

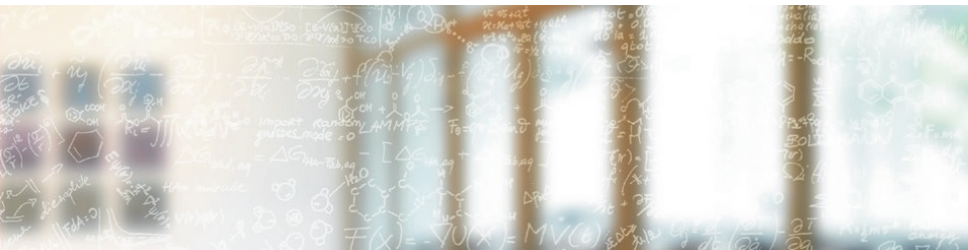


CSCS

Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre



ETH zürich



A PCIe Congestion-Aware Performance Model for Densely Populated Accelerator Servers

Maxime Martinasso^{*}, Grzegorz Kwasniewski[†], Sadaf R. Alam^{*},
Thomas C. Schulthess^{*‡§}, Torsten Hoefler[†]

^{*}Swiss National Supercomputing Centre, ETH Zurich, 6900 Lugano, Switzerland

[†]Department of Computer Science, ETH Zurich, Universitätstr. 6, 8092 Zurich, Switzerland

[‡]Institute for Theoretical Physics, ETH Zurich, 8093 Zurich, Switzerland

[§]Computer Science and Mathematics Division, Oak Ridge National Laboratory, USA

Why more densely populated accelerator servers?

- accelerators are faster and more energy-efficient than CPU
- densely populated accelerator servers are high performance nodes
- reduce space occupancy of the data center

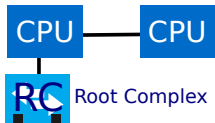
Cray CS Storm – new MeteoSwiss supercomputer

- 2 cabinets
- 12 hybrid computing nodes per cabinet
- 2 Intel Haswell 12-core CPUs per node
- 8 NVIDIA Tesla K80 GPU accelerators per node
- 2 GPU processors per accelerator
- 192 GPU processors in total
- 360 GPU teraflops in total
- **Production system**
- GPUs connected by PCI-Express



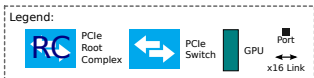
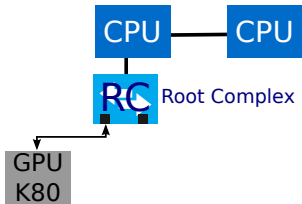
- generation 3, 16 GB/s using x16 wide lane
- dual simplex (a pair of unidirectional links)
- exchange buffer availability between pair of ports of a link
- tree-based topology

Building a densely populated accelerator servers with PCIe:



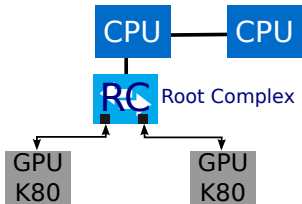
- generation 3, 16 GB/s using x16 wide lane
- dual simplex (a pair of unidirectional links)
- exchange buffer availability between pair of ports of a link
- tree-based topology

Building a densely populated accelerator servers with PCIe:



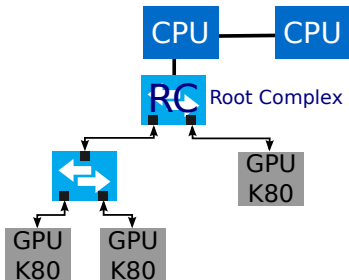
- generation 3, 16 GB/s using x16 wide lane
- dual simplex (a pair of unidirectional links)
- exchange buffer availability between pair of ports of a link
- tree-based topology

Building a densely populated accelerator servers with PCIe:



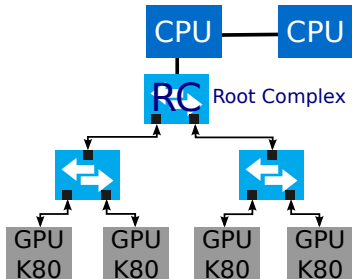
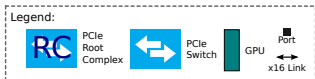
- generation 3, 16 GB/s using x16 wide lane
- dual simplex (a pair of unidirectional links)
- exchange buffer availability between pair of ports of a link
- tree-based topology

Building a densely populated accelerator servers with PCIe:



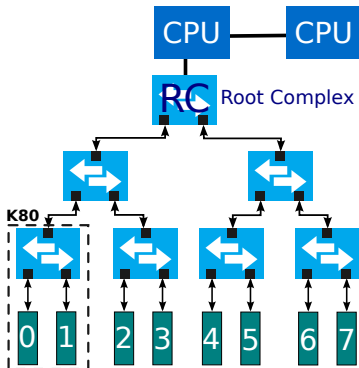
- generation 3, 16 GB/s using x16 wide lane
- dual simplex (a pair of unidirectional links)
- exchange buffer availability between pair of ports of a link
- tree-based topology

Building a densely populated accelerator servers with PCIe:



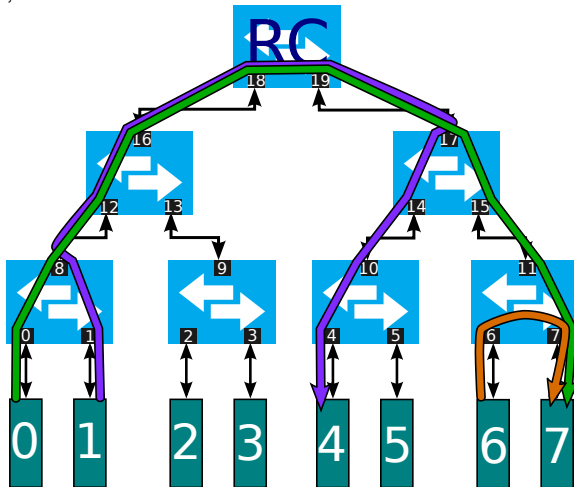
- generation 3, 16 GB/s using x16 wide lane
- dual simplex (a pair of unidirectional links)
- exchange buffer availability between pair of ports of a link
- tree-based topology

Building a densely populated accelerator servers with PCIe:



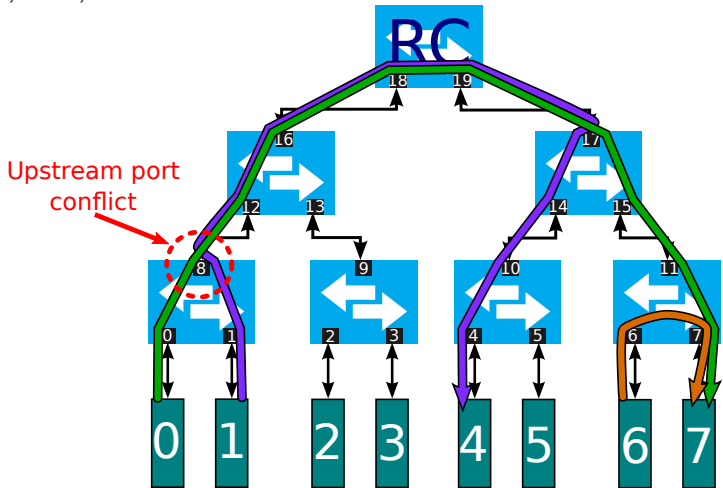
Communication conflicts

0→7, 1→4, 6→7



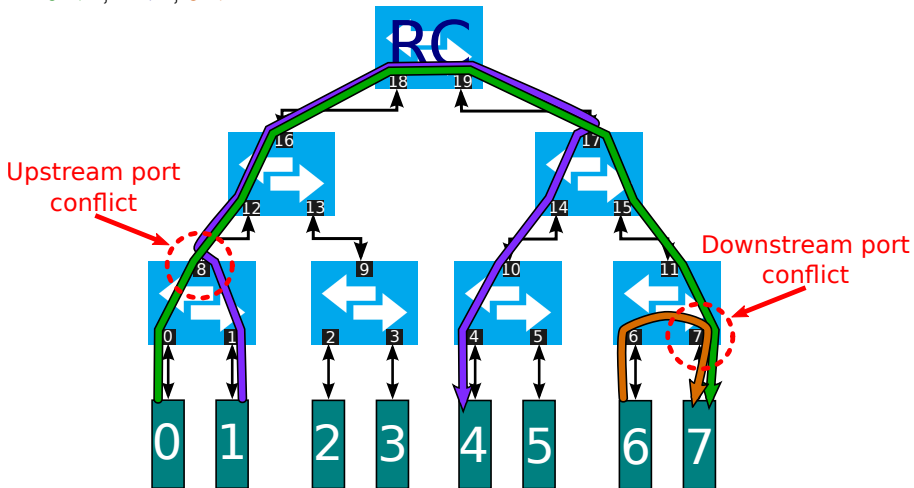
Communication conflicts

0→7, 1→4, 6→7



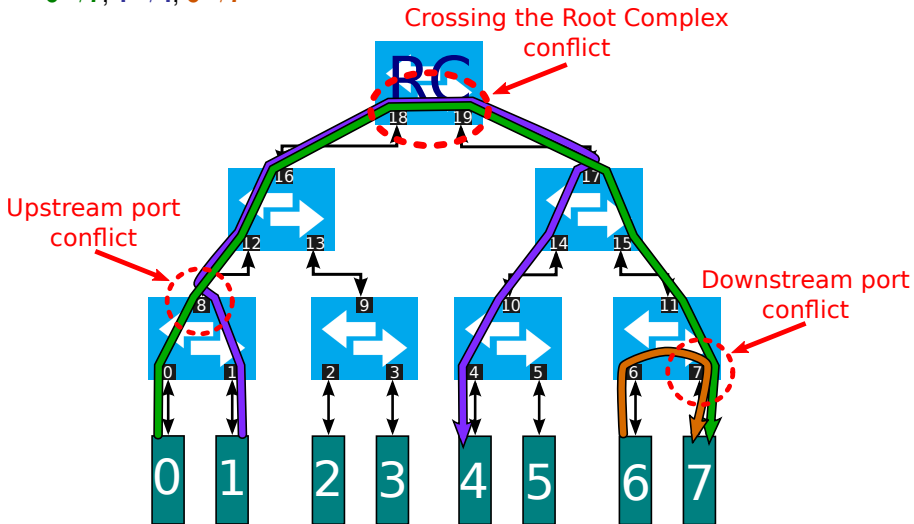
Communication conflicts

0→7, 1→4, 6→7



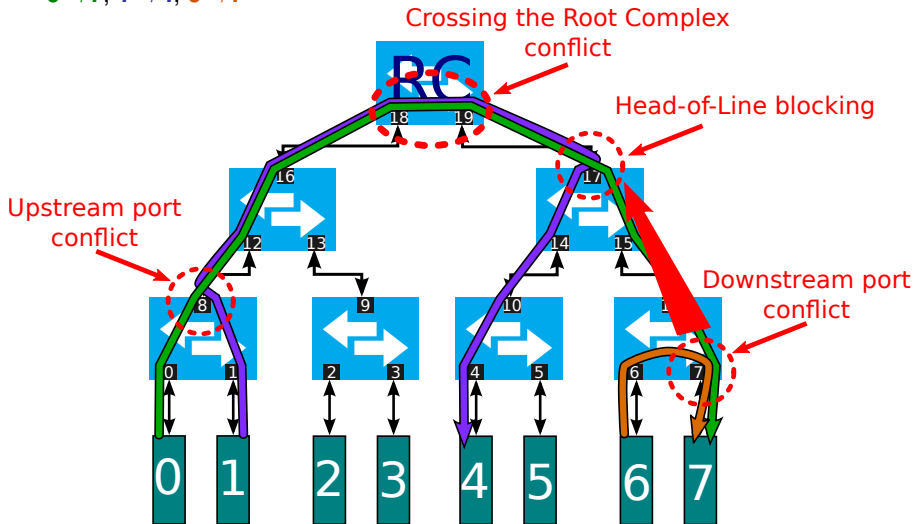
Communication conflicts

0→7, 1→4, 6→7

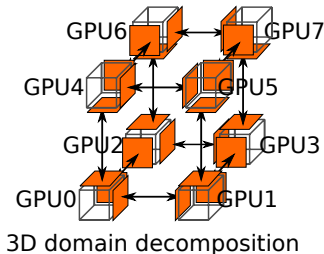
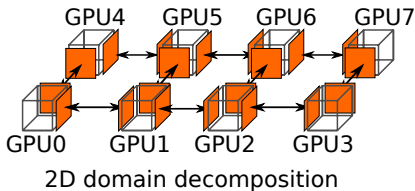


Communication conflicts

0→7, 1→4, 6→7



Motivation – COSMO halo exchange



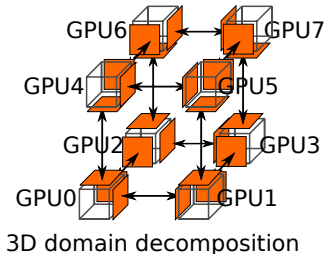
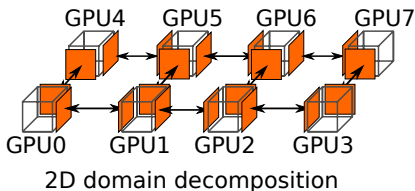
