

Towards Fast, Accurate and Reproducible LU Factorization

Roman Iakymchuk¹, David Defour², and Stef Graillat³

¹KTH Royal Institute of Technology, CSC, CST/PDC

²Université de Perpignan, DALI-LIRMM

³Sorbonne Universités, UPMC Univ Paris VI, UMR 7606, LIP6

riakymch@kth.se

SCAN 2016, Sept 26-29, 2016
Uppsala, Sverige



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

Linear Algebra Libraries

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

LAPACK

FLAME

NAG

Basic Linear Algebra Subprograms (BLAS)

Refer. BLAS

MKL, cuBLAS

OpenBLAS

ATLAS



- To compute BLAS operations with floating-point numbers **fast** and **precise**, ensuring their **numerical reproducibility**, on a wide range of architectures

ExBLAS – Exact BLAS

- **ExBLAS-1**: ExSUM, ExSCAL, ExDOT, ExAXPY, ...
- **ExBLAS-2**: ExGER, ExGEMV, ExTRSV, ExSYR, ...
- **ExBLAS-3**: ExGEMM, ExTRSM, ExSYR2K, ...

- To compute BLAS operations with floating-point numbers **fast** and **precise**, ensuring their **numerical reproducibility**, on a wide range of architectures

ExBLAS – Exact BLAS

- **ExBLAS-1:** ExSUM, ExSCAL, ExDOT, ExAXPY, ...
- **ExBLAS-2:** ExGER, ExGEMV, ExTRSV, ExSYR, ...
- **ExBLAS-3:** ExGEMM, ExTRSM, ExSYR2K, ...

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

- 1 Accuracy and Reproducibility of FP Operations
- 2 Exact Parallel Reduction
- 3 ExBLAS and Reproducible LU
- 4 Performance Results
- 5 Conclusions and Future Work

- 1 Accuracy and Reproducibility of FP Operations
- 2 Exact Parallel Reduction
- 3 ExBLAS and Reproducible LU
- 4 Performance Results
- 5 Conclusions and Future Work

Problems

- Floating-point arithmetic suffers from **rounding errors**
- Floating-point operations (+, ×) are commutative but **non-associative**

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

Problems

- Floating-point arithmetic suffers from **rounding errors**
- Floating-point operations (+, ×) are commutative but **non-associative**

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

Problems

- Floating-point arithmetic suffers from **rounding errors**
- Floating-point operations (+, ×) 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

- **Fix the Order of Computations**

- Sequential mode: intolerably costly at large-scale systems
 - Fixed reduction trees: substantial communication overhead
- Example: Intel **C**onditional **N**umerical **R**eproducibility
($\sim 2x$ for datum, no accuracy guarantees)

● Fix the Order of Computations

- Sequential mode: intolerably costly at large-scale systems
- Fixed reduction trees: substantial communication overhead
- Example: Intel **C**onditional **N**umerical **R**eproducibility
($\sim 2x$ for datum, no accuracy guarantees)

● Eliminate/Reduce the Rounding Errors

- Fixed-point arithmetic: limited range of values
- Fixed FP expansions with Error-Free Transformations (EFT)
- Example: double-double or quad-double (Briggs, Bailey, Hida, Li)
(work well on a set of relatively close numbers)
- “Infinite” precision: reproducible independently from the inputs
- Example: Kulisch accumulator (considered inefficient)

● Fix the Order of Computations

- Sequential mode: intolerably costly at large-scale systems
- Fixed reduction trees: substantial communication overhead
- Example: Intel **Conditional Numerical Reproducibility**
($\sim 2x$ for datum, no accuracy guarantees)

