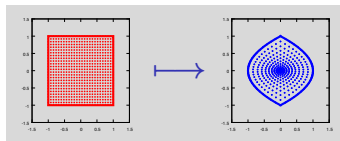
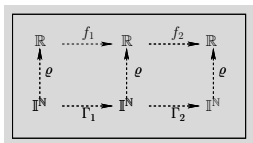


# Introduction to Exact Real Arithmetic \*



Norbert Müller

Universität Trier

Dagstuhl, 2018-04-18

\* DFG: WERA MU 1801/5-1 & CAVER BE 1267/14-1, RFBR: 14-01-91334, EU: CID Marie Skłodowska-Curie 731143

- 1 Introduction
- 2 Computability on real numbers
- 3 Exact real arithmetic
- 4 Examples

Near 1995: **ERA** (Exact Real Arithmetic) starts:

- Real numbers are atomic objects
- Arithmetic is able to deal with arbitrary real numbers ...
- ... usual entrance to  $\mathbb{R}$  is  $\mathbb{Q}$
- ... limits of (certain) sequences  $\rightsquigarrow \sqrt[n]{x}, e^x, \pi \dots$
- Underlying theory: TTE, Type-2-Theory of Effectivity ...
- ... fully(!) consistent with real calculus
- ... implying: computable functions are continuous!
- ... implying: failing tests  $x \leq y$ ,  $x \geq y$ ,  $x = y$  in case of  $x = y$  !
- ... using multi-valued functions instead

## Prototypical implementations near 1995:

- 'Precise computation software' (Oliver Aberth, C++)
- `CRCalc` (Constructive Reals Calculator, Hans Böhm, JAVA)
- XR (eXact Real arithmetic, Keith Briggs, FC++)
- 'Imperial College Reals' (Marko Krznaric, C)
- 'Manchester Reals' (David Lester, HASKELL)
- `iRRAM (M., C++)`

## later:

- `RealLib` (Branimir Lambov, C++)
- `few digits` (Russell O'Connor, HASKELL)
- AERN (Michal Konecny, HASKELL)
- `Mathemagix` (Joris van der Hoeven)
- `Marshall` (Andrej Bauer, Ivo List, HASKELL, OCaml)

- 1 Introduction
- 2 **Computability on real numbers**
- 3 Exact real arithmetic
- 4 Examples

A real number  $x$  is usually represented as follows:

- use open intervals with dyadic endpoints

$$\mathbb{I} := \left\{ \left( \frac{m_1}{2^k}, \frac{m_2}{2^k} \right) \mid m_1, m_2 \in \mathbb{Z}, k \in \mathbb{N} \right\}$$

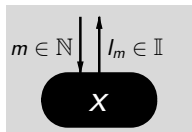
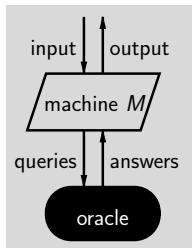
- aiming at oracle Turing machines for sequences

$$[\mathbb{N} \rightarrow \mathbb{I}] = \mathbb{I}^{\mathbb{N}}$$

- define representation  $\varrho : \subseteq \mathbb{I}^{\mathbb{N}} \rightarrow \mathbb{R}$ :

$x \in \mathbb{R}$  is represented by  $(I_m)_{m \in \mathbb{N}}$  iff

$$\lim_{m \rightarrow \infty} \text{diam}(I_m) = 0 \quad \wedge \quad \bigcap_{m \in \mathbb{N}} I_m = \{x\}$$



A real function  $f$  is computed using a machine  $M$  as follows:

- If

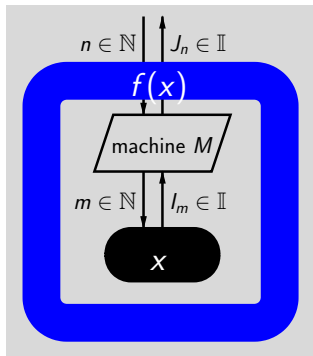
$$(I_m)_{m \in \mathbb{N}} \mapsto \mathbf{x}$$

and

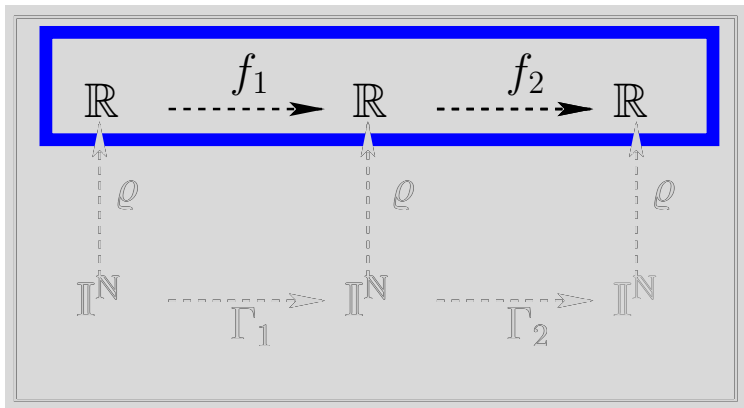
$$(I_m)_{m \in \mathbb{N}} \stackrel{M}{\rightsquigarrow} (J_n)_{n \in \mathbb{N}}$$

then

$$(J_n)_{n \in \mathbb{N}} \mapsto \mathbf{f}(\mathbf{x})$$



Computable analysis (via ‘representations’):



Remember: Computable functions are **continuous!**



- 1 Introduction
- 2 Computability on real numbers
- 3 Exact real arithmetic
- 4 Examples

## Wanted:

## Implementation of real numbers on 'real' computers

- Real numbers as abstract datatype
- Real numbers as (atomic) objects

Close at hand:

- $\mathbf{x} \in \mathbb{R} \iff \lambda n. I_n \in \mathbb{I}^{\mathbb{N}}$
- so just implement  $\lambda n. I_n$  in your favorite language
- with assertion

$$\lim_{n \rightarrow \infty} \text{diam}(I_n) = 0 \quad \wedge \quad \{\mathbf{x}\} = \bigcap_{n \in \mathbb{N}} I_n$$

## Properties of ERA w.r.t TTE:

- Users want to base decisions on the results of programs:
  - ▶ Discrete input must be possible (standard notation of  $\mathbb{N}$ ,  $\mathbb{Q}$ ).
  - ▶ Implementation has to provide human-readable (discrete) output.
  - ▶ Input/output might be (initial segments of) sequences.
- ◊ Equivalence to 'standard' representations!
- Composition is central operator, i.e. interface similar to RealRAM
- Evaluation will be DAG-based  
(although the DAG might be hidden).
- 'Standard' representations too restrictive for efficient composition.

## Common aspects in ERA implementations:

- algebraic approach
- similar to BSS-style RealRAM
- restricted to (TTE-)computability
- complete in matters of (TTE-)computability

## Differences on low level / internal structure:

- Representation of real numbers  
(infinite sequences of signed digits, intervals, Taylor models...)
- Programming paradigm: functional / object-oriented
- Lazy or eager evaluation
- Efficiency (computation time, memory)

## Example: Rump's example (almost polynomial)

```

1 REAL p ( const REAL& a, const REAL& b){
2   return 21*b*b - 2*a*a + 55*b*b*b*b
3           - 10 * a*a*b*b + a/(2*b);
4 }
5
6 void compute(){
7   REAL a = 77617, b = 33096, c = p(a,b);
8
9   cout << "Rump's_example\n";
10
11  cout << setRwidth( 20) << c << "\n";
12  cout << setRwidth( 40) << c << "\n";
13  cout << setRwidth( 60) << c << "\n";
14 }

```

## Result:

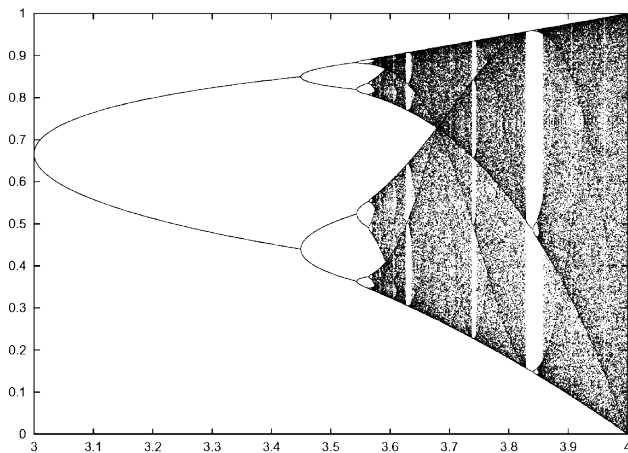
```

1 Rump's_example
2 -.827396059947E+0000
3 -.82739605994682136814116509547982E+0000
4 -.8273960599468213681411650954798162919990331157843848E+0000