Which order of communications is the fastest?

0→4 1→0 2→1 3→7 4→0 5→4 6→2 7→6
0→1 1→2 2→6 3→2 4→5 5→6 6→7 7→3
1→5 2→3 5→1 6→5

0→1 1→0 2→1 3→2 4→0 5→1 6→2 7→3
0→4 1→5 2→6 3→7 4→5 5→6 6→5 7→6
1→2 2→3 5→4 6→7

...

Motivation – COSMO halo exchange



Which order of communications is the fastest?

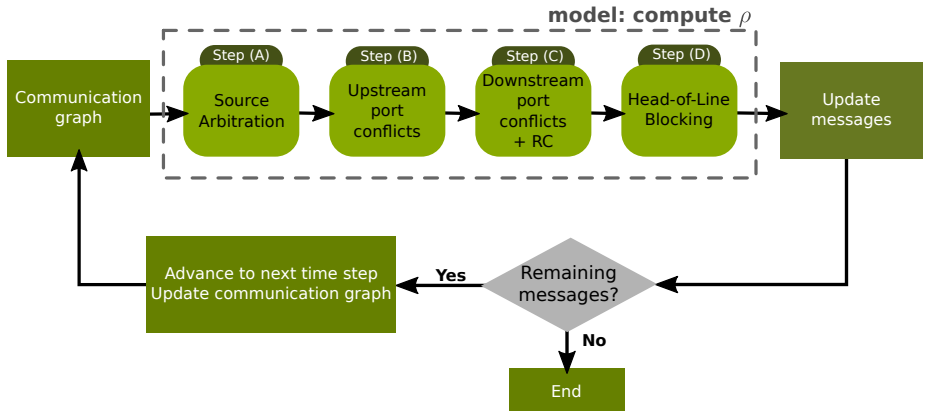
<p>0→4 1→0 2→1 3→7 4→0 5→4 6→2 7→6 0→1 1→2 2→6 3→2 4→5 5→6 6→7 7→3 1→5 2→3 5→1 6→5</p>	<p>0→1 1→0 2→1 3→2 4→0 5→1 6→2 7→3 0→4 1→5 2→6 3→7 4→5 5→6 6→5 7→6 1→2 2→3 5→4 6→7</p>	...
---	---	-----

2D domain decomposition example: 20,376 possibilities

3D domain decomposition has more than 1.6 Million possibilities

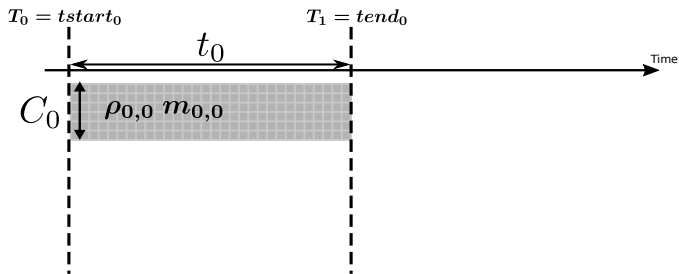
PCIe performance model

We want to identify the congestion factors $\rho \in [0, 1]$ which limit the available bandwidth per communication at each communication phase.



