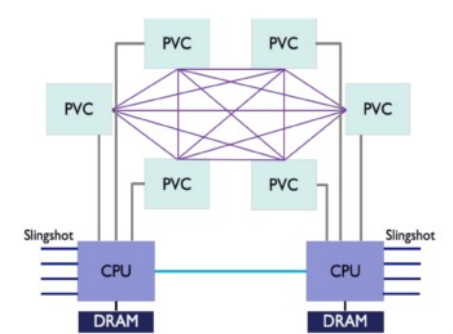
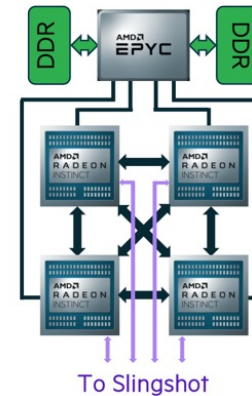
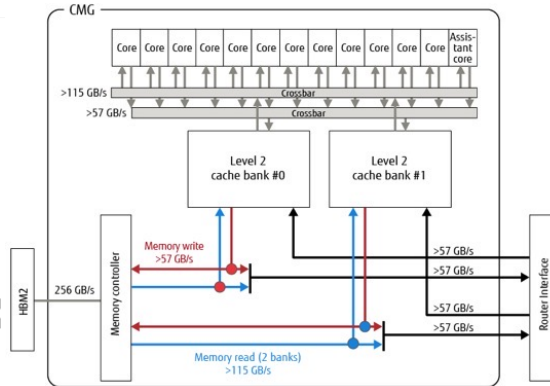
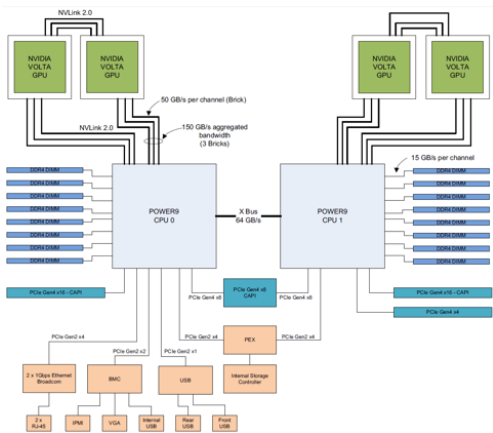


---

# Achieving High Performance with Node-Aware Algorithms

William Gropp  
wgropp.cs.Illinois.edu

# HPC Nodes are Increasingly Complex



## DOE Sierra

- Power 9 with 4 NVIDIA Volta GPU
- 4320 nodes

## DOE Summit similar, but

- 6 NVIDIA GPUs/node
- 4608 nodes

## Fugaku

- Fujitsu A64FX (includes Vector Extensions)
- 158,976 (+) nodes

## DOE Frontier

- AMD with 4 AMD GPU
- 100+ racks

NCSA Delta similar but NVIDIA GPUs and fewer racks 😊

## DOE Aurora

- Intel SR with 6 Intel Ponte Vecchio GPUs
- Being deployed, >9K nodes

---

# Hardware Implications For Programs

- Heterogeneity in many ways
  - Processor – complex compute modes with scalar and vector
  - Many (but not all) include separate accelerators (GPUs and others)
  - Memory – Cache was bad enough; now HBM, other
  - I/O – Burst buffers (often violating POSIX semantics), on node, central, remote (cloud)
- For algorithm developer and programmer, the issue is *Performance Heterogeneity*
  - Whether the implementation uses more than one chip(let) isn't the issue – can you see performance impact of the different elements?
  - Even vectorization counts as performance heterogeneity in this view
    - Compilers still not great at vectorizing code, and often algorithmic changes needed to take full advantage of vectorization (which specializes code, makes it hard to reason about performance)
- Impacts algorithm choice and program realization

