# Programming Heterogeneous & Distributed Architectures using Hierarchical Tasks

Mathieu Faverge, Nathalie Furmento, Abdou Guermouche, Gwenolé Lucas, Thomas Morin, Raymond Namyst, Samuel Thibault, Pierre-André Wacrenier
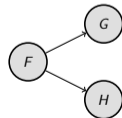
March 2023

## Task Based Programming

- Task-based programming aims to provide portable frameworks capable of exploiting complex architectures.
- Applications are presented as a Directed Acyclic Graph (DAG).
- Runtime systems handle scheduling, communications, . . .

Inría

## Task Based Programming

- Task-based programming aims to provide portable frameworks capable of exploiting complex architectures.
- Applications are presented as a Directed Acyclic Graph (DAG).
- Runtime systems handle scheduling, communications, . . .

## StarPU

- StarPU rely on the *Sequential Task Flow* (STF) to create its DAGs.
- The STF infers dependencies from the order of submission of the tasks and data access modes.

```
F(a)
G(a, b)
H(a, c)
```

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

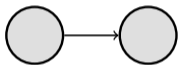## . . . of tasks based programming

- GPUs and CPUs work best on different granularities.
- Some applications are too irregular to fit in a predetermined task-graph.
- Static task graphs limit adaptability during runtime.

## . . . of the STF model

- Runtime overhead induced by a large number of non-ready tasks.
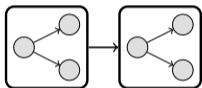- The sequential insertion of tasks can bottleneck the execution of large DAG.

*Inria*

### . . . of tasks based programming

- GPUs and CPUs work best on different granularities.
- Some applications are too irregular to fit in a predetermined task-graph.
- Static task graphs limit adaptability during runtime.

### . . . of the STF model

- Runtime overhead induced by a large number of non-ready tasks.
- The sequential insertion of tasks can bottleneck the execution of large DAG.

$\Rightarrow$ How to create more dymanic task-graphs?

The following runtimes aim at creating more dynamic task-graphs by replacing tasks with an equivalent subgraph:

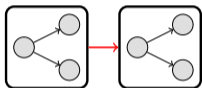| Runtime | Fine-grain Dependencies | Automatic Data Management | Heterogeneity |
|---|---|---|---|
| TaskFlow | | | |
| PaRSEC | | | |
| OmpSs | | | |
| IRIS | | | |
| libtask | | | |
| Our contribution | | | |

Ínría

The following runtimes aim at creating more dynamic task-graphs by replacing tasks with an equivalent subgraph:
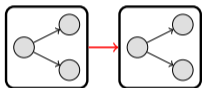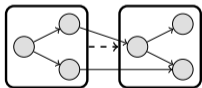


| Runtime | Fine-grain Dependencies | Automatic Data Management | Heterogeneity |
|---|---|---|---|
| TaskFlow | ✗ | ✗ | ✓ |
| PaRSEC | | | |
| OmpSs | | | |
| IRIS | | | |
| libtask | | | |
| Our contribution | | | |

*Inría*

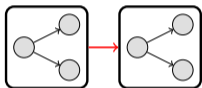The following runtimes aim at creating more dynamic task-graphs by replacing tasks with an equivalent subgraph:
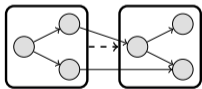


Figure: Barrier between parent tasks

| Runtime | Fine-grain Dependencies | Automatic Data Management | Heterogeneity |
|---|---|---|---|
| TaskFlow | ✗ | ✗ | ✓ |
| PaRSEC | ✗ | ✗ | ✓ |
| OmpSs | | | |
| IRIS | | | |
| libtask | | | |
| Our contribution | | | |

The following runtimes aim at creating more dynamic task-graphs by replacing tasks with an equivalent subgraph:
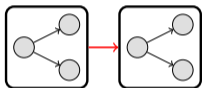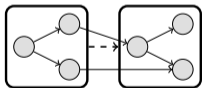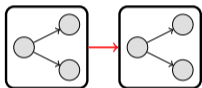


