









Dynamic Task Graph Adaptation with Recursive Tasks

Exa-Soft

Nathalie Furmento, Abdou Guermouche, Gwenolé Lucas, <u>Thomas Morin</u>, Samuel Thibault, Pierre-André Wacrenier

September 2025

















Introduction

#### Introduction - Context

### Task-based Programming

- Motivations:
  - Portable frameworks.
  - Exploit complex architectures.
- Applications: Directed Acyclic Graph (DAG).
- Runtime systems: scheduling, data management, communications, ...





#### Introduction - Context

### Task-based Programming

- Motivations:
  - Portable frameworks.
  - Exploit complex architectures.
- Applications: Directed Acyclic Graph (DAG).
- Runtime systems: scheduling, data management, communications, ...



### STF: Sequential Task Flow

- Dependencies:
  - Automatically inferred.
  - o Order of submission.

```
F(a) submit(F, a:RW)
G(a, b) submit(G, a:R, b:RW)
H(a, c) submit(H, a:R, c:RW)
wait_tasks_completion()
```









#### Submission

- Overhead: large number of non-ready tasks.
- Bottleneck: sequential insertion.
- Adaptability ? static task graphs.





#### Submission

- Overhead: large number of non-ready tasks.
- Bottleneck: sequential insertion.
- Adaptability ? static task graphs.

⇒ How to create more dynamic task-graphs?





#### **Submission**

- Overhead: large number of non-ready tasks.
- Bottleneck: sequential insertion.
- Adaptability ? static task graphs.

 $\Rightarrow$  How to create more dynamic task-graphs ?  $\Rightarrow$  Recursive tasks graphs !





#### Submission

- Overhead: large number of non-ready tasks.
- Bottleneck: sequential insertion.
- Adaptability ? static task graphs.
- $\Rightarrow$  How to create more dynamic task-graphs ?  $\Rightarrow$  Recursive tasks graphs !

### Granularity

- GPUs versus CPUs.
- Lack of parallelism versus Steady State.
- $\Rightarrow$  Steering granularity dynamically ?







#### Overview

- StarPU's Recursive Tasks.
- Introduction of the Splitter component.
- The heterogeneous case:
  - A first policy that relies on Linear Programming.
  - o Improving this policy with a greedy algorithm.









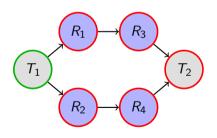






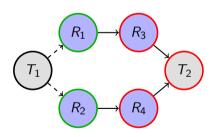
Recursive Tasks





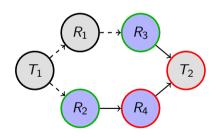


• Recursive task execution:

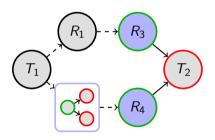




- Recursive task execution:
  - Remain regular task.

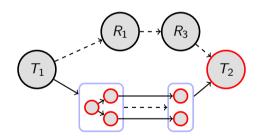


- Recursive task execution:
  - Remain regular task.
  - Insert a subgraph: split.

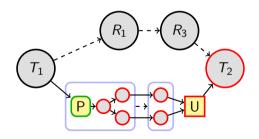




- Recursive task execution:
  - Remain regular task.
  - Insert a subgraph: split.



- Recursive task execution:
  - Remain regular task.
  - o Insert a subgraph: split.













Dynamic task graph adaptation

# Dynamic task graph adaptation: splitting tasks

When do we choose to split task?

Which task should we split?





## Dynamic task graph adaptation: splitting tasks

When do we choose to split task?

Submission, execution, ...

Which task should we split?



## Dynamic task graph adaptation: splitting tasks

When do we choose to split task?

Submission, execution, ...

Which task should we split?

