

Numerical Accuracy and Reproducibility at ExaScale

James Demmel, **Hong Diep Nguyen**

ParLab - EECS - UC Berkeley

ARITH 21 April 7-10, 2013

Floating-point arithmetic: defines a discrete subset of real values and suffers from *rounding errors*.

→ Floating-point operations $(+, \times)$ are commutative but not associative:

$$(-1 + 1) + 2^{-53} \neq -1 + (1 + 2^{-53}).$$

Consequence: results of floating-point computations depend on the order of computation.

Reproducibility: ability to obtain bit-wise identical results from run-to-run on the same input data, with different resources.

Motivations

Demands for reproducible floating-point computations:

- ▶ Debugging: look inside the code step-by-step, and might need to rerun multiple times on the same input data.
- ▶ Understanding the reliability of output. Ex: ¹, Power State Estimation problem (spmv + dot product), after the 5th step the Euclidean norm of the residual vector differs up to 20% from one run to another.
- ▶ Contractual reasons (road type approval, drug design),
- ▶ ...

¹Villa et al, *Effects of Floating-point non-Associativity on Numerical Computations on Massively Multithreaded Systems*, CUG 2009 Proceedings

Sources of non-reproducibility

A performance-optimized floating-point library is prone to inconsistency for various reasons:

- ▶ Changing Data Layouts:
 - ▶ Data alignment,
 - ▶ Data partitioning,
 - ▶ Data ordering,
- ▶ Changing Hardware Resources:
 - ▶ Fused Multiply-Adder support,
 - ▶ Intermediate precision (64 bits, 80 bits, 128 bits, etc),
 - ▶ Data path (SSE, AVX, GPU warp, etc),
 - ▶ Cache line size,
 - ▶ Number of processors,
 - ▶ Network topology,
 - ▶ ???

Sources of non-reproducibility

A performance-optimized floating-point library is prone to inconsistency for various reasons:

- ▶ Changing Data Layouts:
 - ▶ Data alignment,
 - ▶ Data partitioning,
 - ▶ Data ordering,
- ▶ Changing Hardware Resources:
 - ▶ Fused Multiply-Adder support,
 - ▶ Intermediate precision (64 bits, 80 bits, 128 bits, etc),
 - ▶ Data path (SSE, AVX, GPU warp, etc),
 - ▶ Cache line size,
 - ▶ Number of processors,
 - ▶ Network topology,
 - ▶ ???

Reproducibility at ExaScale

ExaScale: ability to execute exaflops (10^{18} floating-point operations) per second, using $\mathcal{O}(10^9)$ processors.

Reproducibility at ExaScale

ExaScale: ability to execute exaflops (10^{18} floating-point operations) per second, using $\mathcal{O}(10^9)$ processors.

- ▶ Highly dynamic scheduling,

Reproducibility at ExaScale

ExaScale: ability to execute exaflops (10^{18} floating-point operations) per second, using $\mathcal{O}(10^9)$ processors.

- ▶ Highly dynamic scheduling,
- ▶ Network heterogeneity: reduction tree shape can vary,

Reproducibility at ExaScale

ExaScale: ability to execute exaflops (10^{18} floating-point operations) per second, using $\mathcal{O}(10^9)$ processors.

- ▶ Highly dynamic scheduling,
- ▶ Network heterogeneity: reduction tree shape can vary,
- ▶ Drastically increased communication time

$$\text{Cost} = \underbrace{\text{Arithmetic}}_{\text{FLOPs}} + \underbrace{\text{Communication}}_{\# \text{words moved} + \# \text{messages}}$$

Reproducibility at ExaScale

ExaScale: ability to execute exaflops (10^{18} floating-point operations) per second, using $\mathcal{O}(10^9)$ processors.

- ▶ Highly dynamic scheduling,
- ▶ Network heterogeneity: reduction tree shape can vary,
- ▶ Drastically increased communication time

$$\text{Cost} = \underbrace{\text{Arithmetic}}_{\text{FLOPs}} + \underbrace{\text{Communication}}_{\# \text{words moved} + \# \text{messages}}$$

- ▶ Communication-Avoiding algorithms change the order of computation **on purpose**, for ex. 2.5D Matmult, 2.5D LU, etc,

Reproducibility at ExaScale

ExaScale: ability to execute exaflops (10^{18} floating-point operations) per second, using $\mathcal{O}(10^9)$ processors.

