



Programming heterogeneous, accelerator-based machines with StarPU

Raymond Namyst
University of Bordeaux
RUNTIME INRIA group

RUNTIME
INRIA Group
INRIA Bordeaux Sud-Ouest

PPAM 2011
Torun, Poland, Sept. 11th-14th

Understanding the evolution of parallel machines

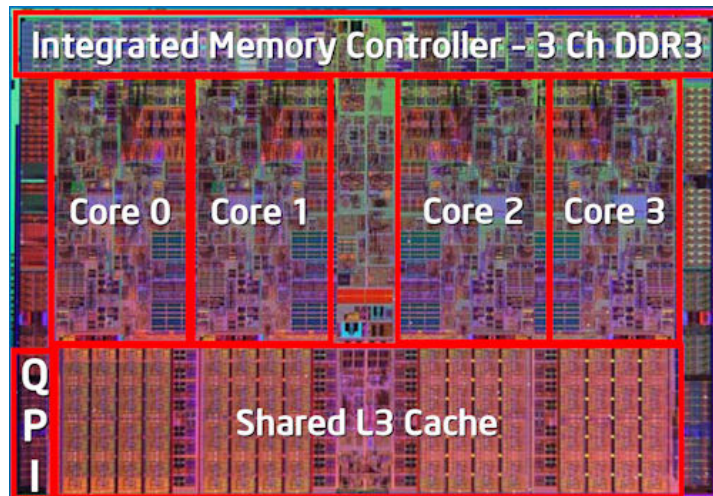
The end of Moore's law?

- The end of single thread performance increase
 - Clock rate is no longer increasing
 - Thermal dissipation
- Processor architecture is already very sophisticated
 - Prediction and Prefetching techniques achieve a very high percentage of success
- Actually, processor complexity is decreasing!
- Question: What circuits should we better add on a die?

Understanding the evolution of parallel machines

Welcome to the multicore era

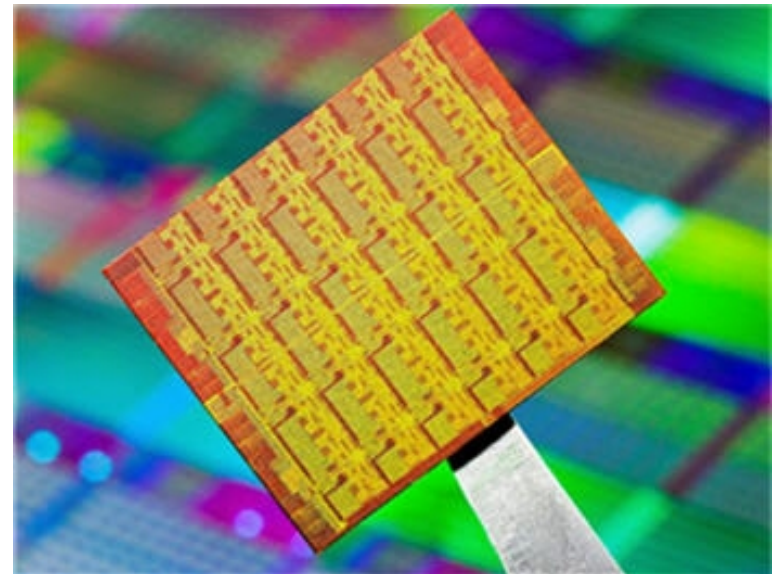
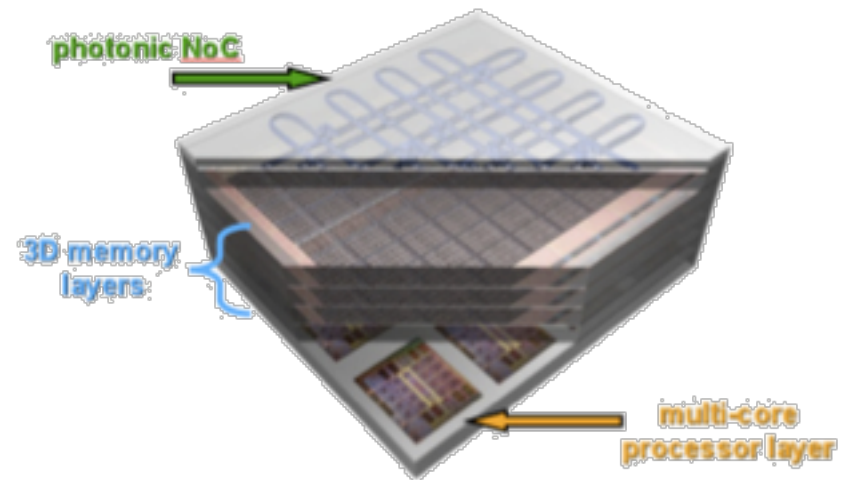
- Answer: Multicore chips
 - Several cores instead of one processor
- Back to complex memory hierarchies
 - Shared caches
 - Organization is vendor-dependent
 - NUMA penalties
- Clusters can no longer be considered as “flat sets of processors”



Understanding the evolution of parallel machines

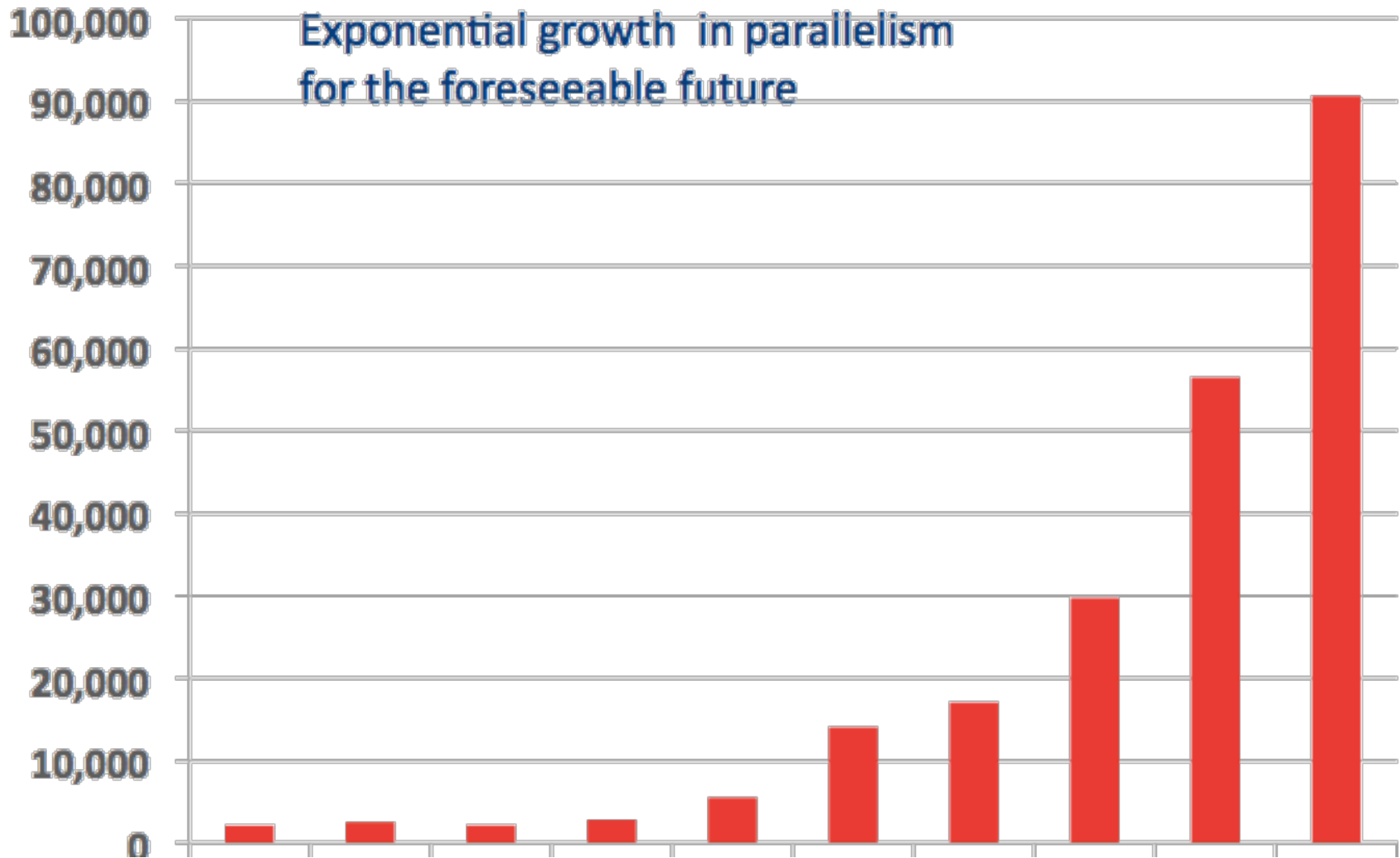
Multicore is a solid trend

- More performance = more cores
 - Toward embarrassingly parallel machines?
- Designing scalable multicore architectures
 - 3D stacked memory
 - Non-coherent cache architectures
 - Intel SCC
 - IBM Cell/BE



Understanding the evolution of parallel machines

Average number of cores per top20 supercomputer



Heterogeneous computing is here

And portable programming is getting harder...

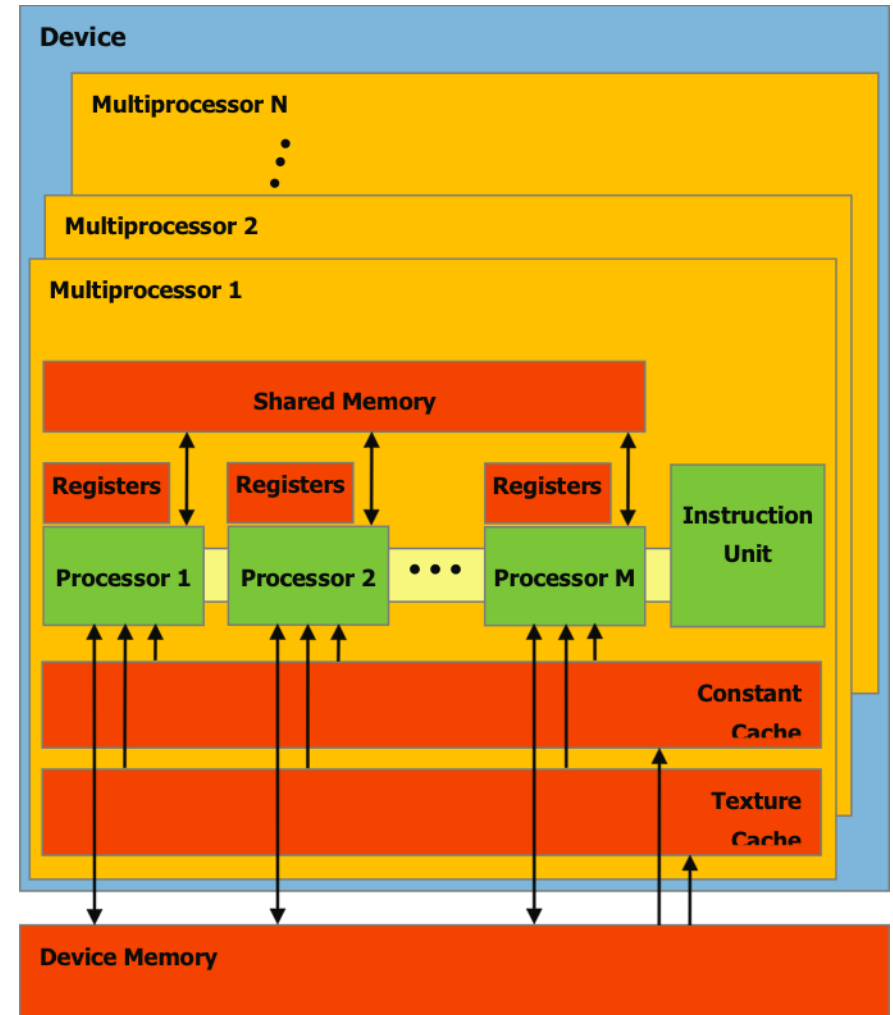
- GPUs are the new kids on the block
 - De facto adoption
 - Concrete success stories
 - Speedups > 50
- Clusters featuring accelerators are already heading the Top500 list
 - Tianhe-1A (#2)
 - Nebulae (#4)
 - Tsubame 2.0 (#5)
 - Roadrunner (#10)



Heterogeneous computing is here

And portable programming is getting harder...

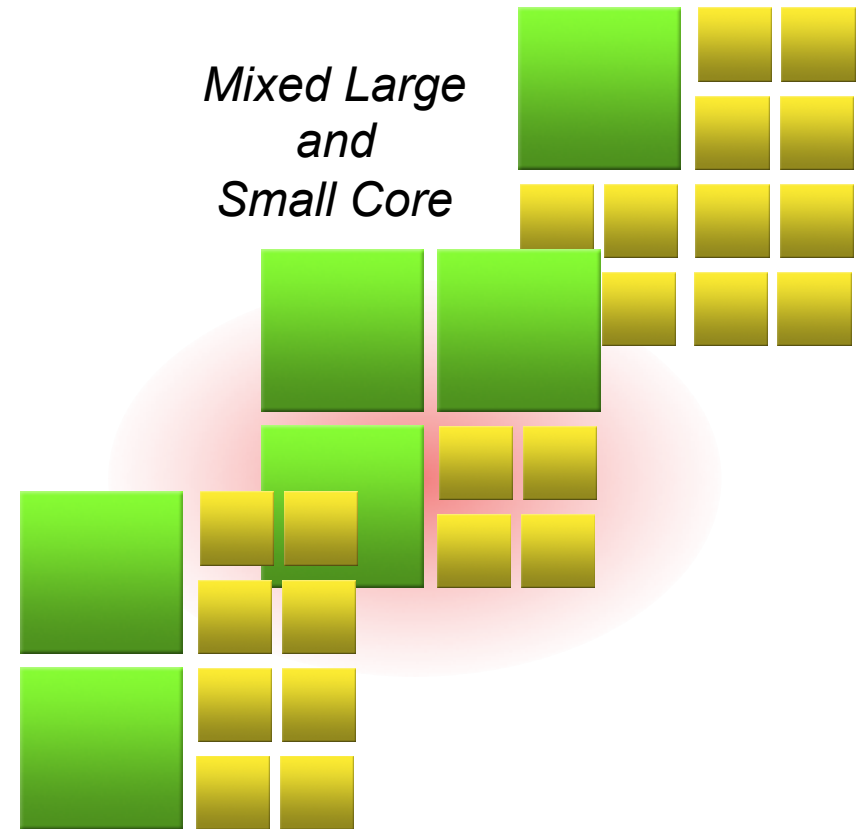
- Programming model
 - Specialized instruction set
 - SIMD execution model
 - Nvidia Fermi GTX 480
512 cores
- Memory
 - Size limitations
 - No hardware consistency
 - Explicit data transfers
- Using GPUs as “side accelerators” is not enough
 - GPU = first class citizens



Heterogeneous computing is here

And it seems to be a solid trend...

- “Future processors will be a mix of general purpose and specialized cores” (anonymous source)
 - One interpretation of “Amdahl’s law”
 - Need powerful, general purpose cores to speed up sequential code
- Accelerators will be more integrated
 - Intel Knights Corner (MIC), SandyBridge
 - AMD Fusion
 - Nvidia Tegra-like
- Are we happy with that?
 - No, but it’s probably unavoidable!



The Quest for programming models



What Programming Models for such machines?

Widely used, standard programming models

- MPI
 - Communication Interface
 - Scalable implementations exist already
 - Was actually designed with scalability in mind
 - Makes programmers “think” scalable algorithms
 - NUMA awareness?
 - Memory consumption
- OpenMP
 - Directive-based, incremental parallelization
 - Shared-memory model
 - Well suited to symmetric machines
 - Portability wrt #cores
 - NUMA awareness?

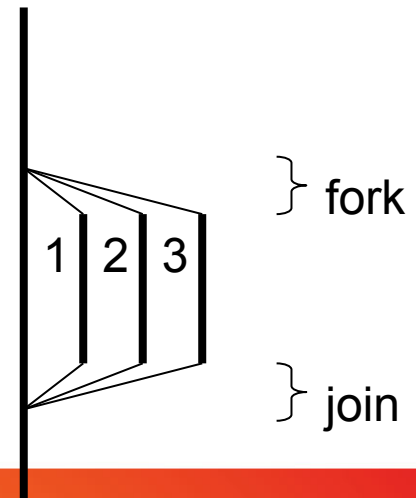
OpenMP (1997)

A portable approach to shared-memory programming

- Extensions to existing languages
 - C, C++, Fortran
 - Set of programming directives
- Fork/join approach
 - Parallel sections
- Well suited to data-parallel programs
 - Parallel loops
- OpenMP 3.0 introduced *tasks*
 - Support for irregular parallelism

```
int matrix[MAX][MAX];
...
#pragma omp parallel for
for (int i; i < 400; i++)
{
    matrix[i][0] += ...
}
```

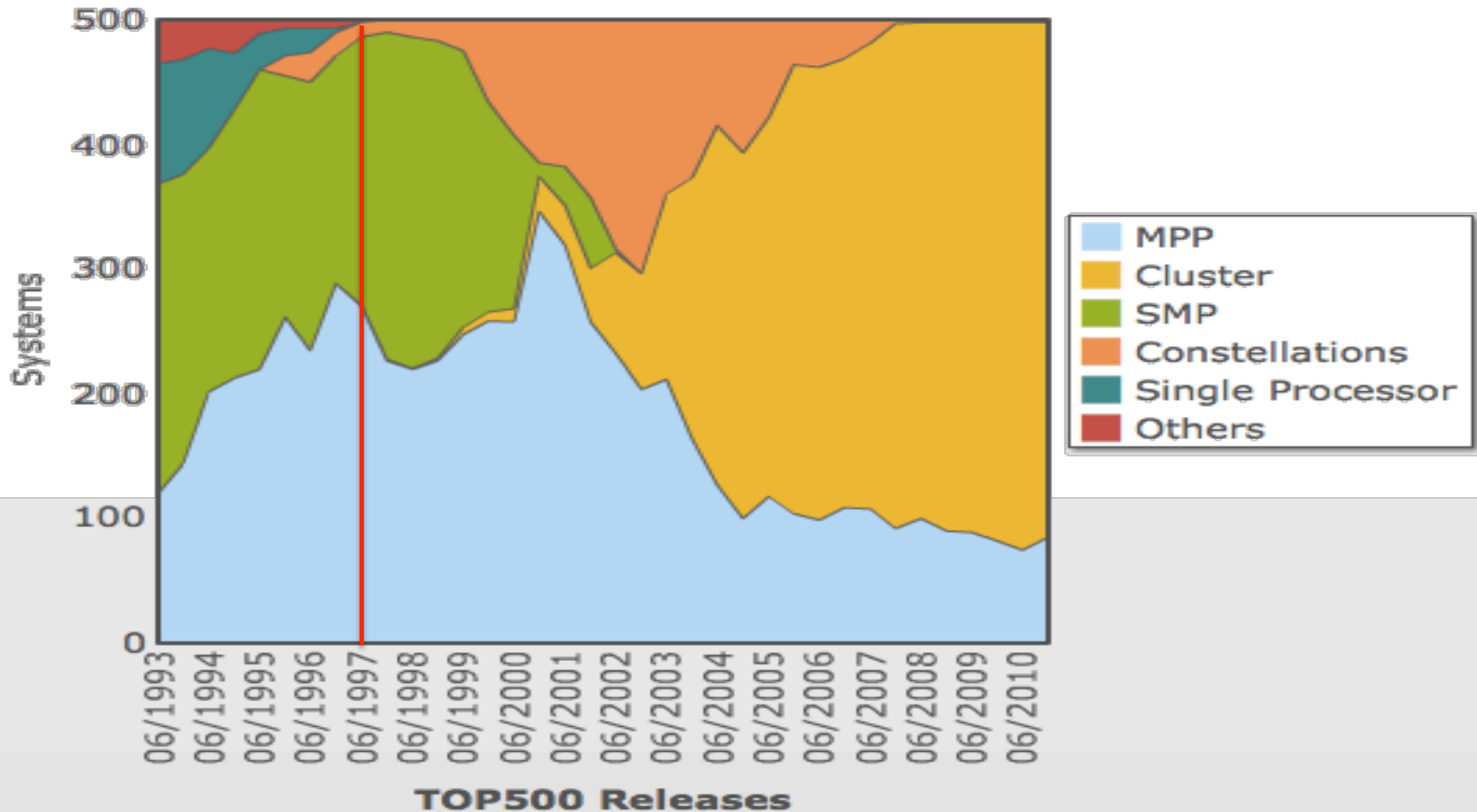
0 (main)



OpenMP (1997)

Multithreading over shared-memory machines

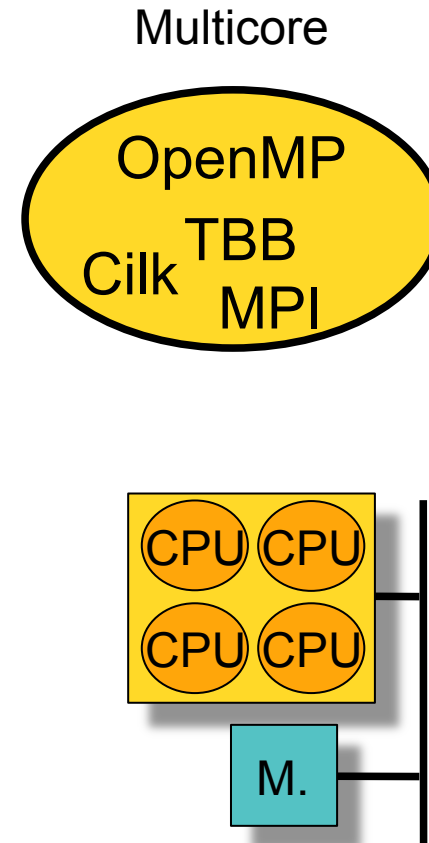
**Architecture Share Over Time
1993-2010**



