Ecole thématique PRECIS 2017, Fréjus, 15-19 mai 2017

Précision et Reproductibilité en Calcul et Informatique Scientifique : garantir la reproductibilité numérique concilier précision et performance

> **Philippe Langlois, David Parello** DALI, University of Perpignan Via Domitia LIRMM, UMR 5506 CNRS-UM, France









Acknowledgement and collaboration

Joint work with (past or current) members of the DALI research team in University of Perpignan:

Prof. B. Goossens, Dr. N. Louvet, Dr. D. Parello, Ch. Chohra, R. Nheili, C. Porada,

and colleagues:

Ch. Denis (EDF R&D, CMLA), J.-M. Hervouet (EDF R&D). Sources of errors when computing the solution of a scientific problem in floating point arithmetic:

- mathematical model,
- truncation errors,
- data uncertainties,
- rounding errors.

Sources of errors when computing the solution of a scientific problem in floating point arithmetic:

- mathematical model,
- truncation errors,
- data uncertainties,
- rounding errors.

Rounding errors may totally corrupt a floating point computation:

- accumulation of billions of floating point operations,
- intrinsic difficulty to solve the problem accurately.

Example: polynomial evaluation

Evaluation of univariate polynomials with floating point coefficients:

- the evaluation of a polynomial suffers from rounding errors
- example : in the neighborhood of a multiple root



 $p(x) = (x - 2)^9$ in expanded form near the multiple root x = 2

Example: polynomial evaluation

Evaluation of univariate polynomials with floating point coefficients:

- the evaluation of a polynomial suffers from rounding errors
- example : in the neighborhood of a multiple root



 $p(x) = (x - 2)^9$ in expanded form near the multiple root x = 2evaluated with the Horner algorithm in IEEE double precision.

How to:

- improve and validate the accuracy of a floating point computation,
- recover the numerical reproducibility of parallel floating-point computation,
- without large computing time overheads ?
 - More hardware precision
 - Software simulation of more computing precision
 - More accurate algorithms
 - As fast as possible accurate algorithms
 - Parallelism is everywhere
 - Reproducibility at least to debug!

Main issues today

Starting point: IEEE-754 floating point arithmetic

- Best possible accuracy for $+, -, \times, /, \sqrt{}$
- Add is not associative

Summing n floating numbers : focus on accuracy

- Core computation, numerous algos, recently some really smart ones
- Computed sum accuracy : doubled or more, faithful or correctly rounded

Summing n floating numbers : focus on running-time and memory print

- Running-time and memory print are discriminant factors when computing the best possible accurate sum
- Appreciating the actual performance of one algo is not an easy task : flop/s ? hardware counters ? compiler options ?

Reliable and significant measure of the time complexity?

The classic way: count the number of flop

- A usual problem: double the accuracy of a computed result
- A usual answer for polynomial evaluation (degree *n*)

Metric	Horner	CompHorner	DDHorner
Flop count	2n	22n + 5	28 <i>n</i> + 4
Flop count ratio	1	pprox 11	pprox 14
Measured #cycles ratio	1	2.8 - 3.2	8.7 – 9.7

Flop count vs. run-time measures

- Flop counts and measured run-times are not proportional
- Run-time measure is a very difficult experimental process
- Which one trust?

Reproducible HPC numerical simulations?

Numerical reproducibility of parallel computation

• Getting bitwise identical results for every *p*-parallel run, $p \ge 1$

One industrial scale simulation code

- Simulation of free-surface flows in 1D-2D-3D hydrodynamics
- 300 000 loc. of open source Fortran 90
- 20 years, 4000 registered users, EDF R&D + international consortium

Telemac 2D [5]

- 2D hydrodynamic: Saint Venant equations
- Finite element method, triangular element mesh, sub-domain decomposition for parallel resolution
- \bullet Mesh node unknowns: water depth (H) and velocity (U,V)

A complex Telemac 2D simulation

The Malpasset dam break (1959)

- A five year old dam break: 433 dead people and huge damage
- Triangular mesh: 26000 elements and 53000 nodes
- \bullet Simulation: $\rightarrow 35 \text{min.}$ after break with a 2sec. time step











How to trust this complex simulation?



A reproducible simulation?

	velocity U	velocity V	depth H
The sequential run	0.4029747E-02	0.7570773E-02	0.3500122E-01
one <mark>64</mark> procs run	0.4 <mark>935279</mark> E-02	0.3422730E-02	0.2748817E-01
one 128 procs run	0.4 <mark>512116</mark> E-02	0.75 <mark>45233</mark> E-02	0.1327634E-01

Reproducibility failure

- Up to $\times 2.5$ uncertainty
- The privileged sequential run?

Telemac2D: the simplest gouttedo simulation

The gouttedo simulation test case

- 2D-simulation of a water drop fall in a square basin
- Unknown: water depth for a 0.2 sec time step
- Triangular mesh: 8978 elements and 4624 nodes



Numerical reproducibility?



Numerical reproducibility?



Numerical reproducibility?



Numerical reproducibility?



Numerical reproducibility?



Numerical reproducibility?



Numerical reproducibility?



Numerical reproducibility?



Numerical reproducibility?



Numerical reproducibility?



Numerical reproducibility?



Numerical reproducibility?



Numerical reproducibility?



Numerical reproducibility?



NO numerical reproducibility!



Telemac2D: gouttedo

NO numerical reproducibility!



Part 1: More accuracy

D Summing *n* floating-point numbers: basic steps

- Some algorithms
- Some accuracy bounds
- 2 Computing sums more accurately
 - Some famous old and magic algorithms
 - Compensation
 - Distillate to understand and go further

Omputing sums as accurately as possible

- iFastSum: distillation again
- Accsum
- Hybridsum and OnLineExactSum
- More EFT and accurate algorithms
 - Some more results



Part 2: More reproducibility

6 Motivation

Reproducibility failure in a finite element simulation

- Sequential and parallel FE assembly
- Sources of non reproducibility in Telemac2D

Recovering reproducibility

- Reproducible parallel FE assembly
- A first application: reproducible Tomawac
- Reproducible algebraic operations
- Reproducible conjugate gradient
- Reproducible Telemac2D

Efficiency



Part 3: More performance

Part I

More accuracy

Summing n floating-point numbers: basic steps

- Some algorithms
- Some accuracy bounds

2 Computing sums more accurately

- Some famous old and magic algorithms
- Compensation
- Distillate to understand and go further

Omputing sums as accurately as possible

- iFastSum: distillation again
- Accsum
- Hybridsum and OnLineExactSum
- More EFT and accurate algorithms
 - Some more results

More accuracy: conclusion

How to manage accuracy and speed?

