# Communication Avoiding and Overlapping for Numerical Linear Algebra

Evangelos Georganas[1], Jorge González-Domínguez[2],
Edgar Solomonik[1], Yili Zheng[3], Juan Touriño[2],
Katherine Yelick[1,3]

[1]Department of Electrical Engineering and Computer Sciences, UC Berkeley
[2]Department of Electronics and Systems, University of A Coruña
[3]Lawrence Berkeley National Laboratory

June 4, 2013

**Introduction**
Linear algebra algorithms
Performance modeling
Conclusions

Communication is Expensive
Communication avoidance vs Overlapping
CA and CO in Linear Algebra

# Table of contents

**Introduction**
Linear algebra algorithms
Performance modeling
Conclusions

**Communication is Expensive**
Communication avoidance vs Overlapping
CA and CO in Linear Algebra

# Communication is Expensive

**Communication has two components:**

- **Bandwidth cost:** # of words moved / bandwidth
- **Latency cost:** # messages × latency

**Introduction**
Linear algebra algorithms
Performance modeling
Conclusions

**Communication is Expensive**
Communication avoidance vs Overlapping
CA and CO in Linear Algebra

# Communication is Expensive

**Communication has two components:**

- **Bandwidth cost:** # of words moved / bandwidth
- **Latency cost:** # messages $\times$ latency

Communication exists in **memory hierarchy** and **network**
Things are bad and **getting worse**:

$$flop\ time \ll 1/bandwidth \ll latency$$

| Annual improvements [FOSC] | | | |
|---|---|---|---|
| Flop time | | Bandwidth | Latency |
| 59% | Network | 26% | 15% |
| | DRAM | 23% | 5% |

Introduction
Linear algebra algorithms
Performance modeling
Conclusions

Communication is Expensive
Communication avoidance vs Overlapping
CA and CO in Linear Algebra

# Communication is Expensive

**Communication has two components:**

- ▶ **Bandwidth cost:** # of words moved / bandwidth
- ▶ **Latency cost:** # messages × latency

Communication exists in **memory hierarchy** and **network**

Things are bad and **getting worse**:

$$flop\ time \ll 1/bandwidth \ll latency$$

| Annual improvements [FOSC] | | | |
|---|---|---|---|
| Flop time | | Bandwidth | Latency |
| 59% | Network | 26% | 15% |
| | DRAM | 23% | 5% |

Communication is also expensive in **energy**

**Introduction**
Linear algebra algorithms
Performance modeling
Conclusions

Communication is Expensive
**Communication avoidance vs Overlapping**
CA and CO in Linear Algebra

# Communication Avoidance vs Overlapping

Two techniques to minimize the impact of communication:

**Introduction**
Linear algebra algorithms
Performance modeling
Conclusions

Communication is Expensive
**Communication avoidance vs Overlapping**
CA and CO in Linear Algebra

# Communication Avoidance vs Overlapping

Two techniques to minimize the impact of communication:

1. **Communication Avoidance (CA):**

**Introduction**
Linear algebra algorithms
Performance modeling
Conclusions

Communication is Expensive
**Communication avoidance vs Overlapping**
CA and CO in Linear Algebra

# Communication Avoidance vs Overlapping

Two techniques to minimize the impact of communication:

1. **Communication Avoidance (CA):**
   - Leads to provably optimal algorithms

**Introduction**
Linear algebra algorithms
Performance modeling
Conclusions

Communication is Expensive
**Communication avoidance vs Overlapping**
CA and CO in Linear Algebra

# Communication Avoidance vs Overlapping

Two techniques to minimize the impact of communication:

1. **Communication Avoidance (CA):**
   - Leads to provably optimal algorithms
   - There might be a trade-off in bandwidth and latency

**Introduction**
Linear algebra algorithms
Performance modeling
Conclusions

Communication is Expensive
**Communication avoidance vs Overlapping**
CA and CO in Linear Algebra

# Communication Avoidance vs Overlapping

Two techniques to minimize the impact of communication:

1. **Communication Avoidance (CA):**
   - ▶ Leads to provably optimal algorithms
   - ▶ There might be a trade-off in bandwidth and latency
2. **Communication Overlapping (CO):**

**Introduction**
Linear algebra algorithms
Performance modeling
Conclusions

Communication is Expensive
**Communication avoidance vs Overlapping**
CA and CO in Linear Algebra

# Communication Avoidance vs Overlapping

Two techniques to minimize the impact of communication:

1. **Communication Avoidance (CA):**
   - ▶ Leads to provably optimal algorithms
   - ▶ There might be a trade-off in bandwidth and latency

2. **Communication Overlapping (CO):**
   - ▶ Reduces the impact of each communication event by overlapping it with computation

**Introduction**
Linear algebra algorithms
Performance modeling
Conclusions

Communication is Expensive
**Communication avoidance vs Overlapping**
CA and CO in Linear Algebra

