



CSCS

Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

ETH zürich



DLAF - Distributed Linear Algebra with Futures

real-world use-cases

Advanced C++

Alberto Invernizzi (alberto.invernizzi@cscs.ch), CSCS

October 12, 2021

HPC and Linear Algebra

The de-facto standard library for distributed linear algebra is ScaLAPACK, a library that has been developed in 1995, when supercomputers were based on nodes which had a single CPU core.

Since then, node architectures have evolved (e.g. multicores and multi-GPU), and new programming paradigms are necessary.

Distributed Linear Algebra with Futures - DLAF

DLAF uses a **task-based approach** aiming to reduce the amount of synchronization necessary between distinct routines.

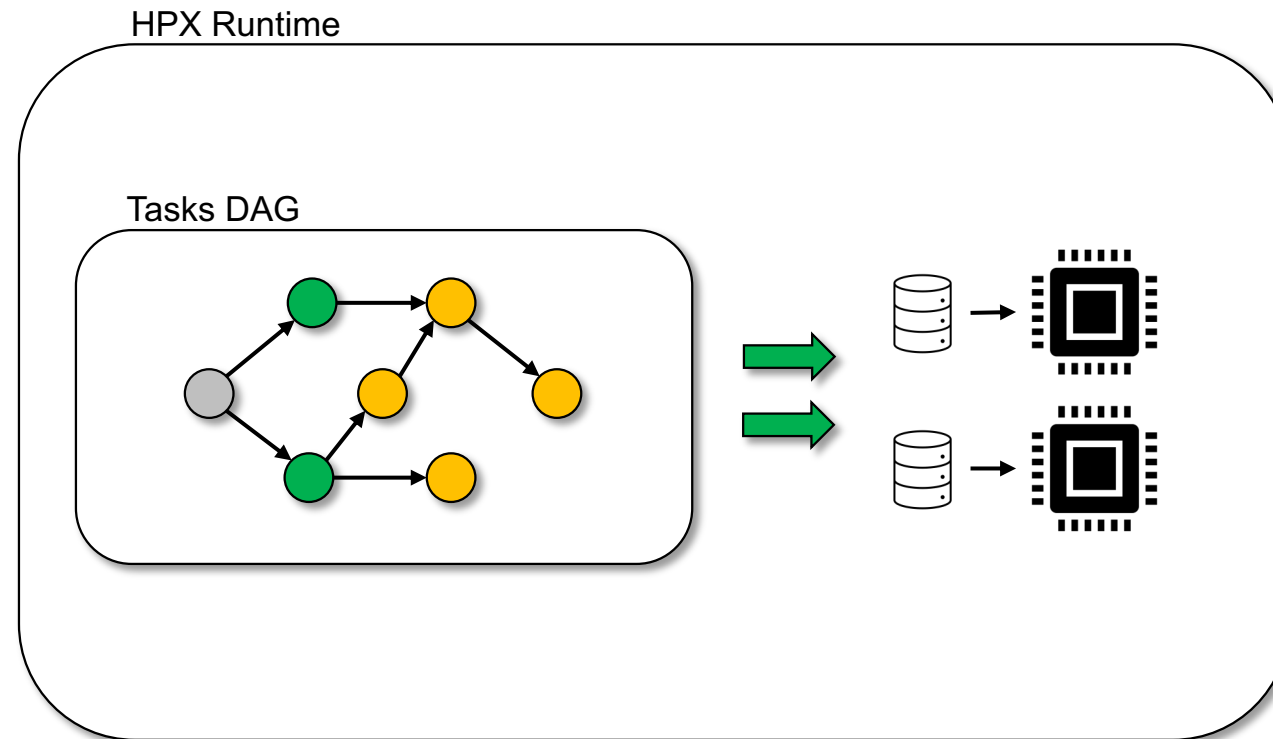
The task-based approach is possible thanks to HPX, a dependency of DLAF for which represents the backbone.

What is *HPX*?

HPX is a C++ Standard Library for **Concurrency and Parallelism**. It implements all of the corresponding facilities as defined by the C++ Standard. Additionally, in *HPX* we implement functionalities proposed as part of the ongoing C++ standardization process. We also extend the C++ Standard APIs to the distributed case. *HPX* is developed by the STE||AR group (see [People](#)).