So many ways ... and a new "best one" every year since 1999 1965 Møller, Ross 1991 Priest 1969 Babuska, Knuth 1992 Clarkson, Priest 1970 Nickel 1993 Higham 1971 Dekker, Malcolm 1997 Shewchuk 1972 Kahan, Pichat 1999 Anderson 1974 Neumaier 2001 Hlavacs/Uberhuber 1975 Kulisch/Bohlender 2002 Li et al. (XBLAS) 1977 Bohlender, Mosteller/Tukey 2003 Demmel/Hida, Nievergelt, 1981 Linnaimaa Zielke/Drygalla 1982 Leuprecht/Oberaigner 2005 Ogita/Rump/Oishi, 1983 Jankowski/Semoktunowicz/-Zhu/Yong/Zeng Wozniakowski 2006 Zhu/Hayes 1985 Jankowski/Wozniakowski 2008 Rump/Ogita/Oishi 1987 Kahan 2009 Rump, Zhu/Hayes 2010 Zhu/Hayes

IEEE-754 floating point arithmetic

Floating-point numbers (normal, non zero)



• The standard model:

 $\mathsf{fl}(a \circ b) = (1 + \varepsilon)(a \circ b), \text{ with } |\varepsilon| \leq \mathbf{u} = 2^{-p}, \text{ or } 2^{1-p}.$

IEEE-754 (1985, 2008)

• formats, rounding modes : +, -, ×, /, \surd are as accurate as possible, exceptions

•
$$u = 2^{-53} \approx 10^{-16}$$
 for b64 in IEEE-754
Using the standard model

Summing *n* floating-point numbers: basic steps

- Some algorithms
- Some accuracy bounds

2 Computing sums more accurately

- Some famous old and magic algorithms
- Compensation
- Distillate to understand and go further

Omputing sums as accurately as possible

- iFastSum: distillation again
- Accsum
- Hybridsum and OnLineExactSum
- More EFT and accurate algorithms
 - Some more results



Floating-point summation: classic results [HFPA,09],[ASNA,06]

Accuracy for backward stable algorithms

- Accuracy of the computed sum $\leq (n-1) imes \mathit{cond} imes \mathbf{u}$
- cond $(\sum x_i) = \sum |x_i|/|\sum x_i|$
- No more significant digit in IEEE-b64 for large cond, *i.e.* > 10¹⁶
- The length also matters for large n



Backward stable summation algorithms

Algorithm (A describes a class of summation algorithms)

Let $S = \{x_1, x_2, ..., x_n\}.$

while ${\mathcal S}$ contains more than one element do

```
Remove two numbers x and y from S and add x \oplus y to S;
```

end while

Return the remaining element of S.

- Recursive summation...
- ... with increasing order sorting (IOS): $|x_1| \le |x_2| \le \cdots \le |x_n|$,
- ... with decreasing order (DOS),
- insertion summation with IOS,
- pairwise summation.

Proposition (Another bound)

 $|\widehat{s_n} - s_n| \leq \mathbf{u} \sum_{i=1}^{n-1} |T_i|$, where T_i is the i-th partial sum in \mathcal{A} .

- Minimize this bound or at least every $|T_i|$.
- All the x_i have the same sign: recursive with IOS (good), insertion with IOS (the best).
- Cancellation when $\sum |x_i| \gg |\sum x_i|$: DOS better than IOS.

Conclusion of the basic steps



It exists no universally better accurate summation algorithm



Computing sums more accurately: step 1

- Summing *n* floating-point numbers: basic steps
 - Some algorithms
 - Some accuracy bounds
- 2 Computing sums more accurately
 - Some famous old and magic algorithms
 - Compensation
 - Distillate to understand and go further
- Computing sums as accurately as possible
 - iFastSum: distillation again
 - Accsum
 - Hybridsum and OnLineExactSum
- 4 More EFT and accurate algorithms
 - Some more results



More accuracy but still conditioning: overview

Arbitrary precision but accuracy still condition dependant

- Simulating more *precision*: double-double, quad-double, ..., MPFR
- Compensated algorithms: Kahan(65), ..., Sum2(05), SumK(05)
- Accuracy $\leq \mathbf{u} + cond \times \mathbf{u}^{K}$
- No more *n* dependency while $n\mathbf{u} \leq 1$.

T. OGITA, S. M. RUMP, AND S. OISHI



Algorithm Dot2, n = 100, 1000 samples

More accuracy but still conditioning: overview

Arbitrary precision but accuracy still condition dependant

- Simulating more *precision*: double-double, quad-double, ..., MPFR
- Compensated algorithms: Kahan(65), ..., Sum2(05), SumK(05)
- Accuracy $\leq \mathbf{u} + cond \times \mathbf{u}^{K}$
- No more *n* dependency while $n\mathbf{u} \leq 1$.





One old and famous solution

Algorithm (Kahan's compensated summation (1965)) $s = x_1;$ c = 0;for i = 2 : n do $y = x_i \ominus c;$ $t = s \oplus y;$ $c = (t \ominus s) \ominus y;$ s = tend for

- Accuracy improvement in practice: -c approximates $s \oplus y$'s rounding error.
- Compensated sum accuracy $\leq 2\mathbf{u} + n \sum |x_i| \mathcal{O}(\mathbf{u}^2)$.

One example

IEEE double precision numbers: $x_1 = 2^{53} - 1$, $x_2 = 2^{53}$ and $x_3 = -(2^{54} - 2)$. Exact sum: $x_1 + x_2 + x_3 = 1$.



Relative error = 1

One example

IEEE double precision numbers: $x_1 = 2^{53} - 1$, $x_2 = 2^{53}$ and $x_3 = -(2^{54} - 2)$. Exact sum: $x_1 + x_2 + x_3 = 1$.



Compensation of the rounding errors



Relative error = 1

The exact result is computed

One example

IEEE double precision numbers: $x_1 = 2^{53} - 1$, $x_2 = 2^{53}$ and $x_3 = -(2^{54} - 2)$. Exact sum: $x_1 + x_2 + x_3 = 1$.



Relative error = 1 The exact result is computed The rounding errors are computed thanks to *error-free transformations*.

Compensation

D Summing *n* floating-point numbers: basic steps

- Some algorithms
- Some accuracy bounds

2 Computing sums more accurately

Some famous old and magic algorithms

Compensation

• Distillate to understand and go further

Omputing sums as accurately as possible

- iFastSum: distillation again
- Accsum
- Hybridsum and OnLineExactSum
- More EFT and accurate algorithms
 - Some more results



The underlined idea of compensated algorithm

Compensation

- Let's compute the generated rounding errors
- **2** and use it further to compensate the whole computation.



Compensating: the later the better?

Approximations in Kahan's compensated sum

• $x_1 = 2^{p+1}, x_2 = 2^{p+1} - 2, x_3 = x_4 = \dots = x_6 = -(2^p - 1)$

computed sum = 1, compensated sum = 3, exact sum = 2 (TP).

 |s| ≥ |y| or "exponent-wise" at least, round-to-nearest, order 2 rounding error in y.

Pichat and Neumaier's compensated summation (72,74)

- Accumulate rounding errors and one final correction.
- Compensated accuracy $\leq \mathbf{u} |\sum x_i| + (0.75n^2 + n)\mathbf{u}^2 \sum |x_i|.$
- No more cond, no initial sort and $n \le n'_{max} = 1/3$ u.

Priest's doubly compensated summation (92)

- Doubly compensated accuracy $\leq 2\mathbf{u} |\sum x_i|$.
- No more cond but initial sort of the summands and $n \le n_{max} = 1/4$ **u**.

Computing sums more accurately: step 2

Summing n floating-point numbers: basic steps

- Some algorithms
- Some accuracy bounds

2 Computing sums more accurately

- Some famous old and magic algorithms
- Compensation
- Distillate to understand and go further

Omputing sums as accurately as possible

- iFastSum: distillation again
- Accsum
- Hybridsum and OnLineExactSum
- More EFT and accurate algorithms
 - Some more results



Next step: towards the best accuracy



- Long accumulator, hardware oriented: Malcolm (71), Kulish (80)
- Cutting the summands: AccSum (SISC-08), FastAccSum (SISC-09)
- Summation at a given exponent: HybridSum (SISC-09), OnLineExact (TOMS-10)

From faithful rounding to correct rounding

• Choosing the "right side": expensive in the break-point neighborhood

• e.g.
$$1 + 2^{-53} \pm 2^{-106}$$

Next step: towards the best accuracy

Distillation: iterate until faithful or correct rounding

- Error free transformation (EFT) $[x] \rightarrow [x^{(1)}] \rightarrow \cdots \rightarrow [x^*]$
 - s.t. $\sum x_i = \sum x_i^*$ et $[x^*]$ returns the expected rounded value.
- Kahan (87), ..., Zhu-Hayes: iFastSum (SISC-09)

More space to keep everything

- Long accumulator, hardware oriented: Malcolm (71), Kulish (80)
- Cutting the summands: AccSum (SISC-08), FastAccSum (SISC-09)
- Summation at a given exponent: HybridSum (SISC-09), OnLineExact (TOMS-10)

From faithful rounding to correct rounding

 \longrightarrow Running-time and memory print are the discriminant factors

 $x + y = x \oplus y$

EFT to sum two floating-point numbers

2Sum (Knuth, 65), Fast2Sum (Dekker, 71) for base \leq 2 and RTN.

$$a + b = x + y$$
, with $a, b, x, y \in \mathbb{F}$ and $x = a \oplus b$.