- ▶ Highly dynamic scheduling,
- ▶ Network heterogeneity: reduction tree shape can vary,
- ▶ Drastically increased communication time

$$\text{Cost} = \frac{\text{Arithmetic}}{\text{FLOPs}} + \text{Communication}$$

$\# \text{words moved} + \# \text{messages}$

- ▶ Communication-Avoiding algorithms change the order of computation **on purpose**, for ex. 2.5D Matmult, 2.5D LU, etc,
- ▶ A little extra arithmetic cost is allowed so long as the **communication cost is controlled**.

State of the art

Source of floating-point non-reproducibility: rounding errors lead to dependence of computed result on order of computations.

To obtain reproducibility ?

State of the art

Source of floating-point non-reproducibility: rounding errors lead to dependence of computed result on **order of computations**.

To obtain reproducibility ?

- ▶ Fix the order of computations:
 - ▶ sequential mode: *intolerably costly at ExaScale*
 - ▶ fixed reduction tree: *substantial communication overhead*

State of the art

Source of floating-point non-reproducibility: **rounding errors** lead to dependence of computed result on **order of computations**.

To obtain reproducibility ?

- ▶ Fix the order of computations:
 - ▶ sequential mode: *intolerably costly at ExaScale*
 - ▶ fixed reduction tree: *substantial communication overhead*
- ▶ Eliminate/Reduce the rounding errors:
 - ▶ exact arithmetic (rounded at the end): *much more expensive in communication and very wide multi-word arithmetic*
 - ▶ fixed-point arithmetic: *limited range of values*
 - ▶ higher precision: *reproducible with high probability (not certain)*.

State of the art

Source of floating-point non-reproducibility: **rounding errors** lead to dependence of computed result on **order of computations**.

To obtain reproducibility ?

- ▶ Fix the order of computations:
 - ▶ sequential mode: *intolerably costly at ExaScale*
 - ▶ fixed reduction tree: *substantial communication overhead*
- ▶ Eliminate/Reduce the rounding errors:
 - ▶ exact arithmetic (rounded at the end): *much more expensive in communication and very wide multi-word arithmetic*
 - ▶ fixed-point arithmetic: *limited range of values*
 - ▶ higher precision: *reproducible with high probability (not certain)*.
- ▶ Our proposed solution: **deterministic errors**.

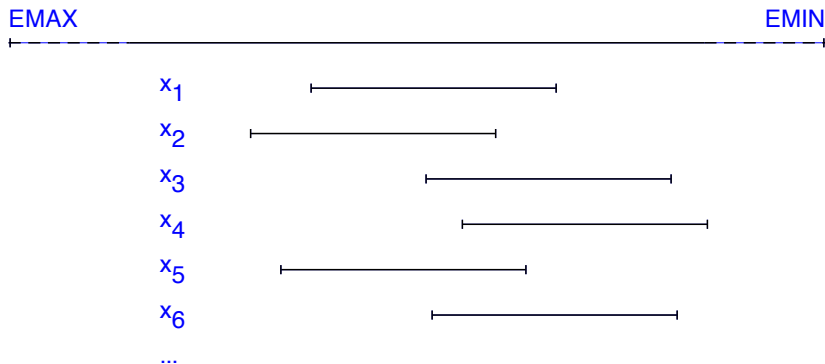
A proposed solution for global sum

Objectives:

- ▶ bit-wise identical results from run-to-run regardless of hardware heterogeneity, # processors, reduction tree shape, ...
- ▶ independent of data ordering,
- ▶ only 1 reduction per sum,
- ▶ no severe loss of accuracy.

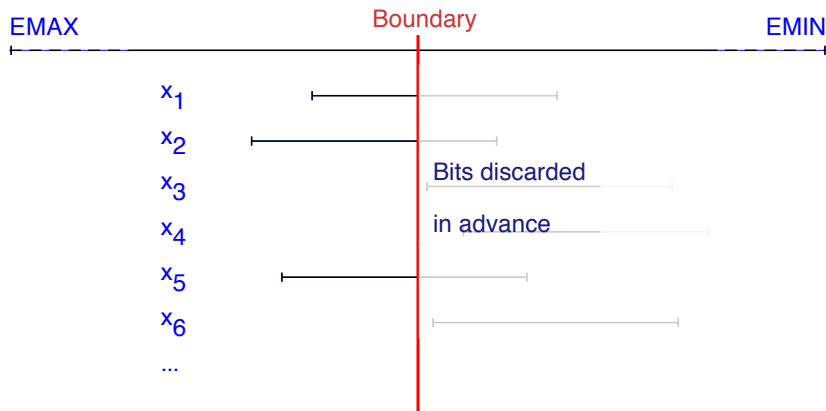
Idea: pre-rounding input values.

