



Heterogeneous Scheduling

IRMA « New Trends in Computing » Summer School

Lionel Eyraud-Dubois

INRIA Bordeaux - Sud-Ouest -- TOPAL Team

Coming up

- (today) StarPU presentation
- (Wednesday) Hands-on with StarPU
- (Thursday) Theoretical Scheduling course



StarPU tutorial

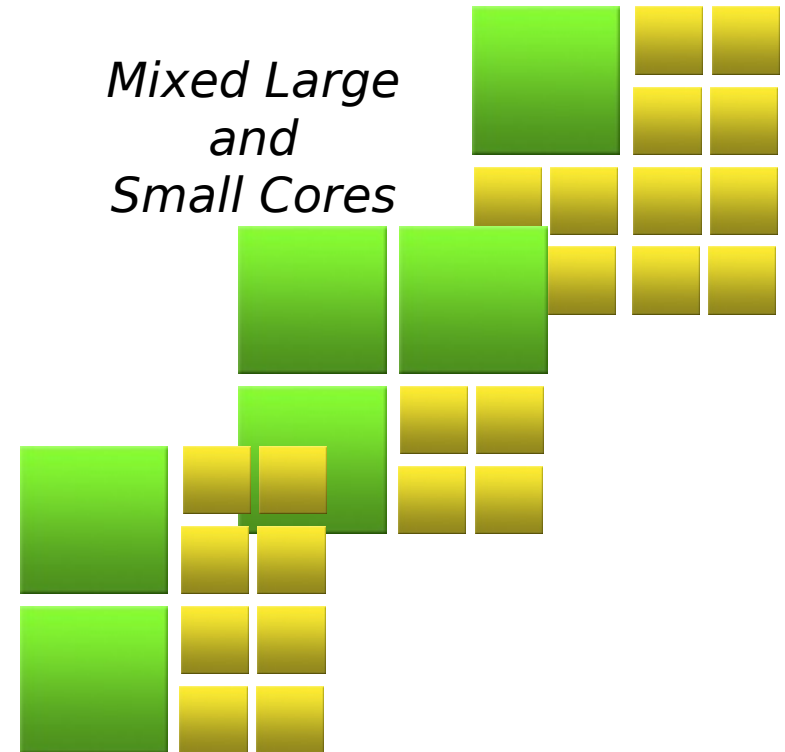
Nathalie Furmento - Samuel Thibault

INRIA Bordeaux - Sud-Ouest -- STORM Team

Introduction

Toward heterogeneous multi-core architectures

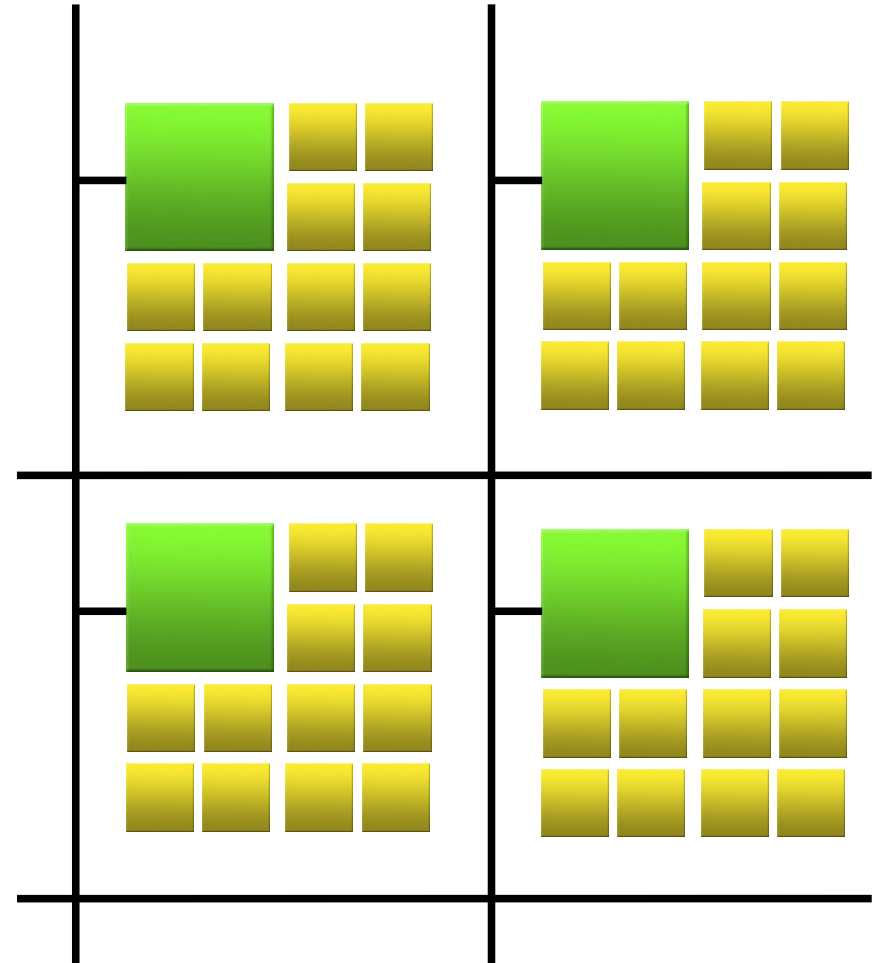
- Multicore is here
 - Hierarchical architectures
 - Manycore
 - Heterogeneous systems
- Architecture specialization
 - Now
 - Accelerators (GPGPUs, FPGAs)
 - Coprocessors (Xeon Phi)
 - All of the above
 - In the near Future
 - Many simple cores
 - A few full-featured cores



Introduction

Toward heterogeneous multi-core clusters

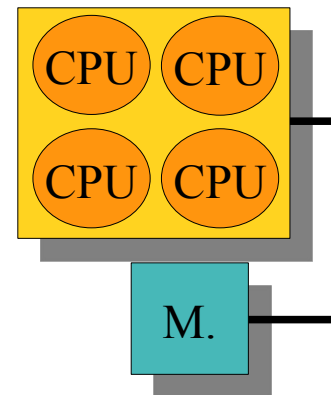
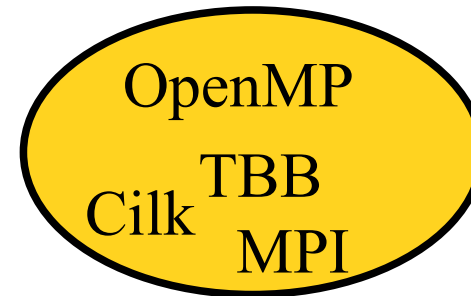
- Multicore is here
 - Hierarchical architectures
 - Manycore
 - Heterogeneous systems
- Clusters thereof
 - High-speed network
 - Network topology
 - Towards exascale



How to program these architectures?

- Multicore programming
 - pthreads, OpenMP, TBB, ...

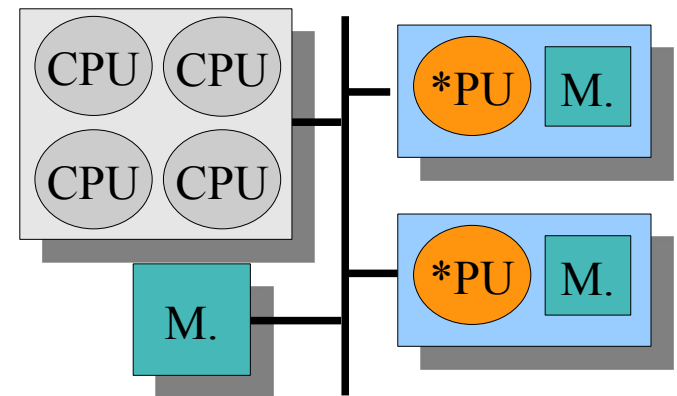
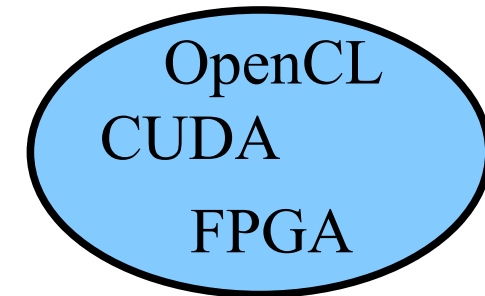
Multicore



How to program these architectures?

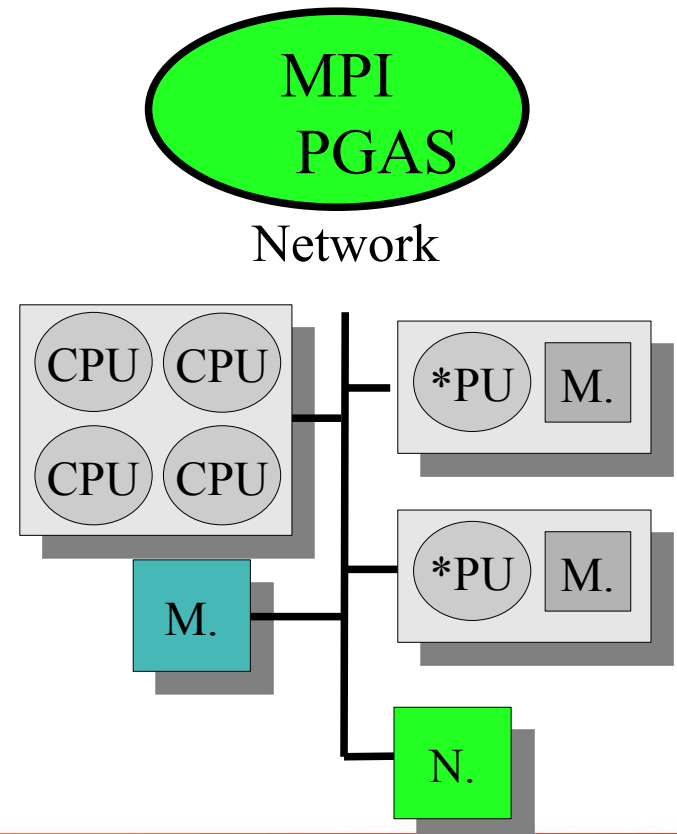
- Multicore programming
 - pthreads, OpenMP, TBB, ...
- Accelerator programming
 - CUDA, OpenCL, FPGA ?
 - OpenMP 5.0?
 - (Often) Pure offloading model

Accelerators



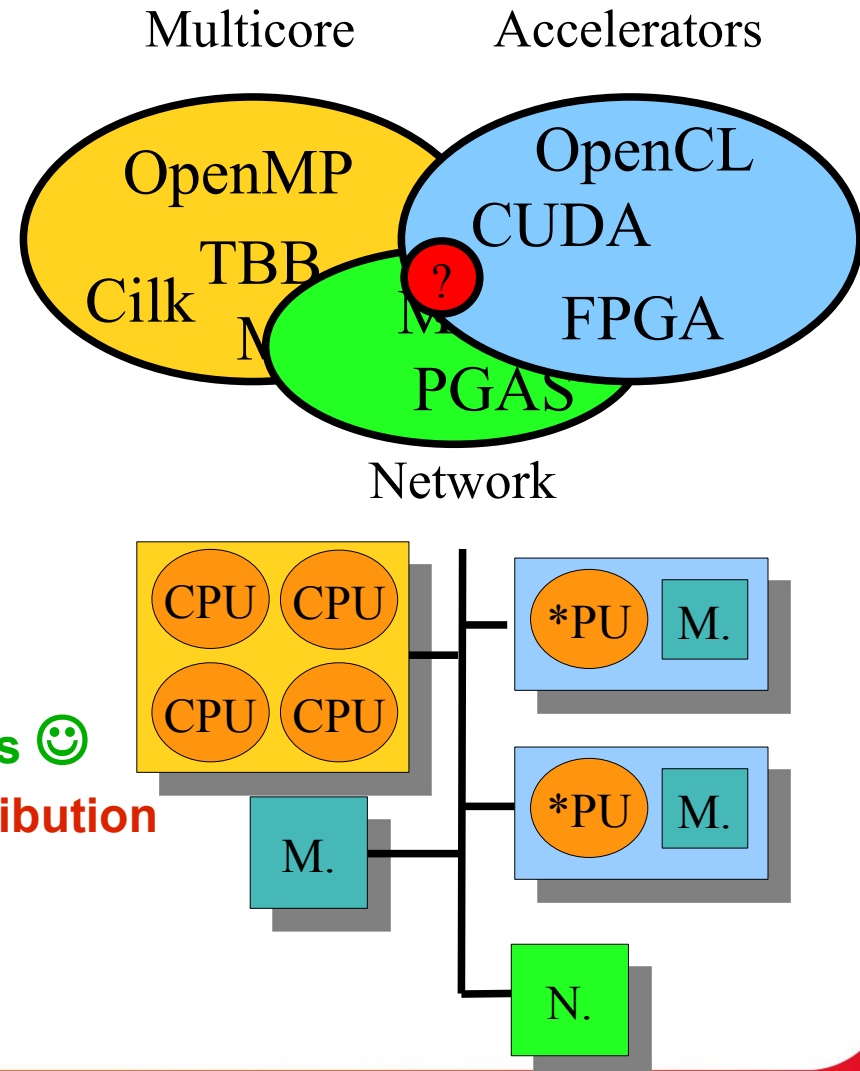
How to program these architectures?

- Multicore programming
 - pthreads, OpenMP, TBB, ...
- Accelerator programming
 - CUDA, OpenCL, FPGA ?
 - OpenMP 5.0?
 - (Often) Pure offloading model
- Network support
 - MPI / PGAS



How to program these architectures?

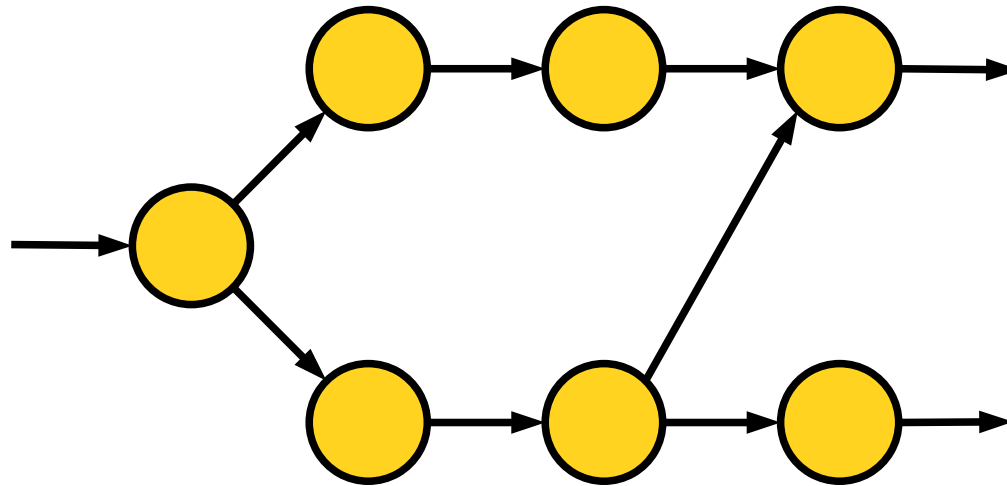
- Multicore programming
 - pthreads, OpenMP, TBB, ...
- Accelerator programming
 - CUDA, OpenCL, FPGA ?
 - OpenMP 5.0?
 - (Often) Pure offloading model
- Network support
 - MPI / PGAS
- Hybrid models?
 - **Take advantage of all resources 😊**
 - **Complex interactions and distribution ☹️**



Task graphs

- Well-studied for scheduling parallelism (since 60's!)
- But only recent trend in HPC
- Departs from usual sequential programming

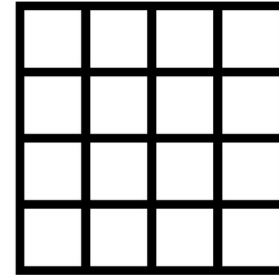
Really ?



Expressing a task graph

Implicit task dependencies

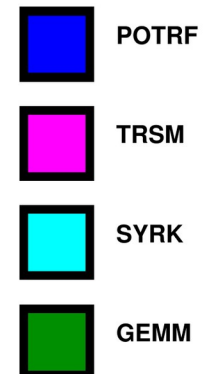
- Right-Looking Cholesky decomposition (from PLASMA)



```

for (j = 0; j < N; j++) {
  POTRF (RW,A[j][j]);
  for (i = j+1; i < N; i++)
    TRSM (RW,A[i][j], R,A[j][j]);
  for (i = j+1; i < N; i++) {
    SYRK (RW,A[i][i], R,A[i][j]);
    for (k = j+1; k < i; k++)
      GEMM (RW,A[i][k],
           R,A[i][j], R,A[k][j]);
  }
}
task_wait_for_all();

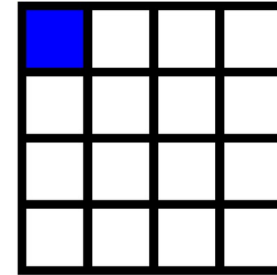
```



Expressing a task graph

Implicit task dependencies

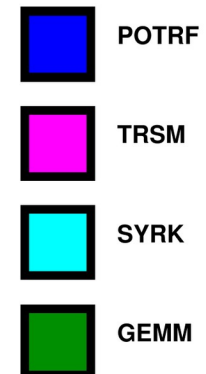
- Right-Looking Cholesky decomposition (from PLASMA)



```

for (j = 0; j < N; j++) {
  POTRF (RW,A[j][j]);
  for (i = j+1; i < N; i++)
    TRSM (RW,A[i][j], R,A[j][j]);
  for (i = j+1; i < N; i++) {
    SYRK (RW,A[i][i], R,A[i][j]);
    for (k = j+1; k < i; k++)
      GEMM (RW,A[i][k],
           R,A[i][j], R,A[k][j]);
  }
}
task_wait_for_all();

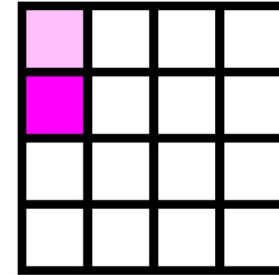
```



Expressing a task graph

Implicit task dependencies

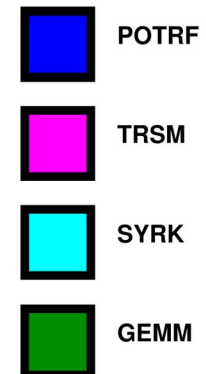
- Right-Looking Cholesky decomposition (from PLASMA)



```

for (j = 0; j < N; j++) {
  POTRF (RW,A[j][j]);
  for (i = j+1; i < N; i++)
    TRSM (RW,A[i][j], R,A[j][j]);
  for (i = j+1; i < N; i++) {
    SYRK (RW,A[i][i], R,A[i][j]);
    for (k = j+1; k < i; k++)
      GEMM (RW,A[i][k],
           R,A[i][j], R,A[k][j]);
  }
}
task_wait_for_all();

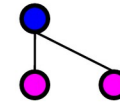
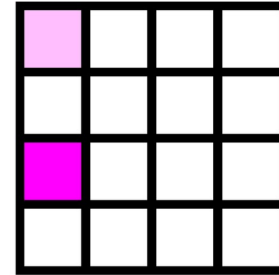
```



Expressing a task graph

Implicit task dependencies

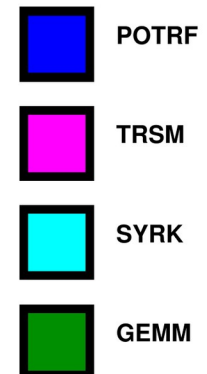
- Right-Looking Cholesky decomposition (from PLASMA)



```

for (j = 0; j < N; j++) {
  POTRF (RW,A[j][j]);
  for (i = j+1; i < N; i++)
    TRSM (RW,A[i][j], R,A[j][j]);
  for (i = j+1; i < N; i++) {
    SYRK (RW,A[i][i], R,A[i][j]);
    for (k = j+1; k < i; k++)
      GEMM (RW,A[i][k],
           R,A[i][j], R,A[k][j]);
  }
}
task_wait_for_all();

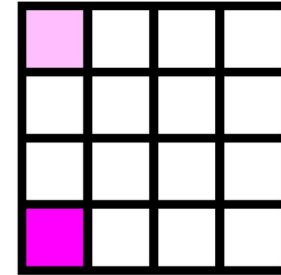
```



Expressing a task graph

Implicit task dependencies

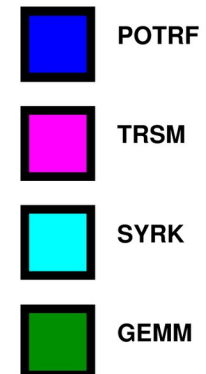
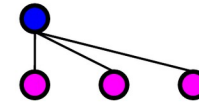
- Right-Looking Cholesky decomposition (from PLASMA)



```

for (j = 0; j < N; j++) {
  POTRF (RW,A[j][j]);
  for (i = j+1; i < N; i++)
    TRSM (RW,A[i][j], R,A[j][j]);
  for (i = j+1; i < N; i++) {
    SYRK (RW,A[i][i], R,A[i][j]);
    for (k = j+1; k < i; k++)
      GEMM (RW,A[i][k],
           R,A[i][j], R,A[k][j]);
  }
}
task_wait_for_all();

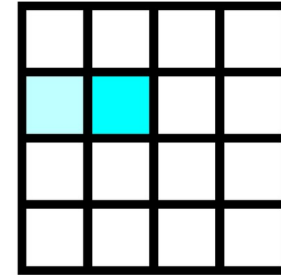
```



Expressing a task graph

Implicit task dependencies

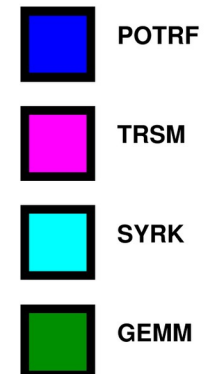
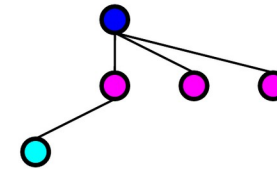
- Right-Looking Cholesky decomposition (from PLASMA)



```

for (j = 0; j < N; j++) {
  POTRF (RW,A[j][j]);
  for (i = j+1; i < N; i++)
    TRSM (RW,A[i][j], R,A[j][j]);
  for (i = j+1; i < N; i++) {
    SYRK (RW,A[i][i], R,A[i][j]);
    for (k = j+1; k < i; k++)
      GEMM (RW,A[i][k],
           R,A[i][j], R,A[k][j]);
  }
}
task_wait_for_all();

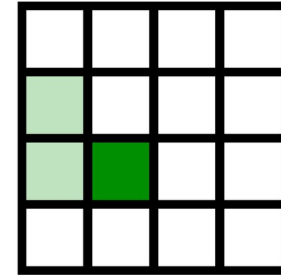
```



Expressing a task graph

Implicit task dependencies

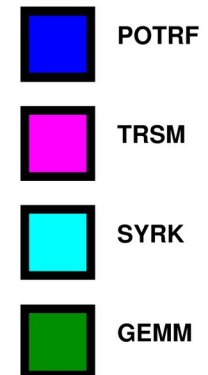
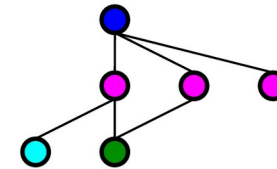
- Right-Looking Cholesky decomposition (from PLASMA)



```

for (j = 0; j < N; j++) {
  POTRF (RW,A[j][j]);
  for (i = j+1; i < N; i++)
    TRSM (RW,A[i][j], R,A[j][j]);
  for (i = j+1; i < N; i++) {
    SYRK (RW,A[i][i], R,A[i][j]);
    for (k = j+1; k < i; k++)
      GEMM (RW,A[i][k],
           R,A[i][j], R,A[k][j]);
  }
}
task_wait_for_all();

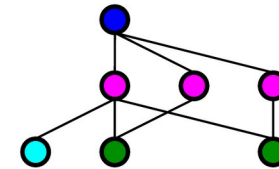
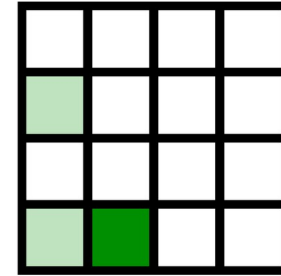
```



Expressing a task graph

Implicit task dependencies

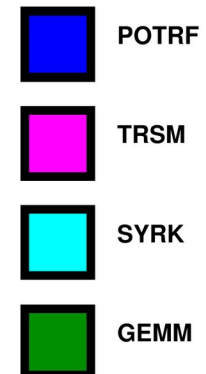
- Right-Looking Cholesky decomposition (from PLASMA)



```

for (j = 0; j < N; j++) {
  POTRF (RW,A[j][j]);
  for (i = j+1; i < N; i++)
    TRSM (RW,A[i][j], R,A[j][j]);
  for (i = j+1; i < N; i++) {
    SYRK (RW,A[i][i], R,A[i][j]);
    for (k = j+1; k < i; k++)
      GEMM (RW,A[i][k],
           R,A[i][j], R,A[k][j]);
  }
}
task_wait_for_all();

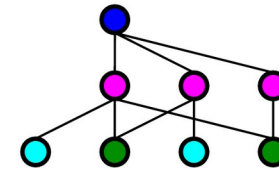
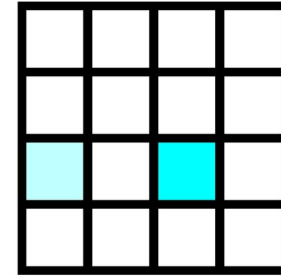
```



Expressing a task graph

Implicit task dependencies

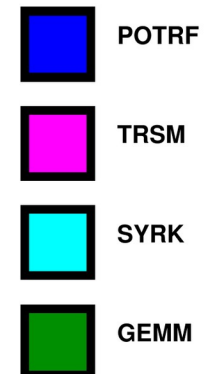
- Right-Looking Cholesky decomposition (from PLASMA)



```

for (j = 0; j < N; j++) {
  POTRF (RW,A[j][j]);
  for (i = j+1; i < N; i++)
    TRSM (RW,A[i][j], R,A[j][j]);
  for (i = j+1; i < N; i++) {
    SYRK (RW,A[i][i], R,A[i][j]);
    for (k = j+1; k < i; k++)
      GEMM (RW,A[i][k],
           R,A[i][j], R,A[k][j]);
  }
}
task_wait_for_all();

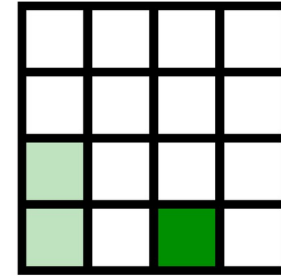
```