The Quest for Programming Models

Dealing with multicore machines

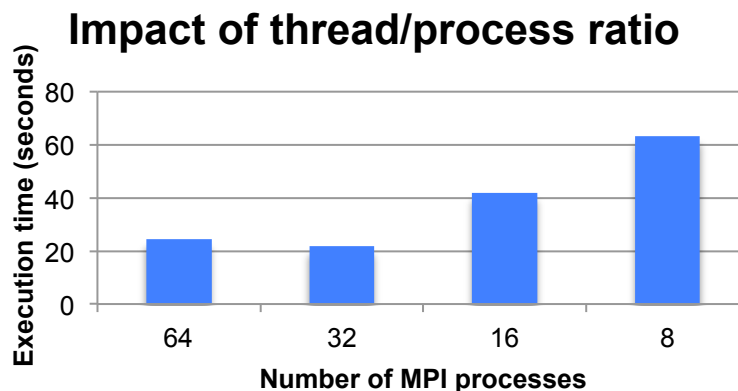
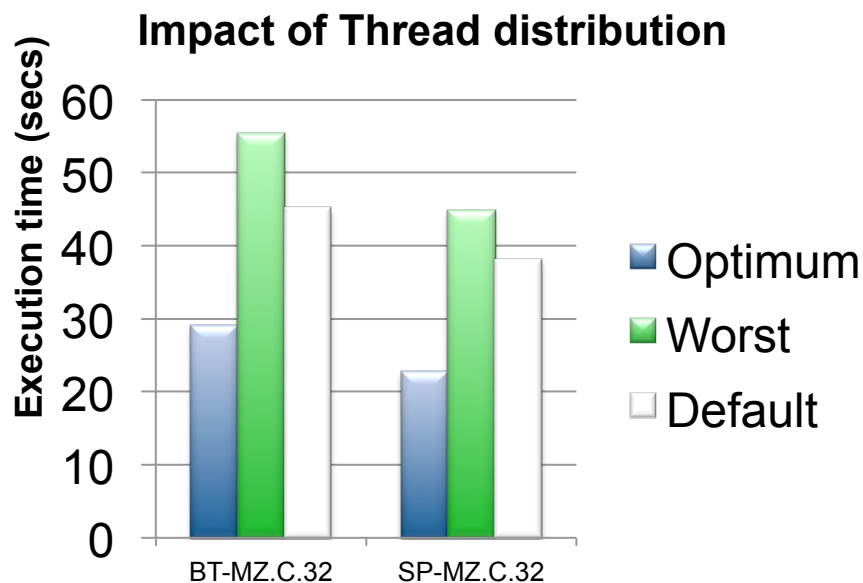
- Several efforts aim at making MPI and OpenMP multicore-ready
 - OpenMP
 - Scheduling in a NUMA context (memory affinity, work stealing)
 - Memory management (page migration)
 - MPI
 - NUMA-aware buffer management
 - Efficient collective operations



Mixing OpenMP with MPI

It makes sense even on shared-memory machines

- MPI address spaces must fit the underlying topology
- Experimental platforms exit to hybrid applications
 - Topology-aware process allocation
 - Customizable core/process ratio
 - # of OpenMP tasks independent from # of cores

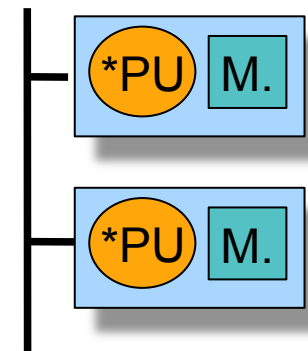
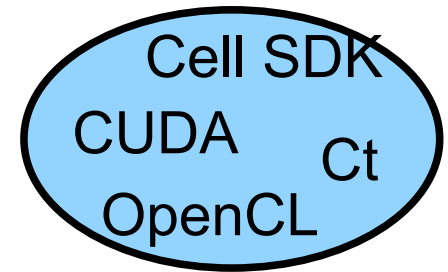


The Quest for Programming Models

Dealing with accelerators

- Software Development Kits and Hardware Specific Languages
 - “Stay close to the hardware and get good performance”
 - Low-level abstractions
 - Compilers generate code for accelerator device
- Examples
 - Nvidia’s CUDA
 - *Compute Unified Device Architecture*)
 - IBM Cell SDK
 - OpenCL

Accelerators



The Quest for Programming Models

The hidden beauty of CUDA

```
__global__ void mykernel(float * A1, float * A2, float * R)
{
    int p = threadIdx.x;
    R[p] = A1[p] + A2[p];
}

int main()
{
    float A1[]={1,2,3,4,5,6,7,8,9}, A2[]={10,20,30,40,50,60,70,80,90}, R[9];
    int size=sizeof(float) * 9;
    float *a1_device, *a2_device, *r_device;
    cudaMalloc ( (void**) &a1_device, size); cudaMalloc ( (void**) &a2_device, size); cudaMalloc ( (void**) &r_device, size);
    cudaMemcpy( a1_device,A1,size,cudaMemcpyHostToDevice); cudaMemcpy( a2_device,A2,size,cudaMemcpyHostToDevice);

    mykernel<<<1,9>>>(a1_device, a2_device, r_device);

    cudaMemcpy(R,r_device,taille_mem,cudaMemcpyDeviceToHost) ;
}
```


The Quest for Programming Models

Are we forced to use such low-level tools?

- Fortunately, well-known kernels are available
 - BLAS routines
 - e.g. CUBLAS
 - FFT kernels
- Implementations are continuously enhanced
 - High Efficiency
- Limitations
 - Data must usually fit accelerators memory
 - Multi-GPU configurations not well supported
- Ongoing efforts
 - Using multi-GPU + multicore
 - MAGMA (Oak Ridge National Lab)

Directive-based approaches

Offloading tasks to accelerators

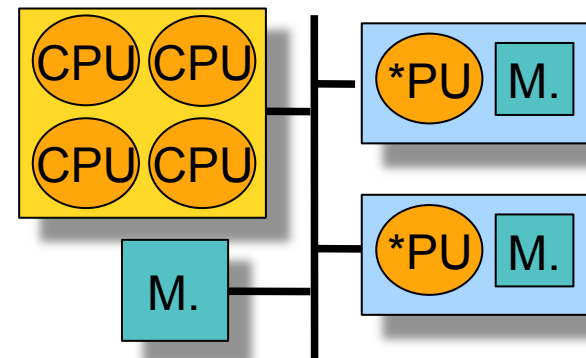
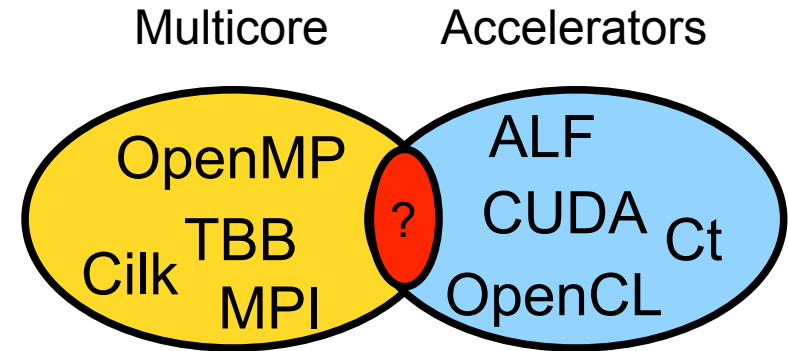
- Idea: use simple directives... and better compilers
 - HMPP (Caps Enterprise)
 - GPU SuperScalar (Barcelona Supercomputing Center)

```
#pragma omp task inout(C[BS][BS])
void matmul( float *A, float *B, float *C) {
// regular implementation
}
#pragma omp target device(cuda) implements(matmul)
copy_in(A[BS][BS] , B[BS][BS] , C[BS][BS])
copy_out(C[BS][BS])
void matmul cuda ( float *A, float *B, float *C) {
// optimized kernel for cuda
}
```

The Quest for Programming Models

How shall we program heterogeneous clusters?

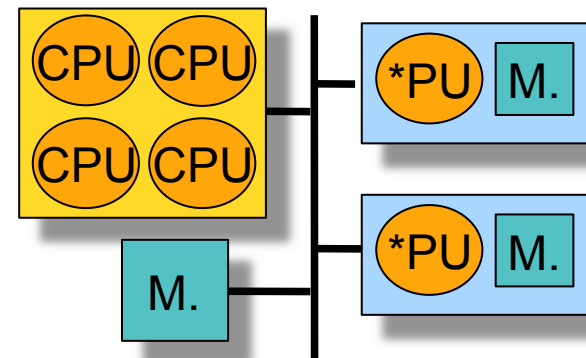
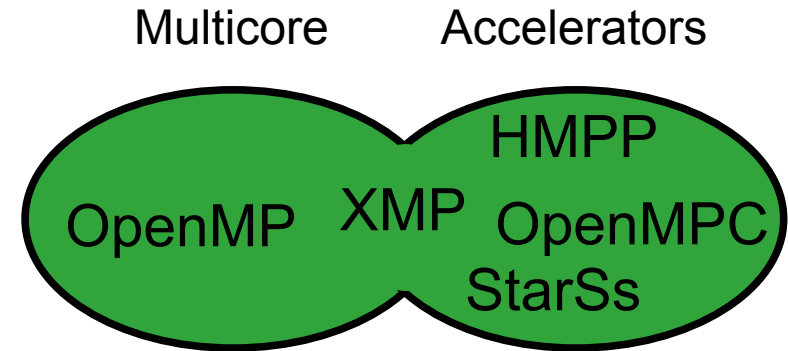
- The ~~hard~~ hybrid way
 - Combine different paradigms by hand
 - MPI + {OpenMP/TBB/???) + {CUDA/OpenCL}
 - Portability is hard to achieve
 - Work distribution depends on #GPU & #CPU per node...
 - Needs aggressive autotuning
 - Currently used for building parallel numerical kernels
 - MAGMA, D-PLASMA, FFT kernels



The Quest for Programming Models

How shall we program heterogeneous clusters?

- The uniform way
 - Use a single (or a combination of) high—level programming language to deal with network + multicore + accelerators
 - Increasing number of directive-based languages
 - Use simple directives... and good compilers!
 - XcalableMP
 - PGAS approach
 - HMPP, OpenMPC, OpenMP 4.0
 - Generate CUDA from OpenMP code
 - StarSs
 - Much better potential for *composability*...
 - If compiler is clever!



**All the things
runtime systems can do for you**

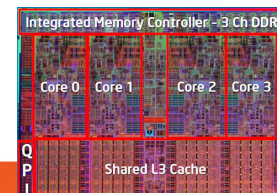
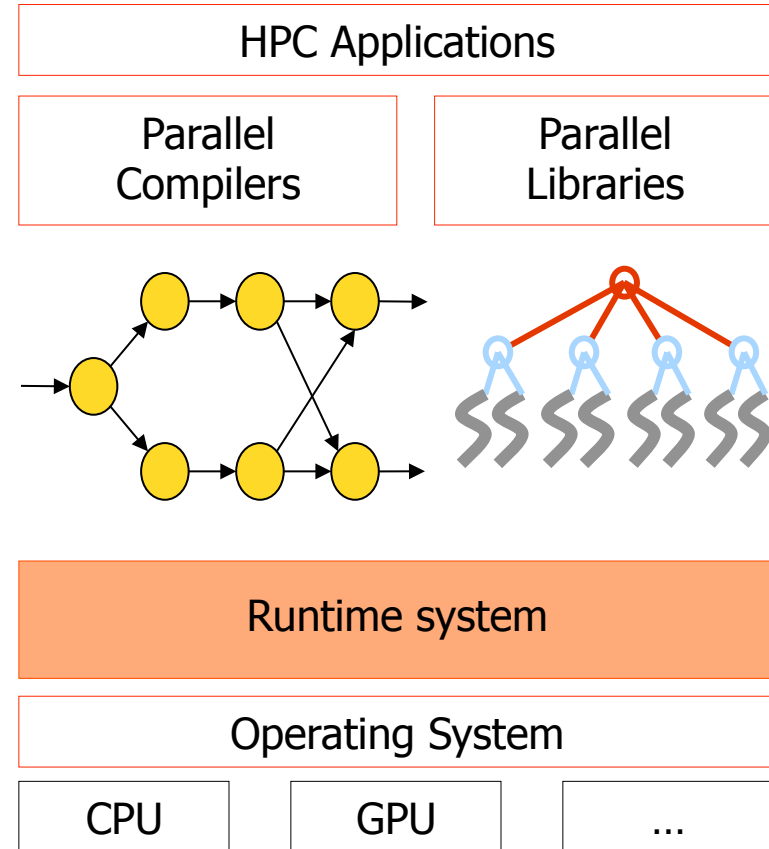
The logo for Inria, featuring the word "inria" in a stylized, cursive font with a color gradient from red to orange. Above the "ria" part, the words "informatics" and "mathematics" are written in a smaller, sans-serif font, separated by a red slash.

informatics mathematics
inria

The role of runtime systems

Toward “portability of performance”

- Do dynamically what can't be done statically
 - Load balance
 - React to hardware feedback
 - Autotuning, self-organization
- We need to put more intelligence into the runtime!



We need new runtime systems!

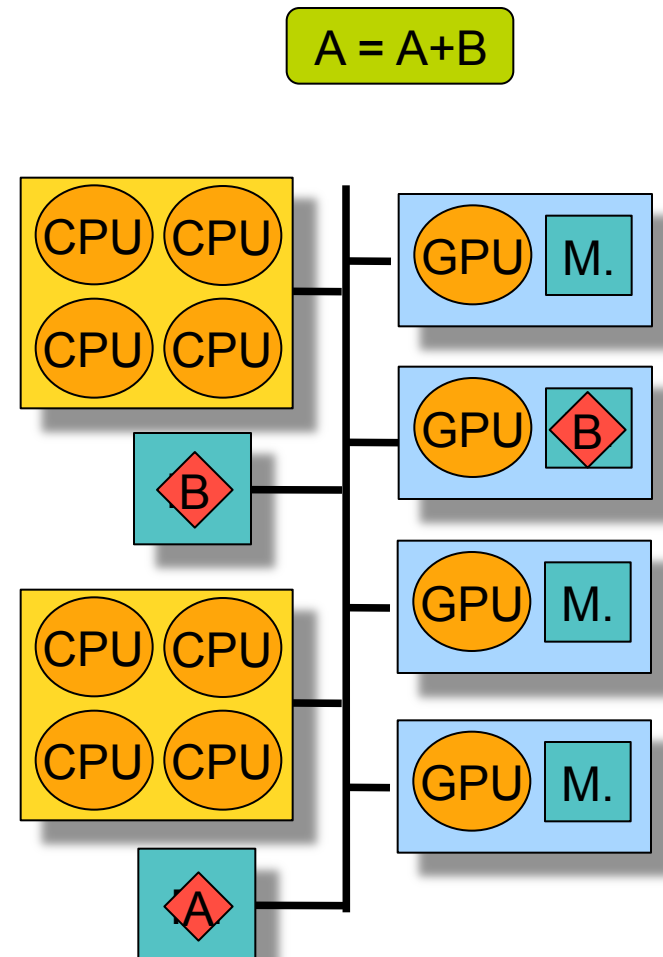
Toward “portability of performance”

- Computations need to exploit accelerators and regular CPUs simultaneously
- Data movements between memory banks
 - Should be minimized
 - Should not be triggered explicitly by application
- Computations need to accommodate to a variable number of processing units
 - Some computations do not scale over a large #cores

Overview of StarPU

A runtime system for heterogeneous architectures

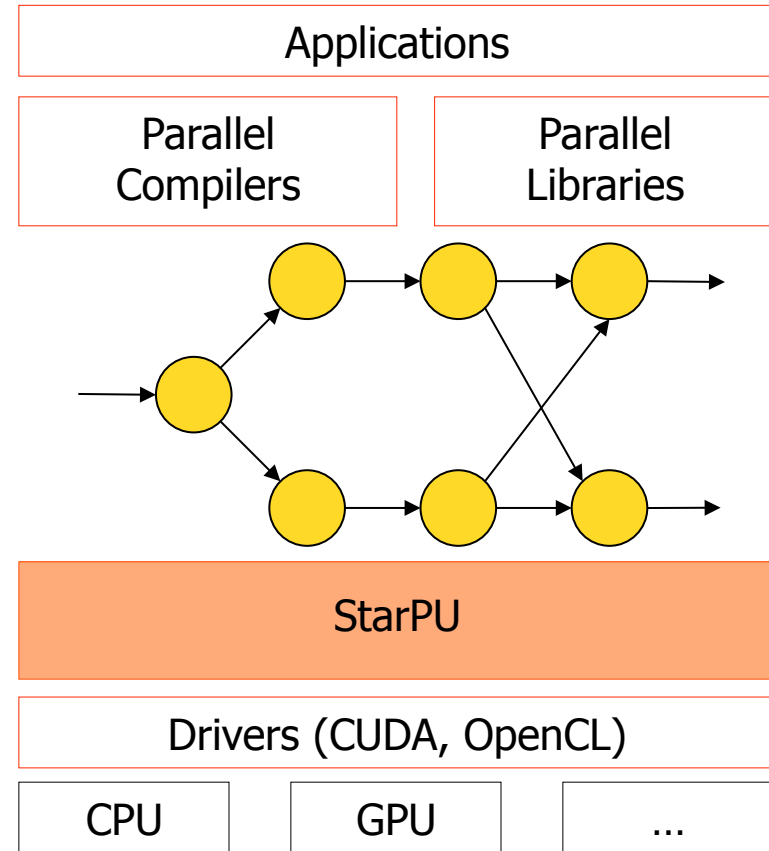
- Rational
 - Dynamically schedule tasks on all processing units
 - See a pool of heterogeneous processing units
 - Avoid unnecessary data transfers between accelerators
 - Software VSM for heterogeneous machines



Overview of StarPU

Maximizing PU occupancy, minimizing data transfers

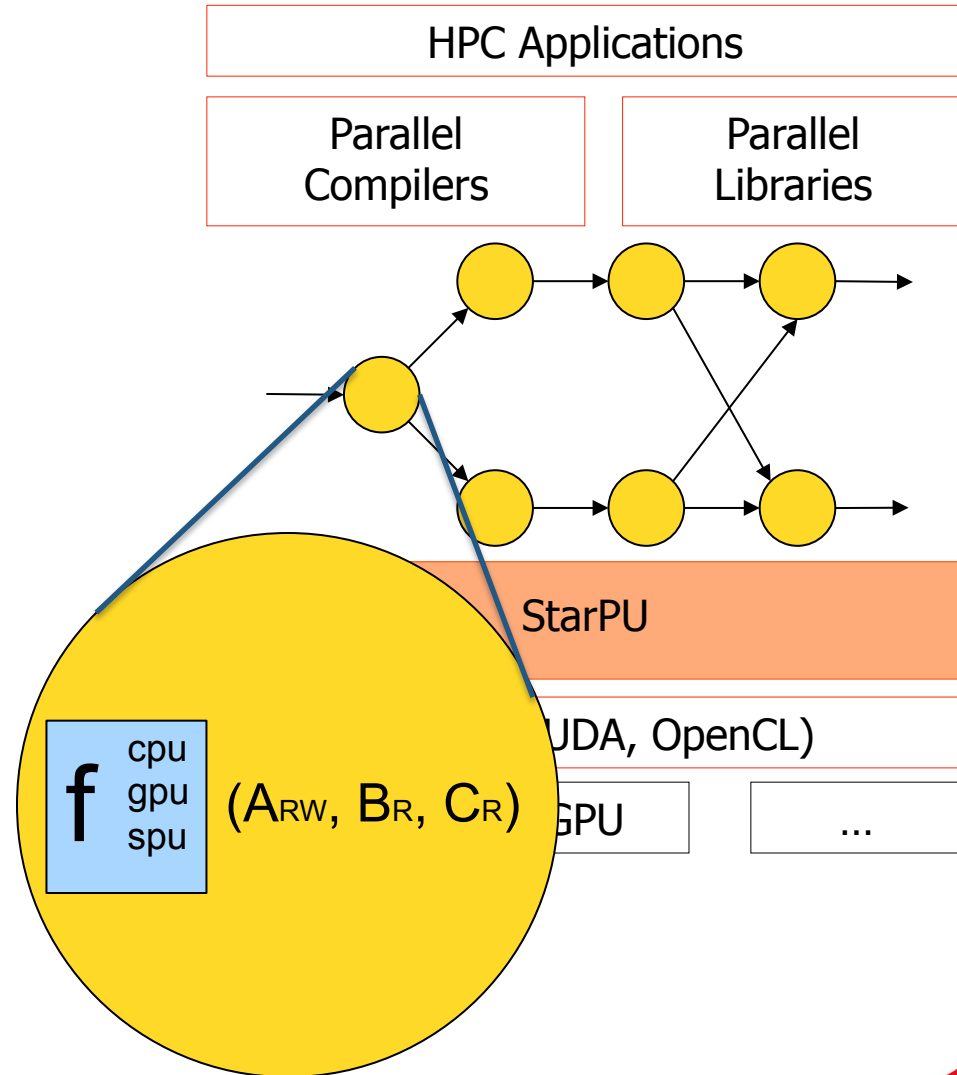
- Ideas
 - Accept tasks that may have multiple implementations
 - Together with potential inter-dependencies
 - Leads to a dynamic acyclic graph of tasks
 - Data-flow approach
 - Provide a high-level data management layer
 - Application should only describe
 - which data may be accessed by tasks
 - How data may be divided



