



**BERKELEY LAB**  
LAWRENCE BERKELEY NATIONAL LABORATORY



# The UPC++ Library for Exascale and Data Intensive Computing

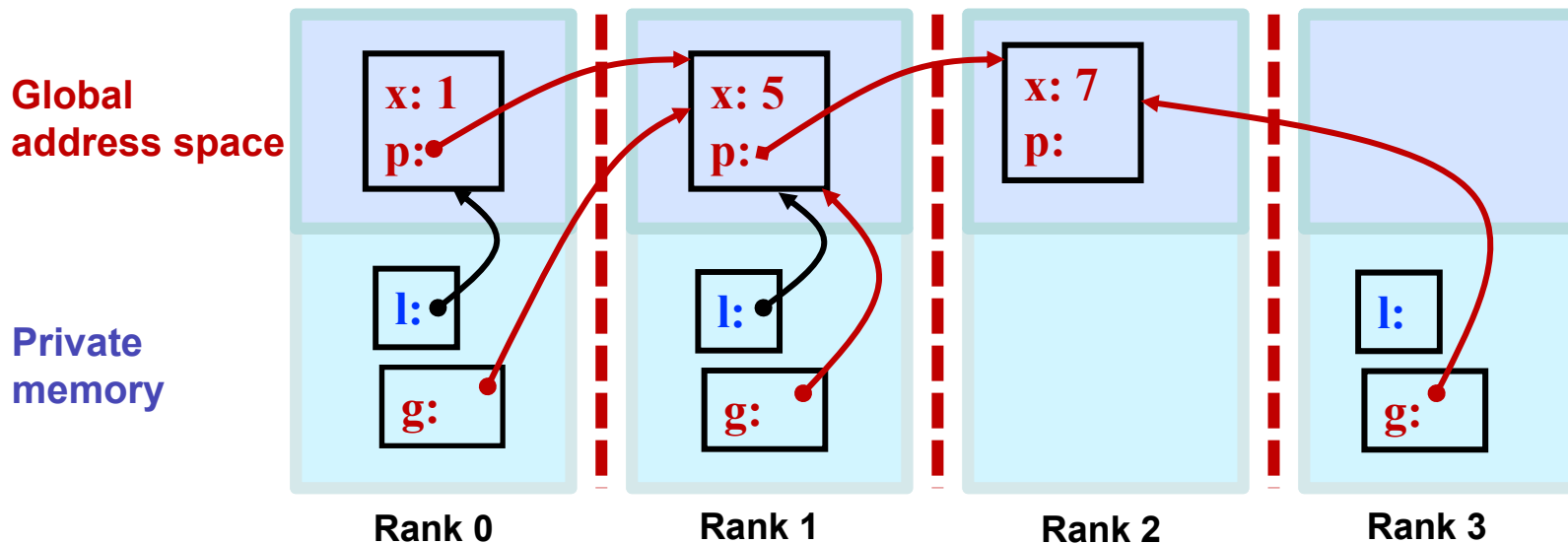
<http://upcxx.lbl.gov>

**Scott B. Baden**

Group Lead, Computer Languages and Systems Software Group

# UPC++: a C++ PGAS Library [iPDPS19]

- Global Address Space (**PGAS**)
  - A portion of the physically distributed address space is visible to all processes. Now generalized to handle GPU memory
- Partitioned (**PGAS**)
  - *Global pointers* to shared memory segments have an *affinity* to a particular rank
  - Explicitly managed by the programmer to optimize for locality



# What does UPC++ offer?

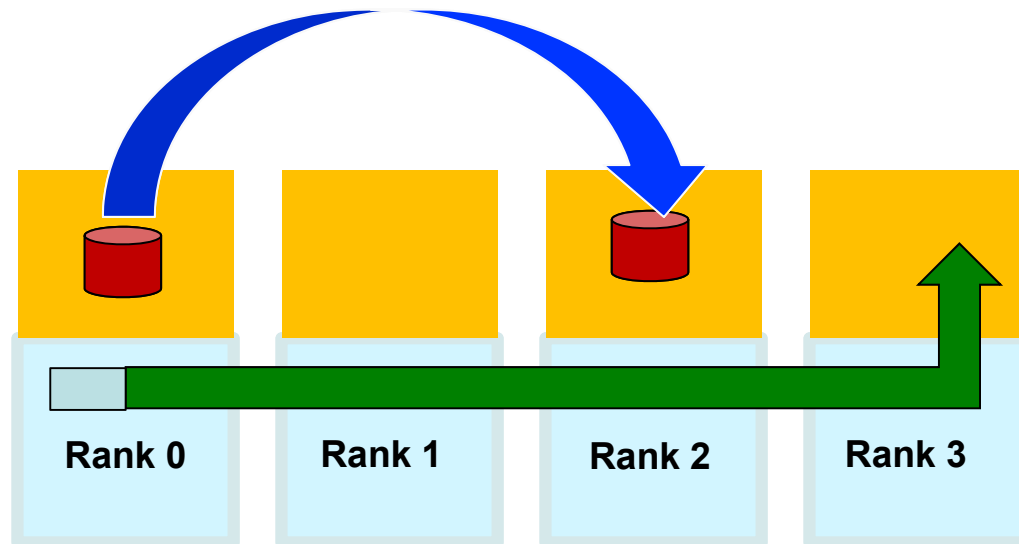
- Asynchronous behavior based on futures/promises
  - RMA: Low overhead, zero-copy 1 sided communication. Get/put to a remote location in another address space
  - RPC: Remote Procedure Call: invoke a function remotely. A higher level of abstraction, though at a cost
- Design principles encourage performant program design
  - All communication is explicit (unlike UPC)
  - All communication is asynchronous: futures and promises
  - Scalability