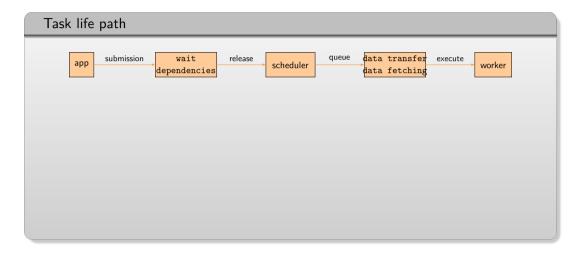
Efficiency, Completion Time, Current Parallelism







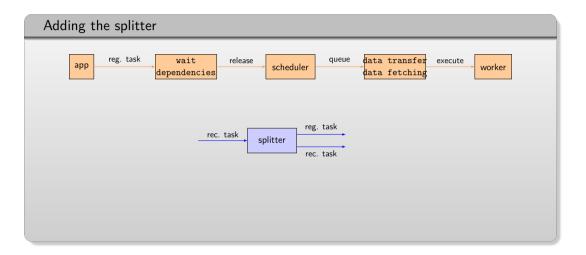
## When do we choose to split tasks







## When do we choose to split tasks

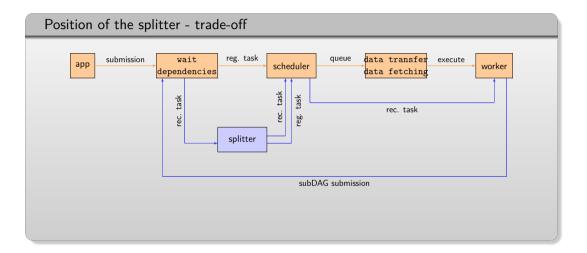








## When do we choose to split tasks



















Heterogeneous

#### General example:

 Heterogeneous platform with 2 GPU Nvidia A100 and 2 AMD Zen3 CPU, each with 32 cores.

#### General example:

- Heterogeneous platform with 2 GPU Nvidia A100 and 2 AMD Zen3 CPU, each with 32 cores.
- Performing a tiled Cholesky factorization.

Kernel	1 gpu 1 core	
GEMM	380	

#### General example:

- Heterogeneous platform with 2 GPU Nvidia A100 and 2 AMD Zen3 CPU, each with 32 cores.
- Performing a tiled Cholesky factorization.

Kernel	$\frac{1 \ gpu}{1 \ core}$	1 gpu StarPU 64 cores (accel.)
GEMM	380	9.4 (40.4)

#### General example:

- Heterogeneous platform with 2 GPU Nvidia A100 and 2 AMD Zen3 CPU, each with 32 cores.
- Performing a tiled Cholesky factorization.

Kernel	1 gpu 1 core	1 gpu StarPU 64 cores (accel.)
GEMM	380	9.4 (40.4)
TRSM	307	6.5 (45.8)
SYRK	343	11.7 (29.3)

#### General example:

- Heterogeneous platform with 2 GPU Nvidia A100 and 2 AMD Zen3 CPU, each with 32 cores.
- Performing a tiled Cholesky factorization.

Kernel	1 gpu	
	1 core	1 gpu StarPU 64 cores (accel.)
		StarPU 64 cores (accel.)
GEMM	380	
		9.4 (40.4)
		9.4 (40.4)
TRSM	307	
		6.5 (45.8)
		0.5 (15.0)
SYRK	343	
	0.10	11.7 (29.3)
		11.1 (23.3)

ullet Q: Steady state  $\Rightarrow$  Which tasks, and how many should be divided ?





 $\bullet$  A task becomes ready  $\Rightarrow$  Recorded by the Splitter.

- A task becomes ready ⇒ Recorded by the Splitter.
- Each 50-task interval ⇒ A Linear Program is called.

- A task becomes ready ⇒ Recorded by the Splitter.
- Each 50-task interval ⇒ A Linear Program is called.
- Linear Program's aim: minimizing execution time by balancing load between the PUs.