# Communication Avoidance vs Overlapping

Two techniques to minimize the impact of communication:

1. **Communication Avoidance (CA):**
   - ▶ Leads to provably optimal algorithms
   - ▶ There might be a trade-off in bandwidth and latency

2. **Communication Overlapping (CO):**
   - ▶ Reduces the impact of each communication event by overlapping it with computation
   - ▶ Hides bandwidth and/or latency cost

**Introduction**
Linear algebra algorithms
Performance modeling
Conclusions

Communication is Expensive
**Communication avoidance vs Overlapping**
CA and CO in Linear Algebra

# Communication Avoidance vs Overlapping

Two techniques to minimize the impact of communication:

1. **Communication Avoidance (CA):**
   - Leads to provably optimal algorithms
   - There might be a trade-off in bandwidth and latency

2. **Communication Overlapping (CO):**
   - Reduces the impact of each communication event by overlapping it with computation
   - Hides bandwidth and/or latency cost
   - Most effective if communication & computation are balanced

**Introduction**
Linear algebra algorithms
Performance modeling
Conclusions

Communication is Expensive
Communication avoidance vs Overlapping
**CA and CO in Linear Algebra**

# CA and CO in Linear Algebra

- ► We studied these techniques in three Linear Algebra routines:
  - ► Matrix Multiplication (SUMMA and Cannon's algorithms)
  - ► Cholesky factorization
  - ► Triangular Solve

**Introduction**
Linear algebra algorithms
Performance modeling
Conclusions

Communication is Expensive
Communication avoidance vs Overlapping
**CA and CO in Linear Algebra**

# CA and CO in Linear Algebra

- We studied these techniques in three Linear Algebra routines:
  - Matrix Multiplication (SUMMA and Cannon's algorithms)
  - Cholesky factorization
  - Triangular Solve
- Prior algorithms but novel implementations

| Optimizations / Algorithm | Overlapping | Avoidance | Overlapping & Avoidance |
|---|---|---|---|
| **SUMMA** | PRIOR | PRIOR | |
| **Cannon's** | PRIOR | PRIOR | |
| **Cholesky** | PRIOR | | |
| **TRSM** | PRIOR | | |

**Introduction**
Linear algebra algorithms
Performance modeling
Conclusions

Communication is Expensive
Communication avoidance vs Overlapping
**CA and CO in Linear Algebra**

# CA and CO in Linear Algebra

- ▶ We studied these techniques in three Linear Algebra routines:
  - ▶ Matrix Multiplication (SUMMA and Cannon's algorithms)
  - ▶ Cholesky factorization
  - ▶ Triangular Solve
- ▶ Prior algorithms but novel implementations

| Optimizations / Algorithm | Overlapping | Avoidance | Overlapping & Avoidance |
|---|---|---|---|
| **SUMMA** | PRIOR | PRIOR | NEW |
| **Cannon's** | PRIOR NEW: One sided communication | PRIOR NEW: One sided communication | NEW |
| **Cholesky** | PRIOR | NEW | NEW |
| **TRSM** | PRIOR | NEW* | NEW* |

*Uses replication but not communication optimal

**Introduction**
Linear algebra algorithms
Performance modeling
Conclusions

Communication is Expensive
Communication avoidance vs Overlapping
**CA and CO in Linear Algebra**

# CA and CO in Linear Algebra

Three major questions arise:

**Introduction**
Linear algebra algorithms
Performance modeling
Conclusions

Communication is Expensive
Communication avoidance vs Overlapping
**CA and CO in Linear Algebra**

# CA and CO in Linear Algebra

Three major questions arise:

1. Which is more important? CA or CO?

**Introduction**
Linear algebra algorithms
Performance modeling
Conclusions

Communication is Expensive
Communication avoidance vs Overlapping
**CA and CO in Linear Algebra**

# CA and CO in Linear Algebra

Three major questions arise:

1. Which is more important? CA or CO? It depends

**Introduction**
Linear algebra algorithms
Performance modeling
Conclusions

Communication is Expensive
Communication avoidance vs Overlapping
**CA and CO in Linear Algebra**

# CA and CO in Linear Algebra

Three major questions arise:

1. Which is more important? CA or CO? It depends
2. Is there a benefit from combining both techniques?

Introduction
Linear algebra algorithms
Performance modeling
Conclusions

Communication is Expensive
Communication avoidance vs Overlapping
CA and CO in Linear Algebra

# CA and CO in Linear Algebra

Three major questions arise:

1. Which is more important? CA or CO? It depends
2. Is there a benefit from combining both techniques? Yes

**Introduction**
Linear algebra algorithms
Performance modeling
Conclusions

Communication is Expensive
Communication avoidance vs Overlapping
**CA and CO in Linear Algebra**

# CA and CO in Linear Algebra

Three major questions arise:

1. Which is more important? CA or CO? It depends
2. Is there a benefit from combining both techniques? Yes
3. Can we explain the behavior of CA and CO under different situations?

