

Round-off Error Analysis of Explicit One-Step Numerical Integration Methods

RAIM 2017, LIP, ÉNS Lyon

Sylvie Boldo¹ Florian Faissole¹ Alexandre Chapoutot²

¹Inria - LRI, Univ. Paris-Sud et CNRS - Univ. Paris-Saclay

²U2IS, ÉNSTA ParisTech

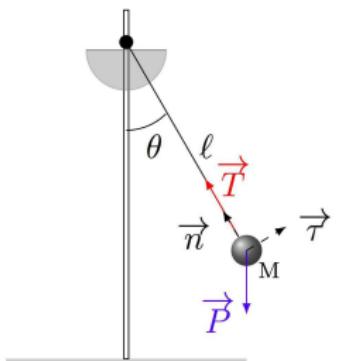


Table of contents

- 1 Motivations and numerical methods
- 2 Roundoff errors of RK methods
 - Local roundoff errors
 - Global roundoff errors of classical methods
- 3 Conclusion and perspectives

Ordinary differential equations (ODEs)

$$y'(t) = f(y, t).$$



Motion of pendulum



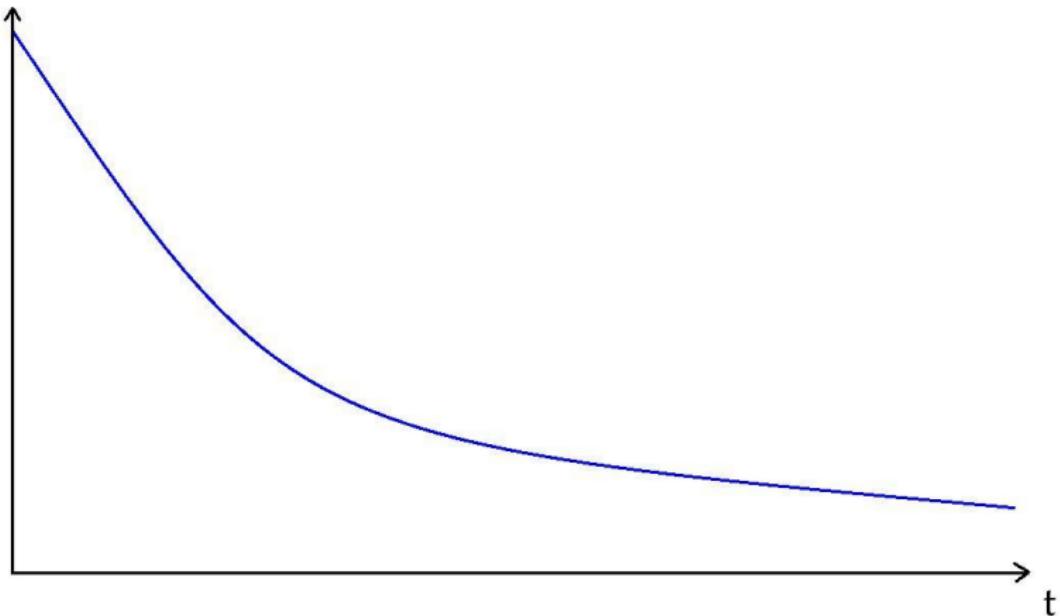
Car shock absorber



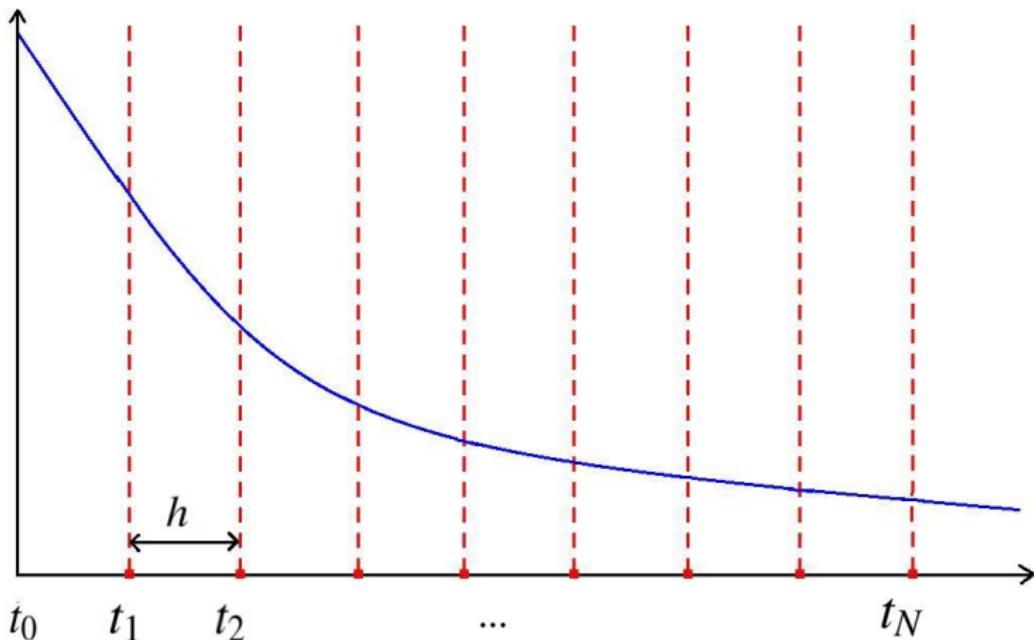
SpaceX Falcon 9 space rocket

Exact resolution is hard \Rightarrow numerical methods.

