Communication Avoiding and Overlapping for Numerical Linear Algebra

Evangelos Georganas¹, Jorge González-Domínguez², Edgar Solomonik¹, Yili Zheng³, Juan Touriño², Katherine Yelick^{1,3}

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

November 15, 2012

Introduction

Linear algebra algorithms Performance modeling Conclusions Communication is Expensive Communication avoidance vs Overlapping CA and CO in Linear Algebra

Table of contents

Introduction

- Communication is Expensive
- Communication avoidance vs Overlapping
- CA and CO in Linear Algebra
- 2 Linear algebra algorithms
 - Matrix Multiplication
 - Cholesky factorization
 - Triangular Solve (TRSM)
- 3 Performance modeling
 - Motivation
 - Methodology
- 4 Conclusions• 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 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 $\ll 1/\textit{bandwidth} \ll \textit{latency}$

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

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 $\ll 1/\textit{bandwidth} \ll \textit{latency}$

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

Communication is also expensive in energy

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:

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):

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

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

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):

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

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

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

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

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		

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

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

CA and CO in Linear Algebra

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?

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

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

CA and CO in Linear Algebra

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

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

CA and CO in Linear Algebra

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

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

CA and CO in Linear Algebra

- 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?

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

CA and CO in Linear Algebra

- 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

Matrix Multiplication Cholesky factorization Triangular Solve (TRSM)

Table of contents

Introduction

- Communication is Expensive
- Communication avoidance vs Overlapping
- CA and CO in Linear Algebra
- 2 Linear algebra algorithms
 - Matrix Multiplication
 - Cholesky factorization
 - Triangular Solve (TRSM)
- 3 Performance modeling
 - Motivation
 - Methodology

Conclusions Conclusions

Matrix Multiplication Cholesky factorization Triangular Solve (TRSM)

2D Matrix Multiplication (SUMMA) [Van De Geijn and Watts 97]



- 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

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× speedup
- Additional communication buffers are needed

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

Matrix Multiplication Cholesky factorization Triangular Solve (TRSM)

Overlapping yields more benefits at medium scale



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

Matrix Multiplication Cholesky factorization Triangular Solve (TRSM)

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



- 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

Matrix Multiplication Cholesky factorization Triangular Solve (TRSM)

Communication avoidance helps a lot at large scale



Matrix Multiplication Cholesky factorization Triangular Solve (TRSM)

Communication avoidance helps a lot at large scale



SUMMA on Hopper (n=16384)

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

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

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

Matrix Multiplication Cholesky factorization Triangular Solve (TRSM)

Cannon's algorithm

In SUMMA we use collective communication operations

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
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

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

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



Matrix Multiplication Cholesky factorization Triangular Solve (TRSM)

- 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

Matrix Multiplication Cholesky factorization Triangular Solve (TRSM)



- 1. Factorize the upper-left (green) block & broadcast it to the column
 - 2. Update via TRSM the (yellow) block column
 - Broadcast the factorized column in two phases
 & update the trailing matrix (white blocks)
 - I. 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

Matrix Multiplication Cholesky factorization Triangular Solve (TRSM)

2D Cholesky factorization



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

2. Update via TRSM the (yellow) block column

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

Matrix Multiplication Cholesky factorization Triangular Solve (TRSM)



- .. Factorize the upper-left (green) block & broadcast it to the column
- . Update via TRSM the (yellow) block column
- Broadcast the factorized column in two phases & update the trailing matrix (white blocks)
 - 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

Matrix Multiplication Cholesky factorization Triangular Solve (TRSM)



- .. Factorize the upper-left (green) block & broadcast it to the column
- . Update via TRSM the (yellow) block column
- Broadcast the factorized column in two phases
 & update the trailing matrix (white blocks)
- 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

Matrix Multiplication Cholesky factorization Triangular Solve (TRSM)