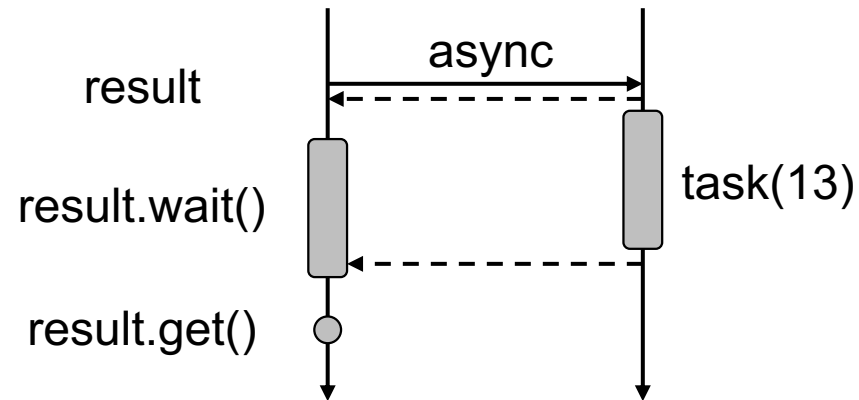
<https://github.com/STELLAR-GROUP/hpx>

HPX overview

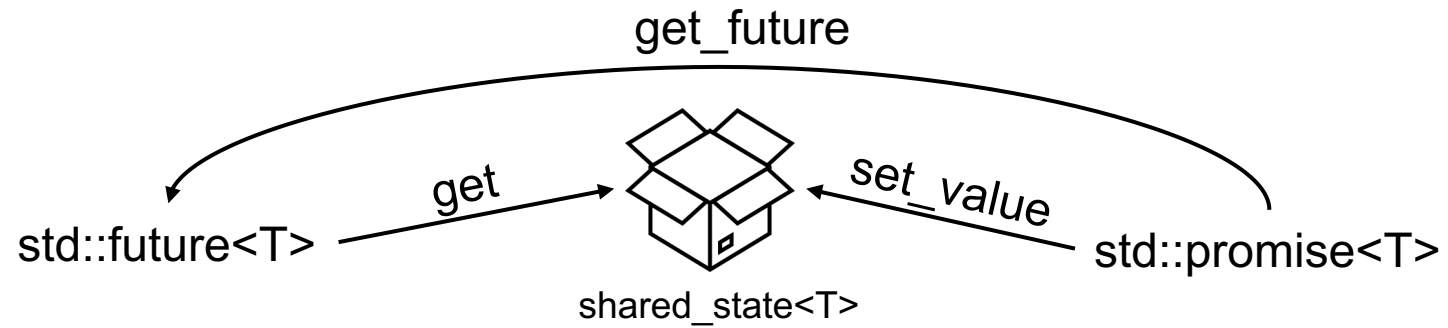


Task: a trivial example

```
float task(int input);  
  
// ...  
  
std::future<float> result = std::async(task, 26);  
  
result.wait();  
float output = result.get();
```

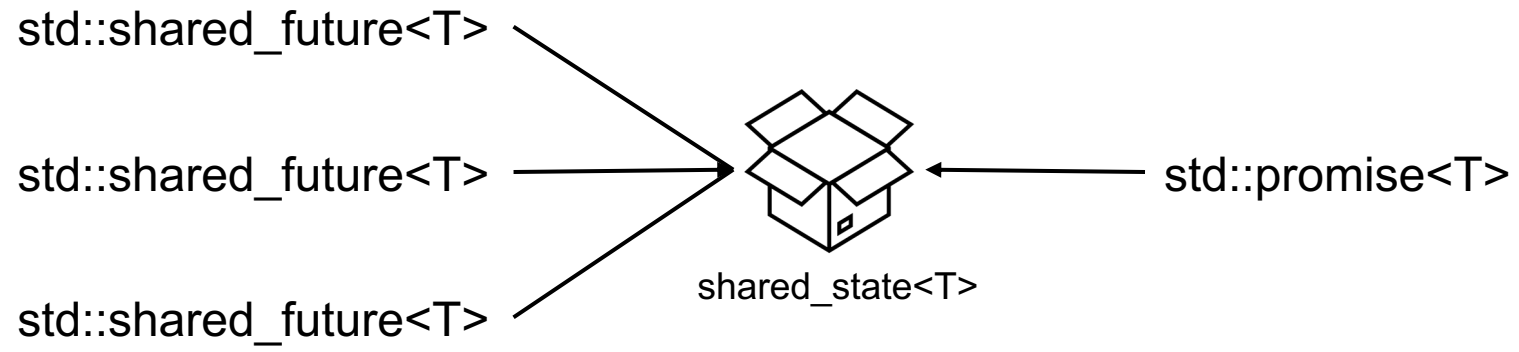


std::promise + std::future



```
1 #include <future>
2
3 int main() {
4     std::promise<int> promise;
5     std::future<int> future = promise.get_future();
6
7     promise.set_value(26);
8
9     int value = future.get();
10 }
```

