinterval
Longreal	*	rrGcinterval	rrGcinterval
rrG	*	rrGcinterval	rrGcinterval
rrGcinterval	/	rrGcinterval	rrGcinterval
rrGcinterval	/	real	rrGcinterval
rrGcinterval	/	Longreal	rrGcinterval
rrGcinterval	/	rrG	rrGcinterval

Als Vergleichsoperatoren existieren lediglich die Operatoren  $= <>$ .

### 6.2.3.4 Überladungen des Zuweisungsoperators

Die Wandlungen zwischen den verschiedenen Datentypen werden in Form überladener Zuweisungen nach folgender Tabelle bereitgestellt:

Zuweisung	Bedeutung
rrGci := R	rrGci = reelle Zahl
rrGci := C	rrGci = komplexe Zahl
rrGci := RI	rrGci = reelles Intervall
rrGci := CI	rrGci = komplexes Intervall
LC := rrGci	rrGci $\subseteq$ LC

Tabelle 6.15: Überladungen des Zuweisungsoperators

rrGci : rrGcinterval; R : rrG/Longreal/real; C : rrGcomplex/Lcomplex/complex;  
 RI : rrGinterval/Linterval/interval; CI : Lcinterval/cinterval; LC : Lcinterval;



## 6.2.3.5 Standardfunktionen

Es stehen folgende Standardfunktionen mit Argumenten  $x, y$ : *rrGinterval* zur Verfügung:

Funktion	Ergb.-Typ	Bedeutung	
<code>sqr (x)</code>	<i>rrGinterval</i>	$x^2 = x \cdot x$	Quadrat
<code>sqrt (x)</code>	<i>rrGinterval</i>	$\sqrt{x}$	Quadratwurzel
<code>exp (x)</code>	<i>rrGinterval</i>	$e^x$	Exponentialfunktion
<code>ln (x)</code>	<i>rrGinterval</i>	$\ln(x)$	Natürlicher Logarithmus
<code>ln1p (x)</code>	<i>rrGinterval</i>	$\ln(1 + x)$	Natürlicher Logarithmus
<code>power (x,y)</code>	<i>rrGinterval</i>	$x^y$	Potenzfunktion
<code>xp1powy (x,y)</code>	<i>rrGinterval</i>	$(1 + x)^y$	Potenzfunktion
<code>sin (x)</code>	<i>rrGinterval</i>	$\sin(x)$	Sinus
<code>cos (x)</code>	<i>rrGinterval</i>	$\cos(x)$	Kosinus
<code>sinh (x)</code>	<i>rrGinterval</i>	$\sinh(x)$	Hyperbolischer Sinus
<code>cosh (x)</code>	<i>rrGinterval</i>	$\cosh(x)$	Hyperbolischer Kosinus
<code>conj (x)</code>	<i>rrGinterval</i>	$\bar{x} = u - i \cdot v$	Konjugation von $x = u + i \cdot v$
<code>arg (x)</code>	<i>rrGinterval</i>	$\varphi$	Argumentintervall für das $x$ -Rechteck
<code>arg1p (x)</code>	<i>rrGinterval</i>	$\varphi$	Argumentintervall für das $(1 + x)$ -Rechteck
<code>abs (x)</code>	<i>rrGinterval</i>	$\sqrt{x.re^2 + x.im^2}$	Absolutbetrag von $x$
<code>diam (x)</code>	<i>rrG</i>	$d$	Durchmesser von $x$

Tabelle 6.16: Standardfunktionen des Moduls LCLARI

$$x, y = \textit{rrGinterval}\text{-Ausdruck}, \quad i = \sqrt{-1}$$

**Anmerkungen zu den Standardfunktionen:**

- Für die Funktionen `sqrt`, `arg`, `arg1p`, `ln`, `ln1p`, `power`, `xp1powy` vom Ergebnistyp *rrGinterval* gelten bezüglich der Verzweigungsschnitte in Abhängigkeit von der Lage des komplexen Argumentintervalls die gleichen Aussagen wie für die entsprechenden *cinterval*-Funktionen von Seite 131.
- $d := \text{diam}(x)$  liefert zum komplexen Intervall  $x$ : *rrGinterval* eine Obergrenze  $d$ : *rrG* für die Länge  $L$  der entsprechenden Rechteckdiagonalen mit:  $0 \leq L \leq d$ ;

Es stehen folgende Standardfunktionen mit Argumenten  $x, y: Lcinterval$  zur Verfügung:

Funktion	Ergb.-Typ	Bedeutung	
<code>sqr_rrG (x)</code>	<code>rrGcinterval</code>	$x^2 = x \cdot x$	Quadrat
<code>sqrt_rrG (x)</code>	<code>rrGcinterval</code>	$\sqrt{x}$	Quadratwurzel
<code>exp_rrG (x)</code>	<code>rrGcinterval</code>	$e^x$	Exponentialfunktion
<code>ln_rrG (x)</code>	<code>rrGcinterval</code>	$\ln(x)$	Natürlicher Logarithmus
<code>ln1p_rrG (x)</code>	<code>rrGcinterval</code>	$\ln(1+x)$	Natürlicher Logarithmus
<code>power_rrG (x,y)</code>	<code>rrGcinterval</code>	$x^y$	Potenzfunktion
<code>xp1powy_rrG (x,y)</code>	<code>rrGcinterval</code>	$(1+x)^y$	Potenzfunktion
<code>sin_rrG (x)</code>	<code>rrGcinterval</code>	$\sin(x)$	Sinus
<code>cos_rrG (x)</code>	<code>rrGcinterval</code>	$\cos(x)$	Kosinus
<code>sinh_rrG (x)</code>	<code>rrGcinterval</code>	$\sinh(x)$	Hyperbolischer Sinus
<code>cosh_rrG (x)</code>	<code>rrGcinterval</code>	$\cosh(x)$	Hyperbolischer Kosinus
<code>arg_rrG (x)</code>	<code>rrGinterval</code>	$\varphi$	Argumentintervall für das $x$ -Rechteck
<code>arg1p_rrG (x)</code>	<code>rrGinterval</code>	$\varphi$	Argumentintervall für das $(1+x)$ -Rechteck
<code>abs_rrG (x)</code>	<code>rrGinterval</code>	$\sqrt{x.re^2 + x.im^2}$	Absolutbetrag von $x$
<code>diam_rrG (x)</code>	<code>rrG</code>	$d$	Durchmesser von $x$

Tabelle 6.17: Standardfunktionen des Moduls LCLARI

$$x, y = Lcinterval\text{-Ausdruck}, \quad i = \sqrt{-1}$$

#### Anmerkungen zu den Standardfunktionen:

1. Für die Funktionen `sqrt_rrG`, `arg_rrG`, `arg1p_rrG`, `ln_rrG`, `ln1p_rrG`, `power_rrG`, `xp1powy_rrG` vom Ergebnistyp `rrGcinterval` gelten bezüglich der Verzweigungsschnitte in Abhängigkeit von der Lage des komplexen Argumentintervalls die gleichen Aussagen wie für die entsprechenden `cinterval`-Funktionen von Seite 131.
2.  $d := \text{diam\_rrG}(x)$  liefert zum komplexen Intervall  $x: Lcinterval$  eine Obergrenze  $d: rrG$  für die Länge  $L$  der entsprechenden Rechteckdiagonalen mit:  $0 \leq L \leq d$ ;

**6.2.3.6 Wertebereich separierbarer Funktionen**

Zur Berechnung einer Wertebereichs-Einschließung einer reellwertigen Funktion  $u(z)$  mit komplexen Punktargumenten  $z$  dient die **PASCAL-XSC** Funktion

```
function BOUNDS (
  mi, Ma : rrGcomplex; { Kompl. Punkte aus dem }
                    { Def.-Bereich von u(z) }
  function U(z: rrGcomplex): rrG; { reellwertig }
  eps : real { eps := 1.00001 *> eps(u) }
          { eps(u) := Fehlerschr. von u(z) }
) : rrGinterval;
```

$U(z)$  sei Näherungsfunktion für die reellwertige Funktion  $u(z)$ :

$$U(z) = u(z) \cdot (1 + \varepsilon_u), \quad |\varepsilon_u| < \varepsilon(u)$$

Der Eingabeparameter `eps` ist definiert durch:

$$\text{eps} := 1.00001 * > \varepsilon(u);$$

Mit der Anweisung

```
y := BOUNDS( mi, Ma, U, eps );    { y: rrGinterval }
```

erhält man unter der Voraussetzung  $u(\text{mi}) \leq u(\text{Ma})$  die folgenden Ergebnisse:

$$y.\text{Inf} \leq u(\text{mi}) \leq u(\text{Ma}) \leq y.\text{Sup}$$

**Anmerkungen:**

- Beim Aufruf von  $U(z)$  muß in der globalen Variablen `FW_exakt`, die vom Standardmodul **STDMOD** bereitgestellt wird, festgehalten werden, ob  $u(z)$  exakt berechnet wurde ( $U(z) = u(z)$ ) oder nicht.
- Im Fall  $u(\text{mi}) \leq u(z) \leq u(\text{Ma})$ ;  $z \in D$ : `rrGinterval` liefert `y: rrGinterval` eine Wertebereichseinschließung bzgl.  $u(z)$  für  $z \in D$ .  
Ist  $U(z)$  der Realteil einer gegebenen, komplexwertigen Näherungsfunktion

$$F(z) \equiv U(z) + i \cdot V(z) \approx f(z) \equiv u(z) + i \cdot v(z)$$

so liefert `y` eine Einschließung des Realteils  $u(z)$ . Ganz analog läßt sich dann auch der Wertebereich des Imaginärteils von  $f(z)$  einschließen.

Ist die komplexwertige Funktion  $f(z) := u(z) + i \cdot v(z)$  separierbar und sind die Minimum- und Maximumpunkte `mi, Ma` jeweils für  $u(z)$  und  $v(z)$  bekannt, so kann man durch Anwendung der Funktion **BOUNDS** auf die reellwertigen Näherungsfunktionen  $U(z)$ ,  $V(z)$  Intervall-Einschließungen für Real- und Imaginärteil von  $f(z)$  berechnen.

**6.2.3.7 Ein-/Ausgabeprozeduren**

Es stehen die bekannten Prozeduren

```
procedure read (var f: text; var RRGCI: rrGcinterval);
procedure write (var f: text; RRGCI: rrGcinterval);
```

mit optionalem Fileparameter, beliebig vielen Ein-/Ausgabeparametern, jedoch ohne Formatspezifikationen zur Verfügung.

Die Eingabe eines komplexen Intervalls  $\text{RRGCI} = [x, y] + i \cdot [v, w]$  kann mit dem Prozeduraufruf **read**(RRGCI) in der Form

$([x, y], [v, w])$	allgemeines komplexes Intervall
oder in der Form	
$(x, [v, w])$	mit Punktintervall als Realteil, d.h. $x = y$
oder in der Form	
$([x, y], v)$	mit Punktintervall als Imaginärteil, d.h. $v = w$
oder in der Form	
$[x, y]$	rein reelles Intervall, d.h. $v = w = 0$
oder in der Form	
$(x, v)$	komplexes Punktintervall, d.h. $x = y$ und $v = w$
oder in der Form	
$x$	rein reelles Punktintervall, d.h. $x = y$ und $v = w = 0$

erfolgen.

Die Vorschriften zur Eingabe einer *rrG-Zahl*  $x$  findet man dabei auf Seite 56.

Die Ausgabe eines komplexen Intervalls RRGCI erfolgt durch den Prozeduraufruf **write**(RRGCI) stets in der Form

$$([x, y], [v, w])$$

Die Ausgabe der Intervallgrenzen  $x, y$  und  $v, w$  kann natürlich auch erfolgen durch:

```
write(RRGCI.re.inf, RRGCI.re.sup) und write(RRGCI.im.inf, RRGCI.im.sup);
```

# Kapitel 7

## Matrix/-Vektorarithmetik

### 7.1 Übersicht

Das Rechnen mit Matrizen und Vektoren bildet die Grundlage zur Lösung vieler Probleme, z.B.:

lineare und nichtlineare Gleichungssysteme, lineare Optimierung, Polynomauswertungen, globale Optimierung, ...

Für Matrix- und Vektoroperanden stehen die Grundoperationen  $+$ ,  $-$ ,  $*$  mit den drei verschiedenen Rundungsarten zur Verfügung. Mit Hilfe des Akkumulators können die Operationen maximalgenau durchgeführt werden, d.h. das wahre Ergebnis z.B. eines Matrixelements aus einem Matrizenprodukt  $A*B$  wird zur nächstgelegenen Maschinenzahl gerundet. Bei gerichteten Rundungen wird das exakte Ergebnis zur nächstkleineren bzw. nächstgrößeren Rasterzahl gerundet. Erst durch solche maximalgenauen Rechnungen wird z.B. die maximalgenaue Einschließung eines Defekts bei linearen Gleichungssystemen oder die maximalgenaue Einschließung eines reellen Polynomwertes ermöglicht! Die genannten Grundoperationen beziehen sich dabei auf reelle oder komplexe Matrix/-Vektoroperanden in Punkt- oder Intervallform.

Im Gegensatz zu den Modulen

STDMOD, LI\_ARI, LC\_ARI oder LCI\_ARI

steht bei Matrizenrechnungen entweder nur der Datentyp *real* oder nur der Datentyp *Longreal* zur Verfügung. Beachten Sie, daß in der vorliegenden BCD-Version die in der Binärversion oft lästigen Konversionsfehler bei Matrix- und Vektorkomponenten **nicht** auftreten!

Die nachfolgende Tabelle liefert eine Zusammenstellung der vorhandenen Module, in denen die für das Rechnen mit Vektoren und Matrizen notwendigen Operatoren, Funktionen und Prozeduren bereitgestellt werden.

Name	Inhalt	Basistyp
MV_ARI	Reelle Matrizen und Vektoren	<i>real</i>
MVL_ARI	Reelle Intervall-Matrizen und -Vektoren	<i>real</i>
MVC_ARI	Komplexe Matrizen und Vektoren	<i>real</i>
MVCL_ARI	Komplexe Intervall-Matrizen und -Vektoren	<i>real</i>
LMV_ARI	Reelle Matrizen und Vektoren	<i>Longreal</i>
LMVL_ARI	Reelle Intervall-Matrizen und -Vektoren	<i>Longreal</i>
LMVC_ARI	Komplexe Matrizen und Vektoren	<i>Longreal</i>
LMVCL_ARI	Komplexe Intervall-Matrizen und -Vektoren	<i>Longreal</i>

Tabelle 7.1: Module zur Matrix/-Vektorarithmetik

Die Namen der Module zum Basistyp **Longreal** ergeben sich also aus den Namen der Module zum Basistyp **real** durch Voranstellen des Buchstaben 'L', wie es auch bei den entsprechenden Typ-Namen gehandhabt wurde.

## 7.2 Das Modul MV\_ARI

In diesem Modul werden die für das Rechnen mit reellen Vektoren und Matrizen notwendigen Operatoren, Funktionen und Prozeduren auf der Grundlage des Datentyps `real` bereitgestellt.

### 7.2.1 Datentypen

Die dynamischen Datentypen zur Darstellung von Vektoren und Matrizen sind entsprechend der Vereinbarung

```
type rvector = dynamic array [*] of real;
      rmatrix = dynamic array [*] of rvector;
```

im Sprachkern von PASCAL-XSC enthalten. Die aktuellen Indexgrenzen werden erst bei der Vereinbarung von Variablen dieser Typen festgelegt.

### 7.2.2 Operatoren

Viele der bekannten, grundlegenden Matrix/Vektor-Operationen sind in diesem Modul vordefiniert. Als arithmetische Operatoren stehen die monadischen Operatoren `+`, `-` und die vier Grundoperationen `+`, `-`, `*`, `/` mit den drei verschiedenen Rundungsarten zur Verfügung, wobei jedoch nur bestimmte Operandenkombinationen zulässig sind. Während die Operationen `+` und `-` für Vektoren und Matrizen komponentenweise erklärt sind gemäß

$$\begin{aligned} c &:= a \pm b & \text{mit} & & c[i] &:= a[i] \pm b[i] \\ C &:= A \pm B & \text{mit} & & C[i, j] &:= A[i, j] \pm B[i, j] \end{aligned}$$

mit `a`, `b`, `c` vom Typ `rvector` und `A`, `B`, `C` vom Typ `rmatrix`, sind die Operationen `*` und `/` definiert durch

$$\begin{aligned} s &:= a * b & \text{mit} & & s &:= \#* (\text{for } i := \text{lb}(a) \text{ to } \text{ub}(a) \\ & & & & & \quad \text{sum } (a[i] * b[i])) \quad \ddagger \\ c &:= r * a & \text{mit} & & c[i] &:= r * a[i] \\ c &:= a * r & \text{mit} & & c[i] &:= a[i] * r \\ c &:= a / r & \text{mit} & & c[i] &:= a[i] / r \\ c &:= A * b & \text{mit} & & c[i] &:= A[i] * b \quad \ddagger \\ C &:= r * A & \text{mit} & & C[i, j] &:= r * A[i, j] \\ C &:= A * r & \text{mit} & & C[i, j] &:= A[i, j] * r \\ C &:= A / r & \text{mit} & & C[i, j] &:= A[i, j] / r \\ C &:= A * B & \text{mit} & & C[i, j] &:= A[i] * B[*, j] \quad \ddagger \end{aligned}$$


---

<sup>‡</sup>Maximalgenaues Skalarprodukt mit Ergebnistyp `real`

wobei  $r, s$  vom Typ *real*,  $a, b, c$  vom Typ *rvector* und  $A, B, C$  vom Typ *rmatrix* sind. Die Operationen mit gerichteten Rundungen sind entsprechend definiert.

Die Vergleichsoperatoren  $=, <>, <, <=, >, >=$  wirken auf alle Komponenten eines Vektors oder einer Matrix. Für  $a, b$  vom Typ *rvector* und  $A, B$  vom Typ *rmatrix* gilt:

$$\begin{aligned} a <= b &\iff a[i] <= b[i] \quad \text{für alle } i \\ A <= B &\iff A[i, j] <= B[i, j] \quad \text{für alle } i, j \end{aligned}$$

Die Operatoren auf der rechten Seite der Äquivalenzbeziehungen sind dabei die Operatoren für Operanden vom Typ *real*.

rechter Operand	integer real	rvector	rmatrix
linker Operand			
monadisch		+, -	+, -
integer real		*, *<, *>	*, *<, *>
rvector	*, *<, *> /, /<, />	◦ ∨	
rmatrix	*, *<, *> /, /<, />	*, *<, *>	◦ ∨

Tabelle 7.2: Operatoren des Moduls MV\_ARI

$$\circ \in \{+, +<, +>, -, -<, ->, *, *<, *>\} \quad \vee \in \{=, <>, <, <=, >, >=\}$$

### 7.2.3 Überladungen des Zuweisungsoperators

Die komponentenweise Initialisierung von *rvector*- und *rmatrix*-Variablen wird in Form überladener Zuweisungen bereitgestellt:

Zuweisung	Bedeutung
$rv := r$	$rv[j] := r; \quad j = lb(rv), \dots, ub(rv)$
$rM := r$	$rM[j, k] := r; \quad j = lb(rM, 1), \dots, ub(rM, 1)$ $k = lb(rM, 2), \dots, ub(rM, 2)$

Tabelle 7.3:  $r = \text{real}$ -Ausdruck,  $rv = \text{rvector-Variable}$ ,  $rM = \text{rmatrix-Variable}$



### 7.2.4 Standardfunktionen

Es stehen die Funktionen *id* und *null* zur Erzeugung einer Einheitsmatrix und einer Nullmatrix bzw. eines Nullvektors sowie die Funktion *transp* zur Berechnung der transponierten Matrix zur Verfügung.

Funktion	Ergebnistyp	Bedeutung
null (v)	<i>rvector</i>	Nullvektor mit dem aktuellen Indexbereich von $v$
vnull (n)	<i>rvector</i>	Nullvektor mit dem Indexbereich [1..n]
null (M)	<i>rmatrix</i>	Nullmatrix mit den aktuellen Indexbereichen von $M$
null (M1,M2)	<i>rmatrix</i>	Nullmatrix mit den aktuellen Indexbereichen der Produktmatrix $M1 \cdot M2$
null (n)	<i>rmatrix</i>	Nullmatrix mit Indexbereich [1..n, 1..n]
null (n1,n2)	<i>rmatrix</i>	Nullmatrix mit Indexbereich [1..n1, 1..n2]
id (M)	<i>rmatrix</i>	Einheitsmatrix mit den aktuellen Indexbereichen von $M$
id (M1,M2)	<i>rmatrix</i>	Einheitsmatrix mit den aktuellen Indexbereichen der Produktmatrix $M1 \cdot M2$
id (n)	<i>rmatrix</i>	Einheitsmatrix mit den Indexbereichen [1..n, 1..n]
id (n1,n2)	<i>rmatrix</i>	Einheitsmatrix mit den Indexbereichen [1..n1, 1..n2]
transp (M)	<i>rmatrix</i>	Transponierte Matrix $Mt$ von $M$ mit $Mt[i, j] = M[j, i]$

Tabelle 7.4:  $n, n1, n2 = \text{integer}$ -Ausdruck,  $v = \text{rvector}$ -Ausdruck,  $M, M1, M2 = \text{rmatrix}$ -Ausdruck

#### Beispiel:

Ist  $E$  die Einheitsmatrix und  $R \approx A^{-1}$  eine Näherungsinverse der quadratischen Matrix  $A$ , so kann der in der Numerik häufig verwendete Defekt

$$D = E - R \cdot A$$

in PASCAL-XSC durch die Anweisung

$$D := \text{id}(A) - R * A$$

bzw. unter Verwendung eines Lattenkreuzausdrucks durch die Anweisung

$$D := \# * (\text{id}(A) - R * A)$$

bestimmt werden, wobei die zweite Form die Defektmatrix mit nur einer einzigen Rundung pro Komponente berechnet.

### 7.2.5 Ein-/AusgabeprozEDUREN

Es stehen die Prozeduren

```
procedure read (var f : text; var a : rvector);
procedure read (var f : text; var A : rmatrix);
procedure write (var f : text; a : rvector);
procedure write (var f : text; A : rmatrix);
```

mit optionalem Fileparameter, beliebig vielen Ein-/Ausgabeparametern, jedoch ohne Formatspezifikationen zur Verfügung.

Die Eingabe eines Vektors bzw. einer Matrix erfolgt komponentenweise entsprechend der Eingabe von *real*-Werten, dabei wird eine Matrix stets zeilenweise eingelesen. Die Ausgabe eines Vektors bzw. einer Matrix erfolgt wie bei der Eingabe komponentenweise mit dem Standardformat für *real*-Größen aus Tabelle 3.4.

#### Beispiel:

Durch die Anweisung

```
read (b, A, x); { b,x : rvector; A : rmatrix; }
```

können der Vektor *b*, die Matrix *A* und der Vektor *x* hintereinander eingelesen werden.

## 7.3 Das Modul LMV\_ARI

In diesem Modul werden die für das Rechnen mit reellen Vektoren und Matrizen notwendigen Operatoren, Funktionen und Prozeduren auf der Grundlage des Datentyps **Longreal** bereitgestellt.

### 7.3.1 Datentypen

Die dynamischen Datentypen zur Darstellung von Vektoren und Matrizen sind entsprechend der Vereinbarung

```
type Lrvector = dynamic array [*] of Longreal;
   Lrmatrix  = dynamic array [*] of Lrvector;
```

definiert. Die aktuellen Indexgrenzen werden erst bei der Vereinbarung von Variablen dieser Typen festgelegt.

### 7.3.2 Operatoren

Viele der bekannten, grundlegenden Matrix/Vektor-Operationen sind in diesem Modul vordefiniert. Als arithmetische Operatoren stehen die monadischen Operatoren  $+$ ,  $-$  und die vier Grundoperationen  $+$ ,  $-$ ,  $*$ ,  $/$  mit den drei verschiedenen Rundungsarten zur Verfügung, wobei jedoch nur bestimmte Operandenkombinationen zulässig sind. Während die Operationen  $+$  und  $-$  für Vektoren und Matrizen komponentenweise erklärt sind gemäß

$$\begin{aligned} c &:= a \pm b & \text{mit} & & c[i] &:= a[i] \pm b[i] \\ C &:= A \pm B & \text{mit} & & C[i, j] &:= A[i, j] \pm B[i, j] \end{aligned}$$

mit  $a, b, c$  vom Typ *Lrvector* und  $A, B, C$  vom Typ *Lrmatrix*, sind die Operationen  $*$ , *times* und  $/$  definiert durch

$$\begin{aligned} s &:= a * b & \text{mit} & & s &:= \text{Longnext} \left( \sum_i a[i] \cdot b[i] \right) \ddagger \\ d &:= a \text{ times } g & \text{mit} & & d &:= \sum_i a[i] \cdot g[i] \dagger \\ c &:= r * a & \text{mit} & & c[i] &:= r * a[i] \\ c &:= a * r & \text{mit} & & c[i] &:= a[i] * r \\ c &:= a / r & \text{mit} & & c[i] &:= a[i] / r \\ c &:= A * b & \text{mit} & & c[i] &:= A[i] * b \ddagger \\ C &:= r * A & \text{mit} & & C[i, j] &:= r * A[i, j] \\ C &:= A * r & \text{mit} & & C[i, j] &:= A[i, j] * r \\ C &:= A / r & \text{mit} & & C[i, j] &:= A[i, j] / r \\ C &:= A * B & \text{mit} & & C[i, j] &:= A[i] * B[* , j] \ddagger \end{aligned}$$


---

$\ddagger$ Maximalgenaues Skalarprodukt mit Ergebnistyp *Longreal*

$\dagger$ Exaktes Skalarprodukt mit Ergebnistyp *dotprecision*

wobei  $r, s$  vom Typ *Longreal*,  $a, b, c$  vom Typ *Lrvector*,  $g$  vom Typ *rvector/Lrvector*,  $d$  vom Typ *dotprecision* und  $A, B, C$  vom Typ *Lrmatrix* sind. Die Operationen mit gerichteten Rundungen sind ganz entsprechend definiert.

Die Vergleichsoperatoren  $=, <>, <, <=, >, >=$  wirken auf alle Komponenten eines Vektors oder einer Matrix. Für  $a, b$  vom Typ *Lrvector* und  $A, B$  vom Typ *Lrmatrix* gilt:

$$\begin{aligned} a <= b &\iff a[i] <= b[i] \quad \text{für alle } i \\ A <= B &\iff A[i, j] <= B[i, j] \quad \text{für alle } i, j \end{aligned}$$

Die Operatoren auf der rechten Seite der Äquivalenzbeziehungen sind dabei die Operatoren für Operanden vom Typ *Longreal*.

rechter Operand	Longreal	rvector	Lrvector	Lrmatrix
linker Operand				
monadisch			+, -	+, -
Longreal			*, *<, *>	*, *<, *>
rvector		times	times	
Lrvector	*, *<, *> /, /<, />	times	o, times ∨	
Lrmatrix	*, *<, *> /, /<, />		*, *<, *>	o ∨

Tabelle 7.5: Operatoren des Moduls LMV\_ARI

$$o \in \{+, +<, +>, -, -<, ->, *, *<, *>\} \quad \vee \in \{=, <>, <, <=, >, >=\}$$

### 7.3.3 Überladungen des Zuweisungsoperators

Zuweisung	Bedeutung
Lv := Lr	Lv[j] := Lr; j = lb(Lv), ... ub(Lv)
LM := Lr	LM[j, k] := Lr; j = lb(LM,1), ... ub(LM,1) k = lb(LM,2), ... ub(LM,2)
LM := M	LM[j, k] := M[j, k]; j = lb(LM,1), ... ub(LM,1) k = lb(LM,2), ... ub(LM,2)
Lv := v	Lv[j] := v[j]; j = lb(Lv), ... ub(Lv)

Lr = *Longreal*-Ausdruck, Lv = *Lrvector-Variable*,  
LM = *Lrmatrix-Variable*, M = *matrix*-Ausdruck, v = *rvector*-Ausdruck

### 7.3.4 Standardfunktionen

Es stehen die Funktionen *lid* und *lnull* zur Erzeugung einer Einheitsmatrix und einer Nullmatrix bzw. eines Nullvektors sowie die Funktion *ltransp* zur Berechnung der transponierten Matrix zur Verfügung. Die *short*-Funktionen erzeugen *rmatrix*- bzw. *rvector*-Werte.

Funktion	Ergebnistyp	Bedeutung
<i>lnull</i> (lv)	<i>Lrvector</i>	Nullvektor mit dem aktuellen Indexbereich von lv
<i>lvnull</i> (n)	<i>Lrvector</i>	Nullvektor mit dem Indexbereich [1..n]
<i>lnull</i> (IM)	<i>Lrmatrix</i>	Nullmatrix mit den aktuellen Indexbereichen von IM
<i>lnull</i> (IM1,IM2)	<i>Lrmatrix</i>	Nullmatrix mit den aktuellen Indexbereichen der Produktmatrix $IM1 \cdot IM2$
<i>lnull</i> (n)	<i>Lrmatrix</i>	Nullmatrix mit Indexbereich [1..n, 1..n]
<i>lnull</i> (n1,n2)	<i>Lrmatrix</i>	Nullmatrix mit Indexbereich [1..n1, 1..n2]
<i>lid</i> (IM)	<i>Lrmatrix</i>	Einheitsmatrix mit den aktuellen Indexbereichen von IM
<i>lid</i> (IM1,IM2)	<i>Lrmatrix</i>	Einheitsmatrix mit den aktuellen Indexbereichen der Produktmatrix $IM1 \cdot IM2$
<i>lid</i> (n)	<i>Lrmatrix</i>	Einheitsmatrix mit den Indexbereichen [1..n, 1..n]
<i>lid</i> (n1,n2)	<i>Lrmatrix</i>	Einheitsmatrix mit den Indexbereichen [1..n1, 1..n2]
<i>ltransp</i> (IM)	<i>Lrmatrix</i>	Transponierte Matrix $IM^t$ von IM mit $IM^t[i, j] = IM[j, i]$
<i>short</i> (IM)	<i>rmatrix</i>	$M[i, j] := short(IM[i, j])$
<i>shortup</i> (IM)	<i>rmatrix</i>	$M[i, j] := shortup(IM[i, j])$
<i>shortdown</i> (IM)	<i>rmatrix</i>	$M[i, j] := shortdown(IM[i, j])$
<i>short</i> (lv)	<i>rvector</i>	$v[j] := short(lv[j])$
<i>shortup</i> (lv)	<i>rvector</i>	$v[j] := shortup(lv[j])$
<i>shortdown</i> (lv)	<i>rvector</i>	$v[j] := shortdown(lv[j])$

Standardfunktionen des Moduls LMV\_ARI, Teil 1

n, n1, n2 = *integer*-Ausdruck, lv = *Lrvector*-Ausdruck, v = *rvector*-**Variable**  
M = *rmatrix*-**Variable**; IM, IM1, IM2 = *Lrmatrix*-Ausdruck;

Funktion	Typ	Bedeutung
short_test (lM,Err)	<i>rmatrix</i>	$M[i, j] := \text{short\_test}(lM[i, j], Err)$
shortup_test (lM,Err)	<i>rmatrix</i>	$M[i, j] := \text{shortup\_test}(lM[i, j], Err)$
shortdown_test (lM,Err)	<i>rmatrix</i>	$M[i, j] := \text{shortdown\_test}(lM[i, j], Err)$
short_test (lv,Err)	<i>rvector</i>	$v[j] := \text{short\_test}(lv[j], Err)$
shortup_test (lv,Err)	<i>rvector</i>	$v[j] := \text{shortup\_test}(lv[j], Err)$
shortdown_test (lv,Err)	<i>rvector</i>	$v[j] := \text{shortdown\_test}(lv[j], Err)$

Tabelle 7.6: Standardfunktionen des Moduls LMV\_ARI, Teil 2

n, n1, n2 = *integer*-Ausdruck, lv = *Lrvector*-Ausdruck, v = *rvector-Variable*  
M = *rmatrix-Variable*; lM, lM1, lM2 = *Lrmatrix*-Ausdruck;  
Err = *boolean*-Variable

### 7.3.5 Ein-/Ausgabeprozeduren

Es stehen die Prozeduren

```

procedure read (var f : text; var a : Lrvector);
procedure read (var f : text; var A : Lrmatrix);
procedure write (var f : text; a : Lrvector);
procedure write (var f : text; A : Lrmatrix);

```

mit optionalem Fileparameter, beliebig vielen Ein-/Ausgabeparametern, jedoch ohne Formatspezifikationen zur Verfügung.