Algorithm (Knuth)

function [x,y] = 2Sum(a,b) $x = a \oplus b$ $z = x \ominus a$ $y = (a \ominus (x \ominus z)) \oplus (b \ominus z)$

Algorithm (|a| > |b|, Dekker) function [x,y] = Fast2Sum(a,b) $x = a \oplus b$ $z = x \ominus a$ $y = b \ominus z$

<i>x</i> ₁	x ₂	x 3	<i>x</i> ₄		x_{n-1}	x _n
-----------------------	-----------------------	------------	-----------------------	--	-----------	----------------



<i>e</i> ₁	s 1	<i>x</i> ₃	<i>x</i> ₄		<i>x</i> _{<i>n</i>-1}	x _n
-----------------------	------------	-----------------------	-----------------------	--	--------------------------------	----------------

$$\sum_{1}^{n} x_{i} = e_{1} + s_{1} + \sum_{3}^{n} x_{i}$$



$$\sum_{1}^{n} x_{i} = e_{1} + e_{2} + s_{2} + \sum_{4}^{n} x_{i}$$

e ₁	<i>e</i> ₂ <i>e</i> ₃	<i>e</i> ₄		<i>e</i> _{<i>n</i>-2}	s _{n-1}	x _n
-----------------------	---	-----------------------	--	--------------------------------	-------------------------	----------------

$$\sum_{1}^{n} x_{i} = \sum_{1}^{n-2} e_{i} + s_{n-1} + x_{n}$$



e ₁	e ₂	e ₃	<i>e</i> ₄		<i>e</i> _{<i>n</i>-2}	<i>e</i> _{<i>n</i>-1}	s _n
-----------------------	-----------------------	----------------	-----------------------	--	--------------------------------	--------------------------------	----------------

$$\sum_{1}^{n} x_{i} = \sum_{1}^{n-1} e_{i} + s_{n} = \sum_{1}^{n} x_{i}^{(1)}$$

and one iterate within this new vector $[x^{(1)}]$...

Let's distillate!



Theorem (Zhu-Hayes,09)

Terminating for IEEE-754 : Iterate the distillation $[x] \rightarrow [x^{(1)}] \rightarrow \cdots \rightarrow [x^{(k)}] \rightarrow \cdots$ converges towards a stable state $[x^*]$ such that $|x_1^*| < |x_2^*| < \cdots < |x_n^*|$, $x_i^* \oplus x_{i+1}^* = x_{i+1}^*$ et $\sum x_i = \sum x_i^*$.

Rmk : the convergence proof uses the round-to-even tie breaking rule.

We need a good stopping criteria!

EFT to sum n floating-point numbers

Rien ne se perd, rien ne se crée, tout se transforme (Anaxagore-Lavoisier)

EFT vector transformation: (Ogita-Rump-Oishi, 05) $[p_1, p_2, \cdots, p_n] \longmapsto [q_2, q_3, \cdots, q_n, \pi_n]$



For $(p_i)_{1 \le i \le n} \in \mathbb{F}$, $p_1 + p_2 = \pi_2 + q_2$, $p_2 + p_3 = \pi_3 + q_3$, \cdots ,

$$\sum_{i=1}^n p_i = \pi_n + \sum_{i=2}^n q_i, \text{ with } \pi_n = \mathsf{fl}(\sum_{i=1}^n p_i)$$

A priori stopping criteria: go away conditioning!

Sum2, SumK (Ogita-Rump-Oishi,05)



FIG. 4.2. Outline of Algorithm 4.8 for n = 5 and K = 4.

Theorem (OgRO,05)

For every p_i in \mathbb{F} , $\sum p_i = \sum p'_i = \dots = \sum p'_i = \dots$, $cond(\sum p^K_i) = \mathcal{O}(\mathbf{u}^K) \times cond(\sum p_i)$

Accuracy $\lesssim \mathbf{u} + \textit{cond} \times \mathbf{u}^{K}$



T. OGITA, S. M. RUMP, AND S. OISHI

FIG. 6.1. Test results for Algorithm 5.3 (Dot2), n = 100, 1000 samples.

Accuracy $\lesssim \mathbf{u} + \mathbf{cond} \times \mathbf{u}^{K}$



T. OGITA, S. M. RUMP, AND S. OISHI

Step 2: towards faithfully or correctly rounded sums

Summing n floating-point numbers: basic steps

- Some algorithms
- Some accuracy bounds
- 2 Computing sums more accurately
 - Some famous old and magic algorithms
 - Compensation
 - Distillate to understand and go further

Omputing sums as accurately as possible

- iFastSum: distillation again
- Accsum
- Hybridsum and OnLineExactSum
- 4 More EFT and accurate algorithms
 - Some more results

No more cond: new solutions

Distillation iFastSum (09)

- Rounding errors quickly vanish in practice.
- So distillation soon iterates on shorter vectors.
- Recursive calls only for difficult cases, e.g. $1 + 2^{-53} + 2^{-106}$.

How to compute exact partial sums?

Exact summation: cut the summands

- Faithfully rounded sum: AccSum (08), FastAccSum (09)
- Correctly rounded sum: NearSum (09)

Exact summation: partial sum by exponant + distillation

• Faithfully rounded sum: HybridSum (09), OnLineExact (10)

Computing sums as accurately as possible

D Summing *n* floating-point numbers: basic steps

- Some algorithms
- Some accuracy bounds

2 Computing sums more accurately

- Some famous old and magic algorithms
- Compensation
- Distillate to understand and go further

Omputing sums as accurately as possible

- iFastSum: distillation again
- Accsum
- Hybridsum and OnLineExactSum
- More EFT and accurate algorithms
 - Some more results



iFastSum (09)

iFastSum = distallation SumK ($K \ge 2$) with a dynamic control of the residual error until the faithfully or correctly rounded value.

Memory overhead = $\mathcal{O}(1)$

No more cond limitation: always faithfully or correctly rounded sum.

- Rounding errors quickly vanish in practice.
- So distillation soon iterates on shorter vectors.
- Recursive calls only for difficult cases, e.g. $1 + 2^{-53} + 2^{-106}$.








Bounding the residual error $e_m = (\text{count}) \times 1/2ulp(\max |s_t|)$









Stopping criteria in 2 steps

- Stop the distillation when e_m = 0 or e_m < HalfUlp(s)</p>
- Orrect rounding: effect of e_m to s and s_t

▶ 2 "recursive" calls of the iFastSum's main part if $s_t \pm e_m$ modifies s or Round3 $(s, (s_t \pm e_m)_h, (s_t \pm e_m)_l) \neq s$.



Principe

Sharp approximation of the exact sum of 3 non overlapping values: 2 floating-points and 1 sign (the value does not matter).



 $\mathsf{Round3}(s_0, s_1, \operatorname{sign}(\widehat{s_2})) = \mathsf{fl}(s_0 + s_1 + \widehat{s_2})$

- One test of s_1
- exponent manipulation
- but no change of the rounding mode

Computing sums as accurately as possible

D Summing *n* floating-point numbers: basic steps

- Some algorithms
- Some accuracy bounds

2 Computing sums more accurately

- Some famous old and magic algorithms
- Compensation
- Distillate to understand and go further

Omputing sums as accurately as possible

• iFastSum: distillation again

Accsum

- Hybridsum and OnLineExactSum
- More EFT and accurate algorithms
 - Some more results











AccSum uses EFT ExtractVector (Fast2Sum)



Steps for sequential Compute max.

