



# What's New in Proxy Apps

ECP Annual Meeting 2019

David Richards  
LLNL

Houston, TX  
15 January 2019

# ECP Proxy App Project: Objectives and Scope

- Assemble and curate a proxy app suite that represents the most important features (especially performance) of exascale applications.
- Improve the quality of proxies created by ECP and maximize the benefit received from their use.
  - Set standards for documentation, build and test systems, performance models, etc.
- Assess the representativeness of proxies
- Project team members at LLNL, LANL, SNL, LBNL, ORNL, ANL



## Over 50 proxy apps can be installed using Spack

- Builds multiple versions, configurations, and compilers
- Encourages “standard” build recipes that conform to common build systems (e.g. Autotools, CMake, Makefile)
- Ease of aggregation
  - `spack fetch <package_name>` to just download the vanilla source
  - `spack install --source <package_name>` to just download the possibly patched source
- Proxies can be build as a whole or separately the “Spack” way
  - Different compilers and with or without system libraries
- Currently using Spack for continuous integration



# Spack



# We maintain a catalog of over 50 proxy apps

The screenshot shows the 'Catalog' page of the proxy apps website. It features a search bar at the top left and a row of filter buttons for programming languages and frameworks: All, Fortran, MPI, C++, CUDA, OpenMP, Python, C, and C and MPI/SHMEM. Below the filters is a grid of application cards, each with a title, description, and status icons. The visible cards are:

- AMG | C**: AMG is a parallel algebraic multigrid solver for linear systems arising from problems on unstructured grids. Status: two green checkmarks.
- ASPA | C++**: ASPA Proxy Application, Multi-scale, adaptive sampling, materials science proxy. Status: two green checkmarks and a person icon.
- CANDLE Benchmarks | Python**: CANDLE Benchmarks. Status: two green checkmarks.
- CloverLeaf | Fortran**: A miniapp that solves the compressible Euler equations on a Cartesian grid. Status: two green checkmarks.
- CLAMR | C++**: CLAMR is a cell-based adaptive mesh refinement mini-app developed as a testbed for hybrid algorithm development using MPI and OpenCL. Status: one green checkmark.
- Chatterbug | C++, MPI**: A suite of communication-intensive proxy applications that mimic commonly found communication patterns in HPC codes. These codes can be used as synthetic codes for benchmarking, or for trace generation using OTF2. Status: two green checkmarks and a person icon.
- CloverLeaf3D | Fortran**: 3D version of a miniapp that solves compressible Euler equations on a Cartesian grid. Status: two green checkmarks.
- CoGL | C++**: Analyzes pattern formation in ferroelastic materials and tests in situ visualization. Status: two green checkmarks and a person icon.
- CoHMM | C**: A proxy application for the Heterogeneous Multiscale Method (HMM) augmented with...

- Work underway on collecting/distributing multiple versions of proxies
  - Issues with maintainability, consistency, releasability, etc.
- Contributions are welcome.
  - Submit your proxy information through the proxy app website

<https://proxyapps.exascaleproject.org>

## We released version 2.0 of the ECP Proxy App Suite in Oct 2018

Proxy	Language	GPU		Proxy	Language	GPU
AMG	C	No public version		miniVite	C++	
CANDLE Benchmarks	Python			NEKbone	Fortran	OpenACC, OpenACC+CUDA
Ember	MPI/SHMEM	N/A		PICSARlite	Fortran	
ExaMiniMD	C++	Kokkos		SW4lite	C, C++, Fortran	CUDA, RAJA
Lagos	C++	CUDA, RAJA, OCCA		SWFFT	C	
MACSio	C	N/A		thornado-mini	Fortran	
miniAMR	C	No public version		XSbench	C	OpenACC, OpenCL, AMP
miniQMC	C++	Kokkos OpenMP				

## Acknowledgements

Slides & information for discussion of new mini-apps provided by:

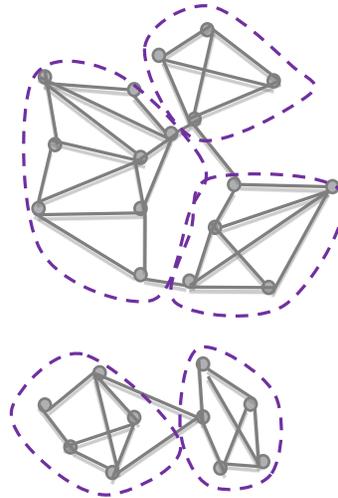
- miniVite: Mahantesh Halappanavar, Sayan Ghosh
- miniQMC: Paul Kent
- thornado-mini: Bronson Messer
- PICSARlite: Jean-Luc Vay
- Ember: Si Hammond

## Graph Clustering (community detection)

- Problem: Given  $G(V, E, \omega)$ , identify tightly knit groups of vertices that **strongly** correlate to one another within their group, and **sparsely** so, outside.

### Input :

- $V = \{1, 2, \dots, n\}$
- $E$ : a set of  $M$  edges
- $\omega(e)$ : weight of edge  $e$   
(non-negative)
- $m = \sum_{\forall e \in E} \omega(e)$

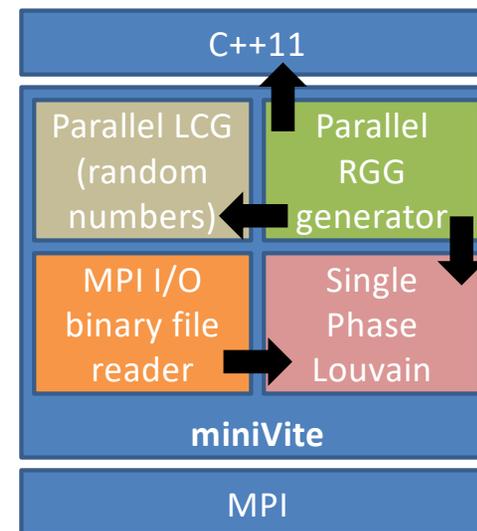


### Output :

- A **partitioning** of  $V$  into  $k$  **mutually disjoint** clusters  $P = \{C_1, C_2, \dots, C_k\}$  such that: ... ?

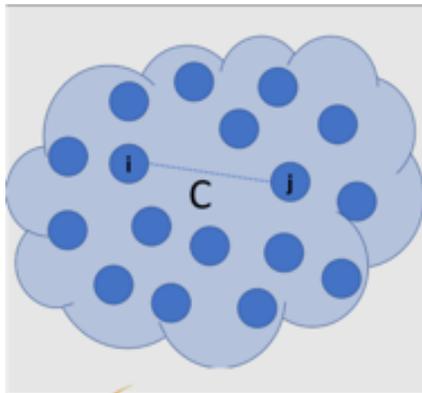
# miniVite (/ˈvi:te/)

- Implements a single phase of Louvain method (without rebuilding the graph)
- Capable of generating synthetic Random Geometric Graphs (RGG) in parallel (needs random numbers)
  - Can also add random edges across processes
- Can also use real world graphs as input (have to convert to a binary format first)
- Parts of code has multiple communication options (can be selected at compile time) – Sendrecv, NB Isend/Irecv (default), MPI RMA and Collectives
- About 3K LoC



# Graph clustering OR Community detection

- **Communities:** Nodes/vertices of most real-world network/graphs tend to be organized into tightly-knit modules known as *communities* or *clusters*
- In 2008, Blondel, et al. introduced a multi-phase, iterative heuristic for modularity optimization, called the **Louvain method**
- Goodness of partitioning into communities is typically measured using a global metric called **modularity**
  - Depends on sum of edge weights between vertices within a community ( $e_{ij}$ ), and sum of weights incident upon a community ( $a_c$ )



The diagram shows a cluster of blue circular nodes enclosed in a light blue cloud-like boundary. Two nodes are labeled 'i' and 'j', and they are connected by a thin blue line representing an edge. The entire cluster is labeled 'C'.

$$Q = \sum_{c \in C} \left[ \frac{e_{ij}}{2m} - \left( \frac{a_c}{2m} \right)^2 \right]$$

where:

$$e_{ij} = \sum w_{ij} : \forall i, j \in c, \text{ and } \{i, j\} \in E$$
$$a_c = \sum_{i \in c} k_i$$

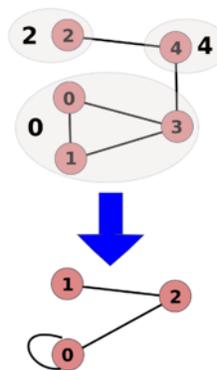
# Louvain algorithm

**Input:** Graph  $G = (V, E)$ , threshold  $\tau$ , Initial community assignment,  $C_{init}$

```

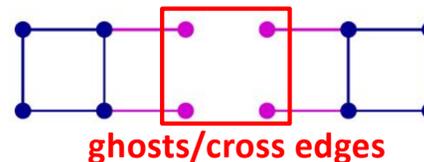
1:  $Q_{prev} \leftarrow -\infty$ 
2:  $C_{prev} \leftarrow$  Initialize each vertex in its own community
3: while true do
4:   for all  $v \in V$  do
5:      $N(v) \leftarrow$  neighboring communities of  $v$ 
6:      $targetComm \leftarrow \arg \max_{t \in N_v} \Delta Q(v \text{ moving to } t)$ 
7:     if  $\Delta Q > 0$  then
8:       Move  $v$  to  $targetComm$  and update  $C_{curr}$ 
9:    $Q_{curr} \leftarrow ComputeModularity(V, E, C_{curr})$ 
10:  if  $Q_{curr} - Q_{prev} \leq \tau$  then
11:    break
12:  else
13:     $Q_{prev} \leftarrow Q_{curr}$ 
  
```

- Initially each vertex assigned to a separate community
- Within each iteration
  - Compute delta Q when a vertex migrates
  - Move vertex from current community to one that yields max delta Q
- Phase continues until delta Q between successive iterations is below a threshold
- At the end of each phase, graph is rebuilt



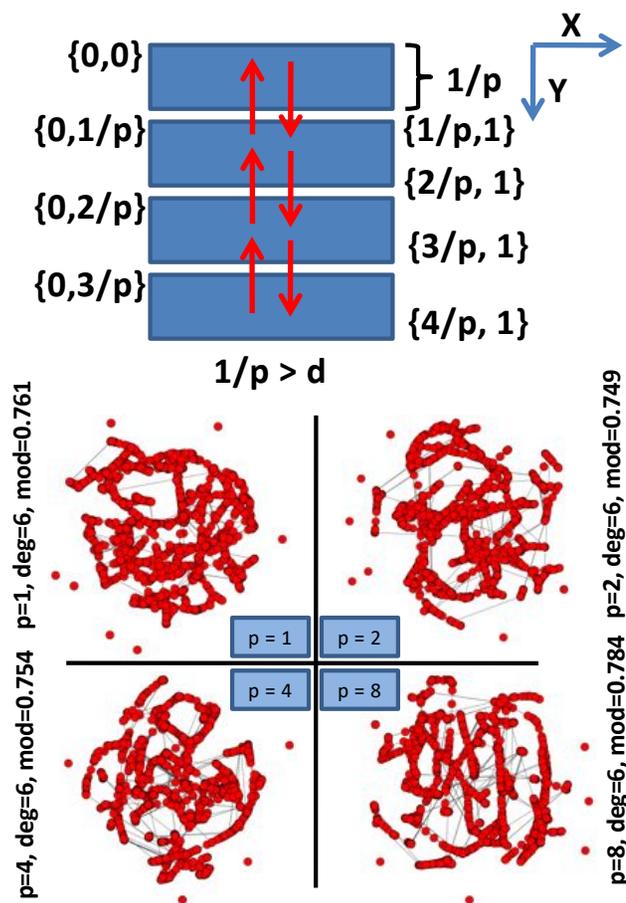
**Rebuilding is nontrivial, may takes ~1-10% of the overall phase time depending on input graph and volume of vertex movement**

**Simple vertex-based data distribution**



# Parallel Random Geometric Graph (RGG) generator

- Constructed by randomly placing  $N$  nodes within a unit square – only add an edge between two vertices if their distance is within a range  $d$ 
  - Basically, generate random numbers within  $(0,1)$  and compute euclidean distance between two points, if it is less than  $d$ , add an edge
  - RGG is known to demonstrate community structure (high modularity)
- We distribute equal number of vertices across processes, each process may have (cross) edges with its up or down neighbor
- Option to add random number of (cross) edges across processes that are far apart



# Results

FIRST PHASE OF LOUVAIN METHOD VS THE LAST PHASE FOR REAL WORLD INPUTS ON 1K PROCESSES OF NERSC CORI.