Remote procedure call (RPC)

Global address space (Shared segments)

One sided communication

Private memory



# Memory Kinds

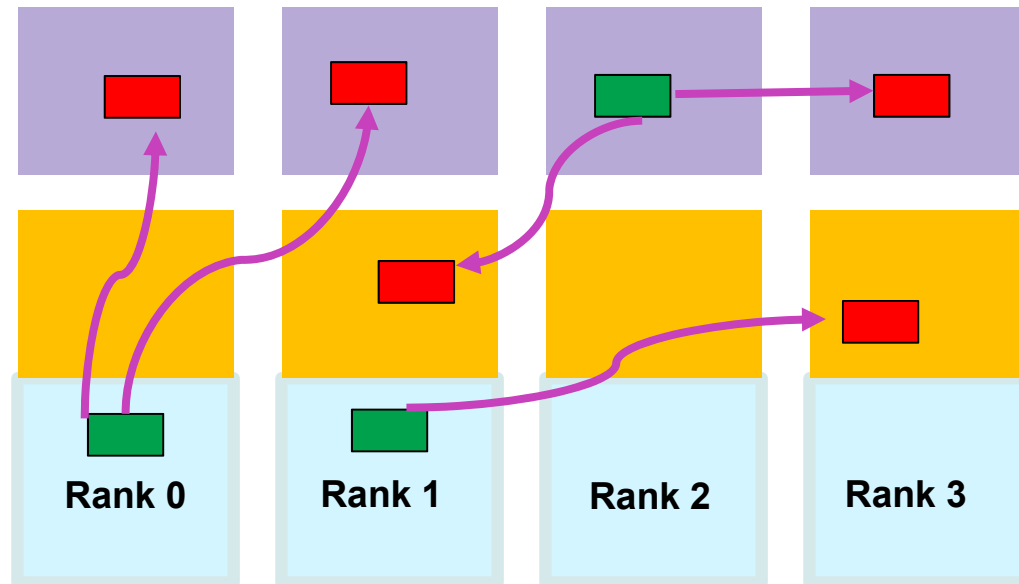
- UPC++ now supports memory kinds, which provides a uniform RMA interface for moving data between DRAM and GPU memory (or between GPU memories) transparently
- Avoids the need for different access methods

Global address space  
(Shared segments)

Device Memory Kind

Host

Private memory

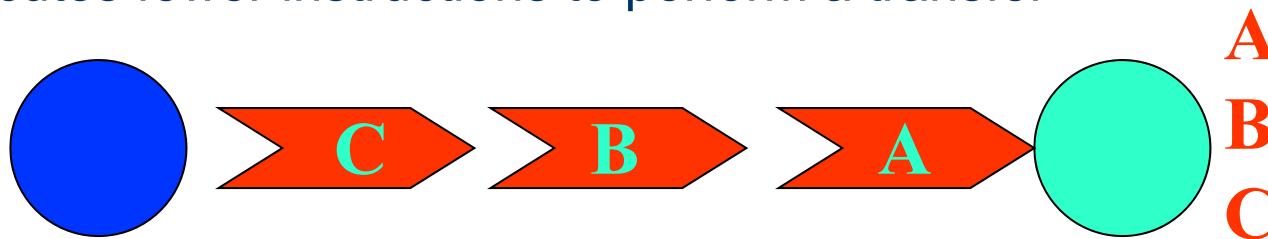


# Why is PGAS attractive?

- The overheads are low  
Multithreading can't speed up overhead
- Memory per core is dropping, requiring reduced communication granularity
- Irregular applications exacerbate granularity problem
- Software managed memories are becoming more common, with different access methods.  
We need a unified method for accessing them
- Current and future HPC networks use one-sided transfers at their lowest level and the PGAS model matches this hardware with very little overhead

# Where does message passing overhead come from?

- Matching sends to receives
  - Messages have an associated context that needs to be matched to handle incoming messages correctly
  - Data movement and synchronization are coupled
- Ordering guarantees are not semantically matched to the hardware
- UPC++ avoids these factors that increase the overhead
  - No matching overhead between source and target
  - Executes fewer instructions to perform a transfer



# How does UPC++ deliver the PGAS model?

- A “Compiler-Free” approach
  - Need only a standard C++ compiler, leverage C++ standards
  - UPC++ is a C++ template library
- Relies on GASNet-EX for low overhead communication
  - Efficiently utilizes the network, whatever that network may be, including any special purpose offload support
- Designed to allow interoperation with existing programming systems
  - 1-to-1 mapping between MPI and UPC++ ranks
  - OpenMP and CUDA can be mixed with UPC++ in the same way as MPI+X

