



Chapel Aggregation Library (CAL)

November 12, 2018

Louis Jenkins

Marcin Zalewski (Pacific Northwest National Lab.),
Michael Ferguson (Cray Inc.)

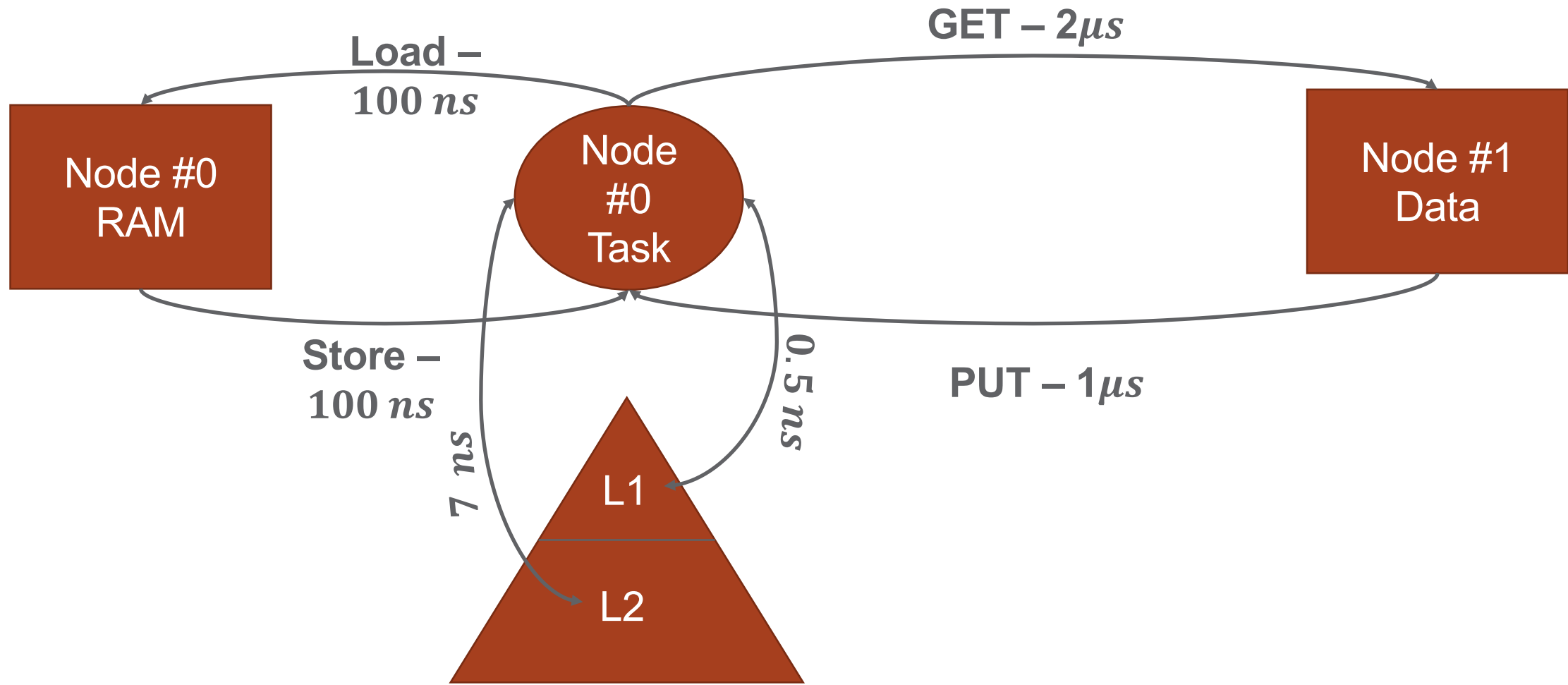


PNNL is operated by Battelle for the U.S. Department of Energy



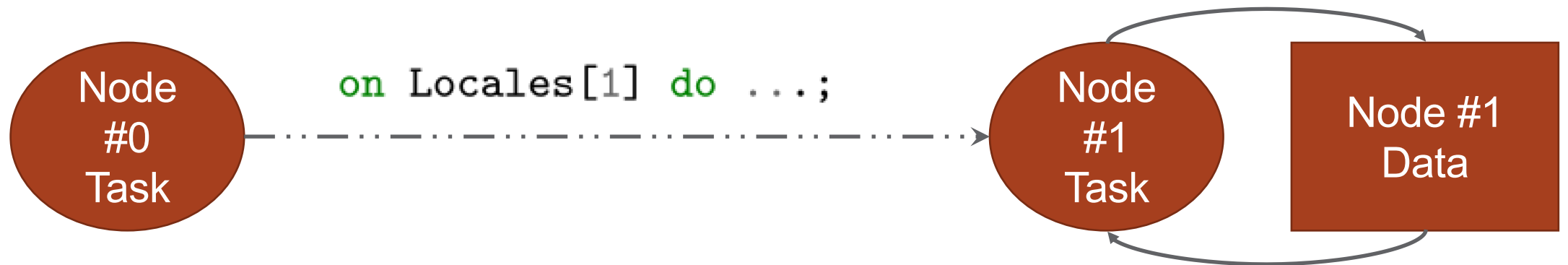
The Problem

- Accessing remote data is slow
 - Multiple orders of magnitude slower to access than local memory