Overview of StarPU

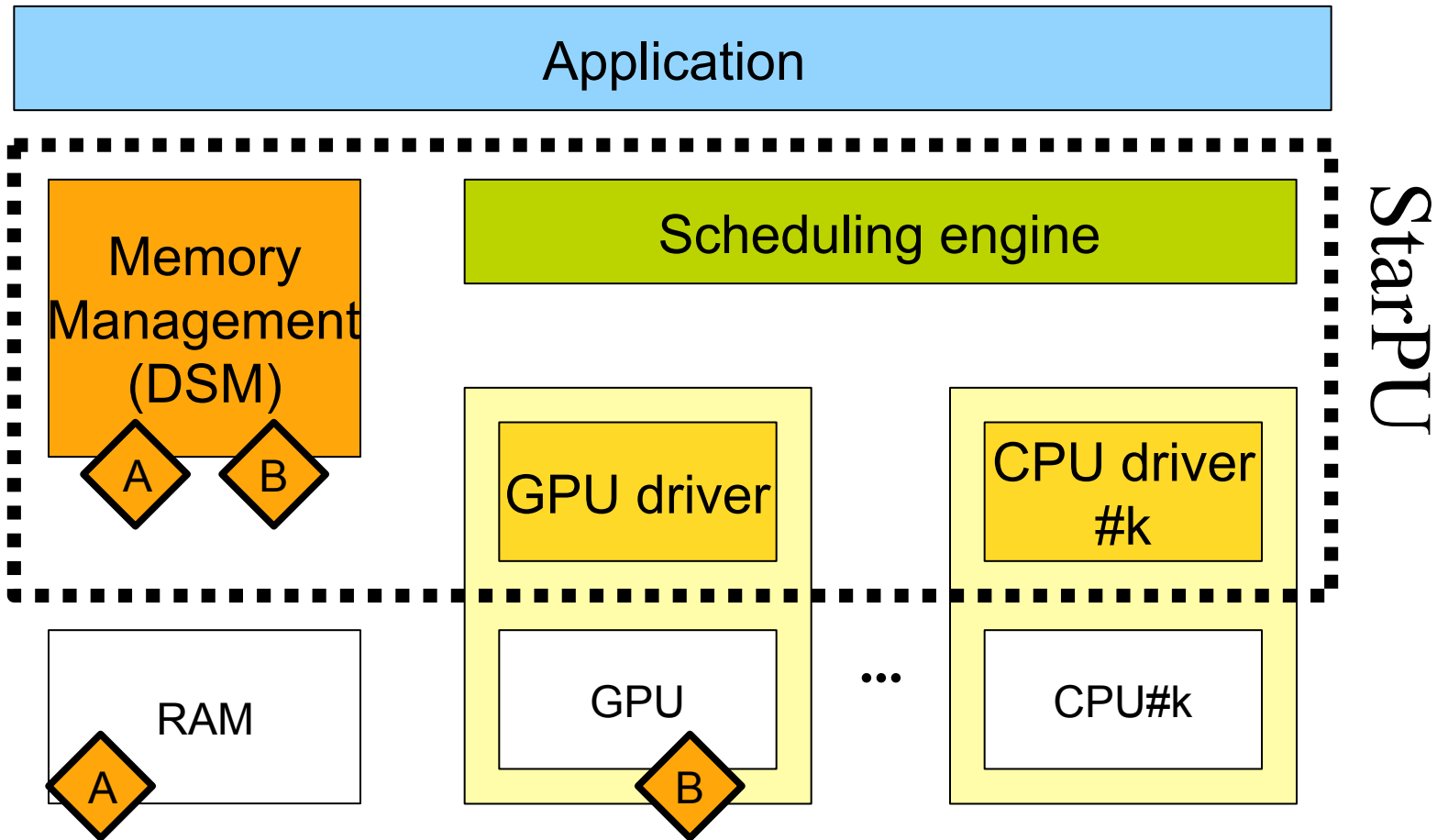
Dealing with heterogeneous hardware accelerators

- Tasks =
 - Data input & output
 - Dependencies with other tasks
 - Multiple implementations
 - E.g. CUDA + CPU implementation
 - Scheduling hints
- StarPU provides an **Open Scheduling platform**
 - Scheduling algorithm = plugins



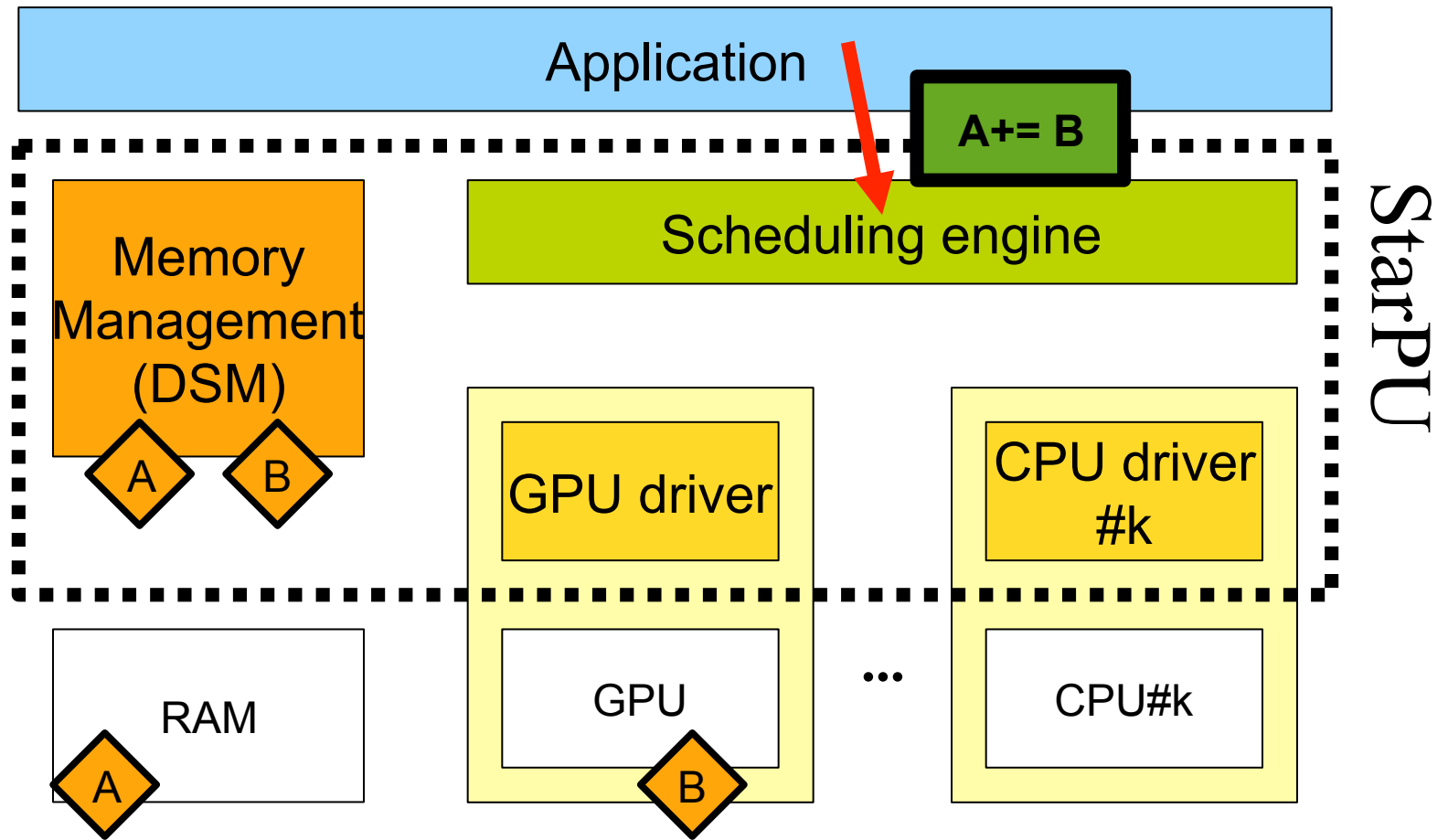
Overview of StarPU

Execution model



Overview of StarPU

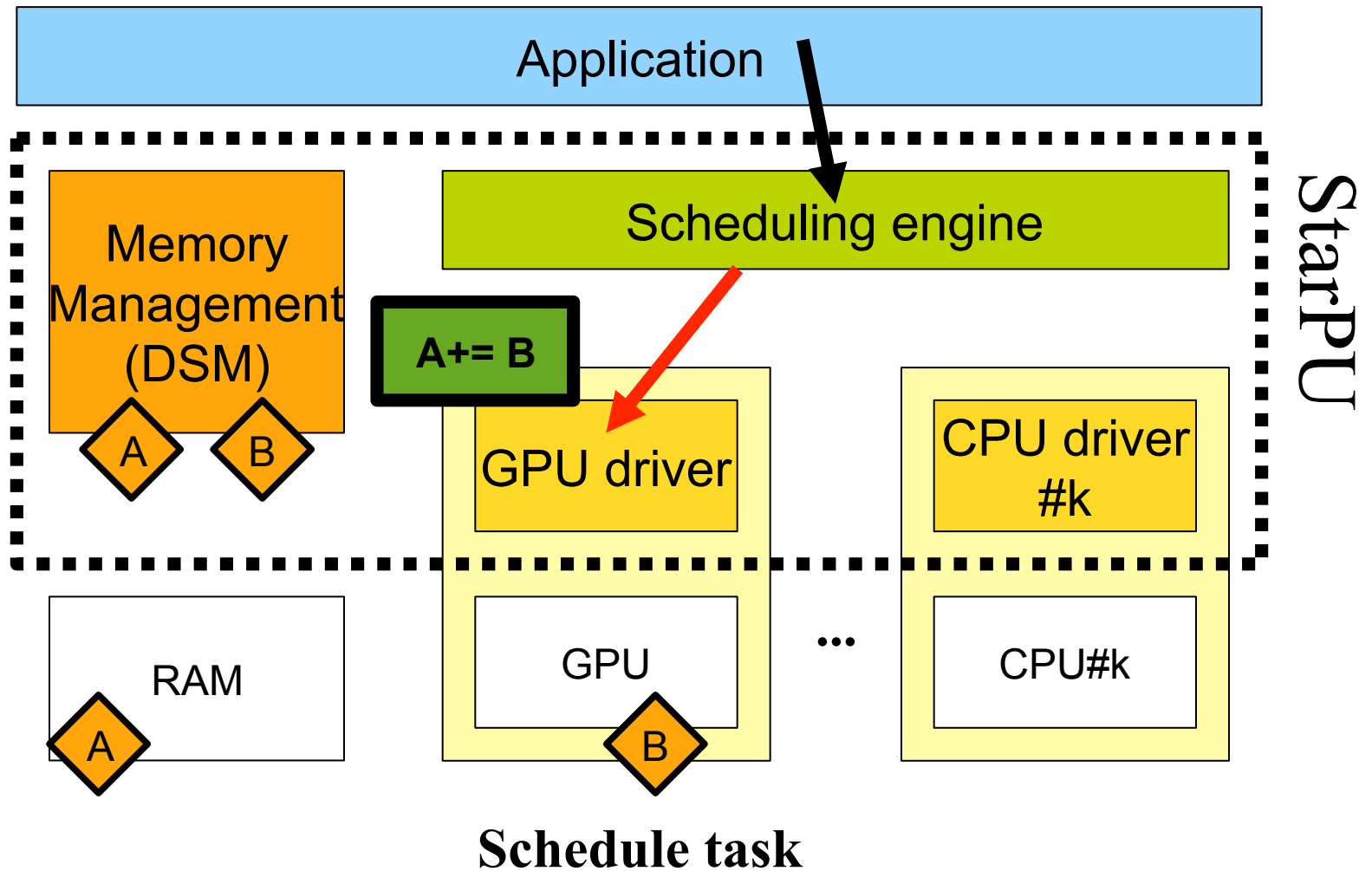
Execution model



Submit task « $A += B$ »