# A simple example of asynchronous execution

By default, all communication ops are split-phased

- **Initiate** operation
- **Wait** for completion

A future holds a value and a state: ready/not ready

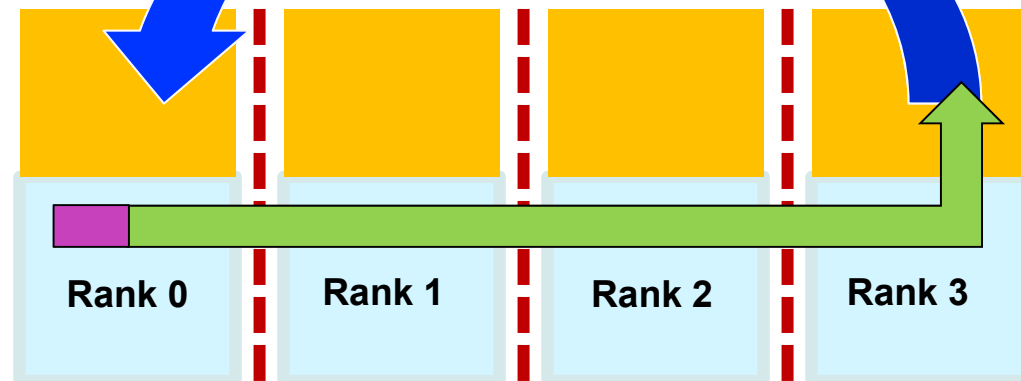
```
global_ptr<T> gptr1 = . . . ;  
future<T> f1 = rget(gptr1);  
// unrelated work..  
T t1 = f1.wait();
```

**Wait returns with result  
when rget completes**

**Global address space**

**Start the get**

**Private memory**

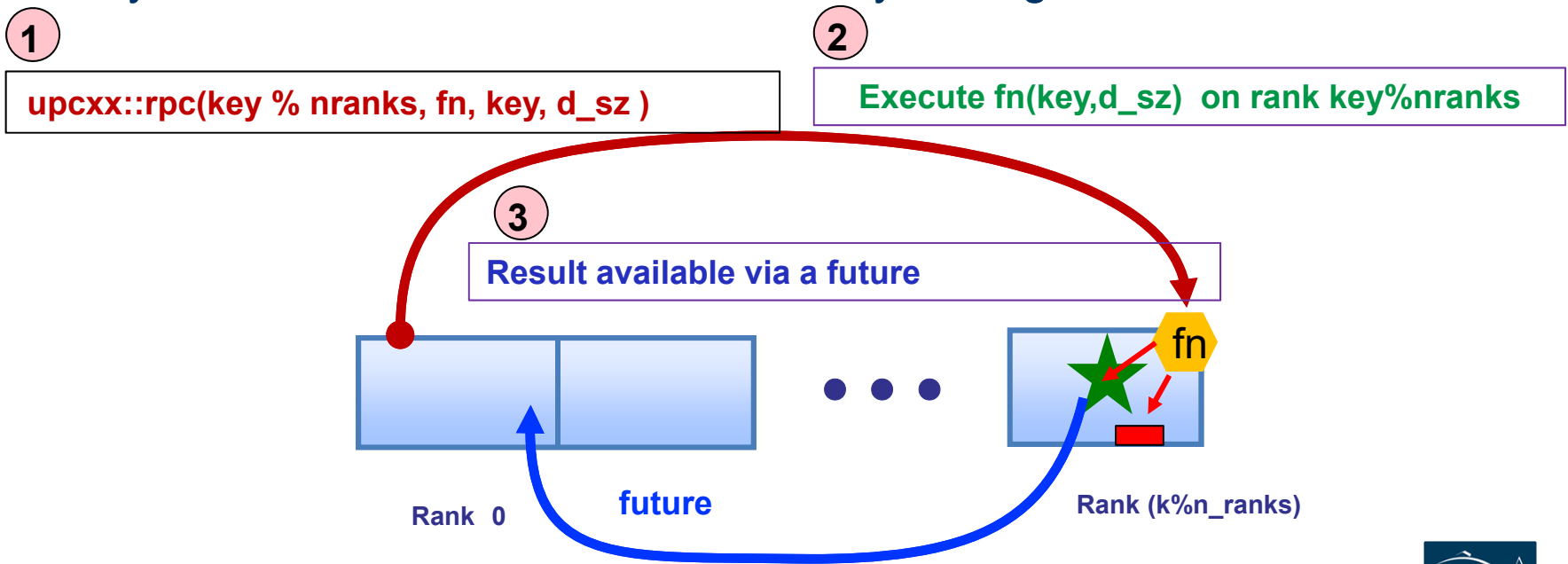


# Simple example of remote procedure call

Execute a function on another rank, sending arguments and returning an optional result

1. Injects the RPC to the *target rank*  $key \% n\_ranks$
2. Executes  $fn(key, d\_sz)$  on target rank at some future time determined at the target
3. Result becomes available to the caller via the future