Expressing a task graph

Implicit task dependencies

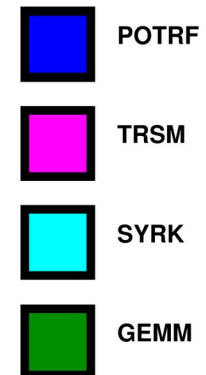
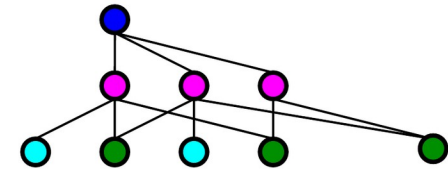
- Right-Looking Cholesky decomposition (from PLASMA)



```

for (j = 0; j < N; j++) {
  POTRF (RW,A[j][j]);
  for (i = j+1; i < N; i++)
    TRSM (RW,A[i][j], R,A[j][j]);
  for (i = j+1; i < N; i++) {
    SYRK (RW,A[i][i], R,A[i][j]);
    for (k = j+1; k < i; k++)
      GEMM (RW,A[i][k],
           R,A[i][j], R,A[k][j]);
  }
}
task_wait_for_all();

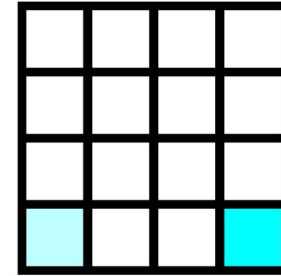
```



Expressing a task graph

Implicit task dependencies

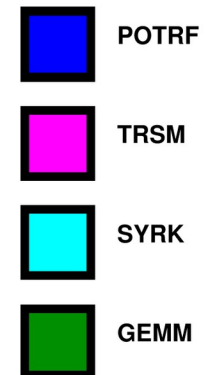
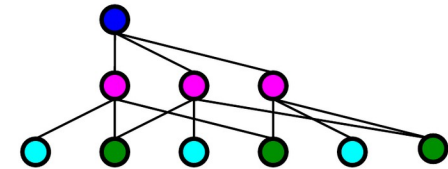
- Right-Looking Cholesky decomposition (from PLASMA)



```

for (j = 0; j < N; j++) {
  POTRF (RW,A[j][j]);
  for (i = j+1; i < N; i++)
    TRSM (RW,A[i][j], R,A[j][j]);
  for (i = j+1; i < N; i++) {
    SYRK (RW,A[i][i], R,A[i][j]);
    for (k = j+1; k < i; k++)
      GEMM (RW,A[i][k],
            R,A[i][j], R,A[k][j]);
  }
}
task_wait_for_all();

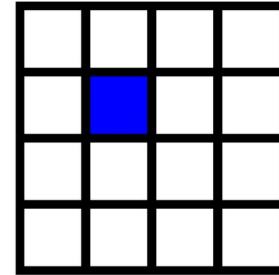
```



Expressing a task graph

Implicit task dependencies

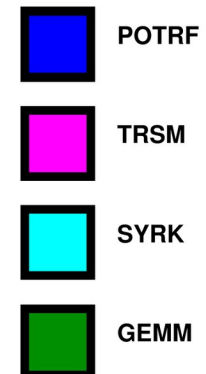
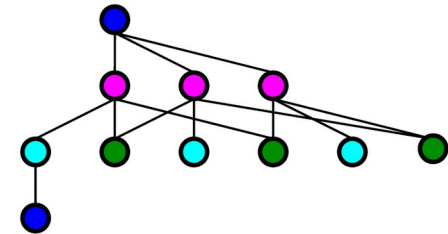
- Right-Looking Cholesky decomposition (from PLASMA)



```

for (j = 0; j < N; j++) {
  POTRF (RW,A[j][j]);
  for (i = j+1; i < N; i++)
    TRSM (RW,A[i][j], R,A[j][j]);
  for (i = j+1; i < N; i++) {
    SYRK (RW,A[i][i], R,A[i][j]);
    for (k = j+1; k < i; k++)
      GEMM (RW,A[i][k],
           R,A[i][j], R,A[k][j]);
  }
}
task_wait_for_all();

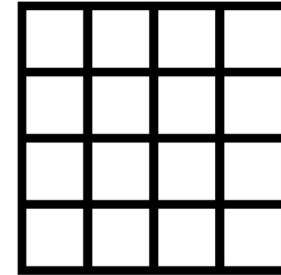
```



Expressing a task graph

Implicit task dependencies

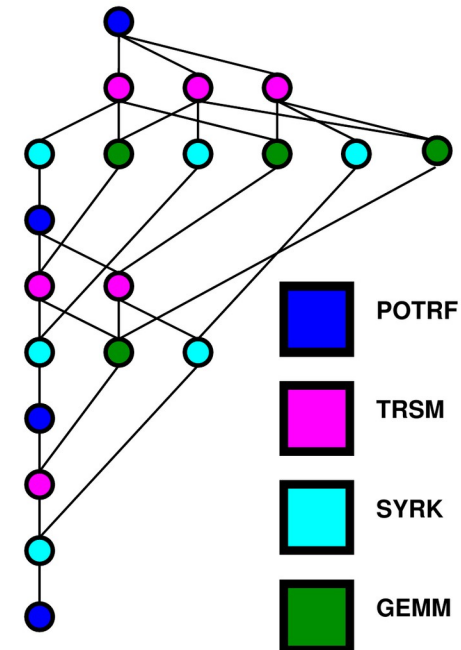
- Right-Looking Cholesky decomposition (from PLASMA)



```

for (j = 0; j < N; j++) {
  POTRF (RW,A[j][j]);
  for (i = j+1; i < N; i++)
    TRSM (RW,A[i][j], R,A[j][j]);
  for (i = j+1; i < N; i++) {
    SYRK (RW,A[i][i], R,A[i][j]);
    for (k = j+1; k < i; k++)
      GEMM (RW,A[i][k],
           R,A[i][j], R,A[k][j]);
  }
}
task_wait_for_all();

```



Write your application as a task graph

Even if using a sequential-looking source code

➔ Portable performance

Sequential Task Flow (STF)

- Algorithm remains the same on the long term
- Can debug the sequential version.
- Only kernels need to be rewritten
 - BLAS libraries, multi-target compilers
- Runtime will handle parallel execution

Task-based programming

- Needs code restructuring
 - Split computation into tasks
 - BLAS, typically
 - Supposed to have “stable” performance
- Constraining
 - No global variables
 - Mandatory for GPUs
- Actually... functional programming

So a good move, in the end 😊

- Have to accept constraints and losing control

Just like we did when moving from assembly to high-level languages

Overview of StarPU

Overview of StarPU

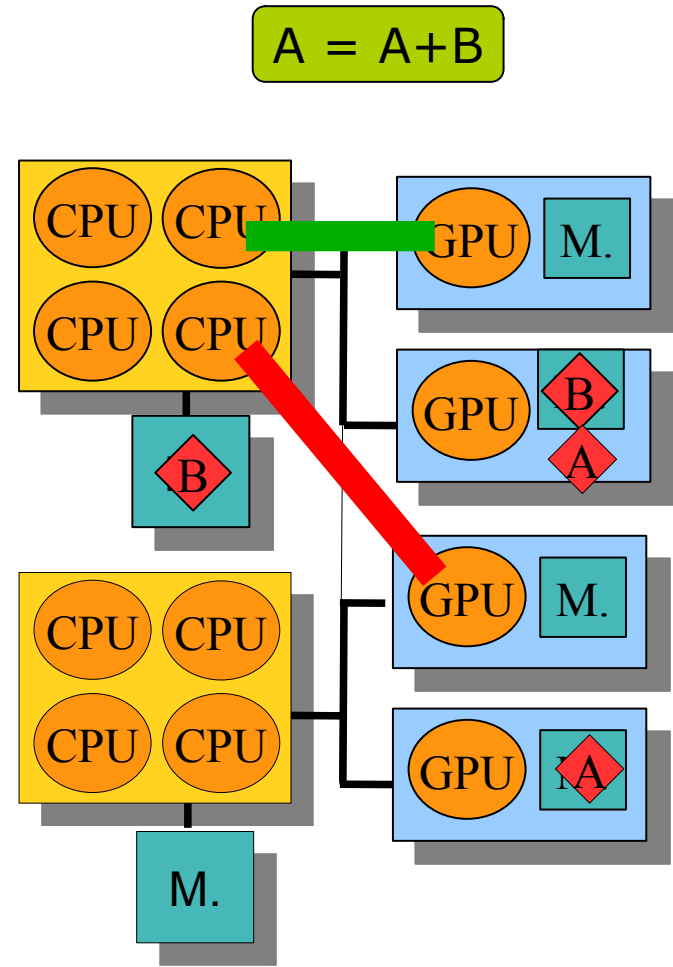
Rationale

Task scheduling

- Dynamic
- On all kinds of PU
 - General purpose
 - Accelerators/specialized

Memory transfer

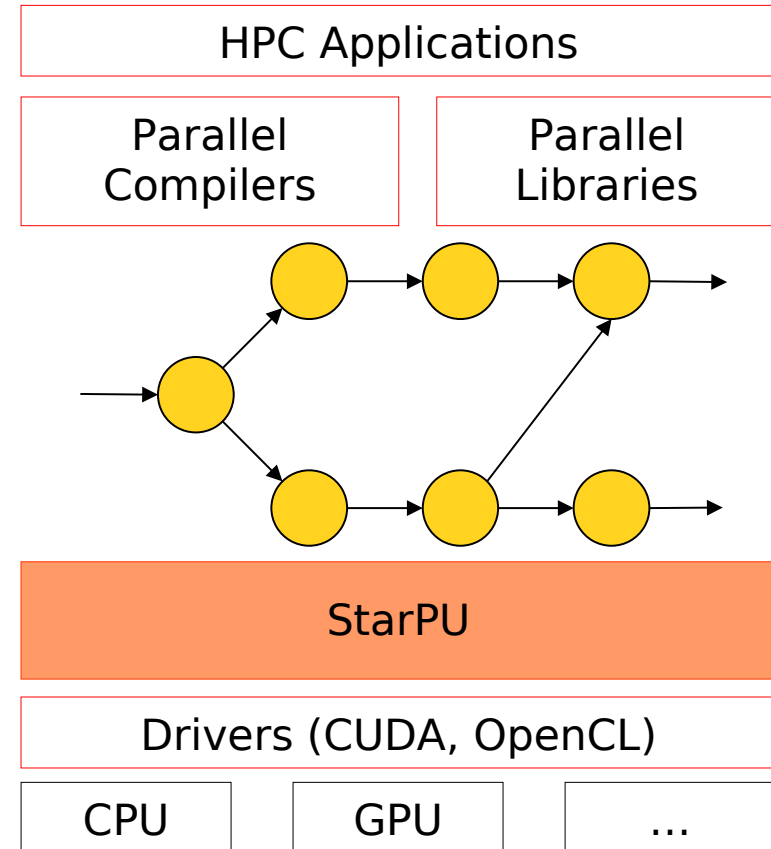
- Eliminate redundant transfers
- Software VSM (Virtual Shared Memory)



The StarPU runtime system

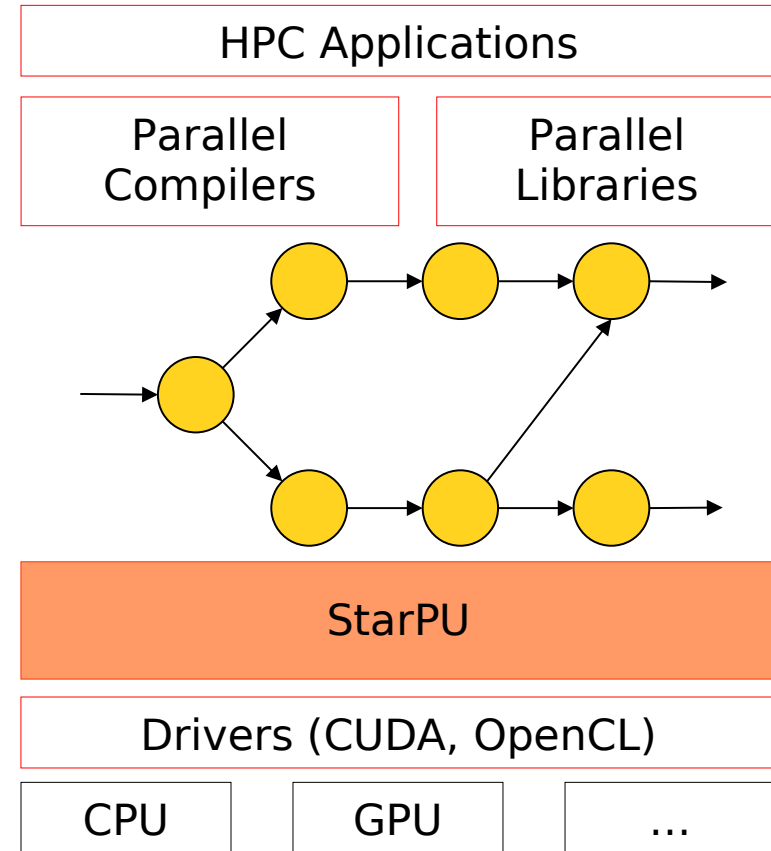
The need for runtime systems

- “do dynamically what can’t be done statically anymore”
- Compilers and libraries generate (graphs of) tasks
 - Additional information is welcome!
- StarPU provides
 - Task scheduling
 - Memory management



Data management

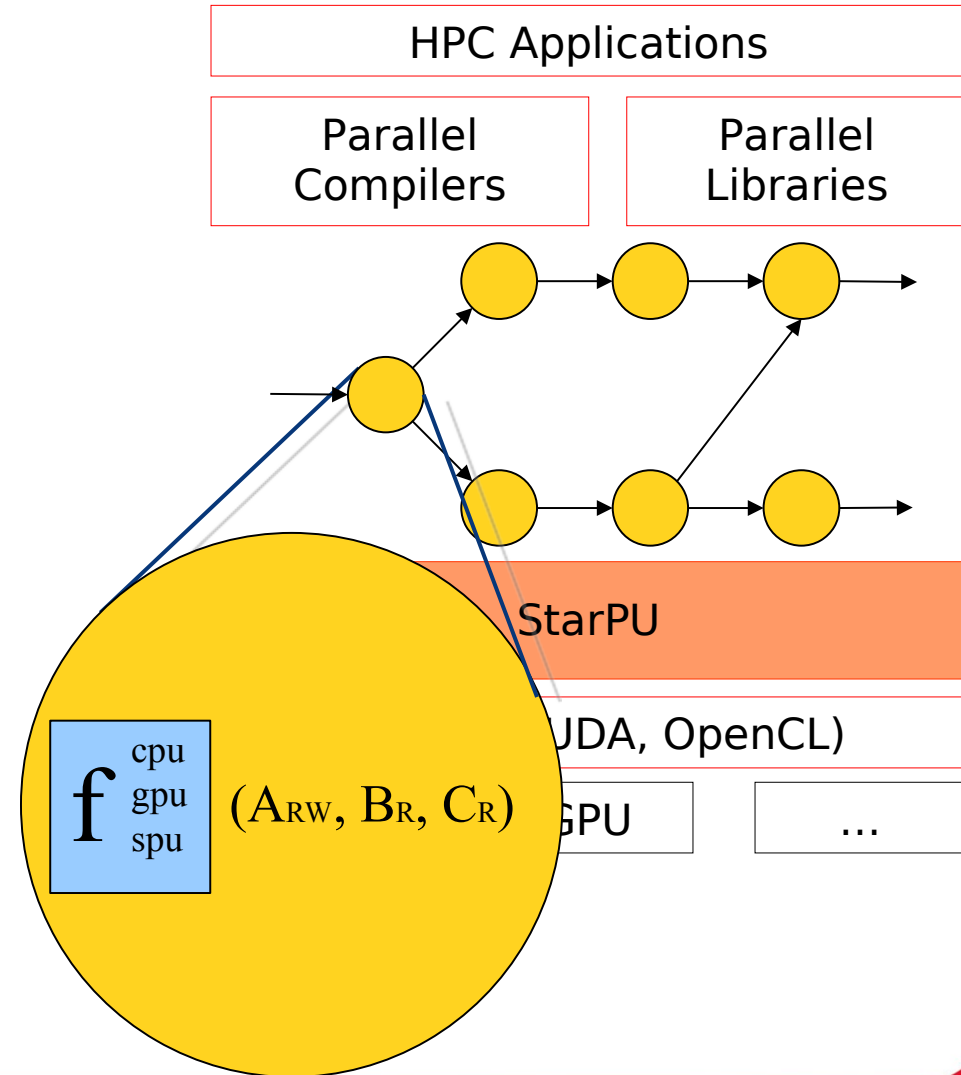
- StarPU provides a **Virtual Shared Memory (VSM)** subsystem (aka DSM)
 - Replication
 - Consistency
 - Single writer
 - Or reduction, ...
- Input & output of tasks = reference to VSM data



The StarPU runtime system

Task scheduling

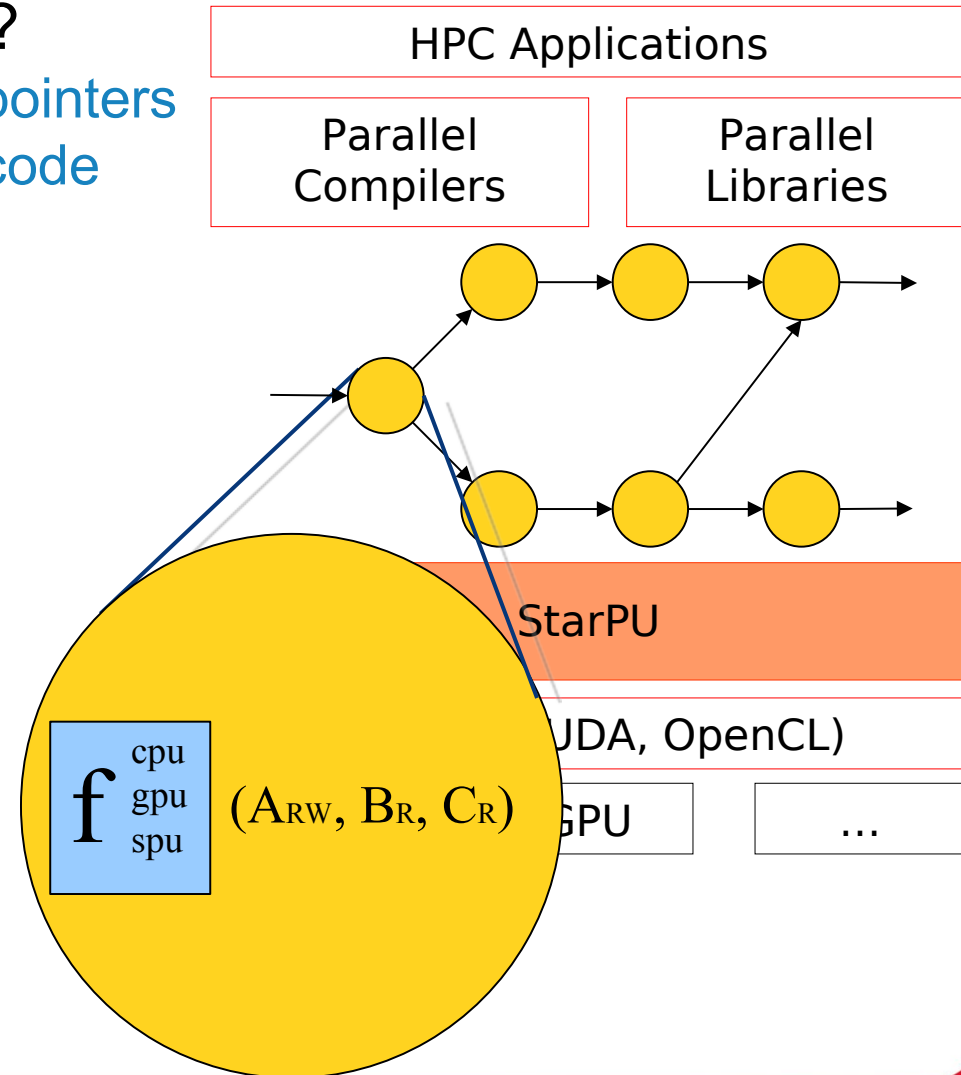
- Tasks =
 - Data input & output
 - Reference to VSM data
 - Multiple implementations
 - E.g. CUDA + CPU implementation
 - Non-preemptible
 - Dependencies with other tasks
- StarPU provides an **Open Scheduling platform**
 - Scheduling algorithm = plug-ins



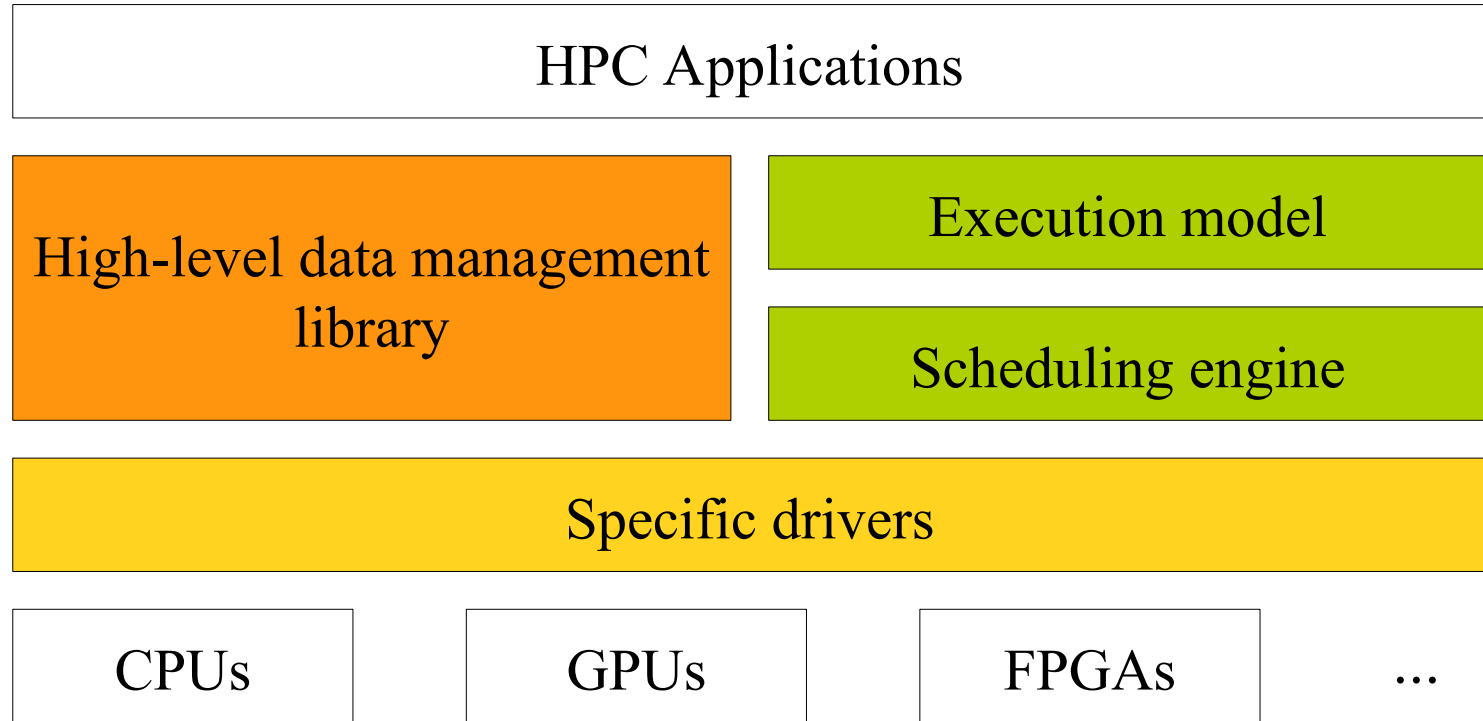
The StarPU runtime system

Task scheduling

- Who generates the code ?
 - StarPU Task \sim function pointers
 - StarPU doesn't generate code
- Libraries era
 - PLASMA + MAGMA
 - FFTW + CUFFT...
 - Variants management
- Rely on compilers



