

# 25+ years of scheduling at ICL

Yves Robert

ENS Lyon  
University of Tennessee Knoxville

Dongarra 70 — July 7, 2021

# Overview

- A few results in parallel linear algebra and resilience that were achieved at ICL under Jack's guidance
- The story started in 1996. It continues today

# Jack's co-authors according to DBLP

- 1 Stanimire Tomov (149)
- 2 Piotr Luszczek (119)
- 3 George Bosilca (91)
- 4 Jakub Kurzak (76)
- 5 Azzam Haidar (69)
- 6 Peter M. A. Sloot (62)
- 7 Thomas Héault (51)
- 8 Hartwig Anzt (51)
- 9 Ichitaro Yamazaki (51)
- 10 Julien Langou (43)
- 11 Aurelien Bouteiller (43)
- 12 Yves Robert (41)
- 13 Graham E. Fagg (40)
- 14 Mark Gates (32)
- 15 Ahmad Abdelfattah (32)
- 16 Hatem Ltaief (31)
- 17 ...

Redistribution  
oooooooo

Linear algebra kernels  
oooooooooooooooooooo

Resilience  
oo

Life at ICL  
oooooooo

# Outline

1 Redistribution

2 Linear algebra kernels

3 Resilience

4 Life at ICL

# Outline

1 Redistribution

2 Linear algebra kernels

3 Resilience

4 Life at ICL

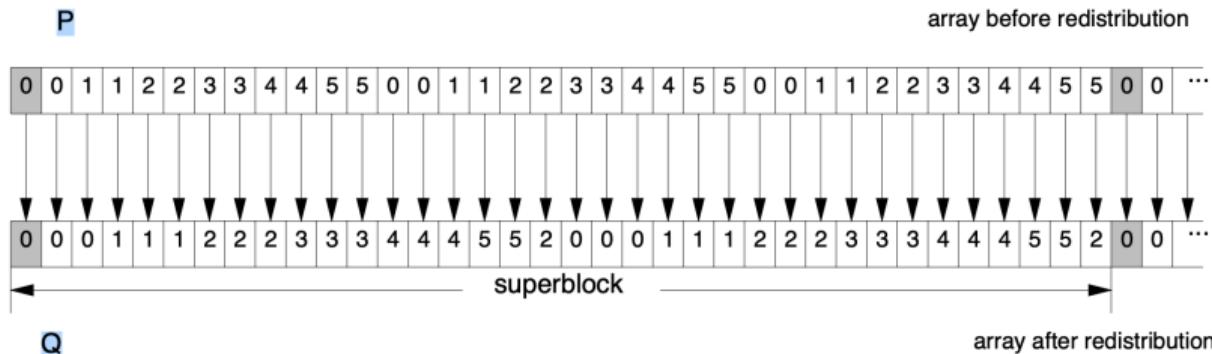
# 1996-2000 Redistribution algorithms

$A_{11}$	$A_{12}$	$A_{13}$	$A_{14}$	$A_{15}$	$A_{16}$	$A_{17}$	$A_{18}$
$A_{21}$	$A_{22}$	$A_{23}$	$A_{24}$	$A_{25}$	$A_{26}$	$A_{27}$	$A_{28}$
$A_{31}$	$A_{32}$	$A_{33}$	$A_{34}$	$A_{35}$	$A_{36}$	$A_{37}$	$A_{38}$
$A_{41}$	$A_{42}$	$A_{43}$	$A_{44}$	$A_{45}$	$A_{46}$	$A_{47}$	$A_{48}$
$A_{51}$	$A_{52}$	$A_{53}$	$A_{54}$	$A_{55}$	$A_{56}$	$A_{57}$	$A_{58}$
$A_{61}$	$A_{62}$	$A_{63}$	$A_{64}$	$A_{65}$	$A_{66}$	$A_{67}$	$A_{68}$
$A_{71}$	$A_{72}$	$A_{73}$	$A_{74}$	$A_{75}$	$A_{76}$	$A_{77}$	$A_{78}$
$A_{81}$	$A_{82}$	$A_{83}$	$A_{84}$	$A_{85}$	$A_{86}$	$A_{87}$	$A_{88}$

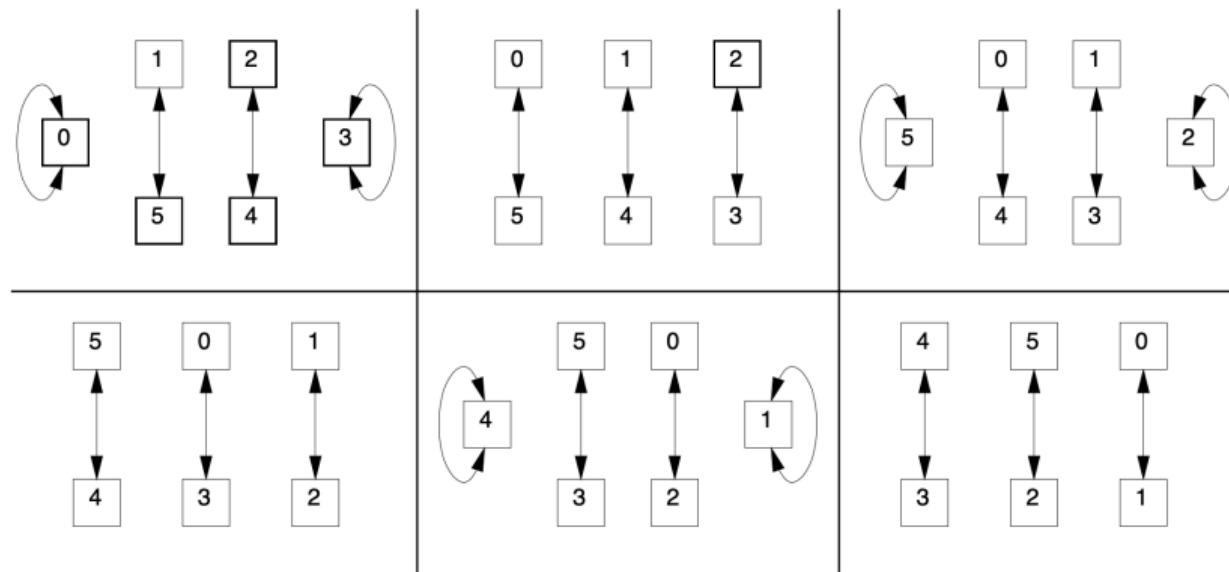
	0	1	2					
0	$A_{11}$	$A_{14}$	$A_{17}$	$A_{12}$	$A_{15}$	$A_{18}$	$A_{13}$	$A_{16}$
1	$A_{31}$	$A_{34}$	$A_{37}$	$A_{32}$	$A_{35}$	$A_{38}$	$A_{33}$	$A_{36}$
2	$A_{51}$	$A_{54}$	$A_{57}$	$A_{52}$	$A_{55}$	$A_{58}$	$A_{53}$	$A_{56}$
3	$A_{71}$	$A_{74}$	$A_{77}$	$A_{72}$	$A_{75}$	$A_{78}$	$A_{73}$	$A_{76}$
4	$A_{21}$	$A_{24}$	$A_{27}$	$A_{22}$	$A_{25}$	$A_{28}$	$A_{23}$	$A_{26}$
5	$A_{41}$	$A_{44}$	$A_{47}$	$A_{42}$	$A_{45}$	$A_{48}$	$A_{43}$	$A_{46}$
6	$A_{61}$	$A_{64}$	$A_{67}$	$A_{62}$	$A_{65}$	$A_{68}$	$A_{63}$	$A_{66}$
7	$A_{81}$	$A_{84}$	$A_{87}$	$A_{82}$	$A_{85}$	$A_{88}$	$A_{83}$	$A_{86}$

Global (left) and distributed (right) views of matrix

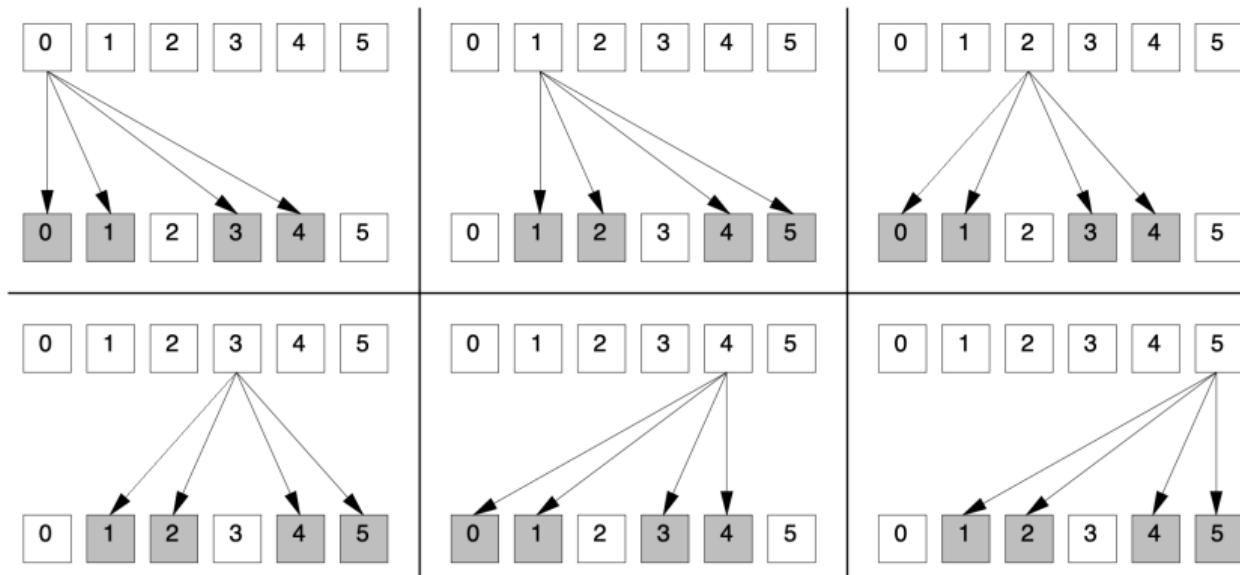
## From CYCLIC(2) to CYCLIC(3)