Many invocations can run simultaneously, hiding data movement



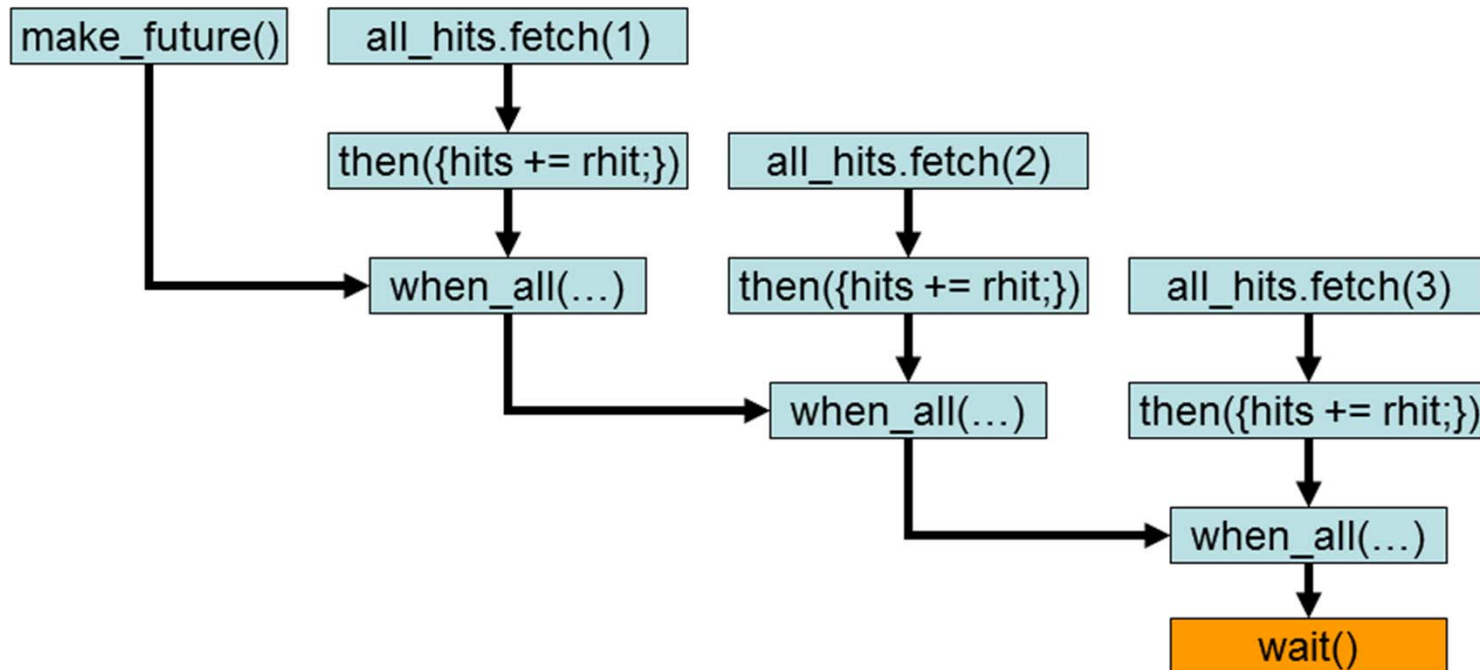
# Composing asynchronous operations

- Rput, rget and rpc return a future
- We can build a DAG of futures, synchronize on the whole rather than on the individual operations
  - Attach a callback: **.then(Foo)**
  - **Foo** is the completion handler, a function or  $\lambda$ 
    - runs locally when the **rget** completes
    - receives arguments containing result associated with the future

```
double Foo(int x){ return sqrt(2*x); }
global_ptr<int> gptr1;
// ... gptr1 initialized
future<int> f1 = rget(gptr1);
future<double> f2 = f1.then(Foo);
// DO SOMETHING ELSE
double y = f2.wait();
```

# Conjoining futures

- We can join futures using `when_all()`
- (Example taken from *UPC++ Programmer's Guide*)



# Road Map

An application

Memory Kinds (GPU Memory)

Inside UPC++

Distributed data structures

UPC++ in Context

# Application: *De Novo* Genome Assembly

Construct a genome (chromosome) from a pool of short fragments, called *reads*, produced by sequencers

Analogy: shred many copies of a book, and reconstruct the book by examining the pieces

Complications: shreds of other books may be intermixed, can also contain errors

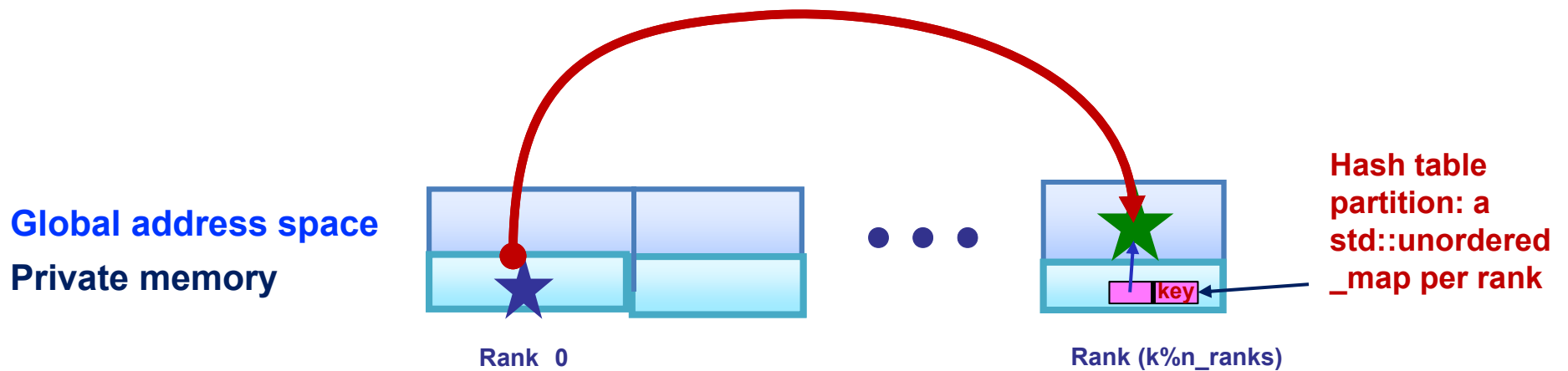
Chop the reads into fixed-length fragments (k-mers)