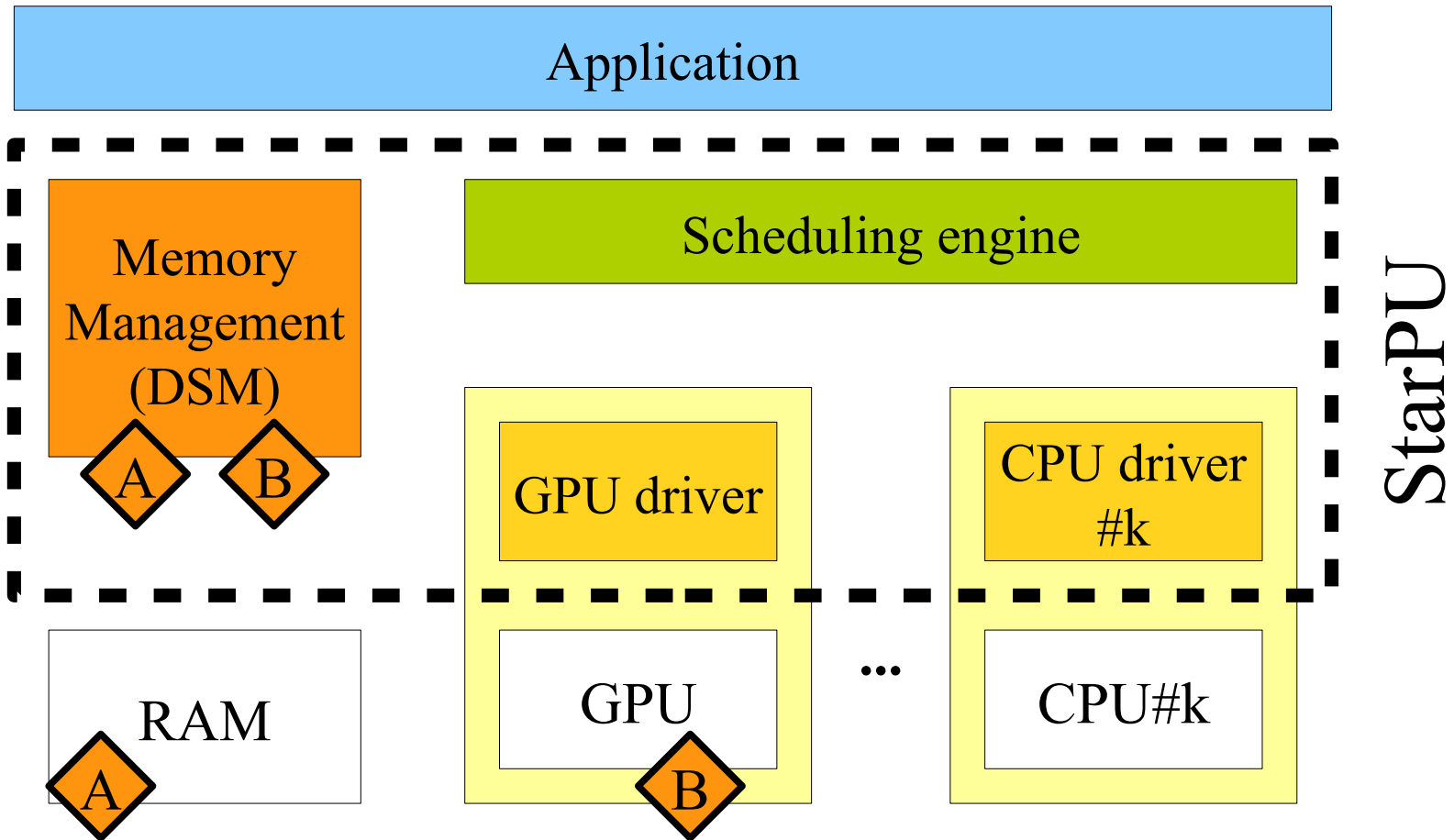
The StarPU runtime system



Mastering CPUs, GPUs, FPGAs ... ***PUs** → **StarPU**

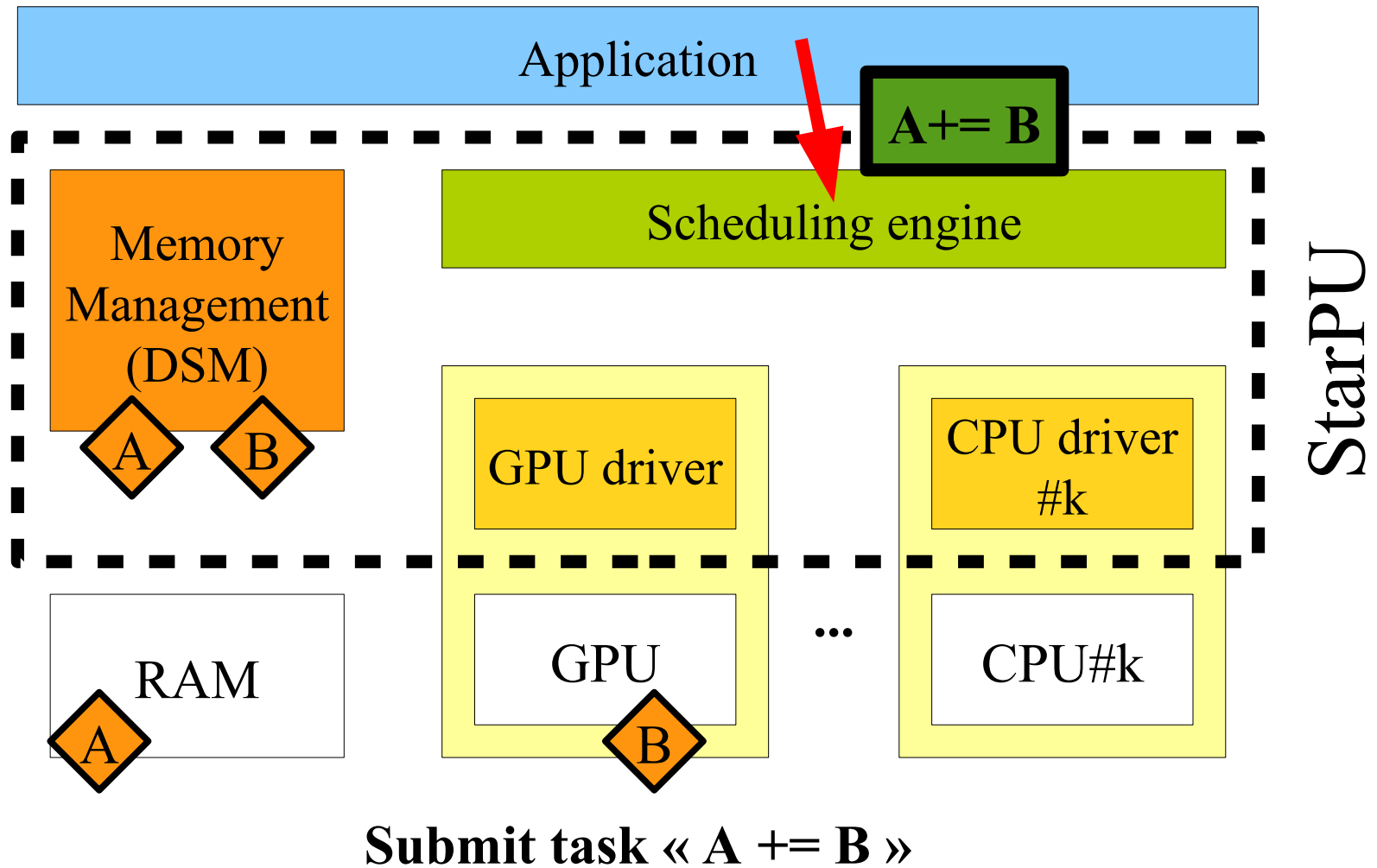
The StarPU runtime system

Execution model



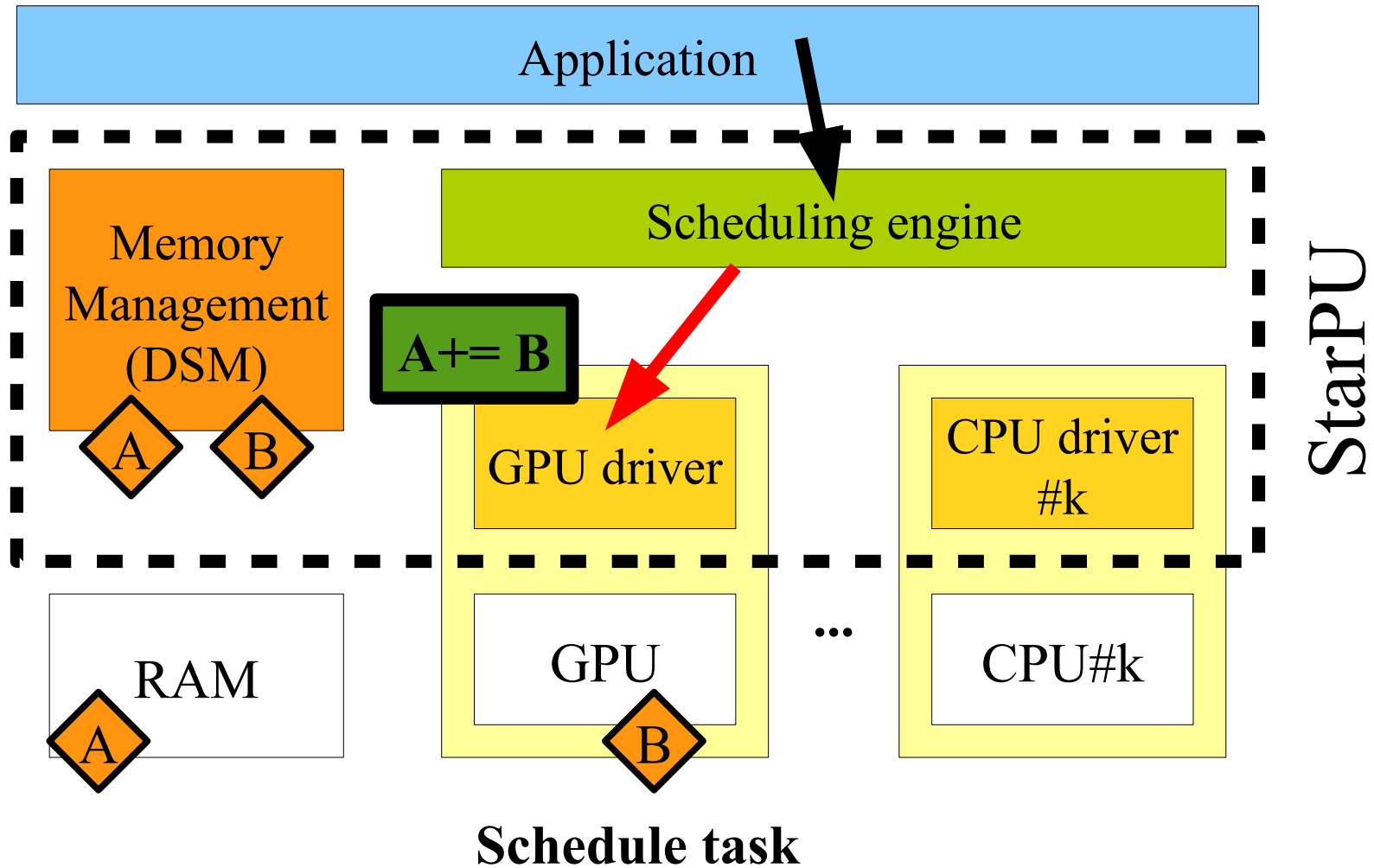
The StarPU runtime system

Execution model



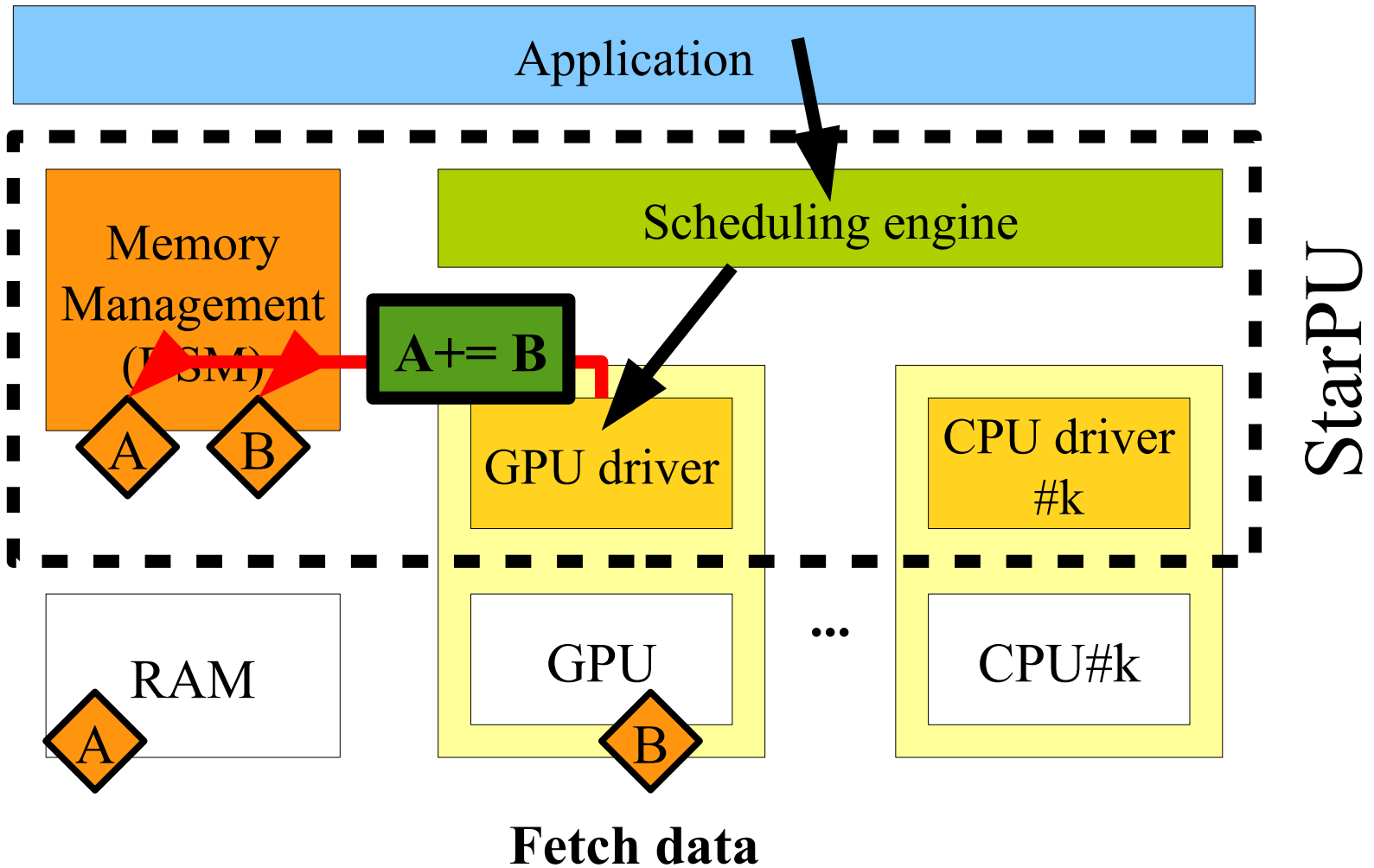
The StarPU runtime system

Execution model



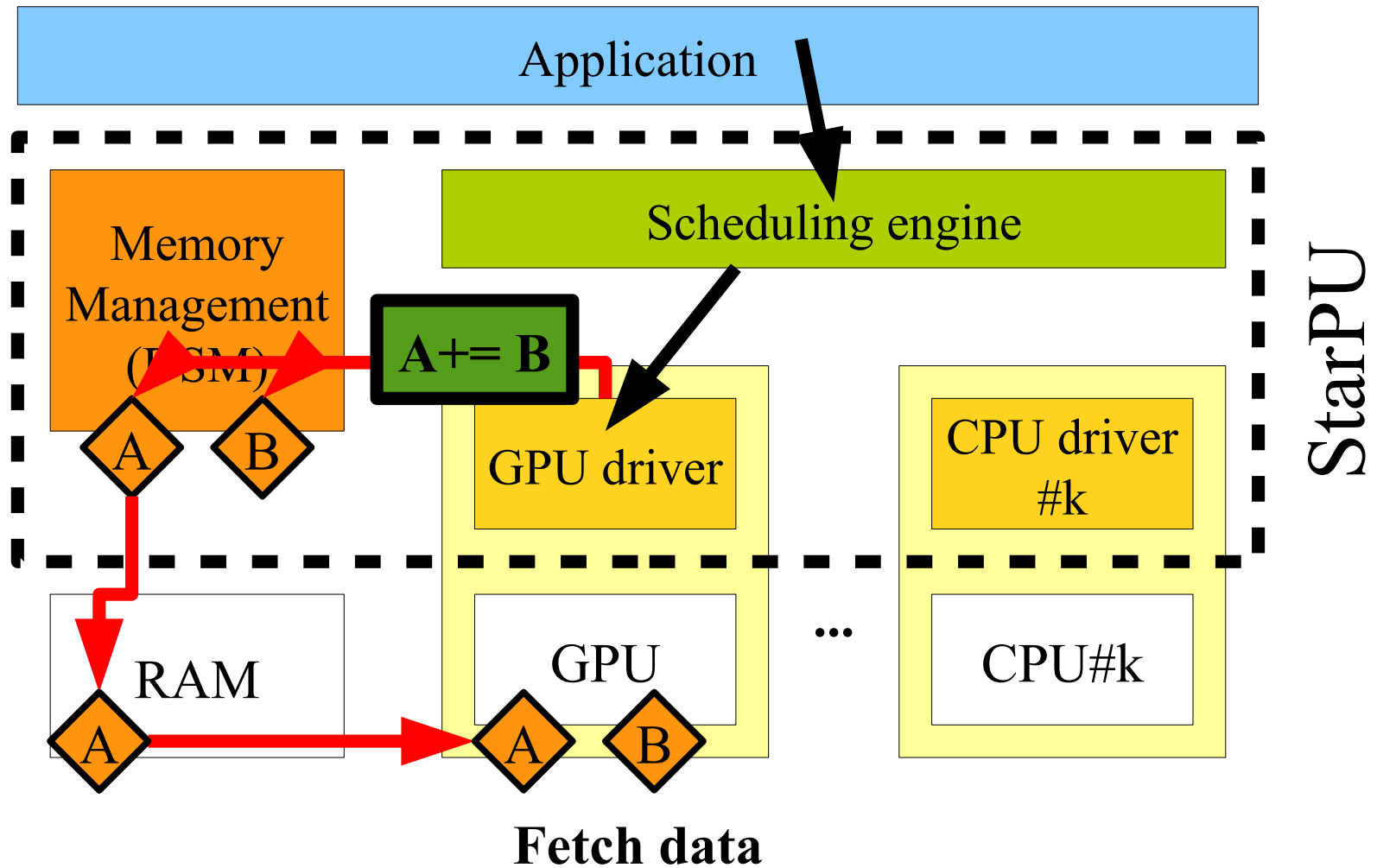
The StarPU runtime system

Execution model



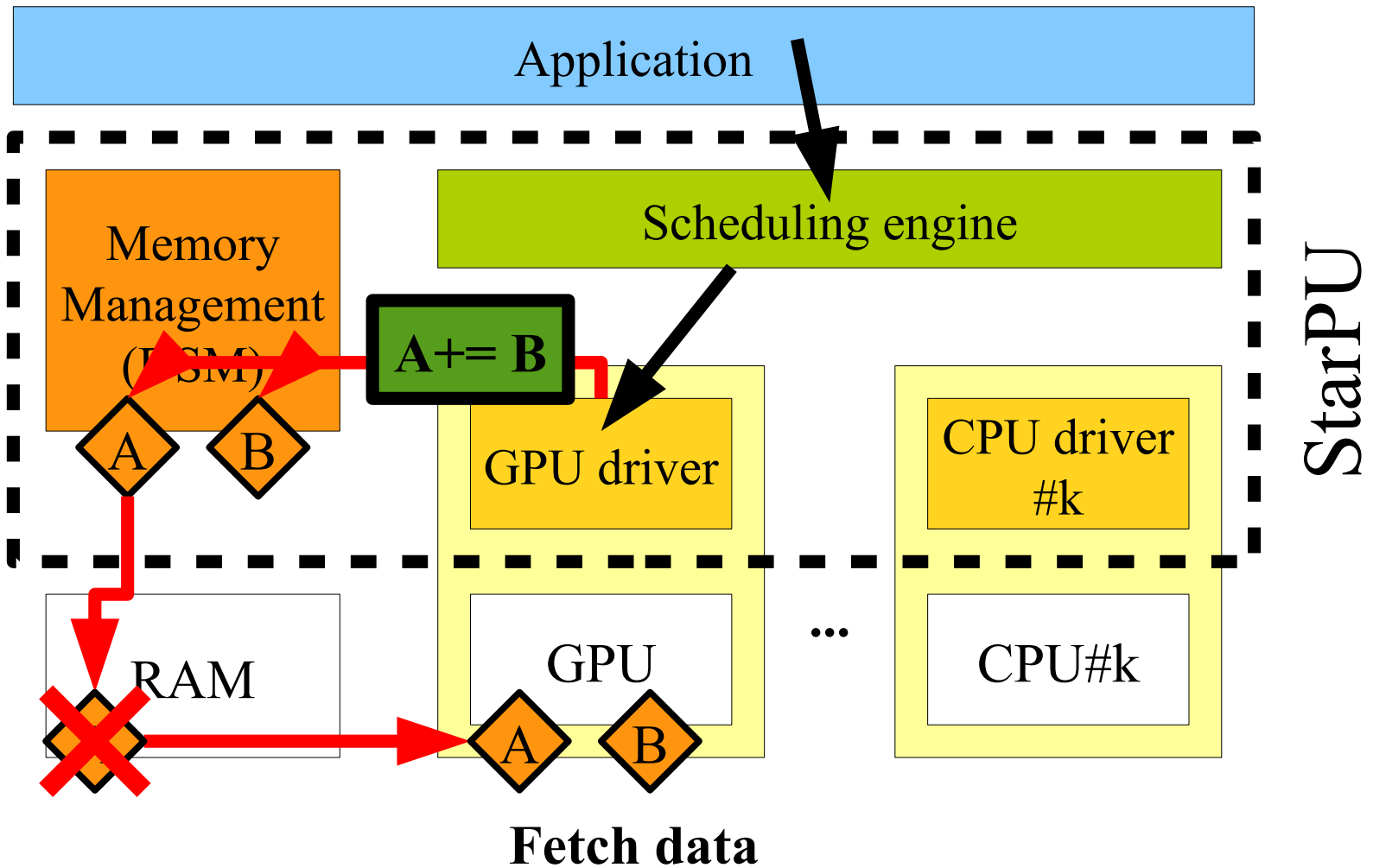
The StarPU runtime system

Execution model



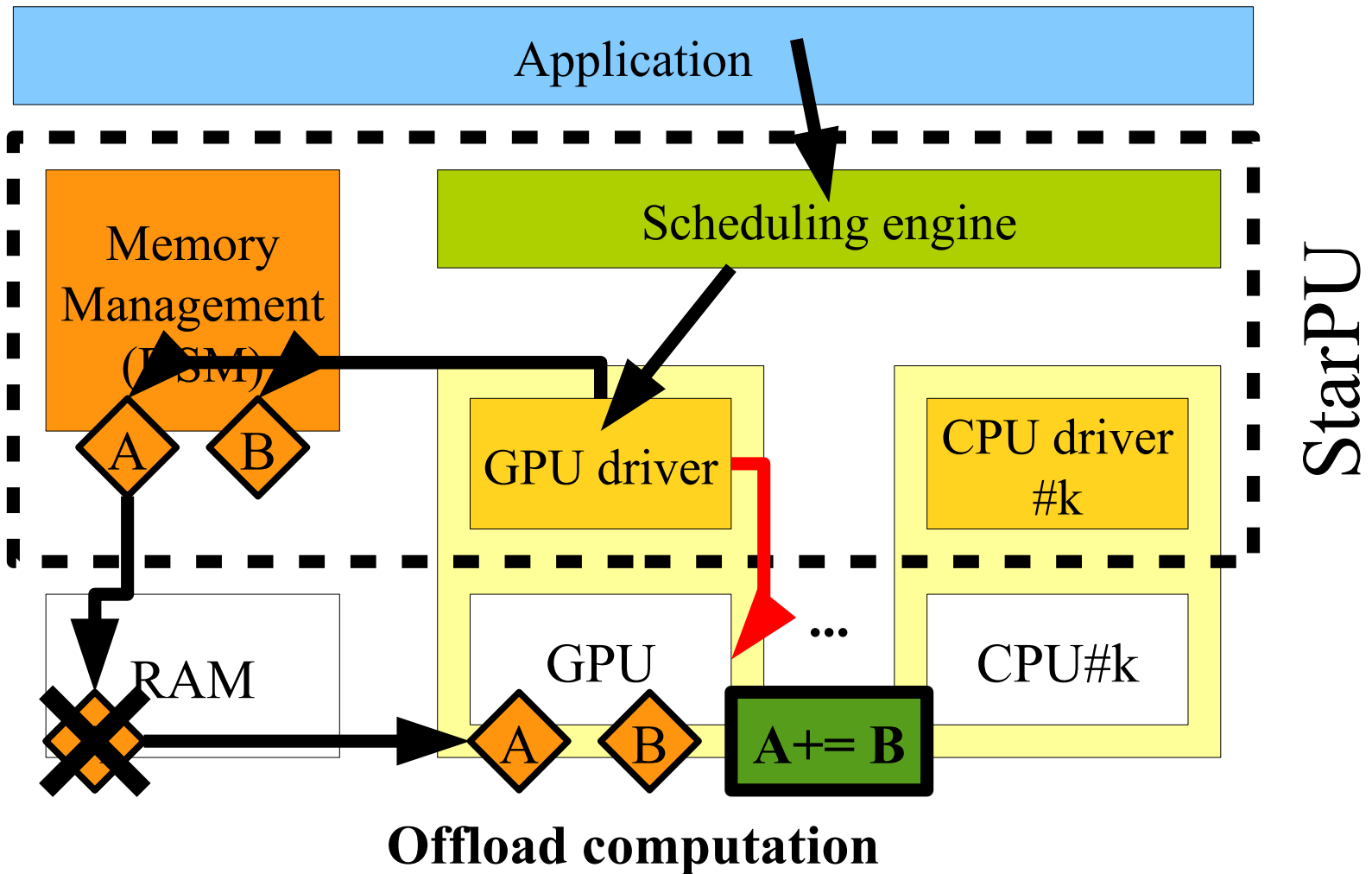
The StarPU runtime system

Execution model



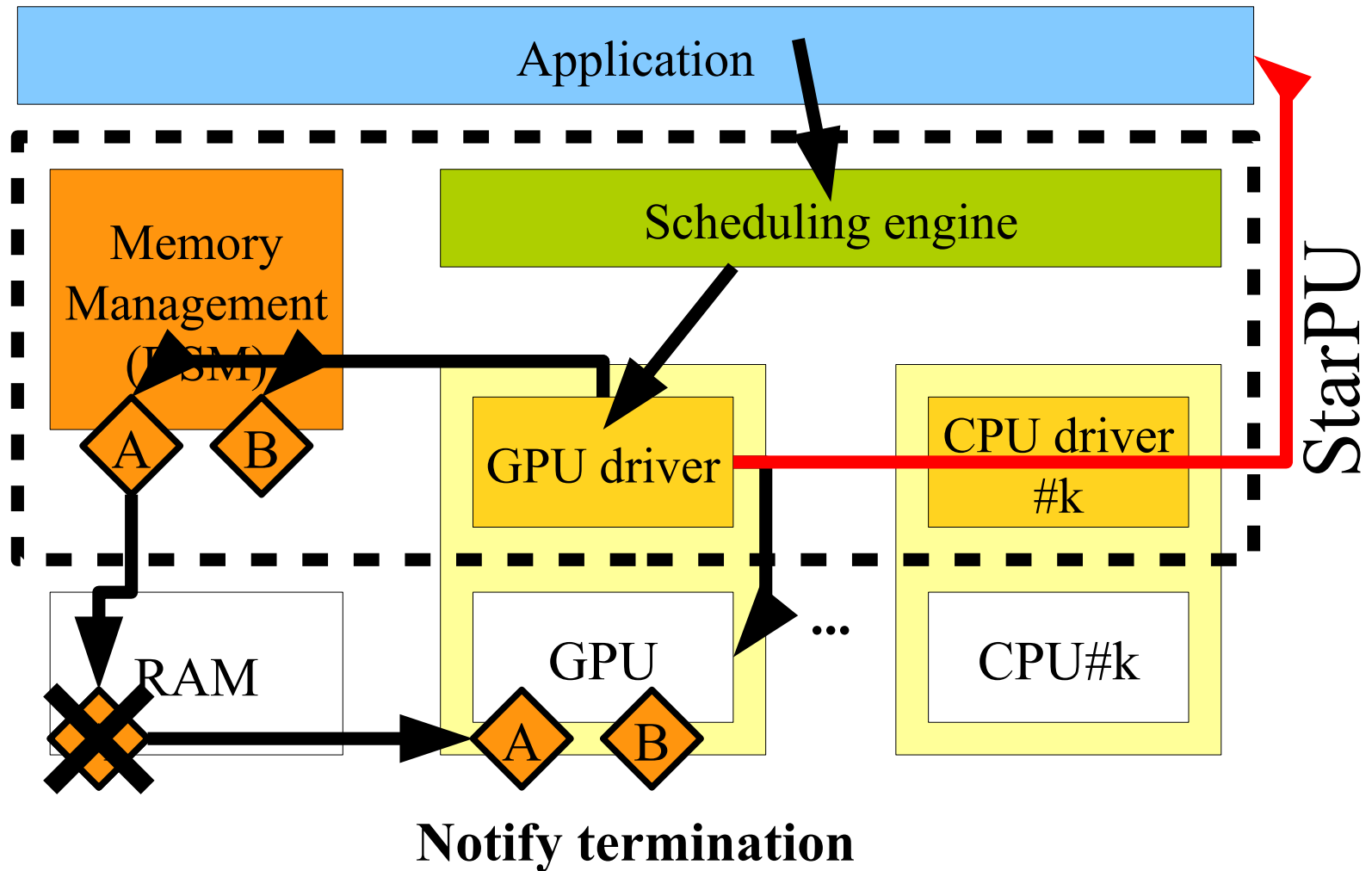
The StarPU runtime system

Execution model



The StarPU runtime system

Execution model



The StarPU runtime system

Development context

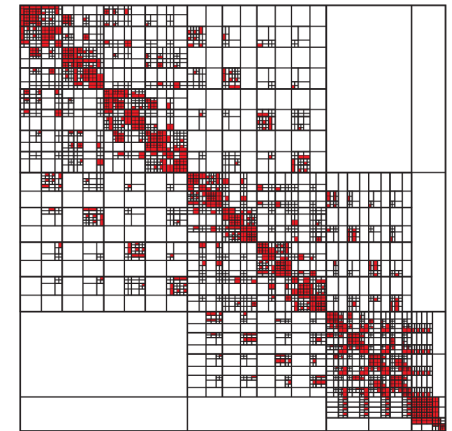
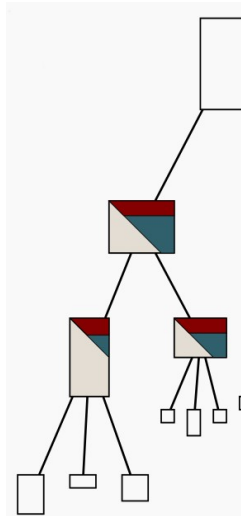
- History
 - Started in 2008
 - PhD Thesis of Cédric Augonnet
 - StarPU main core \approx 70k lines of code
 - Written in C
- Open Source
 - Released under LGPL
 - Sources freely available
 - git repository and nightly tarballs
 - See <https://starpu.gitlabpages.inria.fr/>
 - Open to external contributors
- [HPPC'08]
- [Europar'09] – [CCPE'11],... >1500 citations

The StarPU runtime system

Success stories

Task-based programming actually makes things easier!

