

# Harnessing inexactness in scientific computing

## Lecture 2: summation

Theo Mary (CNRS)

[theo.mary@lip6.fr](mailto:theo.mary@lip6.fr)

<https://perso.lip6.fr/Theo.Mary/>

Elisa Riccietti (ENS Lyon)

[elisa.riccietti@ens-lyon.fr](mailto:elisa.riccietti@ens-lyon.fr)

<https://perso.ens-lyon.fr/elisa.riccietti/>

$$\sum_{i=1}^n x_i$$

**M2 course at ENS Lyon, 2024–2025**

Slides available on course webpage

## Introduction

Dealing with accumulation

Dealing with cancellation

Adaptive precision summation

Conclusion

$$y = \sum_{i=1}^n x_i \quad \dots \text{an ubiquitous and fundamental task!}$$

- **Dot products:**

$$a, b \in \mathbb{R}^n \Rightarrow a^T b = \sum_{i=1}^n a_i b_i$$

- **Matrix–vector products:**

$$A \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^n \Rightarrow (Ab)_j = \sum_{i=1}^n a_{ji} b_i, \quad j = 1: m$$

- **Matrix–matrix products:**

$$A \in \mathbb{R}^{m \times n}, B \in \mathbb{R}^{n \times p} \Rightarrow (AB)_{jk} = \sum_{i=1}^n a_{ji} b_{ik}, \quad j = 1: m, k = 1: p$$

- **Gaussian elimination (LU factorization):**

$$A \in \mathbb{R}^{n \times n}, A = LU \Rightarrow \begin{cases} \ell_{jk} &= \left( a_{jk} - \sum_{i=1}^{k-1} \ell_{ji} u_{ik} \right) / u_{kk} \\ u_{kj} &= a_{kj} - \sum_{i=1}^{k-1} \ell_{ki} u_{ij} \end{cases}$$

Summation suffers from the accumulation of rounding errors

Standard model of FP arithmetic:

$$\text{fl}(x \text{ op } y) = (x \text{ op } y)(1 + \delta), \quad |\delta| \leq u, \quad \text{for op} \in \{+, -, \times, \div\}$$

Consider the computation of  $y = \sum_{i=1}^n x_i$  by recursive summation:

$$\begin{aligned} y_2 &= x_1 + x_2 & \Rightarrow & \hat{y}_2 = (x_1 + x_2)(1 + \delta_1) \\ y_3 &= \hat{y}_2 + x_3 & \Rightarrow & \hat{y}_3 = (\hat{y}_2 + x_3)(1 + \delta_2) \\ & & & = (x_1 + x_2) \underbrace{(1 + \delta_1)(1 + \delta_2)}_{\delta_1 \text{ and } \delta_2 \text{ accumulate!}} + x_3(1 + \delta_2) \end{aligned}$$

$$y_4 = \dots \text{etc.}$$

How can we measure the accumulated effect of all rounding errors?

- Let  $y = f(x)$  be computed in finite precision and let  $\hat{y}$  be the computed result
- **Forward error** analysis measures

$$|\hat{y} - y| \text{ (absolute)} \quad \text{or} \quad \frac{|\hat{y} - y|}{|y|} \text{ (relative)}$$

- **Backward error** analysis computes the smallest perturbation  $\Delta x$  such that

$$\hat{y} = f(x + \Delta x)$$

and measures  $|\Delta x|$  (absolute) or  $|\Delta x|/|x|$  (relative).

- Backward error analysis recasts the rounding errors as perturbations of the input data
- An algorithm is **backward stable** if it yields a small backward error, where “small” usually means  $O(u)$

- **Forward error**

$$\eta_{\text{fwd}} = \frac{|\hat{y} - y|}{|y|}$$

- **Backward error**

$$\eta_{\text{bwd}} = \min \left\{ \varepsilon > 0 : \exists \delta x_i, \hat{y} = \sum_{i=1}^n x_i + \delta x_i, |\delta x_i| \leq \varepsilon |x_i| \right\}.$$

Two questions:

- Find a **formula** for  $\eta_{\text{bwd}}$
- Find **bounds** for  $\eta_{\text{bwd}}$  and  $\eta_{\text{fwd}}$  when  $\hat{y}$  is computed in floating-point arithmetic

$$\eta_{\text{bwd}} = \min \left\{ \varepsilon > 0 : \exists \delta x_i, \hat{y} = \sum_{i=1}^n x_i + \delta x_i, |\delta x_i| \leq \varepsilon |x_i| \right\}.$$

We have the formula

$$\eta_{\text{bwd}} = \frac{|\hat{y} - y|}{\sum_{i=1}^n |x_i|}.$$

*Proof:*

- $\frac{|\hat{y} - y|}{\sum_{i=1}^n |x_i|} \leq \eta_{\text{bwd}}$
- $\eta_{\text{bwd}} \leq \frac{|\hat{y} - y|}{\sum_{i=1}^n |x_i|}$  (using  $\delta x_i = (\hat{y} - y) \frac{|x_i|}{\sum_{i=1}^n |x_i|}$ )

As a result we also obtain the formula

$$\kappa = \frac{\eta_{\text{fwd}}}{\eta_{\text{bwd}}} = \frac{\sum_{i=1}^n |x_i|}{\left| \sum_{i=1}^n x_i \right|}.$$

- $\kappa$  is large if  $\sum |x_i| \gg \left| \sum x_i \right| \Rightarrow$  **cancellation**



$$\begin{aligned}y_2 &= x_1 + x_2 \\ \Rightarrow \hat{y}_2 &= (x_1 + x_2)(1 + \delta_1) = x_1(1 + \delta_1) + x_2(1 + \delta_1) \\ y_3 &= \hat{y}_2 + x_3 \\ \Rightarrow \hat{y}_3 &= (\hat{y}_2 + x_3)(1 + \delta_2) \\ &= x_1(1 + \delta_1)(1 + \delta_2) + x_2(1 + \delta_1)(1 + \delta_2) + x_3(1 + \delta_2) \\ &\dots \\ \Rightarrow \hat{y}_n &= \sum_{i=1}^n \left[ x_i \prod_{k=k_i}^n (1 + \delta_k) \right]\end{aligned}$$

$$\begin{aligned}y_2 &= x_1 + x_2 \\ \Rightarrow \hat{y}_2 &= (x_1 + x_2)(1 + \delta_1) = x_1(1 + \delta_1) + x_2(1 + \delta_1) \\ y_3 &= \hat{y}_2 + x_3 \\ \Rightarrow \hat{y}_3 &= (\hat{y}_2 + x_3)(1 + \delta_2) \\ &= x_1(1 + \delta_1)(1 + \delta_2) + x_2(1 + \delta_1)(1 + \delta_2) + x_3(1 + \delta_2) \\ \dots \\ \Rightarrow \hat{y}_n &= \sum_{i=1}^n \left[ x_i \prod_{k=k_i}^n (1 + \delta_k) \right]\end{aligned}$$