... and `std::shared_future`




```
#include <future>

float task(int value) {
    return value * 1.3;
}

int main() {
    std::future<float> result = std::async(task, 26);

    result.wait();

    assert(result.valid());
    float output = result.get();
    assert(not result.valid());
}
```

```
#include <future>

float task(int value) {
    return value * 2.33;
}

int main() {
    std::shared_future<float> result = std::async(task, 26);

    result.wait();

    assert(result.valid());
    const float& output = result.get();
    assert(result.valid());
}
```



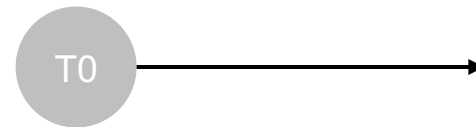
CSCS

Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

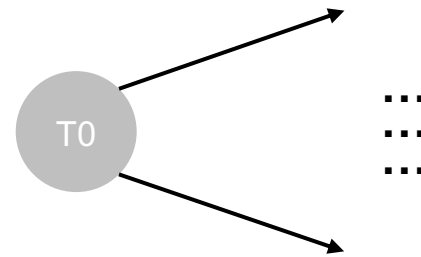
ETH zürich

HPX 101

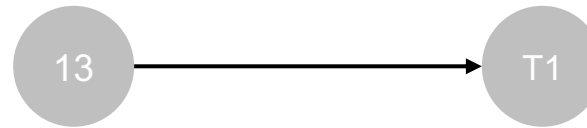
DAG: future



DAG: shared_future



.then

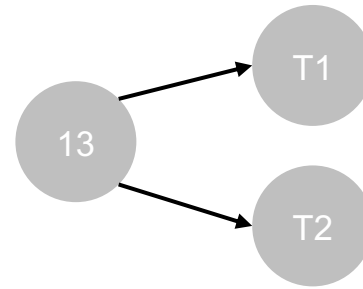


```
1 #include <hpx/future.hpp>
2 #include <hpx/hpx_main.hpp>
3
4 int main() {
5     using namespace hpx;
6
7     future<int> future_value = make_ready_future<int>(13);
8
9     future_value.then([](future<int> value) {});
10 }
```

It returns another future

Inside here the value is ready to be retrieved

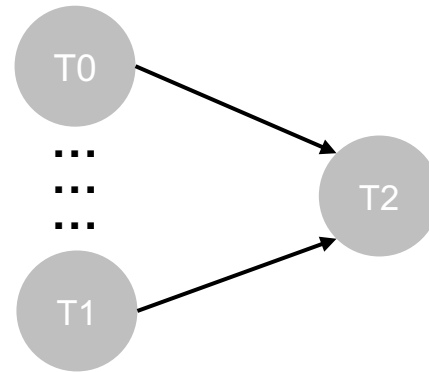
.then



```
1 #include <hpx/future.hpp>
2 #include <hpx/hpx_main.hpp>
3
4 int main() {
5     using namespace hpx;
6
7     shared_future<int> future_value = make_ready_future<int>(13);
8
9     future_value.then([](shared_future<int> value) {});
10    future_value.then([](shared_future<int> value) {});
11 }
```

When will these
be executed?

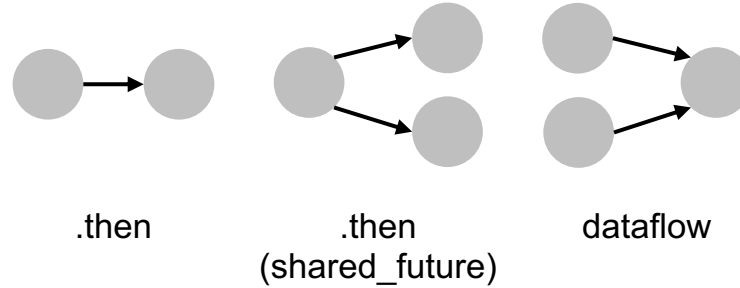
dataflow



```
1 hpx::future<int>    task0_result = hpx::async([]() { return 26; });
2 hpx::future<float>  task1_result = hpx::async([]() { return 13.26f; });
3
4 hpx::dataflow(
5     [](hpx::future<int> data0, hpx::future<> data1) { ... },
6     task0_result, task1_result);
```

