Problem
oooooo

Resolution Methods
ooooooooooo

Evaluation
ooooo

Perspectives
ooo

# Data Localisation for Distributed Applications
## Compressed Cholesky Case Study

Olivier Beaumont
Lionel Eyraud-Dubois
Mathieu Vérité

INRIA Bordeaux - HiePACS Team

Solharis Project Kickoff Meeting
February 5th, 2020

*Inria*

Problem
oooooo

Resolution Methods
ooooooooooo

Evaluation
ooooo

Perspectives
ooo

*Inria*

**Problem**
oooooo

Resolution Methods
ooooooooooo

Evaluation
ooooo

Perspectives
ooo

# 1. Problem

*Inria*

## 1.1 - Data Localisation

Distributed application: static data allocation to computation nodes

### Objectives
- completion time
- result reliability
- execution replicability

### Issues
- load balancing
- communication management

### Strategies
- optimise static allocation
- on the run re-allocation

**Case study: Cholesky decomposition**

*Inria*

# 1.2 - Problem Description

Formal problem:

- $\mathcal{P}$: set of $P$ computation **resources** / **processes**
- tasks **dependencies**: application DAG
- for each task: known constant proportionality
  **input "size"** $\leftrightarrow$ **task execution time**
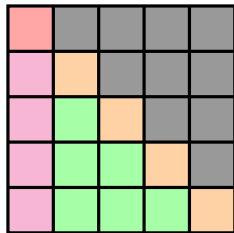
## Objective

minimise *makespan*

## Methodology

- solve approximate problem: *load balancing*
- evaluate *makespan* in execution

**Problem**
○○●○○○

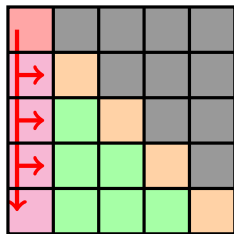Resolution Methods
○○○○○○○○○○○

Evaluation
○○○○○

Perspectives
○○○

# 1.3 - Cholesky Decomposition

- Physical problems modeling: $\mathbf{A} * x = \mathbf{b}$
  $\mathbf{A} \in \mathbb{R}^{n \times n}$ typically **symetric positive definite**
- Steps for resolution: $\mathbf{A} = \mathbf{L} * \mathbf{L}^t \rightarrow \mathbf{L} * y = \mathbf{B} \rightarrow \mathbf{L}^t * x = y$

### *Right-looking* algorithm:
### $\mathcal{O}(N^3)$ complexity

- ▷ **Input** : $\mathbf{A}$
- ▷ **Initialisation** : $\mathbf{L} = \mathbf{A}_{\text{tri.inf.}}$
- ▷ **For** $k = 1 \rightarrow N$ : POTRF($k$)
  - ▷ **For** $i = k + 1 \rightarrow N$ : TRSM($i, k$)
  - ▷ **For** $j = k + 1 \rightarrow N$ : SYRK($j, k$)
    - ▷ **For** $i = j + 1 \rightarrow N$
      GEMM($i, j, k$)
- ▷ **Output** : $\mathbf{L}$
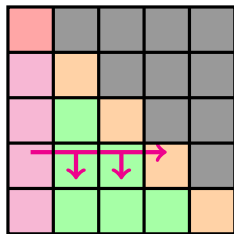
**Problem**
○○●○○○

Resolution Methods
○○○○○○○○○○○

Evaluation
○○○○○

Perspectives
○○○

# 1.3 - Cholesky Decomposition

- Physical problems modeling: $\mathbf{A} * x = \mathbf{b}$
  $\mathbf{A} \in \mathbb{R}^{n \times n}$ typically **symetric positive definite**
- Steps for resolution: $\mathbf{A} = \mathbf{L} * \mathbf{L}^t \rightarrow \mathbf{L} * y = \mathbf{B} \rightarrow \mathbf{L}^t * x = y$

