# Smart-DASH Update and Status

**SPPEXA Annual Plenary Meeting 2019**
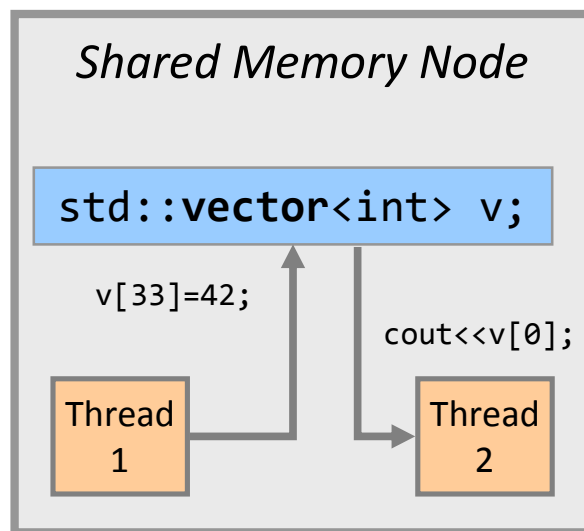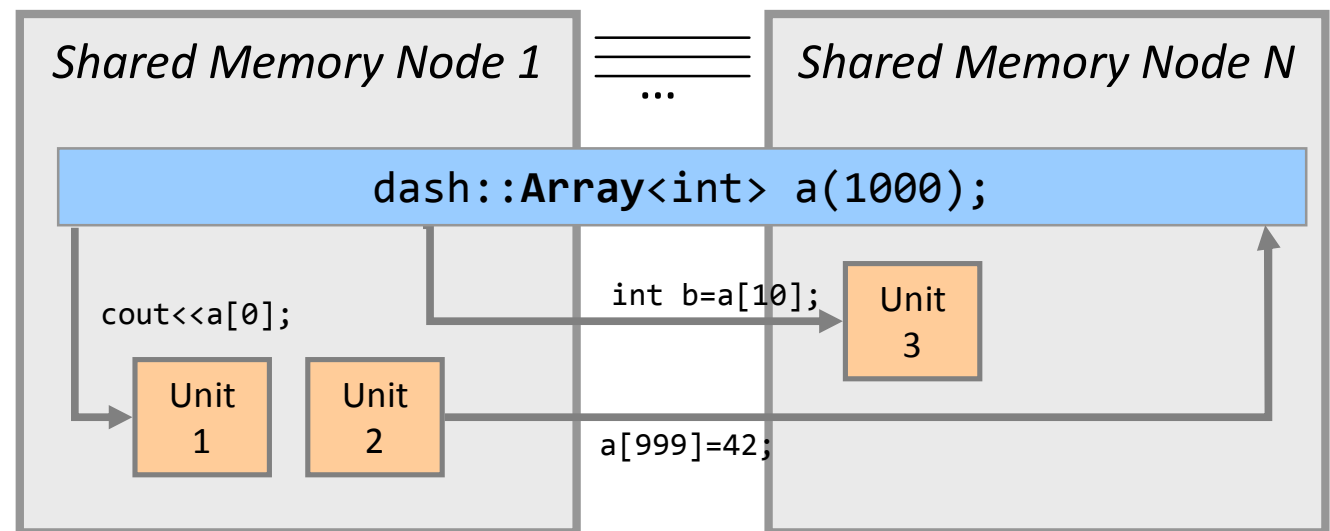**22-23 January 2019, Garching**

www.dash-project.org

Karl Fürlinger
Ludwig-Maximilians-Universität München

- DASH is a C++ template library, that offers
  - Distributed data structures, e.g., `dash::`**`Array`**`<int>`
  - Parallel algorithms, e.g., `dash::`**`sort`**`()`

- Generalizes shared memory programming to distributed memory systems:



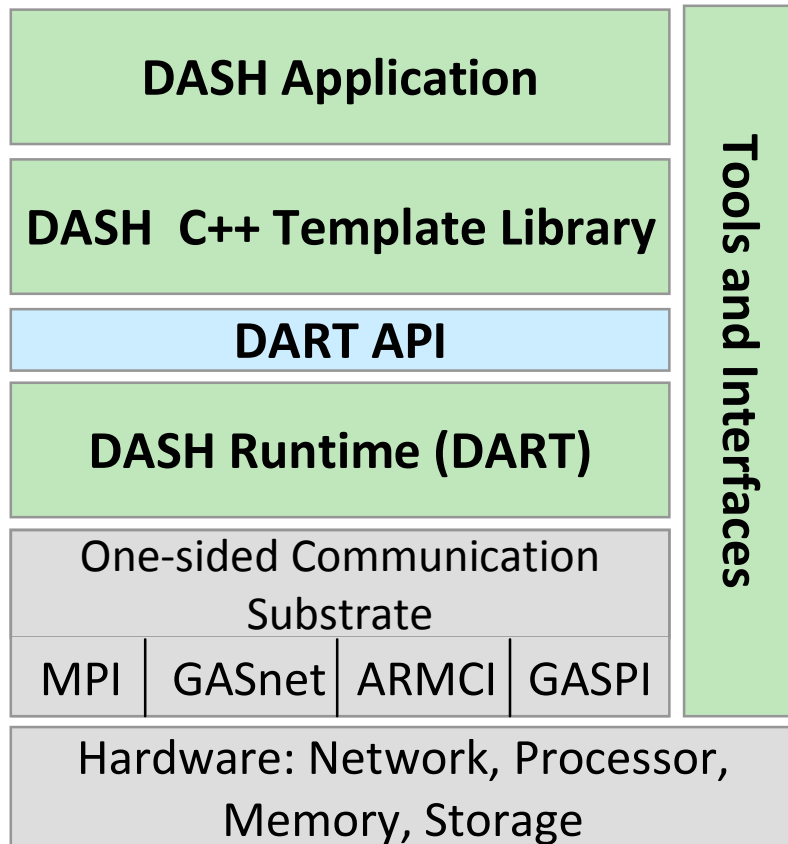| Shared Memory Node | Shared Memory Node 1 ... Shared Memory Node N |
|---|---|
| `std::`**`vector`**`<int> v;` | `dash::`**`Array`**`<int> a(1000);` |
| `v[33]=42;` `cout<<v[0];` Thread 1 Thread 2 | `cout<<a[0];` Unit 1 Unit 2 `int b=a[10];` Unit 3 `a[999]=42;` |