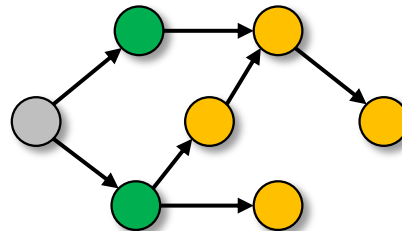
DAG Grammar

These allow you to describe any DAG



!!! IT'S NOT ABOUT FORK/JOIN !!!

these operations describe dependencies between tasks



DLAF

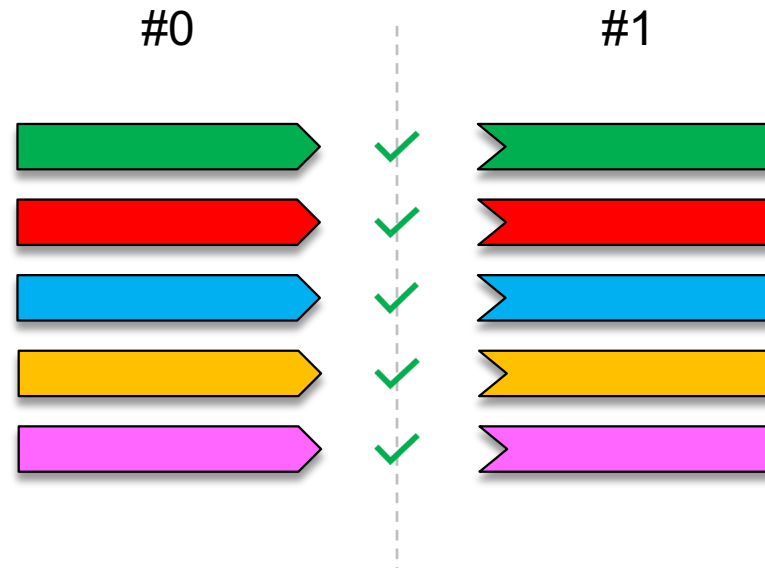
MPI Communication

The D in DLAF stands for DISTRIBUTED, which implies communication between nodes.

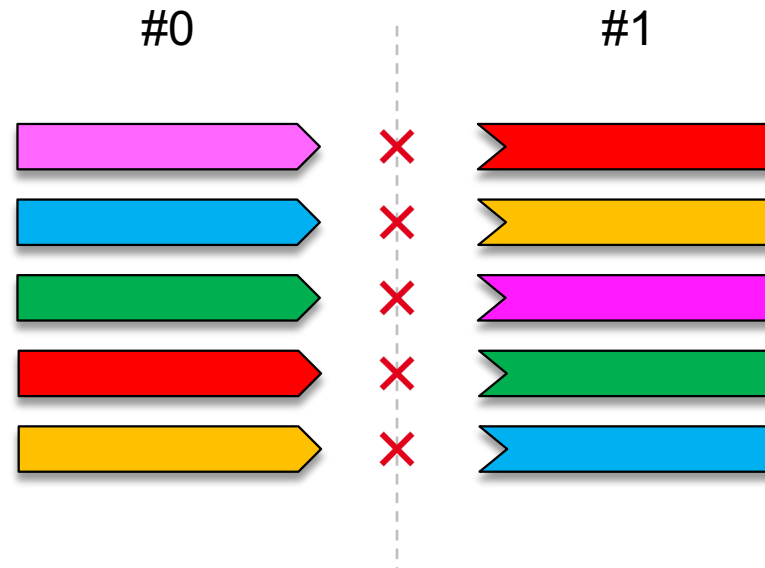
For the communication, our requirement is to use MPI, because DLAF aims at replacing ScaLAPACK in existing Fortran/C/C++ applications.

But using MPI poses a few challenges, especially in a task-based programming model...

Problem: Communication Order



Problem: Communication Order



Pipeline

What about adding a fictitious dependency that describe the ordering requirement?

What do they share?



the data/tile communicated is likely not shared



the communicator does!

