# Reproducibility of Linear Algebra Operations

**Roman Iakymchuk**[1]

joint work with

Sylvain Collange, David Defour, Erwin Laure,
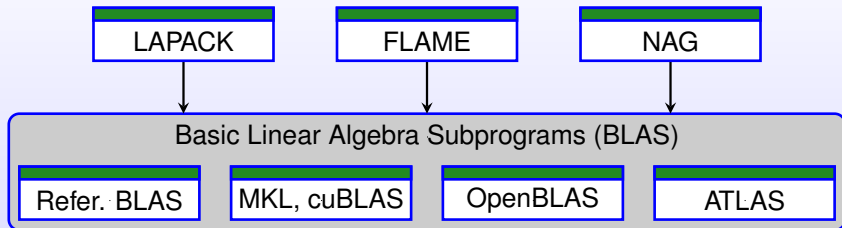Enrique S. Quitana Ortí and Stef Graillat

[1]KTH Royal Institute of Technology, CSC, CST/PDC
`riakymch@kth.se`

MS237: Algorithmic Revolution in Post Moore's Era:
Auto-Tuning and Accuracy Assurance
SIAM CSE17, Feb 27th - Mar 3rd, 2017, Atlanta, GA, USA

# Linear Algebra Libraries



LAPACK    FLAME    NAG

Basic Linear Algebra Subprograms (BLAS)

Refer. BLAS    MKL, cuBLAS    OpenBLAS    ATLAS

$\Downarrow$

| BLAS-1 [1979]: | $y := y + \alpha x$ | $\alpha \in \mathbb{R}; x, y \in \mathbb{R}^n$ | $2/3$ |
|---|---|---|---|
| | $\alpha := \alpha + x^T y$ | | |
| BLAS-2 [1988]: | $A := A + xy^T$ | $A \in \mathbb{R}^{n \times n}; x, y \in \mathbb{R}^n$ | $2$ |
| | $y := A^{-1} x$ | | |
| BLAS-3 [1990]: | $C := C + AB$ | $A, B, C \in \mathbb{R}^{n \times n}$ | $n/2$ |
| | $C := A^{-1} B$ | | |

# Goals

- To ensure BLAS kernels yield precise and numerical reproducible results with comparable performance on a wide range of architectures

## ExBLAS – **Ex**act **BLAS**

- **Ex**BLAS-1: ExSUM, ExSCAL, ExDOT, ExAXPY, ...

- **Ex**BLAS-2: ExGER, ExGEMV, ExTRSV, ExSYR, ...

- **Ex**BLAS-3: ExGEMM, ExTRSM, ExSYR2K, ...

- Use the ExBLAS kernels to construct exact higher-level operations such as matrix factorization

# Outline

# Outline

# Computer Arithmetic

## Problems

- Floating-point arithmetic suffers from rounding errors
- Floating-point operations $(+, \times)$ are commutative but non-associative

$$(-1 + 1) + 2^{-53} \neq -1 + (1 + 2^{-53}) \quad \text{in double precision}$$

# Computer Arithmetic

## Problems

- Floating-point arithmetic suffers from rounding errors
- Floating-point operations $(+, \times)$ are commutative but non-associative

$$2^{-53} \neq 0 \quad \text{in double precision}$$

# Computer Arithmetic

## Problems

- Floating-point arithmetic suffers from rounding errors
- Floating-point operations $(+, \times)$ are commutative but non-associative

  $$(-1 + 1) + 2^{-53} \neq -1 + (1 + 2^{-53}) \quad \text{in double precision}$$

- Consequence: results of floating-point computations depend on the order of computation
- Results computed by performance-optimized parallel floating-point libraries may be often inconsistent: each run returns a different result

- **Reproducibility** – ability to obtain bit-wise identical results from run-to-run on the same input data on the same or different architectures

# Sources of Non-Reproducibility

- Changing Data Layouts:
  - Data partitioning
  - Data alignment

- Changing Hardware Resources
  - Number of threads
  - Fused Multiply-Add support
  - Intermediate precision (64 bits, 80 bits, 128 bits, etc)
  - Data path (SSE, AVX, GPU warp, etc)
  - Number of processors
  - Network topology
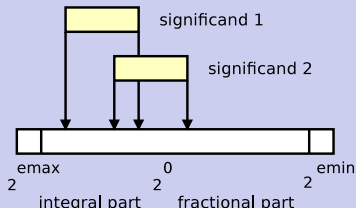
# Exact Multi-Level Parallel Reduction