Communication phase – update messages

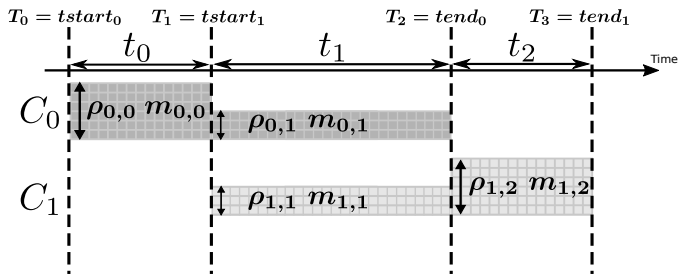
Elapsed time L_C , message size M_C , set of communication phases S_C :



$$L_{C_0} = t_0 = \frac{1}{B} \cdot \frac{m_{0,0}}{\rho_{0,0}} \text{ and } M_{C_0} = m_{0,0}$$

Communication phase – update messages

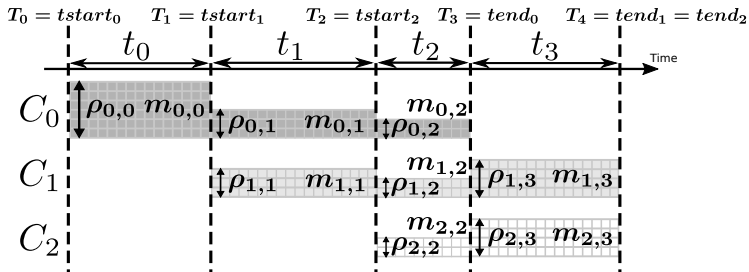
Elapsed time L_C , message size M_C , set of communication phases S_C :



$$L_{C_0} = t_0 + t_1 = \frac{1}{B} \cdot \left(\frac{m_{0,0}}{\rho_{0,0}} + \frac{m_{0,1}}{\rho_{0,1}} \right) \text{ and } M_{C_0} = m_{0,0} + m_{0,1}$$

Communication phase – update messages

Elapsed time L_C , message size M_C , set of communication phases S_C :



$$L_C = \sum_{i \in S_C} t_i = \frac{1}{B} \sum_{i \in S_C} \frac{m_{c,i}}{\rho_{c,i}} \text{ and } M_C = \sum_{i \in S_C} m_{c,i}$$

- (1) start time and M_C are known
 - (2) $\rho_{c,i}$ are given by the model
- with (1) and (2) L_C are computable

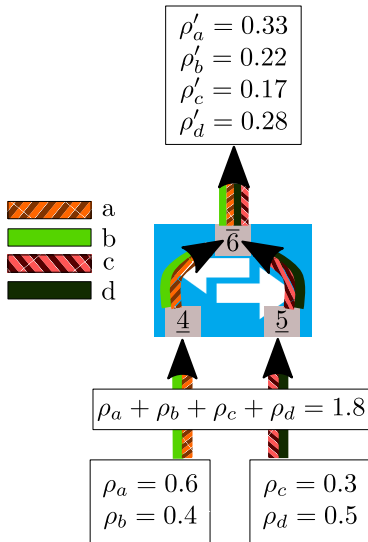
Model conflicts on switch

We want to identify the congestion factors $\rho \in [0, 1]$ which limit the available bandwidth per communication at each communication phase.

- Each communication enters a switch with a congestion factor ρ and leaves with a congestion factor ρ'
- If $\sum_i \rho_i > 1$ then an arbitration policy is required, $\rho' = \rho$ otherwise

Upstream port conflict

Proportional sharing of available bandwidth



Downstream port conflict

- Round-robin policy
- Performance reduction for crossing the root complex

C^R – set of communications crossing the root complex;

n – number of grouped communication sets;

R – congestion factor of a grouped communication set;

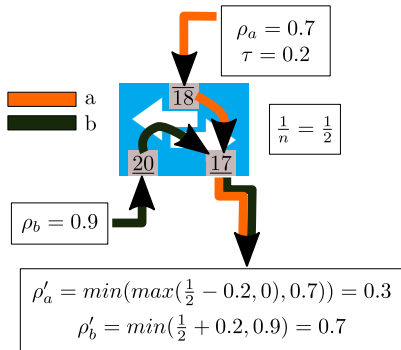
τ – congestion factor for crossing the root complex;

– if $C^R = \emptyset$ then

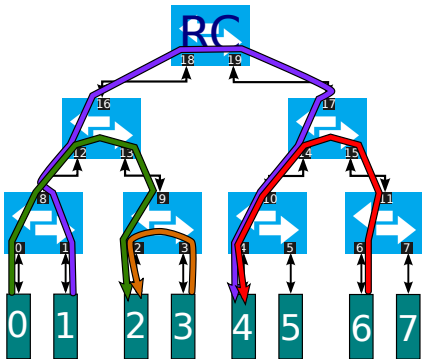
$$R' = \frac{1}{n}$$

– if $C^R \neq \emptyset$ then

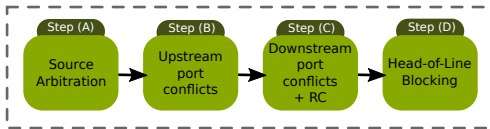
$$R' = \begin{cases} \min(\max(\frac{1}{n} - \tau, 0), R) & \text{if } R \text{ contains comm. } \in C^R \\ \min(\frac{1}{n} + \tau, R) & \text{otherwise} \end{cases}$$