- Sum each shrunk

AccSum uses EFT ExtractVector (Fast2Sum)



Steps for sequential • Compute max.

• Sum each shrunk

AccSum uses EFT ExtractVector (Fast2Sum)



Steps for sequential • Compute max.

• Sum each shrunk





AccSum uses EFT ExtractVector (Fast2Sum)



Steps for sequential

- Compute max.
- Sum each shrunk

AccSum uses EFT ExtractVector (Fast2Sum)



Steps for sequential • Compute max.

• Sum each shrunk



AccSum proportional to $log_2(cond)$.

In practice: how to compute σ and how to extract q_i from p_i :

•
$$M = \lceil \log_2(n+2) \rceil$$

•
$$\sigma = 2^M \cdot max\{|p_i|\}, \text{ for } n+2 \leq 2^M$$

• Extraction of the higher part q_i of p_i : $[q_i, p'_i] = Fast2Sum(\sigma, p_i)$

Properties

- Number of extractions: proportional to log₂(cond)
- Number of summands limitations: $n \lesssim 1/\sqrt{u} \approx 7 \cdot 10^7$ (alternative exist for larger *n*).

AccSum proportional to $log_2(cond)$.

In practice: number of extractions and flop

214

S. RUMP, T. OGITA, AND S. OISHI

TABLE 4.3 Minimum treatable length n for condition number eps^{-1} in IEEE 754 double precision.

m=2	m = 3	m = 4
4094	$1.0\cdot 10^6$	$6.7\cdot 10^7$

TABLE 4.4 Floating-point operations needed for different dimension and condition number.

n	cond	AccSum	Sum2	XBLAS
1000	10^{6}	7n	7n	10n
1000	10^{16}	11n	7n	10n
10^{6}	10^{16}	15n	7n	10n

Computing sums as accurately as possible: step 2

D Summing *n* floating-point numbers: basic steps

- Some algorithms
- Some accuracy bounds

2 Computing sums more accurately

- Some famous old and magic algorithms
- Compensation
- Distillate to understand and go further

Omputing sums as accurately as possible

- iFastSum: distillation again
- Accsum
- Hybridsum and OnLineExactSum
- More EFT and accurate algorithms
 - Some more results



Partial sums by exponent (sorting) + distillation

Zhu-Hayes : HybridSum (09), OnLineExact (10)

Exponent extraction and sorting, then distillation

Principle :

- With no loss, transforming the entry vector to one or two (short) vectors which size equals the exponent range: 2048 for IEEE-b64 ;
 - HybridSum : split + 1 short vector
 - OnLineExact : 2Sum+ 2 vectors concatenated in a very short vector (of non zero entries)
- 2 then distillate this very short vector.

HybridSum (09) while $n \le 7 \cdot 10^7$ in IEEE-b64

Idea : split and exact accumulation with \bigoplus











More EFT and accurate algorithms

Summing *n* floating-point numbers: basic steps

- Some algorithms
- Some accuracy bounds

Computing sums more accurately

- Some famous old and magic algorithms
- Compensation
- Distillate to understand and go further

3 Computing sums as accurately as possible

- iFastSum: distillation again
- Accsum
- Hybridsum and OnLineExactSum
- More EFT and accurate algorithms
 - Some more results



More accuracy: more results

Other EFT to compensate

- Multiplication: $a \times b = (a \otimes b) + e$
- FMA: $a \times b + c = FMA(a, b, c) + e_1 + e_2$, and $FMA(a, b, c) = RN(a \times b + c)$
- Polynomial evaluation

Other compensated algorithms

- Dot product
- Horner and derivative (dsynthetic div.), de Casteljau (Bernstein), Clenshaw (Chebychev)

Sterbenz's lemma: when subtraction is exact

Theorem (Sterbenz, 72)

In a radix- β floating-point system with subnormal numbers, if x and y are finite floating-point numbers such that

$$\frac{y}{2} \le x \le 2y,$$

then x - y is a floating-point number.

Comments

- Exact subtraction for the four IEEE-754 rounding-modes
- Exact subtraction vs. catastrophic cancellation?
- Used in the previous EFT sum proofs

EFT: the multiplication case

Theorem (Dekker)

The multiply rounding-error $xy - (x \otimes y)$ is a floating-point number when no overflow occurs.

When one FMA is available, next algorithm computes the multiply-EFT:

 $xy = r_1 + r_2$.

Algorithm (2ProdFMA)

$$r_1 = FMA(x, y, 0); //r_1 = x \otimes y r_2 = FMA(x, y, -r_1); //r_2 = xy - r_1$$

The multiplication case when no available FMA

Let $C = \beta^s + 1$ where s < p, the next algo splits the *p*-digit floating-point number *x* in two parts x_h, x_l of p - s and *s* digits.

Algorithm (split - Veltkamp, 68)		
$r_shift = C \otimes x$		
$head = x \ominus r_shift$		
$x_h = r_shift \oplus head$		
$x_l = x \ominus x_h$		

Theorem (Dekker, 72; Boldo, 06)

When no overflow occurs,

$$x = x_h + x_l$$
.

The multiplication case when no available FMA

When no overflow occurs, next algorithm computes the multiply-EFT:

 $xy = r_1 + r_2$.

Algorithm (2prod - Dekker, 72) $(x_h, x_l) = split(x);$ $(y_h, y_l) = split(y);$ $r_1 = x \otimes y;$ $high = -r_1 \oplus (x_h \otimes y_h);$ $mid_1 = high \oplus (x_h \otimes y_l);$ $mid_2 = mid_1 \oplus (x_l \otimes y_h);$ $r_2 = mid_2 \oplus (x_l \otimes y_l);$ //low part

The compensated dot product (Ogita-Rump-Oishi, 05)

EFT-DotProd:
$$\sum_{i=1}^{n} x_i y_i = p + \sum_{i=1}^{n} \pi_i + \sum_{i=1}^{n-1} \sigma_i = \sum_{i=1}^{2n} z_i$$
.

Algorithm (Compensated DotProd)

$$(s_1, c_1) = 2Prod(x_1, y_1); \\ for i = 2 : n do \\ (p_i, \pi_i) = 2Prod(x_i, y_i); \\ (s_i, \sigma_i) = 2Sum(p_i, s_{i-1}); \\ c_i = c_{i-1} \oplus (\pi_i \oplus \sigma_i); \\ end for$$

return $s_n \oplus c_n$;

Theorem (Ogita-Rump-Oishi, 05)

Compensated dot product accuracy $\leq \mathbf{u} |\sum_{i=1}^{n} x_i y_i| + \mathcal{O}(n^2 \mathbf{u}^2) \sum |x_i y_i|.$

The compensated dot product (Ogita-Rump-Oishi, 05)





FIG. 6.1. Test results for Algorithm 5.3 (Dot2), n = 100, 1000 samples.
The compensated dot product (Ogita-Rump-Oishi, 05)



T. OGITA, S. M. RUMP, AND S. OISHI

The Horner polynomial evaluation (Graillat-Langlois-Louvet, 07)

EFT-Horner:
$$\sum_{i=0}^{n} a_i x^i = p + \sum_{i=0}^{n-1} \pi_i x^i + \sum_{i=0}^{n-1} \sigma_i y^i = p + \pi(x) + \sigma(x).$$

Algorithm (Compensated Horner) $r_n = a_n;$ for i = n - 1: (-1): 0 do $(p_i, \pi_i) = 2Prod(r_{i+1}, x);$ $(r_i, \sigma_i) = 2Sum(p_i, a_i);$ $q_i = \pi_i \oplus \sigma_i; //correcting polynomial coefficient$ end for; $r = r_0 \oplus Horner(q, x);$ return r;

Theorem (Langlois-Louvet, 07)

Compensated Horner accuracy $\leq \mathbf{u}|p(x)| + \mathcal{O}(4n^2\mathbf{u}^2)\sum_{i=0}^n |a_i||x^i|$.