- Fixed FP expansions (FPE) with Error-Free Transformations
- → Example: double-double or quad-double (Briggs, Bailey, Hida, Li) (work well on a set of relatively close numbers)

---

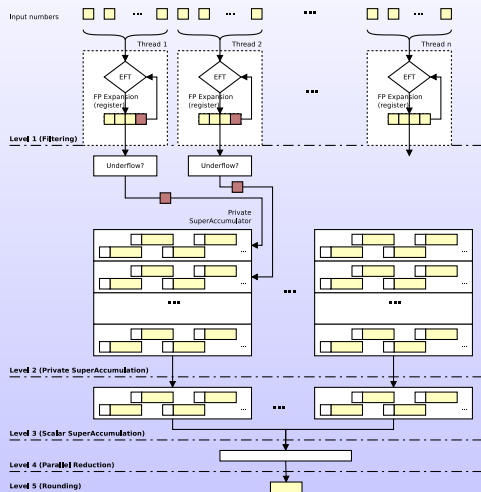**Algorithm 1** FPE of size 2 (Dekker and Knuth)

Function$[r, s]$ = TwoSum$(a, b)$

1: $r \leftarrow a + b$
2: $z \leftarrow r - a$
3: $s \leftarrow (a - (r - z)) + (b - z)$

---

# Exact Multi-Level Parallel Reduction

- Fixed FP expansions (FPE) with Error-Free Transformations
- → Example: double-double or quad-double (Briggs, Bailey, Hida, Li) (work well on a set of relatively close numbers)

---

**Algorithm 1** FPE of size 2 (Dekker and Knuth)

Function$[r, s]$ = TwoSum$(a, b)$

1: $r \leftarrow a + b$
2: $z \leftarrow r - a$
3: $s \leftarrow (a - (r - z)) + (b - z)$

---

- "Infinite" precision: reproducible independently from the inputs
- → Example: Kulisch accumulator (considered inefficient)
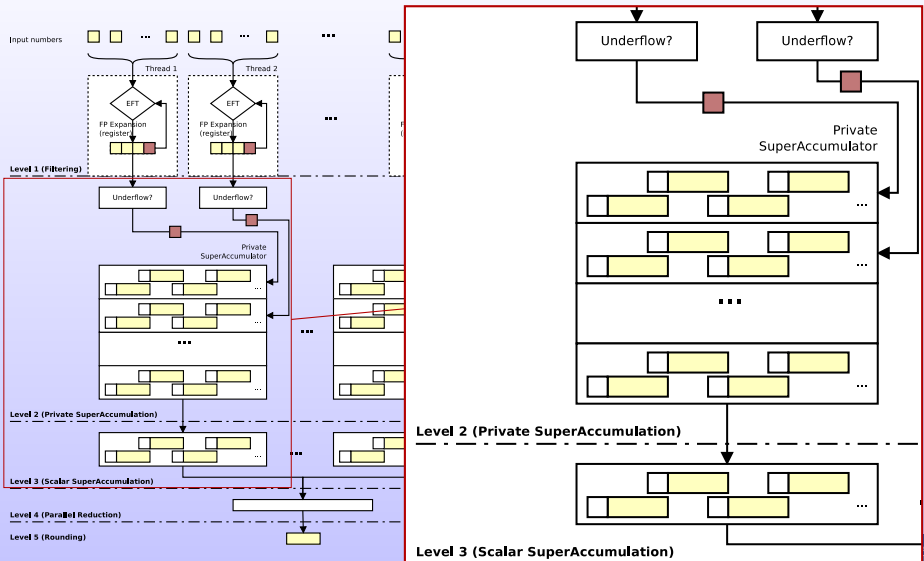
- Parallel algorithm with 5-levels