Graphs	#Vertices	#Edges	Modularity	First phase		Final		
				Iterations	Time (in s)	Phases	Iterations	Time (in s)
friendster	65.6M	1.8B	0.624	143	565.201	3	440	567.173
it-2004	41.3M	1.15B	0.973	14	45.064	4	91	45.849
nlpkkt	27.9M	401.2M	0.939	3	3.57	5	832	21.084
sk-2005	50.6M	1.9B	0.971	11	71.562	4	83	72.94
orkut	3M	117.1M	0.472	89	59.5	3	281	59.64
sinaweibo	58.6M	261.3M	0.482	3	270.254	4	108	281.216
twitter-2010	21.2M	265M	0.478	3	209.385	4	184	386.483
uk2007	105.8M	3.3B	0.972	9	35.174	6	139	37.9887
web-cc12-pay/vladmin	42.8M	1.2B	0.687	31	140.493	4	159	146.92
webbase-2001	118M	1B	0.983	14	14.702	7	239	24.455

For a number of real world graphs, first phase of Louvain method does the most work (little difference between first and final phase)

Processes	NB Send/Recv	Collectives	Sendrecv	RMA
512	7.48492	7.35221	11.2029	7.39827
1024	6.52832	5.56177	13.2942	5.93101

RGG of 134.2M vertices and 1.6B edges on NERSC Cori

# Performance analysis (I)

main	6.40e+11	100.0%
loop at main.cpp: 212	6.34e+11	99.0%
230: distLouvainMethod(int, int, DistGraph const&, std::vector<long, std::allocator<long> >&, double, double)	6.29e+11	98.2%
loop at distLouvainMethodNew.cpp: 47	6.24e+11	97.5%
88: [I] _INTERNAL_24_distLouvainMethodNew_cpp_2c8537cd::distComputeModularity(Graph const&, std::vector<long, std::allocator<long> >&, double, double)	2.67e+11	41.6%
58: _INTERNAL_24_distLouvainMethodNew_cpp_2c8537cd::fillRemoteCommunities(DistGraph const&, int, int, std::vector<long, std::allocator<long> >&)	2.16e+11	33.7%
68: __kmpc_fork_call	1.38e+11	21.6%
125: _INTERNAL_24_distLouvainMethodNew_cpp_2c8537cd::updateRemoteCommunities(DistGraph const&, std::vector<long, std::allocator<long> >&, int, int)	3.30e+09	0.5%
106: __kmpc_fork_call	5.62e+07	0.0%
46: _INTERNAL_24_distLouvainMethodNew_cpp_2c8537cd::exchangeVertexReqs(DistGraph const&, int, int)	3.84e+09	0.6%
138: [I] std::unordered_map<long, long, std::hash<long>, std::equal_to<long>, std::allocator<std::pair<long, long> > >::operator[](const long&)	2.63e+08	0.0%
34: _INTERNAL_24_distLouvainMethodNew_cpp_2c8537cd::distInitLouvain(DistGraph const&, std::vector<long, std::allocator<long> >&, double, double)	2.46e+08	0.0%
245: distbuildNextLevelGraph(int, int, DistGraph*&, std::vector<long, std::allocator<long> >&)	5.06e+09	0.8%
145: loadDistGraphMPIO(int, int, DistGraph*&, std::string&)	6.11e+09	1.0%
365: MPI_Finalize	2.64e+08	0.0%

1. HPCToolkit profiling shows over 60% of time spent in managing and communicating vertex-community information
2. About 40% is spent on global communication (MPI\_Allreduce) for computing modularity

# QMCPACK <https://qmcpack.org>

Developed and released via  
[github.com/QMCPACK/qmcpack](https://github.com/QMCPACK/qmcpack)

Pull request reviews, continuous  
integration & integration testing.

UIUC/NCSA open source license.

~330K code lines, <100K “core”

C++11 (14 soon), MPI, HDF5, XML, FFTW.

OpenMP (walker level only) and CUDA  
(old) on node. **Different vectorization /  
parallelization scheme for GPU & CPU.**  
Major development + maintenance  
headache.

Miniapps: miniqmc, miniafqmc under  
[github.com/QMCPACK](https://github.com/QMCPACK) . Use for  
prototyping, experimenting with individual  
kernels etc.

IOP Publishing Journal of Physics: Condensed Matter  
J. Phys.: Condens. Matter **30** (2018) 195901 (29pp) <https://doi.org/10.1088/1361-648X/aab9c3>

## QMCPACK: an open source *ab initio* quantum Monte Carlo package for the electronic structure of atoms, molecules and solids

Jeongnim Kim<sup>1</sup>, Andrew T Baczewski<sup>2</sup>, Todd D Beaudet<sup>3</sup>, Anouar Benali<sup>4,5</sup>,  
M Chandler Bennett<sup>6</sup>, Mark A Berrill<sup>7</sup>, Nick S Blunt<sup>8</sup>, Edgar Josué  
Landinez Borda<sup>9</sup>, Michele Casula<sup>10</sup>, David M Ceperley<sup>11</sup>, Simone Chiesa<sup>11</sup>,  
Bryan K Clark<sup>11</sup>, Raymond C Clay III<sup>2</sup>, Kris T Delaney<sup>12</sup>, Mark Dewing<sup>5</sup>,  
Kenneth P Esler<sup>13</sup>, Hongxia Hao<sup>14</sup>, Olle Heinonen<sup>15,16</sup>, Paul R C Kent<sup>17,18</sup>,  
Jaron T Krogel<sup>19</sup>, Ilkka Kylänpää<sup>19</sup>, Ying Wai Li<sup>20</sup>, M Graham Lopez<sup>7</sup>,  
Ye Luo<sup>4,5</sup>, Fionn D Malone<sup>9</sup>, Richard M Martin<sup>11</sup>, Amrita Mathuriya<sup>1</sup>,  
Jeremy McMinis<sup>9</sup>, Cody A Melton<sup>6</sup>, Lubos Mitas<sup>6</sup>, Miguel A Morales<sup>9</sup>,  
Eric Neuscamman<sup>21,22</sup>, William D Parker<sup>23</sup>, Sergio D Pineda Flores<sup>21</sup>,

Citation paper with high level method & algorithm descriptions  
J. Kim et al, Journal of Physics: Condensed Matter **30** 195901 (2018)  
<https://doi.org/10.1088/1361-648X/aab9c3>

Open access this year



## High-level Algorithms

Today, GPU operates on walkers (Markov chains) in a “batch” to reduce kernel launch overhead, particularly for small problems ( $N < \text{few } 100$ ). CUDA kernels operate on multiple walkers at once.

Drawbacks: Multiple drivers, Extra APIs through entire app, Code divergence, Memory bandwidth still stressed...

### CPU Algorithm

```
do time step i [ 1K-100K per MPI task ]
  do walker j [ 1 per core, OpenMP ]
    do electron k [ 0.1-10K ]
      do component l [ 3-4 ]
        advance WF
      end l
    end k
  end j
  evaluate Hamiltonian
end j
spawn/kill walkers, load balance
end i
```

### GPU Algorithm: 1 MPI task/GPU

```
do time step i
  do electron k (lock step)
    do component l
      advance WF of all the walkers [10-1000]
    end l
  end k
  evaluate Hamiltonian of all the walkers
  spawn/kill walkers, load balance
end i
```

# miniQMC

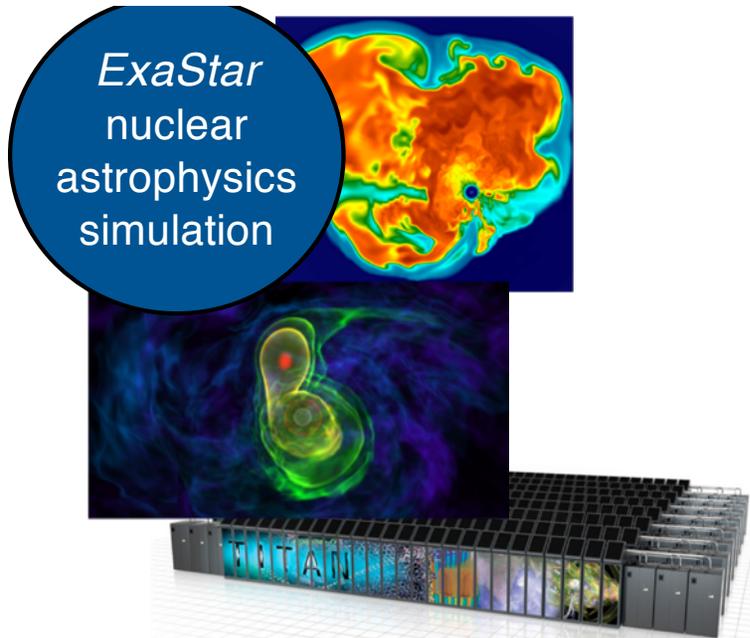
- miniQMC is a mini-app for QMCPack, which is a Quantum Monte Carlo code
  - The goal of QMCPack is to provide highly accurate calculations of the properties of complex materials
- Computational motifs include: particle methods, dense and sparse linear algebra, and Monte Carlo
- The version we are using is based on version 3.1.0 of QMCPack

## miniQMC computation

- QMCPack utilizes an ensemble of walkers which represent particle positions and are moved through space by a drift-diffusion process
- The miniQMC calculation utilizes one walker per rank with a predefined problem to focus on single node performance
- The size of the problem space for the walker in miniQMC can be scaled at runtime by replicating the space by tiling
- There is no MPI communication in miniQMC since the mini-app is focused on the computational portion of the app

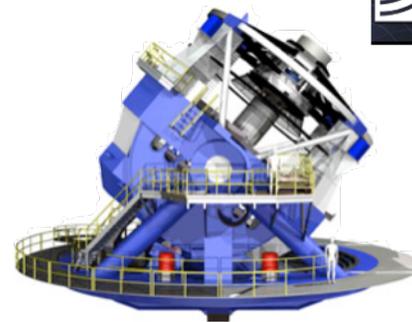
## Significant kernels

- Update – Sherman-Morrison rank-1  $N \times N$  matrix inverse update  
Strongly memory bandwidth limited. Source of  $N^3$  scaling.
- Spline – 3D spline value & gradients  $4 \times 4 \times 4 \times N$  stencil, membw limited
- Jastrow functions – Small classical MD force-field like polynomial evaluations using distances+cutoffs from electron positions. One and two-body forms most common.
- Distance tables – Interparticle distances with minimum image convention/periodic boundary conditions applied.
- Others, depending on threadability of above.