## Worst-case fundamental lemma

Let  $\delta_k$ ,  $k = 1 : n$ , such that  $|\delta_k| \leq u$  and  $nu < 1$ . Then

$$\prod_{k=1}^n (1 + \delta_k) = 1 + \theta_n, \quad |\theta_n| \leq \gamma_n := \frac{nu}{1 - nu}.$$

## General algorithm

$$\mathbb{S} = \{x_1, \dots, x_n\}$$

Repeat

Choose any pair  $(x_i, x_j) \in \mathbb{S}^2$  ( $i \neq j$ )

$$\mathbb{S} \leftarrow \mathbb{S} \setminus \{x_i, x_j\}$$

$$\mathbb{S} \leftarrow \mathbb{S} \cup \{x_i + x_j\}$$

until  $\mathbb{S} = \{y\}$

No matter the summation order we have the bound

$$\eta_{\text{bwd}} \leq \gamma_{n-1} = (n-1)u + O(u^2)$$

Consider the computation

$$y = \sum_{i=1}^n x_i$$

In floating-point arithmetic, the forward error  $\eta_{\text{fwd}}$  is bounded by

$$\eta_{\text{fwd}} \leq \eta_{\text{bwd}} \kappa, \quad \eta_{\text{bwd}} \leq \gamma_{n-1} = (n-1)u + O(u^2), \quad \kappa = \frac{\sum |x_i|}{|\sum x_i|}$$

Thus  $\eta_{\text{fwd}}$  can be large when

- The **unit roundoff**  $u$  is large (**low precision**)
- The **dimension**  $n$  is large (**accumulation**)
- The **condition number**  $\kappa$  is large (**cancellation**)

Introduction

**Dealing with accumulation**

Dealing with cancellation

Adaptive precision summation

Conclusion

No matter the summation order we have the bound

$$\eta_{\text{bwd}} \leq \gamma_{n-1} = (n-1)u + O(u^2)$$

⇒ However, for specific orders, we can get much better bounds, and much smaller errors!

Given a summation order to compute  $y = \sum_{i=1}^n x_i$ , we define its associated summation tree as a **binary tree** such that:

- the  $n$  **leaf nodes** are the  $n$  summands  $x_i$
- any **inner node** is equal to the sum of its two children
- the **root node** is the final sum  $y$

Example: recursive summation is a **comb tree**

- For any summation tree, we have the bound:

$$\eta_{\text{bwd}} \leq \gamma_h = hu + O(u^2)$$

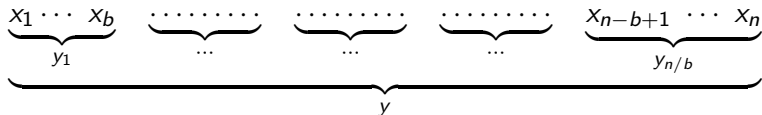
where  $h$  is the height of the tree

- The minimal bound is therefore attained for a **balanced binary tree**, for which  $h = \lceil \log_2 n \rceil$ . This is called **pairwise summation**.
- While it achieves the minimal bound, pairwise summation is not efficient on modern computers.



Blocked summation algorithm:

```
for  $i = 1 : n/b$  do  
  Compute  $y_i = \sum_{j=(i-1)b+1}^{ib} x_j$ .  
end for  
Compute  $y = \sum_{i=1}^{n/b} y_i$ .
```



- Widely used in NLA libraries (BLAS, LAPACK, ...)
- $\eta_{\text{bwd}} \leq \gamma_h$  with  $h = b + n/b - 2$
- With optimal  $b = \sqrt{n}$  :  $h = 2(\sqrt{n} - 1)$

**for**  $i = 1: n/b$  **do**

    Compute  $y_i = \sum_{j=(i-1)b+1}^{ib} x_j$ .

**end for**

Compute  $y = \sum_{i=1}^{n/b} y_i$ .

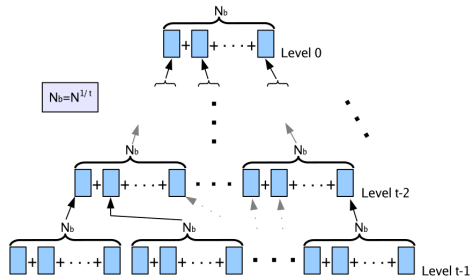
$$\hat{y}_i = \sum_{j=(i-1)b+1}^{ib} \left[ x_j \underbrace{\prod_{k=k_j}^b (1 + \delta_k^{(i)})}_{\text{at most } b-1 \text{ terms}} \right]$$

$$\hat{y} = \sum_{i=1}^{n/b} \left[ \hat{y}_i \underbrace{\prod_{k=k'_i}^{n/b} (1 + \delta'_k)}_{\text{at most } n/b-1 \text{ terms}} \right]$$

$$= \sum_{j=1}^n \left[ x_j \underbrace{\prod_{k=k_j}^b (1 + \delta_k^{(i)}) \prod_{k=k'_j}^{n/b} (1 + \delta'_k)}_{\text{at most } b + n/b - 2 \text{ terms}} \right]$$

# Superblock summation

- **Superblocked summation:** tree summation with  $t$  levels, block size at level  $t$  :  
 $b_t = n^{1/t}$ 
  - $t = 1 \Rightarrow$  standard recursive summation
  - $t = 2 \Rightarrow$  optimal blocked summation
  - $t = \log_2 n \Rightarrow$  pairwise summation
  - $\eta_{\text{bwd}} \leq \gamma_h$  with  $h = t(n^{1/t} - 1)$
  - [Castaldo et al. \(2009\)](#)



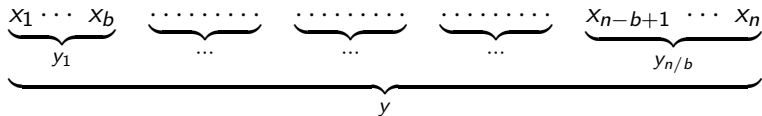
Fast Accurate Blocked summation algorithm (FABsum) [Blanchard, Higham, M. \(2020\)](#)

**for**  $i = 1: n/b$  **do**

    Compute  $y_i = \sum_{j=(i-1)b+1}^{ib} x_j$  with **FastSum**.