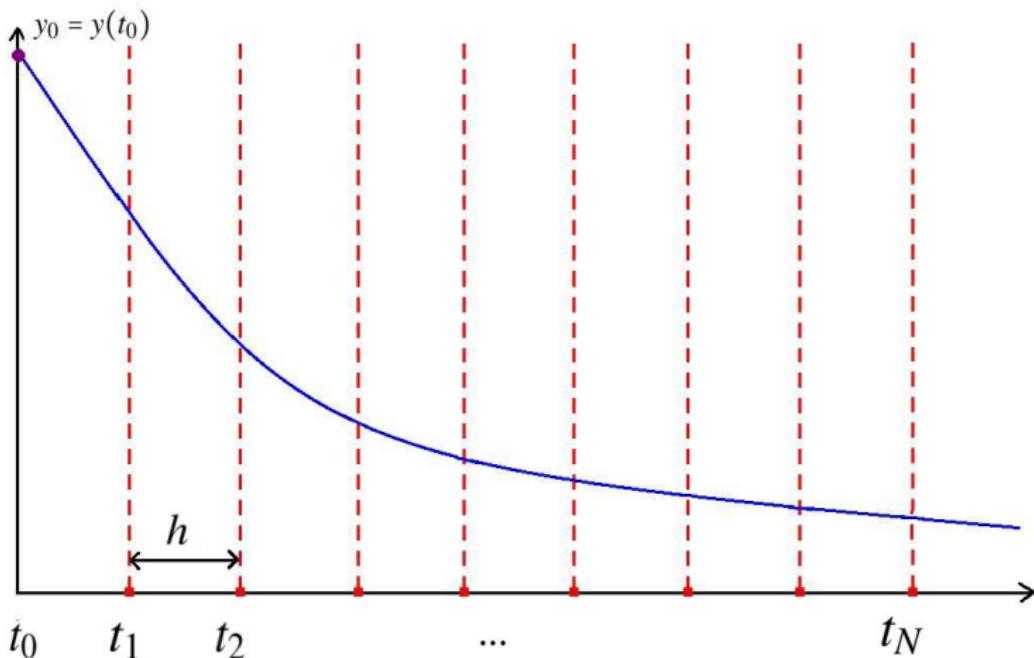
Numerical integration



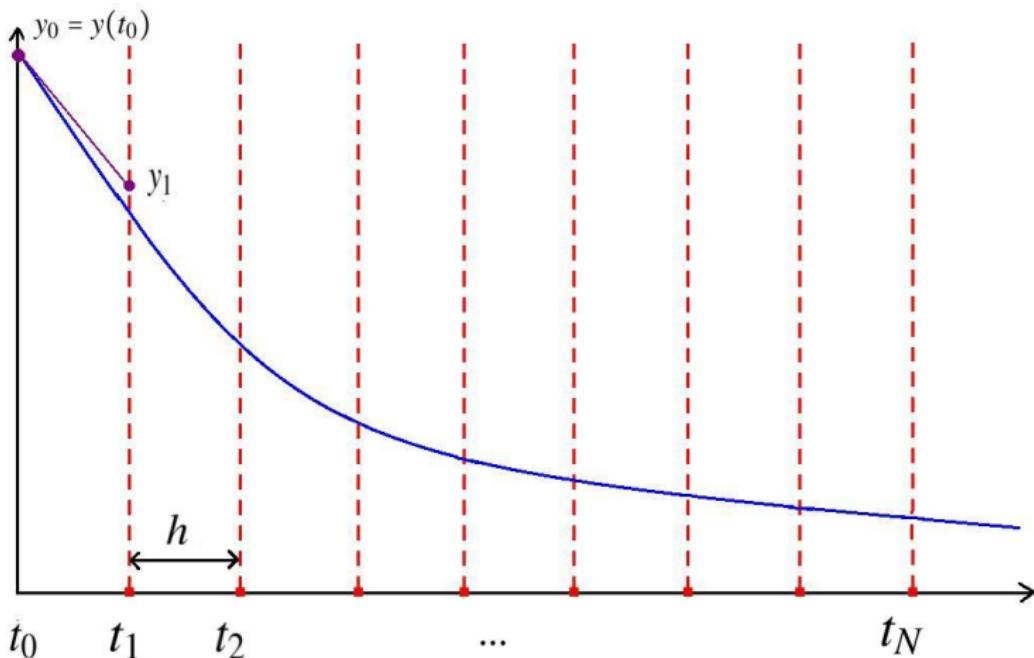
Numerical integration



Numerical integration

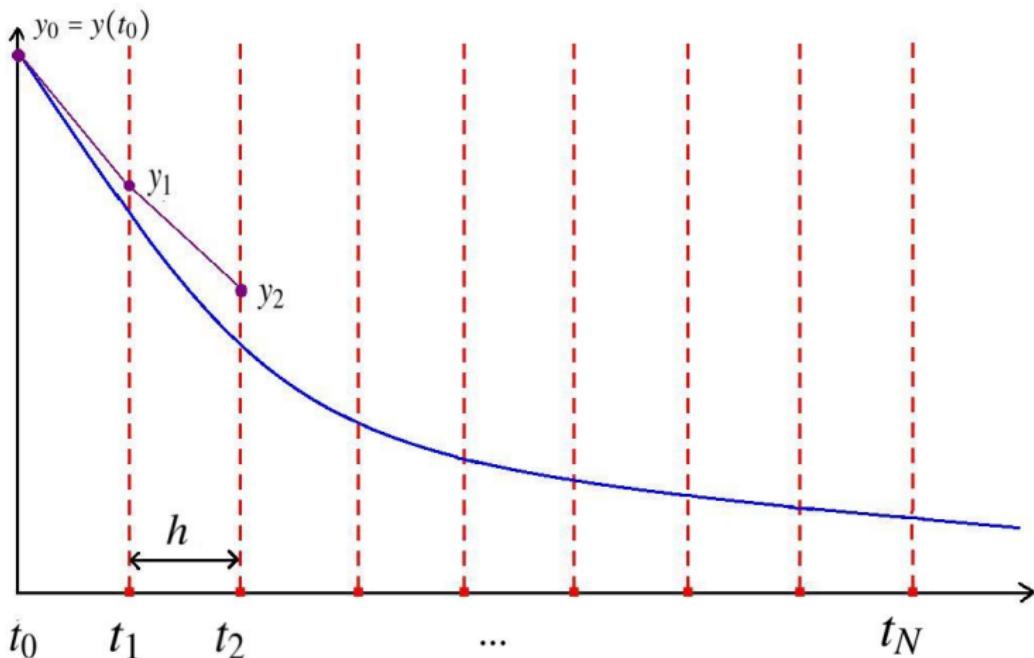


Numerical integration

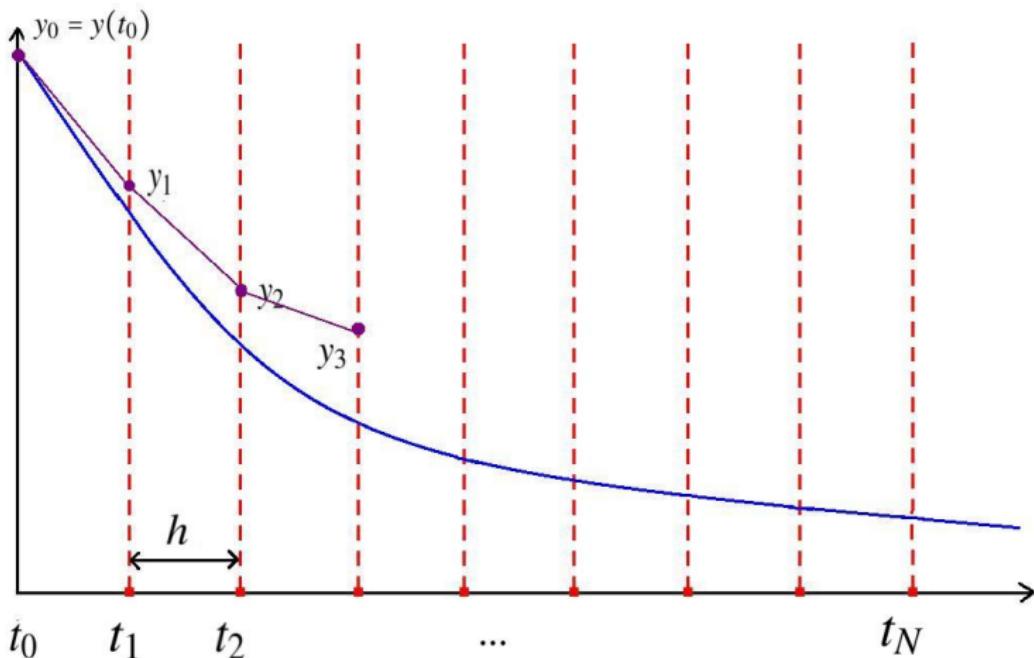


$N = 8$

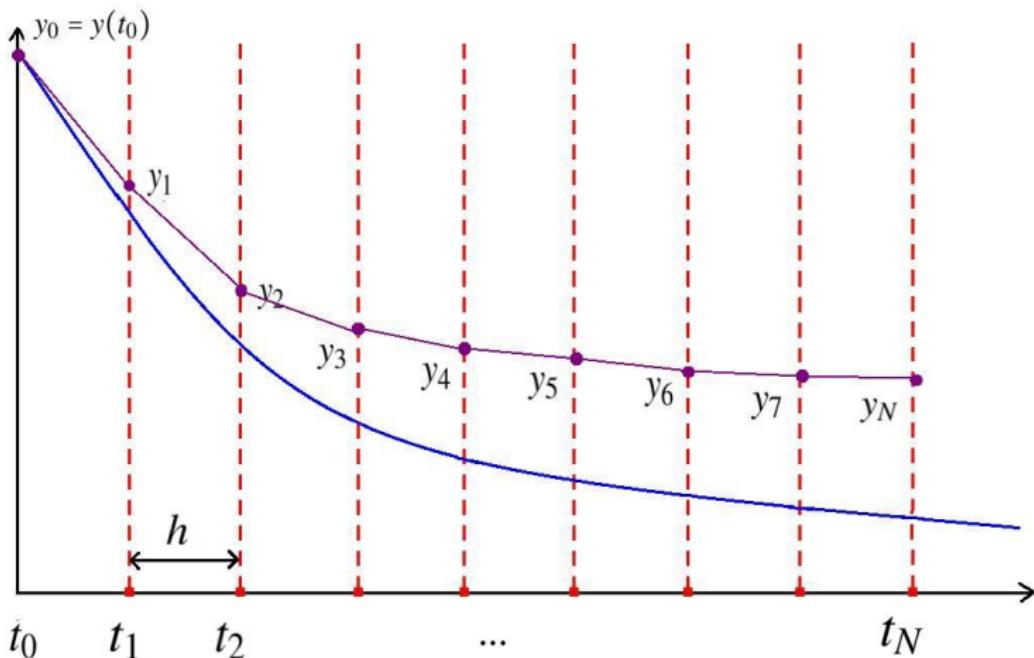
Numerical integration



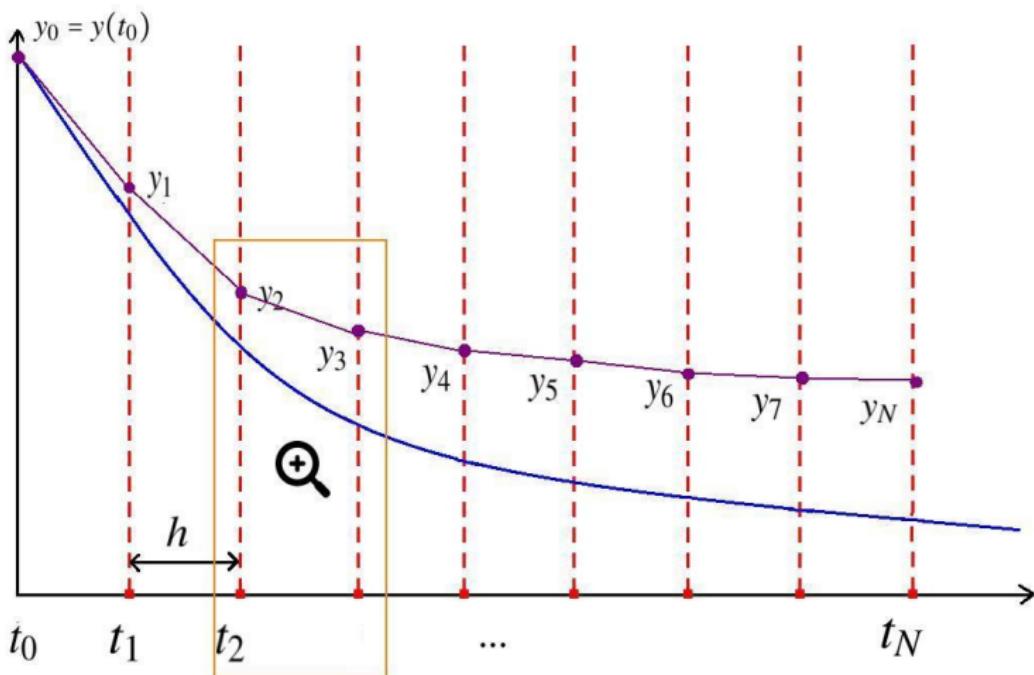
Numerical integration



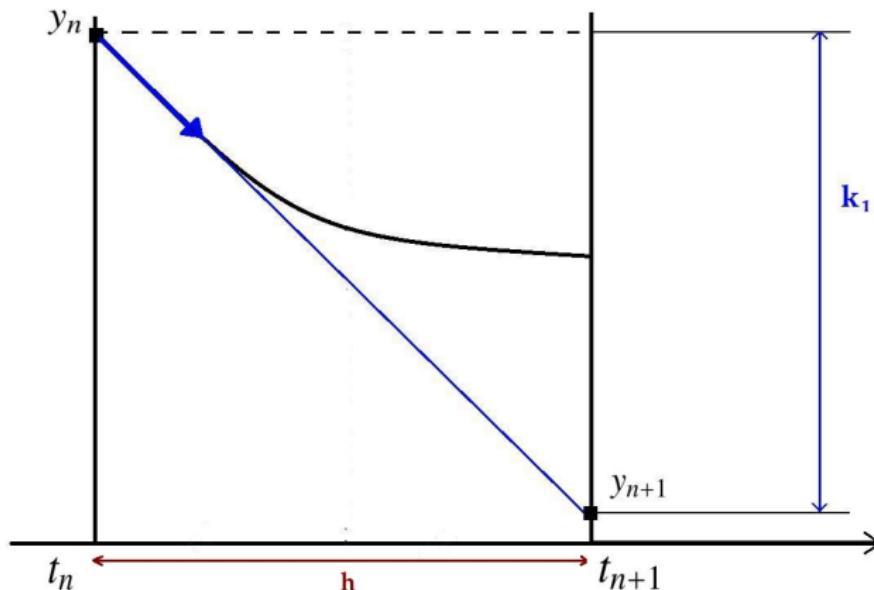
Numerical integration

 $N = 8$

Numerical integration

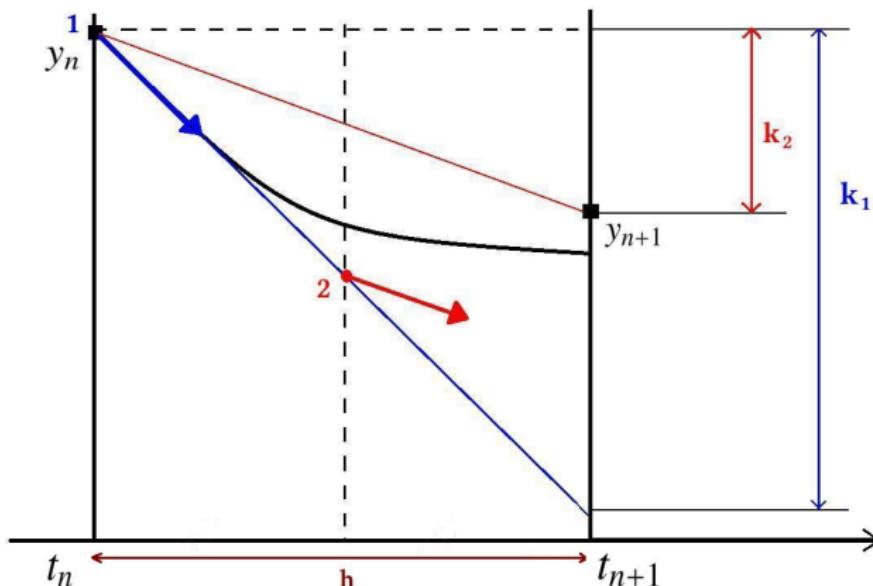
 $N = 8$

Euler method



$$k_1 = h f(t_n, y_n)$$
$$y_{n+1} = y_n + k_1 + O(h^2)$$

RK2 method



$$\begin{aligned}k_1 &= h \times f(t_n, y_n) & k_2 &= h \times f\left(t_n + \frac{h}{2}, y_n + \frac{k_1}{2}\right) \\y_{n+1} &= y_n + k_2 + O(h^3)\end{aligned}$$

Working assumptions

- FP arithmetic:
 - neither underflow nor overflow,
 - radix 2 double precision,

$$u = 2^{-53}$$

Working assumptions

- FP arithmetic;

- neither underflow nor overflow,
 - radix 2 double precision,

$$u = 2^{-53}$$

- ODEs;

- first-order,
 - linear,
 - $y : \mathbb{R} \rightarrow \mathbb{R}$,

$$y' = \lambda y$$

Working assumptions

- FP arithmetic;
 - neither underflow nor overflow,
 - radix 2 double precision, $u = 2^{-53}$
- ODEs;
 - first-order,
 - linear, $y' = \lambda y$
 - $y : \mathbb{R} \rightarrow \mathbb{R}$,
- Methods.
 - explicit,
 - one step,
 - constant step.

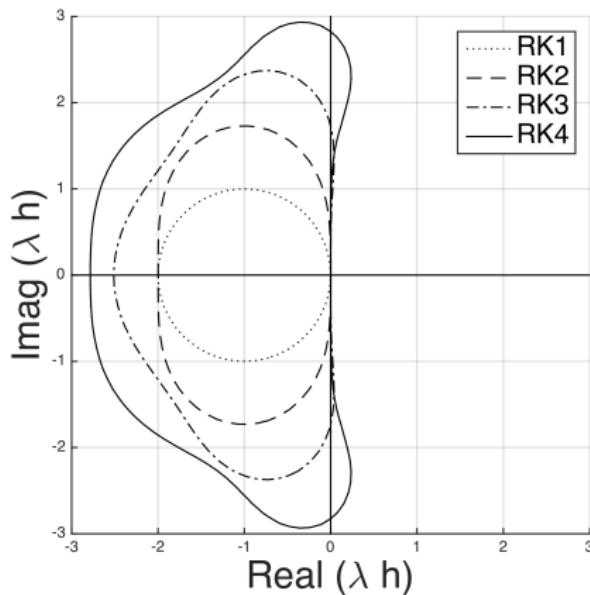
RK methods on linear problems: linear stability

$$\begin{cases} y_0 \in \mathbb{R} \\ y_{n+1} = R(h, \lambda) y_n \\ \text{(} R \text{ polynomial in } h\lambda \text{)} \end{cases}$$

RK methods on linear problems: linear stability