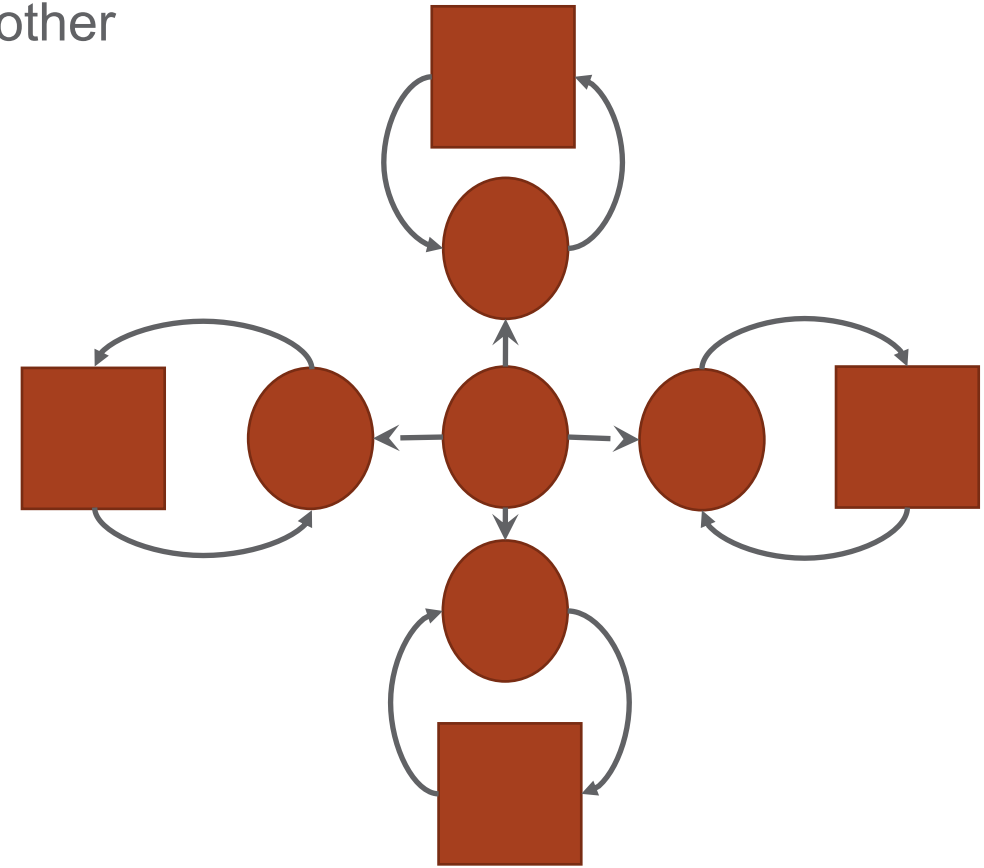
The Problem

- Accessing remote data is slow
 - Multiple orders of magnitude slower to access than local memory
- “Moving the computation to the data” not always the best solution
 - Using an *on* statement requires migrating tasks to another locale