- Multiple threads access **physically** shared memory
- Multiple nodes connected by a high-speed network
- Multiple threads ("units") access **logically** shared memory

# DASH – Project Overview

**DASH Application**

**DASH  C++ Template Library**

**DART API**

**DASH Runtime (DART)**

One-sided Communication Substrate

| MPI | GASnet | ARMCI | GASPI |

Hardware: Network, Processor, Memory, Storage

**Tools and Interfaces**

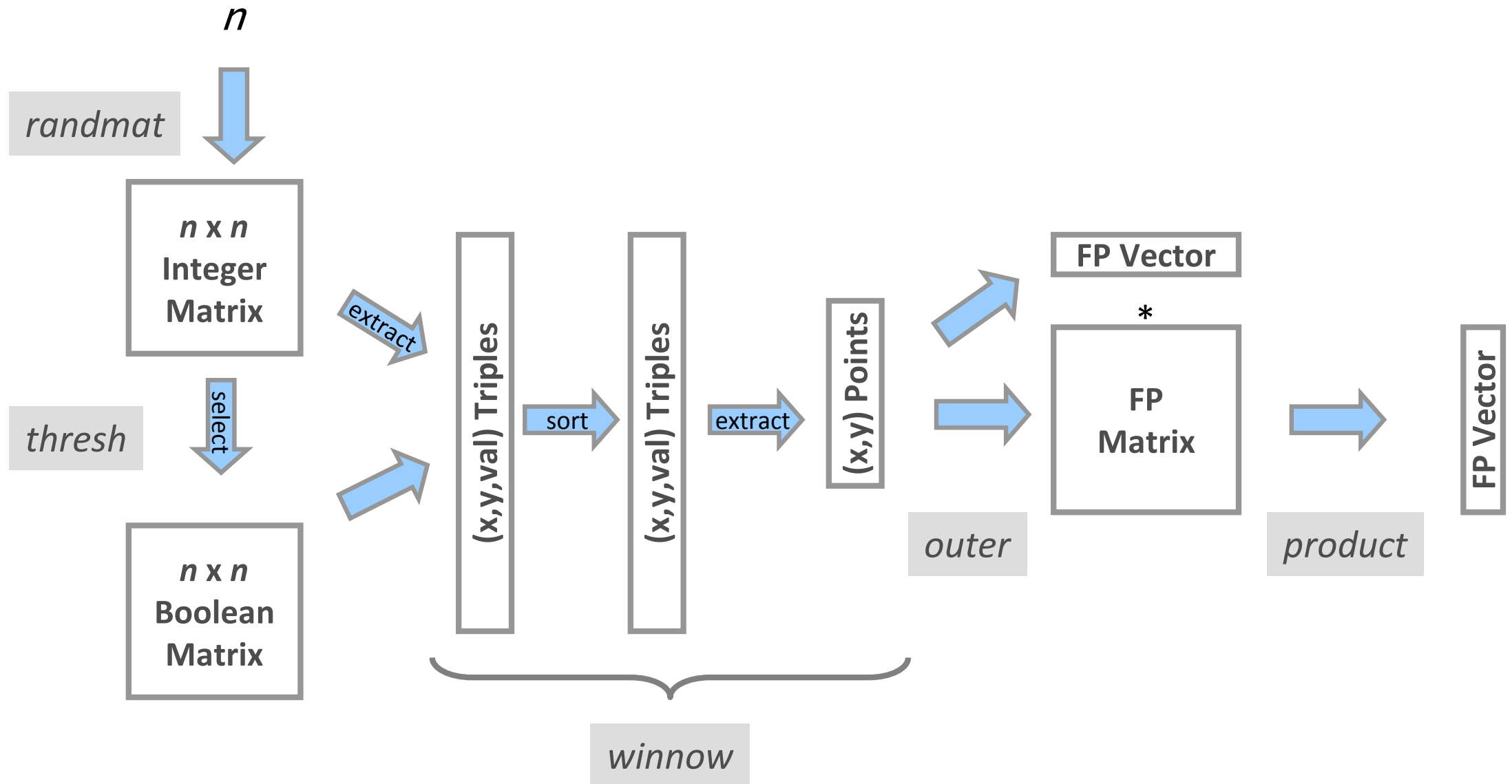|  | Phase I (2013-2015) | Phase II (2016-2018) |
|---|---|---|
| **LMU Munich** | Project management, C++ template library | Project management, C++ template library, DASH data dock |
| **TU Dresden** | Libraries and interfaces, tools support | Smart data structures, resilience |
| **HLRS Stuttgart** | DART runtime | DART runtime, Tasking |
| **KIT Karlsruhe** | Application case studies | |
| **IHR Stuttgart** | | Smart deployment, Application case studies |

www.dash-project.org

DASH is one of 16 SPPEXA projects

- Performance and Productivity Evaluation

- DASH Algorithms – Distributed Parallel Sort

- DASH and task-based execution

# Performance and Productivity

- Cowichan problems
  - A **benchmark suite** designed to investigate the usability of parallel programming systems (1990s)
  - 13 "toy" problems, quick implementation, composable by chaining [1]

- Previous work by Nanz et al. [2] selected **five benchmarks** to evaluate the usability of multicore languages
  - Four programming systems compared:
    - **Go**, **Cilk**, **TBB**, **Chapel**
  - Metrics:
    - **Usability**: LOC, development time
    - **Performance**: execution time and scalability

[1] Wilson, Gregory V., and R. Bruce Irvin. *"Assessing and comparing the usability of parallel programming systems."* University of Toronto. Computer Systems Research Institute, 1995.