$$\begin{cases} y_0 \in \mathbb{R} \\ y_{n+1} = R(h, \lambda) y_n \\ (R \text{ polynomial in } h\lambda) \end{cases}$$

Stable $\Leftrightarrow |R(h, \lambda)| < 1$:



RK methods on linear problems: FP implementation

$$\begin{cases} y_0 \in \mathbb{R} \\ y_{n+1} = R(h, \lambda) y_n \end{cases}$$

$$\begin{cases} \tilde{y}_0 \simeq y_0 \\ \widetilde{y_{n+1}} = \widetilde{R}(\tilde{h}, \tilde{\lambda}, \tilde{y}_n) \end{cases}$$

RK methods on linear problems: FP implementation

$$\begin{cases} y_0 \in \mathbb{R} \\ y_{n+1} = R(h, \lambda) y_n \end{cases}$$

$$\begin{cases} \tilde{y}_0 \simeq y_0 \\ \widetilde{y_{n+1}} = \widetilde{R}(\tilde{h}, \tilde{\lambda}, \tilde{y}_n) \end{cases}$$

Euler:

- $R(h, \lambda) = 1 + h\lambda;$
- $\widetilde{R}(\tilde{h}, \tilde{\lambda}, \tilde{y}_n) = \tilde{y}_n \oplus \tilde{h} \otimes \tilde{\lambda} \otimes \tilde{y}_n.$