*Right-looking* algorithm:
$\mathcal{O}(N^3)$ complexity

▷ **Input** : **A**

▷ **Initialisation** : $\mathbf{L} = \mathbf{A}_{\mathrm{tri.inf.}}$

▷ **For** $k = 1 \rightarrow N$ : POTRF($k$)

    ▷ **For** $i = k + 1 \rightarrow N$ : TRSM($i, k$)

    ▷ **For** $j = k + 1 \rightarrow N$ : SYRK($j, k$)

        ▷ **For** $i = j + 1 \rightarrow N$
        GEMM($i, j, k$)

▷ **Output** : L



Column broadcast
($k = 1$)

*Inria*

**Problem**
○○●○○○

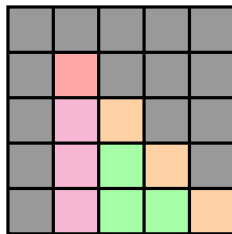**Resolution Methods**
○○○○○○○○○○○

**Evaluation**
○○○○○

**Perspectives**
○○○

# 1.3 - Cholesky Decomposition

- Physical problems modeling: $\mathbf{A} * x = \mathbf{b}$
  $\mathbf{A} \in \mathbb{R}^{n \times n}$ typically **symetric positive definite**
- Steps for resolution: $\mathbf{A} = \mathbf{L} * \mathbf{L}^t \rightarrow \mathbf{L} * y = \mathbf{B} \rightarrow \mathbf{L}^t * x = y$

*Right-looking* algorithm:
$\mathcal{O}(N^3)$ complexity

▷ **Input** : **A**

▷ **Initialisation** : $\mathbf{L} = \mathbf{A}_{\mathrm{tri.inf.}}$

▷ **For** $k = 1 \rightarrow N$ : POTRF($k$)

    ▷ **For** $i = k + 1 \rightarrow N$ : TRSM($i, k$)

    ▷ **For** $j = k + 1 \rightarrow N$ : SYRK($j, k$)

        ▷ **For** $i = j + 1 \rightarrow N$
        GEMM($i, j, k$)

▷ **Output** : L



Row broadcast

($k = 1$)

*Inría*

**Problem**
○○●○○○

Resolution Methods
○○○○○○○○○○○

Evaluation
○○○○○

Perspectives
○○○

# 1.3 - Cholesky Decomposition

- Physical problems modeling: $\mathbf{A} * x = \mathbf{b}$
  $\mathbf{A} \in \mathbb{R}^{n \times n}$ typically **symetric positive definite**
- Steps for resolution: $\mathbf{A} = \mathbf{L} * \mathbf{L}^t \rightarrow \mathbf{L} * y = \mathbf{B} \rightarrow \mathbf{L}^t * x = y$

**_Right-looking_ algorithm:**
$\mathcal{O}(N^3)$ complexity

▷ **Input** : **A**

▷ **Initialisation** : $\mathbf{L} = \mathbf{A}_{\mathrm{tri.inf.}}$

▷ **For** $k = 1 \rightarrow N$ : POTRF($k$)

    ▷ **For** $i = k + 1 \rightarrow N$ : TRSM($i, k$)

    ▷ **For** $j = k + 1 \rightarrow N$ : SYRK($j, k$)

        ▷ **For** $i = j + 1 \rightarrow N$
          GEMM($i, j, k$)

▷ **Output** : **L**



Next iteration

($k = 2$)

*Inria*

**Problem**
oooo● oo

Resolution Methods
ooooooooooo

Evaluation
ooooo

Perspectives
ooo

# 1.4 - Working Assumptions

## Communications scheme:
- data transfer on same position
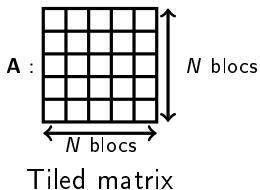- broadcast on same row / column

## Specific solutions
- unmodified association: **tile** $\leftrightarrow$ **proc.**
- max. number of different proc.
  on each row/column: $m^{row}$, $m^{col}$

## Assumption
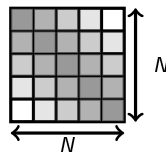- **limited communication** $\Rightarrow$ simultaneous with calculation