K-mers form a De Bruijn graph, traverse the graph to construct longer sequences

Graph is stored in a *distributed hash table*

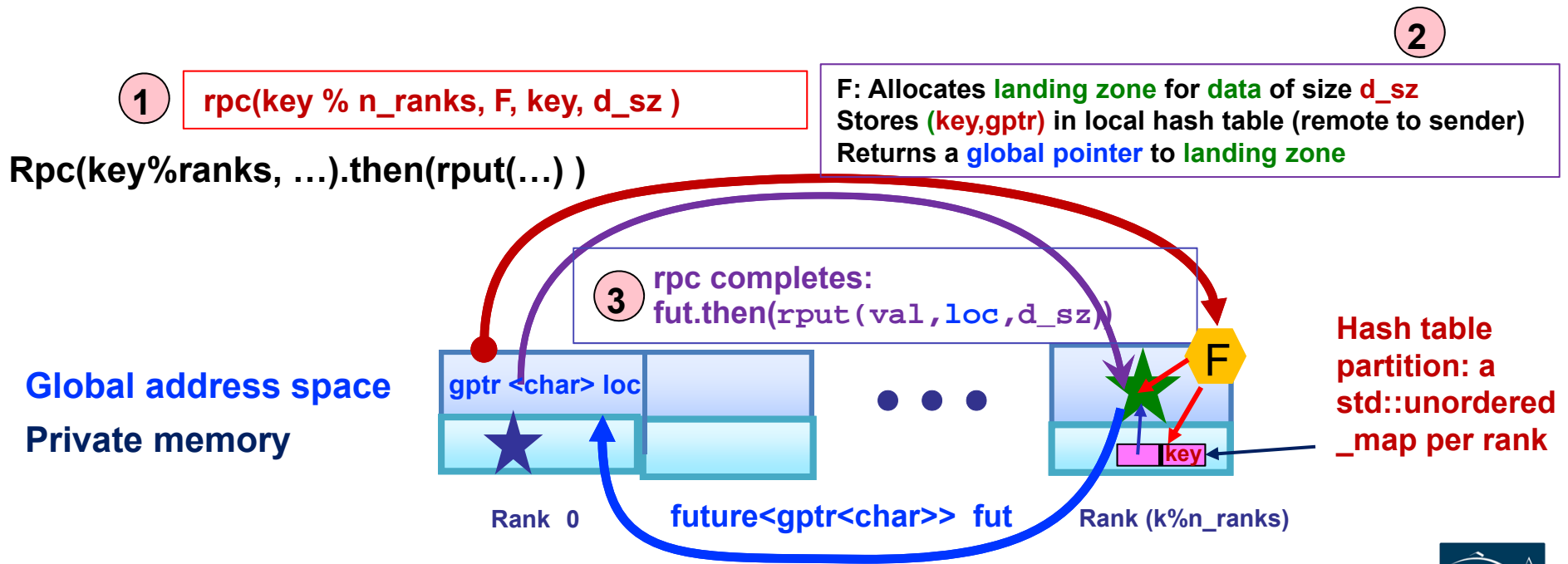
# Distributed hash table implementation

- Used in de-novo Genome Assembly
- This example motivates Remote Procedure Call (RPC)
- RPC simplifies the distributed hash table design
- Store value in a distributed hash table, at a remote location



# Distributed hash table implementation

- Store value in a distributed hash table, at a remote location
- RPC inserts the key @ target and obtains a landing zone pointer
- Once the RPC completes, an attached callback (`.then`) uses `rput` to store the associated data
- We use futures to build a small chain of dependent operations
- The returned future represents the whole operation



# The hash table code

- We use lambda for the RPC function in this example
- RPC inserts key meta data at the target, & allocates the landing zone
  - Leverage implicit synchronization of rpc execution
- Once the RPC completes, the callback (.then()) that was attached to the RPC uses a zero copy rput to store the associated data
  - Exploits the power of rput for high performance RMA where available

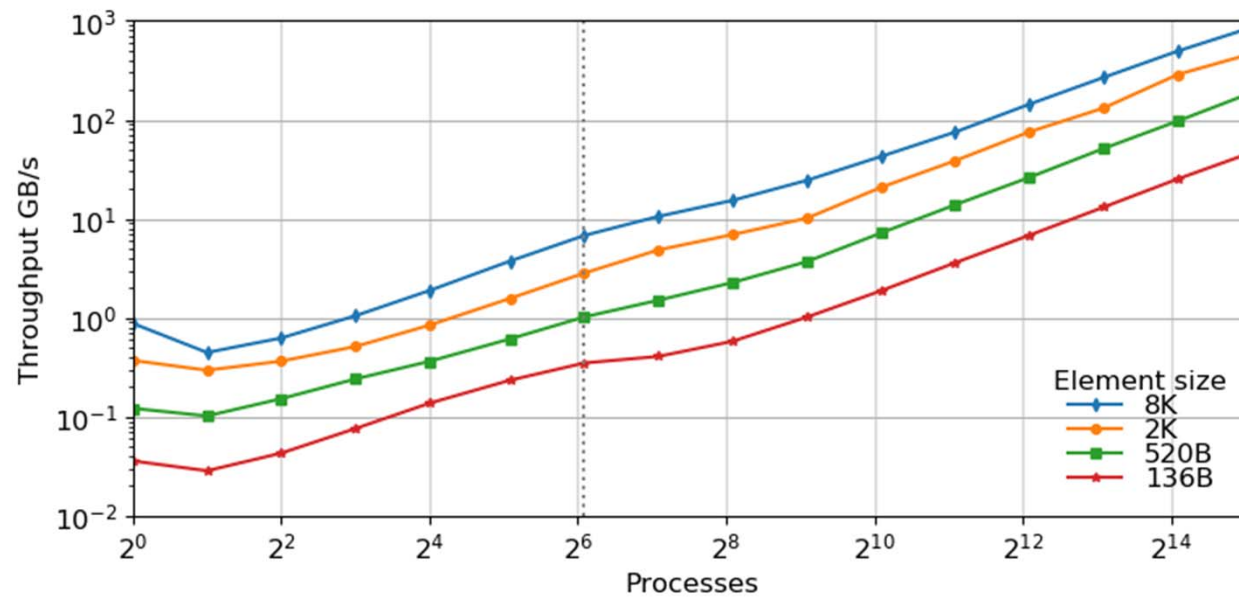
```
// C++ global variables correspond to rank-local state
std::unordered_map<uint64_t, global_ptr<char> > local_map;
// insert a key-value pair and return a future
future<> dht_insert(uint64_t key, char *val, size_t d_sz) {
    auto f1 = rpc(key % rank_n(), // RPC obtains location for the data
        [] (uint64_t key, size_t d_sz) -> global_ptr<char> {
            global_ptr<char> gptr = new_array<char>(d_sz);
            local_map[key] = gptr; // insert in local map
            return gptr;
        }, key, d_sz);
    return f1.then( // callback executes when RPC completes
        [val, d_sz] (global_ptr<char> loc) -> future<> { // λ: RMA put
            return rput(val, loc, d_sz); }
    );
}
```