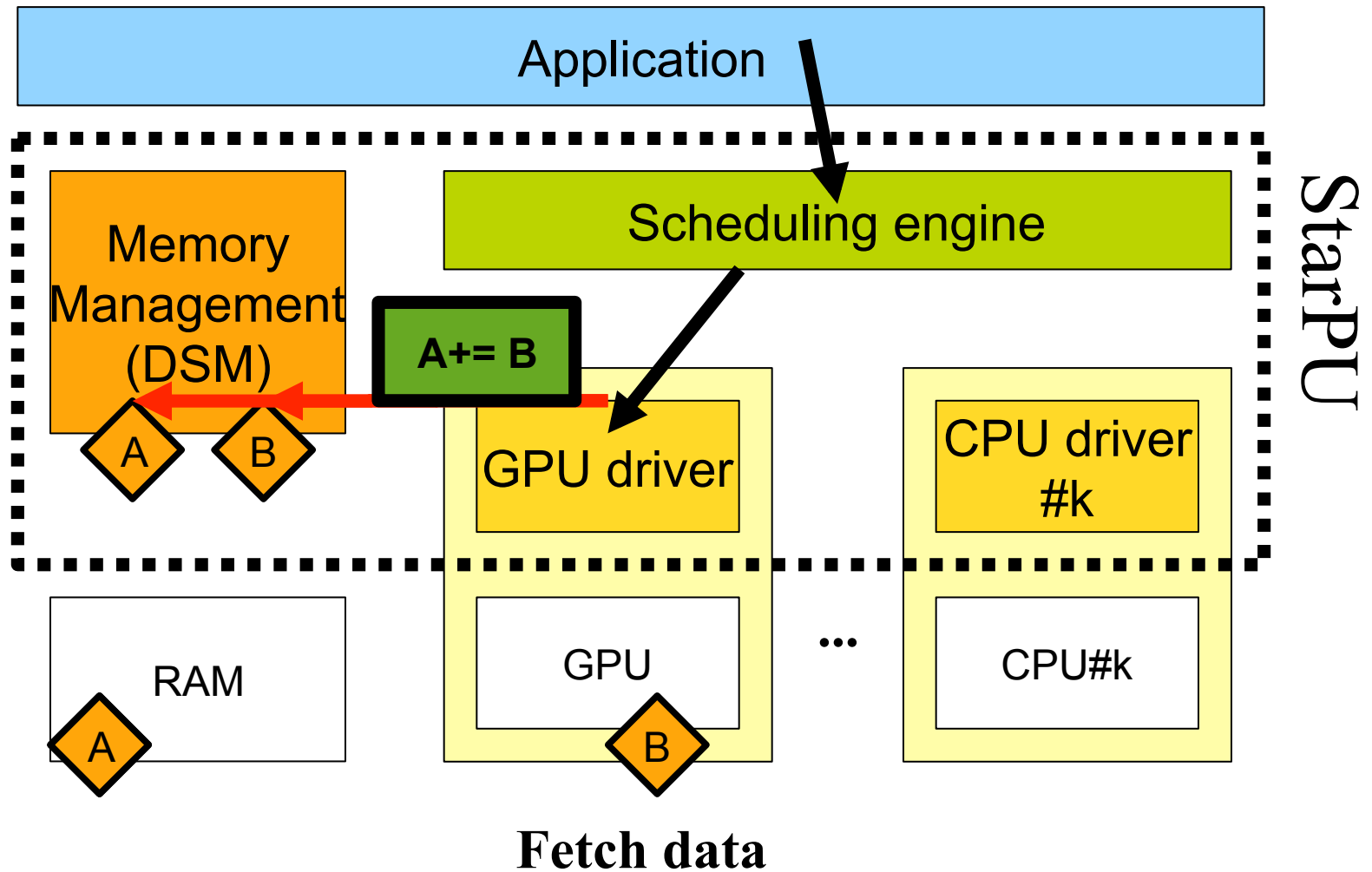
Overview of StarPU

Execution model



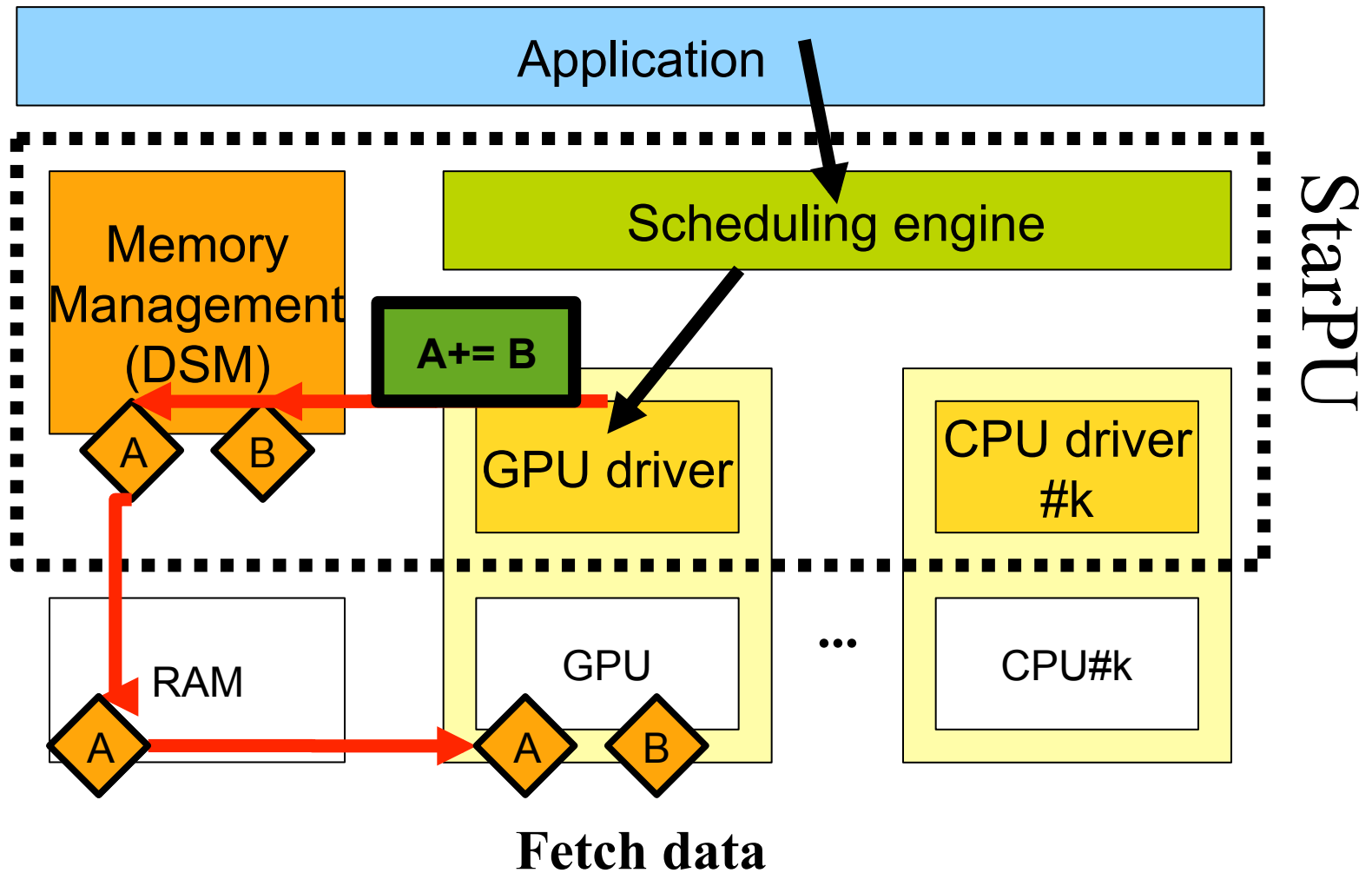
Overview of StarPU

Execution model



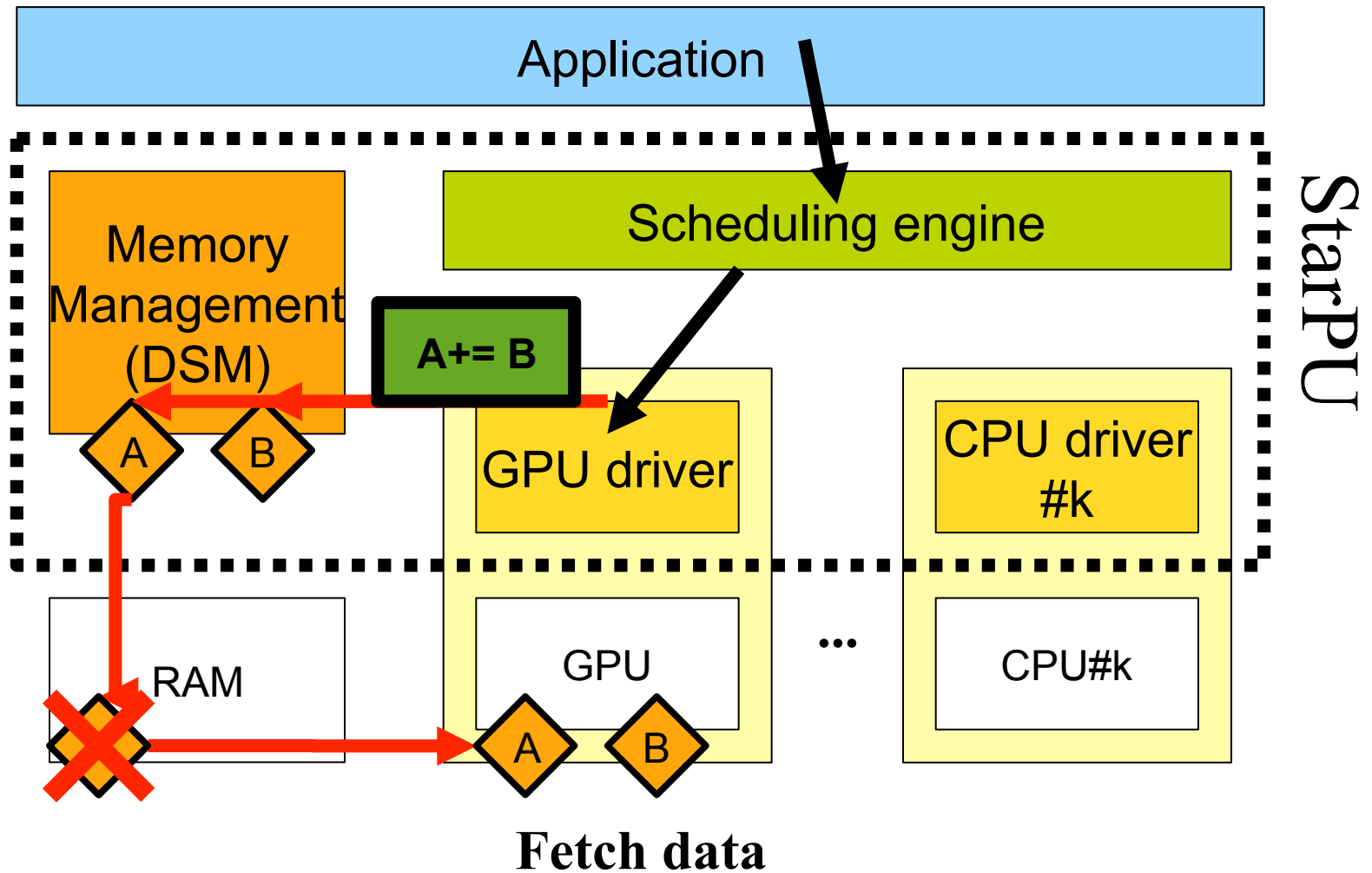
Overview of StarPU

Execution model



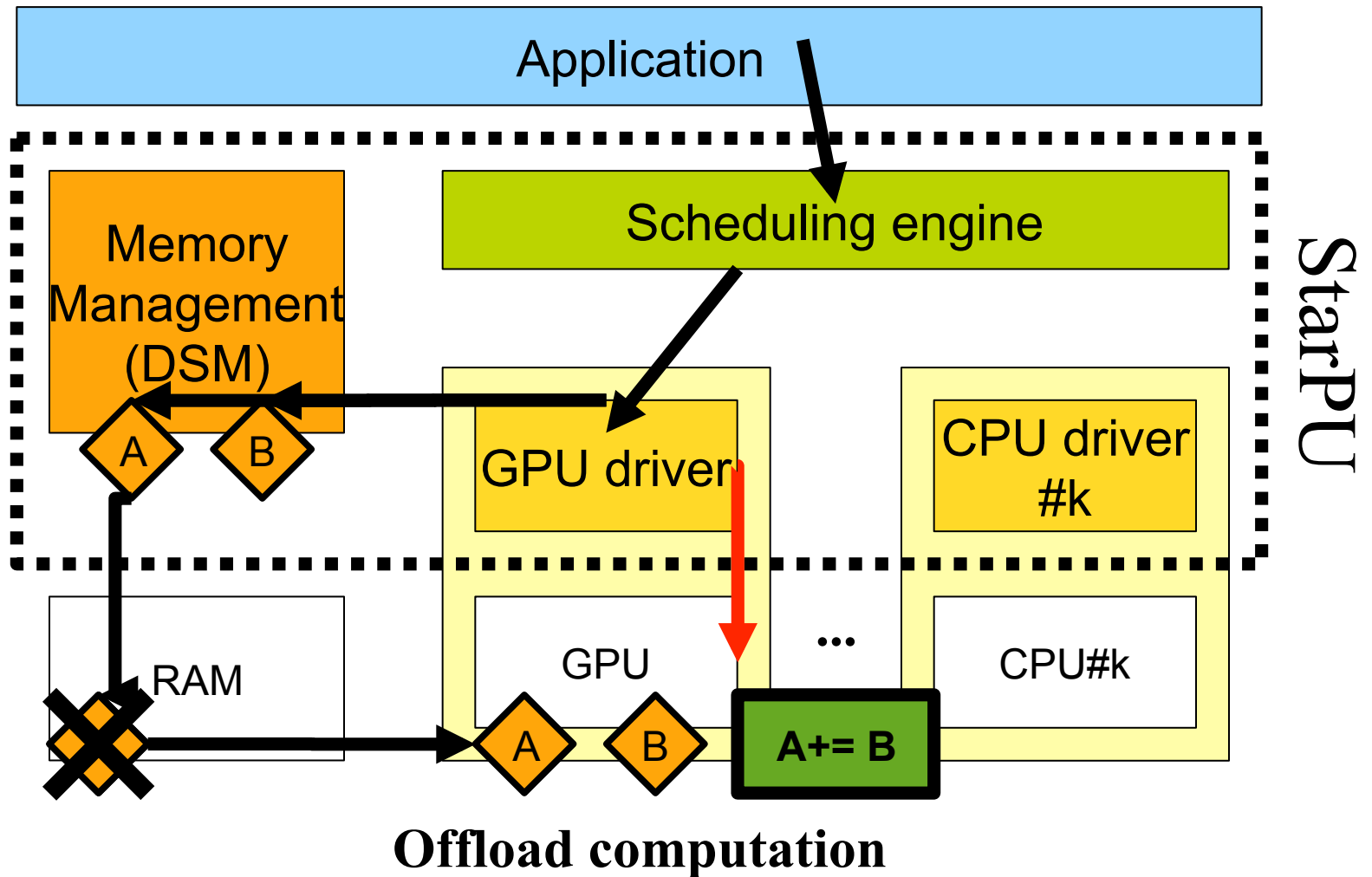
Overview of StarPU

Execution model



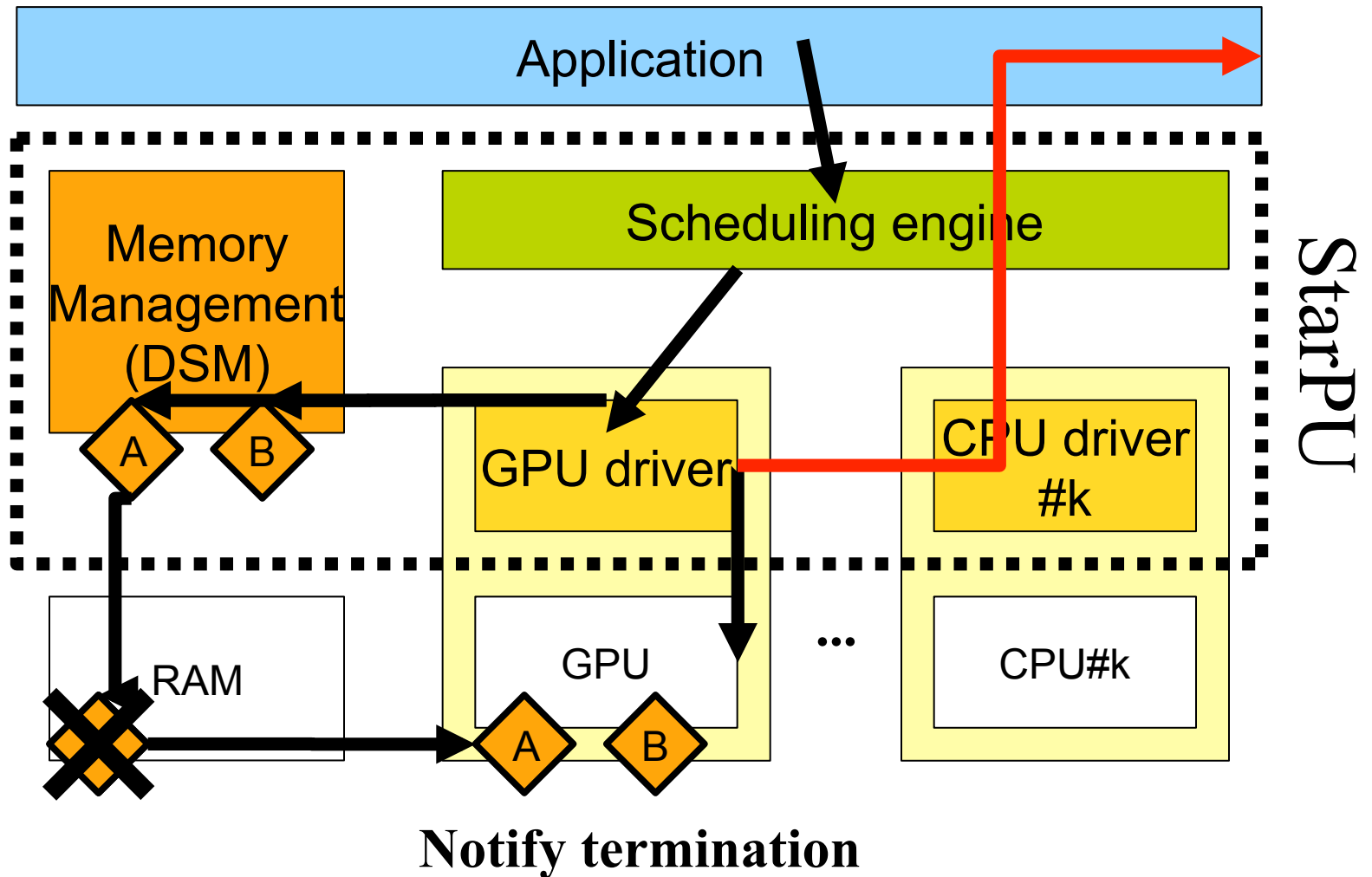
Overview of StarPU

Execution model



Overview of StarPU

Execution model



StarPU's Programming Interface



Scaling vector example

Our first starPU program: vector scaling

Sequential version

```
int main(int argc, char **argv)
{
    float *vector = malloc(NX * sizeof(float)),
          factor = 3.14;
    unsigned i;

    for (i = 0; i < NX; i++)
        vector[i] = (i+1.0f);

    ...

    /* scale the vector */
    for (i = 0; i < NX; i++)
        vector[i] *= factor;

    ...
}
```

Vector scaling

It's all about data

- Registering data will allow the Virtual Shared Memory subsystem to manage
 - Data transfers
 - Data replication
 - Data consistency

```
float *vector = ... ;  
starpu_data_handle vector_handle;
```

```
starpu_vector_data_register(&vector_handle, 0,  
                           (uintptr_t)vector, NX,  
                           sizeof(vector[0]));
```

...

```
starpu_data_unregister(vector_handle);
```

Vector scaling

It's all about data

- Registered data must only be accessed through *handles*
 - To avoid access to dirty copies
 - To retrieve the correct address in local memory
 - GPU embedded memory spaces
 - Data may be copied between cores to avoid NUMA penalties

```
float *vector = ... ;
starpu_data_handle vector_handle;

starpu_vector_data_register(&vector_handle, 0,
                           (uintptr_t)vector, NX,
                           sizeof(vector[0]));

...

starpu_data_unregister(vector_handle);
```

Vector scaling

It's all about data

- StarPU provides a variety of data interfaces to deal with
 - Raw data
 - Single variables
 - Vectors
 - 2D/3D Matrices, BCSR, CSR for sparse matrices
- User-defined interfaces can be added

```
static struct starpu_data_interface_ops_t
interface_vector_ops = {
    .register_data_handle = register_vector_handle,
    .allocate_data_on_node=allocate_vector_buffer_on_node,
    .handle_to_pointer = vector_handle_to_pointer,
    .free_data_on_node = free_vector_buffer_on_node,
    .copy_methods = &vector_copy_data_methods_s,
    .get_size = vector_interface_get_size,
    .compare = vector_compare, ...
};
```

Vector scaling

Before we can create a task...

- StarPU must be able to
 - Choose between multiple implementations at run time
 - Transfer the needed input data to the processing unit if needed
- Because several similar tasks maybe spawned, a “*task type*” is provided to share common information

Vector scaling

Providing multiple implementations

- A StarPU *codelet* is a sort of task type:
 - It contains pointers to each available implementation
 - It stores the number of input+output parameters
 - Only data managed by the VSM, so
nbuffers = 1 vector (the scalar “factor” is ignored)

```
static starpu_codelet cl = {  
    .where = STARPU_CPU | STARPU_CUDA,  
  
    .cpu_func = scal_cpu_func,  
    .cuda_func = scal_cuda_func,  
  
    .nbuffers = 1,  
};
```

Vector scaling

Spawning a task

- The vector scaling task only uses one data from the VSM

```
struct starpu_task *task = starpu_task_create();

task->cl = &cl;

task->buffers[0].handle = vector_handle;
task->buffers[0].mode = STARPU_RW;

float factor = 3.14;
task->cl_arg = &factor;
task->cl_arg_size = sizeof(factor);

starpu_task_submit(task);

starpu_task_wait(task);
```

Vector scaling

Spawning a task

- If no special attribute has to be set, a helper can be used

```
float factor = 3.14;
```

```
    starpu_insert_task(  
        &cl,  
        STARPU_RW, vector_handle,  
        STARPU_VALUE, &factor, sizeof(factor),  
        0);
```

```
starpu_task_wait_for_all();
```

Vector scaling

CPU implementation

```
void scal_cpu_func(void *buffers[], void *cl_arg)
{
    unsigned i;
    float *factor = cl_arg;

    starpu_vector_interface_t *vector = buffers[0];

    unsigned n = STARPU_VECTOR_GET_NX(vector);

    float *val = (float *)STARPU_VECTOR_GET_PTR(vector);

    /* scale the vector */
    for (i = 0; i < n; i++)
        val[i] *= *factor;
}
```

Vector scaling

GPU implementation

```
static __global__ void vector_mult_cuda(float *val, unsigned n,
                                       float factor)
{
    unsigned i = blockIdx.x*blockDim.x + threadIdx.x;
    if (i < n)
        val[i] *= factor;
}

extern "C" void scal_cuda_func(void *buffers[], void *_args)
{
    float *factor = (float *)_args;

    unsigned n = STARPU_VECTOR_GET_NX(buffers[0]);

    float *val = (float *)STARPU_VECTOR_GET_PTR(buffers[0]);
    unsigned threads_per_block = 64;
    unsigned nblocks = (n + threads_per_block - 1) /
                       threads_per_block;

    vector_mult_cuda<<<nblocks, threads_per_block, 0,
                    starpu_cuda_get_local_stream()>>>(val, n, *factor);

    cudaStreamSynchronize(starpu_cuda_get_local_stream());
}
```

Vector scaling

Compiling the program

- Makefile

```
CFLAGS += $(shell pkg-config --cflags libstarpu)
LDFLAGS += $(shell pkg-config --libs libstarpu)
```

```
vector_scal: vector_scal.o \  
             vector_scal_cpu.o \  
             vector_scal_cuda.o
```

```
%.o: %.cu  
     nvcc $(CFLAGS) $< -c $
```

```
clean:  
     rm -f vector_scal *.o
```

Vector scaling

Running the program

- $NX = 1024 * 1024$

```
rnamyst@attila:~/starpu-0.9.1/examples/vec-scal$ ./vector_scal
BEFORE: First element was 1.000000
BEFORE: Last element was 1048576.000000
```

Data transfer statistics:

```
0 -> 1 0.00 MB0.00MB/s      (transfers : 0 - avg -nan MB)
1 -> 0 0.00 MB0.00MB/s      (transfers : 0 - avg -nan MB)
0 -> 2 0.00 MB0.00MB/s      (transfers : 0 - avg -nan MB)
2 -> 0 0.00 MB0.00MB/s      (transfers : 0 - avg -nan MB)
0 -> 3 0.00 MB0.00MB/s      (transfers : 0 - avg -nan MB)
3 -> 0 0.00 MB0.00MB/s      (transfers : 0 - avg -nan MB)
```

Total transfers: 0.00 MB

AFTER: First element is 3.140000

AFTER: Last element is 3292528.750000

Vector scaling

Running the program

- Forcing tasks to run only on GPUs

```
rnamyst@attila:~/starpu-0.9.1/examples/ray-vec-scal$  
STARPU_NCPUS=0 ./vector_scal  
BEFORE: First element was 1.000000  
BEFORE: Last element was 1048576.000000
```

Data transfer statistics:

```
0 -> 1 0.00 MB0.00MB/s (transfers : 0 - avg -nan MB)  
1 -> 0 0.00 MB0.00MB/s (transfers : 0 - avg -nan MB)  
0 -> 2 0.00 MB0.00MB/s (transfers : 0 - avg -nan MB)  
2 -> 0 0.00 MB0.00MB/s (transfers : 0 - avg -nan MB)  
0 -> 3 4.00 MB1443.69MB/s (transfers : 1 - avg 4.00 MB)  
3 -> 0 0.00 MB0.00MB/s (transfers : 0 - avg -nan MB)
```

Total transfers: 4.00 MB

```
AFTER: First element is 3.140000  
AFTER: Last element is 3292528.750000
```


Vector Scaling

Generating execution traces

- StarPU uses the FXT library to generate traces
 - <http://download.savannah.gnu.org/releases/fkt/fxt-0.2.2.tar.gz>
`cd starpu ; configure --with-fxt ; make install`
- Starting/stopping trace generation

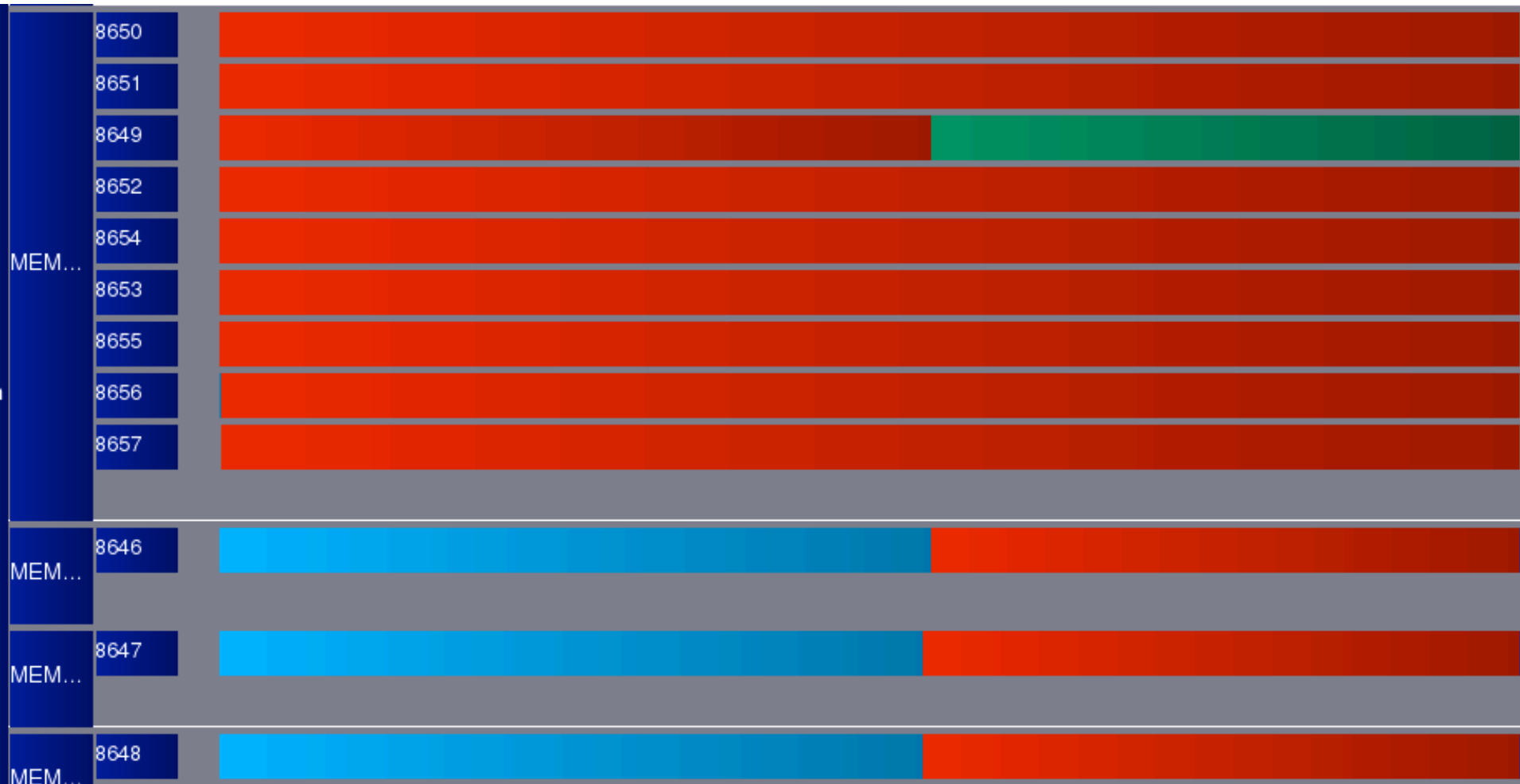
```
starpu_profiling_status_set(STARPU_PROFILING_ENABLE);
starpu_profiling_status_set(STARPU_PROFILING_DISABLE);
```
- A `/tmp/prof_file_xx_yy` is generated

```
starpu_fxt_tool -i /tmp/prof_file_xx_yy
```

 - This translates the raw trace into a “paje” trace
- The paje trace can be displayed using ViTE
 - <http://vite.gforge.inria.fr/>

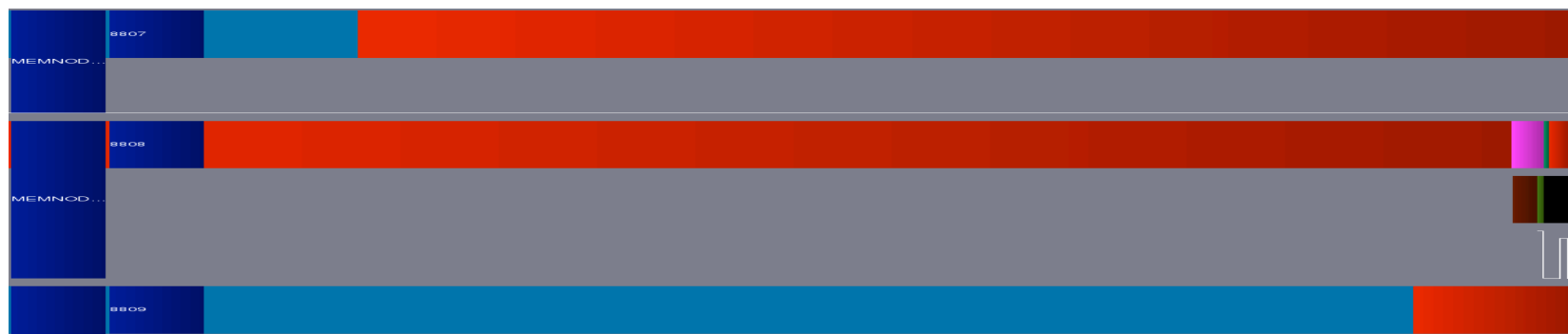
Vector Scaling

Default execution



Vector Scaling

Forcing the execution on a GPU



Divide and Conquer

inria
informatics mathematics

Vector scaling

Splitting the computation

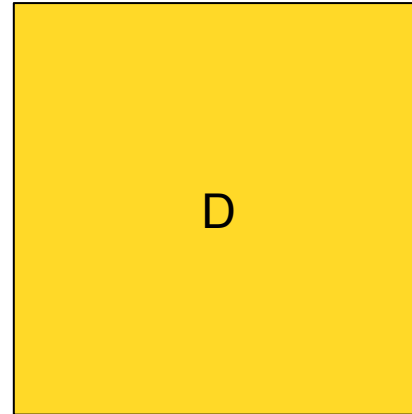
- The VSM subsystem provides *filters* to divide registered data into smaller chunks
 - E.g. vertical/horizontal block filtering
 - Filters produce multiple children handles from a single data handle
- Filters can be applied recursively
- Current limitation: parent data handles cannot be used
 - While a data is divided, only leaves of the tree can be used by tasks
 - Homogeneous granularity