**Introduction**
Linear algebra algorithms
Performance modeling
Conclusions

Communication is Expensive
Communication avoidance vs Overlapping
**CA and CO in Linear Algebra**

# CA and CO in Linear Algebra

Three major questions arise:

1. Which is more important? CA or CO? It depends
2. Is there a benefit from combining both techniques? Yes
3. Can we explain the behavior of CA and CO under different situations? Yes - Performance models

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

Matrix Multiplication
Cholesky factorization
Triangular Solve (TRSM)

# Table of contents

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

**Matrix Multiplication**
Cholesky factorization
Triangular Solve (TRSM)

# 2D Matrix Multiplication (SUMMA)

[Van De Geijn and Watts 97]



16 CPUs (4x4)

- ▶ Outer product form of Mat Mul
- ▶ Partitions $A$, $B$ and $C$ in 2 dimensions
- ▶ Row and column broadcast on 2D grid
- ▶ Costs:
  - ▶ $O(n^3/p)$ flops
  - ▶ $O(n^2/\sqrt{p})$ words moved
  - ▶ $O(\sqrt{p})$ messages

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

**Matrix Multiplication**
Cholesky factorization
Triangular Solve (TRSM)

# 2D Matrix Multiplication with CO (SUMMA)

- ▶ We overlap the broadcasts of next iteration with the local Mat.Mul computation of current iteration
- ▶ Theoretically the execution time becomes
  $t_{exec} = O(\max(t_{computation}, t_{communication}))$
- ▶ If communication and computation are balanced we can achieve up to $2\times$ speedup
- ▶ **Additional communication buffers are needed**

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

**Matrix Multiplication**
Cholesky factorization
Triangular Solve (TRSM)

# Experimental setup

- Experiments on Hopper, a Cray XE6 system (153,216 cores)
- We will focus on communication-limited problems (i.e. small problems on large machine configurations)
- Strong scaling experiments

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

**Matrix Multiplication**
Cholesky factorization
Triangular Solve (TRSM)

# Overlapping yields more benefits at medium scale



SUMMA on Hopper (n=16384)

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

**Matrix Multiplication**
Cholesky factorization
Triangular Solve (TRSM)

# Overlapping yields more benefits at medium scale



SUMMA on Hopper (n=16384)

- At medium scale where communication and computation are balanced we observe larger benefits (1.34× speedup)
- At large scale overlapping does not help

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

**Matrix Multiplication**
Cholesky factorization
Triangular Solve (TRSM)

# 2.5D Matrix Multiplication (SUMMA)
[McColl and Tiskin 99], [Solomonik and Demmel 11]



32 CPUs (4x4x2)

2 copies of matrices

- ▶ The 2.5D algorithm uses extra memory to reduce communication
- ▶ Each one of the $c$ layers of processors computes a different contribution to the matrix $C$
- ▶ Works for $c$ copies, $c \in [1, p^{1/3}]$
- ▶ Costs:
  - ▶ $O(n^3/p)$ flops
  - ▶ $O(n^2/\sqrt{c \cdot p})$ words moved
  - ▶ $O(\sqrt{p/c^3})$ messages

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

**Matrix Multiplication**
Cholesky factorization
Triangular Solve (TRSM)

# Communication avoidance helps a lot at large scale



SUMMA on Hopper (n=16384)

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

**Matrix Multiplication**
Cholesky factorization
Triangular Solve (TRSM)

# Communication avoidance helps a lot at large scale



- At large scale avoidance helps a lot ($2.08\times$ speedup) (there is a lot of communication to avoid!)
- At medium scale communication avoidance may yield slowdown

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

**Matrix Multiplication**
Cholesky factorization
Triangular Solve (TRSM)

# 2.5D Matrix Multiplication with CO (SUMMA)

- ▶ We overlap the broadcasts of the next iteration with the local Mat.Mul computation of current iteration **on each** of the $c$ processor layers

- ▶ **Additional communication buffers are needed on top of the extra memory needed for replication**

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

**Matrix Multiplication**
Cholesky factorization
Triangular Solve (TRSM)

# Putting everything together



SUMMA on Hopper (n=16384)

- ▶ **At medium scale overlapping itself yields best performance**
- ▶ **At large scale combining both techniques pays off**

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

**Matrix Multiplication**
Cholesky factorization
Triangular Solve (TRSM)

# Cannon's algorithm

- In SUMMA we use collective communication operations

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

**Matrix Multiplication**
Cholesky factorization
Triangular Solve (TRSM)

# Cannon's algorithm

- ▶ In SUMMA we use collective communication operations
- ▶ In Cannon's algorithm the communication needed is shifting

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

**Matrix Multiplication**
Cholesky factorization
Triangular Solve (TRSM)

# Cannon's algorithm

- ▶ In SUMMA we use collective communication operations
- ▶ In Cannon's algorithm the communication needed is shifting
- ▶ The same techniques can be applied for Cannon's algorithm