# Algorithm Considerations

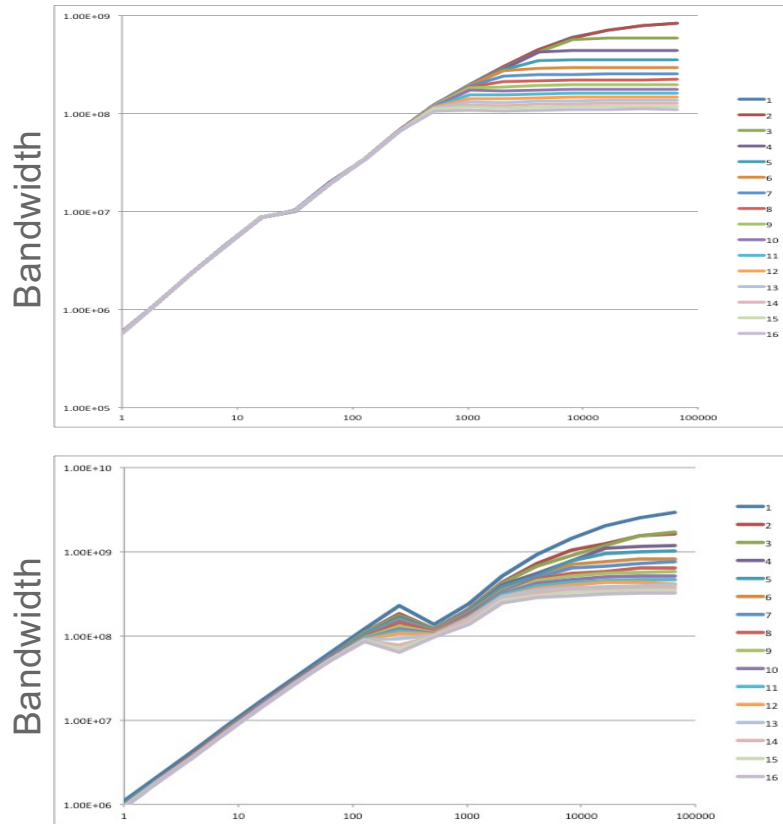
- Start with the choice of mathematical model/numerical method
  - E.g., higher-order approximations for finite difference/element/volume trade floating point operations, data motion, and data size
  - Higher level choices can provide better locality
    - E.g., nonlinear Schwarz, with “local” nonlinear solves
- Performance models needed to guide algorithm design/choice
  - Model does *not* need to be precise – just good enough to guide
  - This is fortunate, as highly accurate performance models are very difficult to create and validate
  - But they need to be accurate enough – and many models haven’t kept up with the evolution of architectures
- One Example: Node-aware algorithms
  - Performance model captures basic system hierarchy at node level
  - Avoid redundant data copies; optimize data motion for HW characteristics
  - Suggests a different approach for process topology mapping...

---

# MPI On Multicore Nodes

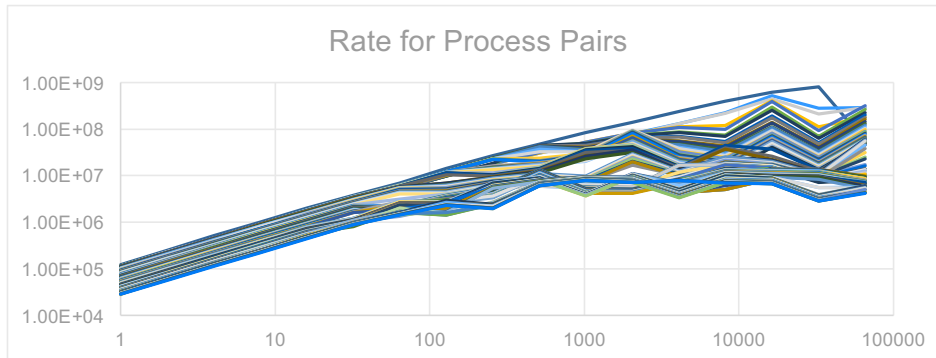
- MPI Everywhere (single core/single thread MPI processes) still common
  - Easy to think about
  - We have good performance models (or do we?)
- In reality, there are issues
  - Memory per core declining
    - Need to avoid large regions for data copies, e.g., halo cells
    - MPI implementations could share internal table, data structures
      - May only be important for extreme scale systems
  - MPI Everywhere implicitly assume uniform communication cost model
    - Limits algorithms explored, communication optimizations used
- Even here, there is much to do for
  - Algorithm designers
  - Application implementers
  - MPI implementation developers
- One example: Can we use the single core performance model for MPI?
  - $T = s + r n$
  - Widely used and effective for designing parallel algorithms