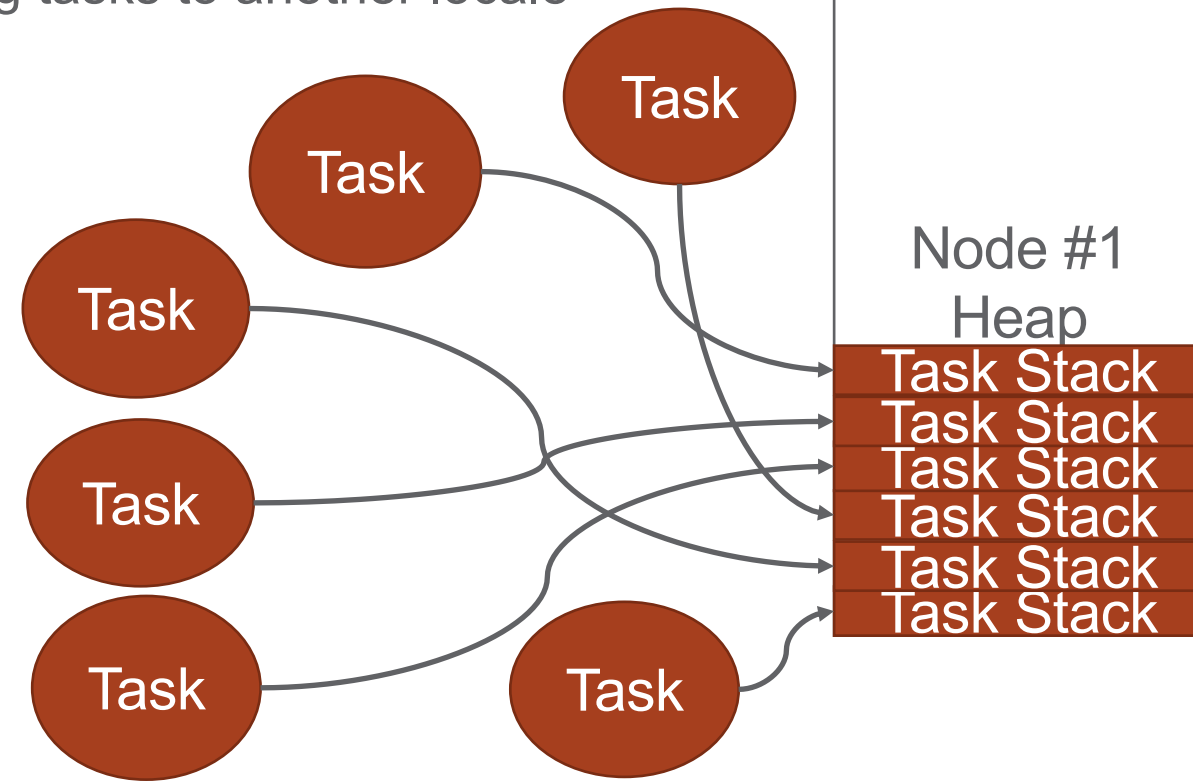
The Problem

- Accessing remote data is slow
 - Multiple orders of magnitude slower to access than local memory
- “Moving the computation to the data” not always the best solution
 - Using an *on* statement requires migrating tasks to another locale
 - ✓ Can become bottleneck if fine-grained



The Problem

- Accessing remote data is slow
 - Multiple orders of magnitude slower to access than local memory
- “Moving the computation to the data” not always the best solution
 - Using an *on* statement requires migrating tasks to another locale
 - ✓ Can become bottleneck if fine-grained
 - ✓ Task creation is relatively expensive
 - Tasks are too large to spawn in a fire-and-forget manner (issue #9984)
 - Migrating tasks require individual active messages (issue #9727)



A Solution

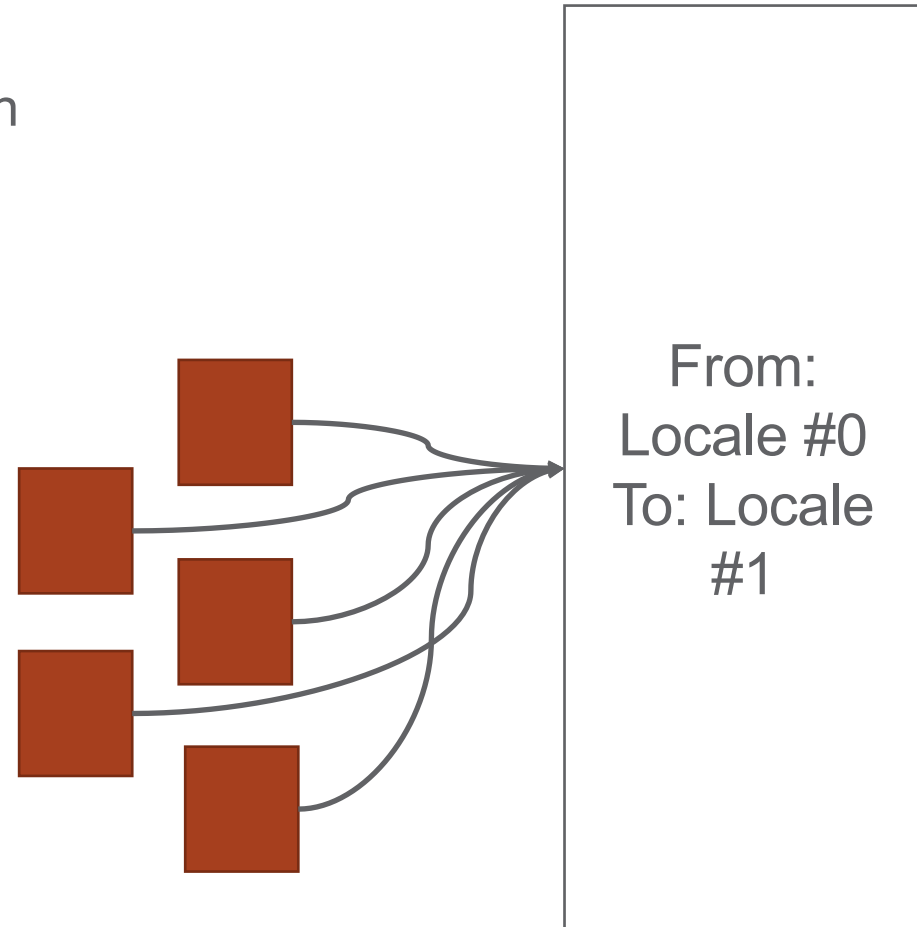
- Coarsen the granularity of the data
 - Buffer units of data to be sent to a locale in *destination buffers*



From:
Locale #0
To: Locale
#1

A Solution

- Coarsen the granularity of the data
 - Buffer units of data to be sent to a locale in *destination buffers*