# ScalAPACK caterpillar



## From CYCLIC(2) to CYCLIC(3)

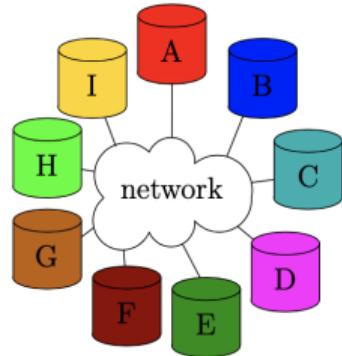


## From CYCLIC(2) to CYCLIC(3)

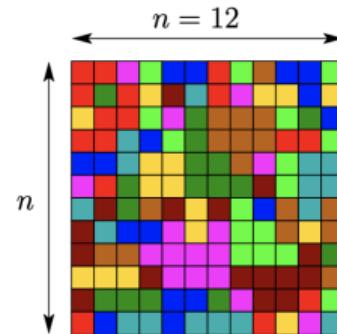
Sender/Recv.	0	1	2	3	4	5	Nbr of msg.
0	2	-	2	-	2	-	3
1	1	1	1	1	1	1	6
2	-	2	-	2	-	2	3
3	2	-	2	-	2	-	3
4	1	1	1	1	1	1	6
5	-	2	-	2	-	2	3
Nbr of msg.	4	4	4	4	4	4	

Message lengths are for a single-slice vector of  $L = \text{lcm}(Pr, Qs) = 36$  components

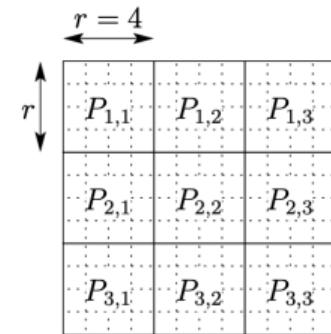
# Redistribution followed by a computational kernel



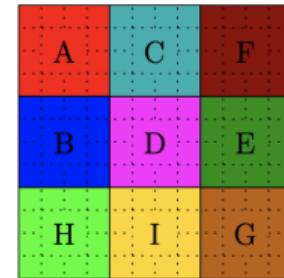
(a) processors holding the data



(b) initial data distribution



(c) target data partition

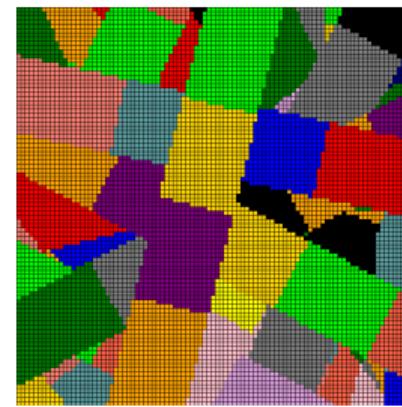


(d) final data distribution

Each color in the data distributions corresponds to a processor,  
e.g., all red data items reside on processor A.

# Redistribution followed by a computational kernel

- Data elements stored on different processors
- Computation kernel must be applied to data
- Initial data distribution may be inefficient for the computation kernel



Distribution of a tiled matrix

# Outline

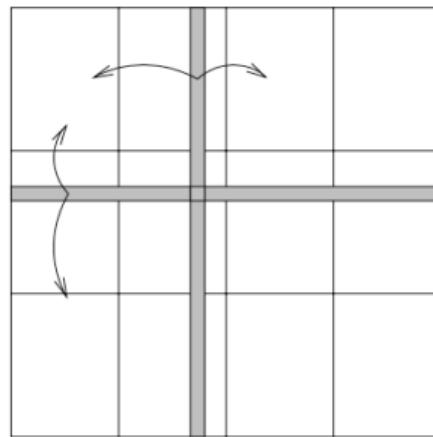
1 Redistribution

2 Linear algebra kernels

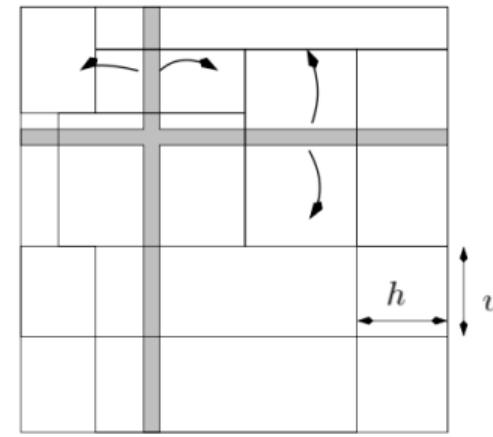
3 Resilience

4 Life at ICL

# 1999-2003 Matrix product on heterogeneous platforms

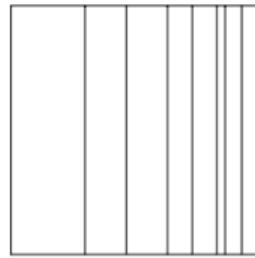


Homogeneous  $3 \times 4$  grid

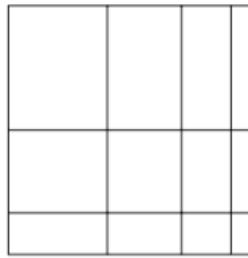


Heterogeneous  $3 \times 4$  grid

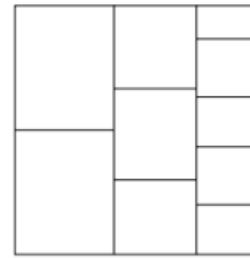
# Matrix product on heterogeneous platforms



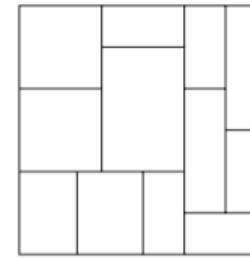
1D



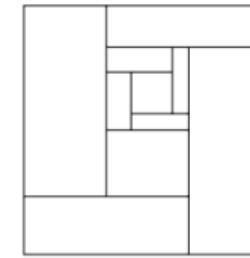
2D



Column

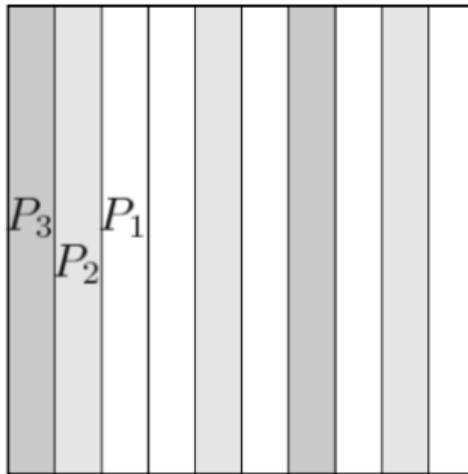


Recursive



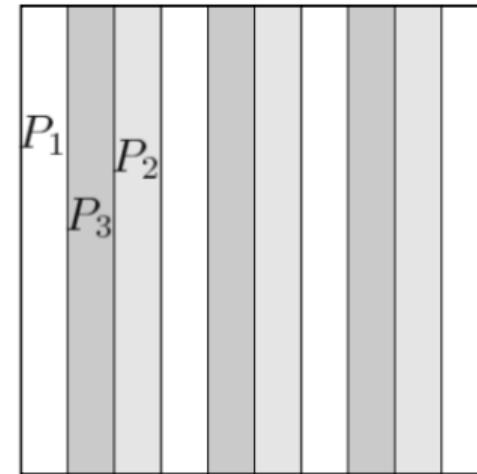
Unconstrained

# LU decomposition on heterogeneous platforms



*Heterogeneous allocation*

Cycle-times:  $c(P_1) = 3$ ,  $c(P_2) = 5$ ,  $c(P_3) = 8$



*Block-cyclic allocation*

# 2008 Matrix product on master-worker platforms

	$\mathcal{B}_{11}$	$\mathcal{B}_{12}$	$\mathcal{B}_{13}$	$\mathcal{B}_{14}$
$\mathcal{A}_{11}$	$\mathcal{C}_{11}$	$\mathcal{C}_{12}$	$\mathcal{C}_{13}$	$\mathcal{C}_{14}$
	$\mathcal{C}_{21}$	$\mathcal{C}_{22}$	$\mathcal{C}_{23}$	$\mathcal{C}_{24}$
	$\mathcal{C}_{31}$	$\mathcal{C}_{32}$	$\mathcal{C}_{33}$	$\mathcal{C}_{34}$
	$\mathcal{C}_{41}$	$\mathcal{C}_{42}$	$\mathcal{C}_{43}$	$\mathcal{C}_{44}$

	$\mathcal{B}_{11}$	$\mathcal{B}_{12}$	$\mathcal{B}_{13}$	$\mathcal{B}_{14}$
	$\mathcal{C}_{11}$	$\mathcal{C}_{12}$	$\mathcal{C}_{13}$	$\mathcal{C}_{14}$
$\mathcal{A}_{21}$	$\mathcal{C}_{21}$	$\mathcal{C}_{22}$	$\mathcal{C}_{23}$	$\mathcal{C}_{24}$
	$\mathcal{C}_{31}$	$\mathcal{C}_{32}$	$\mathcal{C}_{33}$	$\mathcal{C}_{34}$
	$\mathcal{C}_{41}$	$\mathcal{C}_{42}$	$\mathcal{C}_{43}$	$\mathcal{C}_{44}$

	$\mathcal{B}_{11}$	$\mathcal{B}_{12}$	$\mathcal{B}_{13}$	$\mathcal{B}_{14}$
	$\mathcal{C}_{11}$	$\mathcal{C}_{12}$	$\mathcal{C}_{13}$	$\mathcal{C}_{14}$
$\mathcal{A}_{31}$	$\mathcal{C}_{31}$	$\mathcal{C}_{32}$	$\mathcal{C}_{33}$	$\mathcal{C}_{34}$
	$\mathcal{C}_{41}$	$\mathcal{C}_{42}$	$\mathcal{C}_{43}$	$\mathcal{C}_{44}$

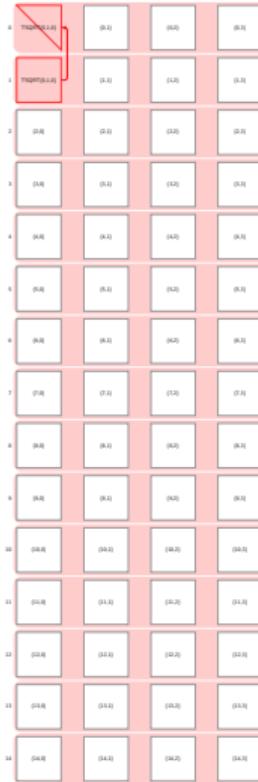
	$\mathcal{B}_{11}$	$\mathcal{B}_{12}$	$\mathcal{B}_{13}$	$\mathcal{B}_{14}$
	$\mathcal{C}_{11}$	$\mathcal{C}_{12}$	$\mathcal{C}_{13}$	$\mathcal{C}_{14}$
$\mathcal{A}_{21}$	$\mathcal{C}_{21}$	$\mathcal{C}_{22}$	$\mathcal{C}_{23}$	$\mathcal{C}_{24}$
	$\mathcal{C}_{31}$	$\mathcal{C}_{32}$	$\mathcal{C}_{33}$	$\mathcal{C}_{34}$
	$\mathcal{C}_{41}$	$\mathcal{C}_{42}$	$\mathcal{C}_{43}$	$\mathcal{C}_{44}$

