

Tackling UQ in DARMA, a Programming Model for Task-Based Execution at Extreme-Scale

F. Rizzi , E. Phipps, D. Hollman, J. Lifflander, J. Wilke, A. Markosyan, H. Kolla,
N. Slattengren, K. Teranishi, J. Stewart, R. Clay and J. Bennett

Sandia National Laboratories

QUIET17 - SISSA - Italy

Follow on your device: <http://fnrizzi.github.io/quiet17>



Motivation

fnrizzi.github.io/quiet17

1 exaFlops: (10e18) calculations per second, supposedly arriving by 2023-2024

As of June 2017:

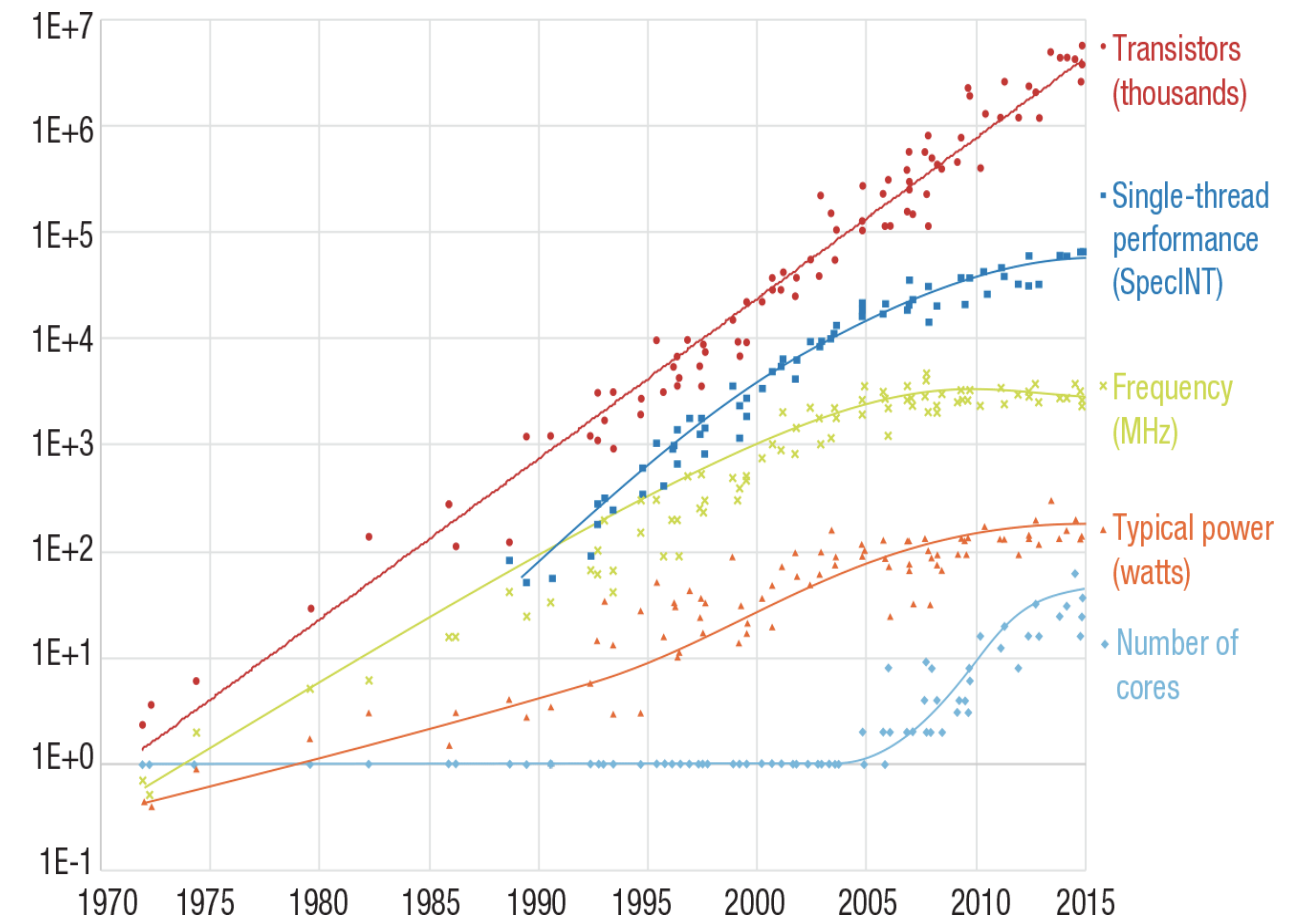
1. Sunway TaihuLight (China): 10,649,600 cores -- 125 PFlops -- 15.4 MW
2. Tianhe-2 (MilkyWay-2) (China): 3,120,000 cores -- 54 PFlops -- 17.8 MW
3. Piz Daint (Switzerland): 361,760 cores -- 25 PFlops -- 2.2 MW
4. Titan (USA): 560,640 cores -- 27 PFlops -- 8.2 MW

Challenges:

- Power consumption
- Complex architectures: heterogeneity
- Unpredictable machines: resilience
- Managing communication/computation
- Increasingly more dynamic workloads and machine performance

Can we ride the wave of current technology?

- Moore's observation: number of transistors on a chip doubles (nearly) every two years
- Dennard scaling: power density remains constant as transistors get smaller
- Dennard scaling broke down ~2005-2007
- Moore's trend is alive and well
- The ability to drop the voltage and current needed to operate reliably has broken, NOT the ability to make smaller transistors.
- Clock speeds are plateauing due to power and thermal limitations





When you can't build outward any longer, build upward!

Clock frequency stalled, so performance growth is/will be achieved by exponential growth in the number of processing elements per chip and hardware threading per core.

Increasing number of cores/chip: expected to double every 18/24 months.

Need for new programming abstractions that:

- virtualize the notion of a core
- threading APIs with expanded semantics for thread control, placement, launching
- synchronization as well as scalable runtimes to manage massive numbers of threads

Locality: Systems bound by communication infrastructure and power dissipation. Management of data locality is a first order concern. Move computation to data, not viceversa.

Heterogeneity: accelerators, implicit data movement, memory hierarchies.

Asynchrony: SPMD/bulk-synchronous programming models presume homogeneous performance across massively parallel systems. Numerous sources of performance inhomogeneity are emerging that will challenge this.

Fault Tolerance: larger, more complex machines. Hundreds of millions of cores, circuits with feature sizes as small as 7 nm, and lower voltages than today. Soft/hard errors more likely.

DARMA

fnrizzi.github.io/quiet17

DARMA: Distributed Asynchronous Resilient Models and Applications

C++ abstraction layer for asynchronous many-task (AMT) runtimes

Provides a set of abstractions to facilitate the expression of tasking that map to a variety of underlying AMT runtime system technologies.

Goals:

- Enables exploration of a variety of underlying runtime system technologies without changing application code.
- Facilitate the expression of coarse-grained tasking.

Decompose into small, transferable units of work (many tasks) with associated inputs (dependencies) rather than simply decomposing at the process level (MPI ranks).

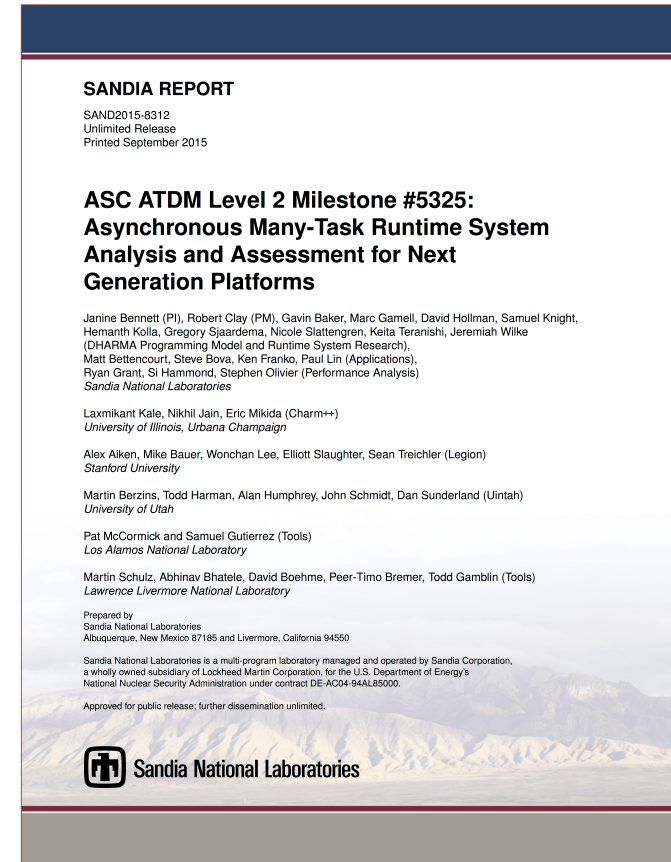
Rather than in a well-defined order, tasks execute when inputs become available.

Is the coarse-grained, distributed memory analog of instruction-level parallelism, extending the concepts of data pre-fetching, out-of-order task execution based on dependency analysis, and even branch prediction (speculative execution).