# Rates Per MPI Process

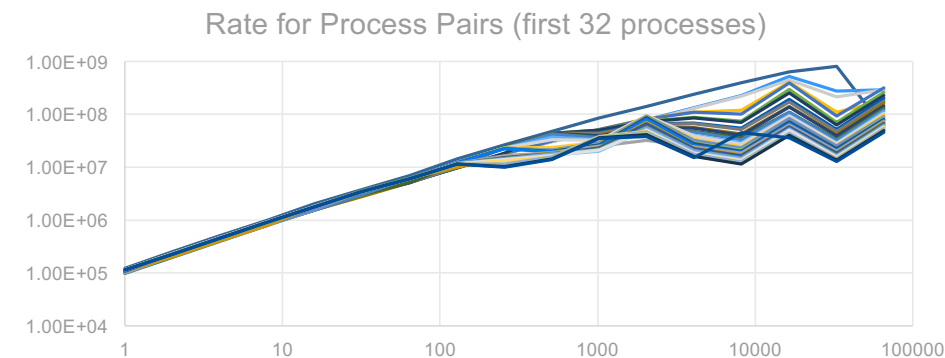


- Ping-pong between 2 nodes using 1-16 cores on each node
- Top is BG/Q, bottom Cray XE6
- “Classic” model predicts a single curve – rates independent of the number of communicating processes

# Rates Per MPI Process: 128 cores

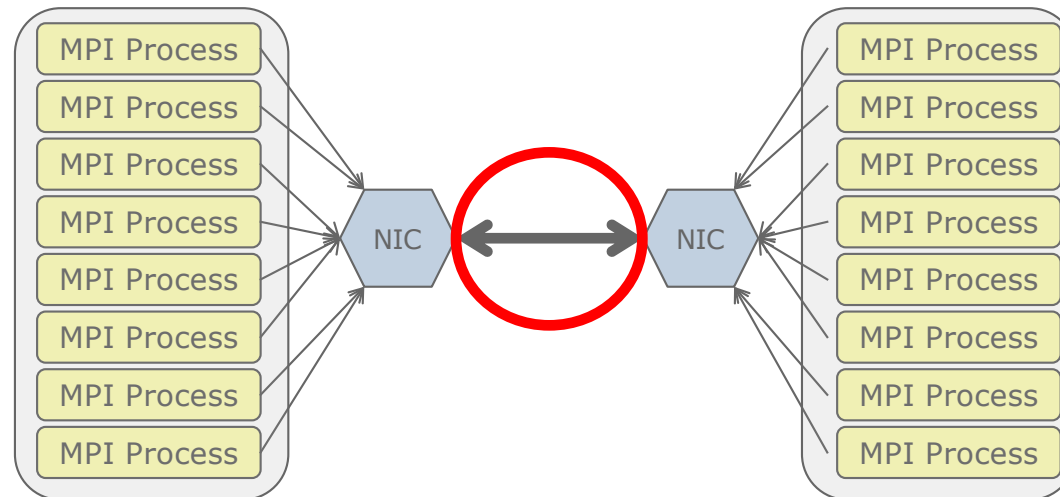


- Increasing core count makes the situation more complex
- Note roughly similar behavior for first 32 processes
  - 1 process / core
  - 64 cores/socket
- As before, classic model predicts a single curve – rate depends only on length, independent of number of communicating processes



# Why this Behavior?

- The  $T = s + r n$  model predicts the *same* performance independent of the number of communicating processes
  - What is going on?
  - How should we model the time for communication?