**end for**

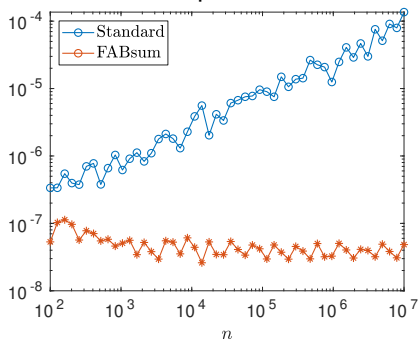
Compute  $y = \sum_{i=1}^{n/b} y_i$  with **AccurateSum**.



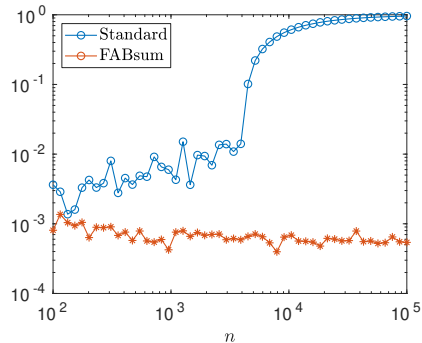
- Cost:  $C(n, b) = \frac{n}{b} C_f(b) + C_a(\frac{n}{b}) \approx C_f(n) + \frac{1}{b} C_a(n)$
- Error:  $\epsilon(n, b) = \epsilon_f(b) + \epsilon_a(n/b) + \epsilon_f(b)\epsilon_a(n/b)$   
 $\Rightarrow$  If  $\epsilon_a(p) = pu^2$  (recursive summation in precision  $u^2$ ), then  $\epsilon(n, b) = bu + O(u^2)$  is independent of  $n$  to first order

Backward error for summing random uniform  $[0, 1]$  data

fp32



fp16



```
for  $i = 1: n/b$  do  
  Compute  $y_i = \sum_{j=(i-1)b+1}^{ib} x_j$ .  
end for  
Compute  $y = \sum_{i=1}^{n/b} y_i$ .
```

- If implemented as is, requires storing  $n/b$  intermediate  $y_i$  values, which requires extra memory and is likely to slow down computation
- Better to implement as follows:

```
 $y = 0$   
for  $i = 1: n/b$  do  
  Compute  $z = \sum_{j=(i-1)b+1}^{ib} x_j$ .  
  Compute  $y = y + z$   
end for
```

Introduction

Dealing with accumulation

**Dealing with cancellation**

Adaptive precision summation

Conclusion

$$[x,y] = \text{Fast2Sum}(a,b)$$

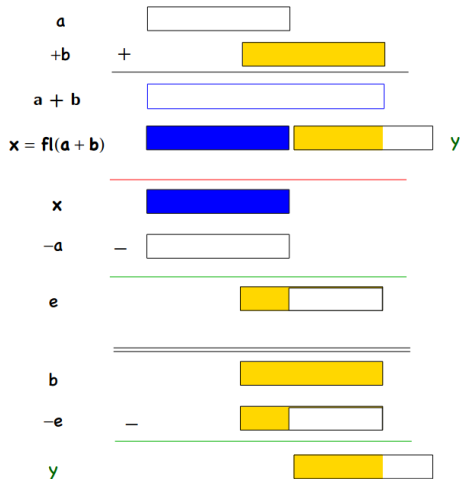
**Input:**  $a, b \in \mathbb{F}$  such that  $|a| \geq |b|$

**Output:**  $x = \text{fl}(a + b), y \in \mathbb{F}$  such that  
 $x + y = a + b$

$$x = a + b$$

$$e = x - a$$

$$y = b - e$$





**Input:**  $x_i \in \mathbb{F}$ ,  $i = 1:n$

**Output:**  $y \approx \sum_{i=1}^n x_i$

$y = 0$

$z = 0$

**for**  $i = 1:n$  **do**

$t = x_i + z$

$[y, z] = \text{Fast2Sum}(y, t)$

**end for**

- Kahan's summation reinjects the errors at each step in the sum
- It satisfies the bound  $\eta_{\text{bwd}} \leq 2u + O(nu^2)$  (proof is quite complicated)

$$\sum_{i=1}^n x_i \xrightarrow{\text{distillation}} \sum_{i=1}^n d_i, \quad \text{where } \kappa(d_i) \ll \kappa(x_i)$$

Fast2Sum:  $\text{fl}(a + b) = a + b + e$ , where  $e \in \mathbb{F}$

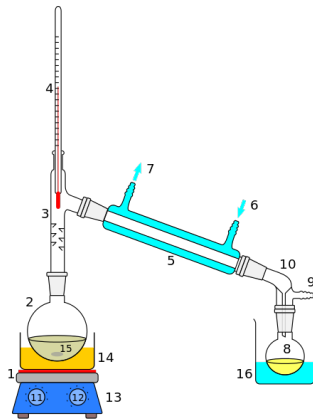
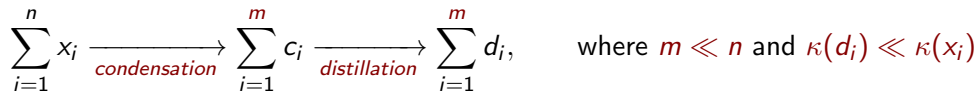
AccSum: repeatedly replace  $(a, b)$  by  $(\text{fl}(a + b), e)$  until the sum is sufficiently well conditioned (higher  $\kappa \Rightarrow$  more iterations)

 Rump, Ogita, Oishi (2008)