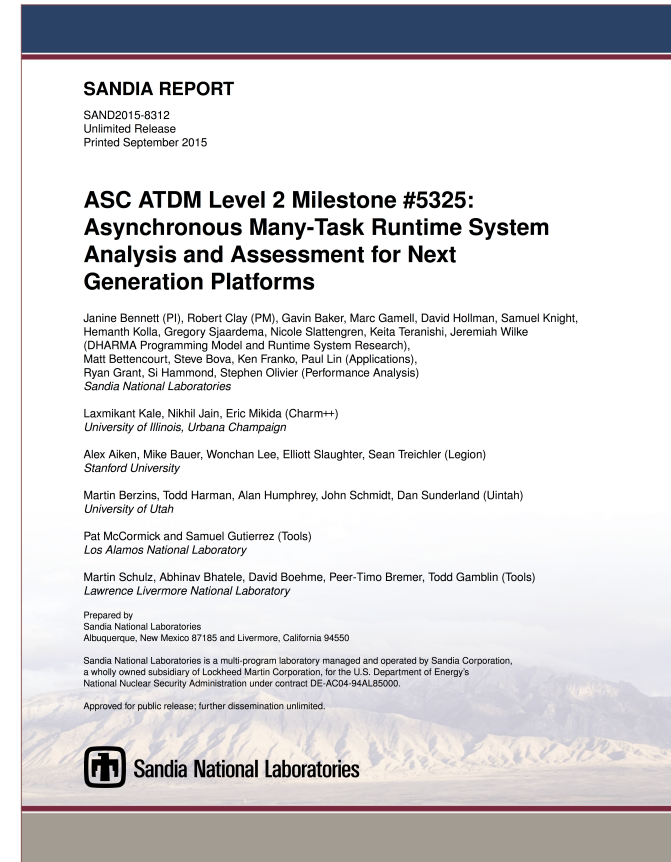
An AMT model aims to leverage all available task and pipeline parallelism, rather just relying on basic data parallelism for concurrency.

Enables communication/computation overlap, as well as reasoning about asynchronous load balancing strategies.

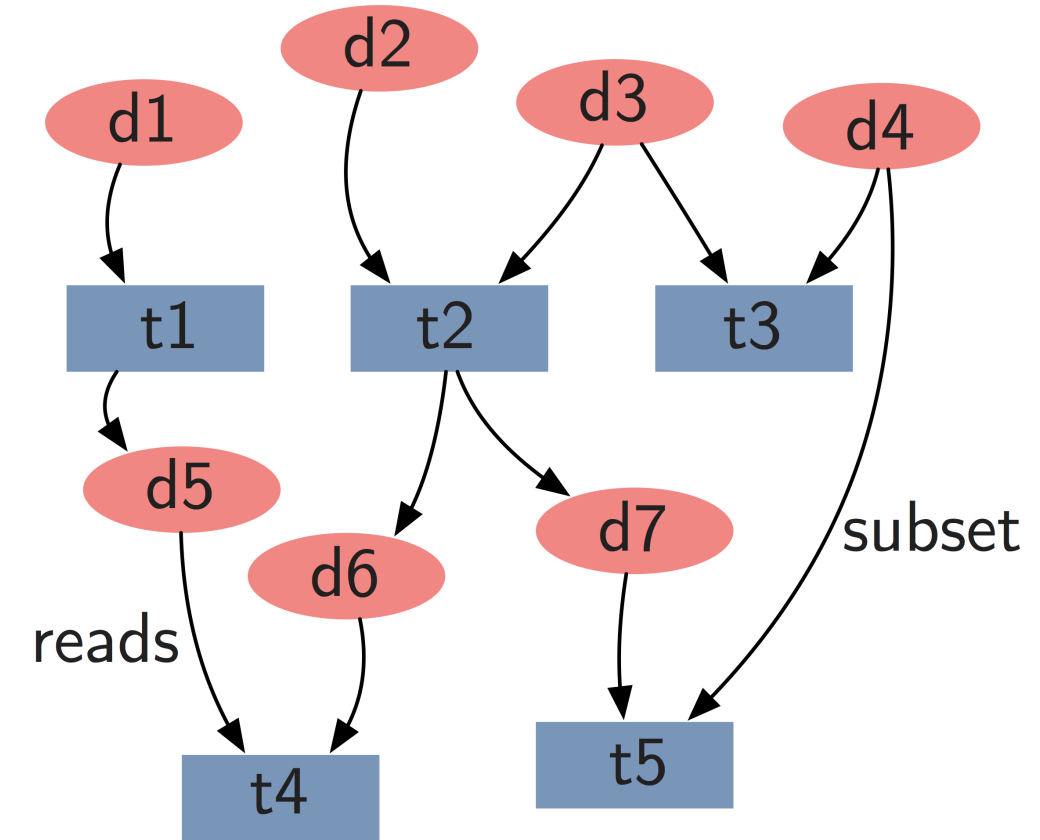
- Broad survey of many AMT runtime systems
- Deep dive on Charm++, Legion, Uintah
- *Programmability:*
Does this runtime enable efficient expression of workloads?
- *Performance:*
How performant is this runtime on current platforms and how well suited is it to address future architecture challenges?
- *Mutability:*
What is the ease of adopting this runtime and modifying it to suit users' needs?



- AMT systems show great promise
- No common user-level APIs
- Need for best practices and standards
- Survey recommendations led to DARMA
- C++ abstraction layer for AMT runtimes
- A single user-level API
- Support multiple AMT runtimes to begin identification of best practices



- AMT runtimes operate with a directed acyclic graph (DAG): captures relationships between application data and inter-dependent tasks
- DAGs can be annotated to capture additional information
 - Tasks' read/write usage of data
 - Task needs a subset of data
- Additional information enables runtime to reason more completely about
 - When and where to execute a task
 - Whether to load balance
- Existing runtimes leverage DAGs with varying degrees of annotation, also leveraging user-provided information about dependencies



Serial code

```
void get_name(string& val){
    /*...*/
    val = "quiet"; }
void get_year(int& val){
    /*...*/
    val = 17; }
void print(string a, int b){
    cout << a << " "
        << b << endl;
}

int main() {
    string name; int year;
    get_name(name);
    get_year(year);
    print(name, year);
}
```

Output: quiet 17

Explicit threads

```
static string name; static int year;
void get_name() {
    /*...*/
    name = "quiet"; }
void get_year() {
    /*...*/
    year = 17; }
void print() {
    cout << name << " " << year << endl;
}

int main() {
    auto thr_n = std::thread(get_name)
    auto thr_y = std::thread(get_year)
    thr_n.join();
    thr_y.join();
    print();
}
```

Output: quiet 17

Using async-future:

```
string get_name() {
    /*...*/
    return "quiet"; }
int get_year() {
    /*...*/
    return 17; }
void print(future<string> a,
           future<int> b) {
    cout << a.get() << " "
        << b.get() << endl;
}

int main() {
    auto n = std::async(get_name);
    auto y = std::async(get_year);
    auto done = std::async(print,
                           move(n), move(y));
    done.wait();
}
```

Output: quiet 17

- Direct extraction of conservative concurrency based on the sequence of data usage
- Conservative because it is "safe by default"
- Enabling runtime-based approaches rather than auto-magic compilers
- There is existing related research (e.g., Legion, OpenMP 4.5)

The function signature itself (from the sequential implementation) can serve as a concurrency specification!

Serial code

```
void get_name(string& val) { val = "quiet"; }
void get_year(int& val) { val = 17; }
void print(string a, int b) {
    cout << a << " " << b << endl;
}

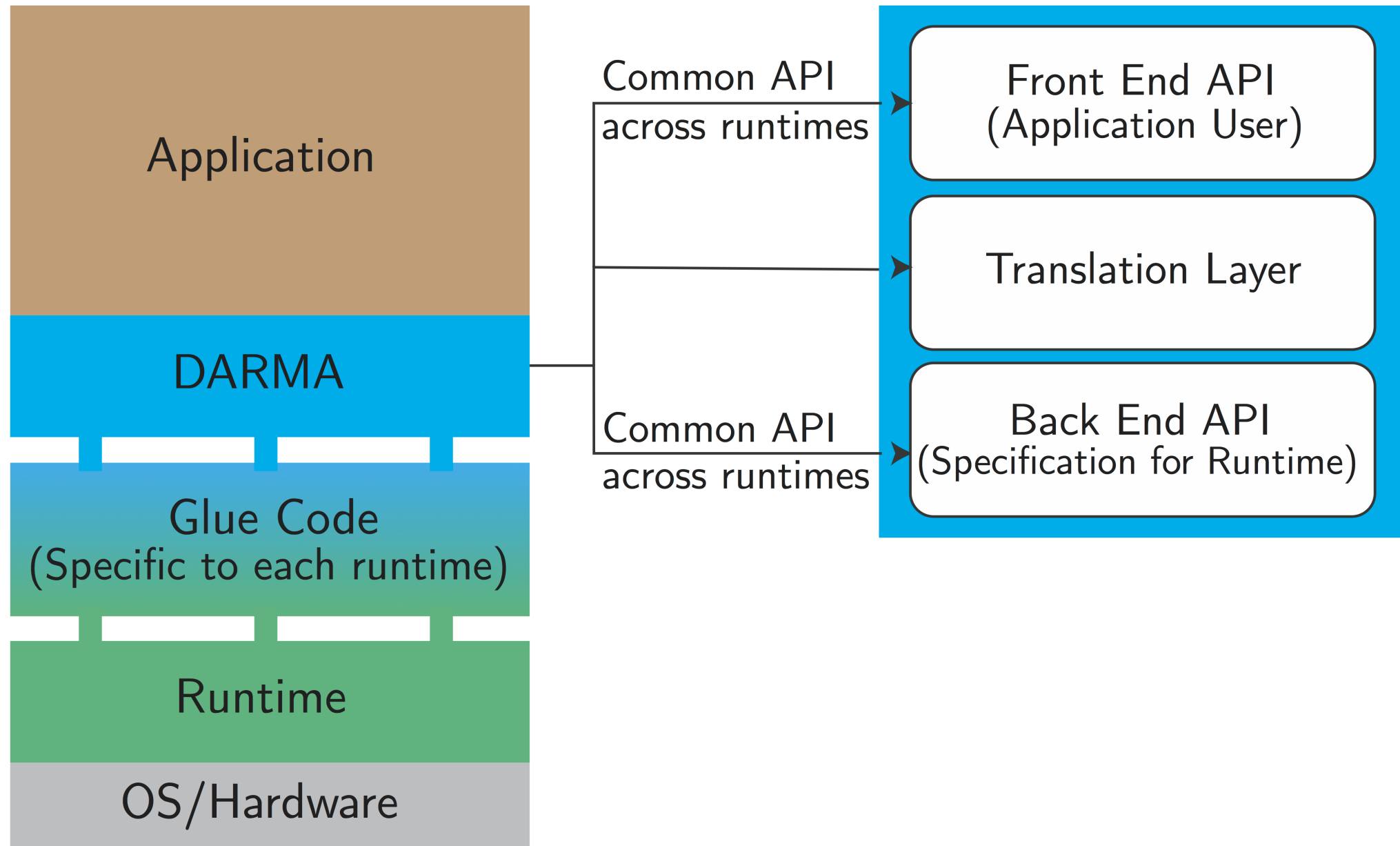
int main() {
    string name; int year;
    get_name(name);
    get_year(year);
    print(name, year);
}
```