- A task becomes ready ⇒ Recorded by the Splitter.
- Each 50-task interval ⇒ A Linear Program is called.
- Linear Program's aim: minimizing execution time by balancing load between the PUs.
- Task Splitting for PUs is required by ensuring minimal parallelism for each type of PU.

- A task becomes ready ⇒ Recorded by the Splitter.
- Each 50-task interval ⇒ A Linear Program is called.
- Linear Program's aim: minimizing execution time by balancing load between the PUs.
- Task Splitting for PUs is required by ensuring minimal parallelism for each type of PU.
- The LP provides ratio ⇒ used by the Splitter.

## Which task do we split - Heterogeneous case - Linear Program

**Parameters** 

 $\mathcal{T}, N_{t,I}^{tot}$   $t \in \mathcal{T}$  Set of task types, Number of task not split of type t at level I.

 $\mathcal{R}$  ,  $\mathcal{R}^{\mathcal{U}}$   $u \in \mathcal{R}$  Set of processing unit types, Number of PU of type u.

C Maximum level of recursion.

 $\mathit{MinN}_{\mathit{U}}$   $u \in \mathcal{R}$  Minimal wanted number of tasks on each PU of type u

Variables

exT Total execution time.  $Ns_{i}^{t} \qquad t \in \mathcal{T}, I < \mathcal{L} \qquad \qquad \text{Number of split task of type}$ 

t and level I.

 $Ne_{I,\,u}^t$   $t\in\mathcal{T},\,l\leq\mathcal{L},\,u\in\mathcal{R}$  Number of task of type t

Minimize

exT

#### Subject to

Task number splitting

$$\sum_{u \in \mathcal{R}} \mathit{Ne}_{l,u}^{t} + \mathit{Ns}_{l}^{t} - \sum_{p \; \in \; \mathit{par}(t)} \mathit{nch}_{p,l}^{t} \cdot \mathit{Ns}_{l-1}^{p} \geq \mathit{N}_{t,l}^{tot} \quad \ t \in \mathcal{T}, \, l \leq \mathcal{L}$$

No last-level splitting

$$\textit{Ns}_{\mathcal{L}}^t = 0 \qquad \forall \ t \in \mathcal{T}$$

Completion time when executing tasks

$$\sum_{\substack{t \in \mathcal{T} \\ 0 \le l \le C}} Ne_{l,u}^t \cdot Ex_{t,l}^u - R^u \cdot exT \le 0 \qquad u \in \mathcal{R}$$

Minimal number of tasks on PU type.

$$\sum_{\substack{t \in \mathcal{T} \\ 0 \leq l \leq \mathcal{L}}} \mathit{Ne}_{l,\,u}^{t} - \mathit{R}^{u} \cdot \mathit{MinN}_{u} \leq 0 \qquad \quad u \in \mathcal{R}$$



#### Benchmarks - Heterogeneous

# The tests were run on PlaFRIM's sirocco nodes:

- 2x 32-core AMD Zen3 EPYC 7513 @ 2.6 GHz
- 2x NVIDIA A100 (40GB)
- 512 GB (8 GB/core) (@3200 MHz)
- Scheduler : Deque Model Data-Aware Ready (DMDAR)

#### Benchmarks - Heterogeneous

#### The tests were run on PlaFRIM's sirocco nodes:

- 2x 32-core AMD Zen3 FPYC 7513 @ 2 6 GHz
- 2x NVIDIA A100 (40GB)
- 512 GB (8 GB/core) (@3200 MHz)
- Scheduler : Deque Model Data-Aware Ready (DMDAR)

#### Tile sizes choosen:

- 3840 : "big" : the most efficient.
- 480 : "small": parallelism, for CPU cores.
- 1920, 2560 : "mid": trade-off.



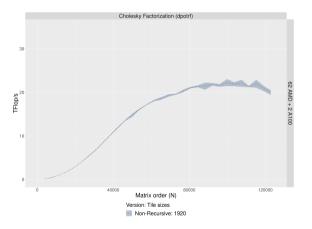


Figure 1: Performance comparison between different versions for Cholesky Factorization.