# Condensation methods

$$\sum_{i=1}^n x_i \xrightarrow{\text{condensation}} \sum_{i=1}^m c_i \xrightarrow{\text{distillation}} \sum_{i=1}^m d_i, \quad \text{where } m \ll n \text{ and } \kappa(d_i) \ll \kappa(x_i)$$

# Condensation methods



## Conceptual algorithm

$$\mathbb{S} = \{x_1, \dots, x_n\}$$

Repeat for all pairs  $(x_i, x_j) \in \mathbb{S}^2$  ( $i \neq j$ ) such that  $x_i + x_j$  is exact

$$\mathbb{S} \leftarrow \mathbb{S} \setminus \{x_i, x_j\}$$

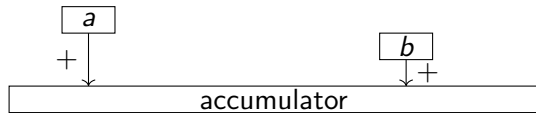
$$\mathbb{S} \leftarrow \mathbb{S} \cup \{x_i + x_j\}$$

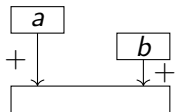
until no such pair remains

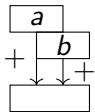
Distill  $\mathbb{S}$

- Can we easily determine when  $x_i + x_j$  is exact?
- Can we bound the maximum number of leftover summands?

# Demmel–Hida method

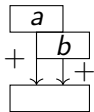






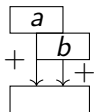
Consider arithmetic with  $f$ -bit mantissa and  $e$ -bit exponent ( $e = 11$  for fp64).





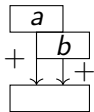
Consider arithmetic with  $f$ -bit mantissa and  $e$ -bit exponent ( $e = 11$  for fp64).

- One big accumulator: [Kulisch](#) method  
... need one accumulator of  $2^e + \log_2 n$  bits



Consider arithmetic with  $f$ -bit mantissa and  $e$ -bit exponent ( $e = 11$  for fp64).

- One big accumulator: **Kulisch** method  
... need one accumulator of  $2^e + \log_2 n$  bits
- One accumulator per exponent:  
**Malcolm** method ... need  $2^e$  accumulators of  $f + \log_2 n$  bits



Consider arithmetic with  $f$ -bit mantissa and  $e$ -bit exponent ( $e = 11$  for fp64).

- One big accumulator: **Kulisch** method ... need one accumulator of  $2^e + \log_2 n$  bits
- One accumulator per exponent: **Malcolm** method ... need  $2^e$  accumulators of  $f + \log_2 n$  bits
- **Demmel–Hida**: general method, balance the number and size of accumulators.

**Input:**  $n$  summands  $x_i$ , number of exponent bits  $m$  to extract

**Output:**  $y = \sum_{j=1}^{2^m} A_j$

Initialize  $A_j = 0$  for  $j = 1, \dots, 2^m$

**for**  $i = 1: n$  **do**

$j \leftarrow m$  leading bits of exponent( $x_i$ )

$A_j \leftarrow A_j + x_i$

**end for**

With  $2^m$  accumulators, need  $F$ -bit mantissa with

$$F \geq f + \lceil \log_2 n \rceil + 2^{e-m} - 1$$

		number of bits				
		signif.	( $t$ )	exp.	range	$u = 2^{-t}$
fp128	quadruple	113		15	$10^{\pm 4932}$	$1 \times 10^{-34}$
fp64	double	53		11	$10^{\pm 308}$	$1 \times 10^{-16}$

Numerical example with fp64 and fp128 arithmetics:

- Assume  $\log_2 n \leq 29$  ( $n \lesssim 0.5 \times 10^9$ )
- $F = 113, f = 53, e = 11 \Rightarrow m$  must thus satisfy

$$\begin{aligned}
 F &\geq f + \lceil \log_2 n \rceil + 2^{e-m} - 1 \\
 \Rightarrow 2^{11-m} &\leq 32 \\
 \Rightarrow 6 &\leq m
 \end{aligned}$$

## Distillation methods (AccSum, etc.)

- 😊 Entirely in the working precision
- 😊 Only uses standard arithmetic operations
- 😞 Strongly dependent on the conditioning
- 😞 Limited parallelism

## Condensation methods (Demmel–Hida, etc.)

- 😊 Independent on the conditioning
- 😊 High level of parallelism
- 😞 Requires access to the exponent
- 😞 Requires extended precision arithmetic

## Distillation methods (AccSum, etc.)

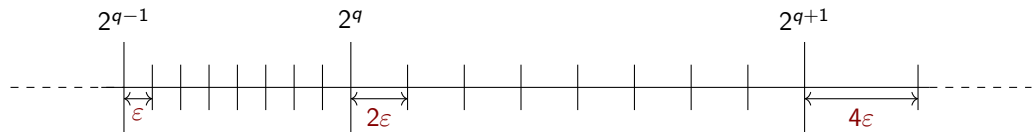
- 😊 Entirely in the working precision
- 😊 Only uses standard arithmetic operations
- 😞 Strongly dependent on the conditioning
- 😞 Limited parallelism

## Condensation methods (Demmel–Hida, etc.)

- 😊 Independent on the conditioning
- 😊 High level of parallelism
- 😞 Requires access to the exponent
- 😞 Requires extended precision arithmetic

**Can we avoid the use of extended precision arithmetic?**

# When is $x + y$ exact? Intuition 1



Let  $x, y \in \mathbb{F} \cap [2^{q-1}, 2^q]$  such that

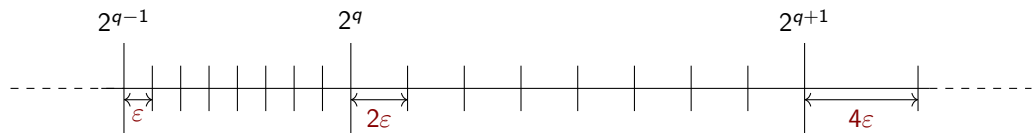
$$x = 2^{q-1} + k_x \epsilon$$

$$y = 2^{q-1} + k_y \epsilon$$

Then

$$\begin{aligned} x + y &= 2^{q-1} + k_x \epsilon + 2^{q-1} + k_y \epsilon \\ &= 2^q + (k_x + k_y) \epsilon \in \mathbb{F} \text{ iff } k_x + k_y \equiv 0 \pmod{2} \end{aligned}$$