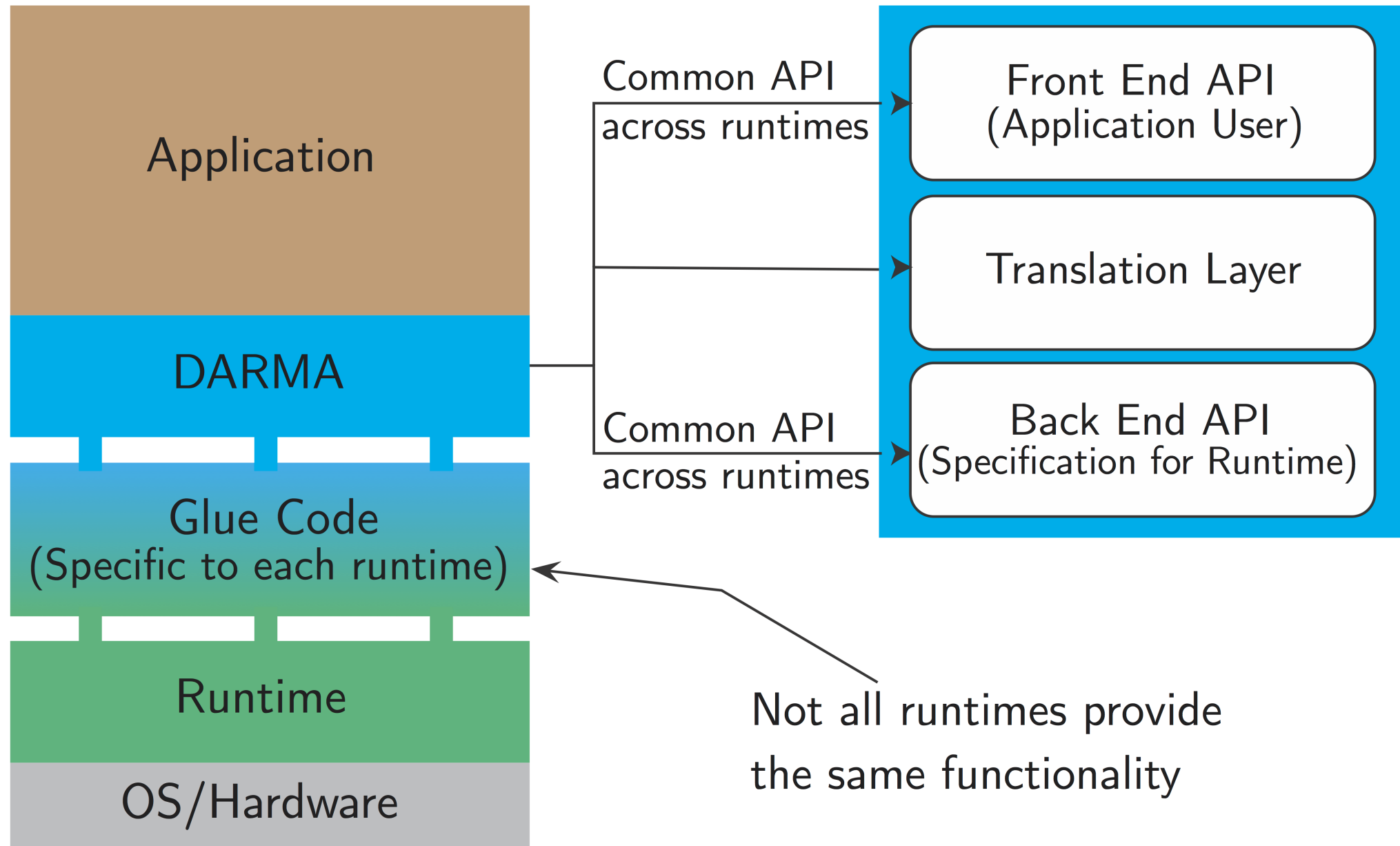
(quasi) DARMA code

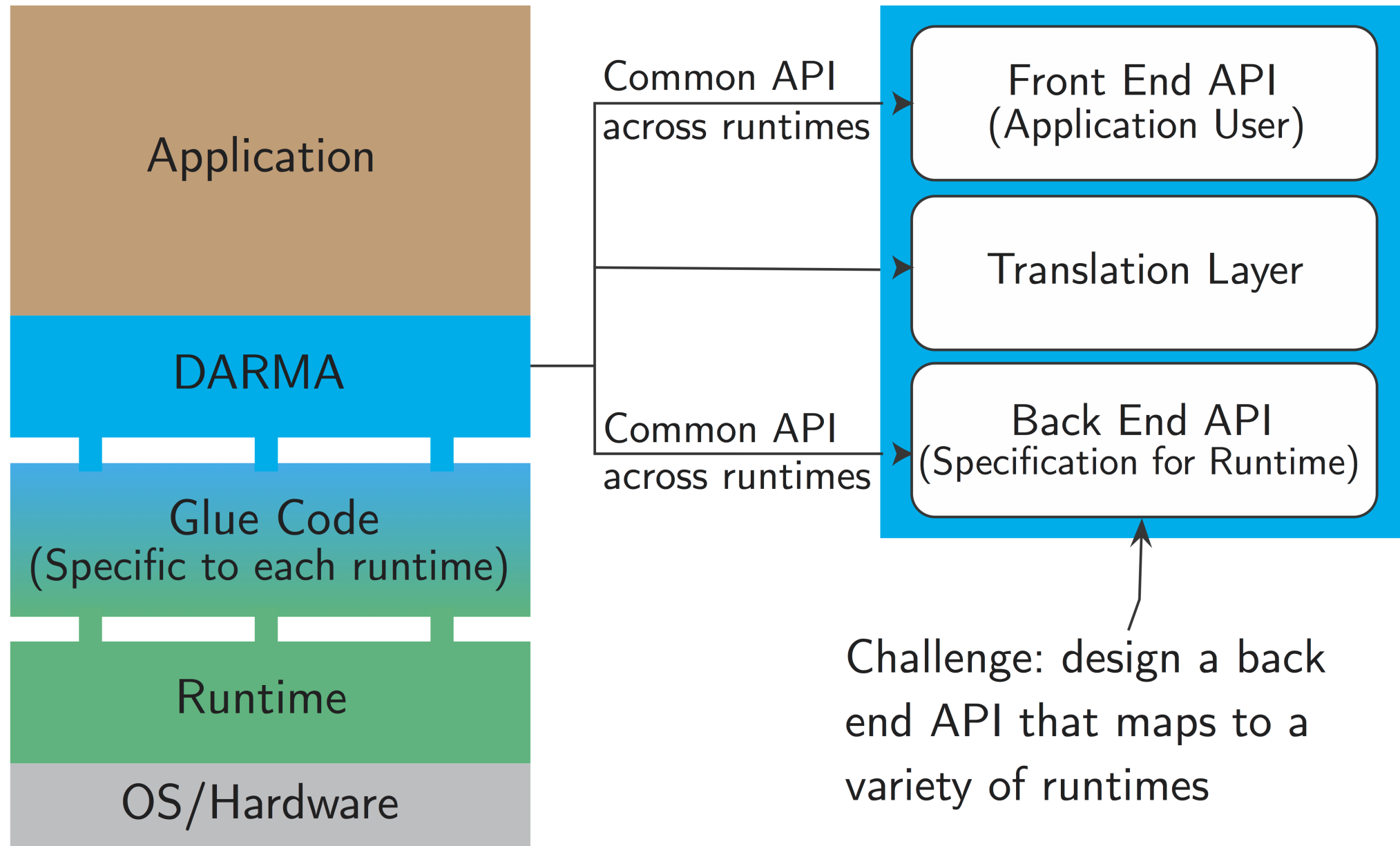
```
void get_name(string& val) { val = "quiet"; }
void get_year(int& val) { val = 17; }
void print(string a, int b) {
    cout << a << " " << b << endl;
}

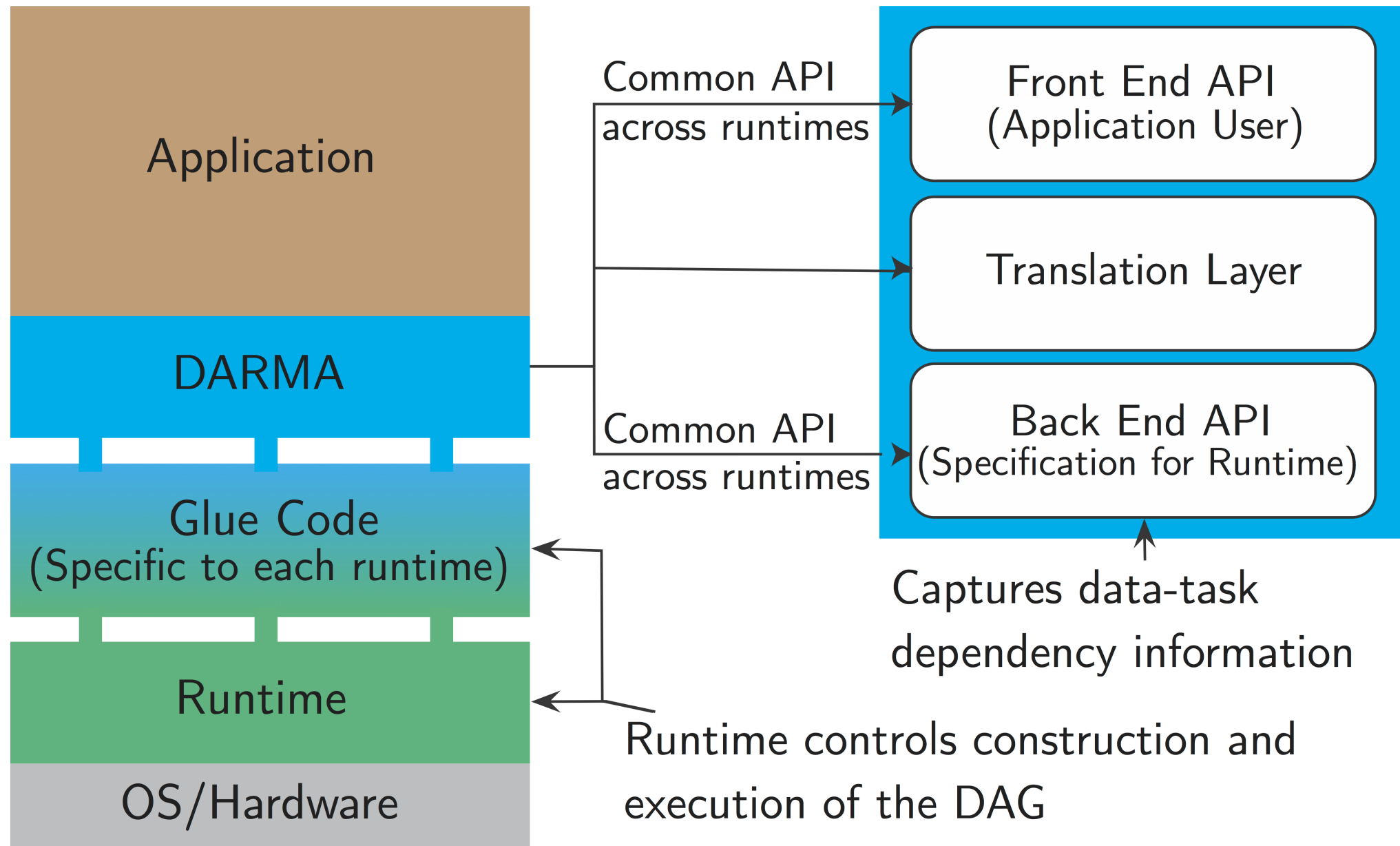
int main() {
    async_ptr<string> name; async_ptr<int> year;
    mpass::async(get_name, name);
    mpass::async(get_year, year);