λ function

λ for callback

# Benefits of UPC++: distributed hash table

- Randomly distributed keys
- Excellent weak scaling up to 32K cores
- RPC leads to simplified and more efficient design
  - Key insertion and storage allocation handled at target
  - Without RPC, complex updates would require explicit synchronization and the need to use global storage, e.g. OpenSHMEM and MPI one-sided

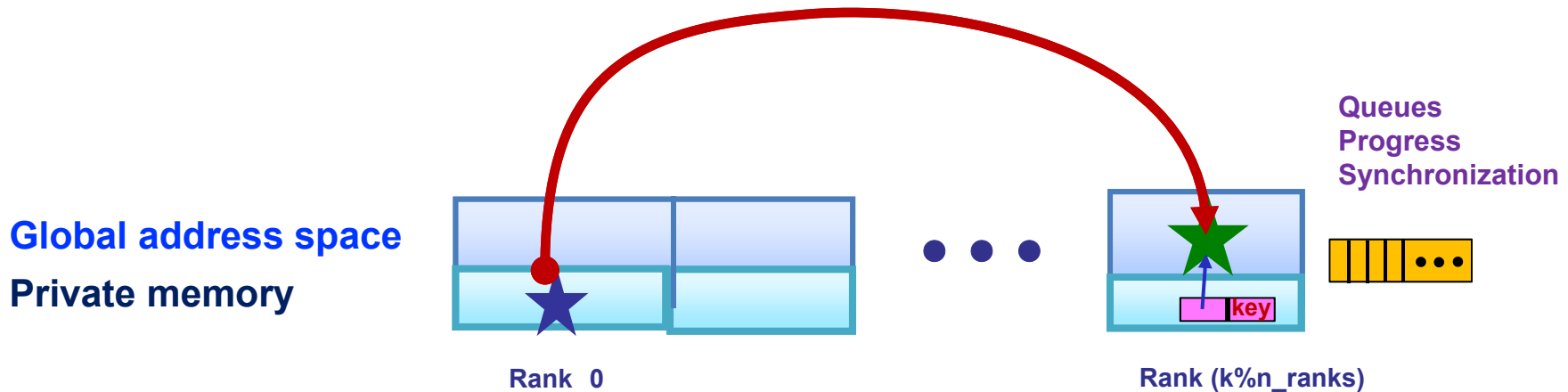


**Cori @ NERSC  
(KNL)**

**Cray XC40**

# The productivity benefit of RPC

- More natural to express hash table insertion with RPC than with one sided communication or message passing
- RPC encapsulates argument passing, queue management and progress, factoring them out of the application code
- More generally, RPC simplifies the coding in updating complicated distributed data structures



# Road Map

An application

Memory Kinds (GPU Memory)

Inside UPC++

Distributed data structures

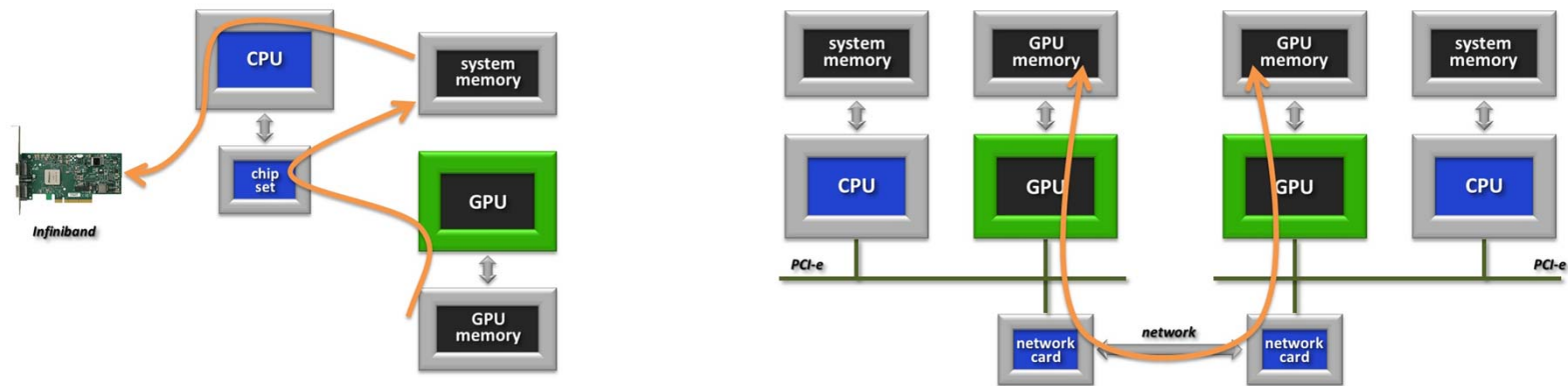
UPC++ in Context

# Memory kinds

In heterogeneous systems, we can have regions of memory requiring different access methods or performance properties (e.g. HBM, NVRAM)