[2] Nanz, Sebastian, Scott West, Kaue Soares Da Silveira, and Bertrand Meyer. *"Benchmarking usability and performance of multicore languages."* In Empirical Software Engineering and Measurement, 2013 ACM/IEEE International Symposium on, pp. 183-192. IEEE, 2013.

| | DASH | go | Chapel | TBB | Cilk |
|---|---|---|---|---|---|
| randmat | 16 | 29 | 14 | 15 | 12 |
| thresh | 31 | 63 | 30 | 56 | 52 |
| winnow | 53 | 94 | 31 | 74 | 78 |
| outer | 23 | 38 | 15 | 19 | 15 |
| product | 20 | 27 | 11 | 14 | 10 |

- **DASH is not the most concise approach, but not much worse than the best solution**
  - DASH is the only case where the same code can be run on shared memory and distributed memory systems!

*"Investigating the Performance and Productivity of DASH Using the Cowichan Problems",* K. Fürlinger, R. Kowalewski, T. Fuchs, and B. Lehmann; Proc. of the International Conference on High Performance Computing in Asia-Pacific Region, Tokyo Jan. 2018

- Platform: **Single node** of SuperMUC Phase 2 (Haswell)
  - Haswell Xeon E5-2697, 2.6 GHz, 28 cores per node, 64 GB mem
  - 30k x 30k matrix
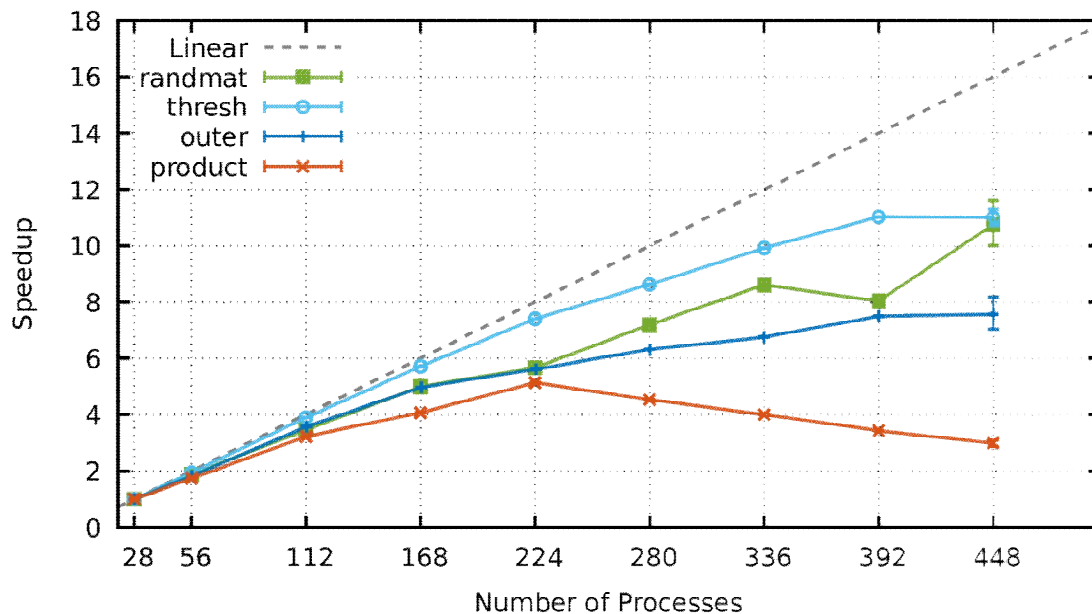  - Intel Compiler (icc) v. 18.0.2 used for all programming systems