# When is $x + y$ exact? Intuition 1



Similarly if

$$x = 2^{q-1} + k_x \epsilon$$

$$y = 2^q + k_y 2\epsilon$$

then  $x + y \in \mathbb{F}$  iff

$$\begin{cases} x + y \leq 2^{q+1} \text{ and } k_x \equiv 0 \pmod{2} \\ x + y > 2^{q+1} \text{ and } k_x + 2k_y \equiv 0 \pmod{4} \end{cases}$$



## When is $x + y$ exact? Intuition 2

$$2^q \times 101 + 2^q \times 111 = 2^q \times 1100 = 2^{q+1} \times 110.0 \in \mathbb{F}$$

$$2^q \times 101 + 2^q \times 110 = 2^q \times 1011 = 2^{q+1} \times 101.1 \notin \mathbb{F}$$

$$2^q \times 101 + 2^{q-1} \times 111 = 2^{q+1} \times 100.01 \notin \mathbb{F}$$

$$2^q \times 101 + 2^{q-1} \times 110 = 2^{q+1} \times 100.00 \in \mathbb{F}$$

## Theorem (Graillat and M.)

Let  $x, y \in \mathbb{F}$  of the same sign  $\sigma = \pm 1$  such that

$$x = \sigma(\beta^{e_x} + k_x \varepsilon_{e_x}),$$

$$y = \sigma(\beta^{e_y} + k_y \varepsilon_{e_y}).$$

Assuming (without loss of generality) that  $|x| \leq |y|$ , then  $x + y \in \mathbb{F}$ , and thus the addition is exact, iff one of the following conditions is met:

- (i)  $x = 0$ ;
- (ii)  $|x + y| < \beta^{e_y+1}$ ,  $e_y - e_x \leq t - 1$ , and  $k_x \equiv 0 \pmod{\beta^{e_y - e_x}}$ ;
- (iii)  $|x + y| = \beta^{e_y+1}$ ,  $e_y + 1 \leq e_{\max}$ ,  $e_y - e_x \leq t - 1$ , and  $k_x \equiv 0 \pmod{\beta^{e_y - e_x}}$ ;
- (iv)  $|x + y| > \beta^{e_y+1}$ ,  $e_y + 1 \leq e_{\max}$ ,  $e_y - e_x \leq t - 2$ , and  $k_x + k_y \beta^{e_y - e_x} \equiv 0 \pmod{\beta^{e_y - e_x + 1}}$ .

# When is $x + y$ exact? Corollary

$$k_x + k_y \beta^{e_y - e_x} \equiv 0 \pmod{\beta^{e_y - e_x + 1}} \xrightarrow{\beta=2, e_x=e_y} k_x + k_y \equiv 0 \pmod{2}$$

## Corollary

If  $x, y \in \mathbb{F}$  with  $\beta = 2$  have the same sign, exponent, and least significant bit, then barring overflow their addition is exact.

Consider the toy example

$$y = 0.25 + 0.3125 + 0.375 + 0.375 + 0.4375 + 0.4375 + 0.625 + 0.625 + 0.75 + 0.75 + 0.875$$

computed with 3-bit arithmetic:

$$\mathbb{F} = \{0.25, 0.3125, 0.375, 0.4375, 0.5, 0.625, 0.75, 0.875, 1, 1.25, 1.5, 1.75, 2, 2.5, 3\}$$

○ LSB=0

□ LSB=1

$e = 1$

$e = 0$

□ 0.625 □ 0.625 ○ 0.75 ○ 0.75 □ 0.875  $e = -1$

○ 0.25 □ 0.3125 ○ 0.375 ○ 0.375 □ 0.4375 □ 0.4375  $e = -2$

# Graillat–Mary method

Consider the toy example

$$y = 0.25 + 0.3125 + 0.375 + 0.375 + 0.4375 + 0.4375 + 0.625 + 0.625 + 0.75 + 0.75 + 0.875$$

computed with 3-bit arithmetic:

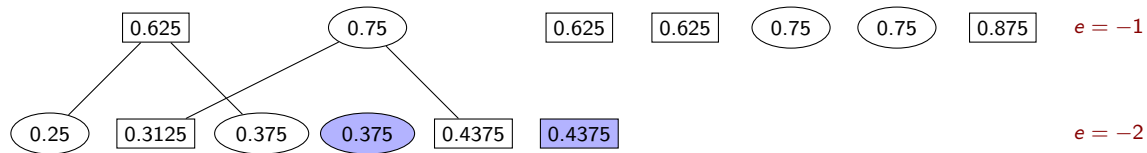
$$\mathbb{F} = \{0.25, 0.3125, 0.375, 0.4375, 0.5, 0.625, 0.75, 0.875, 1, 1.25, 1.5, 1.75, 2, 2.5, 3\}$$

○ LSB=0

□ LSB=1

$e = 1$

$e = 0$



# Graillat–Mary method

Consider the toy example

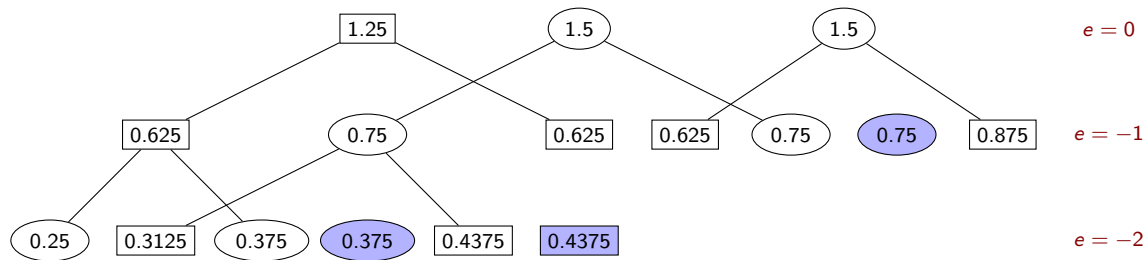
$$y = 0.25 + 0.3125 + 0.375 + 0.375 + 0.4375 + 0.4375 + 0.625 + 0.625 + 0.75 + 0.75 + 0.875$$

computed with 3-bit arithmetic:

$$\mathbb{F} = \{0.25, 0.3125, 0.375, 0.4375, 0.5, 0.625, 0.75, 0.875, 1, 1.25, 1.5, 1.75, 2, 2.5, 3\}$$

