SLATE Project: Objectives, Design, and Results

# SLATE: Software for Linear Algebra Targeting Exascale

# ScaLAPACK Legacy

1991 — Linux

1992 —

1993 — NCSA Mosaic

1994 — MPI 1.0

1995 — **ScaLAPACK**

1996 —

1997 — OpenMP 1.0 for FORTRAN

1998 — OpenMP 1.0 for C/C++

**February 28, 1995:**
- initial release of ScaLAPACK
- Denver International Airport opens

Fujitsu Numerical Wind Tunnel is the fastest machine on the TOP500 list, with **170 Gflop/s** of LINPACK performance.

My 2016 MacBook Pro gets **166 Gflop/s**

# SLATE Objectives

**<u>can be built:</u>**
- serial
- OpenMP multithreading
- MPI message passing
- GPU acceleration

- Coverage            ScaLAPACK and beyond
- Modern Hardware     DOE CORAL (pre Exascale) → DOE Exascale
- Portability           Intel Xeon (&Phi), IBM POWER, ARM, NVIDIA, AMD, …
- Modern Language      C++11/14/17 (templates, STL, overloading, polymorphism, …)
- Modern Standards      MPI 3, OpenMP 4/5 (&omp target)
- Performance         80-90% of peak (asymptotic)
- Scalability           full machine (tens of thousands of nodes)
- Productivity          ca. 4 full time developers
- Maintainability       part time developers + community

# SLATE Software Stack



| molecular dynamics | computational chemistry | quantum mechanics | quantum chemistry | sparse solvers |
|---|---|---|---|---|
| **EXAALT** | **NWChemEx** | **QMCPACK** | **GAMESS** | **FBSS** |

**SLATE** — PARALLEL DENSE LINEAR ALGEBRA ROUTINES, DISTRIBUTED MEMORY, MULTICORE, ACCELERATORS

| **PaRSEC** | **MPI** | **OpenMP** | **BLAS++** | **LAPACK++** | **batch BLAS++** |
|---|---|---|---|---|---|
| | **Exa MPI** | **SOLLVE** | **MKL** | **ESSL** | **cuBLAS** | **ACML** |
| | **OMPI-X** | | | | | |

ECP   standards   vendor   SLATE

# SLATE Resources

- main ECP website:            https://www.exascaleproject.org/
- main SLATE website:         http://icl.utk.edu/slate/
- main SLATE repository:     https://bitbucket.org/icl/slate
- BLAS++ repository:          https://bitbucket.org/icl/blaspp
- LAPACK++ repositry:       https://bitbucket.org/icl/lapackpp
- SLATE Working Notes:      http://www.icl.utk.edu/publications/series/swans
- Research Gate project:     https://www.researchgate.net/project/ECP-SLATE
- SLATE User                https://groups.google.com/a/icl.utk.edu/forum/#!forum/slate-user

# SWAN Catalog

- SWAN 1 Prospectus
- SWAN 2 BLAS++ & LAPACK++
- SWAN 3 Design
- SWAN 4 Batched BLAS++ API
- SWAN 5 Parallel BLAS Performance

- SWAN 6 Parallel Norms Performance
- SWAN 7 Batched BLAS++ Implementation
- SWAN 8 Linear Systems Performance
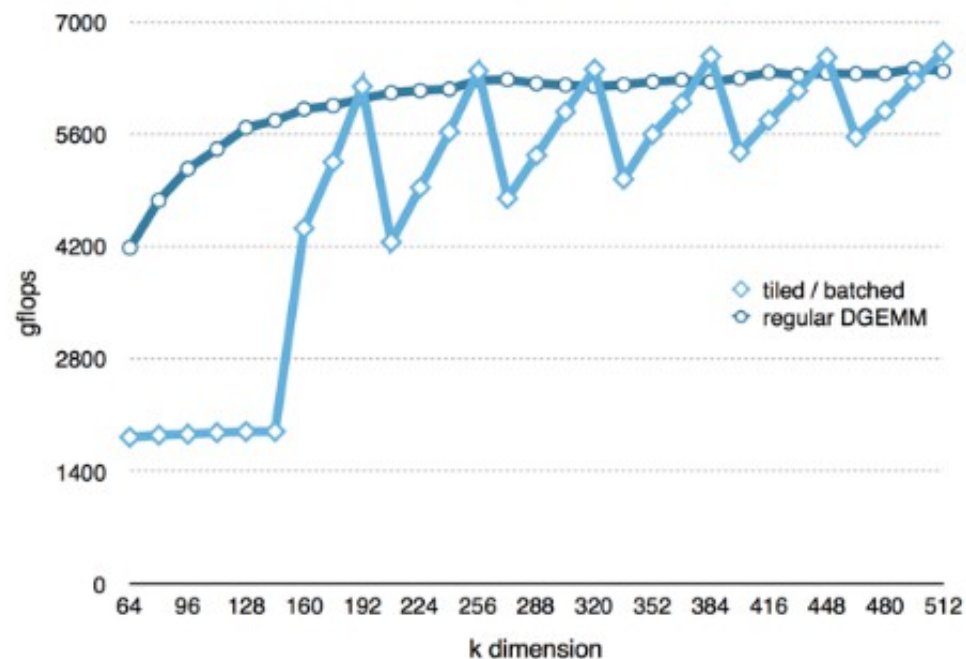- SWAN 9 Least Squares Performance

# SLATE Working Notes

- Main SWANs page
  - http://www.icl.utk.edu/publications/series/swans

- Designing SLATE: Software for Linear Algebra Targeting Exascale
  - http://www.icl.utk.edu/publications/swan-003

- C++ API for BLAS and LAPACK
  - http://www.icl.utk.edu/publications/swan-002
  - https://bitbucket.org/icl/blaspp
  - https://bitbucket.org/icl/lapackpp

- Roadmap for the Development of a Linear Algebra Library for Exascale Computing:
  - SLATE: Software for Linear Algebra Targeting Exascale
  - http://www.icl.utk.edu/publications/swan-001

# GEMM Efficiency



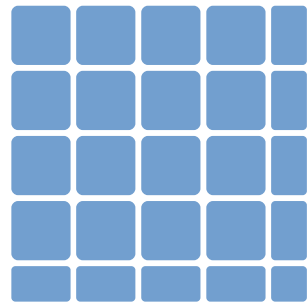Schur complement performance on NVIDIA Pascal

Schur complement performance on NVIDIA Volta

C = C – A × B with small k, i.e., the DGEMM called in LU factorization
The matrix fills out the GPU memory. The X axis shows the k dimension.

# SLATE Matrix

not allocated

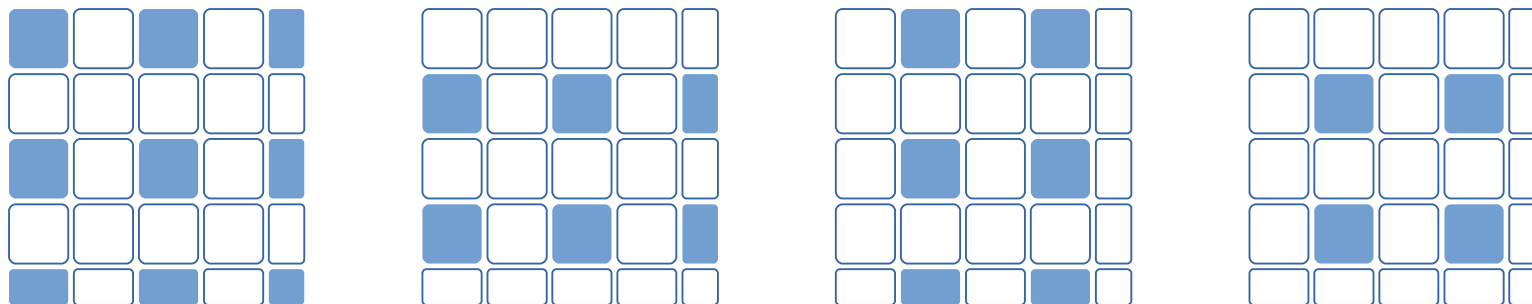std::map<std::tuple<*int64_t*, *int64_t*, *int*>, Tile<FloatType>*> *tiles_;

row

column

host & devices

- collection of tiles
- **individually allocated**
- only allocate what is needed
- accommodates: symmetric, triangular, band, …

While in the PLASMA library the matrix is also stored in tiles, the tiles are laid out contiguously in memory.