Actually, we want to serialize the communications inside the same communicator, because MPI already preserves the independency between different communicators.

So, if we “serialize” the access to the communicator resource, we can impose an order on them!

just movable, enforce
unique ownership

object and promise
ownership

```
1 template <class T> class PromiseGuard {
2 public:
3     PromiseGuard(T object, hpx::lcos::local::promise<T> next)
4         : object_(std::move(object)), promise_(std::move(next)) {}
5
6     PromiseGuard(PromiseGuard &&) = default;
7     PromiseGuard &operator=(PromiseGuard &&) = default;
8
9     PromiseGuard(const PromiseGuard &) = delete;
10    PromiseGuard &operator=(const PromiseGuard &) = delete;
11
12    ~PromiseGuard() {
13        if (promise_.valid())
14            promise_.set_value(std::move(object_));
15    }
16
17    T &ref() { return object_; }
18    const T &ref() const { return object_; };
19
20 private:
21     T object_;
22     hpx::lcos::local::promise<T> promise_;
23 };
```

On destruction it sets the promise
by moving in the value owned

Access to the object via ref()

Get the ownership of the resource

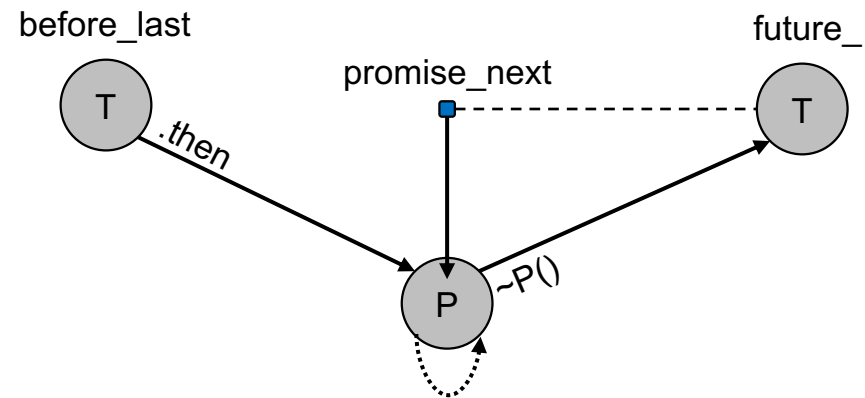
The first element in the pipeline is ready

```
1 template <class T> class Pipeline {
2 public:
3     Pipeline(T object) : future_(hpx::make_ready_future(std::move(object))) {}
4
5     hpx::future<PromiseGuard<T>> operator>()() {
6         auto before_last = std::move(future_);
7
8         hpx::lcos::local::promise<T> promise_next;
9         future_ = promise_next.get_future();
10
11         return before_last.then(
12             hpx::launch::sync,
13             hpx::unwrapping(
14                 [promise_next = std::move(promise_next)](T &&object) mutable {
15                     return PromiseGuard<T>{std::move(object), std::move(promise_next)};
16                 }));
17     }
18
19 private:
20     hpx::future<T> future_;
21 };
```