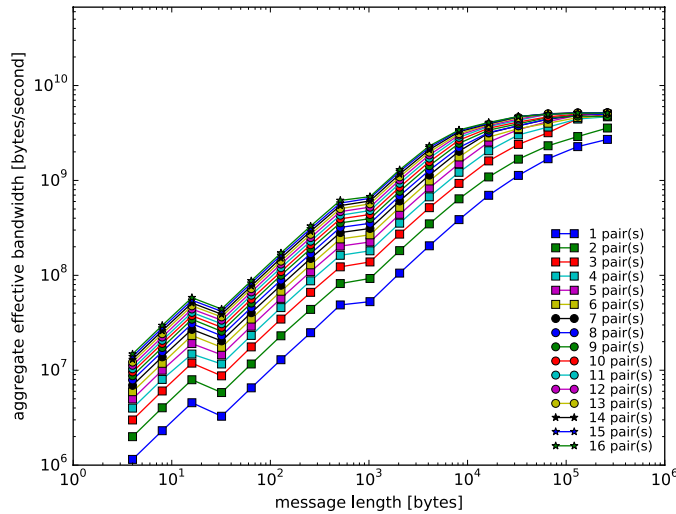




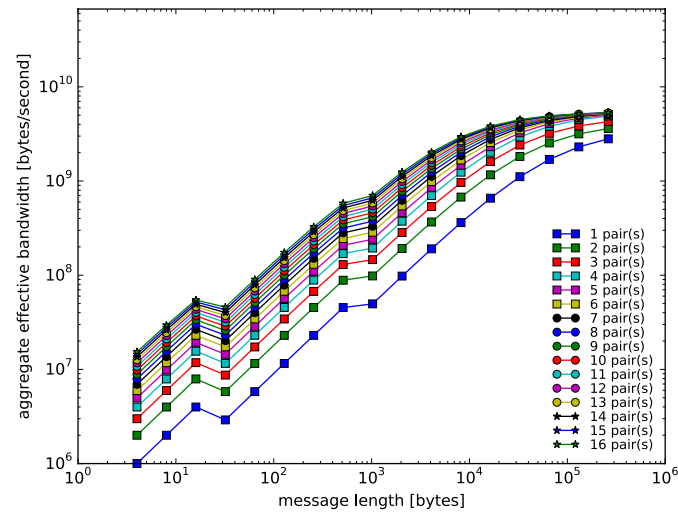
# A Slightly Better Model

- For  $k$  processes sending messages, the sustained rate is
  - $\min(R_{\text{NIC-NIC}}, k R_{\text{CORE-NIC}})$
- Thus
  - $T = s + k n / \min(R_{\text{NIC-NIC}}, k R_{\text{CORE-NIC}})$
- Note if  $R_{\text{NIC-NIC}}$  is very large (very fast network), this reduces to
  - $T = s + k n / (k R_{\text{CORE-NIC}}) = s + n / R_{\text{CORE-NIC}}$
- This model is approximate; additional terms needed to capture effect of shared data paths in node, contention for shared resources, etc.
- But this new term is by far the dominant one
- This is the *max-rate* model (for performance limited by the maximum available bandwidth)
  - Logp model has a similar limitation and needs a similar modification