Absolute performance, using all 28 cores, 30k x 30k Matrix

Performance relative to DASH, using all 28 cores, 30k x 30k Matrix

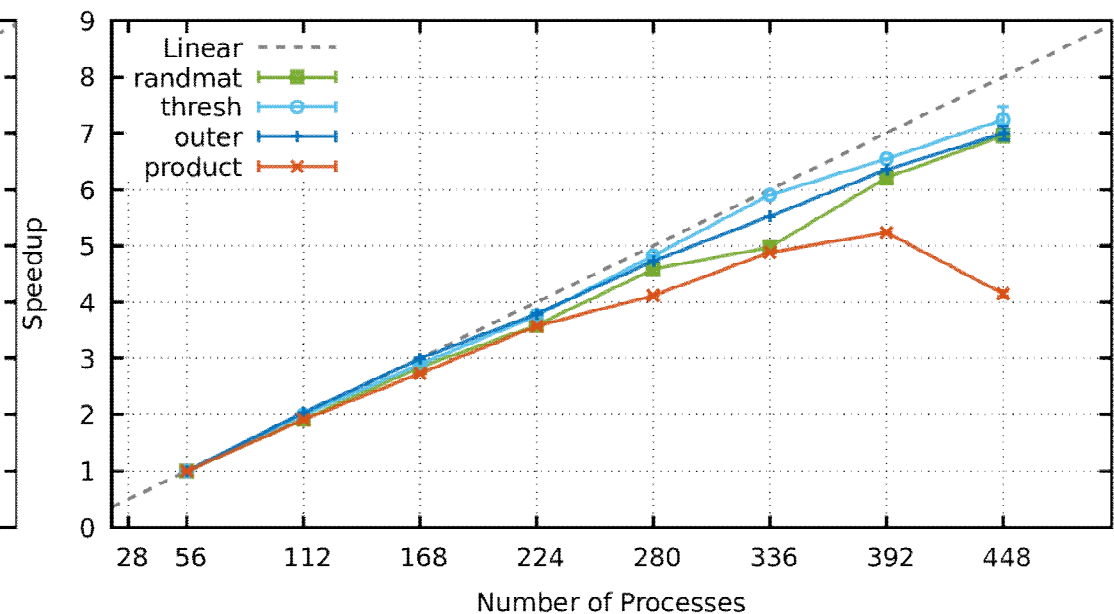- Platform: **Up to 16 nodes** of SuperMUC
    - Haswell Xeon E5-2697, 2.6 GHz, 28 cores per node
    - 64 GB of main memory
    - DASH is the **only** approach that can also use distributed memory machines (the **same source code**)
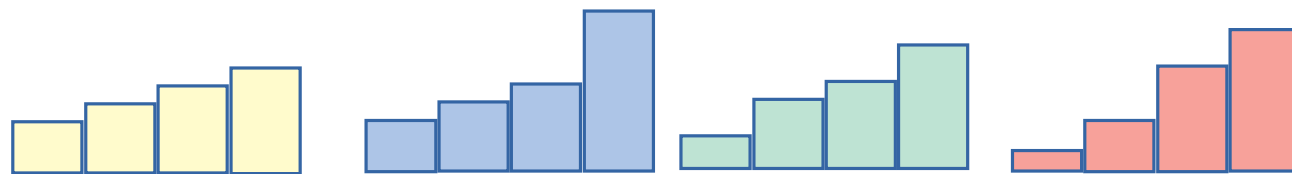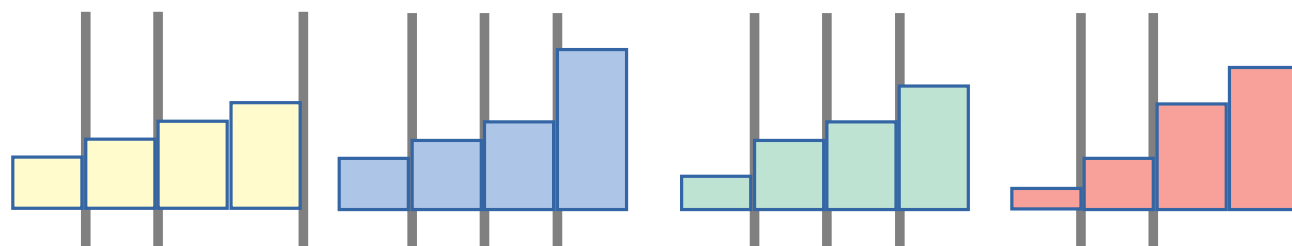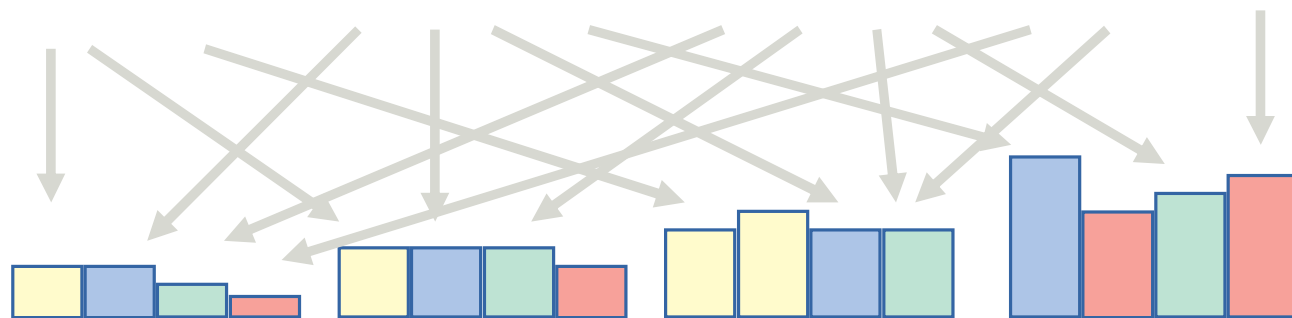