    mpass::async(print, name, year);
}
```

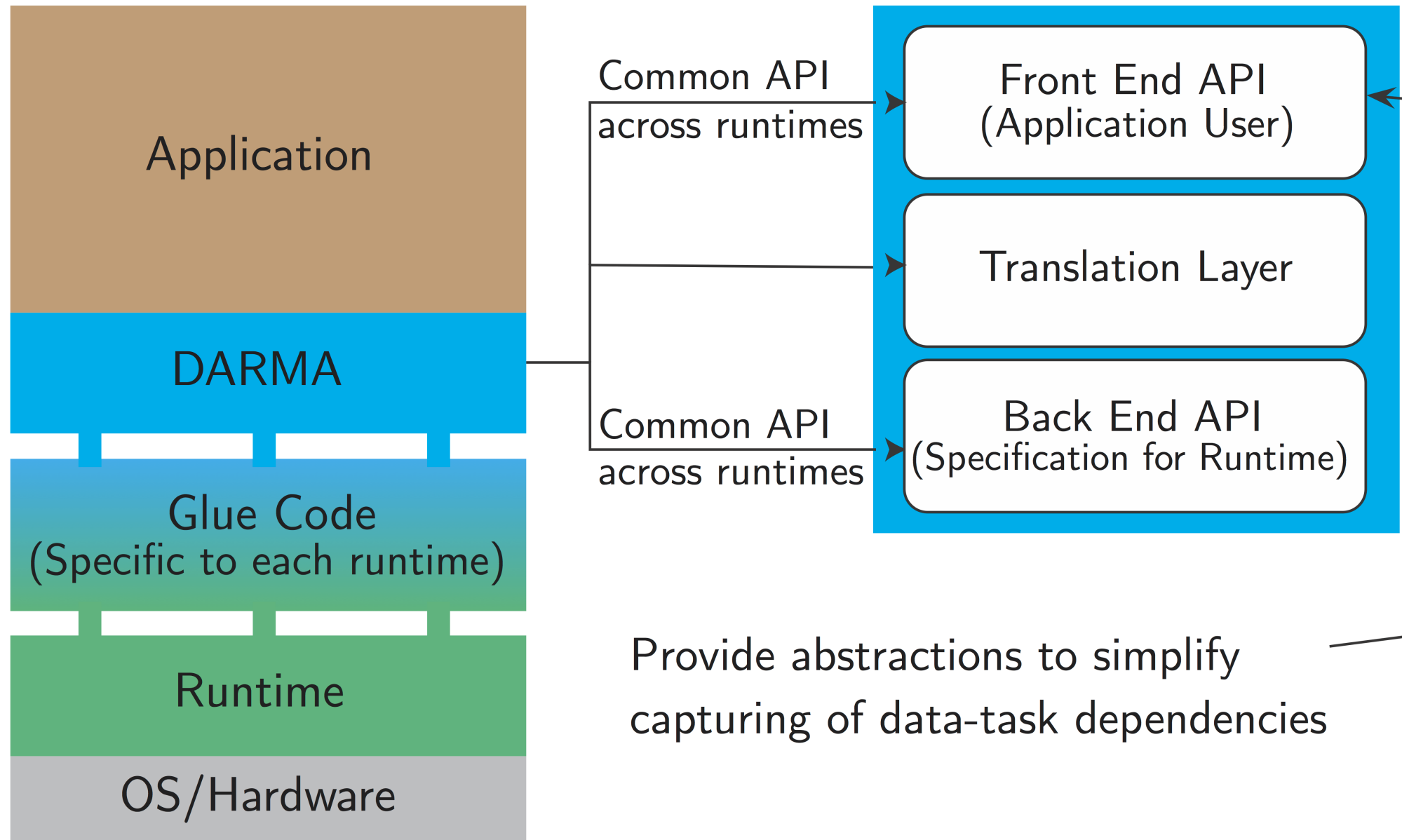
- `mpass::async()` detects dependencies of a task and their use (i.e., read or modify).
- Concurrency with other tasks is implicitly specified by how the data is used.
- A backend task scheduler and runtime layer is needed to execute the DAG.
- Alleviates (some) burden on programmers. Various degrees of user-defined information.