○ LSB=0

□ LSB=1



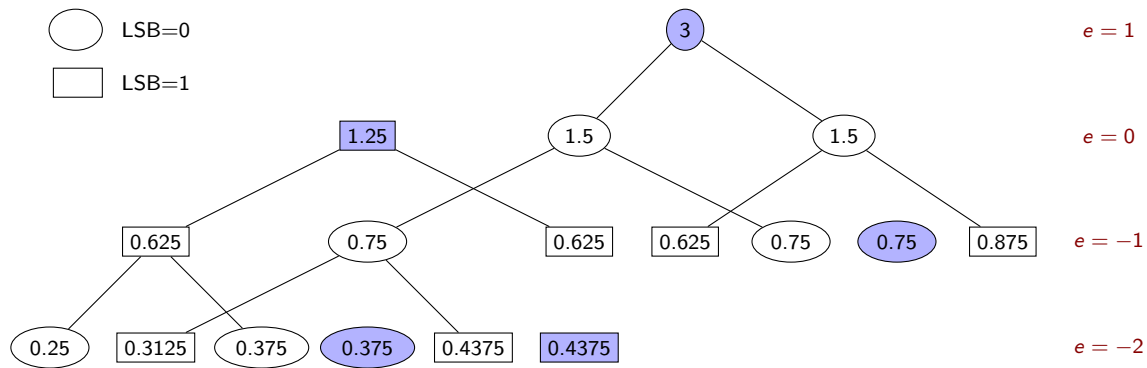
# Graillat–Mary method

Consider the toy example

$$y = 0.25 + 0.3125 + 0.375 + 0.375 + 0.4375 + 0.4375 + 0.625 + 0.625 + 0.75 + 0.75 + 0.875$$

computed with 3-bit arithmetic:

$$\mathbb{F} = \{0.25, 0.3125, 0.375, 0.4375, 0.5, 0.625, 0.75, 0.875, 1, 1.25, 1.5, 1.75, 2, 2.5, 3\}$$



**Input:**  $n$  summands  $x_i$  and a distillation method `distill`

**Output:**  $y = \sum_{i=1}^n x_i$

Initialize  $\text{Acc}(e, s, b)$  to 0 for  $e = e_{\min} : e_{\max}$ ,  
 $s \in \{-1, 1\}$ ,  $b \in \{0, 1\}$ .

**for** all  $x_i$  in any order **do**

$e = \text{exponent}(x_i)$

$s = \text{sign}(x_i)$

$b = \text{LSB}(x_i)$

`insert` ( $\text{Acc}, x_i, e, s, b$ )

**end for**

$x_{\text{condensed}} = \text{gather}(\text{Acc})$

$y = \text{distill}(x_{\text{condensed}})$

**function** `insert` ( $\text{Acc}, x, e, s, b$ )

**if**  $\text{Acc}(e, s, b) = 0$  **then**

$\text{Acc}(e, s, b) = x$

**else**

$x' = \text{Acc}(e, s, b) + x$

$\text{Acc}(e, s, b) = 0$

$b' = \text{LSB}(x')$

`insert`( $\text{Acc}, x', e + 1, s, b'$ )

**end if**

**end function**

**function**  $x_{\text{condensed}} = \text{gather}(\text{Acc})$

$i = 0$

**for** all nonzero  $\text{Acc}(e, s, b)$  **do**

$i = i + 1$

$x_{\text{condensed}}(i) = \text{Acc}(e, s, b)$

**end for**

**end function**



## Conceptual algorithm

$$\mathbb{S} = \{x_1, \dots, x_n\}$$

Repeat for all pairs  $(x_i, x_j) \in \mathbb{S}^2$  ( $i \neq j$ ) such that  $x_i + x_j$  is exact

$$\mathbb{S} \leftarrow \mathbb{S} \setminus \{x_i, x_j\}$$

$$\mathbb{S} \leftarrow \mathbb{S} \cup \{x_i + x_j\}$$

until no such pair remains

Distill  $\mathbb{S}$

- Can we easily determine when  $x_i + x_j$  is exact? **YES!** It suffices to check the sign, exponent, and LSB of  $x_i$  and  $x_j$
- Can we bound the maximum number of leftover summands? **YES!** At most  $4L$  summands where  $L$  is the depth of the tree

$$L \leq \lceil \log_2 n \rceil + d$$

where  $d$  is independent of  $n$  and depends on the range of the values (at most 2047 in binary64)

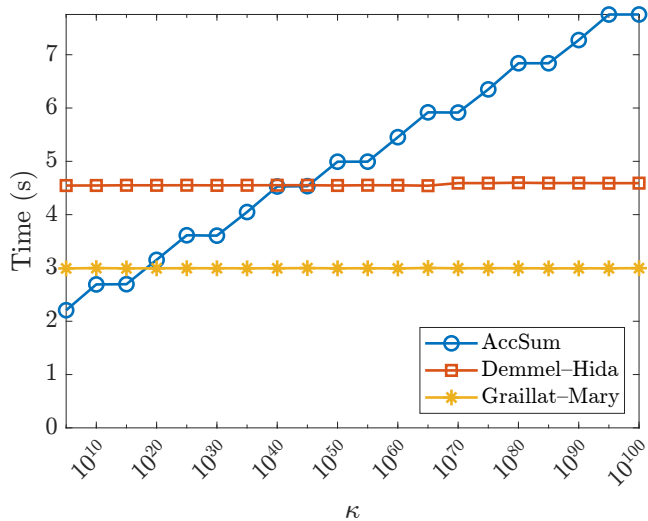
## Distillation methods (AccSum, etc.)