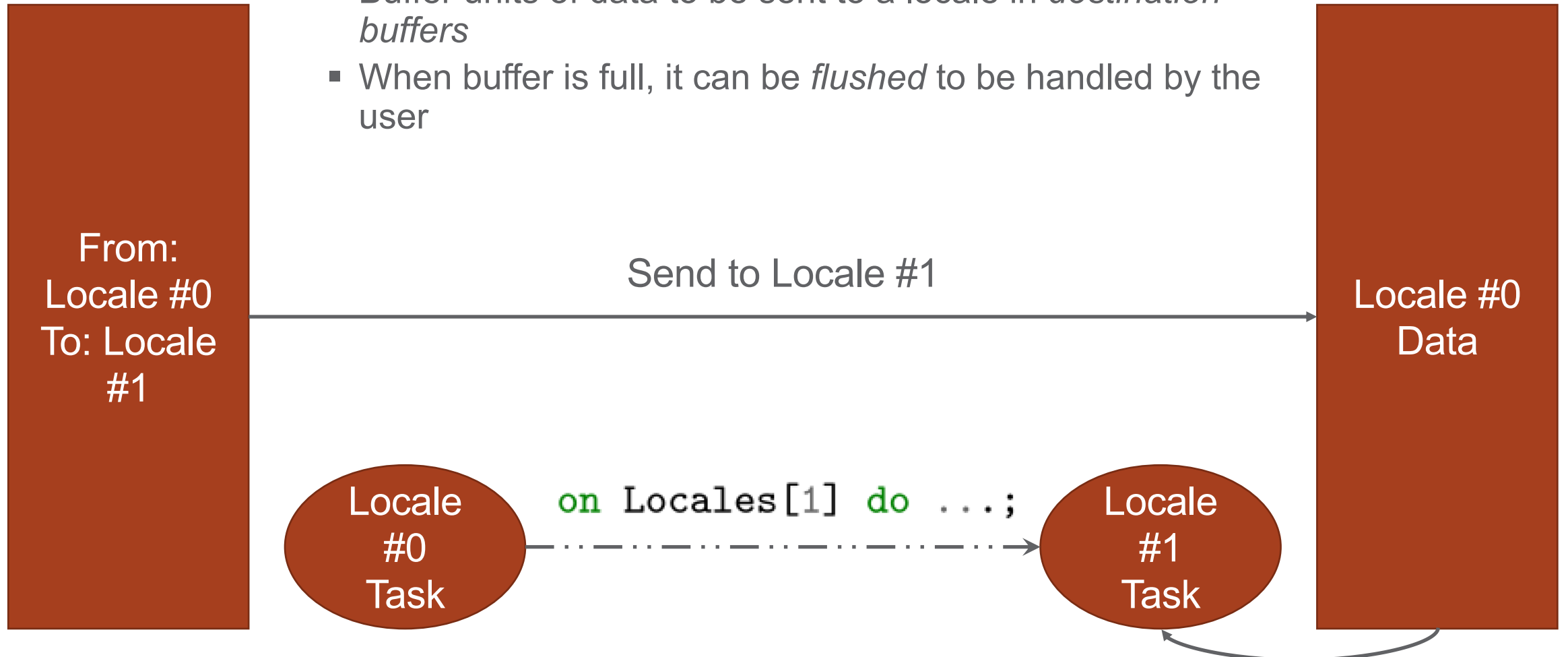
A Solution

- Coarsen the granularity of the data
 - Buffer units of data to be sent to a locale in *destination buffers*

From:
Locale #0
To: Locale
#1

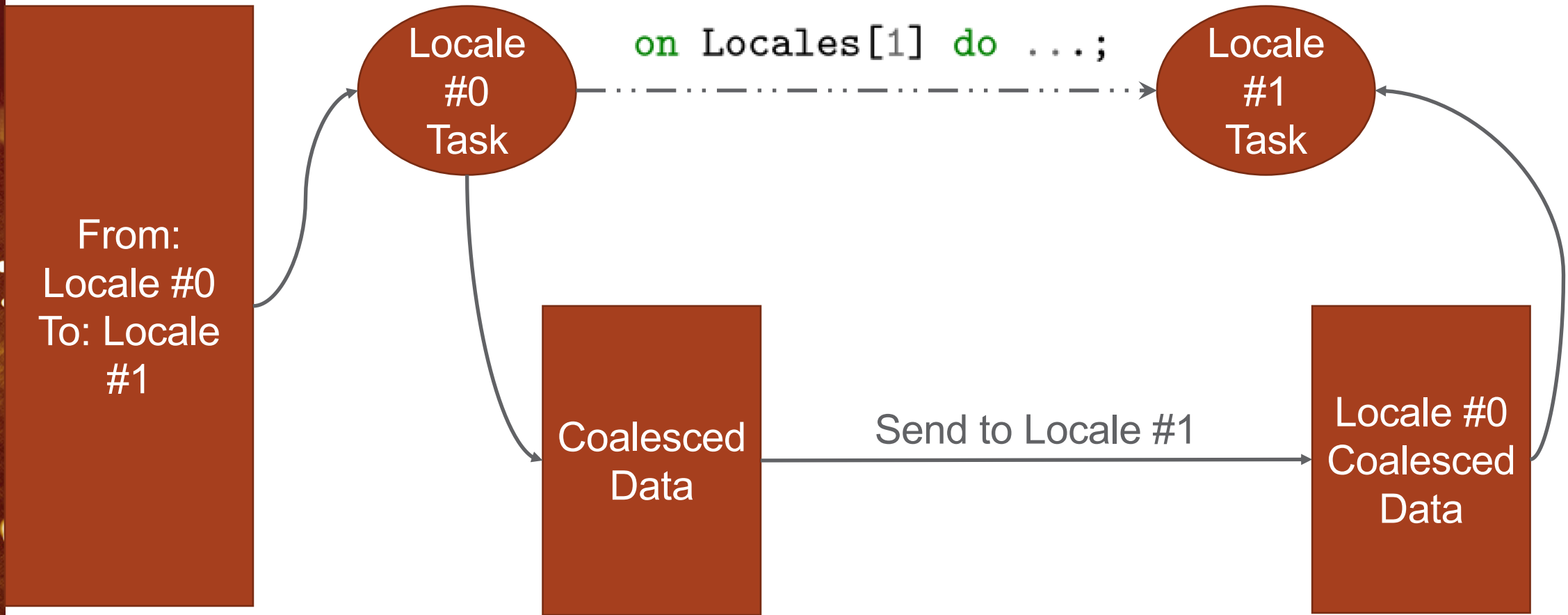
A Solution

- Coarsen the granularity of the data
 - Buffer units of data to be sent to a locale in *destination buffers*
 - When buffer is full, it can be *flushed* to be handled by the user



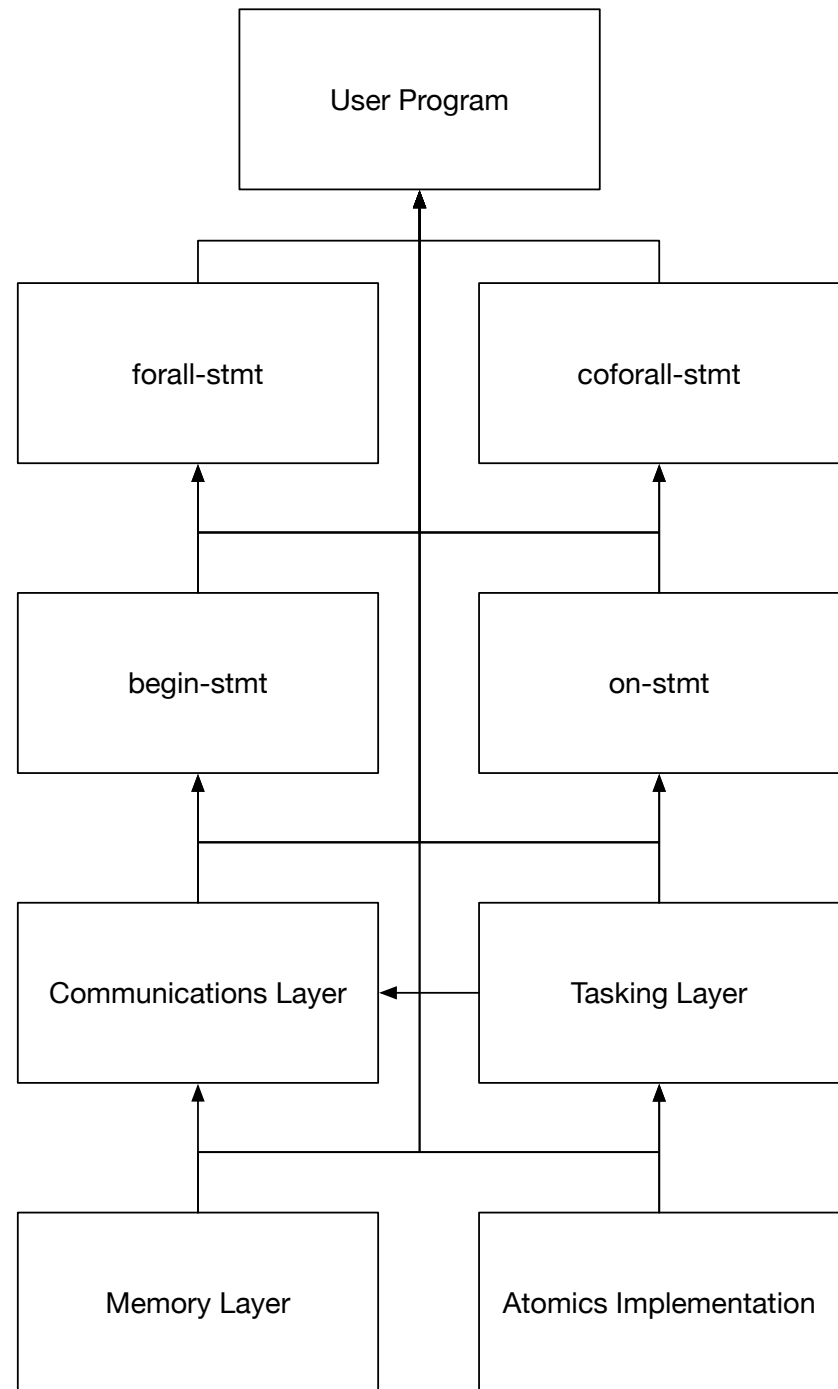
A Solution

- Coarsen the granularity of the data
 - Buffer units of data to be sent to a locale in *destination buffers*
 - When buffer is full, it can be *flushed* to be handled by the user
 - User can perform *coalescing* to combine aggregated data