- Suitable for today's parallel architectures

- Based on FPE with EFT and Kulisch accumulator

- Guarantees "inf" precision
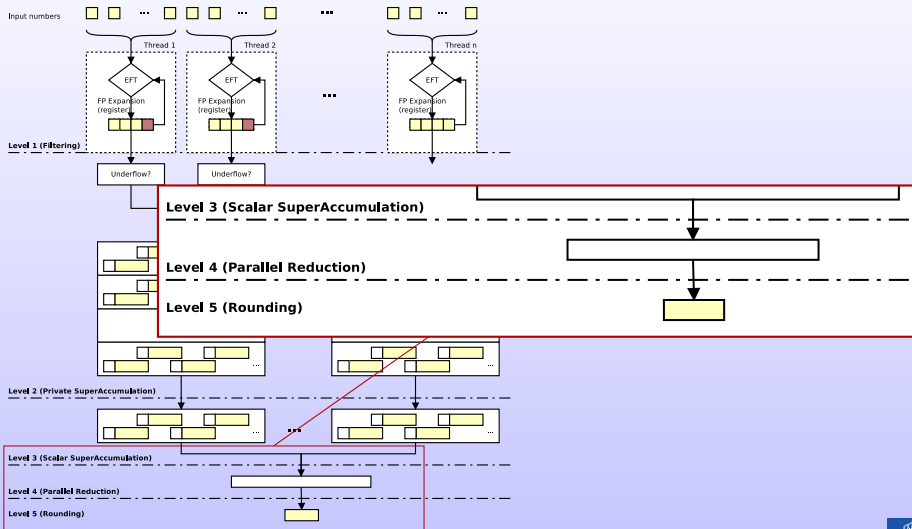
$\rightarrow$ **bit-wise reproducibility**

# Outline

# ExBLAS Highlights

## ExBLAS Status

- ExBLAS-1: ExSUM[a], ExSCAL, ExDOT, ExAXPY, ...

- ExBLAS-2: ExGER, ExGEMV, ExTRSV, ExSYR, ...

- ExBLAS-3: ExGEMM, ExTRSM, ExSYR2K, ...

---
[a]Routines in blue are already in ExBLAS

## ExSCAL

- $x := \alpha * x \rightarrow$ correctly rounded and reproducible

# ExBLAS Highlights

## ExBLAS Status

- ExBLAS-1: `ExSUM`[a], `ExSCAL`, `ExDOT`, `ExAXPY`, ...

- ExBLAS-2: `ExGER`, `ExGEMV`, `ExTRSV`, `ExSYR`, ...

- ExBLAS-3: `ExGEMM`, `ExTRSM`, `ExSYR2K`, ...

---
[a]Routines in blue are already in ExBLAS

## ExSCAL

- $x := \alpha * x \rightarrow$ correctly rounded and reproducible
- Within LU: $x := 1/\alpha * x \rightarrow$ not correctly rounded

# ExBLAS Highlights

## ExBLAS Status

- ExBLAS-1: `ExSUM`[a], `ExSCAL`, `ExDOT`, `ExAXPY`, ...

- ExBLAS-2: `ExGER`, `ExGEMV`, `ExTRSV`, `ExSYR`, ...

- ExBLAS-3: `ExGEMM`, `ExTRSM`, `ExSYR2K`, ...

---
[a]Routines in blue are already in ExBLAS

## ExSCAL

- $x := \alpha * x \rightarrow$ correctly rounded and reproducible
- Within LU: $x := 1/\alpha * x \rightarrow$ not correctly rounded
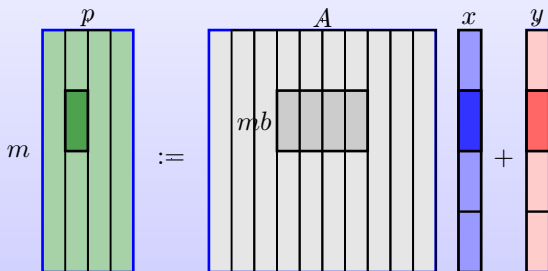- ExInvSCAL: $x := x/\alpha \rightarrow$ correctly rounded and reproducible

## ExGER

- General case: $A := \alpha * x * y^T + A$

# ExBLAS Highlights

## ExBLAS Status

- ExBLAS-1: ExSUM[a], ExSCAL, ExDOT, ExAXPY, ...

- ExBLAS-2: ExGER, ExGEMV, ExTRSV, ExSYR, ...

- ExBLAS-3: ExGEMM, ExTRSM, ExSYR2K, ...