RK methods on linear problems: FP implementation

$$\begin{cases} y_0 \in \mathbb{R} \\ y_{n+1} = R(h, \lambda) y_n \end{cases}$$

$$\begin{cases} \tilde{y}_0 \simeq y_0 \\ \widetilde{y_{n+1}} = \widetilde{R}(\tilde{h}, \tilde{\lambda}, \tilde{y}_n) \end{cases}$$

Euler:

- $R(h, \lambda) = 1 + h\lambda;$
- $\widetilde{R}(\tilde{h}, \tilde{\lambda}, \tilde{y}_n) = \tilde{y}_n \oplus \tilde{h} \otimes \tilde{\lambda} \otimes \tilde{y}_n.$

RK4:

- $R(h, \lambda) = 1 + h\lambda + \frac{1}{2}(h\lambda)^2 + \frac{1}{6}(h\lambda)^3 + \frac{1}{24}(h\lambda)^4;$
- $\widetilde{R}(\tilde{h}, \tilde{\lambda}, \tilde{y}_n) =$

$$\tilde{y}_n \oplus \tilde{h} \otimes 6 \otimes \tilde{\lambda} \otimes \widetilde{y_n} \oplus \tilde{h} \otimes 3 \otimes \tilde{\lambda} \otimes \widetilde{y_n} \oplus \tilde{h} \otimes \tilde{h} \otimes 6 \otimes \tilde{\lambda} \otimes \tilde{\lambda} \otimes \widetilde{y_n} \oplus \tilde{h} \otimes 3 \otimes \tilde{\lambda} \otimes \widetilde{y_n} \otimes$$

$$\tilde{h} \otimes \tilde{h} \otimes 6 \otimes \tilde{\lambda} \otimes \widetilde{y_n} \oplus \tilde{h} \otimes \tilde{h} \otimes \tilde{h} \otimes 12 \otimes \tilde{\lambda} \otimes \tilde{\lambda} \otimes \tilde{\lambda} \otimes \widetilde{y_n} \oplus \tilde{h} \otimes 6 \otimes \tilde{\lambda} \otimes \widetilde{y_n}$$

$$\oplus \tilde{h} \otimes \tilde{h} \otimes 6 \otimes \tilde{\lambda} \otimes \widetilde{y_n} \oplus \tilde{h} \otimes \tilde{h} \otimes \tilde{h} \otimes 12 \otimes \tilde{\lambda} \otimes \tilde{\lambda} \otimes \widetilde{y_n}$$

$$\oplus \tilde{h} \otimes \tilde{h} \otimes \tilde{h} \otimes \tilde{h} \otimes 24 \otimes \tilde{\lambda} \otimes \tilde{\lambda} \otimes \tilde{\lambda} \otimes \widetilde{y_n}.$$

(> 60 flops!)

State-of-the-art

Roundoff errors in numerical methods (N = nb of iterations):

- probabilistic result: error in \sqrt{N} [Henrici, 1963];
- in practice (implicit RK): error in N [Hairer&al, 2008];
- interval analysis [Bouissou-Martel, 2006];
- numerical integration (fine-grained): Newton-Cotes, Gauss-Legendre, ... [Fousse, 2006].

State-of-the-art

Roundoff errors in numerical methods (N = nb of iterations):

- probabilistic result: error in \sqrt{N} [Henrici, 1963];
- in practice (implicit RK): error in N [Hairer&al, 2008];
- interval analysis [Bouissou-Martel, 2006];
- numerical integration (fine-grained): Newton-Cotes, Gauss-Legendre, ... [Fousse, 2006].

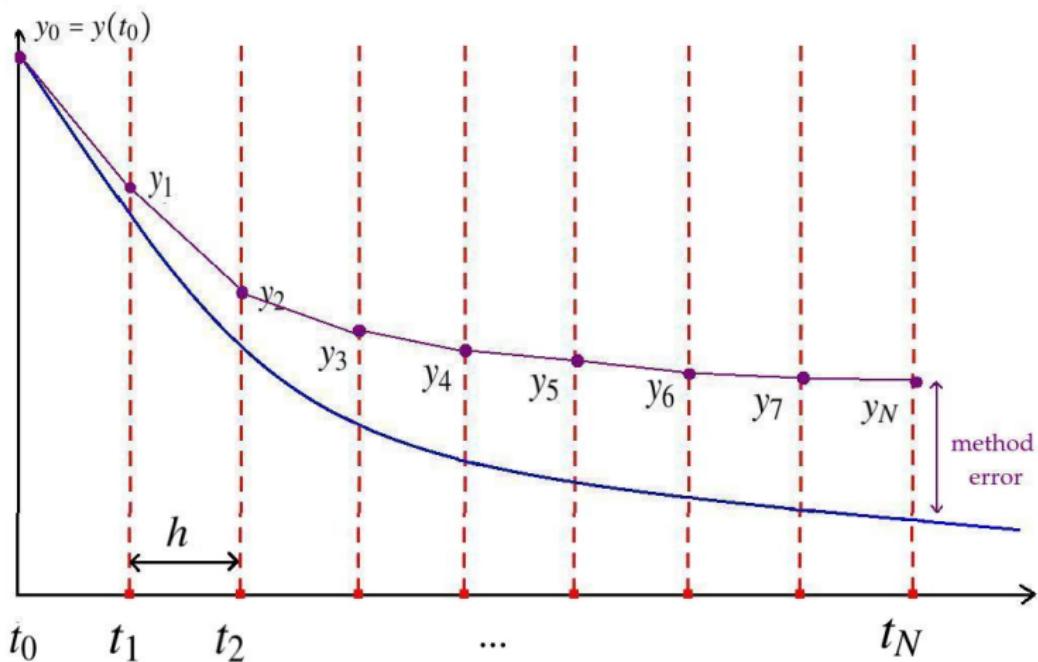
Our approach:

- fine-grained analysis;
- use of mathematical properties of the methods (stability).