Chapel's Multiresolution Design Philosophy

- Higher Level composed of Lower Level abstractions, features, and language constructs
 - Changes to lower level propagate up to higher level
 - User free to use either
 - ✓ High-Level for convenience
 - ✓ Low-Level for performance



Global-View Programming

- Abstracts locality for the user
 - No need to think: “What portion of the array does this task own?”
 - Array can be accessed from any locale, even if it is not distributed over that locale...
 - ✓ Remote references are resolved into remote PUT/GET *implicitly*

Chapel

```
1  var sum : float;  
2  forall a in arr with (+ reduce sum) {  
3    sum += a;  
4  }
```

MPI

```
1  float globalSum = 0;  
2  float localSum = 0;  
3  for (int i = localStart; i < localEnd; i++) {  
4    localSum += arr[i];  
5  }  
6  MPI_REDUCE(&localSum, &globalSum, ...);
```

Global-View Programming

- Abstracts locality for the user
 - No need to think: “What portion of the array does this task own?”
 - Array can be accessed from any locale, even if it is not distributed over that locale...
 - ✓ Remote references are resolved into remote PUT/GET *implicitly*
- Multiresolution: More Abstraction

Chapel

```
1  var sum = + reduce arr;
```

MPI

```
1  float globalSum = 0;  
2  float localSum = 0;  
3  for (int i = localStart; i < localEnd; i++) {  
4      localSum += arr[i];  
5  }  
6  MPI_REDUCE(&localSum, &globalSum, ...);
```


Global-View Programming

- Abstracts locality for the user
 - No need to think: “What portion of the array does this task own?”
 - Array can be accessed from any locale, even if it is not distributed over that locale...
 - ✓ Remote references are resolved into remote PUT/GET *implicitly*
- Multiresolution: Less Abstraction

Chapel

```
1 var sum : float;
2 coforall loc in Locales with (+ reduce sum) do on loc {
3   coforall tid in 0..#here.maxTaskPar with (+ reduce sum) {
4     for i in computeRange(arr.domain.localSubdomain(), tid) {
5       sum += arr[i];
6     }
7   }
8 }
```

MPI

```
1 float globalSum = 0;
2 float localSum = 0;
3 for (int i = localStart; i < localEnd; i++) {
4   localSum += arr[i];
5 }
6 MPI_REDUCE(&localSum, &globalSum, ...);
```

Chapel Aggregation Library (CAL)

- Written in Chapel, for Chapel
 - **Minimal** and User-Friendly
 - ✓ Unassuming of how data is handled
 - ✓ Designed specifically for Chapel
 - **Distributed, Scalable, and Parallel-Safe**
 - ✓ Supports Global-View Programming
 - ✓ Usable with Chapel's parallel and locality constructs
 - **Modular, Reusable, and Generic**
 - ✓ Generic on user-defined type
 - ✓ Easy to use and 'plug in'

Minimalism

- CAL is an aggregation library
 - Processing of the aggregated data is deferred to the user
 - Buffer is returned to the last task that filled it

```
1  const msg = "From Locale#0 to Locale#1";
2  const loc = Locales[1];
3  var aggregator = new Aggregator(string);
4  var buffer = aggregator.aggregate(msg, loc);
5  if buffer != nil then handleBuffer(buffer);
6  [(buf, loc) in aggregator.flush()] on loc do handleBuffer(buf);
```

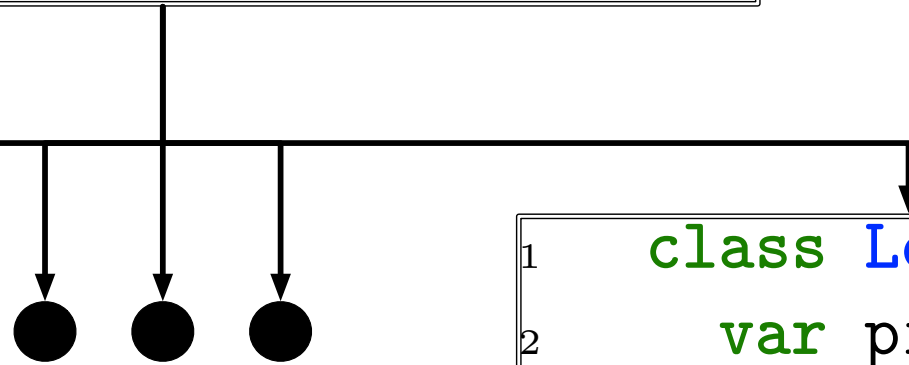
Distributed Object Pattern

- Use privatization to enable global-view programming
 - GlobalClass forwards access to per-locale LocalClass *privatized instances*
 - Each privatized instance can communicate and coordinate with others

```
1  pragma "always RVF"  
2  record GlobalClass {  
3      type classType;  
4      var pid : int;  
5  
6      forwarding chpl_getPrivatizedCopy(pid, classType);  
7  }
```

```
1  class LocalClass {  
2      var pid : int;  
3  }
```