With  $1 + \mu + \mu^2$  blocks in memory, improving Toledo's algorithm

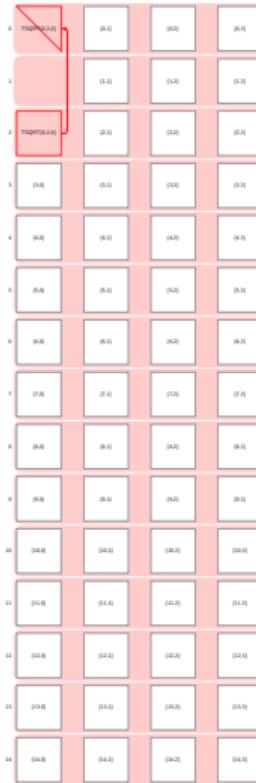
# 2012 QR factorization – Flat Tree in action



# 2012 QR factorization – Flat Tree in action



# 2012 QR factorization – Flat Tree in action



# 2012 QR factorization – Flat Tree in action



# 2012 QR factorization – Flat Tree in action



# 2012 QR factorization – Flat Tree in action



# 2012 QR factorization – Flat Tree in action



# 2012 QR factorization – Flat Tree in action



# 2012 QR factorization – Flat Tree in action



# 2012 QR factorization – Flat Tree in action



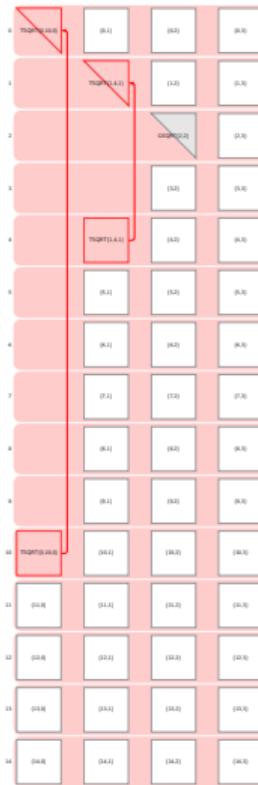
# 2012 QR factorization – Flat Tree in action



# 2012 QR factorization – Flat Tree in action



# 2012 QR factorization – Flat Tree in action



# 2012 QR factorization – Flat Tree in action



# 2012 QR factorization – Flat Tree in action



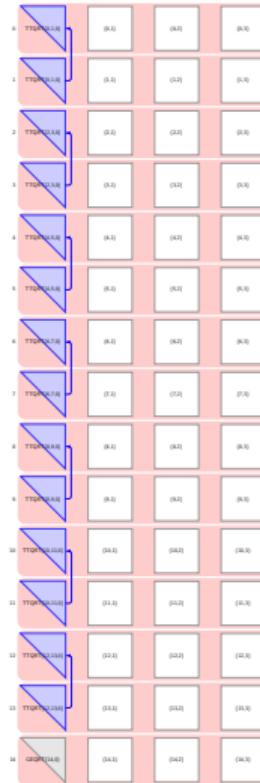
# 2012 QR factorization – Flat Tree in action



# 2012 QR factorization – Flat Tree in action



## QR factorization – Binary Tree in action



## QR factorization – Binary Tree in action



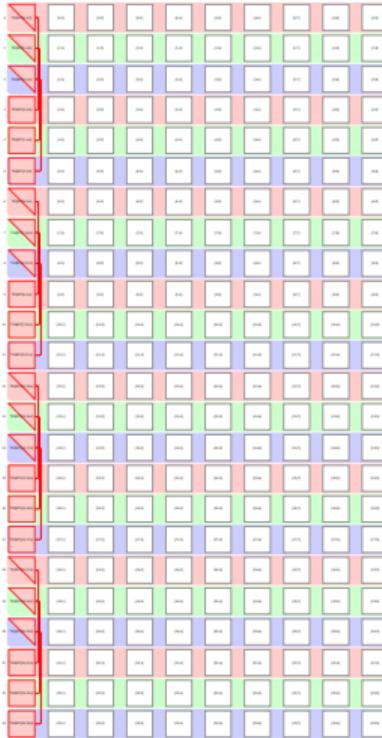
## QR factorization – Binary Tree in action



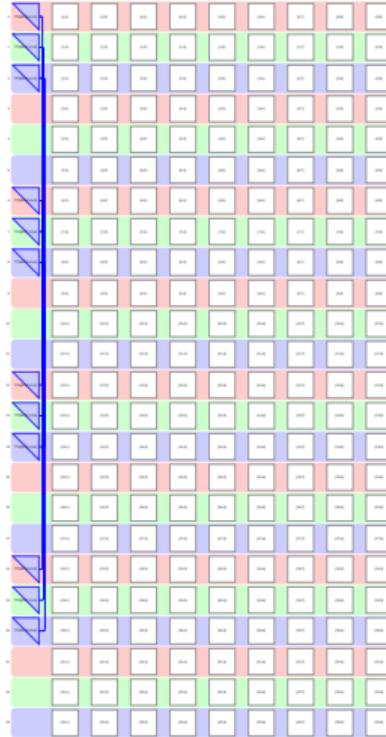
## QR factorization – Binary Tree in action



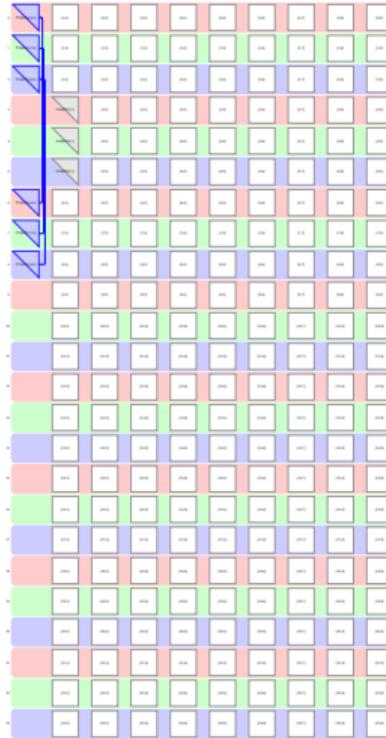
## QR factorization – Hierarchical algorithm in action



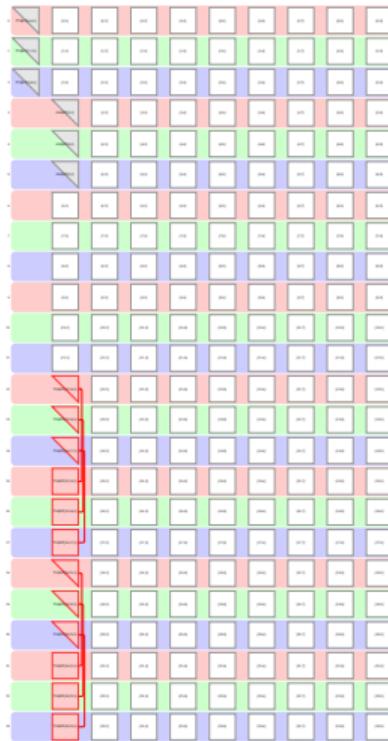
# QR factorization – Hierarchical algorithm in action



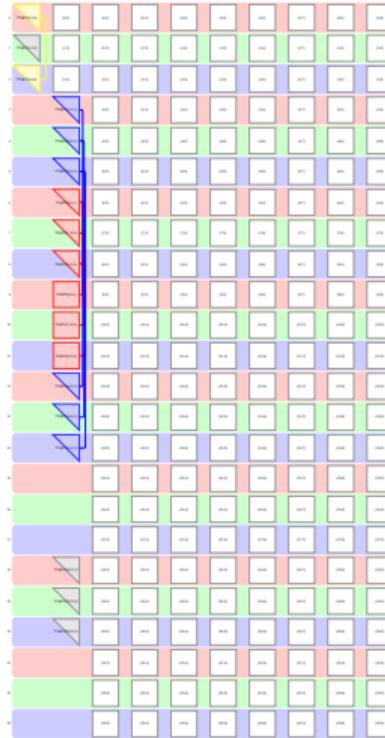
# QR factorization – Hierarchical algorithm in action



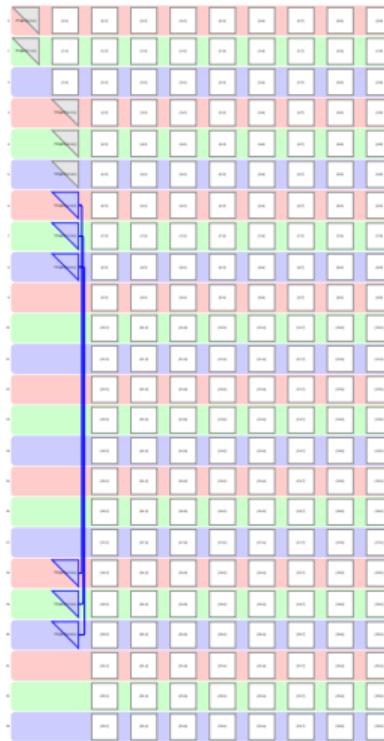
## QR factorization – Hierarchical algorithm in action



# QR factorization – Hierarchical algorithm in action



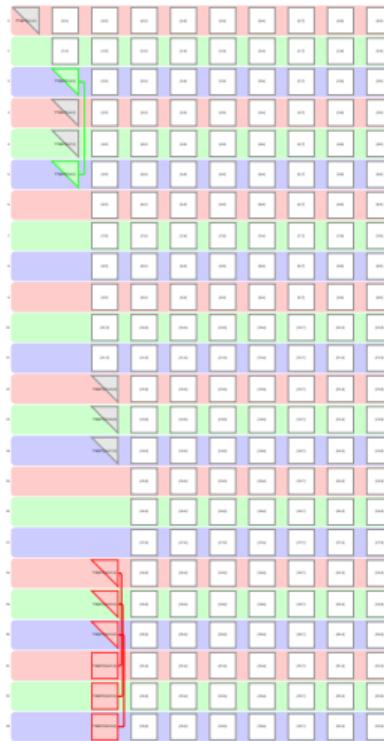
# QR factorization – Hierarchical algorithm in action



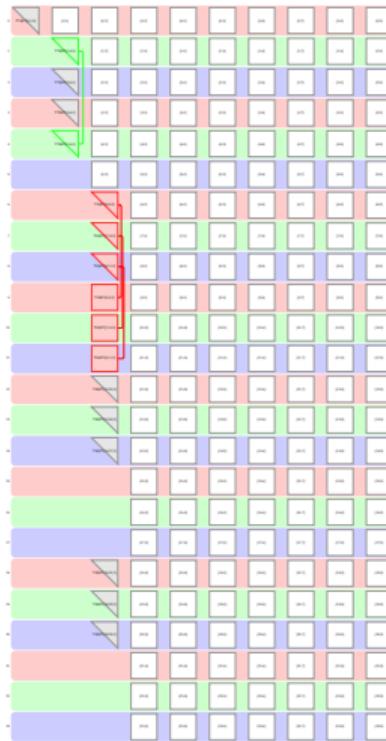
# QR factorization – Hierarchical algorithm in action