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

**Matrix Multiplication**
Cholesky factorization
Triangular Solve (TRSM)

# Cannon's algorithm

- ▶ In SUMMA we use collective communication operations
- ▶ In Cannon's algorithm the communication needed is shifting
- ▶ The same techniques can be applied for Cannon's algorithm
- ▶ Use fast **one-sided communication** provided by UPC

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

**Matrix Multiplication**
Cholesky factorization
Triangular Solve (TRSM)

# Cannon's algorithm

- ▶ In SUMMA we use collective communication operations
- ▶ In Cannon's algorithm the communication needed is shifting
- ▶ The same techniques can be applied for Cannon's algorithm
- ▶ Use fast **one-sided communication** provided by UPC



Cannon on Hopper (n=16384)

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

Matrix Multiplication
**Cholesky factorization**
Triangular Solve (TRSM)

# 2D Cholesky factorization

- Factorize a symmetric positive definite matrix $A$ into $A = L \cdot L^T$, where $L$ is lower triangular
- Take advantage of symmetry and store only half of matrix $A$
- Employ block-cyclic layout for load-balance

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

Matrix Multiplication
**Cholesky factorization**
Triangular Solve (TRSM)

# 2D Cholesky factorization



1. Factorize the upper-left (green) block & broadcast it to the column

2. Update via TRSM the (yellow) block column

3. Broadcast the factorized column in two phases & update the trailing matrix (white blocks)

4. Continue with the factorization of the second block column and repeat previous steps until all matrix is factorized

- Communication involved is row and column broadcasts
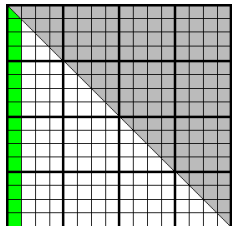- There are dependencies among rows and column during factorization

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

Matrix Multiplication
**Cholesky factorization**
Triangular Solve (TRSM)

# 2D Cholesky factorization



1. Factorize the upper-left (green) block & broadcast it to the column

2. Update via TRSM the (yellow) block column

3. Broadcast the factorized column in two phases & update the trailing matrix (white blocks)

4. Continue with the factorization of the second block column and repeat previous steps until all matrix is factorized
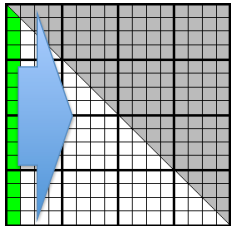
▶ Communication involved is row and column broadcasts

▶ There are dependencies among rows and column during factorization

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

Matrix Multiplication
**Cholesky factorization**
Triangular Solve (TRSM)

# 2D Cholesky factorization



1. Factorize the upper-left (green) block & broadcast it to the column

2. Update via TRSM the (yellow) block column

3. Broadcast the factorized column in two phases & update the trailing matrix (white blocks)

4. Continue with the factorization of the second block column and repeat previous steps until all matrix is factorized

- Communication involved is row and column broadcasts

- There are dependencies among rows and column during factorization

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

Matrix Multiplication
Cholesky factorization
Triangular Solve (TRSM)

# 2D Cholesky factorization



1. Factorize the upper-left (green) block &
   broadcast it to the column

2. Update via TRSM the (yellow) block column

3. Broadcast the factorized column in two phases
   & update the trailing matrix (white blocks)

4. Continue with the factorization of the second
   block column and repeat previous steps until all
   matrix is factorized

▶ Communication involved is row and column broadcasts

▶ There are dependencies among rows and column during
  factorization

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

Matrix Multiplication
**Cholesky factorization**
Triangular Solve (TRSM)

# 2D Cholesky factorization



1. Factorize the upper-left (green) block & broadcast it to the column

2. Update via TRSM the (yellow) block column

3. Broadcast the factorized column in two phases & update the trailing matrix (white blocks)

4. Continue with the factorization of the second block column and repeat previous steps until all matrix is factorized

▶ Communication involved is row and column broadcasts

▶ There are dependencies among rows and column during factorization

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

Matrix Multiplication
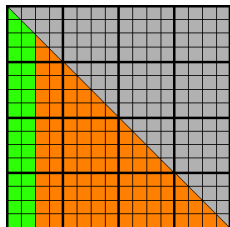**Cholesky factorization**
Triangular Solve (TRSM)

# 2D Cholesky factorization with overlapping



1. Factorize the $1^{st}$ block-column
2. Broadcast factorized column
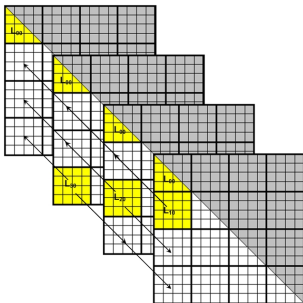3. Update & factorize **only** the $2^{nd}$ block-column
4. **Overlap**
   - broadcast of the $2^{nd}$ block-column
   - update of the rest trailing matrix (using $1^{st}$ block-column)