- QR-Mumps (sparse linear algebra)
 - Non-task version: only 1D decomposition
 - Task version: 2D decomposition, flurry of parallelism
 - With seamless memory control
- H-Matrices (compressed linear algebra, Airbus)
 - Out-of-core support
 - Could run cases unachievable before
 - e.g. 1600 GB matrix with 256 GB memory
 - Shipped to Airbus customers
- Implemented CFD, FMM, CG, stencils, ...



The StarPU runtime system

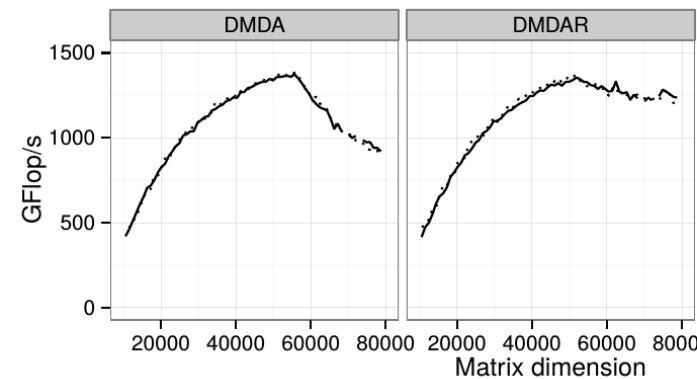
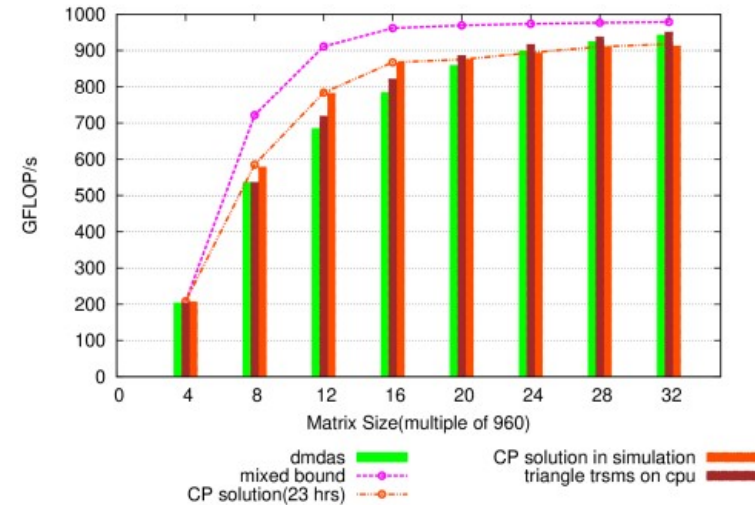
Supported platforms

- Supported architectures
 - Multicore CPUs (x86, PPC, ...)
 - NVIDIA GPUs
 - OpenCL devices (eg. AMD cards)
 - HIP
 - FPGA (ongoing)
 - Old Intel Xeon Phi (MIC), SCC, Kalray MPPA, Cell (decommissioned)
- Supported Operating Systems
 - Linux
 - Mac OS
 - Windows

Task-based support

Then all of this comes “for free” :

- Task/data scheduling
 - Pipelining
 - Load balancing
 - GPU memory limitation management
 - Data prefetching
- Performance bounds
- Distributed execution through MPI
- High-level performance analysis
- Out-of-core : optimized swapping to disk
- Debugging sequential execution
- Reproducible performance simulation



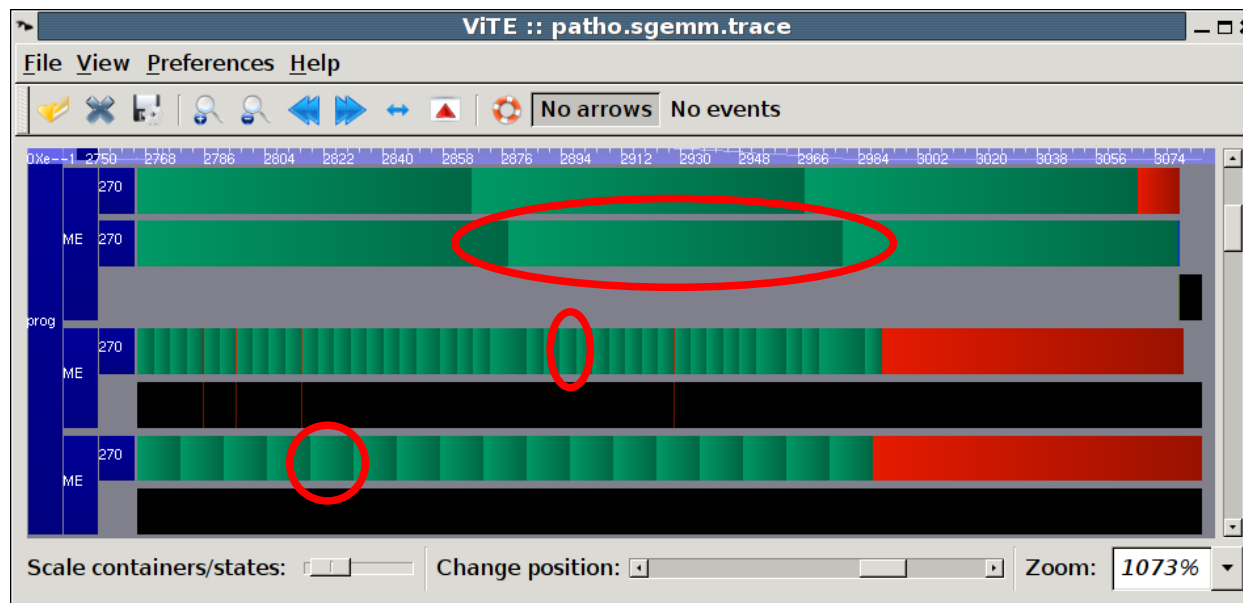
Task Scheduling

Why do we need task scheduling ?

Blocked Matrix multiplication

Things can go (really) wrong even on trivial problems !

- Static mapping ?
 - Not portable, too hard for real-life problems
- Need Dynamic Task Scheduling
 - Performance models



2 Xeon cores

Quadro FX5800

Quadro FX4600

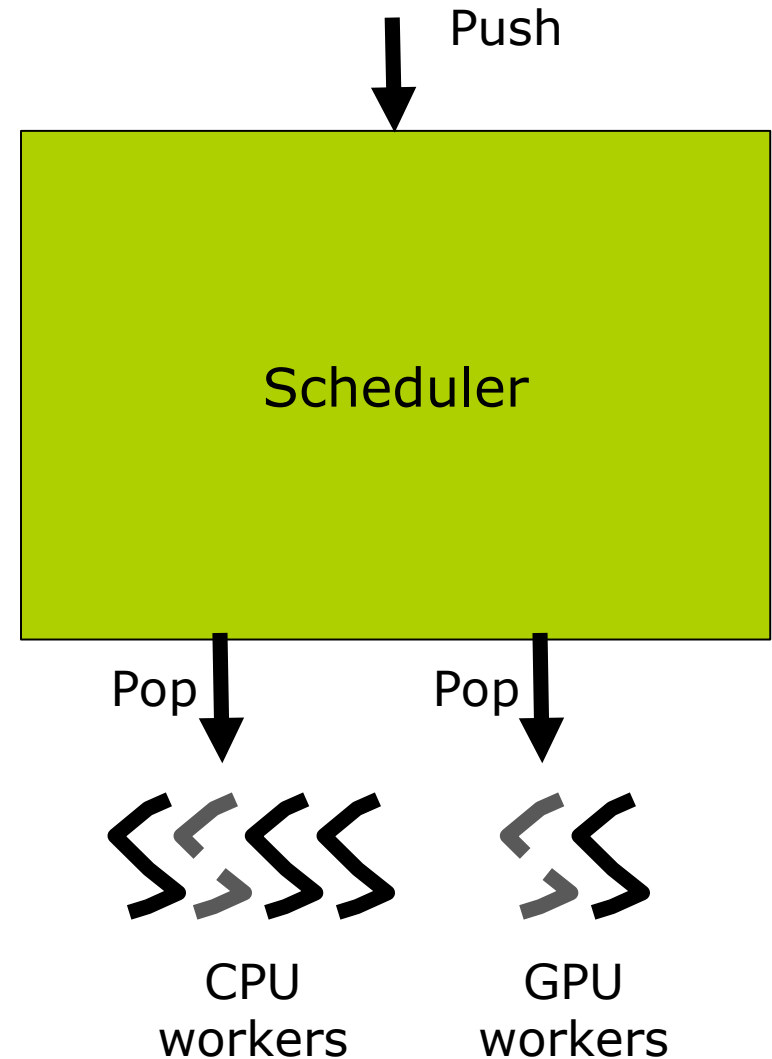
Runtime-based task scheduling

When a task is submitted, it first goes into a pool of “frozen tasks” until all dependencies are met

Then, the task is “pushed” to the scheduler

Idle processing units poll for work (“pop”)

Various scheduling policies, can even be user-defined



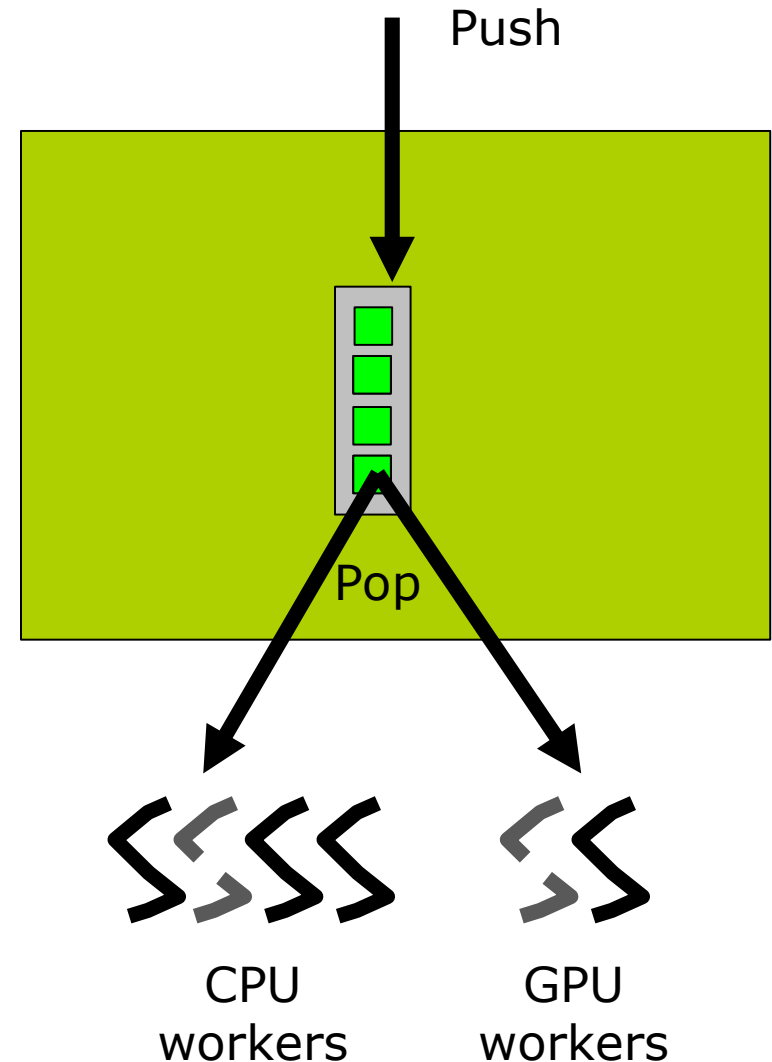
Runtime-based task scheduling

When a task is submitted, it first goes into a pool of “frozen tasks” until all dependencies are met

Then, the task is “pushed” to the scheduler

Idle processing units poll for work (“pop”)

Various scheduling policies, can even be user-defined



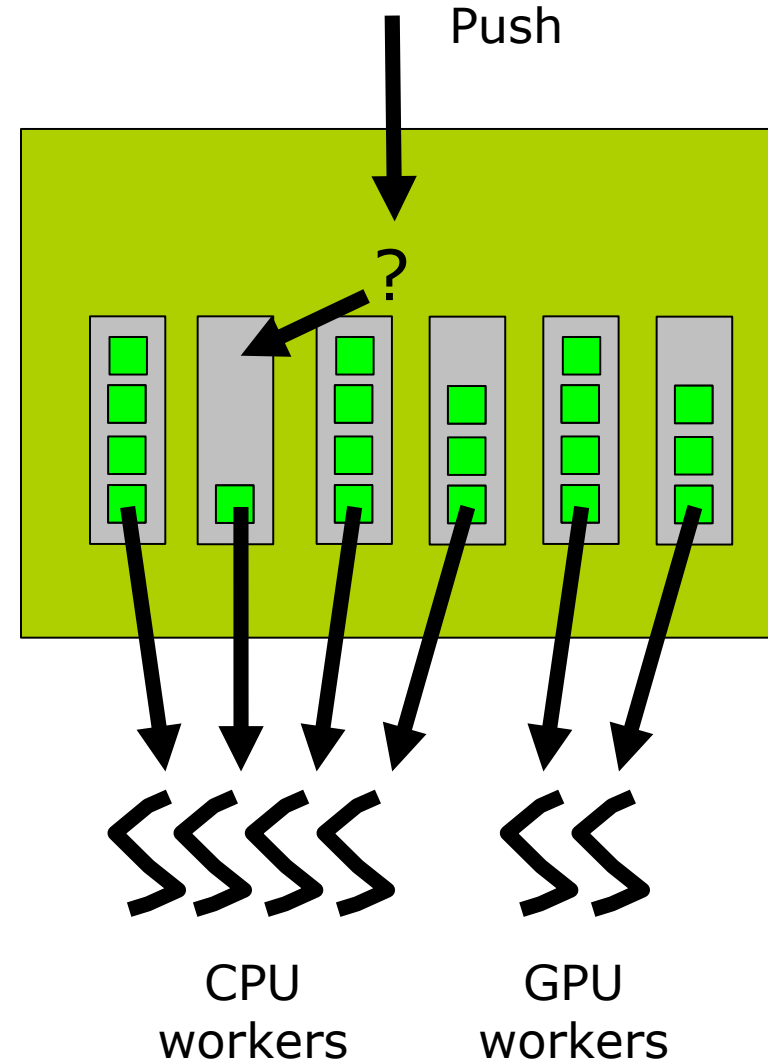
Runtime-based task scheduling

When a task is submitted, it first goes into a pool of “frozen tasks” until all dependencies are met

Then, the task is “pushed” to the scheduler

Idle processing units poll for work (“pop”)

Various scheduling policies, can even be user-defined



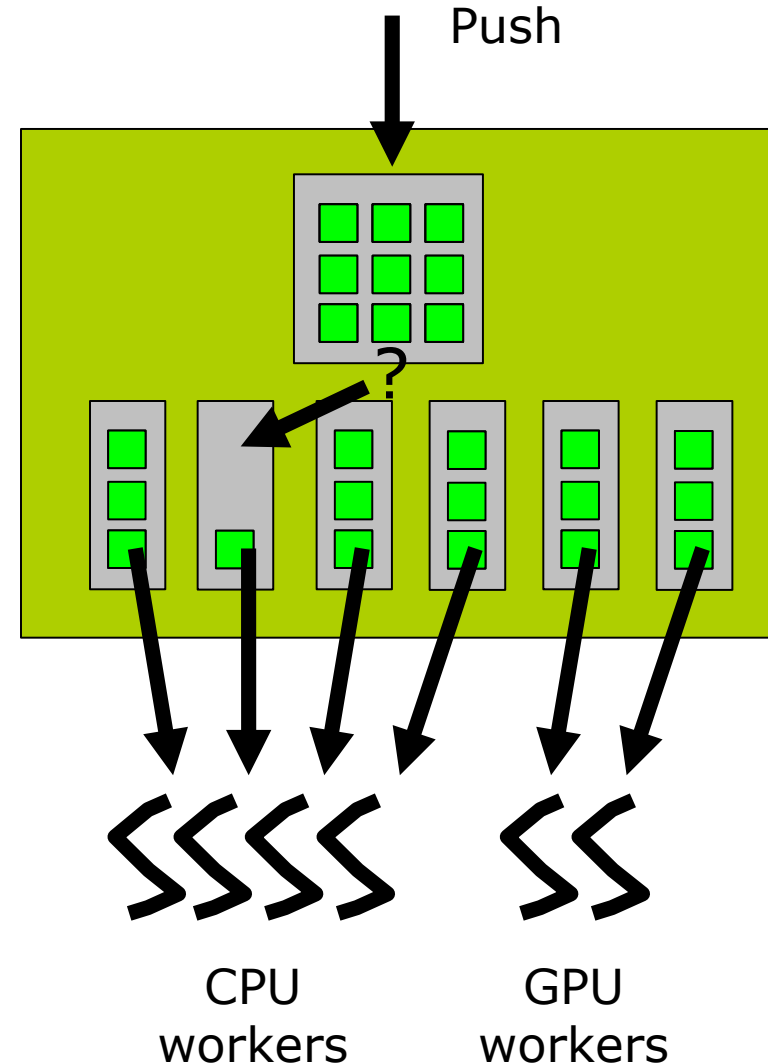
Runtime-based task scheduling

When a task is submitted, it first goes into a pool of “frozen tasks” until all dependencies are met

Then, the task is “pushed” to the scheduler

Idle processing units poll for work (“pop”)

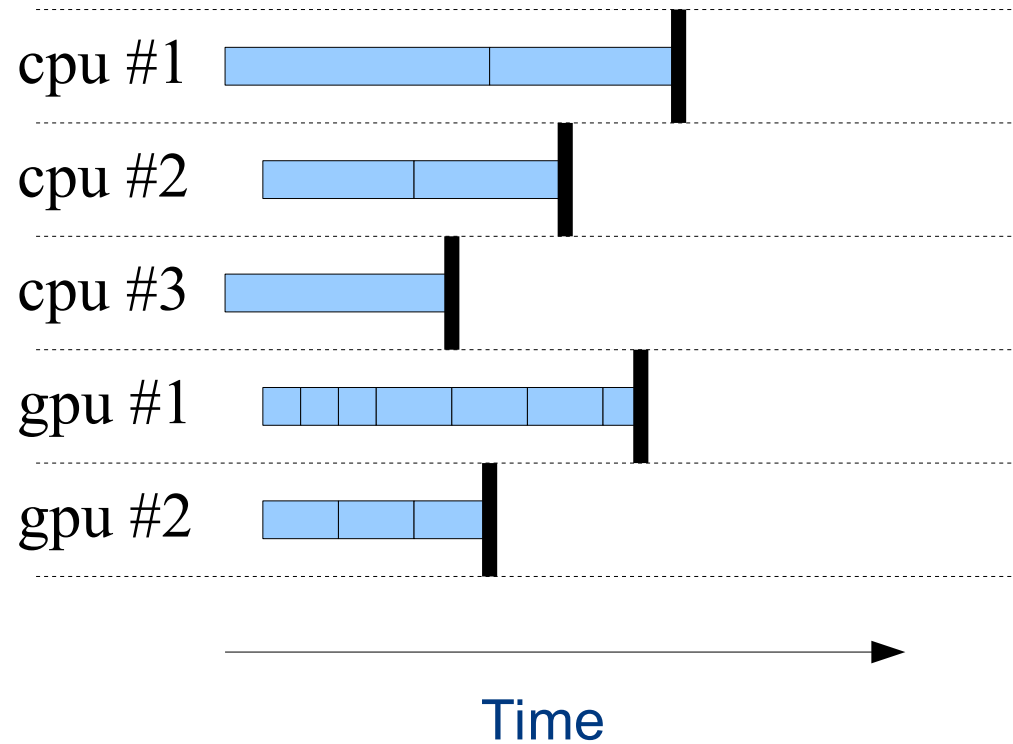
Various scheduling policies, can even be user-defined



Prediction-based scheduling

Load balancing

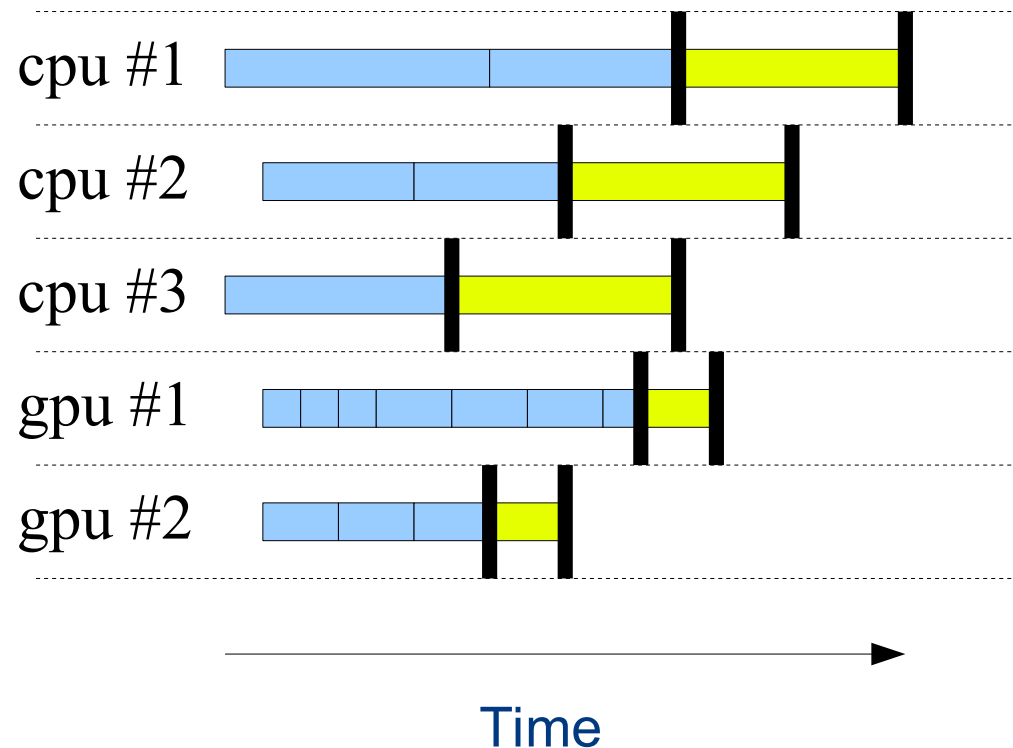
- Task completion time estimation
 - History-based
 - User-defined cost function
 - Parametric cost model
 - [HPPC'09]
- Can be used to implement scheduling
 - E.g. Heterogeneous Earliest Finish Time