In contrast, in SLATE, the tiles are individually allocated, with no correlation of their locations in the matrix to their addresses in memory.

std::map<std::tuple<*int64_t*, *int64_t*, *int*>, Tile<FloatType>*> *tiles_;

- distributed matrix
- global indexing of tiles
- only allocate the local part
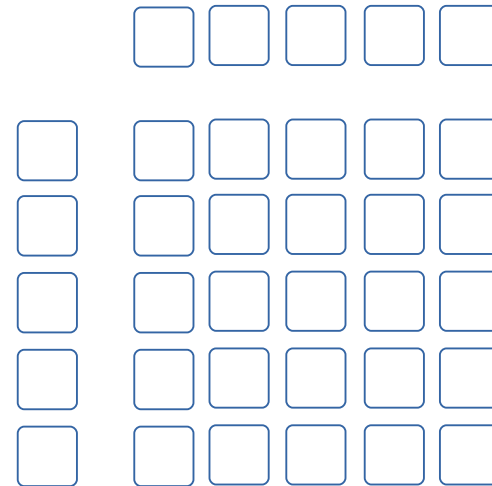- any distribution is possible (2D block cyclic by default)

The same structure, used for single node representation, naturally supports distributed memory representation.
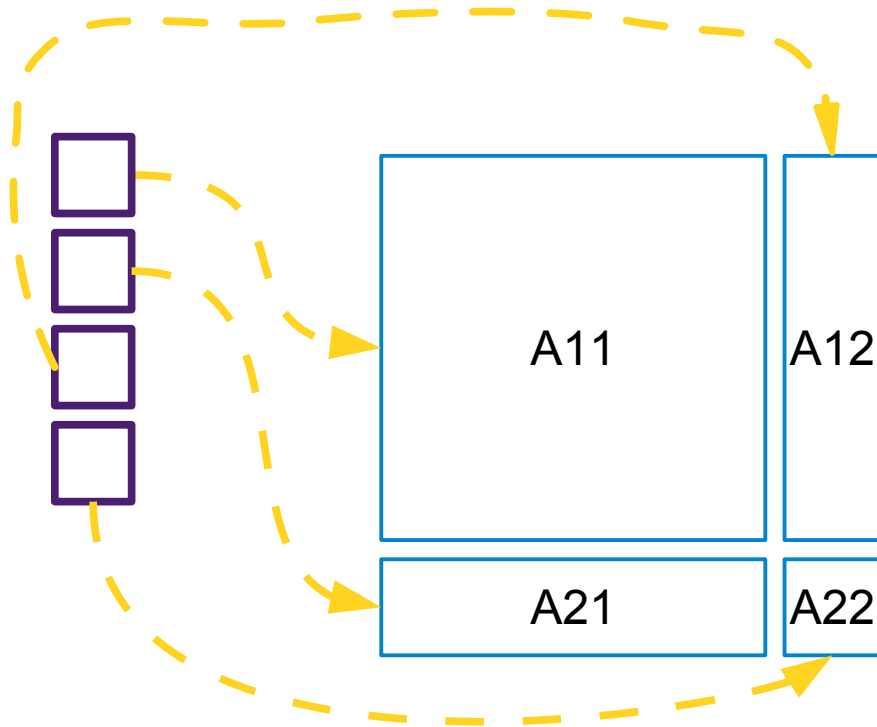
# Data Storage Comparison
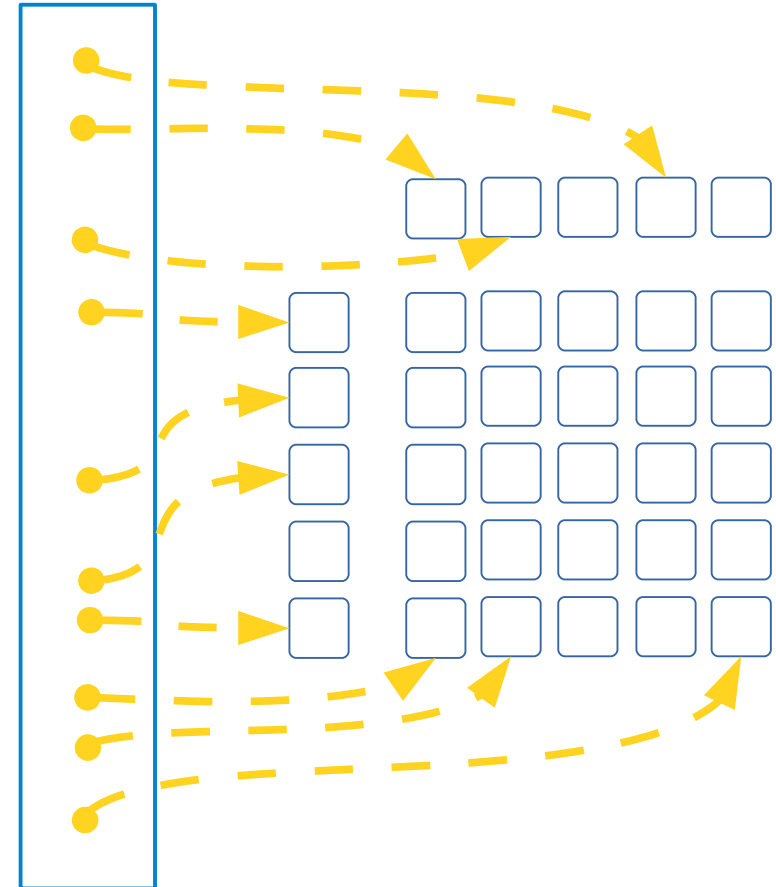
**LAPACK**
**MAGMA**