- .. Factorize the upper-left (green) block & broadcast it to the column
- 2. Update via TRSM the (yellow) block column
- Broadcast the factorized column in two phases
 & update the trailing matrix (white blocks)
- 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

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

Matrix Multiplication Cholesky factorization Triangular Solve (TRSM)

2D Cholesky factorization with overlapping



1. Factorize the 1st block-column

2. Broadcast factorized column

Update & factorize only the 2nd block-column
 Overlap

- broadcast of the 2nd block-column
- update of the rest trailing matrix (using 1st block-column)

Matrix Multiplication Cholesky factorization Triangular Solve (TRSM)

2D Cholesky factorization with overlapping



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

2. Broadcast factorized column

Update & factorize only the 2nd block-column
 Overlap

- broadcast of the 2nd block-column
- update of the rest trailing matrix (using 1st block-column)

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 2nd block-column
- update of the rest trailing matrix (using 1st block-column)

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

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

Matrix Multiplication Cholesky factorization Triangular Solve (TRSM)

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



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



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

Matrix Multiplication Cholesky factorization Triangular Solve (TRSM)

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



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

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

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
- Broadcast different subpanels within each layer & update trailing matrices

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

At step 2 we can overlap computation and communication similarly to the 2D version

Matrix Multiplication Cholesky factorization Triangular Solve (TRSM)

Performance results on Hopper (Cray XE6)



- 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

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

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

Motivation Methodology

Table of contents

Introduction

- Communication is Expensive
- Communication avoidance vs Overlapping
- CA and CO in Linear Algebra
- 2 Linear algebra algorithms
 - Matrix Multiplication
 - Cholesky factorization
 - Triangular Solve (TRSM)
- ③ Performance modeling
 - Motivation
 - Methodology



Motivation Methodology

Can we explain the behavior of CO and CA?

Communication avoiding and overlapping:

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

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

- 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

- 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

- 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 🛧	¥	↓	1
replication 🛧	^/↓	^/↓	NA
fat panel size∱	¥	4	NA

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 🛧	¥	↓	1
replication 🛧	^/↓	^/↓	NA
fat panel size ↑	¥	\	NA

Communication performance depends on number of processors

- 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 🛧	↓	V	1
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

Motivation Methodology

Methodology for constructing performance models

We construct detailed performance models

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

- 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 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

- 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

- 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
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
- We take into account possible idle times

Motivation Methodology

Models make correct qualitative predictions



- The models predict correctly the performance ranking of the four variants
- They encapsulate the complicated interactions
- They can take into account memory limitations
- Ongoing work shows that we can predict absolute performance accurately through integration of congestion

Table of contents

- Communication is Expensive
- Communication avoidance vs Overlapping
- CA and CO in Linear Algebra
- - Matrix Multiplication
 - Cholesky factorization
 - Triangular Solve (TRSM)
- - Motivation
 - Methodology



4 Conclusions Conclusions

Conclusions

Conclusions

1. Which is more important? CA or CO?

Conclusions

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× speedup)
 - ► For larger core counts CA is more beneficial (up to 2.08× 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× speedup)
 - ► For larger core counts CA is more beneficial (up to 2.08× 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× speedup)
 - ► For larger core counts CA is more beneficial (up to 2.08× 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× speedup)
 - ► For larger core counts CA is more beneficial (up to 2.08× speedup)
- 2. Is there a benefit from combining both techniques? Yes
 - CO helps further hide the minimized communication by CA (up to 2.33× 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× speedup)
 - ► For larger core counts CA is more beneficial (up to 2.08× speedup)
- 2. Is there a benefit from combining both techniques? Yes
 - CO helps further hide the minimized communication by CA (up to 2.33× 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× speedup)
 - ► For larger core counts CA is more beneficial (up to 2.08× speedup)
- 2. Is there a benefit from combining both techniques? Yes
 - CO helps further hide the minimized communication by CA (up to 2.33× 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 ranking

Conclusions

Thank you!