Figure: Barrier between parent tasks



Figure: Fine-grain dependencies

| Runtime | Fine-grain Dependencies | Automatic Data Management | Heterogeneity |
|---|---|---|---|
| TaskFlow | ✗ | ✗ | ✓ |
| PaRSEC | ✗ | ✗ | ✓ |
| OmpSs | ✓ | ✗ | ✗ |
| IRIS | | | |
| libtask | | | |
| Our contribution | | | |

The following runtimes aim at creating more dynamic task-graphs by replacing tasks with an equivalent subgraph:
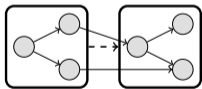


Figure: Barrier between parent tasks
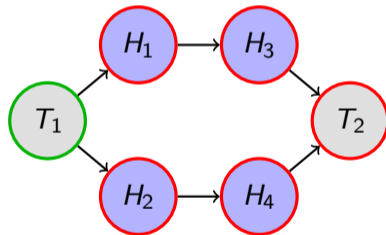


Figure: Fine-grain dependencies

| Runtime | Fine-grain Dependencies | Automatic Data Management | Heterogeneity |
|---|---|---|---|
| TaskFlow | ✗ | ✗ | ✓ |
| PaRSEC | ✗ | ✗ | ✓ |
| OmpSs | ✓ | ✗ | ✗ |
| IRIS | ✗ | ✓ | ✓ |
| libtask | | | |
| Our contribution | | | |

The following runtimes aim at creating more dynamic task-graphs by replacing tasks with an equivalent subgraph:



Figure: Barrier between parent tasks



Figure: Fine-grain dependencies

| Runtime | Fine-grain Dependencies | Automatic Data Management | Heterogeneity |
|---|---|---|---|
| TaskFlow | ✗ | ✗ | ✓ |
| PaRSEC | ✗ | ✗ | ✓ |
| OmpSs | ✓ | ✗ | ✗ |
| IRIS | ✗ | ✓ | ✓ |
| libtask | ✗ | ✓ | ✓/✗ |
| Our contribution | | | |

The following runtimes aim at creating more dynamic task-graphs by replacing tasks with an equivalent subgraph:



Figure: Barrier between parent tasks



Figure: Fine-grain dependencies

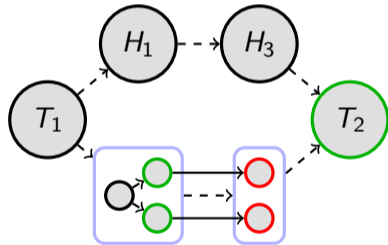| Runtime | Fine-grain Dependencies | Automatic Data Management | Heterogeneity |
|---|---|---|---|
| TaskFlow | ✗ | ✗ | ✓ |
| PaRSEC | ✗ | ✗ | ✓ |
| OmpSs | ✓ | ✗ | ✗ |
| IRIS | ✗ | ✓ | ✓ |
| libtask | ✗ | ✓ | ✓/✗ |
| Our contribution | ✓ | ✓ | ✓ |

## Objectives

- Adapt task granularity to devices
- Reduce the amount of active tasks in StarPU
- Dynamically adapt task implementation at runtime

## Principles

1. No limit for the hierarchy depth
2. Data management is transparent to the programmer
3. Dependencies connect tasks at the finest level possible
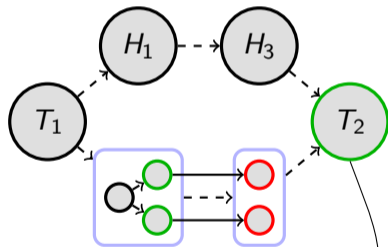
*Inria*

## Objectives

- Adapt task granularity to devices
- Reduce the amount of active tasks in StarPU
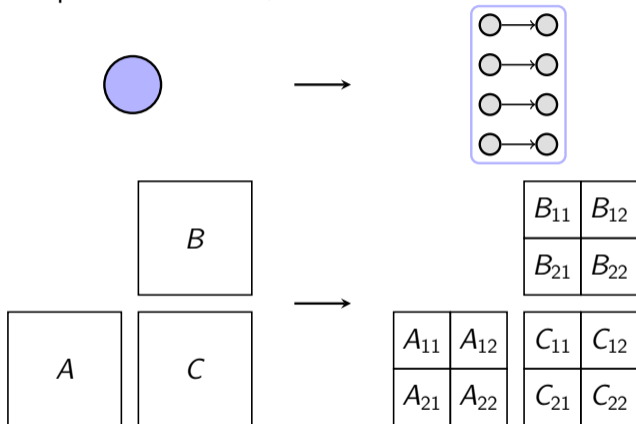- Dynamically adapt task implementation at runtime