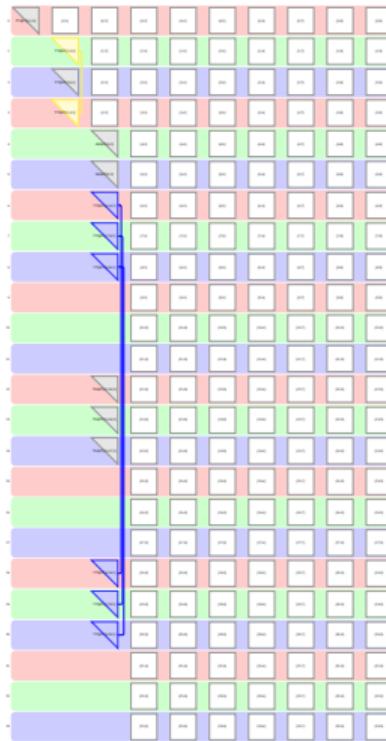
# QR factorization – Hierarchical algorithm in action



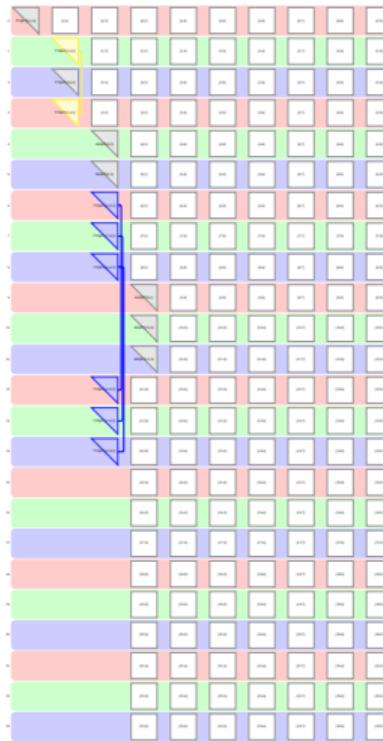
## QR factorization – Hierarchical algorithm in action



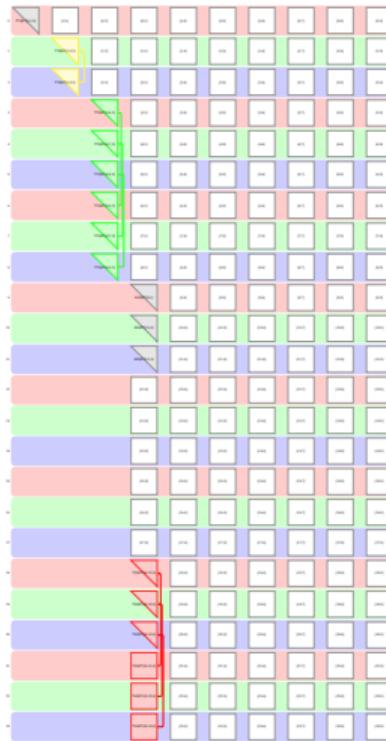
## QR factorization – Hierarchical algorithm in action



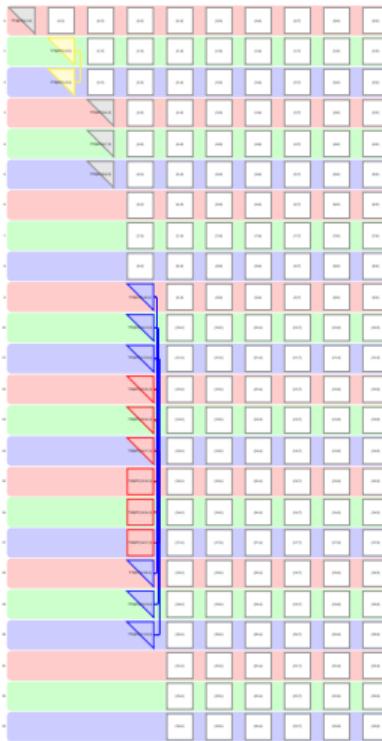
# QR factorization – Hierarchical algorithm in action



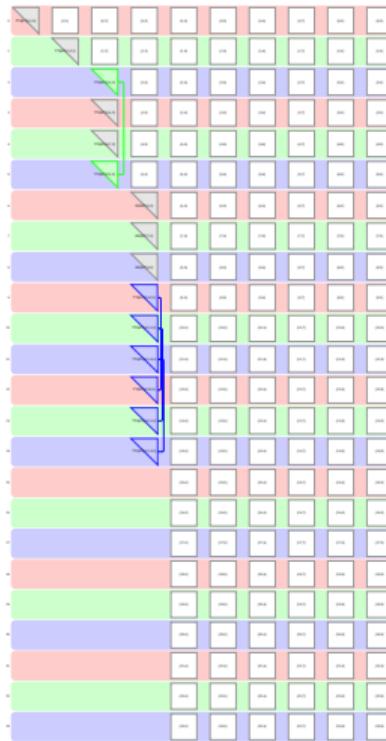
# QR factorization – Hierarchical algorithm in action



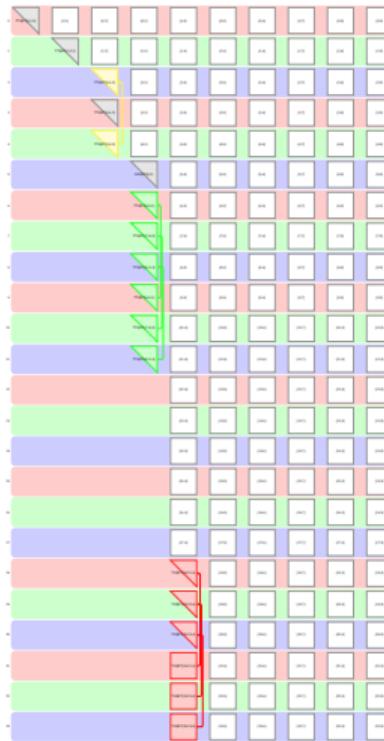
# QR factorization – Hierarchical algorithm in action



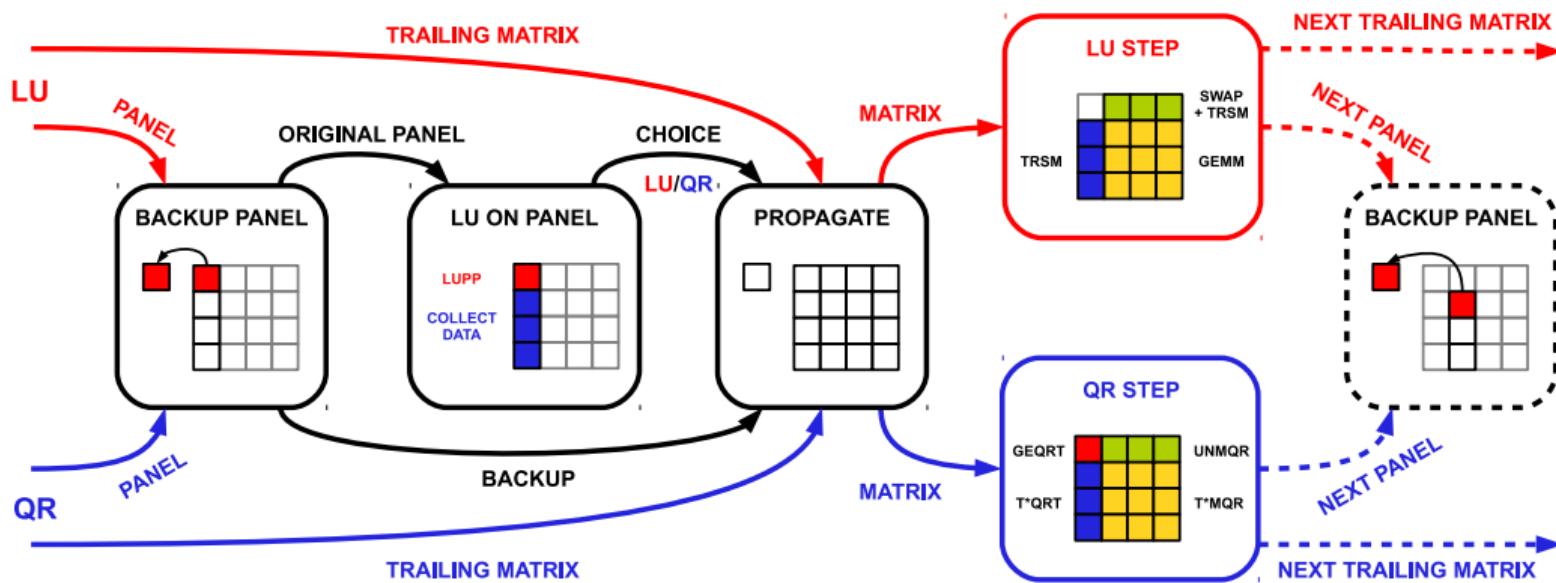
## QR factorization – Hierarchical algorithm in action



# QR factorization – Hierarchical algorithm in action



# 2014 Hybrid LU-QR solvers

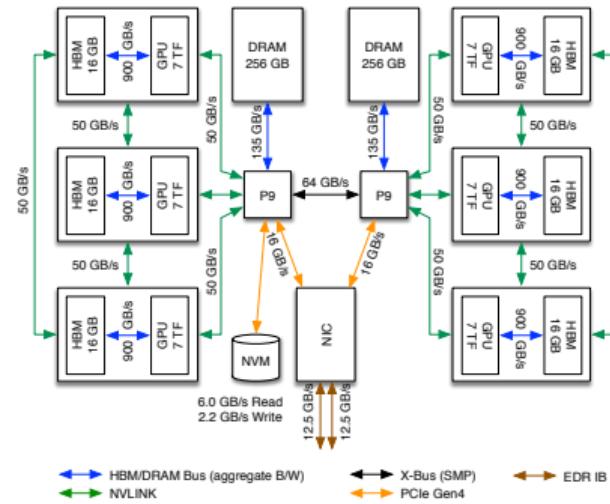


# 2019-2021 Matrix product on multi-GPU nodes

## SUMMIT architecture

- Oak Ridge National Laboratory Leadership Computing Facility
- 9,216 compute nodes
  - POWER9 22-core CPUs with 512 GB / node
  - 6 Nvidia Tesla V100 GPUs with 16GB / card