class invariant:
this stores the last element of the pipeline

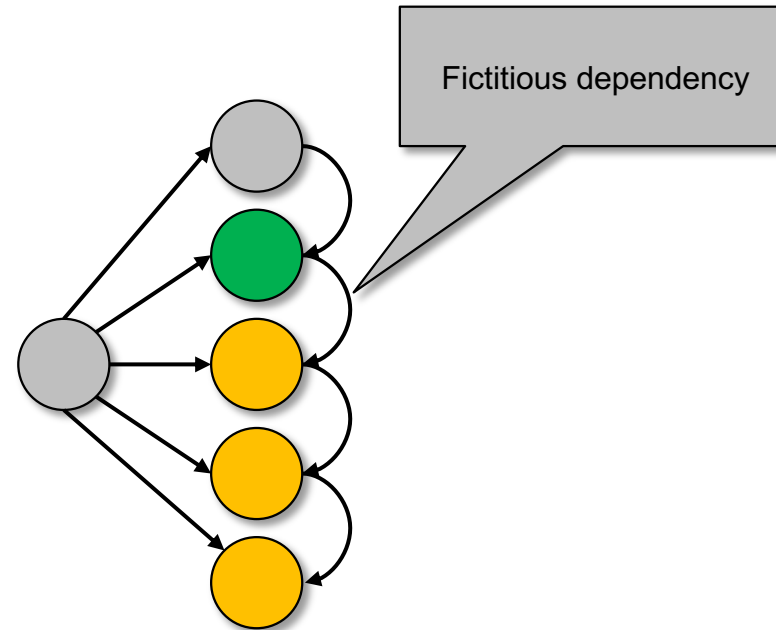
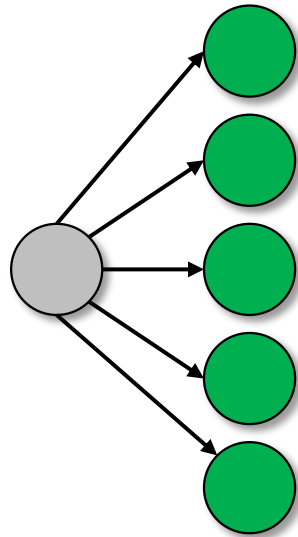


```
1 hpx::future<PromiseGuard<T>> operator()() {  
2     auto before_last = std::move(future_);  
3  
4     hpx::lcos::local::promise<T> promise_next;  
5     future_ = promise_next.get_future();  
6  
7     return before_last.then(  
8         hpx::launch::sync,  
9         hpx::unwrapping(  
10             [promise_next = std::move(promise_next)](T &&object) mutable {  
11                 return PromiseGuard<T>{std::move(object), std::move(promise_next)};  
12             }));  
13 }
```



Note

Pipeline is just a partial solution, since communications that can in theory run in parallel, are forced to run in a specific order, and this may cause performance degradation.



Parallelism

Exploiting new node architectures, which are capable of running multiple operations in parallel, requires a finer grained parallelism.

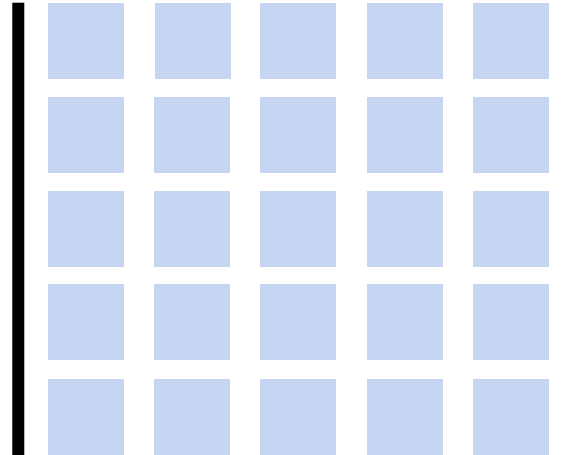
For what concerns linear algebra this is possible thanks to the “tiled” version of algorithms, i.e. divide the problem in smaller blocks.

These blocks can then be run in parallel, respecting the dependencies defined. It sounds like a good use-case for HPX and task-based programming, and it is...