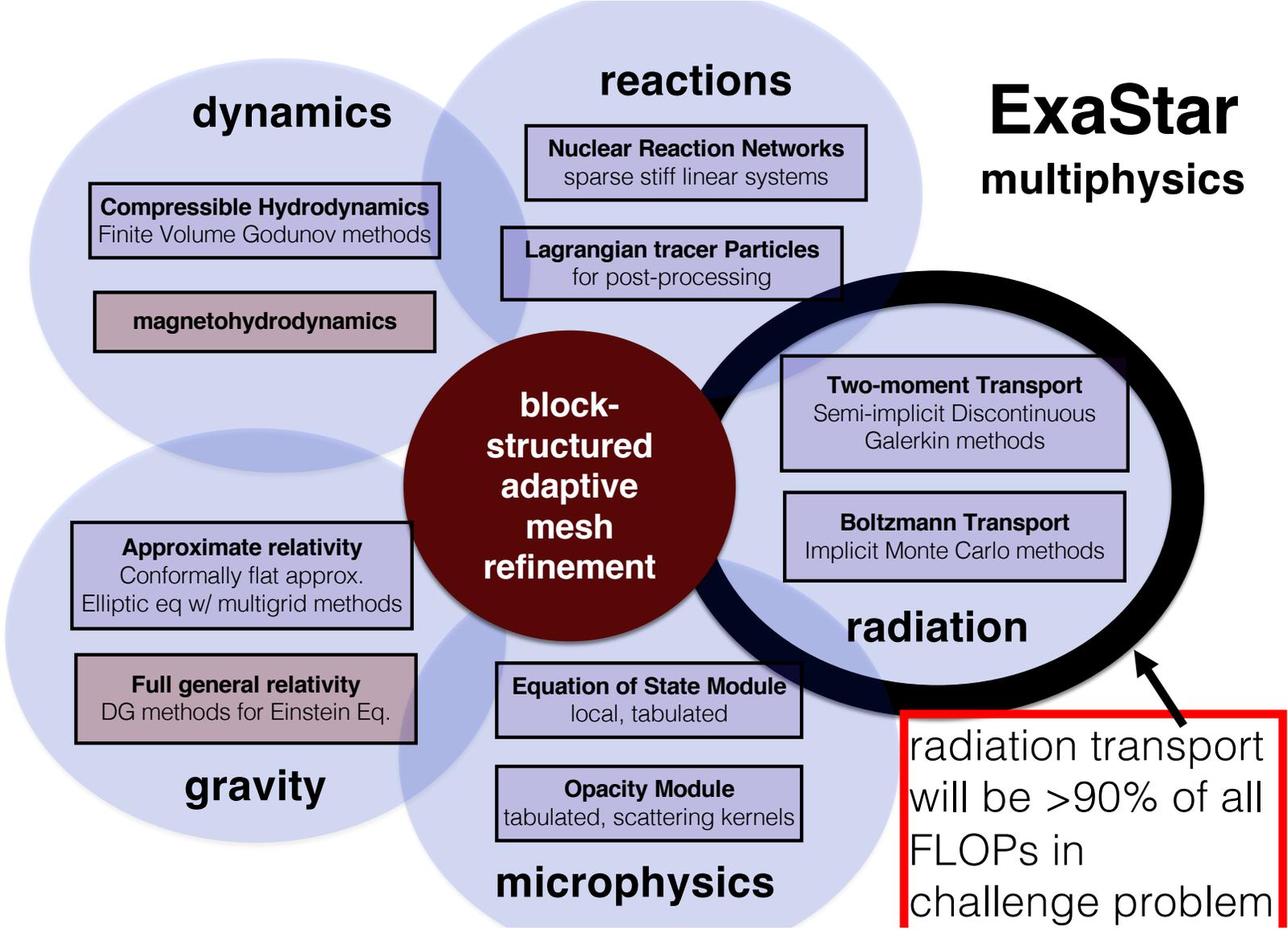


- ExaStar simulations will incorporate:
    - experimental nuclear physics data
    - satellite observations of astrophysical phenomena
    - GW detections
    - neutrino experimental data, including solar and reactor experiments
- to improve predictive power

- ExaStar simulations are essential to:
  - Guide future nuclear physics experimental programs
    - siting the r-process directly impacts which rates are most important to measure
  - Provide reliable templates for gravitational wave and neutrino detectors
    - Low signal-to-noise requires templates for matching
  - Interpret X-ray and gamma-ray observations

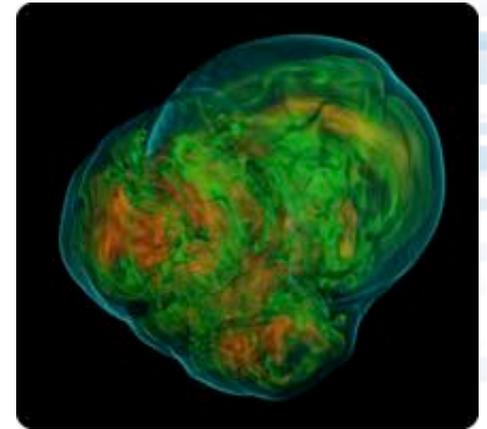


# ExaStar multiphysics

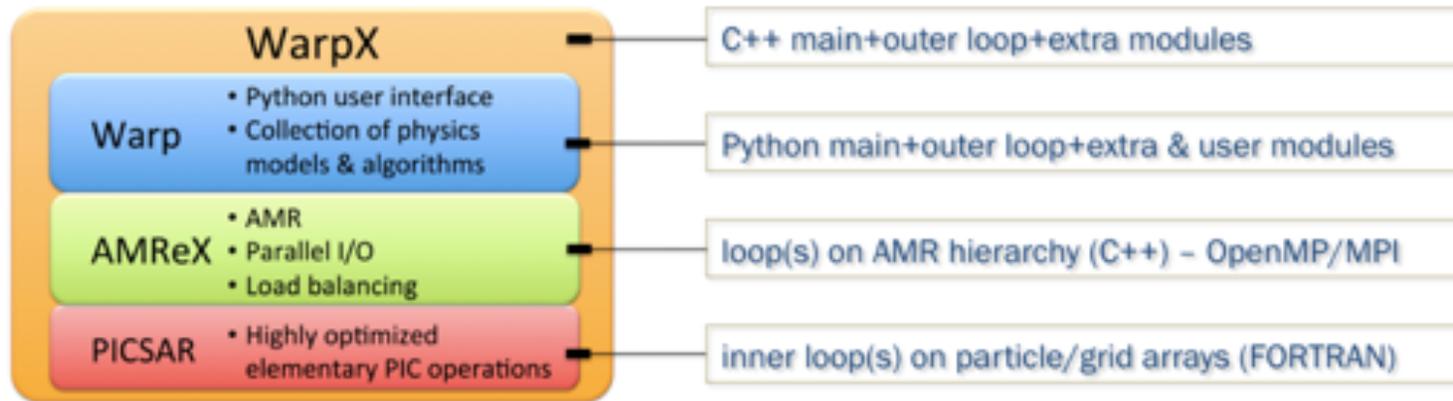


## thornado-mini represents neutrino radiation transport (80-90% of FLOPS in ExaStar challenge problem)

- Key computational motifs are structured mesh, dense linear algebra
- The problem size is defined by the number of elements used, the number of energy groups, and the resolution of the interaction table used.
  - The problem solved in the mini-app (Deleptonization) uses production values for elements and energy groups (2 and 20, respectively), but the table is about quarter-resolution in the three dimensions.
- The mini-app is, in fact, the development vehicle for the finite-element transport algorithms, the mini-app is slowly updated with new implementations. GPU-enabled version (1.1) will be available in March, 2019.



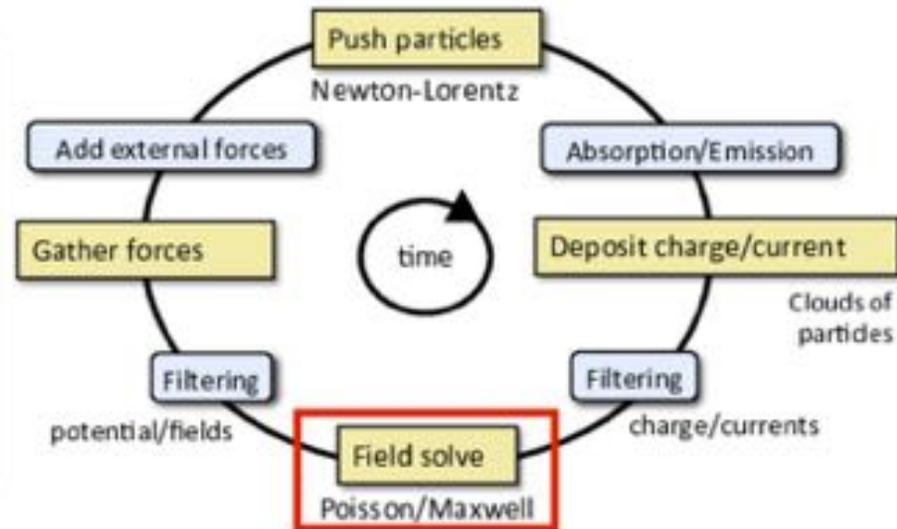
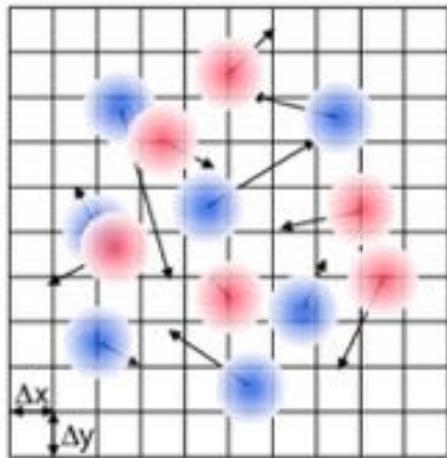
# WarpX: Exascale Modeling of Advanced Particle Accelerators



**WarpX** has code connecting the 3 components

- **Warp**: original Python code (no longer needed!)
- **AMReX**: data and communication
- **PIC SAR**: kernels, also a self-contained mini-app

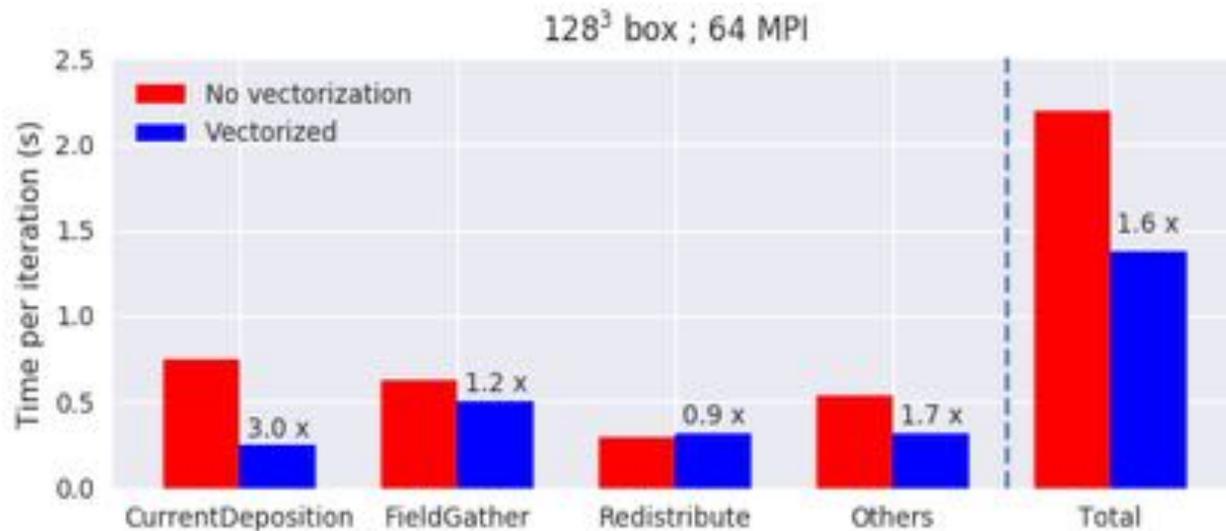
# The Particle-In-Cell (PIC) method



Different regimes:

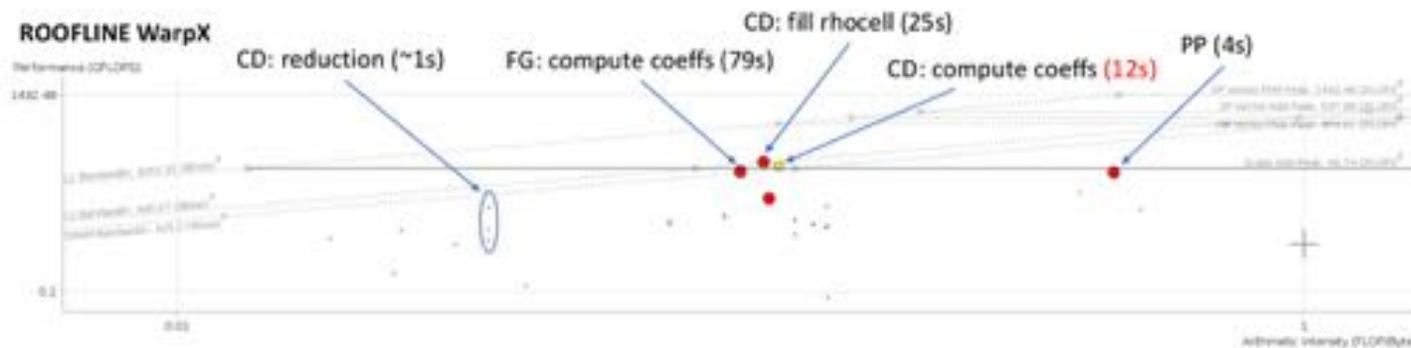
- Particle-dominated,  $\gg 1$  particle/cell
- Commensurate,  $\sim 1$  particle/cell
- Field-dominated,  $\ll 1$  particle/cell

# Time breakdown on 1 KNL node, 50 particles/cell



**Figure 9:** Time taken by the most time-consuming routines (current deposition, field gather, and particle redistribution), as well as by the lumped sum of the other routines (“Others”). The numbers above each blue bar indicates the speedup due to vectorization.