Die Eingabe eines Vektors bzw. einer Matrix erfolgt komponentenweise entsprechend der Eingabe von *Longreal*-Werten, dabei wird eine Matrix stets zeilenweise eingelesen. Die Ausgabe eines Vektors bzw. einer Matrix erfolgt wie bei der Eingabe komponentenweise mit dem Standardformat für *Longreal*-Größen aus Tabelle 3.4.

#### Beispiel:

Durch die Anweisung

```
read (b, A, x); { b,x : Lrvector; A : Lrmatrix }
```

können der Vektor *b*, die Matrix *A* und der Vektor *x* hintereinander eingelesen werden.

### 7.3.6 Exakte Auswertung von Ausdrücken

Mit Hilfe des Moduls **MV\_ARI** lassen sich *rvector*-und *rmatrix*-Ausdrücke mit den drei verschiedenen Rundungsarten **maximalgenau** auswerten. So liefert z.B. die Anweisung

$$D := \# * (id(A) - R * A)$$

die Defektmatrix  $D$  vom Typ *rmatrix*, wobei bzgl. jeder Matrixkomponente nur eine einzige Rundung zur nächsten Rasterzahl erfolgt, vergl. Seite 163.

Mit  $a, rv$  vom Typ *rvector* und  $B$  vom Typ *rmatrix* liefert die Anweisung

$$rv := \# < (3 * a - B * a)$$

für die Variable  $rv$  bzgl. jeder Komponente den zur nächstkleineren *real*-Zahl gerundeten **exakten** Vektorausdruck  $3 \cdot a - B \cdot a$ , wobei  $B \cdot a$  das bekannte Produkt einer Matrix  $B$  mit einem Spaltenvektor  $a$  bedeutet.

Da die obigen #-Ausdrücke nicht anwendbar sind, wenn nur einer der Operanden vom Typ *Longreal*, *Lrvector* oder *Lrmatrix* ist, soll die maximalgenaue Auswertung von *Lrvector*- und *Lrmatrix*-Ausdrücken simuliert werden! Aus Speicherplatzgründen lassen sich aber jetzt keine Operatoren für *Lrmatrix*-Operanden definieren, wenn für jede Ergebniskomponente eine *dotprecision*-Variable erforderlich wird. Man muß vielmehr zur Auswertung eines *Lrmatrix*-Ausdrucks jede Ergebniskomponente **einzeln** in einer *dotprecision*-Variablen rundungsfehlerfrei berechnen und ihren Wert dann anschließend mit einer der drei verschiedenen Rundungsarten in die entsprechende Ergebnis-Komponente vom Typ *Longreal* abspeichern. Dadurch benötigt man für die Auswertung eines beliebigen *Lrmatrix*-Ausdrucks nur eine einzige *dotprecision*-Variable; man muß jedoch die Ausdruckauswertung in einer Laufanweisung programmieren. Neben den Operatoren *plus*, *minus*, *times* aus dem Standardmodul *STDMOD* benötigt man dazu aus dem Modul *LMV\_ARI* lediglich die vier zusätzlichen Operatoren *times* mit *Lrvector*/*rvector*-Operanden und dem Ergebnistyp *dotprecision*:

```
OPERATOR times (var a,b: Lrvector) res: dotprecision;
OPERATOR times (var a,b: rvector) res: dotprecision;
OPERATOR times (var a: Lrvector; var b: rvector) res: dotprecision;
OPERATOR times (var a: rvector; var b: Lrvector) res: dotprecision;
```

Mit diesen Operatoren lassen sich dann auch *Lrmatrix*-Ausdrücke maximalgenau auswerten, die zusätzlich noch *rmatrix*-Operanden enthalten. Für die maximalgenaue Auswertung eines *Lrvector*-Ausdrucks gelten ganz entsprechende Aussagen. Die folgenden Programmbeispiele zeigen die typischen Vorgehensweisen bei der maximalgenauen Auswertung von *Lrvector*- und *Lrmatrix*-Ausdrücken.

```
program Beispiel1 (input, output);
```

```
{* Ist  $R \approx A^{-1}$  eine Näherungsinverse, so liefert die Matrix D: Lrmatrix *}
{* den maximalgenauen Wert des Defekts  $\|E - R \cdot A\| \approx D$ , wobei E *}
{* die Einheitsmatrix mit dem Indexbereich von A ist. *}
```

```
use lmv_ari;
var E, R, A, D : Lrmatrix [1..3, 1..3];
    i, k       : integer;
    Akku       : dotprecision;
```

```

begin
  E := lid (A);   { E = Einheitsmatrix mit Indexbereich von A }
  For i := lb(E,1) to ub(E,1) do
    For k := lb(E,2) to ub(E,2) do
      begin
        Akku := E [ i, k ] - R [ i ] times LRVECTOR( A [ *, k ] );
        D [ i, k ] := LONGNEXT( Akku );
      end;
    end;
  end.

```

**Anmerkungen:**

1. Die Konstruktion `LRVECTOR ( A [ *, k ] )` liefert dem Compiler die Information, das Teilfeld `A [ *, k ]` als Vektor vom Typ *Lrvector* zu behandeln.
2. Um die Matrix `D` maximalgenau berechnen zu können, müssen alle Operatoren in der Anweisung für die *dotprecision*-Variable `Akku` vom Ergebnistyp *dotprecision* sein!
3. Der Operator `times` mit Ergebnistyp *dotprecision* und *Lrvector*-Operanden ist im Modul `LMV_ARI` definiert; vergleiche Tabelle 7.5
4. Der Operator `-` mit *Longreal*- und *dotprecision*- Operanden liefert den Ergebnistyp *dotprecision* und könnte auch durch den Operator `minus` ersetzt werden. Beide Operatoren sind im Modul `STDMOD` definiert, vergl. Tabelle 3.9

Im nächsten Programmbeispiel wird mit  $a, lv$  vom Typ *Lrvector* und  $B$  vom Typ *rmatrix* der *Lrvector*-Ausdruck

$$3 \cdot a - B \cdot a$$

maximalgenau ausgewertet, wobei das exakte Ergebnis komponentenweise zur nächstkleineren *Longreal*-Zahl gerundet und in die Variable  $lv$  geschrieben wird.

```

program Beispiel2 (input, output);
use lmv_ari;
var B      : rmatrix [1..3, 1..3];
    i      : integer;
    a, lv  : Lrvector [1..3];
    Akku   : dotprecision;
begin
  For i := lb(a) to ub(a) do
    begin
      Akku := 3 times a [ i ] minus B [ i ] times a;
      lv [ i ] := LONGDOWN( Akku );
    end;
  end.

```



**Anmerkungen:**

1. Um den Vektor *lv* maximalgenau berechnen zu können, müssen alle Operatoren in der Anweisung für die *dotprecision*-Variable Akku vom Ergebnistyp *dotprecision* sein!
2. Der erste Operator *times* mit *integer/Longreal*-Operanden und der Operator *minus* mit *dotprecision*-Operanden haben den gemeinsamen Ergebnistyp *dotprecision*; vergl. Tabelle 3.9
3. Der zweite Operator *times* mit *rvector/Lrvector*-Operanden hat den Ergebnistyp *dotprecision*; vergl. Tabelle 7.5 und Seite 169.

## 7.4 Das Modul MVI\_ARI

In diesem Modul werden die für das Rechnen mit reellen Intervallvektoren und Intervallmatrizen notwendigen Operatoren, Funktionen und Prozeduren auf der Grundlage des Datentyps `real` bereitgestellt.

### 7.4.1 Datentypen

Die dynamischen Datentypen zur Darstellung von Intervallvektoren und Intervallmatrizen sind entsprechend der Vereinbarung

```
type ivector = dynamic array [*] of interval;
   imatrix  = dynamic array [*] of ivector;
```

im Sprachkern von PASCAL-XSC enthalten. Die aktuellen Indexgrenzen werden erst bei der Vereinbarung von Variablen dieser Typen festgelegt.

### 7.4.2 Operatoren

Viele der bekannten, grundlegenden Matrix/Vektor-Operationen sind in diesem Modul vordefiniert. Als arithmetische Operatoren stehen die monadischen Operatoren  $+$ ,  $-$  und die vier Grundoperationen  $+$ ,  $-$ ,  $*$ ,  $/$  mit der komponentenweisen Rundung zum kleinsten einschließenden Intervall (auch für gemischte Typen) zur Verfügung, wobei jedoch nur bestimmte Operandenkombinationen zulässig sind. Während die Operationen  $+$  und  $-$  für Intervallvektoren und Intervallmatrizen komponentenweise erklärt sind gemäß

$$\begin{aligned} c &:= a \pm b & \text{mit} & & c[i] &:= a[i] \pm b[i] \\ C &:= A \pm B & \text{mit} & & C[i, j] &:= A[i, j] \pm B[i, j] \end{aligned}$$

mit  $a, b, c$  vom Typ *ivector* und  $A, B, C$  vom Typ *imatrix*, sind die Operationen  $*$  und  $/$  definiert durch

$$\begin{aligned} s &:= a * b & \text{mit} & & s &:= \#\# \left( \text{for } i := \text{lb}(a) \text{ to } \text{ub}(a) \right. \\ & & & & & \left. \text{sum } (a[i] * b[i]) \right) \ddagger \\ c &:= r * a & \text{mit} & & c[i] &:= r * a[i] \\ c &:= a * r & \text{mit} & & c[i] &:= a[i] * r \\ c &:= a / r & \text{mit} & & c[i] &:= a[i] / r \\ c &:= A * b & \text{mit} & & c[i] &:= A[i] * b \ddagger \\ C &:= r * A & \text{mit} & & C[i, j] &:= r * A[i, j] \\ C &:= A * r & \text{mit} & & C[i, j] &:= A[i, j] * r \\ C &:= A / r & \text{mit} & & C[i, j] &:= A[i, j] / r \\ C &:= A * B & \text{mit} & & C[i, j] &:= A[i] * B[*, j] \ddagger \end{aligned}$$


---

$\ddagger$ Maximalgenaues Skalarprodukt mit Ergebnistyp *interval*

wobei  $r, s$  vom Typ *interval*,  $a, b, c$  vom Typ *ivector* und  $A, B, C$  vom Typ *imatrix* sind. Die Operationen mit gemischten Operandentypen sind entsprechend definiert.

rechter Operand	integer real	interval	rvector	ivector	rmatrix	imatrix
linker Operand						
monadisch				+, -		+, -
integer real				*		*
interval			*	*	*	*
rvector		*, /	+*	$\diamond$ =, <>, <b>in</b> +*		
ivector	*, /	*, /	$\diamond$ =, <> +*	$\diamond$ <b>in</b> , $\vee$ , >< +*, **		
rmatrix		*, /		*	+*	$\diamond$ =, <>, <b>in</b> +*
imatrix	*, /	*, /	*	*	$\diamond$ =, <> +*	$\diamond$ <b>in</b> , $\vee$ , >< +*, **

Tabelle 7.7: Die Operatoren des Moduls MVLARI

$$\vee \in \{=, <>, <, <=, >, >=\}; \quad \diamond \in \{+, -, *\}$$

Die Vergleichsoperatoren  $=, <>, <, <=, >, >=$  sind wie im Falle von Intervallen mengentheoretisch definiert und wirken auf alle Komponenten eines Intervallvektors oder einer Intervallmatrix. Für  $a, b$  vom Typ *ivector* und  $A, B$  vom Typ *imatrix* gilt:

$$\begin{aligned} a <= b &\iff a[i] <= b[i] \quad \text{für alle } i \\ A <= B &\iff A[i, j] <= B[i, j] \quad \text{für alle } i, j \end{aligned}$$

Die Operatoren auf der rechten Seite der Äquivalenzbeziehungen sind dabei die Operatoren für Operanden vom Typ *interval*.

Zusätzlich stehen noch die Operatoren **in** für die Relation „liegt in“ zwischen einem *rvector*- und einem *ivector*-Operanden bzw. zwischen einem *rmatrix*- und einem *imatrix*-Operanden und für die Relation „enthalten im Innern“ zwischen

zwei *ivector*-Operanden bzw. *imatrix*-Operanden sowie  $><$  für den Test auf Disjunktheit zweier Intervallvektoren bzw. Intervallmatrizen zur Verfügung. Diese Operatoren sind jeweils komponentenweise definiert. Die Verbandsoperatoren  $+*$  und  $**$  bezeichnen die komponentenweise Bildung der Intervall-Hülle bzw. des Durchschnitts, wie es bereits für den Typ *interval* auf Seite 76 beschrieben wurde.

### 7.4.3 Transferfunktionen

Zur Wandlung zwischen den Typen *rvector* und *ivector* werden folgende Transferfunktionen bereitgestellt:

Funktion	Ergebnistyp	Bedeutung
intval (rv1,rv2)	<i>ivector</i>	Intervallvektor <i>iv</i> mit $iv[i] = \text{intval}(rv1[i], rv2[i])$
intval (rv)	<i>ivector</i>	Punktintervallvektor <i>iv</i> mit $iv[i] = \text{intval}(rv[i])$
inf (iv)	<i>rvector</i>	Vektor <i>rv</i> der Untergrenzen mit $rv[i] = \text{inf}(iv[i])$
sup (iv)	<i>rvector</i>	Vektor <i>rv</i> der Obergrenzen mit $rv[i] = \text{sup}(iv[i])$

$rv, rv1, rv2 = rvector$ -Ausdruck mit:  $rv1 \leq rv2$ ;  $iv = ivector$ -Ausdruck

Zur Wandlung zwischen den Typen *rmatrix* und *imatrix* werden folgende Transferfunktionen bereitgestellt:

Funktion	Ergebnistyp	Bedeutung
intval (rM1,rM2)	<i>imatrix</i>	Intervallmatrix <i>iM</i> mit $iM[i, j] = \text{intval}(rM1[i, j], rM2[i, j])$
intval (rM)	<i>imatrix</i>	Punktintervallmatrix <i>iM</i> mit $iM[i, j] = \text{intval}(rM[i, j])$
inf (iM)	<i>rmatrix</i>	Matrix <i>rM</i> der Untergrenzen mit $rM[i, j] = \text{inf}(iM[i, j])$
sup (iM)	<i>rmatrix</i>	Matrix <i>rM</i> der Obergrenzen mit $rM[i, j] = \text{sup}(iM[i, j])$

$rM, rM1, rM2 = rmatrix$ -Ausdruck mit:  $rM1 \leq rM2$ ;  $iM = imatrix$ -Ausdruck

### 7.4.4 Überladungen des Zuweisungsoperators

Die komponentenweise Initialisierung von *ivector*- und *imatrix*-Variablen sowie die Wandlung von *rvector* nach *ivector* und *rmatrix* nach *imatrix* werden in Form überladener Zuweisungen bereitgestellt:

Zuweisung	Bedeutung
$iv := r$	$iv[j] := \text{intval}(r); \quad j = \text{lb}(iv), \dots, \text{ub}(iv)$
$iv := i$	$iv[j] := i; \quad j = \text{lb}(iv), \dots, \text{ub}(iv)$
$iv := rv$	$iv := \text{intval}(rv);$
$iM := r$	$iM[j, k] := \text{intval}(r); \quad j = \text{lb}(iM, 1), \dots, \text{ub}(iM, 1)$ $k = \text{lb}(iM, 2), \dots, \text{ub}(iM, 2)$
$iM := i$	$iM[j, k] := i; \quad j = \text{lb}(iM, 1), \dots, \text{ub}(iM, 1)$ $k = \text{lb}(iM, 2), \dots, \text{ub}(iM, 2)$
$iM := rM$	$iM := \text{intval}(rM);$

Tabelle 7.8: Überladungen des Zuweisungsoperators

$i$  = *interval*-Ausdruck,  $iv$  = *ivector-Variable*,  $iM$  = *imatrix-Variable*  
 $r$  = *real*-Ausdruck,  $rv$  = *rvector*-Ausdruck,  $rM$  = *rmatrix*-Ausdruck

### 7.4.5 Standardfunktionen

Es stehen die Funktionen *id* und *null* zur Erzeugung einer Einheitsmatrix und einer Nullmatrix bzw. eines Nullvektors, die Funktionen *mid* und *diam* für die komponentenweise Berechnung von Mittelpunkt und Durchmesser sowie die Funktion *transp* zur Berechnung der transponierten Intervallmatrix zur Verfügung. Darüberhinaus wird die komponentenweise definierte Funktion *blow* für die Epsilonaufblähung bereitgestellt.

Funktion	Ergebnistyp	Bedeutung
$\text{null}(iv)$	<i>rvector</i>	Nullvektor mit dem aktuellen Indexbereich von <i>iv</i>
$\text{null}(iM)$	<i>rmatrix</i>	Nullmatrix mit den aktuellen Indexbereichen von <i>iM</i>
$\text{null}(iM1, iM2)$	<i>rmatrix</i>	Nullmatrix mit den aktuellen Indexbereichen der Produktmatrix $iM1 \cdot iM2$

Funktion	Ergebnistyp	Bedeutung
id (iM)	<i>rmatrix</i>	Einheitsmatrix mit den aktuellen Indexbereichen von $iM$
id (iM1,iM2)	<i>rmatrix</i>	Einheitsmatrix mit den aktuellen Indexbereichen der Produktmatrix $iM1 \cdot iM2$
mid (iv)	<i>rvector</i>	Mittelpunktvektor $pv$ mit $pv[i] = mid(iv[i])$
diam (iv)	<i>rvector</i>	Durchmesservektor $dv$ mit $dv[i] = diam(iv[i])$
mid (iM)	<i>rmatrix</i>	Mittelpunktmatrix $pM$ mit $pM[i, j] = mid(iM[i, j])$
diam (iM)	<i>rmatrix</i>	Durchmessermatrix $dM$ mit $dM[i, j] = diam(iM[i, j])$
transp (iM)	<i>imatrix</i>	Transponierte Matrix $iMt$ von $iM$ mit $iMt[i, j] = iM[j, i]$
blow (iv,r)	<i>ivector</i>	Vektorielle Epsilonaufblähung $ev$ mit $ev[i] = blow(iv[i], r)$
blow (iM,r)	<i>imatrix</i>	Matrix-Epsilonaufblähung $eM$ mit $eM[i, j] = blow(iM[i, j], r)$

Tabelle 7.9: Standardfunktionen des Moduls MVL\_ARI

$r$  = *real*-Ausdruck,  $iv$  = *ivector*-Ausdruck,  
 $iM, iM1, iM2$  = *imatrix*-Ausdruck

**Beispiel:**

$A, D$  und  $R$  seien Variable vom Typ *imatrix*.  $A$  sei die quadratische Punktintervallmatrix der Punkt-Koeffizienten eines linearen Gleichungssystems.  $R \approx A^{-1}$  sei die quadratische Punktintervallmatrix einer Näherungsinversen zu  $A$ . Dann läßt sich mit der Anweisung

$$D := id(A) - R * A;$$

eine garantierte Intervalleinschließung für den Defekt

$$DEF = E - R \cdot A, \quad (E = \text{Einheitsmatrix})$$

berechnen, die wegen Intervallaufblähungen für die Praxis aber meist unbrauchbar ist!

Unter Verwendung eines Lattenkreuzausdrucks erhält man jedoch mit der Anweisung

$$D := \#\#(id(A) - R * A);$$

durch die Variable  $D$  die **engstmögliche** Einschließung für den Defekt  $DEF$ . In der letzten Anweisung können für  $R$  und  $A$  auch die entsprechenden Matrizen vom Typ *rmatrix* gewählt werden!

### 7.4.6 Ein-/Ausgabeprozeduren

Es stehen die Prozeduren

```
procedure read (var f : text; var a : ivector);
procedure read (var f : text; var A : imatrix);
procedure write (var f : text; a : ivector);
procedure write (var f : text; A : imatrix);
```

mit optionalem Fileparameter, beliebig vielen Ein-/Ausgabeparametern, jedoch ohne Formatspezifikationen zur Verfügung.

Die Eingabe eines Intervallvektors bzw. einer Intervallmatrix erfolgt komponentenweise entsprechend der Eingabe von *interval*-Werten, dabei wird eine Matrix stets zeilenweise eingelesen. Die Ausgabe eines Intervallvektors bzw. einer Intervallmatrix erfolgt wie bei der Eingabe komponentenweise mit dem Standardformat für *real*-Größen aus Tabelle 3.4.

## 7.5 Das Modul LMVI\_ARI

In diesem Modul werden die für das Rechnen mit reellen Intervallvektoren und Intervallmatrizen notwendigen Operatoren, Funktionen und Prozeduren auf der Grundlage des Datentyps `Longreal` bereitgestellt.

### 7.5.1 Datentypen

Die dynamischen Datentypen zur Darstellung von Intervallvektoren und Intervallmatrizen sind entsprechend der Vereinbarung

```
type Livector = dynamic array [*] of Linterval;
   Limatrix  = dynamic array [*] of Livector;
```

definiert. Die aktuellen Indexgrenzen werden erst bei der Vereinbarung von Variablen dieser Typen festgelegt.

### 7.5.2 Operatoren

Viele der bekannten, grundlegenden Matrix/Vektor-Operationen sind in diesem Modul vordefiniert. Als arithmetische Operatoren stehen die monadischen Operatoren  $+$ ,  $-$  und die vier Grundoperationen  $+$ ,  $-$ ,  $*$ ,  $/$  mit der komponentenweisen Rundung zum kleinsten einschließenden Intervall (auch für gemischte Typen) zur Verfügung, wobei jedoch nur bestimmte Operandenkombinationen zulässig sind. Während die Operationen  $+$  und  $-$  für Intervallvektoren und Intervallmatrizen komponentenweise erklärt sind gemäß

$$\begin{aligned} c &:= a \pm b & \text{mit} & & c[i] &:= a[i] \pm b[i] \\ C &:= A \pm B & \text{mit} & & C[i, j] &:= A[i, j] \pm B[i, j] \end{aligned}$$

mit  $a, b, c$  vom Typ *ivector* und  $A, B, C$  vom Typ *imatrix*, sind die Operationen  $*$  und  $/$  definiert durch

$$\begin{aligned} s &:= a * b & \text{mit} & & s &:= \sum_i a[i] \cdot b[i] \quad \ddagger \\ c &:= r * a & \text{mit} & & c[i] &:= r * a[i] \\ c &:= a * r & \text{mit} & & c[i] &:= a[i] * r \\ c &:= a / r & \text{mit} & & c[i] &:= a[i] / r \\ c &:= A * b & \text{mit} & & c[i] &:= A[i] * b \quad \ddagger \\ C &:= r * A & \text{mit} & & C[i, j] &:= r * A[i, j] \\ C &:= A * r & \text{mit} & & C[i, j] &:= A[i, j] * r \\ C &:= A / r & \text{mit} & & C[i, j] &:= A[i, j] / r \\ C &:= A * B & \text{mit} & & C[i, j] &:= A[i] * B[*, j] \quad \ddagger \end{aligned}$$


---

$\ddagger$ Maximalgenaues Skalarprodukt mit Ergebnistyp *Linterval*



wobei  $r, s$  vom Typ *Linterval*,  $a, b, c$  vom Typ *Livector* und  $A, B, C$  vom Typ *Limatrix* sind. Die Operationen mit gemischten Operandentypen sind entsprechend definiert.

rechter Operand	Longreal	Linterval	Lrvector	Livector	Lrmatrix	Limatrix
linker Operand						
monad.				+, -		+, -
Longreal				*		*
Linterval			*	*	*	*
Lrvector		*, /	+*	◇ =, <>, in +*		
Livector	*, /	*, /	◇ =, <> +*	◇ in, V, >< +*, **		
Lrmatrix		*, /		*	+*	◇ =, <>, in +*
Limatrix	*, /	*, /	*	*	◇ =, <> +*	◇ in, V, >< +*, **

Tabelle 7.10: Operatoren aus LMVL\_ARI mit Ergebnistypen: *Linterval*, *Livector*, *Limatrix* oder *boolean*;  $\forall \in \{=, <>, <, <=, >, >=\}$ ;  $\diamond \in \{+, -, *\}$

rechter Operand	rvector	Lrvector	ivector	Livector
linker Operand				
rvector				times_i
Lrvector			times_i	times_i
ivector		times_i		
Livector	times_i	times_i		times_i

Tabelle 7.11: Operator times\_i mit Ergebnistyp *idotprecision*

Zusätzlich zu den Operatoren aus Tabelle 7.10 steht zur exakten Auswertung von *Livector*- und *Limatrix*-Ausdrücken noch der Operator **times\_i** mit dem Ergebnistyp *idotprecision* zur Verfügung. Die erlaubten Operandenkombinationen findet man in Tabelle 7.11.

Die Vergleichsoperatoren  $=, <>, <, <=, >, >=$  sind wie im Falle von Intervallen mengentheoretisch definiert und wirken auf alle Komponenten eines Intervallvektors oder einer Intervallmatrix. Für  $a, b$  vom Typ *Livector* und  $A, B$  vom Typ *Limatrix* gilt:

$$\begin{aligned} a <= b &\iff a[i] <= b[i] \quad \text{für alle } i \\ A <= B &\iff A[i, j] <= B[i, j] \quad \text{für alle } i, j \end{aligned}$$

Die Operatoren auf der rechten Seite der Äquivalenzbeziehungen sind dabei die Operatoren für Operanden vom Typ *Linterval*.

Zusätzlich stehen die Operatoren **in** für die Relation „liegt in“ zwischen einem *Lrvector*- und einem *Livector*-Operanden bzw. zwischen einem *Lrmatrix*- und einem *Limatrix*-Operanden und für die Relation „enthalten im Innern“ zwischen zwei *Livector*-Operanden bzw. *Limatrix*-Operanden sowie  $><$  für den Test auf Disjunktheit zweier Intervallvektoren bzw. Intervallmatrizen zur Verfügung. Diese Operatoren sind jeweils komponentenweise definiert. Die Verbandsoperatoren  $+*$  und  $**$  bezeichnen die komponentenweise Bildung der Intervall-Hülle bzw. des Durchschnitts, wie es bereits für den Typ *interval* auf Seite 76 beschrieben wurde.

### 7.5.3 Transferfunktionen

Zur Wandlung zwischen den Typen *Lrvector* und *Livector* bzw. *Livector* und *ivector* werden folgende Transferfunktionen bereitgestellt:

Funktion	Ergebnistyp	Bedeutung
intval (lv1,lv2)	<i>Livector</i>	Intervallvektor <i>liv</i> mit $liv[i] = \text{intval}(lv1[i], lv2[i])$
intval (lv)	<i>Livector</i>	Punktintervallvektor <i>liv</i> mit $liv[i] = \text{intval}(lv[i])$
inf (liv)	<i>Lrvector</i>	Vektor <i>lv</i> der Untergrenzen mit $lv[i] = \text{inf}(liv[i])$
sup (liv)	<i>Lrvector</i>	Vektor <i>lv</i> der Obergrenzen mit $lv[i] = \text{sup}(liv[i])$
short (liv)	<i>ivector</i>	$iv := \text{short}(liv); \quad liv[i] \subset iv[i]$

$lv, lv1, lv2 = \text{Lrvector}$ -Ausdruck mit:  $lv1 <= lv2$ ;  $liv = \text{Livector}$ -Ausdruck  
 $iv = \text{ivector}$ -Variable

Zur Wandlung zwischen den Typen *Lrmatrix* und *Limatrix* bzw. *Limatrix* und *imatrix* werden folgende Transferfunktionen bereitgestellt:

Funktion	Ergebnistyp	Bedeutung
intval (lM1,lM2)	<i>Limatrix</i>	Intervallmatrix <i>liM</i> mit $liM[i, j] = \text{lintval}(lM1[i, j], lM2[i, j])$
intval (lM)	<i>Limatrix</i>	Punktintervallmatrix <i>liM</i> mit $liM[i, j] = \text{lintval}(lM[i, j])$
inf (liM)	<i>Lrmatrix</i>	Matrix <i>lM</i> der Untergrenzen mit $lM[i, j] = \text{inf}(liM[i, j])$
sup (liM)	<i>Lrmatrix</i>	Matrix <i>lM</i> der Obergrenzen mit $lM[i, j] = \text{sup}(liM[i, j])$
short (liM)	<i>imatrix</i>	$iM := \text{short}(liM); liM[i, j] \subset iM[i, j]$

lM, lM1, lM2 = *Lrmatrix*-Ausdruck mit: lM1 ≤ lM2; liM = *Limatrix*-Ausdruck  
liM = *Limatrix*-Ausdruck; iM = *imatrix*-Variable

#### 7.5.4 Überladungen des Zuweisungsoperators

Zuweisung	Bedeutung
liv := lr	liv[j] := lintval(lr); j = lb(liv),... ,ub(liv)
liv := li	liv[j] := li; j = lb(liv),... ,ub(liv)
liv := lv	liv := intval(lv);
liM := lr	liM[j, k] := lintval(lr); j = lb(liM,1),... ,ub(liM,1) k = lb(liM,2),... ,ub(liM,2)
liM := li	liM[j, k] := li; j = lb(liM,1),... ,ub(liM,1) k = lb(liM,2),... ,ub(liM,2)
liM := lM	liM := intval(lM);
liM := iM	liM[j, k] := iM[j, k]; j = lb(liM,1),... ,ub(liM,1) k = lb(liM,2),... ,ub(liM,2)
liv := iv	liv[j] := iv[j]; j = lb(liv),... ,ub(liv)

Tabelle 7.12: Überladungen des Zuweisungsoperators

li = *Linterval*-Ausdruck, liv = *Livector-Variable*, liM = *Limatrix-Variable*  
lr = *Longreal*-Ausdruck, lv = *Lrvector*-Ausdruck, lM = *Lrmatrix*-Ausdruck  
iM = *imatrix*-Ausdruck, iv = *ivector-Variable*

## 7.5.5 Standardfunktionen

Es stehen die Funktionen *id* und *null* zur Erzeugung einer Einheitsmatrix und einer Nullmatrix bzw. eines Nullvektors, die Funktionen *mid* und *diam* für die komponentenweise Berechnung von Mittelpunkt und Durchmesser sowie die Funktion *transp* zur Berechnung der transponierten Intervallmatrix zur Verfügung. Darüberhinaus werden die komponentenweise definierten Funktionen *blow* und *blow\_z* für die Epsilonaufblähung bereitgestellt.