*Inría*

**Problem**
○○○○●○

Resolution Methods
○○○○○○○○○○○

Evaluation
○○○○○

Perspectives
○○○

# 1.5 - Input Data Handling (1/2): Compression

$\mathbf{A}$ :

$N$ blocs

$N$ blocs

Tiled matrix

### Density distribution

"Pseudo-rank" : $r_{i,j}$

Density: $d(A_{i,j}) = \frac{r_{i,j}}{m}$

$N$

$N$

### Block Low Rank compression

$m$

$m$ $A_{i,j}$

$\approx$

$m$ $r_{i,j}$

$\mathbf{U}_{i,j}$

$m$

$\mathbf{V}_{i,j}$

*Inria*

## 1.5 - Input Data Handling (2/2): Modification

**Allocation:** tile $\leftrightarrow$ proc.

# 1.5 - Input Data Handling (2/2): Modification

**Allocation:** tile $\leftrightarrow$ proc.

Position $(i, j)$

Weighted sum over all iterations

$$\bar{\mathbf{A}}_{i,j} = d(A_{i,j}) \times \sum_{k=1}^{j} \text{task}[i, j, k]$$



TRSM

GEMM

GEMM

$k = 3$

$k = 2$

$k = 1$

$\Rightarrow$ aggregated matrix $\bar{\mathbf{A}}$

*Inria*

Problem
oooooo

Resolution Methods
ooooooooooo

Evaluation
ooooo

Perspectives
ooo

# 2. Resolution Methods

# 2.1 - Problem Simplification

## Assumptions

set of **independent tasks**

## Objective

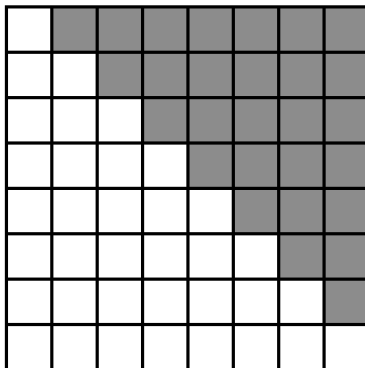$$\min\{makespan\} \Leftrightarrow \min\{\max_{p \in \mathcal{P}}\{load_p\}\}$$

$\rightarrow$ *load balancing* problem

*Inría*

Problem
○○○○○○

Resolution Methods
○●○○○○○○○○○○

Evaluation
○○○○○

Perspectives
○○○

## 2.2 - Block Cyclic Method

Repeated block



$r$

$c$

proc. 1    proc. 4

proc. 2    proc. 5

proc. 3    proc. 6

Matrix $\bar{\mathbf{A}}$



*Ínría*

Problem
○○○○○○

Resolution Methods
○●○○○○○○○○○

Evaluation
○○○○○

Perspectives
○○○

## 2.2 - Block Cyclic Method

Repeated block

Matrix $\bar{\mathbf{A}}$



$r$

$c$

$$\begin{cases} \triangleright & r \times c = P \\ \\ \triangleright & r \neq c \\ \\ \triangleright & \min(|r - c|) \end{cases}$$

*Ínría*

Problem
○○○○○○

Resolution Methods
○●○○○○○○○○○

Evaluation
○○○○○

Perspectives
○○○

## 2.2 - Block Cyclic Method



Repeated block

Matrix $\bar{\mathbf{A}}$

Problem
○○○○○○

Resolution Methods
○○●○○○○○○○○○

Evaluation
○○○○○

Perspectives
○○○

# 2.3 - Extended Block Cyclic (1/3)

Initial block    Extended block



$m^{col} = \alpha \times r$

$m^{row} = \alpha \times c$

$\alpha^2 \times$ more positions

Communication constraints

$\Leftrightarrow$ **parameter**:
$\alpha \in [1; +\infty[$

*Inria*

Problem
oooooo

Resolution Methods
ooooo●oooooo

Evaluation
ooooo

Perspectives
ooo

# 2.3 - Extended Block Cyclic (2/3)

$W$:



| $18^{th}$ | $5^{th}$ | $12^{th}$ | $4^{th}$ | $20^{th}$ | $11^{th}$ |
|---|---|---|---|---|---|
| $17^{th}$ | $24^{th}$ | $22^{nd}$ | $19^{th}$ | $1^{st}$ | $7^{th}$ |
| $10^{th}$ | $6^{th}$ | $23^{rd}$ | $15^{th}$ | $2^{nd}$ | $16^{th}$ |
| $9^{th}$ | $3^{rd}$ | $8^{th}$ | $13^{th}$ | $21^{st}$ | $14^{th}$ |

$m^{col}$

$m^{row}$

### Greedy procedure

- tiles: $\searrow$ load

$$W_{i,j} = \sum_{\substack{u = m^{col} \times i \in [\![1;N]\!] \\ v = m^{row} \times j \in [\![1;N]\!]}} \bar{A}_{u,v}$$

*Inría*

Problem
○○○○○○

Resolution Methods
○○○○●○○○○○○○

Evaluation
○○○○○

Perspectives
○○○

## 2.3 - Extended Block Cyclic (2/3)

$W$:



$m^{col}$

$m^{row}$

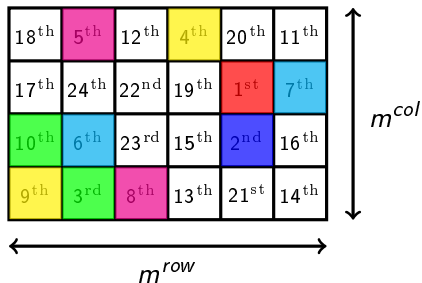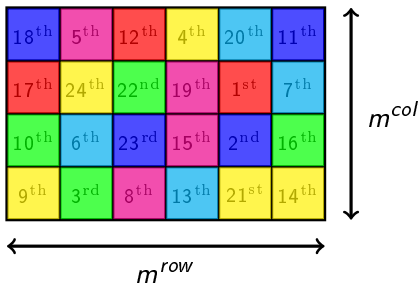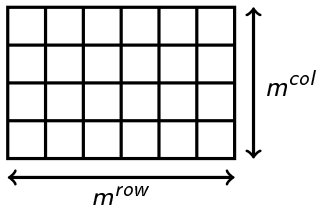Greedy procedure
- tiles: ↘ load
- alloc. least loaded

$$W_{i,j} = \sum_{\substack{u=m^{col} \times i \in [\![1;N]\!] \\ v=m^{row} \times j \in [\![1;N]\!]}} \bar{A}_{u,v}$$

Inría

Problem
oooooo

Resolution Methods
oooo●oooooooo

Evaluation
ooooo

Perspectives
ooo

# 2.3 - Extended Block Cyclic (2/3)

$W$:



Greedy procedure

- tiles: $\searrow$ load
- alloc. least loaded

$$W_{i,j} = \sum_{\substack{u=m^{col} \times i \in [\![1;N]\!] \\ v=m^{row} \times j \in [\![1;N]\!]}} \bar{A}_{u,v}$$

Problem
○○○○○○

Resolution Methods
○○○○●○○○○○○○

Evaluation
○○○○○

Perspectives
○○○

# 2.3 - Extended Block Cyclic (2/3)

$W$:



**Greedy procedure**

- tiles: ↘ load
- alloc. least loaded

$$W_{i,j} = \sum_{\substack{u = m^{col} \times i \in [\![1;N]\!] \\ v = m^{row} \times j \in [\![1;N]\!]}} \bar{A}_{u,v}$$

*Inría*

Problem
oooooo

Resolution Methods
ooooo●oooooo

Evaluation
ooooo

Perspectives
ooo

# 2.3 - Extended Block Cyclic (3/3)

$W$:



$m^{col}$

$m^{row}$

$$W_{i,j} = \sum_{\substack{u=m^{col}\times i\in[\![1;N]\!] \\ v=m^{row}\times j\in[\![1;N]\!]}} \bar{A}_{u,v}$$

Decision variables:

$$\left[ \begin{array}{l} \mathbf{x}_{i,j}^{(p)} = \left\{ \begin{array}{ll} 1 & \text{if proc. } p \text{ on } (i,j) \\ 0 & \text{otherwise} \end{array} \right. \\ \\ \mathbf{w} = \begin{array}{l} \text{maximum working time} \\ \text{among all processes} \end{array} \end{array} \right.$$

*Inria*

## 2.3 - Extended Block Cyclic (3/3)

Integer linear program:

$$\min\{\mathbf{w}\}$$

$$
\begin{cases}
(1) & \forall p \in [\![1;P]\!] & \displaystyle\sum_{\substack{i \in [\![1;m^{col}]\!] \\ j \in [\![1;m^{row}]\!]}} \mathbf{x}_{i,j}^{(p)} \times W_{i,j} \leqslant \mathbf{w} & \text{max. time} \\[3em]
(2) & \forall (i,j) \in [\![1;m^{col}]\!] \times [\![1;m^{row}]\!] & \displaystyle\sum_{p \in [\![1;P]\!]} \mathbf{x}_{i,j}^{(p)} = 1 & \text{allocation}
\end{cases}
$$

*Inria*

Problem
oooooo

Resolution Methods
oooooo●ooooo

Evaluation
ooooo

Perspectives
ooo

# 2.4 - Random Subsets Methods (1/3)

**Idea**: randomly generate subsets of proc. for rows $(R_1, ..., R_Q)$ / columns $(C_1, ..., C_Q)$

- limited number of proc.: $m^{row}$; $m^{col}$
- each $(R_i, C_j)$ pair **compatible** : $\mathrm{Card}(I_{i,j} = R_i \bigcap C_j) \geqslant K$

## Advantages

- **independent** of $N$
- managing sampling difficulty: $Q$, $K$
- degree of freedom for tile allocation: $\mathrm{Card}(I)$



$R_1 \rightarrow \{$ ▢ $\}$
$R_2 \rightarrow \{$ ▢ $\}$ Row subsets
$R_3 \rightarrow \{$ ▢ $\}$

$C_1$ $C_2$ $C_3$

Column subsets

Problem
oooooo

Resolution Methods
ooooooooo●ooo

Evaluation
ooooo

Perspectives
ooo

## 2.4 - Random Subsets Methods (2/3): Two Steps



$8^{th}$
$7^{th}$
$6^{th}$
$5^{th}$
$4^{th}$
$3^{rd}$
$2^{nd}$
$1^{st}$

$7^{th}$ $5^{th}$ $4^{th}$ $3^{rd}$ $1^{st}$ $2^{nd}$ $6^{th}$ $8^{th}$

### Step 1

- rows / columns:
  ↘ sum load

*Inria*

Problem
○○○○○○

Resolution Methods
○○○○○○○●○○○○

Evaluation
○○○○○

Perspectives
○○○

# 2.4 - Random Subsets Methods (2/3): Two Steps



## Step 1

- rows / columns:
  ↘ sum load

- alloc. least
  loaded subset

Problem
○○○○○○

Resolution Methods
○○○○○○○●○○○

Evaluation
○○○○○

Perspectives
○○○

# 2.4 - Random Subsets Methods (2/3): Two Steps



### Step 1

- rows / columns:
  ↘ sum load

- alloc. least
  loaded subset

### Step 2

- tiles: ↘ load

Problem
oooooo

Resolution Methods
oooooooo●oooo

Evaluation
ooooo

Perspectives
ooo

# 2.4 - Random Subsets Methods (2/3): Two Steps



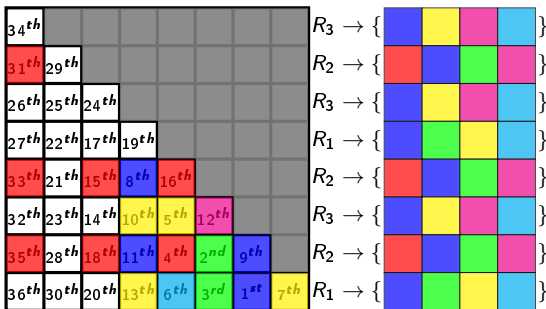## Step 1

- rows / columns: ↘ sum load

- alloc. least loaded subset

## Step 2

- tiles: ↘ load

- alloc *no choice* positions

Problem
○○○○○○

Resolution Methods
○○○○○○○●○○○

Evaluation
○○○○○

Perspectives
○○○

# 2.4 - Random Subsets Methods (2/3): Two Steps



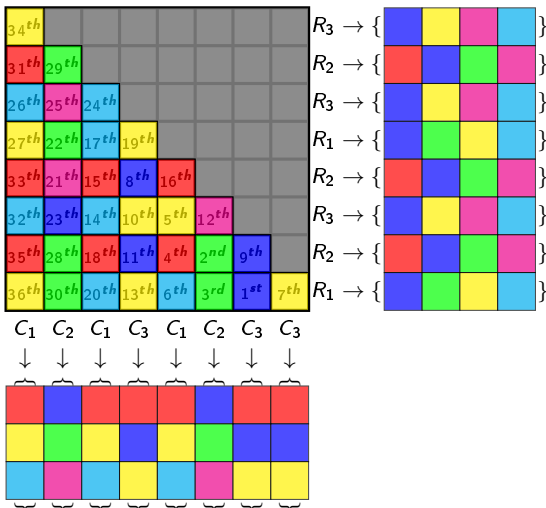$$I_{8,4} = R_1 \bigcap C_3$$
$$= \{ \quad \}$$

## Step 1

- rows / columns: ↘ sum load
- alloc. least loaded subset

## Step 2

- tiles: ↘ load
- alloc *no choice* positions
- available proc.: $I_{i,j} = R_i \bigcap C_j$

*Inria*

Problem
○○○○○○

Resolution Methods
○○○○○○○●○○○

Evaluation
○○○○○

Perspectives
○○○

# 2.4 - Random Subsets Methods (2/3): Two Steps



Step 1

- rows / columns:
  ↘ sum load
- alloc. least
  loaded subset

Step 2

- tiles: ↘ load
- alloc *no choice*
  positions
- available proc.:
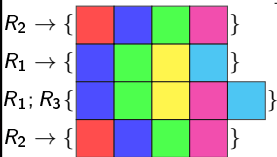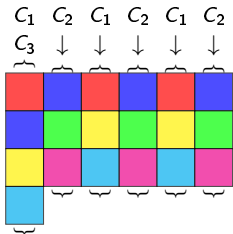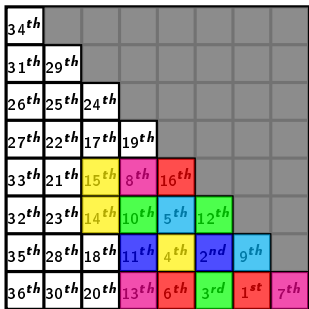  $I_{i,j} = R_i \bigcap C_j$
- alloc. least
  loaded proc.

$$I_{8,4} = R_1 \bigcap C_3$$
$$= \{ \quad \}$$

Problem
○○○○○○

Resolution Methods
○○○○○○○●○○○

Evaluation
○○○○○

Perspectives
○○○

## 2.4 - Random Subsets Methods (2/3): Two Steps



Step 1

- rows / columns: ↘ sum load
- alloc. least loaded subset

Step 2

- tiles: ↘ load
- alloc *no choice* positions
- available proc.: $I_{i,j} = R_i \bigcap C_j$
- alloc. least loaded proc.

Problem
oooooo

Resolution Methods
oooooooo●oo

Evaluation
ooooo

Perspectives
ooo

# 2.4 - Random Subsets Methods (3/3): Direct



**Steps**
- tiles: ↘ load

Problem
oooooo

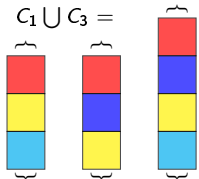Resolution Methods
ooooooooo●oo

Evaluation
ooooo

Perspectives
ooo

# 2.4 - Random Subsets Methods (3/3): Direct



Steps

- tiles: ↘ load
- available proc.

Problem
ooooo

Resolution Methods
ooooooooo●oo

Evaluation
ooooo

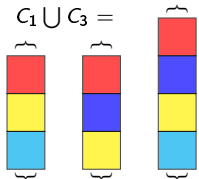Perspectives
ooo

# 2.4 - Random Subsets Methods (3/3): Direct



## Steps

- tiles: ↘ load
- available proc.
- alloc. least loaded

Problem
oooooo

Resolution Methods
ooooooooo●oo

Evaluation
ooooo

Perspectives
ooo

# 2.4 - Random Subsets Methods (3/3): Direct



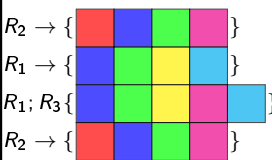### Steps

- tiles: ↘ load
- available proc.
- alloc. least loaded
  ⇒ subsets
  ⇒ *no choice*

Problem
○○○○○○

Resolution Methods
○○○○○○○○●○○

Evaluation
○○○○○

Perspectives
○○○

# 2.4 - Random Subsets Methods (3/3): Direct



Steps

- tiles: ↘ load
- available proc.
- alloc. least loaded
  $\Rightarrow$ subsets
  $\Rightarrow$ *no choice*

Problem
○○○○○○

Resolution Methods
○○○○○○○○○○●○

Evaluation
○○○○○

Perspectives
○○○

## 2.5 - By-Column Greedy Algorithm

Column allocation using integer linear program:



Allocate positions of
the current column

Problem
○○○○○○

Resolution Methods
○○○○○○○○○○○●

Evaluation
○○○○○

Perspectives
○○○

## 2.6 - Load Balancing: Some Results



Figure: Maximum load VS input matrix size
($\alpha = 2$; $N = 20$ to $40$; $\frac{N^2}{P} \approx$ cste)

Problem
○○○○○○

Resolution Methods
○○○○○○○○○○○●

Evaluation
○○○○○

Perspectives
○○○

## 2.6 - Load Balancing: Some Results



Figure: Maximum load VS input matrix size
($\alpha = 2$; $N = 60$ to $80$; $\frac{N^2}{P} \approx$ cste)

Problem
oooooo

Resolution Methods
ooooooooooo

Evaluation
ooooo

Perspectives
ooo

# 3. Evaluation

Problem
oooooo

Resolution Methods
ooooooooooo

Evaluation
●oooo

Perspectives
ooo

## 3.1 - Simulated Execution

**Assumption**

taking into account tasks dependencies

**Simulate execution: task based scheduler**

- DAG + tile allocation
- prioritised queues of ready tasks
- execution at proc. scale $\Rightarrow$ **preemption** allowed
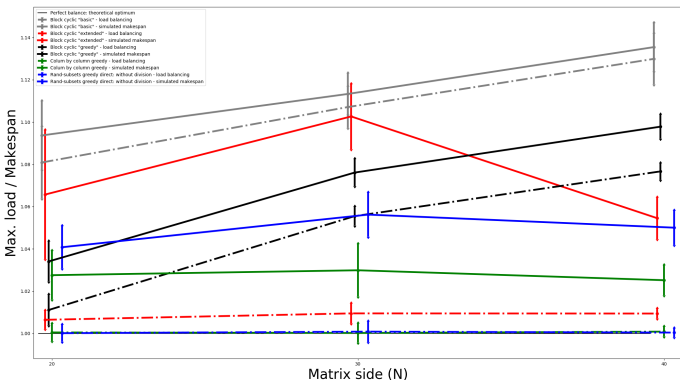- **no communication**

*Inría*

Problem
oooooo

Resolution Methods
ooooooooooo

Evaluation
o●oooo

Perspectives
ooo

## 3.2 - Makespan: Some Results



Figure: Makespan VS input matrix size
($\alpha = 2$; $N = 20$ to $40$; $\frac{N^2}{P} \approx$ cste)

Problem
○○○○○○○

Resolution Methods
○○○○○○○○○○○

Evaluation
○●○○○○

Perspectives
○○○

## 3.2 - Makespan: Some Results



Figure: Makespan VS input matrix size
($\alpha = 2$; $N = 60$ to $80$; $\frac{N^2}{P} \approx$ cste)

Problem
000000

Resolution Methods
00000000000

Evaluation
00●00

Perspectives
000

# 3.3 - Random Subsets: Improved Version (1/2)



- start from end
- unlimited number of proc.
- As Last As Possible
⇒ time threshold:
  first use of unavailable proc.

POTRF
TRSM
SYRK
GEMM

Problem
oooooo

Resolution Methods
ooooooooooo

Evaluation
oooo○●oo

Perspectives
ooo

# 3.3 - Random Subsets: Improved Version (2/2)

Tasks before / after threshold
$\rightarrow$ tiles splitting $+$ reordering

Problem
oooooo

Resolution Methods
ooooooooooo

Evaluation
ooo●oo

Perspectives
ooo

# 3.3 - Random Subsets: Improved Version (2/2)

Tasks before / after threshold
$\rightarrow$ tiles splitting + reordering



splitting + reordering

Problem
○○○○○○

Resolution Methods
○○○○○○○○○○○

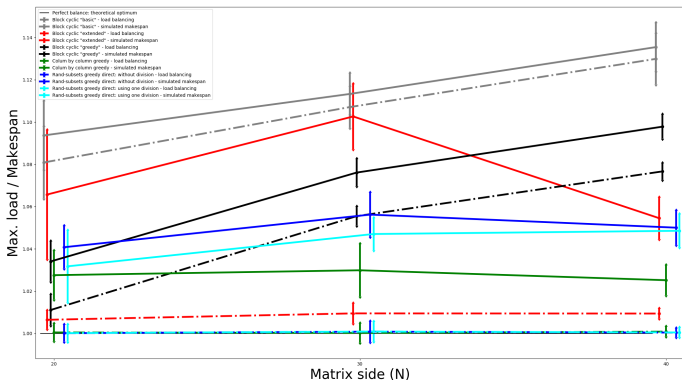Evaluation
○○○○●○

Perspectives
○○○

# 3.4 - Additional Results



Figure: Makespan VS input matrix size
($\alpha = 2$; $N = 20$ to $40$; $\frac{N^2}{P} \approx$ cste)

## 3.4 - Additional Results



Figure: Makespan VS input matrix size
($\alpha = 2$; $N = 60$ to $80$; $\frac{N^2}{P} \approx$ cste)

Problem
oooooo

Resolution Methods
ooooooooooo

Evaluation
ooooo

Perspectives
●oo

# 4. Perspectives

*Inria*

## So far

- efficient methods for load balancing
- many options for strategies

*Inria*

Problem
oooooo

Resolution Methods
ooooooooooo

Evaluation
ooooo

Perspectives
o●o

## So far

- efficient methods for load balancing
- many options for strategies

## Future work

- secure results: larger scale, parameters sets, real data
- improve strategies
- explore new ones: hybrid, relaxed constraints
- dig in scheduling aspect
- other applications / use cases

*Inría*

Problem
oooooo

Resolution Methods
ooooooooooo

Evaluation
ooooo

Perspectives
o●o

## So far
- efficient methods for load balancing
- many options for strategies

## Future work
- secure results: larger scale, parameters sets, real data
- improve strategies
- explore new ones: hybrid, relaxed constraints
- dig in scheduling aspect
- other applications / use cases

## Tools improvement
- evaluation: introduce communications

*Inría*

Problem
oooooo

Resolution Methods
ooooooooooo

Evaluation
ooooo

Perspectives
oo●

# Thank you

*Inria*