## Principles

1. No limit for the hierarchy depth
2. Data management is transparent to the programmer
3. Dependencies connect tasks at the finest level possible

## Objectives

- Adapt task granularity to devices
- Reduce the amount of active tasks in StarPU
- Dynamically adapt task implementation at runtime

## Principles

1. No limit for the hierarchy depth
2. Data management is transparent to the programmer
3. Dependencies connect tasks at the finest level possible

- When executed, a hierarchical task *can* insert a subgraph.

*Inria*

## Objectives

- Adapt task granularity to devices
- Reduce the amount of active tasks in StarPU
- Dynamically adapt task implementation at runtime

## Principles

1. No limit for the hierarchy depth
2. Data management is transparent to the programmer
3. Dependencies connect tasks at the finest level possible

- When executed, a hierarchical task *can* insert a subgraph.
- The dependencies in the subgraphs are infered from the data.

## Objectives

- Adapt task granularity to devices
- Reduce the amount of active tasks in StarPU
- Dynamically adapt task implementation at runtime

## Principles

1. No limit for the hierarchy depth
2. Data management is transparent to the programmer
3. Dependencies connect tasks at the finest level possible

- When executed, a hierarchical task *can* insert a subgraph.
- The dependencies in the subgraphs are infered from the data.



How to ensure the correctness of the DAG?

*Inria*

Matrix-matrix multiplication $C = C + A \times B$:

Matrix-matrix multiplication $C = C + A \times B$:



$\Rightarrow$ How to adapt data partitioning to suit hierarchical tasks?

$A$

```
[Data registration]
register(A)
partition_plan(A, vfilter, V)
partition_plan(A, hfilter, H)




.
```

- The user describes data partitioning through *plan* operations.

```
[Data registration]
register(A)
partition_plan(A, vfilter, V)
partition_plan(A, hfilter, H)


.
```
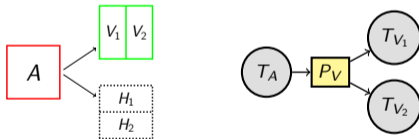
$A \longrightarrow V_1 \; V_2$
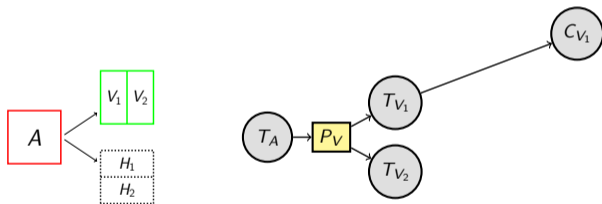
- The user describes data partitioning through *plan* operations.
- Multiple *plans* for the same data piece can coexist.



```
[Data registration]
register(A)
partition_plan(A, vfilter, V)
partition_plan(A, hfilter, H)



.
```

- The user describes data partitioning through *plan* operations.
- Multiple *plans* for the same data piece can coexist.



```
[Data registration]
register(A)
partition_plan(A, vfilter, V)
partition_plan(A, hfilter, H)

[Modification]
submit(T, A:RW)
for (i = 1 to NPARTS)
    submit(T, V[i]:RW)



.
```

- The user describes data partitioning through *plan* operations.
- Multiple *plans* for the same data piece can coexist.
- Partitioning tasks are inserted to ensure that each task will use a valid data.
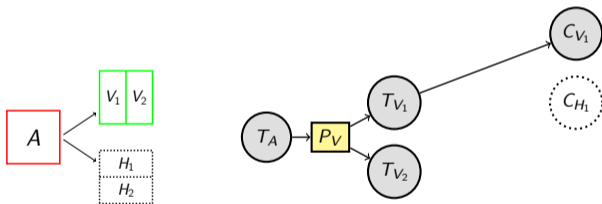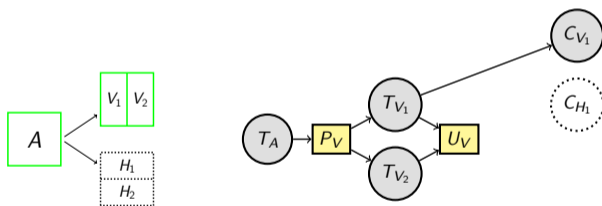


```
[Data registration]
register(A)
partition_plan(A, vfilter, V)
partition_plan(A, hfilter, H)

[Modification]
submit(T, A:RW)
for (i = 1 to NPARTS)
    submit(T, V[i]:RW)



.
```