Funktion	Ergebnistyp	Bedeutung
<i>null</i> ( <i>liv</i> )	<i>Lrvector</i>	Nullvektor mit dem aktuellen Indexbereich von <i>liv</i>
<i>null</i> ( <i>liM</i> )	<i>Lrmatrix</i>	Nullmatrix mit den aktuellen Indexbereichen von <i>liM</i>
<i>null</i> ( <i>liM1</i> , <i>liM2</i> )	<i>Lrmatrix</i>	Nullmatrix mit den aktuellen Indexbereichen der Produktmatrix $liM1 \cdot liM2$
<i>id</i> ( <i>liM</i> )	<i>Lrmatrix</i>	Einheitsmatrix mit den aktuellen Indexbereichen von <i>liM</i>
<i>id</i> ( <i>liM1</i> , <i>liM2</i> )	<i>Lrmatrix</i>	Einheitsmatrix mit den aktuellen Indexbereichen der Produktmatrix $liM1 \cdot liM2$
<i>mid</i> ( <i>liv</i> )	<i>Lrvector</i>	Mittelpunktvektor <i>pv</i> mit $pv[i] = mid(liv[i])$
<i>diam</i> ( <i>liv</i> )	<i>Lrvector</i>	Durchmesservektor <i>dv</i> mit $dv[i] = diam(liv[i])$
<i>mid</i> ( <i>liM</i> )	<i>Lrmatrix</i>	Mittelpunktmatrix <i>pM</i> mit $pM[i, j] = mid(liM[i, j])$
<i>diam</i> ( <i>liM</i> )	<i>Lrmatrix</i>	Durchmessermatrix <i>dM</i> mit $dM[i, j] = diam(liM[i, j])$
<i>transp</i> ( <i>liM</i> )	<i>Limatrix</i>	Transponierte Matrix <i>liMt</i> von <i>liM</i> mit $liMt[i, j] = liM[j, i]$
<i>blow</i> ( <i>liv</i> , <i>r</i> )	<i>Livector</i>	Vektorielle Epsilonaufblähung <i>ev</i> mit $ev[i] = blow(liv[i], r)$ [0, 0] $\rightarrow$ $[-10^{-531}, +10^{-531}]$
<i>blow_z</i> ( <i>liv</i> , <i>r</i> )	<i>Livector</i>	Wie <i>blow</i> ( <i>liv</i> , <i>r</i> ); [0, 0] $\rightarrow$ $[-10^{-255}, +10^{-255}]$
<i>blow</i> ( <i>liM</i> , <i>r</i> )	<i>Limatrix</i>	Matrix-Epsilonaufblähung <i>eM</i> mit $eM[i, j] = blow(liM[i, j], r)$ [0, 0] $\rightarrow$ $[-10^{-531}, +10^{-531}]$
<i>blow_z</i> ( <i>liM</i> , <i>r</i> )	<i>Limatrix</i>	Wie <i>blow</i> ( <i>liM</i> , <i>r</i> ); [0, 0] $\rightarrow$ $[-10^{-255}, +10^{-255}]$

Tabelle 7.13: Standardfunktionen des Moduls LMVL\_ARI

*r* = *real*-Ausdruck, *liv* = *Livector*-Ausdruck,  
*liM*, *liM1*, *liM2* = *Limatrix*-Ausdruck

### 7.5.6 Ein-/AusgabeprozEDUREN

Es stehen die Prozeduren

```

procedure read (var f : text; var a : Livector);
procedure read (var f : text; var A : Limatrix);
procedure write (var f : text; a : Livector);
procedure write (var f : text; A : Limatrix);

```

mit optionalem Fileparameter, beliebig vielen Ein-/Ausgabeparametern, jedoch ohne Formatspezifikationen zur Verfügung.

Die Eingabe eines Intervallvektors bzw. einer Intervallmatrix erfolgt komponentenweise entsprechend der Eingabe von *Linterval*-Werten, dabei wird eine Matrix stets zeilenweise eingelesen. Die Ausgabe eines Intervallvektors bzw. einer Intervallmatrix erfolgt wie bei der Eingabe komponentenweise mit dem Standardformat für *Longreal*-Größen aus Tabelle 3.4.

### 7.5.7 Exakte Auswertung von Ausdrücken

Mit Hilfe des Moduls *MV\_LARI* lassen sich *ivector*- und *imatrix*-Ausdrücke **maximalgenau** auswerten. Bedeutet z.B.  $A$  vom Typ *imatrix* die quadratische Koeffizientenmatrix eines linearen Gleichungssystems  $A \cdot x = b$ ,  $R$  vom Typ *imatrix* die Näherungsinverse  $R \approx A^{-1}$  und ist  $id(A)$  die Einheitsmatrix mit dem Indexbereich von  $A$ , so liefert die Anweisung

$$D := \#\#(id(A) - R * A)$$

die Defektmatrix  $D$  vom Typ *imatrix*, wobei bzgl. jeder Matrixkomponente nur eine einzige Rundung pro Intervallgrenze zur nächstkleineren bzw. nächstgrößeren Rasterzahl erfolgt.

Da der obige *##*-Ausdruck mit  $A$  vom Typ *Limatrix* nicht anwendbar ist, soll die maximalgenaue Auswertung von *Livector*- und *Limatrix*-Ausdrücken simuliert werden! Aus Speicherplatzgründen lassen sich aber jetzt z.B. keine Operatoren für *Limatrix*-Operanden definieren, wenn für jede Ergebniskomponente eine *idotprecision*-Variable erforderlich wird. Man muß vielmehr zur Auswertung eines *Limatrix*-Ausdrucks jede Ergebniskomponente **einzeln** in einer *idotprecision*-Variablen rundungsfehlerfrei berechnen und ihren Wert dann anschließend mit Hilfe der Funktion *LONG* (vergl. Seite 89) mit nur einer einzigen Rundung pro Intervallgrenze durch eine *Linterval*-Variable einschließen. Dadurch benötigt man für die Auswertung eines beliebigen *Limatrix*-Ausdrucks nur eine einzige *idotprecision*-Variable; man muß jedoch die Ausdrucksauswertung in einer Laufanweisung programmieren. Neben den Operatoren *plus\_i*, *minus\_i*, *times\_i* aus dem Modul *LLARI* benötigt man dazu aus dem Modul *LMV\_LARI* lediglich die sieben zusätzlichen Operatoren *times\_i* aus Tabelle 7.11 mit *rvector*/*Lrvector*/*ivector*/*Livector*-Operanden und dem Ergebnistyp *idotprecision*. Mit diesen Operatoren lassen sich dann auch *Limatrix*-Ausdrücke maximalgenau auswerten, die zusätzlich noch *imatrix*- oder *rmatrix*-Operanden enthalten. Für die maximalgenaue Auswertung eines *Livector*-Ausdrucks gelten ganz entsprechende Aussagen.

Die beiden folgenden Programmbeispiele zeigen die typischen Vorgehensweisen bei der maximalgenauen Auswertung eines *Limatrix*-Ausdrucks:

```

program Beispiel1 (input, output);

{ * Defekteinschließung eines linearen Gleichungssystems  $A \cdot x = b$  * }
{ * Ist  $R \approx A^{-1}$  eine Näherungsinverse, so liefert die Matrix D: Limatrix * }
{ * die maximalgenaue Einschließung des Defekts  $E - R \cdot A$ , wobei * }
{ * E die Einheitsmatrix mit dem Indexbereich von A ist. * }

use lmvi_ari, lmv_ari, li_ari;
var A, D, R : Limatrix [1..3, 1..3];
    E : Lrmatrix [1..3, 1..3];
    i, k : integer;
    idot : idotprecision;

begin
  E := id (A); { E = Einheitsmatrix mit Indexbereich von A }
  For i := lb(E,1) to ub(E,1) do
    For k := lb(E,2) to ub(E,2) do
      begin
        idot := E [i, k] minus_i R [i] times_i LIVECTOR ( A [*, k] );
        D [i, k] := LONG ( idot );
      end;
    end;
  end.

```

#### Anmerkungen:

1. In diesem Beispiel wird angenommen, daß die Koeffizienten der Matrix  $A$  **Intervalle** sind, so daß  $A, R$  vom Typ *Limatrix* zu wählen sind.
2. Die Konstruktion `LIVECTOR ( A [*, k] )` liefert dem Compiler die Information, das Teilfeld  $A [*, k]$  als Vektor vom Typ *Livector* zu behandeln.
3. Um die Matrix  $D$  maximalgenau berechnen zu können, müssen alle Operatoren in der Anweisung für die *idotprecision*-Variable `idot` vom Ergebnistyp *idotprecision/idotprecision* sein!
4. Der Operator `times_i` mit dem Ergebnistyp *idotprecision* und *Livector*-Operanden ist im Modul `LMVI_ARI` definiert; vergleiche Tabelle 7.11
5. Der Operator `minus_i` mit *Longreal*- und *idotprecision*- Operanden liefert den Ergebnistyp *idotprecision*; vergl. Tabelle 4.5

Im nächsten Programmbeispiel wird wieder der Defekt  $E - R \cdot A$  eines linearen Gleichungssystems  $A \cdot x = b$  maximalgenau eingeschlossen, wenn  $R \approx A^{-1}$  eine berechnete Näherungsinverse ist. Jetzt wird jedoch angenommen, daß die Koeffizienten der Matrix  $A$  Dezimalzahlen sind, d.h.  $A, R$  müssen vom Typ *Lrmatrix* gewählt werden.

```

program Beispiel2 (input, output);

{* Defekteinschließung eines linearen Gleichungssystems  $A \cdot x = b$  *}
{* Ist  $R \approx A^{-1}$  eine Näherungsinverse, so liefert die Matrix D: Limatrix *}
{* die maximalgenaue Einschließung des Defekts  $E - R \cdot A$ , wobei *}
{* E die Einheitsmatrix mit dem Indexbereich von A ist. *}

use lmvi_ari, lmv_ari, li_ari;
var D : Limatrix [1..3, 1..3];
    A,E,R : Lrmatrix [1..3, 1..3];
    i, k : integer;
    idot : idotprecision;

begin
    E := lid (A); { E = Einheitsmatrix mit Indexbereich von A }
    For i := lb(E,1) to ub(E,1) do
        For k := lb(E,2) to ub(E,2) do
            begin
                idot := E [i, k] - R [i] times LVECTOR ( A [*, k] );
                D [i, k] := LONG ( idot );
            end;
        end;
    end.

```

**Anmerkungen:**

1. In diesem Beispiel wird angenommen, daß die Koeffizienten der Matrix  $A$  *Longreal*-Zahlen sind, so daß  $A, R$  vom Typ *Lrmatrix* zu wählen sind.
2. Die Konstruktion `LVECTOR ( A [*, k] )` liefert dem Compiler die Information, das Teilfeld  $A [*, k]$  als Vektor vom Typ *Lrvector* zu behandeln.
3. Um die Matrix  $D$  maximalgenau berechnen zu können, müssen alle Operatoren in der Anweisung für die *idotprecision*-Variable `idot` vom Ergebnistyp *idotprecision/idotprecision* sein!
4. Der Operator `times` mit dem Ergebnistyp *idotprecision* und *Lrvector*-Operanden ist im Modul LMV\_LARI definiert; vergleiche Tabelle 7.5
5. Der Operator `-` mit *Longreal*- und *idotprecision*- Operanden liefert den Ergebnistyp *idotprecision*; vergl. Tabelle 3.9
6. Die Anweisung `D [i, k] := LONG ( idot )` liefert die **bestmögliche**, d.h. maximalgenaue Einschließung des Intervalls *idot* vom Typ *idotprecision* durch ein Intervall  $D [i, k]$  vom Typ *Linterval*; vergl. Tabelle 4.7 auf Seite 89.

## 7.6 Das Modul MVC\_ARI

In diesem Modul werden die für das Rechnen mit komplexen Vektoren und Matrizen notwendigen Operatoren, Funktionen und Prozeduren auf der Grundlage des Datentyps `real` bereitgestellt.

### 7.6.1 Datentypen

Die dynamischen Datentypen zur Darstellung von komplexen Vektoren und Matrizen sind entsprechend der Vereinbarung

```
type cvector = dynamic array [*] of complex;
      cmatrix  = dynamic array [*] of cvector;
```

im Sprachkern von PASCAL-XSC enthalten. Die aktuellen Indexgrenzen werden erst bei der Vereinbarung von Variablen dieser Typen festgelegt.

### 7.6.2 Operatoren

Viele der bekannten grundlegenden komplexen Matrix/Vektor-Operationen sind in diesem Modul vordefiniert. Als arithmetische Operatoren stehen die monadischen Operatoren `+`, `-` und die vier Grundoperationen `+`, `-`, `*`, `/` mit den drei verschiedenen Rundungsarten zur Verfügung, wobei jedoch nur bestimmte Operandenkombinationen zulässig sind. Während die Operationen `+` und `-` für Vektoren und Matrizen komponentenweise erklärt sind gemäß

$$\begin{aligned} c & := a \pm b \quad \text{mit} \quad c[i] := a[i] \pm b[i] \\ C & := A \pm B \quad \text{mit} \quad C[i, j] := A[i, j] \pm B[i, j] \end{aligned}$$

mit `a`, `b`, `c` vom Typ `cvector` und `A`, `B`, `C` vom Typ `cmatrix`, sind die Operationen `*` und `/` definiert durch

$$\begin{aligned} s & := a * b \quad \text{mit} \quad s := \#* (\text{for } i := \text{lb}(a) \text{ to } \text{ub}(a) \\ & \quad \quad \quad \text{sum } (a[i] * b[i])) \quad \ddagger \\ c & := r * a \quad \text{mit} \quad c[i] := r * a[i] \\ c & := a * r \quad \text{mit} \quad c[i] := a[i] * r \\ c & := a / r \quad \text{mit} \quad c[i] := a[i] / r \\ c & := A * b \quad \text{mit} \quad c[i] := A[i] * b \quad \ddagger \\ C & := r * A \quad \text{mit} \quad C[i, j] := r * A[i, j] \\ C & := A * r \quad \text{mit} \quad C[i, j] := A[i, j] * r \\ C & := A / r \quad \text{mit} \quad C[i, j] := A[i, j] / r \\ C & := A * B \quad \text{mit} \quad C[i, j] := A[i] * B[*, j] \quad \ddagger \end{aligned}$$


---

<sup>‡</sup>Maximalgenaues Skalarprodukt mit Ergebnistyp `real`



wobei  $r, s$  vom Typ *complex*,  $a, b, c$  vom Typ *cvector* und  $A, B, C$  vom Typ *cmatrix* sind. Die Operationen mit gerichteten Rundungen sind entsprechend definiert.

Die Vergleichsoperatoren  $=, <>, <, <=, >, >=$  wirken auf alle Komponenten eines Vektors oder einer Matrix. Für  $a, b$  vom Typ *cvector* und  $A, B$  vom Typ *cmatrix* gilt:

$$\begin{aligned} a <= b &\iff a[i] <= b[i] \quad \text{für alle } i \\ A <= B &\iff A[i, j] <= B[i, j] \quad \text{für alle } i, j \end{aligned}$$

Die Operatoren auf der rechten Seite der Äquivalenzbeziehungen sind dabei die Operatoren für Operanden vom Typ *complex*. In nachfolgender Tabelle sind alle Operatoren des Moduls MVC\_ARI zusammengestellt:

rechter Operand	integer real	complex	rvector	cvector	rmatrix	cmatrix
linker Operand						
monadisch				+, -		+, -
integer real				*, *<, *>		*, *<, *>
complex			*, *<, *>	*, *<, *>	*, *<, *>	*, *<, *>
rvector		*, *<, *> /, /<, />		o v		
cvector	*, *<, *> /, /<, />	*, *<, *> /, /<, />	o v	o v		
rmatrix		*, *<, *> /, /<, />		*, *<, *>		o v
cmatrix	*, *<, *> /, /<, />	*, *<, *> /, /<, />	*, *<, *>	*, *<, *>	o v	o v

Tabelle 7.14: Operatoren des Moduls MVC\_ARI

$$\begin{aligned} o &\in \{+, +<, +>, -, -<, ->, *, *<, *>\} \\ v &\in \{=, <>, <, <=, >, >=\} \end{aligned}$$

### 7.6.3 Transferfunktionen

Zur Wandlung zwischen den Typen *rvector* und *cvector* werden folgende Transferfunktionen bereitgestellt:

Funktion	Ergb.-Typ	Bedeutung
<code>compl (rv1,rv2)</code>	<i>cvector</i>	Komplexer Vektor <i>cv</i> mit $cv[i] = compl (rv1[i], rv2[i])$
<code>compl (rv)</code>	<i>cvector</i>	Rein reeller komplexer Vektor <i>cv</i> mit $cv[i] = compl (rv[i])$
<code>re (cv)</code>	<i>rvector</i>	Realteilvektor <i>rv</i> , mit $rv[i] = re (cv[i])$
<code>im (cv)</code>	<i>rvector</i>	Imaginärteilvektor <i>rv</i> , mit $rv[i] = im (cv[i])$

$rv, rv1, rv2 = rvector$ -Ausdruck,  $cv = cvector$ -Ausdruck

Zur Wandlung zwischen den Typen *rmatrix* und *cmatrix* werden folgende Transferfunktionen bereitgestellt:

Funktion	Ergb.-Typ	Bedeutung
<code>compl (rM1,rM2)</code>	<i>cmatrix</i>	Komplexe Matrix <i>cM</i> mit $cM[i, j] = compl (rM1[i, j], rM2[i, j])$
<code>compl (rM)</code>	<i>cmatrix</i>	Rein reelle komplexe Matrix <i>cM</i> mit $cM[i, j] = compl (rM[i, j])$
<code>re (cM)</code>	<i>rmatrix</i>	Realteilmatrix <i>rM</i> mit $rM[i, j] = re (cM[i, j])$
<code>im (cM)</code>	<i>rmatrix</i>	Imaginärteilmatrix <i>rM</i> mit $rM[i, j] = im (cM[i, j])$

$rM, rM1, rM2 = rmatrix$ -Ausdruck,  $cM = cmatrix$ -Ausdruck

### 7.6.4 Überladungen des Zuweisungsoperators

Die komponentenweise Initialisierung von *cvector*- und *cmatrix*-Variablen sowie die Wandlungen von *rvector* nach *cvector* und *rmatrix* nach *cmatrix* werden in Form überladener Zuweisungen bereitgestellt:

Zuweisung	Bedeutung
$cv := r$	$cv[j] := \text{compl}(r)$ $j = \text{lb}(cv), \dots, \text{ub}(cv)$
$cv := c$	$cv[j] := c$ $j = \text{lb}(cv), \dots, \text{ub}(cv)$
$cv := rv$	$cv := \text{compl}(rv)$
$cM := r$	$cM[j, k] := \text{compl}(r)$ $j = \text{lb}(cM,1), \dots, \text{ub}(cM,1)$ $k = \text{lb}(cM,2), \dots, \text{ub}(cM,2)$
$cM := c$	$cM[j, k] := c$ $j = \text{lb}(cM,1), \dots, \text{ub}(cM,1)$ $k = \text{lb}(cM,2), \dots, \text{ub}(cM,2)$
$cM := rM$	$cM := \text{compl}(rM)$

$c = \text{complex-Ausdruck}$ ,  $cv = \text{cvector-Variable}$   
 $cM = \text{cmatrix-Variable}$ ,  $r = \text{real-Ausdruck}$   
 $rv = \text{rvector-Ausdruck}$ ,  $rM = \text{rmatrix-Ausdruck}$

### 7.6.5 Standardfunktionen

Es stehen die Funktionen *id* und *null* zur Erzeugung einer Einheitsmatrix und einer Nullmatrix bzw. eines Nullvektors, die Funktion *conj* für die Konjugation von Vektoren und Matrizen sowie die Funktionen *transp* und *herm* zur Berechnung der transponierten und der hermiteschen Matrix zur Verfügung.

Funktion	Ergb.-Typ	Bedeutung
$\text{null}(cv)$	<i>rvector</i>	Nullvektor mit dem aktuellen Indexbereich von $cv$
$\text{null}(cM)$	<i>rmatrix</i>	Nullmatrix mit den aktuellen Indexbereichen von $cM$
$\text{null}(cM1, cM2)$	<i>rmatrix</i>	Nullmatrix mit den aktuellen Indexbereichen der Produktmatrix $cM1 \cdot cM2$
$\text{id}(cM)$	<i>rmatrix</i>	Einheitsmatrix mit den aktuellen Indexbereichen von $cM$
$\text{id}(cM1, cM2)$	<i>rmatrix</i>	Einheitsmatrix mit den aktuellen Indexbereichen der Produktmatrix $cM1 \cdot cM2$

$cv = \text{cvector-Ausdruck}$ ,  $cM, cM1, cM2 = \text{cmatrix-Ausdruck}$

Funktion	Ergb.-Typ	Bedeutung
conj (cv)	<i>cvector</i>	Konjugiert komplexer Vektor <i>cvk</i> mit $cvk[i] = conj(cv[i])$
conj (cM)	<i>cmatrix</i>	Konjugiert komplexe Matrix <i>cMk</i> mit $cMk[i, j] = conj(cM[i, j])$
transp (cM)	<i>cmatrix</i>	Transponierte Matrix <i>cMt</i> von <i>cM</i> mit $cMt[i, j] = cM[j, i]$
herm (cM)	<i>cmatrix</i>	Hermitesche Matrix <i>cMh</i> von <i>cM</i> mit $cMh[i, j] = conj(cM[j, i])$

*cv* = *cvector*-Ausdruck, *cM* = *cmatrix*-Ausdruck

### Beispiel:

Für die komplexen Matrizen *M*, *M1*, *M2* vom Typ *cmatrix* ergibt, nach Ausführung der Anweisungen

```
M1 := conj( transp( M ) );
M2 := herm( M );
```

der logische Ausdruck

```
M1 = M2
```

den Wert *true*.

### 7.6.6 Ein-/AusgabeprozEDUREN

Es stehen die Prozeduren

```
procedure read (var f : text; var a : cvector);
procedure read (var f : text; var A : cmatrix);
procedure write (var f : text; a : cvector);
procedure write (var f : text; A : cmatrix);
```

mit optionalem Fileparameter, beliebig vielen Ein-/Ausgabeparametern, jedoch ohne Formatspezifikationen zur Verfügung.

Die Eingabe eines komplexen Vektors bzw. einer komplexen Matrix erfolgt komponentenweise entsprechend der Eingabe von *complex*-Größen, dabei wird eine Matrix stets zeilenweise eingelesen. Die Ausgabe eines komplexen Vektors bzw. einer komplexen Matrix erfolgt wie bei der Eingabe komponentenweise mit dem Standardformat für *real*-Größen aus Tabelle 3.4.

## 7.7 Das Modul LMVC\_ARI

In diesem Modul werden die für das Rechnen mit komplexen Vektoren und Matrizen notwendigen Operatoren, Funktionen und Prozeduren auf der Grundlage des Datentyps `real` bereitgestellt.

### 7.7.1 Datentypen

Die dynamischen Datentypen zur Darstellung von komplexen Vektoren und Matrizen sind definiert durch:

```

type Lcvector = dynamic array [*] of Lcomplex;
          Lcmatrix = dynamic array [*] of Lcvector;

```

Die aktuellen Indexgrenzen werden erst bei der Vereinbarung von Variablen dieser Typen festgelegt.

### 7.7.2 Operatoren

Viele der bekannten grundlegenden komplexen Matrix/Vektor-Operationen sind in diesem Modul vordefiniert. Als arithmetische Operatoren stehen die monadischen Operatoren `+`, `-` und die vier Grundoperationen `+`, `-`, `*`, `/` mit den drei verschiedenen Rundungsarten zur Verfügung, wobei jedoch nur bestimmte Operandenkombinationen zulässig sind. Während die Operationen `+` und `-` für Vektoren und Matrizen komponentenweise erklärt sind gemäß

$$\begin{aligned}
 c &:= a \pm b & \text{mit} & & c[i] &:= a[i] \pm b[i] \\
 C &:= A \pm B & \text{mit} & & C[i, j] &:= A[i, j] \pm B[i, j]
 \end{aligned}$$

mit  $a, b, c$  vom Typ *Lcvector* und  $A, B, C$  vom Typ *Lcmatrix*, sind die Operationen `*` und `/` definiert durch

$$\begin{aligned}
 s &:= a * b & \text{mit} & & s &:= \sum_{i=lb(a)}^{ub(a)} a[i] \cdot b[i] & \ddagger \\
 c &:= r * a & \text{mit} & & c[i] &:= r * a[i] \\
 c &:= a * r & \text{mit} & & c[i] &:= a[i] * r \\
 c &:= a / r & \text{mit} & & c[i] &:= a[i] / r \\
 c &:= A * b & \text{mit} & & c[i] &:= A[i] * b & \ddagger \\
 C &:= r * A & \text{mit} & & C[i, j] &:= r * A[i, j] \\
 C &:= A * r & \text{mit} & & C[i, j] &:= A[i, j] * r \\
 C &:= A / r & \text{mit} & & C[i, j] &:= A[i, j] / r \\
 C &:= A * B & \text{mit} & & C[i, j] &:= A[i] * B[*, j] & \ddagger
 \end{aligned}$$


---

$\ddagger$ Maximalgenaues Skalarprodukt mit Rundung zur nächsten *Longreal*-Zahl

wobei  $r, s$  vom Typ  $Lcomplex$ ,  $a, b, c$  vom Typ  $Lvector$  und  $A, B, C$  vom Typ  $Lmatrix$  sind. Die Operationen mit gerichteten Rundungen sind entsprechend definiert.

Die Vergleichsoperatoren  $=, <>, <, <=, >, >=$  wirken auf alle Komponenten eines Vektors oder einer Matrix. Für  $a, b$  vom Typ  $Lvector$  und  $A, B$  vom Typ  $Lmatrix$  gilt:

$$\begin{aligned} a <= b &\iff a[i] <= b[i] \quad \text{für alle } i \\ A <= B &\iff A[i, j] <= B[i, j] \quad \text{für alle } i, j \end{aligned}$$

Die Operatoren auf der rechten Seite der Äquivalenzbeziehungen sind dabei die Operatoren für Operanden vom Typ  $Lcomplex$ . In nachfolgender Tabelle sind alle Operatoren des Moduls LMVC\_ARI zusammengestellt:

rechter Operand	Longreal	Lcomplex	Lrvector	Lcvector	Lrmatrix	Lcmatrix
linker Operand						
monadisch				+, -		+, -
Longreal				*, *<, *>		*, *<, *>
Lcomplex			*, *<, *>	*, *<, *>	*, *<, *>	*, *<, *>
Lrvector		*, *<, *> /, /<, />		o v		
Lcvector	*, *<, *> /, /<, />	*, *<, *> /, /<, />	o v	o v		
Lrmatrix		*, *<, *> /, /<, />		*, *<, *>		o v
Lcmatrix	*, *<, *> /, /<, />	*, *<, *> /, /<, />	*, *<, *>	*, *<, *>	o v	o v

Tabelle 7.15: Operatoren des Moduls LMVC\_ARI

$$\begin{aligned} o &\in \{+, +<, +>, -, -<, ->, *, *<, *>\} \\ v &\in \{=, <>, <, <=, >, >=\} \end{aligned}$$

### 7.7.3 Transferfunktionen

Zur Wandlung zwischen den Typen *Lrvector* und *Lcvector* werden folgende Transferfunktionen bereitgestellt:

Funktion	Ergb.-Typ	Bedeutung
compl (Lv1,Lv2)	<i>Lcvector</i>	Komplexer Vektor <i>Lcv</i> mit $Lcv[i] = lcompl(Lv1[i], Lv2[i])$
compl (Lv)	<i>Lcvector</i>	Rein reeller komplexer Vektor <i>Lcv</i> mit $Lcv[i] = lcompl(Lv[i])$
re (Lcv)	<i>Lrvector</i>	Realteilvektor <i>Lv</i> , mit $Lv[i] = re(Lcv[i])$
im (Lcv)	<i>Lrvector</i>	Imag.-teilvektor <i>Lv</i> , mit $Lv[i] = im(Lcv[i])$

$Lv, Lv1, Lv2 = Lrvector$ -Ausdruck,  $Lcv = Lcvector$ -Ausdruck

Zur Wandlung zwischen den Typen *Lrmatrix* und *Lcmatrix* werden folgende Transferfunktionen bereitgestellt:

Funktion	Ergb.-Typ	Bedeutung
compl (LM1,LM2)	<i>Lcmatrix</i>	Komplexe Matrix <i>LcM</i> mit $LcM[i, j] = lcompl(LM1[i, j], LM2[i, j])$
compl (LM)	<i>Lcmatrix</i>	Rein reelle komplexe Matrix <i>LcM</i> mit $LcM[i, j] = lcompl(LM[i, j])$
re (LcM)	<i>Lrmatrix</i>	Realteilmatrix <i>LM</i> mit $LM[i, j] = re(LcM[i, j])$
im (LcM)	<i>Lrmatrix</i>	Imaginärteilmatrix <i>LM</i> mit $LM[i, j] = im(LcM[i, j])$

$LM, LM1, LM2 = Lrmatrix$ -Ausdruck,  $LcM = Lcmatrix$ -Ausdruck

Während dem Anwender im Modul **LC\_ARI** alle gemischten Operationen mit den Datentypen *complex* und *Lcomplex* zur Verfügung stehen, sind im Modul **LMVC\_ARI** die gemischten Operationen z.B. mit den Datentypen *cvector* und *Lcvector* nicht definiert.