Pre-rounding technique



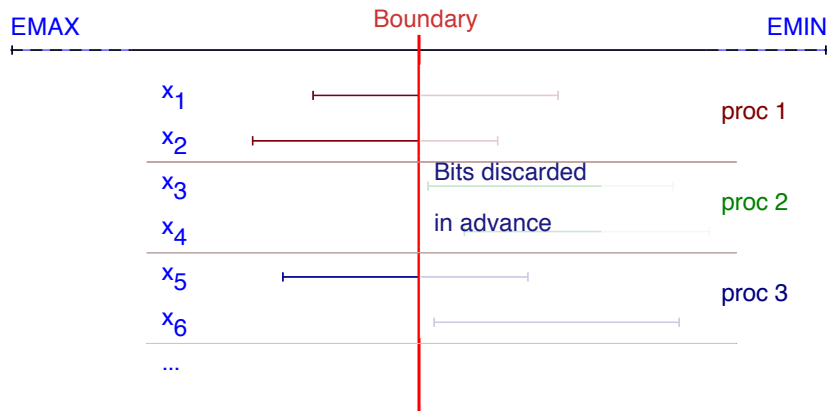
Rounding occurs at each addition. Computation's error depends on the intermediate results, which depend on the order of computation.

Pre-rounding technique



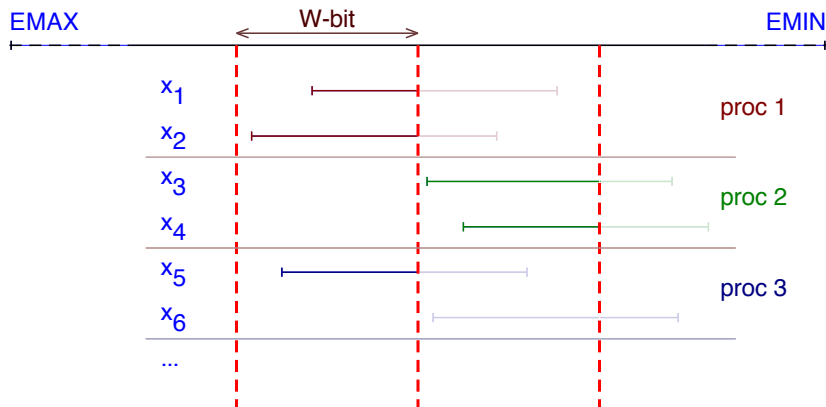
No rounding error at each addition. Computation's error depends on the $Boundary$, which depends on $\max |x_i|$, not on the ordering

Pre-rounding technique



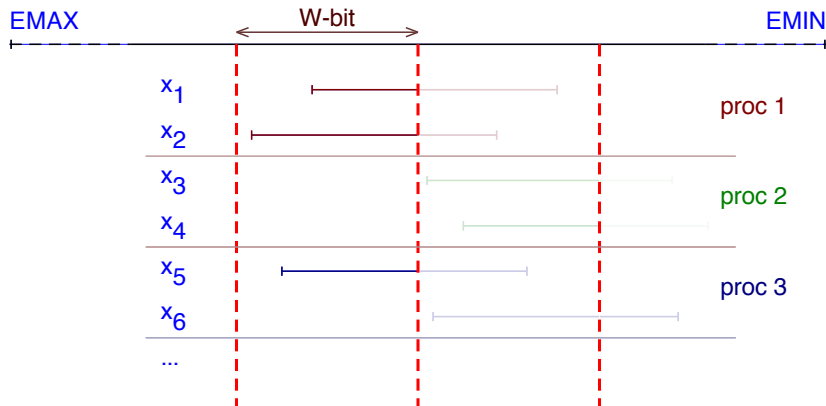
No rounding error at each addition. Computation's error depends on the Boundary, which depends on $\max |x_i|$, not on the ordering \Rightarrow extra communication among processors.