Prediction-based scheduling

Load balancing

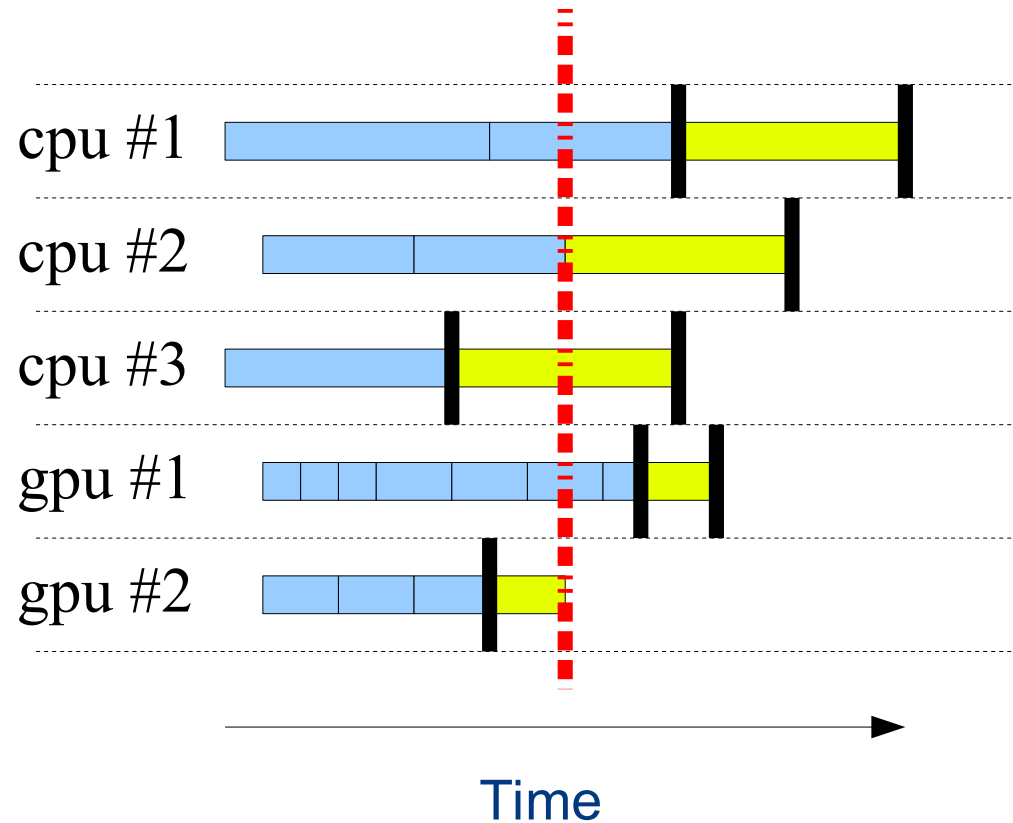
- Task completion time estimation
 - History-based
 - User-defined cost function
 - Parametric cost model
 - [HPPC'09]
- Can be used to implement scheduling
 - E.g. Heterogeneous Earliest Finish Time



Prediction-based scheduling

Load balancing

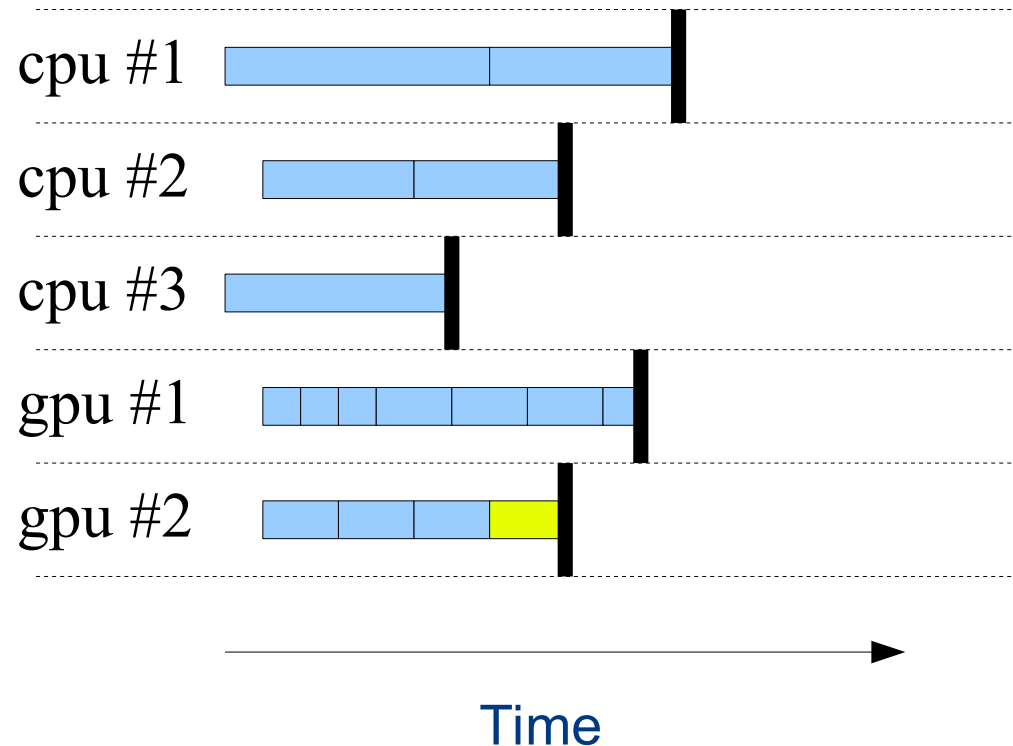
- Task completion time estimation
 - History-based
 - User-defined cost function
 - Parametric cost model
 - [HPPC'09]
- Can be used to implement scheduling
 - E.g. Heterogeneous Earliest Finish Time



Prediction-based scheduling

Load balancing

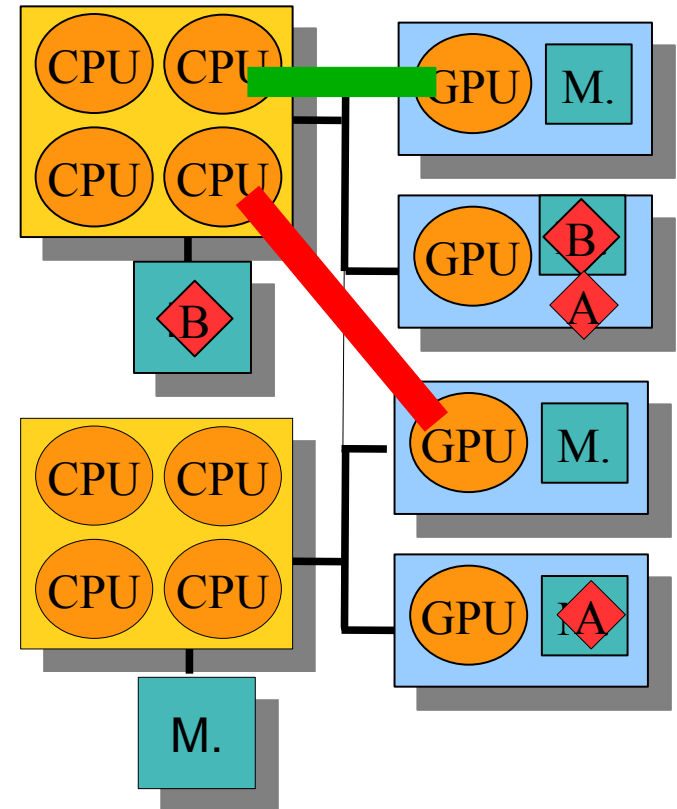
- Task completion time estimation
 - History-based
 - User-defined cost function
 - Parametric cost model
 - [HPPC'09]
- Can be used to implement scheduling
 - E.g. Heterogeneous Earliest Finish Time



Predicting data transfer overhead

Motivations

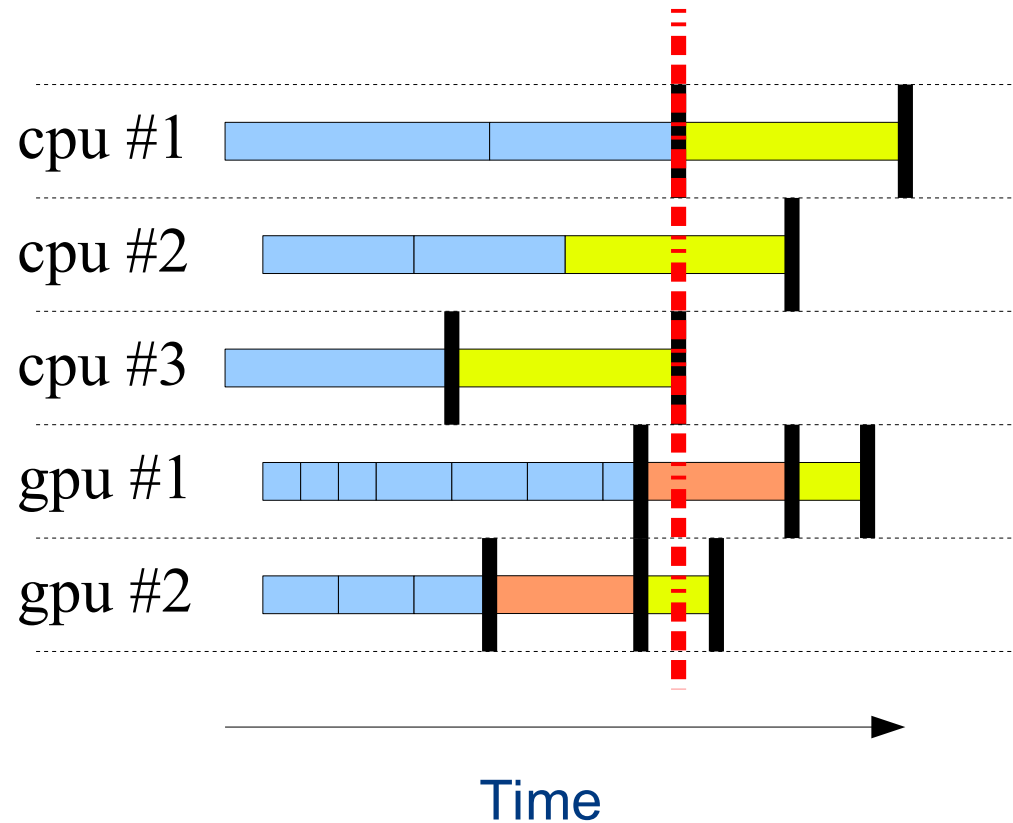
- Hybrid platforms
 - Multicore CPUs and GPUs
 - PCI-e bus is a precious resource
- Data locality vs. Load balancing
 - Cannot avoid all data transfers
 - Minimize them
- StarPU keeps track of
 - data replicates
 - on-going data movements



Prediction-based scheduling

Load balancing

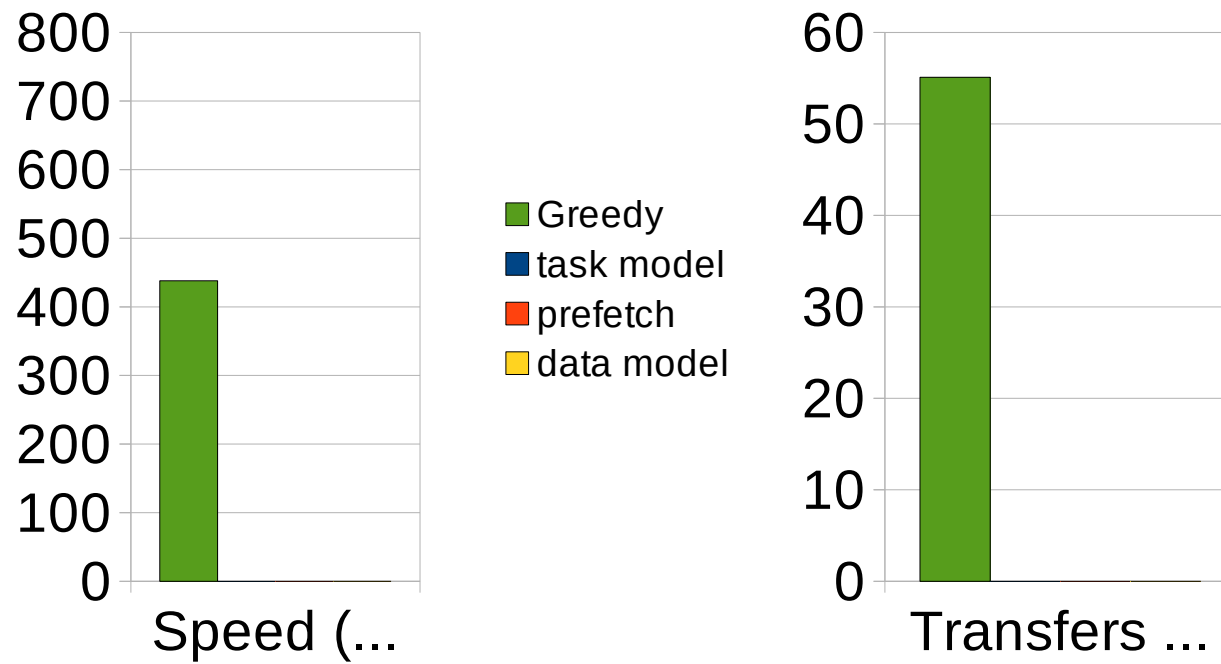
- Data transfer time
 - Sampling based on off-line calibration
- Can be used to
 - Better estimate overall exec time
 - Minimize data movements
- Further
 - Power overhead
- **dmda** [ICPADS'10]



Scheduling in a hybrid environment

Performance models

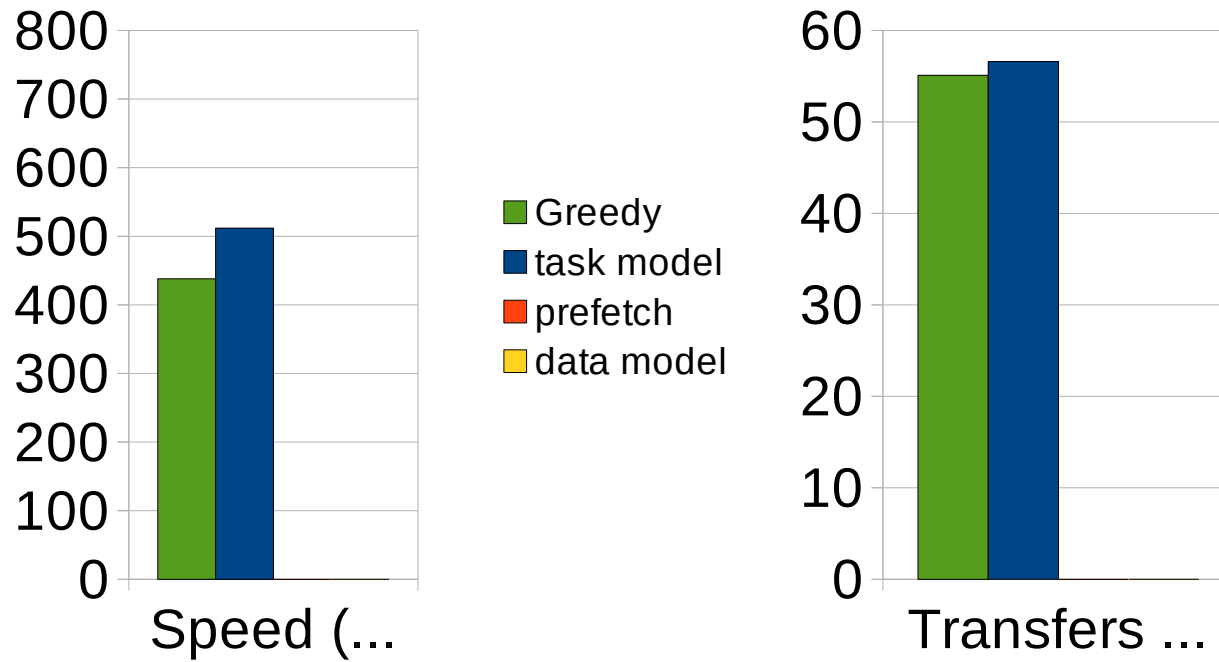
- LU without pivoting (16GB input matrix)
 - 8 CPUs (nehalem) + 3 GPUs (FX5800)



Scheduling in a hybrid environment

Performance models

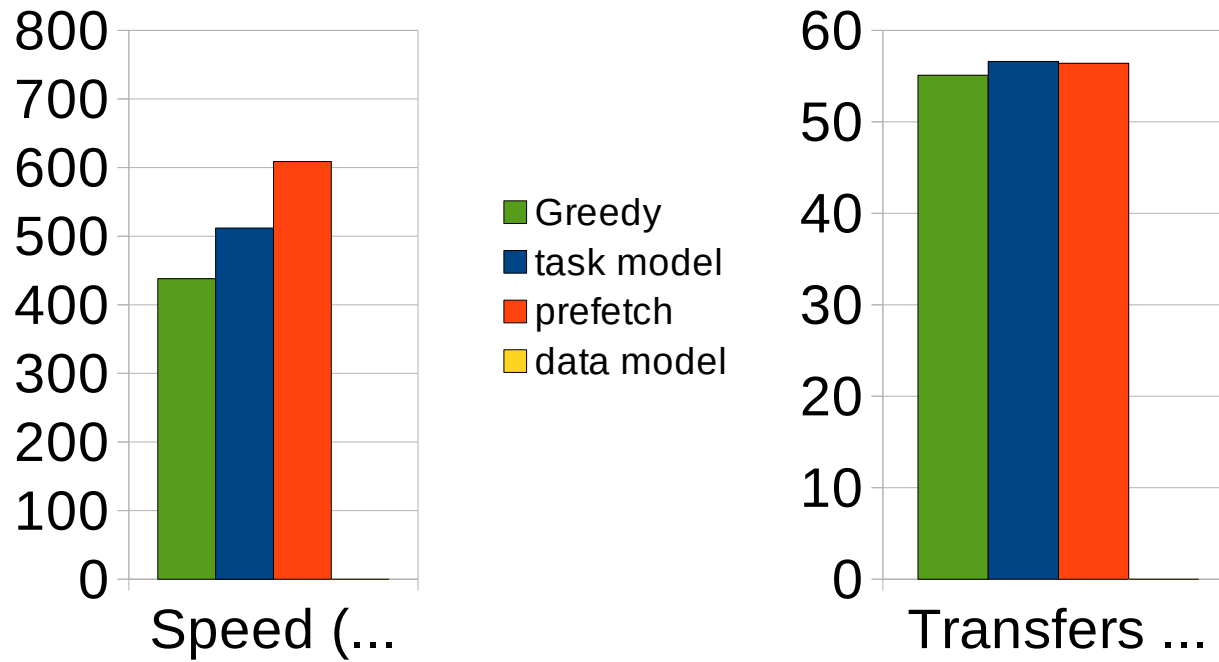
- LU without pivoting (16GB input matrix)
 - 8 CPUs (nehalem) + 3 GPUs (FX5800)



Scheduling in a hybrid environment

Performance models

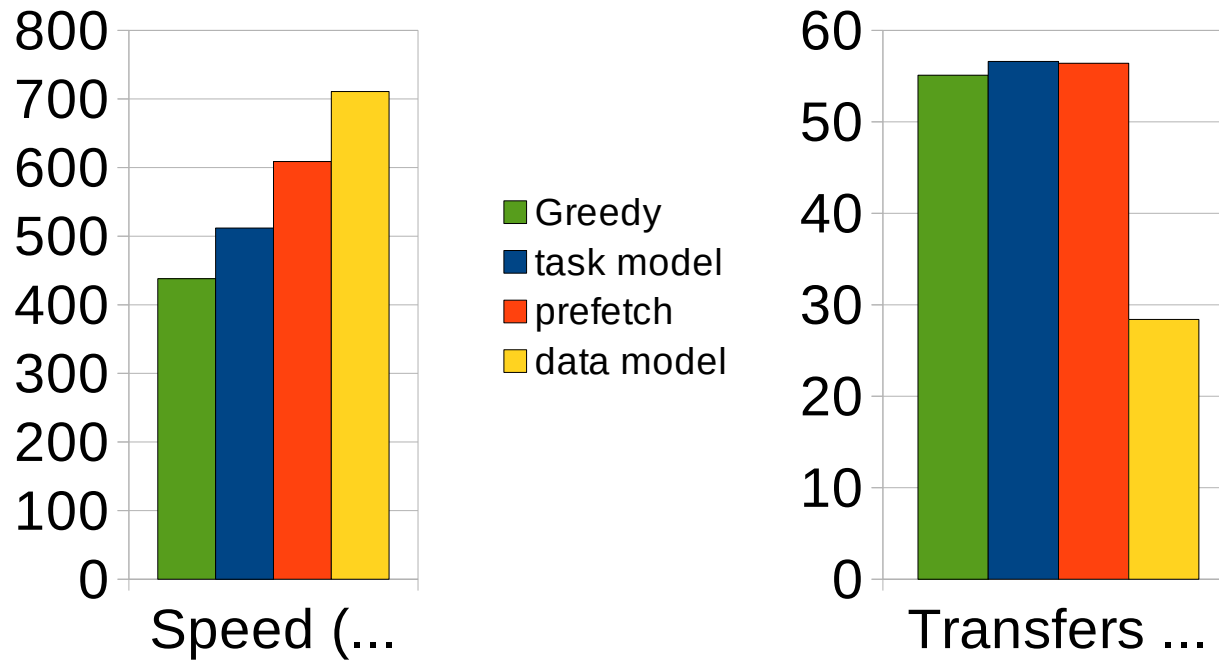
- LU without pivoting (16GB input matrix)
 - 8 CPUs (nehalem) + 3 GPUs (FX5800)



Scheduling in a hybrid environment

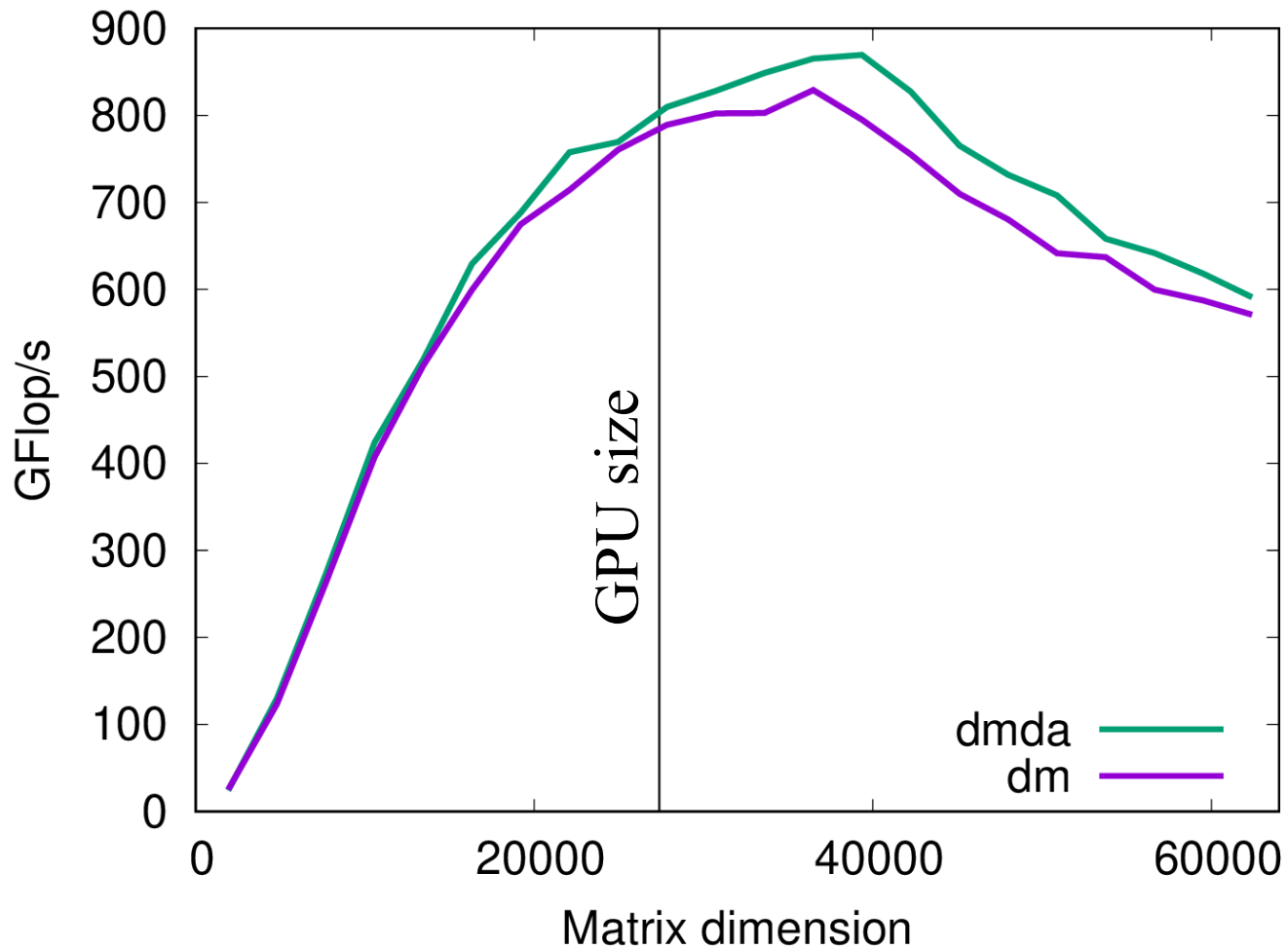
Performance models

- LU without pivoting (16GB input matrix)
 - 8 CPUs (nehalem) + 3 GPUs (FX5800)



LU factorization (no pivoting)