# Comparison on Cray XE6



Measured Data

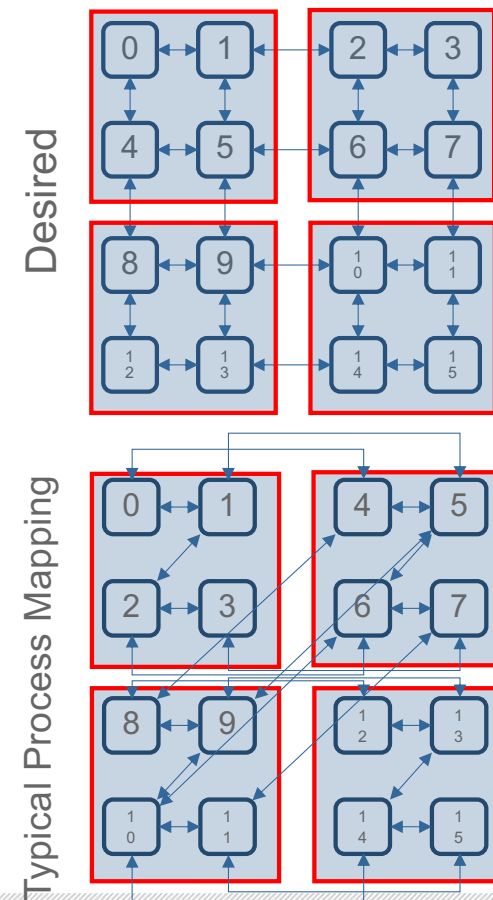


Max-Rate Model

*Modeling MPI Communication Performance on SMP Nodes: Is it Time to Retire the Ping Pong Test*, W Gropp, L Olson, P Samfuss, Proceedings of EuroMPI 16, <https://doi.org/10.1145/2966884.2966919>

# Performance Model to Algorithm

- Performance measurements of halo exchange show poor communication performance
  - Bandwidth per process low relative to “ping pong” measurements
  - Easy target – blame contention in the network
- But common default mapping of processes to nodes leads to more off-node communication
  - The max rate model predicts reduced performance once  $R_{\text{NIC-NIC}}$  limit reached
- We can use this to create a better, and *simpler*, implementation of `MPI_Cart_create`

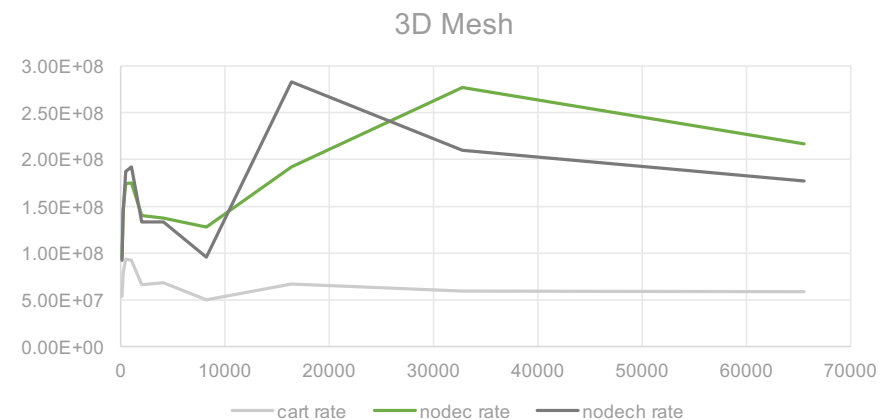
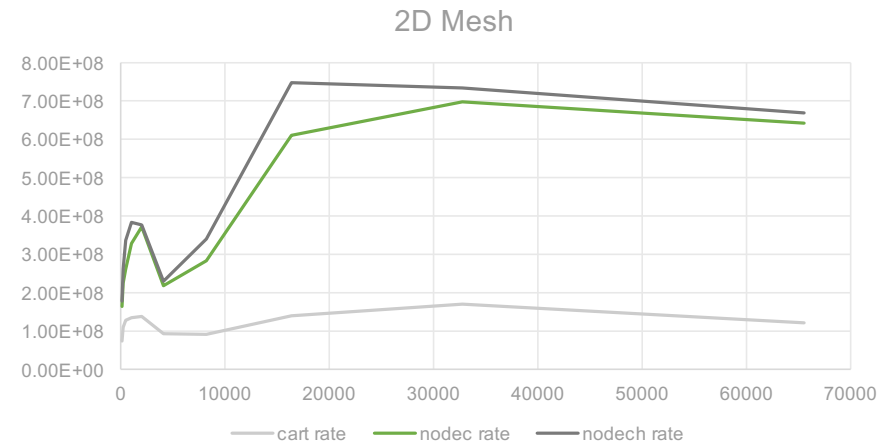