● Eliminate/Reduce the Rounding Errors

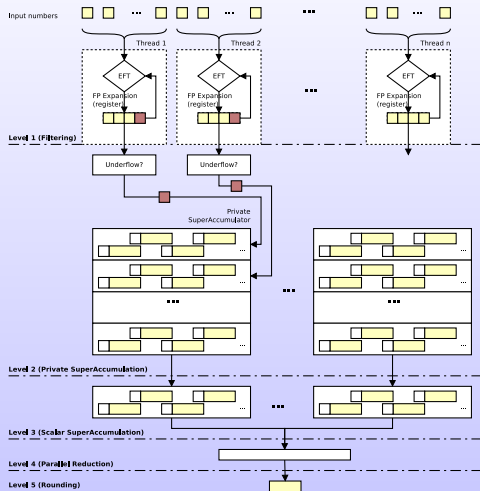
- Fixed-point arithmetic: limited range of values
- Fixed FP expansions with **Error-Free Transformations** (EFT)
- Example: double-double or quad-double (Briggs, Bailey, Hida, Li)
(work well on a set of relatively close numbers)
- “Infinite” precision: reproducible independently from the inputs
- Example: **Kulisch accumulator** (considered inefficient)

● Libraries

- **ReproBLAS**: Reproducible BLAS (Demmel, Nguyen, Ahrens)
- For BLAS-1, GEMV, and GEMM on CPUs
- **RARE-BLAS**: Repr. Accur. Rounded and Eff. BLAS (Chohra, Langlois, Parello) → For BLAS-1 and GEMV on CPUs

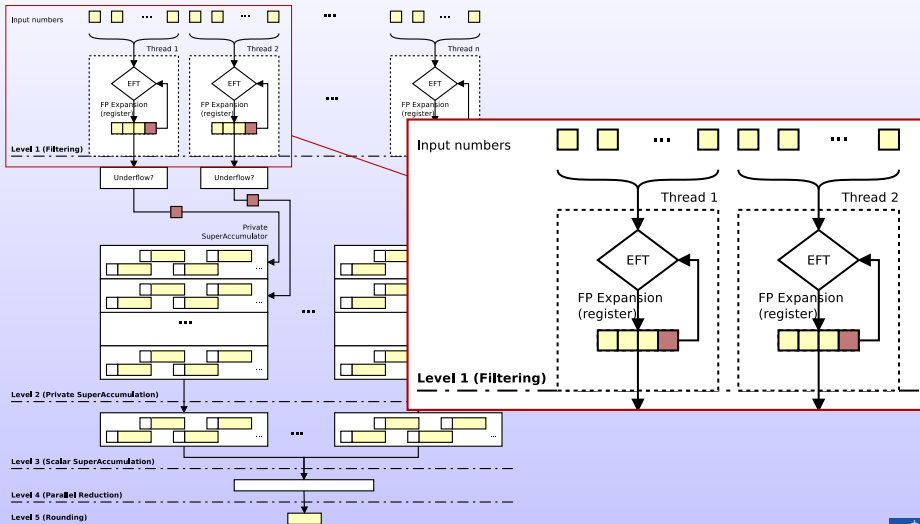
- 1 Accuracy and Reproducibility of FP Operations
- 2 Exact Parallel Reduction**
- 3 ExBLAS and Reproducible LU
- 4 Performance Results
- 5 Conclusions and Future Work