```

## Example: Logistic map

$$x_{n+1} = c \cdot x_n \cdot (1 - x_n) \quad \text{for } x_0 \in (0, 1), c \in (3, 4)$$



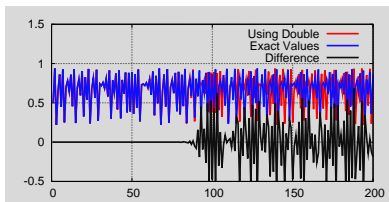
## Example: Logistic map

$$x_{n+1} = c \cdot x_n \cdot (1 - x_n) \quad \text{for } x_0 = 0.5, c = 3.75$$

```

1 void itsyst(int i){
2
3     REAL x, c;
4
5     x = 0.5; c = 3.75;
6
7     for ( int n = 1; n <= i; n++ ){
8         x = c * x * (1-x);
9     }
10    cout << x ;
11 }

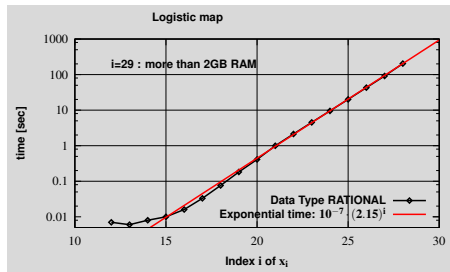
```



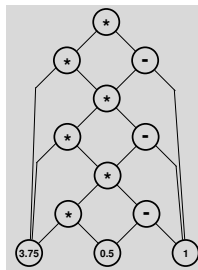
lossless representation of **rational/algebraic/real numbers**, e.g.:

- rational: store nominator/denominator  $\in \mathbb{Z}$
- algebraic: store rational polynomials + choice of root
- DAGs as data structure

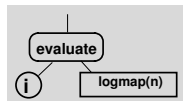
Example: logistic map  $x_{i+1} = c \cdot x_i \cdot (1 - x_i)$  with  $c = 3.75$ ,  $x_0 = 0.5$



memory:  $\exp(i)$



lin( $i$ )



log( $i$ )







## exact real arithmetic: 'constructing' approach using DAGs

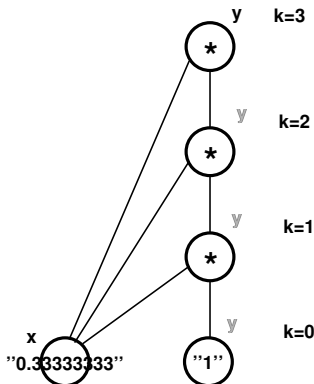
```

REAL z = power( "0.33333333" , 3);

REAL power(const REAL& x, int n) {
  REAL y=1;
  for (int k=0; k<n; k=k+1)
    { y=x*y; }
  return y;
}

```

- 'Lazy Evaluation' (ideally: using MP-intervals)
- values are approximated, but only on demand
- evaluation bottom-up or top-down
- memory requirements!!!



data structures behind REAL variables ...