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

Matrix Multiplication
**Cholesky factorization**
Triangular Solve (TRSM)

# 2D Cholesky factorization with overlapping



1. Factorize the $1^{st}$ block-column
2. Broadcast factorized column
3. Update & factorize **only** the $2^{nd}$ block-column
4. **Overlap**
   - broadcast of the $2^{nd}$ block-column
   - update of the rest trailing matrix
     (using $1^{st}$ block-column)

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

Matrix Multiplication
**Cholesky factorization**
Triangular Solve (TRSM)

# 2D Cholesky factorization with overlapping



1. Factorize the $1^{st}$ block-column

2. Broadcast factorized column

3. Update & factorize **only** the $2^{nd}$ block-column

4. **Overlap**
   - broadcast of the $2^{nd}$ block-column
   - update of the rest trailing matrix
     (using $1^{st}$ block-column)

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

Matrix Multiplication
**Cholesky factorization**
Triangular Solve (TRSM)

# 2D Cholesky factorization with overlapping



1. Factorize the $1^{st}$ block-column
2. Broadcast factorized column
3. Update & factorize **only** the $2^{nd}$ block-column
4. **Overlap**
   - broadcast of the $2^{nd}$ block-column
   - update of the rest trailing matrix
     (using $1^{st}$ block-column)

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

Matrix Multiplication
**Cholesky factorization**
Triangular Solve (TRSM)

# 2.5D Cholesky factorization
[McColl and Tiskin 99], [Solomonik and Demmel 11]
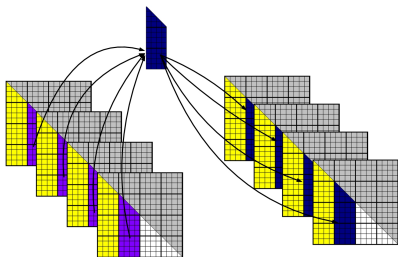
▶ Employ two levels of blocking: "Fat panels" and "blocks"



1. All layers jointly factorize a "fat" panel

2. Broadcast different subpanels within each layer & update trailing matrices

3. All-reduce the next "fat" panel to accumulate the updates

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

Matrix Multiplication
**Cholesky factorization**
Triangular Solve (TRSM)
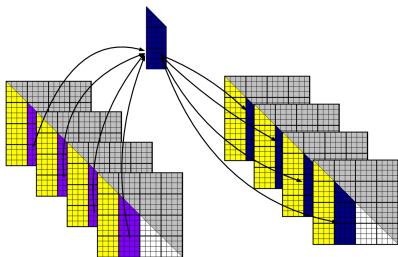
# 2.5D Cholesky factorization
[McColl and Tiskin 99], [Solomonik and Demmel 11]



1. All layers jointly factorize a "fat" panel

2. Broadcast different subpanels within each layer & update trailing matrices

3. All-reduce the next "fat" panel to accumulate the updates

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

Matrix Multiplication
**Cholesky factorization**
Triangular Solve (TRSM)
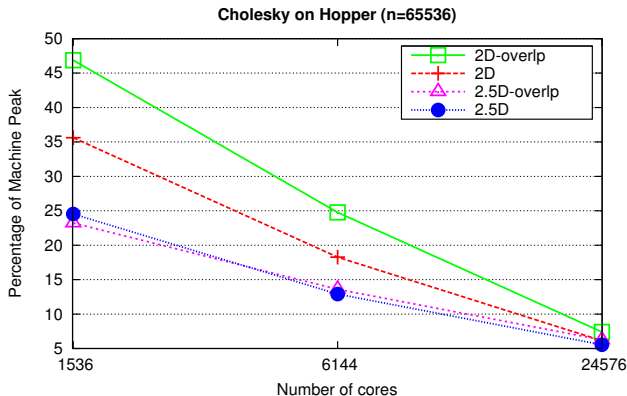
# 2.5D Cholesky factorization
[McColl and Tiskin 99], [Solomonik and Demmel 11]



1. All layers jointly factorize a "fat" panel
2. Broadcast different subpanels within each layer & update trailing matrices

Recall matrix multiplication !!!

3. All-reduce the next "fat" panel to accumulate the updates

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

Matrix Multiplication
**Cholesky factorization**
Triangular Solve (TRSM)

# 2.5D Cholesky factorization
[McColl and Tiskin 99], [Solomonik and Demmel 11]



1. All layers jointly factorize a "fat" panel
2. Broadcast different subpanels within each layer & update trailing matrices
3. All-reduce the next "fat" panel to accumulate the updates

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

Matrix Multiplication
**Cholesky factorization**
Triangular Solve (TRSM)