---
[a]Routines in blue are already in ExBLAS

## ExSCAL

- $x := \alpha * x \to$ correctly rounded and reproducible
- Within LU: $x := 1/\alpha * x \to$ not correctly rounded
- ExInvSCAL: $x := x/\alpha \to$ correctly rounded and reproducible

## ExGER

- General case: $A := \alpha * x * y^T + A$
- Within LU: $A := x * y^T + A$. Using FMA $\to$ correctly rounded and reproducible
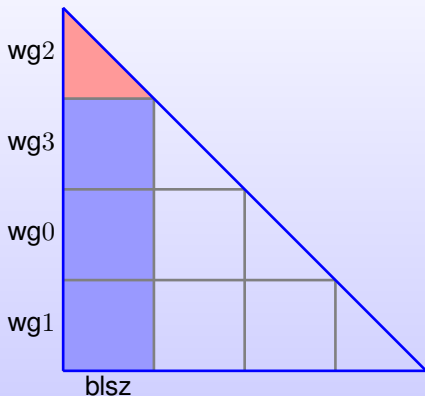
# Matrix-Vector Product

$$\boxed{\text{DGEMV: } y := \alpha A x + \beta y}$$



- Based on ExDOT
- TwoProd$(a, b)$
  1: $r \leftarrow a * b$
  2: $s \leftarrow fma(a, b, -r)$
- $fma(a, b, c) = a * b + c$

## Triangular Solver

```
 1: for i = 0 : blsz : n do
 2:     for k = i : i + blsz do
 3:         for j = 1 : k − 1 do
 4:             [r, e] ← TwoProd(l_kj, −x_j)
 5:             ExpansionAccumulate(r)
 6:             ExpansionAccumulate(e)
 7:         end for
 8:         ExpansionAccumulate(b_k)
 9:         ŝ ← RNDN(acc(k))
10:         x_k ← ŝ/l_kk
11:     end for
12:     for k = i + blsz : n do
13:         for j = i : i + blsz do
14:             [r, e] ← TwoProd(l_kj, −x_j)
15:             ExpansionAccumulate(r)
16:             ExpansionAccumulate(e)
17:         end for
18:     end for
19: end for
```
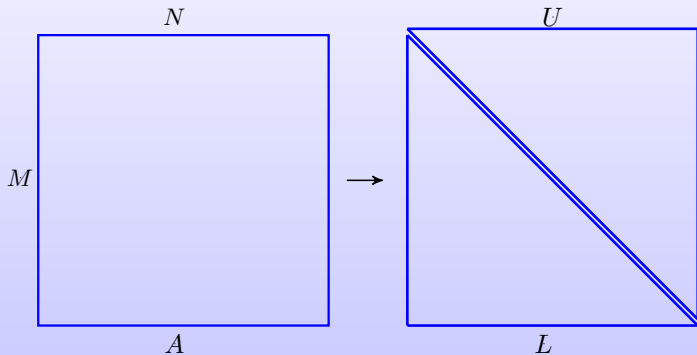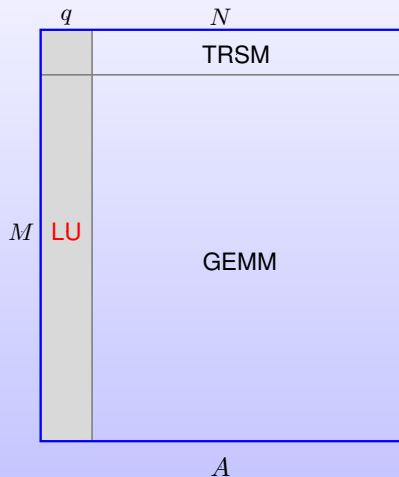


Partitioning of a lower triangular matrix $L$
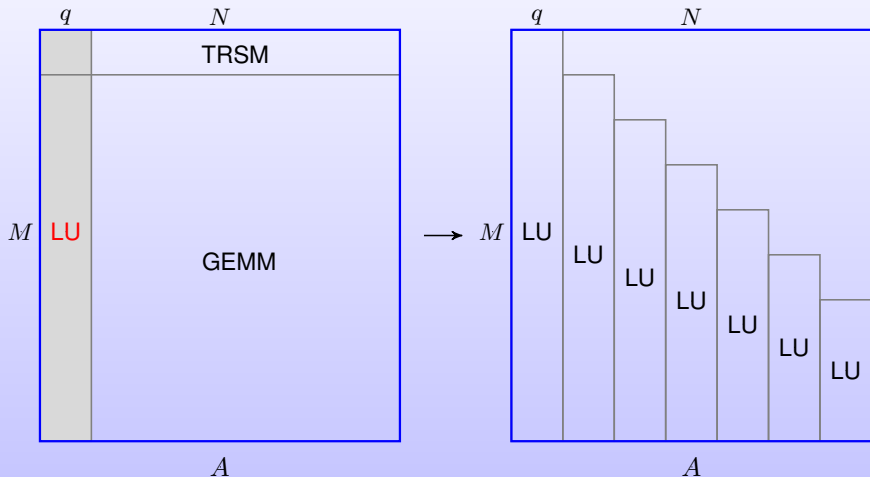
$$\boxed{Ax = b} \Rightarrow \boxed{A = LU}$$