15

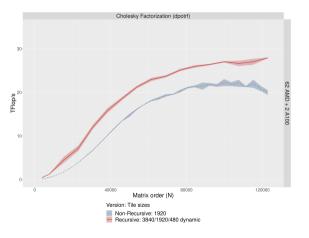


Figure 1: Performance comparison between different versions for Cholesky Factorization.



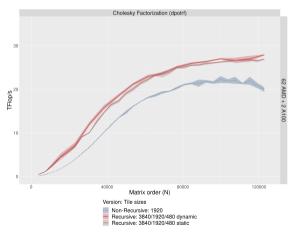


Figure 1: Performance comparison between different versions for Cholesky Factorization.



15

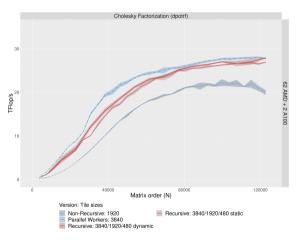


Figure 1: Performance comparison between different versions for Cholesky Factorization.



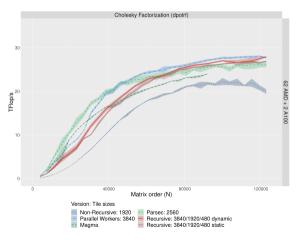


Figure 1: Performance comparison between different versions for Cholesky Factorization.



### Benchmarks - Heterogeneous - without piv

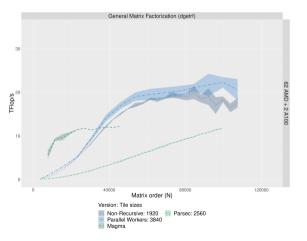


Figure 2: Performance comparison between different versions for LU Factorization without piv.





### Benchmarks - Heterogeneous - without piv

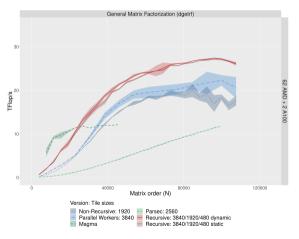


Figure 2: Performance comparison between different versions for LU Factorization without piv.





### Benchmarks - Heterogeneous - All

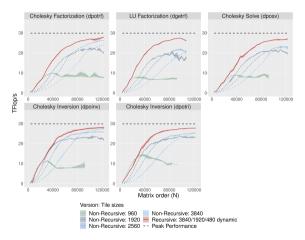


Figure 3: Performance comparison between different kernels and versions.



Exa-Soft, Saclay Thomas Morin

17

### Benchmarks - Heterogeneous Cholesky trace

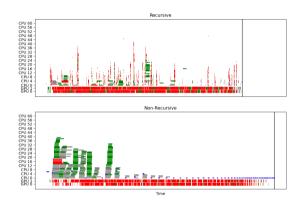


Figure 4: Comparison of traces between Recursive (3840/1920/480) and Non-Recursive (1920) versions for a Cholesky.





- Summary:
  - $\circ$  Automatic + Generic  $\Rightarrow$  Granularity dynamically steered.

- Summary:
  - Automatic + Generic ⇒ Granularity dynamically steered.
  - Dense linear algebra problems:
    - Accross non-recursive Chameleon: Improvements from 5 to 30 %.

- Summary:
  - Automatic + Generic ⇒ Granularity dynamically steered.
  - Dense linear algebra problems:
    - Accross non-recursive Chameleon: Improvements from 5 to 30 %.
    - Accross state-of-the-arts versions: on-par results in a portable way.

- Summary:
  - Automatic + Generic ⇒ Granularity dynamically steered.
  - Dense linear algebra problems:
    - Accross non-recursive Chameleon: Improvements from 5 to 30 %.
    - Accross state-of-the-arts versions: on-par results in a portable way.
- Advantages:
  - o Generic.