30k x 30k, Speedup vs 1 node

80k x 80k, Speedup vs 2 nodes

# Sort

0. Initial Data
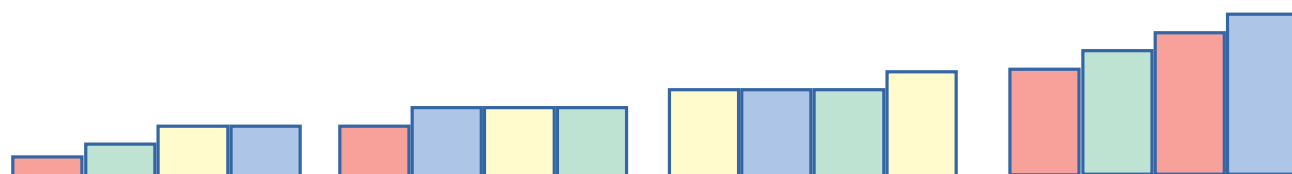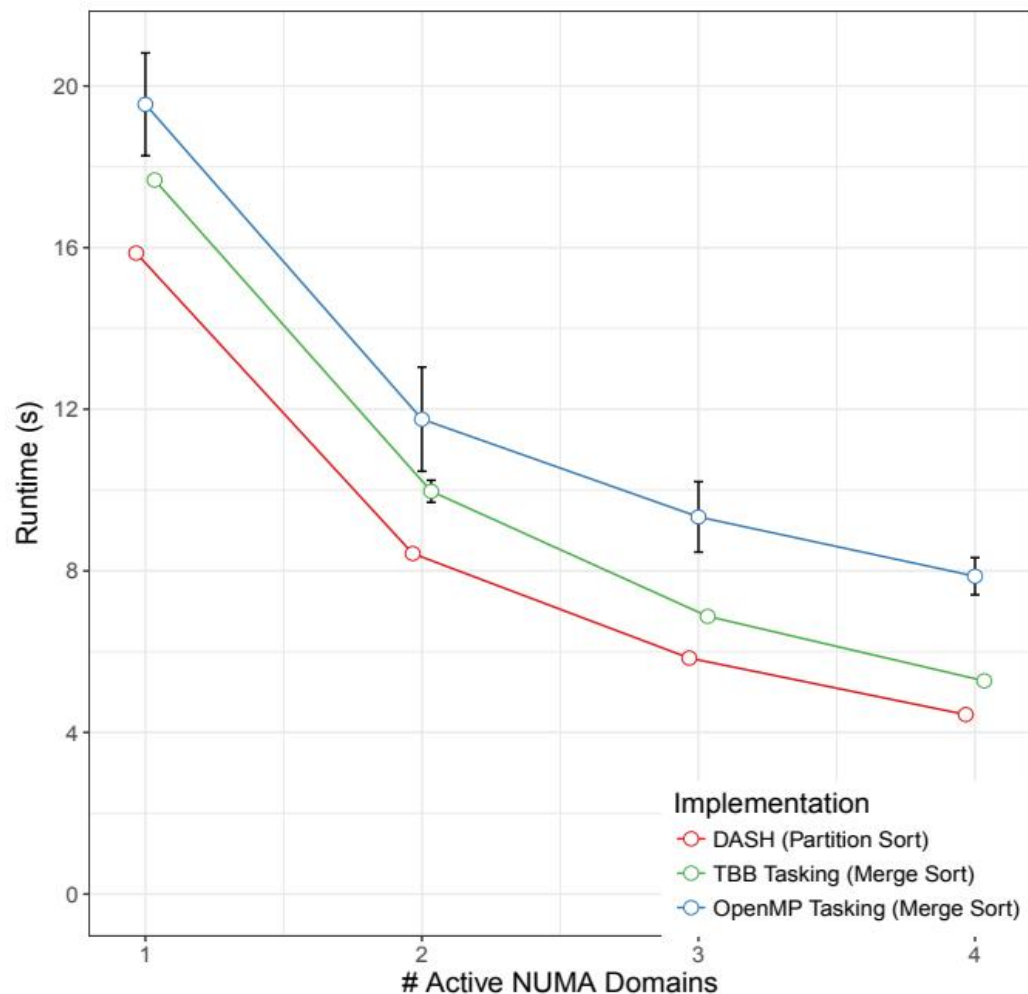
1. Local Sort

2. Determine Partitions

3. All-to-all

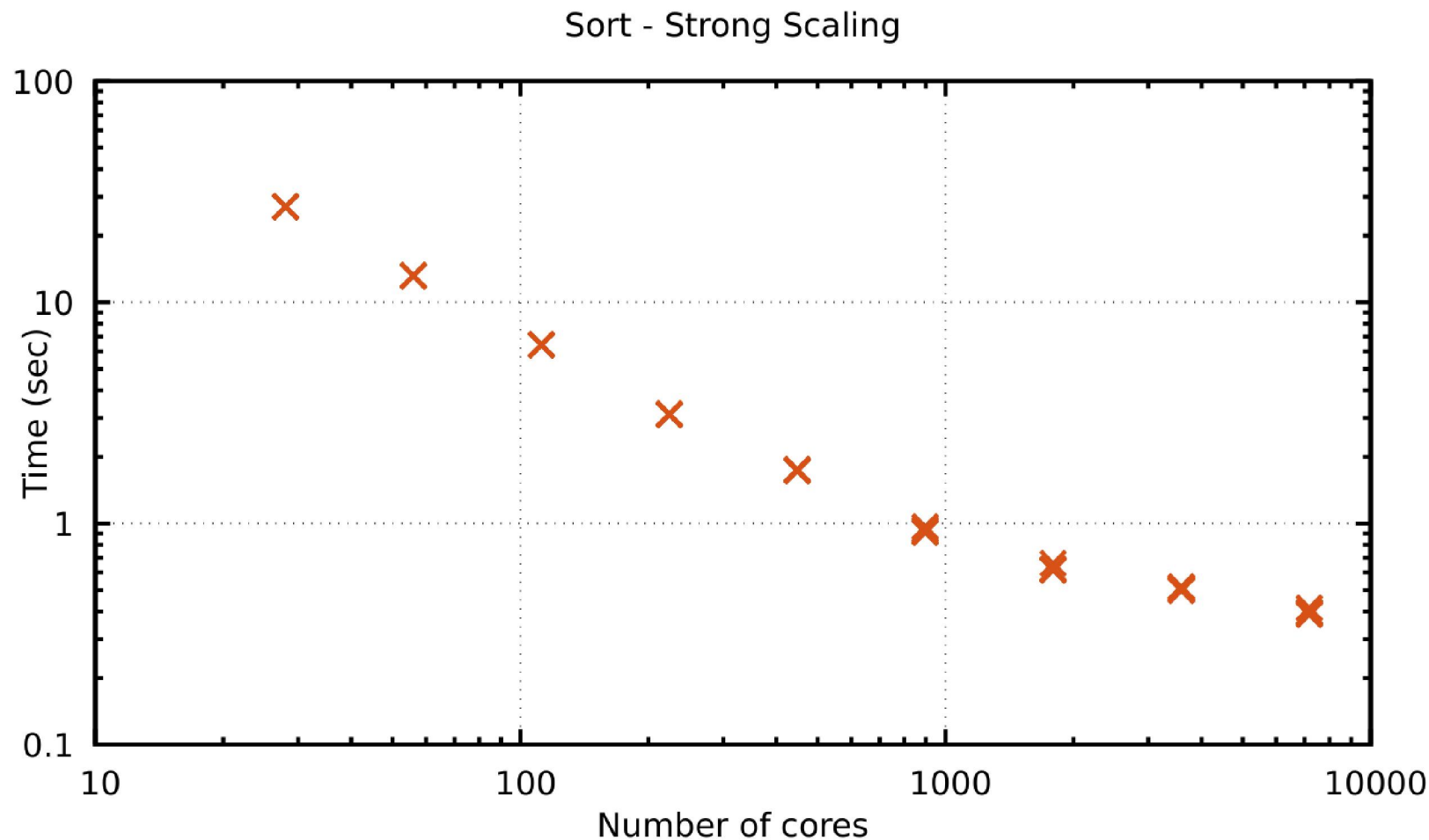4. Final Sort

- Platform: 1 Node of SuperMUC Phase 2
  - 2 sockets, 28 cores, 4 NUMA domains
  - Sorting 4900 MB of data (int)



- PGAS approach beneficial because data locality is primary concern
- Partition-based sort algorithm moves data only once