**SLATE**

$$C = C - A \times B$$

# Data Storage Comparison

SLATE Tile Map<>

PLASMA Descriptor



A11   A12

A21   A22
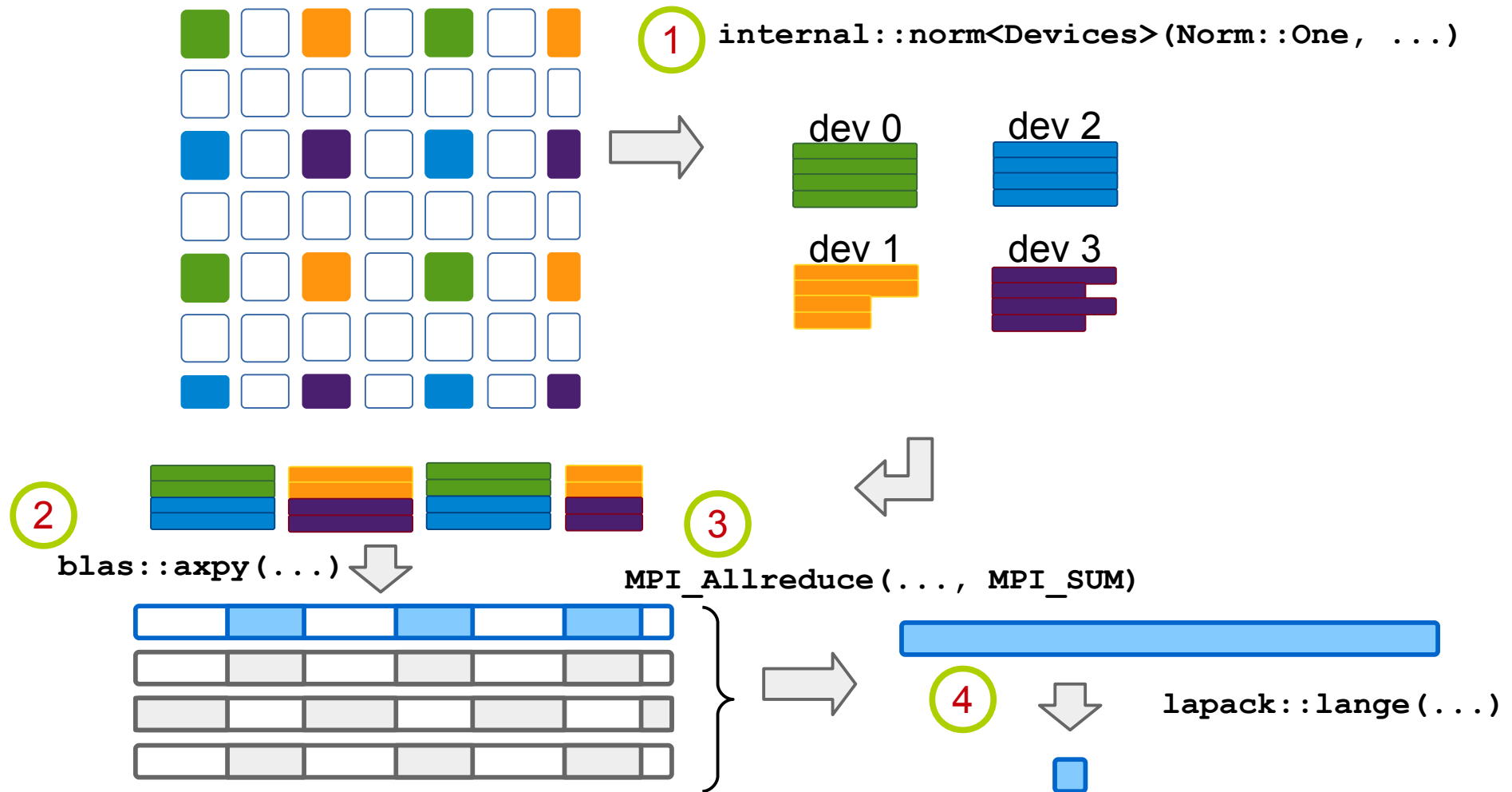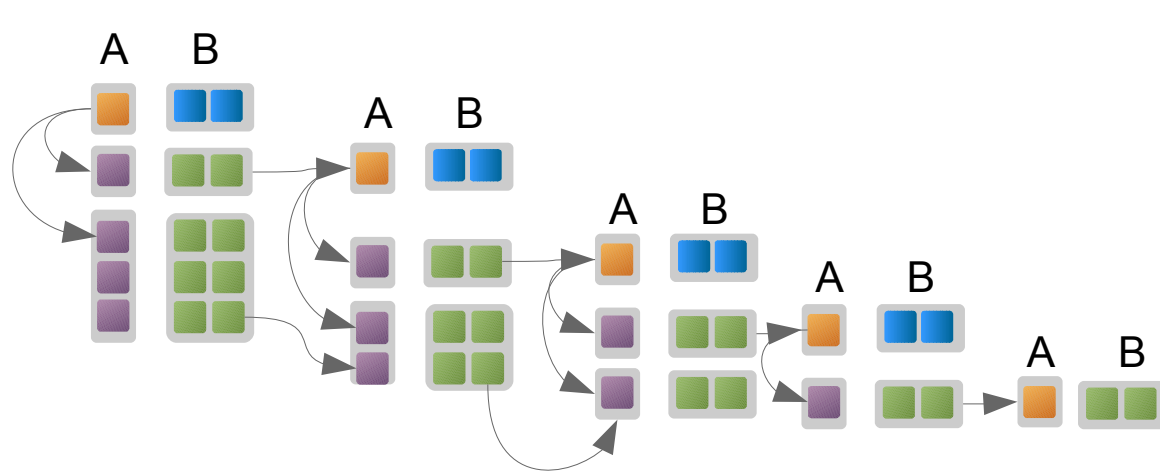
This layout allows in-place translation.

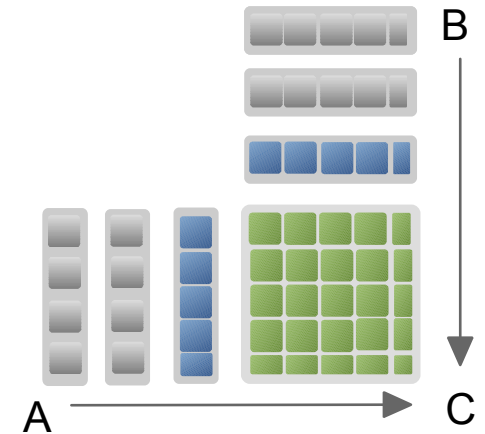# Computing Norms in Parallel in SLATE
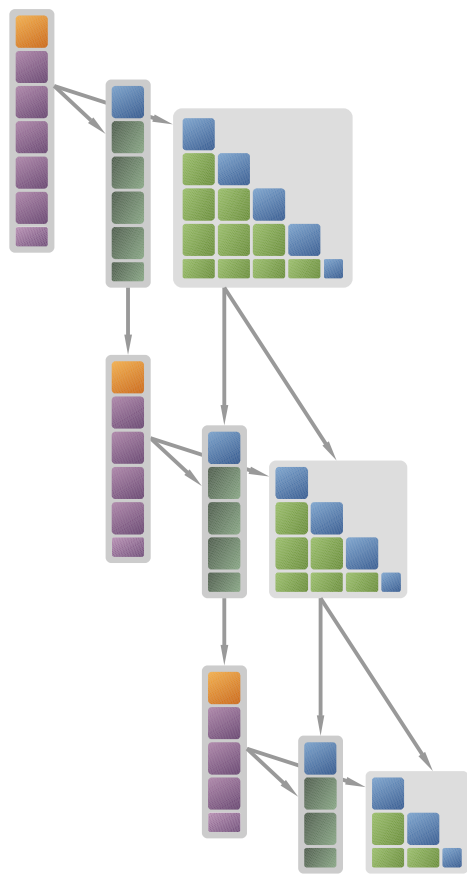
# SLATE Parallel BLAS Scheduling

TRSM

GEMM

- nested parallelism

- top level:

  - #pragma omp task depend
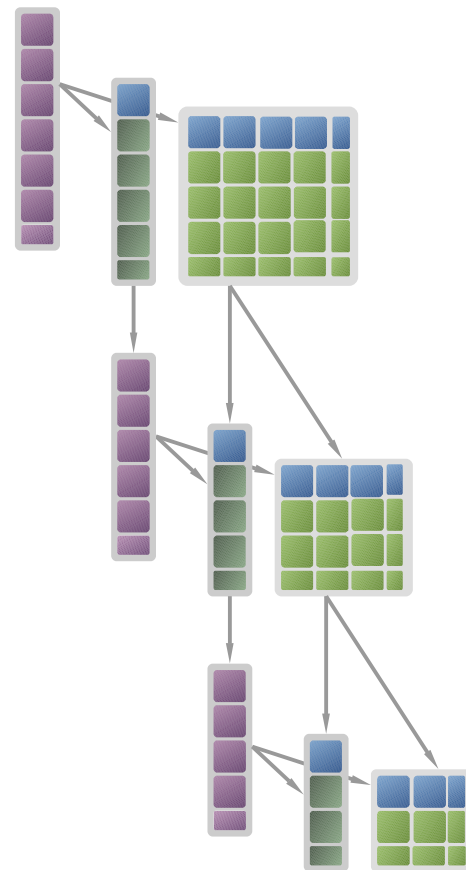
- bottom level:

  - #pragma omp task

  - batch GEMM

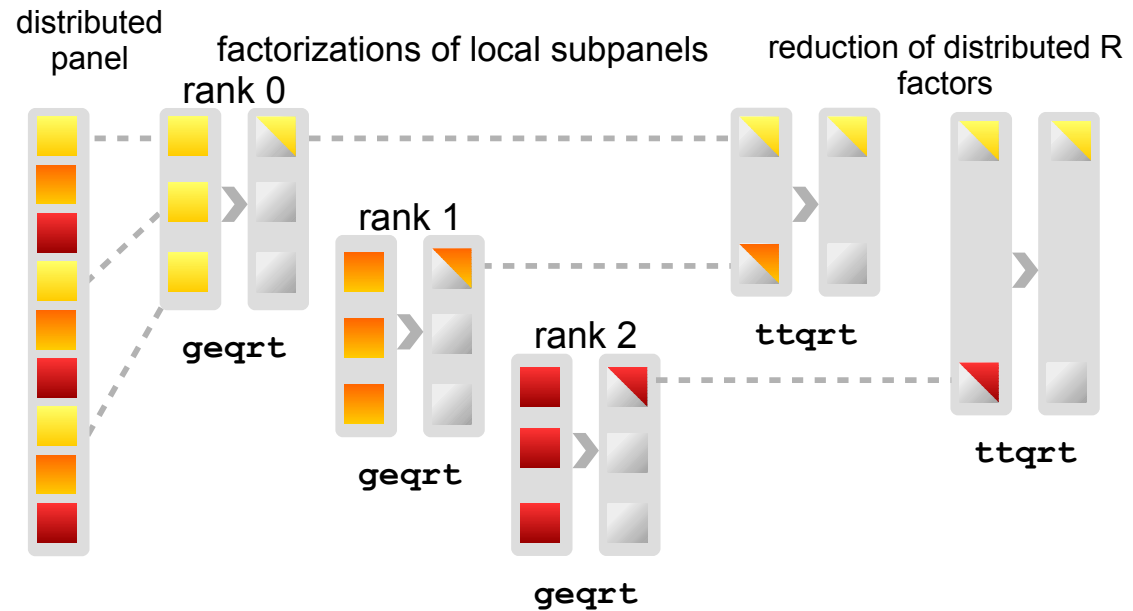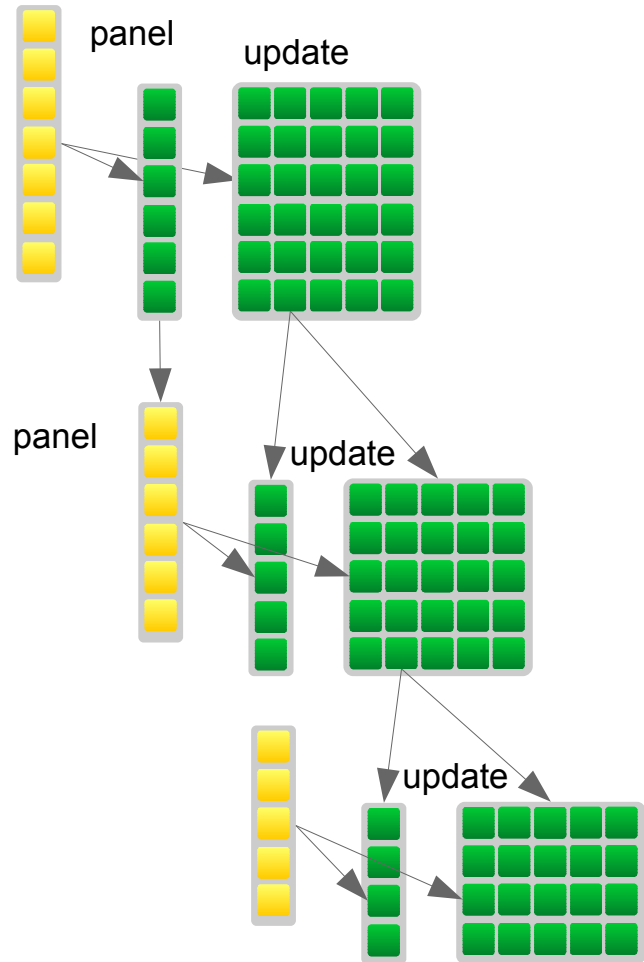# SLATE Scheduling: Linear Solvers



Cholesky factorization

- nested parallelism
- top level:
  - #pragma omp task depend
- bottom level:
  - #pragma omp task
  - batch GEMM
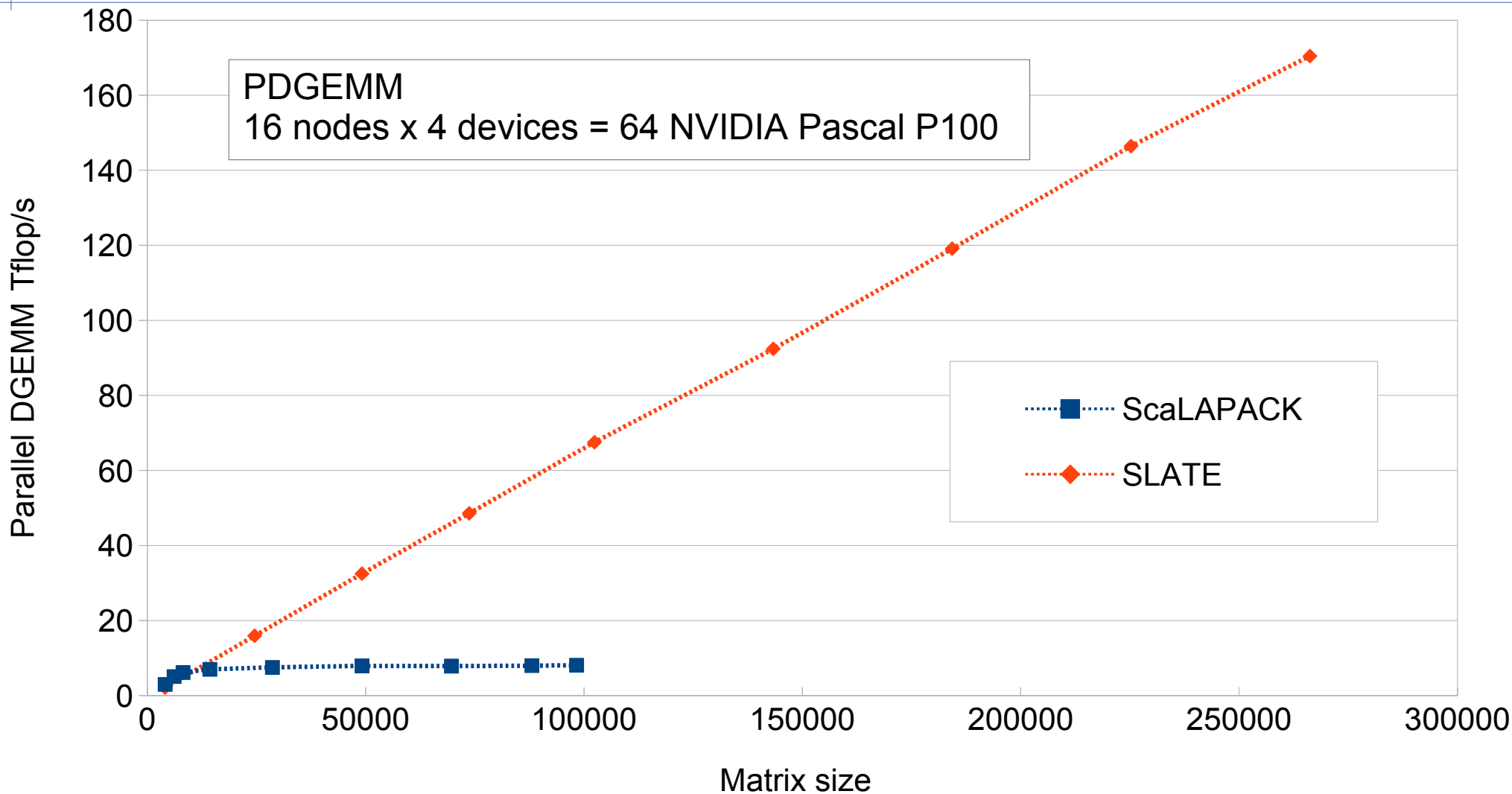
LU factorization