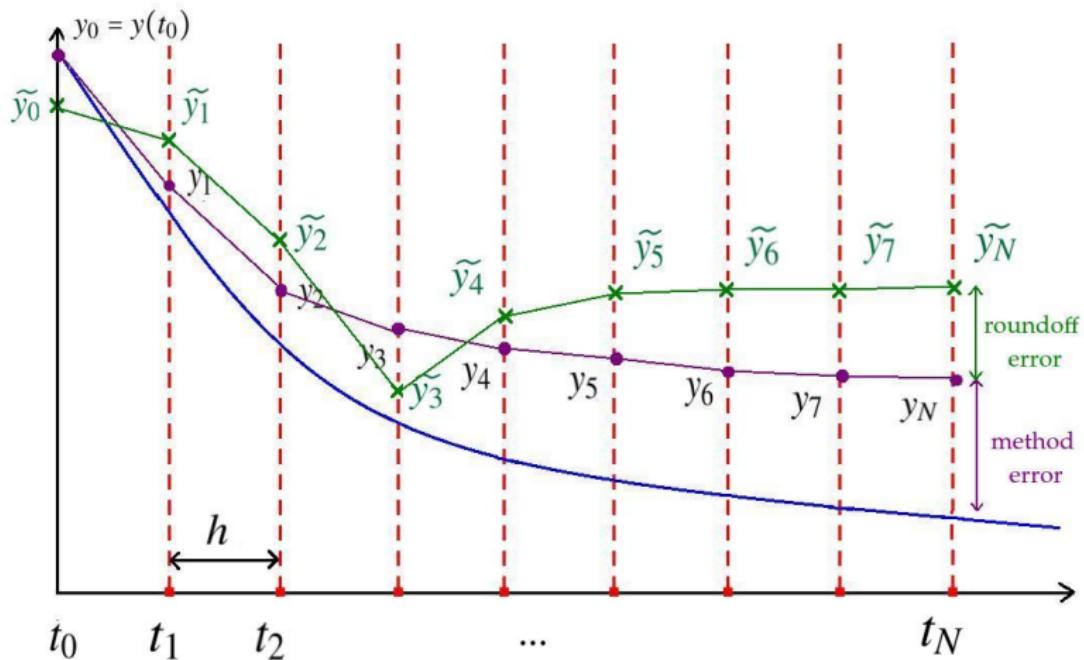
Table of contents

- 1 Motivations and numerical methods
- 2 Roundoff errors of RK methods
 - Local roundoff errors
 - Global roundoff errors of classical methods
- 3 Conclusion and perspectives

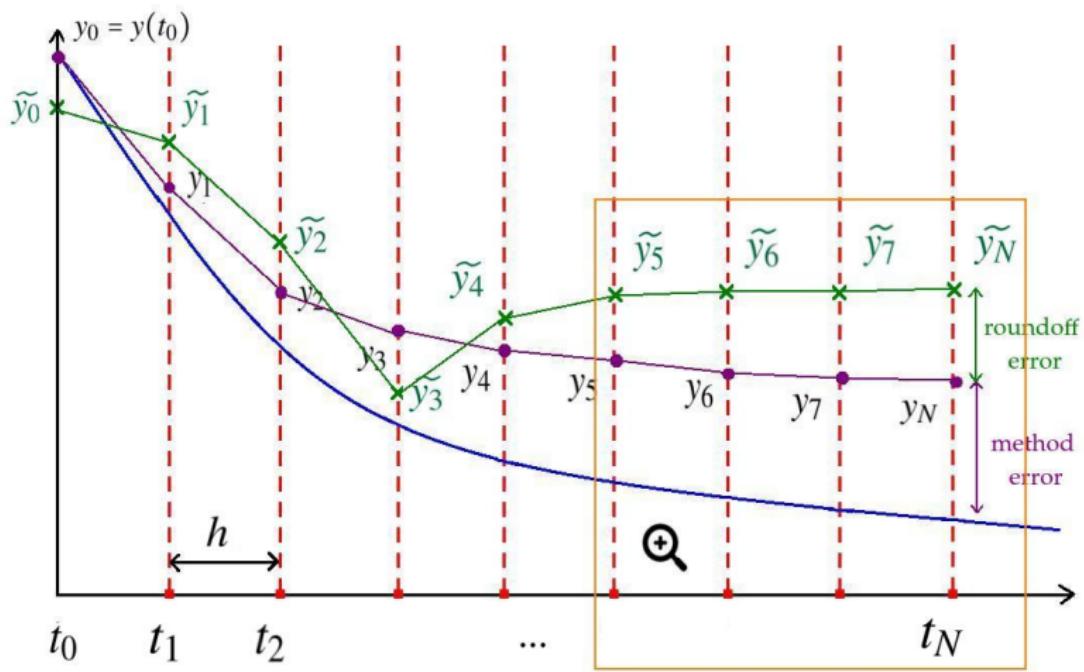
Method error vs roundoff error

 $N = 8$

Method error vs roundoff error

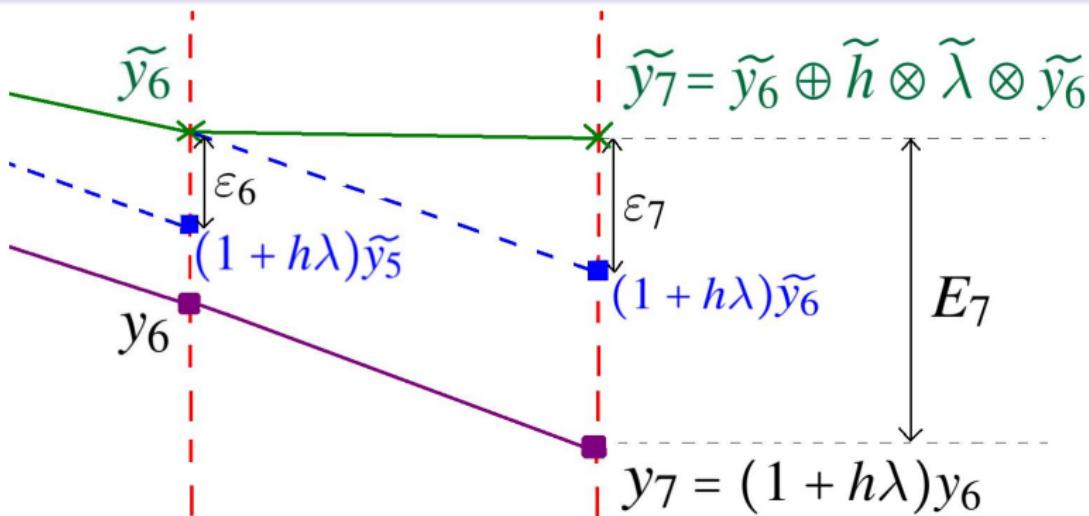
 $N = 8$

Method error vs roundoff error



$N = 8$

Local roundoff error vs global roundoff error



Local error:

$$\varepsilon_0 = |\tilde{y}_0 - y_0|$$

$$\forall n \in \mathbb{N}^*, \varepsilon_n = |\widetilde{R}(\tilde{h}, \tilde{\lambda}, \widetilde{y_{n-1}}) - R(h, \lambda) \widetilde{y_{n-1}}|.$$

Global error:

$$\forall n \in \mathbb{N}, E_n = \tilde{y}_n - y_n.$$

From local to global roundoff error

Local error:

$$\varepsilon_0 = |\tilde{y}_0 - y_0|$$

$$\forall n \in \mathbb{N}^*, \varepsilon_n = |\widetilde{R}(\tilde{h}, \tilde{\lambda}, \widetilde{y_{n-1}}) - R(h, \lambda) \widetilde{y_{n-1}}|.$$