- 😊 Entirely in the working precision
- 😊 Only uses standard arithmetic operations
- 😞 Strongly dependent on the conditioning
- 😞 Limited parallelism

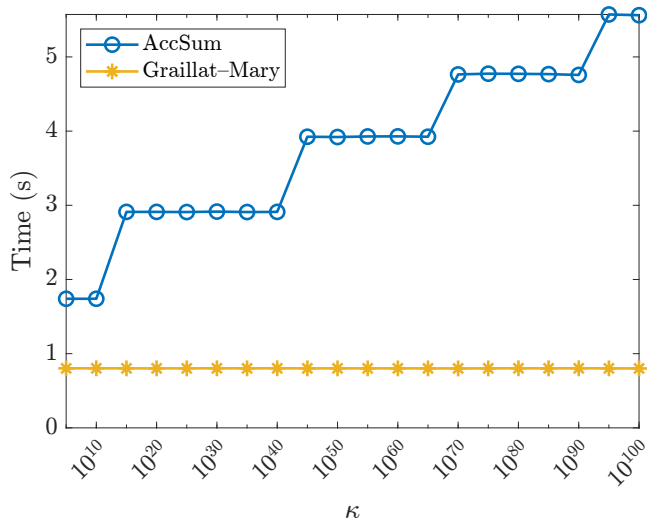
## Condensation methods (Demmel–Hida, Graillat–Mary)

- 😊 Independent on the conditioning
- 😊 High level of parallelism
- 😞 Requires access to the exponent + LSB
- 😞 ~~Requires extended precision arithmetic~~

# Performance comparison



# Quadruple working precision



Introduction

Dealing with accumulation

Dealing with cancellation

**Adaptive precision summation**

Conclusion

- Given an algorithm and a prescribed accuracy  $\varepsilon$ , adaptively select the minimal precision for each instruction depending on the data

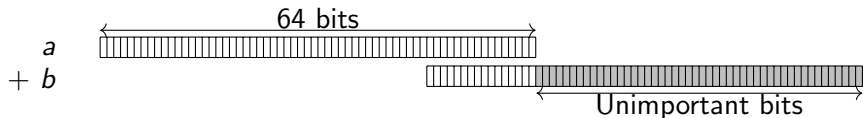
⇒ **First of all, why should the precisions vary?**

- Given an algorithm and a prescribed accuracy  $\varepsilon$ , adaptively select the minimal precision for each instruction depending on the data

⇒ **First of all, why should the precisions vary?**

- Because not all computations are equally “important”!

Example:



⇒ **Opportunity for mixed precision:** adapt the precisions to the data at hand by storing and computing “less important” (which usually means smaller) data in lower precision

# Sparse matrix–vector product (SpMV)

Goal: compute  $y = Ax$ , where  $A$  is a sparse matrix, with a prescribed accuracy  $\varepsilon$

```
for  $i = 1:m$  do  
     $y_i = \sum_{j \in \text{nnz}_i(A)} a_{ij}x_j$   
end for
```

If computed in precision  $\varepsilon$ ,  $\hat{y}$  satisfies

$$|\hat{y}_i - y_i| \leq n_i \varepsilon \sum_{j \in \text{nnz}_i(A)} |a_{ij}x_j|$$

and thus

$$\|\hat{y} - y\| \leq c\varepsilon \|A\| \|x\| \quad (c = \max_i n_i)$$

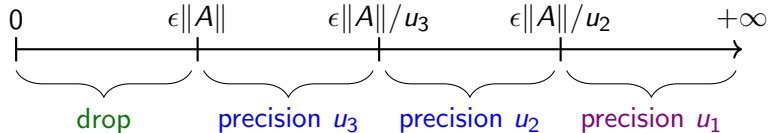
This is a normwise backward error bound:  $\hat{y} = (A + E)x$ ,  $\|E\| \leq c\varepsilon \|A\|$ .



- Given  $p$  available precisions  $u_1 < \epsilon < u_2 < \dots < u_p$ , define partition  $A = \sum_{k=1}^p A^{(k)}$  where

$$a_{ij}^{(k)} = \begin{cases} \text{fl}_k(a_{ij}) & \text{if } |a_{ij}| \in (\epsilon \|A\| / u_k, \epsilon \|A\| / u_{k+1}] \\ 0 & \text{otherwise} \end{cases}$$

$\Rightarrow$  the precision of each element is chosen **inversely proportional to its magnitude**



$$\begin{pmatrix} \times & & \times \\ \times & \times & \\ & \times & \times \end{pmatrix} = \begin{pmatrix} d & & \\ & d & \\ & & d \end{pmatrix} + \begin{pmatrix} & s \\ & & s \\ s & & \end{pmatrix} + \begin{pmatrix} & & \\ h & & \\ & & \end{pmatrix}$$

```

for  $i = 1:m$  do
  for  $k = 1:p$  do
     $y_i^{(k)} = \sum_{j \in \text{nnz}_i(A^{(k)})} a_{ij}^{(k)} x_j$  in precision  $u_k$ 
  end for
   $y_i = \sum_{k=1}^p y_i^{(k)}$  in precision  $u_1$ 
end for

```

- Compute  $y^{(k)} = A^{(k)}x$  in precision  $u_k$ . The computed  $\hat{y}^{(k)}$  satisfies

$$|\hat{y}_i^{(k)} - y_i^{(k)}| \leq (n_i^{(k)})^2 \varepsilon \|A\| \|x\|$$

- Compute  $y = \sum_{k=1}^p y^{(k)}$  in precision  $u_1$ . The computed  $\hat{y}$  satisfies

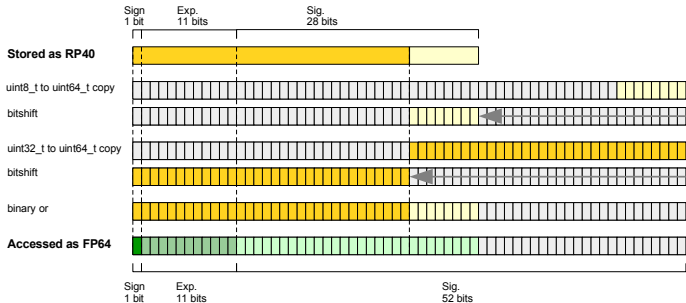
$$\begin{aligned} \hat{y}_i &= \sum_{k=1}^p \hat{y}_i^{(k)} + e_i, & |e_i| &\leq pu_1 \|A\| \|x\| \\ &= y_i + f_i, & |f_i| &\leq c\varepsilon \|A\| \|x\| \end{aligned}$$

The more precisions we have, the more we can reduce storage  $\Rightarrow$  can we exploit custom precision formats?