Locale#0



```
1  class LocalClass {  
2      var pid : int;  
3  }
```

Locale#N

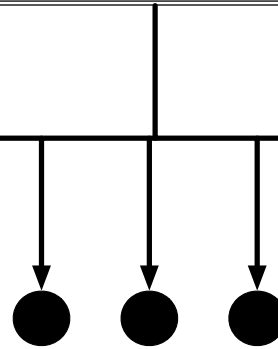
Aggregator

- Aggregator forwards all accesses to per-locale privatized instances
- Distributed and parallel access is abstracted
 - Supports global-view programming

```
1  pragma "always RVF"  
2  record Aggregator {  
3    type bufType;  
4    var pid : int;  
5  
6    forwarding chpl_getPrivatizedCopy(pid, bufType);  
7  }
```

```
1  class LocalBuffer {  
2    type t;  
3    var pid : int;  
4    var buffers : [0..#numLocales] BufferPool(t);  
5  }
```

Locale#0

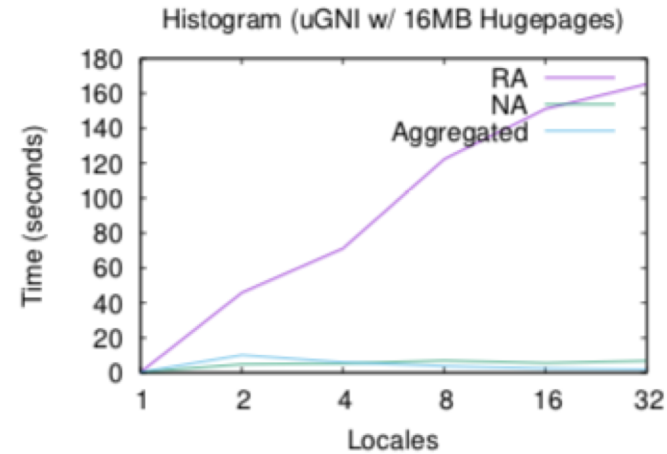


```
1  class LocalBuffer {  
2    type t;  
3    var pid : int;  
4    var buffers : [0..#numLocales] BufferPool(t);  
5  }
```

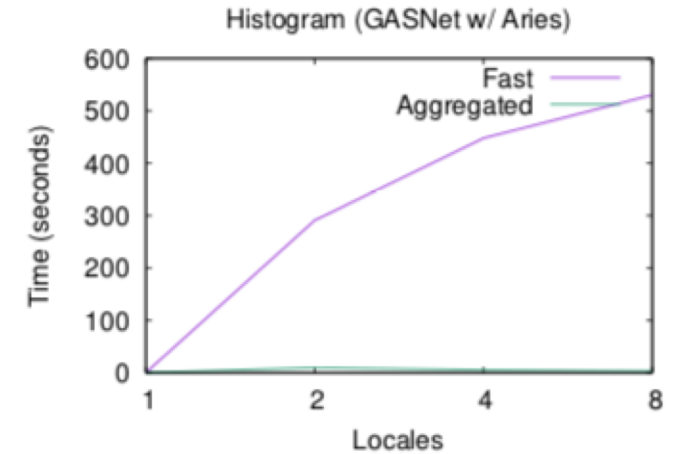
Locale#N

Aggregator - Performance

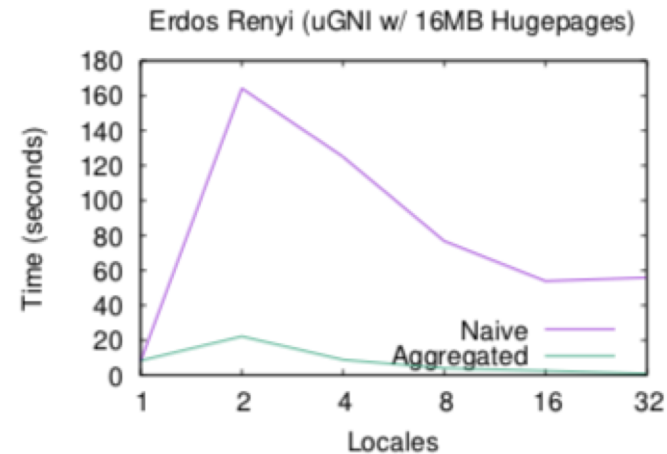
- 10x – 100x speedup at 32 nodes
 - Histogram
 - Hypergraph Generation