- The user describes data partitioning through *plan* operations.
- Multiple *plans* for the same data piece can coexist.
- Partitioning tasks are inserted to ensure that each task will use a valid data.



```
[Data registration]
register(A)
partition_plan(A, vfilter, V)
partition_plan(A, hfilter, H)

[Modification]
submit(T, A:RW)
for (i = 1 to NPARTS)
    submit(T, V[i]:RW)



.
```

- The user describes data partitioning through *plan* operations.
- Multiple *plans* for the same data piece can coexist.
- Partitioning tasks are inserted to ensure that each task will use a valid data.
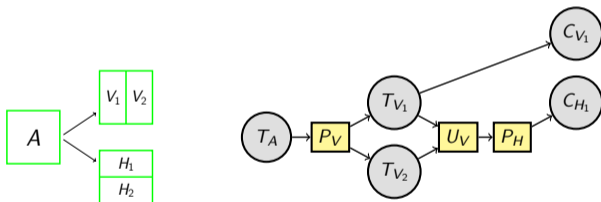


```
[Data registration]
register(A)
partition_plan(A, vfilter, V)
partition_plan(A, hfilter, H)

[Modification]
submit(T, A:RW)
for (i = 1 to NPARTS)
    submit(T, V[i]:RW)


.
```

- The user describes data partitioning through *plan* operations.
- Multiple *plans* for the same data piece can coexist.
- Partitioning tasks are inserted to ensure that each task will use a valid data.
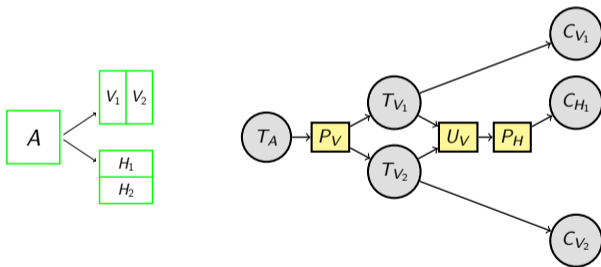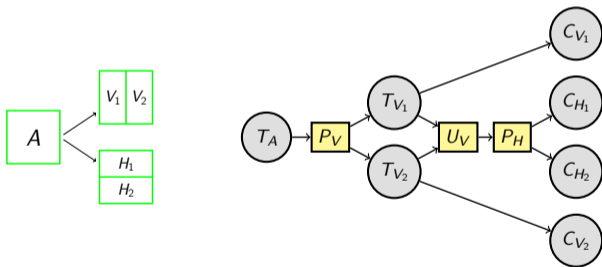


```
[Data registration]
register(A)
partition_plan(A, vfilter, V)
partition_plan(A, hfilter, H)

[Modification]
submit(T, A:RW)
for (i = 1 to NPARTS)
    submit(T, V[i]:RW)
[Verification]
for (i = 1 to NPARTS)
    submit(C, V[i]:R)
    submit(C, H[i]:R)

wait_tasks_completion()
```