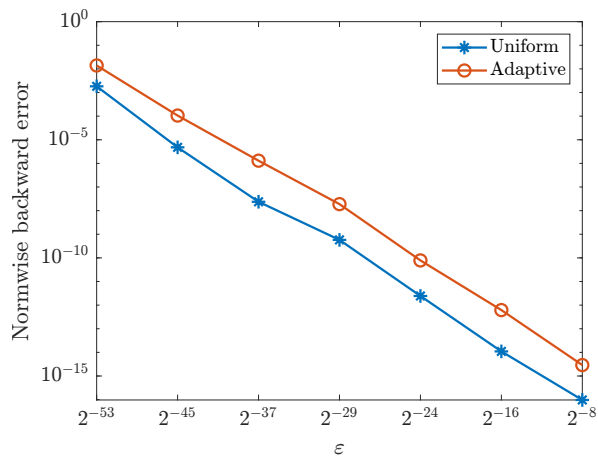
Emulated formats				
Format	Signif. ( $t$ )	Exponent	Range	$u = 2^{-t}$
bf16	8	8	$10^{\pm 38}$	$4 \times 10^{-3}$
fp24	16	8	$10^{\pm 38}$	$2 \times 10^{-5}$
fp32	24	8	$10^{\pm 38}$	$6 \times 10^{-8}$
fp40	29	11	$10^{\pm 308}$	$2 \times 10^{-9}$
fp48	37	11	$10^{\pm 308}$	$8 \times 10^{-12}$
fp56	45	11	$10^{\pm 308}$	$3 \times 10^{-14}$
fp64	53	11	$10^{\pm 308}$	$1 \times 10^{-16}$

How to efficiently implement custom precision storage?

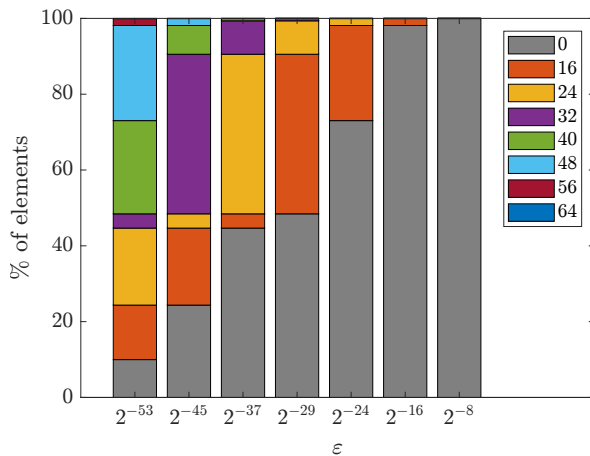
```
union union64 {  
  uint64_t i;  
  double f;  
};  
double RpToFp (rp40 rp, size_t i){  
  union union64 u64;  
  uint64_t i64h, i64l;  
  i64h = (uint64_t)rp.i32[i];  
  i64h = i64h << 32;  
  i64l = (uint64_t)rp.i8[i];  
  i64l = i64l << 24;  
  u64.i = i64h | i64l;  
  return u64.f;  
}
```



[Grailat, Jézéquel, M., Molina, Mukunoki \(2024\)](#)

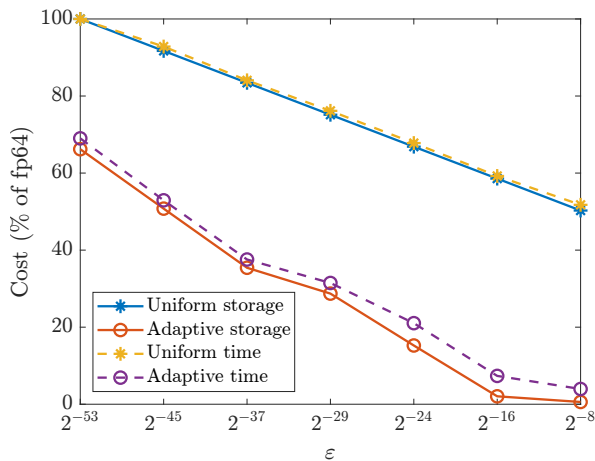


- Controlled accuracy



- Controlled accuracy

# Experimental results (Long\_Coup\_dt6 matrix, $n \approx 1.5M$ )



- Controlled accuracy
- Storage reduced by at least 30% and potentially much more for larger  $\epsilon$ .
- **Time cost matches storage.**

Introduction

Dealing with accumulation

Dealing with cancellation

Adaptive precision summation

**Conclusion**



$$\eta_{\text{fwd}} \leq \eta_{\text{bwd}} \kappa, \quad \eta_{\text{bwd}} \leq \gamma_{n-1} = (n-1)u + O(u^2), \quad \kappa = \frac{\sum |x_i|}{|\sum x_i|}$$

- We have seen various summation methods with different properties/objectives: handling error accumulation, cancellation, using mixed precision. . .
- A common theme has been the **reordering of the summands by grouping them into blocks/buckets**,
  - either fixed-size groups of arbitrary summands
  - or groups of summands of similar magnitude.
- We have seen several possible uses of **mixed precision arithmetic**:
  - Mixed precision blocked summation (FABsum): reduce accumulation  
 $\Rightarrow \eta_{\text{bwd}}$  independent of  $n$
  - Bucket summation with extended precision (Demmel-Hida): reduce cancellation  
 $\Rightarrow \eta_{\text{fwd}}$  independent of  $\kappa$
  - Bucket summation with adaptive precision: exploit lower precisions while controlling  $\eta_{\text{bwd}}$

- You are given a mysterious sum to evaluate as accurately and efficiently as possible.  
**Goal:** achieve close to  $10^{-16}$  accuracy while maintaining a time cost comparable to recursive summation.
- Use of MATLAB's `sum` is obviously forbidden!
- Suggestions:
  - Implement Kahan's summation (slide 24).
  - Implement blocked summation (slide 16). How should you choose the block size  $b$ ?
  - Implement FABsum (slide 19) with Kahan's summation as `AccurateSum`. How should you choose the block size  $b$ ?
  - Compare performance–accuracy tradeoffs.
  - Remember slide number 21.