That's what we are doing with DLAF
DATA-DRIVEN TASK BASED IMPLEMENTATION

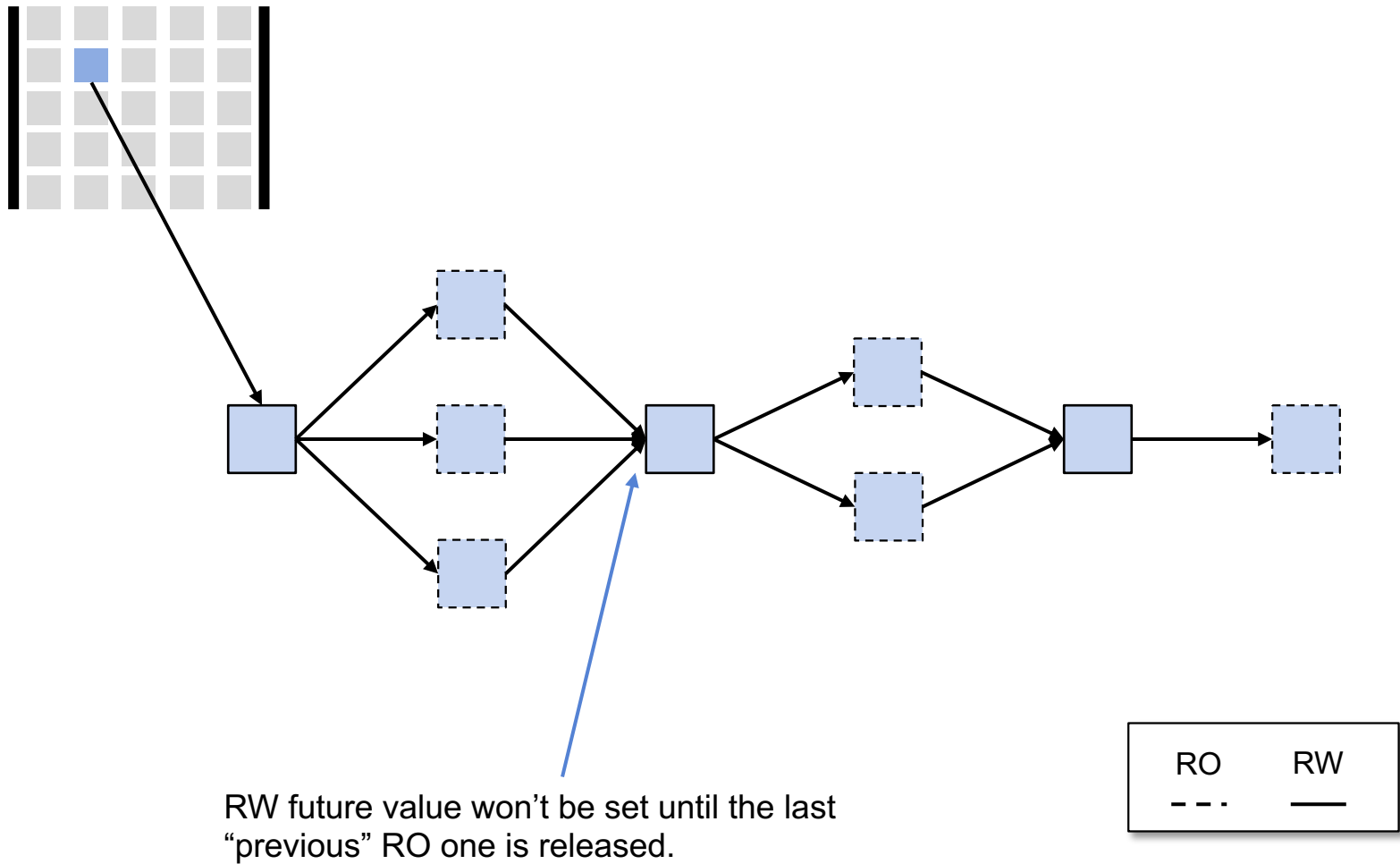
The Matrix

In a linear algebra library, the matrix plays a central role by definition, since it is the data structure on which algorithms runs.