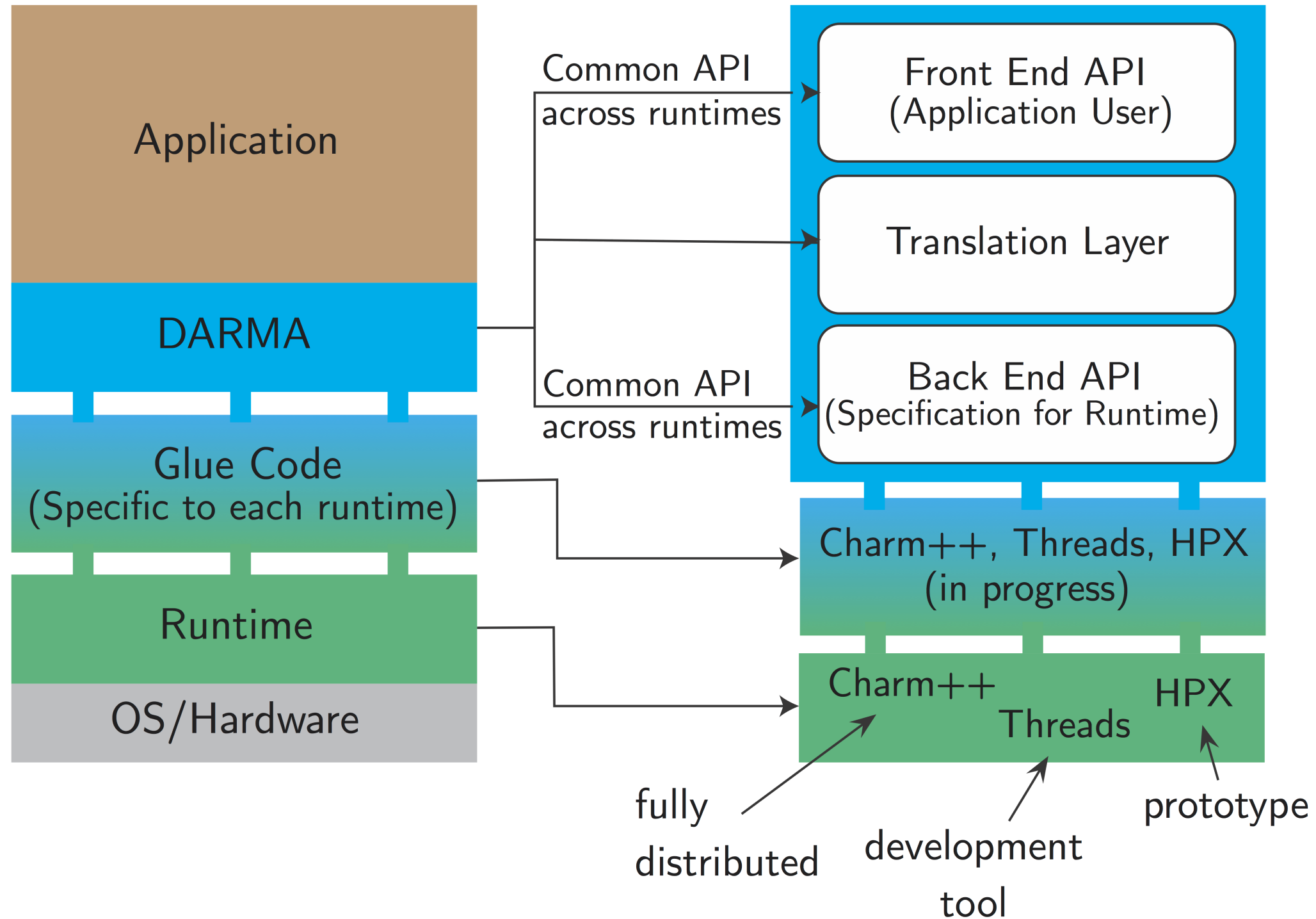












Asynchronous smart pointers wrap application data

- Track meta-data used to build and annotate the DAG
 - Current permissions information (e.g. read-only, read/write)
 - Subsetting information under development

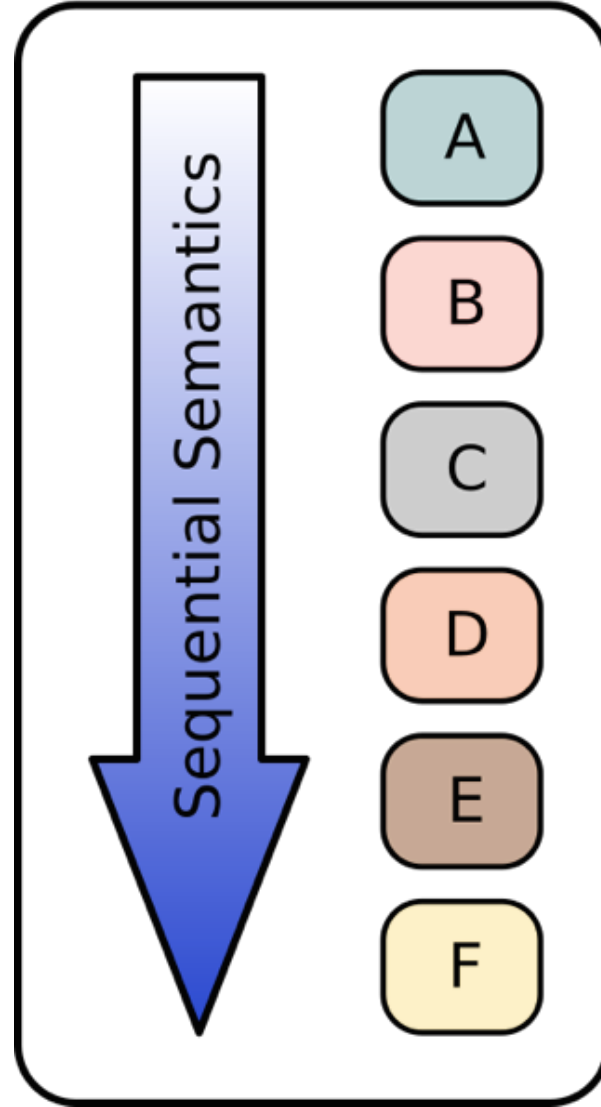
Data partitioning and distribution expressed with explicit, hierarchical, logical decomposition of data using:

- `AccessHandle<T>`
 - Does not span multiple memory spaces
 - Must be serialized to be transferred between memory spaces
- `AccessHandleCollection<T, R>`
 - Expresses a collection of data of type `T`, mapped to range `R`
 - Can be mapped across memory spaces in a scalable manner

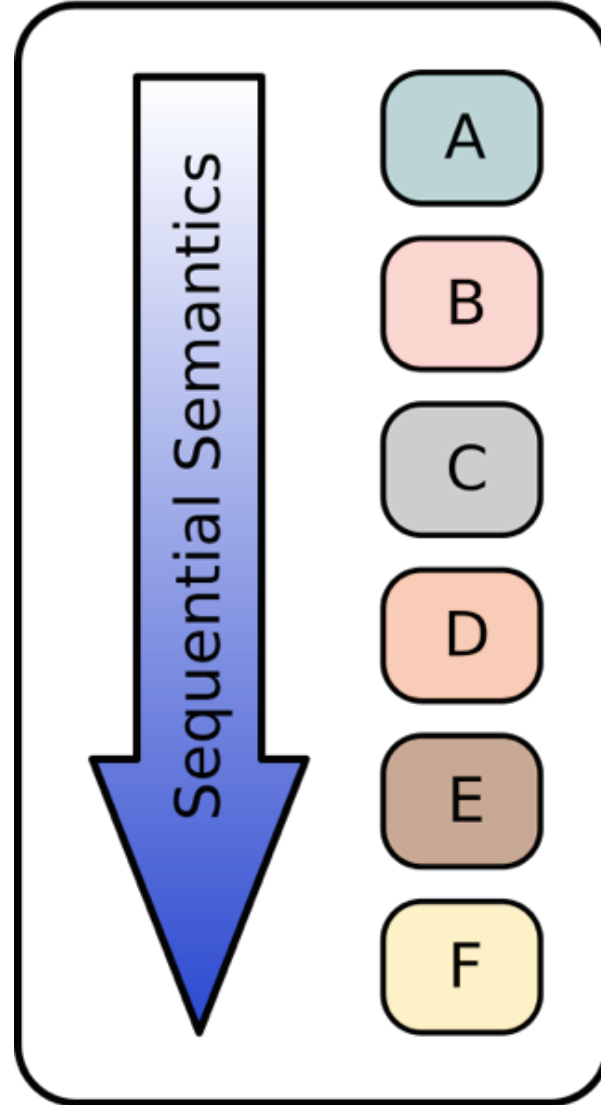
Distribution of data is up to individual backend runtime

- `create_work`
 - A task that doesn't span multiple execution spaces
 - Sequential semantics: the order and manner (e.g., read, write) in which data (`AccessHandle`) is used determines what tasks may be run in parallel
- `create_concurrent_work`
 - Scalable abstraction to launch across distributed systems
 - A collection of tasks that must make simultaneous forward progress
 - Sequential semantics supported across different task collections based on order and manner of `AccessHandleCollection` usage
- How is synchronization expressed?
 - DARMA does not provide explicit temporal synchronization abstractions
 - DARMA does provide data coordination abstractions
 - publish/fetch semantics between participants in a task collection
 - Asynchronous collectives between participants in a task collection