Programming models can hide (or not hide) details

- MPI + CUDA (minimum of 2 additional copies)
- CUDA aware MPI (may or may not avoid 2 additional copies)

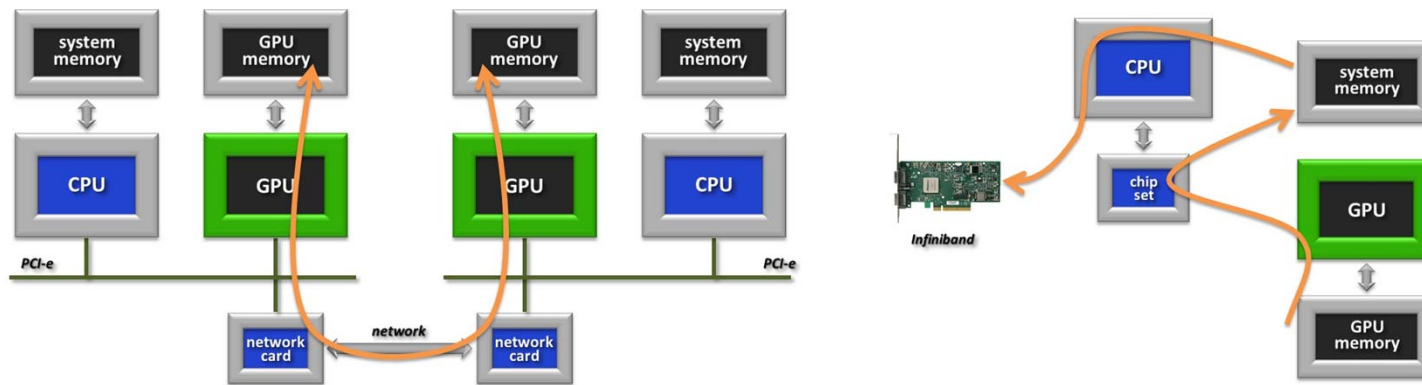


Figures courtesy Cheng, Grossman, and McKercher. *Professional Cuda C Programming*. John Wiley & Sons, 2014.

# Memory kinds

## UPCxx adds GPU memory kinds support (v2019.3.0)

- RMA transfers handled transparently between DRAM and Device Memory using global pointers
- Current implementation uses the extra copy model under the hood - so the user doesn't see the copying
- Future implementations (FY20) will use hardware support for direct communication



Figures courtesy Cheng, Grossman, and McKercher. *Professional Cuda C Programming*. John Wiley & Sons, 2014.

# The API

We are required to construct an allocator for GPU memory with a specified segment size

Global pointers can refer to host or device storage

A single `copy( )` method moves the data

```
auto gpu_device = cuda_device( 0 ); // Open device 0
// Construct and allocator with a 256 MB GPU segment
size_t segsize = 256*1024*1024;
auto gpu_alloc = device_allocator<cuda_device>(gpu_device, segsize);

// Use the allocator allocate an array of 1024 doubles on GPU
global_ptr<double,memory_kind::cuda_device> gpu_array =
gpu_alloc.allocate<double>(1024);

// Allocate an array of 1024 doubles on host
global_ptr<double> host_array = new_array<double>(1024);
// Transfer data from host buffer to device buffer
copy(host_array, gpu_array, 1024).wait();
```

# Road Map

An application

Memory Kinds (GPU Memory)

**Inside UPC++**

Distributed data structures

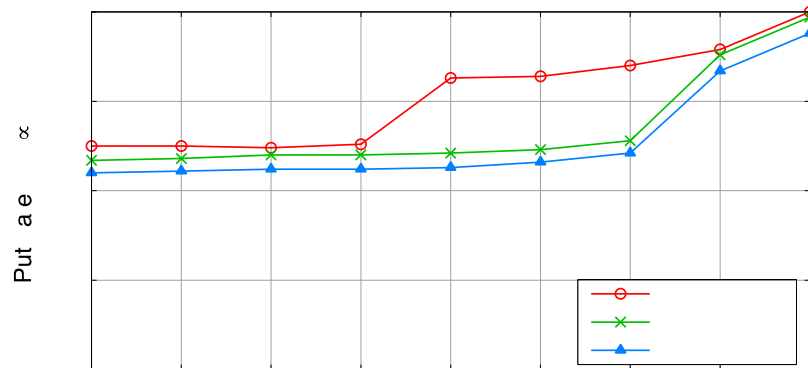
UPC++ in Context

# A look under the hood of UPC++

- Relies on GASNet-EX to provide low overhead communication
  - Efficiently utilizes the network, whatever that network may be, including any special purpose support - low overheads
  - Get/put map directly onto the network hardware's global address support, when available
- RPC uses an active message (AM) to enqueue the function handle remotely.
  - Any return result is also transmitted via an AM
- RPCs can only make progress inside a call to a UPC++ method (Also a distinguished progress() method)
  - Thus, RPCs are serialized at the target, and this attribute can be used to avoid explicit synchronization