Mit Hilfe der Funktionen **SHORT, SHORTUP, SHORNDOWN** lassen sich jedoch Vektoren bzw. Matrizen vom Typ *Lcvector* bzw. *Lcmatrix* durch entsprechende Rundungen umwandeln in Vektoren bzw. Matrizen vom Typ *cvector* bzw. *cmatrix*. Mit den entsprechenden ... \_TEST-Versionen dieser Funktionen kann überprüft werden, ob durch Rundungen in den denormalisierten Zahlenbereich bei wenigstens einer Komponente hintere Mantissenstellen auf Null gesetzt werden.

Durch Überladungen des Zuweisungsoperators können umgekehrt Vektoren und Matrizen vom Typ *cvector* bzw. *cmatrix* rundungsfehlerfrei umgewandelt werden in Vektoren und Matrizen vom Typ *Lcvector* bzw. *Lcmatrix*.

Funktion	Erg.-Typ	Bedeutung
short (Lcv)	<i>cvector</i>	$cv[i] := \text{short}(Lcv[i])$
shortup (Lcv)	<i>cvector</i>	$cv[i] := \text{shortup}(Lcv[i])$
shortdown (Lcv)	<i>cvector</i>	$cv[i] := \text{shortdown}(Lcv[i])$
short_test (Lcv,err)	<i>cvector</i>	$cv[i] := \text{short\_test}(Lcv[i], err)$
shortup_test (Lcv,err)	<i>cvector</i>	$cv[i] := \text{shortup\_test}(Lcv[i], err)$
shortdown_test (Lcv,err)	<i>cvector</i>	$cv[i] := \text{shortdown\_test}(Lcv[i], err)$
short (LcM)	<i>cmatrix</i>	$cM[i] := \text{short}(LcM[i])$
shortup (LcM)	<i>cmatrix</i>	$cM[i] := \text{shortup}(LcM[i])$
shortdown (LcM)	<i>cmatrix</i>	$cM[i] := \text{shortdown}(LcM[i])$
short_test (LcM,err)	<i>cmatrix</i>	$cM[i] := \text{short\_test}(LcM[i], err)$
shortup_test (LcM,err)	<i>cmatrix</i>	$cM[i] := \text{shortup\_test}(LcM[i], err)$
shortdown_test (LcM,err)	<i>cmatrix</i>	$cM[i] := \text{shortdown\_test}(LcM[i], err)$

Lcv = *Lcvector*-Ausdruck, cv = *cvector*-**Variable**, err = *boolean*-**Variable**  
 LcM = *Lcmatrix*-Ausdruck; cM = *cmatrix*-**Variable**

#### 7.7.4 Überladungen des Zuweisungsoperators

Die komponentenweise Initialisierung von *Lcvector*- und *Lcmatrix*-Variablen sowie die Wandlungen von *Lrvector* nach *Lcvector* und *Lrmatrix* nach *Lcmatrix* werden in Form überladener Zuweisungen bereitgestellt:

Zuweisung	Bedeutung
Lcv := Lr	$Lcv[j] := \text{lcompl}(Lr) \quad j = \text{lb}(Lcv), \dots, \text{ub}(Lcv)$
Lcv := Lc	$Lcv[j] := Lc \quad j = \text{lb}(Lcv), \dots, \text{ub}(Lcv)$
Lcv := Lv	$Lcv := \text{lcompl}(Lv)$

Lc = *Lcomplex*-Ausdruck, Lcv = *Lcvector*-**Variable**  
 Lv = *Lrvector*-Ausdruck, Lr = *Longreal*-Ausdruck



Zuweisung	Bedeutung
$LcM := Lr$	$LcM[j, k] := lcompl(Lr)$ $j = lb(LcM,1), \dots, ub(LcM,1)$ $k = lb(LcM,2), \dots, ub(LcM,2)$
$LcM := Lc$	$LcM[j, k] := Lc$ $j = lb(LcM,1), \dots, ub(LcM,1)$ $k = lb(LcM,2), \dots, ub(LcM,2)$
$LcM := LM$	$LcM := compl(LM)$
$Lcv := cv$	$Lcv[j] = cv[j]$ $j = lb(Lcv), \dots, ub(Lcv)$
$LcM := cM$	$LcM[j, k] = cM[j, k]$ $j = lb(LcM,1), \dots, ub(LcM,1)$ $k = lb(LcM,2), \dots, ub(LcM,2)$

$Lc = Lcomplex$ -Ausdruck,  $Lr = Longreal$ -Ausdruck  
 $cv = cvector$ -Ausdruck,  $Lcv = Lcvector$ -**Variable**  
 $cM = cmatrix$ -Ausdruck,  $LcM = Lcmatrix$ -**Variable**,  
 $LM = Lmatrix$ -Ausdruck

### 7.7.5 Standardfunktionen

Es stehen die Funktionen *id* und *null* zur Erzeugung einer Einheitsmatrix und einer Nullmatrix bzw. eines Nullvektors, die Funktion *conj* für die Konjugation von Vektoren und Matrizen sowie die Funktionen *transp* und *herm* zur Berechnung der transponierten und der hermiteschen Matrix zur Verfügung.

Funktion	Ergb.-Typ	Bedeutung
$null(Lcv)$	<i>Lrvector</i>	Nullvektor mit dem aktuellen Indexbereich von <i>Lcv</i>
$null(LcM)$	<i>Lrmatrix</i>	Nullmatrix mit den aktuellen Indexbereichen von <i>LcM</i>
$null(LcM1, LcM2)$	<i>Lrmatrix</i>	Nullmatrix mit den aktuellen Indexbereichen der Produktmatrix $LcM1 \cdot LcM2$
$id(LcM)$	<i>Lrmatrix</i>	Einheitsmatrix mit den aktuellen Indexbereichen von <i>LcM</i>
$id(LcM1, LcM2)$	<i>Lrmatrix</i>	Einheitsmatrix mit den aktuellen Indexbereichen der Produktmatrix $LcM1 \cdot LcM2$
$conj(Lcv)$	<i>Lcvector</i>	Konjugiert komplexer Vektor <i>cvk</i> mit $cvk[i] = conj(Lcv[i])$

$Lcv = Lcvector$ -Ausdruck,  $LcM, LcM1, LcM2 = Lcmatrix$ -Ausdruck

Funktion	Ergb.-Typ	Bedeutung
conj (LcM)	Lcmatrix	Konjugiert komplexe Matrix $lcMk$ mit $lcMk[i, j] = conj(LcM[i, j])$
transp (LcM)	Lcmatrix	Transponierte Matrix $lcMt$ von $LcM$ mit $lcMt[i, j] = LcM[j, i]$
herm (LcM)	Lcmatrix	Hermitesche Matrix $lcMh$ von $LcM$ mit $lcMh[i, j] = conj(LcM[j, i])$

LcM = Lcmatrix-Ausdruck

### Beispiel:

Für die komplexen Matrizen  $M, M1, M2$  vom Typ *Lcmatrix* ergibt, nach Ausführung der Anweisungen

```
M1 := conj( transp ( M ) );
M2 := herm( M );
```

der logische Ausdruck

```
M1 = M2
```

den Wert *true*.

## 7.7.6 Ein-/Ausgabeprozeduren

Es stehen die Prozeduren

```
procedure read (var f : text; var a : Lcvector);
procedure read (var f : text; var A : Lcmatrix);
procedure write (var f : text; a : Lcvector);
procedure write (var f : text; A : Lcmatrix);
```

mit optionalem Fileparameter, beliebig vielen Ein-/Ausgabeparametern, jedoch ohne Formatspezifikationen zur Verfügung.

Die Eingabe eines komplexen Vektors bzw. einer komplexen Matrix erfolgt komponentenweise entsprechend der Eingabe von *Lcomplex*-Größen, dabei wird eine Matrix stets zeilenweise eingelesen. Die Ausgabe eines komplexen Vektors bzw. einer komplexen Matrix erfolgt wie bei der Eingabe komponentenweise mit dem Standardformat für *Longreal*-Größen aus Tabelle 3.4.

### 7.7.7 Exakte Auswertung von Ausdrücken

Benutzt man das Modul `MVC_ARI`, so lassen sich *cvector*- und *cmatrix*-Ausdrücke mit Hilfe des Lattenkreuz-Konzepts maximalgenau auswerten. Da jedoch #-Ausdrücke z.B. nicht mit *Lcmatrix*-Operanden verträglich sind, wird hier an einem Programmbeispiel gezeigt, wie man mit den *Lcmatrix*-Operanden  $A, B, C$  z.B. den Ausdruck  $A + B \cdot C$  ebenfalls maximalgenau berechnen kann:

```

program beispiel (input,output);
use lc_ari, lmvc_ari;
var cdot      : cdotprecision;
    A,B,C,D   : Lcmatrix[1..5,1..5];
    i,k       : integer;

begin { Haupt }
  writeln('A = ?'); read(A);           { Die Eingabe einer Matrix }
  writeln('B = ?'); read(B);           { erfolgt stets zeilenweise }
  writeln('C = ?'); read(C);
  For i := lb(A,1) to ub(A,1) do
    for k := lb(A,2) to ub(A,2) do
      begin
        cdot := A[i,k];
        PADDNACCU( cdot,LCVECTOR(B[i,*]),LCVECTOR(C[*],k) );
        D[i,k] := LONGNEXT(cdot)
      end;
    writeln(D);                        { Maximalgen. Ergebnis in D }
  end.

```

#### Anmerkungen:

1. Eine Matrix wird stets zeilenweise eingegeben, wobei jede Zeile durch ein `< return >` abzuschließen ist.
2. Die Konstruktion `LCVECTOR(B[i,*])` liefert dem Compiler die Information, den  $i$ -ten Zeilenvektor  $B[i,*]$  als Vektor vom Typ *Lcvector* zu behandeln.
3. Die Konstruktion `LCVECTOR(B[*],k)` liefert dem Compiler die Information, den  $k$ -ten Spaltenvektor  $B[*],k$  als Vektor vom Typ *Lcvector* zu behandeln.
4. Die Prozedur `PADDNACCU(..)` aus dem Modul `LC_ARI` addiert zur Variablen `cdot = A[i,k]` das Produkt  $B[i,*] \cdot C[*],k$  rundungsfehlerfrei; vergl. Seite 119, Nr. 12.
5. Die Funktion `LONGNEXT(cdot)` rundet den exakten `cdot`-Wert zur nächsten *Lcomplex*-Zahl  $D[i,k]$ ; vergl. Seite 120.

## 7.8 Das Modul MVCL\_ARI

In diesem Modul werden die für das Rechnen mit komplexen Intervallvektoren und Intervallmatrizen notwendigen Operatoren, Funktionen und Prozeduren auf der Grundlage des Datentyps `real` bereitgestellt.

### 7.8.1 Datentypen

Die dynamischen Datentypen zur Darstellung von komplexen Intervallvektoren und Intervallmatrizen sind entsprechend der Vereinbarung

```
type civector = dynamic array [*] of cinterval;
      cimatrix = dynamic array [*] of civector;
```

im Sprachkern von PASCAL-XSC enthalten. Die aktuellen Indexgrenzen werden erst bei der Vereinbarung von Variablen dieser Typen festgelegt.

### 7.8.2 Operatoren

Viele der bekannten grundlegenden komplexen, intervallmäßigen Matrix/Vektor-Operationen sind in diesem Modul vordefiniert. Als arithmetische Operatoren stehen die monadischen Operatoren `+`, `-` und die vier Grundoperationen `+`, `-`, `*`, `/` mit der komponentenweisen Rundung zum kleinsten einschließenden komplexen Intervall (auch für gemischte Typen) zur Verfügung, wobei jedoch nur bestimmte Operandenkombinationen zulässig sind. Während die Operationen `+` und `-` für Intervallvektoren und Intervallmatrizen komponentenweise erklärt sind gemäß

$$\begin{aligned} c & := a \pm b & \text{mit} & & c[i] & := a[i] \pm b[i] \\ C & := A \pm B & \text{mit} & & C[i, j] & := A[i, j] \pm B[i, j] \end{aligned}$$

mit  $a, b, c$  vom Typ *civector* und  $A, B, C$  vom Typ *cimatrix*, sind die Operationen `*` und `/` definiert durch

$$\begin{aligned} s & := a * b & \text{mit} & & s & := \#\# \left( \text{for } i := \text{lb}(a) \text{ to } \text{ub}(a) \right. \\ & & & & & \left. \text{sum } (a[i] * b[i]) \right) \ddagger \\ c & := r * a & \text{mit} & & c[i] & := r * a[i] \\ c & := a * r & \text{mit} & & c[i] & := a[i] * r \\ c & := a / r & \text{mit} & & c[i] & := a[i] / r \\ c & := A * b & \text{mit} & & c[i] & := A[i] * b \ddagger \\ C & := r * A & \text{mit} & & C[i, j] & := r * A[i, j] \\ C & := A * r & \text{mit} & & C[i, j] & := A[i, j] * r \\ C & := A / r & \text{mit} & & C[i, j] & := A[i, j] / r \\ C & := A * B & \text{mit} & & C[i, j] & := A[i] * B[*, j] \ddagger \end{aligned}$$


---

$\ddagger$ Maximalgenaues Skalarprodukt mit Ergebnistyp *real*

wobei  $r, s$  vom Typ *cinterval*,  $a, b, c$  vom Typ *civector* und  $A, B, C$  vom Typ *cimatrix* sind. Die Operationen mit gemischten Operandentypen sind entsprechend definiert.

Die wie im Falle von Intervallen und komplexen Intervallen mengentheoretisch zu verstehenden Vergleichsoperatoren  $=, <>, <, <=, >, >=$  sind auf der Basis von  $=$  und  $<=$  realisiert, wobei für  $a, b$  vom Typ *civector* und  $A, B$  vom Typ *cimatrix* gilt:

$$\begin{aligned} a <= b &\iff a[i] <= b[i] \quad \text{für alle } i \\ A <= B &\iff A[i, j] <= B[i, j] \quad \text{für alle } i, j \end{aligned}$$

Die Operatoren auf der rechten Seite der Äquivalenzbeziehungen sind dabei die Operatoren für Werte vom Typ *cinterval*.

Zusätzlich stehen die Operatoren **in** für die Relation „liegt in“ zwischen einem *rvector*- oder *cvector*- und einem *civector*-Operanden bzw. zwischen einem *rmatrix*- oder *cmatrix*- und einem *cimatrix*-Operanden sowie für die Relation „enthalten im Innern“ zwischen zwei *civector*-Operanden bzw. zwei *cimatrix*-Operanden zur Verfügung.  $><$  für den Test auf Disjunktheit zweier komplexer Intervallvektoren bzw. Intervallmatrizen steht ebenfalls zur Verfügung. Diese Operatoren sind jeweils komponentenweise definiert.

Die Verbandsoperatoren  $+*$  und  $**$  bezeichnen die komponentenweise Bildung der Intervall-Hülle bzw. des Durchschnitts, wie es bereits für den Typ *cinterval* im Modul CLARI beschrieben wurde.

Die Übersicht über alle im Modul MVCLARI definierten Operatoren wird, bedingt durch die große Zahl von Operatoren, in den beiden nachfolgenden Tabellen gegeben:

rechter Operand	integer real	complex	interval	cinterval	rvector	cvector
linker Operand						
monadisch						
integer real						
complex						
interval						*
cinterval					*	*
rvector				*, /		+*
cvector			*, /	*, /	+*	+*
ivector		*, /		*, /		◇, =, <>, +*
civector	*, /	*, /	*, /	*, /	◇, =, <>, +*	◇, =, <>, +*
rmatrix				*, /		
cmatrix			*, /	*, /		
imatrix		*, /		*, /		*
cimatrix	*, /	*, /	*, /	*, /	*	*

Tabelle 7.16: Operatoren des Moduls MVCL\_ARI, Teil 1

$$\diamond \in \{+, -, *\} \quad \forall \in \{=, <>, <, <=, >, >=\}$$

rechter Operand	ivector	civector	rmatrix	cmatrix	imatrix	cimatrix
linker Operand						
monadisch		+, -				+, -
integer real		*				*
complex	*	*			*	*
interval		*		*		*
cinterval	*	*	*	*	*	*
rvector		$\diamond,$ $=, <>, \mathbf{in}$ $+*$				
cvector	$\diamond,$ $=, <>, \mathbf{in}$ $+*$	$\diamond,$ $=, <>, \mathbf{in}$ $+*$				
ivector		$\diamond,$ $\mathbf{in}, V, ><$ $+*, **$				
civector	$\diamond,$ $V, ><$ $+*, **$	$\diamond,$ $\mathbf{in}, V, ><$ $+*, **$				
rmatrix		*		$+*$		$\diamond,$ $=, <>, \mathbf{in}$ $+*$
cmatrix	*	*	$+*$	$+*$	$\diamond,$ $=, <>, \mathbf{in}$ $+*$	$\diamond,$ $=, <>, \mathbf{in}$ $+*$
imatrix		*		$\diamond,$ $=, <>,$ $+*$		$\diamond,$ $\mathbf{in}V, ><,$ $+*, **$
cimatrix	*	*	$\diamond,$ $=, <>$ $+*$	$\diamond,$ $=, <>,$ $+*$	$\diamond,$ $V, ><,$ $+*, **$	$\diamond,$ $\mathbf{in}V, ><,$ $+*, **$

Tabelle 7.17: Operatoren des Moduls MVCLARI, Teil 2

$$\diamond \in \{+, -, *\} \quad \forall \in \{=, <>, <, <=, >, >=\}$$

### 7.8.3 Transferfunktionen für komplexe Intervallvektoren

Zur Wandlung zwischen den Typen *rvector*, *cvector*, *ivector* und *civector* werden folgende Transferfunktionen bereitgestellt:

Funktion	Ergb.-Typ	Bedeutung
$\text{compl}(iv1, iv2)$	<i>civector</i>	Komplexer Intervallvektor <i>civ</i> mit $civ[i] = \text{compl}(iv1[i], iv2[i])$
$\text{compl}(rv, iv)$	<i>civector</i>	Komplexer Intervallvektor <i>civ</i> mit $civ[i] = \text{compl}(rv[i], iv[i])$
$\text{compl}(iv, rv)$	<i>civector</i>	Komplexer Intervallvektor <i>civ</i> mit $civ[i] = \text{compl}(iv[i], rv[i])$
$\text{compl}(iv)$	<i>civector</i>	Rein reeller komplexer Intervallvektor <i>civ</i> mit $civ[i] = \text{compl}(iv[i])$
$\text{intval}(cv1, cv2)$	<i>civector</i>	Komplexer Intervallvektor <i>civ</i> mit $civ[i] = \text{intval}(cv1[i], cv2[i])$
$\text{intval}(rv, cv)$	<i>civector</i>	Komplexer Intervallvektor <i>civ</i> mit $civ[i] = \text{intval}(rv[i], cv[i])$
$\text{intval}(cv, rv)$	<i>civector</i>	Komplexer Intervallvektor <i>civ</i> mit $civ[i] = \text{intval}(cv[i], rv[i])$
$\text{intval}(cv)$	<i>civector</i>	Punktintervallvektor <i>civ</i> mit $civ[i] = \text{intval}(cv[i])$
$\text{re}(civ)$	<i>ivector</i>	Realteilvektor <i>iv</i> , mit $iv[i] = \text{re}(civ[i])$
$\text{im}(civ)$	<i>ivector</i>	Imaginärteilvektor <i>iv</i> , mit $iv[i] = \text{im}(civ[i])$
$\text{inf}(civ)$	<i>cvector</i>	Komplexer Vektor <i>cv</i> der Untergrenzen mit $cv[i] = \text{inf}(civ[i])$
$\text{sup}(civ)$	<i>cvector</i>	Komplexer Vektor <i>cv</i> der Obergrenzen mit $cv[i] = \text{sup}(civ[i])$

$rv = rvector$ -Ausdruck,  $cv, cv1, cv2 = cvector$ -Ausdruck,  
 $iv, iv1, iv2 = ivector$ -Ausdruck,  $civ = civector$ -Ausdruck



### 7.8.4 Transferfunktionen für komplexe Intervallmatrizen

Zur Wandlung zwischen den Typen *rmatrix*, *cmatrix*, *imatrix* und *cimatrix* werden folgende Transferfunktionen bereitgestellt:

Funktion	Ergb.-Typ	Bedeutung
<code>compl (iM1,iM2)</code>	<i>cimatrix</i>	Komplexe Intervallmatrix <i>ciM</i> mit $ciM[i, j] = compl(iM1[i, j], iM2[i, j])$
<code>compl (rM,iM)</code>	<i>cimatrix</i>	Komplexe Intervallmatrix <i>ciM</i> mit $ciM[i, j] = compl(rM[i, j], iM[i, j])$
<code>compl (iM,rM)</code>	<i>cimatrix</i>	Komplexe Intervallmatrix <i>ciM</i> mit $ciM[i, j] = compl(iM[i, j], rM[i, j])$
<code>compl (iM)</code>	<i>cimatrix</i>	Rein reelle komplexe Intervallmatrix <i>ciM</i> mit $ciM[i, j] = compl(iM[i, j])$
<code>intval (cM1,cM2)</code>	<i>cimatrix</i>	Komplexe Intervallmatrix <i>ciM</i> mit $ciM[i, j] = intval(cM1[i, j], cM2[i, j])$
<code>intval (rM,cM)</code>	<i>cimatrix</i>	Komplexe Intervallmatrix <i>ciM</i> mit $ciM[i, j] = intval(rM[i, j], cM[i, j])$
<code>intval (cM,rM)</code>	<i>cimatrix</i>	Komplexe Intervallmatrix <i>ciM</i> mit $ciM[i, j] = intval(cM[i, j], rM[i, j])$
<code>intval (cM)</code>	<i>cimatrix</i>	Punktintervallmatrix <i>ciM</i> mit $ciM[i, j] = intval(cM[i, j])$
<code>re (ciM)</code>	<i>imatrix</i>	Realteilmatrix <i>iM</i> mit $iM[i, j] = re(ciM[i, j])$
<code>im (ciM)</code>	<i>imatrix</i>	Imaginärteilmatrix <i>iM</i> mit $iM[i, j] = im(ciM[i, j])$
<code>inf (ciM)</code>	<i>cmatrix</i>	Komplexe Matrix <i>cM</i> der Untergrenzen mit $cM[i, j] = inf(ciM[i, j])$
<code>sup (ciM)</code>	<i>cmatrix</i>	Komplexe Matrix <i>cM</i> der Obergrenzen mit $cM[i, j] = sup(ciM[i, j])$

$rM = rmatrix$ -Ausdruck,  $cM, cM1, cM2 = cmatrix$ -Ausdruck  
 $iM, iM1, iM2 = imatrix$ -Ausdruck,  $ciM = cimatrix$ -Ausdruck

### 7.8.5 Überladungen des Zuweisungsoperators

Die komponentenweise Initialisierung von *civector*- und *cimatrix*-Variablen sowie die Wandlungen von *rvector* bzw. *cvector* bzw. *ivector* nach *civector* und *rmatrix* bzw. *cmatrix* bzw. *imatrix* nach *cimatrix* werden in Form überladener Zuweisungen bereitgestellt:

Zuweisung	Bedeutung
$\text{civ} := r$	$\text{civ}[j] := \text{compl}(\text{intval}(r)) \quad j = \text{lb}(\text{civ}), \dots, \text{ub}(\text{civ})$
$\text{civ} := c$	$\text{civ}[j] := \text{intval}(c) \quad j = \text{lb}(\text{civ}), \dots, \text{ub}(\text{civ})$
$\text{civ} := i$	$\text{civ}[j] := \text{compl}(i) \quad j = \text{lb}(\text{civ}), \dots, \text{ub}(\text{civ})$
$\text{civ} := \text{ci}$	$\text{civ}[j] := \text{ci} \quad j = \text{lb}(\text{civ}), \dots, \text{ub}(\text{civ})$
$\text{civ} := \text{rv}$	$\text{civ} := \text{compl}(\text{intval}(\text{rv}))$
$\text{civ} := \text{cv}$	$\text{civ} := \text{intval}(\text{cv})$
$\text{civ} := \text{iv}$	$\text{civ} := \text{compl}(\text{iv})$
$\text{ciM} := r$	$\text{ciM}[j, k] := \text{compl}(\text{intval}(r)) \quad j = \text{lb}(\text{ciM}, 1), \dots, \text{ub}(\text{ciM}, 1)$ $k = \text{lb}(\text{ciM}, 2), \dots, \text{ub}(\text{ciM}, 2)$
$\text{ciM} := c$	$\text{ciM}[j, k] := \text{intval}(c) \quad j = \text{lb}(\text{ciM}, 1), \dots, \text{ub}(\text{ciM}, 1)$ $k = \text{lb}(\text{ciM}, 2), \dots, \text{ub}(\text{ciM}, 2)$
$\text{ciM} := i$	$\text{ciM}[j, k] := \text{compl}(i) \quad j = \text{lb}(\text{ciM}, 1), \dots, \text{ub}(\text{ciM}, 1)$ $k = \text{lb}(\text{ciM}, 2), \dots, \text{ub}(\text{ciM}, 2)$
$\text{ciM} := \text{ci}$	$\text{ciM}[j, k] := \text{ci} \quad j = \text{lb}(\text{ciM}, 1), \dots, \text{ub}(\text{ciM}, 1)$ $k = \text{lb}(\text{ciM}, 2), \dots, \text{ub}(\text{ciM}, 2)$
$\text{ciM} := \text{rM}$	$\text{ciM} := \text{compl}(\text{intval}(\text{rM}))$
$\text{ciM} := \text{cM}$	$\text{ciM} := \text{intval}(\text{cM})$
$\text{ciM} := \text{iM}$	$\text{ciM} := \text{compl}(\text{iM})$

$\text{ci}$  = *cinterval*-Ausdruck,  $\text{civ}$  = *civector*-Variable,  $\text{ciM}$  = *cimatrix*-Variable  
 $i$  = *interval*-Ausdruck,  $\text{iv}$  = *ivector*-Ausdruck,  $\text{iM}$  = *imatrix*-Ausdruck  
 $c$  = *complex*-Ausdruck,  $\text{cv}$  = *cvector*-Ausdruck,  $\text{cM}$  = *cmatrix*-Ausdruck  
 $r$  = *real*-Ausdruck,  $\text{rv}$  = *rvector*-Ausdruck,  $\text{rM}$  = *rmatrix*-Ausdruck

### 7.8.6 Standardfunktionen

Es stehen die Funktionen *id* und *null* zur Erzeugung einer Einheitsmatrix und einer Nullmatrix bzw. eines Nullvektors, die Funktionen *mid* und *diam* für die komponentenweise Berechnung von Mittelpunkt und Durchmesser, die Funktion *conj* für die Konjugation von komplexen Intervallvektoren und Intervallmatrizen, sowie die Funktionen *transp* und *herm* zur Berechnung der transponierten und der hermiteschen Matrix zur Verfügung. Darüber hinaus wird die Funktion *blow* für die komponentenweise Berechnung der Epsilonaufblähung bereitgestellt.

Funktion	Ergb.-Typ	Bedeutung
<i>null</i> ( <i>civ</i> )	<i>rvector</i>	Nullvektor mit dem aktuellen Indexbereich von <i>civ</i>
<i>null</i> ( <i>ciM</i> )	<i>rmatrix</i>	Nullmatrix mit den aktuellen Indexbereichen von <i>ciM</i>
<i>null</i> ( <i>ciM1</i> , <i>ciM2</i> )	<i>rmatrix</i>	Nullmatrix mit den aktuellen Indexbereichen der Produktmatrix $ciM1 \cdot ciM2$
<i>id</i> ( <i>ciM</i> )	<i>rmatrix</i>	Einheitsmatrix mit den aktuellen Indexbereichen von <i>ciM</i>
<i>id</i> ( <i>ciM1</i> , <i>ciM2</i> )	<i>rmatrix</i>	Einheitsmatrix mit den aktuellen Indexbereichen der Produktmatrix $ciM1 \cdot ciM2$
<i>mid</i> ( <i>civ</i> )	<i>cvector</i>	Mittelpunktvektor <i>pv</i> mit $pv[i] = mid(civ[i])$
<i>diam</i> ( <i>civ</i> )	<i>rvector</i>	Durchmesservektor <i>dv</i> mit $dv[i] = diam(civ[i])$
<i>mid</i> ( <i>ciM</i> )	<i>cmatrix</i>	Mittelpunktmatrix <i>pM</i> mit $pM[i, j] = mid(ciM[i, j])$
<i>diam</i> ( <i>ciM</i> )	<i>rmatrix</i>	Durchmessermatrix <i>dM</i> mit $dM[i, j] = diam(ciM[i, j])$
<i>conj</i> ( <i>civ</i> )	<i>civector</i>	Konjugiert komplexer Intervallvektor <i>civk</i> mit $civk[i] = conj(civ[i])$
<i>conj</i> ( <i>ciM</i> )	<i>cimatrix</i>	Konjugiert komplexe Intervallmatrix <i>ciMk</i> mit $ciMk[i, j] = conj(ciM[i, j])$
<i>transp</i> ( <i>ciM</i> )	<i>cimatrix</i>	Transponierte Matrix <i>ciMt</i> von <i>ciM</i> mit $ciMt[i, j] = ciM[j, i]$
<i>herm</i> ( <i>ciM</i> )	<i>cimatrix</i>	Hermitesche Matrix <i>ciMh</i> von <i>ciM</i> mit $ciMh[i, j] = conj(ciM[j, i])$
<i>blow</i> ( <i>civ</i> , <i>r</i> )	<i>civector</i>	Vektorielle Epsilonaufblähung <i>ev</i> mit $ev[i] = blow(civ[i], r)$
<i>blow</i> ( <i>ciM</i> , <i>r</i> )	<i>cimatrix</i>	Vektorielle Epsilonaufblähung <i>eM</i> mit $eM[i, j] = blow(ciM[i, j], r)$

*civ* = *civector*-Ausdruck; *ciM*, *ciM1*, *ciM2* = *cimatrix*-Ausdruck, *r* = *real*-Ausdruck

**Beispiel:**

Für die komplexen Matrizen  $M$ ,  $M1$ ,  $M2$  vom Typ *cimatrix* ergibt, nach Ausführung der Anweisungen

```
M1 := conj( transp (M) );
M2 := herm( M );
```

der logische Ausdruck

```
M1 = M2
```

den Wert *true*.

**7.8.7 Ein-/Ausgabeprozeduren**

Es stehen die Prozeduren

```
procedure read (var f : text; var a : civector);
procedure read (var f : text; var A : cimatrix);
procedure write (var f : text; a : civector);
procedure write (var f : text; A : cimatrix);
```

mit optionalem Fileparameter, beliebig vielen Ein-/Ausgabeparametern, jedoch ohne Formatspezifikationen zur Verfügung.