# LU Factorization

$$A = LU$$

# LU Factorization

$$A = LU$$

# An unblocked LU Factorization

## LU Factorization

$$\left(\frac{a_{01}}{\frac{\alpha_{11}}{a_{21}}}\right) := P(p_0)\left(\frac{a_{01}}{\frac{\alpha_{11}}{a_{21}}}\right) \quad (\textbf{swap})$$

$$a_{01} := L_{00}^{-1}a_{01} \qquad\qquad (\textbf{trsv})$$

$$\alpha_{11} := \alpha_{11} - a_{10}^T a_{01} \qquad (\textbf{dot})$$

$$a_{21} := a_{21} - A_{20}a_{01} \qquad (\textbf{gemv})$$

$$\pi_1 := PivIndex\left(\frac{\alpha_{11}}{a_{21}}\right) \qquad (\textbf{max})$$

$$\left(\frac{\alpha_{11}}{a_{21}}\right) := P(\pi_1)\left(\frac{\alpha_{11}}{a_{21}}\right) \quad (\textbf{swap})$$

$$a_{21} := a_{21}/\alpha_{11} \qquad\qquad\qquad (\textbf{scal})$$

|       | $i$      | $1$           | $p$        |
|-------|----------|---------------|------------|
| $i$   | $A_{00}$ | $a_{01}$      | $A_{02}$   |
| $1$   | $a_{10}^T$ | $\alpha_{11}$ | $a_{12}^T$ |
| $p$   | $A_{20}$ | $a_{21}$      | $A_{22}$   |

$3 \times 3$ partitioning of $A$

# Outline

# Parallel Reduction
Performance Scaling on Intel Xeon Phi

$n = 67e06$
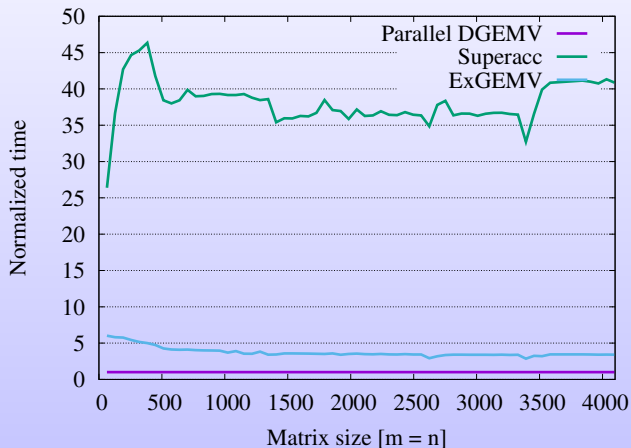
$$\text{DDOT: } \alpha := x^T y = \sum_i^N x_i y_i$$



- Based on `TwoProd` and ExSUM

- `TwoProd`$(a, b)$
  - 1: $r \leftarrow a * b$
  - 2: $s \leftarrow fma(a, b, -r)$

- $fma(a, b, c) = a * b + c$
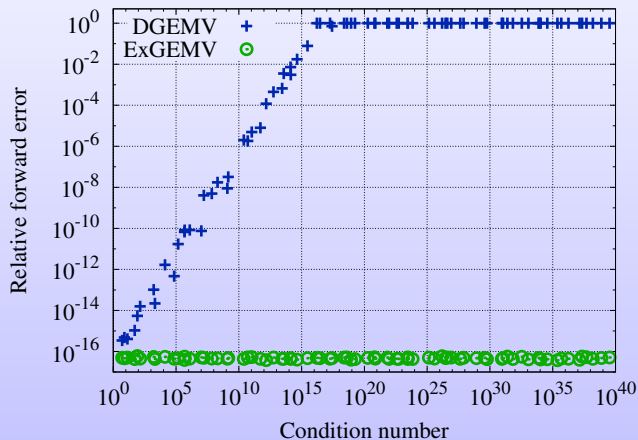
$$\boxed{\text{GEMV: } y := \alpha A x + \beta y}$$