Global error:

$$\forall n \in \mathbb{N}, E_n = \tilde{y}_n - y_n.$$

Theorem 1: Global absolute error of RK methods

Let $C \in \mathbb{R}_+^*$. Suppose $\forall n \in \mathbb{N}^*, \varepsilon_n \leq C|\widetilde{y_{n-1}}|$. Then, $\forall n \in \mathbb{N}$,

$$|E_n| \leq (C + |R(h, \lambda)|)^n \left(\varepsilon_0 + n \frac{C|y_0|}{C + |R(h, \lambda)|} \right).$$

Relative roundoff errors

Relative error:

$$\left| \frac{\tilde{y}_n - y_n}{y_n} \right| \leq \left(\frac{C + |R(h, \lambda)|}{|R(h, \lambda)|} \right)^n \left(\frac{\varepsilon_0}{|y_0|} + n \frac{C}{C + |R(h, \lambda)|} \right).$$

Relative roundoff errors

Relative error:

$$\left| \frac{\tilde{y}_n - y_n}{y_n} \right| \leq \left(\frac{C + |R(h, \lambda)|}{|R(h, \lambda)|} \right)^n \left(\frac{\varepsilon_0}{|y_0|} + n \frac{C}{C + |R(h, \lambda)|} \right).$$

If $C \ll |R(h, \lambda)|$, then:

$$\left| \frac{\tilde{y}_n - y_n}{y_n} \right| \lesssim \frac{\varepsilon_0}{|y_0|} + n \frac{C}{|R(h, \lambda)|}.$$

In practice (Euler, RK2, RK4): $C \leq 200u$ and $200u \ll |R(h, \lambda)|$.

Table of contents

- 1 Motivations and numerical methods**
- 2 Roundoff errors of RK methods**
 - Local roundoff errors
 - Global roundoff errors of classical methods
- 3 Conclusion and perspectives**

Technical lemma for local roundoff errors

Local error of **stable** Euler's method ($-2 \leq h\lambda < 0$):

$$\varepsilon_{n+1} = |\tilde{y}_n \oplus (\tilde{h} \otimes \tilde{\lambda} \otimes \tilde{y}_n) - (1 + h\lambda)y_n|.$$

Technical lemma for local roundoff errors

Local error of **stable** Euler's method ($-2 \leq h\lambda < 0$):

$$\varepsilon_{n+1} = \left| \underbrace{\widetilde{y}_n}_{X_1} \oplus \left(\underbrace{\widetilde{h} \otimes \widetilde{\lambda}}_{X_2} \otimes \underbrace{\widetilde{y}_n}_y \right) - \left(\underbrace{1}_{\alpha_1} + \underbrace{h\lambda}_{\alpha_2} \right) \underbrace{\widetilde{y}_n}_y \right|.$$

Lemma 2

Let $y \in \mathbb{R}$. Let $C_1, C_2 \in \mathbb{R}_+$. Let $\alpha_1, \alpha_2 \in \mathbb{R}$. Let $X_1, X_2 \in \mathbb{F}$ s.t.
 $|X_1 - \alpha_1 y| \leq C_1 |y|$ and $|X_2 - \alpha_2| \leq C_2$.

Then:

$$\begin{aligned} & |X_1 \oplus (X_2 \otimes y) - (\alpha_1 + \alpha_2)y| \\ & \leq |y| (C_1 + C_2 + u (|\alpha_1| + 2|\alpha_2| + C_1 + 2C_2) + u^2 (C_2 + |\alpha_2|)). \end{aligned}$$

Technical lemma for local roundoff errors

Local error of **stable** Euler's method ($-2 \leq h\lambda < 0$):

$$\varepsilon_{n+1} = \left| \underbrace{\widetilde{y}_n}_{X_1} \oplus \left(\underbrace{\widetilde{h} \otimes \widetilde{\lambda}}_{X_2} \otimes \underbrace{\widetilde{y}_n}_y \right) - \left(\underbrace{1}_{\alpha_1} + \underbrace{h\lambda}_{\alpha_2} \right) \underbrace{\widetilde{y}_n}_y \right|.$$

Lemma 2

Let $y \in \mathbb{R}$. Let $C_1, C_2 \in \mathbb{R}_+$. Let $\alpha_1, \alpha_2 \in \mathbb{R}$. Let $X_1, X_2 \in \mathbb{F}$ s.t.
 $|X_1 - \alpha_1 y| \leq C_1 |y|$ and $|X_2 - \alpha_2| \leq C_2$.

Then:

$$\begin{aligned} & |X_1 \oplus (X_2 \otimes y) - (\alpha_1 + \alpha_2)y| \\ & \leq |y| (C_1 + C_2 + u(|\alpha_1| + 2|\alpha_2| + C_1 + 2C_2) + u^2 (C_2 + |\alpha_2|)). \end{aligned}$$

$$\left| \underbrace{\widetilde{y}_n}_{X_1} - \underbrace{\frac{1}{\alpha_1} \underbrace{\widetilde{y}_n}_y}_{C_1} \right| \leq \underbrace{0}_{C_1} \left| \underbrace{\widetilde{y}_n}_y \right|, \quad \left| \underbrace{\widetilde{h} \otimes \widetilde{\lambda}}_{X_2} - \underbrace{h\lambda}_{\alpha_2} \right| \leq \underbrace{6u}_{C_2}$$

(Gappa [Melquiond]).

Local roundoff errors of higher-order methods

Stable RK4 method:

	C
\tilde{y}_n	0
$\circ \left[\tilde{h} \frac{1}{6} \tilde{\lambda} \right]$	$2u$ (Gappa)
$\circ \left[\tilde{y}_n + \tilde{h} \frac{1}{6} \tilde{\lambda} \tilde{y}_n \right] = \tilde{y}_n \oplus \tilde{h} \oslash 6 \otimes \tilde{\lambda} \otimes \tilde{y}_n$	$4u$ (Lemma 2)
$\circ \left[\tilde{h} \frac{1}{3} \tilde{\lambda} \right]$	$4u$ (Gappa)
$\circ \left[\tilde{h} \tilde{h} \frac{1}{6} \tilde{\lambda} \tilde{\lambda} \right]$	$12u$ (Gappa)
$\circ \left[\tilde{h} \tilde{h} \tilde{h} \frac{1}{12} \tilde{\lambda} \tilde{\lambda} \tilde{\lambda} \right]$	$28u$ (Gappa)
$\circ \left[\tilde{h} \tilde{h} \tilde{h} \tilde{h} \frac{1}{24} \tilde{\lambda} \tilde{\lambda} \tilde{\lambda} \tilde{\lambda} \right]$	$53u$ (Gappa)
...	... ($7 \times$ Lemma 2)
$\circ \left[\tilde{y}_n + \dots + \tilde{h} \tilde{h} \tilde{h} \tilde{h} \frac{1}{24} \tilde{\lambda} \tilde{\lambda} \tilde{\lambda} \tilde{\lambda} \tilde{y}_n \right]$	$194u$ (Lemma 2)