The compensated Horner evaluation (Louvet's PhD, 07)



The compensated Horner evaluation: FMA variations



Recent extensions

- derivative Horner evaluation synthetic division (Jiang et al., 12)
- Clenshaw algorithm for Chebychev basis (Jiang et al., 11)
- de Casteljau algorithm for Bernstein basis (Jiang et al., 10)

To conclude this Part 1

Summing *n* floating-point numbers: basic steps

- Some algorithms
- Some accuracy bounds

Computing sums more accurately

- Some famous old and magic algorithms
- Compensation
- Distillate to understand and go further

3 Computing sums as accurately as possible

- iFastSum: distillation again
- Accsum
- Hybridsum and OnLineExactSum
- 4 More EFT and accurate algorithms
 - Some more results



An "ultimately accurate" floating-point summation algorithm cannot be very simple. J.M. Muller et al. [HFPA,09]

- A lot of algorithms!
- Condition dependency and accuracy bounds for classic ones
- Arbitrarily accurate compensated ones: Sum2, SumK
- Faithfully or correctly rounded recent ones: AccSum, FastAccSum, iFastSum, HybdridSum, OnLineExact
- "In place" solutions vs. exponent manipulations

How to choose? Running time performance!

Part II

More reproducibility

6 Motivation

Reproducibility failure in a finite element simulation

- Sequential and parallel FE assembly
- Sources of non reproducibility in Telemac2D

Recovering reproducibility

- Reproducible parallel FE assembly
- A first application: reproducible Tomawac
- Reproducible algebraic operations
- Reproducible conjugate gradient
- Reproducible Telemac2D

9 Efficiency



Reproducibility failure in Telemac2D simulation







• Reproducibility: bitwise identical results for every *p*-parallel run, $p \ge 1$



• Reproducibility: bitwise identical results for every *p*-parallel run, $p \ge 1$



- Reproducibility: bitwise identical results for every *p*-parallel run, $p \ge 1$
- Reproducibility \neq Accuracy



- Reproducibility: bitwise identical results for every *p*-parallel run, $p \ge 1$
- Reproducibility \neq Accuracy
- Failures reported in numerical simulation for energy [21], dynamic weather science [4], dynamic molecular [20], dynamic fluid [17]



- Reproducibility: bitwise identical results for every *p*-parallel run, $p \ge 1$
- Reproducibility \neq Accuracy
- Failures reported in numerical simulation for energy [21], dynamic weather science [4], dynamic molecular [20], dynamic fluid [17]
- How to debug? How to test? How to validate? How to receive legal agreements?



Today's issues

Case study

- Industrial scale software: openTelemac-Mascaret
- Finite element simulation, domain decomposition, linear system solving
- 2 modules: Tomawac, Telemac2D

Today's issues

Case study

- Industrial scale software: openTelemac-Mascaret
- Finite element simulation, domain decomposition, linear system solving
- 2 modules: Tomawac, Telemac2D

Feasibility

- How to recover reproducibility?
- Sources of non-reproducibility?
- Do existing techniques apply? how easily?
- Compensation yields reproducibility here!

Today's issues

Case study

- Industrial scale software: openTelemac-Mascaret
- Finite element simulation, domain decomposition, linear system solving
- 2 modules: Tomawac, Telemac2D

Feasibility

- How to recover reproducibility?
- Sources of non-reproducibility?
- Do existing techniques apply? how easily?
- Compensation yields reproducibility here!

Efficiency

- How much to pay for reproducibility?
- $\times 1.2 \leftrightarrow \times 2.3$ extra-cost which decreases as the problem size increases
- OK to debug, to validate and even to simulate!

6 Motivation

Reproducibility failure in a finite element simulation

- Sequential and parallel FE assembly
- Sources of non reproducibility in Telemac2D

Recovering reproducibility

- Reproducible parallel FE assembly
- A first application: reproducible Tomawac
- Reproducible algebraic operations
- Reproducible conjugate gradient
- Reproducible Telemac2D

9 Efficiency



Non associative floating point addition

• The computed value depends on the operation order

Parallel reduction and compensation techniques

Non associative floating point addition

- The computed value depends on the operation order
- Parallel reduction of undefined order generates reproducibility failure



Parallel reduction and compensation techniques

Non associative floating point addition

- The computed value depends on the operation order
- Parallel reduction of undefined order generates reproducibility failure
- Compensate rounding errors with error free transformations



 $\begin{array}{l} ((a\oplus b)\oplus(c\oplus d))\oplus((e_1\oplus e_2)\oplus e_3)=(((a\oplus b)\oplus c)\oplus d)\oplus((f_1\oplus f_2)\oplus f_3)\\ \\ \text{Should be repeted for too ill-conditionned sums} \end{array}$

Finite element assembly: the sequential case

The assembly step: $V(i) = \sum_{elements} W_{el}(i)$

- compute the inner node values V(i)
- accumulating local W_{el} for every el that contains i



Finite element assembly: the sequential case

The assembly step: $V(i) = \sum_{elements} W_{el}(i)$

- compute the inner node values V(i)
- accumulating local W_{el} for every el that contains i



Finite element assembly: the parallel case





Exact arithmetic

Finite element assembly: the parallel case





Floating point arithmetic

Basic ingredients

Motivation

Reproducibility failure in a finite element simulation

- Sequential and parallel FE assembly
- Sources of non reproducibility in Telemac2D

Recovering reproducibility

- Reproducible parallel FE assembly
- A first application: reproducible Tomawac
- Reproducible algebraic operations
- Reproducible conjugate gradient
- Reproducible Telemac2D

9 Efficiency



Culprits: theory

- Building step: interface point assembly
- **②** Solving step (conjugate gradient): parallel matrix-vector and dot products

Culprits: theory

- Building step: interface point assembly
- Solving step (conjugate gradient): parallel matrix-vector and dot products

Culprits: practice = optimization

- Interface point assembly and linear system solving are merged
- Element-by-element storage (EBE) of the FE matrix
 - EBE parallel matrix-vector product: no BLAS
 - Everything is vector, no matrix!
- Wave equation, "mass-lumping" and associated algebraic transformations
 - System decoupling and many diagonal matrices
 - Everything is vector, no matrix!