- Summary:
  - Automatic + Generic ⇒ Granularity dynamically steered.
  - Dense linear algebra problems:
    - Accross non-recursive Chameleon: Improvements from 5 to 30 %.
    - Accross state-of-the-arts versions: on-par results in a portable way.
- Advantages:
  - o Generic.
  - Automatic.



- Summary:
  - Automatic + Generic ⇒ Granularity dynamically steered.
  - Dense linear algebra problems:
    - Accross non-recursive Chameleon: Improvements from 5 to 30 %.
    - Accross state-of-the-arts versions: on-par results in a portable way.
- Advantages:
  - o Generic.
  - Automatic.
  - Requirements: Recursive implementations + Performance models.

Exa-Soft, Saclay Thomas Morin

19

- Summary:
  - Automatic + Generic ⇒ Granularity dynamically steered.
  - Dense linear algebra problems:
    - Accross non-recursive Chameleon: Improvements from 5 to 30 %.
    - Accross state-of-the-arts versions: on-par results in a portable way.
- Advantages:
  - o Generic.
  - Automatic.
  - Requirements: Recursive implementations + Performance models.
  - At least same performance as non-recursive versions.



19

- Summary:
  - Automatic + Generic ⇒ Granularity dynamically steered.
  - Dense linear algebra problems:
    - Accross non-recursive Chameleon: Improvements from 5 to 30 %.
    - Accross state-of-the-arts versions: on-par results in a portable way.
- Advantages:
  - o Generic.
  - Automatic.
  - Requirements: Recursive implementations + Performance models.
  - At least same performance as non-recursive versions.
- Inconvenients:
  - Decisions taken by looking at past. Cf. call interval.





- Summary:
  - $\circ$  Automatic + Generic  $\Rightarrow$  Granularity dynamically steered.
  - Dense linear algebra problems:
    - Accross non-recursive Chameleon: Improvements from 5 to 30 %.
    - Accross state-of-the-arts versions: on-par results in a portable way.
- Advantages:
  - o Generic.
  - Automatic.
  - Requirements: Recursive implementations + Performance models.
  - At least same performance as non-recursive versions.
- Inconvenients:
  - Decisions taken by looking at past. Cf. call interval.
  - Opaque decision.





- Summary:
  - $\circ$  Automatic + Generic  $\Rightarrow$  Granularity dynamically steered.
  - Dense linear algebra problems:
    - Accross non-recursive Chameleon: Improvements from 5 to 30 %.
    - Accross state-of-the-arts versions: on-par results in a portable way.
- Advantages:
  - o Generic.
  - Automatic.
  - Requirements: Recursive implementations + Performance models.
  - At least same performance as non-recursive versions.
- Inconvenients:
  - Decisions taken by looking at past. Cf. call interval.
  - o Opaque decision.
  - CPU usage can be increased.





#### • Summary:

- Automatic + Generic ⇒ Granularity dynamically steered.
- Dense linear algebra problems:
  - Accross non-recursive Chameleon: Improvements from 5 to 30 %.
  - Accross state-of-the-arts versions: on-par results in a portable way.

#### Advantages:

- o Generic.
- Automatic.
- Requirements: Recursive implementations + Performance models.
- At least same performance as non-recursive versions.

#### Inconvenients:

- Decisions taken by looking at past. Cf. call interval.
- o Opaque decision.
- CPU usage can be increased.
- Turning a global view into a decision with local constraints.





















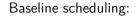
Baseline scheduling:





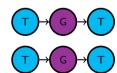








#### Potential split:







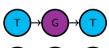


Baseline scheduling:



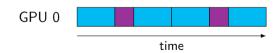


#### Potential split:





#### Potential scheduling:









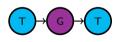
Baseline scheduling:

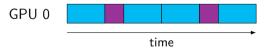


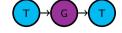


Potential split:

Potential scheduling:







Not interesting: longer than baseline.