# SLATE QR Scheduling

## Lookahead



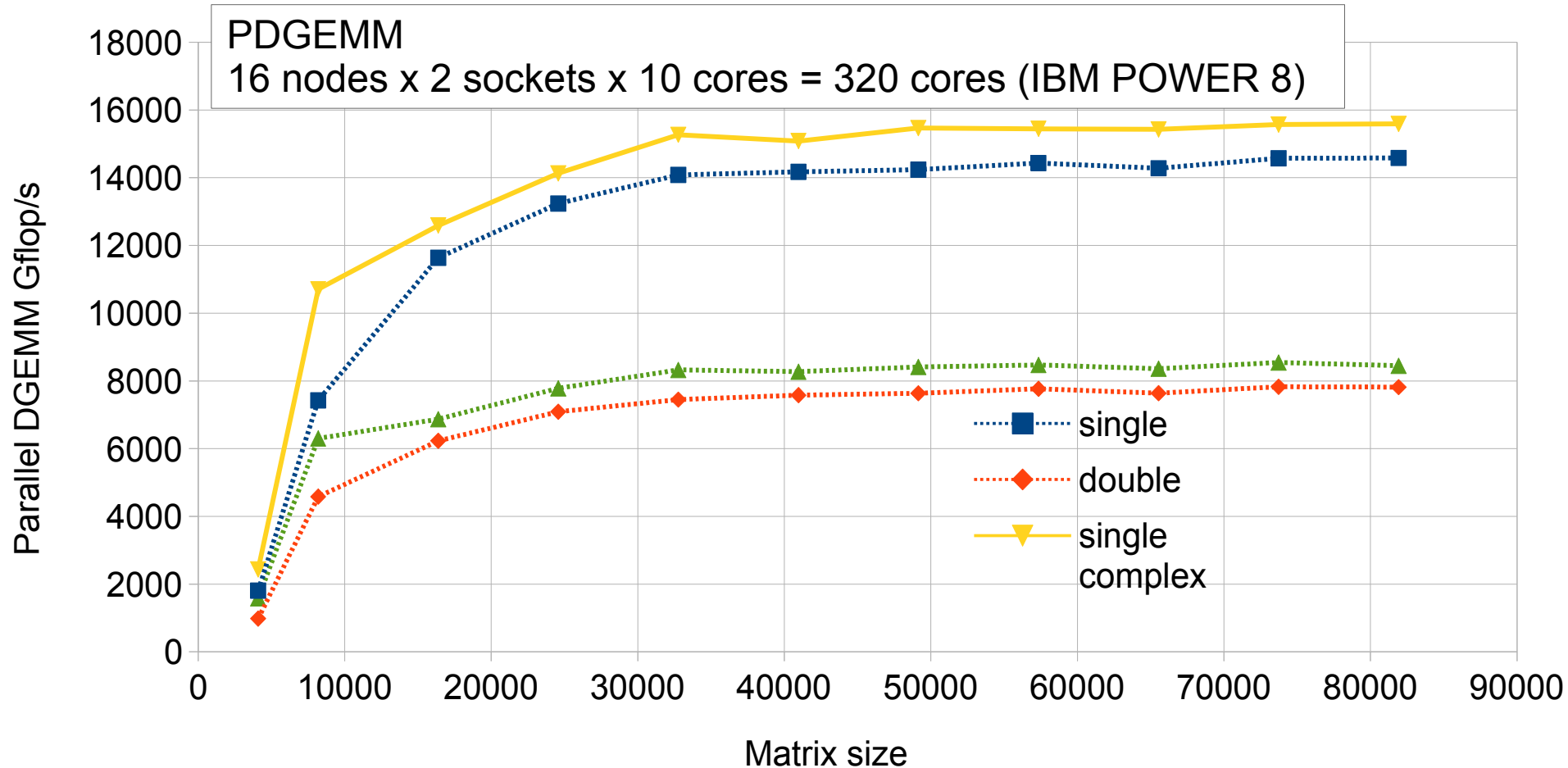distributed panel — factorizations of local subpanels — reduction of distributed R factors

rank 0

rank 1

rank 2

geqrt

geqrt

geqrt

ttqrt

ttqrt

Preliminary SLATE Performance Results

# SummitDev @ OLCF

- 3×18 = **54** nodes (IBM S822LC)
- 2×10 = **20** cores (IBM POWER8)           ca. 0.5 TFLOPS (2.5%)
- **4** GPUs (NVIDIA P100)                ca. 20 TFLOPS (97.5%)
- **256** GB DDR4
- 4×16 = **64** GB **HBM2**
- NVLink 1.0                       80 Gbps (advertised)
- NVLink 2.0                       ~200 Gbps
- GCC 7.1.0
- ESSL 5.5.0
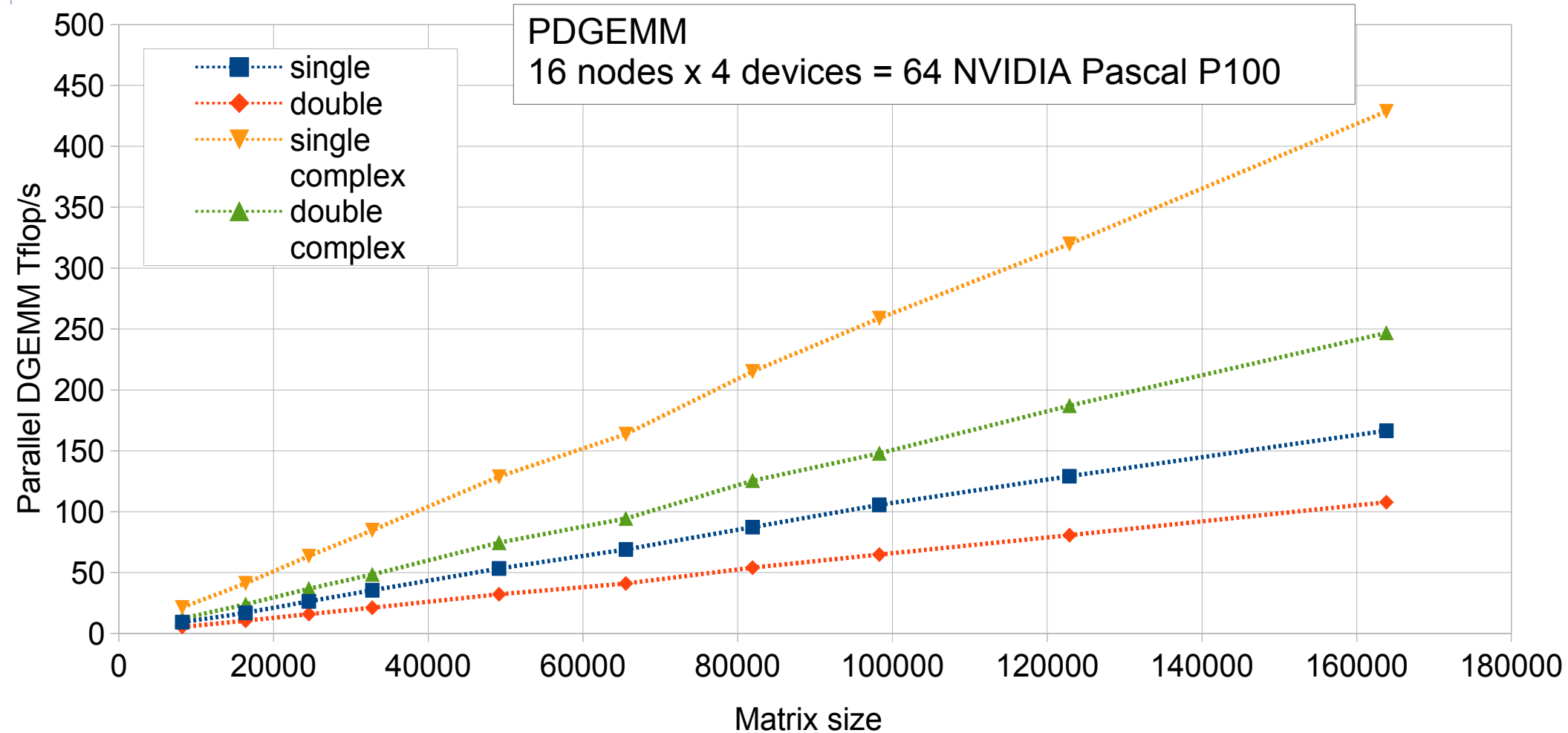- CUDA 8.0.54
- Spectrum MPI 10.1.0.3.

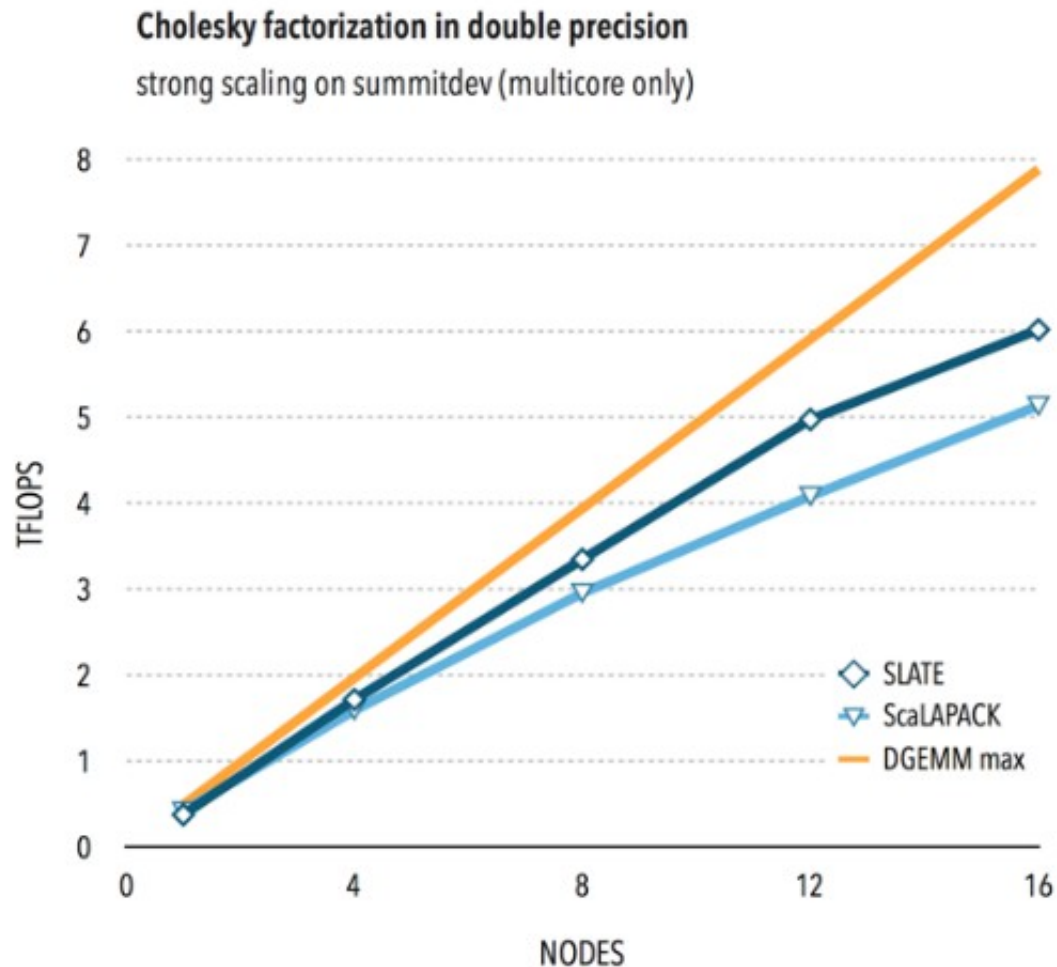# GPU-Accelerated Parallel Matrix Multiply: SLATE vs. ScaLAPACK