Die Eingabe eines komplexen Intervallvektors bzw. einer komplexen Intervallmatrix erfolgt komponentenweise entsprechend der Eingabe von *cinterval*-Größen, dabei wird eine Matrix stets zeilenweise eingelesen. Die Ausgabe eines komplexen Intervallvektors bzw. einer komplexen Intervallmatrix erfolgt wie bei der Eingabe komponentenweise mit dem Standardformat für *real*-Größen aus Tabelle 3.4.

## 7.9 Das Modul LMVCL\_ARI

In diesem Modul werden die für das Rechnen mit komplexen Intervallvektoren und Intervallmatrizen notwendigen Operatoren, Funktionen und Prozeduren auf der Grundlage des Datentyps **Longreal** bereitgestellt.

### 7.9.1 Datentypen

Die dynamischen Datentypen zur Darstellung von komplexen Intervallvektoren und Intervallmatrizen sind definiert durch:

```
type Lcivector = dynamic array [*] of Lcinterval;
   Lcimatrix  = dynamic array [*] of Lcivector;
```

Die aktuellen Indexgrenzen werden erst bei der Vereinbarung von Variablen dieser Typen festgelegt.

### 7.9.2 Operatoren

Viele der bekannten grundlegenden komplexen, intervallmäßigen Matrix/Vektor-Operationen sind in diesem Modul vordefiniert, dabei sind gemischte Operanden mit den Basistypen **real** und **Longreal** nicht realisiert. Als arithmetische Operatoren stehen die monadischen Operatoren  $+$ ,  $-$  und die vier Grundoperationen  $+$ ,  $-$ ,  $*$ ,  $/$  mit der komponentenweisen Rundung zum kleinsten einschließenden komplexen Intervall zur Verfügung, wobei jedoch nur bestimmte Operandenkombinationen zulässig sind. Während die Operationen  $+$  und  $-$  für Intervallvektoren und Intervallmatrizen komponentenweise erklärt sind gemäß

$$\begin{aligned} c &:= a \pm b \quad \text{mit} \quad c[i] := a[i] \pm b[i] \\ C &:= A \pm B \quad \text{mit} \quad C[i, j] := A[i, j] \pm B[i, j] \end{aligned}$$

mit  $a, b, c$  vom Typ *Lcivector* und  $A, B, C$  vom Typ *Lcimatrix*, sind die Operationen  $*$  und  $/$  definiert durch

$$\begin{aligned} s &:= a * b \quad \text{mit} \quad s := \sum_{i=lb(a)}^{ub(a)} a[i] \cdot b[i] \quad \ddagger \\ c &:= r * a \quad \text{mit} \quad c[i] := r * a[i] \\ c &:= a * r \quad \text{mit} \quad c[i] := a[i] * r \\ c &:= a / r \quad \text{mit} \quad c[i] := a[i] / r \\ c &:= A * b \quad \text{mit} \quad c[i] := A[i] * b \quad \ddagger \\ C &:= r * A \quad \text{mit} \quad C[i, j] := r * A[i, j] \\ C &:= A * r \quad \text{mit} \quad C[i, j] := A[i, j] * r \\ C &:= A / r \quad \text{mit} \quad C[i, j] := A[i, j] / r \\ C &:= A * B \quad \text{mit} \quad C[i, j] := A[i] * B[*, j] \quad \ddagger \end{aligned}$$


---

$\ddagger$ Maximalgenaues Skalarprodukt mit Ergebnistyp *Longreal*

wobei  $r, s$  vom Typ *Lcinterval*,  $a, b, c$  vom Typ *Lcivector* und  $A, B, C$  vom Typ *Lcimatrix* sind. Die Operationen mit gemischten Operandentypen sind entsprechend definiert.

Die wie im Falle von Intervallen und komplexen Intervallen mengentheoretisch zu verstehenden Vergleichsoperatoren  $=, <>, <, <=, >, >=$  sind auf der Basis von  $=$  und  $<=$  realisiert, wobei für  $a, b$  vom Typ *Lcivector* und  $A, B$  vom Typ *Lcimatrix* gilt:

$$\begin{aligned} a <= b &\iff a[i] <= b[i] \quad \text{für alle } i \\ A <= B &\iff A[i, j] <= B[i, j] \quad \text{für alle } i, j \end{aligned}$$

Die Operatoren auf der rechten Seite der Äquivalenzbeziehungen sind dabei die Operatoren für Werte vom Typ *Lcinterval*.

Zusätzlich stehen die Operatoren **in** für die Relation „liegt in“ zwischen einem *Lrvector*- oder *Lcvector*- und einem *Lcivector*- Operanden bzw. zwischen einem *Lrmatrix*- oder *Lcmatrix*- und einem *Lcimatrix*-Operanden sowie für die Relation „enthalten im Innern“ zwischen zwei *Lcivector*-Operanden bzw. zwei *Lcimatrix*-Operanden zur Verfügung.  $><$  für den Test auf Disjunktheit zweier komplexer Intervallvektoren bzw. Intervallmatrizen steht ebenfalls zur Verfügung. Diese Operatoren sind jeweils komponentenweise definiert.

Die Verbandsoperatoren  $+*$  und  $**$  bezeichnen die komponentenweise Bildung der Intervall-Hülle bzw. des Durchschnitts, wie es bereits für den Typ *Lcinterval* im Modul `LCLARI` beschrieben wurde.

Die Übersicht über alle im Modul `LMVCLARI` definierten Operatoren wird, bedingt durch die große Zahl von Operatoren, in den beiden nachfolgenden Tabellen gegeben:

rechter Operand	Longreal	Lcomplex	Linterval	Lcinterval	Lrvector	Lcvector
linker Operand						
monadisch						
Longreal						
Lcomplex						
Linterval						*
Lcinterval					*	*
Lrvector				*, /		+*
Lcvector			*, /	*, /	+*	+*
Livector		*, /		*, /		$\diamond$ , =, <>, +*
Lcivector	*, /	*, /	*, /	*, /	$\diamond$ , =, <>, +*	$\diamond$ , =, <>, +*
Lrmatrix				*, /		
Lcmatrix			*, /	*, /		
Limatrix		*, /		*, /		*
Lcimatrix	*, /	*, /	*, /	*, /	*	*

Tabelle 7.18: Operatoren des Moduls LMVCLARI, Teil 1

$$\diamond \in \{+, -, *\} \quad \forall \in \{=, <>, <, <=, >, >=\}$$

rechter Operand	Livector	Lcivector	Lrmatrix	Lcmatrix	Limatrix	Lcimatrix
linker Operand						
monadisch		$+, -$				$+, -$
Longreal		*				*
Lcomplex	*	*			*	*
Linterval		*		*		*
Lcinterval	*	*	*	*	*	*
Lrvector		$\diamond, =, \langle \rangle, \mathbf{in}, +*$				
Lcvector	$\diamond, =, \langle \rangle, \mathbf{in}, +*$	$\diamond, =, \langle \rangle, \mathbf{in}, +*$				
Livector		$\diamond, \mathbf{in}, V, \rangle \langle, +*, **$				
Lcivector	$\diamond, V, \rangle \langle, +*, **$	$\diamond, \mathbf{in}, V, \rangle \langle, +*, **$				
Lrmatrix		*		$+*$		$\diamond, =, \langle \rangle, \mathbf{in}, +*$
Lcmatrix	*	*	$+*$	$+*$	$\diamond, =, \langle \rangle, \mathbf{in}, +*$	$\diamond, =, \langle \rangle, \mathbf{in}, +*$
Limatrix		*		$\diamond, =, \langle \rangle, +*$		$\diamond, \mathbf{in}V, \rangle \langle, +*, **$
Lcimatrix	*	*	$\diamond, =, \langle \rangle, +*$	$\diamond, =, \langle \rangle, +*$	$\diamond, V, \rangle \langle, +*, **$	$\diamond, \mathbf{in}V, \rangle \langle, +*, **$

Tabelle 7.19: Operatoren des Moduls LMVCLARI, Teil 2

$$\diamond \in \{+, -, *\} \quad \forall \in \{=, \langle \rangle, \langle, \langle =, \rangle, \rangle =\}$$



### 7.9.3 Transferfunktionen für komplexe Intervallvektoren

Zur Wandlung zwischen den Typen *Lrvector*, *Lcvector*, *Livector* und *Lcivector* sowie zwischen *Lcivector* und *civector* werden folgende Transferfunktionen bereitgestellt:

Funktion	Ergb.-Typ	Bedeutung
<code>compl (liv1,liv2)</code>	<i>Lcivector</i>	Komplexer Intervallvektor <i>lciv</i> mit $lciv[i] = lcompl(liv1[i], liv2[i])$
<code>compl (lv,liv)</code>	<i>Lcivector</i>	Komplexer Intervallvektor <i>lciv</i> mit $lciv[i] = lcompl(lv[i], liv[i])$
<code>compl (liv,lrv)</code>	<i>Lcivector</i>	Komplexer Intervallvektor <i>lciv</i> mit $lciv[i] = lcompl(liv[i], lrv[i])$
<code>compl (liv)</code>	<i>Lcivector</i>	Rein reeller komplexer Intervallvektor <i>lciv</i> mit $lciv[i] = lcompl(liv[i])$
<code>intval (lcv1,lcv2)</code>	<i>Lcivector</i>	Komplexer Intervallvektor <i>lciv</i> mit $lciv[i] = intval(lcv1[i], lcv2[i])$
<code>intval (lv,lcv)</code>	<i>Lcivector</i>	Komplexer Intervallvektor <i>lciv</i> mit $lciv[i] = intval(lv[i], lcv[i])$
<code>intval (lcv,lv)</code>	<i>Lcivector</i>	Komplexer Intervallvektor <i>lciv</i> mit $lciv[i] = intval(lcv[i], lv[i])$
<code>intval (lcv)</code>	<i>Lcivector</i>	Punktintervallvektor <i>lciv</i> mit $lciv[i] = intval(lcv[i])$
<code>short (lciv)</code>	<i>civector</i>	$lciv[i] \subseteq \text{short}(lciv[i])$
<code>re (lciv)</code>	<i>Livector</i>	Realteilvektor <i>liv</i> , mit $liv[i] = re(lciv[i])$
<code>im (lciv)</code>	<i>Livector</i>	Imaginärteilvektor <i>liv</i> , mit $liv[i] = im(lciv[i])$
<code>inf (lciv)</code>	<i>Lcvector</i>	Komplexer Vektor <i>lcv</i> der Untergrenzen mit $lcv[i] = inf(lciv[i])$
<code>sup (lciv)</code>	<i>Lcvector</i>	Komplexer Vektor <i>lcv</i> der Obergrenzen mit $lcv[i] = sup(lciv[i])$

*lv* = *Lrvector*-Ausdruck, *lcv*, *lcv1*, *lcv2* = *Lcvector*-Ausdruck,  
*liv*, *liv1*, *liv2* = *Livector*-Ausdruck, *lciv* = *Lcivector*-Ausdruck

### 7.9.4 Transferfunktionen für komplexe Intervallmatrizen

Zur Wandlung zwischen den Typen *Lmatrix*, *Lcmatrix*, *Limatrix* und *Lcimatrix* sowie zwischen *Lcimatrix* und *cimatrix* werden folgende Transferfunktionen bereitgestellt:

Funktion	Ergb.-Typ	Bedeutung
compl (liM1,liM2)	<i>Lcimatrix</i>	Komplexe Intervallmatrix <i>lciM</i> mit $lciM[i, j] = lcompl(liM1[i, j], liM2[i, j])$
compl (LM,liM)	<i>Lcimatrix</i>	Komplexe Intervallmatrix <i>lciM</i> mit $lciM[i, j] = lcompl(LM[i, j], liM[i, j])$
compl (liM,LM)	<i>Lcimatrix</i>	Komplexe Intervallmatrix <i>lciM</i> mit $lciM[i, j] = lcompl(liM[i, j], LM[i, j])$
compl (liM)	<i>Lcimatrix</i>	Rein reelle komplexe Intervallmatrix <i>lciM</i> mit $lciM[i, j] = lcompl(liM[i, j])$
intval (lcM1,lcM2)	<i>Lcimatrix</i>	Komplexe Intervallmatrix <i>lciM</i> mit $lciM[i, j] = lintval(lcM1[i, j], lcM2[i, j])$
intval (LM,lcM)	<i>Lcimatrix</i>	Komplexe Intervallmatrix <i>lciM</i> mit $lciM[i, j] = lintval(LM[i, j], lcM[i, j])$
intval (lcM,LM)	<i>Lcimatrix</i>	Komplexe Intervallmatrix <i>lciM</i> mit $lciM[i, j] = lintval(lcM[i, j], LM[i, j])$
intval (lcM)	<i>Lcimatrix</i>	Punktintervallmatrix <i>lciM</i> mit $lciM[i, j] = lintval(lcM[i, j])$
short (lciM)	<i>cimatrix</i>	$lciM[i, j] \subseteq \text{short}(lciM[i, j])$
re (lciM)	<i>Limatrix</i>	Realteilmatrix <i>liM</i> mit $liM[i, j] = re(lciM[i, j])$
im (lciM)	<i>Limatrix</i>	Imaginärteilmatrix <i>liM</i> mit $liM[i, j] = im(lciM[i, j])$
inf (lciM)	<i>Lcmatrix</i>	Komplexe Matrix <i>lcM</i> der Untergrenzen mit $lcM[i, j] = inf(lciM[i, j])$
sup (lciM)	<i>Lcmatrix</i>	Komplexe Matrix <i>lcM</i> der Obergrenzen mit $lcM[i, j] = sup(lciM[i, j])$

LM = *Lmatrix*-Ausdruck, lcM, lcM1, lcM2 = *Lcmatrix*-Ausdruck  
liM, liM1, liM2 = *Limatrix*-Ausdruck, lciM = *Lcimatrix*-Ausdruck

### 7.9.5 Überladungen des Zuweisungsoperators

Die komponentenweise Initialisierung von *Lcivector*- und *Lcimatrix*-Variablen sowie die Wandlungen von *Lrvector* bzw. *Lcvector* bzw. *Livector* nach *Lcivector* und *Lrmatrix* bzw. *Lcmatrix* bzw. *Limatrix* nach *Lcimatrix* und *civector* bzw. *cimatrix* nach *Lcivector* bzw. *Lcimatrix* werden in Form überladener Zuweisungen bereitgestellt:

Zuweisung	Bedeutung
$\text{lciv} := \text{Lr}$	$\text{lciv}[j] := \text{lcompl}(\text{lintval}(\text{Lr}))$ $j = \text{lb}(\text{lciv}), \dots, \text{ub}(\text{lciv})$
$\text{lciv} := \text{lc}$	$\text{lciv}[j] := \text{lintval}(\text{lc})$ $j = \text{lb}(\text{lciv}), \dots, \text{ub}(\text{lciv})$
$\text{lciv} := \text{li}$	$\text{lciv}[j] := \text{lcompl}(\text{li})$ $j = \text{lb}(\text{lciv}), \dots, \text{ub}(\text{lciv})$
$\text{lciv} := \text{lei}$	$\text{lciv}[j] := \text{lei}$ $j = \text{lb}(\text{lciv}), \dots, \text{ub}(\text{lciv})$
$\text{lciv} := \text{civ}$	$\text{lciv}[j] := \text{Long}(\text{civ}[j])$ $j = \text{lb}(\text{lciv}), \dots, \text{ub}(\text{lciv})$
$\text{lciv} := \text{lv}$	$\text{lciv} := \text{lcompl}(\text{lintval}(\text{lv}))$
$\text{lciv} := \text{lcv}$	$\text{lciv} := \text{lintval}(\text{lcv})$
$\text{lciv} := \text{iv}$	$\text{lciv} := \text{lcompl}(\text{liv})$
$\text{lcim} := \text{Lr}$	$\text{lcim}[j, k] := \text{lcompl}(\text{lintval}(\text{Lr}))$ $j = \text{lb}(\text{lcim}, 1), \dots, \text{ub}(\text{lcim}, 1)$ $k = \text{lb}(\text{lcim}, 2), \dots, \text{ub}(\text{lcim}, 2)$
$\text{lcim} := \text{lc}$	$\text{lcim}[j, k] := \text{lintval}(\text{lc})$ $j = \text{lb}(\text{lcim}, 1), \dots, \text{ub}(\text{lcim}, 1)$ $k = \text{lb}(\text{lcim}, 2), \dots, \text{ub}(\text{lcim}, 2)$
$\text{lcim} := \text{li}$	$\text{lcim}[j, k] := \text{lcompl}(\text{li})$ $j = \text{lb}(\text{lcim}, 1), \dots, \text{ub}(\text{lcim}, 1)$ $k = \text{lb}(\text{lcim}, 2), \dots, \text{ub}(\text{lcim}, 2)$
$\text{lcim} := \text{lei}$	$\text{lcim}[j, k] := \text{lei}$ $j = \text{lb}(\text{lcim}, 1), \dots, \text{ub}(\text{lcim}, 1)$ $k = \text{lb}(\text{lcim}, 2), \dots, \text{ub}(\text{lcim}, 2)$
$\text{lcim} := \text{cim}$	$\text{lcim}[j, k] := \text{Long}(\text{cim}[j, k])$ $j = \text{lb}(\text{lcim}, 1), \dots, \text{ub}(\text{lcim}, 1)$ $k = \text{lb}(\text{lcim}, 2), \dots, \text{ub}(\text{lcim}, 2)$
$\text{lcim} := \text{LM}$	$\text{lcim} := \text{lcompl}(\text{lintval}(\text{LM}))$
$\text{lcim} := \text{lcM}$	$\text{lcim} := \text{lintval}(\text{lcM})$
$\text{lcim} := \text{liM}$	$\text{lcim} := \text{lcompl}(\text{liM})$

$\text{lei} = \text{Lcinterval-Ausdr.}$ ,  $\text{lciv} = \text{Lcivector-Variable}$ ,  $\text{lcim} = \text{Lcimatrix-Variable}$   
 $\text{li} = \text{Linterval-Ausdruck}$ ,  $\text{liv} = \text{Livector-Ausdruck}$ ,  $\text{liM} = \text{Limatrix-Ausdruck}$   
 $\text{lc} = \text{Lcomplex-Ausdruck}$ ,  $\text{lcV} = \text{Lcvector-Ausdruck}$ ,  $\text{lcM} = \text{Lcmatrix-Ausdruck}$   
 $\text{Lr} = \text{Longreal-Ausdruck}$ ,  $\text{lv} = \text{Lrvector-Ausdruck}$ ,  $\text{LM} = \text{Lrmatrix-Ausdruck}$   
 $\text{civ} = \text{civector-Ausdruck}$ ,  $\text{cim} = \text{cimatrix-Ausdruck}$ ,  $\text{iv} = \text{ivector-Ausdruck}$

### 7.9.6 Standardfunktionen

Es stehen die Funktionen *id* und *null* zur Erzeugung einer Einheitsmatrix und einer Nullmatrix bzw. eines Nullvektors, die Funktionen *mid* und *diam* für die komponentenweise Berechnung von Mittelpunkt und Durchmesser, die Funktion *conj* für die Konjugation von komplexen Intervallvektoren und Intervallmatrizen, sowie die Funktionen *transp* und *herm* zur Berechnung der transponierten und der hermiteschen Matrix zur Verfügung. Darüber hinaus wird die Funktion *blow* für die komponentenweise Berechnung der Epsilonaufblähung bereitgestellt.

Funktion	Ergb.-Typ	Bedeutung
<i>null</i> ( <i>lciv</i> )	<i>Lrvector</i>	Nullvektor mit dem aktuellen Indexbereich von <i>lciv</i>
<i>null</i> ( <i>lciM</i> )	<i>Lrmatrix</i>	Nullmatrix mit den aktuellen Indexbereichen von <i>lciM</i>
<i>null</i> ( <i>lciM1</i> , <i>lciM2</i> )	<i>Lrmatrix</i>	Nullmatrix mit den aktuellen Indexbereichen der Produktmatrix $lciM1 \cdot lciM2$
<i>id</i> ( <i>lciM</i> )	<i>Lrmatrix</i>	Einheitsmatrix mit den aktuellen Indexbereichen von <i>lciM</i>
<i>id</i> ( <i>lciM1</i> , <i>lciM2</i> )	<i>Lrmatrix</i>	Einheitsmatrix mit den aktuellen Indexbereichen der Produktmatrix $lciM1 \cdot lciM2$
<i>mid</i> ( <i>lciv</i> )	<i>Lcvector</i>	Mittelpunktvektor <i>pv</i> mit $pv[i] = mid(lciv[i])$
<i>diam</i> ( <i>lciv</i> )	<i>Lrvector</i>	Durchmesservektor <i>dv</i> mit $dv[i] = diam(lciv[i])$
<i>mid</i> ( <i>lciM</i> )	<i>Lcmatrix</i>	Mittelpunktmatrix <i>pM</i> mit $pM[i, j] = mid(lciM[i, j])$
<i>diam</i> ( <i>lciM</i> )	<i>Lrmatrix</i>	Durchmessermatrix <i>dM</i> mit $dM[i, j] = diam(lciM[i, j])$
<i>conj</i> ( <i>lciv</i> )	<i>Lcivector</i>	Konjugiert komplexer Intervallvektor <i>lcivk</i> mit $lcivk[i] = conj(lciv[i])$
<i>conj</i> ( <i>lciM</i> )	<i>Lcimatrix</i>	Konjugiert komplexe Intervallmatrix <i>lciMk</i> mit $lciMk[i, j] = conj(lciM[i, j])$
<i>transp</i> ( <i>lciM</i> )	<i>Lcimatrix</i>	Transponierte Matrix <i>lciMt</i> von <i>lciM</i> mit $lciMt[i, j] = lciM[j, i]$
<i>herm</i> ( <i>lciM</i> )	<i>Lcimatrix</i>	Hermitesche Matrix <i>lciMh</i> von <i>lciM</i> mit $lciMh[i, j] = conj(lciM[j, i])$
<i>blow</i> ( <i>lciv</i> , <i>Lr</i> )	<i>Lcivector</i>	Vektorielle Epsilonaufblähung <i>ev</i> mit $ev[i] = blow(lciv[i], short(Lr))$
<i>blow</i> ( <i>lciM</i> , <i>Lr</i> )	<i>Lcimatrix</i>	Vektorielle Epsilonaufblähung <i>eM</i> mit $eM[i, j] = blow(lciM[i, j], short(Lr))$

*lciv* = *Lcivector*-Ausdruck; *lciM*, *lciM1*, *lciM2* = *Lcimatrix*-Ausdruck,  
*Lr* = *Longreal*-Ausdruck

**Beispiel:**

Für die komplexen Matrizen  $M$ ,  $M1$ ,  $M2$  vom Typ *Lcimatix* ergibt, nach Ausführung der Anweisungen

```
M1 := conj( transp (M) );
M2 := herm( M );
```

der logische Ausdruck

```
M1 = M2
```

den Wert *true*.

**7.9.7 Ein-/Ausgabeprozeduren**

Es stehen die Prozeduren

```
procedure read (var f : text; var a : Lcivector);
procedure read (var f : text; var A : Lcimatix);
procedure write (var f : text; a : Lcivector);
procedure write (var f : text; A : Lcimatix);
```

mit optionalem Fileparameter, beliebig vielen Ein-/Ausgabeparametern, jedoch ohne Formatspezifikationen zur Verfügung.

Die Eingabe eines komplexen Intervallvektors bzw. einer komplexen Intervallmatrix erfolgt komponentenweise entsprechend der Eingabe von *leinterval*-Größen, dabei wird eine Matrix stets zeilenweise eingelesen. Die Ausgabe eines komplexen Intervallvektors bzw. einer komplexen Intervallmatrix erfolgt wie bei der Eingabe komponentenweise mit dem Standardformat für *Longreal*-Größen aus Tabelle 3.4.

Die Eingabe einer dreidimensionalen, komplexen Intervallmatrix

$$\begin{pmatrix} 1 + i & 2 + i \cdot [1, 3] & [1.4, 1.5] + i \cdot 0.1 \\ [1, 2] + i & i \cdot [-1, 1] & i \\ [-2.3, 1.22] + i \cdot [-3, -2.99] & 0 & 1.25E-013 \end{pmatrix}$$

kann erfolgen in der Form:

(1, 1)	(2, [1, 3])	([1.4, 1.5], 0.1)	< return >
([1, 2], 1)	(0, [-1, 1])	(0, 1)	< return >
([-2.3, 1.22], [-3, -2.99])	0	1.25E-013	< return >

(1, 1) bedeutet hier also ein komplexes Punktintervall und (0, [-1, 1]) steht für ein rein imaginäres Intervall.

### 7.9.8 Exakte Auswertung von Ausdrücken

Zu berechnen ist das komplexe Intervall-Skalarprodukt

$$(A \cdot b) \cdot c$$

das im folgenden Programm zunächst **exakt** ausgewertet und anschließend in das nächstgrößere Intervall vom Typ *Lcinterval* gerundet wird:

```

program Beispiel (input,output);
use lci_ari, lmvci_ari;
var cidot,
    cidot_su : cidotprecision;
    A        : Lcimatrix[1..3,1..3]; { Kompl. Intervall-Matrix }
    b,c      : Lcivector[1..3];     { Kompl. Intervall-Vektoren }
    sk_pr    : Lcinterval;          { Skalar-Produkt-Einschlg. }
    i        : integer;

begin
  writeln('A: Lcimatrix = ?'); read(A); { Eingabe: zeilenweise }
  writeln('b: Lcivector = ?'); read(b);
  writeln('c: Lcivector = ?'); read(c);
  cidot_su := 0; { Initialisierung }
  For i := lb(a) to ub(a) do
  begin
    cidot := 0; { Initialisierung }
    PADDNACCU( cidot,LCIVECTOR(A[i,*]),b ); { A * b }
    cidot_su := cidot_su + cidot * c[i]; { Skalar-}
  end; { Produkt aufsummieren }
  sk_pr := LONG(cidot_su); { Rundg. nach sk_pr }
  writeln(sk_pr);
end.

```

#### Anmerkungen:

1. Das obige Skalarprodukt läßt sich für *A: cimatrix; b,c: civector* nicht mit dem **PASCAL-XSC** Lattenkreuzkonzept für *cimatrix-, civector*-Ausdrücke behandeln, es muß ebenfalls in einer analogen Schleife berechnet werden!
2. Die Konstruktion `LCIVECTOR(A[i,*])` liefert dem Compiler die Information, den *i*-ten Zeilenvektor `A[i,*]` als Vektor vom Typ *Lcivector* zu behandeln.
3. Die Prozedur `PADDNACCU(..)` aus dem Modul **LCI\_ARI** addiert zur Variablen `cidot = 0` das Skalarprodukt `A[i,*] · b` rundungsfehlerfrei; vergl. Seite 150, Nr. 17.
4. Die Funktion `LONG(cidot_su)` liefert mit dem komplexen Intervall `sk_pr` die bestmögliche Einschließung des exakten, komplexe Intervall-Skalarprodukts `cidot_su`; vergl. Tabelle 6.12 auf Seite 145.

# Kapitel 8

## PASCAL-XSC Module

Wie jede moderne Programmiersprache besitzt auch diese **PASCAL-XSC** BCD-Version ein Modulkonzept, das dem der **IEEE**-Version von **PASCAL-XSC** genau entspricht. Weitere Einzelheiten zum Gebrauch und zur Entwicklung von Modulen findet man in [5, 6, 7].

Wir unterscheiden zwischen **System**-, **Hilfs**- und **Anwendermodulen**:

- **Systemmodule** sind: `STDMOD`, `TIMER`, ... und bestimmen das zur Verfügung stehende **PASCAL-XSC** BCD-System.
- **Hilfsmodule** sind: `DDF_ARI`, `LSYS`, ... zur Behandlung numerischer Spezialgebiete wie Automatische Differentiation, Lineare Gleichungssysteme, ... . Da diese Module erfahrungsgemäß immer wieder erweitert und verbessert werden, stehen entsprechende Updates zur Verfügung. Alle Problemlösungen der Toolbox-Bände I,II [3, 4] sollen so auch für die BCD-Version durch die **Hilfsmodule** bereitgestellt werden. In den entsprechenden Beschreibungen findet der Anwender weniger den schon in den Toolboxbänden behandelten theoretischen Hintergrund sondern schwerpunktmäßig viele Beispiele und Hinweise zu den Lösungsroutinen, womit eine schnelle und problemlose Einarbeitung ermöglicht werden soll.
- **Anwendermodule** enthalten vom Anwender erstellte Problemlösungen, mit dem das **PASCAL-XSC** System nach Bedarf beliebig erweitert werden kann.

### 8.1 Systemmodule

#### 8.1.1 Modulhierarchie

Durch die `use`-Klausel wird in jedem Modul festgelegt, welche anderen Module zu importieren sind. Dabei spielt das Basis-Modul **STDMOD** eine Ausnahmestelle, weil es automatisch von jedem Modul oder Hauptprogramm benutzt wird und daher nicht in die `use`-Klausel aufgenommen werden muß.