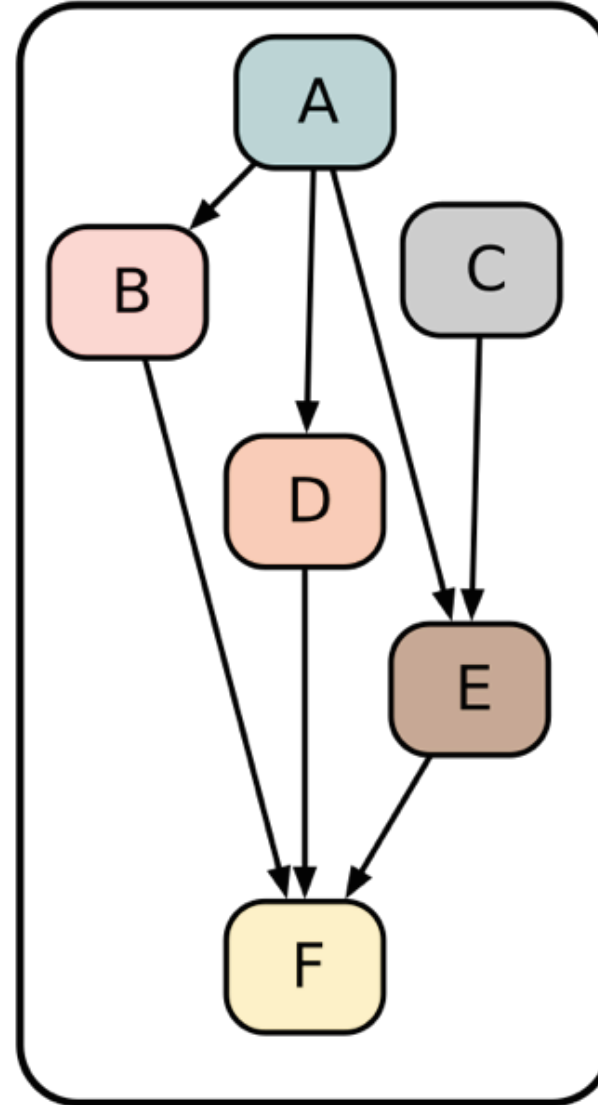
App in DARMA



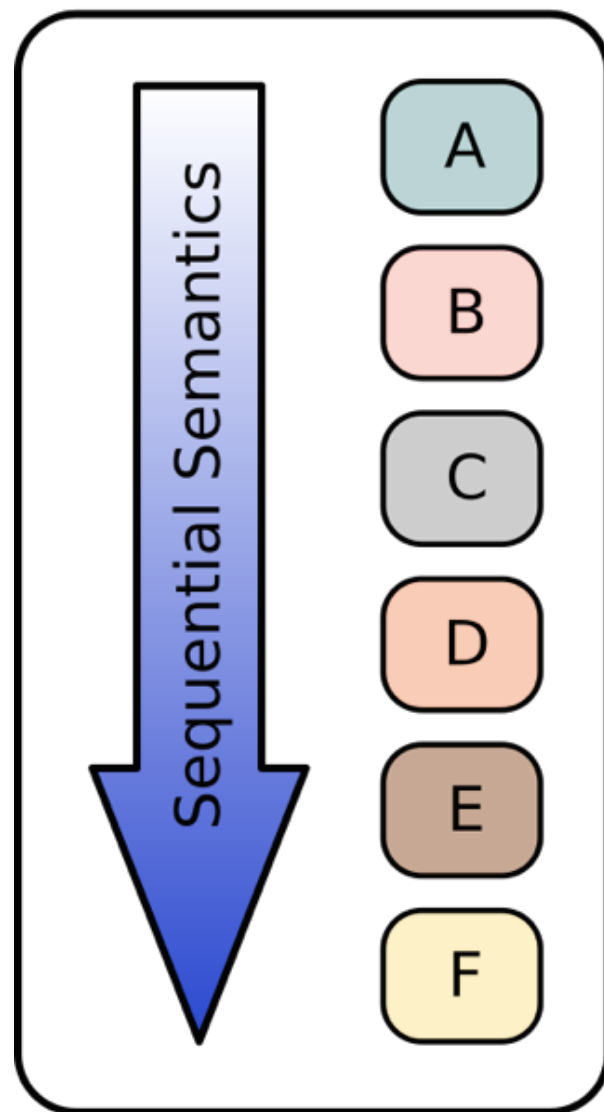
App in DARMA



App in a runtime

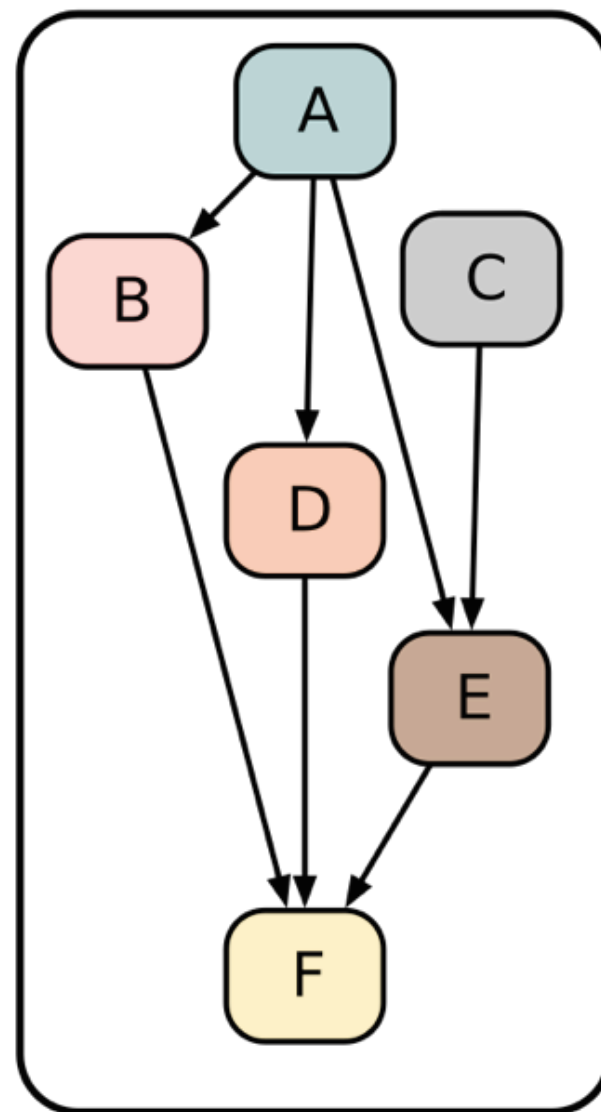


App in DARMA



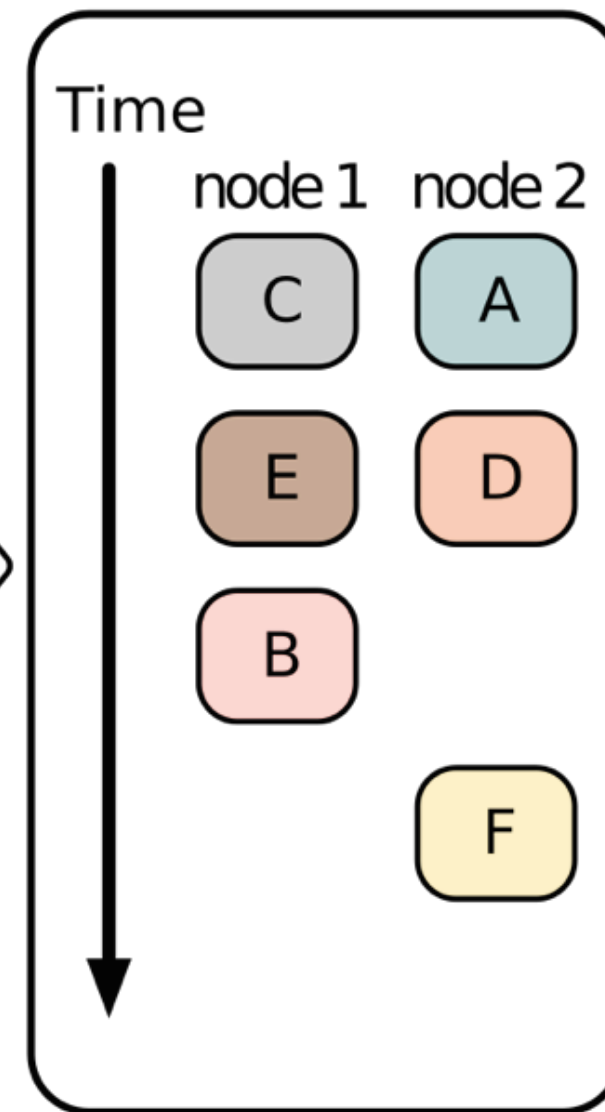
DARMA

App in a runtime



Runtime

App on a hardware



Example Program

```
AccessHandle<int> my_data;
```

```
darma::create_work([=]{  
    my_data.set_value(29);  
});
```