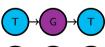


#### Baseline scheduling:





#### Potential split:





#### Potential scheduling:







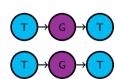




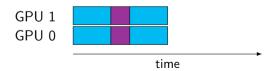
Baseline scheduling:



#### Potential split:



#### Potential scheduling:



gpu-time increased, end-time decreased.





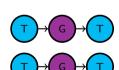




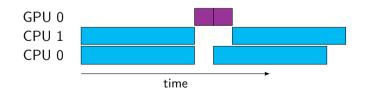
#### Baseline scheduling:



#### Potential split:



#### Potential scheduling:

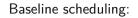






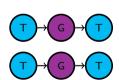




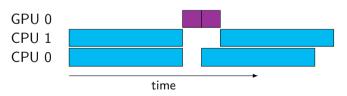




#### Potential split:



#### Potential scheduling:



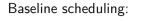
End-time increased, GPU time decreased.





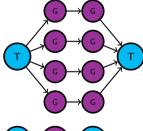








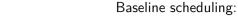
Potential split:







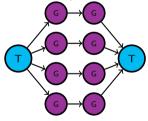


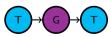






#### Potential split:





Idea: for each type of task, create potential interesting splittings and schedulings:

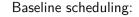
- Decrease GPU time.
- Occupy CPU cores.





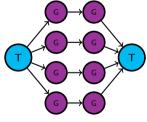








Potential split:



Idea: for each type of task, create potential interesting splittings and schedulings:

- Decrease GPU time.
- Occupy CPU cores.

Record, for each scheduling:

- Sequential GPU time.
- Global finishing time.
- Mean of CPU cores used.







# A new algorithm to improve previous limitations

• LP  $\Rightarrow$  global decision **vs** greedy  $\Rightarrow$  incremental decision with a **global target**.

- LP  $\Rightarrow$  global decision **vs** greedy  $\Rightarrow$  incremental decision with a **global target**.
- Global target: when possible, splitting a task to occupy CPU cores and constraint the scheduling.

- LP  $\Rightarrow$  global decision **vs** greedy  $\Rightarrow$  incremental decision with a **global target**.
- Global target: when possible, splitting a task to occupy CPU cores and constraint the scheduling.
- At start, for each task generation of "all" recursive versions and "all" schedulings.

- LP  $\Rightarrow$  global decision **vs** greedy  $\Rightarrow$  incremental decision with a **global target**.
- Global target: when possible, splitting a task to occupy CPU cores and constraint the scheduling.
- At start, for each task generation of "all" recursive versions and "all" schedulings.
  - All interesting recursive versions and schedulings.



- LP  $\Rightarrow$  global decision **vs** greedy  $\Rightarrow$  incremental decision with a **global target**.
- Global target: when possible, splitting a task to occupy CPU cores and constraint the scheduling.
- At start, for each task generation of "all" recursive versions and "all" schedulings.
  - All **interesting** recursive versions and schedulings.
- When a task is about to be pushed, look at GPU end.



- LP  $\Rightarrow$  global decision **vs** greedy  $\Rightarrow$  incremental decision with a **global target**.
- Global target: when possible, splitting a task to occupy CPU cores and constraint the scheduling.
- At start, for each task generation of "all" recursive versions and "all" schedulings.
  - All **interesting** recursive versions and schedulings.
- When a task is about to be pushed, look at GPU end.
  - If no task can be split with a part executed on CPU cores, push it to GPUs.