1-Reduction technique



Boundaries are precomputed. Special Reduction Operator:
(**MAX** of boundaries combined **SUM** of corresponding partial sums)

1-Reduction technique



Boundaries are precomputed. Special Reduction Operator:
(**MAX** of boundaries combined **SUM** of corresponding partial sums)

Reproducibility vs. Accuracy

Reproducibility does not necessarily mean Accuracy:

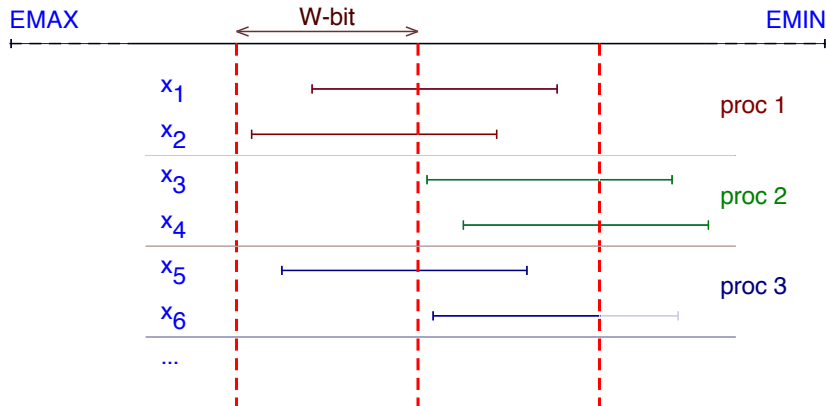
- ▶ A very accurate result (for ex. computed using very high precision) might not be reproducible,
- ▶ A reproducible result might be of much less accuracy.

The accuracy of the pre-rounding technique with 1 bin:

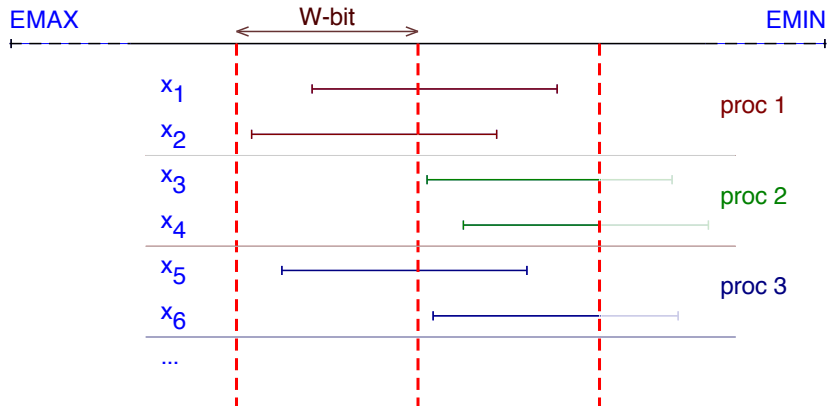
$$\text{absolute error} \leq N \cdot \text{Boundary} < N \cdot \max |x_i|.$$

Solution: increase the number of bins to improve the accuracy.

k-fold Algorithm



k-fold Algorithm



k-fold Algorithm: Accuracy

k-fold algorithm has an error bound:

$$\text{absolute error} \leq N \cdot \text{Boundary}_k < N \cdot 2^{-(k-1) \cdot W} \cdot \max |x_i|.$$

In practice: $k = 3$, $W = 40$.

$$\text{absolute error} < N \cdot 2^{-80} \cdot \max |x_i| = 2^{-27} \cdot N \cdot \epsilon \cdot \max |x_i|$$