... exactly represent the exact values

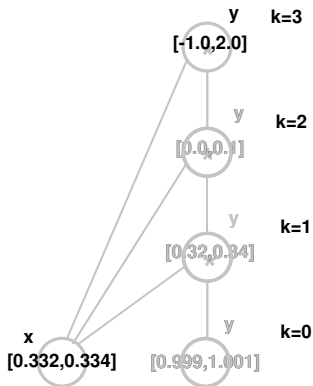
## exact real arithmetic: 'approximating' approach

```

REAL z = power( "0.33333333" , 3);

REAL power(const REAL& x, int n) {
  REAL y=1;
  for (int k=0; k<n; k=k+1)
    { y=x*y; }
  return y;
}

```



## exact real arithmetic: 'approximating' approach

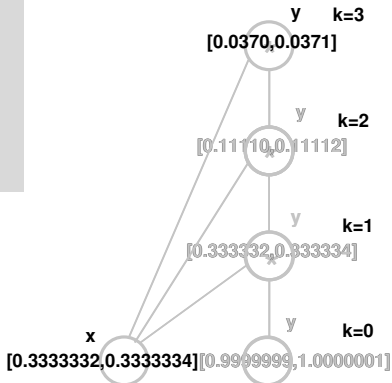
```

REAL z = power( "0.33333333" , 3);

REAL power(const REAL& x, int n) {
  REAL y=1;
  for (int k=0; k<n; k=k+1)
    { y=x*y; }
  return y;
}

```

- iteration of computations!
- 'Exceptions' are the rule...



data structures behind REAL variables ...

... represent only approximations

Cost of reevaluating a DAG:

If  $2t(n) \leq t(2n) \leq c \cdot t(n)$ , then  $\sum_{k=0}^{\lceil \log_2 n \rceil} t(2^k) \in \mathcal{O}(t)$ .

Example:

$$t(n) = \lfloor n^\alpha (\log n)^\beta (\log \log n)^\gamma \rfloor$$

für  $\alpha \geq 1, \beta, \gamma \geq 0, \alpha, \beta, \gamma \in \mathbb{R}$

→ successful evaluation dominates

Some details on the sequences  $(I_m)_{m \in \mathbb{N}} \in \mathbb{I}^{\mathbb{N}}$  in iRRAM:

- Indices  $m \in \mathbb{N}$  are associated with ‘effort’:
  - ▶ Iterations in evaluation correspond to index  $m$
  - ▶ Effort restricts the precision of operations
  - ▶  $m = 0$ : Use double precision intervals
  - ▶  $m > 0$ : Compute with precision of at most  $2^{-p_m}$  with  $p_m \approx 1.1^m$
  - ▶ dynamic change between absolute and relative precision possible
- Simplified intervals:  $\mathbb{I} = (\mathbf{c} \pm \mathbf{r})$  where
  - ▶  $\mathbf{c}$  is MPFR
  - ▶  $\mathbf{r} = m \cdot 2^e$  for 32-bit  $m, e$
  - ▶  $e$  is chosen with  $m \approx 2^{30}$
  - ▶  $\mathbf{c}$  is truncated to absolute precision  $2^e$ .
  - ▶ Precision decreases during computations
- Non-naive interval arithmetic is applicable, e.g. Taylor models

iRRAM uses a list of control precisions as effort:

```

1 Basic precision bounds:
2 double[1] -70[2] -91[3] -113[4] -136[5] -160[6] -186[7]
3 -213[8] -242[9] -273[10] -453[15] -692[20] -1008[25]
4 -1425[30] -1976[35] -2704[40] -3668[45] -4941[50]
5 -6624[55] -8848[60] -11787[65] -15673[70]
6 -20809[75] -27596[80] -36568[85] -48426[90] -64099[95]
7 -84814[100] -112194[105] -148382[110] -196211[115]
8 ...
9 -1118475546[270] -1478304970[275] -1953896587[280]

```

- start program with parameter **-d** (debugging) or **-h** (help):
- precision is currently implemented as **int32**



## Basic Data Types of the iRRAM:

- standard C++ types (**int**, **double**, ...)
- additional data types:

**INTEGER** $\mathbb{Z}$  $\leq 500$  MB per number (wrapper for **GMP**)**RATIONAL** $\mathbb{Q}$  $\leq 500$  MB per nom./denom. (wrapper for **GMP**)**DYADIC** $\{m \cdot 2^e \mid m \in \mathbb{Z}, e \in \mathbb{Z}\}$ exponent **4** B, mantissa  $\leq 500$  MB (wrapper for **MPFR**)**REAL** $\mathbb{R}$ intervals, exponent **4** B, mantissa  $\leq 500$  MB**LAZY\_BOOLEAN** $\{T, F, \perp\}$ 

exact, finite dcpo with non-strict functions

## elementary operators

- **INTEGER / RATIONAL:**  
exact versions of  $+$ ,  $-$ ,  $*$ ,  $/$ ,  $=$ ,  $<$ ,  $==$ ,...
- **DYADIC :**  
approximating versions of  $+$ ,  $-$ ,  $*$ ,  $/$ ,  
exact versions of  $=$ ,  $<$ ,  $==$ ,...
- **REAL:**  
exact versions of  $+$ ,  $-$ ,  $*$ ,  $/$ ,  $=$ ,  
lazy versions of  $<$ ,  $<=$ ,  $==$ ,...

( $\exists$  conversions between all numeric datatypes!)

## special functions

`limit`, `limit_lip`, `limit_mv`, `lipschitz`, `taylor`,...

## derived data types

**COMPLEX, INTERVAL, REALMATRIX, SPARSEREALMATRIX,...**

## non-elementary functions

**sqrt, power, maximum, minimum,...**

**exp, log, sin, cos, asin, acos, sinh, cosh, asinh, acosh,...**

**mag, mig, interval versions of {+, -, \*, /, exp, log, sin, cos},...**

**eye, zeroes, solve, matrix versions of {+, -, \*, /, exp},...**

## direct type conversions:

| from↓        | to→ | string | int32 | double | INTEGER | DYADIC | RATIONAL | REAL | COMPLEX |
|--------------|-----|--------|-------|--------|---------|--------|----------|------|---------|
| char*/string |     |        |       |        | ✓       |        | ✓        | ✓    |         |
| int32        |     |        |       |        | ✓       | ✓      | ✓        | ✓    | ✓       |
| double       |     |        |       |        | ✓       | ✓      | ✓        | ✓    | ✓       |
| INTEGER      | +   | +      |       |        | ✓       | ✓      | ✓        | ✓    | ✓       |
| DYADIC       | +   |        |       |        | +       | ✓      |          | ✓    | ✓       |
| RATIONAL     | +   |        |       |        |         |        | ✓        | ✓    | ✓       |
| REAL         | +   |        | +     | +      | +       | +      |          | ✓    | ✓       |
| COMPLEX      |     |        |       |        |         |        |          | +    | ✓       |

✓: ‘widening’, using explicit constructor

+: ‘narrowing’, using (member) functions like `x.as_double()` or `swrite(x,w)`

explicitly overloaded operators  $x \circ y$ 

| <b>REAL</b> $\circ$       | <b>int32</b> | <b>double</b> | <b>INTEGER</b> | <b>DYADIC</b> | <b>RATIONAL</b> | <b>REAL</b> | <b>COMPLEX</b> |
|---------------------------|--------------|---------------|----------------|---------------|-----------------|-------------|----------------|
| <b>+, -, *, /</b>         | ✓            | ✓             |                |               |                 | ✓           |                |
| <b>&lt;&lt;, &gt;&gt;</b> | ✓            |               |                |               |                 |             |                |
| <b>=</b>                  |              |               |                |               |                 | ✓           |                |
| <b>+=, -=</b>             |              |               |                |               |                 | ✓           |                |
| <b>*=, /=</b>             | ✓            |               |                |               |                 | ✓           |                |

missing combinations possible through (implicit) type conversion ...  
(but with additional overhead)

e.g.  $1 + \text{COMPLEX}(2)$

## Programming in ERA:

- only continuous functions

↪ hence no total test for equality, only

$$\mathit{smaller}(x, y) = \begin{cases} T, & x < y \\ F, & x > y \\ \mathit{undef.} & x = y \end{cases}$$

- instead: multivalued tests, e.g.

$$\text{bound}(\mathbf{x}, \mathbf{k}) = \begin{cases} \mathbf{T} & , \quad |\mathbf{x}| \leq 2^{\mathbf{k}} \\ \mathbf{F} & , \quad |\mathbf{x}| \geq 2^{\mathbf{k}-1} \end{cases}$$

usefull in loops:

```

1 REAL sqrt_approx( long k, REAL x ) {
2   REAL approx = 1, error;
3   do {
4     approx = (approx + x / approx )/2;
5     error  = approx - x / approx;
6   } while ( ! bound( error, k ) );
7   return approx;
8 }
```