Heavy nodes with many accelerators per node

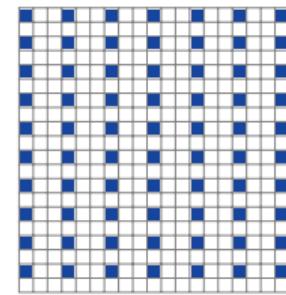
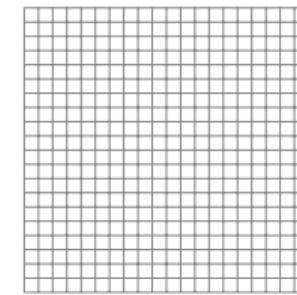
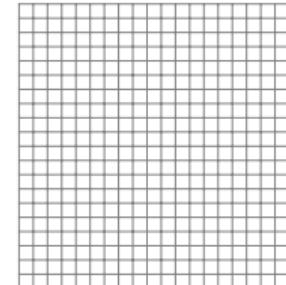


"Summit Interconnection Network Spectrum MPI and InfiniBand", Christopher Zimmer, Dec. 2018

# GEMM Algorithm for GPUs

## Node-level Task Distribution

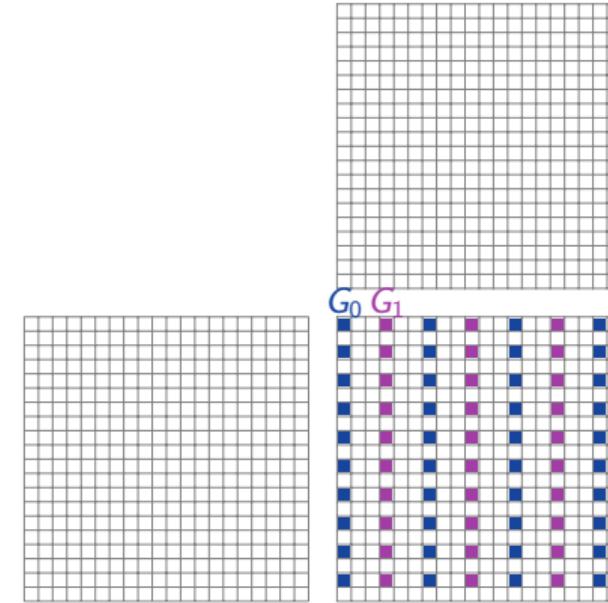
- data affinity
- owner computes heuristic
- 2D block cyclic distribution



# GEMM Algorithm for GPUs

## GPU-level Task Distribution

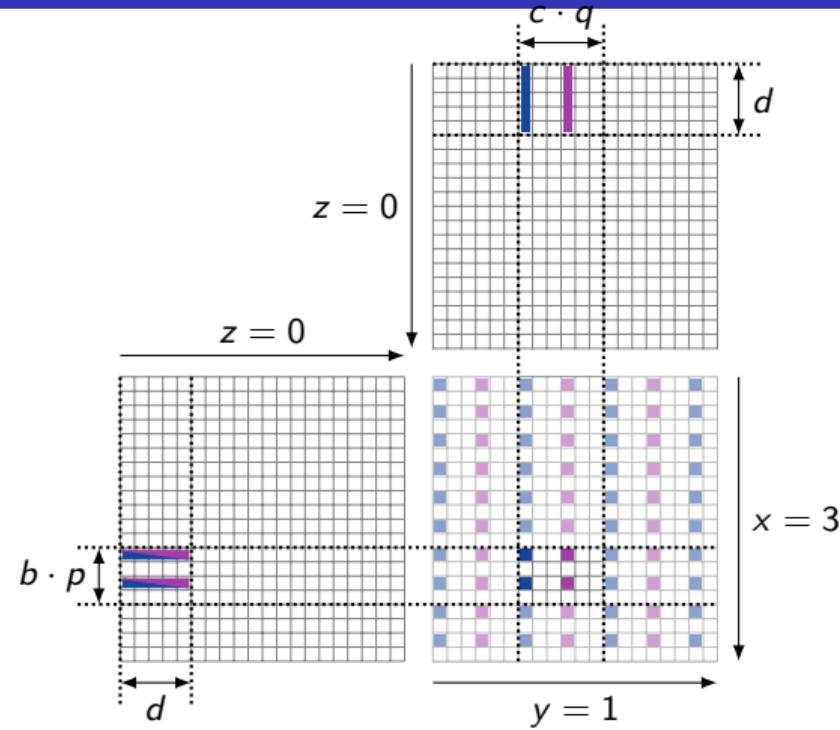
- Fixed, using PaRSEC memory advise API
- Round-robin assignment of tile-columns to the different GPUs



# GEMM Algorithm for GPUs

## Node-Level Blocking

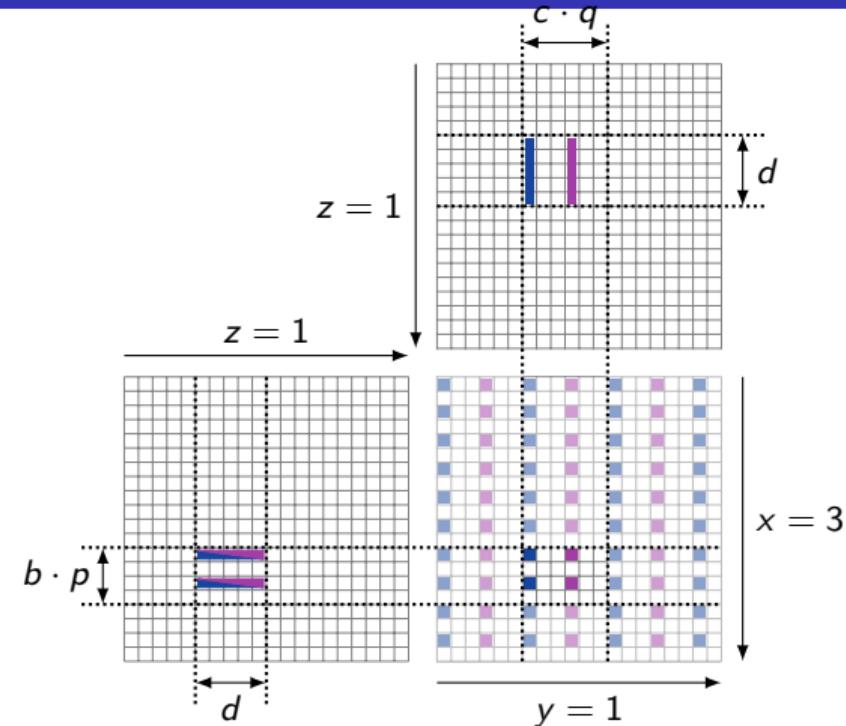
- Active set of GEMMs is defined locally by a block of coordinate  $(x, y, z)$
- $(x, y)$  defines a block in  $C$ , of *local* size  $b \times c$
- $(x, z)$  defines a block in  $A$ , of *global* size  $d \times bp$
- $(z, y)$  defines a block in  $B$ , of *global* size  $d \times cq$
- Order of progress follows  $z$ , then  $x$  and  $y$ .



# GEMM Algorithm for GPUs

## Node-Level Blocking

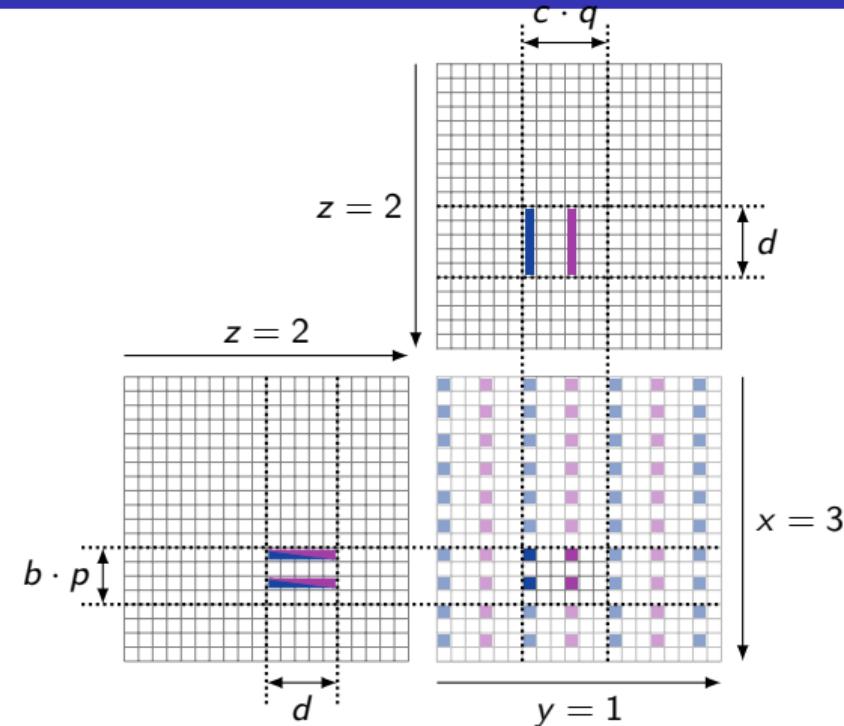
- Controls the size of the active set on GPUs
- if  $z < MAX_z$  :  
 $(x, y, z) \rightarrow (x, y, z + 1)$
- if  $z = MAX_z \wedge x < MAX_x$  :  
 $(x, y, z) \rightarrow (x + 1, y, 0)$
- if  $z = MAX_z \wedge x = MAX_x$  :  
 $(x, y, z) \rightarrow (0, y + 1, 0)$