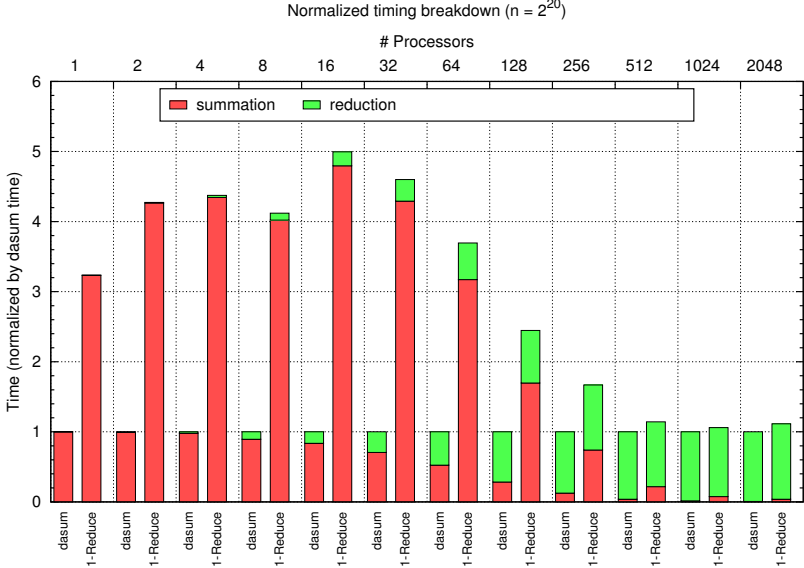
$$\text{Standard sum's error bound} \leq (N - 1) \cdot \epsilon \cdot \sum |x_i|$$

Experimental results: Accuracy

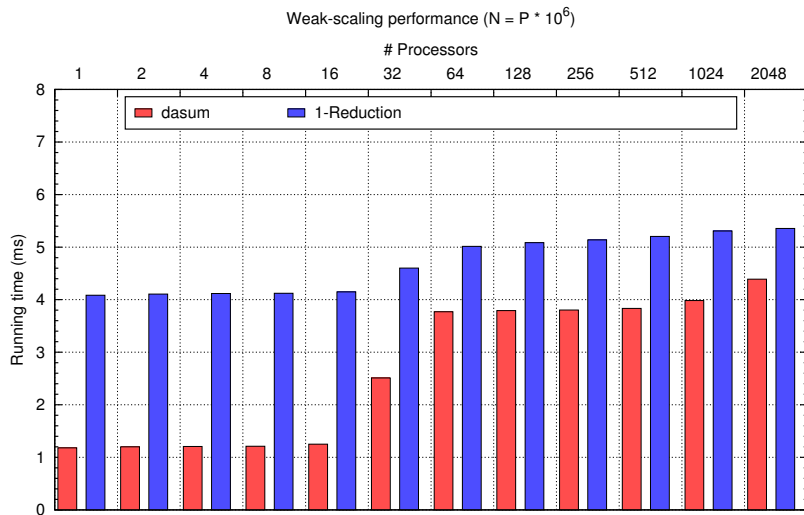
Summation of $n = 10^6$ floating-point numbers. Computed results of both reproducible summation and standard summation (with different ordering: ascending value, descending value, ascending magnitude, descending magnitude) are compared with result computed using quad-double precision.

| Generator x_i | reproducible | standard |
|---|-----------------------|---|
| <code>drand48()</code> | 0 | $-8.5 \times 10^{-15} \div 1.0 \times 10^{-14}$ |
| <code>drand48() - 0.5</code> | 1.5×10^{-16} | $-1.7 \times 10^{-13} \div 1.8 \times 10^{-13}$ |
| <code>sin(2.0 * pi * i/n)</code> | 1.5×10^{-15} | $-1.0 \div 1.0$ |
| <code>sin(2.0 * pi * i/n) * 2⁻⁵</code> | 1.0 | $-1.0 \div 1.0$ |

Experimental results: strong-scaling



Experimental results: weak-scaling



Conclusions

The proposed 1-Reduction pre-rounding technique

- ▶ provides **bit-wise identical reproducibility**, regardless of
 - ▶ data permutation, data assignment,
 - ▶ processor count, reduction tree shape,
 - ▶ hardware heterogeneity, etc.
- ▶ obtains better error bound than the standard sum's,
- ▶ can be done in on-the-fly mode,
- ▶ requires only **ONE** reduction for the global parallel summation,
- ▶ is suitable for very large scale systems (ExaScale),
- ▶ can be applied to Cloud computing environment,
- ▶ can be applied to other operations which use summation as the reduction operator.

Future works

In Progress

- ▶ Parallel Prefix Sum,
- ▶ Matrix-vector / Matrix-matrix multiplication,

TODO

- ▶ Higher level driver routine: `trsm`, factorizations like LU, ...
- ▶ n.5D algorithms (2.5D Matmult, 2.5D LU),
- ▶ `spMV`,
- ▶ Other associative operations.

1-Reduction's reduction operator