# 2.5D Cholesky factorization
[McColl and Tiskin 99], [Solomonik and Demmel 11]



1. All layers jointly factorize a "fat" panel
2. Broadcast different subpanels within each layer & update trailing matrices
3. All-reduce the next "fat" panel to accumulate the updates
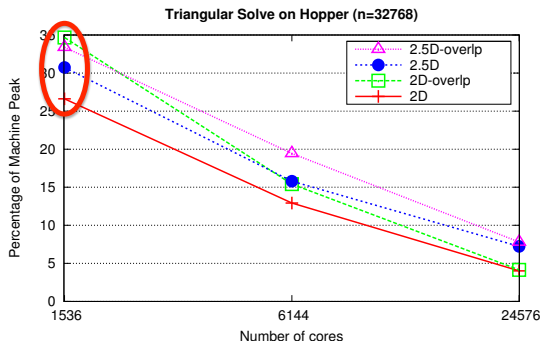
▶ At step 2 we can overlap computation and communication similarly to the 2D version

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

Matrix Multiplication
**Cholesky factorization**
Triangular Solve (TRSM)

# Performance results on Hopper (Cray XE6)



Cholesky on Hopper (n=65536)

- ▶ CO helps more at the smallest scale (1,536 cores)
- ▶ Have not reached yet the cross-point of CA and 2D
- ▶ Future work: Aggregate updates to improve performance

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

Matrix Multiplication
Cholesky factorization
**Triangular Solve (TRSM)**

# Triangular Solve (TRSM)

- Computes $X$, such that $X \cdot U = B$ with upper-triangular $U$
- Similar parallelization to Cholesky



Triangular Solve on Hopper (n=32768)

- At small scale overlapping outperforms other versions

Introduction
**Linear algebra algorithms**
Performance modeling
Conclusions

Matrix Multiplication
Cholesky factorization
**Triangular Solve (TRSM)**

# Triangular Solve (TRSM)

- ▶ Computes $X$, such that $X \cdot U = B$ with upper-triangular $U$
- ▶ Similar parallelization to Cholesky



- ▶ At small scale overlapping outperforms other versions
- ▶ At large scale the combining optimizations is the best option

**Introduction**
**Linear algebra algorithms**
**Performance modeling**
**Conclusions**

Motivation
Methodology

# Table of contents

Introduction
Linear algebra algorithms
**Performance modeling**
Conclusions

**Motivation**
Methodology

# Can we explain the behavior of CO and CA?

- ► Communication avoiding and overlapping:

Introduction
Linear algebra algorithms
**Performance modeling**
Conclusions

**Motivation**
Methodology

# Can we explain the behavior of CO and CA?

- ▶ Communication avoiding and overlapping:
  - ▶ Introduce four different algorithmic variants for each routine

Introduction
Linear algebra algorithms
Performance modeling
Conclusions

Motivation
Methodology

# Can we explain the behavior of CO and CA?

- ▶ Communication avoiding and overlapping:
  - ▶ Introduce four different algorithmic variants for each routine
  - ▶ Have different memory requirements

Introduction
Linear algebra algorithms
**Performance modeling**
Conclusions

**Motivation**
Methodology

# Can we explain the behavior of CO and CA?

- ▶ Communication avoiding and overlapping:
  - ▶ Introduce four different algorithmic variants for each routine
  - ▶ Have different memory requirements
  - ▶ Lead to performance gains dependent on the problem size, the algorithm and the machine configuration

Introduction
Linear algebra algorithms
Performance modeling
Conclusions

Motivation
Methodology

# Can we explain the behavior of CO and CA?

- ▶ Communication avoiding and overlapping:
  - ▶ Introduce four different algorithmic variants for each routine
  - ▶ Have different memory requirements
  - ▶ Lead to performance gains dependent on the problem size, the algorithm and the machine configuration
  - ▶ Tunable parameters with complicated interactions

Introduction
Linear algebra algorithms
**Performance modeling**
Conclusions

**Motivation**
Methodology

# Can we explain the behavior of CO and CA?

- ▶ Communication avoiding and overlapping:
  - ▶ Introduce four different algorithmic variants for each routine
  - ▶ Have different memory requirements
  - ▶ Lead to performance gains dependent on the problem size, the algorithm and the machine configuration
  - ▶ Tunable parameters with complicated interactions

| Effect<br>Parameter | # messages | load<br>balance | computation<br>efficiency |
|---|---|---|---|
| **block size ↑** | ↓ | ↓ | ↑ |
| **replication ↑** | ↑/↓ | ↑/↓ | NA |
| **fat panel size↑** | ↓ | ↓ | NA |

Introduction
Linear algebra algorithms
Performance modeling
Conclusions

**Motivation**
Methodology

# Can we explain the behavior of CO and CA?

- Communication avoiding and overlapping:
  - Introduce four different algorithmic variants for each routine
  - Have different memory requirements
  - Lead to performance gains dependent on the problem size, the algorithm and the machine configuration
  - Tunable parameters with complicated interactions

| Effect<br>Parameter | # messages | load<br>balance | computation<br>efficiency |
|:---:|:---:|:---:|:---:|
| **block size ↑** | ↓ | ↓ | ↑ |
| **replication ↑** | ↑/↓ | ↑/↓ | NA |
| **fat panel size↑** | ↓ | ↓ | NA |