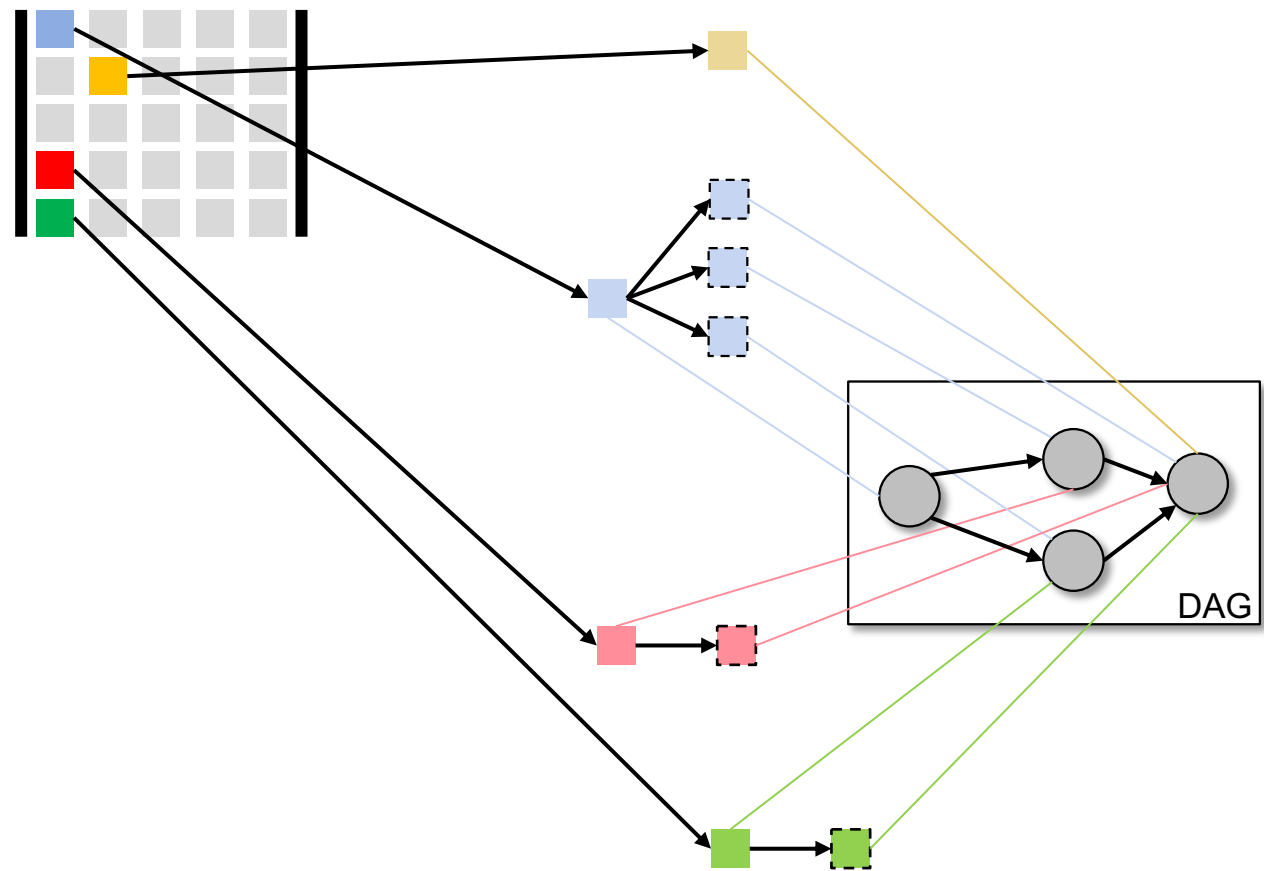


Algorithms are defined in terms of small tasks that run on tiles.

Tile as orchestrator of tasks



Tile as orchestrator of tasks



What's next?

The future of futures: sender-receiver

- `std::future` is largely considered a mistake; it requires heap allocation, does not support continuations, and is largely tied to `std::thread`
 - `hpx::future` inherits many of the downsides of `std::future`
- senders and receivers are a generalization of futures and promises based on concepts, proposed for standardization
- concept-based design allows different asynchronous libraries to interoperate; aims to be for asynchrony what iterators are for iteration
- allows for low-level customization based on the execution context, predecessor, and successor operations

*Thanks to Mikael Simberg (github @msimberg)
HPX Maintainer*

For more information and insights, you can:

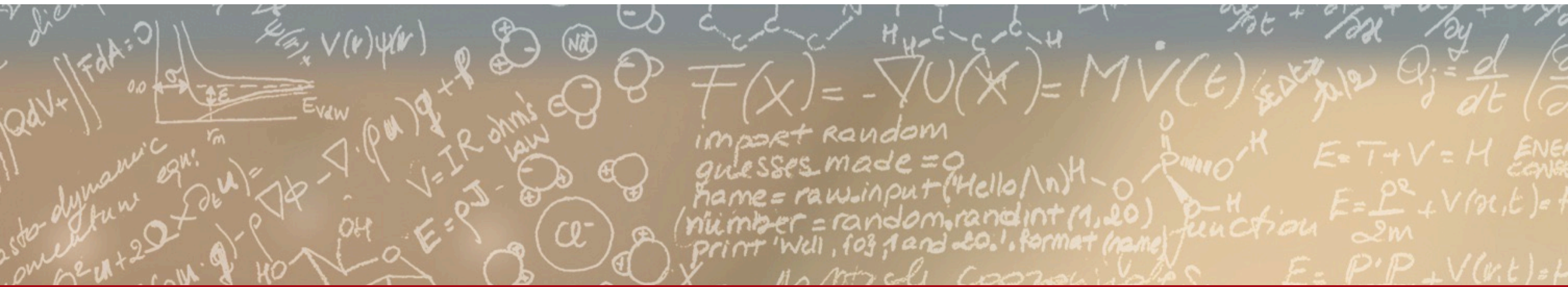
- ✓ see the details in the proposal at wg21.link/p2300
- ✓ look at the presentation given by E.Niebler @ CppCon19
<https://www.youtube.com/watch?v=tF-Nz4aRWAM&t=765s>



CSCS

Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

ETH zürich



Thank you for your attention.