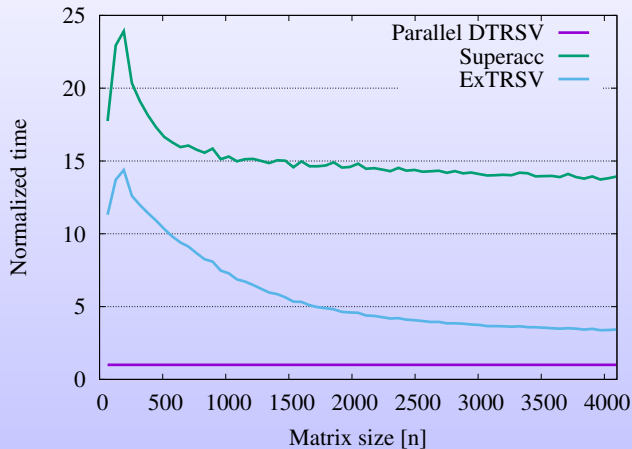- Blocked ExGEMV
- Based on ExDOT

$$\boxed{\text{GEMV: } y := Ax}$$



- Preserve every bit of information
- Correctly-rounded
- $\text{cond}(A, x) = \frac{\||A| \cdot |x|\|}{\|A \cdot x\|}$

# Triangular Solver

$$\boxed{\text{DTRSV: } Ax = b}$$



- Blocked ExTRSV
- Internal ExGEMV
- Based on ExDOT

# Triangular Solver

Accuracy

$$\boxed{\text{TRSV: } Ax = b}$$



1: $x_1 \leftarrow fl(b_1/l_{11})$
2: **for** $i = 2 \rightarrow n$ **do**
3: $\quad s \leftarrow b_i$
4: $\quad$ **for** $j = 1 \rightarrow i - 1$ **do**
5: $\quad\quad s \leftarrow s - l_{ij}x_j$
6: $\quad$ **end for**
7: $\quad x_i \leftarrow fl(RNDN(s)/l_{ii})$
8: **end for**

$$\text{cond}(A, x) = \frac{\||A^{-1}||A||x|\|_\infty}{\|x\|_\infty}$$

# LU Factorization
Performance Scaling on NVIDIA Tesla K420
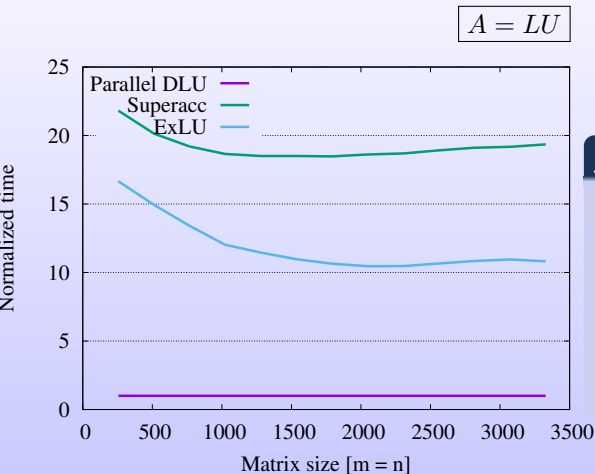
$$A = LU$$



### $jik$ variant of LU

$$\texttt{swap()}$$
$$a_{01} \leftarrow L_{00}^{-1} a_{01} \qquad \texttt{trsv}$$
$$\alpha_{11} \leftarrow \alpha_{11} - a_{10}^T a_{01} \quad \texttt{dot}$$
$$a_{21} \leftarrow a_{21} - A_{20} a_{01} \quad \texttt{gemv}$$
$$\texttt{max()}$$
$$\texttt{swap()}$$
$$a_{21} \leftarrow a_{21} / \alpha_{11} \qquad \texttt{scal}$$

$$A = LU$$

- Slightly better accuracy than DLU

- But, always reproducible

# Outline

# Conclusions and Future Work

## Conclusions

- Compute the results with no errors due to rounding
- Provide bit-wise reproducible results independently from
  - Data permutation, data assignment
  - Warps/threads scheduling
  - Partitioning/blocking
  - Reduction trees

# Conclusions and Future Work

## Conclusions

- Compute the results with no errors due to rounding

- Provide bit-wise reproducible results independently from
  - Data permutation, data assignment
  - Warps/threads scheduling
  - Partitioning/blocking
  - Reduction trees

- Deliver comparable performance to the classic implementations of the memory-bound operations

# Conclusions and Future Work

## Conclusions

- Compute the results with no errors due to rounding
- Provide bit-wise reproducible results independently from
  - Data permutation, data assignment
  - Warps/threads scheduling
  - Partitioning/blocking
  - Reduction trees
- Deliver comparable performance to the classic implementations of the memory-bound operations
- Reproducible underlying kernels → reproducible LU

# Conclusions and Future Work

## Conclusions

- Compute the results with no errors due to rounding
- Provide bit-wise reproducible results independently from
  - Data permutation, data assignment
  - Warps/threads scheduling
  - Partitioning/blocking
  - Reduction trees
- Deliver comparable performance to the classic implementations of the memory-bound operations
- Reproducible underlying kernels $\rightarrow$ reproducible LU

## Future directions

- Lightweight approach for compute-intensive operations
- Performance portability
- Applicability in real-world codes

# Thank you for your attention!

```
https://exblas.lip6.fr
```

## ExBLAS

- ExBLAS-1: ExSUM[a], ExSCAL, ExDOT, ExAXPY, ...

- ExBLAS-2: ExGER, ExGEMV, ExTRSV, ExSYR, ...

- ExBLAS-3: ExGEMM, ExTRSM, ExSYR2K, ...

[a]Routines in blue are already in ExBLAS

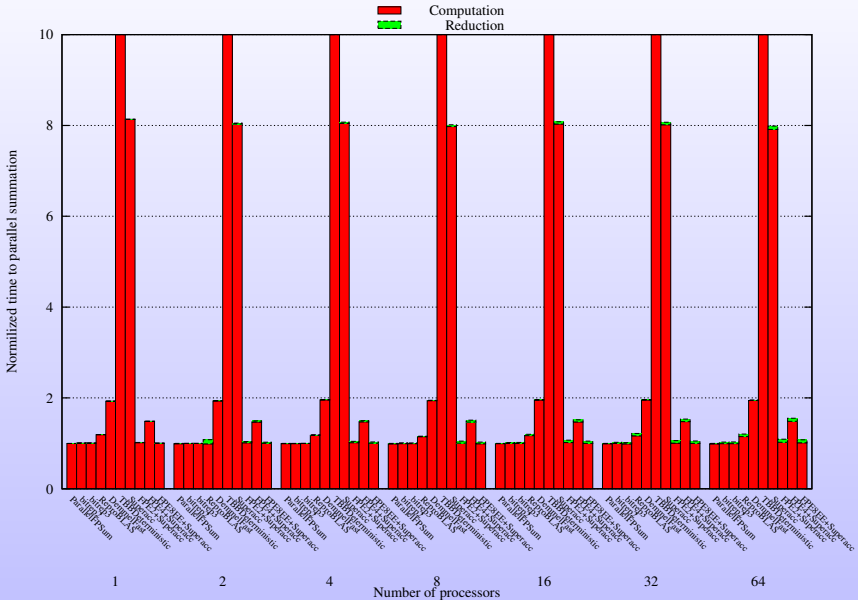## Higher Level Operations

- Unblocked LU factorization variants, including GER+SCAL

- Towards blocked LU factorization

# Efforts on Reproducibility

## Standardalization

- IEEE 754-2018: `TwoSum` and `TwoProd`
- Jim Demmel et. al.: "A Proposal for a Next-Generation BLAS"

## Journals

- TOMS: Replicated Computational Results

## Conferences, Workshops, and Minisymposiums

- ARITH 2016-17: Arithmetic challenges in HPC and exascale computing (accuracy, reproducibility, ...)
- SC 2015-16: Workshop on Numerical Reproducibility at Exascale
- IPDPS 2017 and Euro-par 2014-16: Workshop on Reproducibility in Parallel Computing
- SIAM PP 2016: Numerical Reproducibility for High-Performance Computing
- SIAM CSE 2017: Algorithmic evolution in Post Moore's Era: Auto-Tuning and Accuracy Assurance