- Communication performance depends on number of processors

Introduction
Linear algebra algorithms
Performance modeling
Conclusions

**Motivation**
Methodology

# Can we explain the behavior of CO and CA?

- ▶ Communication avoiding and overlapping:
  - ▶ Introduce four different algorithmic variants for each routine
  - ▶ Have different memory requirements
  - ▶ Lead to performance gains dependent on the problem size, the algorithm and the machine configuration
  - ▶ Tunable parameters with complicated interactions

| Effect / Parameter | # messages | load balance | computation efficiency |
|---|---|---|---|
| block size ↑ | ↓ | ↓ | ↑ |
| replication ↑ | ↑/↓ | ↑/↓ | NA |
| fat panel size↑ | ↓ | ↓ | NA |

- ▶ Communication performance depends on number of processors
- ▶ **Given a problem instance and a machine configuration would like to predict the optimal variant**

Introduction
Linear algebra algorithms
**Performance modeling**
Conclusions

Motivation
**Methodology**

# Methodology for constructing performance models

▶ We construct detailed performance models

Introduction
Linear algebra algorithms
**Performance modeling**
Conclusions

Motivation
**Methodology**

# Methodology for constructing performance models

- We construct detailed performance models
  - Inputs: matrix size, # of processors, BLAS efficiency, LogGP parameters, block sizes, replication factors, fat panel sizes

Introduction
Linear algebra algorithms
**Performance modeling**
Conclusions

Motivation
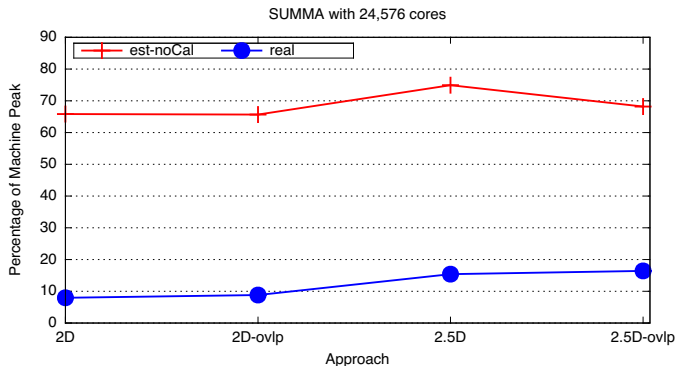**Methodology**

# Methodology for constructing performance models

- We construct detailed performance models
  - Inputs: matrix size, # of processors, BLAS efficiency, LogGP parameters, block sizes, replication factors, fat panel sizes
  - Output: Estimate for execution time of algorithm

Introduction
Linear algebra algorithms
**Performance modeling**
Conclusions

Motivation
**Methodology**

# Methodology for constructing performance models

- ▶ We construct detailed performance models
  - ▶ Inputs: matrix size, # of processors, BLAS efficiency, LogGP parameters, block sizes, replication factors, fat panel sizes
  - ▶ Output: Estimate for execution time of algorithm
- ▶ We track the execution flow of each algorithm and estimate completion time for encountered operation

Introduction
Linear algebra algorithms
**Performance modeling**
Conclusions

Motivation
**Methodology**
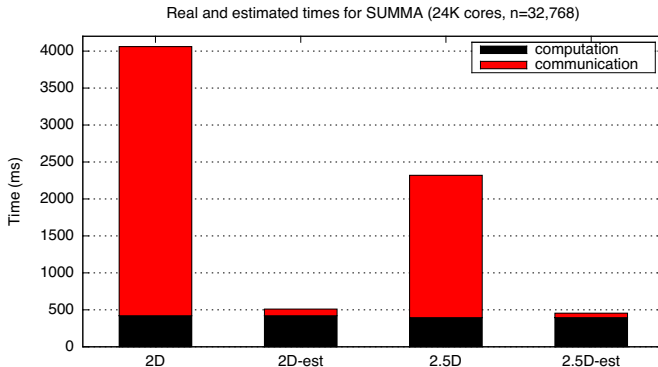
# Methodology for constructing performance models

- ▶ We construct detailed performance models
  - ▶ Inputs: matrix size, # of processors, BLAS efficiency, LogGP parameters, block sizes, replication factors, fat panel sizes
  - ▶ Output: Estimate for execution time of algorithm
- ▶ We track the execution flow of each algorithm and estimate completion time for encountered operation
  - ▶ Estimate computation times through BLAS microbenchmarks

Introduction
Linear algebra algorithms
**Performance modeling**
Conclusions

Motivation
**Methodology**

# Methodology for constructing performance models

- ▶ We construct detailed performance models
  - ▶ Inputs: matrix size, # of processors, BLAS efficiency, LogGP parameters, block sizes, replication factors, fat panel sizes
  - ▶ Output: Estimate for execution time of algorithm
- ▶ We track the execution flow of each algorithm and estimate completion time for encountered operation
  - ▶ Estimate computation times through BLAS microbenchmarks
  - ▶ Estimate communication times through LogGP model (collective models from [Thakur, Rabenseifner, Gropp 2005])