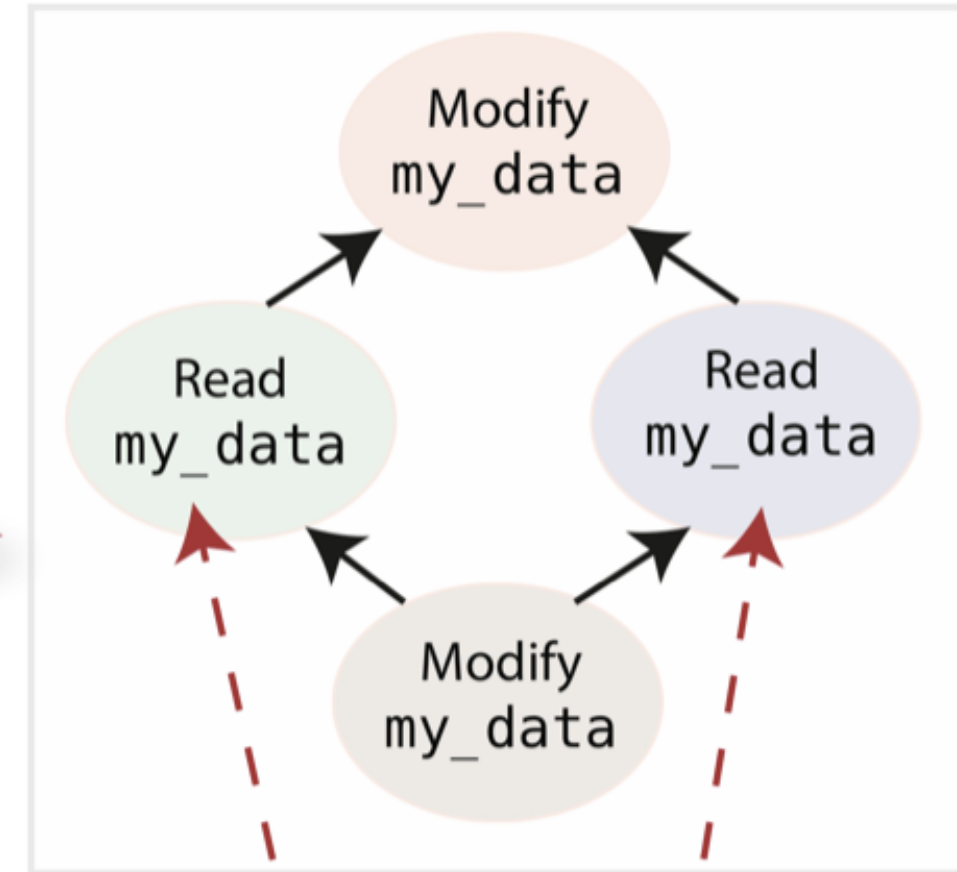
```
darma::create_work(  
    reads(my_data), [=]{  
        cout << my_data.get_value();  
    }  
);
```

```
darma::create_work(  
    reads(my_data), [=]{  
        cout << my_data.get_value();  
    }  
);
```

```
darma::create_work([=]{  
    my_data.set_value(31);  
});
```

DAG (Directed Acyclic Graph)

*Sequential
Semantics*



These two tasks are concurrent
and can be run in parallel by a
DARMA backend!

```

void darma_main_task(std::vector<std::string> args) {

    auto answer = initial_access<int>();

    //set value of answer - must run first
    create_work([=]{ answer.set_value(42); });

    //read-only, can run in parallel with check below
    create_work(reads(answer), [=]{
        std::cout << "The answer is" << *answer << std::endl;
    });

    //read-only, can run in parallel with print above
    create_work(reads(answer), [=]{
        if (*answer != 42){
            darma_runtime::abort("the answer is incorrect");
        }
    });
}

DARMA_REGISTER_TOP_LEVEL_FUNCTION(darma_main_task);

```


UQ with DARMA

fnrizzi.github.io/quiet17

- Uncertainties in inputs propagated to outputs:
 - Moments, reliability, PDFs of the outputs
- Techniques:
 - Sampling methods: ex. Monte Carlo, Multi-level MC, Importance sampling.
 - Functional expansion-based methods: ex. PCe.
- Need multiple evaluation of forward model (e.g. PDE).
- Why is DARMA (AMT) good for UQ?
 - (Dynamic) parallelism: heterogeneity among samples
 - AMT model is a natural fit
 - Nested UQ evaluations
 - Adaptive UQ algorithms
 - Performance portability, expressiveness and productivity

Multiple Solves per Rank

```
using vecD = vector<double>;

struct RunSamples {
    void operator()(
        Index1D<size_t> index,
        AccessHandleCollection<vecD, Range1D> ahcdata, /*...*/
    ) {
        ahcdata[index].local_access().resize(solves_per_rank, 0.0);
        for (uint i = 0; i < solves_per_rank; ++i) {
            create_work([=] {
                // generate sample diffusivity
                // solve PDE for current germ sample
                // independently store QoI from this sample
            });
        } // for
    } // op
};

//
void darma_main_task(std::vector<std::string> args) {

    const uint solves_per_rank = ...; // # of PDE solves per rank
    const uint n_ranks         = ...; // # of ranks

    auto data = initial_access_collection<vecD>(Range1D(n_ranks));

    create_concurrent_work<RunSamples>(data, ..., Range1D(n_ranks))
    create_concurrent_work<Collect>(data, ..., Range1D(n_ranks));

}
DARMA_REGISTER_TOP_LEVEL_FUNCTION(darma_main_task);
```

Single Solve per Rank

```
using vecD = vector<double>;

struct RunSamples {
    void operator()(
        Index1D<size_t> index,
        AccessHandleCollection<double, Range1D> ahcdata, /*...*/
    ) {
        // generate sample diffusivity
        // ...
        // solve PDE for current germ sample
        // store QoI from this sample
        // ...
    }
};

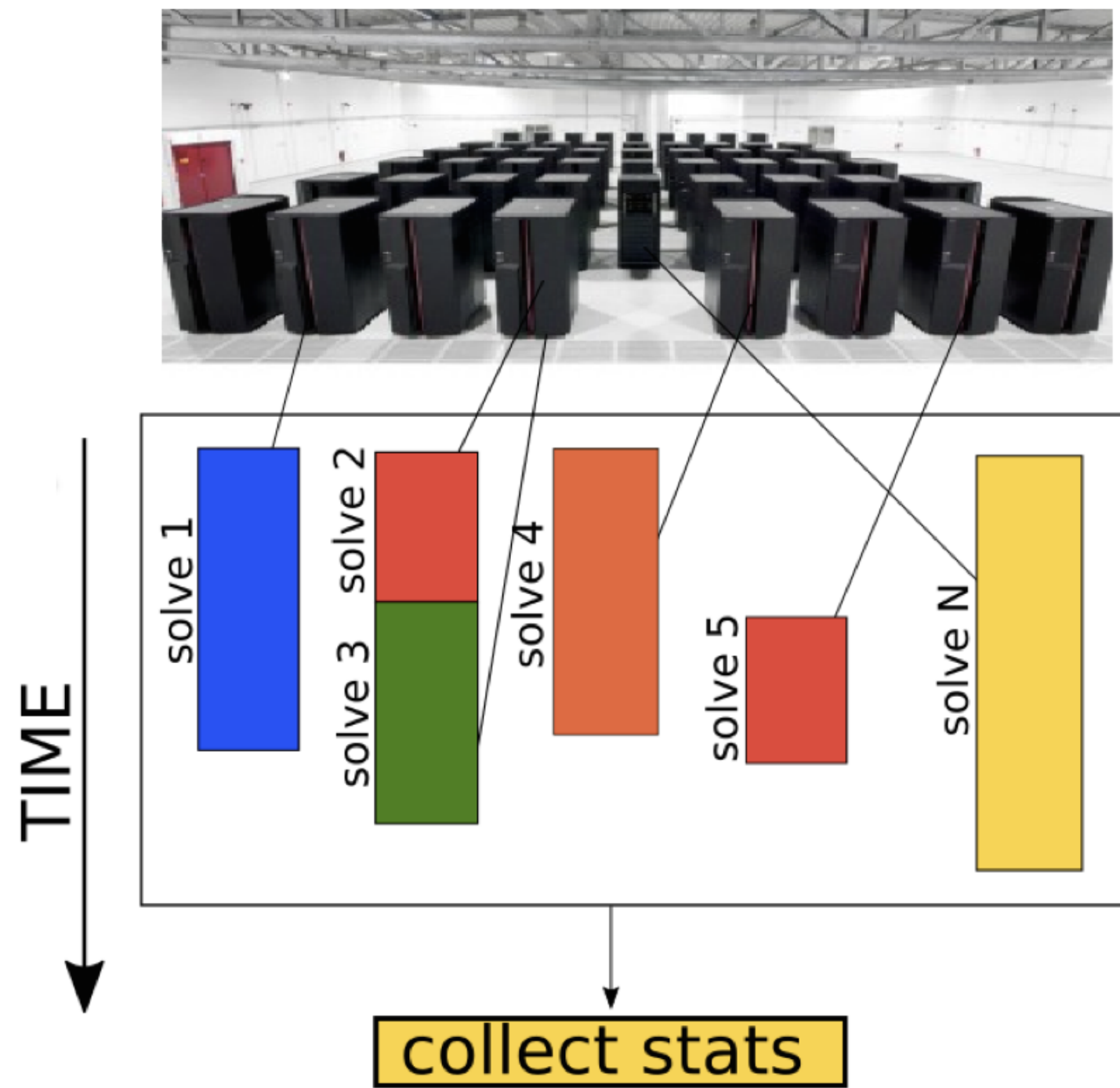
void darma_main_task(std::vector<std::string> args) {

    const uint n_ranks = ...; // # of ranks

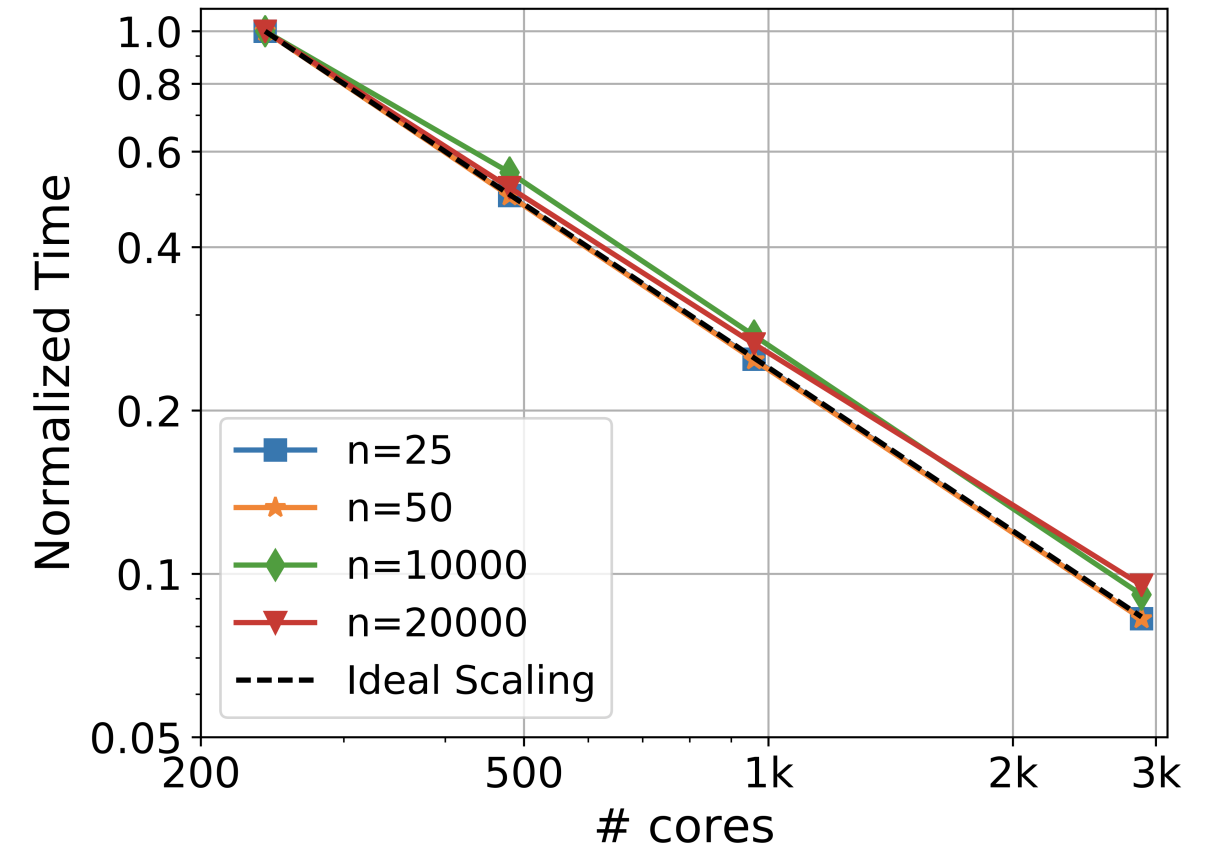
    auto data = initial_access_collection<double>(Range1D(n_ranks))

    create_concurrent_work<RunSamples>(data, ..., Range1D(n_ranks))
    create_concurrent_work<Collect>(data, ..., Range1D(n_ranks));

}
DARMA_REGISTER_TOP_LEVEL_FUNCTION(darma_main_task);
```