# Performance of UPC++ - Latency

- Ping pong microbenchmark using blocking RMA
  - ~20% latency improvement from 16 to 256 bytes
  - Lines never cross at long message sizes not shown
- UPC++ rput, GASNet-Ex *testsmall* and publicly available MPI/IMB-RMA benchmark suite
  - Long sequence of blocking operations:  
    issue put, wait for remote completion, repeat...
  - Latency = average time



Lower is better

**Cori I @ NERSC  
(Haswell)**  
**Cray XC40**

# Road Map

An application

Memory Kinds (GPU Memory)

Inside UPC++

**Distributed data structures**

UPC++ in Context

# Toward distributed data structures

- Other models (UPC, CAF, OpenSHMEM) implement shared arrays via a symmetric heap
  - Scalable and portable implementation is difficult
  - Requires globally collective allocation that does not compose with subset teams
- UPC++ doesn't have a symmetric heap: heap allocation is not collective
- Initially, each rank only knows the locations of shared objects that it allocated
- How does a rank learn the locations of shared objects allocated by other ranks?

# Distributed objects in UPC++

- How does a rank learn the locations of shared objects allocated by other ranks?
- UPC++ provides the *distributed object* (like co-arrays)
  - Globally unique name for each distributed object
  - Each entry holds a rank-specific value
  - Retrieve a remote value using a rank ID
  - Can be used to build a scalable, globally visible directory
- Distributed objects can be used to build shared distributed arrays, among other data structures
- UPC++ does not prescribe solutions, rather it provides building blocks for constructing them

# Distributed 1D Arrays over UPC++

- Design is a work in progress
  - Likely similar to UPC pure-blocked shared array
  - Will support dynamic length and block size
  - Will be built over distributed objects, so it will have a scalable representation and support subset teams
- Exposes the cost of moving distributed data across the network. Compare with UPC, which supports off rank subscripting, but at a cost

```
dist_array<double> array(N, some_team);  
// fetch ith element of array  
future<global_ptr<double>> fut1 = array.pointer_to(i);  
future<double> fut2 = fut1.then(rget);  
// ... other work  
double val = fut2.wait();
```

# Road Map

An application

Memory Kinds (GPU Memory)

Inside UPC++

Distributed data structures

**UPC++ in Context**

# UPC++ in context

- Only existing library with PGAS support that also offers RPC (X10, Chapel and Habanero are languages)
- OpenSHMEM is considering adding RPC
- Besides DASH, only model that support subset teams
- UPC++ supports distributed objects, a generalization of distributed arrays
  - Can construct an object over a subset team
  - Avoids a symmetric heap, which is not scalable
  - Distributed arrays: UPC, OpenSHMEM, DASH, X10, Chapel, co-array C++ via co-arrays

# Persistent storage

- Idea: let's provide a unified global address space that spans storage and memories
- We map objects on disk into the shared segment
- Copy operators works with global pointers
- But not all persistent data can live in memory simultaneously, so we have to manage let the run time cache the data (custom paging)
- A useful operation is remote method invocation, which can leverage hardware offload support to carry out some filtering in the device controller:  
“find me all the data that meets a certain property”

# Summary

- UPC++ provides future and continuation-based completion handling; RPC; one-sided communication
  - Embraces communication networks that use one-sided transfers at their lowest level: low overhead reduces the cost of fine-grained communication
  - Combines explicit locality control with shared memory
  - Overlap communication via asynchronous execution
- Productivity
  - Develop incrementally, enhance code selectively
  - Interoperates with MPI, OpenMP and CUDA
- More advanced constructs (not discussed)
  - Remote atomics, Distributed objects, teams and collectives
  - Promises, personas (end points), generalized completion
  - Serialization, non-contiguous transfers

# The Pagoda Team

- Scott B. Baden (PI)
- Paul H. Hargrove (co-PI)
- John Bachan
- Dan Bonachea
- Mathias Jacquelin
- Amir Kamil
- Hadia Ahmed
- Alumni:  
Brian van Straalen,  
Steve Hofmeyr, Khaled Ibrahim



Code and documentation at <http://upcxx.lbl.gov>

# Acknowledgements

Early work with UPC++ involved Yili Zheng, Amir Kamil, Kathy Yelick, and others [IPDPS '14]

This research was supported in part by the Exascale Computing Project (17-SC-20-SC), funded by the U.S. Department of Energy

ECP collaborators: Kathy Yelick, Sherry Li, Pieter Ghysels, John Bell and Tan Nguyen (Lawrence Berkeley National Laboratory)

Academic collaborators: Alex Pöpl and Michael Bader (TUM) ,  
Niclas Jansson and Johann Hoffman (KTH),  
Sergio Martin (ETH-Z)

<http://upcxx.lbl.gov>

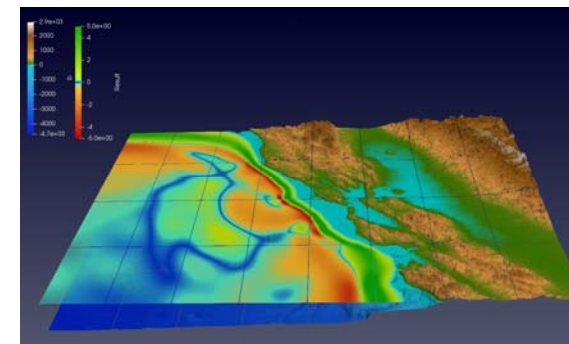


Figure courtesy Alexander Pöpl