Introduction
Linear algebra algorithms
**Performance modeling**
Conclusions

Motivation
**Methodology**

# Methodology for constructing performance models

- ▶ We construct detailed performance models
  - ▶ Inputs: matrix size, # of processors, BLAS efficiency, LogGP parameters, block sizes, replication factors, fat panel sizes
  - ▶ Output: Estimate for execution time of algorithm
- ▶ We track the execution flow of each algorithm and estimate completion time for encountered operation
  - ▶ Estimate computation times through BLAS microbenchmarks
  - ▶ Estimate communication times through LogGP model (collective models from [Thakur, Rabenseifner, Gropp 2005])
- ▶ We take into account possible idle times

# First approach: Ignore network congestion



SUMMA with 24,576 cores

- We predict correctly the relative performance of the CA algorithms
- Prediction of absolute performance is inaccurate

Introduction
Linear algebra algorithms
**Performance modeling**
Conclusions

Motivation
**Methodology**

# First approach: Ignore network congestion



Real and estimated times for SUMMA (24K cores, n=32,768)

- ▶ We predict accurately the computation time
- ▶ We ignore congestion $\rightarrow$ optimistic communication time

Introduction
Linear algebra algorithms
**Performance modeling**
Conclusions

Motivation
Methodology

# Quantify degradation due to congestion

- **Calibration factor**: $\frac{ideal\ BW}{real\ BW}$ when several processes use the network simultaneously

Introduction
Linear algebra algorithms
**Performance modeling**
Conclusions

Motivation
**Methodology**

# Quantify degradation due to congestion

- **Calibration factor**: $\frac{ideal\ BW}{real\ BW}$ when several processes use the network simultaneously
- Extract calibration factors via microbenchmarks

# Quantify degradation due to congestion

▶ **Calibration factor**: $\frac{ideal\ BW}{real\ BW}$ when several processes use the network simultaneously

▶ Extract calibration factors via microbenchmarks



Calibration Factors on Hopper

# Including calibration factors in the models



SUMMA with 24,576 cores

- ▶ **We predict accurately the absolute performance**

# Including calibration factors in the models



SUMMA with 24,576 cores

- **We predict accurately the absolute performance**
- Similar results for the rest algorithms

# Table of contents

## Conclusions

1. Which is more important? CA or CO?

# Conclusions

1. Which is more important? CA or CO? It depends

# Conclusions

1. Which is more important? CA or CO? It depends
   - ▶ For core counts where communication and computation are balanced CO helps more (up to $1.82\times$ speedup)
   - ▶ For larger core counts CA is more beneficial (up to $2.08\times$ speedup)

# Conclusions

1. Which is more important? CA or CO? It depends
   ▶ For core counts where communication and computation are balanced CO helps more (up to $1.82\times$ speedup)
   ▶ For larger core counts CA is more beneficial (up to $2.08\times$ speedup)
2. Is there a benefit from combining both techniques?

# Conclusions

1. Which is more important? CA or CO? It depends
   - For core counts where communication and computation are balanced CO helps more (up to $1.82\times$ speedup)
   - For larger core counts CA is more beneficial (up to $2.08\times$ speedup)
2. Is there a benefit from combining both techniques? Yes

# Conclusions

1. Which is more important? CA or CO? It depends
   - For core counts where communication and computation are balanced CO helps more (up to $1.82\times$ speedup)
   - For larger core counts CA is more beneficial (up to $2.08\times$ speedup)
2. Is there a benefit from combining both techniques? Yes
   - CO helps further hide the minimized communication by CA (up to $2.33\times$ speedup)
3. Can we explain the behavior of CA and CO under different situations?

## Conclusions

1. Which is more important? CA or CO? It depends
   - For core counts where communication and computation are balanced CO helps more (up to $1.82\times$ speedup)
   - For larger core counts CA is more beneficial (up to $2.08\times$ speedup)
2. Is there a benefit from combining both techniques? Yes
   - CO helps further hide the minimized communication by CA (up to $2.33\times$ speedup)
3. Can we explain the behavior of CA and CO under different situations? Yes - Performance models

# Conclusions

1. Which is more important? CA or CO? It depends
   - For core counts where communication and computation are balanced CO helps more (up to $1.82\times$ speedup)
   - For larger core counts CA is more beneficial (up to $2.08\times$ speedup)
2. Is there a benefit from combining both techniques? Yes
   - CO helps further hide the minimized communication by CA (up to $2.33\times$ speedup)
3. Can we explain the behavior of CA and CO under different situations? Yes - Performance models
   - We developed detailed performance models
   - They encapsulate complicated interactions between parameters
   - They predict correctly the performance

# Thank you!