Partitioning data

Using filters

```
starpu_data_handle D, subD;  
filter f1, f2;  
  
starpu_matrix_data_register(&D, 0, ptr, ld, nx,  
ny, 4);  
  
f1.func = starpu_block_filter_func;  
f1.arg = 3;  
  
f2.func = starpu_vertical_block_filter_func;  
f2.arg = 3;  
  
starpu_data_map_filters (D, 2, &f1, &f2);
```

```
subD = starpu_data_get_sub_data (D, 2, i, j);  
starpu_data_acquire (subD, RW);  
// Use subD.ptr  
starpu_data_release (subD);
```

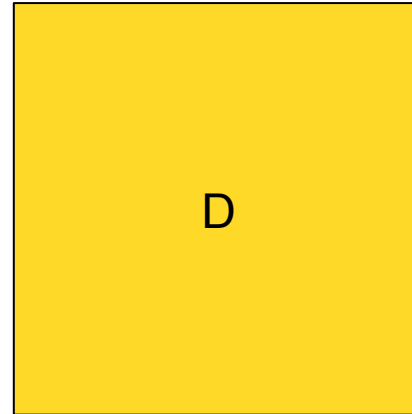


Partitioning data

Using filters

```
starpu_data_handle D, subD;  
filter f1, f2;  
  
starpu_matrix_data_register(&D, 0, ptr, ld, nx,  
ny, 4);  
  
f1.func = starpu_block_filter_func;  
f1.arg = 3;  
  
f2.func = starpu_vertical_block_filter_func;  
f2.arg = 3;  
  
starpu_data_map_filters (D, 2, &f1, &f2);
```

```
subD = starpu_data_get_sub_data (D, 2, i, j);  
starpu_data_acquire (subD, RW);  
// Use subD.ptr  
starpu_data_release (subD);
```

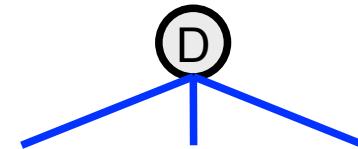
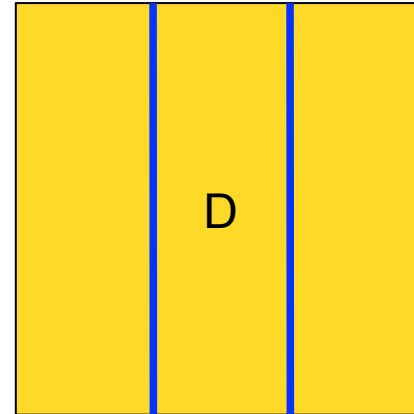


Partitioning data

Using filters

```
starpu_data_handle D, subD;  
filter f1, f2;  
  
starpu_matrix_data_register(&D, 0, ptr, ld, nx,  
ny, 4);  
  
f1.func = starpu_block_filter_func;  
f1.arg = 3;  
  
f2.func = starpu_vertical_block_filter_func;  
f2.arg = 3;  
  
starpu_data_map_filters (D, 2, &f1, &f2);
```

```
subD = starpu_data_get_sub_data (D, 2, i, j);  
starpu_data_acquire (subD, RW);  
// Use subD.ptr  
starpu_data_release (subD);
```

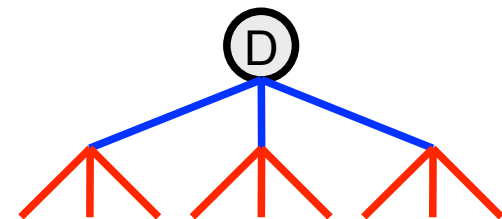
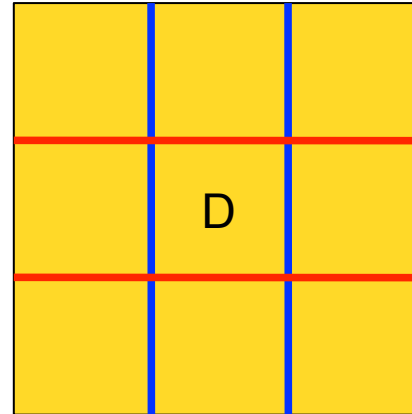


Partitioning data

Using filters

```
starpu_data_handle D, subD;  
filter f1, f2;  
  
starpu_matrix_data_register(&D, 0, ptr, ld, nx,  
ny, 4);  
  
f1.func = starpu_block_filter_func;  
f1.arg = 3;  
  
f2.func = starpu_vertical_block_filter_func;  
f2.arg = 3;  
  
starpu_data_map_filters (D, 2, &f1, &f2);
```

```
subD = starpu_data_get_sub_data (D, 2, i, j);  
starpu_data_acquire (subD, RW);  
// Use subD.ptr  
starpu_data_release (subD);
```



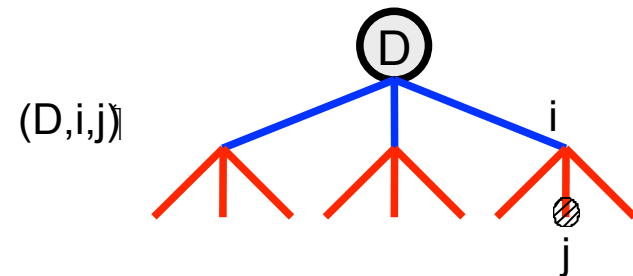
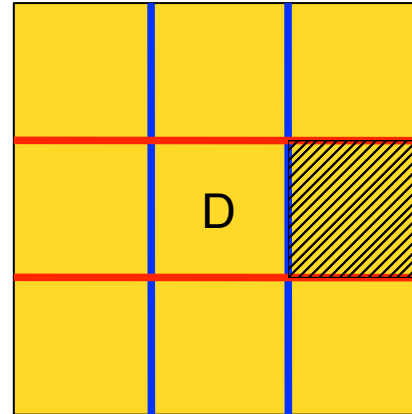
Partitioning data

Using filters

```
starpu_data_handle D, subD;  
filter f1, f2;  
  
starpu_matrix_data_register(&D, 0, ptr, ld, nx,  
ny, 4);  
  
f1.func = starpu_block_filter_func;  
f1.arg = 3;  
  
f2.func = starpu_vertical_block_filter_func;  
f2.arg = 3;  
  
starpu_data_map_filters (D, 2, &f1, &f2);
```

→

```
subD = starpu_data_get_sub_data (D, 2, i, j);  
starpu_data_acquire (subD, RW);  
// Use subD.ptr  
starpu_data_release (subD);
```

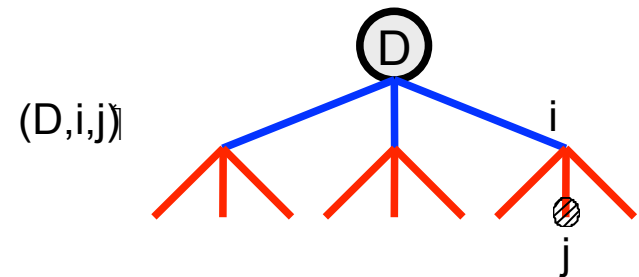
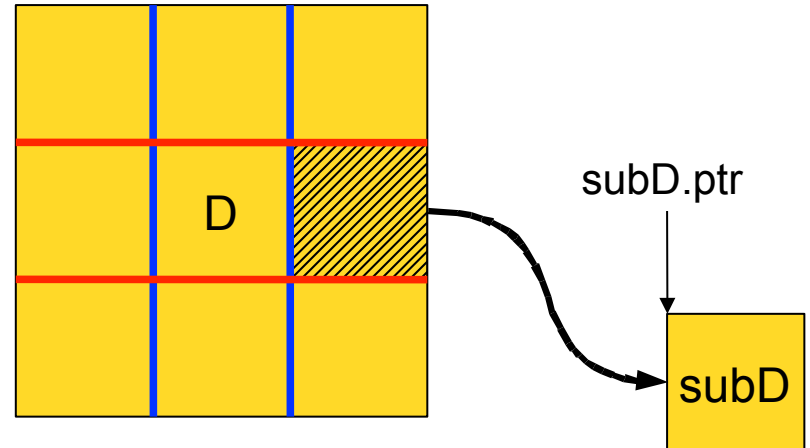


Partitioning data

Using filters

```
starpu_data_handle D, subD;  
filter f1, f2;  
  
starpu_matrix_data_register(&D, 0, ptr, ld, nx,  
ny, 4);  
  
f1.func = starpu_block_filter_func;  
f1.arg = 3;  
  
f2.func = starpu_vertical_block_filter_func;  
f2.arg = 3;  
  
starpu_data_map_filters (D, 2, &f1, &f2);
```

```
subD = starpu_data_get_sub_data (D, 2, i, j);  
starpu_data_acquire (subD, RW);  
// Use subD.ptr  
starpu_data_release (subD);
```



Vector scaling

Back to our example

- Splitting our vector in 8 slices:

```
#define NSLICES 8

starpu_data_handle vector_handle;

    starpu_vector_data_register(&vector_handle, 0
                               (uintptr_t)vector,
                               NX, sizeof(vector[0]));

struct starpu_data_filter vert = {
    .filter_func = starpu_block_filter_func_vector,
    .nchildren = NSLICES,
    .get_nchildren = NULL,
    .get_child_ops = NULL
};

    starpu_data_partition(vector_handle, &vert);
```

Vector scaling

Back to our example

- Spawning 8 tasks:

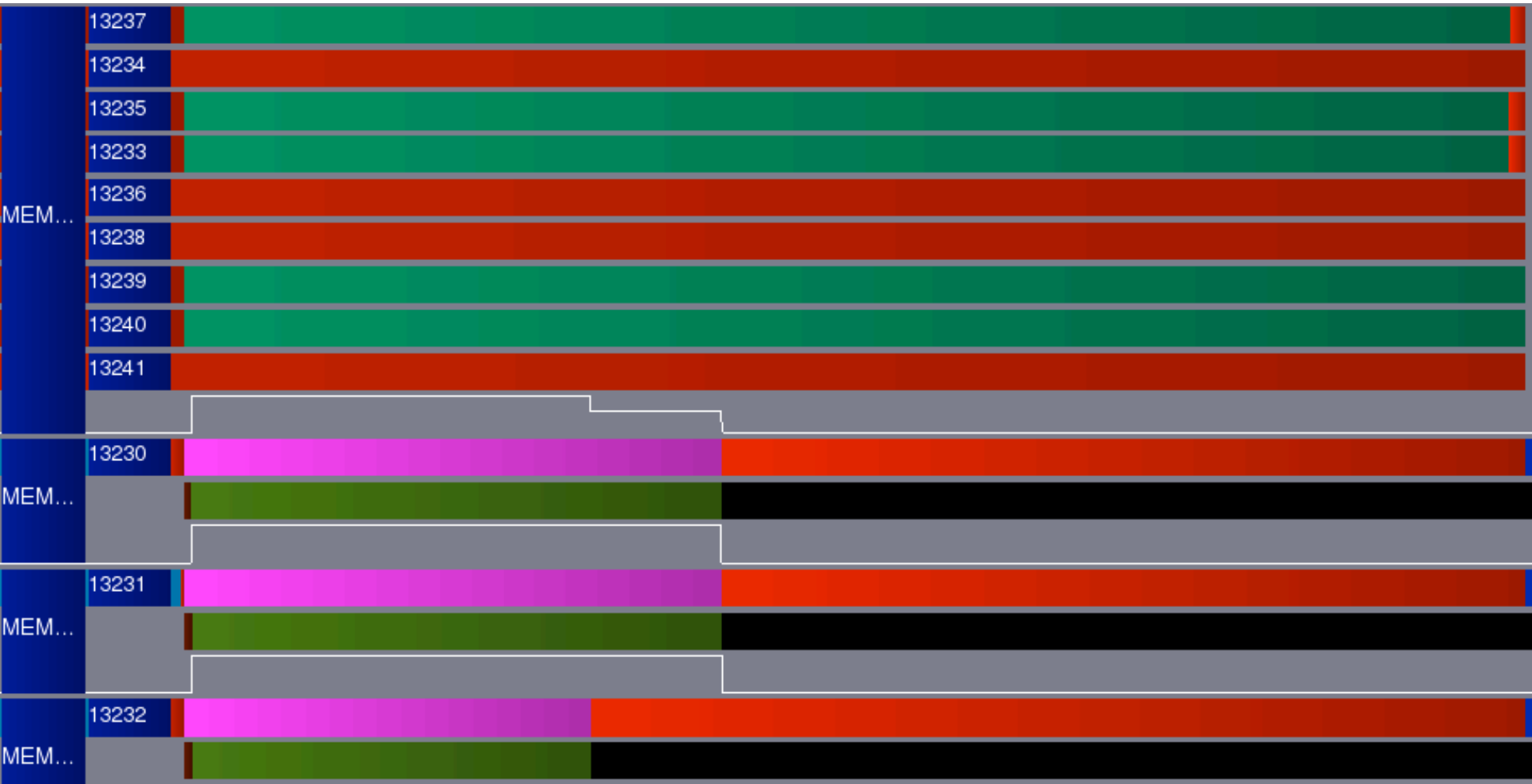
```
for (t = 0; t < NSLICES; t++)
{
    struct starpu_task *task = starpu_task_create();

    starpu_insert_task(
        &cl,
        STARPU_RW, starpu_data_get_sub_data(vector_handle, 1, t),
        STARPU_VALUE, &factor, sizeof(factor),
        0);
}

starpu_task_wait_for_all();
```

Vector Scaling

$NX=1024*1024$, $NSLICES=8$

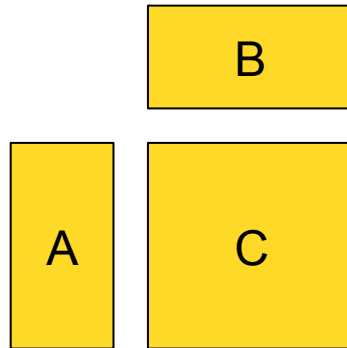


Another example: Blocked Matrix Multiplication



Blocked matrix multiplication

Data partitionning



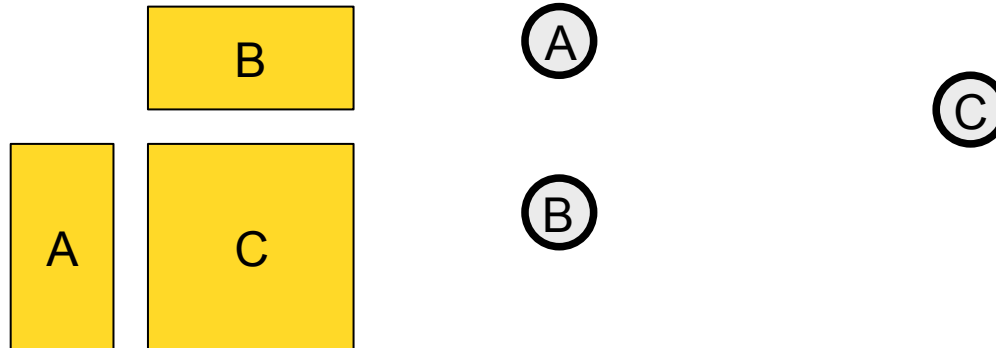
```
filter f1 = {.func = block, .arg = 3};  
filter f2 = {.func = vertical_block, .arg = 3};
```

```
starpu_register_data(&A, pA, ldA, nxA, nyA, sizeof(float));  
starpu_register_data(&B, pB, ldB, nxB, nyB, sizeof(float));  
starpu_register_data(&C, pC, ldC, nxC, nyC, sizeof(float));
```

```
starpu_data_partition(&A, &f2);  
starpu_data_partition(&B, &f1);  
starpu_data_partition(&C, 2, &f1, &f2);
```


Blocked matrix multiplication

Data partitionning



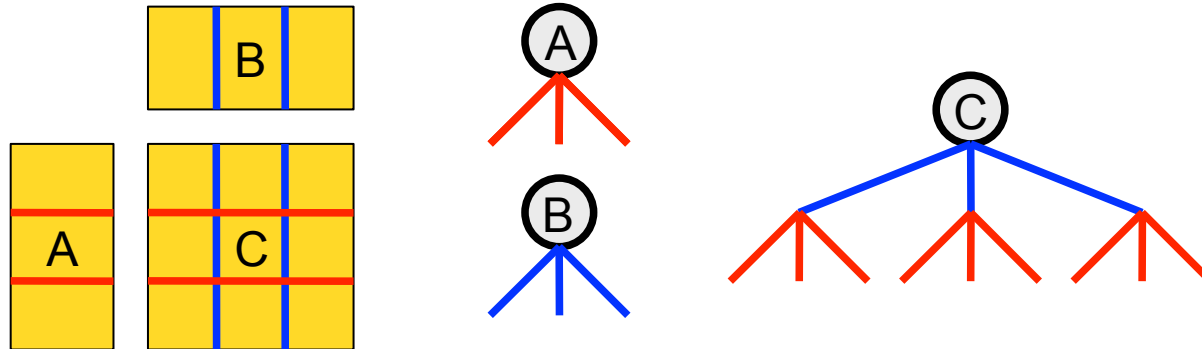
```
filter f1 = {.func = block, .arg = 3};  
filter f2 = {.func = vertical_block, .arg = 3};
```

```
starpu_register_data(&A, pA, ldA, nxA, nyA, sizeof(float));  
starpu_register_data(&B, pB, ldB, nxB, nyB, sizeof(float));  
starpu_register_data(&C, pC, ldC, nxC, nyC, sizeof(float));
```

```
starpu_data_partition(&A, &f2);  
starpu_data_partition(&B, &f1);  
starpu_data_partition(&C, 2, &f1, &f2);
```

Blocked matrix multiplication

Data partitionning



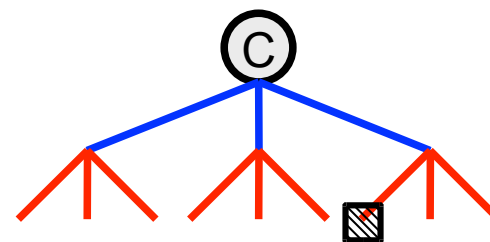
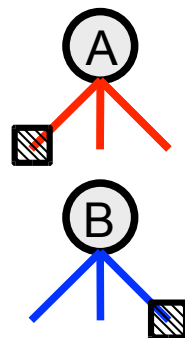
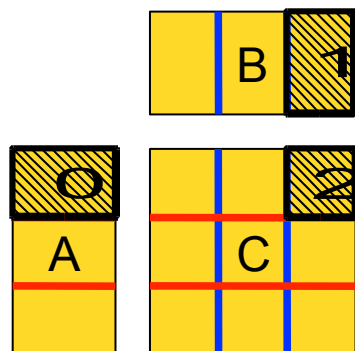
```
filter f1 = {.func = block, .arg = 3};  
filter f2 = {.func = vertical_block, .arg = 3};
```

```
starp_u_register_data(&A, pA, ldA, nxA, nyA, sizeof(float));  
starp_u_register_data(&B, pB, ldB, nxB, nyB, sizeof(float));  
starp_u_register_data(&C, pC, ldC, nxC, nyC, sizeof(float));
```

```
starp_u_data_partition(&A, &f2);  
starp_u_data_partition(&B, &f1);  
starp_u_data_partition(&C, 2, &f1, &f2);
```

Blocked matrix multiplication

Data partitionning



```
task->nbuffers = 3;  
task->buf[0].state = starpu_data_get_sub_data (&A, 1, j);  
task->buf[0].mode = R;  
task->buf[1].state = starpu_data_get_sub_data (&B, 1, i);  
task->buf[1].mode = R;  
task->buf[2].state = starpu_data_get_sub_data (&C, 2, i, j);  
task->buf[2].mode = W;
```

```
task->cublas_func = gpu_mult;  
task->core_func = core_mult;
```

```
starpu_submit_task(task);
```

```
void gpu_mult(buffer_descr *dsc, void *arg)  
{  
    gpu_sgemm(..., dsc[2].ptr, ..);  
}  
  
void core_mult(buffer_descr *dsc, void *arg)  
{  
    cblas_sgemm(..., dsc[2].ptr, ..);  
}
```

Blocked matrix multiplication

Main loop

```
unsigned taskx, tasky;

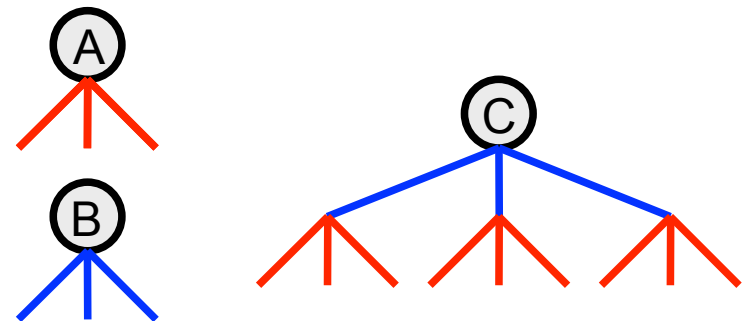
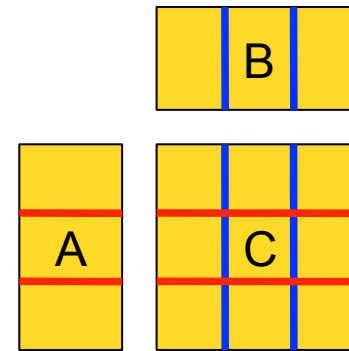
for (taskx = 0; taskx < nslicesx; taskx++)
{
  for (tasky = 0; tasky < nslicesy; tasky++)
  {
    struct starpu_task *task =
      starpu_task_create();

    task->buffers[0].handle =
      starpu_data_get_sub_data(A_handle, 1,
                               tasky);
    task->buffers[0].mode = STARPU_R;

    task->buffers[1].handle =
      starpu_data_get_sub_data(B_handle, 1,
                               taskx);
    task->buffers[1].mode = STARPU_R;

    task->buffers[2].handle =
      starpu_data_get_sub_data(C_handle, 2,
                               taskx, tasky);
    task->buffers[2].mode = STARPU_W;

    starpu_task_submit(task);
  }
}
```



Memory Consistency

Updating main memory

- Most StarPU programs end with a synchronization primitive waiting for all tasks to complete:

```
starpu_task_wait_for_all();
```

- This does not guarantee that data have been updated in main memory!
 - StarPU tries to minimize data transfers...