Our Multi-Level Reproducible Summation

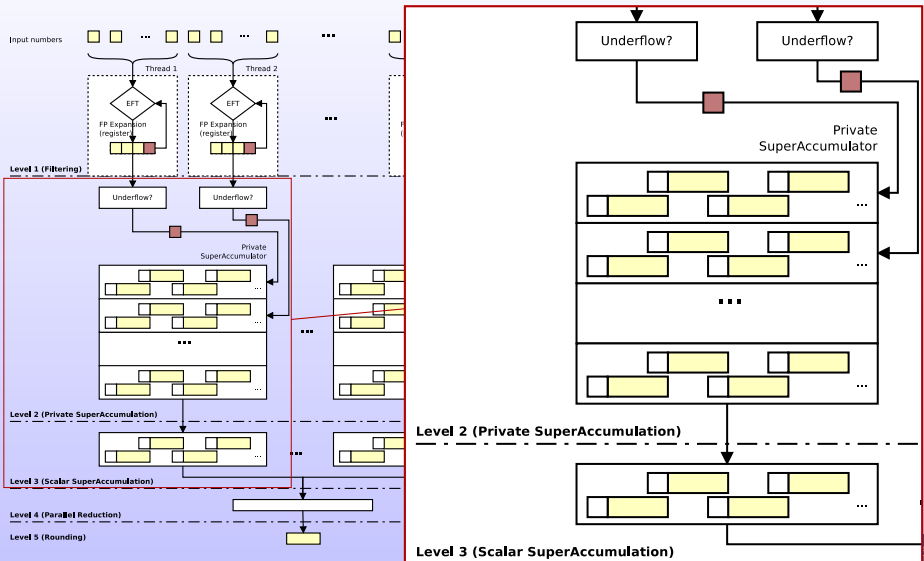


- Parallel algorithm with 5-levels
 - Suitable for today's parallel architectures
 - Based on FPE with EFT and Kulisch accumulator
 - Guarantees “inf” precision
- **bit-wise reproducibility**

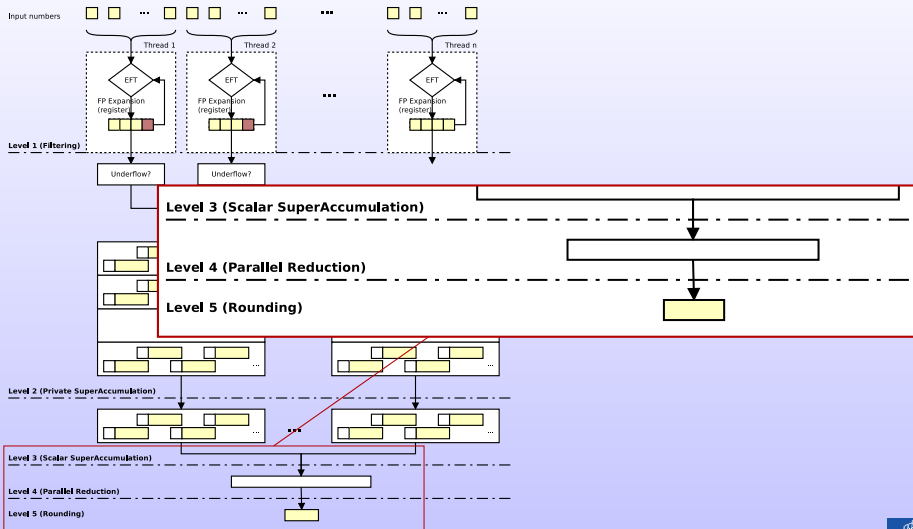
Level 1: Filtering



Level 2 and 3: Scalar Superaccumulator



Level 4 and 5: Reduction and Rounding



- 1 Accuracy and Reproducibility of FP Operations
- 2 Exact Parallel Reduction
- 3 ExBLAS and Reproducible LU**
- 4 Performance Results
- 5 Conclusions and Future Work

ExBLAS Status

- ExBLAS-1: [ExSUM](#), [ExSCAL](#), [ExDOT](#), ExAXPY, ...
- ExBLAS-2: [ExGER](#), [ExGEMV](#), [ExTRSV](#), ExSYR, ...
- ExBLAS-3: [ExGEMM](#), ExTRSM, ExSYR2K, ...

ExBLAS Status

- ExBLAS-1: ExSUM, ExSCAL, ExDOT, ExAXPY, ...
- ExBLAS-2: ExGER, ExGEMV, ExTRSV, ExSYR, ...
- ExBLAS-3: ExGEMM, ExTRSM, ExSYR2K, ...

ExSCAL

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



ExBLAS Status

- ExBLAS-1: ExSUM, ExSCAL, ExDOT, ExAXPY, ...
- ExBLAS-2: ExGER, ExGEMV, ExTRSV, ExSYR, ...
- ExBLAS-3: ExGEMM, ExTRSM, ExSYR2K, ...

ExSCAL

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

ExBLAS Status

- ExBLAS-1: ExSUM, ExSCAL, ExDOT, ExAXPY, ...
- ExBLAS-2: ExGER, ExGEMV, ExTRSV, ExSYR, ...
- ExBLAS-3: ExGEMM, ExTRSM, ExSYR2K, ...

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

ExBLAS Status

- ExBLAS-1: ExSUM, ExSCAL, ExDOT, ExAXPY, ...
- ExBLAS-2: ExGER, ExGEMV, ExTRSV, ExSYR, ...
- ExBLAS-3: ExGEMM, ExTRSM, ExSYR2K, ...