# Roofline model on 1 KNL node, 32 particles/cell



**Figure 12:** Roofline model for WarpX for a simulation with  $128 \times 128 \times 128$  cells and 32 particles per cell on a single Cori KNL node with 64 OpenMP threads. The particle shape factor is 3. This image was obtained using Intel Advisor.

**PICSTAR** kernels mostly memory-bound.  
Shown here:  
Charge Deposition, Field Gather, Particle Push.

## PICSARlite was created to simplify looking a PIC challenges

- Reduced loc from over 80,000 to under 10,000
  - Most steps have multiple options or algorithms. Simplified to only one.
  - Led to script that will automatically reduce PICSAR to a smaller code base containing only the selected methods
- WarpX team is interested in reducing communication:
  - merging messages
  - overlapping communication with computation
- WarpX team has a mini version of WarpX (using only simplest version of each PIC kernel) implemented on SummitDev.

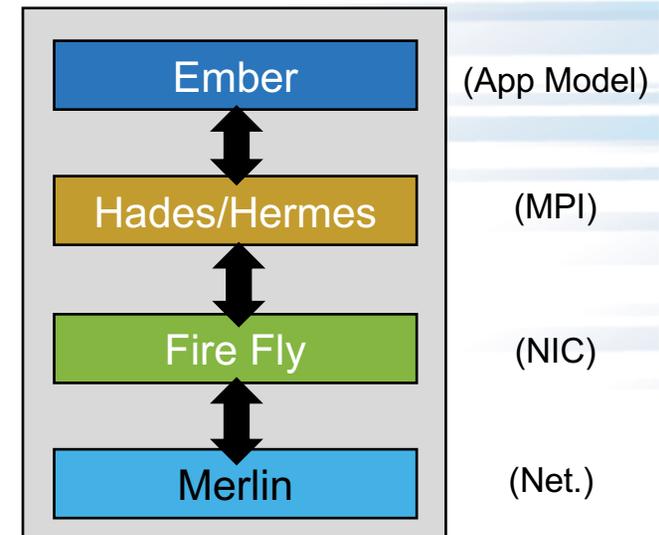
# Modeling Communication Patterns for Exascale

- Exascale systems are likely to be some of the largest HPC machines ever built
  - Implies larger scale and higher performance demands for network interconnects
  - Additional opportunities to add hardware support for some communication operations to reduce latency, improve operation throughput and increase performance

- What we (HPC community, vendors etc) need are:
  - Good models of primitive communication patterns that HPC codes routinely execute (so we can analyze the impact of hardware changes on full system performance)
  - Methods to easily scale communication models so we can evaluate broad range of Exascale system options (traces lock us in to specific rank counts and configurations)
  - Flexible parameterization of communication patterns to reflect choices in decomposition that could be important (e.g. when using a GPU versus a CPU versus a ...)
  - Note: we cannot represent every communication pattern in every application, we need to capture a small subset from which important characteristics of network performance can be modeled in order to support tractable analysis of future systems

# Ember Communication Patterns

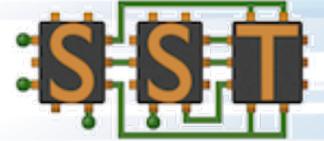
- Ember is a suite of communication patterns that have been developed since 2012 in collaboration with leading industry vendors (Cray, IBM, Intel, HPE, ..)
  - Designed originally to work in the SST Simulator as a scalable model of DOE workloads (<https://github.com/sstsimulator>)
  - Highly parameterizable to easily replicate behavior across range of DOE code bases
  - Flexible enough to scale from small node counts to over a million simulated MPI ranks (much easier to use than communication traces)
  - Can encode complex, dynamic behavior which traces cannot capture
  - Generic enough to work in any simulation environment



- **For ECP Proxy Applications several communication patterns have been turned into simple, equivalent MPI/SHMEM drivers to run on HPC systems**

<http://github.com/sstsimulator/ember>

# Using Ember Communication Patterns



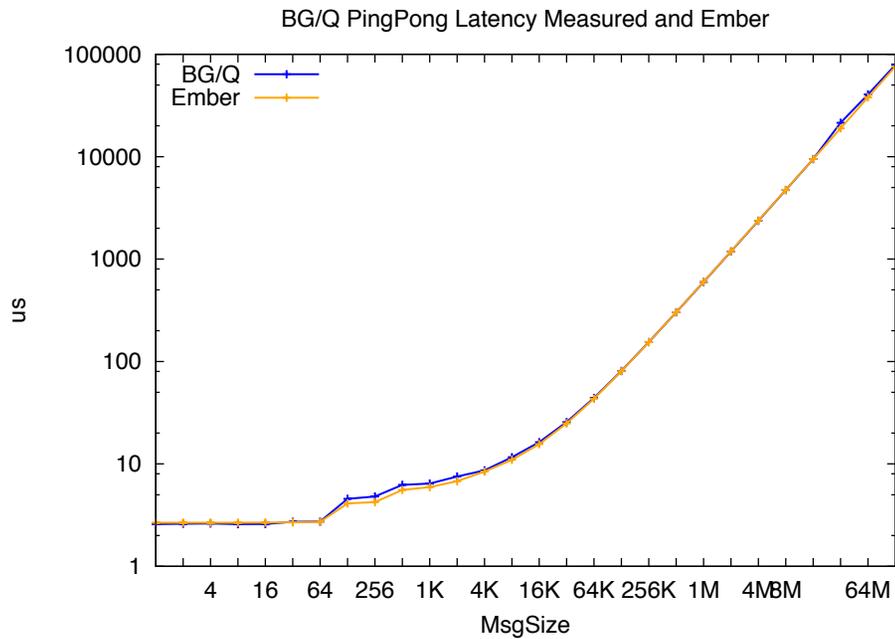
- **How can I use Ember Communication Patterns?**

- Two main choices: (1) Run the MPI/SHMEM implementations, (2) Run the patterns on one of DOE's simulation environments (e.g. SST, CODES..)

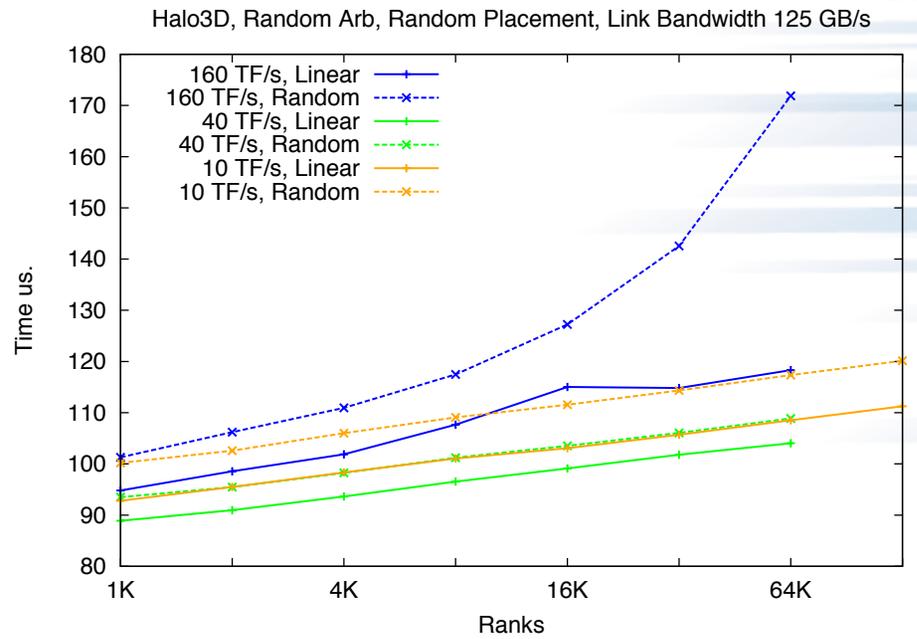
- **What is the SST Simulator?**

- Parallel, Conservative, Discrete-Event Hardware Simulation Environment
- Can simulate cycle-accurate models of single compute nodes through to full-scale models of system interconnects
- Validated against DOE machine installations
- Used by industry vendors (Cray, Intel, IBM, HPE, etc)
- Email: [wg-sst@sandia.gov](mailto:wg-sst@sandia.gov) for more information
- Go to: <https://sst-simulator.org>

# Example Uses of Ember



**Simulator/Performance Model Validation**



**Projection of Communication Patterns For Future Exascale Machine Configurations**



## Communication Patterns included in Ember

- Nearest neighbor (Structured, halo communication)
- Nearest neighbor (Unstructured, halo-like communication, varying message sizes)
- Wavefront sweeps
- MPI Collective operations (reduction, all-gather, scatter, etc)
- Random, fine-grained messaging (closer to large-scale graph analysis)
- FFTs (all-to-all messages)
- Network in-casts (many to few communication flows, also represent I/O transfers)

# Assessment Activities

- Goal

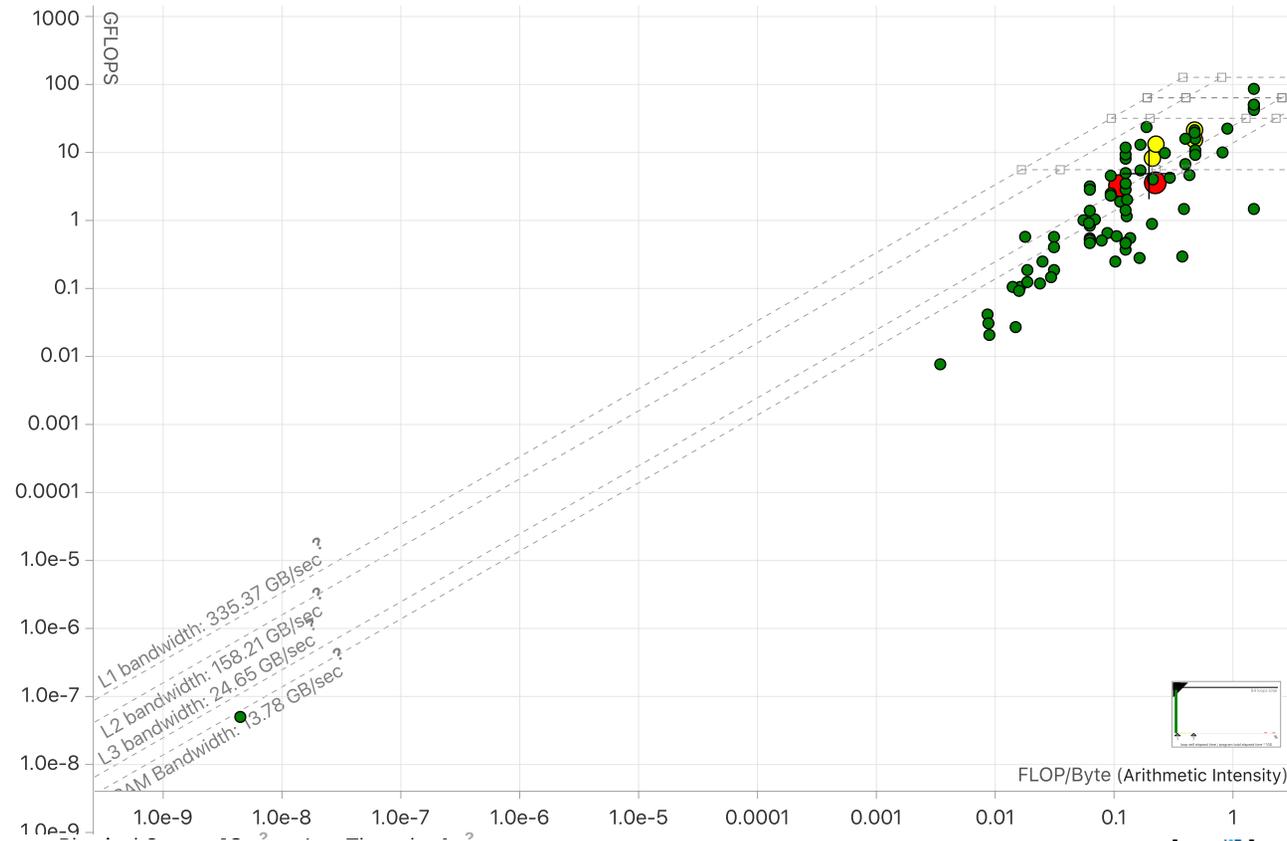
- Performance characterization and identification of underlying hardware bottlenecks that negatively impact performance of both proxy and parent application
  - On current generation platforms
  - Using Exascale challenge problems where possible, otherwise a challenge problem on current systems
- Determine quantitatively if the proxy represents the parent as intended (e.g., memory, communication, computation)
- Target 8 new proxy/parent pairs per year
- Looking at both CPU and GPU-based performance