Complete example

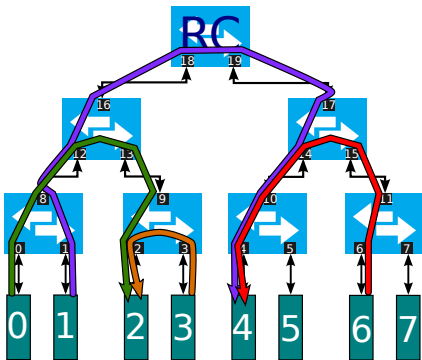


$\tau = 0.2$



comm.	Step (A)	Step (B)	Step (C)	Step (D)
(a) 0 → 2	1	1/2	1/2	3/10 3/10
(b) 1 → 4	1	1/2	3/10	3/10 3/10
(c) 3 → 2	1	1	1/2	1/2 7/10
(d) 6 → 4	1	1	7/10	7/10 7/10

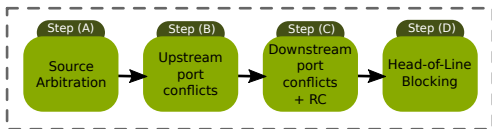
Complete example



$$\tau = 0.2$$

Message size: 300MB

Bandwidth: 11.6 GB/s

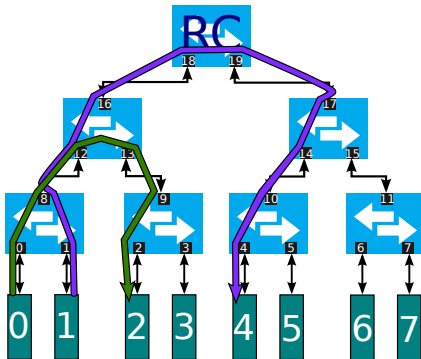


comm.	Step (A)	Step (B)	Step (C)	Step (D)	
(a) 0 → 2	1	1/2	1/2	3/10	3/10
(b) 1 → 4	1	1/2	3/10	3/10	3/10
(c) 3 → 2	1	1	1/2	1/2	7/10
(d) 6 → 4	1	1	7/10	7/10	7/10

Congestion graph step 1

comm.	cong. factor	data remaining	elapsed time
(a) 0 → 2	3/10	128 MB	36 ms
(b) 1 → 4	3/10	128 MB	36 ms
(c) 3 → 2	7/10	0 MB	36 ms
(d) 6 → 4	7/10	0 MB	36 ms

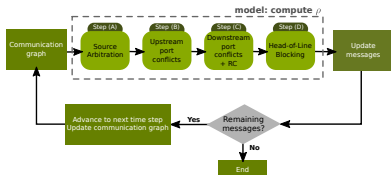
Complete example



$$\tau = 0.2$$

Message size: 300MB

Bandwidth: 11.6 GB/s



Congestion graph step 1

comm.	cong. factor	data remaining	elapsed time
(a) 0→2	3/10	128 MB	36 ms
(b) 1→4	3/10	128 MB	36 ms
(c) 3→2	7/10	0 MB	36 ms
(d) 6→4	7/10	0 MB	36 ms

Congestion graph step 2

comm.	cong. factor	data remaining	elapsed time
(a) 0→2	1/2	0 MB	65 ms
(b) 1→4	1/2	0 MB	65 ms

Model Validation

- Architecture parameters:

$$B = 11.6\text{GB/s}$$

$$\tau = 0.1735$$

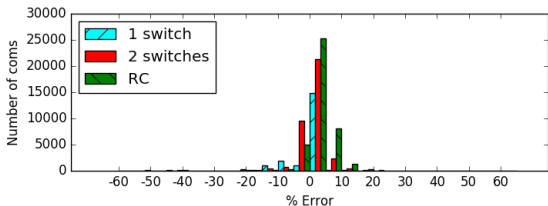
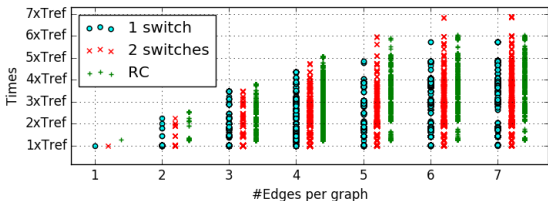
- 22,259 graphs:

- non-isomorphic
- cudaMemcpyAsync*
- Communication pattern: scatter, gather, all-to-all
- Entire set of graphs for subsets of GPUs
- Randomly generated
- $\approx 100\text{K}$ communications

- Message size: 300 MB

- Time no contention:

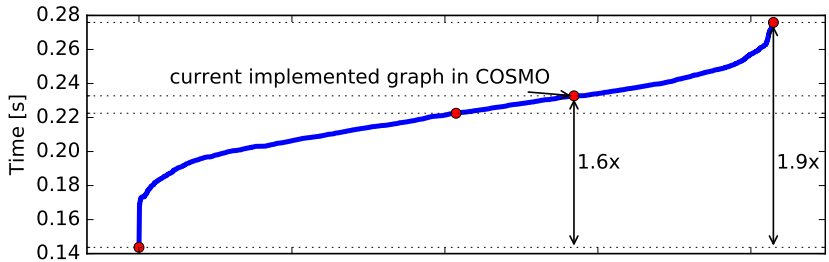
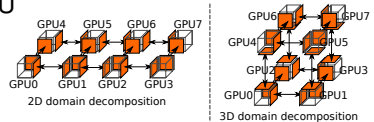
$$T_{ref} = 25.3\text{ms}$$