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 Status

- ExBLAS-1: ExSUM, ExSCAL, ExDOT, ExAXPY, ...
- ExBLAS-2: ExGER, ExGEMV, ExTRSV, ExSYR, ...
- ExBLAS-3: ExGEMM, ExTRSM, ExSYR2K, ...

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

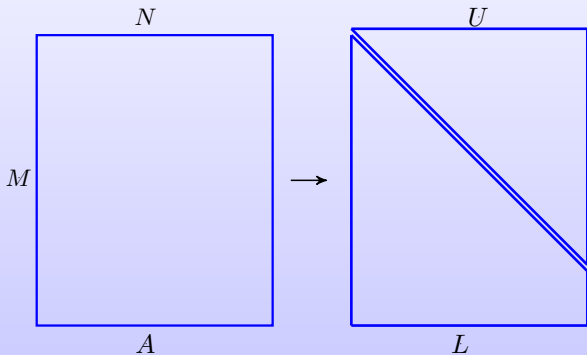
ExGER

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



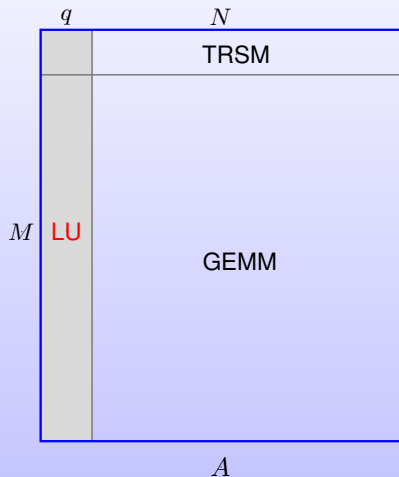
LU Factorization

$$\boxed{Ax = b} \Rightarrow \boxed{A = LU}$$



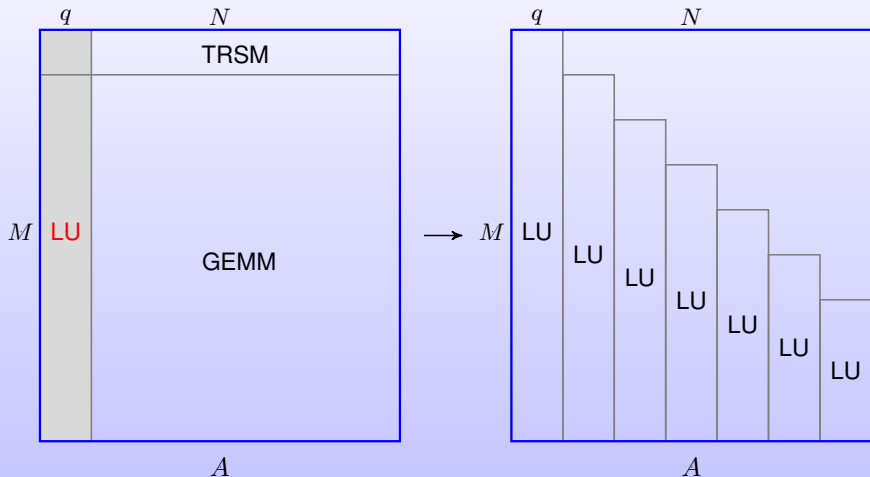
LU Factorization

$$A = LU$$



LU Factorization

$$A = LU$$



An unblocked LU Factorization

LU Factorization

$$a_{10}^T := a_{10}^T U_{00}^{-1} \quad \text{TRSV}$$

$$\alpha_{11} := \alpha_{11} - a_{10}^T a_{01} \quad \text{DOT}$$

$$a_{12}^T := a_{12}^T - a_{10}^T A_{02} \quad \text{GEMV}$$

	i	1	p
i	A_{00}	a_{01}	A_{02}
1	a_{10}^T	α_{11}	a_{12}^T
p	A_{20}	a_{21}	A_{22}

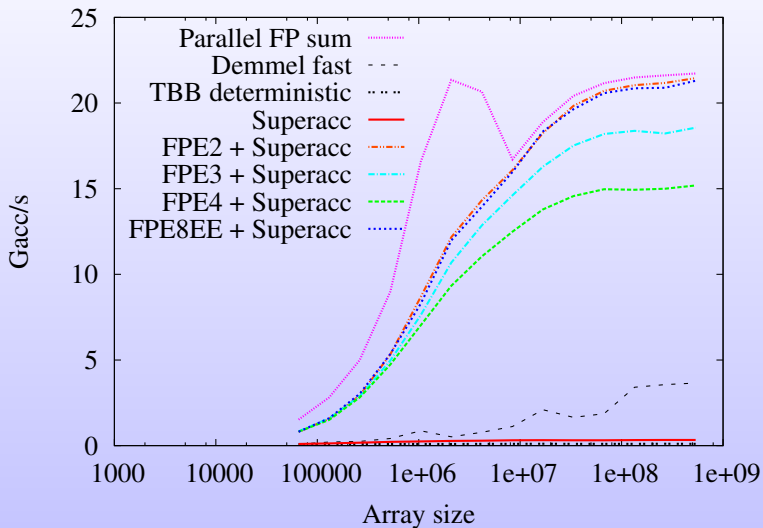
3×3 partitioning of A



- 1 Accuracy and Reproducibility of FP Operations
- 2 Exact Parallel Reduction
- 3 ExBLAS and Reproducible LU
- 4 Performance Results**
- 5 Conclusions and Future Work