The Telemac2D FE steps



algebraic computation

Mesh (elements, nodes)

Discretisation

The Telemac2D FE steps



Recovering reproducibility in a finite element resolution

6 Motivation

- 7 Reproducibility failure in a finite element simulation
 - Sequential and parallel FE assembly
 - Sources of non reproducibility in Telemac2D

Recovering reproducibility

- Reproducible parallel FE assembly
- A first application: reproducible Tomawac
- Reproducible algebraic operations
- Reproducible conjugate gradient
- Reproducible Telemac2D

Efficiency



Recovering reproducibility in Telemac2D

Sources

- FE assembly: diagonal of the matrices and second members
- Resolution: EBE matrix-vector and dot products
- Wave equation: algebraic transformations and diagonal resolutions

Reproducible resolution: principles

- vector $V \rightarrow [V, E_V] \rightarrow V + E_V$
- Computes E_V in the FE assembly of V
- Propagates E_V over each V operation
- Compensates all nodes while assembling the Interface Point
- Compensate MPI parallel dot products that include MPI reduction

Reproducible FE assembly

i: any node; $W_{el}(i)$: contribution for every *el* that contains *i*

Original FE assembly: $V(i) = \sum_{elements} W_{el}(i)$

 $V(i) = W_{el_1}(i) + W_{el_2}(i) + \cdots + W_{el_{ni}}(i)$


Reproducible interface point assembly

i: one interface point of $D_1, D_2, \dots, D_{k-1}, D_k$

Original IP assembly: $V(i) = \sum_{D_k} V(i)$

Communicate V_{D_i} to D_1, D_2, \cdots and compute: $V(i) = V_{D_1}(i) + V_{D_2}(i) + \cdots + V_{D_{k-1}}(i) + V_{D_k}(i)$

Reproducibility: at the IP assembly step

 $[V, E_V] \longrightarrow V + E_V$ compensates every node value

FE assembly: implementation

Main steps

• For every vector value: FE assembly + IP assembly

ASSVEC subroutine

		-	
1	X refers to VEC and ERRX refers to SVEC%E	3 4	CALL ASSVEC(VEC, IKLE, NPT ,NELEM,NELMAX,IELM1, & T,INIT,LV,MSK,MASKEL,NDP)
3	DO TELEM = 1 NELEM	5	
	bo ibbbii - i , wbbbii	6	
4		7	
5	X(IKLE(IELEM,IDP))=X(IKLE(IELEM,IDP))+W(IELEM,IDP)		
6		0	
7		9	Interface point assembly in PARCOM
		10	IF(ASSPAR) CALL PARCOM(SVEC, 2, MESH)
0		11	
9		12	
10		12	
11	ENDDO	13	
12	ENDDO	14	
**	ENDO	15	
_		16	

IFE assembly in ASSVEC

1

Interface point assembly: implementation

```
Receive step
    DO IL=1.NB NEIGHB
 3
      IKA = NB NEIGHB PT(IL)
 4
      IPA = LIST_SEND(IL)
 5
      CALL P IREAD(BUF RECV(1, IL), IAN*IKA*NPLAN*8, IPA, PARACO MSG TAG, RECV REQ(IL))
 6
 7
    ENDDO
 8
    Send step
 9 DO IL=1,NB_NEIGHB
10
       IKA = NB NEIGHB PT(IL)
11
       IPA = LIST SEND(IL)
       Initializes the communication arrays
12
13
       K = 1
14
       DO J=1,NPLAN
15
        DO I=1.IKA
16
          II=NH_COM(I,IL)
17
           BUF_SEND(K,IL) =V1(II,J)
18
19
           K=K+1
20
         ENDDO
21
       ENDDO
       CALL P IWRIT(BUF SEND(1,IL),IAN*IKA*NPLAN*8,IPA,PARACO MSG TAG,SEND REQ(IL))
22
23
24
    ENDDO
25
    | Wait received messages
26 DO IL=1.NB NEIGHB
27
       IKA = NB_NEIGHB_PT(IL)
28
       IPA = LIST_SEND(IL)
29
       CALL P WAIT PARACO(RECV REQ(IL).1)
30
       K=1
31
       DO J=1,NPLAN
32
        DO I=1.IKA
33
           II=NH COM(I,IL)
34
           V1(II,J)=V1(II,J)+ BUF_RECV(K,IL)
35
36
37
38
39
           K=K+1
40
         ENDDO
41
```

Reproducible FE assembly: implementation

Main steps

For every vector value: FE assembly + IP assembly + compensation

ASSVEC subroutine

1	X refers to VEC and ERRX refers to SVEC%E
2	DO IDP = 1 , NDP
3	DO IELEM = 1 , NELEM
4	IF (MODASS.EQ.1)
5	<pre>& X(IKLE(IELEM, IDP))=X(IKLE(IELEM, IDP))+W(IELEM, IDP)</pre>
6	ELSEIF (MODASS.EQ.3) THEN
7	CALL 2SUM(X(IKLE(IELEM, IDP)),
8	& W(IELEM, IDP), X(IKLE(IELEM, IDP)), ERROR)
9	ERRX(IKLE(IELEM, IDP))=ERRX(IKLE(IELEM, IDP))+ERROR
10	ENDIF
11	ENDDO
12	ENDDO

```
Note: VEC is a reference to SVEC%R
    IF(MODASS, EQ. 1) THEN
 3
       CALL ASSVEC(VEC, IKLE, NPT ,NELEM,NELMAX,IELM1,
 Δ
    & T.INIT.LV.MSK.MASKEL.NDP)
 5
    ELSEIF (MODASS, EQ. 3) THEN
       CALL ASSVEC(VEC, IKLE, NPT ,NELEM,NELMAX,IELM1,
 6
    & T.INIT.LV.MSK.MASKEL.NDP.SVEC%E)
 7
 8
    ENDIF
    Implicit modification in PARCOM
 Q
    IF(ASSPAR) CALL PARCOM(SVEC, 2, MESH)
10
    IF (ASSPAR. AND. MODASS. EQ. 3) THEN
11
    The compensation of all the values
12
     DO I = 1 , MESH%NPOIN
13
14
         VEC(I) = VEC(I)+SVEC%E(I)
15
     ENDDO
16 ENDIF
```

Reproducible interface point assembly: implementation

1	! Receive step
2	DO IL=1,NB_NEIGHB
3	IKA = NB_NEIGHB_PT(IL)
4	IPA = LIST_SEND(IL)
5	CALL P_IREAD(BUF_RECV(1,IL),IAN*IKA*NPLAN*8,IPA,PARACO_MSG_TAG,RECV_REQ(IL))
6	CALL P_IREAD(BUF_RECV_ERR(1,IL),IAN*IKA*NPLAN*8,IPA,PARACO_MSG_TAG,RECV_REQ(IL))
7	ENDDO
8	Send step
9	DO IL=1.NB NEIGHB
10	IKA = NB NEIGHB PT(IL)
11	IPA = LIST SEND(IL)
12	Initializes the communication arrays
13	K = 1
14	DO JE1.NPLAN
15	
16	II=NH COM(I II)
17	BUF SEND(K, IL) = V1(IL, I)
18	BUF SEND ERR(K, IL) =ERRX(II)
19	K=K+1
20	ENDDO
21	ENDDO
22	CALL P IWRIT(BUF SEND(1,IL),IAN*IKA*NPLAN*8,IPA,PARACO MSG TAG,SEND REQ(IL))
23	CALL P IWRIT(BUF SEND ERR(1,IL),IAN*IKA*NPLAN*8,IPA,PARACO MSG TAG,SEND REQ(IL))
24	ENDDO
25	! Wait received messages
26	DO IL=1.NB NEIGHB
27	IKA = NB NEIGHB PT(IL)
28	IPA = LIST SEND(IL)
29	CALL P WAIT PARACO(RECV REQ(IL),1)
30	K=1
31	DO J=1,NPLAN
32	DO I=1.IKA
33	II=NH_COM(I,IL)
34	Original version: V1(II,J)=V1(II,J)+ BUF RECV(K,IL)
35	CALL 2SUM(V1(II,J), BUF RECV(K,IL),V1(II,J),ERROR1)
36	CALL 2SUM(ERRV(II), BUF RECV ERR(K, IL), ERRV(II), ERROR2)
37	ERROR=ERROR1+ERROR2
38	ERRV(II)=ERRV(II)+ERROR
39	K=K+1
40	ENDDO
41	ENDDO

Reproducible algebraic operations

6 Motivation

7 Reproducibility failure in a finite element simulation

- Sequential and parallel FE assembly
- Sources of non reproducibility in Telemac2D

Recovering reproducibility

Reproducible parallel FE assembly

• A first application: reproducible Tomawac

- Reproducible algebraic operations
- Reproducible conjugate gradient
- Reproducible Telemac2D