# Standard Performance Characterization

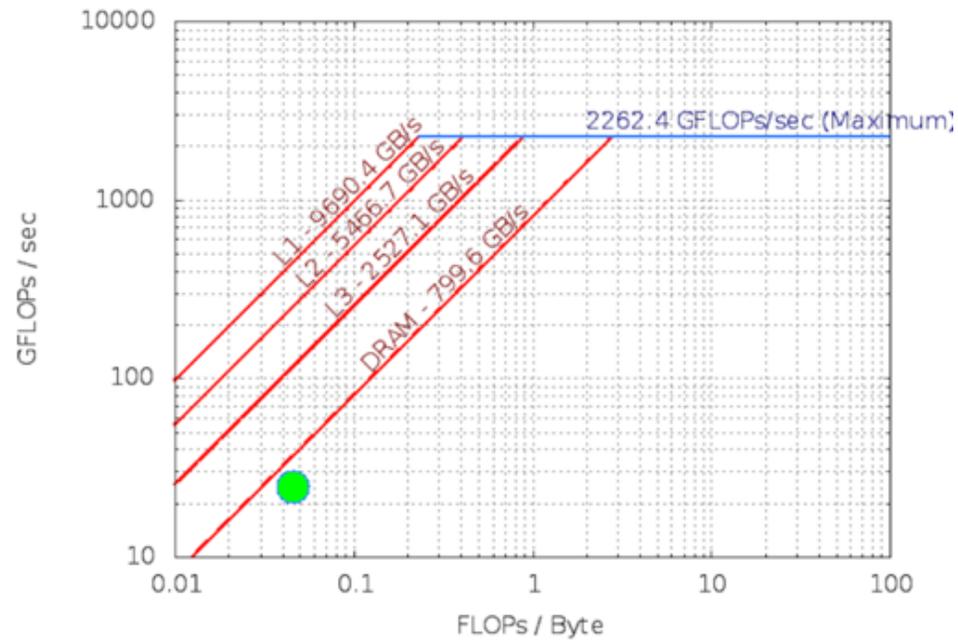
- Dynamic profiling to determine hot functions
- Roofline model to understand upper bound
  - Cache and memory bandwidths
  - FLOPS/arithmetic intensity
- For both proxy and parent and for whole execution and for each of top 10 functions (10 functions that account for largest percentage of total execution time)
  - Cache and memory bandwidths
  - FLOPS/arithmetic intensity

# Baseline using Roofline Model

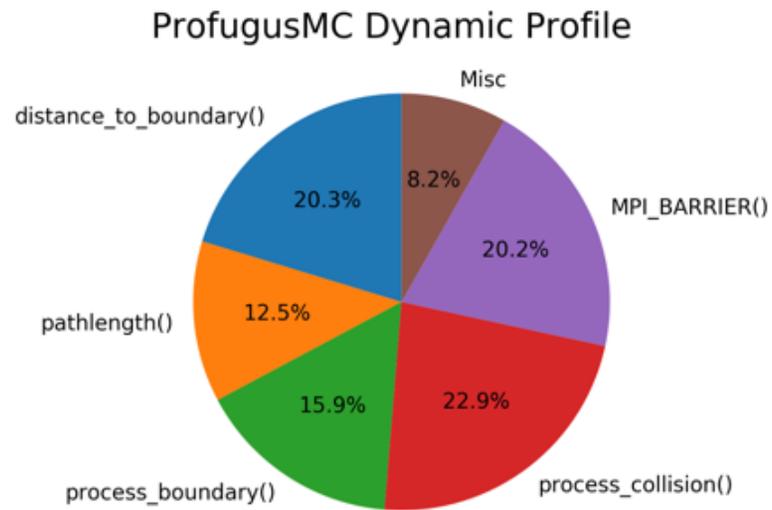
- miniQMC



# Profugus roof line from Brian



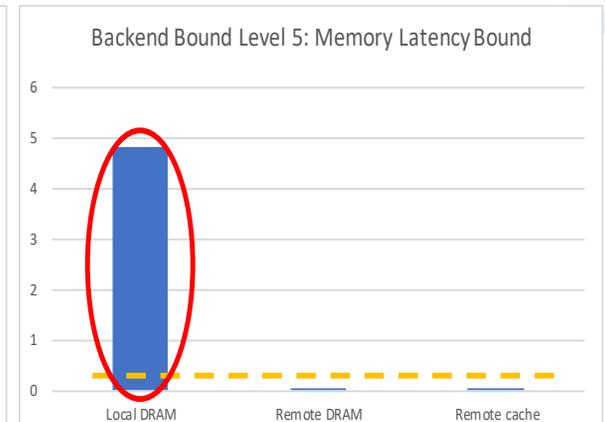
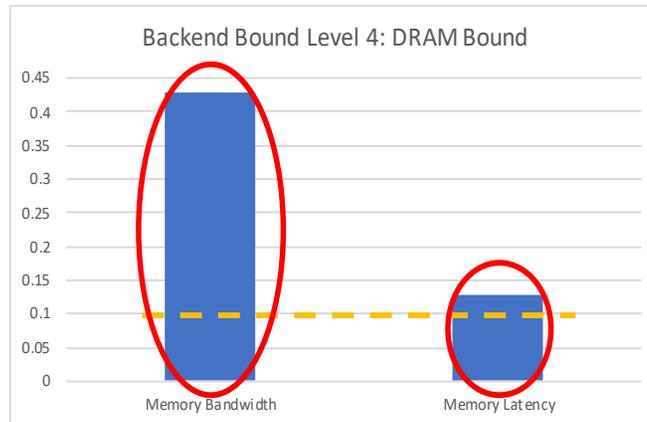
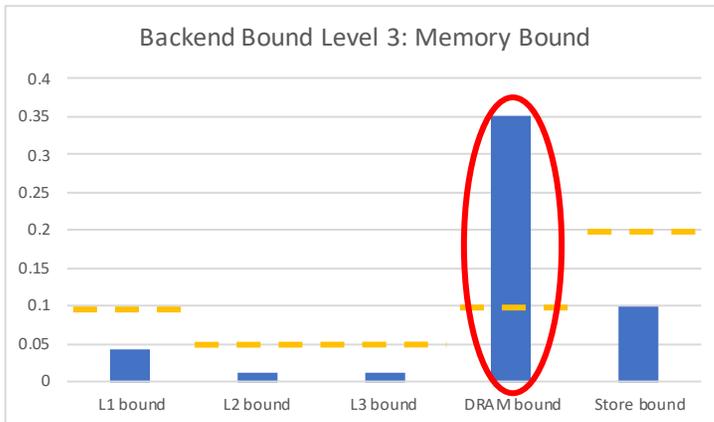
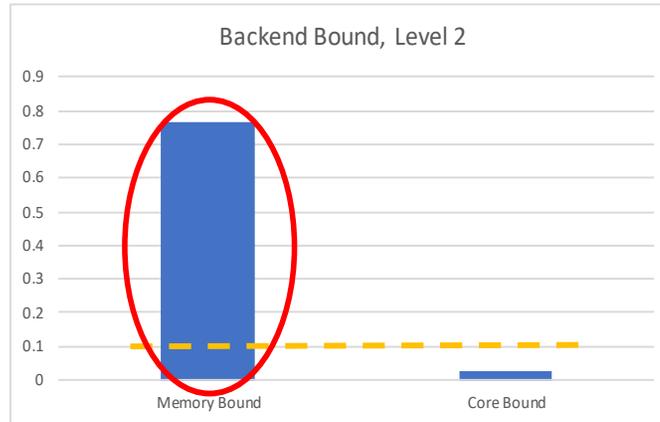
# Profugus Dynamic Profile from Brian



# Identifying Hardware Bottlenecks

- Drill-down on memory bounded-ness
  - Implemented Intel's Top-Down Microarchitecture Analysis (TMA) into LDMS
    - Fast, flexible, and get per-process data
  - Currently validating TMA methodology
    - Does it really identify bottlenecks correctly?

# QMCPack, NiO256 atoms on Skylake

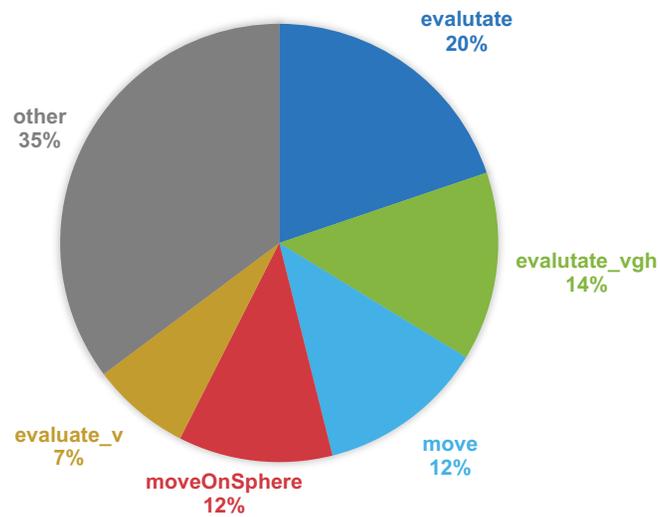


TMA identifies the same bottlenecks for miniQMC

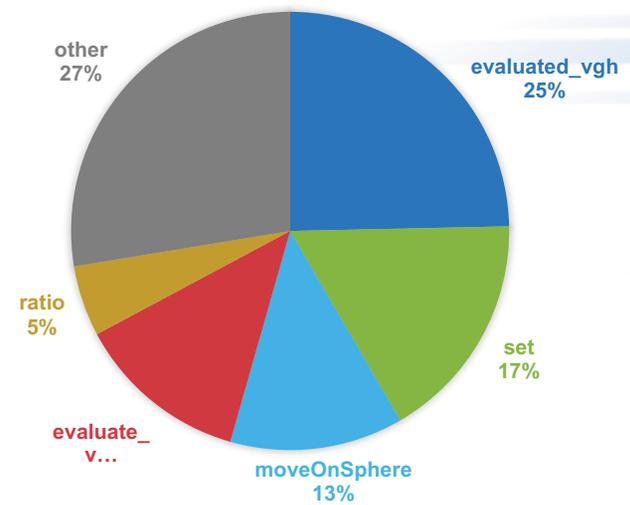
# Function Drill-Down: Dynamic Profiles

- MiniQMC

MINIQMC - 1 X 1 X 32



MINIQMC - 1 X 8 X 1



Instrumented these functions, currently collecting TMA and other hardware performance counters to better understand composite behavior.