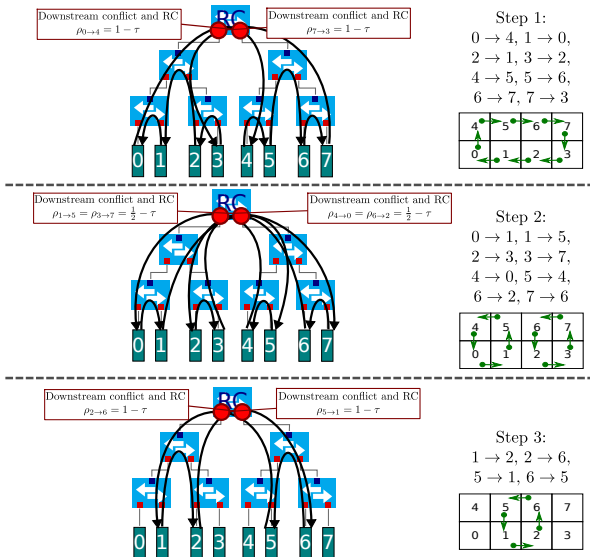
95% of communication are in range +/- 15%

Back to the motivation – COSMO halo exchange

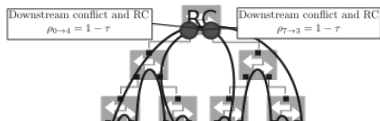
- Upper limit on time to solution, throughput approach
- Running mode one instance per socket (8 GPUs)
- Large domain size 256x256x80 per GPU
- One step triggers 312 halo exchanges
- Message size: 40 KB to 254 KB
- Uses MPI
- $(3!)^4 \times (2!)^4 = 20,736$ communication graphs for 2D domain



Fastest schedule for 2D decomposition

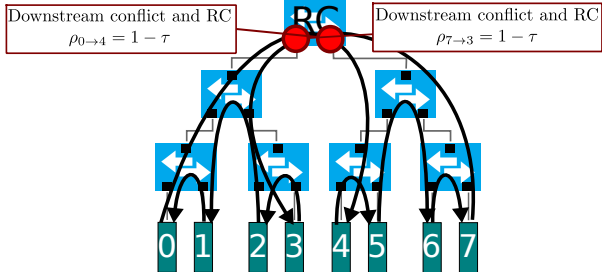


Fastest schedule for 2D decomposition



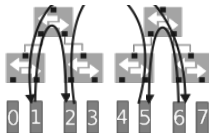
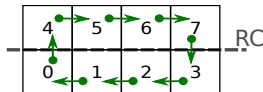
Step 1:

$0 \rightarrow 4, 1 \rightarrow 0,$
 $2 \rightarrow 1, 3 \rightarrow 2,$
 $4 \rightarrow 5, 5 \rightarrow 6,$
 $6 \rightarrow 7, 7 \rightarrow 3$



Step 1:

$0 \rightarrow 4, 1 \rightarrow 0,$
 $2 \rightarrow 1, 3 \rightarrow 2,$
 $4 \rightarrow 5, 5 \rightarrow 6,$
 $6 \rightarrow 7, 7 \rightarrow 3$

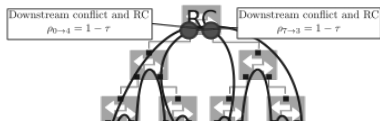


Step 3:

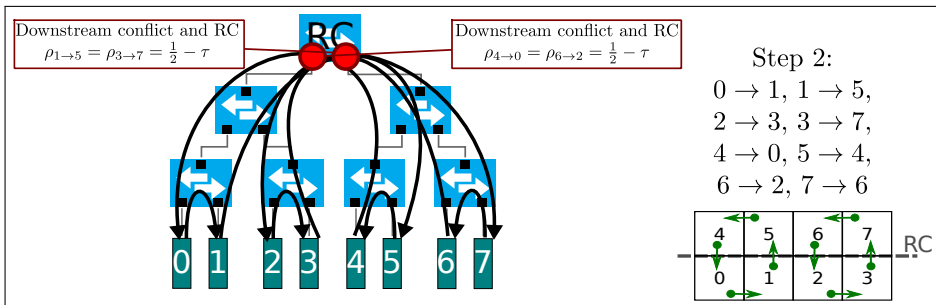
$1 \rightarrow 2, 2 \rightarrow 6,$
 $5 \rightarrow 1, 6 \rightarrow 5$



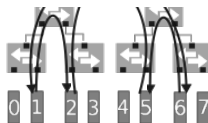
Fastest schedule for 2D decomposition



Step 1:
 $0 \rightarrow 4, 1 \rightarrow 0,$
 $2 \rightarrow 1, 3 \rightarrow 2,$
 $4 \rightarrow 5, 5 \rightarrow 6,$
 $6 \rightarrow 7, 7 \rightarrow 3$



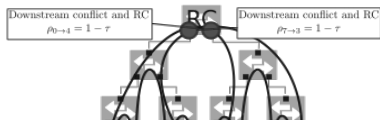
Step 2:
 $0 \rightarrow 1, 1 \rightarrow 5,$
 $2 \rightarrow 3, 3 \rightarrow 7,$
 $4 \rightarrow 0, 5 \rightarrow 4,$
 $6 \rightarrow 2, 7 \rightarrow 6$



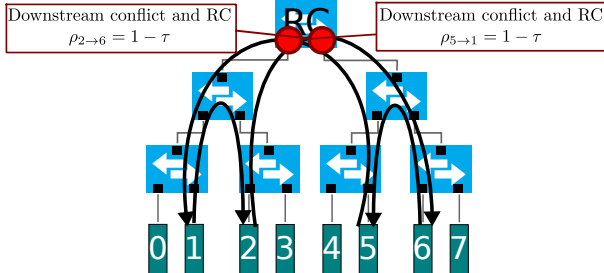
1 \rightarrow 2, 2 \rightarrow 6,
 5 \rightarrow 1, 6 \rightarrow 5