- n : number of samples per DARMA ``stream''
- Total number of samples: $N=n*2880$, where 2880 is tot # of threads (96 nodes * 30 threads/node).
- Each PDE solve involves 4,194,304 points.
- Run on Haswell nodes on in-house machine.
- Grid size of each PDE solve: 4,194,304 points.
- Good scaling (as expected). Similar (or better) runtime of equivalent charm++ code.
- Cache/memory effects appear for larger problems.



Top-level Task

```
void darma_main_task(std::vector<std::string> args){

auto vLevelsH = initial_accesss<vector<Level>>();
create_work<initialize>(vLevelsH, ...);

auto converged = initial_accesss<bool>();
auto iter = initial_accesss<uint>();
create_work([=]{
    converged.set_value(false);
    iter.set_value(1);
});

create_work_while([=]{
    return converged.get_value() == false
    && iter.get_value() <= maxIter;
}).do_([=]
{
    // loop over level and run/collect samples
    for (auto & lev : vLevelsH)
        create_concurrent_work<runFunctor>(lev, ...);

    // compute stats, new # of samples, check convergence
    create_work<checkStats>(vLevelsH, converged, ...);
    iter.get_reference()++;
});

// compute estimator
create_work<MLEstimator>(vLevelsH, ...);
}
DARMA_REGISTER_TOP_LEVEL_FUNCTION(darma_main_task);
```

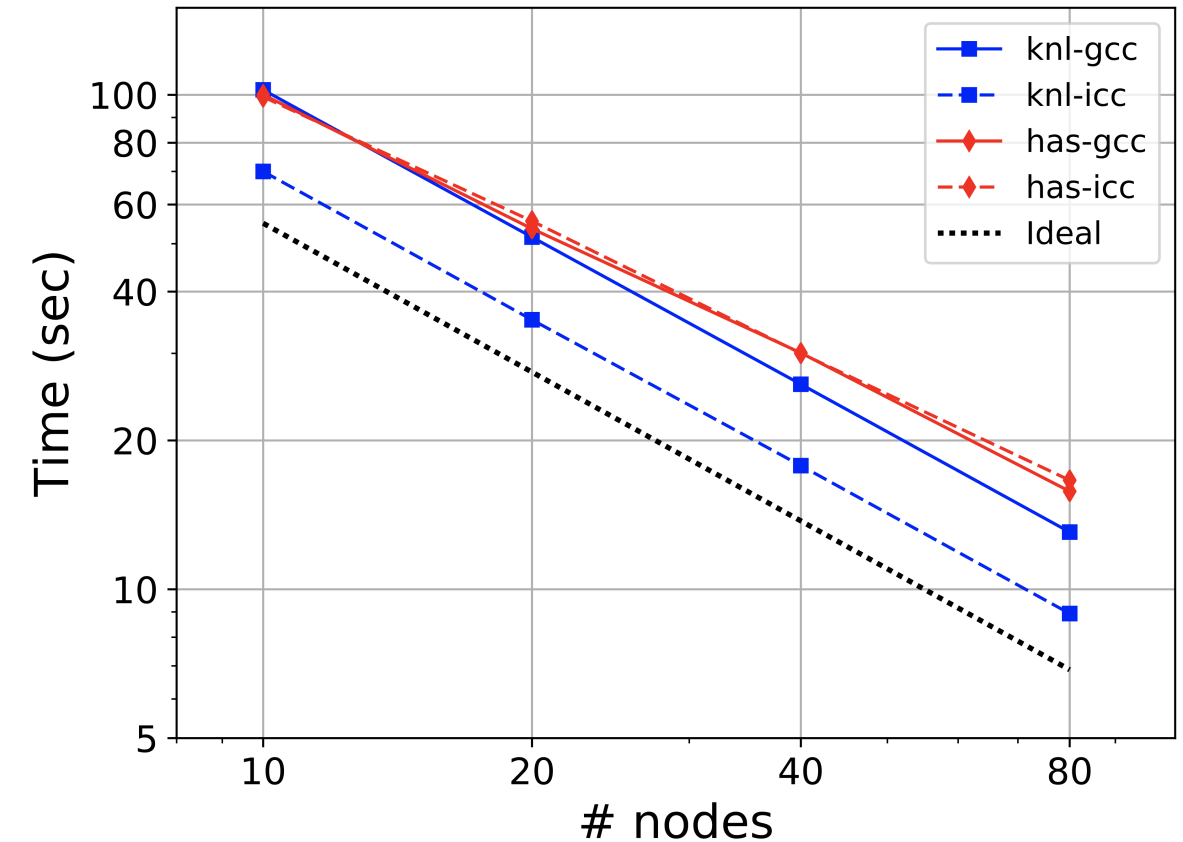
Core Run Functor

```
struct runFunctor{
    void operator()(
        Index1D index, ...) const
    {
        const auto contextSize = index.max_value + 1;

        uint myN = std::ceil(N/contextSize);
        for (uint i = 0; i < myN; ++i){
            create_work([=]
            {
                // generate sample of stochastic diffusivity
                create_work([=]{
                    // PDE solve for l level (fine)
                });
                create_work([=]{
                    // PDE solve for l-1 level (coarser)
                });

                // store target QoI for fine Q_l
                // store target QoI for coarse Q_lm1
                // store target QoI: Y = Q_l - Q_lm1;
            });
        } // for
    } // op
};
```

- Adaptive MLMC: start with fixed initial number of levels (4), add more as needed. Coarsest level has 4096 points.
- E.g. need to handle $\sim 10^7$ tasks.
- Dynamic addition of levels is interesting and challenging for task/data mapping and speculative execution.
- Compare Haswell/KNL using gcc/icc on in-house machine.
- Good scaling. Still investigating Haswell trend.
- Looking now into more heterogeneous problems: impact of load balancing, dynamic parallelism and optimal task mapping.



Take-home message:

- AMT provides a promising framework for exploring UQ on next generation machines.
- DARMA provides a unified interface for AMT: one code, multiple runtimes.
Automatic dependency detection and parallel/concurrency reasoning.

Work in progress:

- Leveraging data reusability
 - Reuse data produced by some tasks to accelerate convergence for other similar tasks.
 - Tradeoff between data movement, execution time (memory access, data locality).
- Benchmarking for distributed backends.
- Optimize load balancing methods for UQ applications.
- Development of abstraction tailored for UQ.

DARMA is about to become publicly available, for info: <https://share-ng.sandia.gov/darma/darma@sandia.gov>

Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC., a wholly owned subsidiary of Honeywell International, Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525.

Questions? Comments?

Thank you for your attention!

fnrizzi@sandia.gov