- To enforce update of main memory:

```
starpu_data_acquire(handle, STARPU_R);  
starpu_data_release(handle);
```

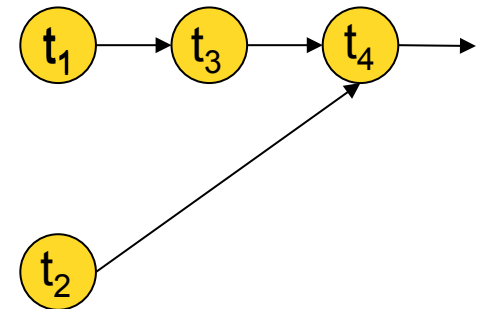
- Note: *unregistering* data is also OK

Memory Consistency

- By default, the StarPU VSM enforces sequential consistency
 - Dependencies between tasks are inferred

```
// t1
starpu_insert_task(...,
                    STARPU_RW, A_handle, 0);
// t2
starpu_insert_task(...,
                    STARPU_R,  B_handle,
                    STARPY_W,  C_handle, 0);
// t3
starpu_insert_task(...,
                    STARPU_R,  A_handle,
                    STARPU_RW, D_handle, 0);
// t4
starpu_insert_task(...,
                    STARPU_R,  C_handle,
                    STARPU_R,  D_handle, 0);
```

Task dependencies



Memory Consistency

Setting dependencies manually

- Disabling automatic dependencies

```
starpu_data_set_default_sequential_consistency_flag(0);
```

- 64-bit tags can be associated to tasks

```
task->use_tag=1  
task->tag_id = <mytag>
```

- Specifying dependencies between tasks is easy

```
// Tag 0x1 depends on tags 0x32 and 0x52  
starpu_tag_t tag_array[2] = {0x32, 0x52};  
starpu_tag_declare_deps_array((starpu_tag_t)0x1,  
                               2, tag_array);
```

- Tags can be unlocked by the application

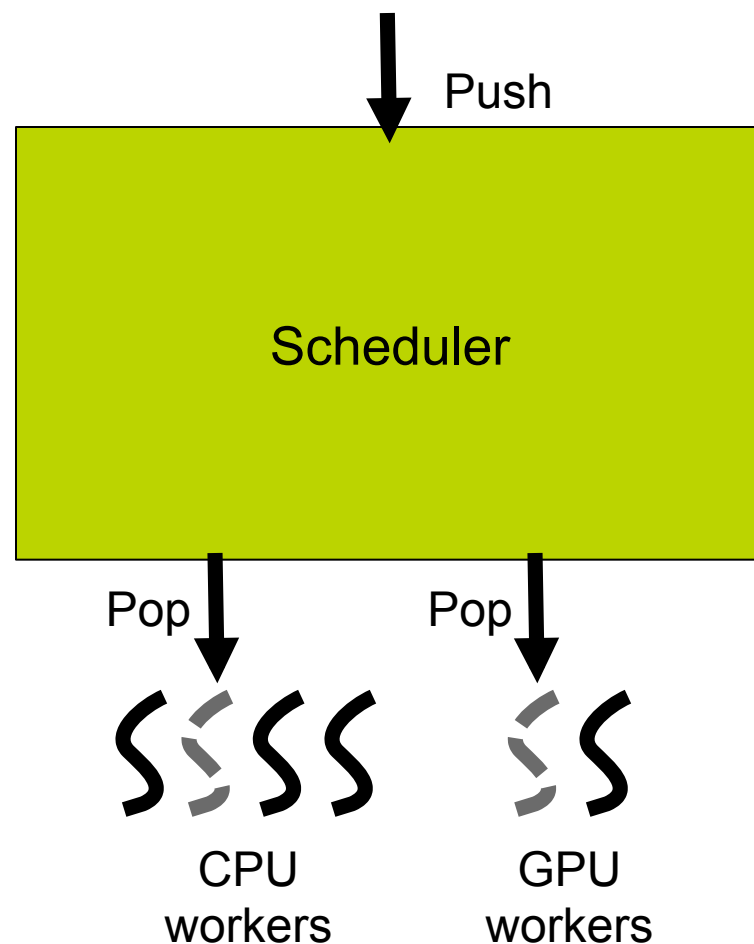
- Tasks can wait for external events (e.g. arrival of a MPI msg)

```
starpu_tag_notify_from_apps (<tag>);
```

Tasks 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 actively poll for work (“pop”)
- What happens inside the scheduler is... up to you!

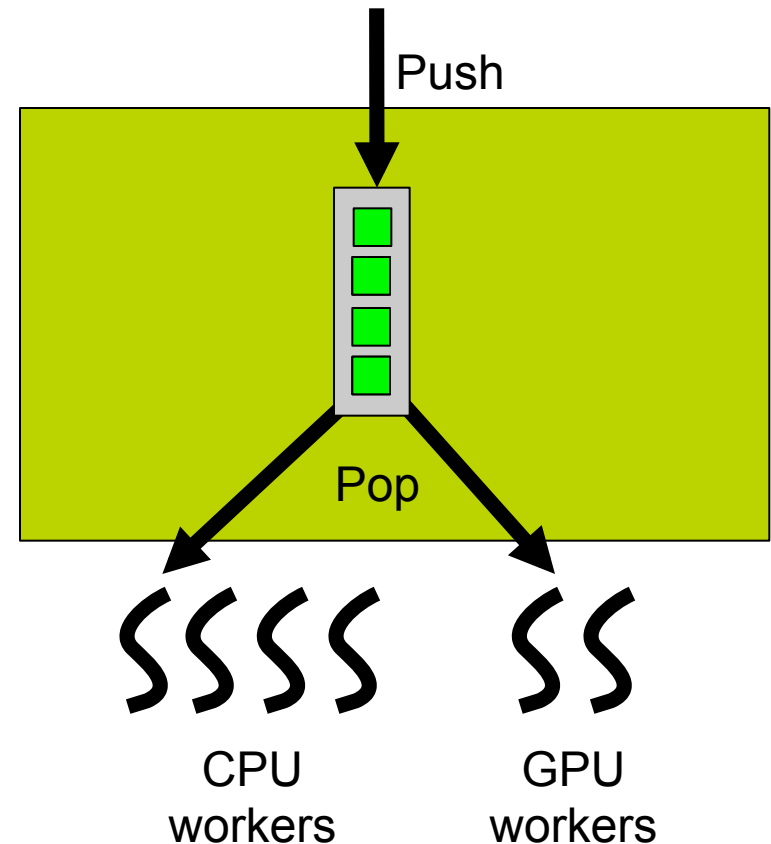
How does it work?



Tasks scheduling

- Queue based scheduler
 - Each worker « pops » task in a specific queue
- Implementing a strategy
 - Easy!
 - Select queue topology
 - Implement « pop » and « push »
 - Priority tasks
 - Work stealing
 - Performance models, ...
- Scheduling algorithms testbed

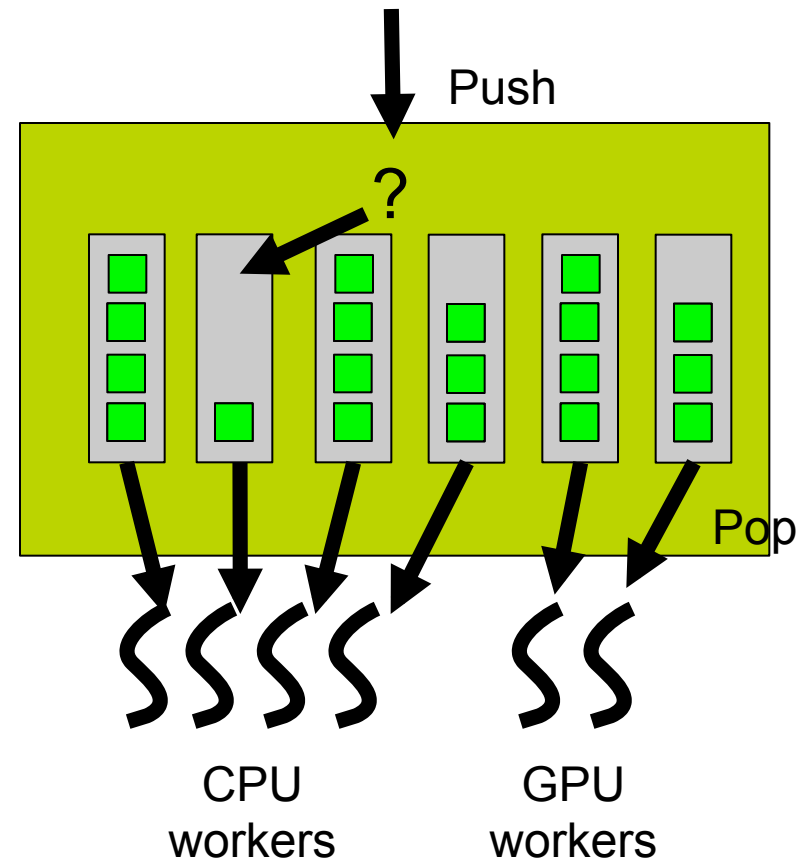
Developing your own scheduler



Tasks scheduling

- Queue based scheduler
 - Each worker « pops » task in a specific queue
- Implementing a strategy
 - Easy!
 - Select queue topology
 - Implement « pop » and « push »
 - Priority tasks
 - Work stealing
 - Performance models, ...
- Scheduling algorithms testbed

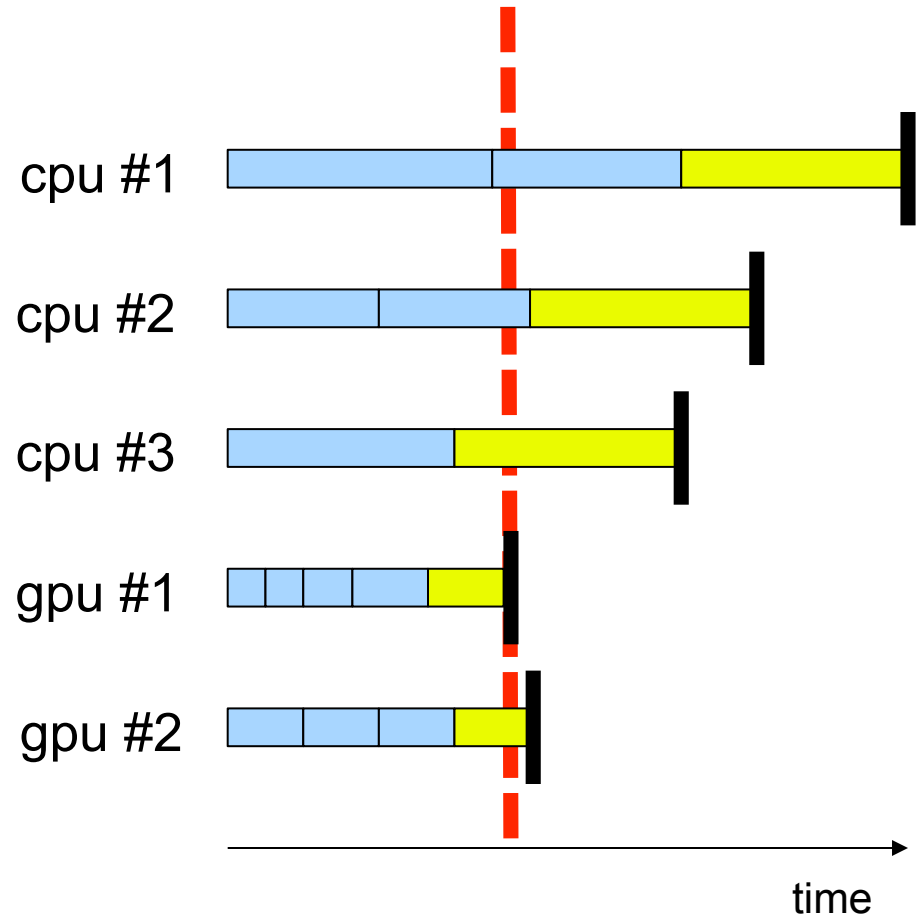
Developing your own scheduler



Dealing with heterogeneous architectures

- Task completion time estimation
 - History-based
 - User-defined cost function
 - Parametric cost model
- Can be used to improve scheduling
 - E.g. Heterogeneous Earliest Finish Time

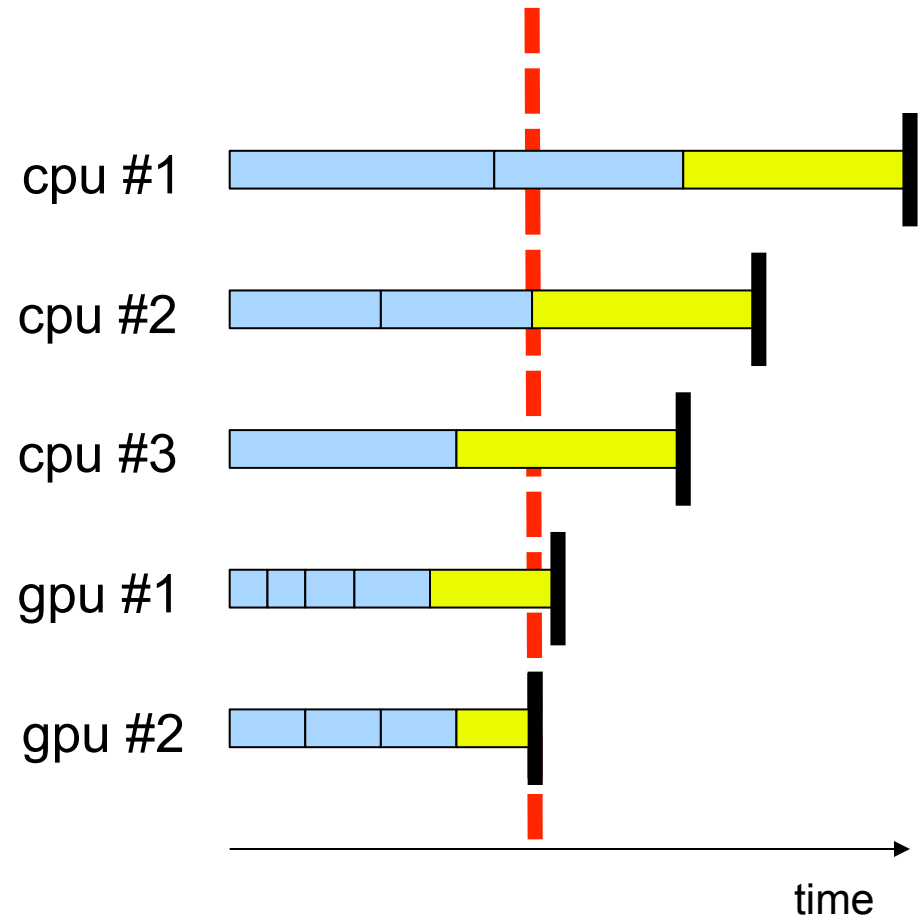
Performance prediction



Dealing with heterogeneous architectures

- Data transfer time estimation
 - Sampling based on off-line calibration
- Can be used to
 - Better estimate overall exec time
 - Minimize data movements

Performance prediction



Controlling scheduling...

...from the command line

- Environment variables
 - STARPU_NCPUS=2
 - STARPU_NCUDA=0
 - STARPU_WORKERS_CPUID = "0 1 4 5"
 - STARPU_WORKERS_CUDAID = "1 3"

 - STARPU_SCHED=heft
 - eager, random, dmdas, etc.
 - STARPU_CALIBRATE=0|1|2

Performance feedback

Getting feedback at run time

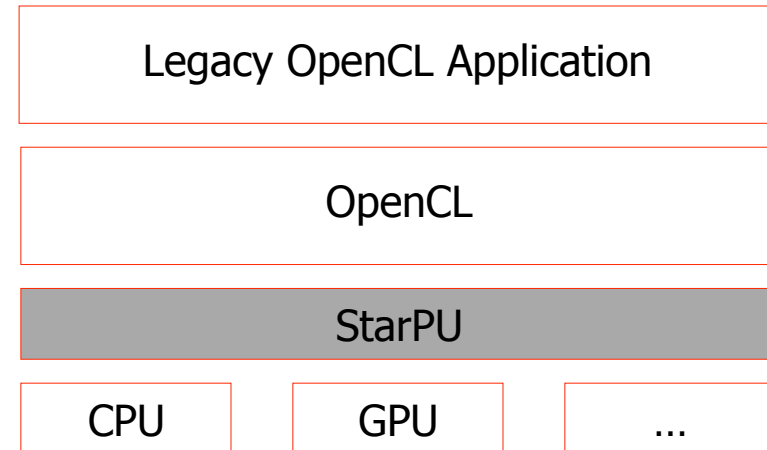
- When profiling is enabled, per-task running information can be collected
 - `profiling_info` field of `starp_u_task` structure
- Information is also cumulated per processing unit
 - See `starp_u_worker_get_profiling_info`

```
struct starpu_task_profiling_info {  
    struct timespec submit_time;  
    struct timespec start_time;  
    struct timespec end_time;  
  
    int workerid;  
  
    uint64_t used_cycles;  
    uint64_t stall_cycles;  
    double power_consumed;  
  
    ...  
};
```

Using StarPU through a standard API

A StarPU driver for OpenCL

- Run legacy OpenCL codes on top of StarPU
 - OpenCL sees a number of starPU devices
- Performance limitations
 - Data transfers performed just-in-time
 - Data replication not managed by StarPU
- Ongoing work
 - We propose light extensions to OpenCL
 - Greatly improves flexibility when used
 - Regular OpenCL behavior if no extension is used



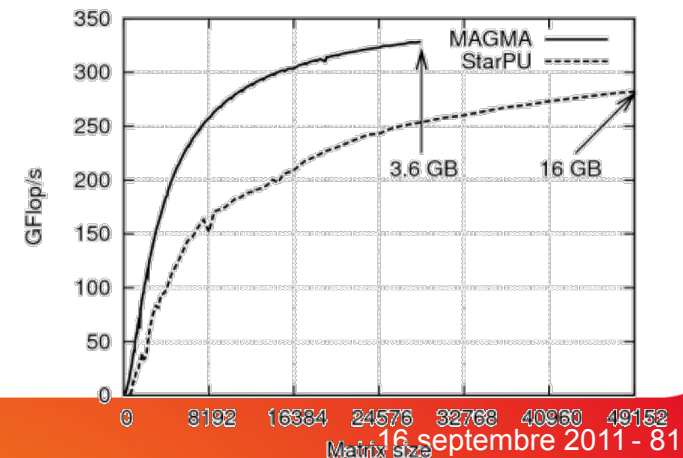
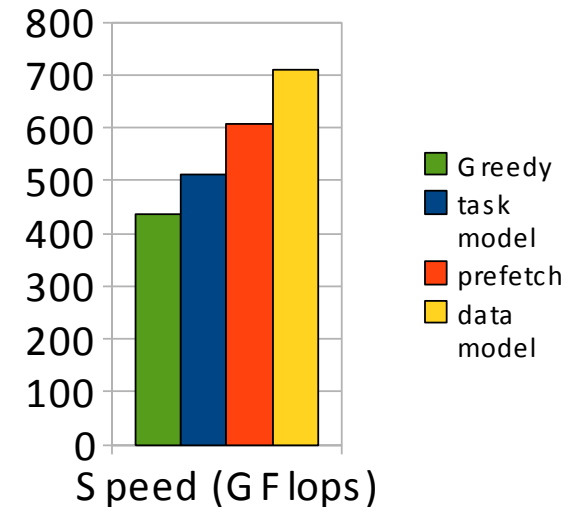
Parallel Dense Linear Algebra over StarPU



Dealing with heterogeneous architectures

Performance

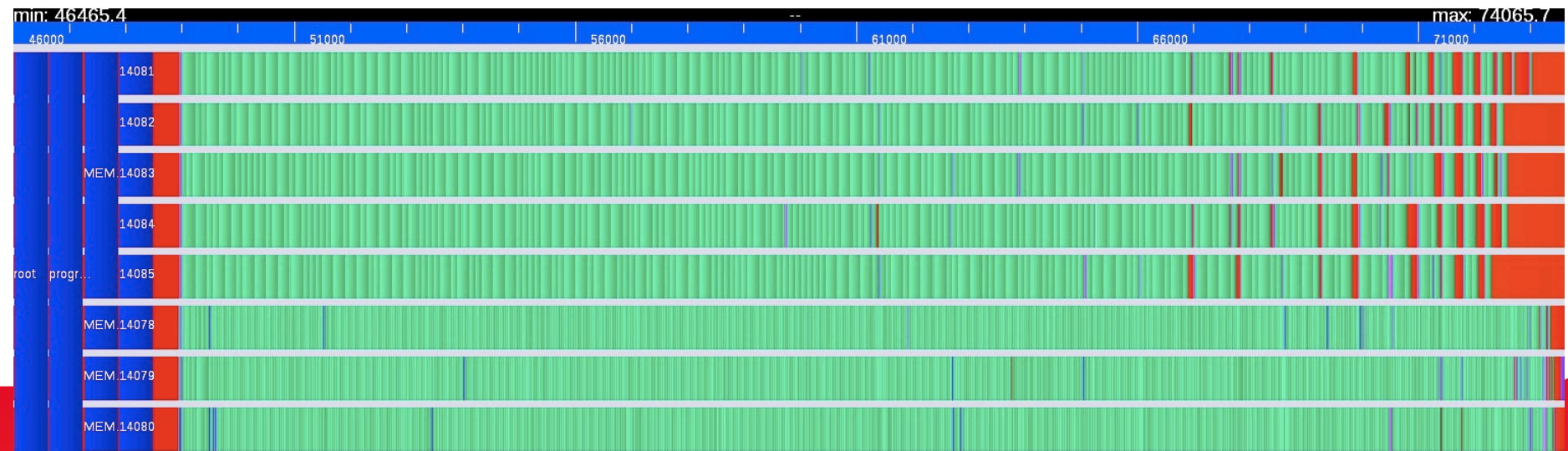
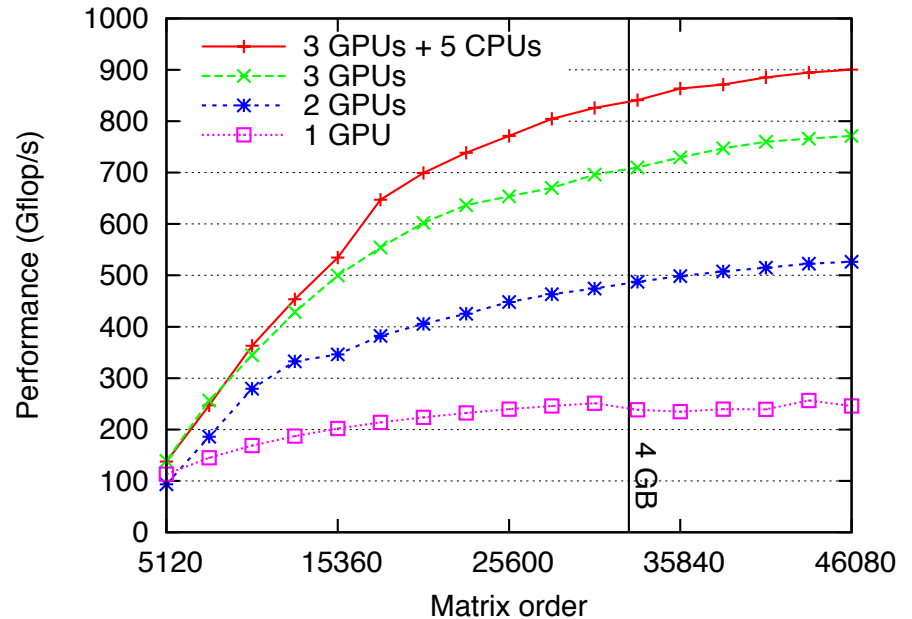
- On the influence of the scheduling policy
 - LU decomposition
 - 8 CPUs (Nehalem) + 3 GPUs (FX5800)
 - 80% of work goes on GPUs, 20% on CPUs
- StarPU exhibits good scalability wrt:
 - Problem size
 - Number of GPUs