9 Efficiency



Tomawac case

Non-reproducible Tomawac

- Wave propagation in coastal areas
- Mesh node unknowns:

height, frequency and direction of the waves







Reproducible Tomawac

- Accurate compensated summation [15]
- Integer conversion [16]
- Demmel-Nguyen's reproducible sums [2]





Reproducible algebraic operations

6 Motivation

7 Reproducibility failure in a finite element simulation

- Sequential and parallel FE assembly
- Sources of non reproducibility in Telemac2D

Recovering reproducibility

- Reproducible parallel FE assembly
- A first application: reproducible Tomawac

• Reproducible algebraic operations

- Reproducible conjugate gradient
- Reproducible Telemac2D

Efficiency



Algebraic operation: the vector case

Reproducible algebraic vector operations

- openTelemac's included library: BIEF
- Entry-wise vector ops: copy, opp/inv., add/sub, Hadamard prod., ...
- Applies also for diagonal matrix
- Propagate rounding errors to compensate while assembling IP

Example: Hadamard product

Original version $X, Y \mapsto V = X \circ Y$

 $V(i) = X(i) \cdot Y(i)$

Modified version $[X, E_X], [Y, E_Y] \mapsto [V, E_V]$ with $(V, e_V) = 2 \operatorname{Prod}(X, Y)$ and $E_V = X \circ E_Y + Y \circ E_X + e_V$

What is reproducible now?

Most of the linear system:

- FE assembly
- algebraic vector operations
- interface point assembly

except:

- the matrix of the H system
- its dependencies: the second members of the U and V systems

Next step:

conjugate gradient

Partially reproducible Telemac2D



Recovering reproducibility in a finite element resolution

6 Motivation

7 Reproducibility failure in a finite element simulation

- Sequential and parallel FE assembly
- Sources of non reproducibility in Telemac2D

Recovering reproducibility

- Reproducible parallel FE assembly
- A first application: reproducible Tomawac
- Reproducible algebraic operations

• Reproducible conjugate gradient

• Reproducible Telemac2D

9 Efficiency



Towards a reproducible conjugate gradient



Non-reproducibility sources

- EBE matrix-vector product
- dot product
 - MPI reduction
 - weighted dot products for IP shared by p sub-domains

The EBE storage

$$M = D + \sum_{el=1}^{nel} X_{el}$$

- nodes $i \in [1, np]$, elements $el \in [1, nel]$, element vertices $j, k, l \in el$
- *M* is decomposed as:
 - **Q** assembled diagonal D[np]: $D = [D(1), \dots, D(np)]$

elementary extra-diagonal X_{el}[6]:

$$X_{el} = \begin{pmatrix} \cdot & X_{jk}(el) & X_{jl}(el) \\ X_{kj}(el) & \cdot & X_{kl}(el) \\ X_{lj}(el) & X_{lk}(el) & \cdot \end{pmatrix} = [X_{el}(1), \cdots, X_{el}(6)]$$

The EBE Matrix-Vector product

$$R = M \cdot V = D \cdot V + \sum_{el=1}^{nel} X_{el} \cdot V_{el}$$

Steps of the EBE matrix-vector product

IP assembly:
$$R(i) = \sum_{D_k} R(i)$$
 for all IP *i*

Reproducible EBE matrix-vector product

Original EBE matrix-vector product

•
$$R = D \cdot V + \sum_{el=1}^{nel} X_{el} \cdot V_{el}$$

• $R(i) = \sum_{D_k} R(i)$

Reproducible EBE matrix-vector product

- $[R, E_R] = [D, E_D] \circ V + \operatorname{ReprodAss}_{el=1}^{nel} X_{el} \cdot V_{el}$
- $R(i) = \operatorname{ReprodAss}_{D_k}[R(i), E_R(i)]$
- Compensation: $R + E_R$

A reproducible conjugate gradient

$$A = [A_{1}, E_{A_{1}}]$$

$$B = C_{1}$$
Initialization: a given d^{0}

$$r^{0} = A\chi^{0} - B;$$

$$\rho^{0} = \frac{(r^{0}, d^{0})}{(Ad^{0}, d^{0})}; X^{1} = X^{0} - \rho^{0}d^{0}$$
Iterations until stopping criteria:
$$r^{m} = r^{m-1} - \rho^{m-1}Ad^{m-1}$$

$$d^{m} = r^{m} + \frac{(r^{m}, r^{m})}{(r^{m-1}, r^{m-1})}d^{m-1}$$

$$\rho^{m} = \frac{(r^{m}, d^{m})}{(d^{m}, Ad^{m})}$$

$$X^{m+1} = X^{m} - \rho^{m}d^{m}$$

Non-reproducibility: sources and solutions

- Reproducible EBE matrix-vector product
- dot product
 - MPI reduction: a parallel compensated dot2
 - weights: $(1/k, 1/k, \dots, 1/k) \rightarrow (1, 0, \dots, 0)$
- Reproducible operations \longrightarrow reproducible results

A reproducible conjugate gradient



Non-reproducibility: sources and solutions

- Reproducible EBE matrix-vector product
- dot product
 - MPI reduction: a parallel compensated dot2
 - weights: $(1/k, 1/k, \dots, 1/k) \rightarrow (1, 0, \dots, 0)$
- Reproducible operations \longrightarrow reproducible results
- Same errors for both sequential and parallel executions

Recovering reproducibility in a finite element resolution

6 Motivation

7 Reproducibility failure in a finite element simulation

- Sequential and parallel FE assembly
- Sources of non reproducibility in Telemac2D

Recovering reproducibility

- Reproducible parallel FE assembly
- A first application: reproducible Tomawac
- Reproducible algebraic operations
- Reproducible conjugate gradient
- Reproducible Telemac2D

9 Efficiency



Reproducible Telemac2D

Reproducible Telemac2D







10

15



2.64859427

2.40000000

95 / 108

2.64859427

10

15



p=4

Version comp: Unknown, H Time step: 2, number of processors: 4





Version comp: Unknown, H Time step: 2, number of processors: 8







2.79880948

2.63361644

2.46842340

2.30323036

2.13803732

.97284429

2.79880948

2.63361644

2 46842340

2.30323036

2.13803732

1.97284429















Time step 7

20

15

10

2.20086447



Version comp: Unknown, H

Time step: 7, number of processors: 8

95 / 108

2.68669468

2.58952864

2.49236260

2.39519655

2.29803051

2.20086447





p=8



p=4









p=4









p=4











Version comp: Unknown, H Time step: 12, number of processors: 8 2.93919137 2.81840066 2 69760994 2.57681923 2,45602852 2.33523781











10

2.20242547

Time step 14







p=4








Efficiency

Motivation

Reproducibility failure in a finite element simulation

- Sequential and parallel FE assembly
- Sources of non reproducibility in Telemac2D

Recovering reproducibility

- Reproducible parallel FE assembly
- A first application: reproducible Tomawac
- Reproducible algebraic operations
- Reproducible conjugate gradient
- Reproducible Telemac2D

9 Efficiency

More reproducibility: conclusion

Runtime extra-cost for reproducible simulations

Measures, test cases and mesh sizes		
hardware cycle counter: rdtsc	#proc.	4624
ý	2	72
gouttedo	4	304
	8	501
• mesh sizes: 4624, 18225, 72361 nodes		Number
$(1, \approx imes 4, \approx imes 16)$		

Hardware and software env.