- # Strong scaling on SuperMUC (HW)
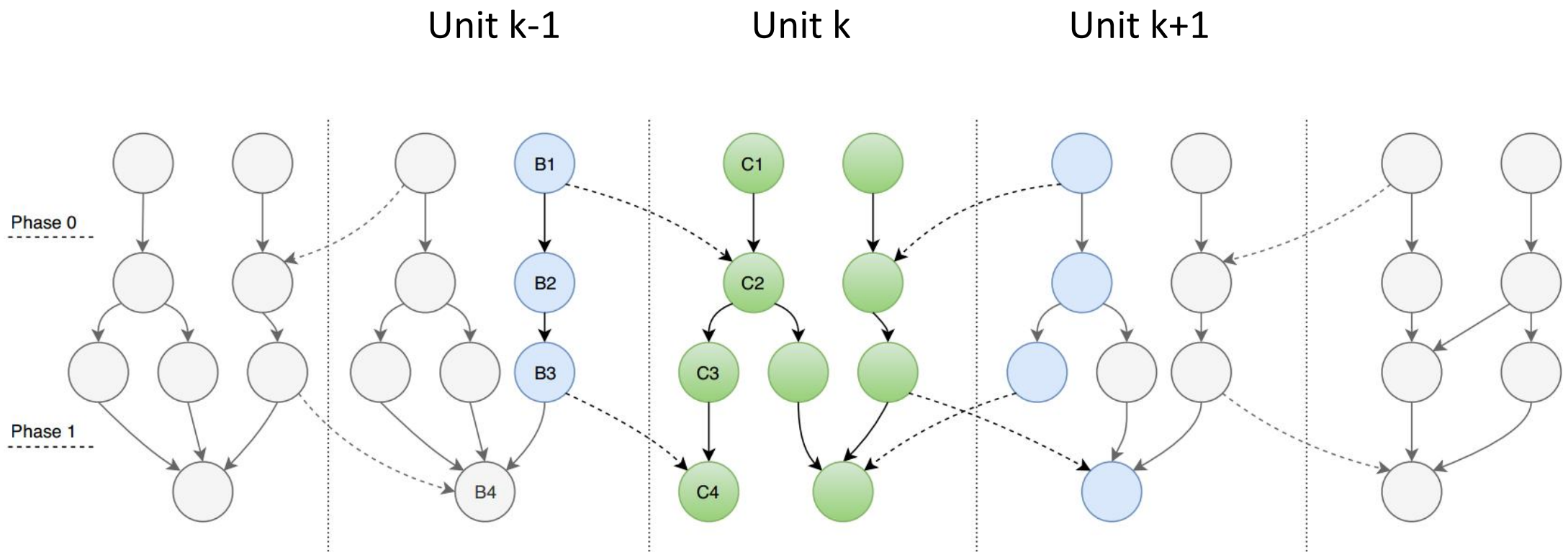  - Sorting ~19 GB data (DP numbers)



*"Engineering a Scalable Histogram Sort",* R. Kowalewski, P. Jungblut, K. Fürlinger, upcoming publication.

# Tasking

- The DASH tasking model is inspired by OpenMP, but can express dependencies in the global address space

| OpenMP Tasks | DASH Tasks |
|---|---|
| #pragma omp task | dash::async(...) w/ lambda and dependency spec. |
| Data dependencies in **node-local memory only** | Data dependencies can **span global address space** (multiple nodes) |
| One scheduler instance | N distributed scheduler instances |
| Dependencies between **previously generated** sibling tasks | Dependencies between sibling tasks in **previous phase** on any node |
| in / out / inout dependencies | in / out / copyin dependencies |

```
 1  for (int k = 0; k < num_blocks; ++k) {
 2    if (mat.block(k,k).is_local()) {
 3      dash::async([&](){ potrf(matrix.block(k,k)); },
 4        dash::out(mat.block(k,k)));
 5    }
 6
 7    // advance to next phase
 8    dash::async_barrier();
 9    for (int i = k+1; i < num_blocks; ++i)
10      if (mat.block(k,i).is_local())
11        dash::async([&](){
12            trsm(cache[k][k], matrix.block(k,i)); },
13          dash::copyin(mat.block(k,k), cache[k][k]),
14          dash::out(mat.block(k,i)));
15
16    // advance to next phase
17    dash::async_barrier();
18    for (int i = k+1; i < num_blocks; ++i) {
19      for (int j = k+1; j < i; ++j) {
20        if (mat.block(j,i).is_local()) {
21          dash::async([&](){
22              gemm(cache[k][i],
23                  cache[k][j], mat.block(j,i)); },
24            dash::copyin(mat.block(k,i), cache[k][i]),
25            dash::copyin(mat.block(k,j), cache[k][j]),
26            dash::out(mat.block(j,i)));
27        }
28      }
29
30      if (mat.block(i,i).is_local()) {
31        dash::async([&](){
32            syrk(cache[k][i], mat.block(i,i)); },
33          dash::copyin(mat.block(k,i), cache[k][i]),
34          dash::out(mat.block(i,i)));
35      }
36    }
37    // advance to next phase
38    dash::async_barrier();
39  }
40  // wait for all tasks to execute
41  dash::complete();
```