- 3 FX5800 GPUs



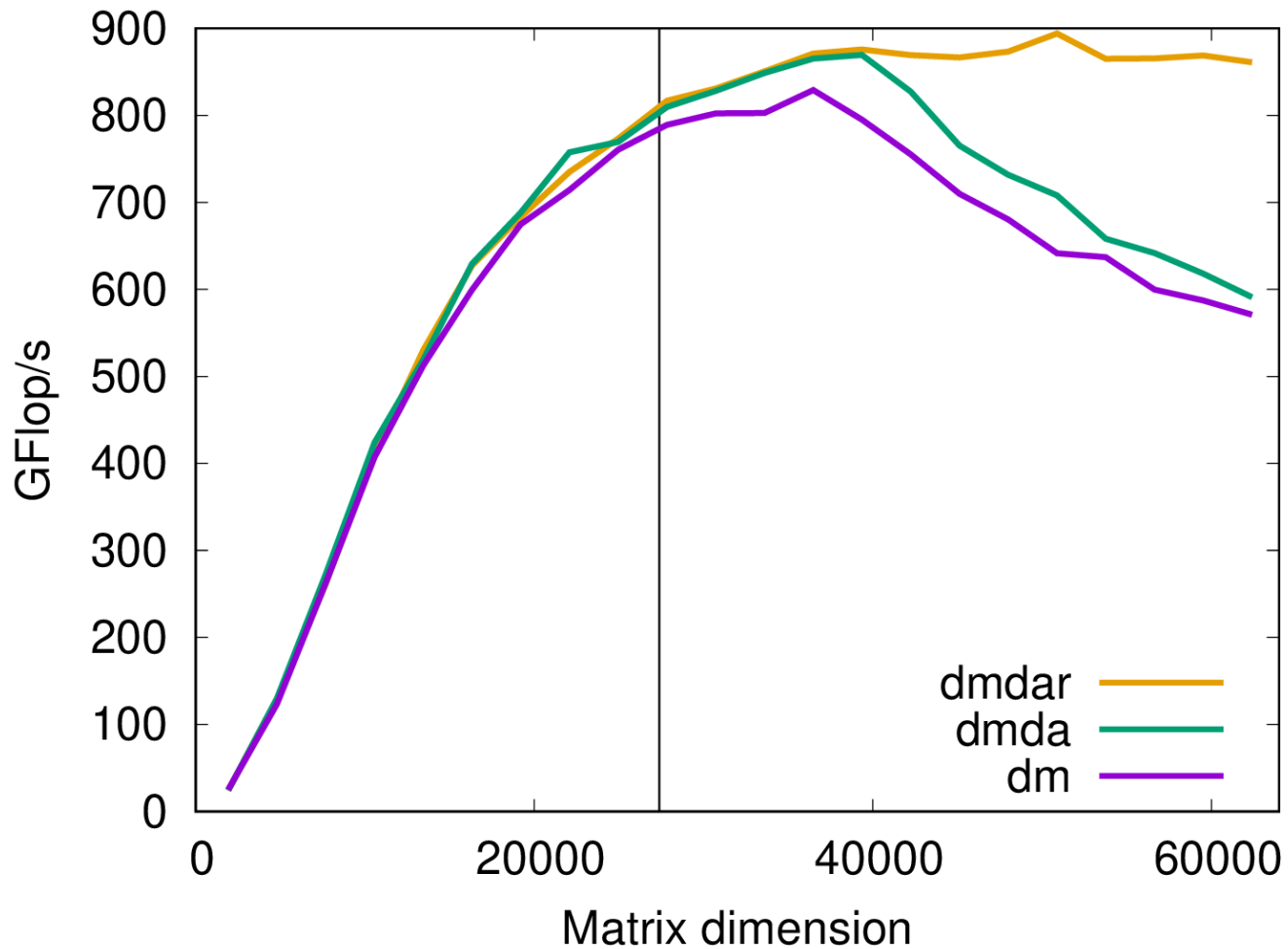
LU factorization (no pivoting)

Locality not so good

- dmda strictly follows task order
 - only decides placement
- dmdar
 - Follows task order
 - But cherry-picks tasks whose data is available (**ready**)
 - Contradict task order with data locality

LU factorization (no pivoting)

- 3 FX5800 GPUs



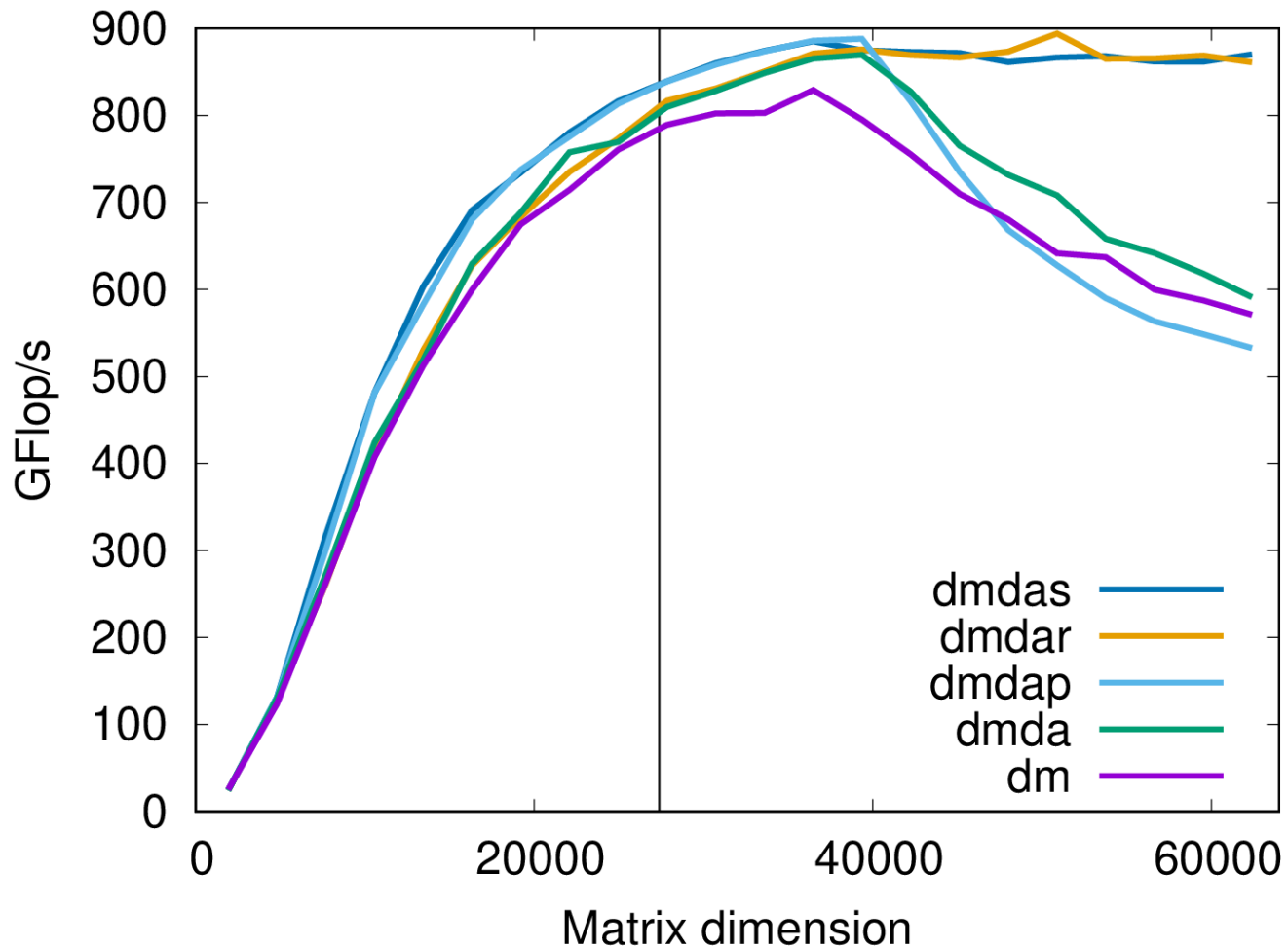
LU factorization (no pivoting)

Not completely heft

- dmda/dmdar do not fast-track tasks with high priority
- dmdap/dmdas keep lists sorted by priority
 - Not contradicted by data locality

LU factorization (no pivoting)

- 3 FX5800 GPUs



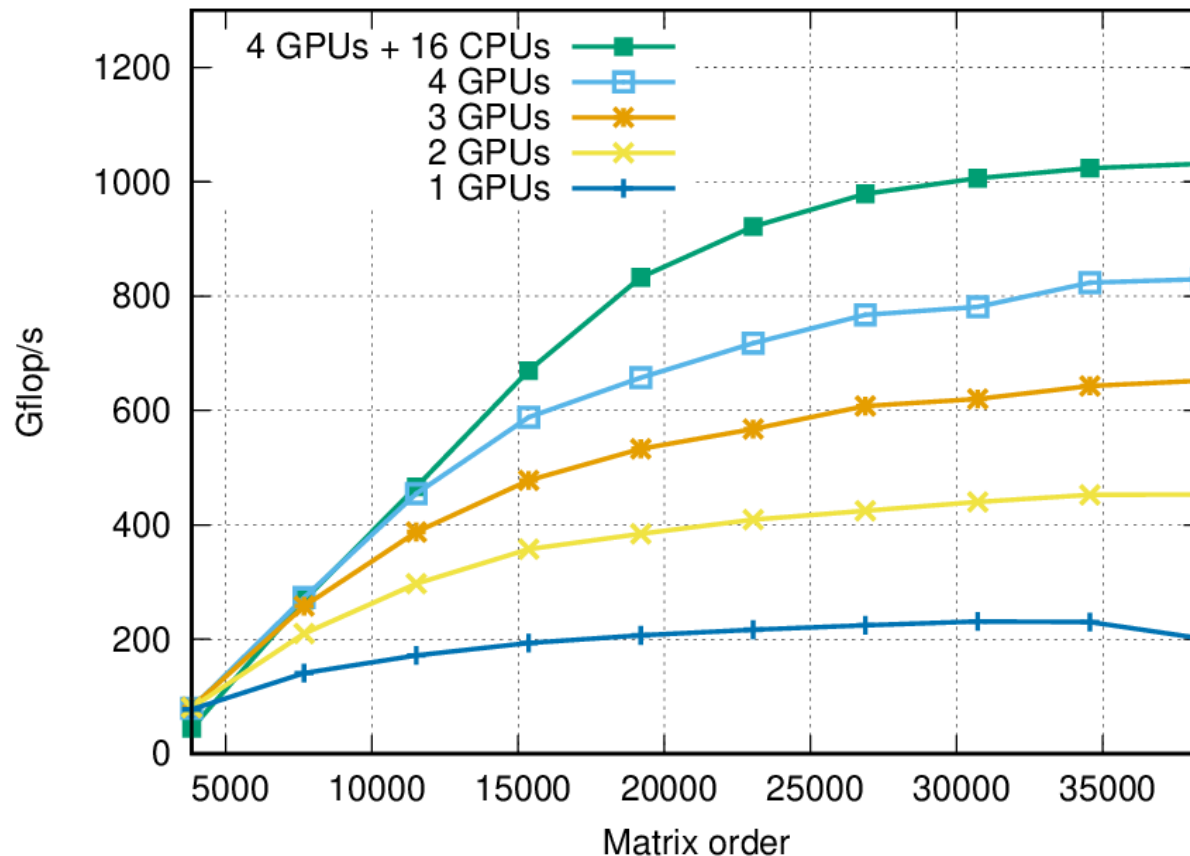
Mixing PLASMA and MAGMA with StarPU

Mixing PLASMA and MAGMA with StarPU

- State of the art algorithms
 - PLASMA (Multicore CPUs)
 - Dynamically scheduled with Quark
 - MAGMA (Multiple GPUs)
 - Hand-coded data transfers
 - Static task mapping
- Design of combination
 - Use PLASMA algorithm with « magnum tiles »
 - PLASMA kernels on CPUs, MAGMA kernels on GPUs
 - Replace the QUARK scheduler with StarPU
- Programmability
 - Cholesky: ~half a week [SAAHPC'10]
 - QR : ~2 days of works [IPDPS'11]
 - Quick algorithmic prototyping
- [GPUgems'10], LU [AICCSA'11]

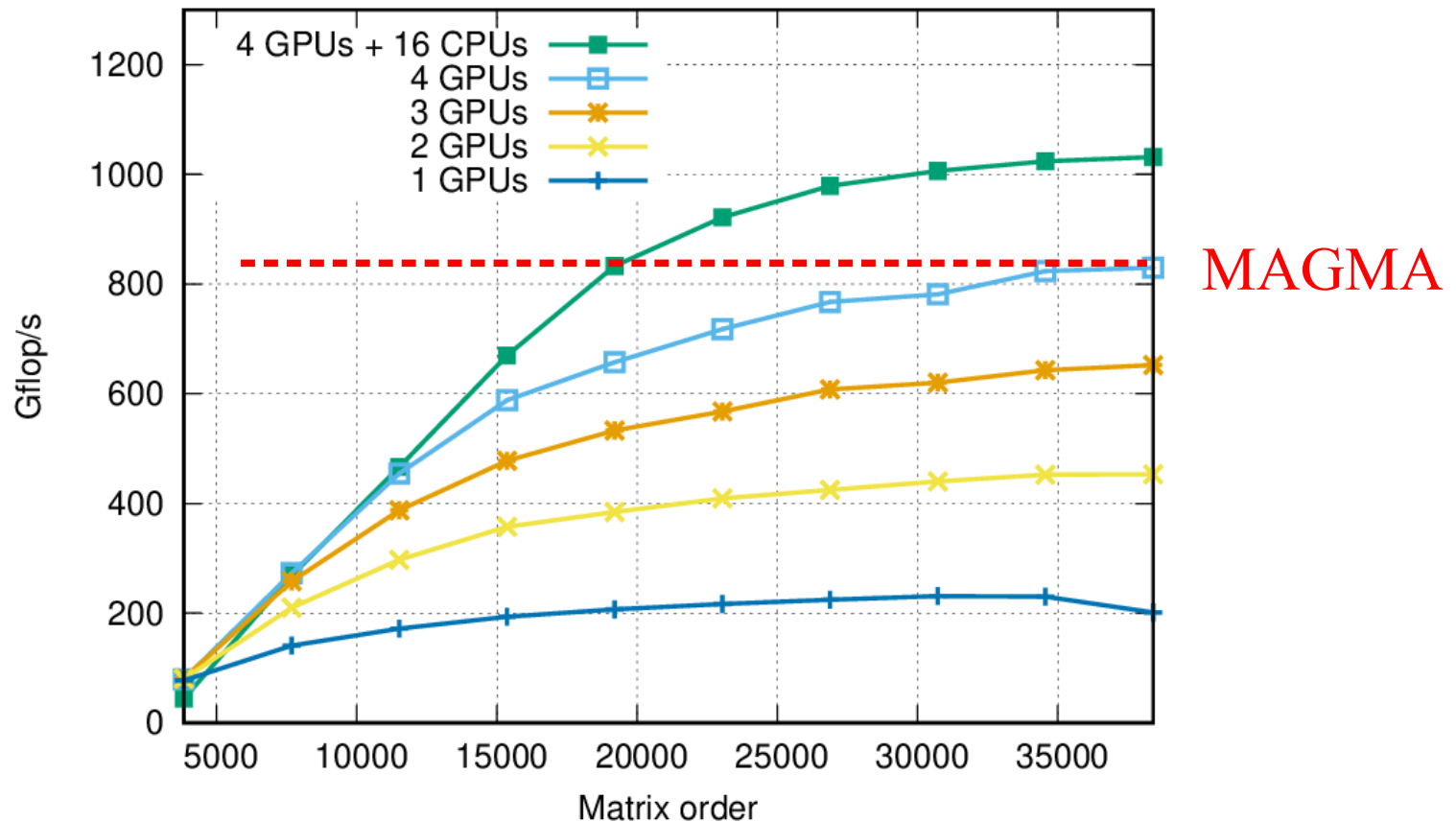
Mixing PLASMA and MAGMA with StarPU

- QR decomposition
 - Mordor8 (UTK) : 16 CPUs (AMD) + 4 GPUs (C1060)



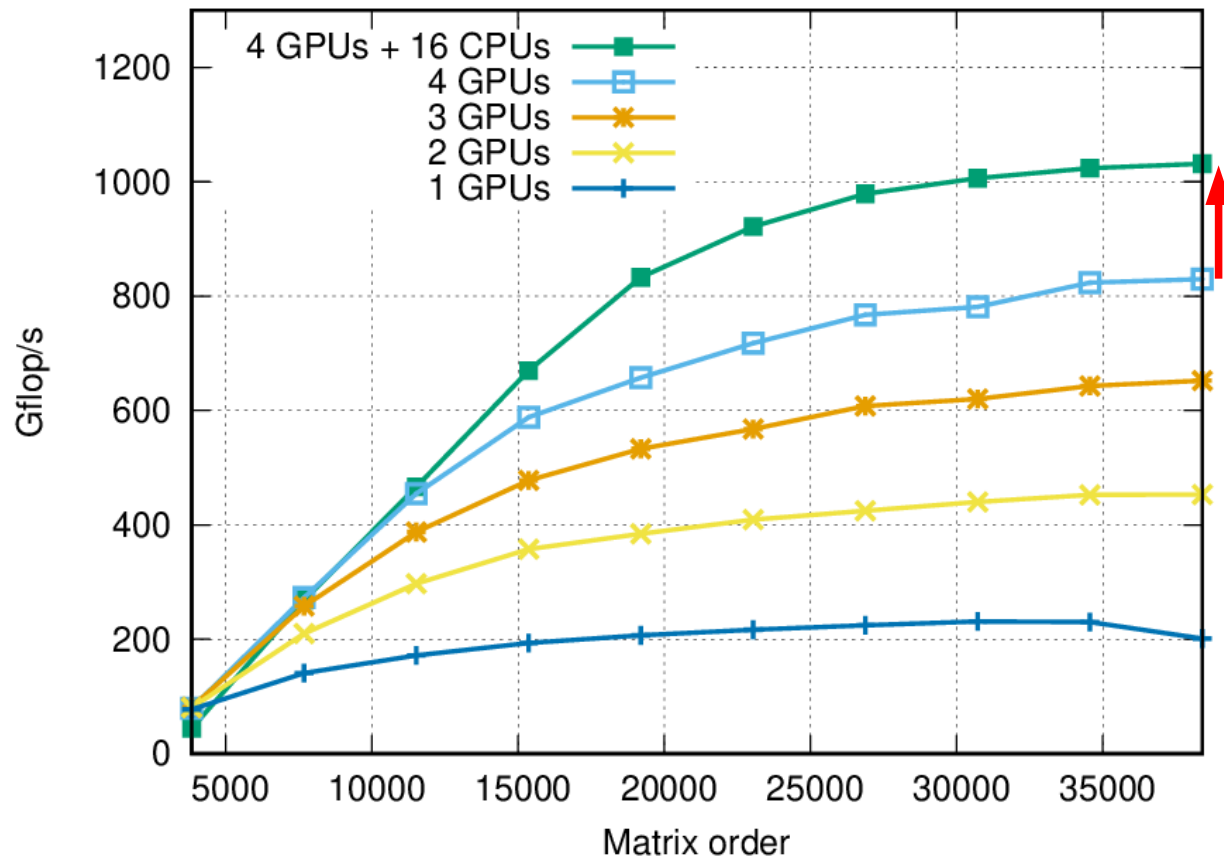
Mixing PLASMA and MAGMA with StarPU

- QR decomposition
 - Mordor8 (UTK) : 16 CPUs (AMD) + 4 GPUs (C1060)



Mixing PLASMA and MAGMA with StarPU

- QR decomposition
 - Mordor8 (UTK) : 16 CPUs (AMD) + 4 GPUs (C1060)



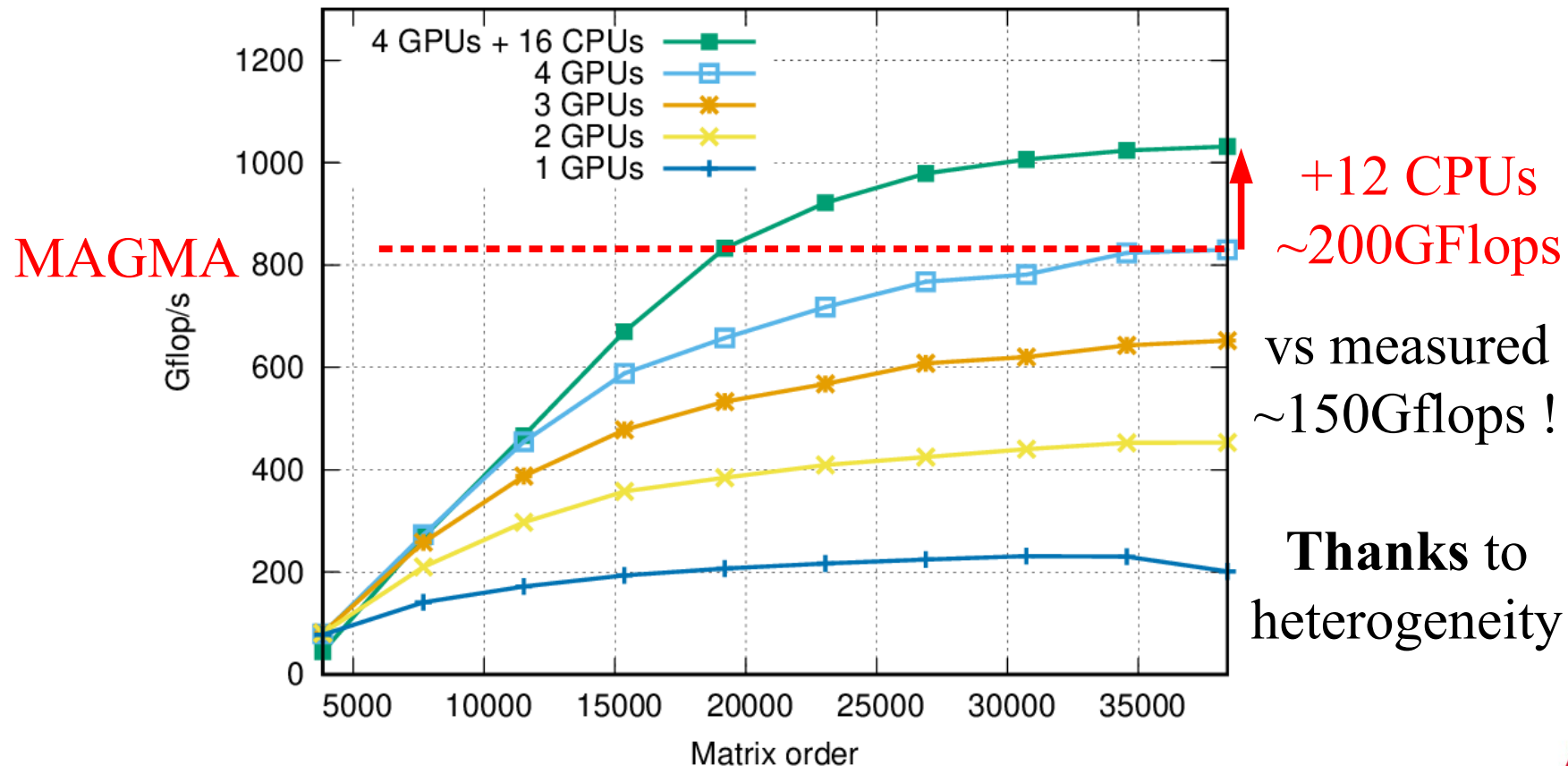
+12 CPUs
~200GFlops

vs measured
~150Gflops !

