

VNODE-LP

A Validated Solver for Initial Value Problems
in Ordinary Differential Equations

Nedialko S. Nedialkov

Department of Computing and Software
McMaster University
Hamilton, Ontario, Canada

Technical Report CAS-06-06-NN

© Nedialko S. Nedialkov, 2006

Contents

Preface	xi
I Introduction, Installation, Use	1
1 Introduction	3
1.1 The problem VNODE-LP solves	3
1.2 On Literate Programming	3
1.3 Applications	4
1.4 Limitations	5
1.5 Prerequisites	5
2 Installation	7
2.1 Prerequisites	7
2.2 Successful installations	7
2.3 Installation process	8
2.3.1 Extracting the source code	8
2.3.2 Preparing a configuration file	8
2.3.3 Building the VNODE-LP library and examples	9
2.3.4 Installing the library files	10
3 Examples	13
3.1 Basic usage	13
3.1.1 Problem definition	13
3.1.2 Main program	14
3.1.3 Files	16
3.1.4 Building an executable	16
3.1.5 Output	17
3.1.6 Standard coding	17
3.2 One-dimensional ODE	19
3.3 Time-dependent ODE	20
3.4 Interval initial conditions	22
3.5 Producing intermediate results	24
3.6 ODE control	26

3.6.1	Passing data to an ODE	26
3.6.2	Integration with parameter change	27
3.7	Integration control	28
3.8	Work versus order	33
3.9	Work versus problem size	35
3.10	Stepsize behavior	36
3.11	Stiff problems	38
4	Interface	41
4.1	Interval data type	41
4.2	Wrapper functions	41
4.3	Interval vector	43
4.4	Solver's public functions	43
4.4.1	Constructor	43
4.4.2	Integrator	43
4.4.3	Set functions	44
4.4.4	Get functions	44
4.5	Constructing an AD object	45
4.6	Some helpful functions	45
5	Testing	47
5.1	General tests	47
5.2	Linear problems	47
5.2.1	Constant coefficient problems	47
5.2.2	Time-dependent problems	48
5.3	Nonlinear problems	49
6	Listings	51
II	Third-party Components	59
7	Packages	61
8	IA package	63
8.1	Functions calling FILIB++	63
8.2	Functions calling PROFIL	65
9	Changing the rounding mode	69
9.1	Changing the rounding mode using FILIB++	69
9.2	Changing the rounding mode using BIAS	70
III	Linear Algebra and Related Functions	71
10	Vectors and Matrices	73

11	Basic functions	75
11.1	Vector operations	75
11.2	Matrix/vector operations	78
11.3	Matrix operations	78
11.4	Get/set column	81
11.5	Conversions	81
12	Interval functions	85
12.1	Inclusion	85
12.2	Interior	85
12.3	Radius	86
12.4	Width	86
12.5	Midpoints	86
12.6	Intersection	87
12.7	Computing h such that $[0, h]\mathbf{a} \subseteq \mathbf{b}$	87
12.7.1	The interval case	87
12.7.2	The interval vector case	88
13	QR factorization	91
14	Matrix inverse	93
14.1	Matrix inverse class	93
14.2	Computing A^{-1}	94
14.3	Enclosing the solution of a linear system	95
14.3.1	Initial box	96
14.3.2	Krawczyk's iteration	96
14.4	Enclosing the inverse of a general point matrix	97
14.5	Enclosing the inverse of an orthogonal matrix	99
14.6	Constructor and destructor	99
IV	Solver Implementation	101
15	Structure	103
16	Solution enclosure representation	105
16.1	Tight enclosure	105
16.2	A priori enclosure	107
17	Taylor coefficient computation	109
17.1	Taylor coefficients for an ODE solution	109
17.2	Taylor coefficients for the solution of the variational equation	110
17.3	AD class	111
18	Control data	113
18.1	Indicator type	113
18.2	Interrupt type	113

18.3	Control data	114
19	Computing a priori bounds	117
19.1	Theory background	117
19.2	The HOE class	118
19.3	Implementation of the HOE method	119
19.3.1	Computing p_j	119
19.3.2	Computing u_j and \tilde{y}_j	120
19.3.3	Computing a stepsize	121
19.3.4	Forming the time interval	123
19.3.5	Selecting a trial stepsize for the next step	125
19.3.6	Computing a priori bounds	125
19.4	Other functions	126
19.4.1	Constructor and destructor	126
19.4.2	Accept a solution	127
19.4.3	Set functions	127
19.4.4	Get functions	128
19.4.5	Enclosing β	128
20	Computing tight bounds on the solution	131
20.1	Theory background	131
20.1.1	Predictor	131
20.1.2	Corrector	132
20.1.3	Computing a solution representation	133
20.1.4	Computing Q_{j+1}	134
20.2	Implementation	134
20.2.1	The IHO class	134
20.2.2	Computing a tight enclosure	135
20.2.3	Initialization	135
20.2.4	Predictor	138
20.2.5	Corrector	140
20.2.6	Enclosure representation	147
20.2.7	Constructor	151
20.2.8	Destructor	152
20.2.9	Accepting a solution	152
20.2.10	Set and get functions	153
20.2.11	Constants	153
20.2.12	Sorting columns of a matrix	155
21	The VNODE class	159
21.1	Declaration	159
21.2	The integrator function	160
21.2.1	Input correctness	160
21.2.2	Determine direction	162
21.2.3	Initialization	162
21.2.4	Methods involved in the initialization	164

21.2.5	Validate existence and uniqueness	166
21.2.6	Check last step	166
21.2.7	Compute a tight enclosure	169
21.2.8	Decide	169
21.3	Constructor/destructor	170
21.4	Get functions	170
21.5	Set parameters	171
21.6	Files	172
21.6.1	Interface	172
21.6.2	Implementation	173
21.7	Interface to the VNODE-LP Package	173
V	AD Implementation	175
22	Using FADBAD++	177
22.1	Computing ODE Taylor coefficients	177
22.1.1	FadbadODE class	177
22.1.2	Function description	178
22.1.3	Files	179
22.2	Computing Taylor coefficients for the variational equation . . .	180
22.2.1	FadbadVarODE class	180
22.2.2	Function description	181
22.3	Files	182
22.4	Encapsulated FADBAD++ AD	183
A	Miscellaneous Functions	185
A.1	Vector output	185
A.2	Check if an interval is finite	185
A.3	Message printing	186
A.4	Check intersection	186
A.5	Timing	188
	Bibliography	189

List of Figures

1	Producing C++ and L ^A T _E X files from cweb files	xi
2.1	Variables of a VNODE-LP configuration file	9
2.2	File <code>config/MacOSXWithProfil</code>	10
2.3	File <code>config/LinuxWithProfil</code>	11
2.4	The first six lines of <code>makefile</code> in <code>vnodelp</code>	12
3.1	<code>makefile</code> in <code>vnodelp/user_program</code>	16
3.2	The “standard” C++ code of <code>basic.cc</code>	18
3.3	Plots generated using <code>integi.cc</code>	24
3.4	Midpoints of the computed bounds with $\beta = 8/3$ from 0 to 20; and with $\beta = 8/3$ changed to 5 at $t = 10$	29
3.5	Plots generated using <code>integctrl.cc</code>	32
3.6	Plots generated using <code>orderstudy.cc</code>	33
3.7	CPU time versus n for Problem 3.2. VNODE-LP takes 8 steps for each n	36
3.8	Plots generated using <code>orbit.cc</code>	37
3.9	Stepsize versus t on (3.8–3.9) for $\mu = 10, 10^2, 10^3, 10^4$	39
6.1	The <code>makefile</code> in the <code>examples</code> directory	52
6.2	The <code>gnuplot</code> file for generating the plot in Figure 3.3	53
6.3	The MATLAB code for the DETEST E1 problem	54
6.4	The <code>gnuplot</code> file for generating the plots in Figure 3.4	54
6.5	The <code>gnuplot</code> file for generating the plots in Figure 3.5	55
6.6	The <code>gnuplot</code> file for generating the plots in Figure 3.6	56
6.7	The <code>gnuplot</code> file for generating the plots in Figure 3.7	57
6.8	The <code>gnuplot</code> file for generating the plots in Figure 3.8	57
6.9	The <code>gnuplot</code> file for generating the plots in Figure 3.9	58
15.1	Classes in VNODE-LP. The triangle arrows denote <i>inheritance</i> relations; the normal arrows denote <i>uses</i> relations.	104
21.1	The case $t_{\text{end}} \subseteq T_j$. We set $t_{j+1} = t_{\text{end}}$	168
21.2	When close to t_{end} , we take the “middle” as the next integration point.	169

Preface

We present VNODE-LP, a C++ solver for computing bounds on the solution of an initial-value problem (IVP) for an ordinary differential equation (ODE). In contrast to traditional ODE solvers, which compute approximate solutions, this solver proves that a unique solution to a problem exists and then computes rigorous bounds that are guaranteed to contain it. Such bounds can be used to help prove a theoretical result, check if a solution satisfies a condition in a safety-critical calculation, or simply to verify the results produced by a traditional ODE solver.

This package is a successor of the VNODE [25], Validated Numerical ODE, package of N. Nedialkov. A distinctive feature of the present solver is that it is developed entirely using Literate Programming (LP) [17]. As a result, the correctness of VNODE-LP's implementation can be examined much easier than the correctness of VNODE—the theory, documentation, and source code of VNODE-LP are interwoven in this manuscript, which can be verified for correctness by a human expert, like in a peer-review process.

Literate programming. With LP, a program (or function) is normally subdivided into pieces of code or *chunks*, and each of them may be subdivided into smaller chunks. How they are divided and put together should be clear from the exposition.

The present document is produced by `cweave` [18] on L^AT_EX-like *cweb* files, which contain both L^AT_EX text and C++ code. The C++ code for VNODE-LP and all the examples are generated by running `ctangle` [18] on those files; see Figure 1.

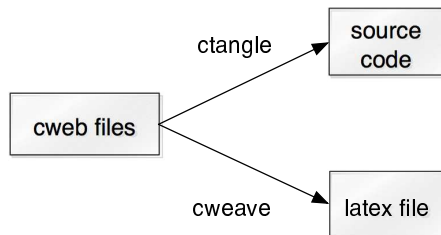


Figure 1. *Producing C++ and L^AT_EX files from cweb files*

Structure. Part I describes the problem VNODE-LP solves, shows how it can

be installed, and illustrates on several examples how VNODE-LP can be used. Parts II–V contain the implementation of this package.

If a reader is interested only in using VNODE-LP, then studying Part I should provide sufficient knowledge for using this package.

This document is *open*: errors found by a reader will be fixed and suggestions on improving it will be incorporated. Such suggestions can be on both exposition and code.

Acknowledgments. This work was supported in part by the Natural Sciences and Engineering Research Council of Canada.

George Corliss has made many valuable comments on this manuscript. Discussions with George Corlis, Baker Kearfott, John Pryce, and Spencer Smith have resulted in various improvements of the presentation.

N. Nedialkov
July 26, 2006

Part I

Introduction, Installation, Use

Chapter 1

Introduction

1.1 The problem VNODE-LP solves

We consider the IVP

$$y'(t) = f(t, y), \quad y(t_0) = y_0, \quad y \in \mathbb{R}^n, \quad t \in \mathbb{R}. \quad (1.1)$$

We denote the set of closed (finite) intervals on \mathbb{R} by

$$\mathbb{IR} = \{ \mathbf{a} = [\underline{\mathbf{a}}, \overline{\mathbf{a}}] \mid \underline{\mathbf{a}} \leq x \leq \overline{\mathbf{a}}, \underline{\mathbf{a}}, \overline{\mathbf{a}} \in \mathbb{R} \}.$$

An interval vector is a vector with interval components. We denote the set of n -dimensional interval vectors by \mathbb{IR}^n .

Given a point $t_{\text{end}} \neq t_0$ ($t_{\text{end}} \in \mathbb{R}$) and $\mathbf{y}_0 \in \mathbb{IR}^n$, the goal of VNODE-LP is to compute $\mathbf{y}_{\text{end}} \in \mathbb{IR}^n$ at t_{end} that contains the solution to (1.1) at t_{end} for all $y_0 \in \mathbf{y}_0$. If VNODE-LP cannot reach t_{end} , bounds on the solution at some t^ between t_0 and t_{end} are returned.*

This package is applicable to ODE problems for which derivatives of the solution $y(t)$ exist to some order; that is, $y(t)$ is sufficiently smooth. As a consequence, the code list of f should not contain functions such as branches, abs, or min.

In practice, t_0 or t_{end} , or both, may not be representable as floating-point numbers; for example the decimal 0.1 has an infinite binary representation. In this case, the user can set a machine-representable interval \mathbf{t}_0 [resp. \mathbf{t}_{end}] containing t_0 [resp. t_{end}].

1.2 On Literate Programming

The VNODE-LP package is a successor of VNODE [23, 25]. Both are written in C++. A major difference is that VNODE-LP is produced entirely, including this manuscript, using Literate Programming (LP) [17] and CWEB [18]. Why LP?

In general, interval methods produce results that can have the power of a mathematical proof. For example, when computing an enclosure of the solution of

an IVP ODE, an interval method first proves that there exists a unique solution to the problem and then produces bounds that contain it. When solving a nonlinear equation, an interval method can prove that a region does not contain a solution or compute bounds that contain a unique solution to the problem.

However, if an interval method is not implemented correctly, it may not produce rigorous results. Furthermore, we cannot claim mathematical rigor if we miss to include even a single roundoff error in a computation. Therefore, it is of paramount importance to ensure that an interval algorithm is encoded correctly in a programming language.

In the author’s opinion, interval software should be written such that it can be readily verified in a human peer-review process, like a mathematical proof is checked for correctness. The main goal of this work is to implement and document an interval solver for IVPs for ODEs such that its correctness can be verified by a reviewer.

To accomplish our goal, we have chosen the LP approach. The author has found LP particularly suitable for ensuring that an implementation of a numerical algorithm is a correct translation of its underlying theory into a programming language. Some of the benefits of employing LP follow.

- We can combine theory, source code, and documentation in a single document; we shall refer to it as an LP document.
- With LP, we can produce nearly “one-to-one” translation of the mathematical theory of a method into a computer program. In particular, we can split the theory into small pieces, translate each of them, and keep mathematical expressions and the corresponding code close together in a unified document. This facilitates verifying the correctness of smaller pieces and of a program as a whole.
- Since theory and implementation are in a single document, it is easier to keep them consistent, compared to having separate theory, source code, and documentation.

The user guide, theory, and source code of VNODE-LP are presented in the remainder of this document. The source code of VNODE-LP is extracted from source, *cweb* files using CWEB’s [18] `ctangle`. This manuscript is produced by running `cweave` on these files and then calling `LATEX`.

If the correctness of this manuscript is confirmed by reviewers in a peer-review-like process, we may trust the correctness of the implementation of VNODE-LP, and accept the bounds it computes as rigorous. When claiming rigor, however, we presume that the operating system, compiler, and the packages VNODE-LP uses do not contain errors.

1.3 Applications

Applications of validated integration include, for example, the solution of Smale’s 14th problem [30] and rigorous computation of asteroid orbits [7]. The (previous)

VNODE package had been employed in applications such as rigorous multibody simulations [5], reliable surface intersection [22, 28], computing bounds on eigenvalues [8], parameter and state estimation [15], rigorous shadowing [11, 12], and theoretical computer science [3].

1.4 Limitations

Generally, VNODE-LP is suitable for computing bounds on the solution of an IVP ODE with point initial conditions, or interval initial conditions with a sufficiently small width, over not very long time intervals. If the initial condition set is not small enough and/or long time integration is desired, the reader is referred to the Taylor models approach of Berz and Makino, and their COSY package. Alternatively, one can subdivide the initial interval vector (box) \mathbf{y}_0 into smaller boxes, perform integrations with them as initial conditions, and build an enclosure of the solution at the desired t_{end} .

1.5 Prerequisites

A user of VNODE-LP does not need to know how the underlying methods work. It is sufficient to know that, if \mathbf{a} and $\mathbf{b} \in \mathbb{IR}$ and $\bullet \in \{+, -, \times, /\}$, then VNODE-LP builds on the interval-arithmetic (IA) operations defined as

$$\mathbf{a} \bullet \mathbf{b} = \{x \bullet y \mid x \in \mathbf{a}, y \in \mathbf{b}\}, \quad (1.2)$$

where division is undefined if $0 \in \mathbf{b}$. This definition can be implemented, for example, as

$$\begin{aligned} \mathbf{a} + \mathbf{b} &= [\underline{\mathbf{a}} + \underline{\mathbf{b}}, \overline{\mathbf{a}} + \overline{\mathbf{b}}], \\ \mathbf{a} - \mathbf{b} &= [\underline{\mathbf{a}} - \overline{\mathbf{b}}, \overline{\mathbf{a}} - \underline{\mathbf{b}}], \\ \mathbf{a} \times \mathbf{b} &= [\min\{\underline{\mathbf{a}}\underline{\mathbf{b}}, \underline{\mathbf{a}}\overline{\mathbf{b}}, \overline{\mathbf{a}}\underline{\mathbf{b}}, \overline{\mathbf{a}}\overline{\mathbf{b}}\}, \max\{\underline{\mathbf{a}}\underline{\mathbf{b}}, \underline{\mathbf{a}}\overline{\mathbf{b}}, \overline{\mathbf{a}}\underline{\mathbf{b}}, \overline{\mathbf{a}}\overline{\mathbf{b}}\}], \quad \text{and} \\ \mathbf{a}/\mathbf{b} &= [\underline{\mathbf{a}}, \overline{\mathbf{a}}] \times [1/\overline{\mathbf{b}}, 1/\underline{\mathbf{b}}], \quad 0 \notin \mathbf{b}. \end{aligned}$$

On a computer, \mathbf{a} and \mathbf{b} are representable machine intervals, and the computed result of an IA operation must contain (1.2), provided no exceptions occur. For example, if intervals are represented by their endpoints, when computing $\mathbf{a} + \mathbf{b}$, the true $\underline{\mathbf{a}} + \underline{\mathbf{b}}$ is rounded towards $-\infty$, and the true $\overline{\mathbf{a}} + \overline{\mathbf{b}}$ is rounded towards $+\infty$.

From a language perspective, we have tried to avoid using advanced C++ techniques; however, basic knowledge of C++ is required.

The installation of VNODE-LP is explained in Chapter 2. Chapter 3 presents various examples of how VNODE-LP can be used. Chapter 4 lists and describes the functions available to a user of VNODE-LP. Chapter 5 contains descriptions of test cases. Various listings are given in Chapter 6.

Chapter 2

Installation

In this Chapter, we list the utilities and packages necessary for installing VNODE-LP, list successful installations, and then describe the installation process.

2.1 Prerequisites

The following utilities are needed:

1. `gunzip` (GNU unzip)
2. `tar` (tape archiver)
3. `ar` (for creating a library archive)
4. C++ compiler
5. GNU `make`
6. `libg2c` run-time library, if the GNU C++ compiler is used

Normally, 1–5 are present on a Unix-based system, while `libg2c` may need to be installed.

The following packages are used by VNODE-LP and must be installed before VNODE-LP is installed:

interval arithmetic: FILIB++ [19] *or* PROFIL/BIAS [16]

linear algebra: LAPACK [2] and BLAS [1]

2.2 Successful installations

To date VNODE-LP has been successfully compiled and installed as follows.

IA package	Operating system	Architecture	Compiler
FILIB++	Linux	x86	gcc
	Solaris	Sparc	gcc
PROFIL	Linux	x86	gcc
	Solaris	Sparc	gcc
	Mac OSX	PowerPC	gcc
	Windows with Cygwin	x86	gcc

Note. At the time of writing this manuscript, the author has not been able to install FILIB++ correctly on Mac OS X. However, VNODE-LP compiles on it.

2.3 Installation process

The installation process consists of the following steps:

1. extracting the source code
2. preparing a configuration file
3. building the VNODE-LP library, examples, and tests
4. installing the library files

2.3.1 Extracting the source code

VNODE-LP can be downloaded from www.cas.mcmaster.ca/~nedialk/vnodelp. The corresponding file is `vnodelp.tar.gz`. To extract the source files, type

```
tar -zxvf vnodelp.tar.gz
```

This will create the directory `vnodelp` and store the VNODE-LP files in it.

2.3.2 Preparing a configuration file

The user has to prepare a *configuration file*, which contains information such as compiler, options, libraries, and various directory paths. There are four such files used by the author: `MacOSXWithFilib`, `MacOSXWithProfil`, `LinuxWithFilib`, and `LinuxWithProfil`, located in `vnodelp/config`. One can modify any of these files or create his own configuration file, where the variables described in Figure 2.1 should be set appropriately. The files `MacOSXWithProfil` and `LinuxWithProfil` are given in Figures 2.2 and 2.3.

variable	stores
CXX	name of C++ compiler
CXXFLAGS	C++ compiler flags
GPP_LIBS	GNU C++ standard library <code>libstdc++</code> and the <code>libg2c</code> run-time library
LDLFLAGS	linker flags
I_PACKAGE	FILIB_VNODE or PROFIL_VNODE
I_INCLUDE	name of the directory containing include files of the interval-arithmetic package
I_LIBDIR	name of the directory containing interval libraries
I_LIBS	names of interval libraries
MAX_ORDER	value for the maximum order VNODE-LP can use
L_LAPACK	name of the directory containing the LAPACK library
L_BLAS	name of the directory containing the BLAS library
LAPACK_LIB	name of the LAPACK library file
BLAS_LIB	name of the BLAS library file

Figure 2.1. *Variables of a VNODE-LP configuration file*

2.3.3 Building the VNODE-LP library and examples

The `makefile` in `vnodelp` (see Figure 2.4) contains two variables that need to be set appropriately:

`CONFIG_FILE` contains the name of the configuration file; and

`INSTALL_DIR` contains the directory, where VNODE-LP should be installed.

After these variables are set appropriately, type

`make`

The library `libvnode.a` will be created in subdirectory `vnodelp/lib`, and the examples will be created in `vnodelp/examples`. Then, several test programs in subdirectory `tests` will be compiled and executed. If VNODE-LP compiles successfully and the tests pass, the following message should appear.

```
*****
***  VNODE-LP has compiled successfully
***  All tests have executed successfully
*****
If you have set the install directory, type
make install
```

```

CXX      = g++
CXXFLAGS = -O2 -g -Wall -pedantic -Wno-deprecated
GPP_LIBS = -lstdc++ /sw/lib/libg2c.a
LD_FLAGS += -bind_at_load -Wno

# interval package
LPACKAGE = PROFIL_VNODE
LINCLUDE =      $(HOME)/NUMLIB/Profil-2.1/src          \
                $(HOME)/NUMLIB/Profil-2.1/src/BIAS      \
                $(HOME)/NUMLIB/Profil-2.1/src/Base

L_LIBDIR =      $(HOME)/NUMLIB/Profil-2.1/src/BIAS      \
                $(HOME)/NUMLIB/Profil-2.1/src/Base      \
                $(HOME)/NUMLIB/Profil-2.1/src/lr

L_LIBS =      -lProfil -lBias -llr

MAX_ORDER = 50

# LAPACK and BLAS
LLAPACK = $(HOME)/NUMLIB/LAPACK
L_BLAS  = $(HOME)/NUMLIB/LAPACK
LAPACK_LIB = -llapack_MACOSX
BLAS_LIB   = -lblas_MACOSX

# — DO NOT CHANGE BELOW —
INCLUDES = $(addprefix -I, $(LINCLUDE))          \
           -I$(PWD)/FADBAD++
LIB_DIRS = $(addprefix -L, $(L_LIBDIR))          \
           $(LLAPACK) $(L_BLAS)
CXXFLAGS += -D${LPACKAGE} \
           -DMAXORDER=$(MAX_ORDER) $(INCLUDES)
LD_FLAGS += $(LIB_DIRS)
LIBS = $(L_LIBS) $(LAPACK_LIB) $(BLAS_LIB)        \
       $(GPP_LIBS)

```

Figure 2.2. *File config/MacOSXWithProfil*

2.3.4 Installing the library files

To install the library and the related include files, type

```
make install
```

This will create a subdirectory `vnodelp` of the directory stored in `INSTALL_DIR` and subdirectories of `vnodelp` as follows:

```

CXX = gcc
CXXFLAGS = -O2 -Wall -Wno-deprecated -DNDEBUG
GPP_LIBS = -lstdc++ -lg2c

# interval package
LPACKAGE = PROFIL_VNODE
LINCLUDE =
    $(HOME)/NUMLIB/Profil-2.0/include \
    $(HOME)/NUMLIB/Profil-2.0/include/BIAS \
    $(HOME)/NUMLIB/Profil-2.0/src/Base
LLIBDIR = $(HOME)/NUMLIB/Profil-2.0/lib
LLIBS = -lProfil -lBias -llr

MAX_ORDER = 50

# LAPACK and BLAS
LLAPACK =
LBLAS =
LAPACK_LIB = -llapack
BLAS_LIB = -lblas

# — DO NOT CHANGE BELOW —
INCLUDES = $(addprefix -I, $(LINCLUDE)) \
    -I$(PWD)/FADBAD++
LIB_DIRS = $(addprefix -L, $(LLIBDIR) \
    $(LLAPACK) $(LBLAS))
CXXFLAGS += -D${LPACKAGE} \
    -DMAXORDER=$(MAX_ORDER) $(INCLUDES)
LDFLAGS += $(LIB_DIRS)
LIBS = $(LLIBS) $(LAPACK_LIB) $(BLAS_LIB) \
    $(GPP_LIBS)

```

Figure 2.3. *File config/LinuxWithProfil*

directory	contains
lib	libvnode.a
include	libvnode.a's include files
config	configuration files
doc	documentation file <code>vnode.pdf</code>

Subsection 3.1.4 contains details about how to build user's programs.

```
# set CONFIG_FILE and INSTALL_DIR
```

```
CONFIG_FILE ?= MacOSXWithProfil  
INSTALL_DIR ?= $(HOME)
```

```
# — DO NOT CHANGE BELOW —
```

Figure 2.4. *The first six lines of `makefile` in `vnode1p`*

Chapter 3

Examples

We start with an example showing how a basic integration with VNODE-LP can be carried out, Section 3.1. In Section 3.2 we examine how VNODE-LP does on a simple scalar ODE. Section 3.3 contains an example of integrating a time-dependent system of ODEs and illustrates how this package can be used to check the correctness of the numerical results produced by a standard ODE method.

Section 3.4 outlines how to integrate with interval initial condition and output intermediate results. We describe how VNODE-LP outputs results at given time points in Section 3.5, and how parameters can be passed to an ODE problem in Section 3.6.

In Section 3.7, we show how an integration can be controlled, and in Section 3.8, we perform a simple study of the computational work versus the order of the method implemented in VNODE-LP. Section 3.9 contains a study of the computational work versus the size of the problem. Section 3.10 illustrates the stepsize behavior when integrating an orbit problem. Finally, Section 3.11 shows the stepsize behavior of VNODE-LP as the stiffness in an ODE increases.

3.1 Basic usage

In VNODE-LP, the user has to specify the right side of an ODE problem and provide a main program.

3.1.1 Problem definition

An ODE must be specified by a template function for evaluating $y' = f(t, y)$ of the form

```
18 <template ODE function 18> ≡  
    template<typename var_type>  
    void ODEName(int n, var_type *yp, const var_type *y, var_type t,  
        void *param)  
    {
```

```

        /* body */
    }

```

Here n is the size of the problem, t is the time variable, y is a pointer to input variables, yp is a pointer to output variables, and $param$ is a pointer to additional parameters that can be passed to this function.

As an example, consider the Lorenz system

$$\begin{aligned} y_1' &= \sigma(y_2 - y_1) \\ y_2' &= y_1(\rho - y_3) - y_2 \\ y_3' &= y_1y_2 - \beta y_3, \end{aligned}$$

where σ , ρ , and β are constants. This system is encoded in the *Lorenz* function below. The constants have values $\sigma = 10$, $\beta = 8/3$, and $\rho = 28$. We initialize *beta* with the interval containing $8/3$: `interval(8.0)` creates an interval with endpoints 8.0, and `interval(8.0)/3.0` is the interval containing $8/3$.¹ The last parameter, *param*, is not used here, but its role is discussed in Section 3.6.

```

19 <Lorenz 19> ≡
    template<typename var_type>
    void Lorenz(int n, var_type *yp, const var_type *y, var_type t,
                void *param)
    {
        interval sigma(10.0), rho(28.0);
        interval beta = interval(8.0)/3.0;
        yp[0] = sigma * (y[1] - y[0]);
        yp[1] = y[0] * (rho - y[2]) - y[1];
        yp[2] = y[0] * y[1] - beta * y[2];
    }

```

This code is used in chunks 20, 45, 48, 61, and 70.

3.1.2 Main program

We give a simple main program and explain its parts.

```

20 <simple main program 20> ≡
    <Lorenz 19>
    int main()
    {
        <set initial condition and endpoint 21>
        <create AD object 22>
        <create a solver 23>
        <integrate (basic) 24>
        <check if success 25>
    }

```

¹The result of this division is the interval with endpoints $8/3$ rounded toward $-\infty$ and $8/3$ rounded towards $+\infty$.

```

    <output results 26>
    return 0;
}

```

This code is used in chunk 27.

The initial condition and endpoint are represented as intervals in VNODE-LP. In this example, they are all point values stored as intervals. The components of **iVector** (interval vector) are accessed like a C/C++ array is accessed.

```

21 <set initial condition and endpoint 21> ≡
    const int n = 3;
    interval t = 0.0, tend = 20.0;
    iVector y(n);
    y[0] = 15.0;
    y[1] = 15.0;
    y[2] = 36.0;

```

This code is used in chunks 20, 42, 48, 58, and 61.

Then we create an AD object of type **FADBAD_AD**. It is instantiated with data types for computing Taylor coefficients (TCs) of the ODE solution and TCs of the solution to the variational equation, respectively [23]. To compute these coefficients, we employ the FADBAD++ package [29]. The first parameter in the constructor of **FADBAD_AD** is the size of the problem. The second and third parameters are the name of the template function, here *Lorenz*.

```

22 <create AD object 22> ≡
    AD *ad = new FADBAD_AD(n, Lorenz, Lorenz);

```

This code is used in chunks 20, 45, 48, 61, and 68.

Now, we create a solver:

```

23 <create a solver 23> ≡
    VNODE *Solver = new VNODE(ad);

```

This code is used in chunks 20, 35, 40, 45, 48, 58, 61, and 68.

The integration is carried out by the *integrate* function. It attempts to compute bounds on the solution at *tend*. When *integrate* returns, either $t = tend$ or $t \neq tend$. In both cases, *y* contains the ODE solution at *t*.

```

24 <integrate (basic) 24> ≡
    Solver->integrate(t, y, tend);

```

This code is used in chunks 20, 35, 40, and 68.

We check if an integration is successful by calling *Solver->successful()*:

```

25 <check if success 25> ≡
    if (!Solver->successful())
        cout << "VNODE-LP could not reach t=" << tend << endl;

```

This code is used in chunks 20, 35, 40, and 68.

Finally, we output the computed enclosure of the solution at t by

```
26 <output results 26> ≡
    cout << "Solution_enclosure_at_t=" << t << endl;
    printVector(y);
```

This code is used in chunks 20, 35, 40, and 47.

3.1.3 Files

The code of VNODE-LP is in the namespace **vnodelp**. The interface to VNODE-LP is stored in the file **vnodelp.h**, which must be included in any file using VNODE-LP. We store our program in the file **basic.cc**.

```
27 <basic.cc 27> ≡
#include <ostream>
#include "vnodelp.h"
using namespace std;
using namespace vnodelp;
<simple main program 20>
```

3.1.4 Building an executable

We describe how **basic.cc** is compiled and linked with the VNODE-LP library. The subdirectory **user_program** of **vnodelp** contains the files **basic.cc** and **makefile**, which is given in Figure 3.1. We consider this file here.

```
INSTALL_DIR = $(HOME)
CONFIG_FILE = $(INSTALL_DIR)/vnodelp/config/MacOSXWithProfil

include $(CONFIG_FILE)

CXXFLAGS += -I$(INSTALL_DIR)/vnodelp/include \
            -I$(INSTALL_DIR)/vnodelp/FADBAD++
LDLAGS    += -L$(INSTALL_DIR)/vnodelp/lib

basic:    basic.o
          $(CXX) $(LDLAGS) -o $@ basic.o -lvnode $(LIBS)

clean:
          @-$(RM) *.o core.* basic
```

Figure 3.1. *makefile in vnodelp/user_program*

The directory where **vnodelp** resides is set in **INSTALL_DIR**, and the configuration file is set in **CONFIG_FILE**. The variables **CXXFLAGS** and **LDLAGS** need not be changed. Finally, the rule for building **basic** is given. To create the executable file **basic**, type **make** in **vnodelp/user_program**.

3.1.5 Output

The output of `basic` is

```
Solution enclosure at t = [20,20]
14.30[38159449608937,44694855332662]
9.5[785941360078012,801274302834650]
39.038[2373597549516,4111043348412]
```

These results are interpreted as

$$y(20) \in \begin{pmatrix} [14.3038159449608937, 14.3044694855332662] \\ [9.5785941360078012, 9.5801274302834650] \\ [39.0382373597549516, 39.0384111043348412] \end{pmatrix}. \quad (3.1)$$

For comparison, if we integrate this problem with MAPLE using `dsolve` with options `method=taylorseries` and `abserr=Float(1,-18)`, and with `Digits := 20`, we obtain

$$y(20) \approx \begin{pmatrix} 14.304146251277895001 \\ 9.5793690774871976695 \\ 39.038325167739731729 \end{pmatrix},$$

which is contained in the bounds (3.1).

Remarks

1. All numerical results in this manuscript are produced with PROFIL on 1.25 GHz PowerPC G4 with MacOS X, 512 MB RAM, and 512KB L2 cache.
2. The output format is due to the PROFIL/BIAS [16] interval-arithmetic package.
3. On different architectures, or with different IA packages on the same architecture, the computed results are likely to differ, but they must contain the true results.

3.1.6 Standard coding

All source-code files in the VNODE-LP distribution, except the test programs in subdirectory `vnodelp/tests`, are generated with `ctangle` from CWEB. Since a user may not use LP, we also give the “standard” C++ code of `basic.cc` in Figure 3.2.

For the remaining examples, we do not explain how they are compiled. For details, see the `makefile` in Figure 6.1. This file is in the directory `vnodelp/examples`. Also, we do not provide “standard” code of the corresponding C++ files; if needed, it can be extracted from the `ctangle` generated files `*.cc` in `vnodelp/examples`.

```

#include <ostream>
#include "vnode.h"

using namespace std;
using namespace vnodelp;

template<typename var_type>
void Lorenz(int n, var_type*yp, const var_type*y, var_type t,
            void*param)
{
    interval sigma(10.0), rho(28.0);
    interval beta = interval(8.0)/3.0;

    yp[0] = sigma*(y[1]-y[0]);
    yp[1] = y[0]*(rho-y[2])-y[1];
    yp[2] = y[0]*y[1]-beta*y[2];
}

int main()
{
    const int n = 3;
    interval t = 0.0, tend = 20.0;
    iVector y(n);
    y[0] = 15.0;
    y[1] = 15.0;
    y[2] = 36.0;

    AD *ad= new FADBAD_AD(n, Lorenz, Lorenz);
    VNODE *Solver= new VNODE(ad);

    Solver->integrate(t, y, tend);
    if (!Solver->successful())
        cout<<"VNODE-LP could not reach t="<<t<<"<<tend<<endl;

    cout<<"Solution enclosure at t="<<t<<endl;
    printVector(y);

    return 0;
}

```

Figure 3.2. The “standard” C++ code of `basic.cc`

3.2 One-dimensional ODE

VNODE-LP is designed to be a general-purpose ODE interval solver. Nevertheless, it should handle scalar ODEs. This example illustrates how VNODE-LP deals with the simple problem

$$y' = -y, \quad y(0) = 1, \quad t_{\text{end}} = 20.$$

We write

```
32 <scalar ODE example 32> ≡
    template<typename var_type>
    void ScalarExample(int n, var_type *yp, const var_type *y, var_type t,
                       void *param)
    {
        yp[0] = -y[0];
    }
```

This code is used in chunk 35.

The initial condition and endpoint are set in

```
33 <set scalar ODE initial condition and endpoint 33> ≡
    const int n = 1;
    interval t = 0.0, tend = 20.0;
    iVector y(n);    /* number of state variables is 1 */
    y[0] = 1.0;
```

This code is used in chunk 35.

To create the necessary AD object, we call

```
34 <create scalar AD object 34> ≡
    AD *ad = new FADBAD_AD(n, ScalarExample, ScalarExample);
```

This code is used in chunk 35.

The main program is

```
35 <scalar.cc 35> ≡
#include <ostream>
#include "vnode.h"
using namespace std;
using namespace vnodelp;
<scalar ODE example 32>
int main()
{
    <set scalar ODE initial condition and endpoint 33>
    <create scalar AD object 34>
    <create a solver 23>
```

```

    <integrate (basic) 24>
    <check if success 25>
    <output results 26>
    return 0;
}

```

The output of this program is 0.000000002061153[6,7], which must enclose e^{-20} . Using MAPLE with a 30-digit computation, we obtain for e^{-20}

$$\underline{0.206115362243855782796594038016} \cdot 10^{-8}$$

(the digits that coincide are underlined), which is contained in the interval computed by VNODE-LP.

3.3 Time-dependent ODE

We show an example of integrating a time-dependent ODE. We choose the E1 problem from the DETEST test set [14]. The problem is

$$\begin{aligned}
 y_1' &= y_2 \\
 y_2' &= -\left(\frac{y_2}{t+1} + \left(1 - \frac{0.25}{(t+1)^2}\right)y_1\right) \\
 y_1(0) &= 0.6713967071418030, \\
 y_2(0) &= 0.09540051444747446, \\
 t_{\text{end}} &= 20.
 \end{aligned}$$

```

37 <DETEST E1 37> ≡
    template<typename var_type>
    void DETEST_E1(int n, var_type *yp, const var_type *y, var_type t, void
        *param)
    {
        var_type t1 = t + 1.0;
        yp[0] = y[1];
        yp[1] = -(y[1]/t1 + (1.0 - 0.25/(t1 * t1)) * y[0]);
    }

```

This code is used in chunk 40.

We store the initial condition as intervals containing the corresponding decimal values. The function *string_to_interval* converts a decimal string to a machine interval that contains the decimal value stored in the string.

```

38 <set E1 initial condition and endpoint 38> ≡
    const int n = 2;
    interval t = 0.0, tend = 20.0;
    iVector y(n);

```



```

y[0] = string_to_interval("0.6713967071418030");
y[1] = string_to_interval("0.09540051444747446");

```

This code is used in chunk 40.

To create an AD object, we call

```

39 <create E1 39> ≡
    AD *ad = new FADBAD_AD(n, DETEST_E1, DETEST_E1);

```

This code is used in chunk 40.

The main program is

```

40 <E1.cc 40> ≡
#include <ostream>
#include "vnode.h"
using namespace std;
using namespace vnodelp;
<DETEST E1 37>
int main()
{
    <set E1 initial condition and endpoint 38>
    <create E1 39>
    <create a solver 23>
    <integrate (basic) 24>
    <check if success 25>
    <output results 26>
    return 0;
}

```

The output of this program is

```

Solution enclosure at t = [20,20]
0.14567236007282[02,87]
-0.0988350019557[410,507]

```

We have also computed a numerical solution using MATLAB's `ode45`. The corresponding programs are in Figure 6.3, and the output is

```

0.14567235996177
-0.09883500182770

```

We have underlined the digits that coincide in the VNODE-LP and MATLAB output.

Since VNODE-LP includes all possible errors in its computation, including conversion errors in the input from decimal to binary, we have bounds on the true solution at $t = 20$. Using these bounds, we can check the accuracy of the numerical solution found by MATLAB's `ode45`.

3.4 Interval initial conditions

Suppose we want to compute bounds on the solution of the Lorenz problem for all

$$y(0) \in \begin{pmatrix} 15 + [-10^{-4}, 10^{-4}] \\ 15 + [-10^{-4}, 10^{-4}] \\ 36 + [-10^{-4}, 10^{-4}] \end{pmatrix}.$$

We set an interval initial condition by

```
42 <set interval initial condition and endpoint 42> ≡
    <set initial condition and endpoint 21>
    interval eps = interval(-1,1)/1 · 104;
    y[0] += eps;
    y[1] += eps;
    y[2] += eps;
```

This code is used in chunk 45.

Before presenting the rest of the code, we introduce some notation. We denote the radius of \mathbf{a} by

$$r(\mathbf{a}) = (\bar{\mathbf{a}} - \underline{\mathbf{a}})/2.$$

The midpoint of \mathbf{a} is denoted by

$$m(\mathbf{a}) = (\bar{\mathbf{a}} + \underline{\mathbf{a}})/2.$$

(In machine arithmetic, the true $\bar{\mathbf{a}} - \underline{\mathbf{a}}$ is rounded upward, and the true $(\bar{\mathbf{a}} + \underline{\mathbf{a}})/2$ is rounded to the nearest.) Radius and midpoint are defined componentwise for interval vectors.

Informally, the *global excess* at a point t^* is the overestimation in the computed bounds over the true solution set at t^* [23]. VNODE-LP computes an estimate, which is a nonnegative number, of the global excess in the computed bounds on each solution component, and also the max norm of these estimates.

When stepping in time from \mathbf{t}_0 to \mathbf{t}_{end} , VNODE-LP computes bounds on the solution at points that are machine numbers, except possibly at \mathbf{t}_{end} , which may be a machine interval with a nonzero radius. In this example, we output intermediate results during the integration. We tell the integrator to return after each step is completed by calling *setOneStep(on)*. This is convenient when we want to access intermediate results, for example, for plotting solutions, stepsize, etc. In the code below, we record such results in a file, which is later used by **gnuplot** to generate the plots in Figure 3.3.

```
43 <indicate single step 43> ≡
    Solver→setOneStep(on);
```

This code is used in chunks 44 and 61.

Now we integrate and record in the file `lorenzi.out` the midpoint and the radius of the computed bounds on $y_1(t)$ for each t selected by the solver. We also output an estimate on the *global excess*. The function `getGlobalExcess` returns the max norm of a vector with estimates on the global excess for each solution component. We exit the **while** loop below if this estimate exceeds 15 (This number is chosen so we can visualize the divergence of the computed bounds; see Figure 3.3(b).)

```
44 <integrate with interval initial condition 44> ≡
    <indicate single step 43>
    ofstream outFile("lorenzi.out", ios::out);
    while (t ≠ tend) {
        Solver→integrate(t, y, tend);
        if (Solver→successful() ∧ Solver→getGlobalExcess() ≤ 15.0) {
            outFile << midpoint(t) << "\t"
                << midpoint(y[0]) << "\t"
                << rad(y[0]) << "\t"
                << Solver→getGlobalExcess() << endl;
        }
        else break;
    }
    outFile.close();
```

This code is used in chunk 45.

The main program is

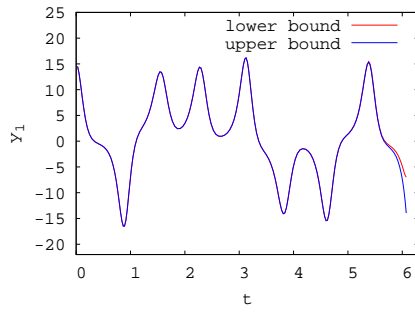
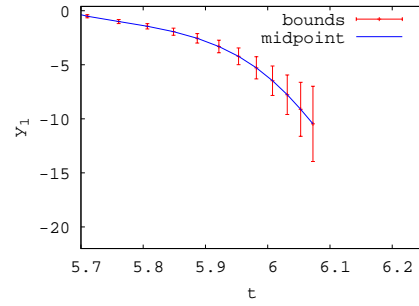
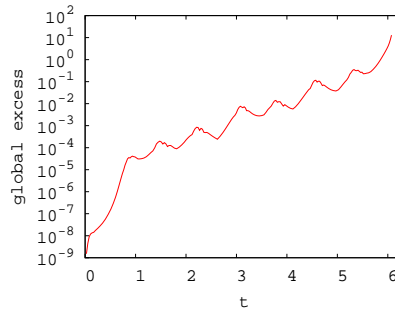
```
45 <integi.cc 45> ≡
#include <fstream>
#include "vnode.h"
using namespace std;
using namespace vnodelp;

<Lorenz 19>

int main()
{
    <set interval initial condition and endpoint 42>
    <create AD object 22>
    <create a solver 23>
    <integrate with interval initial condition 44>
    return 0;
}
```

In Figure 3.3(a) and (b), we plot the lower and upper bounds on y_1 . Their divergence is clearly seen in (b). In (c) we plot the logarithm of the estimate on the global excess versus t . On this problem, the bounds become too wide as t goes beyond ≈ 6.07 .

The `gnuplot` file employed to generate this plot is given in Figure 6.2.

(a) Bounds on $y_1(t)$ versus t (b) Bounds on $y_1(t)$ versus t ; “midpoint” lines connect the midpoints of the computed intervals(c) $\log_{10}(\text{global excess})$ versus t Figure 3.3. Plots generated using `integi.cc`

3.5 Producing intermediate results

We show how we can output enclosures on the solution at given points, for example, at $t = 0.1, 0.2, 0.3$, and $t = 10$. The decimal value 0.1 cannot be stored exactly as an IEEE floating-point number—we store 0.1 as an interval.

```

47 <integrate with intermediate output 47> ≡
    interval step = string_to_interval("0.1");
    tend = 0.0;
    for (int i = 1; i ≤ 3; i++) {
        tend += step;
        Solver→integrate(t, y, tend);
        <output results 26>
    }
    tend = 10;
    Solver→integrate(t, y, tend);
    <output results 26>

```

This code is used in chunk 48.

The main program is

```
48 <intermediate.cc 48> ≡
#include <iostream>
#include "vnode.h"
using namespace std;
using namespace vnodelp;
<Lorenz 19>
int main()
{
    <set initial condition and endpoint 21>
    <create AD object 22>
    <create a solver 23>
    <integrate with intermediate output 47>
    return 0;
}
```

The output is

```
Solution enclosure at t = 0.0999999999999999[9,11]
9.5199890775031[033,833]
1.172296185059[1790,3086]
36.286934318697[3267,4618]
```

```
Solution enclosure at t = 0.1999999999999999[9,11]
2.870955583196[2392,4120]
-1.446088378503[6309,8119]
27.476193552015[1645,3422]
```

```
Solution enclosure at t = 0.2999999999999999[9,11]
0.339508665582[1531,3150]
-0.925287772052[5888,8462]
20.901865632568[5157,7041]
```

```
Solution enclosure at t = [10,10]
-5.909806[3819893544,7254843408]
-11.34140[28468979033,34573302108]
9.08017[76297270730,80129492258]
```

In this output, 0.0999999999999999[9,11] is interpreted as the interval with left point

$$0.0999999999999999 + 00000000000000009 = 0.0999999999999999$$

and right point

$$0.09999999999999999 + 000000000000000011 = 1.00000000000000001.$$

Note. One has to be careful when interpreting the (decimal) output. For example consider the first output point for t , $\mathfrak{t} = 0.09999999999999999[9,11]$. VNODE-LP computes bounds for all t in the **binary interval containing 0.1**, but prints the rounded out decimal interval, which in general contains the binary one (but it is not necessarily the same).

3.6 ODE control

3.6.1 Passing data to an ODE

Suppose that we want to pass the constants to the Lorenz system. We can encapsulate them in the structure

```
51 <constants Lorenz 51> ≡
    struct LorenzConsts {
        interval beta;
        double rho, sigma;
    };

```

This code is used in chunk 58.

We can set σ , β , and ρ in the main programs as

```
52 <set ODE parameters 52> ≡
    LorenzConsts p;
    p.sigma = 10.0;
    p.beta = interval(8.0)/3.0;
    p.rho = 28.0;

```

This code is used in chunk 58.

We can access parameters for the ODE through the **void** pointer *param*. The user has to ensure that such parameters are properly stored and later extracted through *param*.

```
53 <passing parameters to Lorenz 53> ≡
    template<typename var_type>
    void Lorenz2(int n, var_type *yp, const var_type *y, var_type t,
        void *param)
    {
        LorenzConsts *p = (LorenzConsts *) param;
        interval beta = p->beta;
        double sigma = p->sigma;
        double rho = p->rho;
        yp[0] = sigma * (y[1] - y[0]);
        yp[1] = y[0] * (rho - y[2]) - y[1];
    }

```

```

    yp[2] = y[0] * y[1] - beta * y[2];
}

```

This code is used in chunk 58.

To pass parameters to the ODE, we create an AD object *with fourth parameter the address of p*:

```

54 <create problem object with parameters 54> ≡
    AD *ad = new FADBAD_AD(n, Lorenz2, Lorenz2, &p);

```

This code is used in chunk 58.

3.6.2 Integration with parameter change

We illustrate how to integrate with changing β . First, we integrate from t_0 to $tend$ and output $m(\mathbf{y}_j)$ into a file; \mathbf{y}_j is the computed enclosure at t_j .

```

55 <simple integration 55> ≡
    Solver~setOneStep(on);
    ofstream outFile1("odeparam1.out", ios::out);
    while (t ≠ tend) {
        Solver~integrate(t, y, tend);
        outFile1 << midpoint(y[0]) << "\t"
            << midpoint(y[1]) << "\t"
            << midpoint(y[2]) << endl;
    }
    outFile1.close();

```

This code is used in chunk 58.

Now, we

1. integrate from t_0 to $tend/2$,
2. change β , and
3. integrate to $tend$.

Before calling *integrate* again, we call *setFirstEntry*. This call ensures that certain internal data structures are initialized. If *setFirstEntry* is not called, *integrate* would use data corresponding to the last computed solution from the most recent call to *integrate*.

```

56 <integrate from t to tend/2 56> ≡
    t = 0.0;
    y[0] = 15;
    y[1] = 15;
    y[2] = 36;
    interval tend2 = tend/2.0;
    Solver~setFirstEntry();
    while (t ≠ tend2)
        Solver~integrate(t, y, tend2);

```

This code is used in chunk 57.

When changing β , the integration continues with the last computed solution, which is computed with the previous value for β . After a parameter is changed, *ad-`eval`(&p)* must be called. This ensures that some internal data structures are updated, to reflect the change in the ODE.

```
57 <integrate with resetting constants 57> ≡
    <integrate from t to tend/2 56>
    p.beta = 5.0;      /* change  $\beta$  */
    ad-eval(&p);      /* must be called to update internal data structures */
    ofstream outFile2("odeparam2.out", ios::out);
    while (t ≠ tend) {
        Solver→integrate(t, y, tend);
        outFile2 << midpoint(y[0]) << "\t"
            << midpoint(y[1]) << "\t"
            << midpoint(y[2]) << endl;
    }
    outFile2.close();
```

This code is used in chunk 58.

The main program is

```
58 <odeparam.cc 58> ≡
#include <fstream>
#include "vnode.h"
using namespace std;
using namespace vnodelp;
<constants Lorenz 51>
<passing parameters to Lorenz 53>
int main()
{
    <set initial condition and endpoint 21>
    <set ODE parameters 52>
    <create problem object with parameters 54>
    <create a solver 23>
    <simple integration 55>
    <integrate with resetting constants 57>
    return 0;
}
```

In Figure 3.4, we plot $m(\mathbf{y}_j)$ in (y_1, y_2, y_3) coordinates corresponding to $\beta = 8/3$ and $\beta = 5$. The `gnuplot` file for producing this plot is in Figure 6.4.

3.7 Integration control

We start by introducing various facts related to the integration process of VNODE-LP. Then we show ways of controlling it.

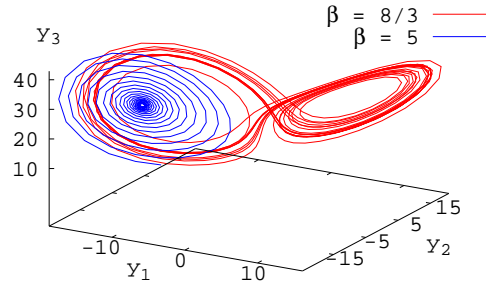


Figure 3.4. *Midpoints of the computed bounds with $\beta = 8/3$ from 0 to 20; and with $\beta = 8/3$ changed to 5 at $t = 10$*

The method implemented in VNODE-LP is a one-step method based on Taylor series and the Hermite-Obreschkoff scheme [23]. These are high-order methods, where typical values for the order, denoted here by p , can be in the range of 20 to 30; see Section 3.8.

VNODE-LP steps in time from t_0 to t_{end} , where the j th integration step, $j \geq 1$, is from t_{j-1} to t_j . The associated stepsize is $h_j = t_j - t_{j-1}$.² Each of these t_j is a representable machine number, except t_0 and t_{end} , which can be machine intervals \mathbf{t}_0 and \mathbf{t}_{end} containing the true t_0 and t_{end} , respectively. For simplicity of the exposition, we assume point values t_0 and t_{end} .

In addition to computing bounds \mathbf{y}_j on the solution at each t_j , VNODE-LP also computes bounds for all $t \in \mathbf{T}_j = [t_{j-1}, t_j]$, or $\mathbf{T}_j = [t_j, t_{j-1}]$ if the integration is in negative direction. We denote such bounds by $\tilde{\mathbf{y}}_j$, and refer to them as a priori bounds [23]; we shall also refer to \mathbf{y}_j as tight bounds.

On the first integration step, VNODE-LP determines initial stepsize and the magnitude of the smallest stepsize that is allowed, h_{\min} . Then, on each step, VNODE-LP automatically selects a stepsize subject to absolute and relative error tolerances, atol and rtol. If the selected stepsize h_j is such that $|h_j| < h_{\min}$, the integration cannot continue, and VNODE-LP returns.

The following parameters can be changed by the user.

parameter	default value
p	20
atol	10^{-12}
rtol	10^{-12}
h_{\min}	computed by VNODE-LP

The functions for changing them are given in Section 4.4. The order p does not

²VNODE-LP selects h_j and then finds t_j , where in computer arithmetic the true $t_{j-1} + h_j$ is rounded toward zero when computing t_j .

vary during an integration. By default $p = 20$, but its value can be changed at the beginning of an integration. As pointed out earlier, p must be between 3 and the value set to `MAX_ORDER` in the `makefile` for building the library; cf. Figures 2.2 and 2.3. If the user sets a value for h_{\min} , then this value is used by *integrate*.

If we wish to set, for example,

$$\begin{aligned} \text{rtol} &= \text{atol} = 10^{-14}, \\ p &= 40, \quad \text{and} \\ h_{\min} &= 10^{-5}, \end{aligned}$$

we proceed with (`cweave` typesets `1e-14` as $1 \cdot 10^{-14}$)

```
60 <set control data for the solver 60> ≡
    Solver~setTols(1 · 10-14, 1 · 10-14);
    Solver~setOrder(40);
    Solver~setHmin(1 · 10-5);
```

This code is used in chunk 61.

We write the main program

```
61 <integctrl.cc 61> ≡
#include <fstream>
#include "vnode.h"
using namespace std;
using namespace vnodelp;
<Lorenz 19>
int main()
{
    <set initial condition and endpoint 21>
    <create AD object 22>
    <create a solver 23>
    <set control data for the solver 60>
    <open files 62>
    <indicate single step 43>
    <output initial condition 64>
    while (t ≠ tend) {
        Solver~integrate(t, y, tend);
        <output solution 65>
    }
    <close files 63>
    return 0;
}
```

Writing into files

We output data into three files, `lorenz.tight`, `lorenz.apriori` and `lorenz.step`. In `lorenz.tight`, each line contains (rounded in decimal to output precision)

$$m(t_j) \quad \underline{y}_{1,j} \quad \overline{y}_{1,j} \quad w(y_{1,j})$$

where the subscripts $1, j$ refers to the j th computed solution for component y_1 ; $w(\mathbf{y}_{1,j}) = 2r(\mathbf{y}_{1,j})$ is the width, or diameter of $\mathbf{y}_{1,j}$. This width can be viewed as the *global excess* in the computed $\mathbf{y}_{1,j}$.

In `lorenz.apriori`, we output the a priori enclosures and the corresponding time intervals in a form suitable for `gnuplot` to produce boxes denoting these a priori bounds; see Figure 3.5(b). The function `getAprioriEncl` returns $\tilde{\mathbf{y}}_j$, and `getT` returns \mathbf{T}_j .

In `lorenz.step`, each line is

$$m(\mathbf{t}_j) \quad h_j$$

The function `getStepsize` returns h_j .

```
62 <open files 62> ≡
    ofstream outFile1("lorenz.tight", ios::out);
    ofstream outFile2("lorenz.step", ios::out);
    ofstream outFile3("lorenz.apriori", ios::out);
```

This code is used in chunk 61.

```
63 <close files 63> ≡
    outFile1.close();
    outFile2.close();
    outFile3.close();
```

This code is used in chunk 61.

In the code below, `inf` returns the left point of an interval, and `sup` returns the right point of an interval. (The output goes through C++'s stream output, so the endpoints are not rounded outward.)

```
64 <output initial condition 64> ≡
    outFile1 << midpoint(t) << "\t"
        << inf(y[0]) << "\t" << sup(y[0]) << "\t" << width(y[0]) << endl;
```

This code is used in chunk 61.

```
65 <output solution 65> ≡
    outFile1 << midpoint(t) << "\t"
        << inf(y[0]) << "\t" << sup(y[0]) << "\t" << width(y[0]) << endl;
    outFile2 << midpoint(t) << "\t" << Solver->getStepsize() << endl;
    iVector Y = Solver->getAprioriEncl();
    interval Tj = Solver->getT();
    outFile3 << inf(Tj) << "\t" << inf(Y[0]) << endl;
    outFile3 << inf(Tj) << "\t" << sup(Y[0]) << endl << endl;
    outFile3 << sup(Tj) << "\t" << inf(Y[0]) << endl;
    outFile3 << sup(Tj) << "\t" << sup(Y[0]) << endl << endl;
    outFile3 << inf(Tj) << "\t" << inf(Y[0]) << endl;
    outFile3 << sup(Tj) << "\t" << inf(Y[0]) << endl << endl;
```

```
outFile3 << inf(Tj) << "\t" << sup(Y[0]) << endl;
outFile3 << sup(Tj) << "\t" << sup(Y[0]) << endl << endl;
```

This code is used in chunk 61.

Plots

To visualize the results in these files, we produce the plots in Figure 3.5. In Figure 3.5(a) and (b), the upper and lower tight bounds cannot be distinguished when plotted. In (b), the a priori bounds are shown as boxes. The `gnuplot` file for generating this figure is displayed in Figure 6.5

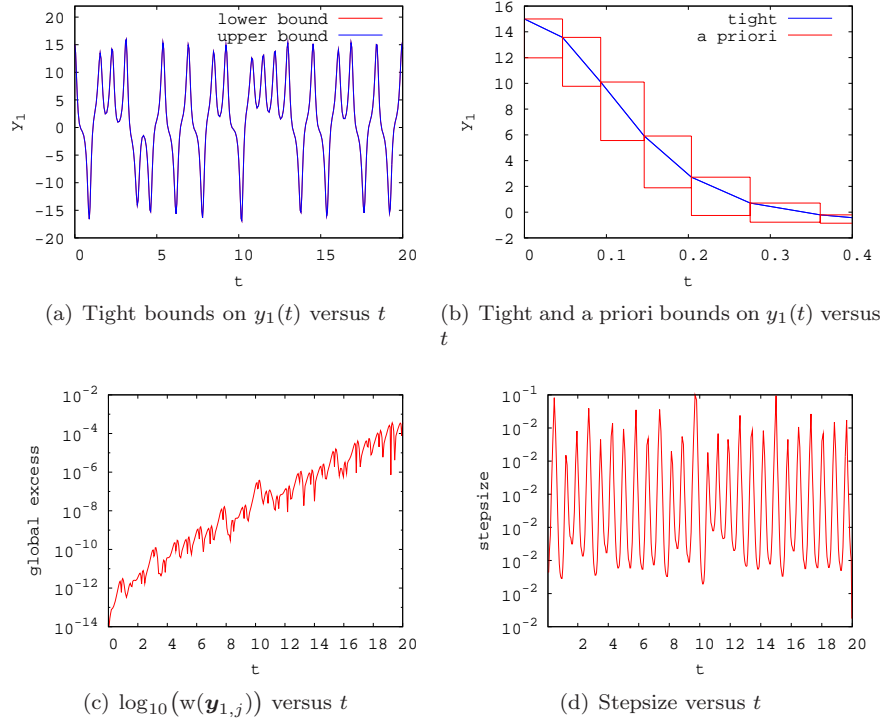


Figure 3.5. Plots generated using `integctrl.cc`

3.8 Work versus order

We show in Figure 3.6 how the computing time depends on the order for various tolerance when integrating the Lorenz system. We consider orders $p = 5, 6, \dots, 40$ and tolerances $\text{atol} = \text{rtol} = 10^{-7}, 10^{-8}, \dots, 10^{-13}$.

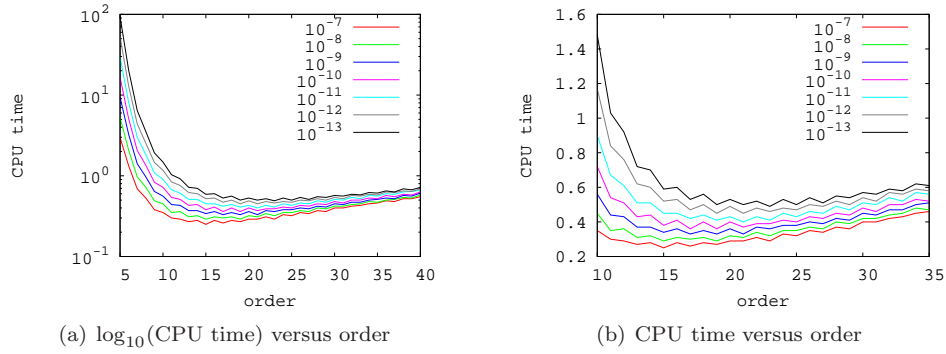


Figure 3.6. *Plots generated using orderstudy.cc*

The gnuplot file employed to generate these plots is given in Figure 6.6.

The main program is

```
68 <main program for order study 68> ≡
    static double tol[] = {1 · 10-7, 1 · 10-8, 1 · 10-9, 1 · 10-10, 1 · 10-11, 1 · 10-12,
        1 · 10-13};
    int main()
    {
        const int n = 3;
        <create AD object 22>
        <create a solver 23>
        iVector y(n);
        for (int i = 0; i < 7; i++)
        {
            Solver→setTols(tol[i]);
            <create file name 69>
            ofstream outFile(file_name.c_str(), ios::out);
            cout << "tol=" << tol[i] <<
                "writing into" << file_name << "... " << endl;
            for (int p = 5; p ≤ 40; p++)
            {
                Solver→setOrder(p);
                interval t = 0.0, tend = 10.0;
```

```

        y[0] = 15.0;
        y[1] = 15.0;
        y[2] = 36.0;
        Solver->setFirstEntry();

        double time = getTime();

        <integrate (basic) 24>
        <check if success 25>
        time = getTotalTime(time, getTime());
        outFile << p << "\t" << time << endl;
    }
    outFile.close();
}
return 0;
}

```

This code is used in chunk 70.

We create a file name for each tolerance value by

```

69 <create file name 69> ≡
    string prefix("order");

    std::stringstream num(std::stringstream::out);

    num << tol[i];

    string file_name = prefix + num.str() + ".out";

```

This code is used in chunk 68.

We store all this into

```

70 <orderstudy.cc 70> ≡
    #include <fstream>
    #include <sstream>
    #include <string>
    #include <cstdlib>
    #include "vnode.h"
    using namespace std;

    using namespace vnodelp;

    <Lorenz 19>
    <main program for order study 68>

```

3.9 Work versus problem size

We investigate how the computing time depends on the size of the problem. We consider the DETEST problem C3 [14]

$$y' = \begin{pmatrix} -2 & 1 & 0 & 0 & \cdots & 0 \\ 1 & -2 & 1 & 0 & \cdots & 0 \\ 0 & 1 & -2 & 1 & \cdots & 0 \\ & & & \vdots & & \\ 0 & \cdots & & 1 & -2 & 1 \\ 0 & \cdots & & 0 & 1 & -2 \end{pmatrix} y \quad (3.2)$$

with $y(0) = (1, 0, \dots, 0)^T$. We integrate with problem sizes $n = 40, 60, \dots, 200$ for $t = 0$ to $t = 5$.

The C++ program is

```

71 <detest_c3.cc 71> ≡
#include <ostream>
#include "vnode.h"
using namespace std;
using namespace vnodelp;
template<typename var_type>
void DETEST_C3(int n, var_type *yp, const var_type *y, var_type t,
               void *param)
{
    yp[0] = -2.0 * y[0] + y[1];
    for (int i = 1; i < n - 1; i++) {
        yp[i] = y[i - 1] - 2.0 * y[i] + y[i + 1];
    }
    yp[n - 1] = y[n - 2] - 2.0 * y[n - 1];
}
int main()
{
    for (int n = 40; n ≤ 200; n += 20) {
        cout << n;

        interval t = 0.0, tend = 5;
        iVector y(n);

        for (int i = 0; i < n; i++) y[i] = 0;
        y[0] = 1;

        AD *ad = new FADBAD_AD(n, DETEST_C3, DETEST_C3);
        VNODE *Solver = new VNODE(ad);
        double time_start = getTime();

        Solver~integrate(t, y, tend);

        double time_end = getTime();
    }
}

```

```

    cout << "___" << getTotalTime(time_start,
        time_end) << "___" << Solver->steps << endl;
    delete Solver;
    delete ad;
}
return 0;
}

```

In Figure 3.7(a) and (b), we display the CPU time versus n . The `gnuplot` file for generating this figure is in Figure 6.7. It is not difficult to see that the computing

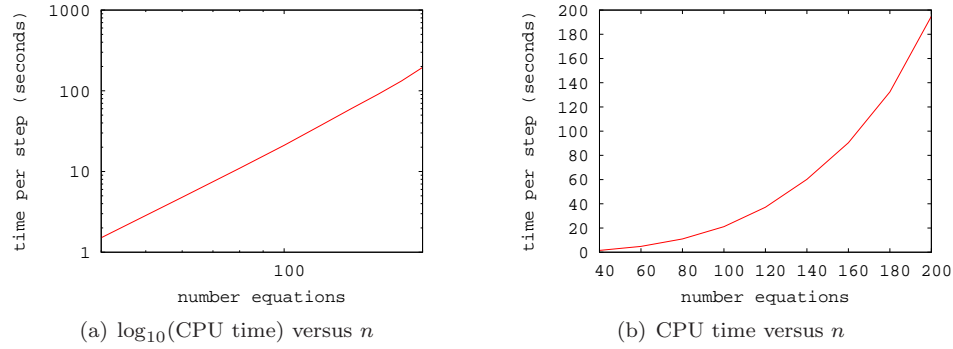


Figure 3.7. CPU time versus n for Problem 3.2. VNODE-LP takes 8 steps for each n .

time grows like n^3 .

3.10 Stepsize behavior

We consider the orbit problem (see for example [4])

$$\begin{aligned}
 y_1'' &= y_1 + 2y_2' - \hat{\mu} \frac{y_1 + \mu}{D_1} - \mu \frac{y_1 - \hat{\mu}}{D_2}, \\
 y_2'' &= y_2 - 2y_1' - \hat{\mu} \frac{y_2}{D_1} - \mu \frac{y_2}{D_2},
 \end{aligned} \tag{3.3}$$

where

$$\mu = 0.012277471, \quad \hat{\mu} = 1 - \mu, \tag{3.4}$$

$$D_1 = ((y_1 + \mu)^2 + y_2^2)^{3/2}, \quad \text{and} \tag{3.5}$$

$$D_2 = ((y_1 - \hat{\mu})^2 + y_2^2)^{3/2}. \tag{3.6}$$

We integrate this problem with

$$\begin{aligned} y_1(0) &= 0.994, \\ y_2(0) &= 0, \\ y_1'(0) &= 0, \\ y_2'(0) &= -2.00158510637908252240537862224, \end{aligned} \tag{3.7}$$

and $t_{\text{end}} = 35$, which corresponds to slightly more than two periods.

In Figure 3.8, we plot y_2 versus y_1 and the stepsize versus t . (The gnuplot file is in Figure 6.8.)

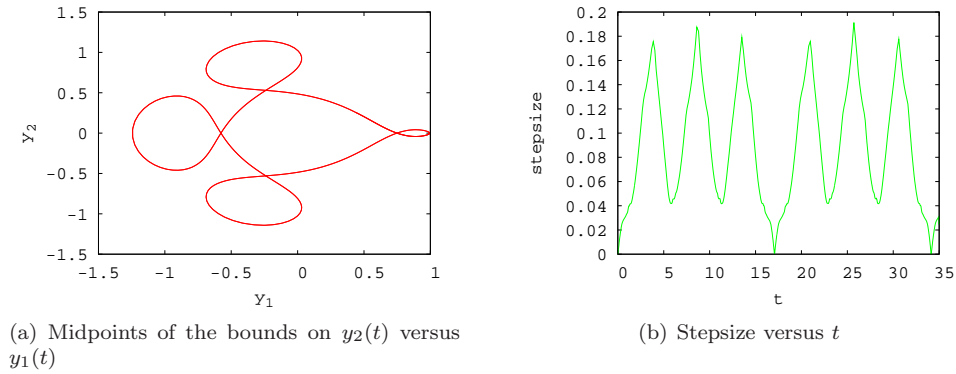


Figure 3.8. *Plots generated using orbit.cc*

The C++ program follows.

```
73 <orbit.cc 73> ≡
#include <fstream>
#include <sstream>
#include <string>
#include <cstdlib>
#include "vnode.h"
using namespace std;
using namespace vnodelp;
template<typename var_type>
void Orbit(int n, var_type *yp, const var_type *y, var_type t,
           void *param)
{
    interval mu = string_to_interval("0.012277471");
    interval mu_h = 1.0 - mu;
    var_type D1 = pow(sqr(y[0] + mu) + sqr(y[1]), interval(1.5));
    var_type D2 = pow(sqr(y[0] - mu_h) + sqr(y[1]), interval(1.5));
    yp[0] = y[2];
    yp[1] = y[3];
}
```

```

    yp[2] = y[0] + 2.0 * y[3] - mu_h * (y[0] + mu)/D1 - mu * (y[0] - mu_h)/D2;
    yp[3] = y[1] - 2.0 * y[2] - mu_h * y[1]/D1 - mu * y[1]/D2;
}
int main()
{
    const int n = 4;
    iVector y(n);
    y[0] = string_to_interval("0.994");
    y[1] = 0;
    y[2] = 0;
    y[3] = string_to_interval("-2.00158510637908252240537862224");
    interval t = 0.0, tend = 35;
    AD *ad = new FADBAD_AD(n, Orbit, Orbit);
    VNODE *Solver = new VNODE(ad);
    ofstream outFileSol("orbit_sol.out", ios::out);
    ofstream outFileStep("orbit_step.out", ios::out);
    outFileSol << midpoint(y[0]) << "\t" << midpoint(y[1]) << endl;
    Solver->setOneStep(on);
    while (t != tend) {
        Solver->integrate(t, y, tend);
        outFileSol << midpoint(y[0]) << "\t" << midpoint(y[1]) << endl;
        outFileStep << midpoint(t) << "\t"
            << Solver->getStepsize() << endl;
    }
    outFileSol.close();
    outFileStep.close();
    return 0;
}

```

3.11 Stiff problems

We illustrate how VNODE-LP behaves when integrating a stiff problem. We integrate Van der Pol's equation (written as a first-order system)

$$\begin{aligned}
 y_1' &= y_2 \\
 y_2' &= \mu(1 - y_1^2)y_2 - y_1
 \end{aligned} \tag{3.8}$$

with

$$y(0) = (2, 0)^T, \quad t_{\text{end}} = 200. \tag{3.9}$$

We perform integrations with $\mu = 10, 10^2, 10^3$, and 10^4 . In Table 3.1, we show the number of steps and CPU time used by VNODE-LP. In Figure 3.9, we plot the stepsizes versus t . As can be seen from this table and figure, VNODE-LP is not efficient when this problem becomes stiff. In general, VNODE-LP works well on

μ	steps	time (secs)
10^1	2377	2.4
10^2	11697	11.6
10^3	126459	124.4
10^4	1180844	1182.7

Table 3.1. Number of steps and CPU time used by VNODE-LP on (3.8–3.9)

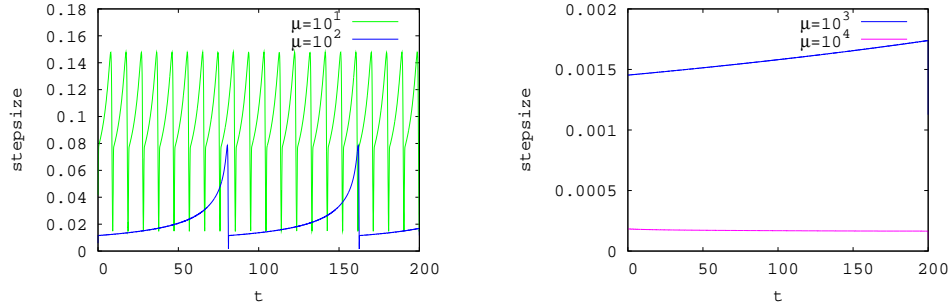


Figure 3.9. Stepsize versus t on (3.8–3.9) for $\mu = 10, 10^2, 10^3, 10^4$

non-stiff and mildly stiff problems. A more detailed study regarding stiff problems can be found in [23].

The program used to produce the numerical results for this problem follows. The gnuplot file for generating the plots is in Figure 6.9.

```

74 <vanderpol.cc 74> ≡
#include <fstream>
#include <iomanip>
#include <sstream>
#include <string>
#include "vnode.h"
using namespace std;
using namespace vnodelp;
template<typename var_type>
void VDP(int n, var_type *yp, const var_type *y, var_type t,
        void *param)
{
    double *MU = (double *) param;
    yp[0] = y[1];
    yp[1] = (*MU) * (1 - sqrt(y[0])) * y[1] - y[0];
}

int main()
{

```

```

const int n = 2;
double MU = 10.0;
AD *ad = new FADBAD_AD(n, VDP, VDP, &MU);
VNODE *Solver = new VNODE(ad);
/* file for storing Table data */
ofstream outSteps("vdp_nosteps.out", ios::out);
outSteps << fixed << showpoint << setprecision(1);
for (int i = 1; i ≤ 4; i++) {
    /* file for storing stepsizes */
    string prefix("vdp_step");
    std::stringstream num(std::stringstream::out);
    num << i;
    string file_name = prefix + num.str() + ".out";
    ofstream outStepSizes(file_name.c_str(), ios::out);
    cout << "└─MU=└─" << MU <<
        "└─writing└─into└─" << file_name << "... " << endl;
    Solver→setFirstEntry();
    Solver→setOneStep(on);
    interval t = 0, tend = 200;
    iVector y(n);
    y[0] = 2.0, y[1] = 0;
    interval t_prev = t;
    double time = getTime();
    while (t ≠ tend) {
        Solver→integrate(t, y, tend);
        if (midpoint(t - t_prev) ≥ 0.01 ∨ t ≡ tend) {
            outStepSizes << midpoint(t) << "\t" << Solver→getStepsize() << endl;
            t_prev = t;
        }
    }
    outStepSizes.close();
    time = getTotalTime(time, getTime());
    outSteps << "$10^{ " << i << " }$" << "\t└─&└─"
        << Solver→getNoSteps() << "\t└─&└─"
        << time << "\t" << "\\\\" << endl;
    MU *= 10.0;
    ad→eval(&MU);
}
outSteps.close();
return 0;
}

```

Chapter 4

Interface

First, we list the data types used by VNODE-LP that are of interest to the user. Then we list and describe briefly the public functions of the VNODE-LP solver.

4.1 Interval data type

If PROFIL/BIAS [16] is employed, **interval** is defined as

```
76 < interval data type (PROFIL) 76 > ≡  
    typedef INTERVAL interval;
```

This code is used in chunk 115.

If FILIB++ [19] is employed, **interval** is defined as

```
77 < interval data type (FILIB++) 77 > ≡  
    typedef filib::interval<double> interval;
```

This code is used in chunk 115.

VNODE-LP does not call the functions of these packages directly. Instead, it implements wrapper functions; see below. To build VNODE-LP on a new IA package, only these functions need to be implemented using this new package.

4.2 Wrapper functions

In the descriptions that follow, ***a***, ***b***, and ***c*** denote corresponding intervals in mathematical notation. We assume that the endpoints of the input intervals are representable machine numbers and no exceptions occur when a result is computed.

double *inf*(const interval &*a*)
returns \underline{a} .

double *sup*(const interval &*a*)
returns \overline{a} .

double *midpoint*(**const interval** &a)
 returns $(\underline{a} + \overline{a})/2$ rounded to the nearest.

double *width*(**const interval** &a)
 returns $\overline{a} - \underline{a}$ rounded to $+\infty$.

double *mag*(**const interval** &a)
 returns $\max\{|\underline{a}|, |\overline{a}|\}$.

bool *subsetq*(**const interval** &a, **const interval** &b)
 returns *true* if $a \subseteq b$; *false* otherwise.

bool *interior*(**const interval** &a, **const interval** &b)
 returns *true* if a is in the interior of b ; *false* otherwise.

bool *disjoint*(**const interval** &a, **const interval** &b)
 returns *true* if $a \cap b = \emptyset$; *false* otherwise.

bool *intersect*(**interval** &c, **const interval** &a, **const interval** &b)
 returns *true* if $a \cap b \neq \emptyset$ and sets $c = a \cap b$; *false* if $a \cap b = \emptyset$ and leaves the input c unchanged.

interval *pi*()
 returns an interval containing π .

interval *pow*(**const interval** &a, **const interval** &b)
 returns an interval containing $\{x^y \mid x \in a, y \in b\}$.

interval *pow*(**const interval** &a, **int** k)
 returns an interval containing $\{x^k \mid x \in a\}$.

Each of the functions e that follows returns $\{e(x) \mid x \in a\}$.

interval *exp*(**const interval** &a)
interval *log*(**const interval** &a)
interval *sqr*(**const interval** &a)
interval *sqrt*(**const interval** &a)
interval *sin*(**const interval** &a)
interval *cos*(**const interval** &a)
interval *tan*(**const interval** &a)
interval *asin*(**const interval** &a)
interval *acos*(**const interval** &a)
interval *atan*(**const interval** &a)

Finally,

interval *string_to_interval*(**const char** *s)
 returns an interval that contains the decimal number that is stored in the character string input.

4.3 Interval vector

VNODE-LP uses interval vector, **iVector**, defined as

```
80 <interval vector 80> ≡
#include <vector>
#include "vnodeinterval.h"
using namespace std;
using namespace v_bias;

typedef vector < interval > iVector;
```

This code is used in chunk 122.

4.4 Solver's public functions

The present solver is implemented by the class **VNODE**. We explain briefly its constructor and public member functions.

4.4.1 Constructor

VNODE(**AD** *ad)
 constructs a **VNODE** object. Here *ad* is a pointer to an object of a class derived from the **AD** class; see Subsection 4.5. Currently, there is only one such class, **FADBAD_AD**, and *ad* is a pointer to an object of this class.

4.4.2 Integrator

void *integrate*(**interval** &t, **iVector** &y, **interval** tend)

By default, *integrate*(*t*, *y*, *tend*) tries to compute an enclosure of the solution to an ODE problem at *tend*. If successful, *y* contains such an enclosure at *t = tend*. If an integration is not successful, *y* is an enclosure of the computed solution at *t ≠ tend*.

The initial and end points, t_0 and t_{end} , can be stored as intervals containing their true values. Normally, the corresponding interval for t_0 [resp. t_{end}] should be of the width of (at most) a few machine numbers. If we denote these intervals by \mathbf{t}_0 and \mathbf{t}_{end} , *integrate* requires that $\mathbf{t}_0 \cap \mathbf{t}_{\text{end}} = \emptyset$. If $\mathbf{t}_0 \cap \mathbf{t}_{\text{end}} \neq \emptyset$, *integrate* returns without performing an integration.

4.4.3 Set functions

void *setTols*(**double** *a*, **double** *r* = 0)

sets *atol* to *a* and *rtol* to *r*. The latter parameter has a default value of 0. For example, *setTols*($1 \cdot 10^{-10}$) sets *atol* = 10^{-10} and *rtol* = 0.

The code roughly controls the “drift” away from the true solution at each integration step to be of size $\text{atol} + \|y\|_{\infty} \cdot \text{rtol}$, where $\|y\|$ is a measure of the size of current solution. This drift is accounted for by the enclosure, thus the size of the enclosure at the end point being roughly proportional to $\text{atol} + \|y\|_{\infty} \cdot \text{rtol}$.

Default values are *atol* = 10^{-12} and *rtol* = 10^{-12} .

void *setOrder*(**int** *p*)

sets the order to *p*. It must be between 3 and the value set to `MAX_ORDER` in the `makefile` for building the library; cf. Figures 2.2 and 2.3. If *p* is not in this range, *integrate* returns without performing an integration.

Experience suggests that order in the range of about 20 to 30 results in efficient integration.

Default value is 20.

void *setHmin*(**double** *h*)

sets the value of the magnitude of minimum stepsize allowed to *h*. If minimum stepsize is not set, or *setHmin*(0) is called, *integrate* computes such.

void *setFirstEntry*()

indicates to *integrate* that this is a first entry into it. If this function is called before *integrate*, the latter will perform various initializations before time-stepping. When *integrate* is called for the first time, *setFirstEntry*() need not be called before *integrate*.

void *setOneStep*(**stepAction** *action*)

tells the integrator whether to stop, *action* \equiv *on* or continue, *action* \equiv *off*, after each step it takes. If *setOneStep*(*on*) is called before *integrate*, the latter will return after each step it takes. To turn off this feature, call *setOneStep*(*off*). In this case, if *integrate* is re-entered, it will not stop after each step it takes (except the last one).

4.4.4 Get functions

bool *successful*() **const**

returns *true* if an integration is successful and *false* otherwise. If *integrate* has not been called, *successful*() returns *true* by default.

int *getMaxOrder*() returns the maximum order allowed in an integration. This is the value set in `MAX_ORDER` in the configuration file.

double *getStepsize*() **const**

returns the value of the most recent stepsize.

double *getNoSteps()* **const**
 returns the number of successful steps taken by VNODE-LP.

const iVector *&getAprioriEncl()* **const**
 returns the last computed $\tilde{\mathbf{y}}_j$.

const interval *&getT()* **const**
 returns the last computed \mathbf{T}_j .

double *getGlobalExcess()* **const**
 returns an estimate of the global excess in the most recent computed enclosure.

double *getGlobalExcess(int i)* **const**
 returns an estimate of the global excess in the i th component of the most recent computed enclosure.

4.5 Constructing an AD object

VNODE-LP computes Taylor coefficients for the ODE and its variational equation. The user has to construct an object for computing such coefficients by

```
AD *ad = new FADBAD_AD(n, function_name, function_name)
or
AD *ad = new FADBAD_AD(n, function_name, function_name, param)
```

Here n is the size of the problem, *function_name* is the name of the template function, as described earlier, and *param* is a pointer to parameters that need to be passed to the ODE problem. After a parameter is changed, **void** *eval(void *p)* must be called to update internal structures in the **AD** object.

4.6 Some helpful functions

template<class **T**> **void** *printVector(const T &v, const char *s = 0)*
 prints the components of a vector v on the standard output. If the second parameter is given, *printVector* prints it before the content of v .

double *getTime()*
 returns the current time measured as user time.

double *getTotalTime(double start_time, double end_time)*
 returns $end_time - start_time$ rounded to the nearest.

Chapter 5

Testing

The code for the test cases is located in subdirectory `tests`. We give a brief description of each test.

5.1 General tests

File `test0.cc`

With the tests in this file, we check if the functions described in Subsection 4.2 compile and execute.

File `test01.cc`

We check if the elementary functions FADBAD++ uses compile and execute.

5.2 Linear problems

5.2.1 Constant coefficient problems

We consider

$$\begin{aligned}y_1' &= y_2 \\ y_2' &= -y_1\end{aligned}\tag{5.1}$$

with

$$y(0) = (1, 1)^T.\tag{5.2}$$

The true solution is

$$y(t) = \begin{pmatrix} \cos(t) & \sin(t) \\ -\sin(t) & \cos(t) \end{pmatrix} y(0).\tag{5.3}$$

File test1.cc

We integrate (5.1, 5.2) from $t = 0$ to $t_{\text{end}} = 10000$. At each integration point t_j , selected by the solver, we evaluate (5.3) in interval arithmetic and check if the resulting enclosure intersects with the computed bounds. This test succeeds if they intersect for all t_j .

File test2.cc

We integrate (5.1, 5.2) and check that at $t = 2k\pi$, for $k = 2, 4, \dots, 1000$, $y(0)$ is contained in the computed bounds. If so, this test is successful.

File test3.cc

We integrate (5.1, 5.2), but with $t_{\text{end}} = -10000$. This test is successful if the true solution (5.3), evaluated in interval arithmetic, and the computed bounds intersect at each point selected by the solver.

File test4.cc

We integrate (5.1, 5.2) and check that at $-2k\pi$, for $k = 2, 4, \dots, 1000$, $y(0)$ is contained in the computed bounds. If so, this test is successful.

File test5.cc

We consider

$$\begin{aligned} y_1' &= y_1 - 2y_2 \\ y_2' &= 3y_1 - 4y_2 \end{aligned} \tag{5.4}$$

with

$$y(0) = (1, -1)^T. \tag{5.5}$$

The true solution is

$$y(t) = \begin{cases} 5e^{-t} - 4e^{-2t} \\ 5e^{-t} - 6e^{-2t} \end{cases}. \tag{5.6}$$

We integrate (5.4, 5.5) and check that, at each point selected by the solver, the true (evaluated in interval arithmetic (5.6)) and computed solutions intersect. If they do, we accept this test as successful.

5.2.2 Time-dependent problems**File test6.cc**

The problem is

$$\begin{aligned}
y'_1 &= \sin(t+10)y_1 - 2y_2 - y_3 - \cos(t) \\
y'_2 &= 3y_1 - 4\cos(t^2)y_2 - \cos(t) \\
y'_3 &= e^{-t^2}y_1 - e^{-t^2}y_2 - \sin(t) \\
y(0) &\in ([0, 5], [-2, 6], [5, 12])^T, \quad t_{\text{end}} = 20.
\end{aligned}$$

We compute an enclosure at $t_{\text{end}} = 20$ with the above initial condition set. Then we compute bounds at t_{end} for each corner of the initial box. These bounds must intersect with the enclosure resulting from $([0, 5], [-2, 6], [5, 12])^T$.

We have also computed (accurate) approximate solutions using MAPLE for each corner of the initial box. At t_{end} , each approximate solution must be inside the bounds computed with the same corner point.

5.3 Nonlinear problems

File test_n1.cc

We integrate the Lorenz system with $y(0) = (15, 15, 36)^T$ from 0 to 1. Denote the enclosure at $t = 1$ by $\mathbf{y}_{(1)}$. We have also computed an approximate solution with MAPLE. Denote it by $\hat{\mathbf{y}}_{(1)}$. We check first if

$$\hat{\mathbf{y}}_{(1)} \in \mathbf{y}_{(1)}. \quad (5.7)$$

Then we integrate this system with $y(1) \in \mathbf{y}_{(1)}$ and $t_{\text{end}} = 0$. Denote the computed enclosure at $t = 0$ by $\mathbf{y}_{(0)}$. We check if

$$y(0) \in \mathbf{y}_{(0)}. \quad (5.8)$$

Finally, we integrate with $y(0) \in \mathbf{y}_{(0)}$ and $t_{\text{end}} = 1$. Denote the computed enclosure by $\mathbf{y}_{(1)}^*$. We check if

$$\hat{\mathbf{y}}_{(1)} \in \mathbf{y}_{(1)}^* \quad \text{and} \quad \mathbf{y}_{(1)}^* \cap \mathbf{y}_{(1)} \neq \emptyset. \quad (5.9)$$

This test is successful if of (5.7), (5.8), and (5.9) hold.

File test_n2.cc

We integrate a three-body problem from 0 to 1 and then from 1 to 0. The initial condition at $t = 0$ must be contained in the computed bounds.

File test_n3.cc

The same three-body problem is integrated from 0 to 1 with orders from 10 to the value set in `MAX_ORDER`. If all the computed enclosures with these orders intersect, we accept this test as successful.

File test_n4.cc

We integrate [20]

$$y'' + cy' + \sin(y) = b \cos(t),$$

written as a first-order system

$$\begin{aligned} y_1' &= y_2 \\ y_2' &= b \cos(t) - cy_2 - \sin(y_1) \end{aligned}$$

with

$$y(0) \in ([0, 0], [1.9999, 2.0001])^T, \quad t_{\text{end}} = 8,$$

and $c = 0$ and $b = 0$.

We compare the computed enclosure by VNODE-LP and AWA [20]. If these enclosure intersect, we accept this test as successful.

File test_n5.cc

This is the restricted three-body test problem from AWA:

$$\begin{aligned} x'' &= x + 2y' - l \frac{x + m}{((x + m)^2 + y^2)^{3/2}} - m \frac{(x - 1)}{((x - 1)^2 + y^2)^{3/2}} \\ x'' &= y - 2x' - l \frac{y}{((x + m)^2 + y^2)^{3/2}} - m \frac{y}{((x - 1)^2 + y^2)^{3/2}}, \end{aligned}$$

where $m = 1/82.45$ and $l = 1 - m$. The initial values are

$$\begin{aligned} x(0) &= 1.2, \\ x'(0) &= 0, \\ y(0) &= 0, \quad \text{and} \\ y'(0) &= -1.04935750983. \end{aligned}$$

We integrate from 0 to 6.192169331396. Again, this test is successful if the computed enclosures by VNODE-LP and AWA intersect.

File test_n6.cc

We integrate the Pleiades problem from the Test Set for IVP Solvers [21]. At the end point, we subtract the reference solution, given in [21], from the computed bounds. If the max norm of the resulting interval vector is $\leq 10^{-2}$, we accept this test as successful.

Chapter 6

Listings

```

# CONFIG_FILE is set in vnodelp/makefile and exported in this
# file. vnodelp/makefile calls this makefile.

include ../config/$(CONFIG_FILE)

CXXFLAGS += -I../include          # compiler flags
LDFLAGS   += -L../lib             # library flags
LIBS      = -lnode $(I_LIBS) $(LAPACK_LIB) \
            $(BLAS_LIB) $(GPP_LIBS) # libraries

EXAMPLES = orbit vanderpol basic E1 scalar basic \
            intermediate integctrl odeparam integri \
            order detest_c3

examples: $(EXAMPLES)

E1:  E1.o
      $(CXX) $(LDFLAGS) -o $@ E1.o $(LIBS)
scalar:  scalar.o
      $(CXX) $(LDFLAGS) -o $@ scalar.o $(LIBS)
basic:  basic.o
      $(CXX) $(LDFLAGS) -o $@ basic.o $(LIBS)
intermediate: intermediate.o
      $(CXX) $(LDFLAGS) -o $@ intermediate.o $(LIBS)
integctrl: integctrl.o
      $(CXX) $(LDFLAGS) -o $@ integctrl.o $(LIBS)
odeparam: odeparam.o
      $(CXX) $(LDFLAGS) -o $@ odeparam.o $(LIBS)
integi: integri.o
      $(CXX) $(LDFLAGS) -o $@ integri.o $(LIBS)
order: orderstudy.o
      $(CXX) $(LDFLAGS) -o $@ orderstudy.o $(LIBS)
detest_c3: detest_c3.o
      $(CXX) $(LDFLAGS) -o $@ detest_c3.o $(LIBS)
vanderpol: vanderpol.o
      $(CXX) $(LDFLAGS) -o $@ vanderpol.o $(LIBS)
orbit: orbit.o
      $(CXX) $(LDFLAGS) -o $@ orbit.o $(LIBS)
clean:
      @-$(RM) *.o *.out core.* $(EXAMPLES)
cleanall:
      @-$(RM) *.o *.cc *.out core.* $(EXAMPLES)

```

Figure 6.1. *The makefile in the examples directory*

```

# file basici.gp
set terminal postscript eps enh color solid "Courier" 28

set xlabel "t"
set ylabel "y_1"

set output 'lorenzi1.eps'
plot [0:6.3][−22:25]\
'lorenzi.out' u 1:($2+$3) \
    title 'lower_bound' w l lt 1 lw 2,\
'lorenzi.out' u 1:($2−$3)\
    title 'upper_bound' w l lt 3 lw 2

set output 'lorenzi2.eps'
set format y "%g"
set xtics 5.7,0.1,6.3
set ylabel "y_1" 0
plot [5.7:6.25][−22:0.4]\
'lorenzi.out' u 1:2:3 \
    title 'bounds' w errorbars lw 2,\
'lorenzi.out' u 1:2 \
    title 'midpoint' w lines lt 3 lw 2

set output 'lorenzi_excess.eps'
set ylabel "global_excess" 2
set logscale y
set xtics 1
set format y "10^{%L}"
plot [0:6.3] 'lorenzi.out' u 1:4 \
    notitle w l lt 1 lw 2

```

Figure 6.2. *The gnuplot file for generating the plot in Figure 3.3*

```

function dy = E1(t,y)
dy = zeros(2,1);
t1 = t+1;
dy(1) = y(2);
dy(2) = -(y(2)/t1 + (1.0 - 0.25/(t1*t1))*y(1));

```

```

clear;
options = odeset('RelTol',1e-10,'AbsTol',1e-10);

y(1)= 0.6713967071418030;
y(2)= 0.09540051444747446;

[T,Y] = ode45(@E1,[0 20],y,options);
format long
Y(end,1:2) '

```

Figure 6.3. The MATLAB code for the DETEST E1 problem

```

# file odeparam.gp
set terminal postscript eps enh color solid "Courier" 28

set xlabel "y_1"
set ylabel "y_2"
set zlabel "y_3"

set xtics -20,10,20
set ytics -25,10,25
set ztics 0,10,45

set output 'odeparam.eps'
splot 'odeparam1.out' u 1:2:3 title '{/Symbol_b}=\_8/3'\
      w l lt 1 lw 2,\
      'odeparam2.out' u 1:2:3 title '{/Symbol_b}=\_5'\
      w l lt 3 lw 2

```

Figure 6.4. The gnuplot file for generating the plots in Figure 3.4

```

# file integctrl.gp
set terminal postscript eps enh color solid "Courier" 28

set ylabel "y_1"
set xlabel "t"

set output 'lorenz.eps'
plot [][-20:22]\
  'lorenz.tight' u 1:2 title 'lower_bound' w l lw 2,\
  'lorenz.tight' u 1:3 title 'upper_bound' w l lt 3 lw 2

set output 'lorenz2.eps'
set xtics 0, 0.1, 0.4
plot [0:0.4]\
  'lorenz.tight' u 1:2 title 'tight' w l lt 3 lw 2,\
  'lorenz.tight' u 1:3 notitle w l lt 3 lw 2,\
  'lorenz.apriori' u 1:2 title 'a_priori' w l lt 1 lw 3

set output 'lorenz_err.eps'
set xtics 0,2,20
set ylabel "global_excess" 2
set logscale y
set format y "10^{%L}"
plot 'lorenz.tight' u 1:4 notitle w l lw 2

set output 'lorenz_step.eps'
set nologscale
set ylabel "stepsize"
plot 'lorenz.step' u 1:2 notitle w l lw 2

```

Figure 6.5. *The gnuplot file for generating the plots in Figure 3.5*

```

# file orderstudy.gp
set terminal postscript eps enh color solid "Courier" 28

set ylabel "CPU_time"
set xlabel "order"

set output 'order.eps'
plot [10:35]\
    'order1e-07.out' u 1:2 title '10^{-7}',
w l lt 1 lw 2,\
    'order1e-08.out' u 1:2 title '10^{-8}',
w l lt 2 lw 2,\
    'order1e-09.out' u 1:2 title '10^{-9}',
w l lt 3 lw 2,\
    'order1e-10.out' u 1:2 title '10^{-10}', w l lt 4 lw 2,\
    'order1e-11.out' u 1:2 title '10^{-11}', w l lt 5 lw 2,\
    'order1e-12.out' u 1:2 title '10^{-12}', w l lt 9 lw 2,\
    'order1e-13.out' u 1:2 title '10^{-13}', w l lt 7 lw 2

set logscale y
set format y "10^{%L}"
set output 'timeorder.eps'
plot 'order1e-07.out' u 1:2 title '10^{-7}' w l lt 1 lw 2,\
    'order1e-08.out' u 1:2 title '10^{-8}' w l lt 2 lw 2,\
    'order1e-09.out' u 1:2 title '10^{-9}' w l lt 3 lw 2,\
    'order1e-10.out' u 1:2 title '10^{-10}' w l lt 4 lw 2,\
    'order1e-11.out' u 1:2 title '10^{-11}' w l lt 5 lw 2,\
    'order1e-12.out' u 1:2 title '10^{-12}' w l lt 9 lw 2,\
    'order1e-13.out' u 1:2 title '10^{-13}' w l lt 7 lw 2

```

Figure 6.6. The gnuplot file for generating the plots in Figure 3.6

```

# file work.gp
set terminal postscript eps enh color solid "Courier" 28

# model of the work
f(x) = a + b*x

fit f(x) "work.out" using (log($1)):(log($2)) via a,b

set xlabel "number_of_equations"
set ylabel "time_per_step_(seconds)"
set xrange [40:200]

set output 'work.eps'
plot 'work.out' using 1:2 notitle with lines

set logscale
set output 'worklog.eps'
plot 'work.out' using 1:2 notitle with lines

```

Figure 6.7. *The gnuplot file for generating the plots in Figure 3.7*

```

# file orbit.gp
set terminal postscript eps enh color solid "Courier" 28

set xlabel "y_1"
set ylabel "y_2"

set output 'orbit_sol.eps'
plot 'orbit_sol.out' u 1:2 notitle w l lt 1 lw 2

set xlabel "t"
set ylabel "stepsize"
set output 'orbit_step.eps'
plot 'orbit_step.out' u 1:2 notitle w l lt 2 lw 2

```

Figure 6.8. *The gnuplot file for generating the plots in Figure 3.8*

```
# file vanderpol.gp
set terminal postscript eps enh color solid "Courier" 28

set xlabel 't'
set ylabel 'stepsize'

set output 'vdp_step1.eps'
plot [0:200][0:0.18] 'vdp_step1.out' u 1:2 \
    title '{/Symbol\m}=10^1' w l lt 2 lw 2,\
    'vdp_step2.out' u 1:2 \
    title '{/Symbol\m}=10^2' w l lt 3 lw 2

set output 'vdp_step2.eps'
plot [0:200][0:0.002] 'vdp_step3.out' u 1:2 \
    title '{/Symbol\m}=10^3' w l lt 3 lw 2,\
    'vdp_step4.out' u 1:2 \
    title '{/Symbol\m}=10^4' w l lt 4 lw 2
```

Figure 6.9. *The gnuplot file for generating the plots in Figure 3.9*

Part II

Third-party Components

Chapter 7

Packages

The VNODE-LP package builds on

- LAPACK [2] and BLAS [1],
- interval-arithmetic (IA) package FILIB++ [19] *or* PROFIL/BIAS [16], and
- automatic differentiation (AD) package FADBAD++ [29].

The interfaces to the IA package are kept as small as possible, which allows a new package to be introduced without substantial programming effort.

Chapter 8 presents the implementation of interfaces to FILIB++ and PROFIL/BIAS. Chapter 9 discusses functions for changing the rounding mode. The AD in VNODE-LP is implemented through abstract classes, which are described in Chapter 17. Implementation of these classes using FADBAD++ is given in Chapter 22.

Chapter 8

IA package

The basic data type in VNODE-LP is **interval**. This package can be built on top of FILIB++ [19] or PROFIL [16], or potentially other packages. The user can select which of these packages to use; for more details see the installation instructions in Section 2.3.

Each of these IA packages can be replaced, provided that the new package supplies an interval data type with overloaded arithmetic operations and elementary functions working with interval arguments. To incorporate a new IA package, the body of the functions described in Section 4, and implemented below, need to be implemented using the new package.

8.1 Functions calling FILIB++

```
113 <functions calling FILIB++ 113> ≡
    inline double inf(const interval &a) {
        return a.inf();
    }
    inline double sup(const interval &a) {
        return a.sup();
    }
    inline double midpoint(const interval &a) {
        return a.mid();
    }
    inline double width(const interval &a) {
        return a.diam();
    }
    inline double mag(const interval &a) {
        return a.mag();
    }
```

```

inline bool subseteq(const interval &a, const interval &b) {
    return filib::subset(a, b);
}

inline bool interior(const interval &a, const interval &b) {
    return filib::interior(a, b);
}

inline bool disjoint(const interval &a, const interval &b) {
    return filib::disjoint(a, b);
}

inline bool intersect(interval &c, const interval &a, const interval &b)
{
    if (filib::disjoint(a, b)) return false;
    c = filib::intersect(a, b);
    return true;
}

inline interval pi() {
    return filib::interval<double>::PI();
}

inline interval pow(const interval &a, const interval &b) {
    return filib::pow(a, b);
}

inline interval pow(const interval &a, int b) {
    return filib::power(a, b);
}

inline interval exp(const interval &a) {
    return filib::exp(a);
}

inline interval log(const interval &a) {
    return filib::log(a);
}

inline interval sqr(const interval &a) {
    return filib::sqr(a);
}

inline interval sqrt(const interval &a) {
    return filib::sqrt(a);
}

inline interval sin(const interval &a) {
    return filib::sin(a);
}

inline interval cos(const interval &a) {
    return filib::cos(a);
}

```

```

inline interval tan(const interval &a) {
    return filib::tan(a);
}
inline interval asin(const interval &a) {
    return filib::asin(a);
}
inline interval acos(const interval &a) {
    return filib::acos(a);
}
inline interval atan(const interval &a) {
    return filib::atan(a);
}
inline interval string_to_interval(const char *s)
{
    std::cerr << "\n\n***_WARNING_***\n";
    std::cerr << "Conversion from a string to an interval containing it\n";
    std::cerr << "has not been implemented in filib++.\n";
    std::cerr << "For this feature, consider using PROFIL/BIAS.\n";
    std::cerr << "The input string" << s <<
        "is converted to double.\n";
    std::cerr << "THE COMPUTED BOUNDS MAY NOT BE CORRECT\n\n";
    double a;
    std::istringstream iss(s);
    iss >> a;
    return interval(a);
}

```

This code is used in chunk 115.

8.2 Functions calling PROFIL

```

114 <functions calling PROFIL 114> ≡
    inline double inf(const interval &a) {
        return Inf(a);
    }
    inline double sup(const interval &a) {
        return Sup(a);
    }
    inline double midpoint(const interval &a) {
        return Mid(a);
    }
    inline double width(const interval &b) {

```

```

    return Diam(b);
}
inline double mag(const interval &a) {
    return Abs(a);
}
inline bool subseq(const interval &a, const interval &b) {
    return  $a \leq b$ ;
}
inline bool interior(const interval &a, const interval &b) {
    return (Inf(a) > (Inf(b))  $\wedge$  (Sup(a) < (Sup(b)));
}
inline bool disjoint(const interval &a, const interval &b) {
    interval c;
    return  $\neg$ Intersection(c, a, b);
}
inline bool intersect(interval &c, const interval &a, const interval &b)
{
    return Intersection(c, a, b);
}
inline interval pi() {
    return ArcCos(-1.0);
}
inline interval pow(const interval &a, int b) {
    return Power(a, b);
}
inline interval pow(const interval &a, const interval &b) {
    return Power(a, b);
}
inline interval exp(const interval &a) {
    return Exp(a);
}
inline interval log(const interval &a) {
    return Log(a);
}
inline interval sqr(const interval &a) {
    return Sqr(a);
}
inline interval sqrt(const interval &a) {
    return Sqrt(a);
}
inline interval sin(const interval &a) {
    return Sin(a);
}

```

```

inline interval cos(const interval &a) {
    return Cos(a);
}
inline interval tan(const interval &a) {
    return Tan(a);
}
inline interval asin(const interval &a) {
    return ArcSin(a);
}
inline interval acos(const interval &a) {
    return ArcCos(a);
}
inline interval atan(const interval &a) {
    return ArcTan(a);
}
inline interval string_to_interval(const char *s)
{
    return Enclosure(s);
}

```

This code is used in chunk 115.

Files

The interface to the IA package is stored in

```

115 <vnodeinterval.h 115> ≡
    #ifndef VNODEINTERVAL_H
    #define VNODEINTERVAL_H
    #ifdef PROFIL_VNODE
    #include <Interval.h>
    #include <Functions.h>
    #include <LongReal.h>
    #include <LongInterval.h>
    namespace v_bias {
        <interval data type (PROFIL) 76>
        <functions calling PROFIL 114>
    }
    #endif
    #ifdef FILIB_VNODE
    #include <interval/interval.hpp>
    #include <iostream>
    #include <sstream>
    #include <string>
    namespace v_bias {
        <interval data type (FILIB++) 77>
        <functions calling FILIB++ 113>
    }

```

```
}  
#endif  
#endif
```


Chapter 9

Changing the rounding mode

For changing the rounding mode, VNODE-LP calls the functions below.

void *round_nearest*()
sets the rounding mode to the nearest.

void *round_down*()
sets the rounding mode to $-\infty$.

void *round_up*()
sets the rounding mode to ∞ .

Depending on the selected IA package, they call corresponding functions either from PROFIL/BIAS or FILIB++. A particular implementation is selected by the value of the variable `I_PACKAGE` in the configuration file; see Subsection 2.3.2.

9.1 Changing the rounding mode using FILIB++

Changing the rounding mode using FILIB++ is done through `round_control` defined as

```
117 <rounding control type (FILIB++) 117> ≡  
    typedef filib::rounding_control < double , true > round_control;  
This code is used in chunk 120.
```

The necessary functions are implemented as follows.

```
118 <changing the rounding mode (FILIB++) 118> ≡  
    inline void round_nearest() {  
        round_control::tonearest();  
    }  
    inline void round_down() {  
        round_control::downward();  
    }
```

```

inline void round_up() {
    round_control::upward();
}

```

This code is used in chunk 120.

9.2 Changing the rounding mode using BIAS

Similarly, we implement functions for changing the rounding mode using BIAS functions.

119 \langle changing the rounding mode (BIAS) 119 $\rangle \equiv$

```

inline void round_nearest() {
    BiasRoundNear();
}

inline void round_down() {
    BiasRoundDown();
}

inline void round_up() {
    BiasRoundUp();
}

```

This code is used in chunk 120.

Files

The above functions are stored in

```

120  $\langle$  vnoderound.h 120  $\rangle \equiv$ 
#ifndef VNODEROUND_H
#define VNODEROUND_H
#ifdef FILIB_VNODE
#include <rounding_control/rounding_control_double.hpp>
    namespace v_bias {
         $\langle$  rounding control type (FILIB++) 117  $\rangle$ 
         $\langle$  changing the rounding mode (FILIB++) 118  $\rangle$ 
    }
#endif
#ifdef PROFIL_VNODE
#include "Bias0.h"
    namespace v_bias {
         $\langle$  changing the rounding mode (BIAS) 119  $\rangle$ 
    }
#endif
#endif

```

Part III

Linear Algebra and Related Functions

Chapter 10

Vectors and Matrices

The VNODE-LP package works with point and interval vectors and matrices, namely **pVector**, **iVector**, **pMatrix**, and **iMatrix**. Here “**p**” is for point and “**i**” is for interval. The **iVector** type was defined in Subsection 4.3; **pVector**, **pMatrix**, and **iMatrix** are defined as

```
122 <vector and matrix types 122> ≡  
    <interval vector 80>  
    typedef vector<double> pVector;  
    typedef vector<vector<double>> pMatrix;  
    typedef vector<vector<interval>> iMatrix;
```

This code is used in chunk 127.

Size and memory allocation

template<class Matrix> void sizeM(Matrix &A, unsigned int n)
allocates space for an $n \times n$ matrix.

template<class Matrix> sizeM(const Matrix &A)
returns the size of a square matrix A .

template<class Vector> unsigned int sizeV(const Vector &a)
allocate space for an n vector.

template<class Vector> unsigned int sizeV(const Vector &a)
returns the size of a vector.

```
124 <size/allocation 124> ≡  
    template<class Matrix> inline void sizeM(Matrix &A, unsigned int n)  
    {  
        A.resize(n);  
        for (unsigned int i = 0; i < A.size(); i++) A[i].resize(n);  
    }
```

```

template<class Matrix> inline unsigned int sizeM(const Matrix &A)
{
    return A.size();
}
template<class Vector> inline void sizeV(Vector &a, unsigned int n)
{
    a.resize(n);
}
template<class Vector> inline unsigned int sizeV(const Vector &a)
{
    return a.size();
}

```

This code is used in chunk 127.

Files

```

127 <vector_matrix.h 127> ≡
    #ifndef VECTOR_MATRIX
    #define VECTOR_MATRIX
    #include <vector>
    #include "vnodeinterval.h"
    using namespace std;
    using namespace v_bias;
    namespace v_blas {
        <vector and matrix types 122>
        <size/allocation 124>
    }
    #endif

```

Chapter 11

Basic functions

We provide “generic” functions for operations involving matrices and vectors. Below, A , B , and C are $n \times n$ matrices, x , y , and z are n vectors, and a is a scalar.

11.1 Vector operations

```
template<class Vector, class scalar>
void setV(Vector &z, scalar a)
    sets each component of  $z$  to  $a$ .
```

```
template<class Vector, class Vector>
void assignV(Vector &z, const Vector &x)
    copies  $x$  to  $z$ 
```

```
template<class Vector, class scalar>
void scaleV(Vector &z, scalar a)
    multiplies each element of  $z$  by  $a$ .
```

```
void addViVi(iVector &z, const iVector &x)
    adds  $z$  and  $x$  and stores the result in  $z$ .
```

```
void addViVp(iVector &z, const pVector &x)
    adds  $z$  and  $x$  and stores the result in  $z$ .
```

```
void subViVp(iVector &z, const pVector &x)
    subtracts  $x$  from  $z$  and stores the result in  $z$ .
```

```
void subViVi(iVector &z, const iVector &x)
    subtracts  $x$  from  $z$  and stores the result in  $z$ .
```

```
void addViVi(iVector &z, const iVector &x, const iVector &y)
    adds  $x$  and  $y$  and stores the result in  $z$ .
```

void *addViVp*(**iVector** &*z*, **const iVector** &*x*, **const pVector** &*y*)
 adds *x* and *y* and stores the result in *z*.

void *subViVp*(**iVector** &*z*, **const iVector** &*x*, **const pVector** &*y*)
 subtracts *y* from *x* and stores the result in *z*.

double *inf_normV*(**const iVector** &*z*)
 returns $\max |z_i|$, where $|a_i| = \max\{|a_i|, |\bar{a}_i|\}$.

double *inf_normV*(**const pVector** &*z*)
 returns $\|z\|$.

template<**class** *scalar*, **class** *Vector1*, **class** *Vector2*>
inline void *dot_product*(*scalar* &*r*, **const Vector1** &*a*, **const Vector2** &*b*)
 computes the dot product of *a* and *b* and stores it in *r*. For input point vectors,
 this dot product is computed in the current rounding mode.

double *norm2*(**const pVector** &*v*)
 returns the two norm of a point vector. For a point vector, this norm is
 computed in the current rounding mode.

129 <vector operations 129> \equiv

```

template<class Vector, class scalar>
  inline void setV(Vector &z, scalar a)
  {
    fill(z.begin(), z.end(), a);
  }

template<class Vector1, class Vector2>
  inline void assignV(Vector1 &z, const Vector2 &x)
  {
    for (unsigned int i = 0; i < sizeV(z); i++) z[i] = x[i];
  }

template<class Vector, class scalar>
  inline void scaleV(Vector &z, scalar a)
  {
    for (unsigned int i = 0; i < sizeV(z); i++) z[i] *= a;
  }

inline void addViVi(iVector &z, const iVector &x)
  {
    transform(z.begin(), z.end(), x.begin(), z.begin(), plus<v_bias::interval>());
  }

inline void addViVp(iVector &z, const pVector &x)
  {
    transform(z.begin(), z.end(), x.begin(), z.begin(), plus<v_bias::interval>());
  }

```



```

inline void subViVp(iVector &z, const pVector &x)
{
    transform(z.begin(), z.end(), x.begin(), z.begin(),
        minus<v_bias::interval>());
}

inline void subViVi(iVector &z, const iVector &x)
{
    transform(z.begin(), z.end(), x.begin(), z.begin(),
        minus<v_bias::interval>());
}

inline void addViVi(iVector &z, const iVector &x, const iVector &y)
{
    transform(x.begin(), x.end(), y.begin(), z.begin(), plus<v_bias::interval>());
}

inline void addViVp(iVector &z, const iVector &x, const pVector &y)
{
    transform(x.begin(), x.end(), y.begin(), z.begin(), plus<v_bias::interval>());
}

inline void subViVp(iVector &z, const iVector &x, const pVector &y)
{
    transform(x.begin(), x.end(), y.begin(), z.begin(),
        minus<v_bias::interval>());
}

inline double inf_normV(const iVector &z)
{
    double s = 0;
    for (unsigned int i = 0; i < sizeV(z); i++)
        if (v_bias::mag(z[i]) > s) s = v_bias::mag(z[i]);
    return s;
}

inline double inf_normV(const pVector &z)
{
    double s = 0;
    for (unsigned int i = 0; i < sizeV(z); i++)
        if (fabs(z[i]) > s) s = fabs(z[i]);
    return s;
}

template<class scalar, class Vector1, class Vector2>
    inline void dot_product(scalar &r, const Vector1 &a, const
        Vector2 &b)
{
    r = 0.0;
    for (unsigned int i = 0; i < a.size(); i++) r += a[i] * b[i];
}

```

```

inline double norm2(const pVector &v) {
    double s;
    dot_product(s, v, v);
    return sqrt(s);
}

```

This code is used in chunk 137.

11.2 Matrix/vector operations

void multMiVi(iVector &z, const iMatrix &A, const iVector &x)
 multiplies A and x and stores the result in z .

void multMpVi(iVector &z, const pMatrix &A, const iVector &x)
 multiplies A and x and stores the result in z .

void multMiVp(iVector &z, const pMatrix &A, const iVector &x)
 multiplies A and x and stores the result in z .

131 \langle matrix times vector 131 $\rangle \equiv$

```

inline void multMiVi(iVector &z, const iMatrix &A, const iVector &x)
{
    for (unsigned int i = 0; i < A.size(); i++) dot_product(z[i], A[i], x);
}
inline void multMpVi(iVector &z, const pMatrix &A, const iVector &x)
{
    for (unsigned int i = 0; i < A.size(); i++) dot_product(z[i], A[i], x);
}
inline void multMiVp(iVector &z, const iMatrix &A, const pVector &x)
{
    for (unsigned int i = 0; i < A.size(); i++) dot_product(z[i], A[i], x);
}

```

This code is used in chunk 137.

11.3 Matrix operations

template<class Matrix, class scalar>
void setM(Matrix &C, scalar a)
 sets each component of C to a .

template<class Matrix> void setId(Matrix &C)
 sets C to the identity matrix.

template<class Matrix1, class Matrix2>
void assignM(Matrix1 &C, const Matrix2 &A)
 copies A to C .

template \langle **class** **Matrix**, **class** **scalar** \rangle **void** *scaleM*(**Matrix** &*C*, **scalar** *a*)
 multiplies each component of *C* by *a*.

template \langle **class** **Matrix** \rangle **void** *transpose*(**Matrix** &*C*, **const** **Matrix** &*A*)
 stores the transpose of *A* in *C*.

template \langle **class** **Matrix** \rangle **void** *addId*(**Matrix** &*C*)
 adds the identity matrix to *C*.

template \langle **class** **Matrix** \rangle **void** *subFromId*(**Matrix** &*C*)
 subtracts *C* from the identity matrix and stores the result in *C*.

void *addMiMi*(**iMatrix** &*C*, **const** **iMatrix** &*A*)
 adds *C* and *A* and stores the result in *C*.

void *subMiMp*(**iMatrix** &*C*, **const** **pMatrix** &*A*)
 subtracts *A* from *C* and stores the result in *C*.

void *multMiMi*(**iMatrix** &*C*, **const** **iMatrix** &*A*, **const** **iMatrix** &*B*)
 multiplies *A* and *B* and stored the result in *C*.

void *multMiMp*(**iMatrix** &*C*, **const** **iMatrix** &*A*, **const** **pMatrix** &*B*)
 multiplies *A* and *B* and stored the result in *C*.

double *inf_normM*(**const** **iMatrix** &*C*)
 computes $\|C\|_{\infty}$. The computation is in the current rounding mode. For example, if an upper bound on $\|C\|_{\infty}$ is desired, the rounding mode must be set to $+\infty$ before calling *inf_normM*.

132 \langle matrix operations 132 $\rangle \equiv$

```

template $\langle$ class Matrix, class scalar $\rangle$ 
    inline void setM(Matrix &C, scalar a)
    {
        for (unsigned int i = 0; i < C.size(); i++) setV(C[i], a);
    }
template $\langle$ class Matrix $\rangle$ 
    inline void setId(Matrix &C)
    {
        setM(C, 0.0);
        for (unsigned int i = 0; i < sizeM(C); i++) C[i][i] = 1.0;
    }
template $\langle$ class Matrix1, class Matrix2 $\rangle$ 
    inline void assignM(Matrix1 &C, const Matrix2 &A)
    {
        for (unsigned int i = 0; i < C.size(); i++) assignV(C[i], A[i]);
    }

```

```

template<class Matrix, class scalar>
    inline void scaleM(Matrix &C, scalar a)
{
    for (unsigned int i = 0; i < C.size(); i++) scaleV(C[i], a);
}

template<class Matrix>
    inline void transpose(Matrix &C, const Matrix &A)
{
    for (unsigned int i = 0; i < C.size(); i++) setColumn(C, A[i], i);
}

template<class Matrix>
    inline void addId(Matrix &C)
{
    for (unsigned int i = 0; i < sizeM(C); i++) C[i][i] += 1.0;
}

template<class Matrix>
    inline void subFromId(Matrix &C)
{
    unsigned int n = sizeM(C);
    for (unsigned int i = 0; i < n; i++)
        for (unsigned int j = 0; j < n; j++) C[i][j] = -C[i][j];
    for (unsigned int i = 0; i < n; i++) C[i][i] += 1.0;
}

inline void addMiMi(iMatrix &C, const iMatrix &A)
{
    for (unsigned int i = 0; i < C.size(); i++) addViVi(C[i], A[i]);
}

inline void subMiMp(iMatrix &C, const pMatrix &A)
{
    for (unsigned int i = 0; i < C.size(); i++) subViVp(C[i], A[i]);
}

inline void multMiMi(iMatrix &C, const iMatrix &A, const iMatrix &B)
{
    unsigned int n = sizeM(A);
    for (unsigned int i = 0; i < n; i++)
        for (unsigned int j = 0; j < n; j++) {
            C[i][j] = 0.0;
            for (unsigned int k = 0; k < n; k++) C[i][j] += A[i][k] * B[k][j];
        }
}

inline void multMiMp(iMatrix &C, const iMatrix &A, const pMatrix &B)
{
    unsigned int n = sizeM(A);

```

```

    for (unsigned int i = 0; i < n; i++)
        for (unsigned int j = 0; j < n; j++) {
            C[i][j] = 0.0;
            for (unsigned int k = 0; k < n; k++) C[i][j] += A[i][k] * B[k][j];
        }
}
template<class Matrix> inline double inf_normM(const Matrix &C)
{
    unsigned int n = sizeM(C);
    double m = 0;
    for (unsigned int i = 0; i < n; i++) {
        double s = 0;
        for (unsigned int j = 0; j < n; j++) s += v_bias::mag(C[i][j]);
        if (s > m) m = s;
    }
    return m;
}

```

This code is used in chunk 137.

11.4 Get/set column

```

template<class Vector, class Matrix>
void getColumn(Vector &z, const Matrix &C, unsigned int j)
    stores the jth column of C in z.

template<class Matrix, class Vector>
void setColumn(Matrix &C, const Vector &z, unsigned int j)
    sets the jth column of C to z.

```

134 <get/set column 134> \equiv

```

template<class Vector, class Matrix>
void getColumn(Vector &z, const Matrix &C, unsigned int j)
{
    for (unsigned int i = 0; i < sizeM(C); i++) z[i] = C[i][j];
}

template<class Matrix, class Vector>
void setColumn(Matrix &C, const Vector &z, unsigned int j)
{
    for (unsigned int i = 0; i < sizeM(C); i++) C[i][j] = z[i];
}

```

This code is used in chunk 137.

11.5 Conversions

The following two functions are convenient when calling LAPACK.

```
template<class scalar, class Matrix>
void matrix2pointer(scalar *M, const Matrix &C)
copies  $C$  into an array pointed to by  $M$ . The matrix at  $M$  is in a column-
major form. That is, the  $(i, j)$  element of the matrix in the array at  $M$  is at
 $jn + i$ , where  $n$  is the size of the matrix.
```

```
template<class Matrix, class scalar>
void pointer2matrix(Matrix &C, const scalar *M)
copies a matrix stored in an array pointed to by  $M$  into a matrix  $C$ . The
matrix at  $M$  is in a column-major form. That is, the  $(i, j)$  element of the
matrix in the array at  $M$  is at  $jn + i$ , where  $n$  is the size of the matrix.
```

```
135 <matrix2pointer 135> ≡
    template<class scalar, class Matrix>
        inline void matrix2pointer(scalar *M, const Matrix &C)
        {
            unsigned int n = sizeM(C);
            for (unsigned int j = 0; j < n; j++)
                for (unsigned int i = 0; i < n; i++) M[j * n + i] = C[i][j];
        }
```

This code is used in chunk 137.

```
136 <pointer2matrix 136> ≡
    template<class Matrix, class scalar>
        inline void pointer2matrix(Matrix &C, const scalar *M)
        {
            unsigned int n = sizeM(C);
            for (unsigned int j = 0; j < n; j++)
                for (unsigned int i = 0; i < n; i++) C[i][j] = M[j * n + i];
        }
```

This code is used in chunk 137.

Files

We store the above functions in

```
137 <basiclinalg.h 137> ≡
    #ifndef BASICLINALG_H
    #define BASICLINALG_H
    #include <algorithm>
    #include "vector_matrix.h"
    namespace v_blas {
        <vector operations 129>
        <matrix times vector 131>
        <matrix operations 132>
        <get/set column 134>
```

```
    <matrix2pointer 135>  
    <pointer2matrix 136>  
    <print vector 366>  
  }  
#endif
```


Chapter 12

Interval functions

12.1 Inclusion

This *subseq* function returns *true* if, for two interval vectors *a* and *b*, $a \subseteq b$, and returns *false* otherwise.

```
139 <check vector inclusion 139> ≡  
    inline bool subseq(const iVector &a, const iVector &b)  
    {  
        return equal(a.begin(), a.end(), b.begin(), v_bias::subseq);  
    }
```

See also chunk 141.

This code is used in chunk 151.

12.2 Interior

The *interior* function returns *true* if, for two interval vectors *a* and *b*, *a* is componentwise in the interior of *b*, and returns *false* otherwise.

```
140 <check if in the interior 140> ≡  
    inline bool interior(const iVector &a, const iVector &b)  
    {  
        return equal(a.begin(), a.end(), b.begin(), v_bias::interior);  
    }
```

This code is used in chunk 151.

The *disjoint* function returns *true* if, for two interval vectors *a* and *b*, $a \cap b = \emptyset$, and returns *false* otherwise.

```
141 <check vector inclusion 139> +≡  
    inline bool disjoint(const iVector &a, const iVector &b)  
    {  
        return equal(a.begin(), a.end(), b.begin(), v_bias::disjoint);  
    }
```

12.3 Radius

The radius of an interval is

```
142 <rad (interval) 142> ≡
    inline double rad(const interval &a)
    {
        return 0.5 * v_bias::width(a);
    }
```

This code is used in chunk 151.

We compute the radius of an interval vector by

```
143 <rad (vector) 143> ≡
    inline void rad(pVector &r, const iVector &v)
    {
        transform(v.begin(), v.end(), r.begin(), v_bias::rad);
    }
```

This code is used in chunk 151.

12.4 Width

We compute the width of an interval vector by

```
144 <width 144> ≡
    inline void width(pVector &r, const iVector &a)
    {
        transform(a.begin(), a.end(), r.begin(), v_bias::width);
    }
```

This code is used in chunk 151.

12.5 Midpoints

```
145 <midpoint of an interval vector 145> ≡
    inline void midpoint(pVector &r, const iVector &a)
    {
        transform(a.begin(), a.end(), r.begin(), v_bias::midpoint);
    }
```

This code is used in chunk 151.

```
146 <midpoint of an interval matrix 146> ≡
    inline void midpoint(pMatrix &R, const iMatrix &A)
    {
        for (unsigned int i = 0; i < A.size(); i++) midpoint(R[i], A[i]);
    }
```

This code is used in chunk 151.

12.6 Intersection

If interval vectors \mathbf{x} and \mathbf{y} intersect, we store their intersection in \mathbf{z} and return *true*. Otherwise, we return *false*.

```

147 < intersection of interval vectors 147 > ≡
    inline bool intersect(iVector &z, const iVector &x, const iVector &y)
    {
        interval c;
        for (unsigned int i = 0; i < sizeV(y); i++) {
            bool b = v_bias::intersect(c, x[i], y[i]);
            if (!b) return false;
            else z[i] = c;
        }
        return true;
    }

```

This code is used in chunk 151.

12.7 Computing h such that $[0, h]\mathbf{a} \subseteq \mathbf{b}$

12.7.1 The interval case

Given intervals \mathbf{a} and \mathbf{b} , where $0 \in \mathbf{b}$, we wish to compute the largest machine-representable $h \geq 0$ such that

$$[0, h]\mathbf{a} \subseteq \mathbf{b}.$$

We consider the following cases.

1. $\underline{\mathbf{a}} = \overline{\mathbf{a}}$.
 - (a) If $\underline{\mathbf{a}} = \overline{\mathbf{a}} = 0$, then we set h to be the largest representable machine number.
 - (b) If $\underline{\mathbf{a}} = \overline{\mathbf{a}} > 0$, then $[0, h]\mathbf{a} = [0, h\overline{\mathbf{a}}] \subseteq \mathbf{b}$ iff $h \leq \overline{\mathbf{b}}/\overline{\mathbf{a}}$.
 - (c) If $\underline{\mathbf{a}} = \overline{\mathbf{a}} < 0$, then $[0, h]\mathbf{a} = [\underline{\mathbf{a}}h, 0] \subseteq \mathbf{b}$ iff $h \leq \underline{\mathbf{b}}/\underline{\mathbf{a}}$.
2. $\underline{\mathbf{a}} \neq \overline{\mathbf{a}}$.
 - (a) $\underline{\mathbf{a}} \geq 0$. Then $\overline{\mathbf{a}} > 0$ and $[0, h]\mathbf{a} = [0, h\overline{\mathbf{a}}] \subseteq \mathbf{b}$ iff $h \leq \overline{\mathbf{b}}/\overline{\mathbf{a}}$.
 - (b) $\overline{\mathbf{a}} \leq 0$. Then $\underline{\mathbf{a}} < 0$ and $[0, h]\mathbf{a} = [h\underline{\mathbf{a}}, 0] \subseteq \mathbf{b}$ iff $h \leq \underline{\mathbf{b}}/\underline{\mathbf{a}}$.

We summarize the above cases in the following procedure:

1. If $\underline{\mathbf{a}} = \overline{\mathbf{a}} = 0$ then set h to be the largest representable machine number;
2. else
 - (a) if $\underline{\mathbf{a}} \geq 0$ then $h = \nabla(\overline{\mathbf{b}}/\overline{\mathbf{a}})$
 - (b) else $h = \nabla(\underline{\mathbf{b}}/\underline{\mathbf{a}})$.

```

149  $\langle h \text{ such that } [0, h]\mathbf{a} \subseteq \mathbf{b} \text{ (intervals)} \ 149 \rangle \equiv$ 
    #include <climits>
    using namespace std;
    using namespace v_bias;

    inline double compH(const v_bias::interval &a, const v_bias::interval
        &b)
    {
        if (inf(a)  $\equiv$  0  $\wedge$  sup(a)  $\equiv$  0) return numeric_limits<double>::max();
        round_down();
        if (inf(a)  $\geq$  0) return sup(b)/sup(a);
        return inf(b)/inf(a);
    }

```

This code is used in chunk 150.

12.7.2 The interval vector case

Given two interval vectors \mathbf{a} and \mathbf{b} , where \mathbf{b} contains the vector with zero component, we compute the largest machine-representable $h \geq 0$ such that

$$[0, h]\mathbf{a} \subseteq \mathbf{b}.$$

```

150  $\langle h \text{ such that } [0, h]\mathbf{a} \subseteq \mathbf{b} \text{ (interval vectors)} \ 150 \rangle \equiv$ 
     $\langle h \text{ such that } [0, h]\mathbf{a} \subseteq \mathbf{b} \text{ (intervals)} \ 149 \rangle$ 
    double compH(const iVector &a, const iVector &b)
    {
        double hmin = compH(a[0], b[0]);
        for (unsigned int i = 1; i < size V(a); i++) {
            double h = compH(a[i], b[i]);
            if (h < hmin) hmin = h;
        }
        return hmin;
    }

```

This code is used in chunk 152.

Files

```

151  $\langle \text{intvfuncs.h} \ 151 \rangle \equiv$ 
    #ifndef INTVFUNC_H
    #define INTVFUNC_H
    #include "vector_matrix.h"
    namespace v_bias {
         $\langle \text{rad (interval)} \ 142 \rangle$ 
    } namespace v_blas {
        using namespace v_bias;

```

```

    <check if in the interior 140>
    <check vector inclusion 139>
    <width 144>
    <midpoint of an interval vector 145>
    <midpoint of an interval matrix 146>
    <intersection of interval vectors 147>
    <rad (vector) 143>
    double compH(const iVector &a, const iVector &b);
  }
#endif

152 <intvfuncs.cc 152> ≡
    #include <limits>
    #include <algorithm>
    #include "vnodeinterval.h"
    #include "vnodearound.h"
    #include "vector_matrix.h"
    namespace v_blas {
      < $h$  such that  $[0, h]\mathbf{a} \subseteq \mathbf{b}$  (interval vectors) 150>
    }

```


Chapter 13

QR factorization

In our implementation of VNODE-LP, we need to compute the QR factorization of an $n \times n$ matrix. We employ the routines DGEQRF and DORGQR from LAPACK. DGEQRF computes a QR factorization of a real $m \times n$ matrix A . DORGQR generates an $m \times n$ real matrix Q with (approximately) orthonormal columns, which is defined from the elementary reflectors returned by DGEQRF. If *computeQR* is successful, *true* is returned; otherwise *false* is returned.

```
153 <compute QR factorization 153> ≡
    extern "C"
    {
        void dgeqrf_(int *m, int *n, double *A, int *lda, double *tau, double
            *work, int *lwork, int *info);
        void dorgqr_(int *m, int *n, int *k, double *A, int *lda, double
            *tau, double *work, int *lwork, int *info);
    }

    bool computeQR(pMatrix &Q, const pMatrix &A)
    {
        int n = sizeM(A);
        int m = n;
        int lda = n;
        int k = n;
        int info;
        int lwork = 10 * n; /* lwork has to be n *(optimal block size). The value
            10 is somewhat random. */
        double *tau = new double[n];
        double *work = new double[lwork];
        double *M = new double[n * n];
        v_bias::round_nearest();
        matrix2pointer(M, A);
        dgeqrf_(&m, &n, M, &lda, tau, work, &lwork, &info);
```

```

    if (info  $\equiv$  0) {
        dorgqr_(&m, &n, &k, M, &lda, tau, work, &lwork, &info);
        if (info  $\equiv$  0) pointer2matrix(Q, M);
    }
    delete[] M;
    delete[] work;
    delete[] tau;
    if (info  $\equiv$  0) return true;
    return false;
}

```

This code is used in chunk 154.

Files

```

154 <qr.cc 154>  $\equiv$ 
    #include "basiclinalg.h"
    #include "vnodearound.h"
    namespace vnodelp {
        using namespace v_blas;
        <compute QR factorization 153>
    }

```


Chapter 14

Matrix inverse

The **MatrixInverse** class provides functions for computing an approximate inverse of a point matrix, enclosing the inverse of a point matrix and enclosing the inverse of a floating-point approximation to an orthogonal matrix.

14.1 Matrix inverse class

```
156 < class MatrixInverse 156 > ≡  
    class MatrixInverse {  
    public:  
        MatrixInverse(int n);  
  
        bool invertMatrix(pMatrix &Ainv, const pMatrix &A);  
        bool encloseMatrixInverse(iMatrix &Ainv, const pMatrix &A);  
        bool orthogonalInverse(iMatrix &Ainv, const pMatrix &A);  
  
        ~MatrixInverse();  
  
        int iterations;  
    private:  
        bool encloseLS(iVector &x, const pMatrix &A, const iVector &b0,  
            const iMatrix &B, double beta);  
        pVector radx;  
        iVector b0, x1, x;  
        iMatrix B, Ci;  
        pMatrix C;  
        double *M;  
        int *ipiv;  
        double *work;  
        int lwork;  
    };
```

This code is used in chunk 172.

invertMatrix tries to compute a floating-point approximation to the inverse of A (if it exists). If successful, *invertMatrix* stores the result in *Ainv* and returns *true*; otherwise, it returns *false*.

encloseMatrixInverse tries to enclose the inverse of A (if it exists). If successful, *encloseMatrixInverse* stores the result in *Ainv* and returns *true*; otherwise, it returns *false*.

orthogonalInverse tries to enclose the inverse of a floating-point approximation to an orthogonal matrix. If successful, *orthogonalInverse* stores the result in *Ainv* and returns *true*; otherwise it returns *false*.

14.2 Computing A^{-1}

First, we compute the LU factorization of A using LAPACK's *dgetrf*. Then, using this LU factorization, we try to compute the inverse of A using LAPACK's *dgetri*.

```

158 <compute  $A^{-1}$  158> ≡
    extern "C"
    {
        void dgetrf_(int *m, int *n, double *A, int *lda,
                     int *ipiv, int *info);
        void dgetri_(int *n, double *A, int *lda,
                     int *ipiv, double *work, int *lwork, int *info);
    }
    bool MatrixInverse::invertMatrix(pMatrix &Ainv, const pMatrix &A)
    {
        int n = sizeM(A);
        int lda = n;
        int info;
        matrix2pointer(M, A);    /*
    •  $n$  number of rows and columns
    •  $M$  the matrix to be factored on input
    •  $lda$  the leading dimension of  $M$ ,  $lda \geq \max(1, n)$ . We set  $lda = n$ 
    •  $ipiv$  pivot indices
    •  $info$  success if  $info \equiv 0$ 

        */
        v_bias::round_nearest();
        dgetrf_(&n, &n, M, &lda, ipiv, &info);
        if (info != 0) {
#ifdef VNODE_DEBUG

```

```

        printMessage("Could_not_invert_a_matrix");
#endif
        return false;
    } /*
    • lwork size of work. We set  $lwork = 2 * n$  in the constructor
    • work array of size lwork
    • info success if  $info \equiv 0$ 
    */
    dgetri_(&n, M, &lida, ipiv, work, &lwork, &info);
    if (info != 0) {
#ifdef VNODE_DEBUG
        printMessage("Could_not_invert_a_matrix");
#endif
        return false;
    }
    pointer2matrix(Ainv, M);
    return true;
}

```

This code is used in chunk 173.

14.3 Enclosing the solution of a linear system

We try to enclose the solution to the linear system $Ax = b$, $A \in \mathbb{R}^{n \times n}$, $b \in \mathbb{R}^n$ in *encloseLS*. It implements Krawczyk's method; see for example [10]. By $C \in \mathbb{R}^{n \times n}$ we denote a preconditioner. This method works if $\|I - CA\|_\infty = \beta < 1$. Let

$$\alpha = \frac{\|Cb\|_\infty}{1 - \beta}.$$

Then

$$x \in ([-\alpha, \alpha], \dots, [-\alpha, \alpha])^T.$$

This shows how to compute an initial box that is guaranteed to contain x . Then Krawczyk's iteration is

$$x^{(i+1)} = (Cb + (I - CA)x^{(i)}) \cap x^{(i)}.$$

The *encloseLS* function follows. In the comment inside it, before the horizontal line, we list the input variables and what they contain. The output is described after this line.

When describing code in this manuscript, we shall mix variable names and math symbols, where appropriate. Such a mixture may be viewed as an abuse of notation, but the author has found it helpful. For example $b0 \ni Cb$ below means that the interval vector $b0$ contains the true Cb . Similarly, $beta \geq \|I - CA\|_\infty$ denotes that the value of the input variable *beta* must be greater or equal the true value of $\|I - CA\|_\infty$.

```

160 <enclose the solution to  $Ax = b$  160>  $\equiv$ 
    bool MatrixInverse::encloseLS(iVector &x,
        const pMatrix &A, const iVector &b0,
        const iMatrix &B, double beta)
    {
        /*
             $A$  is the matrix from  $Ax = b$ 
             $b0 \ni Cb$ 
             $B \ni I - CA$ 
             $beta \geq \|I - CA\|_\infty$ ,  $beta$  must be  $< 1$ 
        */


---


         $x$  contains the solution to  $Ax = b$  if true is returned
        /*
            <compute initial box 161>
            <do Krawczyk's iteration 165>
            return true;
        */
    }

```

This code is used in chunk 173.

14.3.1 Initial box

The input $b0$ contains Cb , and the input B contains $I - CA$. Also, $beta \geq \|I - CA\|_\infty$. We compute $a \geq \|Cb\|_\infty / (1 - \beta)$ and set the initial box containing the solution.

```

161 <compute initial box 161>  $\equiv$ 
    v_bias::round_down();
    double a = 1 - beta; /*  $a \leq 1 - \beta$  */
    double a1 = inf_normV(b0); /*  $a1 \geq \|Cb\|_\infty$  */
    v_bias::round_up();
    a = a1 / a; /*  $a \geq \|Cb\|_\infty / (1 - \beta)$  */
    setV(x, v_bias::interval(-a, a));

```

This code is used in chunk 160.

14.3.2 Krawczyk's iteration

We implement $Cb + (I - CA)x^{(i)}$. That is, we program the expression $x1 = b0 + B*x$.

```

162 <Krawczyk's iteration 162>  $\equiv$ 
    multMiVi(x1, B, x);
    addViVi(x1, b0);

```

See also chunk 163.

This code is used in chunk 165.

Now we can intersect $x1$ and x . The result is stored in x , if $b \equiv true$. Otherwise, we return *false*.

```

163 <Krawczyk's iteration 162> +=
    bool b = intersect(x, x, x1);
    if (b == false) {
#ifdef VNODE_DEBUG
        printMessage("x_and_x1_do_not_intersect");
#endif
        return false;
    }

```

For the **while** loop that follows, we need to compute the sum of the radii of the components of *x*.

```

164 <sum radii 164> ==
    rad(radx, x);
    v_bias::round_nearest();
    sum_radii = accumulate(radx.begin(), radx.end(), 0.0);

```

This code is used in chunk 165.

Now, we compose the whole iteration. Initially, we set *sum_old_radii* to the largest machine number. The **while** loop below iterates as far as *sum_radii* < *mult* * *sum_old_radii*. The factor *mult* = (1 + *beta*)/2 is the same as in [10].

```

165 <do Krawczyk's iteration 165> ==
    double sum_old_radii = numeric_limits<double>::max();
    double sum_radii;
    <sum radii 164>
    round_up();
    int max_iterations = 20;
    int counter = 0;
    double mult = (1 + beta)/2;
    while (sum_radii < mult * sum_old_radii & counter < max_iterations) {
        <Krawczyk's iteration 162>
        sum_old_radii = sum_radii;
        <sum radii 164>
        counter++;
    }
    iterations = counter;

```

This code is used in chunk 160.

14.4 Enclosing the inverse of a general point matrix

To enclose the inverse of a point matrix *A*, we enclose the solution to

$$Ax_i = e_i \quad \text{for } i = 1, \dots, n,$$

where *e_i* is the *i*th unit vector. Then, the *i*th column of the inverse is *x_i*.

First we compute a floating-point inverse of A . If *invertMatrix* fails, *encloseMatrixInverse* returns *false*. Then we compute *beta* that is needed for *encloseLS* and call this function for each e_i .

```

166 <enclose the inverse of a matrix 166> ≡
    bool MatrixInverse :: encloseMatrixInverse(iMatrix &Ainv,
        const pMatrix &A)
    {
        bool b = invertMatrix(C, A);
        if (b ≡ false) return false;
        <find beta 167>
        <enclose each column 169>
        #ifdef VNODE_DEBUG
            iMatrix B = Ainv;
            multMiMp(B, Ainv, A);
            int n = sizeM(A);
            for (int i = 0; i < n; i++)
                for (int j = 0; j < n; j++) {
                    interval b = B[i][j];
                    if (i ≡ j) assert(v_bias :: subseteq(interval(1.0), b));
                    else assert(v_bias :: subseteq(interval(0.0), b));
                }
        #endif
        return true;
    }

```

This code is used in chunk 173.

We compute B such that it contains $I - CA$.

```

167 <find beta 167> ≡
    assignM(Ci, C);
    multMiMp(B, Ci, A);
    subFromId(B);

```

See also chunk 168.

This code is used in chunks 166 and 170.

Now we find $\beta \geq \beta$. If $\beta \geq 1$, we return *false*, as Krawczyk's iteration cannot proceed.

```

168 <find beta 167> +≡
    v_bias :: round_up();
    double beta = inf_normM(B);
    if (beta ≥ 1) return false;

```

Finally, we can call *encloseLS* for each right side e_i . *getColumn* extracts the i th column of C in $b0$. *setColumn* sets the interval vector containing the corresponding solution to $Ax = e_i$.

```

169 <enclose each column 169> ≡
    for (unsigned int i = 0; i < sizeM(A); i++) {
        getColumn(b0, C, i);
        bool b = encloseLS(x, A, b0, B, beta);
        if (b ≡ false) {
#ifdef VNODE_DEBUG
            printMessage("Could not enclose the solution to a linear system");
#endif
            return false;
        }
        setColumn(Ainv, x, i);
    }

```

This code is used in chunks 166 and 170.

14.5 Enclosing the inverse of an orthogonal matrix

We have a floating-point approximation for an orthogonal matrix. Normally, its transpose is not the same as the inverse of this matrix. Hence, we need to enclose it.

```

170 <enclose the inverse of an orthogonal matrix 170> ≡
    bool MatrixInverse::orthogonalInverse(iMatrix &Ainv,
        const pMatrix &A)
    {
        transpose(C, A);
        <find beta 167>
        <enclose each column 169>
        return true;
    }

```

This code is used in chunk 173.

14.6 Constructor and destructor

```

171 <MatrixInverse constructor/destructor 171> ≡
    MatrixInverse::MatrixInverse(int n)
    {
        M = new double[n * n];
        ipiv = new int[n];
        lwork = 2 * n;
        work = new double[lwork];
        sizeV(radix, n);
        sizeV(b0, n);
        sizeV(x1, n);
        sizeV(x, n);
        sizeM(B, n);
        sizeM(Ci, n);
    }

```

```

    sizeM(C,n);
}
MatrixInverse::~MatrixInverse()
{
    delete[] work;
    delete[] ipiv;
    delete[] M;
}

```

This code is used in chunk 173.

Files

```

172 <matrixinverse.h 172> ≡
    #ifndef MATRIXINVERSE_H
    #define MATRIXINVERSE_H
    #include <numeric>
    #include "vnodeinterval.h"
    #include "vnoderound.h"
    #include "vector_matrix.h"
    namespace v_blas {
        using namespace v_bias;
        <class MatrixInverse 156>
    }
    #endif

173 <matrixinverse.cc 173> ≡
    #include <climits>
    #include "matrixinverse.h"
    #include "basiclinalg.h"
    #include "intvfuncs.h"
    #include "debug.h"
    using namespace std;
    namespace v_blas {
        <MatrixInverse constructor/destructor 171>
        <compute  $A^{-1}$  158>
        <enclose the solution to  $Ax = b$  160>
        <enclose the inverse of a matrix 166>
        <enclose the inverse of an orthogonal matrix 170>
    }

```


Part IV

Solver Implementation

Chapter 15

Structure

We list the classes in VNODE-LP along with brief descriptions. These classes are depicted in Figure 15.1.

Solution contains a representation of the solution at each time point.

Apriori contains a representation of an enclosure of the solution over an integration step.

Control stores various data for controlling an integration.

AD_ODE provides functions for generating Taylor coefficients for the solution to an ODE.

AD_VAR provides functions for generating Taylor coefficients for the solution to the variational equation.

AD aggregates objects of **AD_ODE** and **AD_VAR**.

FadbadODE, **FadbadVarODE**, and **FADBAD_AD** are implementations of **AD_ODE**, **AD_VAR**, and **AD**, respectively.

HOE implements the High-Order Enclosure method [27] for enclosing the solution to an ODE and computing a priori bounds.

IHO implements the Interval Hermite-Obreschkoff method [23] for computing tight bounds on the solution.

VNODE implements the overall integrator.

MatrixInverse provides functions for computing the inverse of a matrix.

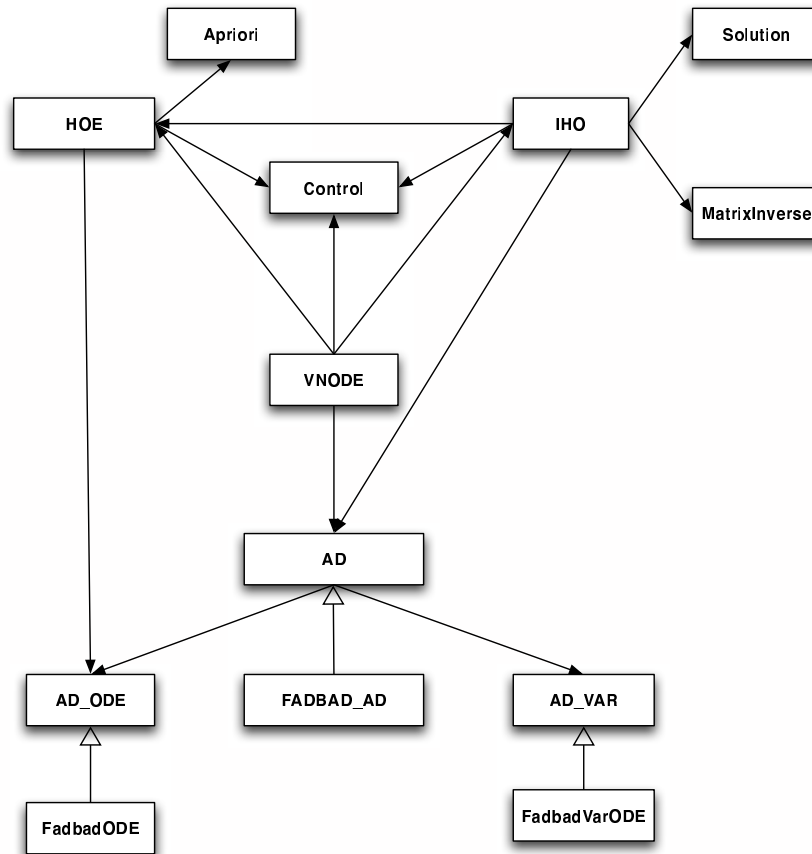


Figure 15.1. *Classes in VNODE-LP. The triangle arrows denote inheritance relations; the normal arrows denote uses relations.*

Chapter 16

Solution enclosure representation

The present solver implements a one-step method. The initial point t_0 and endpoint t_{end} are stored as intervals. For example, if t_0 is representable, the interval is $\mathbf{t}_0 = [t_0, t_0]$; otherwise, \mathbf{t}_0 is an interval containing t_0 . The interval containing t_{end} is denoted by \mathbf{t}_{end} .

Internally, VNODE-LP select points that are machine-representable numbers. For simplicity in the considerations that follow, we denote such a point by the interval \mathbf{t}_j , which contains t_j . At \mathbf{t}_0 , the user provides \mathbf{y}_0 . On each integration step, we maintain two types of representations of an enclosure of an ODE solution: tight enclosure at \mathbf{t}_j and an a priori enclosure over $[\underline{\mathbf{t}}_j, \bar{\mathbf{t}}_{j+1}]$ (or $[\underline{\mathbf{t}}_{j+1}, \bar{\mathbf{t}}_j]$ if the integration is in negative direction).

16.1 Tight enclosure

We maintain three representations as enclosures on the solution at \mathbf{t}_j :

$$\begin{aligned} &\{u_j + S_j \alpha + A_j r \mid \alpha \in \boldsymbol{\alpha}, r \in r_j\} \\ &\{u_j + S_j \alpha + Q_j r \mid \alpha \in \boldsymbol{\alpha}, r \in r_{\text{QR},j}\}, \quad \text{and} \\ &\mathbf{y}_j. \end{aligned}$$

Here, $u_j, \alpha, r \in \mathbb{R}^n$; $S_j, A_j, Q_j \in \mathbb{R}^{n \times n}$; and $\mathbf{r}_j, \mathbf{r}_{\text{QR},j}, \boldsymbol{\alpha}, \mathbf{y}_j \in \mathbb{I}\mathbb{R}^n$. The first set corresponds to the P (parallelepiped) method, and the second set corresponds to the QR method [24]. The use of these sets is discussed in detail in Chapter 20.

178 \langle tight enclosure representation 178 $\rangle \equiv$

```
class Solution {
public:
    Solution(int n);
    void init(const v_bias::interval &t0, const iVector &y0);
    v_bias::interval t;
    pVector u;
```

```

    iVector y;
    iVector alpha, r, rQR;
    pMatrix S, A, Q;
};

```

This code is used in chunk 184.

The constructor of the **Solution** class allocates memory for the vector and matrix members of this class.

```

179 <create Solution 179> ≡
    Solution::Solution(int n)
    {
        sizeV(u, n);
        sizeV(alpha, n);
        sizeV(r, n);
        sizeV(rQR, n);
        sizeV(y, n);
        sizeM(S, n);
        sizeM(A, n);
        sizeM(Q, n);
    }

```

This code is used in chunk 185.

Initially, when $j = 0$, we set

$$\begin{aligned}
 u_0 &= m(\mathbf{y}_0), \\
 \boldsymbol{\alpha} &= \mathbf{y}_0 - u_0, \\
 \mathbf{r}_0 &= 0, \\
 \mathbf{r}_{\text{QR},0} &= 0, \\
 S_0 &= I, \\
 A_0 &= I, \quad \text{and} \\
 Q_0 &= I.
 \end{aligned}$$

When computing the midpoint of $y\theta = \mathbf{y}_0$, we find a point vector u that is the rounded to the nearest true midpoint of $y\theta$.

```

180 <initialize Solution 180> ≡
    void Solution::init(const v_bias::interval &t0, const iVector &y0)
    {
        t = t0;
        y = y0;
        midpoint(u, y0);    /* u ≈ m(y0) */
        subViVp(alpha, y0, u); /* alpha ⊇ y0 - u */
        setV(r, 0.0);
        setV(rQR, 0.0);
        setId(S);
    }

```

```

    setId(Q);
    setId(A);
}

```

This code is used in chunk 185.

16.2 A priori enclosure

On each integration step, we validate existence and uniqueness of a solution and compute a priori bounds $\tilde{\mathbf{y}}_j$ such that

$$y(t; t_j, y_j) \in \tilde{\mathbf{y}}_j$$

for all $t \in \tilde{t}_j := t_j + [0, 1]h_j$ and all $y_j \in \mathbf{y}_j$. Here, $h_j \in \mathbb{R}$ is a stepsize selected by VNODE-LP.

The interval \tilde{t}_j and the interval vector $\tilde{\mathbf{y}}_j$ are stored in

```

181 <a priori enclosure representation 181> ≡
    class Apriori {
    public:
        Apriori(int n);
        void init(const v_bias::interval t0, const iVector y0);
        v_bias::interval t;
        iVector y;
    };

```

This code is used in chunk 184.

The constructor allocates a vector y

```

182 <create Apriori 182> ≡
    Apriori::Apriori(int n) {
        sizeV(y, n);
    }

```

This code is used in chunk 185.

When initializing an **Apriori** object, we set t and y .

```

183 <initialize Apriori 183> ≡
    void Apriori::init(const v_bias::interval t0, const iVector y0) {
        t = t0;
        y = y0;
    }

```

This code is used in chunk 185.

Files

```

184 <solution.h 184> ≡
    #ifndef SOLUTION_H

```

```
#define SOLUTION_H
    namespace vnodelp {
        using namespace v_bias;
        using namespace v_blas;

        ⟨tight enclosure representation 178⟩
        ⟨a priori enclosure representation 181⟩
    }
#endif

185 ⟨solution.cc 185⟩ ≡
    #include "vnodeinterval.h"
    #include "basiclinalg.h"
    #include "solution.h"
    #include "intvfuncs.h"
    namespace vnodelp {
        ⟨create Solution 179⟩
        ⟨initialize Solution 180⟩
        ⟨create Apriori 182⟩
        ⟨initialize Apriori 183⟩
    }
```


Chapter 17

Taylor coefficient computation

We provide an abstract class, **AD_ODE**, for generating Taylor coefficients (TCs) for the solution to an ODE and an abstract class, **AD_VAR**, for generating TCs for the solution to the ODE's variational equation. Concrete implementations of these classes using FADBAD++ [6] are discussed in Chapter 22.

17.1 Taylor coefficients for an ODE solution

```
187 <AD_ODE 187> ≡  
    class AD_ODE {  
    public:  
        virtual void set(const v_bias::interval &t0, const iVector &y0,  
            const v_bias::interval &h, int k) = 0;  
        virtual void compTerms() = 0;  
        virtual void sumTerms(iVector &sum, int m) = 0;  
        virtual void getTerm(iVector &term, int i) const = 0;  
        virtual v_bias::interval getStepsize() const = 0;  
        virtual void eval(void *param) = 0;  
        virtual ~AD_ODE() {}  
    };
```

This code is used in chunk 194.

Brief descriptions of the above functions follow. The k th TC of the solution to (1.1) at (t^*, y^*) is denoted by $f^{[i]}(t^*, y^*)$; cf. [23].

set initializes a TC computation by setting a point of expansion \mathbf{t}_0 , \mathbf{y}_0 , enclosure on the stepsize \mathbf{h} , and order k .

compTerms encloses

$$\mathbf{h} f^{[1]}(\mathbf{t}_0, \mathbf{y}_0), \mathbf{h}^2 f^{[2]}(\mathbf{t}_0, \mathbf{y}_0), \dots, \mathbf{h}^k f^{[k]}(\mathbf{t}_0, \mathbf{y}_0).$$

sumTerms encloses $\sum_{i=0}^m \mathbf{h}^i f^{[i]}(\mathbf{t}_0, \mathbf{y}_0)$, where $m \leq k$. The result is stored in the parameter *sum*.

getTerm obtains the *i*th term, $\mathbf{h}^i f^{[i]}(\mathbf{t}_0, \mathbf{y}_0)$, where $i \leq k$.

getStepsize returns *h*.

eval evaluates $f(t, y)$ and rebuilds the computational graph. *param* is a pointer to parameters that can be passed to *f*.

17.2 Taylor coefficients for the solution of the variational equation

```
189 <AD_VAR 189> ≡
    class AD_VAR {
    public:
        virtual void set(const v_bias::interval &t0, const iVector &y0,
            const v_bias::interval &h, int k) = 0;
        virtual void compTerms() = 0;
        virtual void sumTerms(iMatrix &sum, int m) = 0;
        virtual void getTerm(iMatrix &term, int i) const = 0;
        virtual void eval(void *param) = 0;
        virtual ~AD_VAR() {}
    };

```

This code is used in chunk 195.

Brief descriptions of the above functions follow.

set initializes a TC computation by setting a point of expansion $\mathbf{t}_0, \mathbf{y}_0$, stepsize *h*, and order *k*.

compTerms encloses

$$\mathbf{h} \frac{\partial f^{[1]}}{\partial \mathbf{y}}(\mathbf{t}_0, \mathbf{y}_0), \mathbf{h}^2 \frac{\partial f^{[2]}}{\partial \mathbf{y}}(\mathbf{t}_0, \mathbf{y}_0), \dots, \mathbf{h}^k \frac{\partial f^{[k]}}{\partial \mathbf{y}}(\mathbf{t}_0, \mathbf{y}_0).$$

sumTerms encloses

$$\sum_{i=0}^m \mathbf{h}^i \frac{\partial f^{[i]}}{\partial \mathbf{y}}(\mathbf{t}_0, \mathbf{y}_0),$$

where $m \leq k$. The result is stored in the parameter *sum*.

getTerm obtains the *i*th term, $\mathbf{h}^i \frac{\partial f^{[i]}}{\partial \mathbf{y}}(\mathbf{t}_0, \mathbf{y}_0)$, where $i \leq k$.

eval evaluates $f(t, y)$ and rebuilds the computational graph. *param* is a pointer to parameters that can be passed to *f*.

17.3 AD class

It is convenient to encapsulate the above classes into

```
191 <encapsulated AD 191> ≡
    class AD {
    public:
        AD(int n, AD_ODE *a, AD_VAR *av);
        void eval(void *p);
        virtual int getMaxOrder() const = 0;
        /* maximum order that is allowed */
    public:
        int size;
        AD_ODE *tayl_coeff_ode;
        AD_VAR *tayl_coeff_var;
    };

```

This code is used in chunk 196.

The constructor of **AD** sets the size of the problem and pointers to objects of **AD_ODE** and **AD_VAR**.

```
192 <implementation of encapsulated AD 192> ≡
    inline AD::AD(int n, AD_ODE *a, AD_VAR *av)
        : size(n), tayl_coeff_ode(a), tayl_coeff_var(av) {}

```

See also chunk 193.

This code is used in chunk 196.

The *eval* function calls the corresponding *eval* functions of **AD_ODE** and **AD_VAR**.

```
193 <implementation of encapsulated AD 192> +≡
    inline void AD::eval(void *p)
    {
        tayl_coeff_ode->eval(p);
        tayl_coeff_var->eval(p);
    }

```

Files

The above classes are stored in **ad_ode.h**, **ad_var.h**, and **allad.h**, respectively.

```
194 <ad_ode.h 194> ≡
    #ifndef AD_ODE_H
    #define AD_ODE_H
    #include "vnodeinterval.h"
    #include "vector_matrix.h"
    namespace vnodelp {
        using namespace v_bias;
        using namespace v_blas;
    }

```

```

    <AD_ODE 187>
  }
#endif

195 <ad_var.h 195> ≡
    #ifndef AD_VAR_H
    #define AD_VAR_H
    #include "vnodeinterval.h"
    #include "vector_matrix.h"
    namespace vnodelp {
        using namespace v_bias;
        using namespace v_blas;
        <AD_VAR 189>
    }
#endif

196 <allad.h 196> ≡
    #ifndef ALLAD_H
    #define ALLAD_H
    #include "ad_ode.h"
    #include "ad_var.h"
    namespace vnodelp {
        <encapsulated AD 191>
        <implementation of encapsulated AD 192>
    }
#endif

```

Chapter 18

Control data

We store various data needed to control an integration in a **Control** class. An integration in VNODE-LP is carried out by the *integrate* function of **VNODE**; see Chapter 21.

18.1 Indicator type

First, we introduce an enumerated **Ind** data type, where a variable of this type can take the following values:

value	description
<i>first_entry</i>	indicates a first entry into <i>integrate</i>
<i>success</i>	<i>integrate</i> has reached t_{end} successfully
<i>failure</i>	an error has occurred in <i>integrate</i>

```
198 <indicator type 198> ≡  
    typedef enum {  
        first_entry, success, failure  
    } Ind;
```

This code is used in chunk 201.

18.2 Interrupt type

Similarly, we have an **Interrupt** type, where a variable can take values as described below.

value	description
<i>no</i>	<i>integrate</i> tries to reach t_{end}
<i>before_accept</i>	<i>integrate</i> takes a step and returns before accepting this step

```

199 <interrupt type 199> ≡
    typedef enum {
        no, before_accept
    } Interrupt;

```

This code is used in chunk 201.

18.3 Control data

In the **Control** class, we store the following data:

name	default value	description
<i>ind</i>	<i>first_entry</i>	indicator variable
<i>interrupt</i>	<i>no</i>	indicates if interrupts are requested
<i>order</i>	20	order of the method
<i>atol</i>	10^{-12}	absolute error tolerance
<i>rtol</i>	10^{-12}	relative error tolerance
<i>hmin</i>	0	magnitude of the minimum stepsize allowed. If a positive value is set by the user, this value for <i>hmin</i> will be used in an integration. Otherwise, the solver computes a minimum stepsize as discussed in Subsection 21.2.3.

```

200 <control class 200> ≡
    class Control {
    public:
        Ind ind;
        Interrupt interrupt;
        unsigned int order;
        double atol, rtol;
        double hmin;
        Control() :
            ind(first_entry),
            interrupt(no),
            order(20),
            atol( $1 \cdot 10^{-12}$ ), rtol( $1 \cdot 10^{-12}$ ),
            hmin(0) {}
    };

```

This code is used in chunk 201.

Files

We store **Ind**, **Interrupt**, and **Control** in

```
201 <control.h 201> ≡  
    #ifndef CONTROL_H  
    #define CONTROL_H  
        namespace vndelp {  
            <indicator type 198>  
            <interrupt type 199>  
            <control class 200>  
        }  
    #endif
```


Chapter 19

Computing a priori bounds

We have implemented the HOE method [27] to compute a priori bounds on the solution of an ODE problem. In Section 19.1 we summarize the relevant theory. The **HOE** class is given in Section 19.2. The implementation of the HOE method is in Section 19.3. The rest of the functions of this class is in Section 19.4.

19.1 Theory background

The HOE method is based on the following two results; cf. [9, 27].

1. If h_j and $\tilde{\mathbf{y}}_j$ are such that $\mathbf{y}_j \subseteq \text{int}(\tilde{\mathbf{y}}_j)$ and

$$\sum_{i=0}^{k-1} (t - t_j)^i f^{[i]}(t_j, y_j) + (t - t_j)^k f^{[k]}(t_j + [0, 1]h_j, \tilde{\mathbf{y}}_j) \subseteq \tilde{\mathbf{y}}_j$$

for all $t \in t_j + [0, 1]h_j$ and all $y_j \in \mathbf{y}_j$, then

$$y' = f(t, y), \quad y(t_j) = y_j \in \mathbf{y}_j \tag{19.1}$$

has a unique solution

$$y(t; t_j, y_j) \in \tilde{\mathbf{y}}_j$$

for all $t \in t_j + [0, 1]h_j$ and all $y_j \in \mathbf{y}_j$.

2. Let $h_{j,0} \neq 0$ and let \mathbf{p}_j be an interval vector enclosing the set

$$\mathcal{P}_j = \left\{ \sum_{i=0}^{k-1} (t - t_j)^i f^{[i]}(t_j, y_j) \mid t \in t_j + [0, 1]h_{j,0}, y_j \in \mathbf{y}_j \right\}. \tag{19.2}$$

That is, $\mathcal{P}_j \subseteq \mathbf{p}_j$.

Let \mathbf{u}_j be such that

$$\begin{aligned}\tilde{\mathbf{y}}_j &= \mathbf{p}_j + \mathbf{u}_j, \quad \text{and} \\ \mathbf{y}_j &\subseteq \text{int}(\tilde{\mathbf{y}}_j).\end{aligned}\tag{19.3}$$

If $h_{j,1} \neq 0$ is such that

$$[0, 1]h_{j,1}^k f^{[k]}(t_j + [0, 1]h_{j,0}, \tilde{\mathbf{y}}_j) \subseteq \mathbf{u}_j,\tag{19.4}$$

and

$$h_j = \text{sign}(h_{j,0}) \cdot \min\{|h_{j,0}|, |h_{j,1}|\},\tag{19.5}$$

then there exists a unique solution $y(t; t_j, y_j)$ to (19.1) for all $t \in t_j + [0, 1]h_j$ and all $y_j \in \mathbf{y}_j$. Moreover,

$$y(t; t_j, y_j) \in \tilde{\mathbf{y}}_j \quad \text{for all } t \in t_j + [0, 1]h_j \quad \text{and all } y_j \in \mathbf{y}_j.$$

19.2 The HOE class

The class implementing the HOE method is

```
204 < class HOE 204 > ≡
    class HOE {
    public:
        HOE(int n);
        void compAprioriEnclosure(const interval &t0, const iVector &y0,
            bool &info);
        void acceptSolution();
        < set functions HOE 221 >
        < get functions HOE 222 >
        ~HOE();
    private:
        Apriori *apriori_trial, *apriori;
        Control *control;
        AD_ODE *tayl_coeff;
        double h, h_next, h_trial, t_trial;
        int order_trial;
        iVector term, p, u, v;
        const interval one; /* one = [0, 1] */
        interval comp_beta(const iVector &v, const iVector &u, int k);
    };
```

This code is used in chunk 226.

19.3 Implementation of the HOE method

19.3.1 Computing p_j

We have to enclose \mathcal{P}_j in (19.2). On the first step, if t_0 is not a representable machine number, an interval \mathbf{t}_0 containing t_0 would be given. We assume that, in general, $t_j \in \mathbf{t}_j$. Then

$$t \in \mathbf{t}_j + [0, 1]h_{j,0} \quad \text{and} \quad t - t_j \in \mathbf{t}_j - \mathbf{t}_j + [0, 1]h_{j,0}.$$

The width of an interval \mathbf{a} is $w(\mathbf{a}) = \overline{\mathbf{a}} - \underline{\mathbf{a}}$. (Width of an interval vector is defined componentwise.) Denote $\mathbf{t}_j^* = (\mathbf{t}_j - \mathbf{t}_j)/h_{j,0} + [0, 1]$. Then

$$\mathbf{t}_j^* = \begin{cases} (\mathbf{t}_j - \mathbf{t}_j)/h_{j,0} + [0, 1] & \text{if } w(\mathbf{t}_j) > 0, \\ [0, 1] & \text{if } w(\mathbf{t}_j) = 0. \end{cases}$$

Hence $t - t_j \in \mathbf{t}_j^* h_{j,0}$, and we can bound \mathcal{P}_j as

$$\mathcal{P}_j \subseteq \sum_{i=0}^{k-1} \mathbf{t}_j^{*i} h_{j,0}^i f^{[i]}(\mathbf{t}_j, \mathbf{y}_j).$$

We enclose first $h_{j,0}^i f^{[i]}(\mathbf{t}_j, \mathbf{y}_j)$ for $i = 1, \dots, k-1$. Then, we use an interval form of Horner's rule to compute

$$\mathbf{p}_j := \mathbf{y}_j + \mathbf{t}_j^* \left(h_{j,0} f^{[1]}(\mathbf{t}_j, \mathbf{y}_j) + \dots + \mathbf{t}_j^* (h_{j,0}^{k-2} f^{[k-2]}(\mathbf{t}_j, \mathbf{y}_j) + \mathbf{t}_j^* h_{j,0}^{k-1} f^{[k-1]}(\mathbf{t}_j, \mathbf{y}_j)) \dots \right).$$

206 `< compute \mathbf{p}_j 206 > \equiv /*`

$$t_0 = \mathbf{t}_j$$

$$y0 \supseteq \mathbf{y}_j$$

$$h_trial = h_{j,0}$$

$$order_trial = k$$

$$taylor_coeff \text{ contains enclosures on } h_{j,0}^i f^{[i]}(\mathbf{t}_j, \mathbf{y}_j) \text{ for } i = 0, \dots, k-1$$

$$t_enc \supseteq \mathbf{t}_j^* = (\mathbf{t}_j - \mathbf{t}_j)/h_{j,0} + [0, 1]$$

$$p \supseteq \mathbf{p}_j$$

`*/`

`taylor_coeff \rightarrow set($t0, y0, h_trial, order_trial - 1$);`

`taylor_coeff \rightarrow compTerms();`

interval `t_enc = ($t0 - t0$)/ $h_trial + one$;`

`taylor_coeff \rightarrow getTerm($p, order_trial - 1$);`

for (**int** `i = order_trial - 2; i \geq 0; i--`) {

`scaleV(p, t_enc);`

`taylor_coeff \rightarrow getTerm($term, i$);`

`addViVi($p, term$);`

}

This code is used in chunk 217.

19.3.2 Computing \mathbf{u}_j and $\tilde{\mathbf{y}}_j$

If \mathbf{u}_j is symmetric with no component equal to $[0, 0]$ then (19.3) holds. Denote

$$\text{tol}_j = \text{rtol} \cdot \|\mathbf{y}_j\| + \text{atol}.$$

Throughout this manuscript, we shall assume the infinity norm, unless stated otherwise. For an interval vector \mathbf{a} ,

$$\|\mathbf{a}\| = \max_i \{|\underline{a}_i|, |\overline{a}_i|\}.$$

Let \mathbf{u}_j be the n -vector with each component $h_{j,0}[-\text{tol}_j/2, \text{tol}_j/2]$. Then, we form

$$\tilde{\mathbf{y}}_j = \mathbf{p}_j + \mathbf{u}_j.$$

207 $\langle \text{compute } \mathbf{u}_j \text{ and } \tilde{\mathbf{y}}_j \text{ } 207 \rangle \equiv \quad / *$

$$\begin{array}{c} y0 \supseteq \mathbf{y}_j \\ \text{control-atol} = \text{atol} \\ \text{control-rtol} = \text{rtol} \\ h_trial = h_{j,0} \\ p \supseteq \mathbf{p}_j \\ \hline \text{tol} \geq \text{rtol} \cdot \|\mathbf{y}_j\| + \text{atol} \\ u \supseteq \mathbf{u}_j \\ \text{apriori_trial} \supseteq \tilde{\mathbf{y}}_j = \mathbf{p}_j + \mathbf{u}_j. \end{array}$$

```

*/
round_up();
double tol = inf_norm V(y0) * control-rtol + control-atol;
set V(u, h_trial * interval(-tol/2, tol/2));
add Vi Vi (apriori_trial-y, p, u);
assert (interior(y0, apriori_trial-y));

```

This code is used in chunk 217.

Motivation for the choice of \mathbf{u}_j

Denote

$$\mathbf{z}_j = f^{[k]}(t_j + [0, 1]h_{j,0}\tilde{\mathbf{y}}_j).$$

We can consider $|h_{j,1}|^k \|\mathbf{w}(\mathbf{z}_j)\|$ as an estimate of the *local excess* at \mathbf{t}_{j+1} that we introduce on the step from \mathbf{t}_j to \mathbf{t}_{j+1} .

The magnitude of an interval \mathbf{a} is $|\mathbf{a}| = \max\{|\underline{\mathbf{a}}|, |\overline{\mathbf{a}}|\}$. Magnitude of an interval vector is defined componentwise. From (19.4), (19.5), and the choice for \mathbf{u}_j ,

$$\begin{aligned} |h_{j,1}|^k w(\mathbf{z}_j) &\leq w([0, 1]h_{j,1}^k \mathbf{z}_j) = |h_{j,1}|^k |\mathbf{z}_j| \\ &\leq w(\mathbf{u}_j) = h_{j,0} (\text{tol}_j, \text{tol}_j, \dots, \text{tol}_j)^T. \end{aligned}$$

Hence,

$$|h_{j,1}|^k \|w(\mathbf{z}_j)\| \leq h_{j,0} \cdot \text{tol}_j.$$

If $h_{j,1} = h_{j,0}$, we have a local excess per unit step (LEPUS) control. If $h_{j,1} < h_{j,0}$, then usually $h_{j,1}$ is not much smaller than $h_{j,0}$, and we have almost LEPUS.

19.3.3 Computing a stepsize

We enclose first

$$\mathbf{v}_j = h_{j,0}^k f^{[k]}(\mathbf{t}_j + [0, 1]h_{j,0}\tilde{\mathbf{y}}_j).$$

210 `< compute stepsize 210 > ≡` `/*`

$$\begin{aligned} t0 &= \mathbf{t}_j \\ \text{apriori_trial} &\supseteq \tilde{\mathbf{y}}_j \\ h_trial &= h_{j,0} \\ \text{order_trial} &= k \end{aligned}$$

$$\begin{aligned} \text{taylor_coeff} &\text{ contains enclosures on } h_{j,0}^i f^{[i]}(\mathbf{t}_j + [0, 1]h_{j,0}\tilde{\mathbf{y}}_j) \text{ for } i = 0, \dots, k \\ v &\supseteq h_{j,0}^k f^{[k]}(\mathbf{t}_j + [0, 1]h_{j,0}\tilde{\mathbf{y}}_j) \end{aligned}$$

`*/`

```
taylor\_coeff←set(t0 + one * h\_trial, apriori\_trial+y, h\_trial, order\_trial);
taylor\_coeff←compTerms();
taylor\_coeff←getTerm(v, order\_trial);
```

See also chunks 212 and 214.

This code is used in chunk 217.

Now we consider two cases: $\underline{\mathbf{t}}_j = \bar{\mathbf{t}}_j$ and $\underline{\mathbf{t}}_j < \bar{\mathbf{t}}_j$.

The case $\underline{\mathbf{t}}_j = \bar{\mathbf{t}}_j$

Let $\gamma > 0$ be the largest number such that

$$\begin{aligned} [0, 1]h_{j,1}^k f^{[k]}(\mathbf{t}_j + [0, 1]h_{j,0}\tilde{\mathbf{y}}_j) &= [0, 1](\gamma h_{j,0}^k) f^{[k]}(\mathbf{t}_j + [0, 1]h_{j,0}\tilde{\mathbf{y}}_j) \\ &= [0, \gamma] \mathbf{v}_j \\ &\subseteq \mathbf{u}_j. \end{aligned}$$

Such a γ exists since \mathbf{u}_j is symmetric. Then, we write $\beta = \gamma^{1/k}$ and set

$$h_j = \begin{cases} \beta h_{j,0} & \text{if } \beta < 1 \\ h_{j,0} & \text{if } \beta \geq 1. \end{cases}$$

In practice, we compute $\beta \ni \beta > 0$, with $\underline{\beta} \geq 0$ and find

$$h_j = \begin{cases} \downarrow(\underline{\beta} h_{j,0}) & \text{if } \underline{\beta} < 1, h_{j,0} > 0 \\ \uparrow(\underline{\beta} h_{j,0}) & \text{if } \underline{\beta} < 1, h_{j,0} < 0 \\ h_{j,0} & \text{if } \underline{\beta} \geq 1. \end{cases} \quad (19.6)$$

Here \uparrow denotes rounding up, and \downarrow denotes rounding down.

212 $\langle \text{compute stepsize } 210 \rangle + \equiv \quad / *$

$$\frac{\begin{array}{l} u \supseteq \mathbf{u}_j \\ v \supseteq \mathbf{v}_j = h_{j,0}^k f^{[k]}(\mathbf{t}_j + [0, 1]h_{j,0}, \tilde{\mathbf{y}}_j) \end{array}}{h = h_j}$$

$*/$

```

interval beta = comp_beta(v, u, order_trial);
assert(inf(beta) ≥ 0);
if (inf(beta) < 1) {
  if (h_trial > 0) round_down();
  else round_up();
  h = inf(beta) * h_trial;
}
else h = h_trial;

```

The case $\underline{\mathbf{t}}_j < \bar{\mathbf{t}}_j$

Similar to (19.4), we try to find (the largest) $|h_{j,1}| \leq |h_{j,0}|$ such that

$$\begin{aligned} (t - \mathbf{t}_j)^k f^{[k]}(\mathbf{t}_j + [0, 1]h_{j,0}) &\subseteq (\mathbf{t}_j - \mathbf{t}_j + [0, 1]h_{j,1})^k f^{[k]}(\mathbf{t}_j + [0, 1]h_{j,0}, \tilde{\mathbf{y}}_j) \\ &\subseteq \mathbf{u}_j. \end{aligned} \quad (19.7)$$

Normally $|h_{j,0}| \gg |\mathbf{t}_j - \mathbf{t}_j|$, and therefore,

$$(\mathbf{t}_j - \mathbf{t}_j + [0, 1]h_{j,0})^k \approx [0, 1]h_{j,0}^k.$$

Hence, we may consider setting h_j as in (19.6). However, (19.7) may not hold. To find h_j such that it is likely to hold, we set

$$h_j = 0.9h_{j,0}$$

and check if

$$\begin{aligned}
& ((\mathbf{t}_j - \mathbf{t}_j) + [0, 1]h_j)^k f^{[k]}(\mathbf{t}_j + [0, 1]h_{j,0}, \tilde{\mathbf{y}}_j) \\
&= \left(\frac{\mathbf{t}_j - \mathbf{t}_j + [0, 1]h_j}{h_{j,0}} \right)^k h_{j,0}^k f^{[k]}(\mathbf{t}_j + [0, 1]h_{j,0}, \tilde{\mathbf{y}}_j) \\
&= \left(\frac{\mathbf{t}_j - \mathbf{t}_j + [0, 1]h_j}{h_{j,0}} \right)^k \mathbf{v}_j \subseteq \mathbf{u}_j.
\end{aligned}$$

If the above inclusion test fails, we reduce h_j and repeat as shown below.

```

214 < compute stepsize 210 > +=
    if (inf(t0) < sup(t0)) {
        v ⊇ v_j
        h = h_j
        h_trial = h_{j,0}
        t0 = t_j
        control-hmin = h_min
        -----
        tt ⊇ t_j - t_j
        t_enc ⊇ (t_j - t_j + [0, 1]h_j) / h_{j,0}
        v ⊇ ( (t_j - t_j + [0, 1]h_j) / h_{j,0} )^k h_{j,0}^k f^{[k]}(t_j + [0, 1]h_{j,0}, y_tilde_j)
        -----
        h = h_j
    }
    */
interval tt = t0 - t0;
while (fabs(h) > control-hmin) {
    t_enc = (tt + one * h) / h_trial;
    scale V(v, pow(t_enc, order_trial));
    if (subteq(v, u)) break;
    h = 0.9 * h;
}

```

19.3.4 Forming the time interval

Once h_j is found, we need to determine the next representable integration point and the machine interval over which the a priori bounds hold.

1. If $h_j > 0$, we compute

$$\begin{aligned}
t_{j+1} &= \downarrow (\mathbf{t}_j + h_j) \quad \text{and} \\
\mathbf{T}_j &= [\mathbf{t}_j, t_{j+1}].
\end{aligned}$$

2. If $h_j < 0$, we compute

$$t_{j+1} = \uparrow (\bar{\mathbf{t}}_j + h_j) \quad \text{and} \\ \mathbf{T}_j = [t_{j+1}, \bar{\mathbf{t}}_j].$$

Proposition 19.1.

(i) $\mathbf{T}_j \subseteq \mathbf{t}_j + [0, 1]h_{j,0}$

(ii) If $|h_j| \geq w(\mathbf{t}_j)$, then $\mathbf{t}_j \subseteq \mathbf{T}_j$

Proof. $h_{j,0} > 0$. Then $0 < h_j \leq h_{j,0}$ and

$$\begin{aligned} \mathbf{T}_j &= [\underline{\mathbf{t}}_j, t_{j+1}] \subseteq [\underline{\mathbf{t}}_j, \underline{\mathbf{t}}_j + h_j] \subseteq [\underline{\mathbf{t}}_j, \underline{\mathbf{t}}_j + h_{j,0}] = \underline{\mathbf{t}}_j + [0, 1]h_{j,0} \\ &\subseteq \mathbf{t}_j + [0, 1]h_{j,0}. \end{aligned}$$

If $h_j \geq w(\mathbf{t}_j) = \bar{\mathbf{t}}_j - \underline{\mathbf{t}}_j$, then $\bar{\mathbf{t}}_j \leq \downarrow (\underline{\mathbf{t}}_j + h_j) \leq \underline{\mathbf{t}}_j + h_j$, and

$$[\underline{\mathbf{t}}_j, \bar{\mathbf{t}}_j] \subseteq [\underline{\mathbf{t}}_j, \downarrow (\underline{\mathbf{t}}_j + h_j)] = \mathbf{T}_j.$$

$h_{j,0} < 0$. Then $0 > h_j \geq h_{j,0}$ and

$$\begin{aligned} \mathbf{T}_j &= [t_{j+1}, \bar{\mathbf{t}}_j] \subseteq [\bar{\mathbf{t}}_j + h_j, \bar{\mathbf{t}}_j] \subseteq [\bar{\mathbf{t}}_j + h_{j,0}, \bar{\mathbf{t}}_j] = \bar{\mathbf{t}}_j + [0, 1]h_{j,0} \\ &\subseteq \mathbf{t}_j + [0, 1]h_{j,0}. \end{aligned}$$

If $-h_j \geq \bar{\mathbf{t}}_j - \underline{\mathbf{t}}_j$, then $\underline{\mathbf{t}}_j \geq \uparrow (\bar{\mathbf{t}}_j + h_j) \geq \bar{\mathbf{t}}_j + h_j$, and

$$\mathbf{t}_j = [\underline{\mathbf{t}}_j, \bar{\mathbf{t}}_j] \subseteq [\uparrow (\bar{\mathbf{t}}_j + h_j), \bar{\mathbf{t}}_j] = \mathbf{T}_j.$$

□

215 $\langle \text{form time interval } 215 \rangle \equiv \quad / *$

$$h = h_j$$

$$t0 = \mathbf{t}_j$$

$$\begin{aligned} t_trial &= t_{j+1} \\ \text{apriori_trial-}t &= \mathbf{T}_j \end{aligned}$$

$*/$

double td ;

if $(h > 0)$ {

$td = \inf(t0)$;

$\text{round_down}()$;

$t_trial = td + h$;


```

    apriori_trial_t = interval(td, t_trial);
}
else {
    td = sup(t0);
    round_up();
    t_trial = td + h;
    apriori_trial_t = interval(t_trial, td);
}
assert(subseteq(t0, apriori_trial_t));

```

This code is used in chunk 217.

19.3.5 Selecting a trial stepsize for the next step

We can select for the next step $h_{j+1,0} = \beta h_{j,0}$. This is reasonable since, if $|h_{j+1,0}| > |h_{j,0}|$, we have computed an a priori enclosure with $h_{j,0}$, but we could have possibly done this with a stepsize between $h_{j,0}$ and $h_{j+1,0}$. That is, we assume that we might be successful on the next step with $|h_{j+1,0}| = \beta |h_{j,0}| > |h_{j,0}|$. Similar considerations apply when $|h_{j+1,0}| < |h_{j,0}|$. Hence, we select

$$h_{j+1,0} = \beta h_{j,0}$$

for the next step.

216 $\langle \text{select stepsize 216} \rangle \equiv$ /*

$$\frac{\begin{array}{l} h_trial = h_{j,0} \\ beta \supseteq \beta \end{array}}{h_next = h_{j+1,0}}$$

```

    */
    h_next = inf(beta) * h_trial;

```

This code is used in chunk 217.

19.3.6 Computing a priori bounds

First, we check if the stepsize is not very small. Then, we compute \mathbf{p}_j , $\tilde{\mathbf{y}}_j$, and h_j . If $|h_j| \leq h_{\min}$, we cannot validate existence and uniqueness and return *false*. Otherwise, we form \mathbf{T}_j and select $h_{j+1,0}$ for the next step.

217 $\langle \text{validate existence and uniqueness 217} \rangle \equiv$

```

    void HOE::compAprioriEnclosure(const interval &t0, const iVector &y0,
                                   bool &info)

```

```

{      /*

           $t_0 = \mathbf{t}_j$ 
           $y_0 \supseteq \mathbf{y}_j$ 
           $h\_trial = h_{j,0}$ 
           $order\_trial = k$ 

          -----
           $apriori\_trial \rightarrow t \supseteq \mathbf{T}_j$ 
           $apriori\_trial \rightarrow y \supseteq \tilde{\mathbf{y}}_j$ 
           $tayl\_coeff$  contains enclosures on  $h_{j,0}^i f^{[i]}(\mathbf{t}_j + [0, 1]h_{j,0}, \tilde{\mathbf{y}}_j)$  for  $i = 0, \dots, k$ 
           $h\_trial = h_{j+1,0}$ 

      */
  if (fabs(h_trial) ≤ control-hmin) {
    info = false;
    return;
  }
  { compute  $\mathbf{p}_j$  206 };
  { compute  $\mathbf{u}_j$  and  $\tilde{\mathbf{y}}_j$  207 };
  { compute stepsize 210 };
  if (fabs(h) ≤ control-hmin) {
    info = false;
    return;
  }
  { form time interval 215 };
  { select stepsize 216 };
  info = true;
}

```

This code is used in chunk 227.

19.4 Other functions

19.4.1 Constructor and destructor

```

219 { constructor-destructor HOE 219 } ≡
    HOE::HOE(int n)
    : one(interval(0,1)) {
      sizeV(term, n);
      sizeV(p, n);
      sizeV(u, n);
      sizeV(v, n);
      apriori_trial = new Apriori(n);
      apriori = new Apriori(n);
      assert(apriori ∧ apriori_trial);
      tayl_coeff = 0;
    }

```

```

    control = 0;
}
HOE::~HOE()
{
    delete apriori;
    delete apriori_trial;
}

```

This code is used in chunk 227.

19.4.2 Accept a solution

To accept a trial enclosure, we store it in the **apriori* object.

```

220 <accept solution (HOE) 220> ≡
    void HOE::acceptSolution() { /*

        apriori_trial contains  $\tilde{\mathbf{y}}_j$  and  $\mathbf{T}_j$ 
        apriori contains  $\tilde{\mathbf{y}}_{j-1}$  and  $\mathbf{T}_{j-1}$ 
        h_trial =  $h_{j,0}$ 
        h_next =  $h_{j+1,0}$ 

        -----
        apriori contains  $\tilde{\mathbf{y}}_j$  and  $\mathbf{T}_j$ 
        h_trial =  $h_{j+1,0}$ 

        */
        apriori->t = apriori_trial->t;
        assignV(apriori->y, apriori_trial->y);
        h_trial = h_next;
    }

```

This code is used in chunk 227.

19.4.3 Set functions

```

221 <set functions HOE 221> ≡
    void set(Control *ctrl, AD *ad)
    {
        control = ctrl;
        tayl_coeff = ad->tayl_coeff_ode;
    }

    void init(const interval &t0, const iVector &y0) {
        apriori->init(t0, y0);
    }

    void setTrialStepsize(double h0) {
        h_trial = h0;
    }

```

```

void setTrialOrder(int order0) {
    order_trial = order0;
}

```

This code is used in chunk 204.

19.4.4 Get functions

```

222 <get functions HOE 222> ≡
    double getStepsize() const {
        return h;
    }
    double getTrialStepsize() const {
        return h_trial;
    }
    const interval &getT() const {
        return apriori-t;
    }
    interval getTrialT() const {
        return apriori_trial-t;
    }
    const iVector &getApriori() const {
        return apriori-y;
    }
    const iVector &getTrialApriori() const {
        return apriori_trial-y;
    }

```

See also chunk 223.

This code is used in chunk 204.

Obtain error term

We obtain $h_{j,0}^i f^{[i]}(\mathbf{t}_j + [0, 1]h_{j,0}, \tilde{\mathbf{y}}_j)$ for given $i = 0, 1, \dots, k$ by

```

223 <get functions HOE 222> +≡      /*
    
$$e \supseteq h_{j,0}^i f^{[i]}(\mathbf{t}_j + [0, 1]h_{j,0}, \tilde{\mathbf{y}}_j)$$

    */
    void getErrorTerm(iVector &e, int i) const {
        tayl-coeff→getTerm(e, i);
    }

```

19.4.5 Enclosing β

The function *comp_beta* returns an enclosure of β .

```

224  ⟨compute  $\beta$  224⟩ ≡      /*
       $gamma = \gamma$ , the largest representable  $\gamma > 0$  such that  $\gamma \mathbf{v}_j \subseteq \mathbf{u}_j$ 
       $beta \supseteq \beta = \gamma^{1/k}$ 

      */
      interval HOE::comp_beta(const iVector &v, const iVector &u, int k)
      {
        double gamma = v_blas::compH(v, u);
        interval i_gamma = gamma;
        interval i_pw = interval(1.0)/interval(double(k));
        interval beta = pow(i_gamma, i_pw);
        return beta;
      }

```

This code is used in chunk 227.

Files

```

226  ⟨hoe.h 226⟩ ≡
      #ifndef HOE_H
      #define HOE_H
      namespace vnodelp {
        ⟨class HOE 204⟩
      }
      #endif

227  ⟨hoe.cc 227⟩ ≡
      #include <cmath>
      #include <cassert>
      #include <algorithm>
      #include "vnodeinterval.h"
      #include "vnoderound.h"
      #include "basiclinalg.h"
      #include "intvfuncs.h"
      #include "control.h"
      #include "solution.h"
      #include "allad.h"
      #include "hoe.h"
      using namespace v_blas;
      namespace vnodelp {
        ⟨constructor-destructor HOE 219⟩
        ⟨validate existence and uniqueness 217⟩
        ⟨accept solution (HOE) 220⟩
        ⟨compute  $\beta$  224⟩
      }

```


Chapter 20

Computing tight bounds on the solution

To compute tight bounds on the solution, we employ the interval Hermite-Obreschkoff (IHO) method developed in [23]. It consists of two phases: a predictor and a corrector. We present the relevant theory first and then describe its implementation.

20.1 Theory background

20.1.1 Predictor

For $y_j \in \mathbf{y}_j \subseteq \tilde{\mathbf{y}}_j$, we have

$$y(t_{j+1}; t_j, y_j) \in y_j + \sum_{i=1}^q h_j^i f^{[i]}(t_j, y_j) + h_j^{q+1} f^{[q+1]}(\mathbf{T}_j, \tilde{\mathbf{y}}_j), \quad (20.1)$$

where $h_j = t_{j+1} - t_j$ and $t_j, t_{j+1} \in \mathbf{T}_j$.

Let $J(f^{[i]}; \mathbf{y}_j)$ be the Jacobian of $f^{[i]}$ evaluated at \mathbf{y}_j and denote

$$\mathbf{U}_{j+1} = I + \sum_{i=1}^q h_j^i J(f^{[i]}; \mathbf{y}_j).$$

These Jacobians are computed by generating TCs for the solution of the associated variational equation

$$Y' = \frac{\partial f}{\partial y} Y, \quad Y(t_j) = I,$$

where I is the $n \times n$ identity matrix.

We assume that at t_j the solution $y(t_j; t_0, y_0)$ is contained in

$$\begin{aligned} & \mathbf{y}_j \\ & \{u_j + S_j \alpha + A_j r \mid \alpha \in \boldsymbol{\alpha}, r \in \mathbf{r}_j\}, \quad \text{and} \\ & \{u_j + S_j \alpha + Q_j r \mid \alpha \in \boldsymbol{\alpha}, r \in \mathbf{r}_{\text{QR},j}\}, \end{aligned}$$

where $u_j \in \mathbf{y}_j$.

If we apply the mean-value theorem to the $f^{[i]}$ in (20.1), we obtain that for any

$$\begin{aligned} y_j &\in \{u_j + S_j \alpha + A_j r \mid \alpha \in \boldsymbol{\alpha}, r \in \mathbf{r}_j\} \\ &\cap \{u_j + S_j \alpha + Q_j r \mid \alpha \in \boldsymbol{\alpha}, r \in \mathbf{r}_{\text{QR},j}\} \\ &\cap \mathbf{y}_j, \end{aligned}$$

$$\begin{aligned} y(t_{j+1}; t_j, y_j) &\in u_j + \sum_{i=1}^q h_j^i f^{[i]}(t_j, u_j) + h_j^{q+1} f^{[q+1]}(\mathbf{T}_j, \tilde{\mathbf{y}}_j) \\ &\quad + (U_{j+1} S_j) \boldsymbol{\alpha} + \{(U_{j+1} A_j) \mathbf{r}_j \cap (U_{j+1} Q_j) \mathbf{r}_{\text{QR},j}\}. \end{aligned} \quad (20.2)$$

For brevity, denote

$$\begin{aligned} \hat{u}_{j+1} &= u_j + \sum_{i=1}^q h_j^i f^{[i]}(t_j, u_j) \\ \mathbf{z}_{j+1} &= h_j^{q+1} f^{[q+1]}(\mathbf{T}_j, \tilde{\mathbf{y}}_j), \quad \text{and} \\ \mathbf{x}_{j+1} &= (U_{j+1} S_j) \boldsymbol{\alpha} + \{(U_{j+1} A_j) \mathbf{r}_j \cap (U_{j+1} Q_j) \mathbf{r}_{\text{QR},j}\}. \end{aligned}$$

Then

$$y(t_{j+1}; t_0, y_0) \in \mathbf{y}_{j+1}^* := (\hat{u}_{j+1} + \mathbf{z}_{j+1} + \mathbf{x}_{j+1}) \cap \tilde{\mathbf{y}}_j.$$

20.1.2 Corrector

In the corrector, we compute $\mathbf{y}_{j+1} \subseteq \mathbf{y}_{j+1}^*$. Usually, \mathbf{y}_{j+1} is much tighter than \mathbf{y}_{j+1}^* .

Let

$$y_j = y(t_j; t_0, y_0) \quad \text{and} \quad y_{j+1} = y(t_{j+1}; t_0, y_0).$$

Denote $k = p + q + 1$ ($p, q \geq 0$) and

$$c_i^{q,p} = \frac{q! (q + p - i)!}{(p + q)! (q - i)!} \quad (q, p, \text{ and } i \geq 0). \quad (20.3)$$

Denote also

$$\gamma_{p,q} = \frac{q! p!}{(p + q)!} \quad (20.4)$$

Then [23]

$$\begin{aligned} \sum_{i=0}^q (-1)^i c_i^{q,p} h_j^i f^{[i]}(t_{j+1}, y_{j+1}) &\in \sum_{i=0}^p c_i^{p,q} h_j^i f^{[i]}(t_j, y_j) \\ &\quad + (-1)^q \gamma_{p,q} h_j^{p+q+1} f^{[p+q+1]}(\mathbf{T}_j, \tilde{\mathbf{y}}_j). \end{aligned} \quad (20.5)$$

Denote

$$y_{j+1}^* = m(\mathbf{y}_{j+1}^*), \quad (20.6)$$

$$\mathbf{B}_{j+1} = \sum_{i=0}^q (-1)^i c_i^{q,p} h_j^i J(f^{[i]}; \mathbf{y}_{j+1}^*), \quad (20.7)$$

$$\mathbf{F}_j = \sum_{i=0}^p c_i^{p,q} h_j^i J(f^{[i]}; \mathbf{y}_j), \quad (20.8)$$

$$\mathbf{C}_{j+1} = m(\mathbf{B}_{j+1}), \quad (20.9)$$

$$\mathbf{S}_{j+1} = (\mathbf{C}_{j+1}^{-1} \mathbf{F}_j) \mathbf{S}_j, \quad (20.10)$$

$$\mathbf{A}_{j+1} = (\mathbf{C}_{j+1}^{-1} \mathbf{F}_j) \mathbf{A}_j, \quad (20.11)$$

$$\mathbf{Q}_{j+1} = (\mathbf{C}_{j+1}^{-1} \mathbf{F}_j) \mathbf{Q}_j, \quad (20.12)$$

$$\mathbf{e}_{j+1} = (-1)^q \gamma_{p,q} h_j^{p+q+1} f^{[p+q+1]}(\mathbf{T}_j, \tilde{\mathbf{y}}_j), \quad (20.13)$$

$$g_{j+1} = \sum_{i=0}^p c_i^{p,q} h_j^i f^{[i]}(u_j) - \sum_{i=0}^q (-1)^i c_i^{q,p} h_j^i f^{[i]}(y_{j+1}^*), \quad (20.14)$$

$$\mathbf{d}_{j+1} = g_{j+1} + \mathbf{e}_{j+1}, \quad \text{and} \quad (20.15)$$

$$\mathbf{w}_{j+1} = \mathbf{C}_{j+1}^{-1} \mathbf{d}_{j+1} + (\mathbf{I} - \mathbf{C}_{j+1}^{-1} \mathbf{B}_{j+1})(\mathbf{y}_{j+1}^* - \mathbf{y}_{j+1}^*). \quad (20.16)$$

(For sufficiently small h_j , we can enclose the inverse of \mathbf{C}_{j+1} .)

Since

$$y_{j+1}, y_{j+1}^* \in \mathbf{y}_{j+1}^* \quad \text{and} \quad y_j, u_j \in \mathbf{y}_j,$$

we can apply the mean-value theorem to the two sums in (20.5), and using the above notation derive that

$$y_{j+1} \in y_{j+1}^* + \mathbf{S}_{j+1} \boldsymbol{\alpha} + \mathbf{A}_{j+1} \mathbf{r}_j + \mathbf{w}_{j+1} \quad \text{and} \quad (20.17)$$

$$y_{j+1} \in y_{j+1}^* + \mathbf{S}_{j+1} \boldsymbol{\alpha} + \mathbf{Q}_{j+1} \mathbf{r}_{\text{QR},j} + \mathbf{w}_{j+1}. \quad (20.18)$$

Denoting

$$\mathbf{s}_{j+1} = (\mathbf{A}_{j+1} \mathbf{r}_j) \cap (\mathbf{Q}_{j+1} \mathbf{r}_{\text{QR},j}),$$

we have

$$y(t_{j+1}; t_0, y_0) \in \mathbf{y}_{j+1} := (y_{j+1}^* + \mathbf{S}_{j+1} \boldsymbol{\alpha} + \mathbf{s}_{j+1} + \mathbf{w}_{j+1}) \cap \mathbf{y}_{j+1}^*.$$

20.1.3 Computing a solution representation

We have to determine u_{j+1} , \mathbf{S}_{j+1} , \mathbf{A}_{j+1} , \mathbf{Q}_{j+1} , \mathbf{r}_{j+1} , and $\mathbf{r}_{\text{QR},j+1}$ for the next step.

Let

$$u_{j+1} = m(\mathbf{y}_{j+1}), \quad (20.19)$$

$$S_{j+1} = m(\mathbf{S}_{j+1}), \quad (20.20)$$

$$\mathbf{v}_{j+1} = \mathbf{y}_{j+1}^* - u_{j+1} + (\mathbf{S}_{j+1} - S_{j+1})\boldsymbol{\alpha} + \mathbf{w}_{j+1}. \quad (20.21)$$

$$\mathbf{A}_{j+1} = m(\mathbf{A}_{j+1}), \quad (20.22)$$

$$\mathbf{r}_{j+1} = (\mathbf{A}_{j+1}^{-1} \mathbf{A}_{j+1})\mathbf{r}_j + \mathbf{A}_{j+1}^{-1} \mathbf{v}_{j+1}, \quad (20.23)$$

$$\mathbf{Q}_{j+1} \quad \text{the orthonormal matrix described in Subsection 20.1.4, and} \quad (20.24)$$

$$\mathbf{r}_{\text{QR},j+1} = (\mathbf{Q}_{j+1}^{-1} \mathbf{Q}_{j+1})\mathbf{r}_{\text{QR},j} + \mathbf{Q}_{j+1}^{-1} \mathbf{v}_{j+1}. \quad (20.25)$$

Using (20.19–20.25) in (20.17–20.18), we derive that

$$y_{j+1} \in \{u_{j+1} + S_{j+1}\alpha + A_{j+1}r \mid \alpha \in \boldsymbol{\alpha}, r \in \mathbf{r}_{j+1}\} \quad \text{and}$$

$$y_{j+1} \in \{u_{j+1} + S_{j+1}\alpha + Q_{j+1}r \mid \alpha \in \boldsymbol{\alpha}, r \in \mathbf{r}_{\text{QR},j+1}\}.$$

We assumed above that we can enclose A_{j+1}^{-1} . If we cannot enclose A_{j+1}^{-1} , or if

$$Q_{j+1}\mathbf{r}_{\text{QR},j+1} \subseteq A_{j+1}\mathbf{r}_{j+1},$$

we set

$$A_{j+1} = Q_{j+1} \quad \text{and} \quad \mathbf{r}_{j+1} = \mathbf{r}_{\text{QR},j+1}.$$

20.1.4 Computing Q_{j+1}

Let

$$\tilde{A}_{j+1} = m(\mathbf{Q}_{j+1}) \quad \text{and} \quad D = \text{diag}(w(\mathbf{r}_{\text{QR},j})).$$

Let also P_{j+1} be a permutation matrix such that the columns of $\tilde{A}_{j+1}DP_{j+1}$ are sorted in non-increasing order in the Euclidean norm. Then, we perform the QR factorization of $\tilde{A}_{j+1}DP_{j+1} = Q_{j+1}R_{j+1}$ and use Q_{j+1} in our method.

20.2 Implementation

20.2.1 The IHO class

```

235 <class IHO 235> ≡
    class IHO {
    public:
        IHO(int n);
        <set and get functions 277>
        void compCoeffs();
        void compTightEnclosure(interval &t_next);
        void acceptSolution();
        virtual ~IHO();

```

```

private:
    void compCpq(int p, int q);
    void compCqp(int p, int q);
    interval compErrorConstant(int p, int q);
private:
    unsigned int p, q, order_trial;
    interval h_trial;
    pVector y_pred_point;
    iVector y, y_pred, globalExcess, temp, temp2, x, u_next, predictor_excess,
        corrector_excess, z, w, gj, term, d, s;
    iMatrix Fj, M, Cinv, G, B, S, A, Q, U, V, Ainv;
    pMatrix C, A_point;
    MatrixInverse *matrix_inverse;
    Solution *solution, *trial_solution;
    interval *C_pq, *C_qp;
    HOE *hoe;
    AD *ad;
    Control *control;
    interval errorConstant;
};

```

This code is used in chunk 292.

20.2.2 Computing a tight enclosure

The main functions is

```

236 <compute tight enclosure 236> ≡
    void IHO::compTightEnclosure(interval &t_next)
    {
        <initialize IHO method 238>
        <predictor: compute  $\mathbf{y}_{j+1}^*$  242>
        <corrector: compute  $\mathbf{y}_{j+1}$  247>
        <set  $\mathbf{t}_{j+1}$  264>
        <find solution representation for next step 265>
    }

```

This code is used in chunk 293.

20.2.3 Initialization

Stepsize

We assume we have enclosures on the solution at $t_j \in \mathbf{t}_j$ and now wish to compute enclosures at $t_{j+1} \in \mathbf{t}_{j+1}$. Hence, we have for the stepsize

$$h_j = t_{j+1} - t_j \in \mathbf{h}_j = \mathbf{t}_{j+1} - \mathbf{t}_j.$$

238 $\langle \text{initialize IHO method 238} \rangle \equiv$ $/*$

$$\frac{\begin{array}{l} t_next = t_{j+1} \\ solution_t = t_j \end{array}}{h_trial \supseteq h_j = t_{j+1} - t_j}$$

$*/$
 $h_trial = t_next - solution_t;$

See also chunks 240 and 241.

This code is used in chunk 236.

Order and method coefficients

The order is in *order_trial*. Initially, *order_trial* = 0. If *control_order* \neq *order_trial*, we set *order_trial* = *control_order* and compute the coefficients (20.3) and (20.4) of the method. If the value for the order is *k*, we require that $p + q + 1 = k$ and $p \leq q$. We set

$$p = \lfloor (k - 1)/2 \rfloor \quad \text{and} \quad q = \lceil (k - 1)/2 \rceil.$$

239 $\langle \text{compute IHO method coefficients 239} \rangle \equiv$

```
void IHO::compCoeffs()
{
    if (order_trial  $\neq$  control_order) {
        /* deal with order */
        order_trial = control_order;
        double pq = (order_trial - 1)/2.0;
        p = int(floor((pq)));
        q = int(ceil((pq)));
        assert(p + q + 1  $\equiv$  order_trial);
        /* reallocate memory if necessary */
        if (C_pq) delete[] C_pq;
        C_pq = new interval[p + 1];
        if (C_qp) delete[] C_qp;
        C_qp = new interval[q + 1];
        /* compute coefficients */
        compCpq(p, q);
        compCqp(p, q);
        errorConstant = compErrorConstant(p, q);
    }
}
```

This code is used in chunk 293.

240 $\langle \text{initialize IHO method 238} \rangle + \equiv$

$compCoeffs();$

Local excess

In the HOE method, we enclose

$$h_{j,0}^{q+1} f^{[q+1]}(\mathbf{t}_j + [0, 1]h_{j,0}, \tilde{\mathbf{y}}_j)$$

and

$$h_{j,0}^{p+q+1} f^{[p+q+1]}(\mathbf{t}_j + [0, 1]h_{j,0}, \tilde{\mathbf{y}}_j) = h_{j,0}^k f^{[k]}(\mathbf{t}_j + [0, 1]h_j, \tilde{\mathbf{y}}_j).$$

In the IHO method, we have to enclose

$$h_j^{q+1} f^{[q+1]}(\mathbf{T}_j, \tilde{\mathbf{y}}_j) \quad \text{and} \quad h_j^{p+q+1} f^{[p+q+1]}(\mathbf{T}_j, \tilde{\mathbf{y}}_j),$$

where $\mathbf{T}_j \subseteq \mathbf{t}_j + [0, 1]h_{j,0}$, cf. Proposition 19.1.

We have

$$\begin{aligned} h_j^{q+1} f^{[q+1]}(\mathbf{T}_j, \tilde{\mathbf{y}}_j) &= \left(\frac{h_j}{h_{j,0}} \right)^{q+1} h_{j,0}^{q+1} f^{[q+1]}(\mathbf{T}_j, \tilde{\mathbf{y}}_j) \\ &\subseteq \left(\frac{h_j}{h_{j,0}} \right)^{q+1} h_{j,0}^{q+1} f^{[q+1]}(\mathbf{t}_j + [0, h_{j,0}], \tilde{\mathbf{y}}_j). \end{aligned}$$

Similarly, we have

$$\begin{aligned} h_j^k f^{[k]}(\mathbf{T}_j, \tilde{\mathbf{y}}_j) &= \left(\frac{h_j}{h_{j,0}} \right)^k h_{j,0}^k f^{[k]}(\mathbf{T}_j, \tilde{\mathbf{y}}_j) \\ &\subseteq \left(\frac{h_j}{h_{j,0}} \right)^k h_{j,0}^k f^{[k]}(\mathbf{t}_j + [0, h_{j,0}], \tilde{\mathbf{y}}_j). \end{aligned}$$

Hence, we just re-scale.

```

241 < initialize IHO method 238 > +≡      /*
                                     h_trial ⊃ h_j
      tayl_coeff → getStepsize() returns h_{j,0}
      tayl_coeff contains enclosures on h_{j,0}^i f^{[i]}(t_j + [0, 1]h_{j,0}, y_tilde_j)
                                     for i = 0, ..., k = p + q + 1
      -----
      rescale ⊃ h_j / h_{j,0}
      predictor_excess ⊇ h_{j,0}^{q+1} f^{[q+1]}(t_j + [0, 1]h_{j,0}, y_tilde_j)
      predictor_excess ⊇ h_j^{q+1} f^{[q+1]}(t_j + [0, 1]h_{j,0}, y_tilde_j)
      corrector_excess ⊇ h_{j,0}^{p+q+1} f^{[p+q+1]}(t_j + [0, 1]h_{j,0}, y_tilde_j)
      corrector_excess ⊇ h_j^{p+q+1} f^{[p+q+1]}(t_j + [0, 1]h_{j,0}, y_tilde_j)

      */
interval rescale = h_trial / ad → tayl_coeff_ode → getStepsize();
hoe → getErrorTerm(predictor_excess, q + 1);
scaleV(predictor_excess, v_bias :: pow(rescale, q + 1));
hoe → getErrorTerm(corrector_excess, p + q + 1);
scaleV(corrector_excess, pow(rescale, order_trial));

```

20.2.4 Predictor

We enclose \hat{u}_{j+1} , \mathbf{U}_{j+1} , \mathbf{x}_{j+1} , and \mathbf{y}_{j+1}^* . Below, $t_j \in \mathbf{t}_j$.

242 $\langle \text{predictor: compute } \mathbf{y}_{j+1}^* \text{ 242} \rangle \equiv$
 $\langle \hat{u}_{j+1} = u_j + \sum_{i=1}^q h_j^i f^{[i]}(t_j, u_j) \text{ 243} \rangle$
 $\langle \mathbf{U}_{j+1} = I + \sum_{i=1}^q h_j^i J(f^{[i]}; \mathbf{y}_j) \text{ 244} \rangle$
 $\langle \mathbf{x}_{j+1} = (\mathbf{U}_{j+1} S_j) \boldsymbol{\alpha} + \{(\mathbf{U}_{j+1} A_j) \mathbf{r}_j \cap (\mathbf{U}_{j+1} Q_j) \mathbf{r}_{\text{QR},j}\} \text{ 245} \rangle$
 $\langle \mathbf{y}_{j+1}^* = (\hat{u}_{j+1} + \mathbf{z}_{j+1} + \mathbf{x}_{j+1}) \cap \tilde{\mathbf{y}}_j \text{ 246} \rangle$

This code is used in chunk 236.

243 $\langle \hat{u}_{j+1} = u_j + \sum_{i=1}^q h_j^i f^{[i]}(t_j, u_j) \text{ 243} \rangle \equiv$ $/*$
 $\text{solution} \rightarrow t = \mathbf{t}_j$
 $\text{solution} \rightarrow u = u_j$
 $h_trial \supseteq \mathbf{h}_j$
 $q \text{ order}$

 tayl_coeff contains enclosures on $h_j^i f^{[i]}(t_j, u_j)$ for $i = 0, \dots, q$
 $u_next \ni \hat{u}_{j+1}$
 $*/$
 $\text{assign} V(\text{temp}, \text{solution} \rightarrow u);$ $/* \text{temp stores } u_j \text{ as an interval } */$
 $\text{ad} \rightarrow \text{tayl_coeff_ode} \rightarrow \text{set}(\text{solution} \rightarrow t, \text{temp}, h_trial, q);$
 $\text{ad} \rightarrow \text{tayl_coeff_ode} \rightarrow \text{compTerms}();$
 $\text{ad} \rightarrow \text{tayl_coeff_ode} \rightarrow \text{sumTerms}(u_next, q);$

This code is used in chunk 242.

244 $\langle \mathbf{U}_{j+1} = I + \sum_{i=1}^q h_j^i J(f^{[i]}; \mathbf{y}_j) \text{ 244} \rangle \equiv$ $/*$
 $\text{solution} \rightarrow t = \mathbf{t}_j$
 $\text{solution} \rightarrow y \supseteq \mathbf{y}_j$
 $h_trial \supseteq \mathbf{h}_j$
 $q \text{ order}$

 tayl_coeff_var contains enclosures on $h_j^i J(f^{[i]}; \mathbf{y}_j)$
for $i = 0, 1, \dots, q$

$$U \supseteq \mathbf{U}_{j+1} = I + \sum_{i=1}^q h_j^i J(f^{[i]}; \mathbf{y}_j)$$
 $*/$
 $\text{ad} \rightarrow \text{tayl_coeff_var} \rightarrow \text{set}(\text{solution} \rightarrow t, \text{solution} \rightarrow y, h_trial, q);$
 $\text{ad} \rightarrow \text{tayl_coeff_var} \rightarrow \text{compTerms}();$
 $\text{ad} \rightarrow \text{tayl_coeff_var} \rightarrow \text{sumTerms}(U, q);$

This code is used in chunk 242.

Finally, we compute

```

246  $\langle \mathbf{y}_{j+1}^* = (\hat{u}_{j+1} + \mathbf{z}_{j+1} + \mathbf{x}_{j+1}) \cap \tilde{\mathbf{y}}_j \text{ 246} \rangle \equiv \quad / *$ 

$$u\_next \ni \hat{u}_{j+1}$$


$$predictor\_excess \supseteq \mathbf{z}_{j+1}$$


$$x \supseteq \mathbf{x}_{j+1}$$


$$hoe \rightarrow getTrialApriori() \supseteq \tilde{\mathbf{y}}_j$$



---



$$y\_pred \supseteq u_{j+1} + \mathbf{z}_{j+1}$$


$$y\_pred \supseteq \hat{u}_{j+1} + \mathbf{z}_{j+1} + \mathbf{x}_{j+1}$$


$$y\_pred \supseteq (\hat{u}_{j+1} + \mathbf{z}_{j+1} + \mathbf{x}_{j+1}) \cap \tilde{\mathbf{y}}_j$$


$$*/$$


$$addViVi(y\_pred, u\_next, predictor\_excess);$$


$$addViVi(y\_pred, x);$$


$$\text{bool } b = intersect(y\_pred, hoe \rightarrow getTrialApriori(), y\_pred);$$


$$assert(b);$$


```

This code is used in chunk 242.

20.2.5 Corrector

The computation of \mathbf{y}_{j+1} is given by

```

247  $\langle \text{corrector: compute } \mathbf{y}_{j+1} \text{ 247} \rangle \equiv$ 

$$\langle \mathbf{F}_j = \sum_{i=0}^p c_i^{p,q} h_j^i J(f^{[i]}; \mathbf{y}_j) \text{ 248} \rangle$$


$$\langle \mathbf{B}_{j+1} = \sum_{i=0}^q (-1)^i c_i^{q,p} h_j^i J(f^{[i]}; \mathbf{y}_{j+1}^*) \text{ 249} \rangle$$


$$\langle C_{j+1} = m(\mathbf{B}_{j+1}) \text{ 250} \rangle$$


$$\langle \mathbf{G}_{j+1} = C_{j+1}^{-1} \mathbf{F}_j \text{ 251} \rangle$$


$$\langle \mathbf{S}_{j+1} = \mathbf{G}_{j+1} \mathbf{S}_j \text{ 252} \rangle$$


$$\langle \mathbf{A}_{j+1} = \mathbf{G}_{j+1} \mathbf{A}_j \text{ 253} \rangle$$


$$\langle \mathbf{Q}_{j+1} = \mathbf{G}_{j+1} \mathbf{Q}_j \text{ 254} \rangle$$


$$\langle \mathbf{e}_{j+1} = (-1)^q \gamma_{p,q} \mathbf{h}_j^{p+q+1} f^{[p+q+1]}(\mathbf{T}_j, \tilde{\mathbf{y}}_j) \text{ 255} \rangle$$


$$\langle g_{j+1} = \sum_{i=0}^p c_i^{p,q} h_j^i f^{[i]}(u_j) - \sum_{i=0}^q (-1)^i c_i^{q,p} h_j^i f^{[i]}(y_{j+1}^*) \text{ 256} \rangle$$


$$\langle \mathbf{d}_{j+1} = g_{j+1} + \mathbf{e}_{j+1} \text{ 260} \rangle$$


$$\langle \mathbf{w}_{j+1} = C_{j+1}^{-1} \mathbf{d}_{j+1} + (I - C_{j+1}^{-1} \mathbf{B}_{j+1})(\mathbf{y}_{j+1}^* - \mathbf{y}_{j+1}^*) \text{ 261} \rangle$$


$$\langle \mathbf{s}_{j+1} = (\mathbf{A}_{j+1} \mathbf{r}_j) \cap (\mathbf{Q}_{j+1} \mathbf{r}_{QR,j}) \text{ 262} \rangle$$


$$\langle \mathbf{y}_{j+1} = (\mathbf{y}_{j+1}^* + \mathbf{S}_{j+1} \boldsymbol{\alpha} + \mathbf{s}_{j+1} + \mathbf{w}_{j+1}) \cap \mathbf{y}_{j+1}^* \text{ 263} \rangle$$


```

This code is used in chunk 236.

Computing \mathbf{F}_j

The matrices $h_j^i J(f^{[i]}; \mathbf{y}_j)$ for $i = 1, \dots, q$ have been already enclosed in the predictor. Since $p \leq q$, we have the terms that we need to enclose \mathbf{F}_j in (20.8). In the code below, we obtain the enclosure on $h_j^i J(f^{[i]}; \mathbf{y}_j)$ in M and then scale M such that it contains $c_i^{p,q} h_j^i J(f^{[i]}; \mathbf{y}_j)$.

248 $\langle \mathbf{F}_j = \sum_{i=0}^p c_i^{p,q} h_j^i J(f^{[i]}; \mathbf{y}_j) \text{ 248} \rangle \equiv \quad /*$

taylor_coeff_var contains enclosures on $h_j^i J(f^{[i]}; \mathbf{y}_j)$ for $i = 0, 1, \dots, q$

p order

$$M \supseteq h_j^i J(f^{[i]}; \mathbf{y}_j)$$

$$C_pq[i] \ni c_i^{p,q}$$

$$M \supseteq c_i^{p,q} h_j^i J(f^{[i]}; \mathbf{y}_j)$$

$$Fj \supseteq \mathbf{F}_j$$

$*/$

$\text{setM}(Fj, 0.0);$

for (**int** $i = p; i \geq 1; i--$) {

$\text{ad} \rightarrow \text{taylor_coeff_var} \rightarrow \text{getTerm}(M, i);$

$\text{scaleM}(M, C_pq[i]);$

$\text{addMiMi}(Fj, M);$

}

$\text{addId}(Fj);$

This code is used in chunk 247.

Computing \mathbf{B}_{j+1}

We enclose $h_j^i J(f^{[i]}; \mathbf{y}_{j+1}^*)$ for $i = 1, \dots, q$. We obtain these enclosures in M , which is scaled such that it encloses $(-1)^i c_i^{q,p} h_j^i J(f^{[i]}; \mathbf{y}_{j+1}^*)$.

249 $\langle \mathbf{B}_{j+1} = \sum_{i=0}^q (-1)^i c_i^{q,p} h_j^i J(f^{[i]}; \mathbf{y}_{j+1}^*) \text{ 249} \rangle \equiv \quad /*$

$$t_next = \mathbf{t}_{j+1}$$

$$y_pred \supseteq \mathbf{y}_{j+1}^*$$

$$h_trial \supseteq \mathbf{h}_j$$

q order

taylor_coeff_var contains $h_j^i J(f^{[i]}; \mathbf{y}_{j+1}^*)$ for $i = 0, 1, \dots, q$

$$C_qp[i] \ni (-1)^i c_i^{q,p}$$

$$M \supseteq h_j^i J(f^{[i]}; \mathbf{y}_{j+1}^*)$$

$$M \supseteq (-1)^i c_i^{q,p} h_j^i J(f^{[i]}; \mathbf{y}_{j+1}^*)$$

$$B \supseteq \mathbf{B}_{j+1}$$

$*/$

$\text{ad} \rightarrow \text{taylor_coeff_var} \rightarrow \text{set}(t_next, y_pred, h_trial, q);$

$\text{ad} \rightarrow \text{taylor_coeff_var} \rightarrow \text{compTerms}();$

$\text{setM}(B, 0.0);$

for (**int** $i = q; i \geq 1; i--$) {

```

    ad→taylcoeff→var→getTerm(M, i);
    scaleM(M, C_qp[i]);
    addMiMi(B, M);
  }
  addId(B);

```

This code is used in chunk 247.

```

250  ⟨ Cj+1 = m(Bj+1) 250 ⟩ ≡      /*
      
$$\frac{B \supseteq \mathbf{B}_{j+1}}{C = C_{j+1} = m(\mathbf{B}_{j+1})}$$

      */
    midpoint(C, B);

```

This code is used in chunk 247.

The function *encloseInverse* below encloses the inverse of a point matrix. If this function fails, it returns *false*; otherwise, it returns *true*.

```

251  ⟨ Gj+1 = Cj+1-1Fj 251 ⟩ ≡      /*
      
$$\frac{\begin{array}{l} C = C_{j+1} \\ Fj \supseteq \mathbf{F}_j \end{array}}{C_{inv} \ni C_{j+1}^{-1} \text{ if } ok}$$

      
$$G \supseteq \mathbf{G}_{j+1} = C_{j+1}^{-1} \mathbf{F}_j$$

      */
    bool ok = matrix_inverse→encloseMatrixInverse(C_inv, C);
    if (¬ok)
    {
      control→ind = failure;
    }
    #ifndef VNODE_DEBUG
      printMessage("Could not invert the C matrix.");
    #endif
    return;
  }
  multMiMi(G, C_inv, Fj);

```

This code is used in chunk 247.

```

252  ⟨ Sj+1 = Gj+1Sj 252 ⟩ ≡      /*
      
$$\frac{\begin{array}{l} solution→S = S_j \\ G \supseteq \mathbf{G}_{j+1} \end{array}}{S \supseteq \mathbf{S}_{j+1} = \mathbf{G}_{j+1} S_j}$$


```

*/
 $\text{multMiMp}(S, G, \text{solution} \rightarrow S);$
 This code is used in chunk 247.

$$253 \quad \langle \mathbf{A}_{j+1} = \mathbf{G}_{j+1} A_j \quad 253 \rangle \equiv \quad / * \\
\frac{\text{solution} \rightarrow A = A_j \quad G \supseteq \mathbf{G}_{j+1}}{A \supseteq \mathbf{A}_{j+1} = \mathbf{G}_{j+1} A_j}$$

*/
 $\text{multMiMp}(A, G, \text{solution} \rightarrow A);$
 This code is used in chunk 247.

$$254 \quad \langle \mathbf{Q}_{j+1} = \mathbf{G}_{j+1} Q_j \quad 254 \rangle \equiv \quad / * \\
\frac{\text{solution} \rightarrow Q = Q_j \quad G \supseteq \mathbf{G}_{j+1}}{Q \supseteq \mathbf{Q}_{j+1} = \mathbf{G}_{j+1} Q_j}$$

*/
 $\text{multMiMp}(Q, G, \text{solution} \rightarrow Q);$
 This code is used in chunk 247.

Computing e_{j+1}

corrector_excess encloses $\mathbf{h}_j^{p+q+1} f^{[p+q+1]}(\mathbf{T}_j, \tilde{\mathbf{y}}_j)$. We multiply it by $(-1)^q \gamma_{p,q}$.

$$255 \quad \langle e_{j+1} = (-1)^q \gamma_{p,q} \mathbf{h}_j^{p+q+1} f^{[p+q+1]}(\mathbf{T}_j, \tilde{\mathbf{y}}_j) \quad 255 \rangle \equiv \quad / * \\
\frac{\text{corrector_excess} \supseteq \mathbf{h}_j^{p+q+1} f^{[p+q+1]}(\mathbf{T}_j, \tilde{\mathbf{y}}_j) \quad \text{errorConstant} \supseteq (-1)^q \gamma_{p,q}}{\text{corrector_excess} \supseteq (-1)^q \gamma_{p,q} \mathbf{h}_j^{p+q+1} f^{[p+q+1]}(\mathbf{T}_j, \tilde{\mathbf{y}}_j)}$$

*/
 $\text{scaleV}(\text{corrector_excess}, \text{errorConstant});$
 This code is used in chunk 247.

Enclosing g_{j+1}

$$256 \quad \langle g_{j+1} = \sum_{i=0}^p c_i^{p,q} h_j^i f^{[i]}(u_j) - \sum_{i=0}^q (-1)^i c_i^{q,p} h_j^i f^{[i]}(y_{j+1}^*) \quad 256 \rangle \equiv \\
\langle g_{j+1}^f = \sum_{i=0}^p c_i^{p,q} h_j^i f^{[i]}(u_j) \quad 257 \rangle \\
\langle g_{j+1}^b = \sum_{i=0}^q (-1)^i c_i^{q,p} h_j^i f^{[i]}(y_{j+1}^*) \quad 258 \rangle \\
\langle g_{j+1} = g_{j+1}^f - g_{j+1}^b \quad 259 \rangle$$

This code is used in chunk 247.

We have the terms $h_j^i f^{[i]}(u_j)$ for $i = 1, \dots, q$ computed in the predictor.

$$257 \quad \langle g_{j+1}^f = \sum_{i=0}^p c_i^{p,q} h_j^i f^{[i]}(u_j) \quad 257 \rangle \equiv \quad /*$$

tayl_coeff contains $h_j^i f^{[i]}(u_j)$ for $i = 0, 1, \dots, q$
p order

$$gj = 0$$

$$term \ni h_j^i f^{[i]}(u_j)$$

$$C_pq[i] \ni c_i^{p,q}$$

$$term \ni c_i^{p,q} h_j^i f^{[i]}(u_j)$$

$$gj \ni g_{j+1}^f = \sum_{i=0}^p c_i^{p,q} h_j^i f^{[i]}(u_j)$$

*/

setV(*gj*, 0.0);

for (**int** *i* = *p*; *i* ≥ 0; *i*--) {
 ad→taylor_coeff_ode→getTerm(*term*, *i*);
 scaleV(*term*, *C_pq*[*i*]);
 addViVi(*gj*, *term*);
}

This code is used in chunk 256.

$$258 \quad \langle g_{j+1}^b = \sum_{i=0}^q (-1)^i c_i^{q,p} h_j^i f^{[i]}(y_{j+1}^*) \quad 258 \rangle \equiv \quad /*$$

$$t_next = t_{j+1}$$

$$y_pred \supseteq y_{j+1}^*$$

$$h_trial \supseteq h_j$$

$$q \text{ order}$$

$$temp2 = 0$$

$$y_pred_point = y_{j+1}^* = m(y_{j+1}^*)$$

temp interval vector storing y_{j+1}^*

tayl_coeff contains $h_j^i f^{[i]}(y_{j+1}^*)$ after *compTerms* is called

$$term \ni h_j^i f^{[i]}(y_{j+1}^*)$$

$$C_qp[i] \ni (-1)^i c_i^{q,p}$$

$$term \ni (-1)^i c_i^{q,p} h_j^i f^{[i]}(y_{j+1}^*)$$

$$temp2 \ni g_{j+1}^b = \sum_{i=0}^q (-1)^i c_i^{q,p} h_j^i f^{[i]}(y_{j+1}^*)$$

*/

```

midpoint(y_pred_point, y_pred);
assignV(temp, y_pred_point);
ad_tayl_coeff_ode_set(t_next, temp, h_trial, q);
ad_tayl_coeff_ode_compTerms();

setV(temp2, 0.0);
for (int i = q; i ≥ 0; i--) {
    ad_tayl_coeff_ode_getTerm(term, i);
    scaleV(term, C_qp[i]);
    addViVi(temp2, term);
}

```

This code is used in chunk 256.

259 $\langle g_{j+1} = g_{j+1}^f - g_{j+1}^b \quad 259 \rangle \equiv \quad /*$

$$\frac{\begin{array}{l} gj \ni g_{j+1}^f \\ temp2 \ni g_{j+1}^b \end{array}}{gj \ni g_{j+1} = g_{j+1}^f - g_{j+1}^b}$$

```

*/
subViVi(gj, temp2);

```

This code is used in chunk 256.

Computing d_{j+1}

260 $\langle d_{j+1} = g_{j+1} + e_{j+1} \quad 260 \rangle \equiv \quad /*$

$$\frac{\begin{array}{l} corrector_excess \supseteq e_{j+1} \\ gj \ni g_{j+1} \end{array}}{d \supseteq d_{j+1} = g_{j+1} + e_{j+1}}$$

```

*/
addViVi(d, gj, corrector_excess);

```

This code is used in chunk 247.

Computing w_{j+1}

Now, we compute

261 $\langle \mathbf{w}_{j+1} = C_{j+1}^{-1} \mathbf{d}_{j+1} + (I - C_{j+1}^{-1} \mathbf{B}_{j+1})(\mathbf{y}_{j+1}^* - \mathbf{y}_{j+1}^*) \rangle \equiv \quad /*$

$$\begin{array}{l}
 y_pred_point = y_{j+1}^* \\
 y_pred \supseteq \mathbf{y}_{j+1}^* \\
 B \supseteq \mathbf{B}_{j+1} \\
 Cinv \ni C_j^{-1} \\
 d \supseteq \mathbf{d}_{j+1} \\
 \hline
 term \supseteq \mathbf{y}_{j+1}^* - y_{j+1}^* \\
 M \supseteq C_{j+1}^{-1} \mathbf{B}_{j+1} \\
 M \supseteq I - C_j^{-1} \mathbf{B}_{j+1} \\
 temp \supseteq (I - C_{j+1}^{-1} \mathbf{B}_{j+1})(\mathbf{y}_{j+1}^* - y_{j+1}^*) \\
 w \supseteq \mathbf{w}_{j+1} = C_{j+1}^{-1} \mathbf{d}_{j+1} + (I - C_{j+1}^{-1} \mathbf{B}_{j+1})(\mathbf{y}_{j+1}^* - y_{j+1}^*) \\
 \hline
 */ \\
 subViVp(term, y_pred, y_pred_point); \\
 multMiMi(M, Cinv, B); \\
 subFromId(M); \\
 multMiVi(temp, M, term); \\
 multMiVi(w, Cinv, d); \quad /* w = Cinv * d + temp */ \\
 addViVi(w, temp);
 \end{array}$$

This code is used in chunk 247.

262 $\langle \mathbf{s}_{j+1} = (\mathbf{A}_{j+1} \mathbf{r}_j) \cap (\mathbf{Q}_{j+1} \mathbf{r}_{QR,j}) \rangle \equiv \quad /*$

$$\begin{array}{l}
 solution\text{-}r \supseteq \mathbf{r}_j \\
 A \supseteq \mathbf{A}_{j+1} \\
 solution\text{-}rQR \supseteq \mathbf{r}_{QR,j} \\
 Q \supseteq \mathbf{Q}_{j+1} \\
 \hline
 temp \supseteq \mathbf{A}_{j+1} \mathbf{r}_j \\
 s \supseteq \mathbf{Q}_{j+1} \mathbf{r}_{QR,j} \\
 s \supseteq \mathbf{s}_{j+1} = (\mathbf{A}_{j+1} \mathbf{r}_j) \cap (\mathbf{Q}_{j+1} \mathbf{r}_{QR,j}) \\
 \hline
 */ \\
 multMiVi(temp, A, solution\text{-}r); \\
 multMiVi(s, Q, solution\text{-}rQR); \\
 b = intersect(s, temp, s); \\
 assert(b);
 \end{array}$$

This code is used in chunk 247.

```

263  $\langle \mathbf{y}_{j+1} = (\mathbf{y}_{j+1}^* + \mathbf{S}_{j+1}\boldsymbol{\alpha} + \mathbf{s}_{j+1} + \mathbf{w}_{j+1}) \cap \mathbf{y}_{j+1}^* \rangle \equiv$   $/*$ 

     $trial\_solution\text{-}\alpha \supseteq \boldsymbol{\alpha}$ 
     $S \supseteq \mathbf{S}_{j+1}$ 
     $y\_pred \supseteq \mathbf{y}_{j+1}^*$ 
     $y\_pred\_point = \mathbf{y}_{j+1}^*$ 
     $s \supseteq \mathbf{s}_{j+1}$ 
     $w \supseteq \mathbf{w}_{j+1}$ 

    -----
     $globalExcess \supseteq \mathbf{s}_{j+1} + \mathbf{w}_{j+1}$ 
     $temp \supseteq \mathbf{S}_{j+1}\boldsymbol{\alpha}$ 
     $temp \supseteq \mathbf{S}_{j+1}\boldsymbol{\alpha} + \mathbf{s}_{j+1} + \mathbf{w}_{j+1}$ 
     $temp \supseteq \mathbf{y}_{j+1}^* + \mathbf{S}_{j+1}\boldsymbol{\alpha} + \mathbf{s}_{j+1} + \mathbf{w}_{j+1}$ 
     $trial\_solution\text{-}y \supseteq \mathbf{y}_{j+1} = (\mathbf{y}_{j+1}^* + \mathbf{S}_{j+1}\boldsymbol{\alpha} + \mathbf{x}_{j+1}) \cap \mathbf{y}_{j+1}^*$ 

    */
    addViVi(globalExcess, s, w);
    multMiVi(temp, S, trial_solution-alpha);
    addViVi(temp, globalExcess);
    addViVp(temp, y-pred_point);
    b = intersect(trial_solution-y, y-pred, temp);
    assert(b);

```

This code is used in chunk 247.

```

264  $\langle \text{set } t_{j+1} \rangle \equiv$ 
     $trial\_solution\text{-}t = t\_next;$ 

```

This code is used in chunk 236.

20.2.6 Enclosure representation

```

265  $\langle \text{find solution representation for next step} \rangle \equiv$ 
     $\langle u_{j+1} = m(\mathbf{y}_{j+1}) \rangle$ 
     $\langle S_{j+1} = m(\mathbf{S}_{j+1}) \rangle$ 
     $\langle \mathbf{v}_{j+1} = \mathbf{y}_{j+1}^* - u_{j+1} + (\mathbf{S}_{j+1} - \mathbf{S}_{j+1})\boldsymbol{\alpha} + \mathbf{w}_{j+1} \rangle$ 
     $\langle A_{j+1} = m(\mathbf{A}_{j+1}) \rangle$ 
     $\langle \mathbf{r}_{j+1} = (\mathbf{A}_{j+1}^{-1}\mathbf{A}_{j+1})\mathbf{r}_j + \mathbf{A}_{j+1}^{-1}\mathbf{v}_{j+1} \rangle$ 
     $\langle \text{compute } Q_{j+1} \rangle$ 
     $\langle \mathbf{r}_{QR,j+1} = (\mathbf{Q}_{j+1}^{-1}\mathbf{Q}_{j+1})\mathbf{r}_{QR,j} + \mathbf{Q}_{j+1}^{-1}\mathbf{v}_{j+1} \rangle$ 
     $\langle \text{reset if needed} \rangle$ 

```

This code is used in chunk 236.

266 $\langle u_{j+1} = m(\mathbf{y}_{j+1}) \text{ 266} \rangle \equiv \quad /*$

$$\frac{\text{trial_solution} \rightarrow y \supseteq \mathbf{y}_{j+1}}{\text{trial_solution} \rightarrow u = u_{j+1}}$$

$*/$
 $\text{midpoint}(\text{trial_solution} \rightarrow u, \text{trial_solution} \rightarrow y);$
 This code is used in chunk 265.

267 $\langle S_{j+1} = m(\mathbf{S}_{j+1}) \text{ 267} \rangle \equiv \quad /*$

$$\text{trial_solution} \rightarrow S = S_{j+1} = m(\mathbf{S}_{j+1})$$

$*/$
 $\text{midpoint}(\text{trial_solution} \rightarrow S, S);$
 This code is used in chunk 265.

268 $\langle \mathbf{v}_{j+1} = \mathbf{y}_{j+1}^* - u_{j+1} + (\mathbf{S}_{j+1} - S_{j+1})\boldsymbol{\alpha} + \mathbf{w}_{j+1} \text{ 268} \rangle \equiv \quad /*$

$$\begin{array}{l} y_pred_point = \mathbf{y}_{j+1}^* \\ \text{trial_solution} \rightarrow u = u_{j+1} \\ S \supseteq \mathbf{S}_{j+1} \\ \text{trial_solution} \rightarrow S = S_{j+1} \\ \text{trial_solution} \rightarrow \alpha \supseteq \boldsymbol{\alpha} \\ w \supseteq \mathbf{w}_{j+1} \\ \hline S \supseteq \mathbf{S}_{j+1} - S_{j+1} \\ z \supseteq (\mathbf{S}_{j+1} - S_{j+1})\boldsymbol{\alpha} \\ z \supseteq (\mathbf{S}_{j+1} - S_{j+1})\boldsymbol{\alpha} + \mathbf{w}_{j+1} \\ z \supseteq -u_{j+1} + (\mathbf{S}_{j+1} - S_{j+1})\boldsymbol{\alpha} + \mathbf{w}_{j+1} \\ z \supseteq \mathbf{v}_{j+1} = \mathbf{y}_{j+1}^* - u_{j+1} + (\mathbf{S}_{j+1} - S_{j+1})\boldsymbol{\alpha} + \mathbf{w}_{j+1} \end{array}$$

$*/$
 $\text{subMiMp}(S, \text{trial_solution} \rightarrow S);$
 $\text{multMiVi}(z, S, \text{trial_solution} \rightarrow \alpha);$
 $\text{addViVi}(z, w);$
 $\text{subViVp}(z, \text{trial_solution} \rightarrow u);$
 $\text{addViVp}(z, y_pred_point);$

This code is used in chunk 265.

269 $\langle A_{j+1} = m(\mathbf{A}_{j+1}) \text{ 269} \rangle \equiv \quad /*$

$$\frac{A \supseteq \mathbf{A}_{j+1}}{\text{trial_solution} \rightarrow A = A_{j+1}}$$

*/
midpoint(*trial_solution*→*A*, *A*);

This code is used in chunk 265.

270 $\langle \mathbf{r}_{j+1} = (A_{j+1}^{-1} \mathbf{A}_{j+1}) \mathbf{r}_j + A_{j+1}^{-1} \mathbf{v}_{j+1} \quad 270 \rangle \equiv \quad / *$

trial_solution→*A* = *A*_{*j*+1}

A \supseteq *A*_{*j*+1}

z \supseteq *v*_{*j*+1}

solution→*r* \supseteq *r*_{*j*}

Ainv $\ni A_{j+1}^{-1}$ if *ok*

temp $\supseteq A_{j+1}^{-1} \mathbf{v}_{j+1}$

M $\supseteq A_{j+1}^{-1} \mathbf{A}_{j+1}$

trial_solution→*r* $\supseteq (A_{j+1}^{-1} \mathbf{A}_{j+1}) \mathbf{r}_j$

trial_solution→*r* $\supseteq \mathbf{r}_{j+1} = (A_{j+1}^{-1} \mathbf{A}_{j+1}) \mathbf{r}_j + A_{j+1}^{-1} \mathbf{v}_{j+1}$

*/
ok = *matrix_inverse_encloseMatrixInverse*(*Ainv*, *trial_solution*→*A*);

if (*ok*) {
 multMiVi(*temp*, *Ainv*, *z*);
 multMiMi(*M*, *Ainv*, *A*);
 multMiVi(*trial_solution*→*r*, *M*, *solution*→*r*);
 addViVi(*trial_solution*→*r*, *temp*);
}

This code is used in chunk 265.

271 $\langle \mathbf{r}_{\text{QR},j+1} = (Q_{j+1}^{-1} \mathbf{Q}_{j+1}) \mathbf{r}_{\text{QR},j} + Q_{j+1}^{-1} \mathbf{v}_{j+1} \quad 271 \rangle \equiv \quad / *$

trial_solution→*Q* = *Q*_{*j*+1}

z \supseteq *v*_{*j*+1}

solution→*rQR* \supseteq *r*_{QR,*j*}

Q \supseteq *Q*_{*j*+1}

Ainv $\ni Q_{j+1}^{-1}$ if *ok*

temp $\supseteq Q_{j+1}^{-1} \mathbf{v}_{j+1}$

M $\supseteq Q_{j+1}^{-1} \mathbf{Q}_{j+1}$

trial_solution→*rQR* $\supseteq (Q_{j+1}^{-1} \mathbf{Q}_{j+1}) \mathbf{r}_{\text{QR},j}$

trial_solution→*rQR* $\supseteq \mathbf{r}_{\text{QR},j+1} = (Q_{j+1}^{-1} \mathbf{Q}_{j+1}) \mathbf{r}_{\text{QR},j} + Q_{j+1}^{-1} \mathbf{v}_{j+1}$

*/

```

    b = matrix_inverse_orthogonalInverse(Ainv, trial_solution-Q);
    if (b == false) {
        control-ind = failure;
#ifdef VNODE_DEBUG
        printMessage("Could not invert the Q matrix.");
#endif
    return;
    }
    multMiVi(temp, Ainv, z);
    multMiMi(M, Ainv, Q);
    multMiVi(trial_solution-rQR, M, solution-rQR);
    addViVi(trial_solution-rQR, temp);

```

This code is used in chunk 265.

272 $\langle \text{reset if needed } 272 \rangle \equiv$ $/*$

$$\begin{array}{l}
 \text{trial_solution-Q} = Q_{j+1} \\
 \text{trial_solution-rQR} \supseteq \mathbf{r}_{\text{QR},j+1} \\
 \text{trial_solution-A} = A_{j+1} \\
 \text{trial_solution-r} \supseteq \mathbf{r}_{j+1} \\
 \hline
 \text{temp} \supseteq Q_{j+1} \mathbf{r}_{\text{QR},j+1} \\
 \text{temp2} \supseteq A_{j+1} \mathbf{r}_{j+1} \\
 \text{if (subseteq(temp, temp2) } \vee \neg \text{ok) } \{ \\
 \quad \text{trial_solution-A} = Q_{j+1} \\
 \quad \text{trial_solution-r} \supseteq \mathbf{r}_{\text{QR},j+1} \\
 \}
 \end{array}$$

```

    */
    multMpVi(temp, trial_solution-Q, trial_solution-rQR);
    multMpVi(temp2, trial_solution-A, trial_solution-r);
    if (subseteq(temp, temp2) ∨ ¬ok) {
        assignM(trial_solution-A, trial_solution-Q);
        assignV(trial_solution-r, trial_solution-rQR);
    }

```

This code is used in chunk 265.

Computing Q_{j+1}

```

273  ⟨ compute  $Q_{j+1}$  273 ⟩ ≡      /*
                                      $Q \supseteq Q_{j+1}$ 
                                     -----
                                      $A\_point = m(Q_{j+1})$ 
                                      $A\_point = m(Q_{j+1}) \text{ diag } w(r_{QR,j})$ 
                                      $C = A\_point$  with columns sorted in non-increasing order in  $\|\cdot\|_2$ 
                                      $C = Q_{j+1} R_{j+1}$ 
                                      $trial\_solution-Q = Q_{j+1}$ 

                                     */
midpoint( $A\_point$ ,  $Q$ );
int  $n = sizeM(Q)$ ;
for (int  $i = 0$ ;  $i < n$ ;  $i++$ )
    for (int  $j = 0$ ;  $j < n$ ;  $j++$ )
         $A\_point[i][j] *= v\_bias :: width(solution-rQR[j])$ ;
if ( $\neg(inf\_normM(A\_point) < numeric\_limits<double>::max())$ ) {
#ifdef VNODE_DEBUG
     $printMessage("The\_computed\_enclosures\_are\_too\_wide")$ ;
#endif
     $control-ind = failure$ ;
    return;
}
sortColumns( $C$ ,  $A\_point$ );
 $b = computeQR(trial\_solution-Q, C)$ ;
assert( $b$ );

```

This code is used in chunk 265.

20.2.7 Constructor

```

274  ⟨ constructor IHO 274 ⟩ ≡
    IHO :: IHO(int  $n$ )
    {
         $order\_trial = 0$ ;
         $solution = new\ Solution(n)$ ;
         $trial\_solution = new\ Solution(n)$ ;
         $matrix\_inverse = new\ MatrixInverse(n)$ ;
         $C\_pq = C\_qp = 0$ ;
         $sizeV(y, n)$ ;
         $sizeV(y\_pred, n)$ ;
         $sizeV(globalExcess, n)$ ;
         $sizeV(y\_pred\_point, n)$ ;
         $sizeV(temp, n)$ ;
         $sizeV(temp2, n)$ ;
         $sizeV(x, n)$ ;
    }

```

```

sizeV(u_next, n);
sizeV(predictor_excess, n);
sizeV(corrector_excess, n);
sizeV(z, n);
sizeV(w, n);
sizeV(gj, n);
sizeV(term, n);
sizeV(d, n);
sizeV(s, n);
sizeM(Fj, n);
sizeM(M, n);
sizeM(Cinv, n);
sizeM(G, n);
sizeM(B, n);
sizeM(S, n);
sizeM(A, n);
sizeM(Q, n);
sizeM(U, n);
sizeM(V, n);
sizeM(Ainv, n);
sizeM(C, n);
sizeM(A_point, n);
}

```

This code is used in chunk 293.

20.2.8 Destructor

```

275 < destructor IHO 275 > ≡
    IHO :: ~IHO()
    {
        delete matrix_inverse;
        delete trial_solution;
        delete solution;
        delete[] C_pq;
        delete[] C_qp;
    }

```

This code is used in chunk 293.

20.2.9 Accepting a solution

```

276 < accept solution (IHO) 276 > ≡
    void IHO :: acceptSolution() {
        solution→t = trial_solution→t;
        assignV(solution→y, trial_solution→y);
        assignV(solution→u, trial_solution→u);
    }

```

```

    assignM (solution→S, trial_solution→S);
    assignV (solution→alpha, trial_solution→alpha);
    assignM (solution→A, trial_solution→A);
    assignV (solution→r, trial_solution→r);
    assignM (solution→Q, trial_solution→Q);
    assignV (solution→rQR, trial_solution→rQR);
}

```

This code is used in chunk 293.

20.2.10 Set and get functions

```

277 <set and get functions 277> ≡
    void set(Control *c, HOE *hoem, AD *ad0)
    {
        assert(c ∧ hoem ∧ ad0);
        control = c;
        hoe = hoem;
        ad = ad0;
    }
    void getTightEnclosure(iVector &y) const
    {
        y = solution→y;
    }
    void init(const interval &t, const iVector &y)
    {
        trial_solution→init(t, y);
        solution→init(t, y);
    }
    const iVector &getGlobalExcess() const
    {
        return globalExcess;
    }

```

This code is used in chunk 235.

20.2.11 Constants

Computing $c_i^{p,q}$

We show how to compute

$$c_i^{p,q} = \frac{p!}{(p+q)!} \frac{(q+p-i)!}{(p-i)!}.$$

From

$$\begin{aligned}
 c_{i-1}^{p,q} &= \frac{p!}{(p+q)!} \frac{(q+p-i+1)!}{(p-i+1)!} = \frac{p!}{(p+q)!} \frac{(q+p-i)!}{(p-i)!} \frac{q+p-i+1}{p-i+1} \\
 &= \frac{q+p-i+1}{p-i+1} c_i^{p,q}
 \end{aligned}$$

we compute

$$c_i^{p,q} = \frac{p-i+1}{q+p-i+1} c_{i-1}^{p,q}, \quad \text{where } c_0^{p,q} = 1.$$

```

279 < c_i^{p,q} = \frac{p!}{(p+q)!} \frac{(q+p-i)!}{(p-i)!} 279 > \equiv
    void IHO :: compCpq (int p, int q)
    {
        /*
            C_pq[i] \ni c_i^{p,q}    for i = 1, \dots, p

            */
        int tmp = q + p + 1;
        C_pq[0] = 1.0;
        for (int i = 1; i \leq p; i++) {
            C_pq[i] = (C_pq[i-1] * double(p-i+1))/double(tmp-i);
        }
    }

```

This code is used in chunk 293.

Computing $(-1)^i c_i^{q,p}$

When computing $c_i^{q,p}$, we multiply them by $(-1)^i$, that is $(-1)^i c_i^{q,p}$.

From

$$\begin{aligned}
 c_{i-1}^{q,p} &= \frac{q!}{(q+p)!} \frac{(p+q-i+1)!}{(q-i+1)!} = \frac{q!}{(q+p)!} \frac{(p+q-i)!}{(q-i)!} \frac{p+q-i+1}{q-i+1} \\
 &= \frac{p+q-i+1}{q-i+1} c_i^{q,p}
 \end{aligned}$$

we compute

$$(-1)^i c_i^{q,p} = (-1) \frac{q-i+1}{p+q-i+1} ((-1)^{i-1} c_{i-1}^{q,p}), \quad \text{where } c_0^{q,p} = 1.$$

```

280 < (-1)^i c_i^{q,p} = (-1)^i \frac{q!}{(p+q)!} \frac{(q+p-i)!}{(q-i)!} 280 > \equiv
    void IHO :: compCqp (int p, int q)
    {
        /*
            C_qp[i] \ni (-1)^i c_i^{q,p}    for i = 1, \dots, q

            */
        int tmp = q + p + 1;
        C_qp[0] = 1.0;
        for (int i = 1; i \leq q; i++) {
            C_qp[i] = (-C_qp[i-1] * double(q-i+1))/double(tmp-i);
        }
    }

```

This code is used in chunk 293.

Even number

This function returns *true* if its arguments ($k \geq 0$) is even and *false* otherwise. We check if the last bit of k is 1.

```

281 <check if a number is even 281> ≡
    inline bool isEven(unsigned int k)
    {
        if (k ≡ 0) return true;
        if (k & #01) return false;
        return true;
    }

```

This code is used in chunk 293.

Computing $(-1)^q q! p! / (p + q)!$

For given p and q , *compErrorConst* computes $(-1)^q q! p! / (p + q)!$. We use the relation

$$\frac{q! p!}{(p + q)!} = \frac{p!}{(q + 1)(q + 2) \cdots (p + q)} = \frac{1}{q + 1} \frac{2}{q + 2} \cdots \frac{p}{q + p}.$$

```

282 <compute  $(-1)^q q! p! / (p + q)!$  282> ≡      /*

```

```

    err_const ≥  $(-1)^q q! p! / (p + q)!$ 

    */
    interval IHO :: compErrorConstant(int p, int q)
    {
        interval err_const(1.0);
        for (int i = 1; i ≤ p; i++) {
            err_const *= double(i);
            err_const /= (q + i);
        }
        if (¬isEven(q)) err_const = -err_const;
        return err_const;
    }

```

This code is used in chunk 293.

20.2.12 Sorting columns of a matrix

Given an $n \times n$ point matrix A , we sort its columns such that, in the 2-norm, they are in non-increasing order. In *sortColumns* the columns of the input matrix A are sorted, and the result is stored in B .

```

284 < sort columns 284 > ≡
    void sortColumns(pMatrix &B, const pMatrix &A)
    {
        < compute column norms of A 286 >
        < check if sorting is needed 287 >
        < perform sorting 289 >
    }

```

This code is used in chunk 291.

First, we create a structure to store the index and the norm of the column corresponding to this index.

```

285 < index-norm structure 285 > ≡
    struct index_norm {
        int index;
        double norm;
    };

```

This code is used in chunk 291.

We store in an array b in $b[i].norm$ the 2-norm of column i of A .

```

286 < compute column norms of A 286 > ≡
    int n = v_blas::sizeM(B);
    index_norm *b = new index_norm[n];
    pVector tmp;
    v_blas::sizeV(tmp, n);
    for (int j = 0; j < n; j++) {
        b[j].index = j;
        getColumn(tmp, A, j);
        b[j].norm = v_blas::norm2(tmp);
    }

```

This code is used in chunk 284.

```

287 < check if sorting is needed 287 > ≡
    bool sorting_needed = false;
    for (int i = 0; i < n - 1; i++) {
        if (b[i].norm < b[i + 1].norm) {
            sorting_needed = true;
            break;
        }
    }

```

This code is used in chunk 284.

We use the standard *qsort* function. We need a *compare* function for it.

```

288 < compare function 288 > ≡

```



```

inline int v_compare(const void *a1, const void *b1)
{
    index_norm *a = (index_norm *) a1;
    index_norm *b = (index_norm *) b1;

    if (a->norm > b->norm) return -1;
    if (a->norm < b->norm) return 1;
    return 0;
}

```

This code is used in chunk 291.

We sort the b array based on $b[i].norm$. Then the j th column of B is the $b[j].index$ column of A .

```

289 <perform sorting 289> ≡
    if (sorting_needed) std::qsort((void *) b, n, sizeof(index_norm), v_compare);
    for (int j = 0; j < n; j++)
    {
        getColumn(tmp, A, b[j].index);
        setColumn(B, tmp, j);
    }
    delete[] b;

```

This code is used in chunk 284.

Files

```

291 <sortcolumns.cc 291> ≡
#include <cstdlib>
#include <cassert>
#include "vnodeinterval.h"
#include "basiclinalg.h"
namespace vnodelp {
    using namespace v_bias;
    using namespace v_blas;

    <index-norm structure 285>
    <compare function 288>
    <sort columns 284>
}

```

```

292 <iho.h 292> ≡
#ifndef IHO_H
#define IHO_H
#include <cassert>
#include <cmath>
#include "ad_ode.h"
#include "ad_var.h"

```

```

#include "hoe.h"
#include "intvfuncs.h"
#include "basiclinalg.h"
#include "matrixinverse.h"
    namespace vnodelp {
        using namespace v_bias;
        using namespace v_blas;
        < class IHO 235 >
    }
#endif

293 < iho.cc 293 > ≡
#include <cmath>
#include "vnodeinterval.h"
#include "basiclinalg.h"
#include "solution.h"
#include "control.h"
#include "allad.h"
#include "matrixinverse.h"
#include "iho.h"
#include "debug.h"
#include "miscfuncs.h"

    namespace vnodelp {
        using namespace v_bias;
        using namespace v_blas;
        extern void sortColumns(pMatrix &B, const pMatrix &A);
        extern bool computeQR(pMatrix &B, const pMatrix &A);
        < constructor IHO 274 >
        < destructor IHO 275 >
        < check if a number is even 281 >
        <  $c_i^{p,q} = \frac{p!}{(p+q)!} \frac{(q+p-i)!}{(p-i)!}$  279 >
        <  $(-1)^i c_i^{q,p} = (-1)^i \frac{q!}{(p+q)!} \frac{(q+p-i)!}{(q-i)!}$  280 >
        < compute  $(-1)^q q! p! / (p+q)!$  282 >
        < accept solution (IHO) 276 >
        < compute tight enclosure 236 >
        < compute IHO method coefficients 239 > }

```

Chapter 21

The VNODE class

21.1 Declaration

```
295 < class VNODE 295 > ≡
    typedef enum {
        on, off
    } stepAction;
    class VNODE {
    public:
        int steps;
        VNODE(AD *ad);
        void integrate(interval &t0, iVector &y0, const interval &t_end);
        < set VNODE parameters 331 >
        < get functions (VNODE) 329 >
        ~VNODE();
    private:
        void acceptSolution(interval &t0, iVector &y0);
        double compHstart(const interval &t0, const iVector &y0);
    private:
        int direction;
        interval t_trial, Tj, h_accepted;
        iVector temp;
        double h_start, h_min;
        pVector tp;
        Control *control;
        HOE *hoe;
        IHO *iho;
        AD *ad;
    };

```

This code is used in chunk 333.

21.2 The integrator function

The integrator function consists of the following steps:

1. check input correctness and determine the direction of the integration
2. initialize
3. validate existence and uniqueness and select stepsize
4. check if the end point is reached
5. compute a tight enclosure
6. decide how to proceed

```

296 <integrator 296> ≡
    void VNODE::integrate(interval &t0, iVector &y0, const interval &t_end)
    {
        <check input correctness 298>
        <determine direction 305>
        if (control-ind ≡ first_entry) {
            <initialize integration 308>
        }
        while (t0 ≠ t_end) {
            if (control-ind ≡ first_entry ∨ control-ind ≡ success) {
                <validate and select stepsize 317>
                <check if last step 318>
            }
            <compute enclosure 324>
            <decide 325>
        } /* restore hmin */
        control-hmin = h_min;
    }

```

This code is used in chunk 334.

21.2.1 Input correctness

control-ind must not be *failure*

```

298 <check input correctness 298> ≡
    if (control-ind ≡ failure) {
        vnodeMessage("Previous_call_to_integrate()_failed");
        vnodeMessage("Call_setFirstEntry()_before_calling_integrate");
        return;
    }

```

See also chunks 299, 300, 301, 302, 303, and 304.

This code is used in chunk 296.

The initial condition and the final points must be representable machine intervals.

```

299 <check input correctness 298> +≡
    if (¬v_bias::finite_interval(t0) ∨ ¬v_bias::finite_interval(t_end)) {
        vnodeMessage("t0_and_t_end_must_be_finite");
        control-ind = failure;
        return;
    }
    for (unsigned int i = 0; i < size V(y0); i++)
        if (¬v_bias::finite_interval(y0[i])) {
            vnodeMessage("y0_must_contain_finite_intervals");
            control-ind = failure;
            return;
        }

```

If the integrator is re-entered with the same end point, t_{end} , as the t_{end} from the most recent call, then it should return.

```

300 <check input correctness 298> +≡
    t_trial = t0;
    if (t_trial ≡ t_end ∧ control-ind ≠ first_entry) {
        vnodeMessage("Set_different_t_end");
        return;
    }

```

At least one of the tolerances must be positive.

```

301 <check input correctness 298> +≡
    if (control-atol ≤ 0 ∧ control-rtol ≤ 0) {
        vnodeMessage("Set_nonnegative_tolerances_and_at_least\
            one_positive_tolerance");
        return;
    }

```

The order must be between 3 and $getMaxOrder()$.

```

302 <check input correctness 298> +≡
    if (control-order < 3 | control-order > getMaxOrder()) {
        vnodeMessage("Set_order_>=3_and_order_<=MAX_ORDER");
        control-ind = failure;
        return;
    }

```

The value for the minimum stepsize must be non-negative

```

303 <check input correctness 298> +≡
    if (control-hmin < 0) {
        vnodeMessage("Set_minimum_stepsize_>=0");
        return;
    }

```

The initial and final time intervals may not intersect.

```

304 <check input correctness 298> +≡
    if (¬v_bias::disjoint(t0, t_end)) {
        control-ind = failure;
        vnodeMessage("t0_and_t_end_must_be_disjoint");
        return;
    }

```

21.2.2 Determine direction

The user provides t_0 and t_{end} . The solver decides if the integration is from left to right or right to left.

If a first entry into *integrate*, we initialize *direction* with 0. We also save the direction from the most recent integration.

```

305 <determine direction 305> ≡
    if (control-ind ≡ first_entry) direction = 0;
    int last_direction = direction;

```

See also chunks 306 and 307.

This code is used in chunk 296.

1. If $\bar{t}_0 < \underline{t}_{\text{end}}$, the integration is in positive direction.
2. If $\underline{t}_0 > \bar{t}_{\text{end}}$, the integration is in negative direction.

```

306 <determine direction 305> +≡
    direction = 1;
    if (v_bias::sup(t_end) < v_bias::inf(t0))
        direction = -1;

```

If the direction of the integration is reversed without calling first *setFirstEntry()*, then print a message and return.

```

307 <determine direction 305> +≡
    if (last_direction ≠ 0 ∧ last_direction ≡ -direction) {
        vnodeMessage("Integration_direction_changed_without_r\
            esetting.\n" "Call_setFirstEntry()_before_integrating.");
        control-ind = failure;
        return;
    }

```

21.2.3 Initialization

Initializing an integration includes the following parts discussed below.

```

308 < initialize integration 308 > ≡
    steps = 0;
    < find minimum stepsize 309 >
    < find initial stepsize 310 >
    < set IHO method 312 >
    < set HOE method 311 >

```

This code is used in chunk 296.

Find minimum stepsize

- If *control-hmin* $\equiv 0$, the solver computes minimum stepsize.
- If *control-ind* \equiv *success*, this is a re-entry to the integrator; it computes minimum stepsize.
- If *control-hmin* $> 0 \wedge$ *control-ind* \equiv *first_entry*, the user has specified minimum stepsize. We still have to check if this value is not smaller than what the solver finds as a minimum stepsize.

```

309 < find minimum stepsize 309 > ≡
    h_min = control-hmin;
    if (control-hmin  $\equiv 0 \vee$  control-ind  $\equiv$  success)
        control-hmin = compHmin(t0, t_end);
    if (control-hmin  $> 0 \wedge$  control-ind  $\equiv$  first_entry) {
        double h = compHmin(t0, t_end);
        if (control-hmin  $< h$ ) {
            control-ind = failure;
            vnodeMessage("Set_larger_value_for_hmin");
            return;
        }
        assert(control-hmin  $> 0$ );
    }

```

This code is used in chunk 308.

Find initial stepsize

If *control-ind* \equiv *first_entry*, the solver computes initial stepsize.

```

310 < find initial stepsize 310 > ≡
    if (control-ind  $\equiv$  first_entry) {
        h_start = compHstart(t0, y0);
    }
    if (h_start  $<$  control-hmin) {
        control-ind = failure;
        vnodeMessage("Minimum_stepsize_reached.");
        return;
    }

```

```

    if (direction  $\equiv$  -1) {
        if (control→ind  $\equiv$  first_entry) h_start = -h_start;
    }

```

This code is used in chunk 308.

Setting the HOE method

We set in the HOE method the value for the initial stepsize, order, control structure, and an AD object.

```

311 <set HOE method 311>  $\equiv$ 
    hoe→setTrialStepsize(h_start);
    hoe→setTrialOrder(control→order);
    hoe→set(control, ad);
    hoe→init(t0, y0);

```

This code is used in chunk 308.

Setting the IHO method

In the IHO method, we set the control object, the HOE method, and an AD object. We also compute the coefficients of the method and set the initial point.

```

312 <set IHO method 312>  $\equiv$ 
    iho→set(control, hoe, ad);
    iho→compCoeffs();
    iho→init(t0, y0);

```

This code is used in chunk 308.

21.2.4 Methods involved in the initialization

Minimum stepsize

We determine the minimum magnitude, h_{\min} , for the stepsize allowed as follows:

$$\hat{t} = \max\{|\mathbf{t}_0|, |\mathbf{t}_{\text{end}}|\},$$

$$h_{\min} = \max\{\Delta(\text{next}(\hat{t}) - \hat{t}), w(\mathbf{t}_{\text{end}})\}.$$

Here $\text{next}(a)$ returns the next representable machine number to the right of a .

```

314 <compute  $h_{\min}$  314>  $\equiv$ 
    double compHmin(const interval &t0, const interval &t_end)
    {
        /*
             $t_0 = \mathbf{t}_0$ 
             $t_{\text{end}} = \mathbf{t}_{\text{end}}$ 

```

```

             $t = \max\{|\mathbf{t}_0|, |\mathbf{t}_{\text{end}}|\}$ 
             $t_{\text{next}} =$  next machine number to the right of  $t$ 
             $t = \Delta(\text{next}(\hat{t}) - \hat{t})$ 
             $t = \max\{\Delta(\text{next}(\hat{t}) - \hat{t}), w(\mathbf{t}_{\text{end}})\}$ 

```



```

    */
    double t = std::max(v_bias::mag(t0), v_bias::mag(t_end));
    double t_next = nextafter(t, t + 1);
    v_bias::round_up();
    t = t_next - t;
    t = std::max(t, v_bias::width(t_end));
    return t;
}

```

This code is used in chunk 334.

Initial stepsize

On the first step we compute

$$h_{0,0} = \left(\frac{\text{tol}}{\|(k+1)f^{[k+1]}(\mathbf{y}_0)\|} \right)^{1/k},$$

(see [26]) where

$$\text{tol} = \text{rtol} \cdot \|\mathbf{y}_0\| + \text{atol}.$$

315 \langle compute initial stepsize 315 $\rangle \equiv$

```

double VNODE::compHstart(const interval &t0, const iVector &y0)
{
    /*

```

$$\begin{array}{l} y0 = \mathbf{y}_0 \\ k \quad \text{order} \end{array}$$

tayl_coeff_ode contains $f^{[i]}(\mathbf{t}_0, \mathbf{y}_0)$ for $i = 0, 1, \dots, k+1$
after *compTerm* is called

$$\text{temp} \supseteq f^{[k+1]}(\mathbf{t}_0, \mathbf{y}_0)$$

$$\text{tol} = \text{tol} = \text{rtol} \cdot \|\mathbf{y}_0\| + \text{atol}$$

$$t = (k+1) \|f^{[k+1]}(\mathbf{t}_0, \mathbf{y}_0)\|$$

$$t = \frac{\text{tol}}{(k+1) \|f^{[k+1]}(\mathbf{t}_0, \mathbf{y}_0)\|}$$

$$t = \left(\frac{\text{tol}}{(k+1) \|f^{[k+1]}(\mathbf{t}_0, \mathbf{y}_0)\|} \right)^{1/k}$$

```

    */
    int k = control_order;
    ad_tayl_coeff_ode.set(t0, y0, 1, k + 1);
    ad_tayl_coeff_ode.compTerms();
    ad_tayl_coeff_ode.getTerm(temp, k + 1);
    double tol = control_rtol * inf_norm V(y0) + control_atol;
    double t = (k + 1) * inf_norm V(temp);

```

```

    < check  $t$  316 >
     $t = tol/t$ ;
     $t = \mathbf{std}::pow(t, 1.0/k)$ ;
    return  $t$ ;
}

```

This code is used in chunk 334.

The t that we compute above is normally $t > 0$. To prevent the case that it might be zero, we set t to be the largest of *control- $atol$* and *control- $rtol$* .

```

316 < check  $t$  316 >  $\equiv$ 
    if ( $t \equiv 0$ )  $t = \mathbf{std}::max(\text{control-}atol, \text{control-}rtol)$ ;

```

This code is used in chunk 315.

21.2.5 Validate existence and uniqueness

We try to validate existence and uniqueness with \mathbf{t}_j and \mathbf{y}_j . If *validate* returns *false*, we set *control- ind* = *failure*.

```

317 < validate and select stepsize 317 >  $\equiv$       /*
         $t0 = \mathbf{t}_j$ 
         $y0 = \mathbf{y}_j$ 

        */
    bool info;
    hoe-compAprioriEnclosure( $t0, y0, info$ );
    if (info  $\equiv$  false) {
        control-ind = failure;
        vnodeMessage("Could_not_validate_solution");
        return;
    }

```

This code is used in chunk 296.

21.2.6 Check last step

If *validate* returns successfully, we obtain an interval $[t_j, t_{j+1}]$, or $[t_{j+1}, t_j]$ if the integration is backwards, over which we have verified existence and uniqueness. Denote this interval by \mathbf{T}_j .

We want to prevent the solver from taking a very small last stepsize. To accomplish this, we adapt the idea for taking a last step from [13]. We consider three cases:

1. $\mathbf{t}_{\text{end}} \subseteq \mathbf{T}_j$
2. $\mathbf{t}_{\text{end}} \cap \mathbf{T}_j = \emptyset$
3. $t_{j+1} \in \mathbf{t}_{\text{end}}$

```

318  ⟨check if last step 318⟩ ≡      /*
       $T_j \supseteq \mathbf{T}_j$ 
       $t_{\text{end}} = \mathbf{t}_{\text{end}}$ 
      */
       $T_j = \text{hoe\_getTrialT}()$ ;
      if ( $\mathbf{v\_bias} :: \text{subsetq}(t_{\text{end}}, T_j)$ ) {
        ⟨case 1 320⟩
      }
      else {
        if ( $\mathbf{v\_bias} :: \text{disjoint}(t_{\text{end}}, T_j)$ ) {
          ⟨case 2 321⟩
        }
        else {
          interval tmp;
          if ( $\text{direction} \equiv 1$ )  $\text{tmp} = \mathbf{v\_bias} :: \text{sup}(T_j)$ ;
          else  $\text{tmp} = \text{inf}(T_j)$ ; /*  $\text{tmp} \ni \mathbf{t}_{j+1}$  */
          if ( $\mathbf{v\_bias} :: \text{subsetq}(\text{tmp}, t_{\text{end}})$ ) {
            ⟨case 3 323⟩
          }
          else assert(0);
        }
      }
}

```

This code is used in chunk 296.

Case 1: $\mathbf{t}_{\text{end}} \subseteq \mathbf{T}_j$. We set $\mathbf{t}_{j+1} = \mathbf{t}_{\text{end}}$; Figure 21.1.

```

320  ⟨case 1 320⟩ ≡
       $t_{\text{trial}} = t_{\text{end}}$ ; /*  $t_{\text{end}}$  is inside  $T_j$  */

```

This code is used in chunk 318.

Case 2: $\mathbf{t}_{\text{end}} \cap \mathbf{T}_j = \emptyset$. Denote

$$\Theta = |\mathbf{t}_{\text{end}} - t_j| = \max\{|\underline{\mathbf{t}}_{\text{end}} - t_j|, |\bar{\mathbf{t}}_{\text{end}} - t_j|\}.$$

- (i) If $w(\mathbf{T}_j) < \Theta/2$, then the next point we chose is the already selected t_{j+1} .
- (ii) Otherwise, $w(\mathbf{T}_j) \geq \Theta/2$; see Figure 21.2.

We compute

$$\mathbf{T}_j^* = \begin{cases} \mathbf{T}_j \cap (t_j + [0, \Theta/2]) & \text{if } t_j < t_{j+1} \\ \mathbf{T}_j \cap (t_j - [0, \Theta/2]) & \text{if } t_j > t_{j+1} \end{cases} \quad (21.1)$$

and select for the next step

$$t_{j+1} \leftarrow \begin{cases} \bar{\mathbf{T}}_{j+1}^* & \text{if } t_j < t_{j+1} \\ \underline{\mathbf{T}}_{j+1}^* & \text{if } t_j > t_{j+1} \end{cases} \quad (21.2)$$

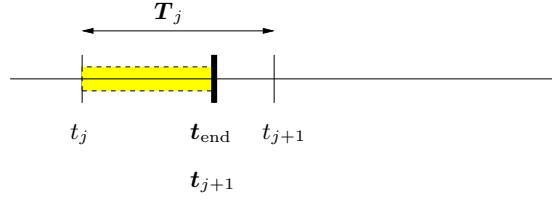


Figure 21.1. The case $t_{\text{end}} \subseteq T_j$. We set $t_{j+1} = t_{\text{end}}$.

```

321  ⟨ case 2 321 ⟩ ≡      /*
      theta ≥ Θ = |t_end - t_j|
      d ≥ w(T_j)

      */
      double theta = v_bias::mag(t_end - t0);
      if (v_bias::width(Tj) ≥ theta/2) {
        double d = theta/2.0;    /* d = Θ/2 */
        interval tmp;
        if (direction ≡ 1) tmp = interval(0, d);
        else tmp = interval(-d, 0);
        bool b;
        b = v_bias::intersect(Tj, Tj, t0 + tmp);
        assert(b);
      }
      if (direction ≡ 1) t_trial = v_bias::sup(Tj);
      else t_trial = v_bias::inf(Tj);

```

This code is used in chunk 318.

Case 3: $t_{j+1} \in t_{\text{end}}$. We expect this situation to occur very rarely. We set

$$t_{j+1} \leftarrow \begin{cases} \underline{t}_{\text{end}} & \text{if } t_j < t_{j+1} \\ \bar{t}_{\text{end}} & \text{if } t_j > t_{j+1}. \end{cases}$$

```

323  ⟨ case 3 323 ⟩ ≡
      if (direction ≡ 1) t_trial = v_bias::inf(t_end);
      else t_trial = v_bias::sup(t_end);

```

This code is used in chunk 318.

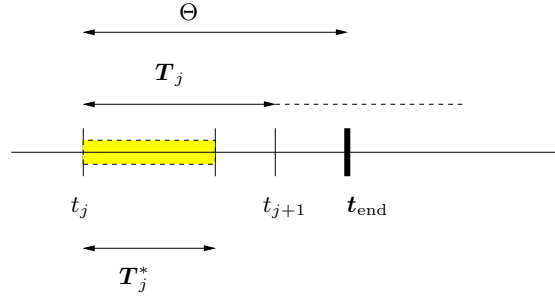


Figure 21.2. When close to t_{end} , we take the “middle” as the next integration point.

21.2.7 Compute a tight enclosure

To compute a tight enclosure at t_{j+1} , we call

```
324 <compute enclosure 324> ≡
    iho→compTightEnclosure(t_trial);
```

This code is used in chunk 296.

21.2.8 Decide

We decide how to proceed as shown below.

1. If *control-ind* \equiv *first_entry* we set *control-ind* \equiv *success*.
2. If *control-ind* \equiv *success*, we consider two cases.
 - (a) If *control-interrupt* \equiv *no*, we continue the integration.
 - (b) If *control-interrupt* \equiv *before_accept*, we return from the integration.
3. If *control-ind* \equiv *failure*, we return from the integration.

```
325 <decide 325> ≡
    bool ret = false;
    switch (control-ind) {
    case first_entry:
        control-ind = success;
    case success:
        if (control-interrupt ≡ before_accept) ret = true;
        break;
    case failure: return;
    }
    acceptSolution(t0, y0);
    if (ret) return;
```

This code is used in chunk 296.

Accept solution

We accept the solution computed in the *hoe* and *iho* objects. We also update *t0* and *y0*.

```
326 <VNODE accept solution 326> ≡
    void VNODE::acceptSolution(interval &t0, iVector &y0)
    {
        hoe->acceptSolution();
        iho->acceptSolution();
        h_accepted = t_trial - t0;
        t0 = t_trial;
        iho->getTightEnclosure(y0);
        steps++;
    }
```

This code is used in chunk 334.

21.3 Constructor/destructor

```
327 <constructor (VNODE) 327> ≡
    VNODE::VNODE(AD *a)
    {
        steps = 0;
        direction = 0;
        int n = a->size;
        sizeV(temp, n);
        sizeV(tp, n);
        hoe = new HOE(n);
        iho = new IHO(n);
        control = new Control;
        assert(a);
        ad = a;
    }
```

This code is used in chunk 334.

```
328 <destructor (VNODE) 328> ≡
    VNODE::~~VNODE()
    {
        delete control;
        delete iho;
        delete hoe;
    }
```

This code is used in chunk 334.

21.4 Get functions

To check if an integration is successful, we call

```

329 <get functions (VNODE) 329> ≡
    unsigned int getMaxOrder() const
    {
        return ad→getMaxOrder();
    }
    bool successful() const
    {
        if (control→ind ≡ success ∨ control→ind ≡ first_entry) return true;
        return false;
    }

```

See also chunk 330.

This code is used in chunk 295.

```

330 <get functions (VNODE) 329> +≡
    double getStepsize() const { /* Returns the last stepsize taken. */
        return midpoint(h_accepted);
    }
    const iVector &getAprioriEncl() const {
        /* Obtains a reference to the a priori enclosure */
        return hoe→getApriori();
    }
    const interval &getT() const { /* Returns  $T_j$ . */
        return hoe→getT();
    }
    double getGlobalExcess() { /* Returns an estimate of the global excess. */
        width(tp, iho→getGlobalExcess());
        return inf_norm V(tp);
    }
    double getGlobalExcess(int i) {
        /* Returns an estimate of the global excess in component  $i$ . */
        width(tp, iho→getGlobalExcess());
        return tp[i];
    }
    int getNoSteps() const {
        /* Returns the number of accepted steps during an integration. */
        return steps;
    }

```

21.5 Set parameters

We set integration parameters by the following functions.

```

331 <set VNODE parameters 331> ≡

```

```

void setTols(double a,double r = 0) {
    /* Sets atol and rtol. By default, rtol is set to 0 */
    control-atol = a;
    control-rtol = r;
}

void setOrder(int p) {    /* Sets order. */
    control-order = p;
}

void setHmin(double h) {    /* Sets a value for the minimum stepsize. */
    control-hmin = h;
}

void setOneStep(stepAction action) {
    /* Indicates an interrupt after each computed solution. */
    if (action  $\equiv$  on) control-interrupt = before_accept;
    else    /* off */
        control-interrupt = no;
}

void setFirstEntry() {    /* Indicates first entry */
    control-ind = first_entry;
}

```

This code is used in chunk 295.

21.6 Files

21.6.1 Interface

```

333 <vnoint.h 333>  $\equiv$ 
    #ifndef VNODEINT_H
    #define VNODEINT_H
    #include <algorithm>
    #include "miscfuncs.h"
    #include "vnointerval.h"
    #include "vnoderound.h"
    #include "vector_matrix.h"
    #include "ad_ode.h"
    #include "ad_var.h"
    #include "allad.h"
    #include "solution.h"
    #include "control.h"
    #include "matrixinverse.h"
    #include "iho.h"
    #include "hoe.h"
    #include "vtiming.h"
    #include "debug.h"
    namespace vnodelp {

```



```

    < class VNODE 295 >
}
#endif

```

21.6.2 Implementation

```

334 < integ.cc 334 > ≡
#include <cmath>
#include <algorithm>
#include <ostream>
using namespace std;
#include "vnodeinterval.h"
#include "vector_matrix.h"
#include "matrixinverse.h"
#include "basiclinalg.h"
#include "vnodeint.h"
namespace vnodelp {
    < constructor (VNODE) 327 >
    < destructor (VNODE) 328 >
    < compute  $h_{\min}$  314 >
    < compute initial stepsize 315 >
    < integrator 296 >
    < VNODE accept solution 326 >
}

```

21.7 Interface to the VNODE-LP Package

The interface to the VNODE-LP package is stored in

```

335 < vnode.h 335 > ≡
#ifndef VNODE_H
#define VNODE_H
#include <algorithm>
#include "vnodeinterval.h"
#include "vnodearound.h"
#include "vector_matrix.h"
#include "ad_ode.h"
#include "ad_var.h"
#include "allad.h"
#include "vnodeint.h"
#include "solution.h"
#include "control.h"
#include "iho.h"
#include "hoe.h"
#include "matrixinverse.h"
#include "basiclinalg.h"

```

```
#include "fadb_ad.h"  
#include "fadb_advar.h"  
#include "fadb_ad.h"  
#endif
```

Part V

AD Implementation

Chapter 22

Using FADBAD++

Currently, VNODE employs the FADBAD++ package [29].

22.1 Computing ODE Taylor coefficients

22.1.1 FadbadODE class

```
339 <class FadbadODE 339> ≡
    typedef T<interval> Tinterval;
    typedef void(*Tfunction)(int n, Tinterval *yp, const Tinterval
        *y, Tinterval t, void *param);
    class FadbadODE : public AD_ODE {
    public:
        FadbadODE(int n, Tfunction, void *param = 0);
        void set(const interval &t0, const iVector &y0, const interval &h, int
            k);
        void compTerms();
        void sumTerms(iVector &sum, int m);
        void getTerm(iVector &term, int i) const;
        interval getStepsize() const;
        void eval(void *param);
        ~FadbadODE();
    private:
        Tfunction fcn;
        Tinterval *y_coeff, *f_coeff, t;
        int size;
        int order;
        interval stepsize;
    };

```

This code is used in chunk 349.

22.1.2 Function description

In the constructor, we allocate the necessary memory, set the ODE problem, and generate the computational graph by calling *fcn*.

```
340 < constructor (FadbadODE) 340 > ≡
    FadbadODE::FadbadODE(int n, Tfunction f, void *param)
    : AD_ODE() {
        size = n;
        y_coeff = new Tinterval[2 * n];
        f_coeff = y_coeff + n;
        fcn = f;
        fcn(size, f_coeff, y_coeff, t, param);
    }
```

This code is used in chunk 350.

```
341 < FadbadODE destructor 341 > ≡
    FadbadODE::~FadbadODE()
    {
        delete y_coeff;
    }
```

This code is used in chunk 350.

```
342 < initialize Taylor coefficients (FadbadODE) 342 > ≡
    void FadbadODE::set(const interval &t0, const iVector &y0, const
        interval &h, int k)
    {
        t[0] = t0;
        t[1] = h;
        stepsize = h;
        order = k;
        for (int eqn = 0; eqn < size; eqn++) y_coeff[eqn][0] = y0[eqn];
    }
```

This code is used in chunk 350.

The Taylor coefficients for the right side of $y' = f(t, y)$ are evaluated by *eval* and stored in *f_coeff*. Then, the *i*th coefficient for *y* is $h/i \times$ (the $(i - 1)$ st coefficient for *f*).

```
343 < compute Taylor coefficients (FadbadODE) 343 > ≡
    void FadbadODE::compTerms()
    {
        for (int eqn = 0; eqn < size; eqn++)
            f_coeff[eqn].reset(); /* reset previously created terms */
        for (int coeff = 1; coeff ≤ order; coeff++) /* compute coefficients */
            for (int eqn = 0; eqn < size; eqn++) {
                f_coeff[eqn].eval(coeff - 1);
            }
```

```

        y_coeff[eqn][coeff] = stepsize * f_coeff[eqn][coeff - 1]/double(coeff);
    }
}

```

This code is used in chunk 350.

```

344  <sum terms (FadbadODE) 344> ≡
    void FadbadODE::sumTerms(iVector &sum, int m)
    {
        interval s;
        for (int eqn = 0; eqn < size; eqn++) {
            s = 0.0;
            for (int coeff = m; coeff ≥ 0; coeff--) s += y_coeff[eqn][coeff];
            sum[eqn] = s;
        }
    }

```

This code is used in chunk 350.

```

345  <get term (FadbadODE) 345> ≡
    void FadbadODE::getTerm(iVector &term, int i) const
    {
        for (int eqn = 0; eqn < size; eqn++) term[eqn] = y_coeff[eqn][i];
    }

```

This code is used in chunk 350.

```

346  <obtain stepsize 346> ≡
    interval FadbadODE::getStepsize() const
    {
        return stepsize;
    }

```

This code is used in chunk 350.

To evaluate the function and rebuild the computational graph, we do

```

347  <rebuild computational graph 347> ≡
    void FadbadODE::eval(void *param)
    {
        fcn(size, f_coeff, y_coeff, t, param);
    }

```

This code is used in chunk 350.

22.1.3 Files

```

349  <fadbad_ad.h 349> ≡
    #ifndef Fadbad_ODE
    #define Fadbad_ODE

```

```

#include "vnodeinterval.h"
#include "basicclinalg.h"
#include "ad_ode.h"
#include "ffadiff.h"
#include "fadbad_intv.inc"
    namespace vnodelp {
        < class FadbadODE 339 >
    }
#endif

350 < fadbad_ad.cc 350 > ≡
#include "fadbad_ad.h"
    namespace vnodelp {
        < constructor (FadbadODE) 340 >
        < FadbadODE destructor 341 >
        < initialize Taylor coefficients (FadbadODE) 342 >
        < compute Taylor coefficients (FadbadODE) 343 >
        < sum terms (FadbadODE) 344 >
        < get term (FadbadODE) 345 >
        < obtain stepsize 346 >
        < rebuild computational graph 347 >
    }

```

22.2 Computing Taylor coefficients for the variational equation

22.2.1 FadbadVarODE class

```

352 < class FadbadVarODE 352 > ≡
    typedef T<F<interval>> TFinterval;
    typedef void(*TFfunction)(int n, TFinterval *yp, const TFinterval
        *y, TFinterval tf, void *param);
    class FadbadVarODE : public AD_VAR {
    public:
        FadbadVarODE(int n, TFfunction f, void *param = 0);
        void set(const interval &t0, const iVector &y0, const interval &h, int
            k);
        void compTerms();
        void sumTerms(iMatrix &sum, int m);
        void getTerm(iMatrix &term, int i) const;
        void eval(void *param) { fcn(size, tf_out, tf_in, tf, param); }
        ~FadbadVarODE();
    private:
        TFinterval *tf_in, *tf_out, tf;

```



```

    TFfunction fcn;
    int size;
    int order;
    interval stepsize;
};

```

This code is used in chunk 360.

22.2.2 Function description

```

353 < constructor (FadbadVarODE) 353 > ≡
    FadbadVarODE :: FadbadVarODE(int n, TFfunction f, void *param)
    {
        size = n;
        tf_in = new TFinterval[2 * n];
        tf_out = tf_in + n;
        fcn = f;
        fcn(size, tf_out, tf_in, tf, param);
    }

```

This code is used in chunk 361.

```

354 < destructor (FadbadVarODE) 354 > ≡
    FadbadVarODE :: ~FadbadVarODE()
    {
        delete[] tf_in;
    }

```

This code is used in chunk 361.

```

355 < initialize Taylor coefficients (FadbadVarODE) 355 > ≡
    void FadbadVarODE :: set(const interval &t0, const iVector &y0, const
        interval &h, int k)
    {
        stepsize = h;
        order = k;
        tf[0].x() = t0;
        tf[1].x() = h;
        for (int eqn = 0; eqn < size; eqn++) {
            tf_in[eqn][0] = y0[eqn];
        }
    }

```

This code is used in chunk 361.

```

356 < compute Taylor coefficients (FadbadVarODE) 356 > ≡
    void FadbadVarODE :: compTerms()
    {

```

```

    for (int eqn = 0; eqn < size; eqn++) tf_out[eqn].reset();
    for (int eqn = 0; eqn < size; eqn++) tf_in[eqn][0].diff(eqn, size);
    for (int coeff = 0; coeff < order; coeff++) {
        for (int eqn = 0; eqn < size; eqn++) {
            tf_out[eqn].eval(coeff);
            tf_in[eqn][coeff + 1] = stepsize * (tf_out[eqn][coeff]/double(coeff + 1));
        }
    }
}

```

This code is used in chunk 361.

```

357 <sum terms (FadbadVarODE) 357> ≡
    void FadbadVarODE::sumTerms(iMatrix &Sum, int k)
    {
        for (int row = 0; row < size; row++)
            for (int col = 0; col < size; col++) {
                interval s = 0.0;
                for (int coeff = k; coeff ≥ 1; coeff--) s += tf_in[row][coeff].d(col);
                Sum[row][col] = s;
            }
        for (int row = 0; row < size; row++) Sum[row][row] += 1.0;
    }

```

This code is used in chunk 361.

```

358 <get term (FadbadVarODE) 358> ≡
    void FadbadVarODE::getTerm(iMatrix &Term, int i) const
    {
        for (int row = 0; row < size; row++)
            for (int col = 0; col < size; col++) Term[row][col] = tf_in[row][i].d(col);
    }

```

This code is used in chunk 361.

22.3 Files

```

360 <fadbad_advar.h 360> ≡
    #ifndef Fadbad_Var_ODE
    #define Fadbad_Var_ODE
    #include "vnodeinterval.h"
    #include "vector_matrix.h"
    #include "ad_var.h"
    #include "ffadiff.h"
    #include "fadbad_intv.inc"
    namespace vnodelp {
        using namespace v_bias;
        using namespace v_blas;
    }

```

```

    < class FadbadVarODE 352 >
  }
#endif

361 < fadbad_advar.cc 361 > ≡
#include "fadbad_advar.h"
namespace vnodelp {
  < constructor (FadbadVarODE) 353 >
  < destructor (FadbadVarODE) 354 >
  < initialize Taylor coefficients (FadbadVarODE) 355 >
  < compute Taylor coefficients (FadbadVarODE) 356 >
  < sum terms (FadbadVarODE) 357 >
  < get term (FadbadVarODE) 358 >
}

```

22.4 Encapsulated FADBAD++ AD

Now we encapsulate all AD using FADBAD++ in

```

362 < encapsulated FADBAD++ AD 362 > ≡
class FADBAD_AD : public AD {

public:
  FADBAD_AD(int n, Tfunction f, TFfunction tf)
  : AD(n, new FadbadODE(n, f), new FadbadVarODE(n, tf)),
    max_order(MaxLength - 2) {}

  FADBAD_AD(int n, Tfunction f, TFfunction tf, void *p)

  : AD(n, new FadbadODE(n, f, p), new FadbadVarODE(n, tf, p)),
    max_order(MaxLength - 2) {}

  virtual int getMaxOrder() const {
    return max_order;
  }

private: const int max_order;
};

```

This code is used in chunk 363.

and store the new class in

```

363 < fadbadad.h 363 > ≡
#ifndef FADBADAD_H
#define FADBADAD_H
#include "fadbad_ad.h"
#include "fadbad_advar.h"
#include "allad.h"
namespace vnodelp {

```

```
    { encapsulated FADBAD++ AD 362 }  
#endif
```

Appendix A

Miscellaneous Functions

A.1 Vector output

We output a vector to the standard output by

```
366 <print vector 366> ≡  
#include <iostream>  
using namespace std;  
template<class T> void printVector(const T &v, const char *s = 0)  
{  
    if (s) cout << s << "␣="␣ << endl;  
    for (unsigned int i = 0; i < v.blas::sizeV(v); i++) cout << v[i] << endl;  
    cout << endl;  
}
```

This code is used in chunk 137.

A.2 Check if an interval is finite

```
367 <check finite 367> ≡  
namespace v_blas {  
    inline bool finite_interval(const interval &a)  
    {  
        return (isfinite(inf(a)) & isfinite(sup(a)));  
    }  
}
```

See also chunk 368.

This code is used in chunk 374.

We check if an interval vector contains finite intervals by

```
368 <check finite 367> +≡  
namespace v_blas {
```

```

inline bool finite_interval(const iVector &a)
{
    for (unsigned int i = 0; i < a.size(); i++)
        if (!v_bias::finite_interval(a[i])) return false;
    return true;
}

```

A.3 Message printing

We would like to print various messages.

```

369 < message printing 369 > ≡
    #ifndef VNODE_DEBUG
    #define printMessage(s)
    {
        cerr << "\n***" << __FILE__ ":" << __LINE__ << " " << s << endl;
    }
    #else
    #define printMessage(s) (0)
    #endif
    #ifndef VNODE_DEBUG
    #define exitOnError(s)
    {
        printMessage(s);
        exit(-1);
    }
    #else
    #define exitOnError(s) (0)
    #endif

```

This code is used in chunk 372.

```

370 < VNODE-LP message 370 > ≡
    void vnodeMessage(const char *s)
    {
        cerr << endl << "***VNODE-LP:" << s << endl;
    }

```

This code is used in chunk 373.

A.4 Check intersection

```

371 < check vector intersection 371 > ≡
    void checkIntersection(const iVector &a, const iVector &b)
    {
        v_bias::interval ai, bi;

```

```

    for (unsigned int i = 0; i < vnodelp::sizeV(a); i++) {
        ai = a[i];
        bi = b[i];
        if (vnodelp::disjoint(ai, bi))
            cout << "i_=" << i << " _=" << endl << ai << endl << bi;
    }
}

```

This code is used in chunk 373.

Files

```

372 <debug.h 372> ≡
    #ifndef DEBUG_H
    #define DEBUG_H
    #include "basiclinalg.h"
    <message printing 369>
    using namespace v_blas;
    void checkIntersection(const iVector &a, const iVector &b);
    #endif

373 <debug.cc 373> ≡
    #include <ostream>
    #include <stdarg.h>
    #include <stdio.h>
    #include <stdlib.h>
    #include <iostream>
    using namespace std;
    #include "vnodeinterval.h"
    #include "basiclinalg.h"
    using namespace std; namespace vnodelp {
        using namespace v_bias;
        using namespace v_blas;
        <check vector intersection 371>
        <VNODE-LP message 370>
    }

374 <miscfuncs.h 374> ≡
    #ifndef MISCFUNS_H
    #define MISCFUNS_H
    #include <cmath>
    #include "basiclinalg.h"
    using namespace std; <check finite 367> namespace vnodelp {
        void vnodeMessage(const char *s);
    }
    #endif

```

A.5 Timing

The *getTime* function returns the current user time. The *getTotalTime* subtracts the end time from the start time and returns the result.

```

375 <vtiming.h 375> ≡
    #ifndef VTIMING_H
    #define VTIMING_H
        double getTime();
        double getTotalTime(double start_time, double end_time);
    #endif

376 <vtiming.cc 376> ≡
    #include <cassert>
    #include <sys/times.h>
    #include <unistd.h>
    #include <ctime>
    #include "vnodeinterval.h"
    #include "vnoderound.h"
    using namespace std;
    static struct tms Tms;
    double getTime()
    {
        times(&Tms);
        v_bias::round_nearest();
        long int ClockTcks = sysconf(_SC_CLK_TCK);
        return (Tms.tms_utime)/(double(ClockTcks));
    }
    double getTotalTime(double start_time, double end_time)
    {
        assert(start_time ≤ end_time);
        v_bias::round_nearest();
        return (end_time - start_time);
    }

```


Bibliography

- [1] *BLAS — Basic Linear Algebra Subprograms*. www.netlib.org/blas/.
- [2] *LAPACK — Linear Algebra PACKage*. www.netlib.org/lapack/.
- [3] D. ACHLIOPTAS, *Setting 2 variables at a time yields a new lower bound for random 3-SAT*, Tech. Rep. MSR-TR-99-96, Microsoft Research, Microsoft Corp., One Microsoft Way, Redmond, WA 98052, December 1999.
- [4] U. M. ASCHER AND L. R. PETZOLD, *Computer Methods for Ordinary Differential Equations and Differential-Algebraic Equations*, SIAM, Philadelphia, 1998.
- [5] E. AUER, A. KECSKEMÉTHY, M. TÄNDL, AND H. TRACZINSKI, *Interval algorithms in modelling of multibody systems*, in Numerical Software with Result Verification, vol. 2991 of Lecture Notes in Computer Science in Computer Science, Springer-Verlag, 2004, pp. 132–159.
- [6] C. BENDSTEN AND O. STAUNING, *FADBAD, a flexible C++ package for automatic differentiation using the forward and backward methods*, Tech. Rep. 1996-x5-94, Department of Mathematical Modelling, Technical University of Denmark, DK-2800, Lyngby, Denmark, August 1996.
- [7] M. BERZ, K. MAKINO, AND J. HOEFKENS, *Verified integration of dynamics in the solar system*, Nonlinear Analysis: Theory, Methods & Applications, 47 (2001).
- [8] B. M. BROWN, M. LANGER, M. MARLETTA, C. TRETTER, AND M. WAGENHOFER, *Eigenvalue bounds for the singular Sturm-Liouville problem with a complex potential*, J. Phys. A: Math. Gen., 36 (2003), pp. 3773–3787.
- [9] G. F. CORLISS AND R. RIHM, *Validating an a priori enclosure using high-order Taylor series*, in Scientific Computing, Computer Arithmetic, and Validated Numerics, G. Alefeld and A. Frommer, eds., Akademie Verlag, Berlin, 1996, pp. 228–238.
- [10] G. I. HARGREAVES, *Interval analysis in MATLAB*, Tech. Rep. No. 416, Department of Mathematics, University of Manchester, UK, December 2002.

- [11] W. HAYES, *Rigorous shadowing of numerical solutions of ordinary differential equations by containment*, PhD thesis, Department of Computer Science, University of Toronto, Toronto, Canada, 2001.
- [12] W. HAYES AND K. R. JACKSON, *Rigorous shadowing of numerical solutions of ordinary differential equations by containment*, SIAM J. Num. Anal., to appear (2005).
- [13] T. E. HULL AND W. H. ENRIGHT, *A structure for programs that solve ordinary differential equations*, Tech. Rep. 66, Department of Computer Science, University of Toronto, May 1974.
- [14] T. E. HULL, W. H. ENRIGHT, B. M. FELLEN, AND A. E. SEDGWICK, *Comparing numerical methods for ordinary differential equations*, SIAM J. on Numerical Analysis, 9 (1972), pp. 603–637.
- [15] M. KIEFFER AND E. WALTER, *Nonlinear parameter and state estimation for cooperative systems in a bounded-error context*, in Numerical Software with Result Verification, vol. 2991 of Lecture Notes in Computer Science, Springer-Verlag, 2004, pp. 107–123.
- [16] O. KNÜPPEL, *PROFIL/BIAS – a fast interval library*, Computing, 53 (1994), pp. 277–287.
- [17] D. E. KNUTH, *Literate Programming*, Center for the Study of Language and Information, Stanford, CA, USA, 1992.
- [18] D. E. KNUTH AND S. LEVY, *The CWEB System of Structured Documentation*, Addison-Wesley, Reading, Massachusetts, 1993.
- [19] M. LERCH, G. TISCHLER, AND J. WOLFF VON GUDENBERG, *FILIB++ – interval library specification and reference manual*, Tech. Rep. 279, Universität Würzburg, Germany, 2001.
- [20] R. J. LOHNER, *Einschließung der Lösung gewöhnlicher Anfangs- und Randwertaufgaben und Anwendungen*, PhD thesis, Universität Karlsruhe, 1988.
- [21] F. MAZZIA AND F. IAVERNARO, *Test set for initial value problem solvers*, Tech. Rep. 40, Department of Mathematics, University of Bari, Italy, 2003. <http://pitagora.dm.uniba.it/~testset/>.
- [22] H. MUKUNDAN, K. H. KO, T. MAEKAWA, T. SAKKALIS, AND N. M. PATRIKALAKIS, *Tracing surface intersections with a validated ODE system solver*, in Proceedings of the Ninth EG/ACM Symposium on Solid Modeling and Applications, G. Elber and G. Taubin, eds., Eurographics Press, June 2004, June 2004.
- [23] N. S. NEDIALKOV, *Computing Rigorous Bounds on the Solution of an Initial Value Problem for an Ordinary Differential Equation*, PhD thesis, Department of Computer Science, University of Toronto, Toronto, Canada, M5S 3G4, February 1999.

- [24] N. S. NEDIALKOV, *An interval method for initial value problems in linear ordinary differential equations*, SIAM Journal Numerical Analysis, (2004). Submitted.
- [25] N. S. NEDIALKOV AND K. R. JACKSON, *The design and implementation of a validated object-oriented solver for IVPs for ODEs*, Tech. Rep. 6, Software Quality Research Laboratory, Department of Computing and Software, McMaster University, Hamilton, Canada, L8S 4L7, 2002.
- [26] N. S. NEDIALKOV, K. R. JACKSON, AND G. F. CORLISS, *Validated solutions of initial value problems for ordinary differential equations*, Applied Mathematics and Computation, 105 (1999), pp. 21–68.
- [27] N. S. NEDIALKOV, K. R. JACKSON, AND J. D. PRYCE, *An effective high-order interval method for validating existence and uniqueness of the solution of an IVP for an ODE*, Reliable Computing, 7 (2001), pp. 449–465.
- [28] N. M. PATRIKALAKIS, T. MAEKAWA, K. H. KO, AND H. MUKUNDAN, *Surface to surface intersection*, in International CAD Conference and Exhibition, CAD'04, L. Piegl, ed., Thailand, May 2004.
- [29] O. STAUNING AND C. BENDTSEN, *FADBAD++ web page*, May 2003. FADBAD++ is available at www.imm.dtu.dk/fadbad.html.
- [30] W. TUCKER, *A rigorous ODE solver and Smale's 14th problem*, Found. Comput. Math., 2 (2002), pp. 53–117.

Index

__FILE__: 369.
 __LINE__: 369.
 _SC_CLK_TCK: 376.
 A: 123, 124, 130, 131, 132, 146,
 153, 156, 158, 160, 166, 170,
 178, 235, 284, 293.
 a: 79, 84, 113, 114, 123, 124, 129,
 132, 139, 140, 141, 142, 144, 145,
 149, 150, 151, 161, 191, 192, 288,
 327, 331, 367, 368, 371, 372.
 A_point: 235, 273, 274.
 Abs: 114.
 accept: 325.
 acceptSolution: 204, 220, 235, 276,
 295, 325, 326.
 accumulate: 164.
 acos: 79, 113, 114.
 action: 84, 331.
 ad: 22, 23, 34, 39, 54, 57, 71, 73,
 74, 82, 86, 221, 235, 241, 243,
 244, 248, 249, 257, 258, 277, 295,
 311, 312, 315, 327, 329.
 AD: 22, 34, 39, 54, 71, 73, 74, 82,
 86, 175, 191, 192, 193, 221, 235,
 277, 295, 327, 362.
 AD_ODE: 175, 186, 187, 191, 192,
 193, 204, 339, 340.
 AD_ODE_H: 194.
 AD_VAR: 175, 186, 189, 191,
 192, 193, 352.
 AD_VAR_H: 195.
 addId: 132, 248, 249.
 addMiMi: 132, 248, 249.
 addViVi: 129, 132, 162, 206, 207,
 245, 246, 257, 258, 260, 261,
 263, 268, 270, 271.
 addViVp: 129, 263, 268.
 ad0: 277.
 ai: 371.
 Ainv: 156, 157, 158, 166, 169, 170,
 235, 270, 271, 274.
 ALLAD_H: 196.
 alpha: 178, 179, 180, 245, 263,
 268, 276.
 apriori: 204, 219, 220, 221, 222.
 Apriori: 175, 181, 182, 183, 204,
 219.
 apriori_trial: 204, 207, 210, 215,
 217, 219, 220, 222.
 ArcCos: 114.
 ArcSin: 114.
 ArcTan: 114.
 asin: 79, 113, 114.
 assert: 166, 207, 212, 215, 219, 239,
 245, 246, 262, 263, 273, 277, 309,
 318, 321, 327, 376.
 assignM: 132, 167, 272, 276.
 assignV: 129, 132, 220, 243, 245,
 258, 272, 276.
 atan: 79, 113, 114.
 atol: 200, 207, 301, 315, 316, 331.
 av: 191, 192.
 Ax: 157.
 a1: 161, 288.
 B: 132, 156, 160, 166, 235, 284, 293.
 b: 79, 113, 114, 129, 139, 140, 141,
 147, 149, 150, 151, 163, 166, 169,
 245, 246, 286, 288, 321, 371, 372.
 BASICLINALG_H: 137.
 before_accept: 199, 325, 331.

- begin*: 129, 139, 140, 141, 143,
 144, 145, 164.
beta: 19, 51, 52, 53, 57, 156, 160,
 161, 165, 166, 168, 169, 212,
 216, 224.
bi: 371.
BiasRoundDown: 119.
BiasRoundNear: 119.
BiasRoundUp: 119.
b0: 156, 160, 161, 162, 169, 171.
b1: 288.
C: 132, 133, 134, 135, 136, 156, 235.
c: 79, 113, 114, 147, 277.
C-pq: 235, 239, 248, 257, 274,
 275, 279.
C-qp: 235, 239, 249, 258, 274,
 275, 280.
c_str: 68, 74.
ceil: 239.
cerr: 113, 369, 370.
checkIntersection: 371, 372.
Ci: 156, 167, 171.
Cinv: 235, 251, 261, 274.
ClockTcks: 376.
close: 44, 55, 57, 63, 68, 73, 74.
coeff: 343, 344, 356, 357.
col: 357, 358.
comp_beta: 204, 212, 224.
compAprioriEnclosure: 204, 217,
 317.
compare: 288.
compCoeffs: 235, 239, 240, 312.
compCpq: 235, 239, 279.
compCqp: 235, 239, 280.
compErrorConst: 282.
compErrorConstant: 235, 239, 282.
compH: 149, 150, 151, 224.
compHmin: 309, 314.
compHstart: 295, 310, 315.
compTerm: 315.
compTerms: 187, 188, 189, 190,
 206, 210, 243, 244, 249, 258, 315,
339, 343, 352, 356.
compTightEnclosure: 235, 236, 324.
computeQR: 153, 273, 293.
contol: 309.
Control: 175, 197, 200, 201, 204,
 221, 235, 277, 295, 327.
control: 204, 207, 214, 217, 219,
 221, 235, 239, 245, 251, 271,
 273, 277, 295, 296, 298, 299, 300,
 301, 302, 303, 304, 305, 307, 309,
 310, 311, 312, 315, 316, 317, 325,
 327, 328, 329, 331.
CONTROL_H: 201.
corrector_excess: 235, 241, 255,
 260, 274.
cos: 79, 113, 114.
Cos: 114.
counter: 165.
cout: 25, 26, 68, 71, 74, 366, 371.
ctrl: 221.
d: 235, 321.
DEBUG_H: 372.
DETEST_C3: 71.
DETEST_E1: 37, 39.
DGEQRF: 153.
dgeqrf-.: 153.
dgetrf: 158.
dgetrf-.: 158.
dgetri: 158.
dgetri-.: 158.
Diam: 114.
diam: 113.
diff: 356.
direction: 295, 305, 306, 307, 310,
 318, 321, 323, 327.
disjoint: 79, 113, 114, 141, 304,
 318, 371.
DORGQR: 153.
dorgqr-.: 153.
dot_product: 129, 131.
downward: 118.
D1: 73.
D2: 73.
e: 223.
encloseInverse: 251.
encloseLS: 156, 159, 160, 166, 169.
encloseLSS: 157.
encloseMatrixInverse: 156, 157,
166, 251, 270.
Enclosure: 114.

- end*: 129, 139, 140, 141, 143, 144, 145, 164.
end_time: 87, 375, 376.
endl: 25, 26, 44, 55, 57, 64, 65, 68, 71, 73, 74, 366, 369, 370, 371.
eps: 42.
eqn: 342, 343, 344, 345, 355, 356.
equal: 139, 140, 141.
err_const: 282.
errorConstant: 235, 239, 255.
eval: 57, 74, 86, 187, 188, 189, 190, 191, 193, 339, 343, 347, 352, 356.
exit: 369.
exitOnError: 369.
exp: 79, 113, 114.
Exp: 114.
f: 352, 353.
f_coeff: 339, 340, 343, 347.
fabs: 129, 214, 217.
FADBAD_AD: 22, 34, 39, 54, 71, 73, 74, 82, 86, 175, 362.
Fadbad_ODE: 349.
Fadbad_Var_ODE: 360.
FADBADAD_H: 363.
FadbadODE: 175, 339, 340, 341, 342, 343, 344, 345, 346, 347, 362.
FadbadVarODE: 175, 352, 353, 354, 355, 356, 357, 358, 362.
failure: 198, 245, 251, 271, 273, 298, 299, 302, 304, 307, 309, 310, 317, 325.
false: 79, 85, 113, 139, 140, 141, 147, 153, 157, 158, 163, 166, 168, 169, 217, 251, 271, 281, 287, 317, 325, 329, 368.
fcn: 339, 340, 347, 352, 353.
file_name: 68, 69, 74.
filib: 77, 113, 117.
FILIB_VNODE: 115, 120.
fill: 129.
finite_interval: 245, 299, 367, 368.
finite_temp: 245.
finite_temp2: 245.
first_entry: 198, 200, 296, 300, 305, 309, 310, 325, 329, 331.
fixed: 74.
Fj: 235, 248, 251, 274.
floor: 239.
function_name: 86.
G: 235.
gamma: 224.
getApriori: 222, 330.
getAprioriEncl: 62, 65, 85, 330.
getColumn: 133, 134, 169, 286, 289.
getErrorTerm: 223, 241.
getGlobalExcess: 44, 85, 277, 330.
getMaxOrder: 85, 191, 302, 329, 362.
getNoSteps: 74, 85, 330.
getStepsize: 62, 65, 73, 74, 85, 187, 188, 222, 241, 330, 339, 346.
getT: 62, 65, 85, 222, 330.
getTerm: 187, 188, 189, 190, 206, 210, 223, 248, 249, 257, 258, 315, 339, 345, 352, 358.
getTightEnclosure: 277, 326.
getTime: 68, 71, 74, 87, 375, 376.
getTotalTime: 68, 71, 74, 87, 375, 376.
getTrialApriori: 222, 246.
getTrialStepsize: 222.
getTrialT: 222, 318.
gj: 235, 257, 259, 260, 274.
globalExcess: 235, 263, 274, 277.
h: 84, 150, 187, 189, 204, 309, 331, 339, 342, 352, 355.
h_accepted: 295, 326, 330.
h_min: 295, 296, 309.
h_next: 204, 216, 220.
h_start: 295, 310, 311.
h_trial: 204, 206, 207, 210, 212, 214, 216, 217, 220, 221, 222, 235, 238, 241, 243, 244, 249, 258.
hmin: 150, 200, 214, 217, 296, 303, 309, 310, 331.
hoe: 235, 241, 246, 277, 295, 311, 312, 317, 318, 326, 327, 328, 330.
HOE: 175, 202, 204, 217, 219, 220, 224, 235, 277, 295, 327.
HOE_H: 226.
hoem: 277.
h0: 221.

- i*: [47](#), [68](#), [71](#), [74](#), [85](#), [124](#), [129](#), [131](#),
[132](#), [134](#), [135](#), [136](#), [146](#), [147](#),
[150](#), [166](#), [169](#), [187](#), [189](#), [206](#), [223](#),
[248](#), [249](#), [257](#), [258](#), [273](#), [279](#), [280](#),
[282](#), [287](#), [299](#), [330](#), [339](#), [345](#), [352](#),
[358](#), [366](#), [368](#), [371](#).
- i_gamma*: [224](#).
- I_PACKAGE*: [116](#).
- i_pw*: [224](#).
- iho*: [295](#), [312](#), [324](#), [326](#), [327](#), [328](#),
[330](#).
- IHO**: [175](#), [235](#), [236](#), [239](#), [274](#), [275](#),
[276](#), [279](#), [280](#), [282](#), [295](#), [327](#).
- IHO_H*: [292](#).
- iMatrix**: [122](#), [130](#), [131](#), [132](#), [146](#),
[156](#), [160](#), [166](#), [170](#), [189](#), [235](#),
[352](#), [357](#), [358](#).
- ind*: [200](#), [245](#), [251](#), [271](#), [273](#), [296](#),
[298](#), [299](#), [300](#), [302](#), [304](#), [305](#), [307](#),
[309](#), [310](#), [317](#), [325](#), [329](#), [331](#).
- Ind**: [198](#), [200](#), [201](#).
- index*: [285](#), [286](#), [289](#).
- index_norm**: [285](#), [286](#), [288](#), [289](#).
- inf*: [64](#), [65](#), [79](#), [113](#), [114](#), [149](#), [212](#),
[214](#), [215](#), [216](#), [306](#), [318](#), [321](#),
[323](#), [367](#).
- Inf*: [114](#).
- inf_normM*: [132](#), [168](#), [273](#).
- inf_normV*: [129](#), [161](#), [207](#), [315](#), [330](#).
- info*: [153](#), [158](#), [204](#), [217](#), [317](#).
- init*: [178](#), [180](#), [181](#), [183](#), [221](#), [277](#),
[311](#), [312](#).
- integrate*: [24](#), [44](#), [47](#), [55](#), [56](#), [57](#), [60](#),
[61](#), [71](#), [73](#), [74](#), [83](#), [84](#), [85](#), [197](#),
[198](#), [199](#), [295](#), [296](#), [305](#).
- interior*: [79](#), [113](#), [114](#), [140](#), [207](#).
- Interrupt**: [199](#), [200](#), [201](#).
- interrupt*: [200](#), [325](#), [331](#).
- intersect*: [79](#), [113](#), [114](#), [147](#), [163](#),
[245](#), [246](#), [262](#), [263](#), [321](#).
- Intersection*: [114](#).
- interval**: [19](#), [21](#), [33](#), [38](#), [42](#), [47](#), [51](#),
[52](#), [53](#), [56](#), [65](#), [68](#), [71](#), [73](#), [74](#), [76](#),
[77](#), [79](#), [80](#), [83](#), [85](#), [112](#), [113](#), [114](#),
[122](#), [129](#), [142](#), [147](#), [149](#), [161](#), [166](#),
[178](#), [180](#), [181](#), [183](#), [187](#), [189](#), [204](#),
[206](#), [207](#), [212](#), [214](#), [215](#), [217](#), [219](#),
[221](#), [222](#), [224](#), [235](#), [236](#), [239](#), [241](#),
[277](#), [282](#), [295](#), [296](#), [314](#), [315](#), [318](#),
[321](#), [326](#), [330](#), [339](#), [342](#), [344](#), [346](#),
[352](#), [355](#), [357](#), [367](#), [371](#).
- INTERVAL**: [76](#).
- INTVFUNC_H*: [151](#).
- invertMatrix*: [156](#), [157](#), [158](#), [166](#).
- ios*: [44](#), [55](#), [57](#), [62](#), [68](#), [73](#), [74](#).
- ipiv*: [156](#), [158](#), [171](#).
- isEven*: [281](#), [282](#).
- isfinite*: [367](#).
- iss*: [113](#).
- istringstream*: [113](#).
- iterations*: [156](#), [165](#).
- iVector**: [21](#), [33](#), [38](#), [65](#), [68](#), [71](#), [73](#),
[74](#), [80](#), [83](#), [85](#), [122](#), [129](#), [130](#),
[131](#), [139](#), [140](#), [141](#), [143](#), [144](#), [145](#),
[147](#), [150](#), [151](#), [156](#), [160](#), [178](#), [180](#),
[181](#), [183](#), [187](#), [189](#), [204](#), [217](#), [221](#),
[222](#), [223](#), [224](#), [235](#), [277](#), [295](#), [296](#),
[315](#), [326](#), [330](#), [339](#), [342](#), [344](#), [345](#),
[352](#), [355](#), [368](#), [371](#), [372](#).
- j*: [132](#), [133](#), [134](#), [135](#), [136](#), [166](#),
[273](#), [286](#), [289](#).
- k*: [79](#), [132](#), [153](#), [187](#), [189](#), [204](#), [224](#),
[281](#), [315](#), [339](#), [342](#), [352](#), [355](#), [357](#).
- last_direction*: [305](#), [307](#).
- lda*: [153](#), [158](#).
- log*: [79](#), [113](#), [114](#).
- Log*: [114](#).
- Lorenz*: [19](#), [22](#).
- LorenzConsts**: [51](#), [52](#), [53](#).
- Lorenz2*: [53](#), [54](#).
- lwork*: [153](#), [156](#), [158](#), [171](#).
- M*: [135](#), [136](#), [153](#), [156](#), [235](#).
- m*: [132](#), [153](#), [158](#), [187](#), [189](#), [339](#),
[344](#), [352](#).
- mag*: [79](#), [113](#), [114](#), [129](#), [132](#), [314](#),
[321](#).
- main*: [20](#), [35](#), [40](#), [45](#), [48](#), [58](#), [61](#),
[68](#), [71](#), [73](#), [74](#).
- Matrix**: [123](#), [124](#), [132](#), [133](#), [134](#),
[135](#), [136](#).
- matrix_inverse*: [235](#), [251](#), [270](#), [271](#),
[274](#), [275](#).

- MatrixInverse:** 155, 156, 158,
 160, 166, 170, 171, 175, 235, 274.
MATRIXINVERSE_H: 172.
Matrix1: 132.
Matrix2: 132.
matrix2pointer: 135, 153, 158.
max: 149, 158, 165, 273, 314, 316.
max_iterations: 165.
max_order: 362.
MaxLength: 362.
mid: 113.
Mid: 114.
midpoint: 44, 55, 57, 64, 65, 73, 74,
79, 113, 114, 145, 146, 180, 250,
 258, 266, 267, 269, 273, 330.
minus: 129.
MISCFUNS_H: 374.
mu: 73.
MU: 74.
mu_h: 73.
mult: 165.
multMiMi: 132, 251, 261, 270, 271.
multMiMp: 132, 166, 167, 245,
 252, 253, 254.
multMiVi: 130, 131, 162, 245, 261,
 262, 263, 268, 270, 271.
multMiVp: 130, 131.
multMpVi: 130, 131, 272.
n: 18, 19, 21, 32, 33, 37, 38, 53,
68, 71, 73, 74, 123, 124, 132,
135, 136, 153, 156, 158, 166, 171,
178, 179, 181, 182, 191, 192, 204,
219, 235, 273, 274, 286, 327, 339,
340, 352, 353, 362.
nextafter: 314.
no: 199, 200, 325, 331.
norm: 285, 286, 287, 288, 289.
norm2: 129, 286.
num: 69, 74.
numeric_limits: 149, 165, 273.
ODE: 353.
ODENAME: 18.
off: 84, 295.
ofstream: 44, 55, 57, 62, 68, 73, 74.
ok: 251, 270, 271, 272.
on: 43, 55, 73, 74, 84, 295, 331.
one: 204, 206, 210, 214, 219.
Orbit: 73.
order: 200, 239, 302, 311, 315, 331,
339, 342, 343, 352, 355, 356.
order_trial: 204, 206, 210, 212, 214,
 217, 221, 235, 239, 241, 274.
order0: 221.
orthogonalInverse: 156, 157, 170,
 271.
out: 44, 55, 57, 62, 68, 69, 73, 74.
outFile: 44, 68.
outFileSol: 73.
outFileStep: 73.
outFile1: 55, 62, 63, 64, 65.
outFile2: 57, 62, 63, 65.
outFile3: 62, 63, 65.
outSteps: 74.
outStepSizes: 74.
p: 52, 53, 68, 79, 84, 86, 191, 193,
204, 235, 279, 280, 282, 331, 362.
param: 18, 19, 32, 37, 53, 71, 73,
74, 86, 187, 188, 189, 190, 339,
340, 347, 352, 353.
pi: 79, 113, 114.
PI: 113.
plus: 129.
pMatrix: 122, 130, 131, 132, 146,
 153, 156, 158, 160, 166, 170,
 178, 235, 284, 293.
pointer2matrix: 135, 136, 153, 158.
pow: 73, 79, 113, 114, 214, 224,
 241, 315.
Power: 114.
power: 113.
pq: 239.
predictor_excess: 235, 241, 246, 274.
prefix: 69, 74.
printMessage: 158, 163, 169, 251,
 271, 273, 369.
printVector: 26, 87, 366.
PROFIL_VNODE: 115, 120.
pVector: 122, 129, 131, 143, 144,
 145, 156, 178, 235, 286, 295.
Q: 153, 178, 235.
q: 235, 279, 280, 282.
qsort: 288, 289.

- R*: [146](#).
r: [84](#), [129](#), [143](#), [144](#), [145](#), [178](#), [331](#).
rad: [44](#), [142](#), [143](#), [164](#).
radx: [156](#), [164](#), [171](#).
rescale: [241](#).
reset: [343](#), [356](#).
resize: [124](#).
ret: [325](#).
rho: [19](#), [51](#), [52](#), [53](#).
round_control: [117](#), [118](#).
round_down: [116](#), [118](#), [119](#), [149](#),
[161](#), [212](#), [215](#).
round_nearest: [116](#), [118](#), [119](#), [153](#),
[158](#), [164](#), [376](#).
round_up: [116](#), [118](#), [119](#), [161](#), [165](#),
[168](#), [207](#), [212](#), [215](#), [314](#).
rounding_control: [117](#).
row: [357](#), [358](#).
rQR: [178](#), [179](#), [180](#), [245](#), [262](#), [271](#),
[272](#), [273](#), [276](#).
rtol: [200](#), [207](#), [301](#), [315](#), [316](#), [331](#).
S: [178](#), [235](#).
s: [79](#), [87](#), [113](#), [114](#), [129](#), [132](#), [235](#),
[344](#), [357](#), [366](#), [370](#), [374](#).
scalar: [129](#), [132](#), [135](#), [136](#).
ScalarExample: [32](#), [34](#).
scaleM: [132](#), [248](#), [249](#).
scaleV: [129](#), [132](#), [206](#), [214](#), [241](#),
[255](#), [257](#), [258](#).
set: [187](#), [188](#), [189](#), [190](#), [206](#), [210](#),
[221](#), [243](#), [244](#), [249](#), [258](#), [277](#), [311](#),
[312](#), [315](#), [339](#), [342](#), [352](#), [355](#).
set_output_digits: [79](#).
setColumn: [132](#), [133](#), [134](#), [169](#), [289](#).
setFirstEntry: [56](#), [68](#), [74](#), [84](#),
[307](#), [331](#).
setHmin: [60](#), [84](#), [331](#).
setId: [132](#), [180](#).
setM: [132](#), [248](#), [249](#).
setOneStep: [43](#), [55](#), [73](#), [74](#), [84](#), [331](#).
setOrder: [60](#), [68](#), [84](#), [331](#).
setprecision: [74](#).
setTols: [60](#), [68](#), [84](#), [331](#).
setTrialOrder: [221](#), [311](#).
setTrialStepsize: [221](#), [311](#).
setV: [129](#), [132](#), [161](#), [180](#), [207](#),
[257](#), [258](#).
showpoint: [74](#).
sigma: [19](#), [51](#), [52](#), [53](#).
sin: [79](#), [113](#), [114](#).
Sin: [114](#).
size: [124](#), [129](#), [131](#), [132](#), [146](#), [191](#),
[192](#), [327](#), [339](#), [340](#), [342](#), [343](#),
[344](#), [345](#), [347](#), [352](#), [353](#), [355](#),
[356](#), [357](#), [358](#), [368](#).
sizeM: [123](#), [124](#), [132](#), [134](#), [135](#), [136](#),
[153](#), [158](#), [166](#), [169](#), [171](#), [179](#),
[273](#), [274](#), [286](#).
sizeV: [123](#), [124](#), [129](#), [147](#), [150](#), [171](#),
[179](#), [182](#), [219](#), [274](#), [286](#), [299](#),
[327](#), [366](#), [371](#).
solution: [235](#), [238](#), [243](#), [244](#), [245](#),
[252](#), [253](#), [254](#), [262](#), [270](#), [271](#), [273](#),
[274](#), [275](#), [276](#), [277](#).
Solution: [175](#), [178](#), [179](#), [180](#),
[235](#), [274](#).
SOLUTION_H: [184](#).
Solver: [23](#), [24](#), [25](#), [43](#), [44](#), [47](#), [55](#),
[56](#), [57](#), [60](#), [61](#), [65](#), [68](#), [71](#), [73](#), [74](#).
sortColumns: [273](#), [284](#), [293](#).
sorting_needed: [287](#), [289](#).
sqr: [73](#), [74](#), [79](#), [113](#), [114](#).
Sqr: [114](#).
sqrt: [79](#), [113](#), [114](#), [129](#).
Sqrt: [114](#).
start_time: [87](#), [375](#), [376](#).
std: [27](#), [35](#), [40](#), [45](#), [48](#), [58](#), [61](#), [69](#),
[70](#), [71](#), [73](#), [74](#), [80](#), [113](#), [127](#), [149](#),
[173](#), [289](#), [314](#), [315](#), [316](#), [334](#),
[366](#), [373](#), [374](#), [376](#).
step: [47](#).
stepAction: [84](#), [295](#), [331](#).
steps: [71](#), [295](#), [308](#), [326](#), [327](#), [330](#).
stepsize: [339](#), [342](#), [343](#), [346](#), [352](#),
[355](#), [356](#).
str: [69](#), [74](#).
string: [69](#), [74](#).
string_to_interval: [38](#), [47](#), [73](#), [79](#),
[113](#), [114](#).
stringstream: [69](#), [74](#).
subFromId: [132](#), [167](#), [261](#).

- subMiMp*: [132](#), [268](#).
subset: [113](#).
subseteq: [79](#), [113](#), [114](#), [139](#), [166](#),
[214](#), [215](#), [272](#), [318](#).
subViVi: [129](#), [259](#).
subViVp: [129](#), [132](#), [180](#), [261](#), [268](#).
success: [198](#), [296](#), [309](#), [325](#), [329](#).
successful: [25](#), [44](#), [85](#), [329](#).
sum: [187](#), [188](#), [189](#), [190](#), [339](#),
[344](#), [352](#).
Sum: [357](#).
sum_old_radui: [165](#).
sum_radui: [164](#), [165](#).
sumTerms: [187](#), [188](#), [189](#), [190](#), [243](#),
[244](#), [339](#), [344](#), [352](#), [357](#).
sumV: [129](#).
sup: [64](#), [65](#), [79](#), [113](#), [114](#), [149](#), [214](#),
[215](#), [306](#), [318](#), [321](#), [323](#), [367](#).
Sup: [114](#).
sysconf: [376](#).
t: [18](#), [19](#), [21](#), [32](#), [33](#), [37](#), [38](#), [53](#), [68](#),
[71](#), [73](#), [74](#), [83](#), [178](#), [181](#), [277](#),
[314](#), [315](#), [339](#).
t_enc: [206](#), [214](#).
t_end: [295](#), [296](#), [299](#), [300](#), [304](#), [306](#),
[309](#), [314](#), [318](#), [320](#), [321](#), [323](#).
t_next: [235](#), [236](#), [238](#), [249](#), [258](#),
[264](#), [314](#).
t_prev: [74](#).
t_trial: [204](#), [215](#), [295](#), [300](#), [320](#), [321](#),
[323](#), [324](#), [326](#).
t_0: [206](#).
tan: [79](#), [113](#), [114](#).
Tan: [114](#).
tau: [153](#).
tayl_coeff: [204](#), [206](#), [210](#), [217](#), [219](#),
[221](#), [223](#), [241](#), [243](#), [257](#), [258](#).
tayl_coeff_ode: [191](#), [192](#), [193](#), [221](#),
[241](#), [243](#), [257](#), [258](#), [315](#).
tayl_coeff_var: [191](#), [192](#), [193](#), [244](#),
[248](#), [249](#).
td: [215](#).
temp: [235](#), [243](#), [245](#), [258](#), [261](#),
[262](#), [263](#), [270](#), [271](#), [272](#), [274](#),
[295](#), [315](#), [327](#).
temp2: [235](#), [245](#), [258](#), [259](#), [272](#), [274](#).
tem2: [245](#).
tend: [21](#), [24](#), [25](#), [33](#), [38](#), [44](#), [47](#), [55](#),
[56](#), [57](#), [61](#), [68](#), [71](#), [73](#), [74](#), [83](#).
tend2: [56](#).
term: [187](#), [189](#), [204](#), [206](#), [219](#), [235](#),
[257](#), [258](#), [261](#), [274](#), [339](#), [345](#), [352](#).
Term: [358](#).
tf: [352](#), [353](#), [355](#), [362](#).
tf_in: [352](#), [353](#), [354](#), [355](#), [356](#),
[357](#), [358](#).
tf_out: [352](#), [353](#), [356](#).
TFfunction: [352](#), [353](#), [362](#).
TFinterval: [352](#), [353](#).
Tfunction: [339](#), [340](#), [362](#).
theta: [321](#).
time: [68](#), [74](#).
time_end: [71](#).
time_start: [71](#).
times: [376](#).
Tinterval: [339](#), [340](#).
Tj: [65](#), [295](#), [318](#), [320](#), [321](#).
tmp: [279](#), [280](#), [286](#), [289](#), [318](#), [321](#).
tms: [376](#).
Tms: [376](#).
tms_etime: [376](#).
tol: [68](#), [69](#), [207](#), [315](#).
tonearest: [118](#).
tp: [295](#), [327](#), [330](#).
transform: [129](#), [143](#), [144](#), [145](#).
transpose: [132](#), [170](#).
trial_solution: [235](#), [263](#), [264](#), [266](#),
[267](#), [268](#), [269](#), [270](#), [271](#), [272](#), [273](#),
[274](#), [275](#), [276](#), [277](#).
true: [79](#), [85](#), [113](#), [117](#), [139](#), [140](#),
[141](#), [147](#), [153](#), [157](#), [158](#), [160](#),
[163](#), [166](#), [170](#), [217](#), [251](#), [281](#),
[287](#), [325](#), [329](#), [368](#).
tt: [214](#).
t0: [55](#), [56](#), [178](#), [180](#), [181](#), [183](#), [187](#),
[189](#), [204](#), [206](#), [210](#), [214](#), [215](#), [217](#),
[221](#), [295](#), [296](#), [299](#), [300](#), [304](#), [306](#),
[309](#), [310](#), [311](#), [312](#), [314](#), [315](#), [317](#),
[321](#), [325](#), [326](#), [339](#), [342](#), [352](#), [355](#).
t1: [37](#).
U: [235](#).
u: [178](#), [204](#), [224](#).

- u_next*: [235](#), [243](#), [246](#), [274](#).
upward: [118](#).
V: [235](#).
v: [87](#), [129](#), [143](#), [204](#), [224](#), [366](#).
v_bias: [80](#), [115](#), [120](#), [127](#), [129](#), [132](#),
[139](#), [140](#), [141](#), [142](#), [143](#), [144](#),
[145](#), [147](#), [149](#), [151](#), [153](#), [158](#),
[161](#), [164](#), [166](#), [168](#), [172](#), [178](#), [180](#),
[181](#), [183](#), [184](#), [187](#), [189](#), [194](#), [195](#),
[241](#), [273](#), [291](#), [292](#), [293](#), [299](#), [304](#),
[306](#), [314](#), [318](#), [321](#), [323](#), [360](#), [367](#),
[368](#), [371](#), [373](#), [376](#).
v_blas: [127](#), [137](#), [151](#), [152](#), [154](#),
[172](#), [173](#), [184](#), [194](#), [195](#), [224](#),
[227](#), [286](#), [291](#), [292](#), [293](#), [360](#),
[366](#), [368](#), [372](#), [373](#).
v_compare: [288](#), [289](#).
validate: [317](#), [318](#).
var_type: [18](#), [19](#), [32](#), [37](#), [53](#), [71](#),
[73](#), [74](#).
VDP: [74](#).
Vector: [123](#), [124](#), [129](#), [133](#), [134](#).
vector: [80](#), [122](#).
VECTOR_MATRIX: [127](#).
Vector1: [129](#).
Vector2: [129](#).
VNODE: [23](#), [71](#), [73](#), [74](#), [81](#), [82](#),
[175](#), [197](#), [295](#), [296](#), [315](#), [326](#),
[327](#), [328](#).
VNODE_DEBUG: [158](#), [163](#), [166](#), [169](#),
[251](#), [271](#), [273](#), [369](#).
VNODE_H: [335](#).
VNODEINT_H: [333](#).
VNODEINTERVAL_H: [115](#).
vnodelp: [27](#), [35](#), [40](#), [45](#), [48](#), [58](#),
[61](#), [70](#), [71](#), [73](#), [74](#), [154](#), [184](#), [185](#),
[194](#), [195](#), [196](#), [201](#), [226](#), [227](#), [291](#),
[292](#), [293](#), [333](#), [334](#), [349](#), [350](#), [360](#),
[361](#), [363](#), [371](#), [373](#), [374](#).
vnodelpMessage: [298](#), [299](#), [300](#), [301](#),
[302](#), [303](#), [304](#), [307](#), [309](#), [310](#),
[317](#), [370](#), [374](#).
VNODEROUND_H: [120](#).
void: [339](#), [352](#).
VTIMING_H: [375](#).
w: [235](#).
width: [64](#), [65](#), [79](#), [113](#), [114](#), [142](#),
[144](#), [273](#), [314](#), [321](#), [330](#).
work: [153](#), [156](#), [158](#), [171](#).
x: [129](#), [130](#), [131](#), [147](#), [156](#), [160](#), [235](#).
x1: [156](#), [162](#), [163](#), [171](#).
Y: [65](#).
y: [18](#), [19](#), [21](#), [32](#), [33](#), [37](#), [38](#), [53](#), [68](#),
[71](#), [73](#), [74](#), [83](#), [129](#), [147](#), [178](#),
[181](#), [235](#), [277](#), [339](#), [352](#).
y_coeff: [339](#), [340](#), [341](#), [342](#), [343](#),
[344](#), [345](#), [347](#).
y_pred: [235](#), [246](#), [249](#), [258](#), [261](#),
[263](#), [274](#).
y_pred_point: [235](#), [258](#), [261](#), [263](#),
[268](#), [274](#).
yp: [18](#), [19](#), [32](#), [37](#), [53](#), [71](#), [73](#), [74](#),
[339](#), [352](#).
y0: [178](#), [180](#), [181](#), [183](#), [187](#), [189](#),
[204](#), [206](#), [207](#), [217](#), [221](#), [295](#), [296](#),
[299](#), [310](#), [311](#), [312](#), [315](#), [317](#), [325](#),
[326](#), [339](#), [342](#), [352](#), [355](#).
z: [129](#), [130](#), [131](#), [133](#), [134](#), [147](#), [235](#).

List of Refinements

- $\langle (-1)^i c_i^{q,p} = (-1)^i \frac{q!}{(p+q)!} \frac{(q+p-i)!}{(q-i)!} \text{ 280} \rangle$ Used in chunk 293.
- $\langle A_{j+1} = m(\mathbf{A}_{j+1}) \text{ 269} \rangle$ Used in chunk 265.
- $\langle C_{j+1} = m(\mathbf{B}_{j+1}) \text{ 250} \rangle$ Used in chunk 247.
- $\langle S_{j+1} = m(\mathbf{S}_{j+1}) \text{ 267} \rangle$ Used in chunk 265.
- $\langle \mathbf{A}_{j+1} = \mathbf{G}_{j+1} A_j \text{ 253} \rangle$ Used in chunk 247.
- $\langle \mathbf{G}_{j+1} = C_{j+1}^{-1} \mathbf{F}_j \text{ 251} \rangle$ Used in chunk 247.
- $\langle \mathbf{Q}_{j+1} = \mathbf{G}_{j+1} Q_j \text{ 254} \rangle$ Used in chunk 247.
- $\langle \mathbf{S}_{j+1} = \mathbf{G}_{j+1} S_j \text{ 252} \rangle$ Used in chunk 247.
- $\langle \mathbf{y}_{j+1}^* = (\hat{u}_{j+1} + \mathbf{z}_{j+1} + \mathbf{x}_{j+1}) \cap \tilde{\mathbf{y}}_j \text{ 246} \rangle$ Used in chunk 242.
- $\langle \mathbf{B}_{j+1} = \sum_{i=0}^q (-1)^i c_i^{q,p} h_j^i J(f^{[i]}; \mathbf{y}_{j+1}^*) \text{ 249} \rangle$ Used in chunk 247.
- $\langle \mathbf{F}_j = \sum_{i=0}^p c_i^{p,q} h_j^i J(f^{[i]}; \mathbf{y}_j) \text{ 248} \rangle$ Used in chunk 247.
- $\langle \mathbf{U}_{j+1} = I + \sum_{i=1}^q h_j^i J(f^{[i]}; \mathbf{y}_j) \text{ 244} \rangle$ Used in chunk 242.
- $\langle \mathbf{d}_{j+1} = g_{j+1} + \mathbf{e}_{j+1} \text{ 260} \rangle$ Used in chunk 247.
- $\langle \mathbf{e}_{j+1} = (-1)^q \gamma_{p,q} \mathbf{h}_j^{p+q+1} f^{[p+q+1]}(\mathbf{T}_j, \tilde{\mathbf{y}}_j) \text{ 255} \rangle$ Used in chunk 247.
- $\langle \mathbf{r}_{\text{QR},j+1} = (Q_{j+1}^{-1} \mathbf{Q}_{j+1}) \mathbf{r}_{\text{QR},j} + Q_{j+1}^{-1} \mathbf{v}_{j+1} \text{ 271} \rangle$ Used in chunk 265.
- $\langle \mathbf{r}_{j+1} = (A_{j+1}^{-1} \mathbf{A}_{j+1}) \mathbf{r}_j + A_{j+1}^{-1} \mathbf{v}_{j+1} \text{ 270} \rangle$ Used in chunk 265.
- $\langle \mathbf{s}_{j+1} = (\mathbf{A}_{j+1} \mathbf{r}_j) \cap (\mathbf{Q}_{j+1} \mathbf{r}_{\text{QR},j}) \text{ 262} \rangle$ Used in chunk 247.
- $\langle \mathbf{v}_{j+1} = \mathbf{y}_{j+1}^* - u_{j+1} + (\mathbf{S}_{j+1} - \mathbf{S}_{j+1}) \boldsymbol{\alpha} + \mathbf{w}_{j+1} \text{ 268} \rangle$ Used in chunk 265.
- $\langle \mathbf{w}_{j+1} = C_{j+1}^{-1} \mathbf{d}_{j+1} + (I - C_{j+1}^{-1} \mathbf{B}_{j+1})(\mathbf{y}_{j+1}^* - \mathbf{y}_{j+1}^*) \text{ 261} \rangle$ Used in chunk 247.
- $\langle \mathbf{y}_{j+1} = (\mathbf{y}_{j+1}^* + \mathbf{S}_{j+1} \boldsymbol{\alpha} + \mathbf{s}_{j+1} + \mathbf{w}_{j+1}) \cap \mathbf{y}_{j+1}^* \text{ 263} \rangle$ Used in chunk 247.
- $\langle \hat{u}_{j+1} = u_j + \sum_{i=1}^q h_j^i f^{[i]}(t_j, u_j) \text{ 243} \rangle$ Used in chunk 242.
- $\langle c_i^{p,q} = \frac{p!}{(p+q)!} \frac{(q+p-i)!}{(p-i)!} \text{ 279} \rangle$ Used in chunk 293.
- $\langle g_{j+1} = g_{j+1}^f - g_{j+1}^b \text{ 259} \rangle$ Used in chunk 256.
- $\langle g_{j+1} = \sum_{i=0}^p c_i^{p,q} h_j^i f^{[i]}(u_j) - \sum_{i=0}^q (-1)^i c_i^{q,p} h_j^i f^{[i]}(y_{j+1}^*) \text{ 256} \rangle$ Used in chunk 247.
- $\langle g_{j+1}^b = \sum_{i=0}^q (-1)^i c_i^{q,p} h_j^i f^{[i]}(y_{j+1}^*) \text{ 258} \rangle$ Used in chunk 256.
- $\langle g_{j+1}^f = \sum_{i=0}^p c_i^{p,q} h_j^i f^{[i]}(u_j) \text{ 257} \rangle$ Used in chunk 256.
- $\langle h \text{ such that } [0, h] \mathbf{a} \subseteq \mathbf{b} \text{ (intervals) 149} \rangle$ Used in chunk 150.
- $\langle h \text{ such that } [0, h] \mathbf{a} \subseteq \mathbf{b} \text{ (interval vectors) 150} \rangle$ Used in chunk 152.
- $\langle u_{j+1} = m(\mathbf{y}_{j+1}) \text{ 266} \rangle$ Used in chunk 265.
- $\langle \mathbf{x}_{j+1} = (\mathbf{U}_{j+1} S_j) \boldsymbol{\alpha} + \{(\mathbf{U}_{j+1} A_j) \mathbf{r}_j \cap (\mathbf{U}_{j+1} Q_j) \mathbf{r}_{\text{QR},j}\} \text{ 245} \rangle$ Used in chunk 242.
- $\langle \text{AD_ODE 187} \rangle$ Used in chunk 194.
- $\langle \text{AD_VAR 189} \rangle$ Used in chunk 195.

⟨DETEST E1 37⟩ Used in chunk 40.
 ⟨E1.cc 40⟩
 ⟨FadbadODE destructor 341⟩ Used in chunk 350.
 ⟨Krawczyk's iteration 162, 163⟩ Used in chunk 165.
 ⟨Lorenz 19⟩ Used in chunks 20, 45, 48, 61, and 70.
 ⟨VNODE accept solution 326⟩ Used in chunk 334.
 ⟨VNODE-LP message 370⟩ Used in chunk 373.
 ⟨a priori enclosure representation 181⟩ Used in chunk 184.
 ⟨accept solution (HOE) 220⟩ Used in chunk 227.
 ⟨accept solution (IHO) 276⟩ Used in chunk 293.
 ⟨ad_ode.h 194⟩
 ⟨ad_var.h 195⟩
 ⟨allad.h 196⟩
 ⟨basic.cc 27⟩
 ⟨basiclinalg.h 137⟩
 ⟨case 1 320⟩ Used in chunk 318.
 ⟨case 2 321⟩ Used in chunk 318.
 ⟨case 3 323⟩ Used in chunk 318.
 ⟨changing the rounding mode (BIAS) 119⟩ Used in chunk 120.
 ⟨changing the rounding mode (FILIB++) 118⟩ Used in chunk 120.
 ⟨check finite 367, 368⟩ Used in chunk 374.
 ⟨check if a number is even 281⟩ Used in chunk 293.
 ⟨check if in the interior 140⟩ Used in chunk 151.
 ⟨check if last step 318⟩ Used in chunk 296.
 ⟨check if sorting is needed 287⟩ Used in chunk 284.
 ⟨check if success 25⟩ Used in chunks 20, 35, 40, and 68.
 ⟨check input correctness 298, 299, 300, 301, 302, 303, 304⟩ Used in chunk 296.
 ⟨check vector inclusion 139, 141⟩ Used in chunk 151.
 ⟨check vector intersection 371⟩ Used in chunk 373.
 ⟨check t 316⟩ Used in chunk 315.
 ⟨class FadbadODE 339⟩ Used in chunk 349.
 ⟨class FadbadVarODE 352⟩ Used in chunk 360.
 ⟨class HOE 204⟩ Used in chunk 226.
 ⟨class IHO 235⟩ Used in chunk 292.
 ⟨class MatrixInverse 156⟩ Used in chunk 172.
 ⟨class VNODE 295⟩ Used in chunk 333.
 ⟨close files 63⟩ Used in chunk 61.
 ⟨compare function 288⟩ Used in chunk 291.
 ⟨compute $(-1)^q q! p! / (p + q)!$ 282⟩ Used in chunk 293.
 ⟨compute A^{-1} 158⟩ Used in chunk 173.
 ⟨compute Q_{j+1} 273⟩ Used in chunk 265.
 ⟨compute β 224⟩ Used in chunk 227.
 ⟨compute p_j 206⟩ Used in chunk 217.
 ⟨compute u_j and \tilde{y}_j 207⟩ Used in chunk 217.
 ⟨compute h_{\min} 314⟩ Used in chunk 334.
 ⟨compute IHO method coefficients 239⟩ Used in chunk 293.

<compute QR factorization 153> Used in chunk 154.
 <compute Taylor coefficients (FadbadODE) 343> Used in chunk 350.
 <compute Taylor coefficients (FadbadVarODE) 356> Used in chunk 361.
 <compute column norms of A 286> Used in chunk 284.
 <compute enclosure 324> Used in chunk 296.
 <compute initial box 161> Used in chunk 160.
 <compute initial stepsize 315> Used in chunk 334.
 <compute stepsize 210, 212, 214> Used in chunk 217.
 <compute tight enclosure 236> Used in chunk 293.
 <constants Lorenz 51> Used in chunk 58.
 <constructor (FadbadODE) 340> Used in chunk 350.
 <constructor (FadbadVarODE) 353> Used in chunk 361.
 <constructor (VNODE) 327> Used in chunk 334.
 <constructor IHO 274> Used in chunk 293.
 <constructor-destructor HOE 219> Used in chunk 227.
 <control class 200> Used in chunk 201.
 <control.h 201>
 <corrector: compute y_{j+1} 247> Used in chunk 236.
 <create AD object 22> Used in chunks 20, 45, 48, 61, and 68.
 <create E1 39> Used in chunk 40.
 <create a solver 23> Used in chunks 20, 35, 40, 45, 48, 58, 61, and 68.
 <create file name 69> Used in chunk 68.
 <create problem object with parameters 54> Used in chunk 58.
 <create scalar AD object 34> Used in chunk 35.
 <create **Apriori** 182> Used in chunk 185.
 <create **Solution** 179> Used in chunk 185.
 <debug.cc 373>
 <debug.h 372>
 <decide 325> Used in chunk 296.
 <destructor (FadbadVarODE) 354> Used in chunk 361.
 <destructor (VNODE) 328> Used in chunk 334.
 <destructor IHO 275> Used in chunk 293.
 <determine direction 305, 306, 307> Used in chunk 296.
 <detest_c3.cc 71>
 <do Krawczyk's iteration 165> Used in chunk 160.
 <encapsulated AD 191> Used in chunk 196.
 <encapsulated FADBAD++ AD 362> Used in chunk 363.
 <enclose each column 169> Used in chunks 166 and 170.
 <enclose the inverse of a matrix 166> Used in chunk 173.
 <enclose the inverse of an orthogonal matrix 170> Used in chunk 173.
 <enclose the solution to $Ax = b$ 160> Used in chunk 173.
 <fadbad_ad.cc 350>
 <fadbad_ad.h 349>
 <fadbad_advar.cc 361>
 <fadbad_advar.h 360>
 <fadbadad.h 363>

<find beta 167, 168> Used in chunks 166 and 170.
 <find initial stepsize 310> Used in chunk 308.
 <find minimum stepsize 309> Used in chunk 308.
 <find solution representation for next step 265> Used in chunk 236.
 <form time interval 215> Used in chunk 217.
 <functions calling FILIB++ 113> Used in chunk 115.
 <functions calling PROFIL 114> Used in chunk 115.
 <get functions (VNODE) 329, 330> Used in chunk 295.
 <get functions HOE 222, 223> Used in chunk 204.
 <get term (FadbadODE) 345> Used in chunk 350.
 <get term (FadbadVarODE) 358> Used in chunk 361.
 <get/set column 134> Used in chunk 137.
 <hoe.cc 227>
 <hoe.h 226>
 <iho.cc 293>
 <iho.h 292>
 <implementation of encapsulated AD 192, 193> Used in chunk 196.
 <index-norm structure 285> Used in chunk 291.
 <indicate single step 43> Used in chunks 44 and 61.
 <indicator type 198> Used in chunk 201.
 <initialize IHO method 238, 240, 241> Used in chunk 236.
 <initialize Taylor coefficients (FadbadODE) 342> Used in chunk 350.
 <initialize Taylor coefficients (FadbadVarODE) 355> Used in chunk 361.
 <initialize integration 308> Used in chunk 296.
 <initialize **Apriori** 183> Used in chunk 185.
 <initialize **Solution** 180> Used in chunk 185.
 <integ.cc 334>
 <integctrl.cc 61>
 <integi.cc 45>
 <integrate (basic) 24> Used in chunks 20, 35, 40, and 68.
 <integrate from t to $tend/2$ 56> Used in chunk 57.
 <integrate with intermediate output 47> Used in chunk 48.
 <integrate with interval initial condition 44> Used in chunk 45.
 <integrate with resetting constants 57> Used in chunk 58.
 <integrator 296> Used in chunk 334.
 <intermediate.cc 48>
 <interrupt type 199> Used in chunk 201.
 <intersection of interval vectors 147> Used in chunk 151.
 <interval data type (FILIB++) 77> Used in chunk 115.
 <interval data type (PROFIL) 76> Used in chunk 115.
 <interval vector 80> Used in chunk 122.
 <intvfuncs.cc 152>
 <intvfuncs.h 151>
 <main program for order study 68> Used in chunk 70.
 <matrix operations 132> Used in chunk 137.
 <matrix times vector 131> Used in chunk 137.

<matrix2pointer 135> Used in chunk 137.
 <matrixinverse.cc 173>
 <matrixinverse.h 172>
 <message printing 369> Used in chunk 372.
 <midpoint of an interval matrix 146> Used in chunk 151.
 <midpoint of an interval vector 145> Used in chunk 151.
 <miscfuncs.h 374>
 <obtain stepsize 346> Used in chunk 350.
 <odeparam.cc 58>
 <open files 62> Used in chunk 61.
 <orbit.cc 73>
 <orderstudy.cc 70>
 <output initial condition 64> Used in chunk 61.
 <output results 26> Used in chunks 20, 35, 40, and 47.
 <output solution 65> Used in chunk 61.
 <passing parameters to Lorenz 53> Used in chunk 58.
 <perform sorting 289> Used in chunk 284.
 <pointer2matrix 136> Used in chunk 137.
 <predictor: compute \mathbf{y}_{j+1}^* 242> Used in chunk 236.
 <print vector 366> Used in chunk 137.
 <qr.cc 154>
 <rad (interval) 142> Used in chunk 151.
 <rad (vector) 143> Used in chunk 151.
 <rebuild computational graph 347> Used in chunk 350.
 <reset if needed 272> Used in chunk 265.
 <rounding control type (FILIB++) 117> Used in chunk 120.
 <scalar ODE example 32> Used in chunk 35.
 <scalar.cc 35>
 <select stepsize 216> Used in chunk 217.
 <set \mathbf{t}_{j+1} 264> Used in chunk 236.
 <set E1 initial condition and endpoint 38> Used in chunk 40.
 <set HOE method 311> Used in chunk 308.
 <set IHO method 312> Used in chunk 308.
 <set ODE parameters 52> Used in chunk 58.
 <set VNODE parameters 331> Used in chunk 295.
 <set and get functions 277> Used in chunk 235.
 <set control data for the solver 60> Used in chunk 61.
 <set functions HOE 221> Used in chunk 204.
 <set initial condition and endpoint 21> Used in chunks 20, 42, 48, 58, and 61.
 <set interval initial condition and endpoint 42> Used in chunk 45.
 <set scalar ODE initial condition and endpoint 33> Used in chunk 35.
 <simple integration 55> Used in chunk 58.
 <simple main program 20> Used in chunk 27.
 <size/allocation 124> Used in chunk 127.
 <solution.cc 185>
 <solution.h 184>

`<sort columns 284>` Used in chunk 291.
`<sortcolumns.cc 291>`
`<sum radii 164>` Used in chunk 165.
`<sum terms (FadbadODE) 344>` Used in chunk 350.
`<sum terms (FadbadVarODE) 357>` Used in chunk 361.
`<template ODE function 18>`
`<tight enclosure representation 178>` Used in chunk 184.
`<validate and select stepsize 317>` Used in chunk 296.
`<validate existence and uniqueness 217>` Used in chunk 227.
`<vanderpol.cc 74>`
`<vector and matrix types 122>` Used in chunk 127.
`<vector operations 129>` Used in chunk 137.
`<vector_matrix.h 127>`
`<vnode.h 335>`
`<vnodeint.h 333>`
`<vnodeinterval.h 115>`
`<vnodeinterval.h 120>`
`<vtiming.cc 376>`
`<vtiming.h 375>`
`<width 144>` Used in chunk 151.
`<MatrixInverse constructor/destructor 171>` Used in chunk 173.