- The user describes data partitioning through *plan* operations.
- Multiple *plans* for the same data piece can coexist.
- Partitioning tasks are inserted to ensure that each task will use a valid data.
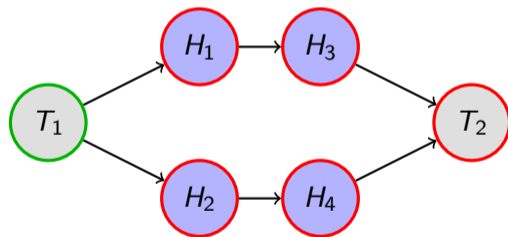


```
[Data registration]
register(A)
partition_plan(A, vfilter, V)
partition_plan(A, hfilter, H)

[Modification]
submit(T, A:RW)
for (i = 1 to NPARTS)
    submit(T, V[i]:RW)
[Verification]
for (i = 1 to NPARTS)
    submit(C, V[i]:R)
    submit(C, H[i]:R)

wait_tasks_completion()
```
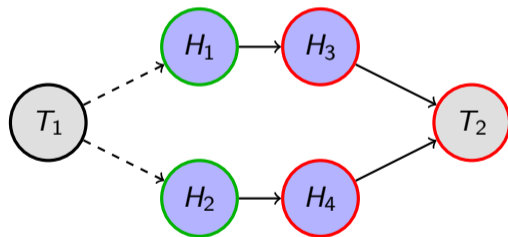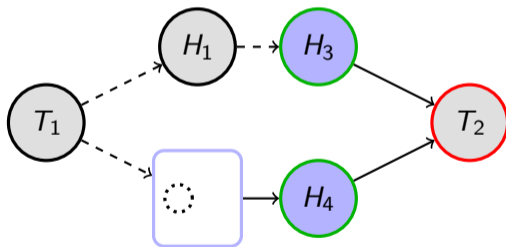
- The user describes data partitioning through *plan* operations.
- Multiple *plans* for the same data piece can coexist.
- Partitioning tasks are inserted to ensure that each task will use a valid data.



```
[Data registration]
register(A)
partition_plan(A, vfilter, V)
partition_plan(A, hfilter, H)

[Modification]
submit(T, A:RW)
for (i = 1 to NPARTS)
    submit(T, V[i]:RW)
[Verification]
for (i = 1 to NPARTS)
    submit(C, V[i]:R)
    submit(C, H[i]:R)

wait_tasks_completion()
```

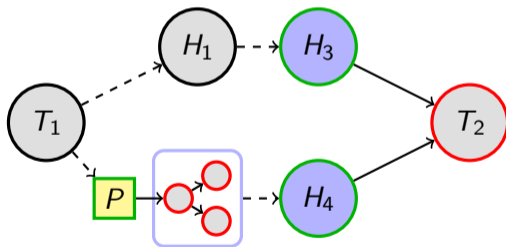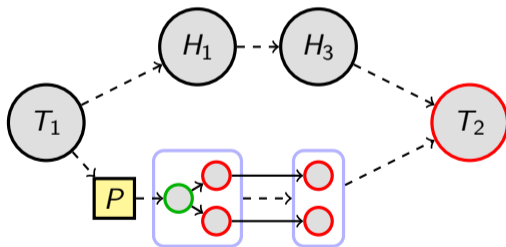- The user describes data partitioning through *plan* operations.
- Multiple *plans* for the same data piece can coexist.
- Partitioning tasks are inserted to ensure that each task will use a valid data.



```
[Data registration]
register(A)
partition_plan(A, vfilter, V)
partition_plan(A, hfilter, H)

[Modification]
submit(T, A:RW)
for (i = 1 to NPARTS)
    submit(T, V[i]:RW)
[Verification]
for (i = 1 to NPARTS)
    submit(C, V[i]:R)
    submit(C, H[i]:R)

wait_tasks_completion()
```

- The user describes data partitioning through *plan* operations.
- Multiple *plans* for the same data piece can coexist.
- Partitioning tasks are inserted to ensure that each task will use a valid data.



```
[Data registration]
register(A)
partition_plan(A, vfilter, V)
partition_plan(A, hfilter, H)

[Modification]
submit(T, A:RW)
for (i = 1 to NPARTS)
    submit(T, V[i]:RW)
[Verification]
for (i = 1 to NPARTS)
    submit(C, V[i]:R)
    submit(C, H[i]:R)

wait_tasks_completion()
```

- The user describes data partitioning through *plan* operations.
- Multiple *plans* for the same data piece can coexist.
- Partitioning tasks are inserted to ensure that each task will use a valid data.



```
[Data registration]
register(A)
partition_plan(A, vfilter, V)
partition_plan(A, hfilter, H)

[Modification]
submit(T, A:RW)
for (i = 1 to NPARTS)
    submit(T, V[i]:RW)
[Verification]
for (i = 1 to NPARTS)
    submit(C, V[i]:R)
    submit(C, H[i]:R)

wait_tasks_completion()
```

- Partitioning tasks are added, if needed, at the submission of a subtask.

Inria

- Partitioning tasks are added, if needed, at the submission of a subtask.