**Thanks to
heterogeneity**

Mixing PLASMA and MAGMA with StarPU

- QR decomposition
 - Mordor8 (UTK) : 16 CPUs (AMD) + 4 GPUs (C1060)



Mixing PLASMA and MAGMA with StarPU

- « Super-Linear » efficiency in QR?
 - Heterogeneous kernel efficiency
 - sgeqrt
 - CPU: 9 Gflops GPU: 30 Gflops (Speedup : ~3)
 - stsqrt
 - CPU: 12Gflops GPU: 37 Gflops (Speedup: ~3)
 - somqr
 - CPU: 8.5 Gflops GPU: 227 Gflops (Speedup: ~27)
 - Sssmqr
 - CPU: 10Gflops GPU: 285Gflops (Speedup: ~28)
 - Task distribution observed on StarPU
 - sgeqrt: 20% of tasks on GPUs
 - Sssmqr: 92.5% of tasks on GPUs
 - Taking advantage of heterogeneity !
 - Only do what you are good for
 - Don't do what you are not good for

HeteroPrio

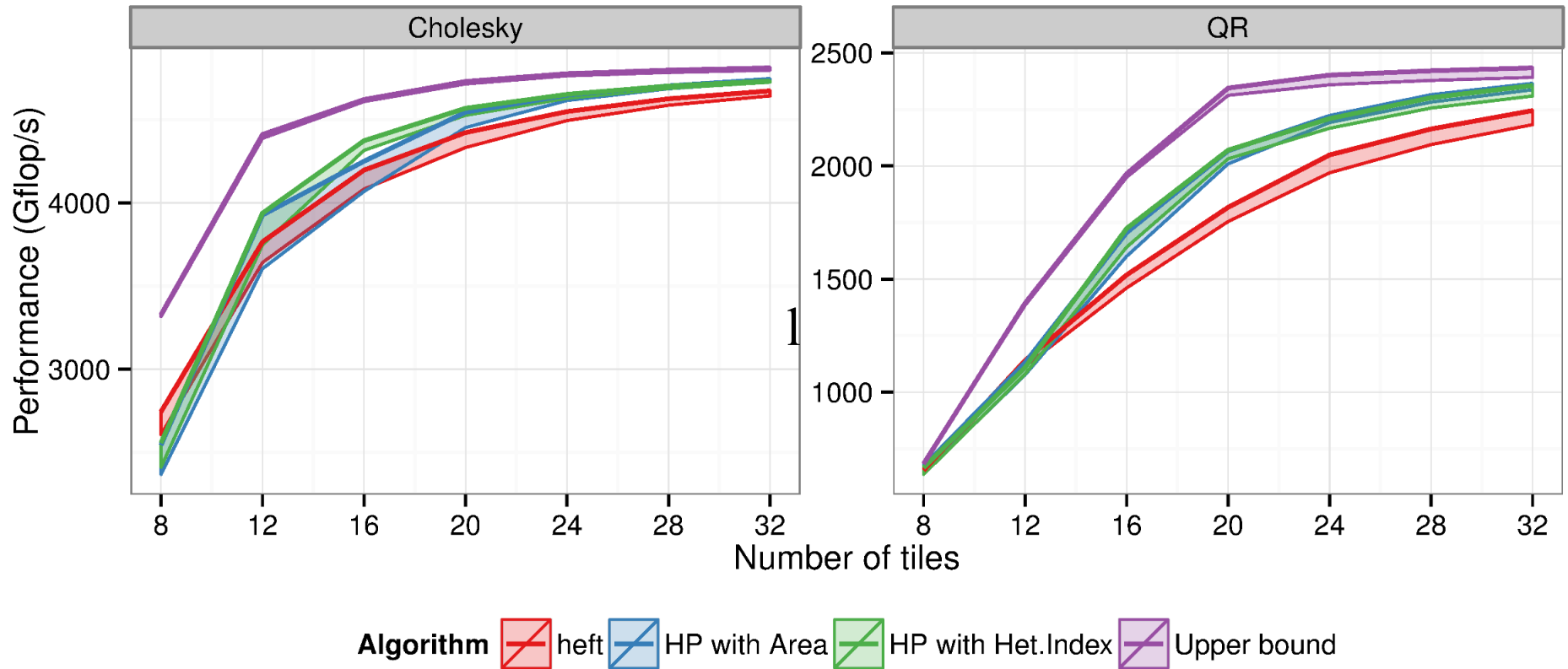
Scheduler initially designed for an FMM application on StarPU

- First just to respect hand-made priorities [Bramas14]
- Then extended to heterogeneous priorities [Bramas15]
- Then studied and generalized by theoreticians [Kumar16]
- And improved [Kumar&Cojean16]

Assumes a small number of types of tasks

- Sorted by acceleration factor
- List of tasks for each type, sorted by priority
 - I.e. contradict upward rank with acceleration factor

HeteroPrio



(From [Kumar&Cojean16])

Performance analysis tools

(see StarPU handbook for details)

Performance analysis tools

Performance models

- Offline
 - RAM/GPU bandwidth, RAM/Disk bandwidth
 - Task completion time linear / non-linear regression
- Online
 - Task completion time history-based average
 - React to performance changes
 - Eliminate outliers

Traces

- Offline analysis
 - Gantt Chart
 - Activity statistics

Bus performance

```
$ ./tools/starpu_machine_display
```

```
5 CPU cores
```

```
  CPU 0
```

```
  ...
```

```
3 CUDA Devices
```

```
  CUDA 0 (Tesla C2050 3.0 GiB 02:00.0)
```

```
  ...
```

from	to RAM	to CUDA 0	to CUDA 1	to CUDA 2
RAM	0.0	5236.89	5236.71	5240.12
CUDA 0	4547.68	0.0	3031.37	3093.99
CUDA 1	4547.62	3030.38	0.0	3093.90
CUDA 2	4537.36	3823.06	3823.17	0.0

Task distribution

```
$ STARPU_WORKER_STATS=1 ./examples/mult/sgemm
```

```
Time: 34.78 ms
```

```
GFlop/s: 24.12
```

```
Worker statistics:
```

```
*****
```

CUDA 0 (Quadro FX 5800)	264 task(s)
CUDA 1 (Quadro FX 5800)	237 task(s)
CUDA 2 (Quadro FX 5800)	237 task(s)
CPU 0	177 task(s)
CPU 1	175 task(s)
CPU 2	168 task(s)
CPU 3	177 task(s)

Bus usage

```
$ STARPU_BUS_STATS=1 ./examples/mult/sgemm
```

```
Time: 35.71 ms
```

```
GFlop/s: 23.49
```

```
Data transfer statistics:
```

```
*****
```

```
0 -> 1  2.52 MB 1.32MB/s          (transfers : 161 - avg 0.02 MB)
```

```
1 -> 0  2.39 MB 1.26MB/s          (transfers : 153 - avg 0.02 MB)
```

```
0 -> 2  3.12 MB 1.64MB/s          (transfers : 200 - avg 0.02 MB)
```

```
2 -> 0  3.00 MB 1.58MB/s          (transfers : 192 - avg 0.02 MB)
```

```
0 -> 3  3.03 MB 1.59MB/s          (transfers : 194 - avg 0.02 MB)
```

```
3 -> 0  2.91 MB 1.53MB/s          (transfers : 186 - avg 0.02 MB)
```

```
Total transfers: 16.97 MB
```

Disk usage

```
$ STARPU_BUS_STATS=1 ./tests/disk/disk_copy
```

```
0 -> 1: 337 MB/s
```

```
1 -> 0: 337 MB/s
```

```
0 -> 1: 1593  $\mu$ s
```

```
1 -> 0: 1593  $\mu$ s
```

```
NUMA 0 -> Disk 0 0.0625 GB 88.6847 MB/s (transfers: 2 - avg 32MB)
```

```
Total transfers: 0.0625 GB
```


Energy consumption

```
$ STARPU_WORKER_STATS=1 STARPU_PROFILING=1 ./examples/stencil/stencil
```

```
OpenCL 0 (Quadro FX 5800)
```

```
773 task(s)
```

```
total: 409.60 ms executing: 340.51 ms sleeping: 0.00
```

```
5040.000000 J consumed
```

```
OpenCL 1 (Quadro FX 5800)
```

```
767 task(s)
```

```
total: 409.62 ms executing: 346.28 ms sleeping: 0.00
```

```
10280.000000 J consumed
```

```
OpenCL 2 (Quadro FX 5800)
```

```
756 task(s)
```

```
total: 409.63 ms executing: 343.72 ms sleeping: 0.00
```

```
14880.000000 J consumed
```

Performance models

```
$ starpu_perfmodel_display -l
```

```
file: <starpu_sgemm_gemm>
```

```
$ starpu_perfmodel_display -s starpu_sgemm_gemm
```

```
performance model for cpu
```

# hash	size	mean	dev	n
880805ba49152		1.233333e+02	1.063576e+01	1612
8bd4e11d2359296		1.331984e+04	6.971079e+02	635

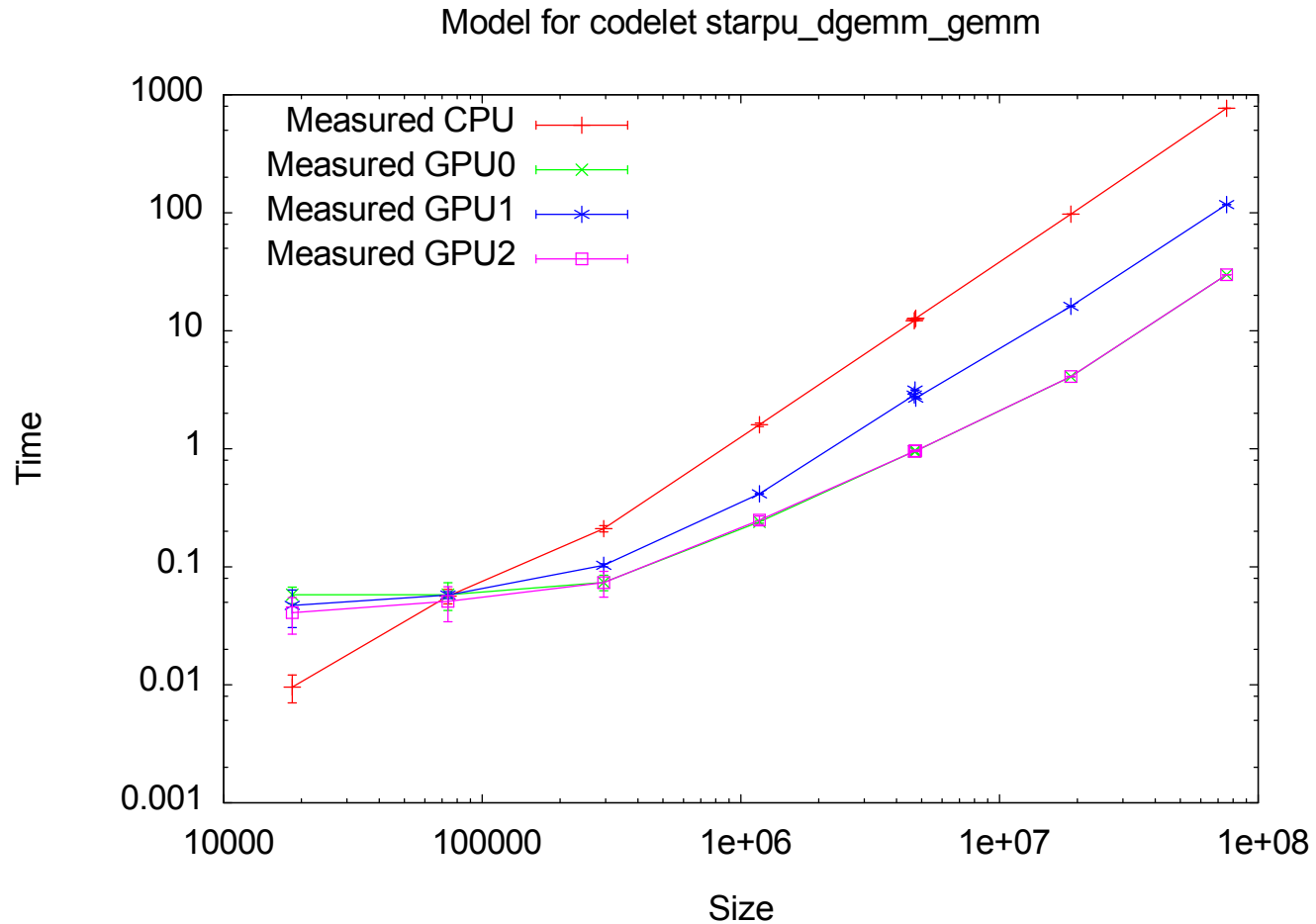
```
performance model for cuda_0
```

# hash	size	mean	dev	n
880805ba49152		2.743658e+01	2.178427e+00	496
8bd4e11d2359296		6.207991e+02	6.941988e+00	307

Performance models plot

```
$ starpu_perfmodel_plot -s starpu_dgemm_gemm
```

```
$ ./starpu_dgemm_gemm_gp
```

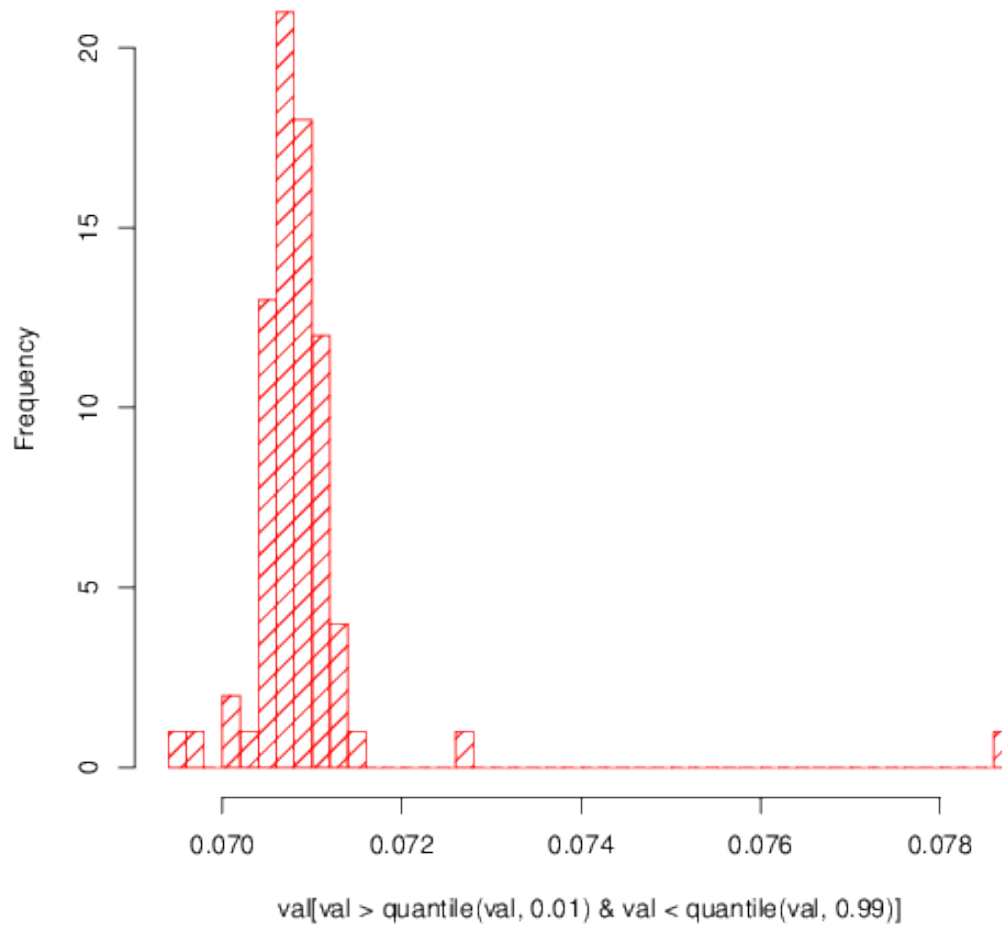


Kernel performance plot

```
$ starpu_fxt_tool -i /tmp/prof_file_user_sthibaul0
```

```
$ starpu_codelet_histo_profile distrib.data
```

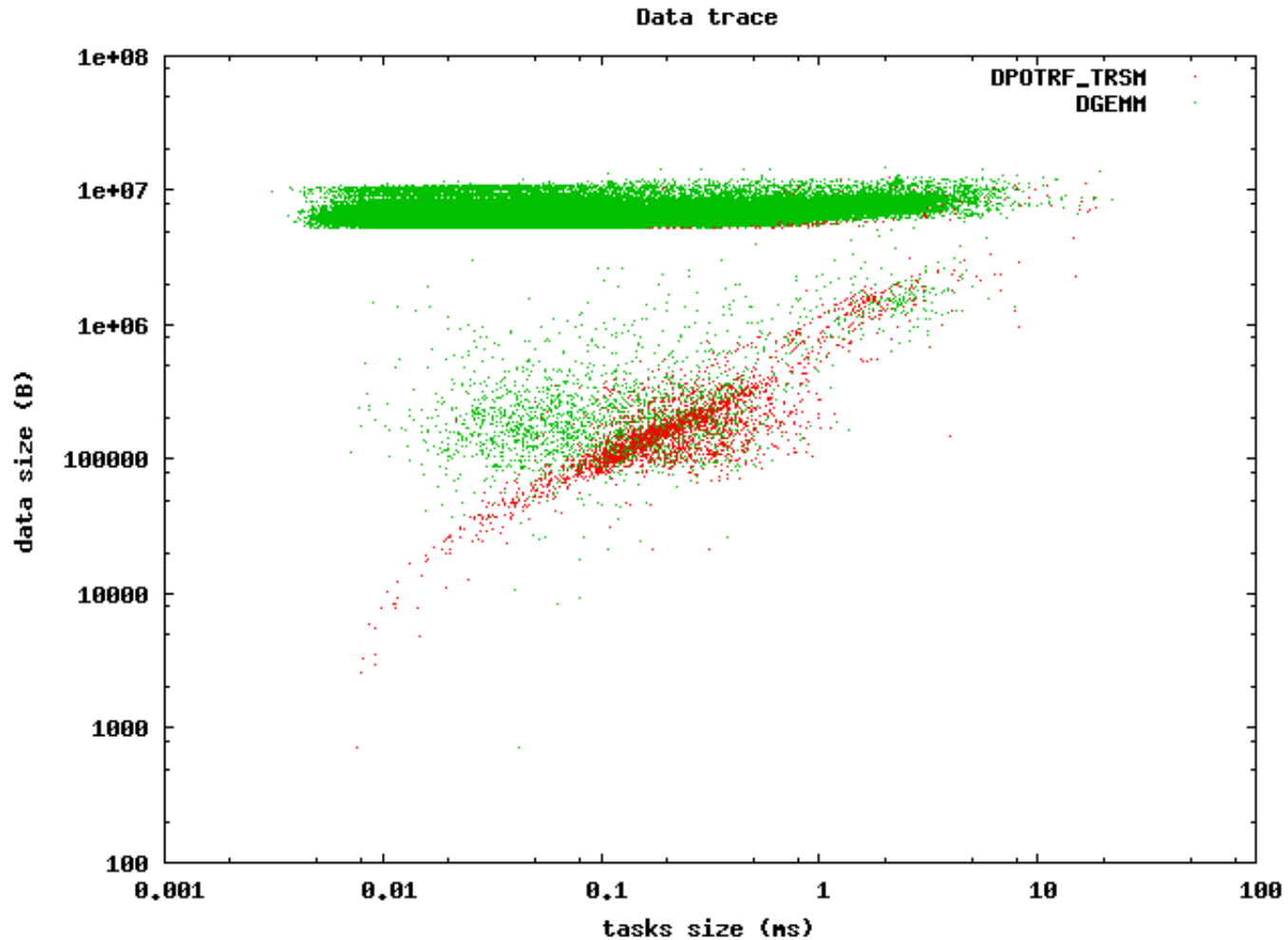
Histogram of `val[val > quantile(val, 0.01) & val < quantile(val, 0.99)]`



Kernel performance plot

```
$ starpu_fxt_data_trace /tmp/prof_file_sthibaul_0
```

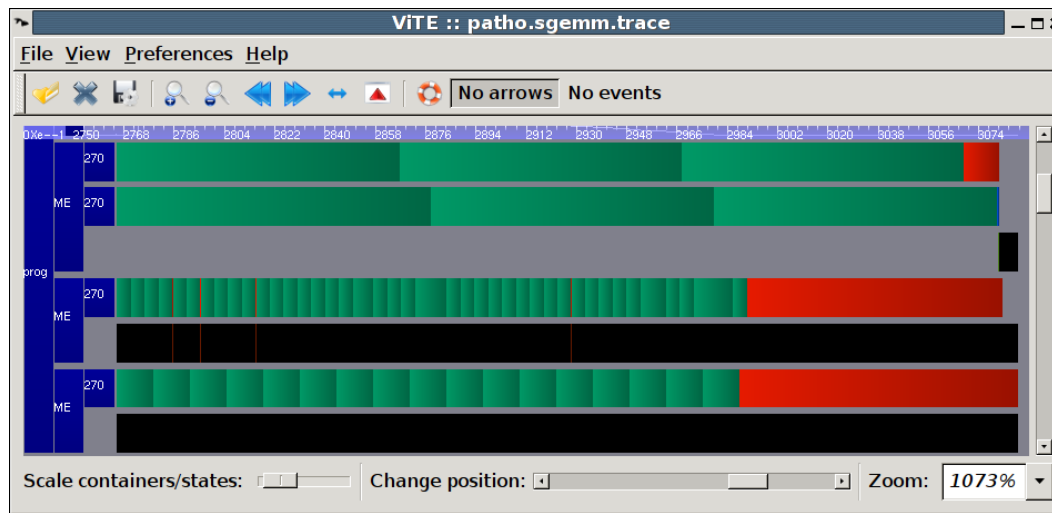
```
$ gnuplot data_trace.gp
```



Offline performance analysis

Visualize execution traces

- Generate a Pajé trace
 - <https://savannah.nongnu.org/projects/fkt>
 - `./configure --with-fxt`
 - `fxt_tool -i /tmp/prof_file_user_yourlogin`
→ `paje.trace`
- Vite trace visualization tool
 - Freely available from <http://vite.gforge.inria.fr/> (open source !)
 - `vite paje.trace`



2 Xeon cores

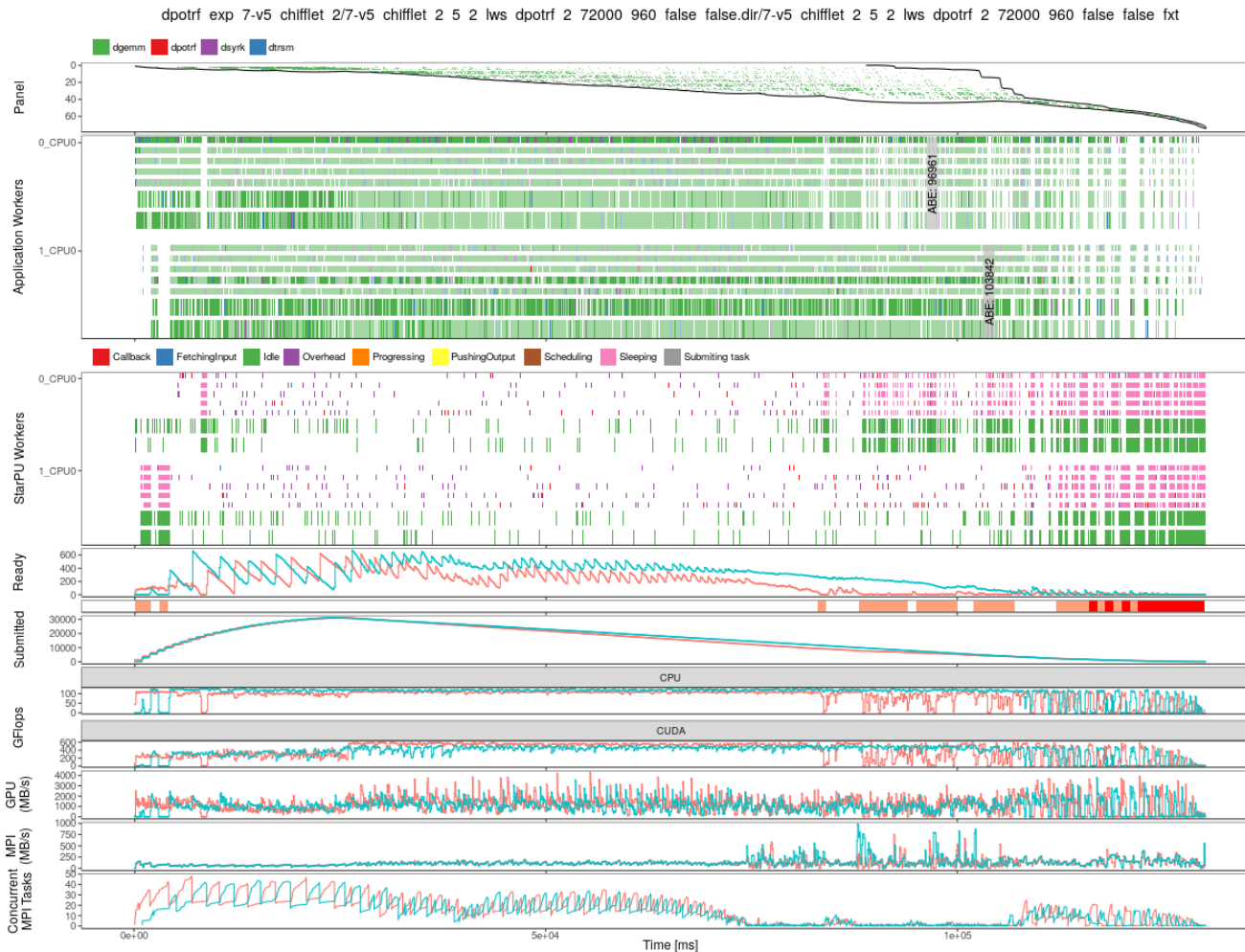
Quadro FX5800

Quadro FX4600

Offline performance analysis

Visualize execution traces

- starvz



Cluster support

How to scale over MPI?

(StarPU handles intra-MPI node scheduling fine)

- Splitting graph by hand
 - Complex, not flexible
 - Master-Slave does not scale
 - Each node should determine its duty by itself
 - Algebraic representation of e.g. Parsec
 - Difficult to write
 - Not flexible enough for any kind of application
 - Recursive task graph unrolling
 - Complex
- Rather just unroll the whole task graph on each node

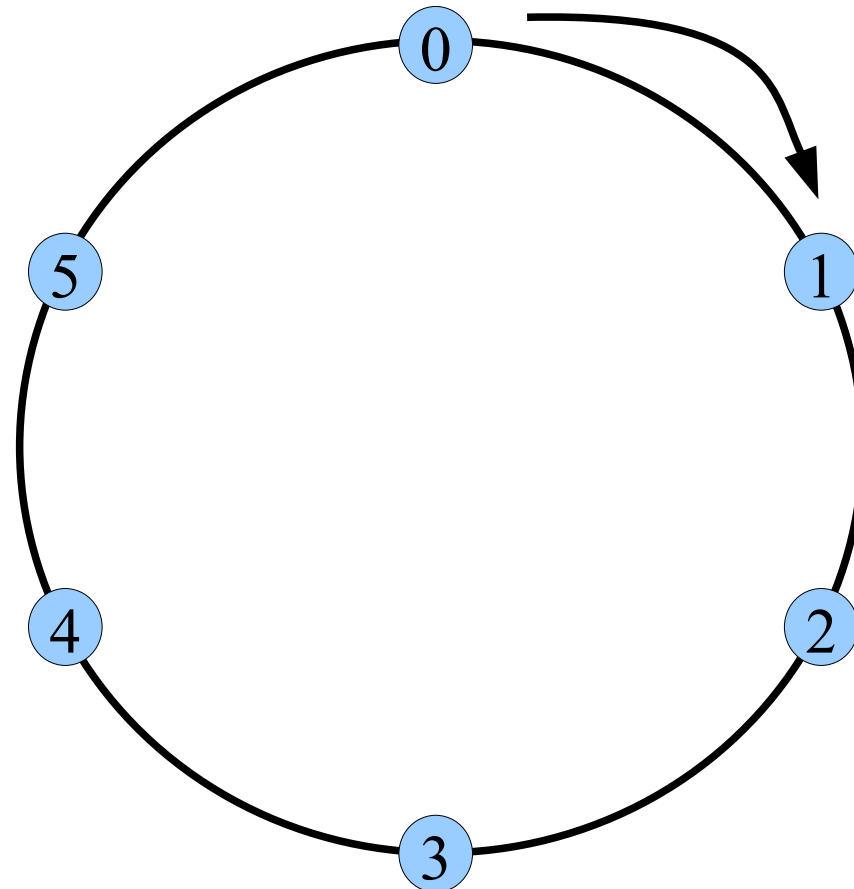
How about MPI + StarPU?

- Save programmers the burden of rewriting their MPI code
 - Keep the same MPI flow
 - Work on StarPU data instead of plain data buffers.
- 1 MPI process per machine, handles all CPUs and GPUs
- StarPU provides support for sending data over MPI
 - `starpu_mpi_send/recv, isend/irecv, ...`
 - Equivalents of `MPI_Send/Recv, Isend/Irecv, ...`
 - ... but working on StarPU data Handles
 - CPU/GPU transfers
 - task/communications dependencies
 - Overlapping everything
- [ICPADS'10]

But StarPU can also automatically generate all of this!