Parallel Reduction

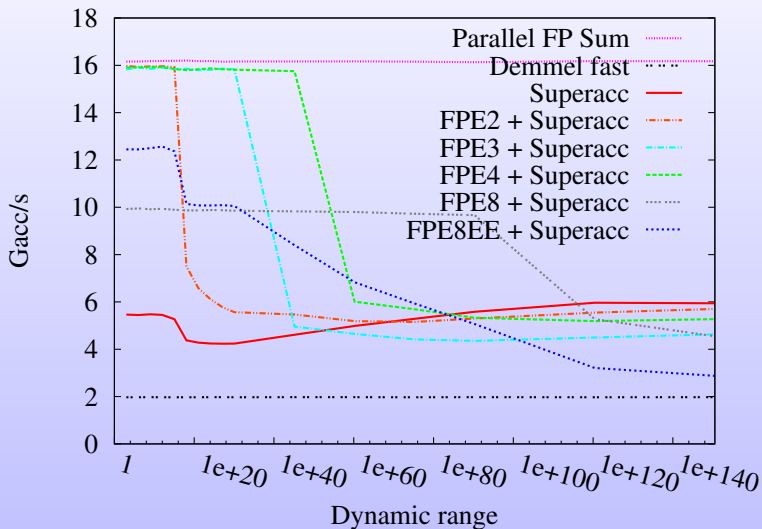
Performance Scaling on Intel Xeon Phi



Parallel Reduction

Data-Dependent Performance on NVIDIA Tesla K20c

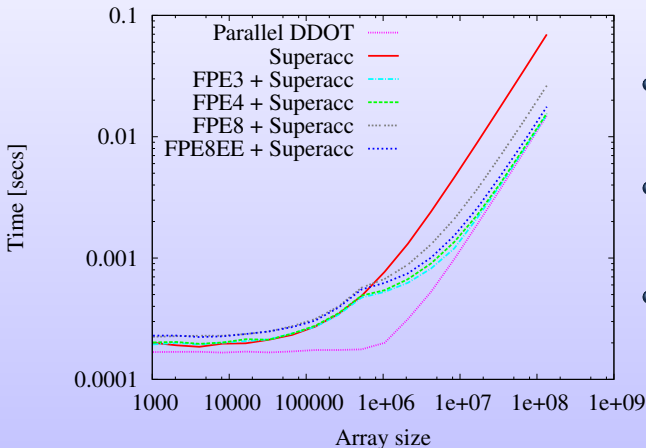
$n = 67e06$



Dot Product

Performance Scaling on NVIDIA Tesla K20c

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



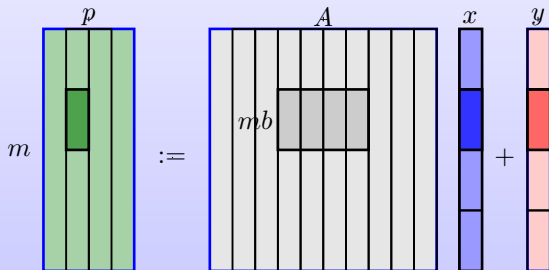
- Based on **TwoProduct** and Reproducible Summation
- **TwoProduct**(a, b)
 - 1: $r \leftarrow a * b$