4	5	6	7	RC
0	1	2	3	

Fastest schedule for 2D decomposition

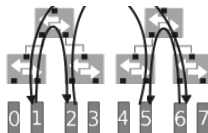
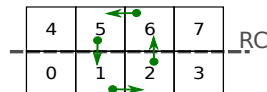


Step 1:
 $0 \rightarrow 4, 1 \rightarrow 0,$
 $2 \rightarrow 1, 3 \rightarrow 2,$
 $4 \rightarrow 5, 5 \rightarrow 6,$
 $6 \rightarrow 7, 7 \rightarrow 3$



Step 3:

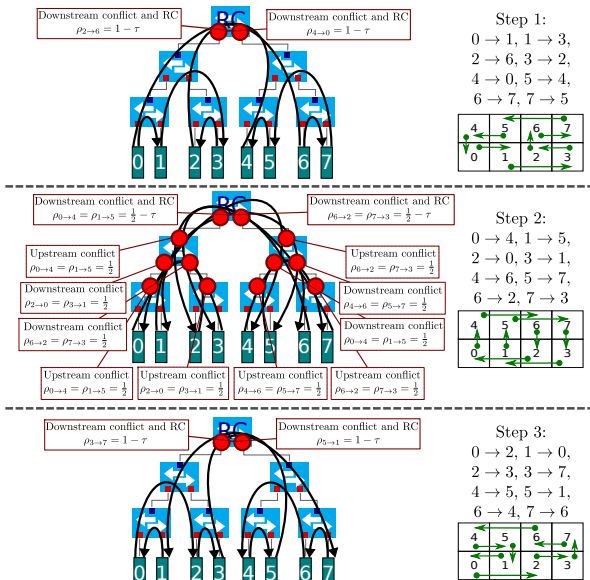
$1 \rightarrow 2, 2 \rightarrow 6,$
 $5 \rightarrow 1, 6 \rightarrow 5$



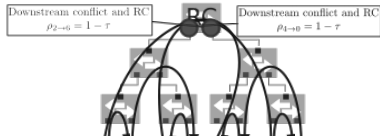
Step 3:
 $1 \rightarrow 2, 2 \rightarrow 6,$
 $5 \rightarrow 1, 6 \rightarrow 5$



Fastest schedule for 3D decomposition



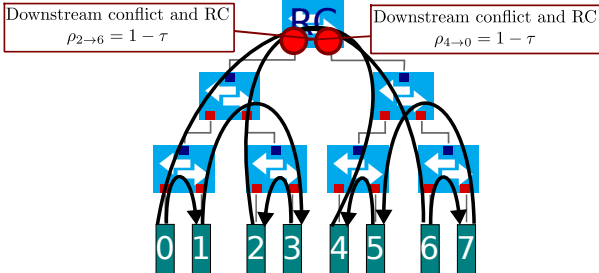
Fastest schedule for 3D decomposition



Step 1:

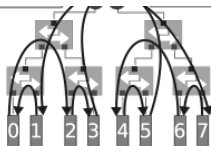
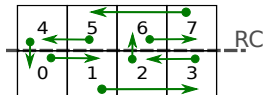
$0 \rightarrow 1, 1 \rightarrow 3,$
 $2 \rightarrow 6, 3 \rightarrow 2,$
 $4 \rightarrow 0, 5 \rightarrow 4,$
 $6 \rightarrow 7, 7 \rightarrow 5$

4	5	6	7
---	---	---	---



Step 1:

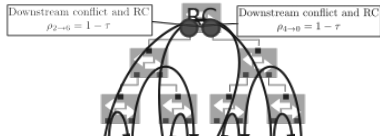
$0 \rightarrow 1, 1 \rightarrow 3,$
 $2 \rightarrow 6, 3 \rightarrow 2,$
 $4 \rightarrow 0, 5 \rightarrow 4,$
 $6 \rightarrow 7, 7 \rightarrow 5$



$0 \rightarrow 2, 1 \rightarrow 0,$
 $2 \rightarrow 3, 3 \rightarrow 7,$
 $4 \rightarrow 5, 5 \rightarrow 1,$
 $6 \rightarrow 4, 7 \rightarrow 6$

4	5	6	7
0	1	2	3

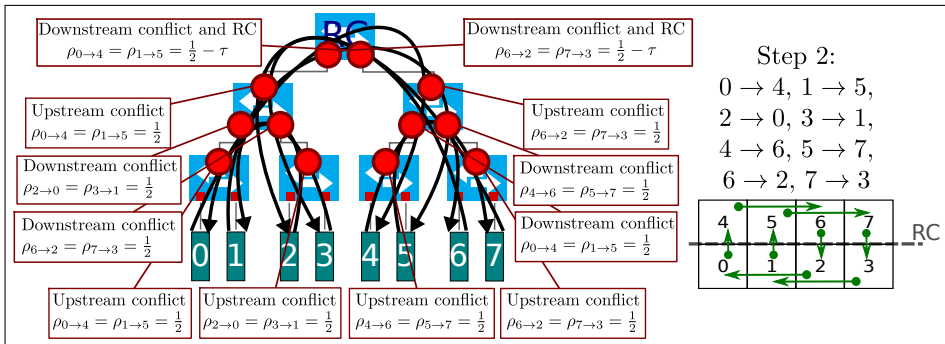
Fastest schedule for 3D decomposition



Step 1:

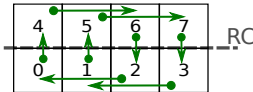
0 → 1, 1 → 3,
2 → 6, 3 → 2,
4 → 0, 5 → 4,
6 → 7, 7 → 5