Mixing PLASMA and MAGMA with StarPU

With University of Tennessee & INRIA HiePACS

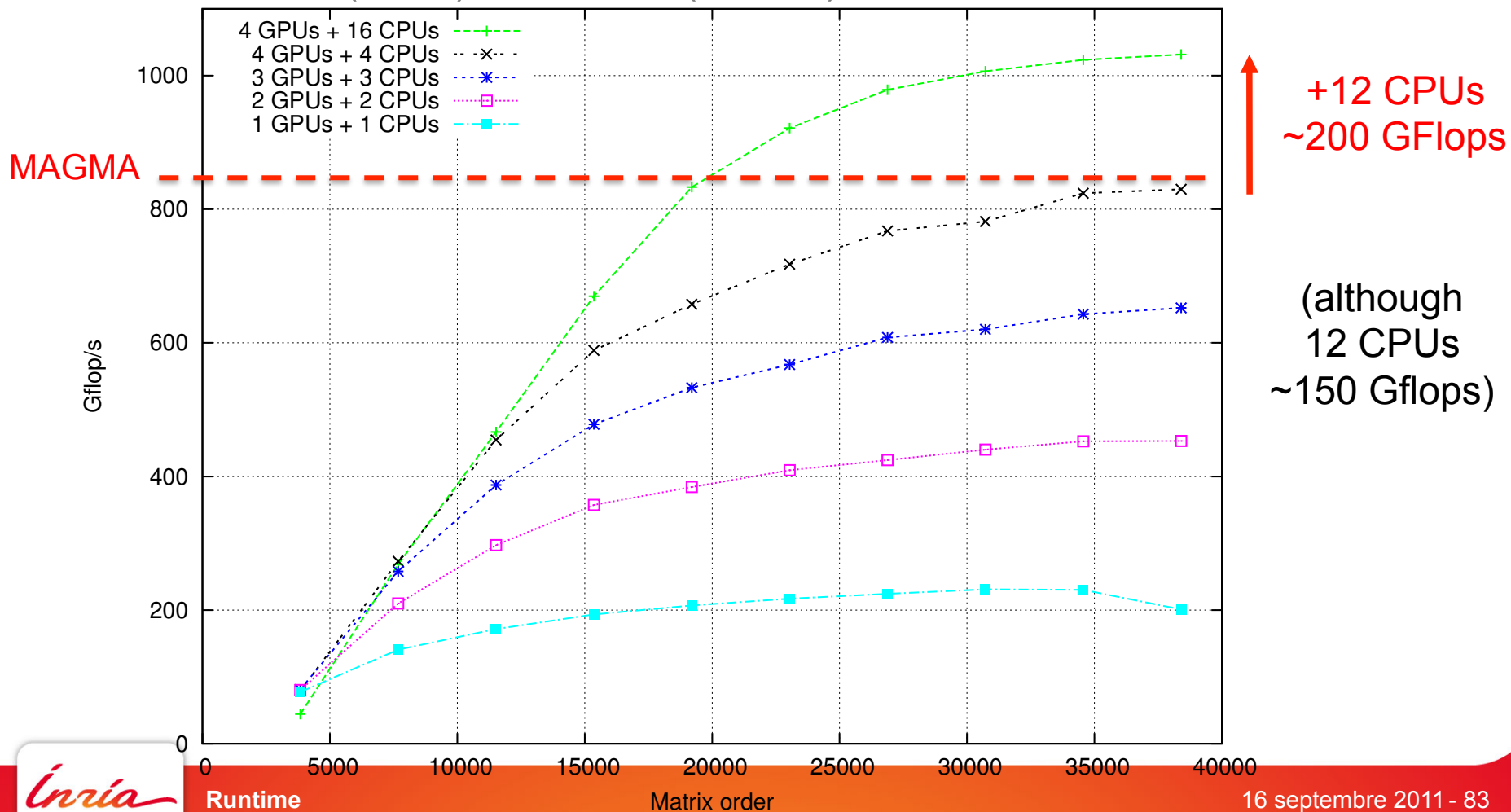
- Cholesky decomposition
 - 5 CPUs (Nehalem) + 3 GPUs (FX5800)
 - Efficiency > 100%



Mixing PLASMA and MAGMA with StarPU

With University of Tennessee & INRIA HiePACS

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



Mixing PLASMA and MAGMA with StarPU

« Super-Linear » efficiency in QR?

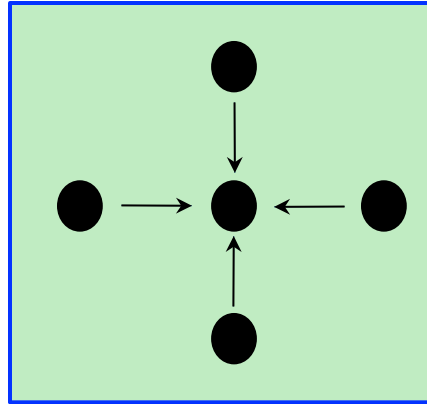
- Kernel efficiency
 - sgeqrt
 - CPU: 9 Gflops GPU: 30 Gflops Ratio: x3
 - stsqrt
 - CPU: 12 Gflops GPU: 37 Gflops Ratio: x3
 - somqr
 - CPU: 8.5 Gflops GPU: 227 Gflops Ratio: x27
 - Sssmqr
 - CPU: 10 Gflops GPU: 285 Gflops Ratio: x28
- Task distribution observed on StarPU
 - sgeqrt: 20% of tasks on GPUs
 - Sssmqr: 92.5% of tasks on GPUs
- Heterogeneous architectures are cool! 😊

About dynamic scheduling...

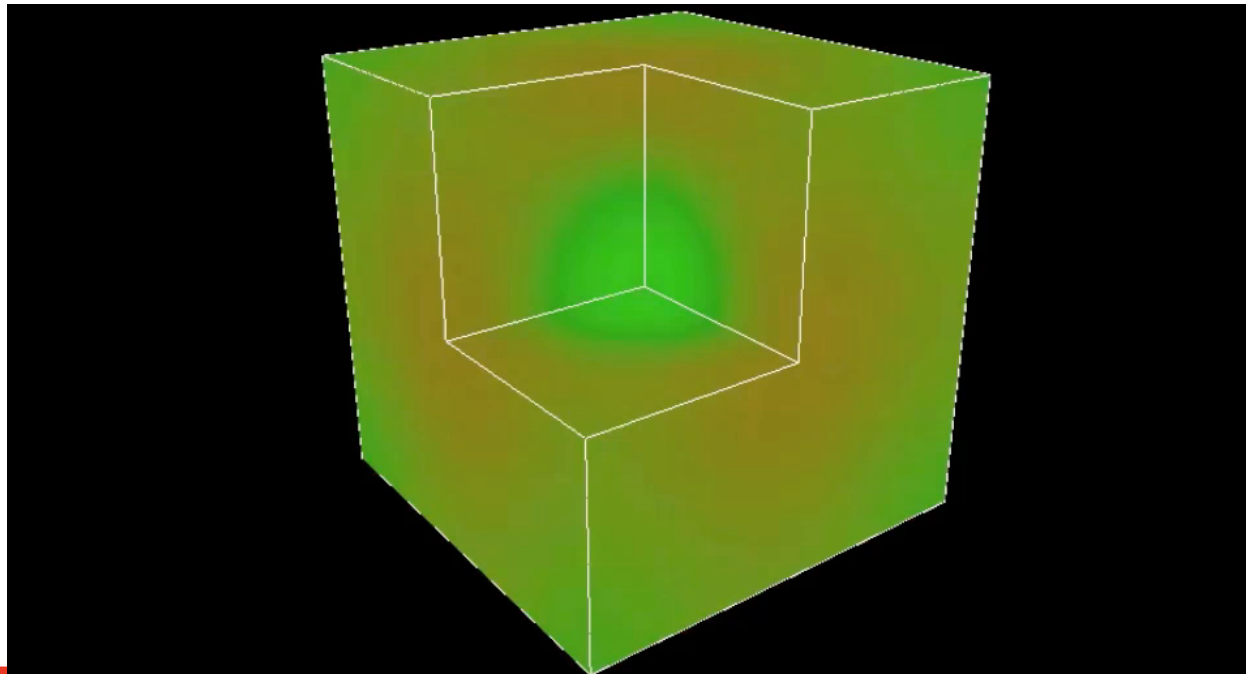
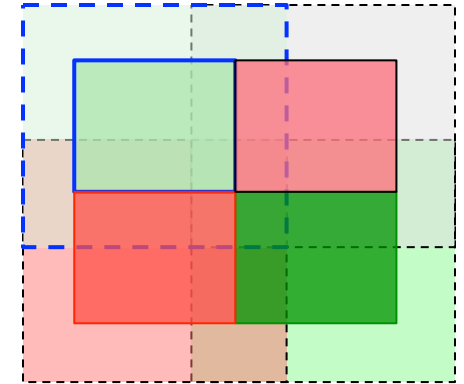


Static vs Dynamic scheduling

- Wave propagation
 - Prefetching
 - Asynchronism



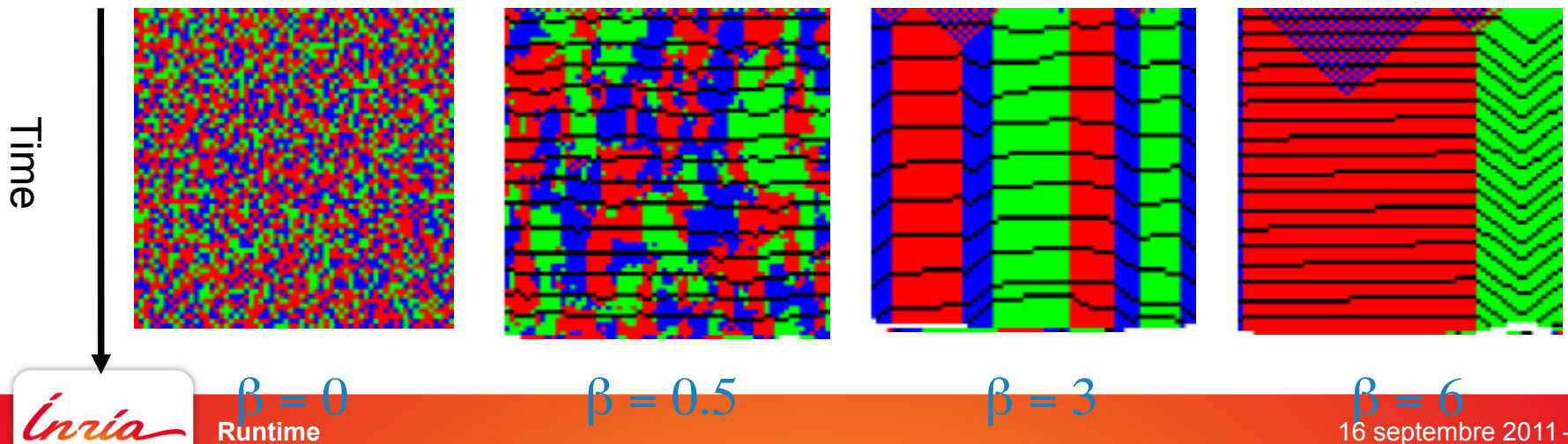
Stencil computation



Static vs Dynamic scheduling

Can a dynamic scheduler compete with a static approach?

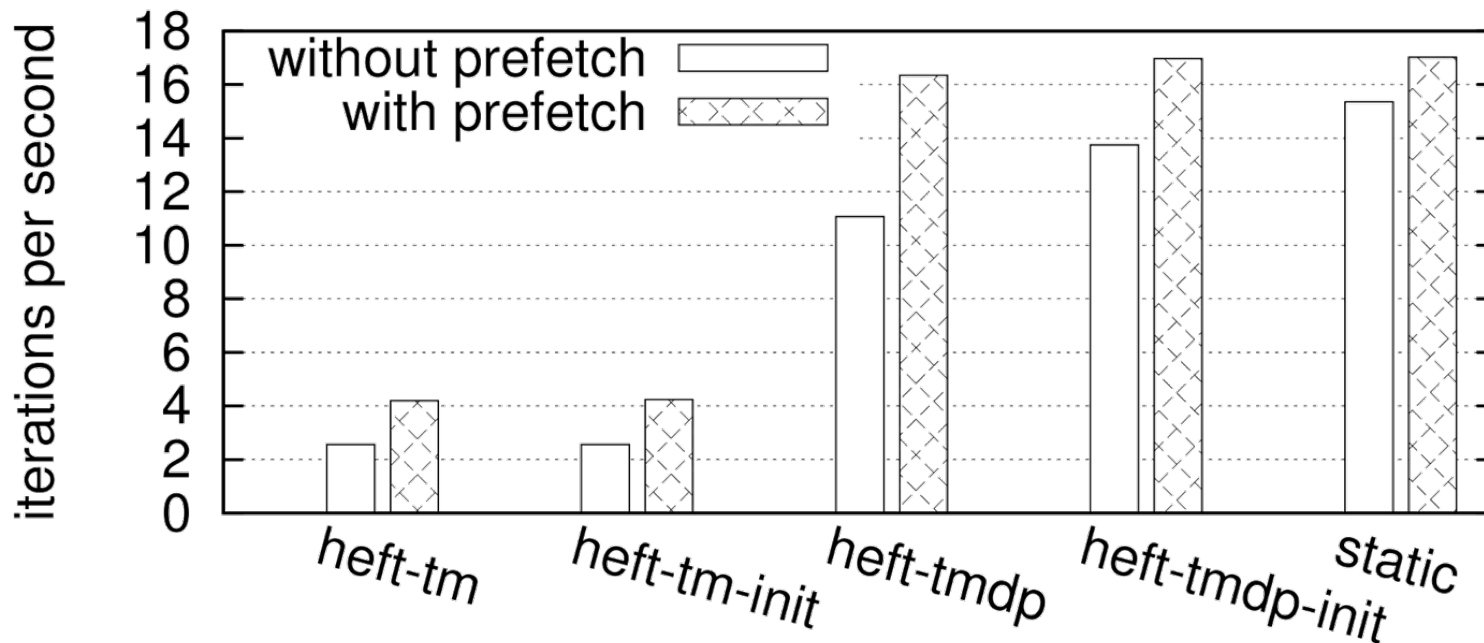
- Load balancing vs data stability
 - We estimate the task cost as $\alpha \text{ compute} + \beta \text{ transfer}$
 - Problem size: 256 x 4096 x 4096, divided into 64 blocks
 - Task distribution (1 color per GPU)
 - Dynamic scheduling can lead to stable configurations



Static vs Dynamic scheduling

Performance

- Impact of scheduling policy
 - 3 GPUs (FX5800) – no CPU used
 - 256 x 4096 x 4096 : 64 blocks
 - Speed up = 2.7 (2 PCI 16x + 1 PCI 8x config)



MPI and StarPU

inria
informatics mathematics

Using MPI and StarPU

Sending and Receiving StarPU data

- Rational: keep an MPI-looking code
 - Work on StarPU data instead of plain data buffers.
- MPI-like primitives to transfer StarPU data
 - `starpu_mpi_send/recv, isend/irecv, ...`
 - Triggers all needed CPU/GPU transfers
 - Respect task/communications dependencies
 - Overlaps MPI communications with *PU communications

```
// unlocks 'tag' upon reception of the message
```

```
starpu_mpi_irecv_detached_unlock_tag (starpu_data_handle data_handle,  
                                       int source, int mpi_tag,  
                                       MPI Comm comm,  
                                       starpu_tag_t tag);
```

Using MPI and StarPU

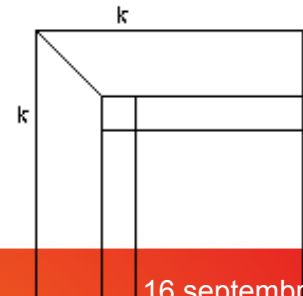
MPI version of `starpu_insert_task`

- Data distribution over MPI nodes decided by application
- Data consistency enforced at the cluster level
 - Automatic `starpu_mpi_send/recv` calls for each task
 - \approx DSM with task-based granularity
- All nodes execute the same algorithm
 - Actual task distribution according to data being written to
 - Owner compute rule
- Sequential-looking code !

MPI version of starpu_insert_task

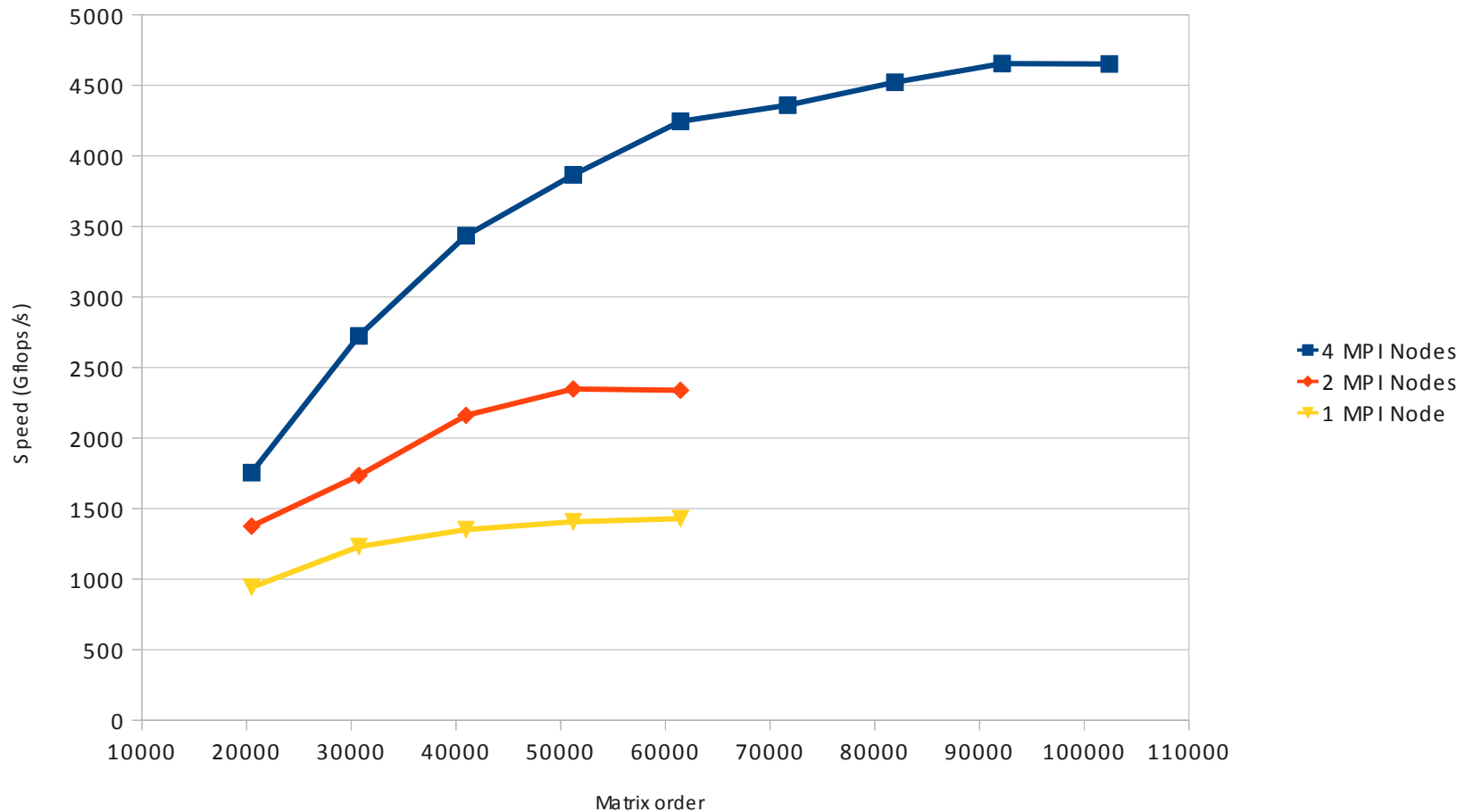
Cholesky decomposition

```
for (k = 0; k < nblocks; k++) {
    starpu_mpi_insert_task(MPI_COMM_WORLD, &c111,
                          STARPU_RW, data_handles[k][k], 0);
    for (j = k+1; j<nblocks; j++) {
        starpu_mpi_insert_task(MPI_COMM_WORLD, &c121,
                              STARPU_R, data_handles[k][k],
                              STARPU_RW, data_handles[k][j], 0);
        for (i = k+1; i<nblocks; i++)
            if (i <= j)
                starpu_mpi_insert_task(MPI_COMM_WORLD, &c122,
                                        STARPU_R, data_handles[k][i],
                                        STARPU_R, data_handles[k][j],
                                        STARPU_RW, data_handles[i][j], 0);
    }
}
starpu_task_wait_for_all();
```