PDGEMM
16 nodes x 4 devices = 64 NVIDIA Pascal P100

# SLATE PDGEMM with Multiple Precisions: CPU only



PDGEMM
16 nodes x 2 sockets x 10 cores = 320 cores (IBM POWER 8)

- single
- double
- single complex

# SLATE PDGEMM with Multiple Precisions: GPUs



PDGEMM
16 nodes x 4 devices = 64 NVIDIA Pascal P100

- single
- double
- single complex
- double complex

Y-axis: Parallel DGEMM Tflop/s
X-axis: Matrix size

# SLATE Cholesky Multicore Performance

**Cholesky factorization in double precision**

strong scaling on summitdev (multicore only)
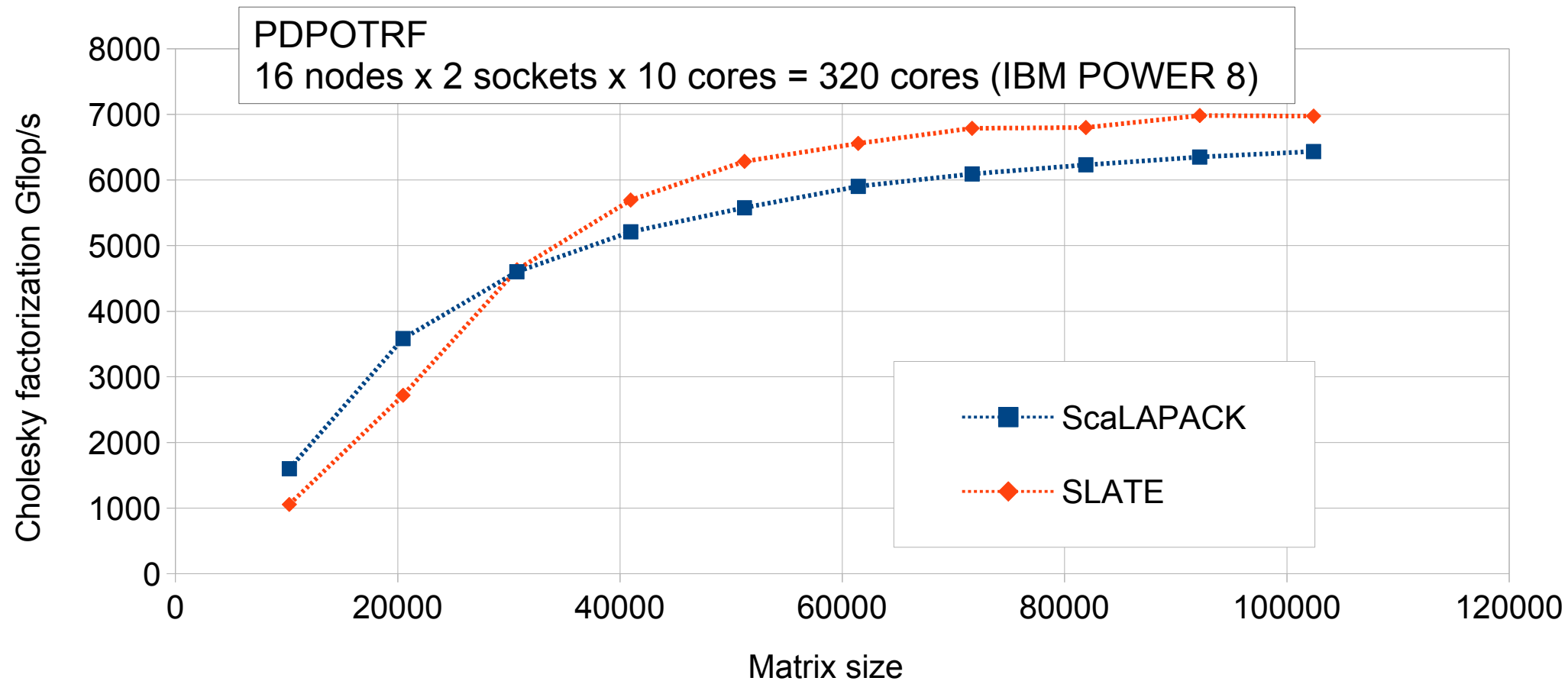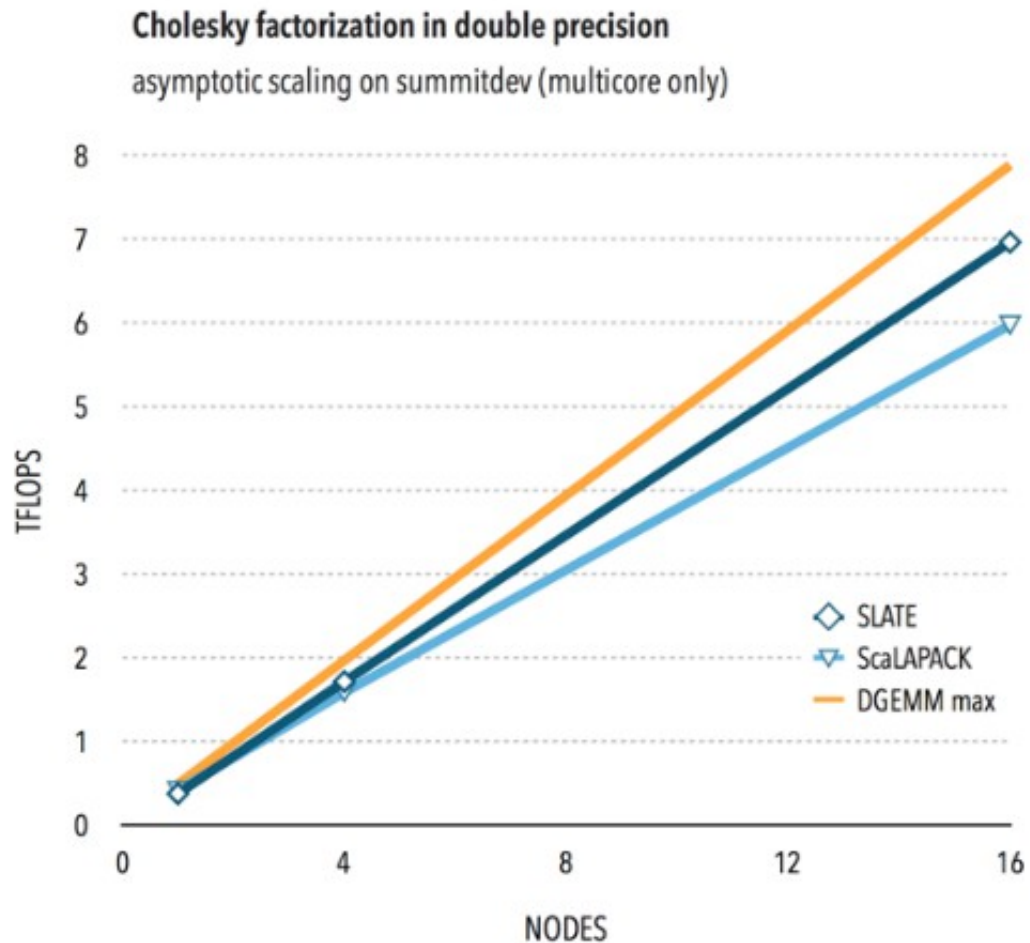


strong scaling
- 40 K × 40 K
- up to 16 nodes / 32 sockets / 320 cores

# SLATE Cholesky Factorization Comparison



PDPOTRF
16 nodes x 2 sockets x 10 cores = 320 cores (IBM POWER 8)