 - 2: $s \leftarrow fma(a, b, -r)$
- $fma(a, b, c) = a * b + c$



Matrix-Vector Product

Performance Scaling on NVIDIA Tesla K80

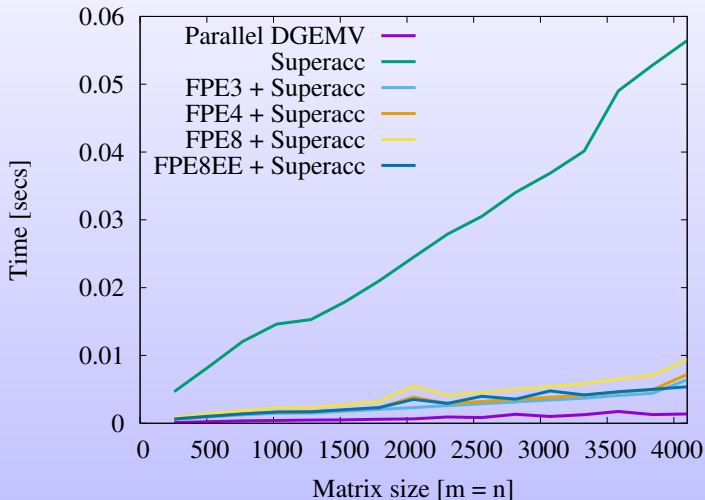
$$\text{DGEMV: } y := \alpha Ax + \beta y$$



Matrix-Vector Product

Performance Scaling on NVIDIA Tesla K80

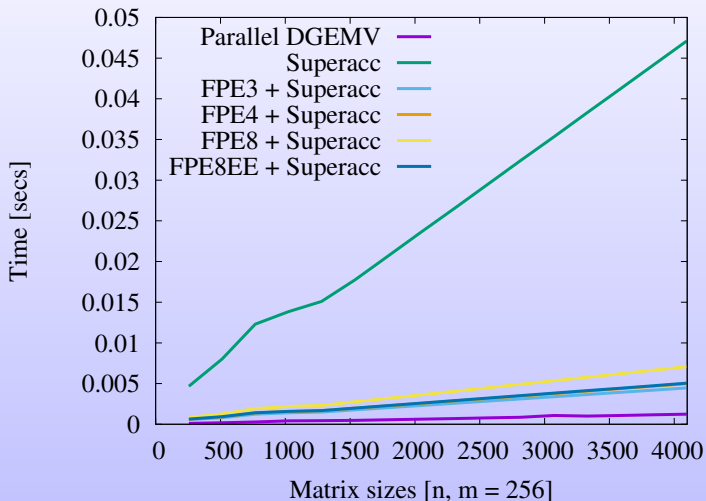
$$\text{DGEMV: } y := \alpha Ax + \beta y$$



Matrix-Vector Product

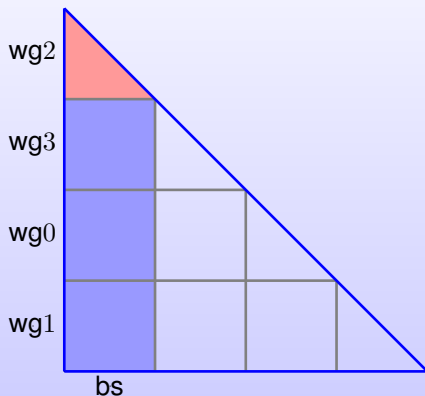
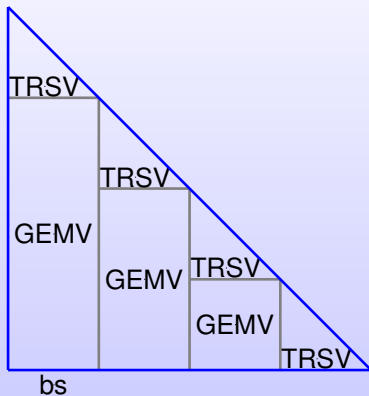
Performance Scaling on NVIDIA Tesla K80