4	5	6	7
---	---	---	---



Step 2:

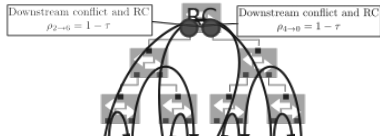
0 → 4, 1 → 5,
2 → 0, 3 → 1,
4 → 6, 5 → 7,
6 → 2, 7 → 3



6 → 4, 7 → 6

4	5	6	7
0	1	2	3

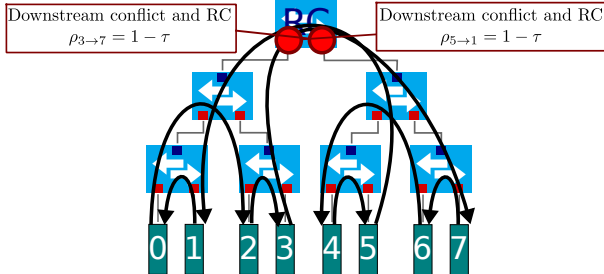
Fastest schedule for 3D decomposition



Step 1:

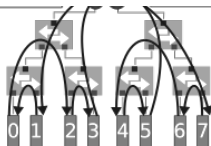
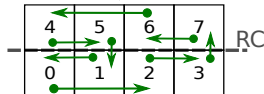
$0 \rightarrow 1, 1 \rightarrow 3,$
 $2 \rightarrow 6, 3 \rightarrow 2,$
 $4 \rightarrow 0, 5 \rightarrow 4,$
 $6 \rightarrow 7, 7 \rightarrow 5$

4	5	6	7
---	---	---	---



Step 3:

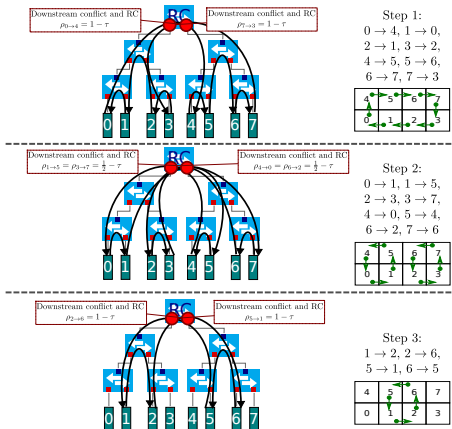
$0 \rightarrow 2, 1 \rightarrow 0,$
 $2 \rightarrow 3, 3 \rightarrow 7,$
 $4 \rightarrow 5, 5 \rightarrow 1,$
 $6 \rightarrow 4, 7 \rightarrow 6$



$0 \rightarrow 2, 1 \rightarrow 0,$
 $2 \rightarrow 3, 3 \rightarrow 7,$
 $4 \rightarrow 5, 5 \rightarrow 1,$
 $6 \rightarrow 4, 7 \rightarrow 6$

4	5	6	7
0	1	2	3

COSMO improvement – fastest schedule



COSMO gain: 5.6% per halo exchange step,
gain is limited by MPI 2-sided overhead.

Conclusion

- Latency not modeled
- MPI 2-sided overhead not modeled (use one-sided?)
- + Captures all PCIe features including congestion
- + Simple model only 2 parameters (B and τ)
- + Precise for large messages
- + Design of topology-aware algorithms
- + COSMO halo exchange performance gain

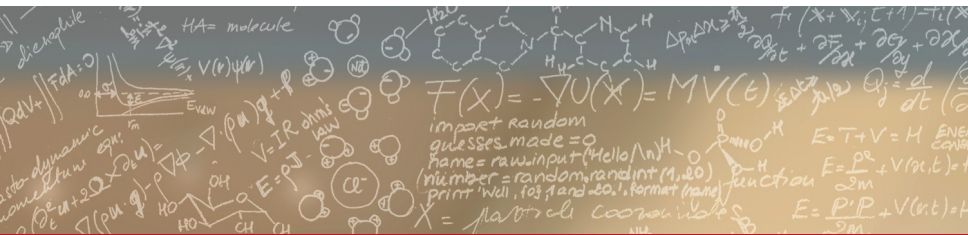


CSCS

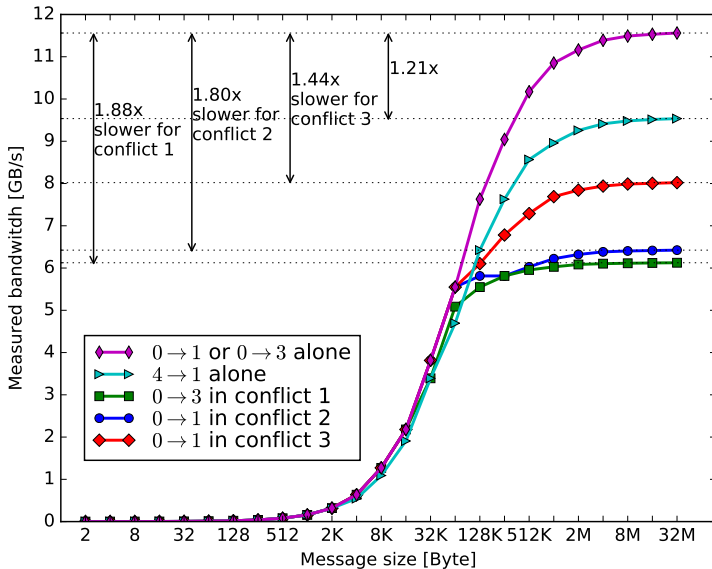
Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre



ETH zürich



Thank you for your attention.



$$\tau = 1 - 1/1.21$$

Legend:

