

# TASK-BASED PARALLEL PROGRAMMING FOR SCALABLE ALGORITHMS AND ILLUSTRATION WITH MATRIX MULTIPLICATION

JCAD 2024

---

Antoine Jégo

Supervised by Alfredo Buttari, Emmanuel Agullo & Abdou Guermouche

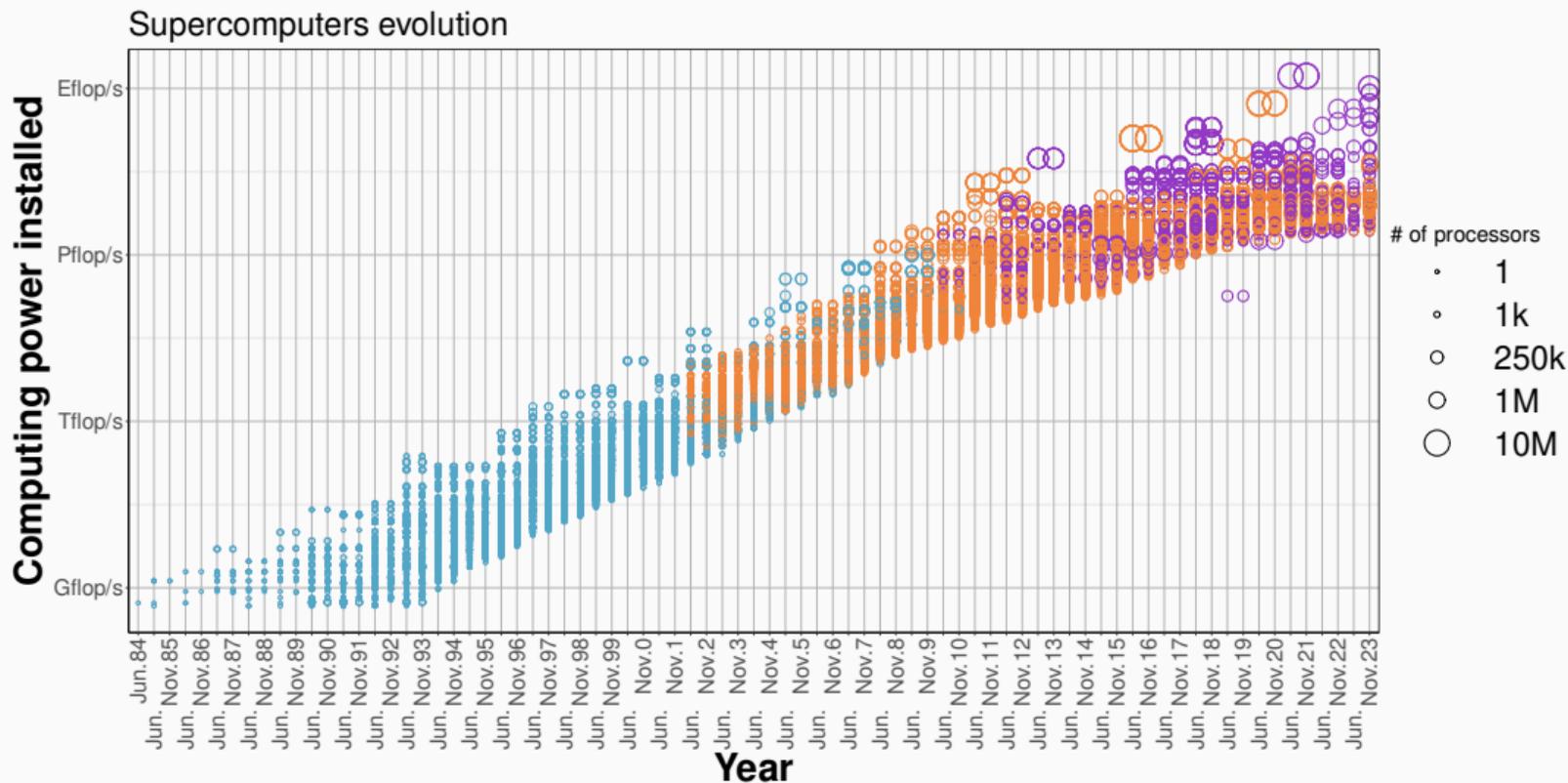


5 novembre 2024



BACKGROUND

---

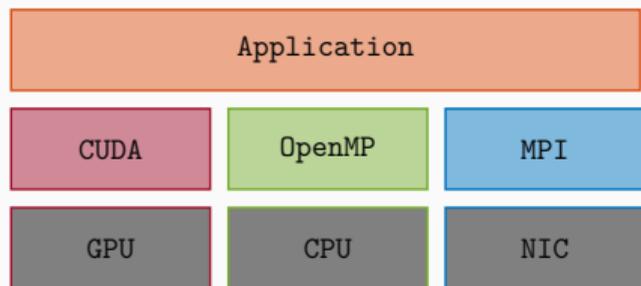


Distributed

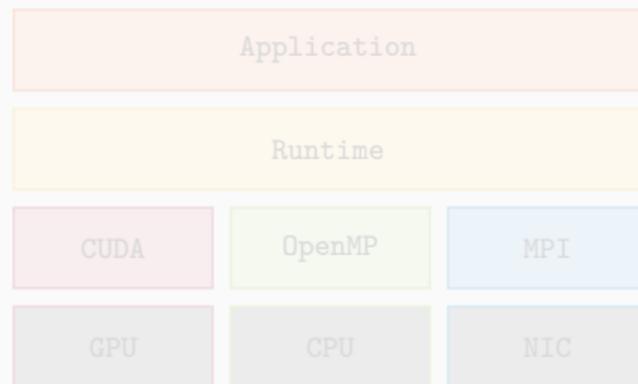
+ Multicore

+ Accelerated

## Classical approach



## Runtime system approach



Complex supercomputers



Advanced scalability techniques

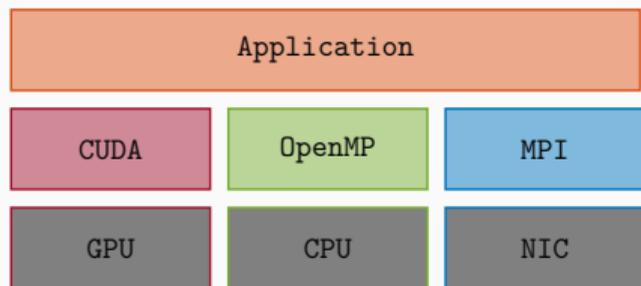


Complex software

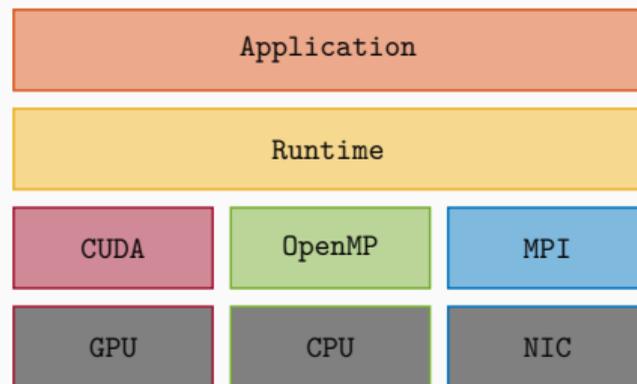
## General-purpose runtimes

- High-level interface
- Data management
- Heterogeneity management

## Classical approach



## Runtime system approach



Complex supercomputers



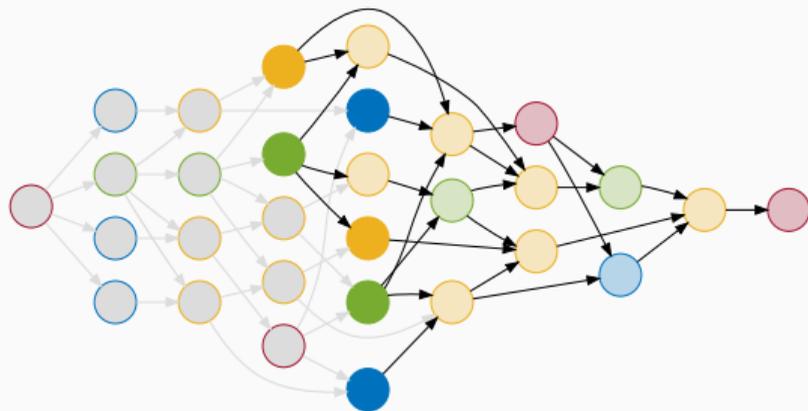
Advanced scalability techniques



Complex software

## General-purpose runtimes

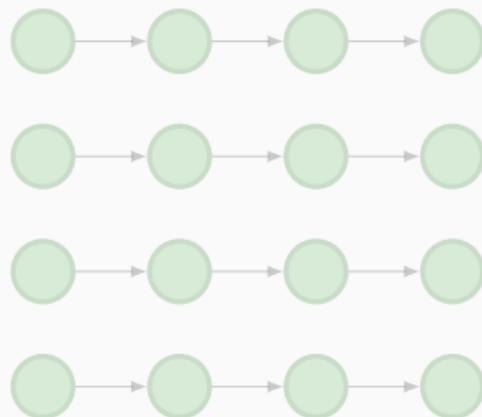
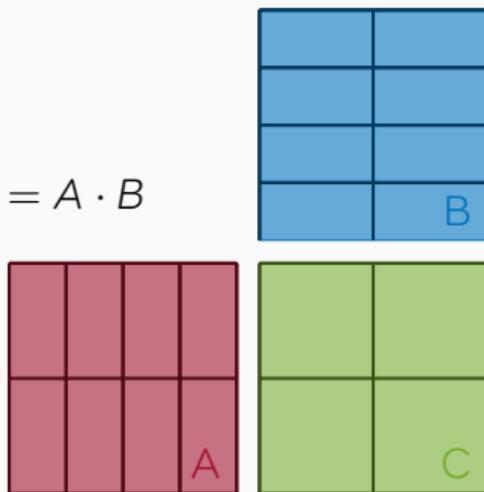
- High-level interface
- Data management
- Heterogeneity management



- ✓ Approach validated on heterogeneous single-node machine
  - SOLHAR among multiple scientific projects
- Multiple ways to program the high-level description of the DAG
  - Parametrized Task Graph (PTG): PaRSEC
  - Sequential Task Flow (**STF**): PaRSEC, StarPU
  - ...

# SEQUENTIAL TASK FLOW MODEL

$$C = A \cdot B$$



DAG of GEMM

---

**Algorithm 1:** Sequential blocked GEMM.

---

```
1 for i = 1...m do
2   for j = 1...n do
3     for l = 1...k do
4       Cij += Ail · Blj
```

---

---

**Algorithm 2:** Shared-memory STF GEMM.

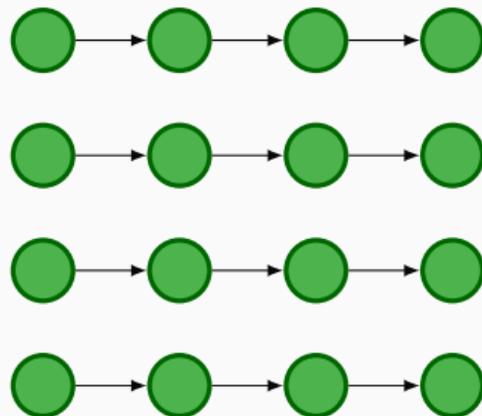
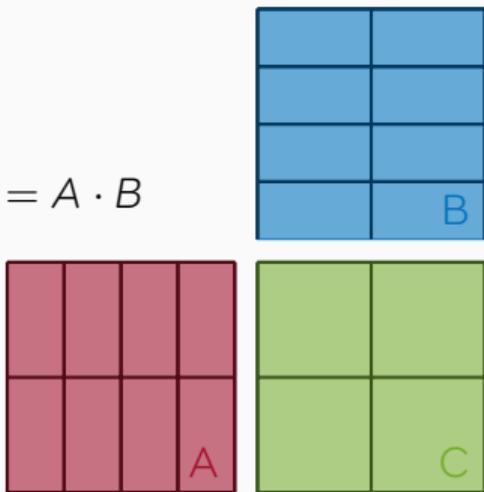
---

```
1 for i = 1...m do
2   for j = 1...n do
3     for l = 1...k do
4       insert_task(gemm, Ail:R, Blj:R, Cij:RW)
5 task_wait()
```

---

# SEQUENTIAL TASK FLOW MODEL

$$C = A \cdot B$$



DAG of GEMM

---

**Algorithm 1:** Sequential blocked GEMM.

---

```
1 for i = 1...m do
2   for j = 1...n do
3     for l = 1...k do
4       Cij += Ail · Blj
```

---

---

**Algorithm 2:** Shared-memory STF GEMM.

---

```
1 for i = 1...m do
2   for j = 1...n do
3     for l = 1...k do
4       insert_task(gemm, Ail:R, Blj:R, Cij:RW)
5 task_wait()
```

---

## Pros of runtime systems:

- Increased portability
- Improved separation of concerns

## Pros of the STF model:

- Productive
- Clear translation from sequential

In practice, scientific computing software use **hand-tuned, application-tailored advanced** communication patterns, data layouts, workload mappings, ...

Is **STF** expressive enough to transparently schedule these techniques and reach existing and new **scalable algorithms** ?



→?

ScaLAPACK GEMM  
 $\mathcal{O}(1k)$  lines of code  
legacy package

---

### Algorithm 3: STF GEMM.

---

```
1 for i = 1...m do
2   for j = 1...n do
3     for l = 1...k do
4       insert_task(gemm, ?)
```

---

**Three portable nested loops**

## Contributions

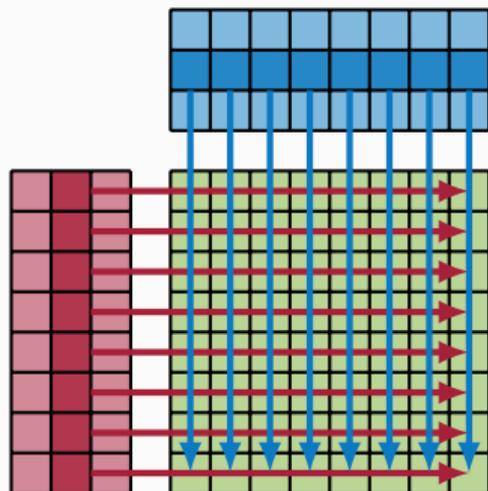
- Extension of the STF programming model
  - Data replication features (reduction, replication)
  - Validated over dense linear algebra operations (GEMM, SYMM, POTRF)

# STF SCALABLE MATRIX-MATRIX MULTIPLICATION

---

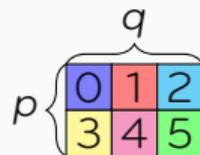
## Scalable **U**niversal **M**atrix **M**ultiplication **A**lgorithms (SUMMA)

**A**-, **B**- and **C**-stationary depending on matrices' dimensions.



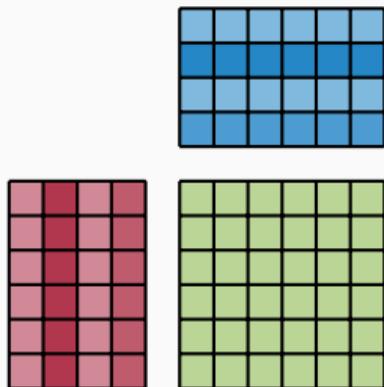
### SUMMA – 2D C-stationary<sup>1</sup>

- 2D block-cyclic distribution



- owner of C computes
- **Communications**
  - row-wise broadcast of A
  - column-wise broadcast of B

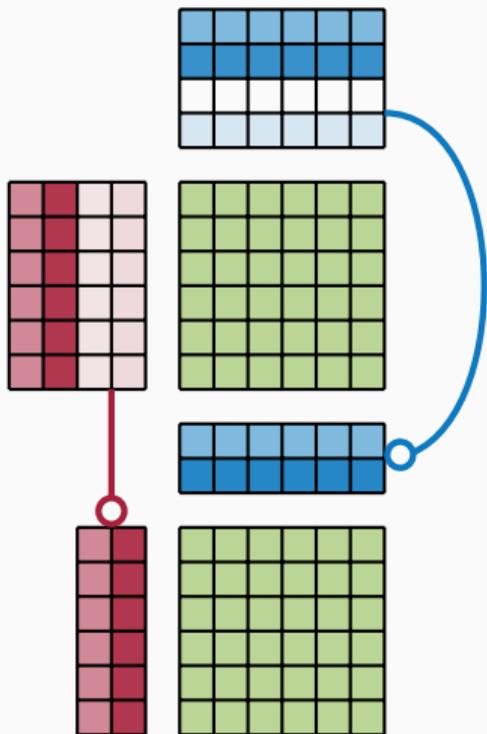
<sup>1</sup>R. v. d. Geijn and Watts, 1997



## SUMMA – 3D C-stationary<sup>2</sup>

- 2D block-cyclic (first layer)
- owner of C's aisle (across layers) compute
- **Communications**
  - row-wise scatter+broadcast of A
  - column-wise scatter+broadcast of B
  - aisle-wise reduce of C

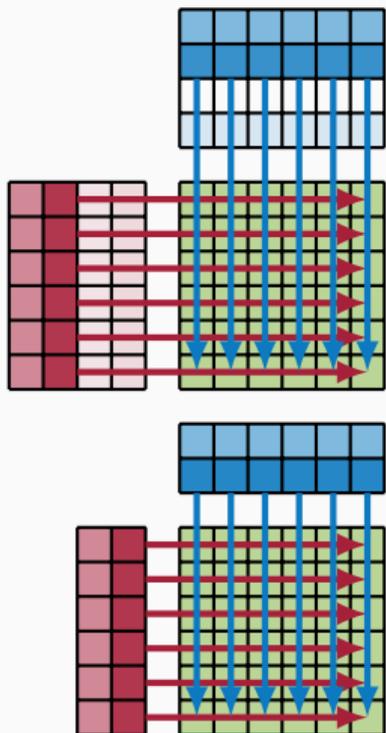
<sup>2</sup>Schatz, R. A. v. d. Geijn, and Poulson, 2016



## SUMMA – 3D C-stationary<sup>2</sup>

- 2D block-cyclic (first layer)
- owner of C's aisle (across layers) compute
- **Communications**
  - row-wise **scatter**+broadcast of A
  - column-wise **scatter**+broadcast of B
  - aisle-wise reduce of C

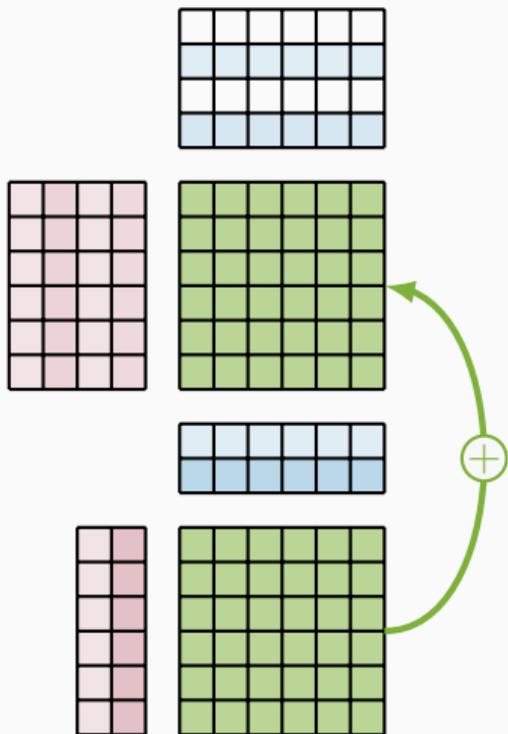
<sup>2</sup>Schatz, R. A. v. d. Geijn, and Poulson, 2016



## SUMMA – 3D C-stationary<sup>2</sup>

- 2D block-cyclic (first layer)
- owner of C's aisle (across layers) compute
- **Communications**
  - row-wise scatter+**broadcast** of A
  - column-wise scatter+**broadcast** of B
  - aisle-wise reduce of C

<sup>2</sup>Schatz, R. A. v. d. Geijn, and Poulson, 2016



## SUMMA – 3D C-stationary<sup>2</sup>

- 2D block-cyclic (first layer)
- owner of C's aisle (across layers) compute
- **Communications**
  - row-wise scatter+broadcast of A
  - column-wise scatter+broadcast of B
  - aisle-wise **reduce** of C

<sup>2</sup>Schatz, R. A. v. d. Geijn, and Poulson, 2016

---

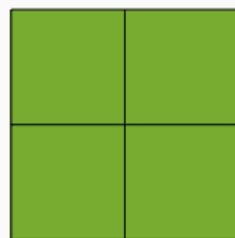
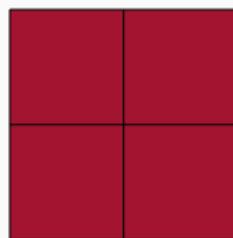
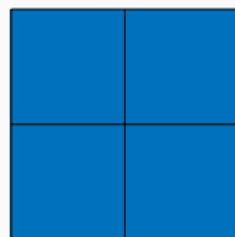
## Algorithm 4: Baseline 2D GEMM.

---

```

1 for  $i = 1 \dots m$  do
2   for  $j = 1 \dots n$  do
3     for  $l = 1 \dots k$  do
4       insert_task (gemm,  $A_{ij}:\mathbb{R}$ ,  $B_{jl}:\mathbb{R}$ ,
                     $C_{ij}:\mathbb{RW}$ )
  
```

---



### Baseline model<sup>a</sup>

1. Data mapping
2. Task mapping **inferred** from data mapping
3. Point-to-point comms. **inferred** from task mapping

---

<sup>a</sup>Agullo, Aumage, Faverge, Furmento, Pruvost, Sergent, and Thibault, 2017.

---

## Algorithm 4: Baseline 2D GEMM.

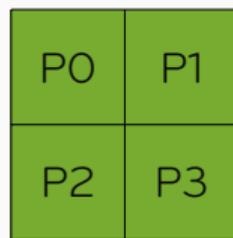
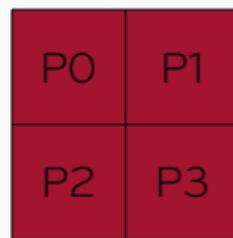
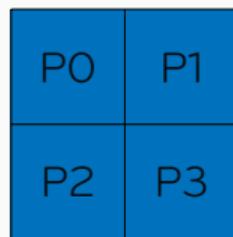
---

```

1 for  $i = 1 \dots m$  do
2   for  $j = 1 \dots n$  do
3     for  $l = 1 \dots k$  do
4       insert_task (gemm,  $A_{ij}:\mathbb{R}$ ,  $B_{lj}:\mathbb{R}$ ,
                     $C_{ij}:\mathbb{RW}$ )

```

---



### Baseline model<sup>a</sup>

1. Data mapping
2. Task mapping **inferred** from data mapping
3. Point-to-point comms. **inferred** from task mapping

---

<sup>a</sup>Agullo, Aumage, Faverge, Furmento, Pruvost, Sergent, and Thibault, 2017.

---

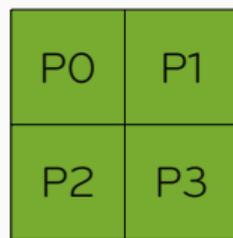
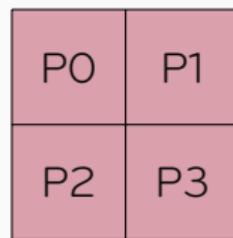
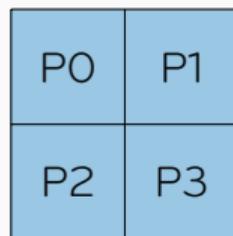
## Algorithm 4: Baseline 2D GEMM.

---

```

1 for  $i = 1 \dots m$  do
2   for  $j = 1 \dots n$  do
3     for  $l = 1 \dots k$  do
4       insert_task (gemm,  $A_{ij}:\mathbb{R}$ ,  $B_{lj}:\mathbb{R}$ ,
                     $C_{ij}:\mathbb{RW}$ )
  
```

---



### Baseline model<sup>a</sup>

1. Data mapping
2. Task mapping **inferred** from data mapping
3. Point-to-point comms. **inferred** from task mapping

---

<sup>a</sup>Agullo, Aumage, Faverge, Furmento, Pruvost, Sergent, and Thibault, 2017.

---

## Algorithm 4: Baseline 2D GEMM.

---

```

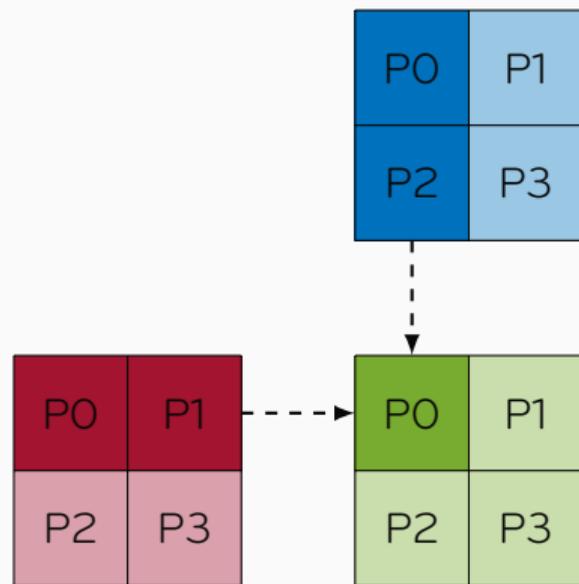
1 for  $i = 1 \dots m$  do
2   for  $j = 1 \dots n$  do
3     for  $l = 1 \dots k$  do
4       insert_task (gemm,  $A_{ij}:\mathbb{R}$ ,  $B_{lj}:\mathbb{R}$ ,
                     $C_{ij}:\mathbb{RW}$ )

```

---

### Baseline model<sup>a</sup>

1. Data mapping
2. Task mapping **inferred** from data mapping
3. Point-to-point comms. **inferred** from task mapping




---

<sup>a</sup>Agullo, Aumage, Faverge, Furmento, Pruvost, Sergent, and Thibault, 2017.

Baseline STF model **cannot productively** provide an expression of SUMMA algorithms

- ✗ Task mapping is bound to data mapping
- ✗ Point-to-point communications patterns are inefficient
- ✗ Reduction patterns would burden the expression

---

**Algorithm 5:** 3D GEMM (I/III).

---

```
1 for  $i = 1 \dots m$  do
2   | for  $j = 1 \dots n$  do
3   |   | for  $l = 1 \dots k$  do
4   |   |   | insert_task (gemm,  $A_{il}:\mathbb{R}$ ,  $B_{lj}:\mathbb{R}$ ,
   |   |   |    $C_{ij}:\mathbb{RW}$  )
```

---

## Advanced model

1. Data mapping
2. Implicit task mapping
3. Point-to-point comms. inferred from task mapping

---

**Algorithm 5:** 3D GEMM (I/III).

---

```
1 for  $i = 1 \dots m$  do
2   | for  $j = 1 \dots n$  do
3   |   | for  $l = 1 \dots k$  do
4   |   |   | rank = map( $i, j, l, \text{stat}, s$ )
5   |   |   | insert_task (gemm,  $A_{ij}:\mathbb{R}, B_{lj}:\mathbb{R},$ 
   |   |   |   |  $C_{ij}:\mathbb{RW}, \text{rank}:\text{ON\_RANK}$ )
```

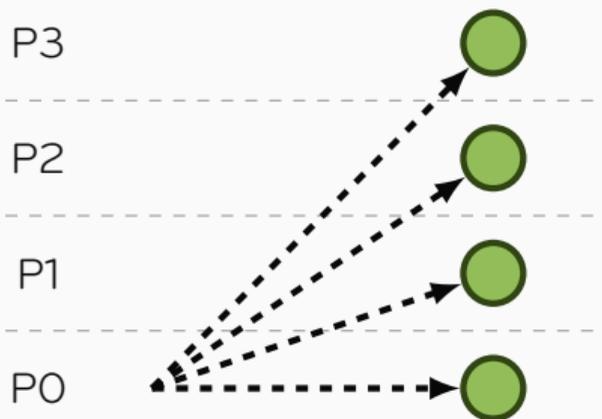
---

## Advanced model

1. Data mapping
2. **Explicit** task mapping
3. Point-to-point comms. inferred from task mapping

## Algorithm 6: 3D GEMM (II/III).

```
1 for  $i = 1 \dots m$  do
2   for  $j = 1 \dots n$  do
3     for  $l = 1 \dots k$  do
4       rank = map( $i, j, l, \text{stat}, s$ )
5       insert_task (gemm,  $A_{ij}:\mathbb{R}$ ,  $B_{lj}:\mathbb{R}$ ,
                     $C_{ij}:\mathbb{RW}$ , rank:ON_RANK)
```

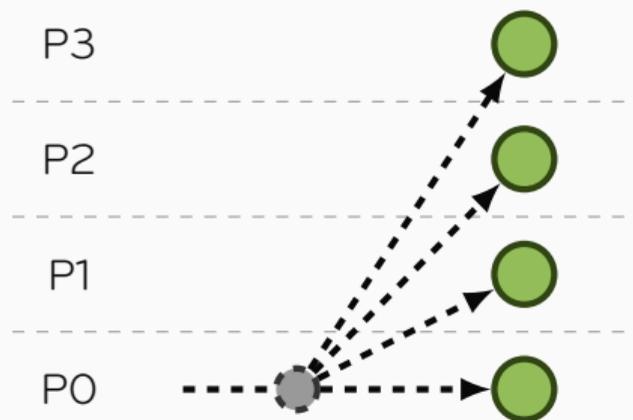


## Advanced model

1. Data mapping
2. **Explicit** task mapping
3. Point-to-point comms. inferred from task mapping

## Algorithm 6: 3D GEMM (II/III).

```
1 for  $i = 1 \dots m$  do
2   for  $j = 1 \dots n$  do
3     for  $l = 1 \dots k$  do
4       rank = map( $i, j, l, \text{stat}, s$ )
5       insert_task (gemm,  $A_{ij}:\mathbb{R}$ ,  $B_{lj}:\mathbb{R}$ ,
                    $C_{ij}:\mathbb{RW}$ , rank:ON_RANK)
```



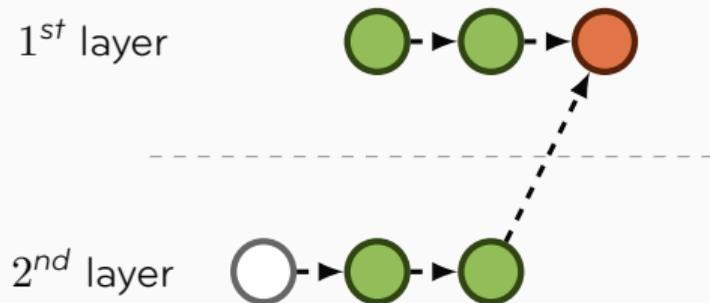
## Advanced model

1. Data mapping
2. **Explicit** task mapping
3. **Collective** communications dynamically detected<sup>a</sup>

<sup>a</sup>Denis, Jeannot, Swartvagher, and Thibault, 2020.

## Algorithm 7: 3D GEMM (III/III).

```
1 submit_initilization_tasks()
2 for  $i = 1 \dots m$  do
3   for  $j = 1 \dots n$  do
4     for  $l = 1 \dots k$  do
5       rank = map( $i, j, l, \text{stat}, s$ )
6       insert_task (gemm,  $A_{ij}:\mathbb{R}$ ,  $B_{lj}:\mathbb{R}$ ,
7                    $C_{ij}^{(\text{rank})}:\mathbb{RW}$ , rank:ON_RANK)
7 submit_reduction_tasks()
```

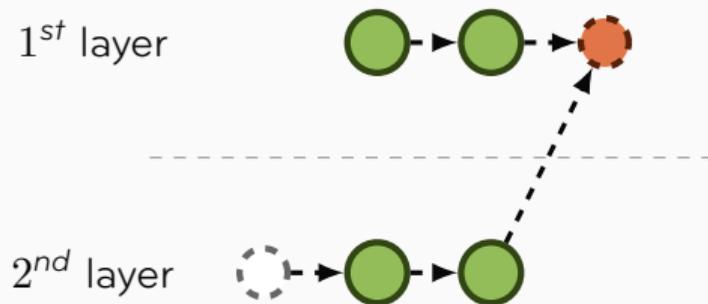


## Advanced model

1. Data mapping
2. **Explicit** task mapping
3. **Collective** communications dynamically detected
4. Explicit reduction pattern

## Algorithm 7: 3D GEMM (III/III).

```
1 bind_methods(C,redux,init)
2 for i = 1 ... m do
3   for j = 1 ... n do
4     for l = 1 ... k do
5       rank = map(i,j,l,stat,s)
6       insert_task (gemm, Aij:R, Blj:R,
                    Cij:RANK_REDUX,
                    rank:ON_RANK)
```



## Advanced model

1. Data mapping
2. **Explicit** task mapping
3. **Collective** communications dynamically detected
4. **Implicit** reduction pattern

**Algorithm 8:** 3D GEMM.

```

1 bind_methods(C,redux,init)
2 for  $i = 1 \dots m$  do
3   for  $j = 1 \dots n$  do
4     for  $l = 1 \dots k$  do
5       rank = map(i,j,l,stat,s)
6       insert_task (gemm,  $A_{ij}:\mathbb{R}$ ,  $B_{lj}:\mathbb{R}$ ,
                     $C_{ij}:\text{RANK\_REDUX}$ ,
                    rank:ON_RANK)

```

## Conclusion I

**Suitably extended** STF model can **portably** express scalable GEMM algorithms.

Is this approach scalable in practice ?

**Advanced model** = STF<sup>♦</sup>

1. Data mapping
2. **Explicit** task mapping
3. **Collective** communications dynamically detected
4. **Implicit** reduction pattern

**Algorithm 8:** 3D GEMM.

```

1 bind_methods(C,redux,init)
2 for  $i = 1 \dots m$  do
3   for  $j = 1 \dots n$  do
4     for  $l = 1 \dots k$  do
5       rank = map(i,j,l,stat,s)
6       insert_task (gemm,  $A_{ij}:\mathbb{R}$ ,  $B_{lj}:\mathbb{R}$ ,
                     $C_{ij}:\text{RANK\_REDUX}$ ,
                    rank:ON_RANK)

```

## Conclusion I

**Suitably extended** STF model can **portably** express scalable GEMM algorithms.

Is this approach scalable in practice ?

**Advanced model** = STF<sup>♦</sup>

1. Data mapping
2. **Explicit** task mapping
3. **Collective** communications dynamically detected
4. **Implicit** reduction pattern

Runtime system used is StarPU – extended with improved reduction features.

STF<sup>♦</sup> GEMM has been implemented in `qr_mumps`.

software package	programming model	2D			3D	GPU
		C	A	B	all	
Chameleon 1.1	STF <sup>X</sup> , StarPU 1.3	✓	X	X	X	✓
SLATE	MPI+OpenMP	✓	✓	X	X	✓
DPlasma	PTG, PaRSEC	✓	✓	✓	X	✓
ScaLAPACK	SPMD, BLACS	✓	✓	✓	X	X
Elemental	SPMD	✓	✓	✓	X	X
this work	STF <sup>♦</sup> , StarPU 1.4	✓	✓	✓	✓	✓

Machines hosted inside the Très Grand Centre de Calcul (TGCC).

A. **Irène-Rome** AMD Rome 7H12 2.6 GHz

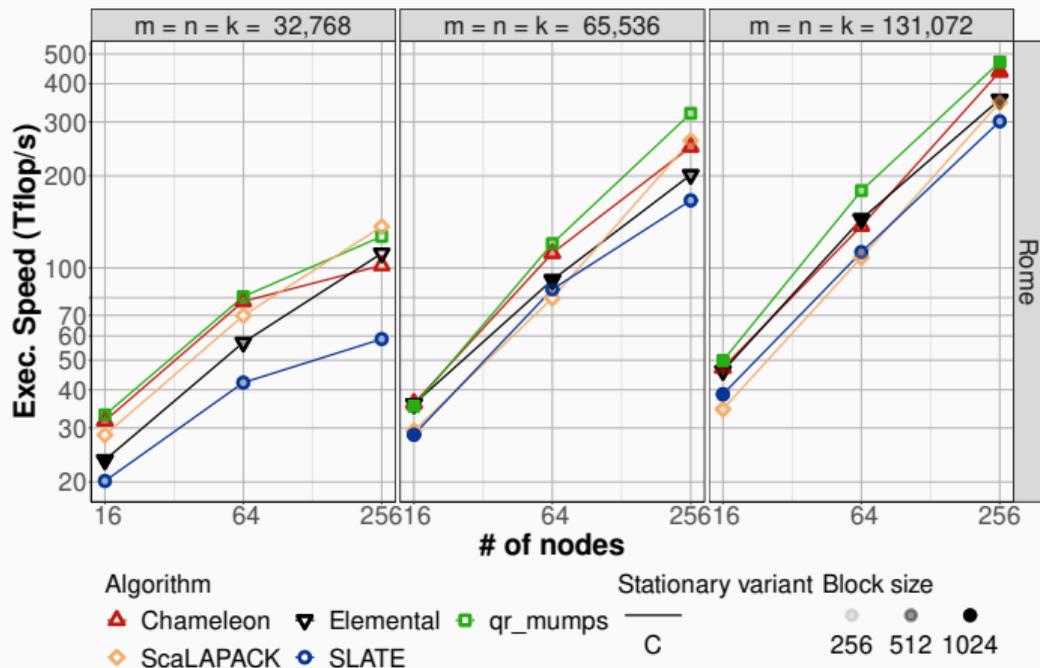
- 2,292 128-cores nodes with Infiniband HDR100 interconnect
- $R_{\text{peak}} = 11.75$  PFlop/s
- 256GB DDR4 RAM per node

B. **Irène-Skylake** Intel Skylake 8168 2.7 GHz

- 1,656 48-cores nodes with Infiniband EDR interconnect
- $R_{\text{peak}} = 6.86$  PFlop/s
- 192GB DDR4 RAM per node

Overall about 1 million CPU hours consumed.

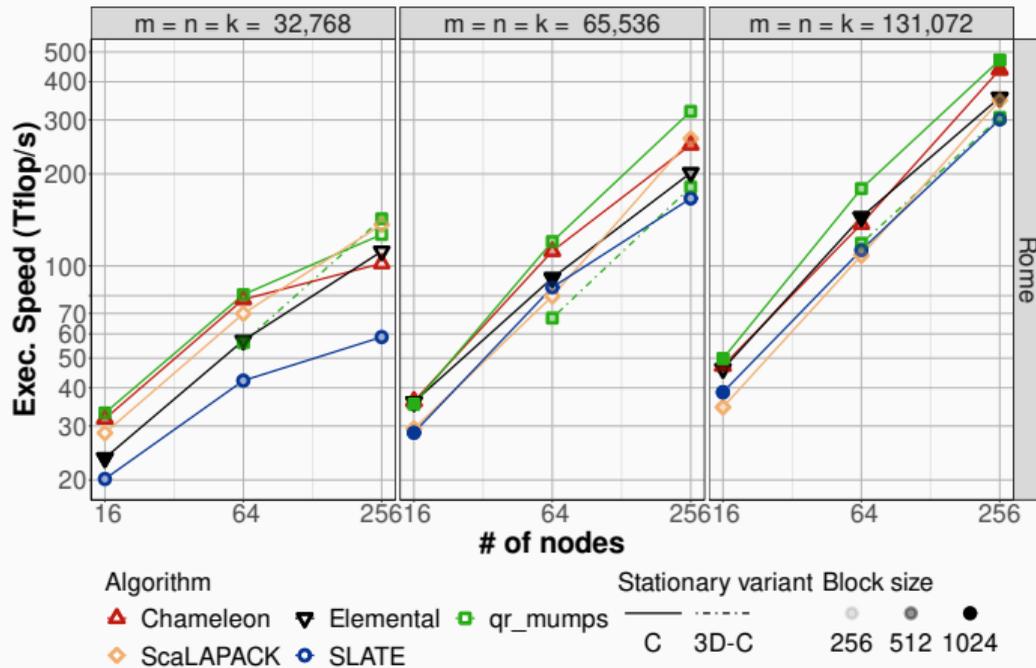
# EXPERIMENTS – PERFORMANCE (C-STAT)



Conclusion II

STF<sup>♦</sup> offers competitive performance w.r.t reference libraries

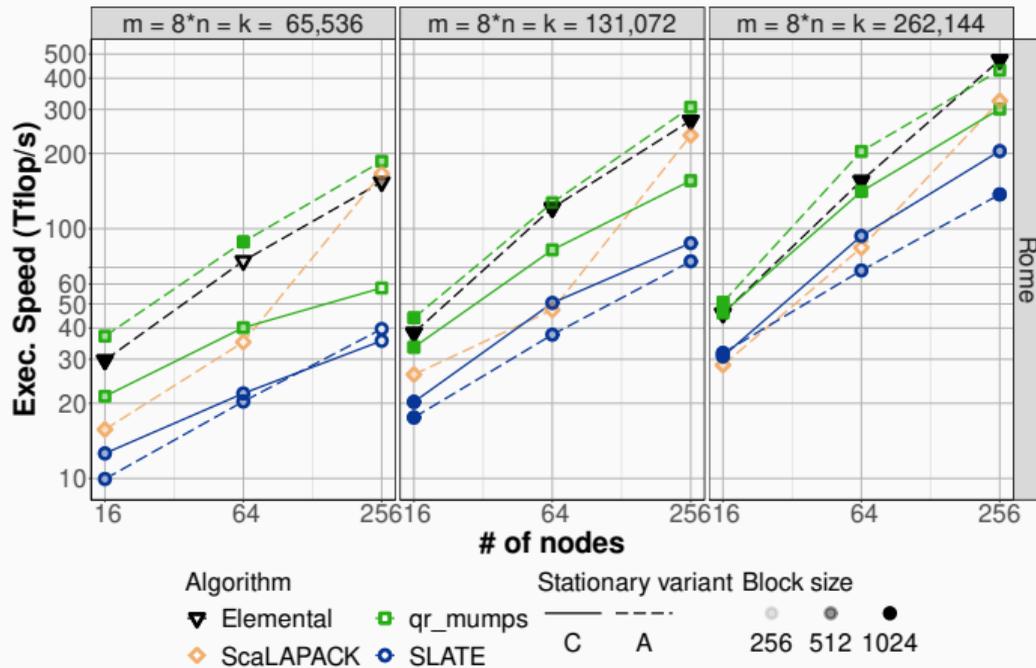
# EXPERIMENTS – PERFORMANCE (C-STAT)



## Conclusion II

STF<sup>♦</sup> offers competitive performance w.r.t reference libraries

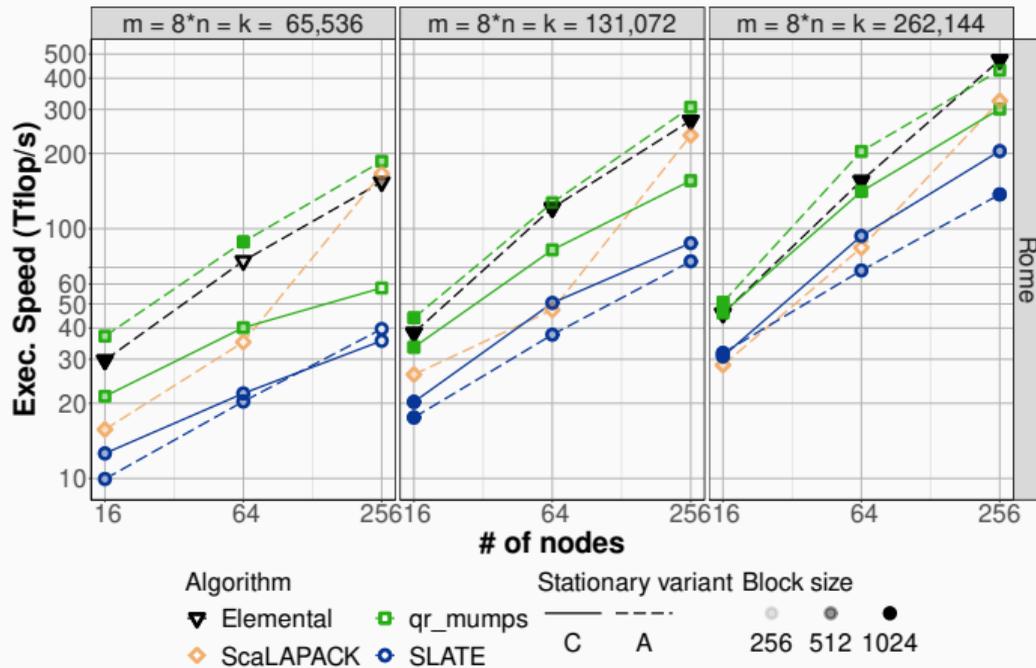
# EXPERIMENTS – PERFORMANCE (A-STAT)



Conclusion III

STF<sup>♦</sup> brings versatility together with performance

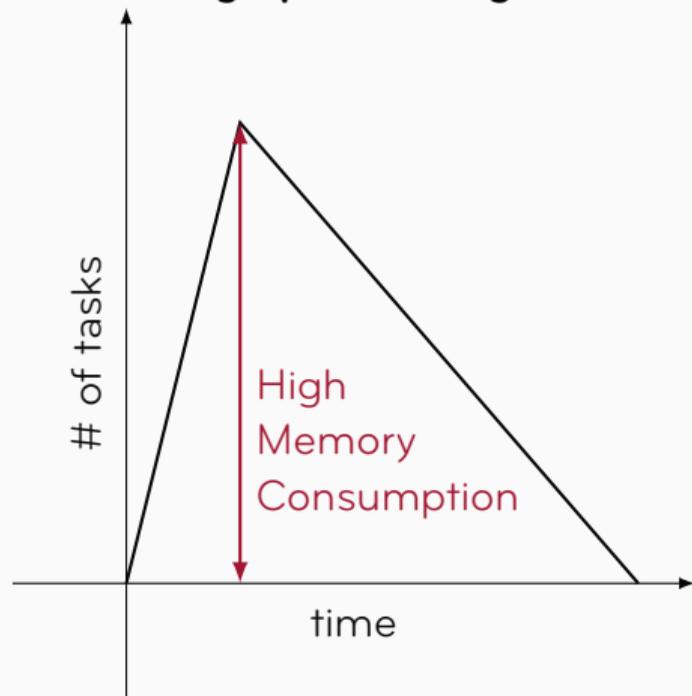
# EXPERIMENTS – PERFORMANCE (A-STAT)



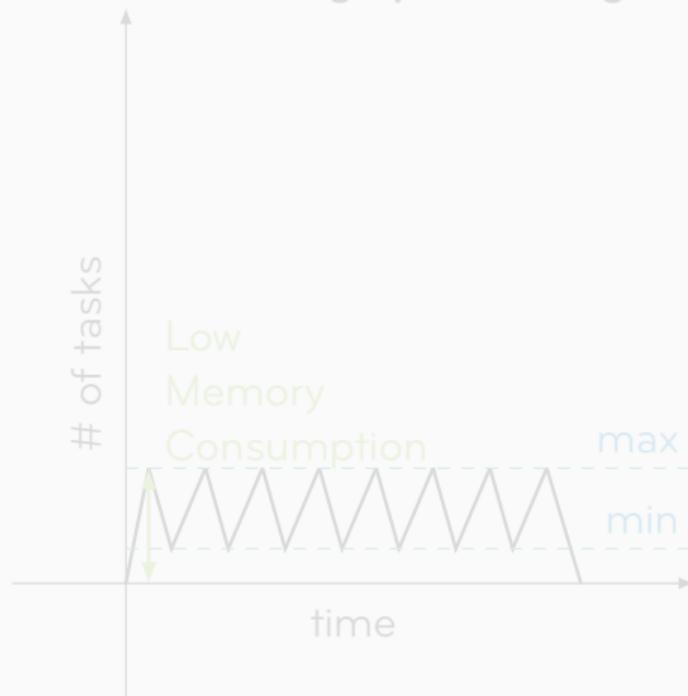
## Conclusion III

STF<sup>♦</sup> brings versatility together with performance

## Unlimited graph unrolling

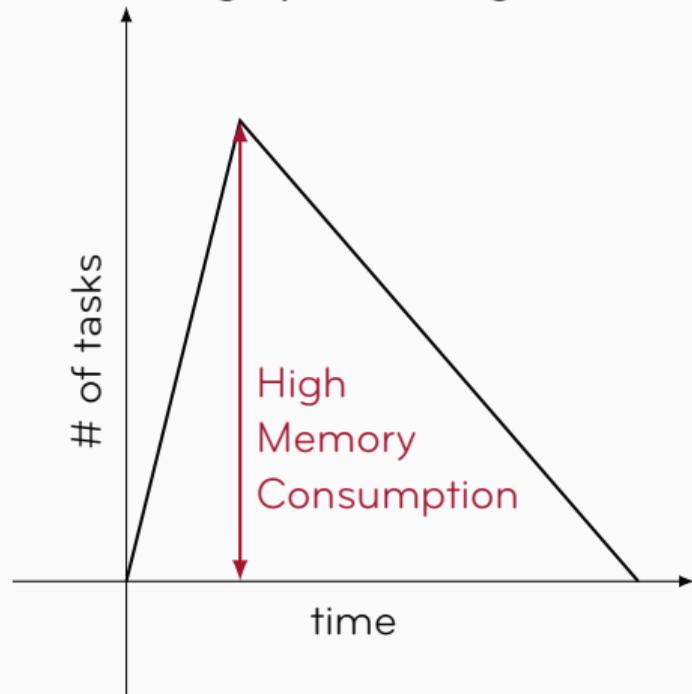


## User-controlled graph unrolling

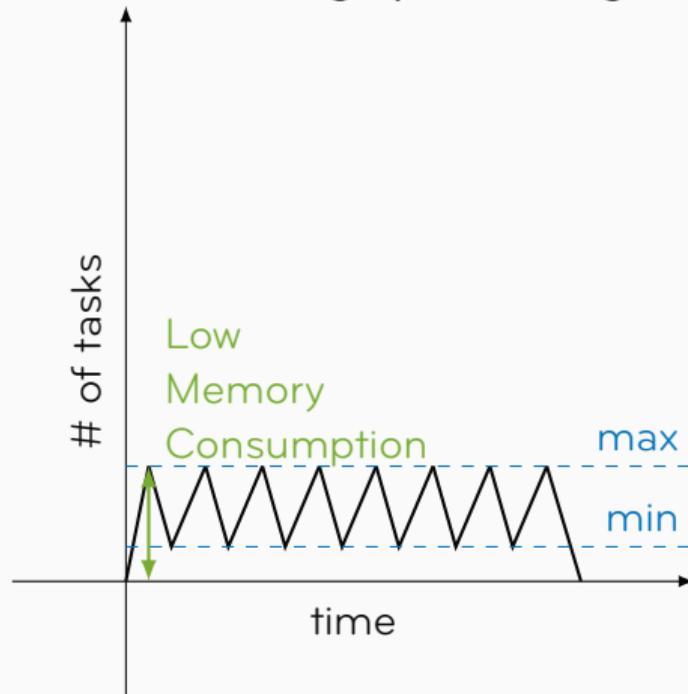


Memory consumption control through limits on tasks' flow.

## Unlimited graph unrolling

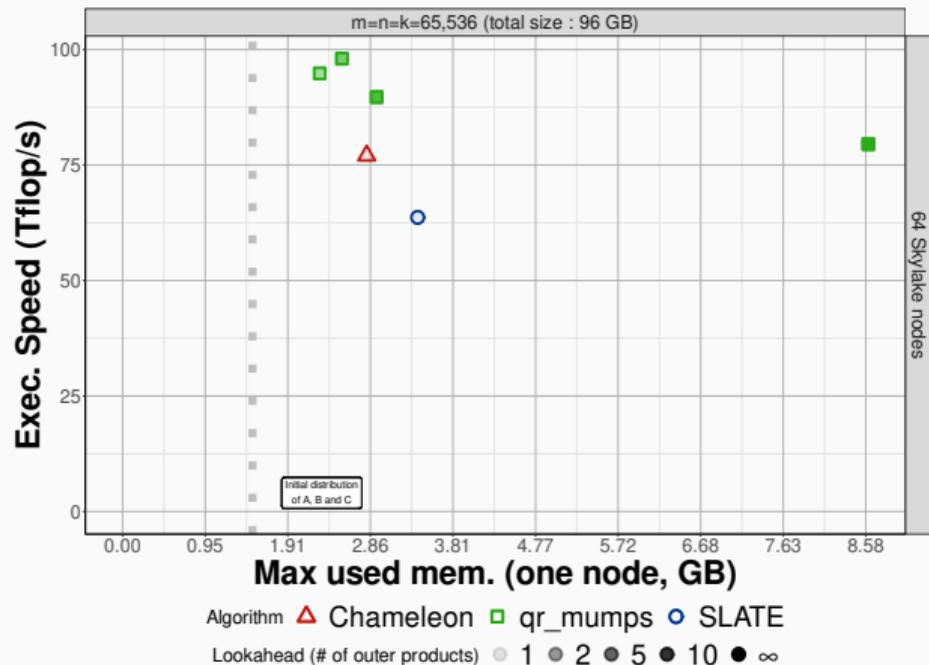


## User-controlled graph unrolling



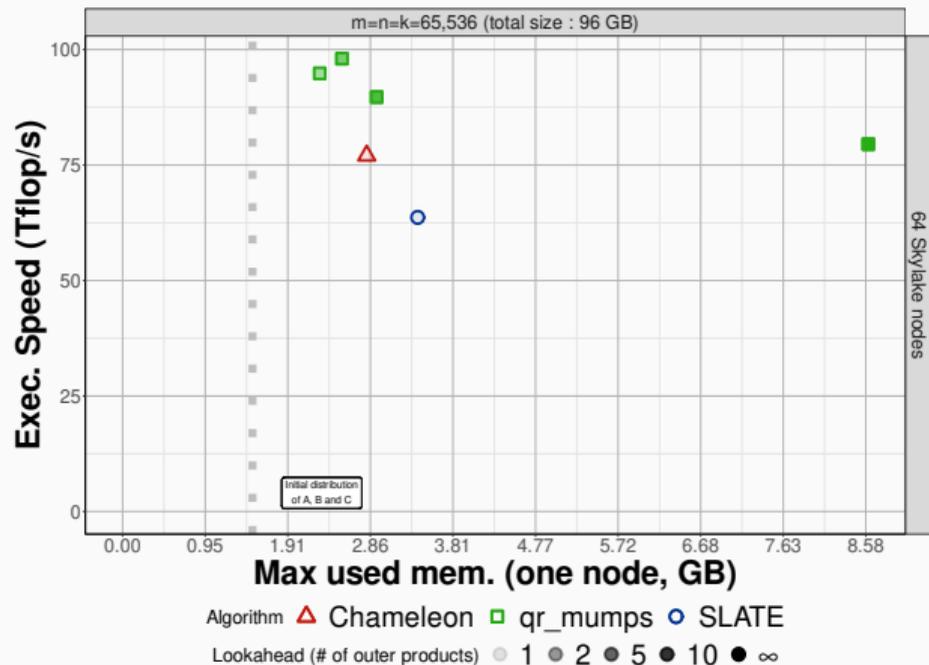
**Memory consumption control** through limits on tasks' flow.

# EXPERIMENTS – MEMORY CONSUMPTION



## Conclusion IV

STF $\diamond$  incorporates tunable indirect method to control memory consumption



## Conclusion IV

STF $\blacklozenge$  incorporates tunable indirect method to control memory consumption

## CONCLUSIONS

---

## STF<sup>♦</sup> assessment over linear algebra operations

- I **Suitably extended** STF model can **portably** express scalable GEMM algorithms
- II STF<sup>♦</sup> offers competitive performance w.r.t reference libraries
- III STF<sup>♦</sup> brings versatility together with performance

## Additional contributions (not discussed here)

- Implementation of scalable SYMM, POTRF algorithms
- Implementation of a communication-avoiding 2D 5-points stencil prototype

## Software contributions

- Implementation of additional features in StarPU
  - Distributed-memory reduction patterns
  - Data write replication (not discussed here, MR 66  and 68 )
- Implementation of dense scalable routines in `qr_mumps`

## Transfer

- Chameleon 1.3
  - Used by INRAE colleagues to analyse biodiversity datasets