Da dem Anwender mit der BCD-Version die verschiedenen Zahlenformate

**integer, real, Longreal, rrG, lrG, dotprecision**

zu Verfügung stehen, existiert eine Vielzahl von Modulen, deren Hierarchie nicht geschlossen auf einer Buchseite dargestellt werden kann. In nachfolgender Tabelle sind daher in der ersten Spalte alle Systemmodule in der Reihenfolge angegeben, wie sie vom Compiler zur Übersetzung benötigt werden. In jeder Zeile findet man in den nachfolgenden Spalten alle diejenigen Module, die vom Modul aus der ersten Spalte durch die **use**-Klausel importiert werden:

Modul	Importierte Module				
STDMOD					
TIMER					
X_INTG					
X_REAL					
X_STRG					
IOSTD					
LI_ARI					
I_ARI	li_ari				
RRGI_ARI	li_ari				
C_ARIAUX	x_real				
LC_HELP	li_ari				
LC_ARI	lc_help	li_ari			
LCI_ARI	i_ari	li_ari	lc_ari	rrGi_ari	
C_ARI	c_ariaux	i_ari	lc_ari		
CI_ARI	c_ari	i_ari	lci_ari	li_ari	
MV_ARI	iostd				
MVI_ARI	i_ari	mv_ari			
MVC_ARI	c_ari	mv_ari			
MVCI_ARI	ci_ari	i_ari	mvc_ari	mvi_ari	mv_ari
LMV_ARI	iostd				
LMVI_ARI	li_ari	lmv_ari	mv_ari		
LMVC_ARI	lc_ari	lmv_ari			
LMVCI_ARI	lci_ari	li_ari	lmvc_ari	lmvi_ari	lmv_ari
	mvci_ari				





Die string-Variablen `format` der Funktion `get_date` bewirkt die in folgender Tabelle zusammengestellten Ausgabeformate:

Bedeutung der string-Variablen <code>format</code> in <code>get_date</code>	
string-Variablen	Ausgabeformat
<code>format := ' '</code>	Wochentag Monat Monatstag-Nr. Uhrzeit Jahr
<code>format := '% a '</code>	Abgekürzter Wochentag
<code>format := '% A '</code>	Ausgeschriebener Wochentag
<code>format := '% b '</code>	Ausgeschriebener Monat
<code>format := '% B '</code>	Abgekürzter Monat
<code>format := '% c '</code>	Datum und Uhrzeit
<code>format := '% d '</code>	Nur Monatstag-Nr.
<code>format := '% H '</code>	Stundenangabe im Bereich: 0 – 23
<code>format := '% I '</code>	Stundenangabe im Bereich: 0 – 12
<code>format := '% j '</code>	Nur die Nummer des Jahrestages
<code>format := '% m '</code>	Nur der Monat
<code>format := '% M '</code>	Nur die Minute der Uhrzeit
<code>format := '% p '</code>	Nur AM / PM
<code>format := '% S '</code>	Nur die Sekundenangabe der Uhrzeit
<code>format := '% w '</code>	Nur die Wochentags-Nr.
<code>format := '% x '</code>	Datum in Zifferform, z.B. 17/02/39
<code>format := '% X '</code>	Uhrzeit in Zifferform, z.B. 21:13:33
<code>format := '% y '</code>	Jahresangabe ohne Jahrhundert, z.B. 97
<code>format := '% Y '</code>	Jahresangabe mit Jahrhundert, z.B. 1997
<code>format := '% % '</code>	Ausgabe des % Zeichens

### 8.1.2.3 X\_REAL

Das Modul `x_real` liefert zusätzliche Konstanten, Funktionen und Prozeduren bzgl. der Datentypen `REAL`, `LONGREAL`. In [7] findet man weitere Informationen zu diesen Modul-Routinen, die für die BCD-Version jedoch nur dann gültig sind, wenn in diesem Abschnitt ausdrücklich auf [7] verwiesen wird. Beachten Sie bitte, daß sich das genannte Handbuch nur auf die `IEEE`-Version bezieht!

## 8.1.2.3.1 Konstanten

Name	Wert	Bedeutung
IEEE_INV_OP	256	Illegale Operation
IEEE_DIV_BY_ZERO	2560	Division durch Null
IEEE_OVERFLOW	2816	Überlauf
IEEE_UNDERFLOW	2816	Unterlauf
IEEE_INEXACT	3328	Nicht exakte Operation
IEEE_CONTINUE	64	Zur Programmfortsetzung
IEEE_ALL	0	Zum Rücksetzen aller IEEE-flags
MAXREAL	9.99 ... 99E+255	Größte positive real-Zahl
MINREAL	1.00 ... 00E-267	Kleinste positive real-Zahl

Die ersten sieben integer-Konstanten aus obiger Tabelle werden zur Fehler- und Ausnahmebehandlung benötigt, die in einem der nächsten Abschnitte beschrieben wird.

## 8.1.2.3.2 Speicherformat für real/Longreal-Konstanten

Wir betrachten zunächst eine real-Konstante , für deren Speicherung 8 Byte mit zusammen 64 bits benötigt werden:

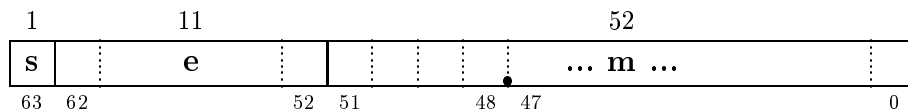


Abbildung 8.1: Speicherformat einer real-Konstanten

Die real-Zahl  $r = 0$  wird gespeichert, indem alle 64 bits auf Null gesetzt werden. Für  $r \neq 0$  bestimmt das 63.-te bit das Vorzeichen von  $r$  :

$$s = \begin{cases} 0 & \text{falls } r > 0 \\ 1 & \text{falls } r < 0 \end{cases}$$

Die vorzeichenlose integer-Zahl  $e$  ist definiert durch das 52. bis 62. bit, und die Mantisse  $m$  wird in den 52 restlichen bits gespeichert, wobei jede BCD-Ziffer der Mantisse 4 bits belegt, so daß insgesamt 13 BCD-Ziffern zur Verfügung stehen. Die Ziffer 9 wird beispielsweise durch die bit-Folge 1001 repräsentiert. Der Dezimalpunkt ist (gedanklich) stets zwischen dem 47. und 48. bit zu setzen. Nach diesen Vereinbarungen gilt für  $r \neq 0$  :

$$r = (-1)^s \cdot m \cdot 10^{e-255}; \quad 0 \leq e \leq 510$$

Für normalisierte real-Zahlen gilt:  $1 \leq m < 10$ ;  
 Für denormalisierte real-Zahlen gilt:  $10^{-12} \leq m \leq 0.999999999999$ ;

Eine **Longreal**-Zahl benötigt 12 Byte und wird ganz analog gespeichert:

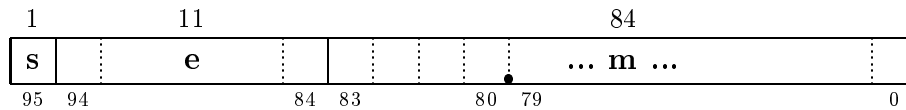


Abbildung 8.2: Speicherformat einer Longreal-Konstanten

Die Longreal-Zahl  $Lr = 0$  wird gespeichert, indem alle 96 bits auf Null gesetzt werden.

Für  $Lr \neq 0$  bestimmt das 95.-te bit das Vorzeichen von  $Lr$  :

$$s = \begin{cases} 0 & \text{falls } Lr > 0 \\ 1 & \text{falls } Lr < 0 \end{cases}$$

Die vorzeichenlose integer-Zahl  $e$  ist definiert durch das 84. bis 94. bit, und die Mantisse  $m$  wird in den 84 restlichen bits gespeichert, wobei jede BCD-Ziffer der Mantisse 4 bits belegt, so daß insgesamt 21 BCD-Ziffern zur Verfügung stehen. Der Dezimalpunkt ist (gedanklich) stets zwischen dem 79. und 80. bit zu setzen. Nach diesen Vereinbarungen gilt für  $Lr \neq 0$  wieder:

$$r = (-1)^s \cdot m \cdot 10^{e-511}; \quad 0 \leq e \leq 1022$$

Für normal. Longreal-Zahlen gilt:  $1 \leq m < 10$ ;

Für denormal. Longreal-Zahlen gilt:  $10^{-20} \leq m \leq 0.999999999999999999$ ;

### 8.1.2.3.3 Hexadezimale Ein-/Ausgabe

Die Darstellung des Speichers für **real**- und **Longreal**-Konstanten in hexadezimaler Form entsprechend der Abbildungen 8.1,8.2 wird ermöglicht durch die Prozeduren

```
procedure read (var f : text; var r : real);
procedure write (var f : text; var r : real);
```

```
procedure read (var f : text; var Lr : Longreal);
procedure write (var f : text; var Lr : Longreal);
```

indem an `r` bzw. `Lr` einen Doppelpunkt mit nachfolgendem 'x' oder 'X' angehängt wird. Mit den Wertzuweisungen

```
Lr := -1.23E-33; { Lr : Longreal }
writeLn(Lr: 'X'); write(Lr: 'x');
```

erhält man die Bildschirmausgabe:

```

9 D E 1 2 3 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
9 d e 1 2 3 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

```

'X' bewirkt demnach die Ausgabe der Hexadezimal-Zahlen in Großbuchstaben. Die Hexadezimalziffer 9 beschreibt die bit-Folge 1001 mit den Nummern 95 bis 92, und D beschreibt die bit-Folge 1101 mit den Nummern 91 bis 88. 9 D E definieren demnach die bit-Folge: 100111011110, woraus sich nach Abb. 8.2 für e der Wert 478 ergibt;  $478 - 511 = -33$  liefert dann den Zehnerexponenten -33 unseres Beispiels.

Ganz entsprechend kann man mit read(Lr:'X') und Eingabe der obigen 24-stelligen Hexadezimalfolge der Longreal-Variablen Lr den Wert  $-1.23 \cdot 10^{-33}$  zuweisen, wobei zwischen 'X' und 'x' nicht unterschieden wird. Für die Ein-/Ausgabe von real-Konstanten in 16-stelliger hexadezimaler Form gelten ganz analoge Aussagen.

#### 8.1.2.3.4 Mantisse/Exponent

Die Standardfunktionen **comp**, **mant**, **expo** beziehen sich auf die Form:

$$x = \text{mant}(x) \cdot 10^{\text{expo}(x)}, \quad \text{mit: } 0.1 \leq \text{mant}(x) < 1;$$

Die nachfolgenden Funktionen

```

function x_comp (m : real; ex : integer) : real;
function x_expo (r : real) : integer;
function x_mant (r : real) : real;

function x_comp (m : Longreal; ex : integer) : Longreal;
function x_expo (Lr : Longreal) : integer;
function x_mant (Lr : Longreal) : Longreal;

```

des Moduls **x\_real** beziehen sich auf eine Zahlendarstellung der Form:

$$x = \text{x\_mant}(x) \cdot 10^{\text{x\_expo}(x)}, \quad \text{mit: } 1 \leq \text{x\_mant}(x) < 10;$$

#### 8.1.2.3.5 Klassifizierung von real/Longreal-Werten

Die Klassifizierung von **real/Longreal**-Konstanten erfolgt mit dem Datentyp **x\_ccode** und den Funktionen:

```

function x_class (r : real) : x_ccode;
function x_class (Lr : Longreal) : x_ccode;
function x_value (c : x_ccode) : real;
function x_Lvalue (c : x_ccode) : Longreal;

```

Der nachfolgend definierte Typ **x\_ccode** unterscheidet dabei 10 verschiedene Klassen:

```

type x_ccode = ( x_sNaN, { signaling NaN          }
  x_qNaN, { quiet NaN              }
  x_minf, { minus infinity         }
  x_mnor, { negative normalized   }
  x_mden, { negative denormalized }
  x_mnul, { minus zero            }
  x_pnul, { plus zero             }
  x_pden, { positive denormalized }
  x_pnor, { positive normalized   }
  x_pinf, { plus infinity         }
);

```

Mit `c`: `x_ccode` liefert die Anweisung `c := x_class(r)` zum `real/Longreal`-Wert `r` den entsprechenden `c`-Wert vom Typ `x_ccode`.

Die Funktion `x_value` liefert zu jedem der 10 Argumente `c`: `x_ccode` einen **speziellen** `real`-Wert; die hexadezimalen Darstellungen dieser `real`-Konstanten sind zusammen mit ihrem jeweiligen Dezimalwert in folgender Tabelle zusammengestellt:

function x_value ( c : x_ccode ) : real		
c	Hexadezimalwert	Dezimalwert
x_sNaN	1 FF 8 0 0 0 0 F F F F F F F F	qNaN
x_qNaN	1 FF 0 0 0 0 0 F F F F F F F F	qNaN
x_minf	9 FF 0 0 0 0 0 0 0 0 0 0 0 0 0 0	-infinity
x_mnor	9 FE 9 9 9 9 9 9 9 9 9 9 9 9 9 9	$-9.99 \dots 99 \cdot 10^{+255}$
x_mden	8 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1	$-1.00 \dots 00 \cdot 10^{-267}$
x_mnul	8 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	$-0.00 \dots 00$
x_pnul	0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	$+0.00 \dots 00$
x_pden	0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1	$+1.00 \dots 00 \cdot 10^{-267}$
x_pnor	1 FE 9 9 9 9 9 9 9 9 9 9 9 9 9 9	$+9.99 \dots 99 \cdot 10^{+255}$
x_pinf	1 FF 0 0 0 0 0 0 0 0 0 0 0 0 0 0	+infinity

Ganz entsprechend liefert die Funktion `x_Lvalue` zu jedem der 10 Argumente `c`: `x_ccode` einen **speziellen** `Longreal`-Wert; die hexadezimalen Darstellungen dieser `Longreal`-Konstanten sind zusammen mit ihrem jeweiligen Dezimalwert in folgender Tabelle zusammengestellt:

function x_Lvalue ( c : x_ccode ) : Longreal		
c	Hexadezimalwert	Dezimalwert
x_sNaN	3 FF8 0000 FFFFFFFF00000000	qNaN
x_qNaN	3 FF0 0000 FFFFFFFF00000000	qNaN
x_minf	BFF0 0000 0000000000000000	-infinity
x_mnor	BFE9 9999 9999999999999999	$-9.99\dots99 \cdot 10^{+511}$
x_mden	8000 0000 0000000000000001	$-1.00\dots00 \cdot 10^{-531}$
x_mnul	8000 0000 0000000000000000	$-0.00\dots00$
x_pnul	0000 0000 0000000000000000	$+0.00\dots00$
x_pden	0000 0000 0000000000000001	$+1.00\dots00 \cdot 10^{-531}$
x_pnor	3 FE9 9999 9999999999999999	$+9.99\dots99 \cdot 10^{+511}$
x_pinf	3 FF0 0000 0000000000000000	+infinity

#### 8.1.2.3.6 IEEE Ausnahmebehandlungs-Routinen

Zusammen mit den im ersten Abschnitt definierten Konstanten stehen die nachfolgenden Funktionen und Prozeduren zur Ausnahmebehandlung zur Verfügung:

```

procedure IEEE_environment( action : integer;
  handler : integer;
  mode : boolean);
procedure IEEE_trap_enable( handler : integer;
  mode : boolean);
function IEEE_test ( handler : integer) : boolean;
procedure IEEE_set ( handler : integer);
procedure IEEE_reset ( handler : integer);
procedure IEEE_save ( var x : integer);
procedure IEEE_restore ( x : integer);

```

Eine genaue Beschreibung findet man in [7]; beachten Sie dort auch das Programmbeispiel im Anhang D.

#### 8.1.2.4 IOSTD, X\_INTG, X\_STRG

Die Beschreibung der Module `iostd`, `x_intg`, `x_strg` findet man in [7].

#### 8.1.2.5 LI\_ARI

Das Modul `Li_ari` stellt für den Datentyp `Longreal` alle für die reelle Intervallrechnung notwendigen Operatoren, Funktionen und Prozeduren zur Verfügung. Darüber hinaus existieren Operatoren und Prozeduren zur exakten Akkumulator-Auswertung von Longreal-Intervallausdrücken. Eine genaue Beschreibung findet man ab Seite 81.

### 8.1.2.6 I\_ARI

Das Modul `i_ari` stellt für den Typ `real` alle die für die reelle Intervallrechnung notwendigen Operatoren, Funktionen und Prozeduren zur Verfügung. Darüber hinaus existiert das komplette Lattenkreuzkonzept zur exakten Auswertung von Intervallausdrücken. Eine genaue Beschreibung findet man ab Seite 75.

### 8.1.2.7 RRGI\_ARI

Das Modul `rrGi_ari` stellt für den 26-stelligen Datentyp `rrG` alle die für die Implementierung reeller Intervall-Standardfunktionen notwendigen Operatoren und Prozeduren zur Verfügung. Darüber hinaus existieren Hilfsfunktionen zur Wertebereichseinschließung von monotonen Funktionen oder Funktionen mit einem oder mehreren lokalen Extrema. Eine genaue Beschreibung findet man ab Seite 93. Die Intervall-Standardfunktionen sind in Tabelle 4.8 auf Seite 94 zusammengestellt.

### 8.1.2.8 C\_ARIAUX

`c_ariaux` wird nur vom System-Modul `c_ari` importiert und exportiert zwei Routinen zur Realisierung der komplexen Division für Punktargumente:

1. procedure **produkt** (a,b,c,d : real; var overfl : boolean;  
var p1 : real; var p2 : interval);

Mit der Definition:

$$ex := \begin{cases} -267, & \text{falls overfl} = +1 \\ +0, & \text{falls overfl} = +0 \\ +267, & \text{falls overfl} = -1 \end{cases} \quad \text{gilt für alle } a,b,c,d : \text{real}$$

$$(a \cdot b + c \cdot d) \cdot 10^{ex} \in (p1 + p2)$$

$p1$  : maximalgenaue Näherung für:  $(a \cdot b + c \cdot d) \cdot 10^{ex}$

$p2$  : maximalgenaue Einschließung für:  $(a \cdot b + c \cdot d - p1) \cdot 10^{ex}$

Die Berechnung von  $p1, p2$  erfolgt wegen des Skalierungsfaktors  $10^{ex}$  ohne Programmabbruch durch internen Über- oder Unterlauf!

2. function **quotient** (z1 : real; z2 : interval; n1 : real; n2 : interval;  
round, zoverfl, noverfl : integer) : real;

Mit der Anweisung

```
y := quotient (z1,z2,n1,n2,rd,zo,no)
```

und  $zo, no \in \{-1, 0, +1\}$  wird intern zunächst durch staggered correction Technik unter der Voraussetzung:  $(n1 + n2) > 0$  eine Einschließung  $Q$  des exakten Intervalls

$$q := \frac{(z1 + z2)}{(n1 + n2)} \cdot 10^{(zo-no) \cdot 267} \subseteq Q$$



berechnet, wobei nur Überlauf entstehen kann. Der Funktionswert  $y$ : real wird dann durch den Rundungsparameter  $rd$  wie folgt bestimmt:

$$y := \begin{cases} \text{Nächstkleinere real-Zahl bzgl. } Q.\text{inf} \leq q.\text{inf}, & \text{falls } rd = -1 \\ \text{Nächste real-Zahl bzgl. 'Mittelpunkt' von } Q, & \text{falls } rd = +0 \\ \text{Nächstgrößere real-Zahl bzgl. } Q.\text{sup} \geq q.\text{sup}, & \text{falls } rd = +1 \end{cases}$$

**Anmerkungen:**

- Da die Einschließung  $Q$  mittels staggered correktion berechnet wird, ist  $Q$  eine sehr enge Einschließung des exakten Quotienten  $q$ . Deshalb liefert z.B.  $rd = -1$  mit  $y$  fast immer die größte Unterschranke von  $q$ , und nur in extremen Ausnahmesituationen liegt zwischen  $y$  und  $q.\text{inf}$  eine weitere real-Zahl. Für  $rd = +1$  gelten ganz entsprechende Aussagen. Da bei der komplexen Division die gerichteten Rundungen mittels  $rd = \pm 1$  berechnet werden, muß deshalb die Fehlerschranke größer sein als  $1 \cdot 10^{-12}$  und wird festgelegt durch  $1.5 \cdot 10^{-12}$ ; vergl. die Tabelle auf Seite 108.
- Wir setzen voraus, daß die Intervalle  $(z1 + z2), (n1 + n2)$  selbst, z.B. mit der Prozedur produkt(...), durch staggered correction Technik berechnet werden, so daß ihre Durchmesser extrem klein gegenüber dem Abstand benachbarter real-Zahlen angenommen werden, und dies gilt dann auch für die Einschließung  $Q$  des exakten Quotienten  $q$ . Enthält nun  $Q$  in einer seltenen Ausnahmesituation den Mittelpunkt seiner umgebenden real-Zahlen, so können unterschiedliche Zahlen aus  $Q$  verschiedene (benachbarte) 'nächste real-Zahlen' besitzen, so daß dann der oben genannte 'Mittelpunkt' von  $Q$  nur eine der beiden möglichen 'nächsten real-Zahlen' festlegt. Da bei der komplexen Division der Operator  $/$  durch die Funktion quotient(..) mit  $rd = 0$  definiert wird, ist seine Fehlerschranke nicht mehr  $5 \cdot 10^{-13}$  sondern  $1 \cdot 10^{-12}$ , d.h. die komplexe Division ist nur hochgenau; vergl. die Tabelle auf Seite 108.

### 8.1.2.9 LC\_HELP

**Lc\_help** wird nur vom System-Modul **Lc\_ari** importiert und exportiert zwei Routinen zur Realisierung der komplexen Division für Punktargumente:

1. procedure **produkt** (a,b,c,d : Longreal; var overfl : boolean;  
var p1 : Longreal; var p2 : Linterval);

Mit der Definition:

$$ex := \begin{cases} -531, & \text{falls overfl} = +1 \\ +0, & \text{falls overfl} = +0 \\ +531, & \text{falls overfl} = -1 \end{cases} \quad \text{gilt für alle } a,b,c,d : \text{Longreal}$$

$$(a \cdot b + c \cdot d) \cdot 10^{ex} \in (p1 + p2)$$

$p1$  : maximalgenaue Näherung für:  $(a \cdot b + c \cdot d) \cdot 10^{ex}$

$p2$  : maximalgenaue Einschließung für:  $(a \cdot b + c \cdot d - p1) \cdot 10^{ex}$

Die Berechnung von  $p1, p2$  erfolgt wegen des Skalierungsfaktors  $10^{ex}$  ohne Programmabbruch durch internen Über- oder Unterlauf!

2. function **quotient** ( $z1 : \text{Longreal}; z2 : \text{Linterval}; n1 : \text{Longreal}; n2 : \text{Linterval};$   
 $\text{round}, \text{zoverfl}, \text{noverfl} : \text{integer}$ ) :  $\text{Longreal}$ ;

Mit der Anweisung

```
y := quotient (z1, z2, n1, n2, rd, zo, no)
```

und  $zo, no \in \{-1, 0, +1\}$  wird intern zunächst durch staggered correction Technik unter der Voraussetzung:  $(n1 + n2) > 0$  eine Einschließung  $Q$  des exakten Intervalls

$$q := \frac{(z1 + z2)}{(n1 + n2)} \cdot 10^{(zo-no) \cdot 531} \subseteq Q$$

berechnet, wobei nur Überlauf entstehen kann. Der Funktionswert  $y : \text{Longreal}$  wird dann durch den Rundungsparameter  $rd$  wie folgt bestimmt:

$$y := \begin{cases} \text{Nächstkleinere Longreal-Zahl bzgl. } Q.\text{inf} \leq q.\text{inf}, & \text{falls } rd = -1 \\ \text{Nächste Longreal-Zahl bzgl. 'Mittelpunkt' von } Q, & \text{falls } rd = +0 \\ \text{Nächstgrößere Longreal-Zahl bzgl. } Q.\text{sup} \geq q.\text{sup}, & \text{falls } rd = +1 \end{cases}$$

#### Anmerkungen:

- Da die Einschließung  $Q$  mittels staggered correction berechnet wird, ist  $Q$  eine sehr enge Einschließung des exakten Quotienten  $q$ . Deshalb liefert z.B.  $rd = -1$  mit  $y$  fast immer die größte Unterschranke von  $q$ , und nur in extremen Ausnahmesituationen liegt zwischen  $y$  und  $q.\text{inf}$  eine weitere Longreal-Zahl. Für  $rd = +1$  gelten ganz entsprechende Aussagen. Da bei der komplexen Division die gerichteten Rundungen mittels  $rd = \pm 1$  berechnet werden, muß deshalb die Fehlerschranke größer sein als  $1 \cdot 10^{-20}$  und wird festgelegt durch  $1.5 \cdot 10^{-20}$ ; vergl. die Tabelle auf Seite 113.
- Wir setzen voraus, daß die Intervalle  $(z1 + z2), (n1 + n2)$  selbst, z.B. mit der Prozedur produkt(...), durch staggered correction Technik berechnet werden, so daß ihre Durchmesser extrem klein gegenüber dem Abstand benachbarter Longreal-Zahlen angenommen werden, und dies gilt dann auch für die Einschließung  $Q$  des exakten Quotienten  $q$ . Enthält nun  $Q$  in einer seltenen Ausnahmesituation den Mittelpunkt seiner umgebenden Longreal-Zahlen, so können unterschiedliche Zahlen aus  $Q$  verschiedene (benachbarte) 'nächste Longreal-Zahlen' besitzen, so daß dann der oben genannte 'Mittelpunkt' von  $Q$  nur eine der beiden möglichen 'nächsten Longreal-Zahlen' festlegt. Da bei der komplexen Division der Operator  $/$  durch die Funktion `quotient(.)` mit  $rd = 0$  definiert wird, ist seine Fehlerschranke nicht mehr  $5 \cdot 10^{-21}$  sondern  $1 \cdot 10^{-20}$ , d.h. die komplexe Division ist nur hochgenau; vergl. die Tabelle auf Seite 113.

#### 8.1.2.10 LC\_ARI

Das Modul `Lc_ari` stellt für den Datentyp `Longreal` alle für komplexe Rechnungen notwendigen Operatoren, Funktionen und Prozeduren zur Verfügung. Darüber hinaus existieren Operatoren und Prozeduren zur exakten Akkumulator-Auswertung von `Lcomplex`-Ausdrücken.

Für den 26-stelligen Datentyp `rrGcomplex` stehen zusätzlich die wichtigsten Operatoren und Transferfunktionen zur Realisierung 26-stelliger Standardfunktionen zur Verfügung. Eine genaue Beschreibung findet man ab Seite 112.

#### 8.1.2.11 LCI\_ARI

Das Modul `Lci_ari` stellt für den Datentyp `Longreal` alle für komplexe Intervallrechnungen notwendigen Operatoren, Funktionen und Prozeduren zur Verfügung. Darüber hinaus existieren Operatoren und Prozeduren zur exakten Akkumulator-Auswertung von `Lcinterval`-Ausdrücken.

Für den 26-stelligen Datentyp `rrGcinterval` stehen zusätzlich die wichtigsten Operatoren und Transferfunktionen zur Realisierung 26-stelliger Intervall-Standardfunktionen zur Verfügung. Eine genaue Beschreibung findet man dazu ab Seite 136.

#### 8.1.2.12 C\_ARI

Das Modul `c_ari` stellt für den Typ `real` alle die für komplexe Rechnungen notwendigen Operatoren, Funktionen und Prozeduren zur Verfügung. Darüber hinaus existiert das komplette Lattenkreuzkonzept zur exakten Auswertung von `complex`-Ausdrücken. Eine genaue Beschreibung findet man ab Seite 107.

#### 8.1.2.13 CI\_ARI

Das Modul `ci_ari` stellt für den Typ `real` alle die für komplexe Intervallrechnungen notwendigen Operatoren, Funktionen und Prozeduren zur Verfügung. Darüber hinaus existiert das komplette Lattenkreuzkonzept zur exakten Auswertung von `cinterval`-Ausdrücken. Genaue Beschreibungen findet man ab Seite 127.

#### 8.1.2.14 MV\_ARI

In diesem Modul werden zum Basis-Typ `real` die für das Rechnen mit reellen Vektoren und Matrizen notwendigen Operatoren, Funktionen und Prozeduren bereitgestellt. Darüber hinaus existiert das komplette Lattenkreuzkonzept zur exakten Auswertung von `rmatrix`- oder `rvector`-Ausdrücken. Genaue Beschreibungen findet man ab Seite 161.

#### 8.1.2.15 MVI\_ARI

In diesem Modul werden zum Basis-Typ `real` die für das Rechnen mit reellen Intervallvektoren und Intervallmatrizen notwendigen Operatoren, Funktionen und Prozeduren bereitgestellt. Darüber hinaus existiert das komplette Lattenkreuzkonzept zur exakten Auswertung von `imatrix`- oder `ivector`-Ausdrücken. Eine genaue Beschreibung findet man ab Seite 172.

#### 8.1.2.16 MVC\_ARI

Das Modul `mvc_ari` stellt zum Basis-Typ **real** die für das Rechnen mit komplexen Vektoren und Matrizen notwendigen Operatoren, Funktionen und Prozeduren zur Verfügung. Darüber hinaus existiert auch das komplette Lattenkreuzkonzept zur exakten Auswertung von `cmatrix`- oder `cvector`-Ausdrücken. Eine genaue Beschreibung findet man ab Seite 186.

#### 8.1.2.17 MVCI\_ARI

In diesem Modul werden zum Basis-Typ **real** die für das Rechnen mit komplexen Intervallvektoren und Intervallmatrizen notwendigen Operatoren, Funktionen und Prozeduren bereitgestellt. Darüber hinaus existiert das komplette Lattenkreuzkonzept zur exakten Auswertung von `cimatrix`- oder `civector`-Ausdrücken. Eine genaue Beschreibung findet man ab Seite 198.

#### 8.1.2.18 LMV\_ARI

In diesem Modul werden zum Basis-Typ **Longreal** die für das Rechnen mit reellen Vektoren und Matrizen notwendigen Operatoren, Funktionen und Prozeduren bereitgestellt. Operanden vom Basis-Typ **real** sind nicht erlaubt. Darüber hinaus existieren alle notwendigen Operatoren zur exakten Auswertung von `Lmatrix`- oder `Lrvector`-Ausdrücken, wobei auch `rmatrix`- bzw. `rvector`-Operanden möglich sind. Eine genaue Beschreibung findet man ab Seite 165.

#### 8.1.2.19 LMVI\_ARI

In diesem Modul werden zum Basis-Typ **Longreal** die für das Rechnen mit reellen Intervallvektoren und Intervallmatrizen notwendigen Operatoren, Funktionen und Prozeduren bereitgestellt. Operanden vom Basis-Typ **real** sind nicht erlaubt. Darüber hinaus existieren alle notwendigen Operatoren zur exakten Auswertung von `Limatrix`- oder `Livector`-Ausdrücken, wobei in diesen Ausdrücken auch `rmatrix`- und `imatrix`- bzw. `rvector`- und `ivector`-Operanden möglich sind. Eine genaue Beschreibung findet man ab Seite 178.

#### 8.1.2.20 LMVC\_ARI

In diesem Modul werden zum Basis-Typ **Longreal** die für das Rechnen mit komplexen Vektoren und Matrizen notwendigen Operatoren, Funktionen und Prozeduren bereitgestellt. Operanden vom Basis-Typ **real** sind nicht erlaubt. Zur exakten Auswertung von `Lcmatrix`- oder `Lcvector`-Ausdrücken stehen die entsprechenden Prozeduren aus dem Modul **Lc\_ari** zur Verfügung. Eine genaue Beschreibung mit Beispielen findet man ab Seite 191.

#### 8.1.2.21 LMVCI\_ARI

In diesem Modul werden zum Basis-Typ **Longreal** die für das Rechnen mit komplexen Intervallvektoren und Intervallmatrizen notwendigen Operatoren, Funktionen und Prozeduren bereitgestellt. Operanden vom Basis-Typ **real** sind nicht erlaubt. Zur exakten Auswertung von `Lcimatrix`- oder `Lcivector`-Ausdrücken stehen die entsprechenden Prozeduren aus dem Modul **Lci\_ari** zur Verfügung. Eine genaue Beschreibung mit Beispielen findet man ab Seite 207.

## 8.2 Hilfsmodule

In Form einer zusätzlichen Modulbibliothek stehen **PASCAL-XSC** -Routinen der BCD-Version für die Lösung häufig auftretender numerischer Probleme zur Verfügung. Die darin verwendeten Lösungsmethoden berechnen eine scharfe Einschließung der wahren Lösung des gestellten Problems und weisen gleichzeitig die Existenz und Eindeutigkeit im angegebenen Intervall nach. Diese Routinen haben daher folgende Vorzüge:

- Die Berechnung einer Lösung erfolgt mit maximaler oder hoher, immer aber mit kontrollierter Genauigkeit, auch in vielen schlecht konditionierten Fällen.
- Die Richtigkeit des Ergebnisses wird verifiziert, d.h. es wird eine Einschließungsmenge berechnet, in der die Existenz und Eindeutigkeit der exakten Lösung gesichert ist.
- Falls keine Lösung existiert oder das Problem zu schlecht konditioniert ist, so wird eine Fehlermeldung ausgegeben.
- Der Benutzer kann kaum etwas falsch machen, und somit sind diese Methoden auch von Nichtspezialisten anwendbar.
- Bei der BCD-Version entfallen die bei Binärversionen oft so lästigen Konversionsprobleme, da die über die Tastatur eingegebenen, dezimalen Gleitpunkt-Zahlen **fehlerfrei** gespeichert werden und so den Anwender-Routinen unverfälscht zur Verfügung stehen. Beachten Sie bitte, daß bei einer Binärversion die Dezimalzahl 0.1 schon nicht mehr fehlerfrei gespeichert wird!

Im einzelnen decken die **PASCAL-XSC** Routinen die folgenden Gebiete ab:

- Lineare Gleichungssysteme
  - Vollbesetzte Systeme
  - Matrixinvertierung
  - Ausgleichsprobleme
  - Berechnung der Pseudoinversen
  - Bandmatrizen
  - Spärlich besetzte Matrizen
- Polynomauswertung
  - in einer Variablen
  - in mehreren Variablen
- Polynomnullstellen
- Eigenwerte und Eigenvektoren
  - Symmetrische Matrizen
  - Beliebige Matrizen

- Anfangs- und Randwertprobleme bei gewöhnlichen Differentialgleichungen
  - linear
  - nichtlinear
- Auswertung von arithmetischen Ausdrücken
- Nichtlineare Gleichungssysteme
- Numerische Quadratur
- Integralgleichungen
- Automatische Differentiation
- Spezielle Funktionen der mathematischen Physik

Näheres zu den einzelnen Routinen und Modulen kann der begleitenden Dokumentation der **PASCAL-XSC** Numerikbibliothek entnommen werden.

Über die jeweils behandelte Grundproblematik hinaus sind die Einsatzmöglichkeiten dieser Routinen vielfältig. Mit ihnen wird die Beantwortung so interessanter und wichtiger Fragestellungen möglich, wie z.B.

- Feststellung der Kondition von Problemen durch Intervalleingabe.
- Feststellung von lokalen Ausschlußgebieten von Lösungen, indem man große und kleine Eindeutigkeitsbereiche bestimmt, deren Differenzmenge dann das Ausschlußgebiet darstellt.
- Verifikation von Aussagen wie Bestimmung des Mindestranges einer Matrix oder Bestimmung einer Kugel bzw. einer Halbebene, in der alle Nullstellen eines komplexen Polynoms liegen. Damit ist es z.B. möglich, die Stabilität elektrotechnischer Schaltungen zu garantieren, soweit das Modell die Realität erfaßt.
- Parameterkontrolle bei Modellbildungen. Es kann festgestellt werden, welche Modelldaten mit welcher Empfindlichkeit in eine Modellformel einwirken und umgekehrt, wie genau Meßdaten sein müssen, damit eine abhängige Größe überhaupt noch eine vorgeschriebene Genauigkeitsrelevanz besitzt.

Die Modulbibliothek soll den Anwender in die Lage versetzen, eigene Problemstellungen so zu lösen, daß seine numerischen Programmergebnisse wirkliche mathematische Aussagekraft besitzen und nicht nur eine Näherung liefern, über deren Genauigkeit allenfalls spekuliert werden kann nach dem Motto:

'Der Rechner liefert eine 16-stellige Lösung, von denen die ersten vier Ziffern  
hoffentlich wieder richtig sind'

Diese heute noch weit verbreitete Denkweise auch bei Mathematikern sollte mit Hilfe der jetzt vorhandenen **PASCAL-XSC** Werkzeuge bald überwunden sein!

# Literaturverzeichnis

- [1] G. Alefeld, J. Herzberger. – *Einführung in die Intervallrechnung*. Reihe Informatik/12; BI, 1974; ISBN 3-411-01466-0;
- [2] F. Blomquist. – *Automatische a priori Fehlerabschätzungen*. Inst. f. Angewandte Mathematik, Universität Karlsruhe, 1992.
- [3] R. Hammer, M. Hocks, U. Kulisch, D. Ratz. – *Numerical Toolbox for Verified Computing I*. Springer Series in Computational Mathematics 21.
- [4] R. Hammer, M. Hocks, U. Kulisch, D. Ratz. – *Numerical Toolbox for Verified Computing II*. Springer Series in Computational Mathematics 21.
- [5] R. Klätte, U. Kulisch, M. Neaga, D. Ratz, Ch. Ullrich **PASCAL-XSC** – *Sprachbeschreibung mit Beispielen*. Springer-Verlag, Heidelberg, 1991.
- [6] R. Klätte, U. Kulisch, M. Neaga, D. Ratz, Ch. Ullrich **PASCAL-XSC** – *Language Reference with Examples*. Springer-Verlag, New York, 1992.
- [7] Numeric Software GmbH: **PASCAL-XSC** : *A PASCAL Extension for Scientific Computation. User's Guide*. Numeric Software GmbH, P.O.Box 2232, D-76492, Baden-Baden, Germany, 1991.

# Stichwortverzeichnis

- a\_DIV\_bc** (Funktion)
  - , Longreal 25
  - , real 14
- ab\_DIV\_c** (Funktion)
  - , Longreal 25
  - , real 14
- ab\_DIV\_cd** (Funktion)
  - , Longreal 25
  - , real 14
- abc** (Funktion)
  - , Longreal 24
  - , real 14
- abcd** (Funktion)
  - , Longreal 25
  - , real 14
- abs** (Funktion)
  - , interval 131
  - , complex 110
  - , interval 78
  - , Linterval 142
  - , Lcomplex 115
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 58
  - , rrGinterval 155
  - , rrGcomplex 124
  - , rrGinterval 94
- abs\_rrG** (Funktion)
  - , Linterval-rrGinterval 156
  - , Lcomplex-rrGcomplex 125
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 60
- add\_11** (Prozedur) 47
- add\_11G** (Prozedur) 49
- add\_21** 67
- add\_21** (Prozedur) 47
- add\_22** 67
- add\_22** (Prozedur) 47
- addaccu** (Prozedur) 38
- addnaccu** (Prozedur) 39, 89, 90, 118, 119, 149
- akku\_times\_10power** (Prozedur) 39
- Akkumulator** 46, 56, 66
  - , zweimaliges Auslesen 46, 66
- Akkumulator-Einschränkung für Longreal-Produkte** 40
  - , Beispiele 41
- Akkuoperationen mit**
  - \* 42, 43, 87, 88, 117, 151
  - + 38, 43, 147, 151
  - - 38, 43, 147, 151
  - addaccu 38
  - addnaccu 39, 89, 90, 118, 119, 149
  - akku\_times\_10power 39
  - minus 43
  - minus\_c 117
  - minus\_i 87
  - paddnaccu 39, 90, 119, 150
  - plus 43
  - plus\_c 117
  - plus\_i 87
  - psubnaccu 39, 90, 119, 120, 150, 151
  - subaccu 38
  - subnaccu 39, 90, 119, 149, 150
  - times 42, 43, 165, 166, 169
  - times\_ci 147
  - times\_c 117
  - times\_i 87, 179
- Anwendermodul** 217
- arccos** (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 59
  - , rrGinterval 94
- arccos\_rrG** (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 61
- arccot** (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 59
  - , rrGinterval 94
- arccot\_rrG** (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 61



- arcosh (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 59
  - , rrGinterval 94
- arcosh\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 61
- arcoshlp (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 59
  - , rrGinterval 94
- arcoshlp\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 61
- arcoth (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 59
  - , rrGinterval 94
- arcoth\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 61
- arcsin (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 59
  - , rrGinterval 94
- arcsin\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 61
- arctan (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 59
  - , rrGinterval 94
- arctan\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 61
- arctan2 (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 59
  - , rrGinterval 94
- arctan2\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 61
- arg (Funktion)
  - , cinterval 130, 132
  - , complex 110
  - , Lcinterval 142
  - , Lcomplex 115
  - , rrGcinterval 155
  - , rrGcomplex 124
- arg\_red (Prozedur) 55
- arg\_rrG (Funktion)
  - , Lcinterval-rrGcinterval 156
  - , Lcomplex-rrGcomplex 125
- arglp (Funktion)
  - , cinterval 131, 133
  - , complex 110
  - , Lcinterval 142
  - , Lcomplex 115
  - , rrGcinterval 155
  - , rrGcomplex 124
- arglp\_rrG (Funktion)
  - , Lcinterval-rrGcinterval 156
  - , Lcomplex-rrGcomplex 125
- arithmetische Operatoren
  - , - (Akku) 38, 43, 147, 151
  - , \* (Akku) 42, 43, 87, 88, 117, 147, 151
  - , + (Akku) 38, 43, 147, 151
  - , cidotprecision-cdotprecision, cidotprecision 147
  - , cidotprecision-complex, Lcomplex, cinterval, Lcinterval 147
  - , cidotprecision-dotprecision, idotprecision 147
  - , cidotprecision-real, Longreal, interval, Linterval 147
  - , cimatrix-real, complex, interval, cinterval 198, 200, 201
  - , cimatrix-rmatrix, cmatrix, imatrix, cimatrix 198, 200, 201
  - , cimatrix-rvector, cvector, ivector, civector 198, 200, 201
  - , cinterval-real, complex, interval, cinterval 128
  - , civector-real, complex, interval, cinterval 198, 200, 201
  - , civector-rvector, cvector, ivector, civector 198, 200, 201
  - , cmatrix-cvector, rmatrix, cmatrix 186, 187
  - , cmatrix-real, complex, rvector 186, 187
  - , cvector-real, complex, rvector, cvector 186, 187
  - , dotprecision-Longreal 38
  - , Ergebnistyp cidotprecision 147

- , Ergebnistyp cimatrix 198, 200, 201
- , Ergebnistyp cinterval 128
- , Ergebnistyp civector 198, 200, 201
- , Ergebnistyp cmatrix 186, 187
- , Ergebnistyp cvector 186, 187
- , Ergebnistyp dotprecision 42, 169
- , Ergebnistyp idotprecision 87, 179
- , Ergebnistyp imatrix 172, 173
- , Ergebnistyp interval 76
- , Ergebnistyp ivector 172, 173
- , Ergebnistyp Lcimatrix 207, 209, 210
- , Ergebnistyp Lcinterval 136
- , Ergebnistyp Lcivector 207, 209, 210
- , Ergebnistyp Lcmatrix 191, 192
- , Ergebnistyp Lcomplex 112
- , Ergebnistyp Lcvector 191, 192
- , Ergebnistyp Limatrix 178, 179
- , Ergebnistyp Linterval 81
- , Ergebnistyp Livector 178, 179
- , Ergebnistyp Longreal 26
- , Ergebnistyp Lrmatrix 165, 166
- , Ergebnistyp Lrvector 165, 166
- , Ergebnistyp real 17
- , Ergebnistyp rmatrix 161, 162
- , Ergebnistyp rrGcinterval 154
- , Ergebnistyp rrGcomplex 123
- , Ergebnistyp rrGinterval 93
- , Ergebnistyp rvector 161, 162
- , imatrix-ivector, rmatrix, imatrix 173
- , imatrix-ivector, rmatrix, imatrix 172
- , imatrix-real, interval, rvector 172, 173
- , integer-real-complex 107
- , ivector-real, interval, rvector, ivector 172, 173
- , Lcimatrix-Longreal, Lcomplex, Linterval, Lcinterval 207, 209, 210
- , Lcimatrix-Lrmatrix, Lcmatrix, Limatrix, Lcimatrix 207, 209, 210
- , Lcimatrix-Lrvector, Lcvector, Livector, Lcivector 207, 209, 210
- , Lcinterval-interval, Linterval, cinterval, Lcinterval 136
- , Lcinterval-real, Longreal, complex, Lcomplex 136
- , Lcivector-Longreal, Lcomplex, Linterval, Lcinterval 207, 209, 210
- , Lcivector-Lrvector, Lcvector, Livector, Lcivector 207, 209, 210
- , Lcmatrix-Lcvector, Lrmatrix, Lcmatrix 191, 192
- , Lcmatrix-Longreal, Lcomplex, Lrvector 191, 192
- , Lcomplex-real, Longreal, complex, Lcomplex 112
- , Lcvector-Longreal, Lcomplex, Lrvector, Lcvector 191, 192
- , Limatrix-Livector, Lrmatrix, Limatrix 178, 179
- , Limatrix-Longreal, Linterval, Lrvector 178, 179
- , Linterval-real, Longreal, interval, Linterval 81
- , Livector-Longreal, Linterval, Lrvector, Livector 178, 179
- , Longreal-rrG 47, 48, 50, 51
- , lrG-lrG 67, 68, 69
- , Lrmatrix-Longreal, Lrvector, Lrmatrix 165, 166
- , Lrvector-Longreal, Lrvector 165, 166
- , minus 42
- , minus\_i 87
- , plus 42
- , plus\_i 87
- , real, complex, interval-cinterval 128
- , real-Longreal 26
- , real-real 17
- , real-rrG 47, 48, 50, 51
- , rmatrix-real, rvector, rmatrix 161, 162
- , rrG-Longreal 48, 51
- , rrG-real 48, 51
- , rrG-rrG 47, 48, 49, 50, 51
- , rrGcinterval-real, Longreal, rrG, rrGcinterval 154
- , rrGcomplex-real, Longreal, rrG, rrGcomplex 123
- , rrGinterval-real, Longreal, rrG, rrGinterval 93
- , rvector-real, rvector 161, 162
- , times 42, 166, 169
- , times\_ci 147
- , times\_i 87, 179
- arsinh (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 59
  - , rrGinterval 94
- arsinh\_rrG (Funktion)

- , Linterval-rrGinterval 95
- , Longreal-rrG 61
- artanh (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 59
  - , rrGinterval 94
- artanh\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 61
- Ausgabe
  - Beispiel: lrG 72
  - Beispiel: rrG 57
  - einer complex-Zahl 111
  - einer hexadezimalen Longreal-Zahl 222
  - einer hexadezimalen real-Zahl 222
  - einer komplexen Intervall-Matrix (cimatrix) 206
  - einer komplexen Intervall-Matrix (Lcimatrix) 215
  - einer komplexen Matrix (cmatrix) 190
  - einer komplexen Matrix (Lcmatrix) 196
  - einer Lcomplex-Zahl 116
  - einer Longreal-Zahl 37
  - einer lrG-Zahl 71
  - einer Matrix (imatrix) 177
  - einer Matrix (Limatrix) 183
  - einer Matrix (Lrmatrix) 168
  - einer Matrix (rmatrix) 164
  - einer real-Zahl 37
  - einer rrG-Zahl 56
  - einer rrGcomplex-Zahl 126
  - eines Intervalls (interval) 79
  - eines Intervalls (Linterval) 85
  - eines Intervalls (rrGinterval) 96
  - eines komplexen Intervall-Vektors (civector) 206
  - eines komplexen Intervall-Vektors (Lcivector) 215
  - eines komplexen Intervalls (cinterval) 135
  - eines komplexen Intervalls (Lcinterval) 143
  - eines komplexen Intervalls (rrGcinterval) 158
  - eines komplexen Vektors (cvector) 190
  - eines komplexen Vektors (Lcvector) 196
  - eines Vektors (ivector) 177
  - eines Vektors (Livector) 183
  - eines Vektors (Lrvector) 168
  - eines Vektors (rvector) 164
- Ausnahmebehandlungen 225
- awf\_bcd (Programm) 56
- awf\_bcds (Programm) 71
- B**asismodul
  - , exportierte numerische Datentypen 219
- Basismodul stdmod 219
- BCD-Format 3
  - , Vorteile 3
  - , Vorteile: Beispiele 4
- BCD-Version
  - , Vorteile 8
- Binär-Format 3
  - , Nachteile 3
  - , Nachteile: Beispiele 4
- blow (Funktion)
  - , cimatrix 205
  - , cinterval 131
  - , civector 205
  - , imatrix 176
  - , interval 78
  - , ivector 176
  - , Lcimatrix 214
  - , Lcinterval 142
  - , Lcivector 214
  - , Limatrix 182
  - , Linterval 83
  - , Livector 182
- blow\_z (Funktion)
  - , Limatrix 182
  - , Livector 182
- bound (Funktion) 37, 54
- bound\_toL (Funktion) 54
- bounds (Funktion) 157
- C**-Compiler 1
- c\_ari (Modul) 107, 229
- c\_ariaux (Modul) 226
- catalan\_rrG (rrG-Konstante) 64, 65
- cdotprecision (Typ) 117
- ceil (Funktion)
  - , Longreal 28
  - , real 18
- ci\_ari (Modul) 127, 229
- cidotprecision (Typ) 144
- cimatrix (Typ) 198
- cinterval (Typ) 127
- civector (Typ) 149, 198
- cmatrix (Typ) 186
- comp (Funktion)
  - , Longreal 30
  - , real 19
- compl (Funktion) 109, 129, 188, 193, 202, 203, 211, 212
- complex (Typ) 107

- conj (Funktion)
  - , cimatrix 205
  - , cinterval 130
  - , civector 205
  - , cmatrix 190
  - , complex 110
  - , cvector 190
  - , Lcimatrix 214
  - , Lcinterval 142
  - , Lcivector 214
  - , Lcmatrix 196
  - , Lcomplex 115
  - , Lcvector 195
  - , rrGcinterval 155
  - , rrGcomplex 124
- conj\_rrG (Funktion)
  - , Lcomplex-rrGcomplex 125
- cos (Funktion)
  - , cinterval 130
  - , complex 110
  - , interval 78
  - , Lcinterval 142
  - , Lcomplex 115
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 59
  - , rrGcinterval 155
  - , rrGcomplex 124
  - , rrGinterval 94
- cos\_rrG (Funktion)
  - , Lcinterval-rrGcinterval 156
  - , Lcomplex-rrGcomplex 125
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 61
- cosh (Funktion)
  - , cinterval 130
  - , complex 110
  - , interval 78
  - , Lcinterval 142
  - , Lcomplex 115
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 59
  - , rrGcinterval 155
  - , rrGcomplex 124
  - , rrGinterval 94
- cosh\_rrG (Funktion)
  - , Lcinterval-rrGcinterval 156
  - , Lcomplex-rrGcomplex 125
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 61
- cospi (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
- , real 16
- , rrG 59
- , rrGinterval 94
- , real 16
- , rrG 59
- , rrGcomplex 124
- , rrGinterval 94
- cospi\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 61
- cvector (Typ) 118, 186
- , real 16
- , rrG 59
- , rrGinterval 94
- cospi\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 61
- cot (Funktion)
  - , complex 110
  - , interval 78
  - , Lcomplex 115
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 59
  - , rrGcomplex 124
  - , rrGinterval 94
- cot\_rrG (Funktion)
  - , Lcomplex-rrGcomplex 125
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 61
- coth (Funktion)
  - , complex 110
  - , interval 78
  - , Lcomplex 115
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 59
  - , rrGcomplex 124
  - , rrGinterval 94
- coth\_rrG (Funktion)
  - , Lcomplex-rrGcomplex 125
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 61
- cotpi (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 59
  - , rrGinterval 94
- cotpi\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 61
- Datenkonversion 3
  - , Beispiele 4
- Datentyp
  - , cdotprecision 117
  - , cidotprecision 144
  - , cimatrix 198
  - , cinterval 127
  - , civector 149, 198
  - , cmatrix 186
  - , complex 107
  - , cvector 118, 186

- , dotprecision 38
  - , idotprecision 86
  - , imatrix 172
  - , interval 75
  - , ivector 89, 172
  - , Lcimatrix 207
  - , Lcinterval 136
  - , Lcivector 149, 207
  - , Lcmatrix 191
  - , Lcomplex 112
  - , Lvvector 118, 191
  - , Limatrix 178
  - , Linterval 81
  - , Livector 89, 178
  - , Longreal 22
  - , lrG 66
  - , Lrmatrix 165
  - , Lrvector 34, 165
  - , real 11
  - , rmatrix 161
  - , rrG 46
  - , rrGcinterval 153
  - , rrGcomplex 123
  - , rrGinterval 93
  - , rvector 34, 161
  - , x\_ccode 223
  - denormalisierte Zahlen
    - , Longreal 23
    - , real 12
  - diam (Funktion)
    - , cimatrix 205
    - , cinterval 131
    - , civector 205
    - , imatrix 176
    - , interval 78
    - , ivector 176
    - , Lcimatrix 214
    - , Lcinterval 142
    - , Lcivector 214
    - , Limatrix 182
    - , Linterval 83
    - , Livector 182
    - , rrGcinterval 155
    - , rrGinterval 94
  - diam\_rrG (Funktion)
    - , Lcinterval-rrGcinterval 156
    - , Linterval-rrGinterval 95
  - div\_22 (Prozedur) 50, 68
  - divdown\_test (Funktion) 108
    - , Lcomplex 113
    - , Longreal 25
    - , real 14
  - divnext\_test (Funktion) 108
    - , Lcomplex 113
    - , Longreal 25
    - , real 14
  - divup\_test (Funktion) 108
    - , Lcomplex 113
    - , Longreal 25
    - , real 14
  - , Lcomplex 113
  - , Longreal 25
  - , real 14
  - dotcompl (Funktion) 118
  - doticmpl (Funktion) 144, 145
  - dotintval (Funktion) 89
  - dotprecision (Typ) 38
- E**ingabe
- Beispiel: lrG 72
  - Beispiel: rrG 57
  - einer complex-Zahl 111
  - einer hexadezimalen Longreal-Zahl 222
  - einer hexadezimalen real-Zahl 222
  - einer komplexen Intervall-Matrix (cimatrix) 206
  - einer komplexen Intervall-Matrix (Lcimatrix) 215
  - einer komplexen Matrix (cmatrix) 190
  - einer komplexen Matrix (Lcmatrix) 196
  - einer Lcomplex-Zahl 116
  - einer Longreal-Zahl 37
  - einer lrG-Zahl 71
  - einer Matrix (imatrix) 177
  - einer Matrix (Limatrix) 183
  - einer Matrix (Lrmatrix) 168
  - einer Matrix (rmatrix) 164
  - einer real-Zahl 37
  - einer rrG-Zahl 56
  - einer rrGcomplex-Zahl 126
  - eines Intervalls (interval) 79
  - eines Intervalls (Linterval) 85
  - eines Intervalls (rrGinterval) 96
  - eines komplexen Intervall-Vektors (civector) 206
  - eines komplexen Intervall-Vektors (Lcivector) 215
  - eines komplexen Intervalls (cinterval) 135
  - eines komplexen Intervalls (Lcinterval) 143
  - eines komplexen Intervalls (rrGcinterval) 158
  - eines komplexen Vektors (cvector) 190
  - eines komplexen Vektors (Lcvector) 196
  - eines Vektors (ivector) 177
  - eines Vektors (Livector) 183
  - eines Vektors (Lrvector) 168
  - eines Vektors (rvector) 164
- Einschließung eines
- exakten Funktionswertes 36
  - Wertebereichs 97, 157

- em\_ganz (Funktion) 32
- entire (Funktion)
  - , Longreal 32
  - , real 21
  - , rrG 55
- eul (Longreal-Konstante) 33
- eul\_rrG (rrG-Konstante) 64, 65
- exakte Auswertung von
  - Lcinterval-Ausdrücken 144
  - Lcivector, Lcmatrix-Ausdrücken, Beispiel 216
  - Lcomplex-Ausdrücken 116
  - Lcomplex-Ausdrücken, Beispiele 121
  - Lcvector, Lcmatrix-Ausdrücken, Beispiel 197
  - Limatrix-Ausdrücken, Beispiele 184
  - Linterval-Ausdrücken 86
  - Livector, Limatrix-Ausdrücken 183
  - Longreal-Ausdrücken 38
  - Lrvector, Lrmatrix-Ausdrücken, Beispiele 169, 170
- exp (Funktion)
  - , cinterval 130
  - , complex 110
  - , interval 78
  - , Lcinterval 141
  - , Lcomplex 115
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 58
  - , rrGcinterval 155
  - , rrGcomplex 124
  - , rrGinterval 94
- exp\_rrG (Funktion)
  - , Lcinterval-rrGcinterval 156
  - , Lcomplex-rrGcomplex 125
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 60
- exp\_x2 (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 58
  - , rrGinterval 94
- exp\_x2\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 60
- expm1 (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
- , rrG 58
- , rrGinterval 94
- expm1\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 60
- expo (Funktion)
  - , Longreal 29
  - , lrG 70
  - , real 18
  - , rrG 53
- expo\_add (Funktion)
  - , real 19
- expo\_0 (Funktion)
  - , Longreal 30
  - , real 19
- Exponent
  - , Longreal 29, 30, 223
  - , lrG 70
  - , real 18, 19, 223
  - , rrG 53
- exp10 (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 58
  - , rrGinterval 94
- exp10\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 60
- exp2 (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 58
  - , rrGinterval 94
- exp2\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 60
- f**loor (Funktion)
  - , Longreal 28
  - , real 18
- frac (Funktion)
  - , Longreal 28
  - , real 18
- g**amma\_c (Longreal-Konstante) 33
- gamma\_c\_rrG (rrG-Konstante) 64, 65
- get\_date (Funktion) 219, 220
- get\_timer (Funktion) 219
- greater\_21 (Prozedur) 51, 69
- h**erm (Funktion)
  - , cimatrix 205
  - , cmatrix 190

- , Lcimatix 214
- , Lcimatix 196
- hexadezimale Ein-/Ausgabe
  - , Longreal-Konstanten 222
  - , real-Konstanten 222
- Hilfsmodul 217
- horn\_ts (Programm) 71
- horner (Prozedur) 56, 71
- horner\_cdot (Funktion) 121
- horner\_cidot (Funktion) 152
- horner\_dot (Funktion) 43
- horner\_idot (Funktion) 91
- horner\_l (Funktion) 34
- horner\_t (Programm) 56
- Hornerschema 34, 43, 55, 70, 91, 121, 152
  
- I**ari (Modul) 75, 226
- id (Funktion)
  - , cimatrix 205
  - , cmatrix 189
  - , imatrix 176
  - , Lcimatix 214
  - , Lcimatix 195
  - , Limatrix 182
  - , rmatrix 163
- idotprecision (Typ) 86
- IEEE\_all (real-Konstante) 221
- IEEE\_continue (real-Konstante) 221
- IEEE\_div\_by\_zero (real-Konstante) 221
- IEEE\_enviroment (Prozedur) 225
- IEEE\_inexact (real-Konstante) 221
- IEEE\_inv\_op (real-Konstante) 221
- IEEE\_overflow (real-Konstante) 221
- IEEE\_reset (Prozedur) 225
- IEEE\_restore (Prozedur) 225
- IEEE\_save (Prozedur) 225
- IEEE\_set (Prozedur) 225
- IEEE\_test (Funktion) 225
- IEEE\_trap\_enable (Prozedur) 225
- IEEE\_underflow (real-Konstante) 221
- ifo\_bcd (Programm) 56
- ifo\_bcds (Programm) 71
- im (Funktion) 109, 114, 118, 123, 129, 140, 145, 153, 188, 193, 202, 203, 211, 212
- image (Funktion)
  - , Longreal 31
  - , real 20
- imatrix (Typ) 172
- in (Operator) siehe Operatoren
- incl\_mami (Funktion) 103
- incl\_max (Funktion) 99
- inf (Funktion) 77, 82, 89, 129, 140, 145, 174, 180, 181, 202, 203, 211, 212
- init\_timer (Prozedur) 219
- int\_nr (Funktion) 36
- integer\_add (Prozedur) 19
- interval (Typ) 75
- Intervallarithmetik
  - , cidotprecision 144
  - , cimatrix 198
  - , cinterval 127
  - , civector 198
  - , idotprecision 86
  - , interval 75
  - , ivector, imatrix 172
  - , Lcimatix 207
  - , Lcinterval 136
  - , Lcivector 207
  - , Linterval 81
  - , Livector, Limatrix 178
  - , rrGcinterval 153
  - , rrGinterval 93
- intval (Funktion) 77, 129, 174, 180, 181, 202, 203, 211, 212
- iostd (Modul) 225
- ivector (Typ) 89, 172
  
- K**omplexe Arithmetik
  - , cdotprecision 117
  - , complex 107
  - , cvector, cmatrix 186
  - , Lcomplex 112
  - , Lcvector, Lcmatrix 191
  - , rrGcomplex 123
- komplexe Division 107, 112, 226, 227, 228
- komplexe Intervallarithmetik
  - , cidotprecision 144
  - , cinterval 127
  - , civector, cimatrix 198
  - , Lcinterval 136
  - , Lcivector, Lcimatix 207
  - , rrGcinterval 153
- komplexes Intervall-Skalarprodukt 216
- Konstanten
  - , Longreal 33
  - , real 15, 221
  - , rrG 64, 65
- Konversionsfehler 44
  - , Vermeidung 7, 8
  
- L**attenkreuzkonzept für cimatrix 216
- Lattenkreuzkonzept für cinterval 144
- Lattenkreuzkonzept für civector 216
- Lattenkreuzkonzept für complex 116
- Lattenkreuzkonzept für cvector, cmatrix 197
- Lattenkreuzkonzept für interval 86
- Lattenkreuzkonzept für ivector, imatrix 183
- Lattenkreuzkonzept für real 38
- Lattenkreuzkonzept für rvector, rmatrix 168
- Lattenkreuzsimulation für