$$\text{DGEMV: } y := \alpha Ax + \beta y$$



Triangular Solver

Matrix Partitioning

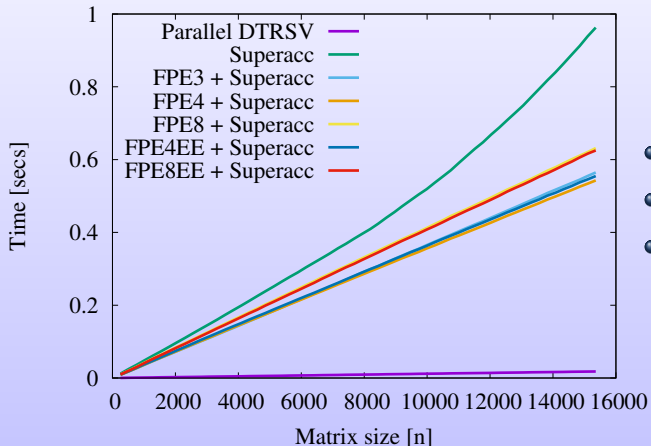


Partitioning of a lower triangular matrix L

Triangular Solver

Performance Scaling on NVIDIA Tesla K20c

DTRS_V: $Ax = b$

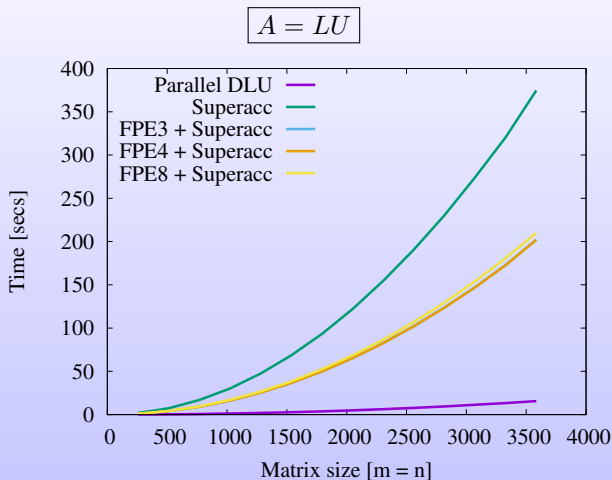


- Blocked ExTRS_V
- Based on ExDOT
- Internal ExGEMV



LU Factorization

Performance Scaling on NVIDIA Tesla K80



$a_{10}^T \leftarrow a_{10}^T U_{00}^{-1}$ trsv
 $\alpha_{11} \leftarrow \alpha_{11} - a_{10}^T a_{01}$ dot
 $a_{12}^T \leftarrow a_{12}^T - a_{10}^T A_{02}$ gemv



- 1 Accuracy and Reproducibility of FP Operations
- 2 Exact Parallel Reduction
- 3 ExBLAS and Reproducible LU
- 4 Performance Results
- 5 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
 - Thread scheduling
 - Partitioning/blocking
 - Reduction trees

Conclusions

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

Conclusions

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

Future directions

- Enhance compute-intensive operations and the LU factorization
- Cover the other variants of the unblocked LU factorization
- Application of our implementations in real-world codes



Thank you for your attention!

URL: <https://exblas.lip6.fr>

ExBLAS Status

- ExBLAS-1: [ExSUM^a](#), [ExSCAL](#), [ExDOT](#), ExAXPY, ...
- ExBLAS-2: [ExGER](#), [ExGEMV](#), [ExTRSV](#), ExSYR, ...
- ExBLAS-3: [ExGEMM](#), ExTRSM, ExSYR2K, ...

^aRoutines in [blue](#) are already in ExBLAS