- openTelemac v7.2
- socket: Intel Xeon E5-2660 2.20GHz (L3 cache = 20 M)
- 2 sockets of 8 cores each
- GNU Fortran 4.6.3, -O3
- OpenMPI 1.5.4
- Linux 3.5.0-54-generic

		#nodes		
#proc.	4624	18225	72361	
2	72	143	280	
4	304	674	1368	
8	501	1152	2020	
Number of IP				

The core runtime extra-cost for reproducible gouttedo

gouttedo core: no input/output steps, just building+solving



98/108

Time to conclude

6 Motivation

Reproducibility failure in a finite element simulation

- Sequential and parallel FE assembly
- Sources of non reproducibility in Telemac2D

Recovering reproducibility

- Reproducible parallel FE assembly
- A first application: reproducible Tomawac
- Reproducible algebraic operations
- Reproducible conjugate gradient
- Reproducible Telemac2D

9 Efficiency



Recovering numerical reproducibility

- Industrial scale software: openTelemac-Mascaret
- Finite element simulation, domain decomposition, linear system solving
- 2 reproducible modules: Tomawac, Telemac2D
- Integration in the next openTelemac version: current work

Recovering numerical reproducibility

- Industrial scale software: openTelemac-Mascaret
- Finite element simulation, domain decomposition, linear system solving
- 2 reproducible modules: Tomawac, Telemac2D
- Integration in the next openTelemac version: current work

Feasibility

- How to recover reproducibility? Sources of non-reproducibility? Do existing techniques apply? how easily?
- Hand-made analysis of the computing workflow
- Compensation yields reproducibility here!
- Fits well to the openTelemac's vector library
- Other existing techniques also apply and more or less easily [11]

Efficiency

- How much to pay for reproducibility?
- $\times 1.2 \leftrightarrow \times 2.3$ extra-cost which decreases as the problem size increases
- OK to debug, to validate and even to simulate!

Efficiency

- How much to pay for reproducibility?
- $\times 1.2 \leftrightarrow \times 2.3$ extra-cost which decreases as the problem size increases
- OK to debug, to validate and even to simulate!

Reproducibility at a larger scale: the whole openTelemac software suite

- Does it still working for complex, large and real-life simulations?
- The two FE test cases are significant enough to validate the methodology
- Localization of the failure sources is difficult to automatize
- but the methodology application is easy for the software developpers

Existing techniques to recover numerical reproducibility in summation

- Accurate compensated summation [15]
- Demmel-Nguyen's reproducible sums [2]
- Integer convertion [16]

Part III

More performance

• Développé en TP.

References I

D. H. Bailey.

Twelve ways to fool the masses when giving performance results on parallel computers. *Supercomputing Review*, pages 54–55, Aug. 1991.

J. W. Demmel and H. D. Nguyen.

Fast reproducible floating-point summation.

In Proc. 21th IEEE Symposium on Computer Arithmetic. Austin, Texas, USA, 2013.

B. Goossens, P. Langlois, D. Parello, and E. Petit.

PerPI: A tool to measure instruction level parallelism.

In K. Jónasson, editor, Applied Parallel and Scientific Computing - 10th International Conference, PARA 2010, Reykjavík, Iceland, June 6-9, 2010, Revised Selected Papers, Part I, volume 7133 of Lecture Notes in Computer Science, pages 270–281. Springer, 2012.

Y. He and C. Ding.

Using accurate arithmetics to improve numerical reproducibility and stability in parallel applications.

```
J. Supercomput., 18:259–277, 2001.
```

References II

J.-M. Hervouet.

Hydrodynamics of free surface flows: Modelling with the finite element method. John Wiley & Sons, 2007.

N. J. Higham.

Accuracy and Stability of Numerical Algorithms. SIAM, 2nd edition, 2002.

H. Jiang, R. Barrio, H. Li, X. Liao, L. Cheng, and F. Su.

Accurate evaluation of a polynomial in chebyshev form.

Applied Mathematics and Computation, 217(23):9702 – 9716, 2011.

- H. Jiang, S. Graillat, C. Hu, S. Li, X. Liao, L. Cheng, and F. Su.
 Accurate evaluation of the k-th derivative of a polynomial and its application.
 J. Comput. Appl. Math., 243:28–47, May 2013.
 - H. Jiang, S. Li, L. Cheng, and F. Su.

Accurate evaluation of a polynomial and its derivative in bernstein form. *Computers & Mathematics with Applications*, 60(3):744 – 755, 2010.

References III

P. Langlois and N. Louvet.

How to ensure a faithful polynomial evaluation with the compensated Horner algorithm? In P. Kornerup and J.-M. Muller, editors, 18th IEEE International Symposium on Computer Arithmetic, number ISBN 0-7695-2854-6, pages 141–149. IEEE Computer Society, June 2007.

P. Langlois, R. Nheili, and C. Denis.

Numerical Reproducibility: Feasibility Issues.

In NTMS'2015: 7th IFIP International Conference on New Technologies, Mobility and Security, pages 1–5, Paris, France, July 2015. IEEE, IEEE COMSOC & IFIP TC6.5 WG.

P. Langlois, R. Nheili, and C. Denis.

Recovering numerical reproducibility in hydrodynamic simulations.

In J. H. P. Montuschi, M. Schulte, S. Oberman, and N. Revol, editors, 23rd IEEE International Symposium on Computer Arithmetic, number ISBN 978-1-5090-1615-0, pages 63–70. IEEE Computer Society, July 2016.

(Silicon Valley, USA. July 10-13 2016).

References IV

N. Louvet.

Algorithmes compensés en arithmétique flottante : précision, validation, performances. PhD thesis, Université de Perpignan Via Domitia, Nov. 2007.

J.-M. Muller, N. Brisebarre, F. de Dinechin, C.-P. Jeannerod, V. Lefèvre, G. Melquiond, N. Revol, D. Stehlé, and S. Torres.

Handbook of Floating-Point Arithmetic.

Birkhäuser Boston, 2010.

T. Ogita, S. M. Rump, and S. Oishi.

Accurate sum and dot product.

SIAM J. Sci. Comput., 26(6):1955-1988, 2005.



Open TELEMAC-MASCARET. v.7.0, Release notes.

www.opentelemac.org, 2014.



R. W. Robey, J. M. Robey, and R. Aulwes.

In search of numerical consistency in parallel programming. *Parallel Comput.*, 37(4-5):217–229, 2011.

References V



S. M. Rump.

Ultimately fast accurate summation.

SIAM J. Sci. Comput., 31(5):3466-3502, 2009.

S. M. Rump, T. Ogita, and S. Oishi.

Accurate floating-point summation – part I: Faithful rounding. SIAM J. Sci. Comput., 31(1):189–224, 2008.

M. Taufer, O. Padron, P. Saponaro, and S. Patel.

Improving numerical reproducibility and stability in large-scale numerical simulations on gpus.

In IPDPS, pages 1-9. IEEE, 2010.

🛛 📄 O. Villa, D. G. Chavarría-Miranda, V. Gurumoorthi, A. Márquez, and S. Krishnamoorthy.

Effects of floating-point non-associativity on numerical computations on massively multithreaded systems.

In CUG 2009 Proceedings, pages 1–11, 2009.

References VI



V. Weaver and J. Dongarra.

Can hardware performance counters produce expected, deterministic results?

In 3rd Workshop on Functionality of Hardware Performance Monitoring, 2010, pages 1–11, Atlanta, USA, 2010.

D. Zaparanuks, M. Jovic, and M. Hauswirth.

Accuracy of performance counter measurements.

In IEEE International Symposium on Performance Analysis of Systems and Software, ISPASS 2009, April 26-28, 2009, Boston, Massachusetts, USA, pages 23–32, 2009.



Y.-K. Zhu and W. B. Hayes.

Correct rounding and hybrid approach to exact floating-point summation. *SIAM J. Sci. Comput.*, **31**(4):2981–3001, 2009.

1

Y.-K. Zhu and W. B. Hayes.

Algorithm 908: Online exact summation of floating-point streams.

ACM Trans. Math. Software, 37(3):37:1-37:13, Sept. 2010.

version: May 22, 2017