potrf

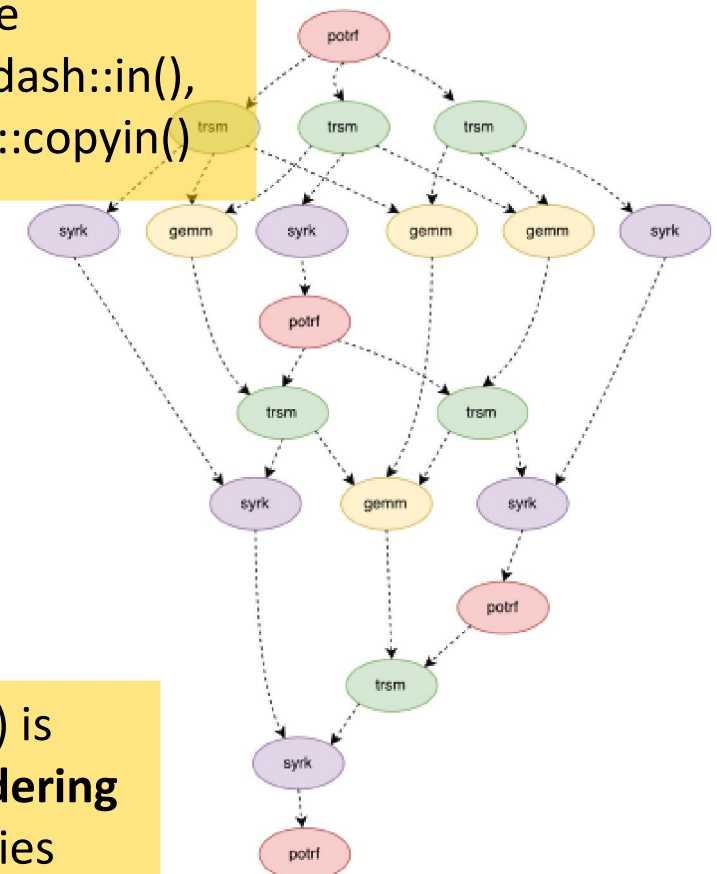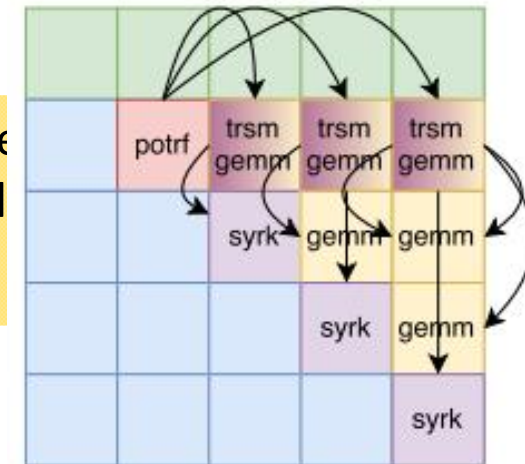trsm

gemm

syrk

sh::async() de...
**sk** using C++ l...
expressions

**dependencies** are
expressed using dash::in(),
dash::out(), dash::copyin()
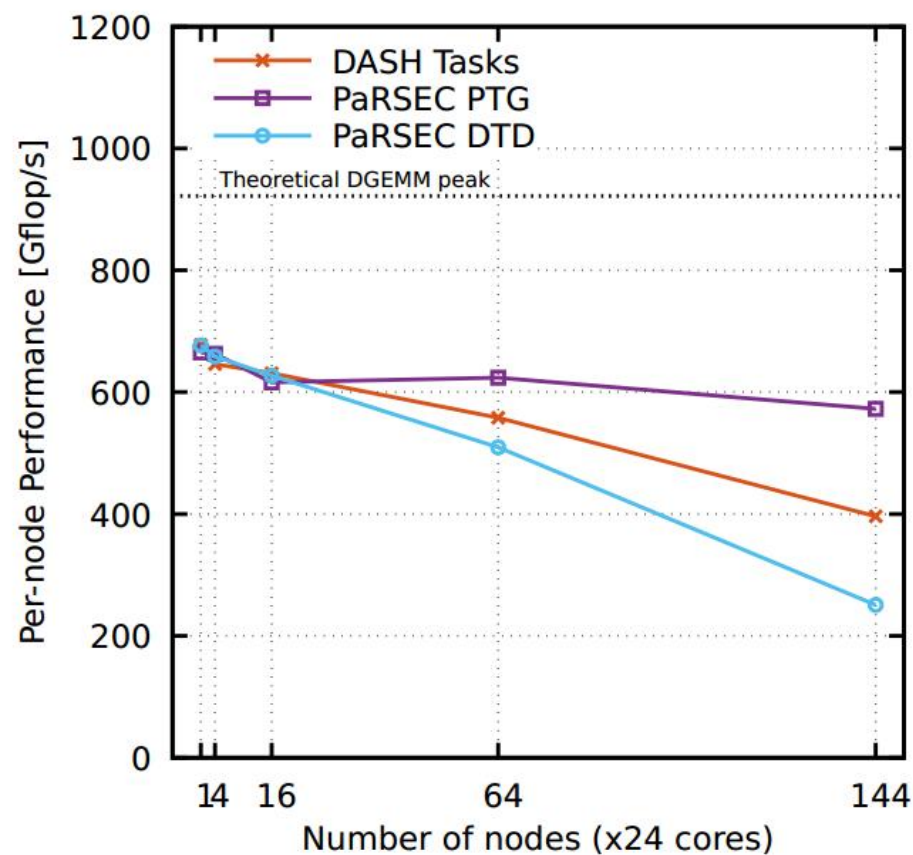
dash::async_barrier() is
used to **establish ordering**
between dependencies