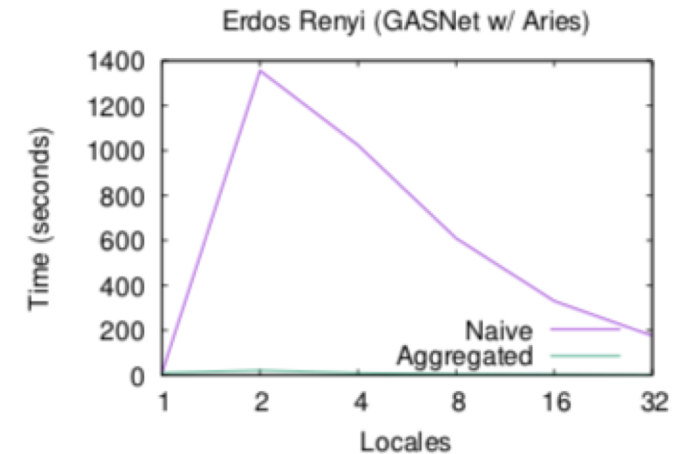
(a) uGNI Histogram



(b) GASNet Histogram



(a) uGNI Erdos Renyi



(b) GASNet Erdos Renyi

Distributed - Example

- Aggregator is allocated on Locale#0, but accessible from Locale#1
 - Accesses are forwarded to Locale#1's privatized instance
 - Global-View Programming
- Implicit parallelism (line 9) vs Explicit parallelism (line 11)

```
1  var aggregator = new Aggregator(int);
2  // Migrate to Locale #1 from Locale #0
3  on Locales[1] {
4      // Aggregate single value to Locale #0
5      var buffer = aggregator.aggregate(0, Locales[0]);
6      // If non-nil, then handle buffer.
7      if buffer != nil then handleBuffer(buffer);
8      // Aggregate multiple units of data via Chapel's implicit parallelism
9      var buffers = aggregator.aggregate(1..1024, Locales[0]);
10     // Check if any of the buffers are nil
11     [buf in buffers] if buf != nil then handleBuffer(buf);
12 }
```

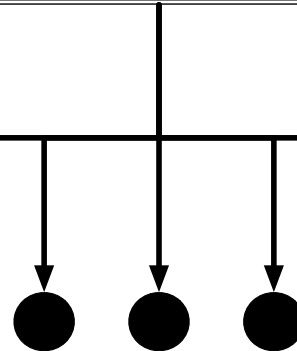
Modularity

- Composition of Distributed Objects
 - Aggregator can be used within other global-view data structures
 - Future of Distributed Object Oriented Programming (?)

```
1  pragma "always RVF"  
2  record GlobalClass {  
3      type classType;  
4      var pid : int;  
5  
6      forwarding chpl_getPrivatizedCopy(pid, classType);  
7  }
```

```
1  class LocalClass {  
2      type t;  
3      var pid : int;  
4      var aggregator : Aggregator(t);  
5  }
```

Locale#0



```
1  class LocalClass {  
2      type t;  
3      var pid : int;  
4      var aggregator : Aggregator(t);  
5  }
```

Locale#N

Future Works

- Software release of CAL
 - Currently only available as module under Chapel HyperGraph Library (CHGL)
 - ✓ github.com/pnnl/chgl
 - Independent release coming soon (?)
- Integration into Chapel
 - Mason package or Standard Module (?)
 - Run-time integration
- Aggregation handlers as first-class functions
 - Once Chapel has better first-class function support

Potential Application Light Weight Tasks (LWT)

- Chapel Tasks are infeasible to use in fire-and-forget manner
 - Stack size of tasks in Chapel are static and large (8MB default)
 - Task migration can be made asynchronous but is not aggregated
- Solution – Make a library for LWT
 - Use Distributed Object pattern for Global-View programming
 - Use *Aggregator* for aggregation
 - Use First-Class Functions (once improved) to represent a lightweight task

```
1  var lwt = new LWT(visit);
2  proc visit(v : Vertex) {
3      for vv in neighbors(v) {
4          if hasProperty(vv) {
5              lwt.spawn(vv, vv.locale);
6          }
7      }
8  }
9  forall v in vertices {
10     if hasProperty(v) {
11         lwt.spawn(v);
12     }
13 }
```

Vertex Degree Distribution

```
1 // Find largest degree of all vertices in distributed graph
2 var N = max reduce [v in graph.getVertices()] graph.degree(v);
3 // Histogram is cyclically distributed over all locales
4 var histogramDomain = {1..N} dmapped Cyclic(startIdx=1);
5 var histogram : [histogramDomain] atomic int;
6
7 // Aggregate increments to histogram
8 var aggregator = new Aggregator(int);
9 forall v in graph.getVertices() {
10     const deg = graph.degree(v);
11     const loc = histogram[deg].locale;
12     var buffer = aggregator.aggregate(deg, loc);
13     if buffer != nil {
14         on loc do [deg in buffer] histogram[deg].add(1);
15         buffer.done();
16     }
17 }
18
19 // Flush
20 forall (buf, loc) in aggregator.flush() {
21     on loc do [deg in buf] histogram[deg].add(1);
22     buffer.done();
23 }
```