MPI ring example

- Token passed and incremented from node to node



MPI ring example

```
for (loop = 0 ; loop < NLOOPS; loop++) {  
    if ( !(loop == 0 && rank == 0))  
  
        MPI_Recv(&data, prev_rank, ...);  
  
    increment(&data);  
  
    if ( !(loop == NLOOPS-1 && rank == size-1))  
  
        MPI_Send(&data, next_rank, ...);  
  
}
```

StarPU-MPI ring example

```
for (loop = 0 ; loop < NLOOPS; loop++) {  
    if ( !(loop == 0 && rank == 0)) {  
        starpu_data_acquire(data_handle, STARPU_W) ;  
        MPI_Recv(&data, prev_rank, ...) ;  
        starpu_data_release(data_handle) ;  
    }  
    starpu_task_insert(&increment_codelet, STARPU_RW, data_handle, 0);  
  
    if ( !(loop == NLOOPS-1 && rank == size-1)) {  
        starpu_data_acquire(data_handle, STARPU_R) ;  
        MPI_Send(&data, next_rank, ...) ;  
        starpu_data_release(data_handle) ;  
    }  
}
```

StarPU-MPI ring example

```
for (loop = 0 ; loop < NLOOPS; loop++) {  
    if ( !(loop == 0 && rank == 0))  
  
        starpu_mpi_irecv_submit(data_handle, prev_rank, ...) ;  
  
    starpu_task_insert(&increment_codelet, STARPU_RW, data_handle, 0);  
  
    if ( !(loop == NLOOPS-1 && rank == size-1))  
  
        starpu_mpi_isend_submit(data_handle, next_rank, ...) ;  
  
}  
starpu_task_wait_for_all() ;
```

StarPU-MPI ring example

```
for (loop = 0 ; loop < N * NLOOPS; loop++) {
```

```
    starpu_mpi_task_insert(&increment_codelet, STARPU_RW, data_handle,  
                          STARPU_ON_NODE, loop % N, 0);
```

```
}
```

```
starpu_task_wait_for_all();
```

Automatic generation of Send/Recv MPI VSM

- Application decides data distribution over MPI nodes
- But data coherency extended to the MPI level
 - Automatic `starpu_mpi_send/recv` calls for each task
- Similar to a DSM, but granularity is whole data and whole task

- All nodes process the whole algorithm
 - Actual task execution according to data being written to

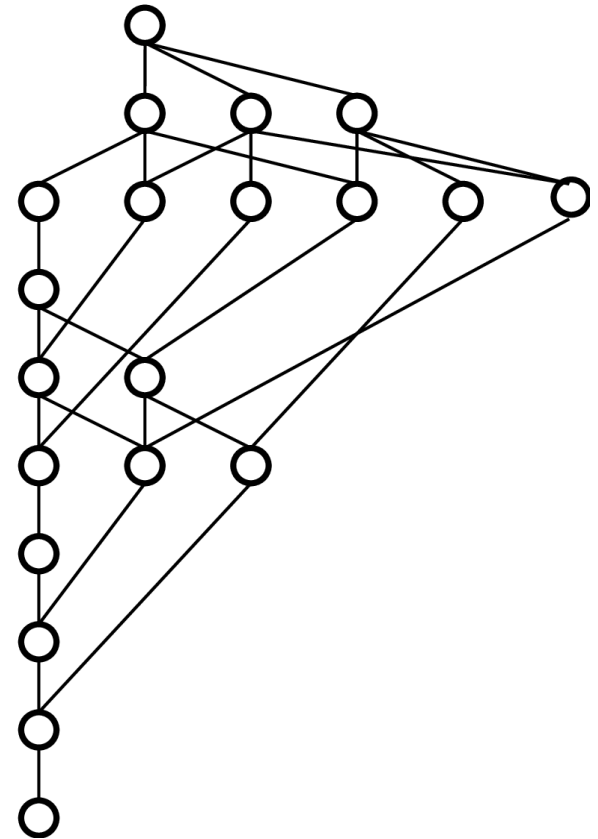
Sequential-looking code !

MPI VSM

```

For (k = 0 .. tiles - 1) {
  POTRF(A[k,k])
  for (m = k+1 .. tiles - 1)
    TRSM(A[k,k], A[m,k])
  for (m = k+1 .. tiles - 1) {
    SYRK(A[m,k], A[m,m])
    for (n = m+1 .. tiles - 1)
      GEMM(A[m,k], A[n,k], A[n,m])
  }
}

```



MPI VSM

- Data mapping (e.g. 2D block-cyclic)

```
int get_rank(int m, int n) { return ((m%p)*q + n%q); }
```

```
For (m = 0 .. tiles - 1)
```

```
  For (n = m .. tiles - 1)
```

```
    set_rank(A[m,n], get_rank(m,n));
```

```
For (k = 0 .. tiles - 1) {
```

```
  POTRF(A[k,k])
```

```
  for (m = k+1 .. tiles - 1)
```

```
    TRSM(A[k,k], A[m,k])
```

```
  for (m = k+1 .. tiles - 1) {
```

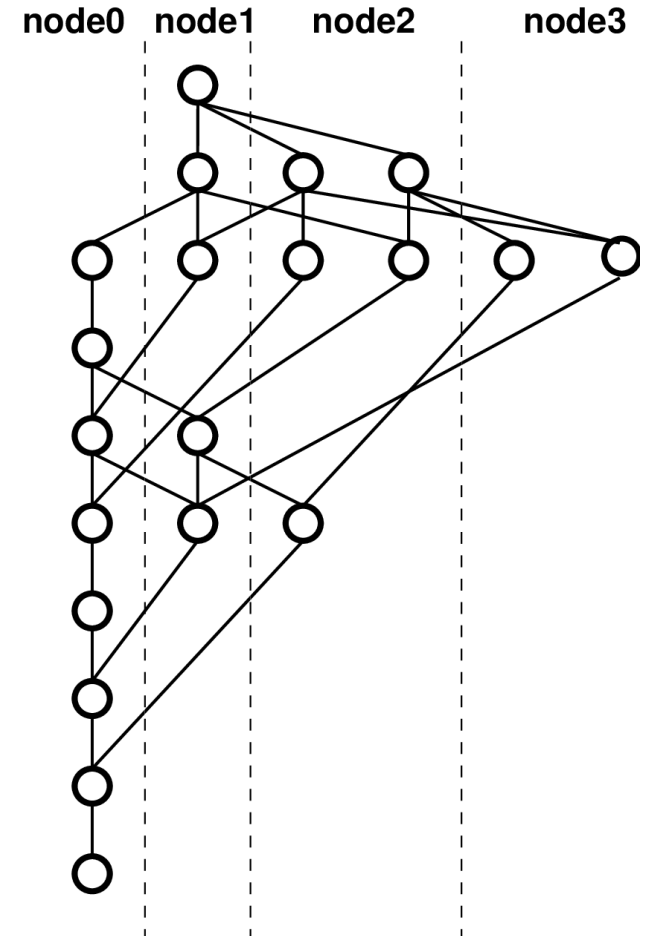
```
    SYRK(A[m,k], A[m,m])
```

```
    for (n = m+1 .. tiles - 1)
```

```
      GEMM(A[m,k], A[n,k], A[n,m])
```

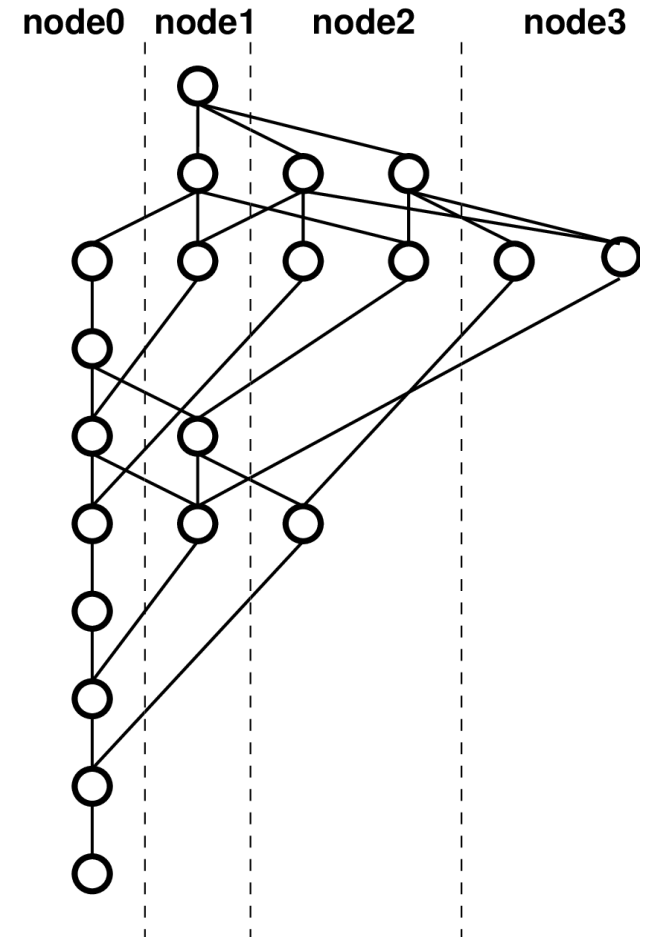
```
  }
```

```
}
```



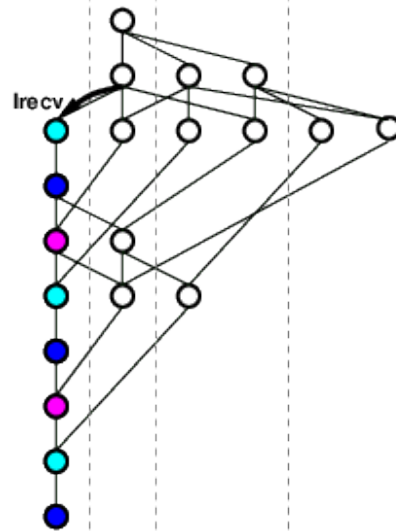
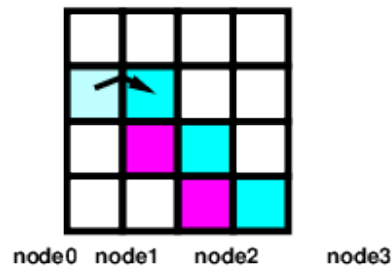
MPI VSM

- Each node unrolls the whole task graph
- Data \leftrightarrow node mapping
 - Provided by the application
 - E.g. 2D block-cyclic
 - Can be modified during submission
 - `starpu_mpi_data_migrate()`
- Task \leftrightarrow node mapping
 - Tasks move to data they modify
- MPI transfers
 - Automatically queued

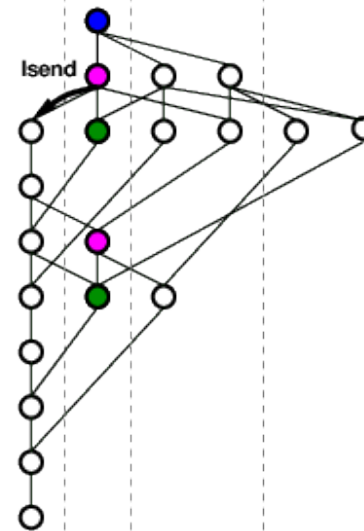
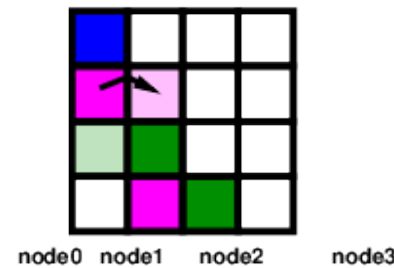


MPI VSM

- Right-Looking Cholesky decomposition (from PLASMA)



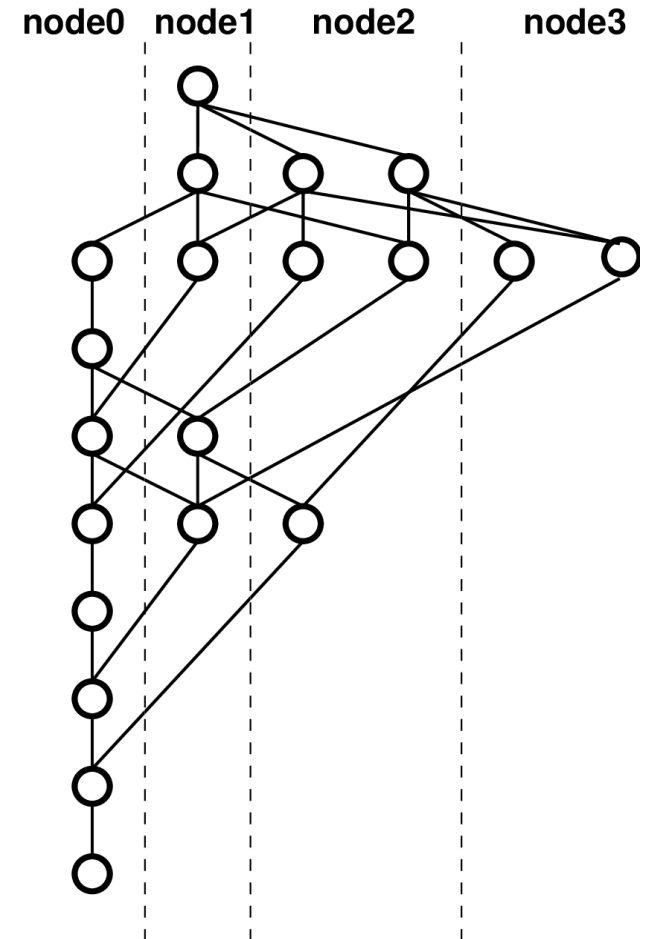
Node 0 execution



Node 1 execution

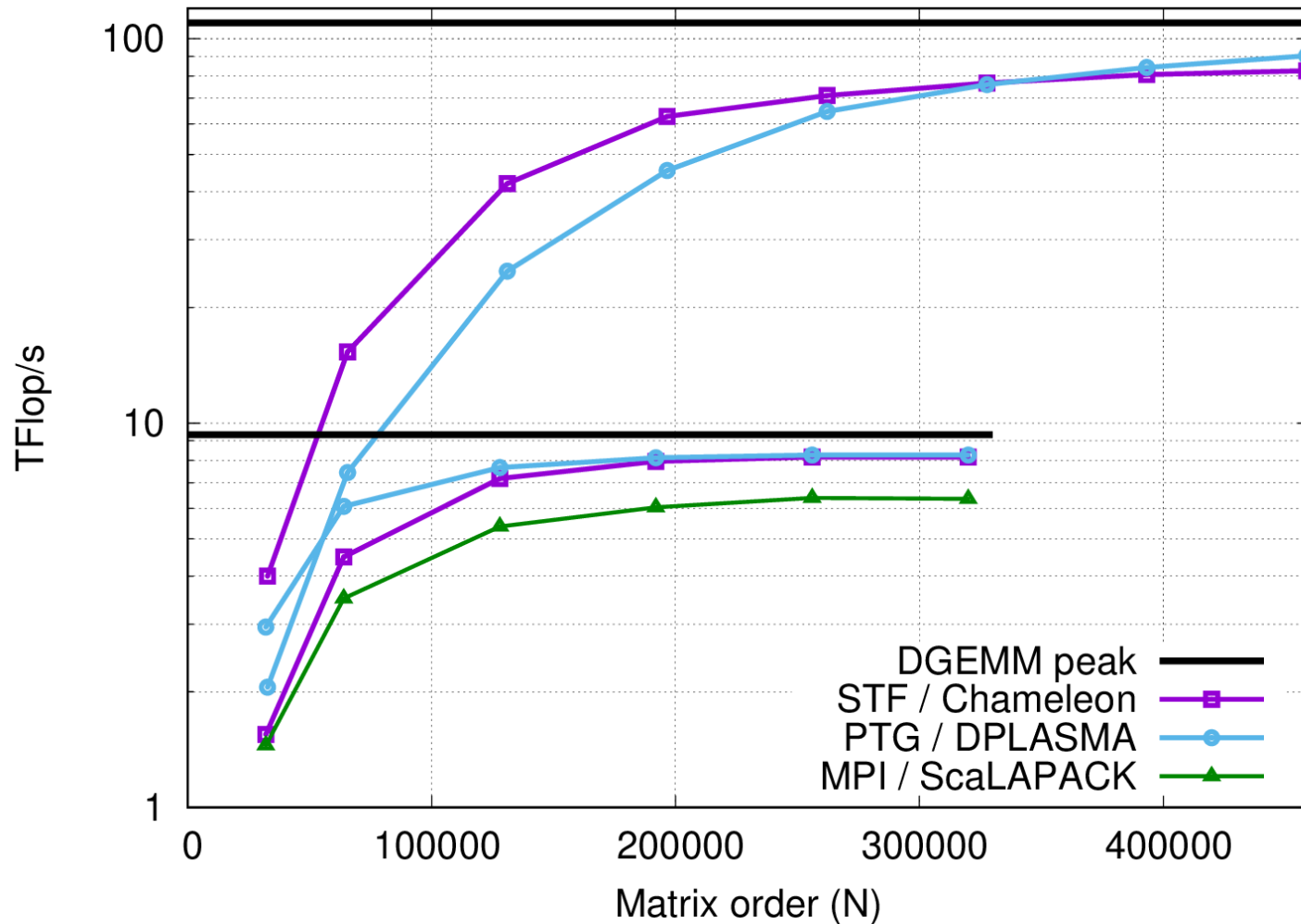
MPI VSM

- Separation of concerns: graph vs mapping
- Local view of the computation
 - No synchronizations
 - No global scheduling



Cholesky cluster performance

@CEA: 144 nodes with 8 CPU cores (E5620) + 2 GPUs (M2090)

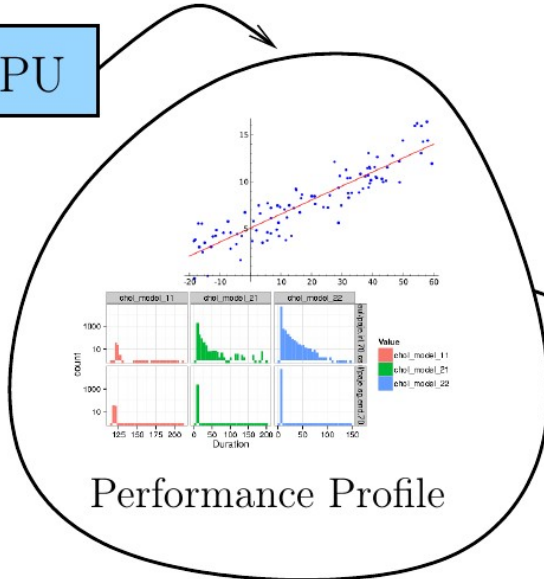


Simulation with SimGrid

Calibration



App StarPU



From A. Legrand
and L. Stanisc

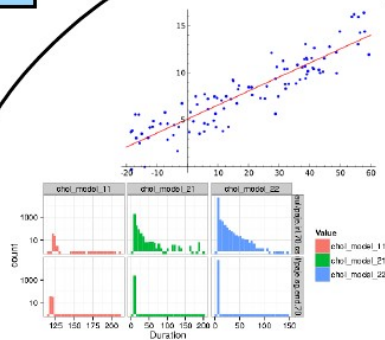
Run once!

Simulation with SimGrid

Calibration



App StarPU

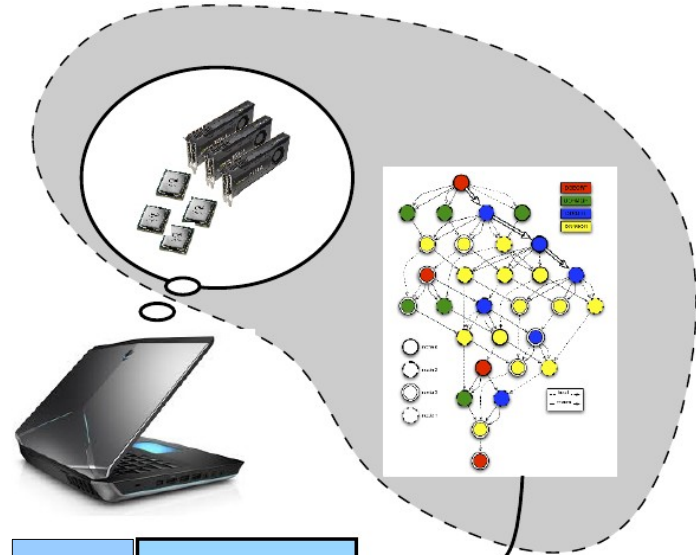


Performance Profile

From A. Legrand
and L. Stanisc

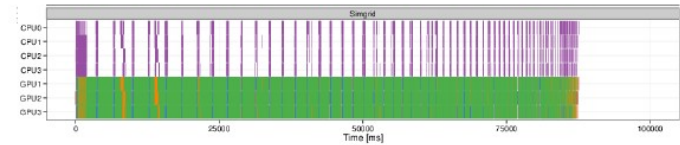
Run once!

Simulation



App StarPU

SimGrid

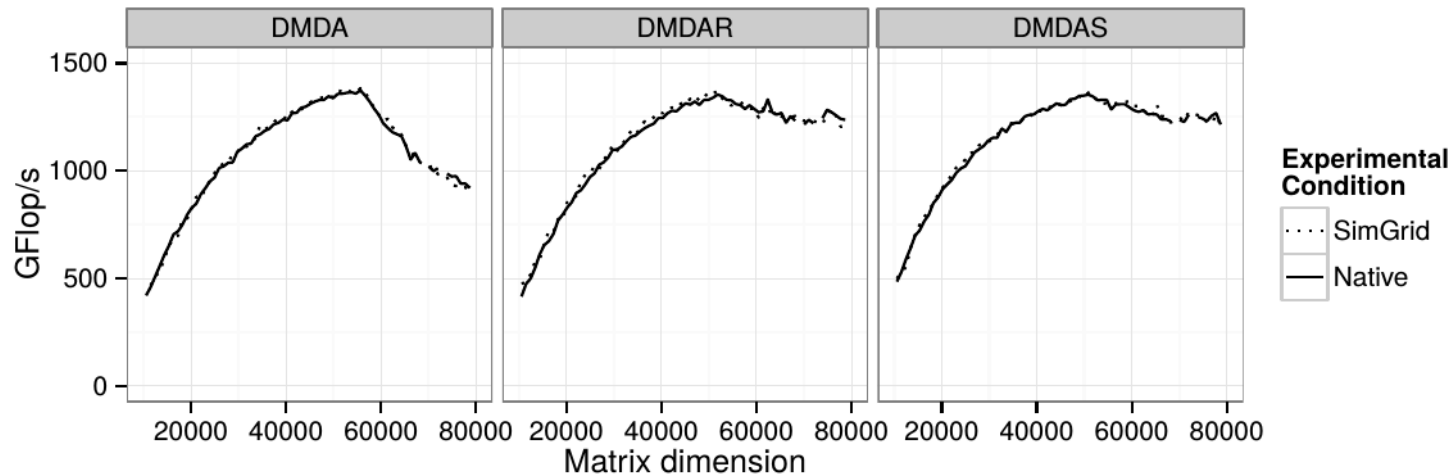


Quickly Simulate Many Times

Simulation with SimGrid

- Run application natively on target system
 - Records performance models
- Rebuild application against simgrid-compiled StarPU
- Run again
 - Uses performance model estimations instead of actually executing tasks
- Way faster execution time
- Reproducible experiments
- No need to run on target system
- Can change system architecture

Simulation with SimGrid



- Way faster execution time
- Reproducible experiments
- No need to run on target system
- Can change system architecture
- Can evaluate and study different scheduling algorithms

Applications on top of StarPU

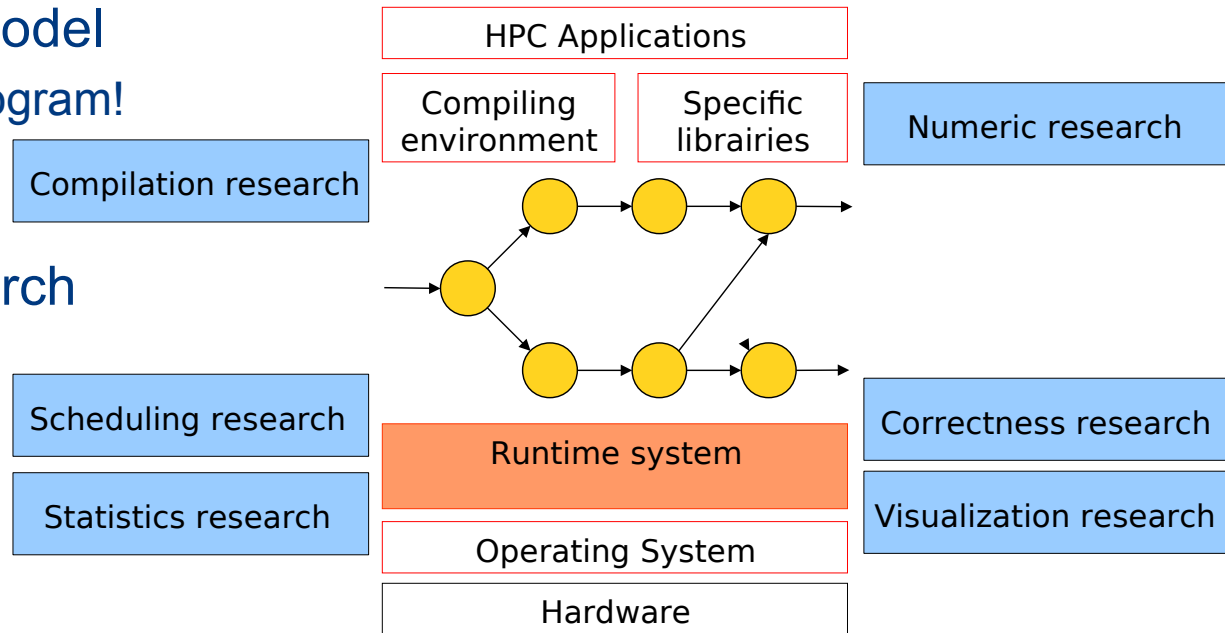
Using CPUs, GPUs, distributed, out of core, ...

- Dense linear algebra
 - Cholesky, QR, LU, ... : Chameleon (based on Plasma/Magma)
- Sparse linear algebra
 - QR_MUMPS
 - PaStiX
- Compressed linear algebra
 - BLR, h-matrices
- Fast Multipole Method
 - ScalFMM
- Conjugate Gradient
- Other programming models : Data flow, skeletons
 - SignalPU, SkePU
- ...

Conclusion

Task graphs

- Nice programming model
 - Keep sequential program!
- Optimized execution
- Playground for research
 - Runtime
 - Scheduling
 - Numeric algorithms
 - Statistics
 - Correctness
- Used for various real-world computations
 - Cholesky/QR/LU (dense/sparse/compressed), stencil, CG, CFD, FMM...



<http://starpu.gitlabpages.inria.fr/tutorials/>