- LP  $\Rightarrow$  global decision **vs** greedy  $\Rightarrow$  incremental decision with a **global target**.
- Global target: when possible, splitting a task to occupy CPU cores and constraint the scheduling.
- At start, for each task generation of "all" recursive versions and "all" schedulings.
  - All interesting recursive versions and schedulings.
- When a task is about to be pushed, look at GPU end.
  - If no task can be split with a part executed on CPU cores, push it to GPUs.
  - Else, take the less accelerated task,



Exa-Soft, Saclay Thomas Morin

- LP  $\Rightarrow$  global decision **vs** greedy  $\Rightarrow$  incremental decision with a **global target**.
- Global target: when possible, splitting a task to occupy CPU cores and constraint the scheduling.
- At start, for each task generation of "all" recursive versions and "all" schedulings.
  - All interesting recursive versions and schedulings.
- When a task is about to be pushed, look at GPU end.
  - If no task can be split with a part executed on CPU cores, push it to GPUs.
  - Else, take the less accelerated task,



Exa-Soft, Saclay Thomas Morin

- LP  $\Rightarrow$  global decision **vs** greedy  $\Rightarrow$  incremental decision with a **global target**.
- Global target: when possible, splitting a task to occupy CPU cores and constraint the scheduling.
- At start, for each task generation of "all" recursive versions and "all" schedulings.
  - All interesting recursive versions and schedulings.
- When a task is about to be pushed, look at GPU end.
  - If no task can be split with a part executed on CPU cores, push it to GPUs.
  - o Else, take the less accelerated task, and search for a recursive scheduling that



Thomas Morin

- LP  $\Rightarrow$  global decision **vs** greedy  $\Rightarrow$  incremental decision with a **global target**.
- Global target: when possible, splitting a task to occupy CPU cores and constraint the scheduling.
- At start, for each task generation of "all" recursive versions and "all" schedulings.
  - All interesting recursive versions and schedulings.
- When a task is about to be pushed, look at GPU end.
  - If no task can be split with a part executed on CPU cores, push it to GPUs.
  - o Else, take the less accelerated task, and search for a recursive scheduling that
    - Can finish before the GPUs.
    - Occupy less CPU cores than the available CPUs.



Exa-Soft, Saclay Thomas Morin

- LP  $\Rightarrow$  global decision **vs** greedy  $\Rightarrow$  incremental decision with a **global target**.
- Global target: when possible, splitting a task to occupy CPU cores and constraint the scheduling.
- At start, for each task generation of "all" recursive versions and "all" schedulings.
  - All interesting recursive versions and schedulings.
- When a task is about to be pushed, look at GPU end.
  - If no task can be split with a part executed on CPU cores, push it to GPUs.
  - o Else, take the less accelerated task, and search for a recursive scheduling that
    - Can finish before the GPUs.
    - Occupy less CPU cores than the available CPUs.
- Warning: Splitting generates data transfers that should be taken into account.



Exa-Soft, Saclay Thomas Morin

# Benchmarks - Heterogeneous - Cholesky

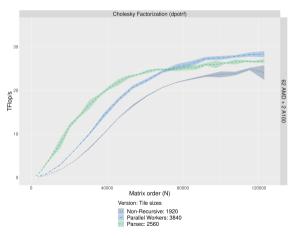


Figure 5: Performance comparison between different versions for Cholesky Factorization with Greedy Algorithm.





# Benchmarks - Heterogeneous - Cholesky

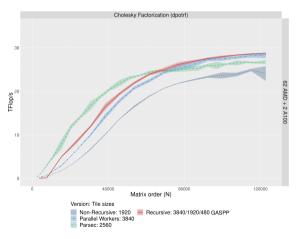


Figure 5: Performance comparison between different versions for Cholesky Factorization with Greedy Algorithm.



# Benchmarks - Heterogeneous - Cholesky

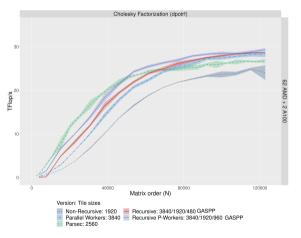


Figure 5: Performance comparison between different versions for Cholesky Factorization with Greedy Algorithm.