# GEMM Algorithm for GPUs

## Node-Level Blocking

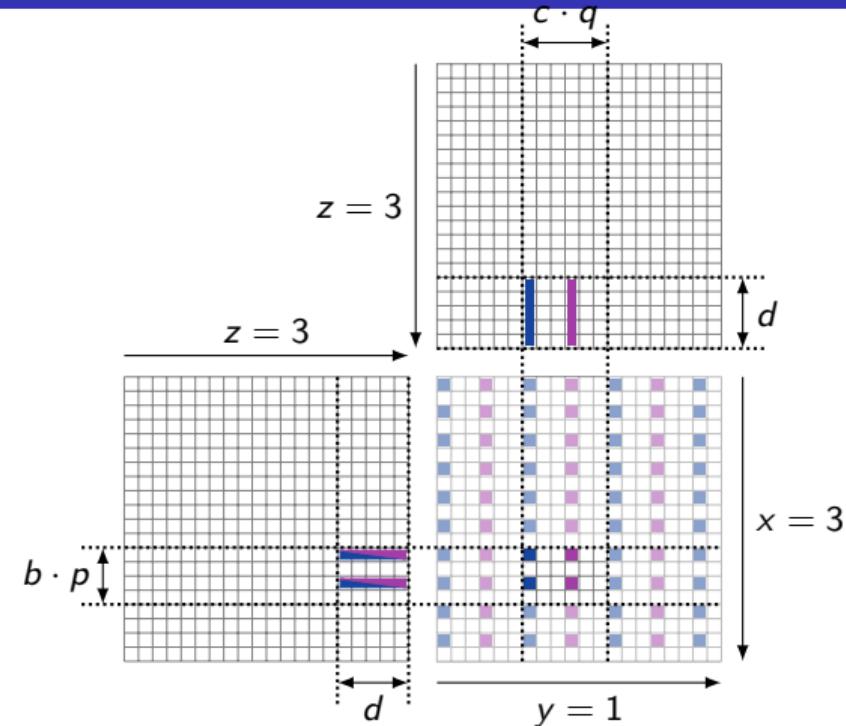
- Controls the size of the active set on GPUs
- if  $z < MAX_z$  :  
 $(x, y, z) \rightarrow (x, y, z + 1)$
- if  $z = MAX_z \wedge x < MAX_x$  :  
 $(x, y, z) \rightarrow (x + 1, y, 0)$
- if  $z = MAX_z \wedge x = MAX_x$  :  
 $(x, y, z) \rightarrow (0, y + 1, 0)$



# GEMM Algorithm for GPUs

## Node-Level Blocking

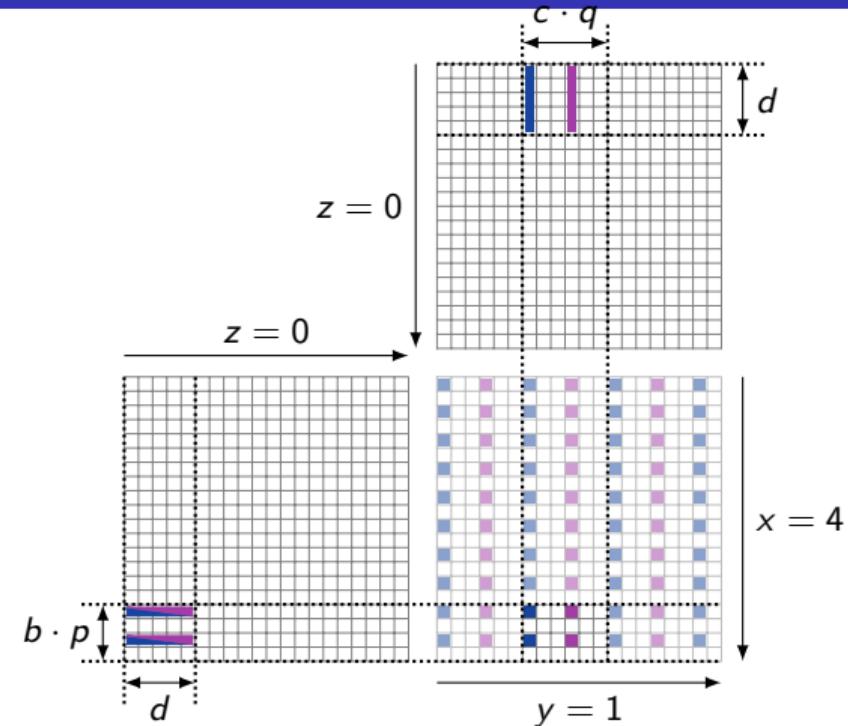
- Controls the size of the active set on GPUs
- if  $z < MAX_z$  :  
 $(x, y, z) \rightarrow (x, y, z + 1)$
- if  $z = MAX_z \wedge x < MAX_x$  :  
 $(x, y, z) \rightarrow (x + 1, y, 0)$
- if  $z = MAX_z \wedge x = MAX_x$  :  
 $(x, y, z) \rightarrow (0, y + 1, 0)$



# GEMM Algorithm for GPUs

## Node-Level Blocking

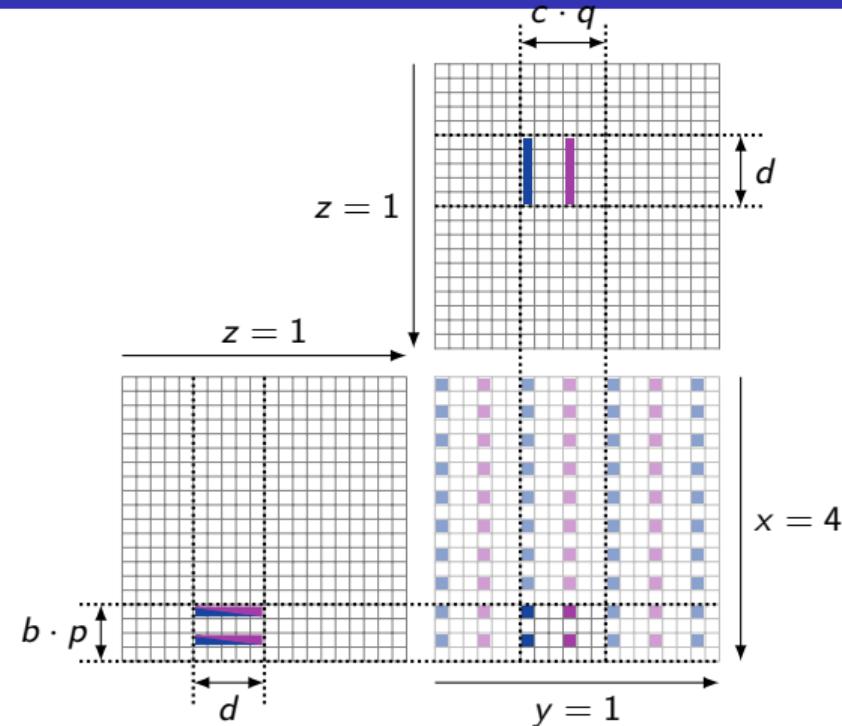
- Controls the size of the active set on GPUs
- if  $z < MAX_z$  :  
 $(x, y, z) \rightarrow (x, y, z + 1)$
- if  $z = MAX_z \wedge x < MAX_x$  :  
 $(x, y, z) \rightarrow (x + 1, y, 0)$
- if  $z = MAX_z \wedge x = MAX_x$  :  
 $(x, y, z) \rightarrow (0, y + 1, 0)$



# GEMM Algorithm for GPUs

## Node-Level Blocking

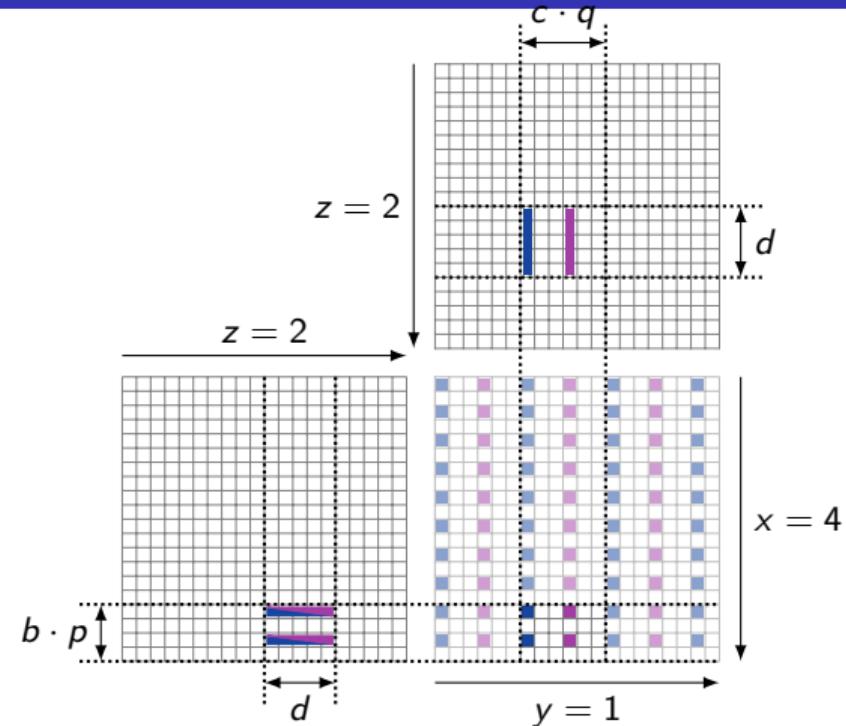
- Controls the size of the active set on GPUs
- if  $z < MAX_z$  :  
 $(x, y, z) \rightarrow (x, y, z + 1)$
- if  $z = MAX_z \wedge x < MAX_x$  :  
 $(x, y, z) \rightarrow (x + 1, y, 0)$
- if  $z = MAX_z \wedge x = MAX_x$  :  
 $(x, y, z) \rightarrow (0, y + 1, 0)$



# GEMM Algorithm for GPUs

## Node-Level Blocking

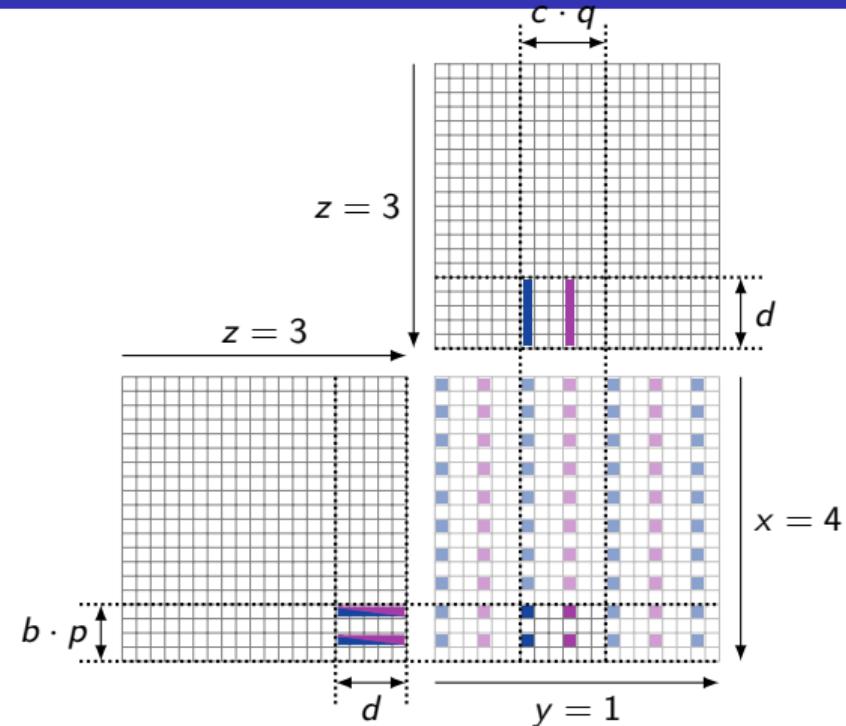
- Controls the size of the active set on GPUs
- if  $z < MAX_z$  :  
 $(x, y, z) \rightarrow (x, y, z + 1)$
- if  $z = MAX_z \wedge x < MAX_x$  :  
 $(x, y, z) \rightarrow (x + 1, y, 0)$
- if  $z = MAX_z \wedge x = MAX_x$  :  
 $(x, y, z) \rightarrow (0, y + 1, 0)$



# GEMM Algorithm for GPUs

## Node-Level Blocking

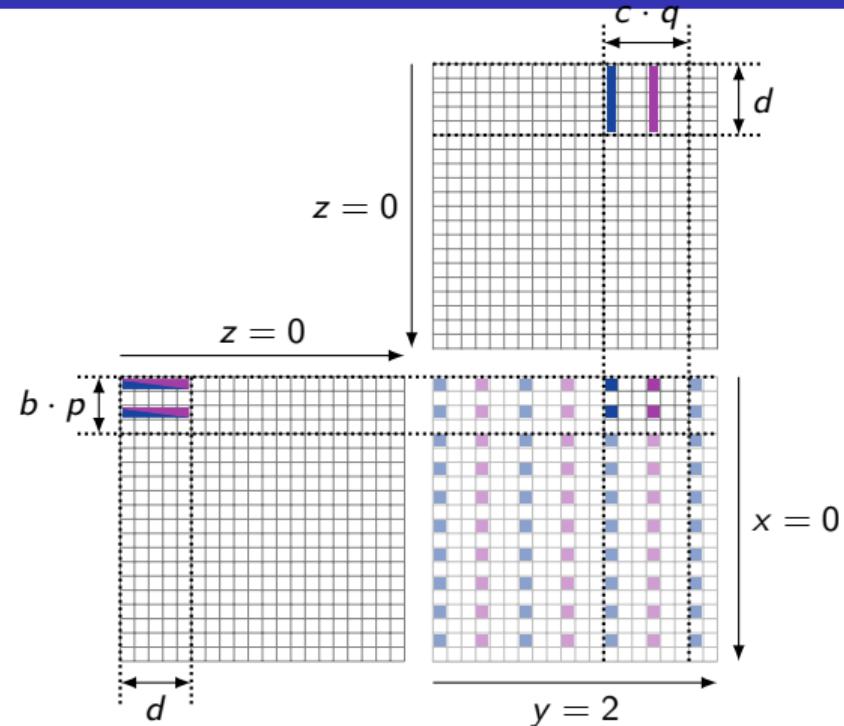
- Controls the size of the active set on GPUs
- if  $z < MAX_z$  :  
 $(x, y, z) \rightarrow (x, y, z + 1)$
- if  $z = MAX_z \wedge x < MAX_x$  :  
 $(x, y, z) \rightarrow (x + 1, y, 0)$
- if  $z = MAX_z \wedge x = MAX_x$  :  
 $(x, y, z) \rightarrow (0, y + 1, 0)$



# GEMM Algorithm for GPUs

## Node-Level Blocking

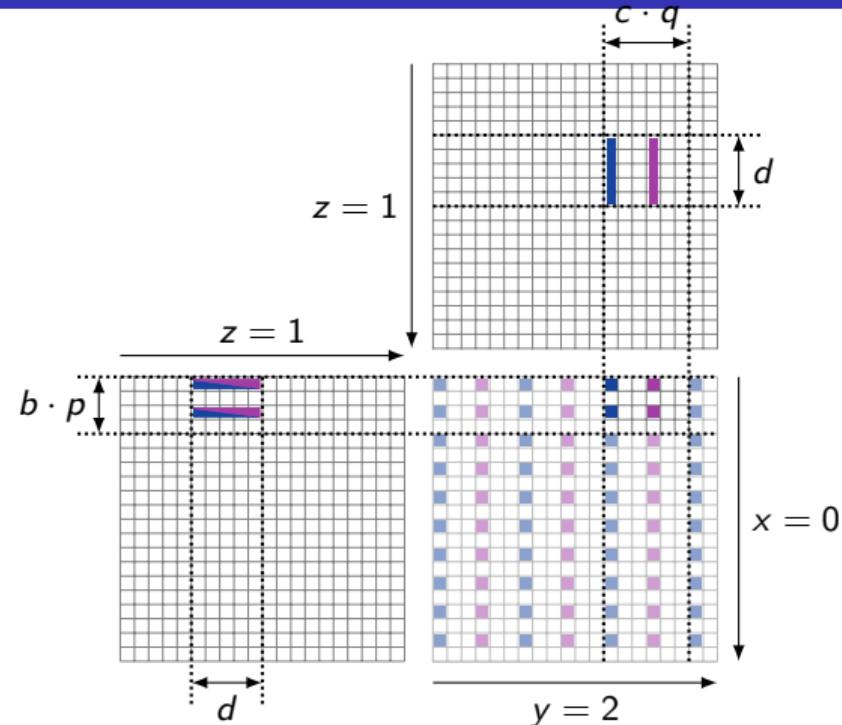
- Controls the size of the active set on GPUs
- if  $z < MAX_z$  :  
 $(x, y, z) \rightarrow (x, y, z + 1)$
- if  $z = MAX_z \wedge x < MAX_x$  :  
 $(x, y, z) \rightarrow (x + 1, y, 0)$
- if  $z = MAX_z \wedge x = MAX_x$  :  
 $(x, y, z) \rightarrow (0, y + 1, 0)$



# GEMM Algorithm for GPUs

## Node-Level Blocking

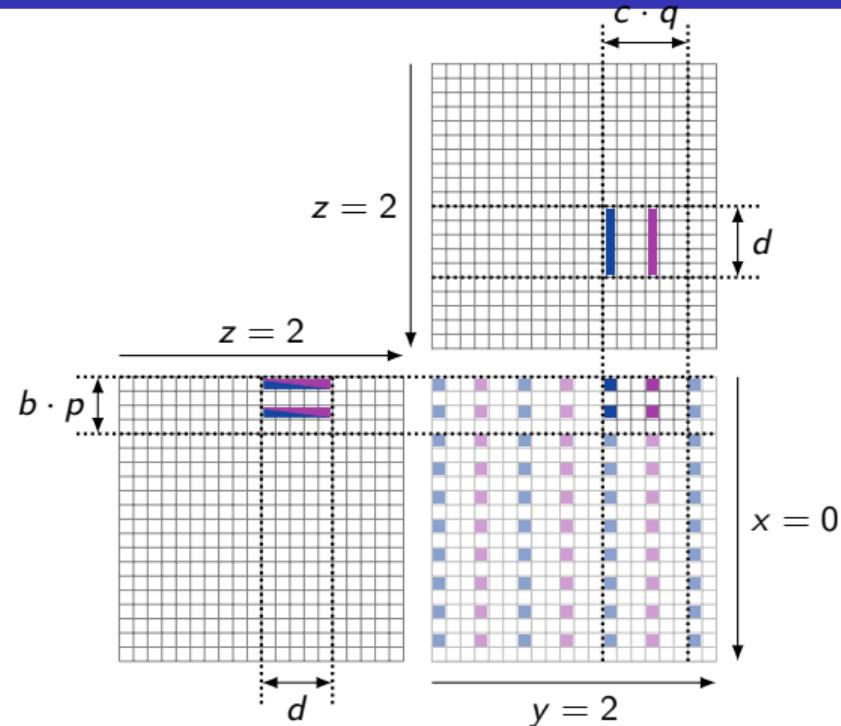
- Controls the size of the active set on GPUs
- if  $z < MAX_z$  :  
 $(x, y, z) \rightarrow (x, y, z + 1)$
- if  $z = MAX_z \wedge x < MAX_x$  :  
 $(x, y, z) \rightarrow (x + 1, y, 0)$
- if  $z = MAX_z \wedge x = MAX_x$  :  
 $(x, y, z) \rightarrow (0, y + 1, 0)$



# GEMM Algorithm for GPUs

## Node-Level Blocking

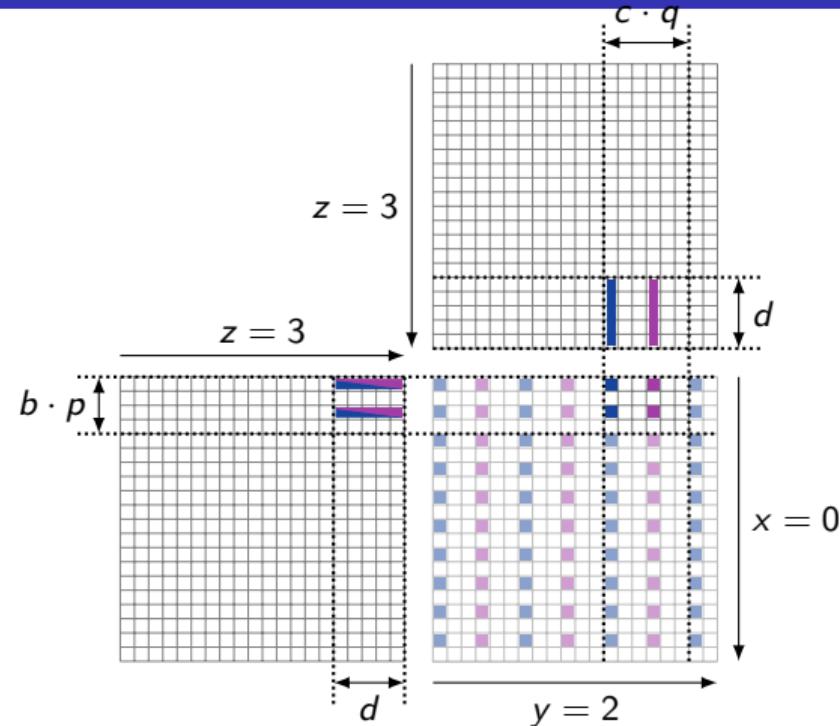
- Controls the size of the active set on GPUs
- if  $z < MAX_z$  :  
 $(x, y, z) \rightarrow (x, y, z + 1)$
- if  $z = MAX_z \wedge x < MAX_x$  :  
 $(x, y, z) \rightarrow (x + 1, y, 0)$
- if  $z = MAX_z \wedge x = MAX_x$  :  
 $(x, y, z) \rightarrow (0, y + 1, 0)$



# GEMM Algorithm for GPUs

## Node-Level Blocking

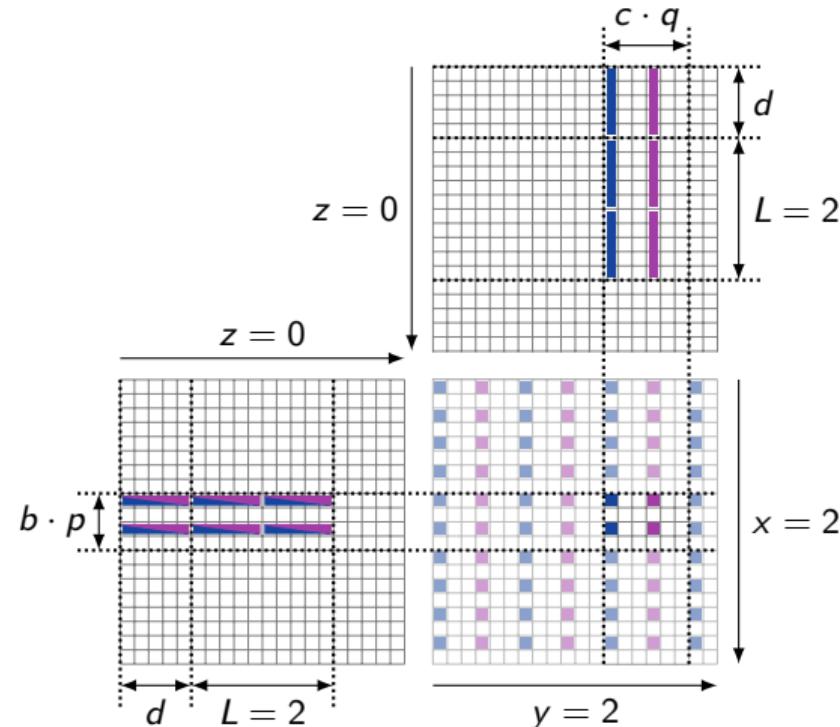
- Controls the size of the active set on GPUs
- if  $z < MAX_z$  :  
 $(x, y, z) \rightarrow (x, y, z + 1)$
- if  $z = MAX_z \wedge x < MAX_x$  :  
 $(x, y, z) \rightarrow (x + 1, y, 0)$
- if  $z = MAX_z \wedge x = MAX_x$  :  
 $(x, y, z) \rightarrow (0, y + 1, 0)$



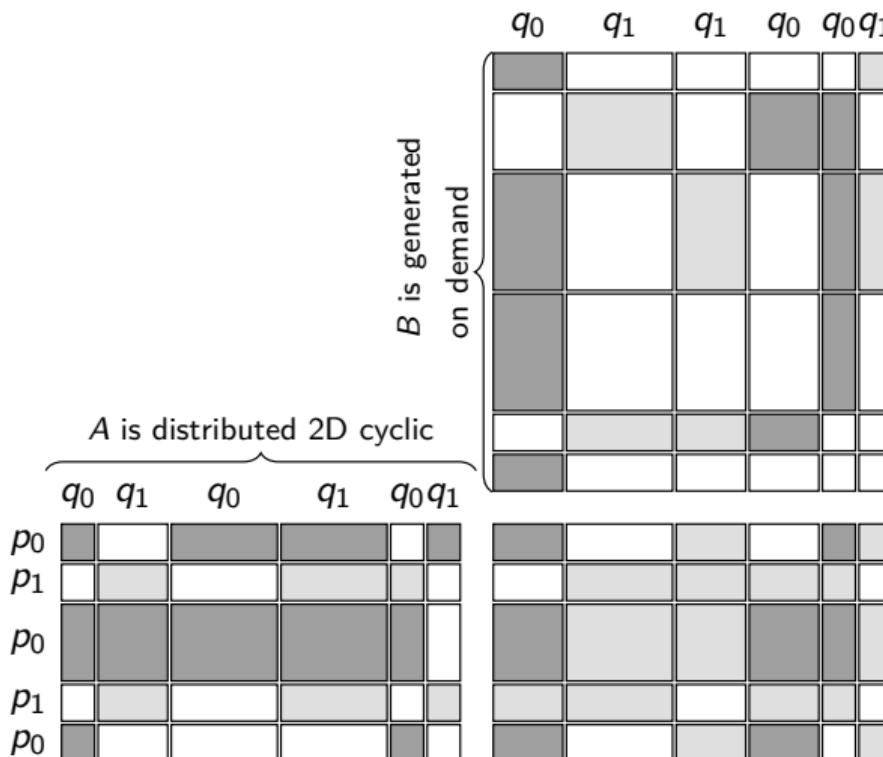
# GEMM Algorithm for GPUs

Machine-level Blocking:

- Main RAM is used as a temporary buffer
- Look ahead parameter  $L$ : #blocks loaded in advance
- Global synchronizations prevent any node to progress more than  $L$  steps than slowest node
  - Prevent overloading a node with download request
  - Control amount of temporary memory



# Tensor product



- $M = K \gg N$ :  $B$  is huge in front of  $A$  or  $C$ .

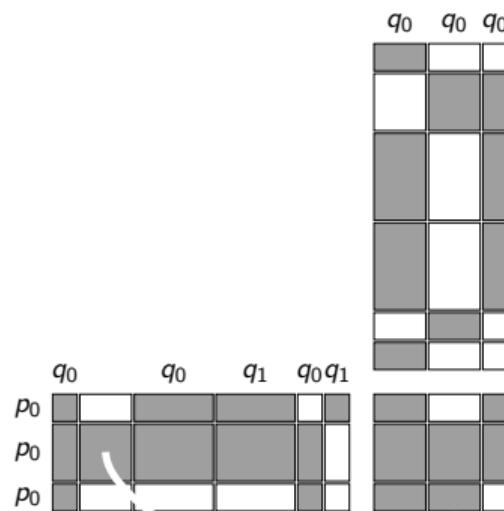
- Tiles of  $B$  are generated on demand
- Generating a tile is non trivial

## Strategy for $B$

Tiles of  $B$  are generated once, when needed, used then discarded

A single node reads a given tile of  $B$   
Load balance flops between nodes

# On a node



## Distribute $B$ on GPUs

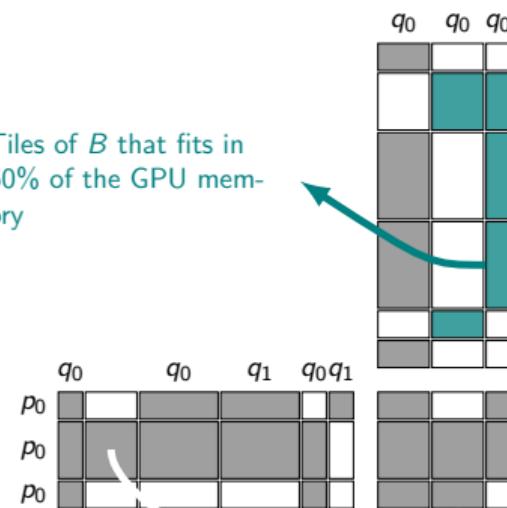
Assign each columns of  $B$  to a single GPU on the node  
Load balance flops between GPUs

## Blocking $B$ for a given GPU

Tiles of  $B$  occupy at most 50% of GPU memory

- Sort local columns of  $B$  by size
- Split execution in phases
- Greedy algorithm to fill each phase with as many (full) columns of  $B$

# On a node



Tiles of  $B$  that fits in 50% of the GPU memory

## Distribute $B$ on GPUs

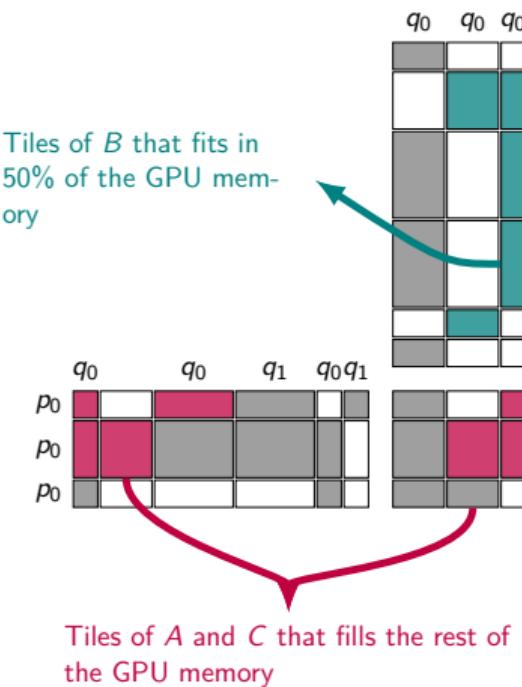
Assign each columns of  $B$  to a single GPU on the node  
Load balance flops between GPUs

## Blocking $B$ for a given GPU

Tiles of  $B$  occupy at most 50% of GPU memory

- Sort local columns of  $B$  by size
- Split execution in phases
- Greedy algorithm to fill each phase with as many (full) columns of  $B$

# On a node



## Blocking A and C for a given block phase

Tiles of A and C fill up the rest

Greedy with heuristic:  $\approx$ rectangles of A of size  $d$  by what fits, and rectangles of C of size  $d$  by #columns of B

- Split previous column phases in block-phases
- In each block-phase, assign as many GEMMs as possible
- Such that tiles of A and C fit in the remaining GPU memory
- Looking at tiles of A vertically until at least  $d$  rows are selected, then adding only tiles of the same rows, from left to right
- if memory is still available, add a row, iterate.

# Outline

1 Redistribution

2 Linear algebra kernels

3 Resilience

4 Life at ICL

# 2013-2021 Resilience

- In-memory checkpointing
- ABFT and composite strategies
- Failure detection
- Cooperative checkpointing
- Distributed termination
- ...

Redistribution  
oooooooo

Linear algebra kernels  
oooooooooooooooooooo

Resilience  
oo

Life at ICL  
●oooooooo

# Outline

1 Redistribution

2 Linear algebra kernels

3 Resilience

4 Life at ICL

## 1996-1997



- Spent a few months with Jack in IBM ECSEC Rome, 1986-87
- Waited for Bernard Tourancheau and Frédéric Desprez to report on Knoxville 😊
- Well, I liked country music **before** coming to Tennessee 😊

Redistribution  
oooooooo

Linear algebra kernels  
oooooooooooooooooooo

Resilience  
oo

Life at ICL  
oo●oooo

# At Jack's



Redistribution  
ooooooooo

Linear algebra kernels  
oooooooooooooooooooo

Resilience  
oo

Life at ICL  
ooo●ooo

# At Jack's



Redistribution  
oooooooo

Linear algebra kernels  
oooooooooooooooooooo

Resilience  
oo

Life at ICL  
oooo●○○

## The dream team



Redistribution  
oooooooo

Linear algebra kernels  
oooooooooooooooooooo

Resilience  
oo

Life at ICL  
oooooo●o

# Leaving ...



# A little theorem for the road

**Theorem** Life is good at ICL

## Proof

- By evidence (look at us)
- By enumeration (enough to count the French mafia)

Thanks to Jack and everybody at ICL 😊