- additional: usefull operator for limits, e.g. to define  $\sqrt{\mathbf{x}}$ :

```

1 REAL sqrt(REAL x) {return limit(sqrt_apprx ,x);}
```

- 1 Introduction
- 2 Computability on real numbers
- 3 Exact real arithmetic
- 4 Examples

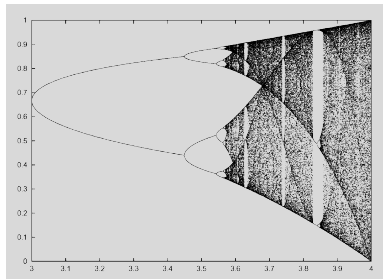


## Example: Logistic map using Taylor models in iRRAM

```

1 void itsyst( REAL& c, int n){
2   TM x;
3   x = REAL(0.125);
4   for ( int i=0; i<=n; i++ ){
5     TM::polish(x);
6     x = x * c * (REAL(1)-x);
7   }
8   cout << REAL(x) ;
9 }

```



| c             | Data type TM |                     |             |                     | Data type REAL |                     |             |                     |
|---------------|--------------|---------------------|-------------|---------------------|----------------|---------------------|-------------|---------------------|
|               | n=10000      |                     | n=100000    |                     | n=10000        |                     | n=100000    |                     |
|               | time<br>[s]  | precision<br>[bits] | time<br>[s] | precision<br>[bits] | time<br>[s]    | precision<br>[bits] | time<br>[s] | precision<br>[bits] |
| 3.125         | 0.09         | double              | 0.90        | double              | 1.08           | 18581               | 266         | 175466              |
| 3.56982421875 | 0.09         | double              | 0.94        | double              | 0.85           | 18581               | 363         | 219405              |
| 3.75          | 0.64         | 5894                | 115         | 57301               | 1.60           | 23299               | 400         | 219405              |
| 3.82          | 0.75         | 7440                | 148         | 71699               | 1.38           | 23299               | 340         | 219405              |
| 3.830078125   | 0.09         | double              | 0.92        | double              | 1.40           | 23299               | 337         | 219405              |
| 3.84          | 0.09         | 136                 | 0.89        | 136                 | 1.46           | 23299               | 354         | 219405              |

## Example: Van der Pol oscillator, discretized

- nonlinear differential equation,  $d = 2$

$$\dot{x} = y$$

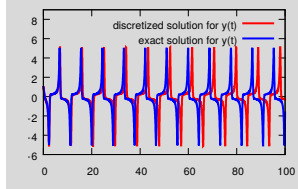
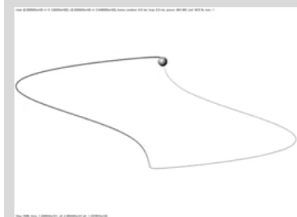
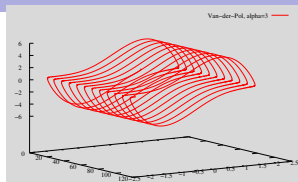
$$\dot{y} = \alpha y - x - \alpha x^2 y$$

- using  $\alpha = 3$
- initial value  $w_0 = (1, 1)$  at  $t_0 = 0$
- discretized with  $\Delta t = 0.01$  to

$$x_{n+1} = x_n + \Delta t \cdot y_n$$

$$y_{n+1} = y_n + \Delta t \cdot (\alpha y_n - x_n - \alpha x_n^2 y_n)$$

| $t_{end}$ | $n$        | Data type TM |                  | Data type REAL |                  |
|-----------|------------|--------------|------------------|----------------|------------------|
|           |            | time [s]     | precision [bits] | time [s]       | precision [bits] |
| 10        | 1 000      | 0.05         | double           | 0.01           | 136              |
| 100       | 10 000     | 0.42         | double           | 0.18           | 1737             |
| 1 000     | 100 000    | 4.6          | 136              | 6.9            | 14807            |
| 10 000    | 1 000 000  | 32           | 136              | 2395           | 175466           |
| 100 000   | 10 000 000 | 305          | 136              | -              | -                |