# Benchmarks - Heterogeneous - LU without piv

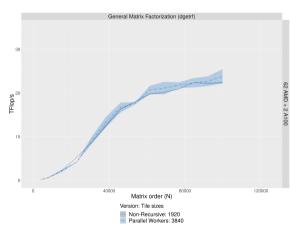


Figure 6: Performance comparison between different versions for LU Factorization without piv with Greedy Algorithm.





# Benchmarks - Heterogeneous - LU without piv

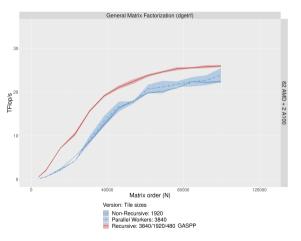


Figure 6: Performance comparison between different versions for LU Factorization without piv with Greedy Algorithm.





# Benchmarks - Heterogeneous - LU without piv

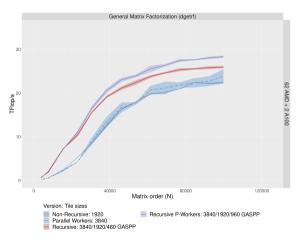


Figure 6: Performance comparison between different versions for LU Factorization without piv with Greedy Algorithm.





# Benchmarks - Heterogeneous - All

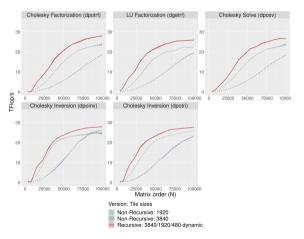


Figure 7: Performance comparison between different kernels and versions with Greedy Algorithm.





# Benchmarks - Heterogeneous - POTRF on GPU

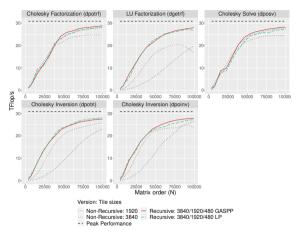


Figure 8: Performance comparison between different kernels and versions with Greedy Algorithm, having dpotrf on GPU.





# Benchmarks - Heterogeneous - Cholesky - POTRF on GPU

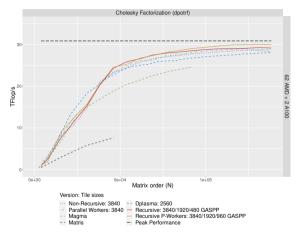


Figure 9: Performance comparison between different versions for Cholesky Factorization with Greedy Algorithm, having POTRF on GPU.















Conclusion





#### Conclusion

- Recursive tasks:
  - Insert subgraph at runtime.
  - More dynamic DAG.
- Splitting task dynamically brings different questions:
  - Which task sould we split.
  - When do we choose to split.

#### Current Work

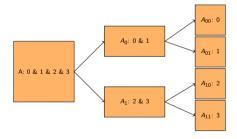
- Distributed recursive tasks, with scheduling.
- Extend to more irregular applications.







### Shared data

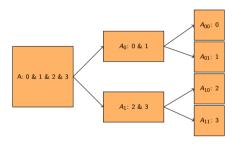








### Shared data



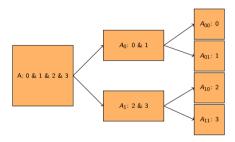
# Auto-pruning







### Shared data



# Auto-pruning Node 0







#### Advantages

- Auto-pruning is important to decrease runtime pressure.
- Communications submitted Just-In-Time, reducing pressure.







# Benchmarks - Homogeneous - LU nopiv

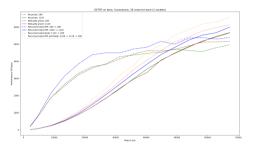


Figure 10: Performance comparison between different versions for LU nopiv.

- No prune: no pruning, all tasks inserted, StarPU remove not submit some.
- Manually prune: The app submits only the needed tasks.
- Recursive tasks MPI: All tasks are submitted, without pruning, but communications are inserted at splitting.