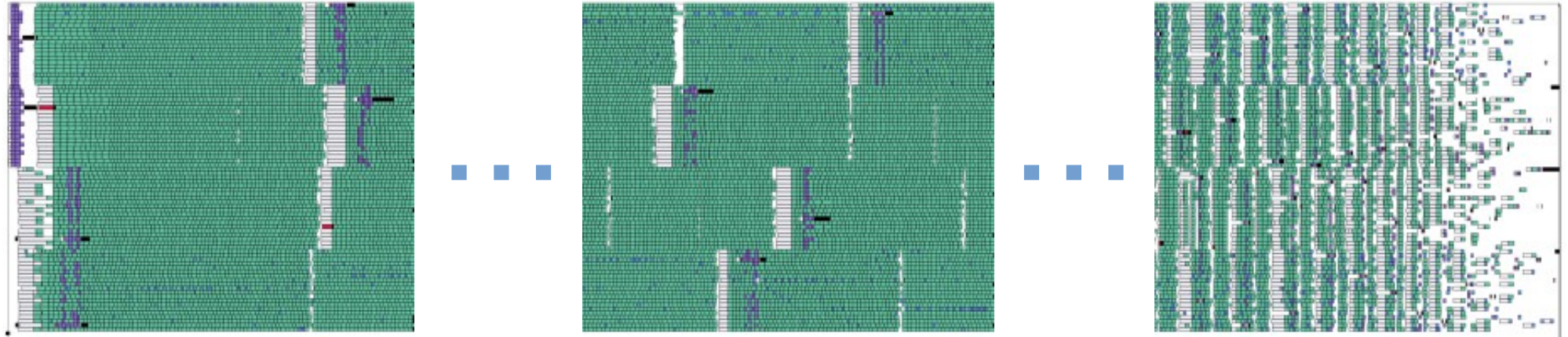**Cholesky factorization in double precision**

asymptotic scaling on summitdev (multicore only)

<u>asymptotic scaling</u>
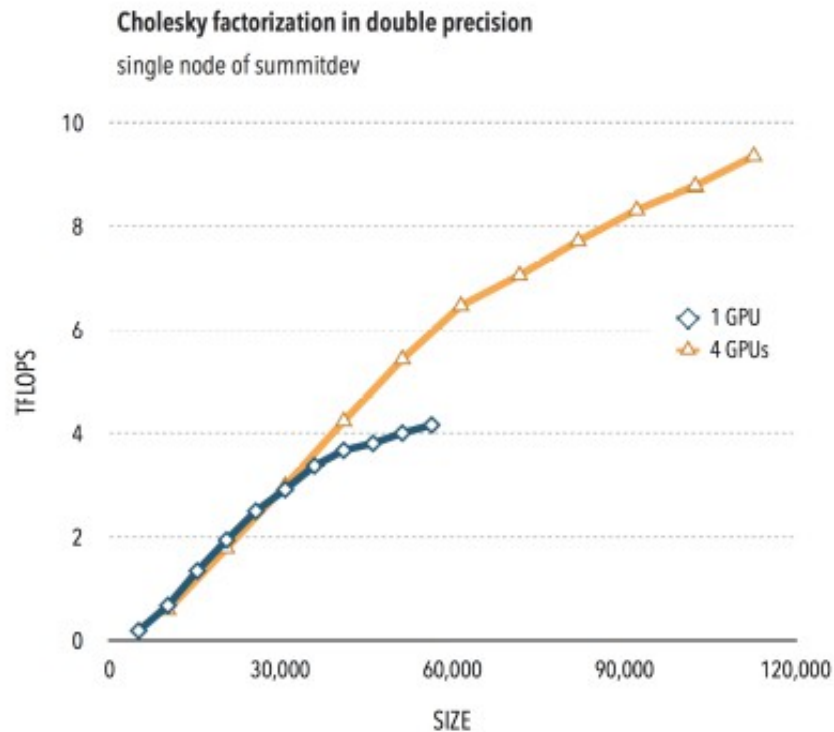
- 20 K × 20 K 1 node   20 cores
- 40 K × 40 K 4 nodes      80 cores
- 80 K × 80 K 16 nodes    320 cores

- Cholesky factorization
- 4 nodes (80 cores) factoring a 25 K × 25 K matrix using a tile size of 256
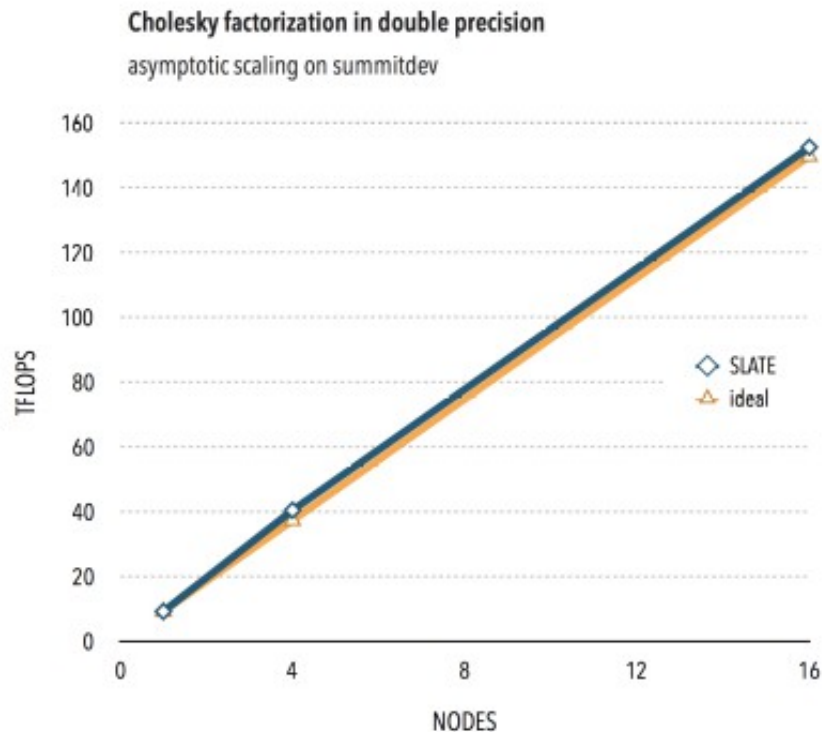- dynamic scheduling – no fork-and-join synchronization

# SLATE Cholesky GPU Performance

**Cholesky factorization in double precision**
single node of summitdev



GPU performance
- up to 56 K × 56 K        1 GPU
- up to 112 K × 112 K      4 GPUs