- Lcinterval-Ausdrücke 144
- Lcivector-Ausdrücke 149, 150
- Lcivector-Ausdrücke: Beispiele 151
- Lcivector, Lcimatix-Ausdrücke, Beispiel 216
- Lcomplex-Ausdrücke 117
- Lcvector, Lcimatix-Ausdrücke, Beispiel 197
- Livector, Limatrix-Ausdrücke 183
- Longreal-Ausdrücke 38
- Longreal-Ausdrücke, Beispiele 43
- Lrvector, Lrmatix-Ausdrücke, Beispiele 169, 170
- Laufzeiten 7, 55, 70
- lc\_ari (Modul) **112**, 229
- lc\_help (Modul) 227
- lci\_ari (Modul) **136**, 229
- Lcimatix (Typ) 207
- Lcinterval (Typ) 136
- Lcivector (Typ) 149, 207
- Lcimatix (Typ) 191
- lcompl (Funktion) 114, 139
- Lcomplex (Typ) 112
- Lcvector (Typ) 118, 191
- less\_21 (Prozedur) 51, 69
- li\_ari (Modul) **81**, 225
- lid (Funktion)
  - , Lrmatix 167
- Limatrix (Typ) 178
- Linterval (Typ) 81
- lintval (Funktion) 82, 139, 140
- lintval\_f (Funktion) 92
- lintval\_frrG (Funktion) 92
- Livector (Typ) 89, 178
- lmv\_ari (Modul) 160, **165**, 230
- lmvc\_ari (Modul) 160, **191**, 230
- lmvci\_ari (Modul) 160, **207**, 230
- lmvi\_ari (Modul) 160, **178**, 230
- ln (Funktion)
  - , cinterval 130, 133
  - , complex 110
  - , interval 78
  - , Lcinterval 141, 142
  - , Lcomplex 115
  - , Linterval 83
  - , Longreal 27
  - , lrG 73
  - , real 16
  - , rrG 58
  - , rrGcinterval 155
  - , rrGcomplex 124
  - , rrGinterval 94
- ln\_ex (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
- , real 16
- , rrG 58
- , rrGinterval 94
- ln\_ex\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 60
- ln\_gam\_c (Longreal-Konstante) 33
- ln\_gam\_c\_rrG (rrG-Konstante) 64, 65
- ln\_pi (Longreal-Konstante) 33
- ln\_pi\_rrG (rrG-Konstante) 64, 65
- ln\_rrG (Funktion)
  - , Lcinterval-rrGcinterval 156
  - , Lcomplex-rrGcomplex 125
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 60
- ln\_10 (Longreal-Konstante) 33
- ln\_10\_rrG (rrG-Konstante) 64, 65
- ln\_2 (Longreal-Konstante) 33
- ln\_2\_rrG (rrG-Konstante) 64, 65
- ln\_sqrtx2y2 (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 58
  - , rrGinterval 94
- ln\_sqrtx2y2\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 60
- ln\_sqrt1px2y2 (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 58
  - , rrGinterval 94
- ln\_sqrt1px2y2\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 60
- lnull (Funktion)
  - , Lrmatix 167
  - , Lrvector 167
- ln1p (Funktion)
  - , cinterval 130, 134
  - , complex 110
  - , interval 78
  - , Lcinterval 141, 142
  - , Lcomplex 115
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 58
  - , rrGcinterval 155
  - , rrGcomplex 124
  - , rrGinterval 94
- ln1p\_rrG (Funktion)
  - , Lcinterval-rrGcinterval 156



- , Lcomplex-rrGcomplex 125
- , Linterval-rrGinterval 95
- , Longreal-rrG 60
- loga (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 58
  - , rrGinterval 94
- loga\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 60
- log1p (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , Longreal-rrG 60
  - , real 16
  - , rrG 58
  - , rrGinterval 94
- log1p\_rrG (Funktion)
  - , Linterval-rrGinterval 95
- log10 (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , lrG 73
  - , real 16
  - , rrG 58
  - , rrGinterval 94
- log10\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 60
- log2 (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 58
  - , rrGinterval 94
- log2\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 60
- long (Funktion) 28, 82, 89, 114, 140, 145
- longdown (Funktion) 39, 53, 120
- longdown\_test (Funktion) 39, 120
- longnext (Funktion) 39, 53, 120
- longnext\_test (Funktion) 39, 120
- Longreal-Konstanten 33
  - , Klassifizierung 223
- Longreal-Konstanten (Speicherformat) 222
- Longreal (Typ) 22
- longup (Funktion) 39, 53, 120
- longup\_test (Funktion) 39, 120
- lrG (Typ) 66
- Lrmatrix (Typ) 165
- Lrvector (Typ) 165
- ltransp (Funktion)
  - , Lrmatrix 167
- lval (Funktion) 30, 31
- lvnull (Funktion)
  - , Lrvector 167
- m**\_10hoch (Operator)
  - , Longreal 32
  - , real 21
- mant (Funktion)
  - , Longreal 29
  - , real 18
- mant\_expo\_product (Prozedur)
  - , Longreal 30
  - , real 21
- mant\_w (real-Konstante) 15
- mant\_wl (real-Konstante) 15
- mant\_26 (Prozedur) 53
- mant\_34 (Prozedur) 70
- Mantisse
  - , Longreal 29, 223
  - , lrG 70
  - , real 18, 223
  - , rrG 53
- Matrix/Vektorarithmetik
  - , Übersicht 159
  - , civector, cimatrix 198
  - , cvector, cmatrix 186
  - , ivector, imatrix 172
  - , Lcivector, Lcimatrix 207
  - , Lcvector, Lcmatrix 191
  - , Livector, Limatrix 178
  - , Lrvector, Lrmatrix 165
  - , module 160
  - , rvector, rmatrix 161
- max\_exp (real-Konstante) 15
- max\_expl (real-Konstante) 15
- maxlongreal (Longreal-Konstante) 33
- maxreal (real-Konstante) 11, 221
- mid (Funktion)
  - , cimatrix 205
  - , cinterval 131
  - , civector 205
  - , imatrix 176
  - , interval 78
  - , ivector 176
  - , Lcimatrix 214
  - , Lcinterval 142
  - , Lcivector 214
  - , Limatrix 182
  - , Linterval 83
  - , Livector 182
- min\_exp (real-Konstante) 15
- min\_exp\_ (real-Konstante) 15
- min\_expl (real-Konstante) 15

- min\_expl\_ (real-Konstante) 15
- min\_longreal (Longreal-Konstante) 33
- min\_Longreal (Longreal-Konstante) 23
- min\_real (real-Konstante) 11, 15
- minlongreal (Longreal-Konstante) 33
- minreal (real-Konstante) 12, 221
- minus (Operator) 42
- minus\_c (Operator) 117
- minus\_i (Operator) 87
- Modul
  - , Anwender- 217
  - , Basis- 219
  - , Hierarchie 217, 218
  - , Hilfs- 217
  - , Import- 218
  - , Konzept 217
  - , System- 217
- Modulbibliothek 231
- Module
  - , c\_ari 107, 229
  - , c\_ariaux 226
  - , ci\_ari 127, 229
  - , i\_ari 75, 226
  - , iostd 225
  - , lc\_ari 112, 229
  - , lc\_help 227
  - , lci\_ari 136, 229
  - , li\_ari 81, 225
  - , lmv\_ari 165, 230
  - , lmvc\_ari 191, 230
  - , lmvci\_ari 207, 230
  - , lmv\_i\_ari 178, 230
  - , mv\_ari 161, 229
  - , mvc\_ari 186, 230
  - , mvci\_ari 198, 230
  - , mvi\_ari 172, 229
  - , rrGi\_ari 93, 226
  - , stdmod 217, 219
  - , x\_intg 225
  - , x\_real 220
  - , x\_strg 225
- monadische Operatoren 47, 67, 76, 81, 128, 136, 154, 161, 165, 172, 178, 179, 186, 187, 191, 192, 198, 200, 201, 2
- mul\_11 (Prozedur) 49
- mul\_12 (Prozedur) 49, 68
- mul\_21 (Prozedur) 49
- mul\_22 (Prozedur) 49, 68
- mul\_22S (Prozedur) 49
- muldown\_test (Funktion) 108
  - , Lcomplex 113
  - , Longreal 24
  - , real 13
- mulnext\_test (Funktion) 108
  - , Lcomplex 113
  - , Longreal 24
- , real 13
- mulup\_test (Funktion) 108
  - , Lcomplex 113
  - , Longreal 24
  - , real 13
- mv\_ari (Modul) 160, 161, 229
- mvc\_ari (Modul) 160, 186, 230
- mvci\_ari (Modul) 160, 198, 230
- mvi\_ari (Modul) 160, 172, 229
- N**ormalisiere (Funktion)
  - , rrG 54
- normalisierte Zahlen
  - , Longreal 22
  - , real 11
- null (Funktion)
  - , cimatrix 205
  - , civector 205
  - , cmatrix 189
  - , cvector 189
  - , imatrix 175
  - , ivector 175
  - , Lcimatrix 214
  - , Lcivector 214
  - , Lcmatrix 195
  - , Lcvector 195
  - , Limatrix 182
  - , Livector 182
  - , rmatrix 163
  - , rvector 163
- Numeric Software GmbH 1
- numerische Datentypen 219
- O**peratoren
  - , +, \* (Akku), siehe arithmetische Operatoren
  - , arithmetische, siehe arithmetische Operatoren
  - , Fehlerschranken (complex) 108
  - , Fehlerschranken (Lcomplex) 113
  - , in 76, 81, 88, 93, 128, 138, 179, 199, 200, 201, 208, 209, 210
  - , m\_10hoch 21, 32
  - , minus 42
  - , minus\_c 117
  - , minus\_i 87
  - , mit Ergebnistyp cidotprecision 117
  - , mit Ergebnistyp cidotprecision 146, 147
  - , mit Ergebnistyp cidotprecision: Beispiele 148
  - , mit Ergebnistyp cidotprecision: Einschränkung bei Addition, Subtraktion 146
  - , mit Ergebnistyp cimatrix 198, 200, 201

- , mit Ergebnistyp `cinterval` 128
- , mit Ergebnistyp `civector` 198, 200, 201
- , mit Ergebnistyp `cmatrix` 186, 187
- , mit Ergebnistyp `complex` 107
- , mit Ergebnistyp `cvector` 186, 187
- , mit Ergebnistyp `dotprecision` 42, 169
- , mit Ergebnistyp `idotprecision` 87, 179
- , mit Ergebnistyp `imatrix` 172, 173
- , mit Ergebnistyp `interval` 76
- , mit Ergebnistyp `ivector` 172, 173
- , mit Ergebnistyp `Lcmatrix` 207, 209, 210
- , mit Ergebnistyp `Lcinterval` 136, 137, 138
- , mit Ergebnistyp `Lcivector` 207, 209, 210
- , mit Ergebnistyp `Lcmatrix` 191, 192
- , mit Ergebnistyp `Lcomplex` 112
- , mit Ergebnistyp `Lcvector` 191, 192
- , mit Ergebnistyp `Limatrix` 178, 179
- , mit Ergebnistyp `Linterval` 81
- , mit Ergebnistyp `Livector` 178, 179
- , mit Ergebnistyp `Longreal` 26
- , mit Ergebnistyp `lrG` 67, 68, 69
- , mit Ergebnistyp `Lrmatrix` 165, 166
- , mit Ergebnistyp `Lrvector` 165, 166
- , mit Ergebnistyp `real` 17
- , mit Ergebnistyp `rmatrix` 161, 162
- , mit Ergebnistyp `rrG` 47, 48, 49, 50, 51
- , mit Ergebnistyp `rrGinterval` 154
- , mit Ergebnistyp `rrGcomplex` 123
- , mit Ergebnistyp `rrGinterval` 93
- , mit Ergebnistyp `rvector` 161, 162
- , `plus` 42
- , `plus_c` 117
- , `plus_i` 87
- , Tabellen 42, 43, 76, 81, 87, 88, 109, 112, 117, 123, 128, 137, 138, 147, 154, 162, 166, 173, 179, 187, 192, 20
- , `times` 42, 165, 166, 169
- , `times_c` 117
- , `times_ci` 147
- , `times_i` 87, 179
- , Vergleichs-, siehe Vergleichsoperatoren
- P**addnaccu (Prozedur) 39, 90, 119, 150
- `pi_div4` (Longreal-Konstante) 33
- `pi_div4_rrG` (rrG-Konstante) 64, 65
- `pi_down` (Longreal-Konstante) 33
- `pi_half` (Longreal-Konstante) 33
- `pi_half_rrG` (rrG-Konstante) 64, 65
- `pi_half1` (Longreal-Konstante) 33
- `pi_rrG` (rrG-Konstante) 64, 65
- `pi_up` (Longreal-Konstante) 33
- `pi_1, pi_2, pi_3, pi_4` (real-Konstanten) 15
- `pi_5, pi_rest` (real-Konstanten) 15
- `pi2` (Longreal-Konstante) 33
- `pi2_rrG` (rrG-Konstante) 64, 65
- `plus` (Operator) 42
- `plus_c` (Operator) 117
- `plus_i` (Operator) 87
- `power` (Funktion)
  - , `cinterval` 130, 134
  - , `Lcinterval` 141, 142
  - , `rrGinterval` 155
- `power_rrG` (Funktion)
  - , `Lcinterval-rrGinterval` 156
- `power_rrG(x,n,rndmode)` (Funktion)
  - , Longreal-rrG 60
- `power_rrG(x,y)` (Funktion)
  - , `Linterval-rrGinterval` 95
  - , Longreal-rrG 60
- `power(x,n,rndmode)` (Funktion)
  - , Longreal 27
  - , real 16
  - , rrG 58
- `power(x,y)` (Funktion)
  - , `interval` 78
  - , `Linterval` 83
  - , Longreal 27
  - , real 16
  - , rrG 59
  - , `rrGinterval` 94
- `pred` (Funktion)
  - , Longreal 30
  - , real 19
- `prod_sum` (Funktion) 55
- `product` (Funktion)
  - , Longreal 24
  - , real 13
- produkt (Prozedur) 226, 227
- `psubnaccu` (Prozedur) 39, 90, 119, 120, 150, 151
- Q**uotient (Funktion) 226, 228
  - , Longreal 25
  - , real 14
- R**\_eul\_rrG (rrG-Konstante) 64, 65

- r\_ln\_10 (Longreal-Konstante) 33
- r\_ln\_10\_rrG (rrG-Konstante) 64, 65
- r\_ln\_2 (Longreal-Konstante) 33
- r\_pi (Longreal-Konstante) 33
- r\_pi\_r (real-Konstante) 15
- r\_pi\_rrG (rrG-Konstante) 64, 65
- r\_sqrtpi (Longreal-Konstante) 33
- r\_sqrtpi\_rrG (rrG-Konstante) 64, 65
- r\_sqrtpi2 (Longreal-Konstante) 33
- r\_sqrtpi2\_rrG (rrG-Konstante) 64, 65
- r\_sqrt2 (Longreal-Konstante) 33
- r\_sqrt2\_rrG (rrG-Konstante) 64, 65
- r\_sqrt2pi (Longreal-Konstante) 33
- r\_sqrt2pi\_rrG (rrG-Konstante) 64, 65
- r\_2pi (Longreal-Konstante) 33
- r\_2pi\_rrG (rrG-Konstante) 64, 65
- re (Funktion) 109, 114, 118, 123, 129, 140, 145, 153, 188, 193, 202, 203, 211, 212
- read (Prozedur) 37, 56, 71, 79, 85, 96, 111, 116, 126, 135, 143, 158, 164, 168, 177, 183, 190, 196, 206, 215, 222
- real-Konstanten 15, 221
  - , Klassifizierung 223
- real-Konstanten (Speicherformat) 221
- real-Konstanten (x\_real) 221
- real (Typ) 11
- realdown\_test (Funktion) 40
- realnext (Funktion) 53
- realnext\_test (Funktion) 40
- realup\_test (Funktion) 40
- rezip (Prozedur) 50
- rmatrix (Typ) 161
- round (Funktion)
  - , integer 18
  - , Longreal 28
  - , real 18
- rrG-Konstanten 64, 65
- rrG (Typ) 46
- rrGcinterval (Typ) 153
- rrGcintval (Funktion) 153
- rrGcompl (Funktion) 123
- rrGcomplex (Typ) 123
- rrGi\_lari (Modul) **93**, 226
- rrGinterval (Typ) 93
- rrGintval (Funktion) 93
- rrgintval\_f (Funktion) 98
- rrGintval\_f\_rrG (Funktion) 97
- rrgintval\_frrG (Funktion) 98
- rval (Funktion) 19, 20
- rvector (Typ) 161
  
- Scale (Funktion)
  - , lrG 70
  - , rrG 54
- short (Funktion) 29, 82, 89, 114, 140, 145, 180, 181, 193, 194, 212
  - , Lrmatrix 167
  - , Lrvector 167
- short\_test (Funktion) 29, 114, 194
  - , Lrmatrix 168
  - , Lrvector 168
- shortdown (Funktion) 29, 114, 193, 194
  - , Lrmatrix 167
  - , Lrvector 167
- shortdown\_test (Funktion) 29, 114, 194
  - , Lrmatrix 168
  - , Lrvector 168
- shortup (Funktion) 29, 114, 193, 194
  - , Lrmatrix 167
  - , Lrvector 167
- shortup\_test (Funktion) 29, 114, 194
  - , Lrmatrix 168
  - , Lrvector 168
- sign (Funktion)
  - , dotprecision 40
  - , interval 78
  - , Linterval 83
  - , Longreal 29
  - , real 18
  - , rrGinterval 94
- sin (Funktion)
  - , cinterval 130
  - , complex 110
  - , interval 78
  - , Lcinterval 142
  - , Lcomplex 115
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 59
  - , rrGcinterval 155
  - , rrGcomplex 124
  - , rrGinterval 94
- sin\_rrG (Funktion)
  - , Lcinterval-rrGcinterval 156
  - , Lcomplex-rrGcomplex 125
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 61
- sinh (Funktion)
  - , cinterval 130
  - , complex 110
  - , interval 78
  - , Lcinterval 142
  - , Lcomplex 115
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 59
  - , rrGcinterval 155
  - , rrGcomplex 124
  - , rrGinterval 94

- sinh\_rrG (Funktion)
  - , Lcinterval-rrGcinterval 156
  - , Lcomplex-rrGcomplex 125
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 61
- sinpi (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 59
  - , rrGinterval 94
- sinpi\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 61
- Speicherformat
  - , Longreal 222
  - , real 221
- sqr (Funktion)
  - , cinterval 130
  - , complex 110
  - , interval 78
  - , Lcinterval 141
  - , Lcomplex 115
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 58
  - , rrGcinterval 155
  - , rrGcomplex 124
  - , rrGinterval 94
- sqr\_rrG (Funktion)
  - , Lcinterval-rrGcinterval 156
  - , Lcomplex-rrGcomplex 125
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 60
- sqrt (Funktion)
  - , cinterval 130, 131
  - , complex 110
  - , interval 78
  - , Lcinterval 141, 142
  - , Lcomplex 115
  - , Linterval 83
  - , Longreal 27
  - , lrG 73
  - , real 16
  - , rrG 58
  - , rrGcinterval 155
  - , rrGcomplex 124
  - , rrGinterval 94
- sqrt\_negim (Funktion)
  - , complex 110
  - , Lcomplex 115
  - , rrGcomplex 124
- sqrt\_negim\_rrG (Funktion)
  - , Lcomplex-rrGcomplex 125
- sqrt\_pi (Longreal-Konstante) 33
- sqrt\_pi\_d2 (Longreal-Konstante) 33
- sqrt\_pi\_d2\_rrG (rrG-Konstante) 64, 65
- sqrt\_pi\_rrG (rrG-Konstante) 64, 65
- sqrt\_rrG (Funktion)
  - , Lcinterval-rrGcinterval 156
  - , Lcomplex-rrGcomplex 125
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 60
- sqrt\_u (Funktion)
  - , complex 110
  - , Lcomplex 115
  - , rrGcomplex 124
- sqrt\_u\_rrG (Funktion)
  - , Lcomplex-rrGcomplex 125
- sqrt\_eul\_rrG (rrG-Konstante) 64, 65
- sqrtx2m1 (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , lrG 73
  - , real 16
  - , rrG 58
  - , rrGinterval 94
- sqrtx2m1\_lrG (Funktion)
  - , Longreal-lrG 73
- sqrtx2m1\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 60
- sqrtx2y2 (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 58
  - , rrGinterval 94
- sqrtx2y2\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 60
- sqrt1mx2 (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , lrG 73
  - , real 16
  - , rrG 58
  - , rrGinterval 94
- sqrt1mx2\_lrG (Funktion)
  - , Longreal-lrG 73
- sqrt1mx2\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 60
- sqrt1pm1 (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 58

- , rrGinterval 94
- sqrt1pm1\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 60
- sqrt1px2 (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , lrG 73
  - , real 16
  - , rrG 58
  - , rrGinterval 94
- sqrt1px2\_lrG (Funktion)
  - , Longreal-lrG 73
- sqrt1px2\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 60
- sqrt10 (Longreal-Konstante) 33
- sqrt2 (Longreal-Konstante) 33
- sqrt2\_rrG (rrG-Konstante) 64, 65
- sqrt3 (Longreal-Konstante) 33
- sqrt5 (Longreal-Konstante) 33
- Standardfunktionen
  - , cinterval 130, 131
  - , civector, cimatrix 205
  - , cvector, cmatrix 189, 190
  - , interval 77, 78, 79
  - , ivector, imatrix 175, 176
  - , Lcinterval 141, 142
  - , Lcinterval-rrGinterval 156
  - , Lcivector, Lcimatrix 214
  - , Lcomplex-rrGcomplex 125
  - , Lcvector, Lcmatrix 195, 196
  - , Linterval 83
  - , Linterval-rrGinterval 95
  - , Livector, Limatrix 182
  - , Longreal 26
  - , Longreal-lrG 73
  - , Longreal-rrG 60, 61
  - , lrG 73
  - , Lrvector, Lrmatrix 167, 168
  - , real 15
  - , rrG 58, 59
  - , rrGcinterval 155
  - , rrGcomplex 124
  - , rrGinterval 94
  - , rvector, rmatrix 163
  - , verschachtelte, Beispiel 62
- stdmod (Modul) 217, 218, 219
- subaccu (Prozedur) 38
- subnaccu (Prozedur) 39, 90, 119, 149, 150
- succ (Funktion)
  - , Longreal 30
  - , real 19
- sup (Funktion) 77, 82, 129, 140, 145, 174, 180, 181, 202, 203, 211, 212
- Systemmodul 217
- Systemmodule 218
- tan** (Funktion)
  - , complex 110
  - , interval 78
  - , Lcomplex 115
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 59
  - , rrGcomplex 124
  - , rrGinterval 94
- tan\_rrG (Funktion)
  - , Lcomplex-rrGcomplex 125
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 61
- tanh (Funktion)
  - , complex 110
  - , interval 78
  - , Lcomplex 115
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 59
  - , rrGcomplex 124
  - , rrGinterval 94
- tanh\_rrG (Funktion)
  - , Lcomplex-rrGcomplex 125
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 61
- tanpi (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 59
  - , rrGinterval 94
- tanpi\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 61
- Teilintervallsuche 35
- test\_26 (real-Konstante) 15
- test\_34 (real-Konstante) 15
- timer (Modul) 219
- times (Operator) 42, 165, 166, 169
- times\_c (Operator) 117
- times\_ci (Operator) 147
- times\_i (Operator) 87, 179
- Toolbox-Bände 217
- trans (Prozedur) 69
- trans\_rrG (Funktion) 53
- Transferfunktionen
  - , cdotprecision-real, Longreal, dotprecision 118
  - , cidotprecision-complex, Lcomplex 144, 145

- , cidotprecision-dotprecision, idotprecision 144, 145
- , cidotprecision-real, Longreal, interval, Linterval 144, 145
- , cimatix-imatrix, rmatrix, cmatrix 203
- , cinterval-real, interval, complex 129
- , civector-ivector, rvector, cvector 202
- , cmatrix-rmatrix 188
- , complex-Lcomplex 114
- , complex-real 109
- , cvector-rvector 188
- , dotprecision-cdotprecision 118
- , dotprecision-Longreal 39
- , dotprecision-real 40
- , idotprecision-real, Longreal, interval, Linterval, dotprecision 89
- , imatrix-rmatrix 174
- , ivector-rvector 174
- , Lcimatix-Limatix, Lrmatrix, Lcmatrix 212
- , Lcinterval-real, Longreal, interval, Linterval, complex, Lcomplex 139, 140
- , Lcivector-Livector, Lrvector, Lcvector 211
- , Lcmatrix-Lrmatrix 193, 194
- , Lcomplex-cdotprecision 120
- , Lcomplex-integer, real, Longreal, complex 114
- , Lcvector-Lrvector 193, 194
- , Limatrix-Lrmatrix, imatrix 181
- , Linterval-real, Longreal, interval 82
- , Livector-Lrvector, ivector 180
- , Longreal-integer 28, 29
- , Longreal-real 29
- , Longreal-rrG 53
- , lrG-rrG 69
- , real-complex 109
- , real-interval 77
- , real-Longreal 28
- , rrG-Longreal 53
- , rrG-lrG 69
- , rrG-real 53
- , rrG-rrGcomplex 123
- , rrGcinterval-rrGinterval 153
- , rrGcomplex-rrG 123
- , rrGinterval-real, Longreal, rrG 93
- transform (Funktion) 55, 70
- transp (Funktion)
  - , cimatix 205
  - , cmatrix 190
  - , imatrix 176
  - , Lcimatix 214
  - , Lcmatrix 196
  - , Limatrix 182
  - , rmatrix 163
- trunc (Funktion)
  - , Longreal 29
  - , real 18
- Use-Klausel** 217
- V**erbandsoperatoren 81, 88, 93, 129, 137, 138, 173, 174, 179, 199, 200, 201, 208, 209, 210
- Vergleiche mit Prozeduren 51, 69
- Vergleichsoperatoren
  - , cimatix-rmatrix, cmatrix, imatrix, cimatix 199, 200, 201
  - , cinterval-interval, cinterval 128
  - , civector-rvector, cvector, ivector, cvector 199, 200, 201
  - , cmatrix-rmatrix, cmatrix 187
  - , complex-real, integer, complex 108
  - , cvector-rvector, cvector 187
  - , dotprecision-integer, real, Longreal, dotprecision 45
  - , idotprecision-interval-Linterval 88
  - , imatrix-rmatrix, imatrix 173
  - , interval-integer, real, interval 76
  - , ivector-rvector, ivector 173
  - , Lcimatix-Lrmatrix, Lcmatrix, Limatrix, Lcimatix 208, 209, 210
  - , Lcinterval-interval, Linterval, cinterval, Lcinterval 136
  - , Lcivector-Lrvector, Lcvector, Livector, Lcivector 208, 209, 210
  - , Lcmatrix-Lrmatrix, Lcmatrix 192
  - , Lcomplex-real, Longreal, complex, Lcomplex 112, 114
  - , Lcvector-Lrvector, Lcvector 192
  - , Limatrix-Lrmatrix, Limatrix 179
  - , Linterval-real, Longreal, interval, Linterval 81
  - , Livector-Lrvector, Livector 179
  - , Lrmatrix-Lrmatrix 166
  - , Lrvector-Lrvector 166
  - , real-Longreal 26
  - , real-real 17
  - , rmatrix-rmatrix 162
  - , rrG-Longreal 53
  - , rrG-real 52, 53
  - , rrG-rrG 52

- , rrGcinterval-rrGcinterval 154
- , rrGcomplex-rrGcomplex 123
- , rrGinterval-rrGinterval 93
- , rvector-rvector 162
- vnull (Funktion)
  - , rvector 163
- W**ertebereich
  - , bei lokalem Maximum 98
  - , bei lokalem Maximum, Minimum 102
  - , bei lokalem Minimum 100
  - , bei lokalem Minimum, Maximum 103
  - , bei mehr als zwei Extrema 105
  - , monotoner Funktionen 92, 97
  - , reellw. Funktion mit komplexen Punktargumenten 157
- Wertzuweisungen
  - , siehe Zuweisungsoperatoren
- write (Prozedur) 37, 56, 71, 79, 85, 96, 111, 116, 126, 135, 143, 158, 164, 168, 177, 183, 190, 196, 206, 215, 222
- X**\_ccode (Typ) 223, 224
- x\_class (Funktion) 223
- x\_comp (Funktion) 223
- x\_expo (Funktion) 223
- x\_intg (Modul) 225
- x\_lvalue (Funktion) 223, 224, 225
- x\_mant (Funktion) 223
- x\_real (Modul) 220
- x\_strg (Modul) 225
- x\_value (Funktion) 223, 224
- xplpowy (Funktion)
  - , interval 130, 134
  - , Lcinterval 141, 142
  - , rrGcinterval 155
- xplpowy\_rrG (Funktion)
  - , Lcinterval-rrGcinterval 156
- x2\_y2 (Funktion)
  - , interval 78
  - , Linterval 83
  - , Longreal 27
  - , real 16
  - , rrG 58
  - , rrGinterval 94
- x2\_y2\_rrG (Funktion)
  - , Linterval-rrGinterval 95
  - , Longreal-rrG 60
- Z**erl (Funktion) 32
- Zuweisungsoperatoren
  - , cidotprecision := real, Longreal, complex, Lcomplex, dotprecision 118

- , cidotprecision := complex, Lcomplex, cidotprecision 146
- , cidotprecision := interval, Linterval, idotprecision, cinterval, Lcinterval 146
- , cidotprecision := real, Longreal, dotprecision 146
- , cimatrix := real, complex, interval, cinterval 204
- , cimatrix := rmatrix, cmatrix, imatrix 204
- , cinterval := real, interval, complex 130
- , civector := real, complex, interval, cinterval 204
- , civector := rvector, cvector, ivector 204
- , cmatrix-real, complex, rmatrix 189
- , complex := real 109
- , cvector-real, complex, rvector 189
- , dotprec. := real, Longreal, dotprec. 45
- , imatrix := real, interval, rmatrix 175
- , interval := real 77
- , ivector := real, interval, rvector 175
- , Lcimatrix := cimatrix, Lrmatrix, Lcmatrix, Limatrix 213
- , Lcimatrix := Longreal, Lcomplex, Linterval, Lcinterval 213
- , Lcinterval := real, Longreal, cinterval, Lcinterval, complex, Lcomplex 141
- , Lcinterval := rrGcinterval 154
- , Lcivector := civector, Lrvector, Lcvector, ivector 213
- , Lcivector := Longreal, Lcomplex, Linterval, Lcinterval 213
- , Lcmatrix := Longreal, Lcomplex, Lrmatrix, cmatrix 195
- , Lcomplex := real, Longreal, complex 114
- , Lcomplex := rrGcomplex 126
- , Lcvector := cvector 195
- , Lcvector := Longreal, Lcomplex, Lrvector 194
- , Limatrix := Longreal, Linterval, Lrmatrix, imatrix 181
- , Linterval := real, Longreal, interval 85
- , Linterval := rrGinterval 96
- , Livector := Longreal, Linterval, Lrvector, ivector 181
- , Longreal := real 34



- , lrG := rrG 69
- , Lrmatrix := Longreal, rmatrix  
166
- , Lrvector := Longreal, rvector  
166
- , real := Longreal 34
- , rmatrix := real 162
- , rrG := Longreal 53
- , rrG := lrG 69
- , rrG := real 53
- , rrGinterval := interval, Linterval,  
rrGinterval, cinterval, Lcinterval  
154
- , rrGinterval := real, Longreal,  
complex, Lcomplex 154
- , rrGcomplex := Lcomplex 126
- , rrGinterval := real, Longreal, rrG,  
Linterval 96
- , rvector := real 162