- Partitioning tasks are added, if needed, at the submission of a subtask.

- Partitioning tasks are added, if needed, at the submission of a subtask.
- Unpartitioning tasks are added, if needed, before a regular task. They enforce the correctness of the DAG.

Inria

Figure: Submission cost of computational tasks for the matrix-matrix multiplication kernel.

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

| Name | Processor | GPU | Memory |
|------|-----------|-----|--------|
| INTEL-V100 | 2 × INTEL XEON GOLD 6142, 16 cores, 2.6GHz | 2 × NVIDIA V100 (16GB) | 384GB |
| AMD-A100 | 2 × AMD ZEN3 EPYC 7513, 32 cores, 2.6GHz | 2 × NVIDIA A100 (40GB) | 512GB |



Figure: Full matrix partitioning.

Figure: Tile size of 960.



Version: Tile sizes

Non-Hierarchical: 960

Figure: Matrix-matrix multiplication kernel with a fixed percentage of hierarchical tasks on INTEL-V100.

Figure: Tile size of 2880.



Figure: Matrix-matrix multiplication kernel with a fixed percentage of hierarchical tasks on INTEL-V100.

Figure: 10% recursive matrix partitioning.



Version: Tile sizes

- Non-Hierarchical: 2880
- Non-Hierarchical: 960
- Hierarchical: 2880 / 960 (10%)
- Hierarchical: 2880 / 960 (15%)

Figure: Matrix-matrix multiplication kernel with a fixed percentage of hierarchical tasks on INTEL-V100.

Figure: Tile size of 320.



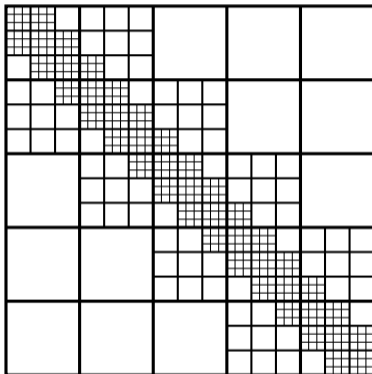Figure: LU factorization kernel with a diagonal repartition of hierarchical tasks on INTEL-V100.

Figure: Tile size of 960.



LU Factorization (dgetrf_nopiv)

Version: Tile sizes
— Non-Hierarchical: 960
— Non-Hierarchical: 320

Figure: LU factorization kernel with a diagonal repartition of hierarchical tasks on INTEL-V100.

*Inría*

Figure: Tile size of 2880.



Figure: LU factorization kernel with a diagonal repartition of hierarchical tasks on INTEL-V100.

Figure: Diagonal matrix partitioning.



Figure: LU factorization kernel with a diagonal repartition of hierarchical tasks on INTEL-V100.

Figure: Diagonal matrix partitioning.



LU Factorization (dgetrf_nopiv)

Version: Tile sizes

- Non-Hierarchical: 2880
- Non-Hierarchical: 960
- Non-Hierarchical: 320
- Hierarchical: 2880 / 960
- Hierarchical: 2880 / 960 / 320

Figure: LU factorization kernel with a diagonal repartition of hierarchical tasks on INTEL-V100.

Ínría

Figure: Cholesky type operations (DPOTRF, DPOSV, DPOINV) kernel with diagonal distribution of the hierarchical tasks on INTEL-V100.

- Non-hierarchical version: Various tile sizes and number of stream (1, 2, 4, 8).
- Hierarchical version: Various levels of partitioning.



Figure: Aggregation of best results for LU factorization on INTEL-V100.

Figure: Aggregation of best results for LU factorization on AMD-A100.
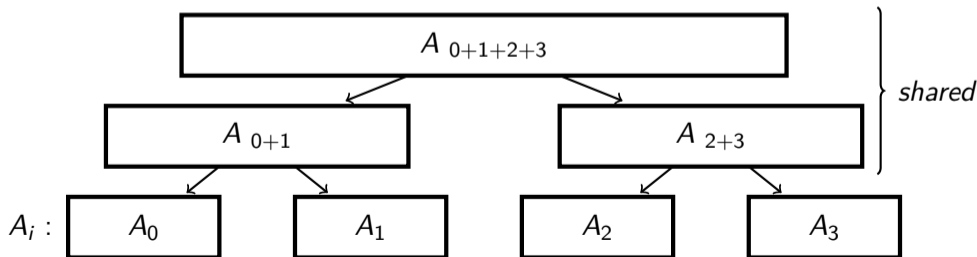
- Thomas Morin extends hierarchical tasks to the distributed context for his intership.