# SLATE GPU Performance



Cholesky factorization in double precision
asymptotic scaling on summitdev

## asymptotic scaling
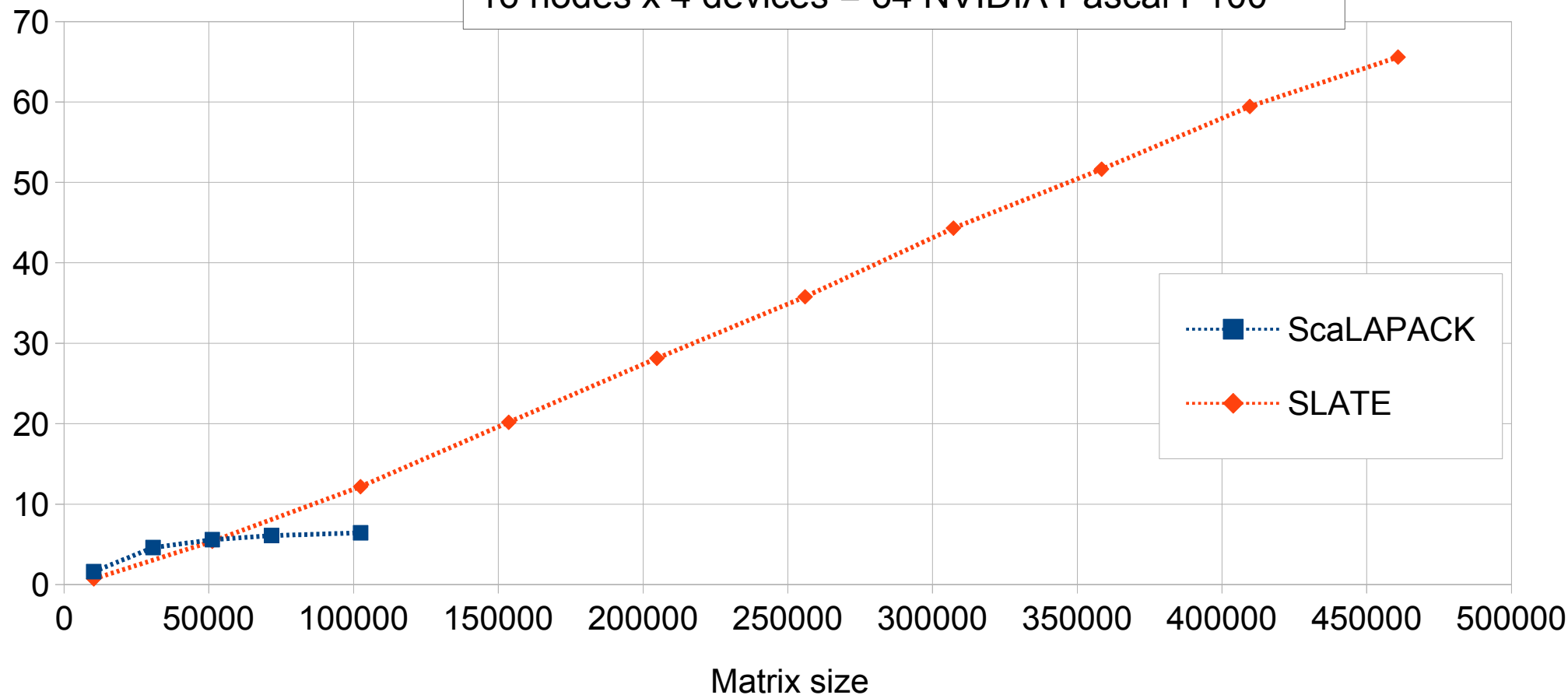- 112 K × 112 K     1 node   4 GPUs
- 225 K × 225 K     4 nodes     16 GPUs
- 450 K × 450 K     16 nodes    64 GPUs

# SLATE Cholesky Factorization: GPUs
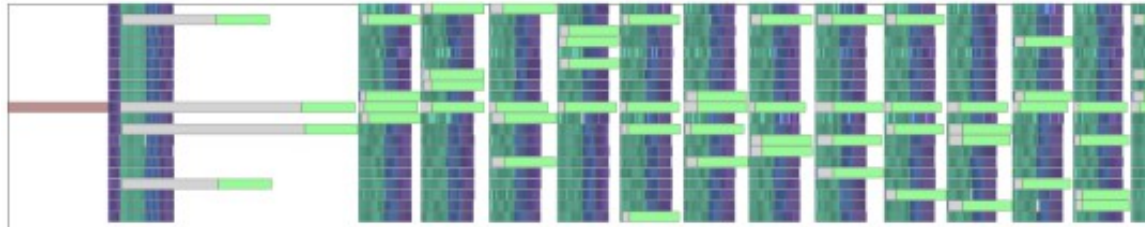


PDPOTRF
16 nodes x 4 devices = 64 NVIDIA Pascal P100

Legend: ScaLAPACK, SLATE

X-axis: Matrix size
Y-axis: Parallel Cholesky Factorization Tflop/s

Cholesky factorization
20 cores + 1 GPU
56 K × 56 K matrix
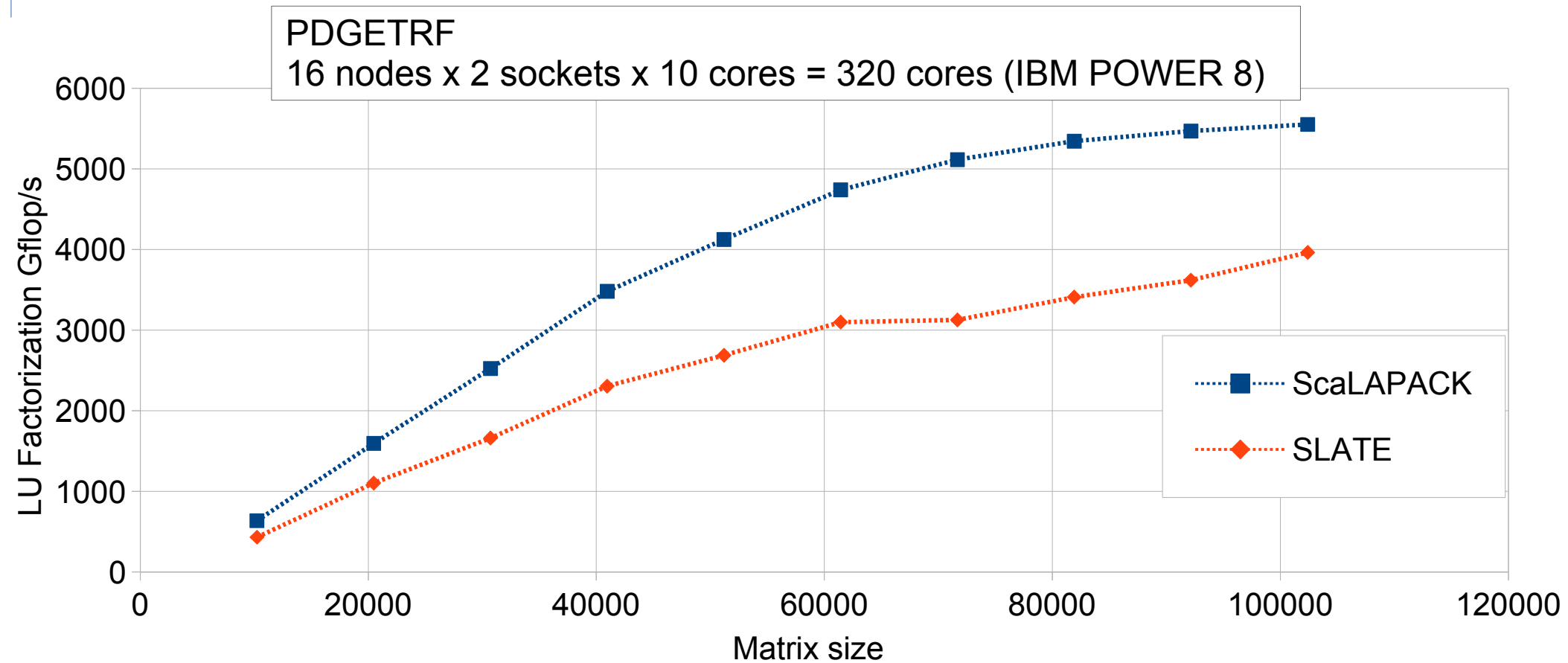tile size of 512

Cholesky factorization
20 cores + 4 GPUs
112 K × 112 K matrix
tile size of 512

# SLATE LU Performance: CPUs



PDGETRF
16 nodes x 2 sockets x 10 cores = 320 cores (IBM POWER 8)

# SLATE LU Performance: GPUs



PDGETRF
16 nodes x 4 devices = 64 NVIDIA Pascal P100

Legend:
- ScaLAPACK
- SLATE

X axis: Matrix size
Y axis: Parallel DGETRF Tflop/s

# SLATE QR Factorization Performance: CPUs



PDGEQRF
16 nodes x 2 sockets x 10 cores = 320 cores (IBM POWER 8)

- ScaLAPACK
- SLATE

X-axis: Matrix size
Y-axis: Parallel DGEQRF Gflop/s

# SLATE QR Factorization Performance: GPUs



PDGEQRF
16 nodes x 4 devices = 64 NVIDIA Pascal P100

ScaLAPACK
SLATE

Parallel DGEQRF Gflop/s

Matrix size

# SLATE Timeline: Past and Future

2016

   Q1 *research*
   Q2    *design*
   Q3      *prototyping*

2017
   Q4  C++ APIs for BLAS and LAPACK
   Q1     parallel BLAS
   Q2       parallel norms
   Q3         linear systems (LU, LLT, LDLT)

2018
   Q4           least squares (CA-QR/LQ)
   Q1            mixed precision linear systems
   Q2             matrix inversion

2019
   Q3              SVD
   Q4               EVP