## Example: Van der Pol oscillator, exact

### Part I: Taylor series

- Consider a sequence of Taylor coefficients  $(a_n)_{n \in \mathbb{N}}$  together with pair  $R, M$  for  $|a_n| \leq M \cdot R^{-n}$
- ↪ operator for *infinite* summation, transparent for Taylor models:
- Use Taylor model arithmetic for partial sums  $S_{n,x}$  and error bounds

$$S_{n,x} := \sum_{k=0}^n a_k x^k \quad E_{n,x} := \frac{M \cdot R}{R - |x|} \cdot \left( \frac{|x|}{R} \right)^{n+1}$$

until  $E_{n,x}$  is ‘small enough’

- then perform eager approximation by  $S_{n,x} + [0 \pm E_{n,x}]$

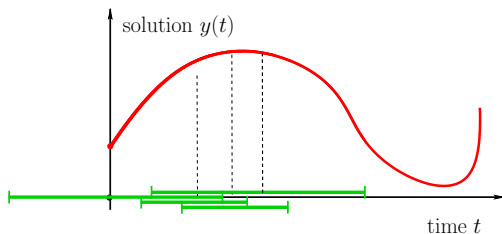
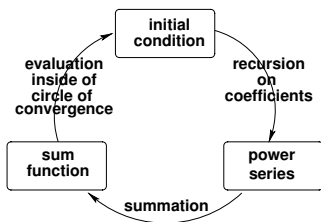
```

1 FUNCTION <TM, int> a = ...;
2 REAL R = ...; REAL M = ...; TM x = ...;
3 FUNCTION<TM, TM> f = taylor_sum(a,R,M);
4 cout << f(x);

```

Example: Van der Pol oscillator, exact

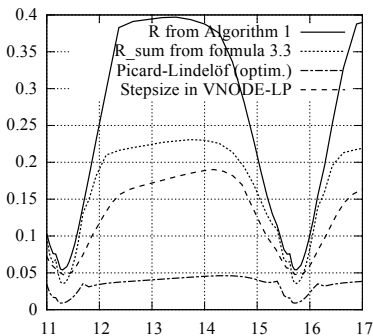
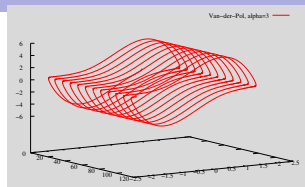
Part II: power series method, iterated:



- radii of convergence are finite (unless system is linear)
- similar to analytic continuation, but finite(!) states  $\mathbf{w}_i$  at  $\mathbf{t}_i$
- again *polish* states  $\mathbf{w}_i$  at times  $\mathbf{t}_i$

## Example: Van der Pol oscillator, exact

compute solution at  $t_{\text{end}}$  with **22** decimals:



| $t_{\text{end}}$ | Taylor models |             | intervals |             |                               |
|------------------|---------------|-------------|-----------|-------------|-------------------------------|
|                  | time [s]      | prec [bits] | time [s]  | prec [bits] | prec/ $t_{\text{end}}$ [bits] |
| 8.25             | 56            | 242         | 33        | 242         | 29.3                          |
| 15.75            | 108           | 242         | 124       | 375         | 23.8                          |
| 23.75            | 153           | 242         | 272       | 541         | 22.8                          |
| 34.00            | 232           | 242         | 1301      | 1008        | 29.6                          |
| 63.00            | 412           | 242         | 2848      | 1332        | 21.1                          |
| 83.50            | 572           | 242         | 5562      | 1737        | 20.8                          |
| 109.25           | 761           | 242         | 10470     | 2242        | 20.5                          |
| 140.75           | 924           | 242         | 21354     | 2876        | 20.4                          |
| 200.00           | 1680          | 242         |           |             |                               |
| 250.00           | 2100          | 242         |           |             |                               |
| 300.00           | 2520          | 242         |           |             |                               |
| 350.00           | 2940          | 242         |           |             |                               |
| 500.00           | 3883          | 242         |           |             |                               |

**VNODE-LP**: 0.2s for  $t_{\text{end}} = 100$ ,  $\sim 12$  decimals

Thank you for your attention!

Questions?

Remarks?