- Thomas Morin extends hierarchical tasks to the distributed context for his intership.
- In StarPU+MPI, data must be distributed among the nodes.

$A$ :

| $A_0$ | $A_1$ | $A_2$ | $A_3$ |
|-------|-------|-------|-------|

Inría

- Thomas Morin extends hierarchical tasks to the distributed context for his intership.
- In StarPU+MPI, data must be distributed among the nodes.

$\Rightarrow$ How to submit a hierarchical task on $A$ that will be processed on different nodes?

$A$ :

| $A_0$ | $A_1$ | $A_2$ | $A_3$ |
|:---:|:---:|:---:|:---:|

- Thomas Morin extends hierarchical tasks to the distributed context for his intership.
- In StarPU+MPI, data must be distributed among the nodes.

$\Rightarrow$ How to submit a hierarchical task on $A$ that will be processed on different nodes?

- He introduced the notion of *shared data*.

- Thomas Morin extends hierarchical tasks to the distributed context for his intership.
- In StarPU+MPI, data must be distributed among the nodes.

$\Rightarrow$ How to submit a hierarchical task on $A$ that will be processed on different nodes?

- He introduced the notion of *shared data*.
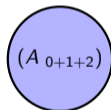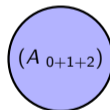- Whether a data is *shared* is automatically deduced when the $A_i$ are registered.

- Using shared data allows us to automatically prune the distributed DAG of tasks irrelevant for certain nodes.

- Using shared data allows us to automatically prune the distributed DAG of tasks irrelevant for certain nodes.
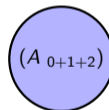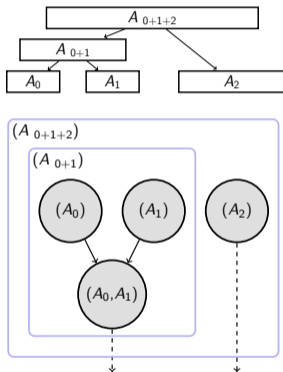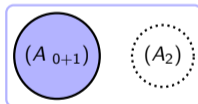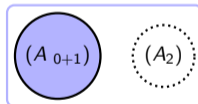


N0:

N1:

N2:

- Using shared data allows us to automatically prune the distributed DAG of tasks irrelevant for certain nodes.

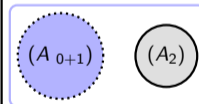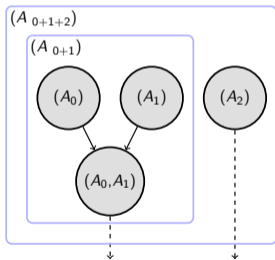- Using shared data allows us to automatically prune the distributed DAG of tasks irrelevant for certain nodes.
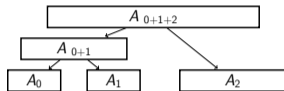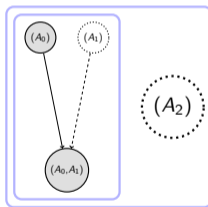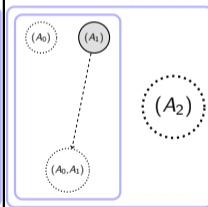
- Using shared data allows us to automatically prune the distributed DAG of tasks irrelevant for certain nodes.
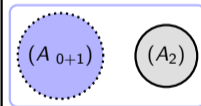
- Hierarchical tasks can insert a subgraph at runtime, resulting in a more dynamic DAG.
- Data management is handled automatically and contributes to the correctness of hierarchical DAGs.
- Hierarchical tasks can be used with StarPU+MPI and allow for automatic pruning of the DAG.

## Future Work

- Scheduling questions:
  - > When should we insert a subgraph ?
  - > Where should we execute it ?
  - > Using which implementation ?
- Testing with applications benefitting more from dynamic task graphs.
- In distributed computing, there are investigations to be made around data transfers (delay them, use higher level data, etc).

*Inria*