Cholesky Using MPI+StarPU + Magma kernels

Early results



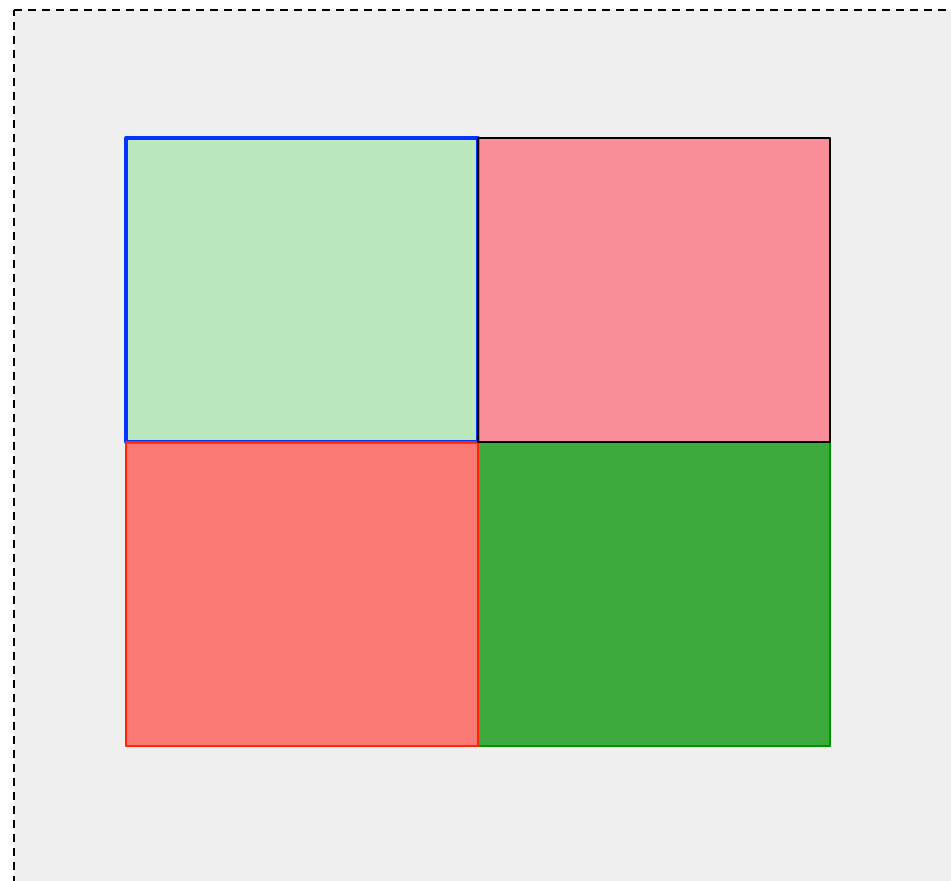
Integrating multithreading and StarPU



Towards parallel tasks on CPUs

Going further

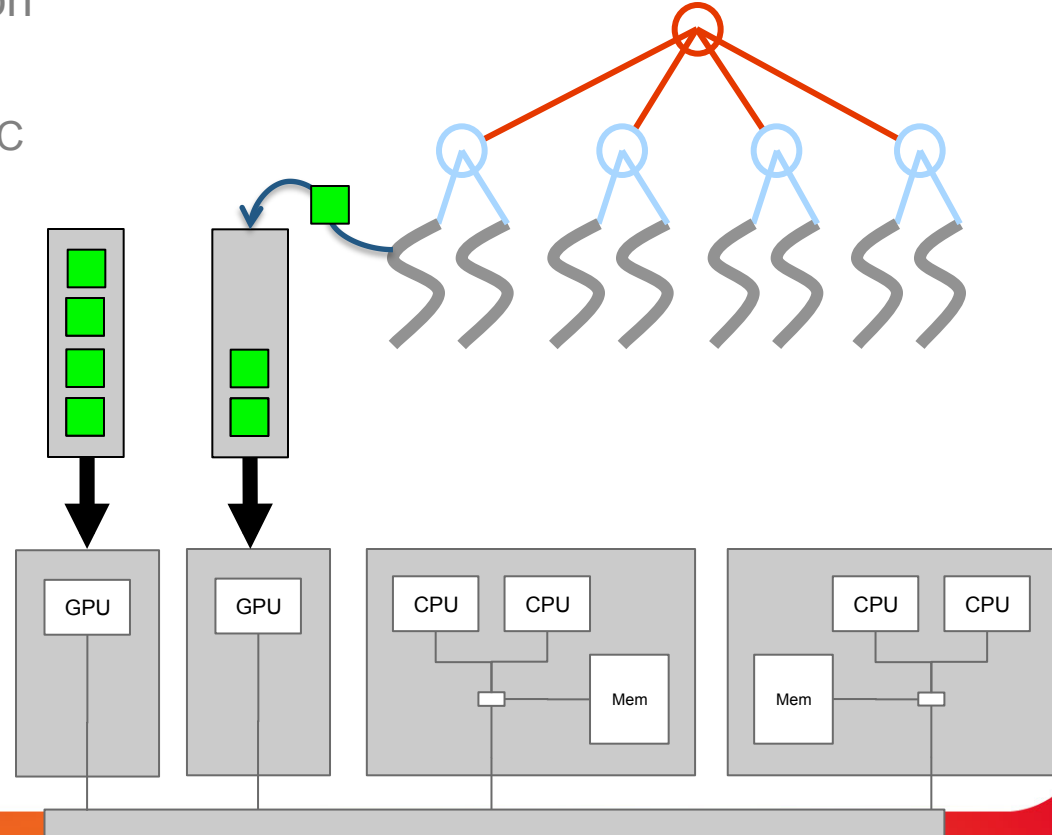
- MPI + StarPU + OpenMP
 - Many algorithms can take advantage of shared memory
 - We can't seriously "*taskify*" the world!
- The Stencil case
 - When neighbor tasks can be scheduled on a single node
 - Just use shared memory!
 - Hence an OpenMP stencil kernel
- We need to mix StarPU with OpenMP/TBB/Phreads/<add your favorite lib here>



Integrating StarPU and Multithreading

Integrating tasks and threads

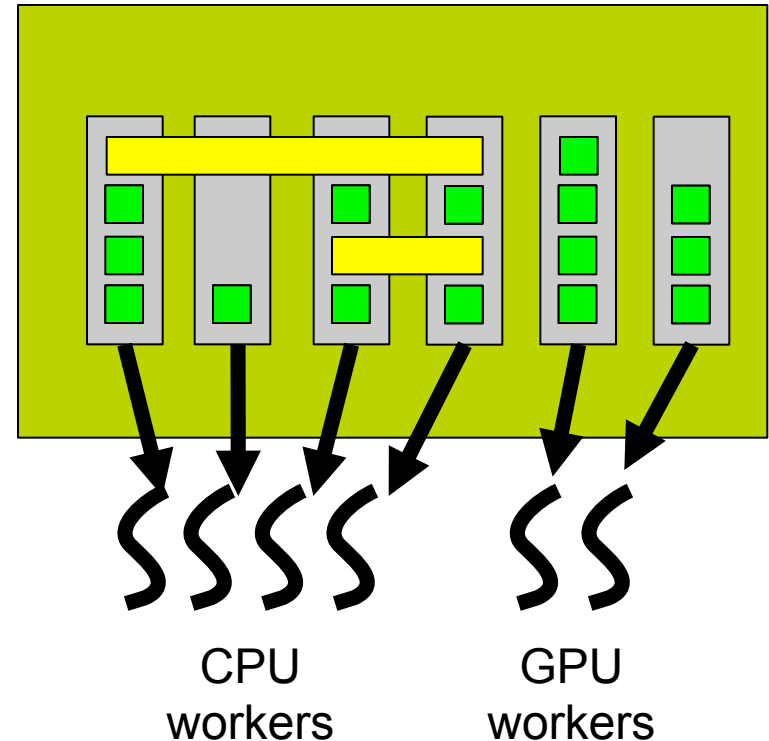
- First approach
 - Use an OpenMP main stream
 - Suggested (?) by recent parallel language extension proposals
 - E.g. Star SuperScalar (UPC Barcelona)
 - HMPP (CAPS Enterprise)
 - Implementing scheduling is difficult
 - Much more than a simple offloading approach...



Integrating StarPU and Multithreading

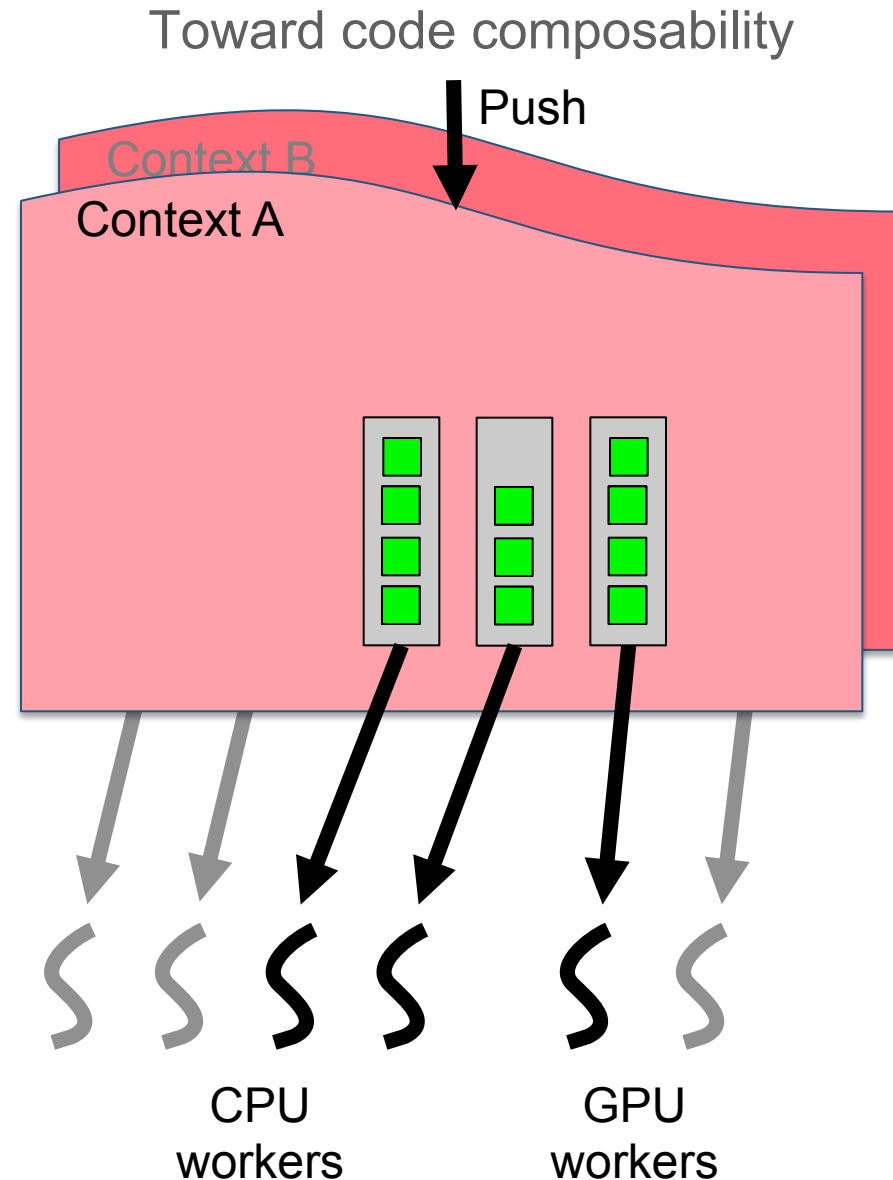
Integrating tasks and threads

- Alternate approach
 - Let StarPU spawn OpenMP tasks
 - Performance modeling would still be valid
 - Would also work with other tools
E.g. Intel TBB
 - How to find the appropriate granularity?
May depend on the concurrent tasks!
 - StarPU tasks = first class citizen
Need to bridge the gap with existing parallel languages



StarPU's Scheduling Contexts

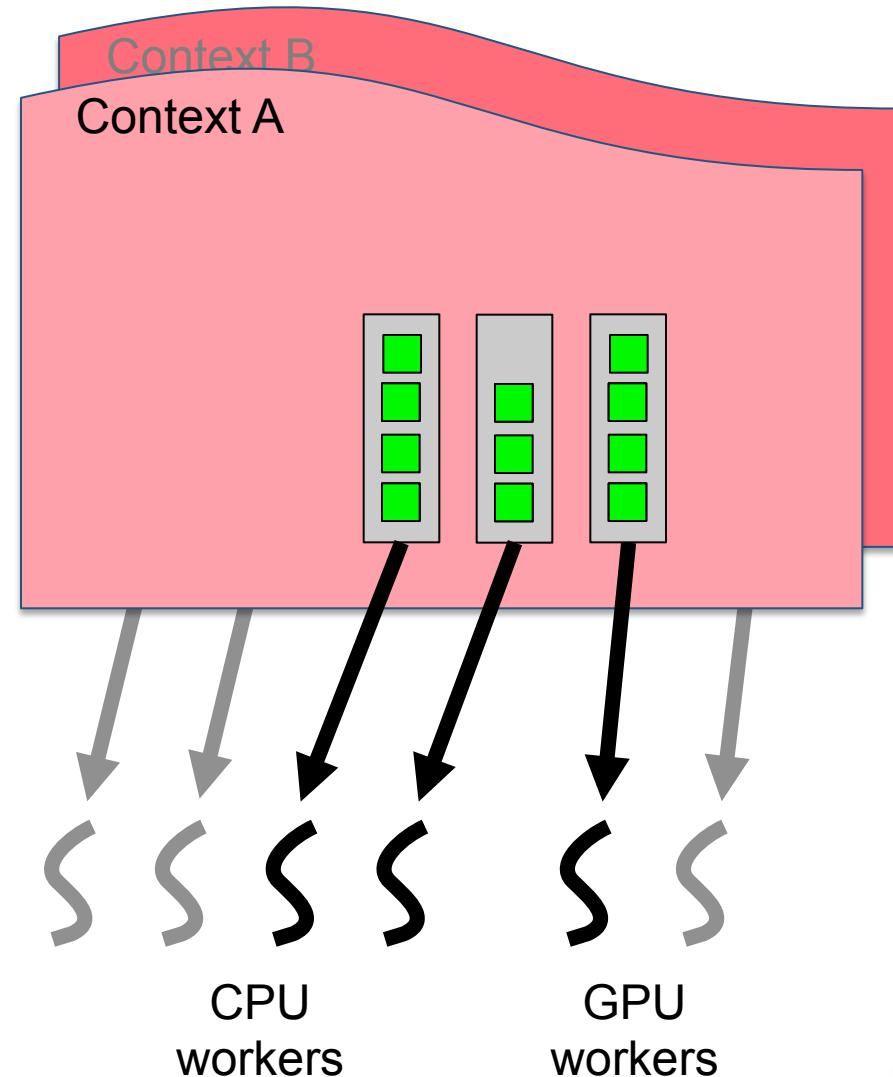
- Similar to OpenCL contexts
 - Except that each context features its own scheduler
- Multiple parallel libraries can run simultaneously
 - Virtualization of resources
 - At minimal overhead
 - Scheduling overhead reduced
 - Scalability workaround



StarPU's Scheduling Contexts

- Contexts may share processing units
 - Avoid underutilized resources
 - Schedulers are aware of each other
- Contexts may expand and shrink
 - Maximize overall throughput
 - Use dynamic feedback both from application and runtime

Toward code composability



What's next?

inria
informatics mathematics

Future parallel machines

Exascale (10^{18} flop/s) systems, by 2018?

- The biggest change comes from node architecture
 - Hybrid systems will be commonplace
 - Multicore chips + accelerators (GPUs?)
 - More integrated design
 - Extreme parallelism
 - Total system concurrency $\sim O(10^9)$!
 - Including $O(10)$ to $O(100)$ to hide latency
- = x 10 000 increase

How will we program these machines?

Let's prepare for serious changes

- Billions of threads will be necessary to occupy exascale machines
 - Exploit every source of (fine-grain) parallelism
 - Not every algorithm/problem resolution can scale that far ☹
 - Multi-scale, Multi-physics applications are welcome!
 - Great opportunity to exploit multiple levels of parallelism
 - Is SIMD the only reasonable approach?
 - Are CUDA & OpenCL our future?
- No global, consistent view of node's state
 - Local algorithms
 - Hierarchical coordination/load balance
- Maybe, this time, we should seriously consider enabling (parallel) code reuse...

Parallel code reuse

Mixing different paradigms leads to several issues

- Can we really use several hybrid parallel kernels simultaneously?
 - Ever tried to mix OpenMP and Intel MKL?
 - Could be helpful in order to exploit millions of cores
- It's all about composability
 - Probably the biggest challenge for runtime systems
 - Hybridization will mostly be indirect (linking libraries)
- And with composability come a lot of related issues
 - Need for autotuning / scheduling hints

International Exascale Software Project (IESP)

“A call to action”

- Build an international plan for coordinating research for the next generation open source software for scientific high-performance computing
 - Hardware is evolving more rapidly than software
 - New hardware trends not handled by existing software
 - Emerging software technologies not yet integrated into a common software stack
 - No global evaluation of key missing components



European Exascale Software Initiative (EESI)

Position of Europe in the international HPC landscape

- WP4: Enabling technologies for Exascale computing
 - Assess novel HW and SW technologies for Exascale challenges
 - Build a European vision and a roadmap
- WG 4.2: Software eco-systems
 - Subtopic: Runtime systems (Raymond Namyst, Jesús Labarta)

European Exascale Software Initiative (EESI)

Runtime systems: Scientific and Technical Hurdles

- Mastering heterogeneity
 - Unified/transparent accelerator models
 - Providing support for adaptive granularity
 - Fine grain parallelism
 - Scheduling for latency/bandwidth
 - (NC)-NUMA
- Supporting multiple programming models
 - Hybrid runtimes
 - MPI + threading model + accelerator model
 - Matching hybrid parallelism on heterogeneous architectures
 - Tuning the number of (processes/threads) per level

European Exascale Software Initiative (EESI)

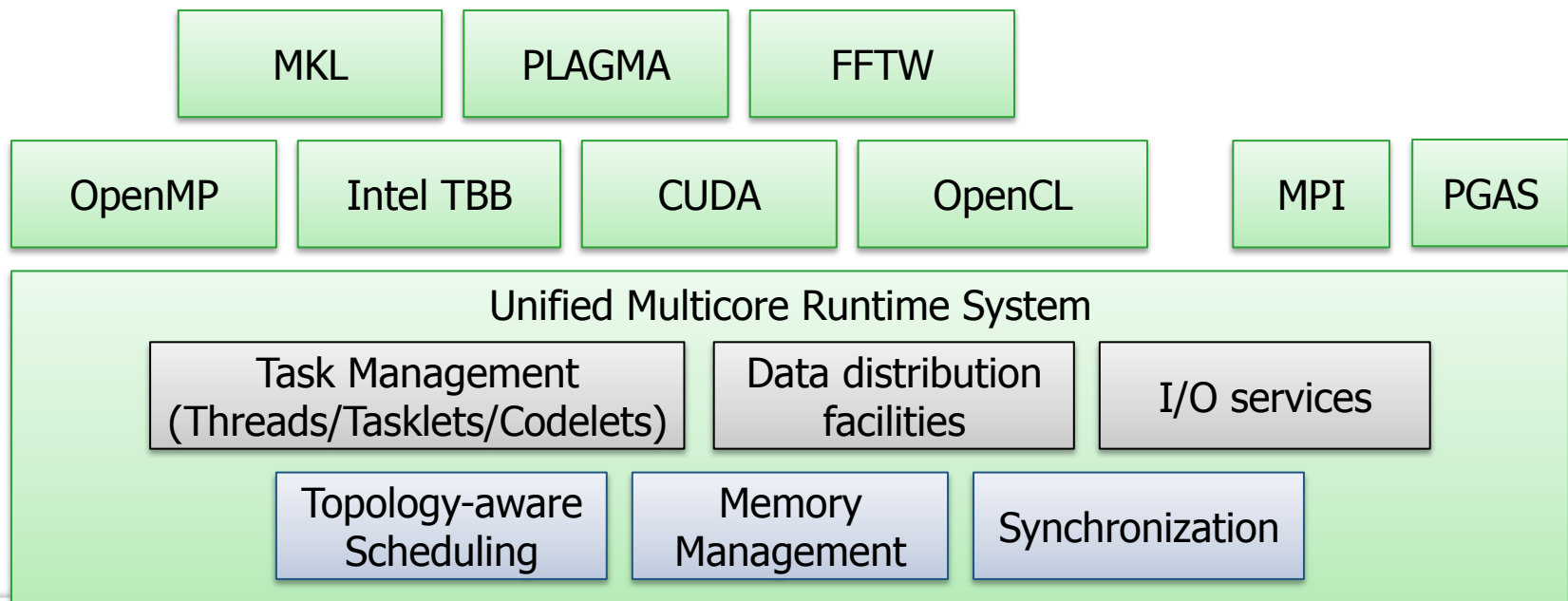
Runtime systems: Scientific and Technical Hurdles

- Dealing with millions of cores/nodes
 - Scheduling
 - Hierarchical scheduling
 - Data-flow task bases approaches
 - Non-coherent architectures
 - Software data prefetching
 - Imbalance detection, prediction and (local) correction
 - Avoid global balancing strategies
 - Work stealing
 - Communication
 - Scalable implementations of MPI/PGAS
 - Minimize memory consumption (per connection)
 - Redesign of collective operations
 - Asynchrony, overlap
 - Discourage use of global synchronization primitives?

Toward a common runtime system?

i.e. Unified Software Stack

- There is currently no consensus on a common runtime system
 - One objective of the Exascale Software Center
 - Coordinated exascale software stack
 - Technically feasible...



Major Challenges are ahead...

We are living exciting times!



more information:

<http://runtime.bordeaux.inria.fr>