# GPU Characterization

- Motivation
  - Large proportion of performance on exascale computers expected to come from accelerators such as GPUs.
  - Tradeoffs between performance and portability depending on GPU program model used.
  - Application may have different characteristics and bottlenecks on GPU vs. CPU architectures.
- Goal is to use proxy apps to answer the following questions:
  - What are the implications of GPU architecture trends for ECP applications?
  - What GPU programming models will enable the best performance portability across future accelerator architectures?
  - What software optimizations are key to achieving good performance on GPU architectures?

## Interaction with ECP AD Teams

Best practices for proxy app development for GPUs (Proxy App team can help as resources allow)

- Provide proxy app configurations and inputs representative of exascale application problems (can be scaled back to reduce runtime and number of nodes)
- Keep proxy app up-to-date with main app development
  - In some cases (e.g., SW4lite, ProfugusMC), proxy app leads GPU development
- Provide feedback on Proxy App team assessment results (maybe we missed something)
- Optimize proxy app GPU implementation for new GPU architectures and features
  - Most time-consuming kernels are the most important.

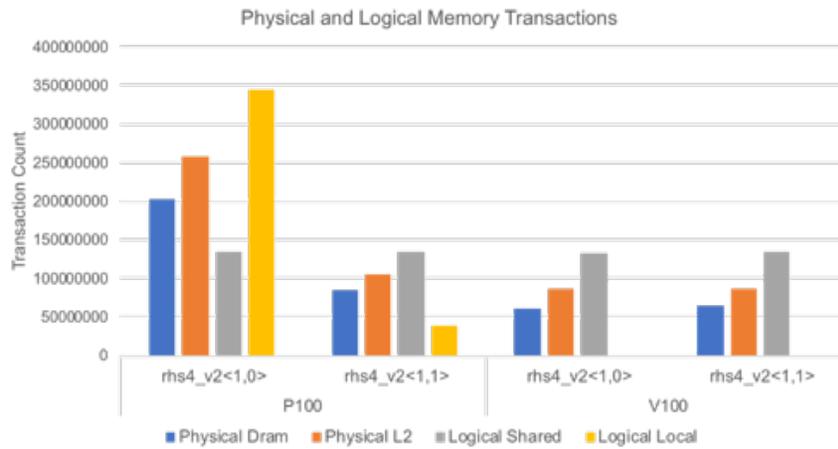
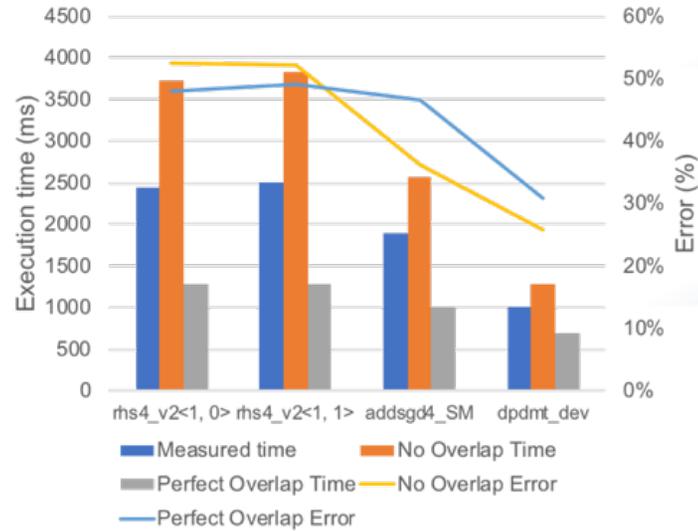
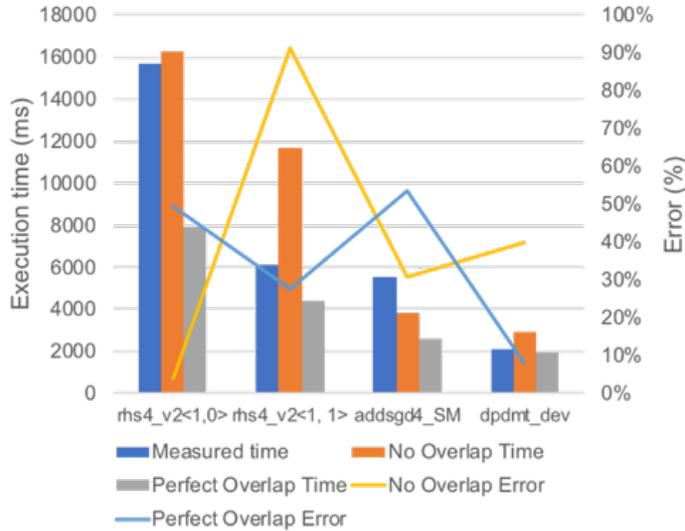
# Why GPU Assessment Matters

- Ask yourselves
  - If a vendor uses my proxy app to help with design decisions for GPU architecture features, will the result benefit performance of my application?
  - If my proxy app is used for procurement decisions for a machine with accelerators, will the resulting machine architecture and configuration match my application requirements?
- If you don't have a representative GPU implementation that is agile enough to be easily built and run (including on simulators and prototypes), your application will not influence these decisions!

# GPU Assessment Procedure

- Similar to CPU assessment that has already started
- Steps
  - Obtain code and inputs for CPU and GPU versions
  - Build and run CPU and GPU versions and verify results
  - Profile execution
  - Collect characterization data
    - Instruction mix
    - Achieved cache and memory bandwidths
    - Communication characteristics
  - Analysis
    - Compare full and proxy app characteristics
    - Roofline models
    - Identify performance and scaling bottlenecks

# SW4lite Performance on P100 vs. V100



LOH.1-h100.in input running on one GPU



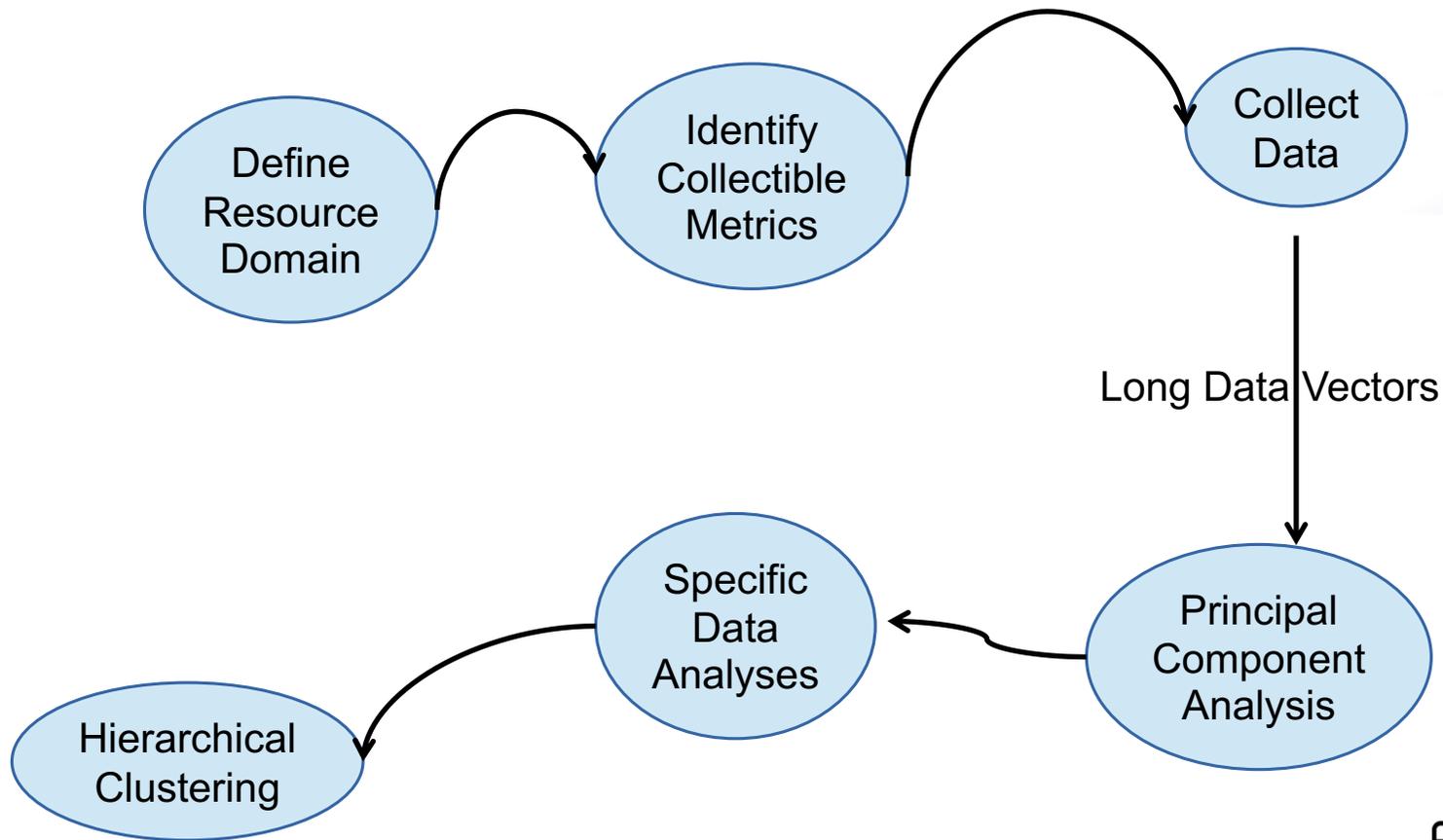
## SW4lite Results on Summit

Test case	Ngp	Nts	Nodes	GPUs	Runtime (sec)
LOH.1-h100	1.55e7	1073	1	1	19.4
LOH.1-h100	1.55e7	1073	1	6	5.16
LOH.1-h50	1.23e8	1073	1	3	55.4
LOH.1-h50	1.23e8	1073	1	6	31.8
LOH.1-h50	1.23e8	1073	2	8	25.4
LOH.1-h25	9.82e8	1073	2	12	115
LOH.1-h25	9.82e8	1073	4	24	69.7