Local roundoff errors of classical methods

Lemma 3: Local error of Euler method

Suppose $-2 \leq h\lambda \leq -2^{-100}$ and $2^{-60} \leq h \leq 1$. Then:

$$\forall n \in \mathbb{N}, \quad \varepsilon_{n+1}^{Euler} \leq 11.01u |\tilde{y}_n|.$$

Lemma 4: Local error of RK2 method

Suppose $-2 \leq h\lambda \leq -2^{-100}$ and $2^{-60} \leq h \leq 1$. Then:

$$\forall n \in \mathbb{N}, \quad \varepsilon_{n+1}^{RK2} \leq 28.01u |\tilde{y}_n|.$$

Lemma 5: Local error of RK4 method

Suppose $-3 \leq h\lambda \leq -2^{-100}$ and $2^{-60} \leq h \leq 1$. Then:

$$\forall n \in \mathbb{N}, \quad \varepsilon_{n+1}^{RK4} \leq 194u |\tilde{y}_n|.$$

Table of contents

- 1 Motivations and numerical methods
- 2 Roundoff errors of RK methods
 - Local roundoff errors
 - Global roundoff errors of classical methods
- 3 Conclusion and perspectives

Global roundoff errors of classical methods

Bounds on **global errors**:

- **Euler:** $C = 11.01u$

$$|E_n| \leq (11.01u + |R(h, \lambda)|)^n \left(\varepsilon_0 + n \frac{11.01u|y_0|}{11.01u + |R(h, \lambda)|} \right);$$

- **RK2:** $C = 28.01u$

$$|E_n| \leq (28.01u + |R(h, \lambda)|)^n \left(\varepsilon_0 + n \frac{28.01u|y_0|}{28.01u + |R(h, \lambda)|} \right);$$

- **RK4:** $C = 194u$

$$|E_n| \leq (194u + |R(h, \lambda)|)^n \left(\varepsilon_0 + n \frac{194u|y_0|}{194u + |R(h, \lambda)|} \right).$$

⇒ no compensation ☹ but reasonable bounds ☺.

Table of contents

- 1 Motivations and numerical methods
- 2 Roundoff errors of RK methods
 - Local roundoff errors
 - Global roundoff errors of classical methods
- 3 Conclusion and perspectives

General methodology

$$y_{n+1} = \sum_{i=0}^M \alpha_i y_n \quad \widetilde{y_{n+1}} = \oplus_{i=0}^M \widetilde{\alpha}_i \otimes \widetilde{y_n} \quad (\alpha_i = h\lambda, \frac{h\lambda}{3}, \frac{h\lambda}{6}, \frac{h^4\lambda^4}{24} \dots)$$

General methodology

$$y_{n+1} = \sum_{i=0}^M \alpha_i y_n \quad \widetilde{y_{n+1}} = \oplus_{i=0}^M \widetilde{\alpha_i} \otimes \widetilde{y_n} \quad (\alpha_i = h\lambda, \frac{h\lambda}{3}, \frac{h\lambda}{6}, \frac{h^4\lambda^4}{24} \dots)$$

Methodology to bound roundoff errors:

General methodology

$$y_{n+1} = \sum_{i=0}^M \alpha_i y_n \quad \widetilde{y_{n+1}} = \oplus_{i=0}^M \widetilde{\alpha_i} \otimes \widetilde{y_n} \quad (\alpha_i = h\lambda, \frac{h\lambda}{3}, \frac{h\lambda}{6}, \frac{h^4\lambda^4}{24} \dots)$$

Methodology to bound roundoff errors:

- 1) bound $\varepsilon_0 = |\widetilde{y}_0 - y_0|$;

General methodology

$$y_{n+1} = \sum_{i=0}^M \alpha_i y_n \quad \widetilde{y_{n+1}} = \bigoplus_{i=0}^M \widetilde{\alpha_i} \otimes \widetilde{y_n} \quad (\alpha_i = h\lambda, \frac{h\lambda}{3}, \frac{h\lambda}{6}, \frac{h^4\lambda^4}{24} \dots)$$

Methodology to bound roundoff errors:

- 1) bound $\varepsilon_0 = |\widetilde{y}_0 - y_0|$;
- 2) bound the error on each term α_i (Gappa + stability);

General methodology

$$y_{n+1} = \sum_{i=0}^M \alpha_i y_n \quad \widetilde{y_{n+1}} = \bigoplus_{i=0}^M \widetilde{\alpha_i} \otimes \widetilde{y_n} \quad (\alpha_i = h\lambda, \frac{h\lambda}{3}, \frac{h\lambda}{6}, \frac{h^4\lambda^4}{24} \dots)$$

Methodology to bound roundoff errors:

- 1) bound $\varepsilon_0 = |\widetilde{y}_0 - y_0|$;
- 2) bound the error on each term α_i (Gappa + stability);
- 3) bound local errors by M applications of Lemma 2;

General methodology

$$y_{n+1} = \sum_{i=0}^M \alpha_i y_n \quad \widetilde{y_{n+1}} = \bigoplus_{i=0}^M \widetilde{\alpha_i} \otimes \widetilde{y_n} \quad (\alpha_i = h\lambda, \frac{h\lambda}{3}, \frac{h\lambda}{6}, \frac{h^4\lambda^4}{24} \dots)$$

Methodology to bound roundoff errors:

- 1) bound $\varepsilon_0 = |\widetilde{y}_0 - y_0|$;
- 2) bound the error on each term α_i (Gappa + stability);
- 3) bound local errors by M applications of Lemma 2;
- 4) bound the global error by instantiating Theorem 1.

Conclusion and perspectives

Conclusion:

- fine and mechanical analysis of roundoff errors;
- results on useful and classical methods (Euler, RK2, RK4);
- linear growth of roundoff errors;
- no compensation (for explicit one-step methods).

Perspectives:

- overflows and underflows;
- formalization in Coq (proof assistant):
 - based on the Flocq library [Boldo-Melquiond],
 - based on the gappa and interval tactics [Melquiond],
- more general ODEs: complex, matricial, non-linear;
- more general methods: multi-step, variable step, implicit.