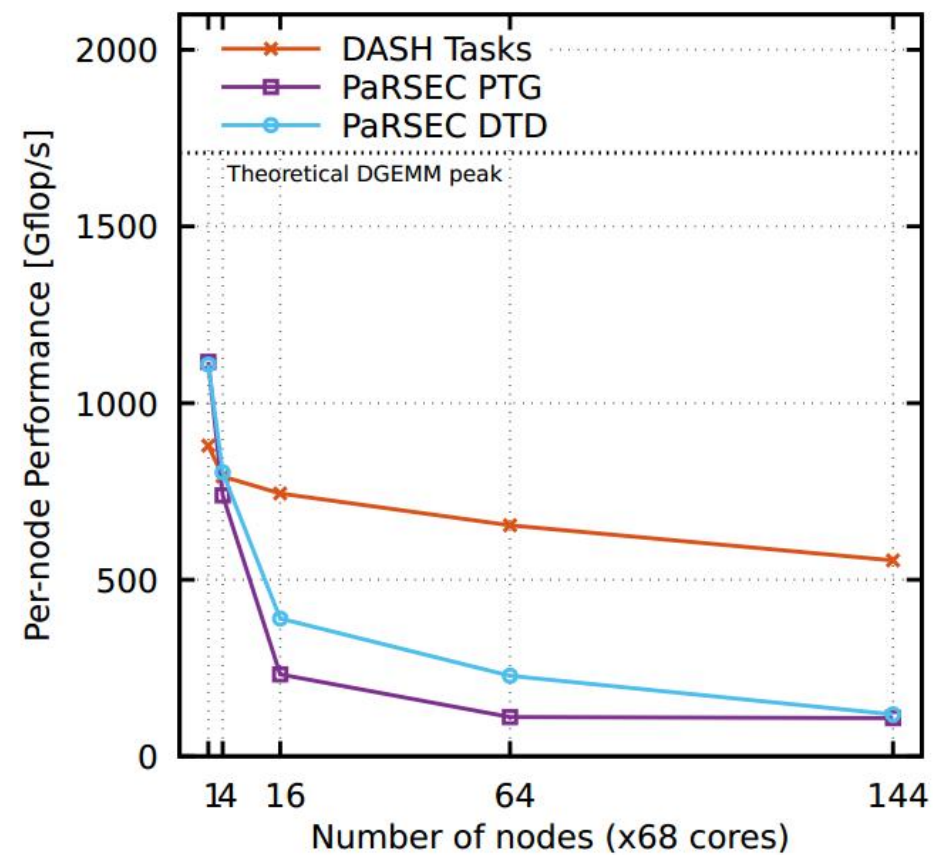- Weak scaling study (20k x 20k matrix per node)
  - Compares with PaRSEC (http://icl.utk.edu/parsec/)
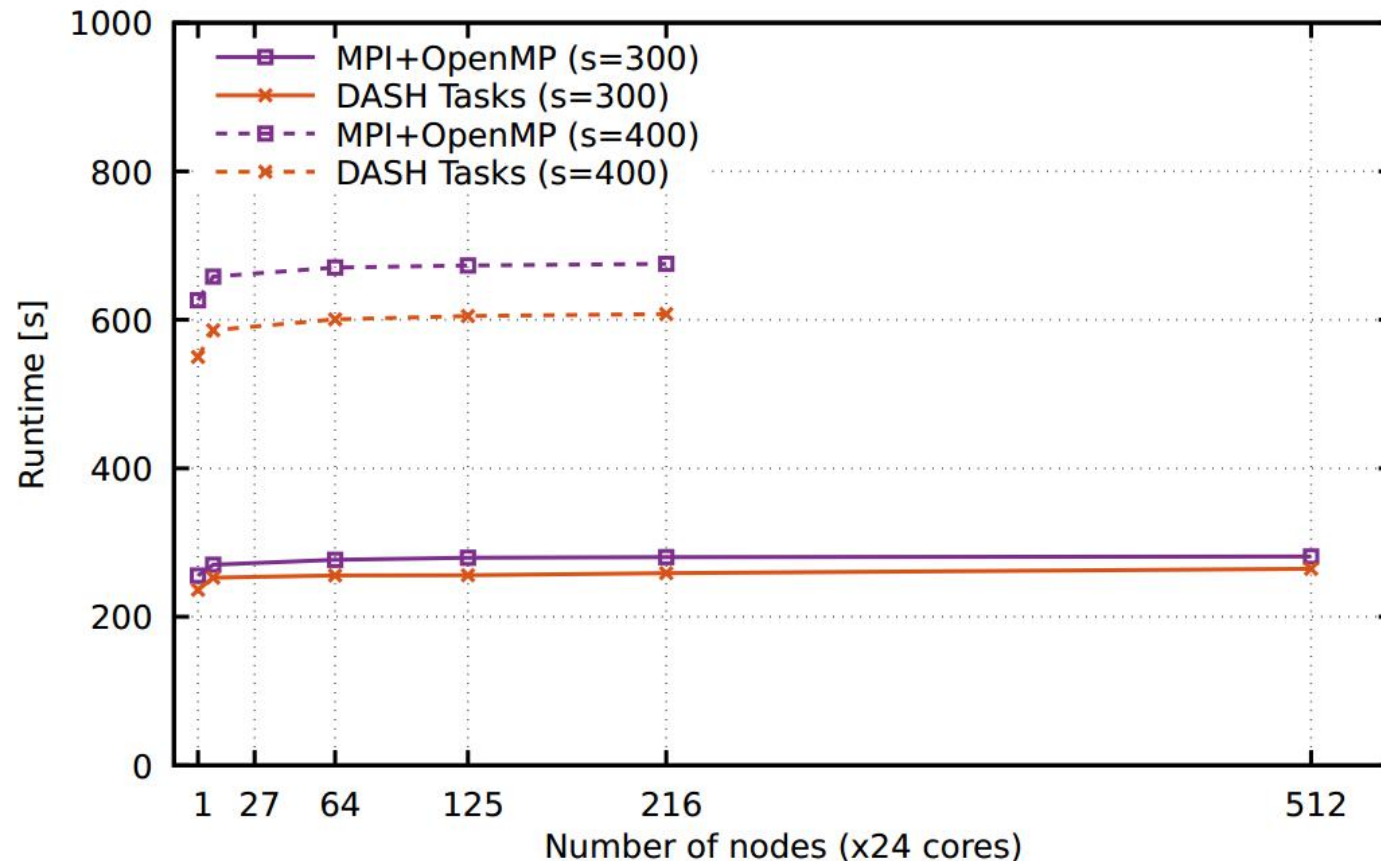


Cray XC40 (Intel SandyBridge)

Oakforest-PACS (KNL)

- Weak scaling on Cray XC40 (Intel SandyBridge)



*"Global Task Data-Dependencies in PGAS Applications",* Joseph Schuchart and Jose Gracia, upcoming publication.

- Collaboration with the **MYX** and **ESSEX-II** projects: Workshops on *"Parallel Programming Models - Productivity and Applications"* and *"Programming and Computing for Exascale and Beyond"*
  - Aachen (15 March, 2018) and Tokyo (30 October, 2018)
  - Plans to use **DART** as a use-case for the **MUST** correctness checking tool

- BMBF Project **MEPHISTO** (02/2017 – 01/2020)
  - TU Dresden (lead) + HZDR + LMU Munich
  - Integration of ALPAKA (portable compute kernel abstraction) with DASH (data structure abstraction)

# Acknowledgements

- Funding

- The DASH Team

- DASH is on GitHub
  - https://github.com/dash-project/dash/