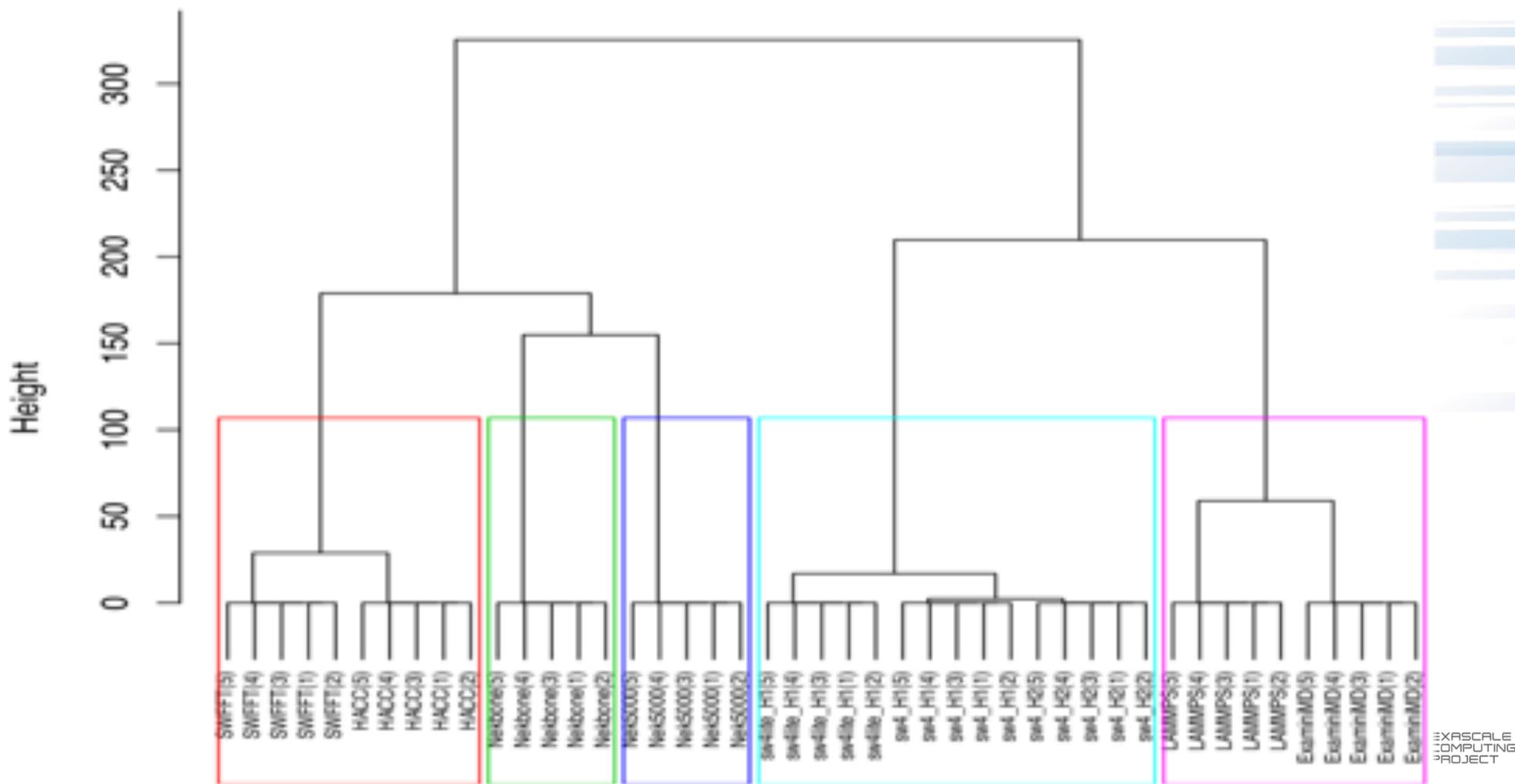
### Remarks:

- Reasonable strong scaling
- Weak scaling is not good (much worse than SW4 on Cori 2 KNL, where runtime increased 6% going from h100 to h50 and increased 30% going from h50 to h25)

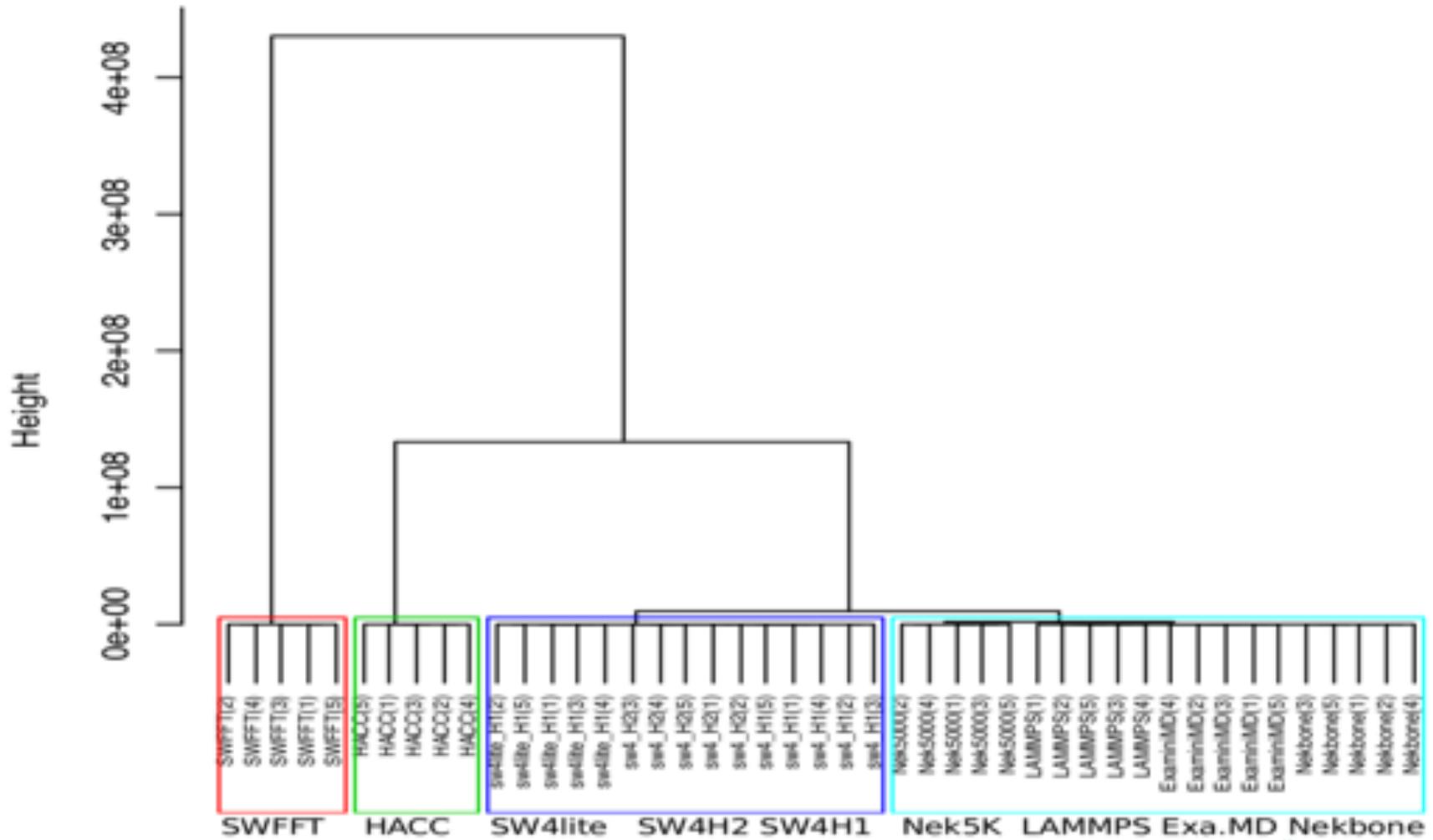
# Clustering Analysis of Proxy Representativeness



# Broadwell Basic Node Hardware Metric Clustering



# Broadwell Communication Clustering



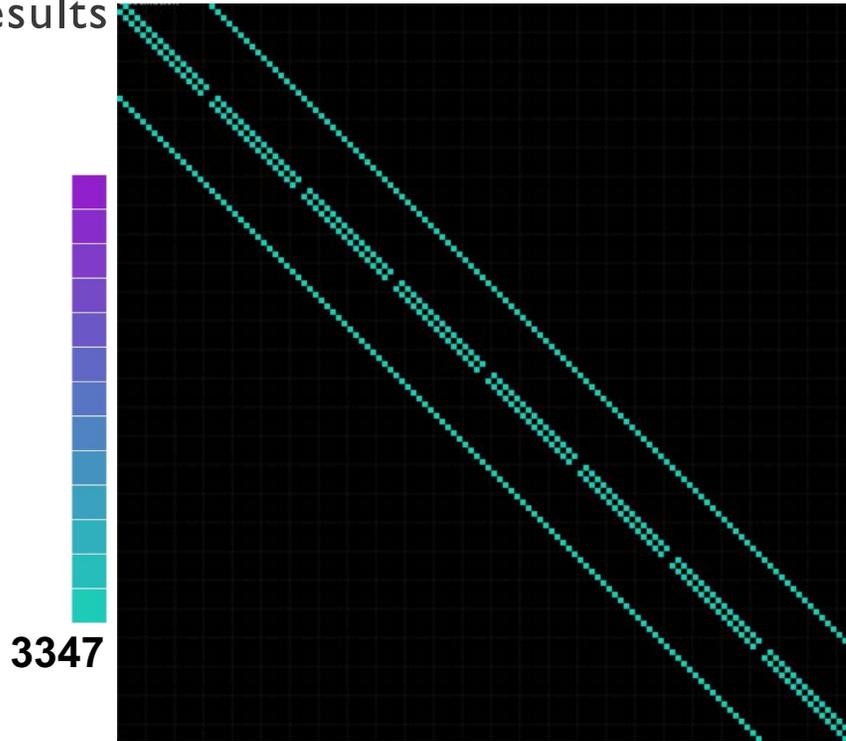
# Proxy/Parent Communication Similarity

- Communication metrics based on behavior (not MPI primitives used)
  - Pairwise communication data analysis
    - Point to point communication patterns (source, destination)
    - Total number of messages sent for each pair
    - Using CrayPat tool
  - Communication vector data clustering
    - KB/sec - Total size of data transferred (KB) / total execution time (sec)
    - MPI KB/sec - Total size of data transferred (KB) / total time spent in MPI (sec)
    - Message size histogram data
    - Using mpiP tool

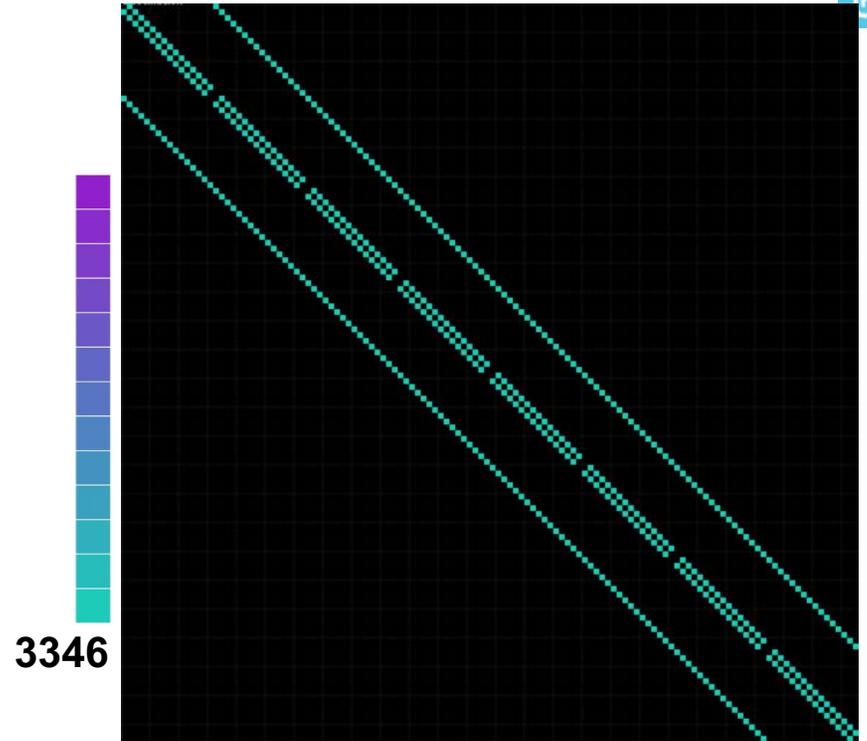
PMBS2018, “Exploring and Quantifying How Communication Behaviors in Proxies Relate to Real Applications”

50 Results

### SW4



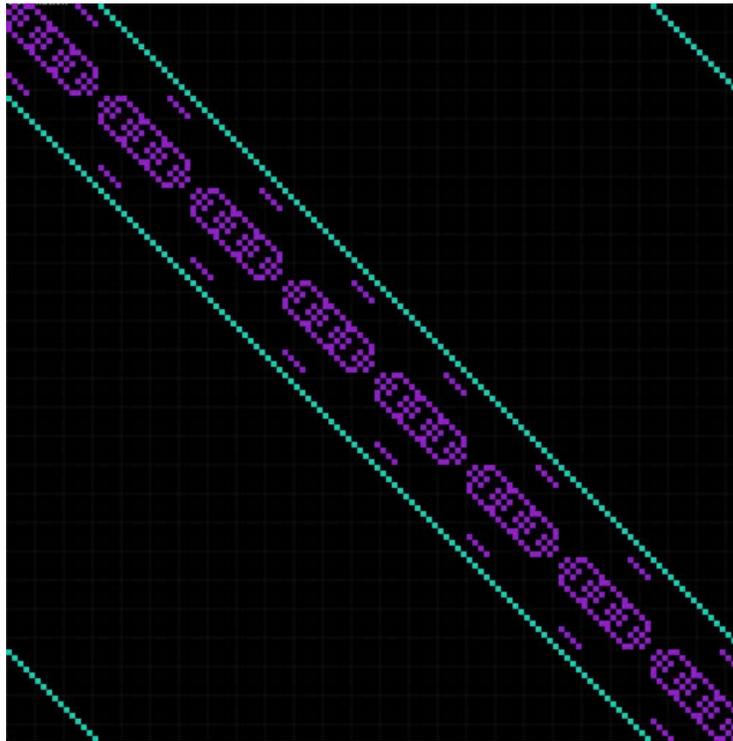
### SW4Lite



Parent/Proxy	Parent in Proxy		Proxy in Parent		Full Set		Parent in Proxy		Proxy in Parent	
	#msg	#pair	#msg	#pair	PCorr	SCorr	PCorr	SCorr	PCorr	SCorr
SW4/ SW4lite	100	100	100	100						

## LAMMPS

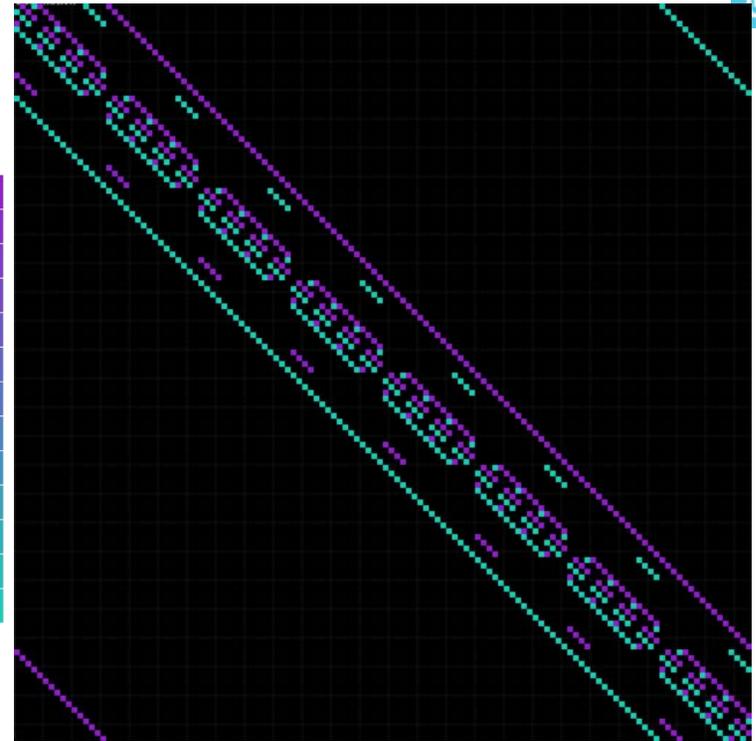
14414



7208

## ExaMiniMD

20703



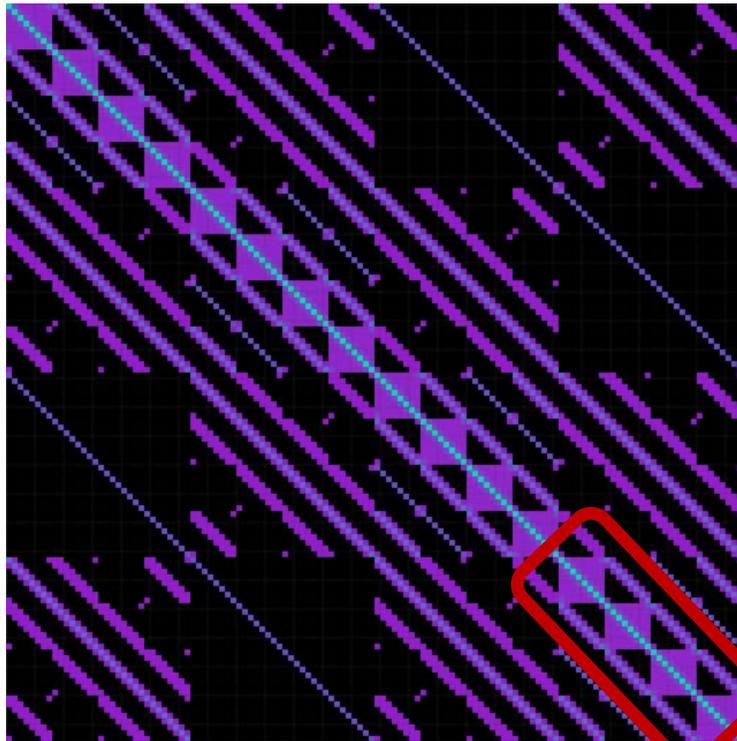
20702

Parent/Proxy	Parent in Proxy		Proxy in Parent		Full Set		Parent in Proxy		Proxy in Parent	
	#msg	#pair	#msg	#pair	PCorr	SCorr	PCorr	SCorr	PCorr	SCorr
LAMMPS/ ExaMMD	100	100	100	100	0	0	0	0	0	0

## HACC

1210

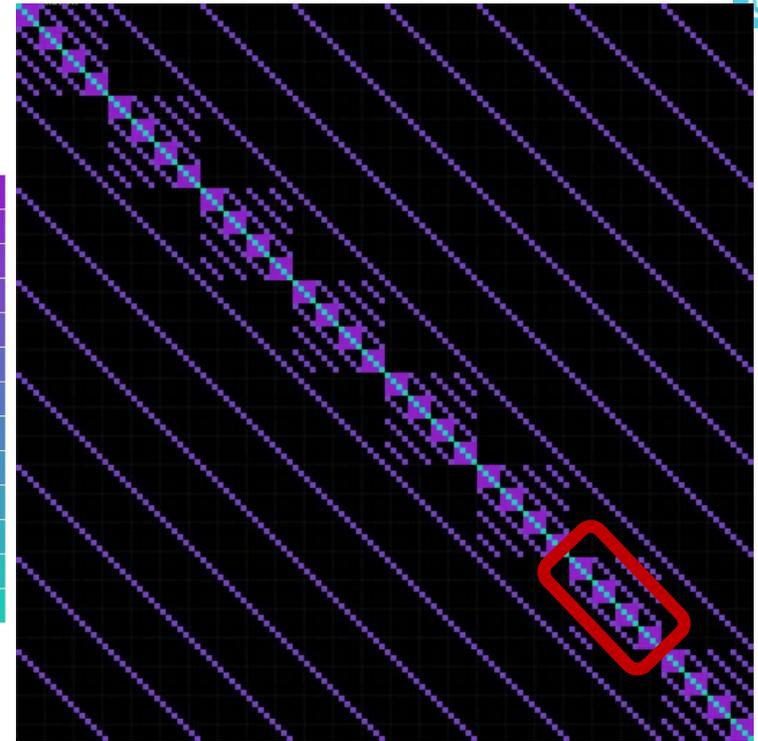
178



## SWFFT

500

200

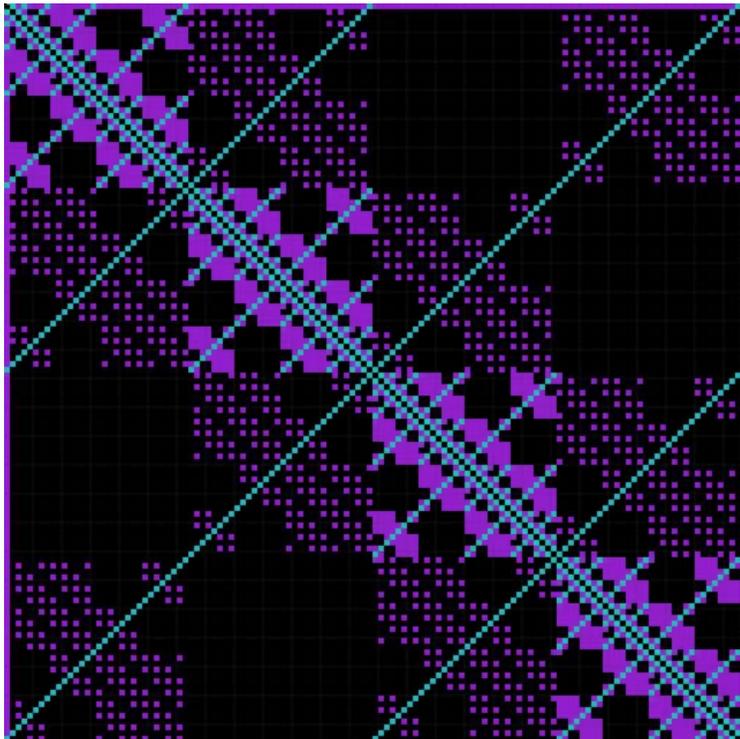


Parent/Proxy	Parent in Proxy		Proxy in Parent		Full Set		Parent in Proxy		Proxy in Parent	
	#msg	#pair	#msg	#pair	PCorr	SCorr	PCorr	SCorr	PCorr	SCorr
HACC/ SWFFT	51.7	29.4	71.4	71.4	0.58	0.31	0.61	0.28	0.87	0.81

Nek5K 3D

35046

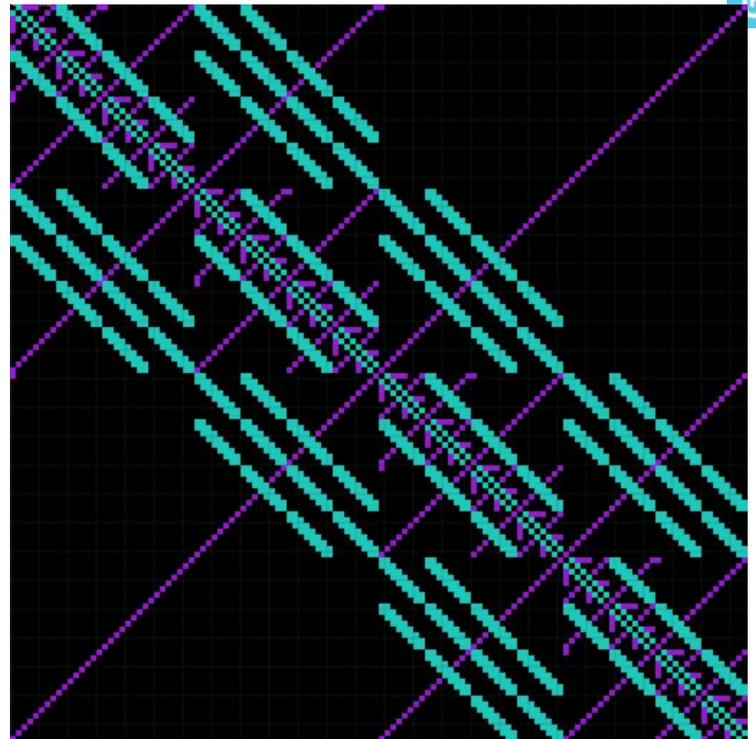
24



Nekbone

15342

651



Parent/Proxy	Parent in Proxy		Proxy in Parent		Full Set		Parent in Proxy		Proxy in Parent	
	#msg	#pair	#msg	#pair	PCorr	SCorr	PCorr	SCorr	PCorr	SCorr
Nek5K 3D/ Nekbone 3D	99.9	51.4	58.0	68.4	-0.1	-0.05	-0.65	-0.23	0.04	0.49

## Acknowledgment

This research was supported by the **Exascale Computing Project (ECP)**, Project Number 17-SC-20-SC, a collaborative effort of two DOE organizations, the Office of Science and the National Nuclear Security Administration, responsible for the planning and preparation of a capable exascale ecosystem including software, applications, hardware, advanced system engineering, and early testbed platforms, to support the nation's exascale computing imperative.

# Questions???



**Thank you!**

[exascaleproject.org](http://exascaleproject.org)



U.S. DEPARTMENT OF  
**ENERGY**

Office of  
Science