# Building A Better MPI\_Cart\_create

- Hypothesis: A better process mapping **within** a node will provide significant benefits
  - **Ignore** the internode network topology
    - Vendors have argued that their network is fast enough that process mapping isn't necessary
    - They may be (almost) right – once data enters the network
- Idea for Cartesian Process Topologies
  - Identify nodes (see MPI\_Comm\_split\_type)
  - Map processes *within* a node to minimize **internode** communication
    - Trading **intranode** for **internode** communication
    - *Using Node Information to Implement MPI Cartesian Topologies*, Gropp, William D., Proceedings of the 25th European MPI Users' Group Meeting, 18:1–18:9, 2018 <https://dl.acm.org/citation.cfm?id=3236377>
    - *Using Node and Socket Information to Implement MPI Cartesian Topologies*, Parallel Computing, 2019 <https://doi.org/10.1016/j.parco.2019.01.001>

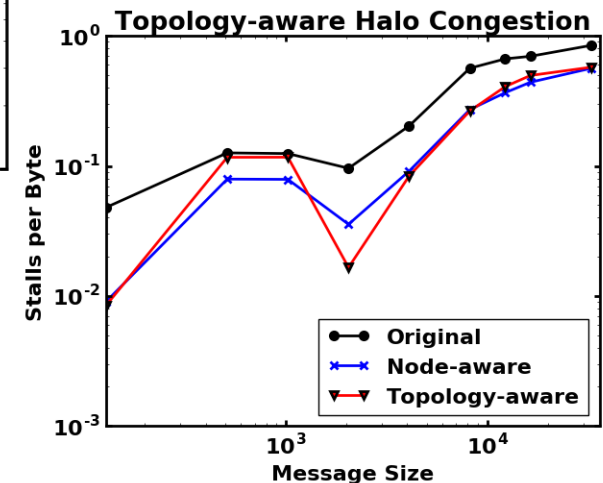
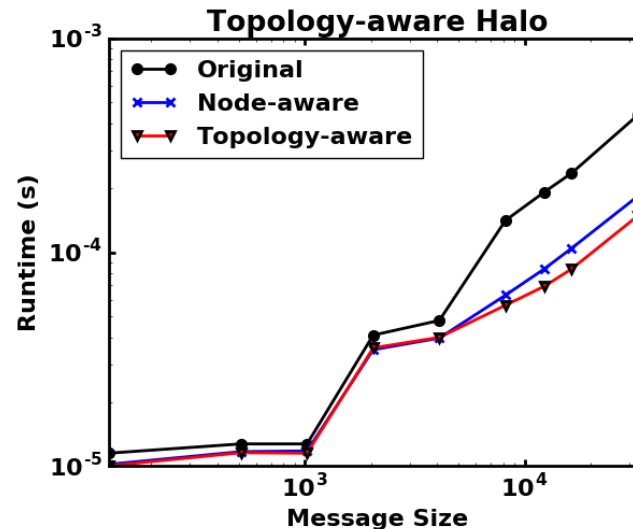
# Increasing Core Count Makes Proper Mapping More Important

- Cartesian mapping on Delta
  - CPU nodes have 2 AMD Milan x 64 cores each (GPU nodes have 1 AMD Milan and 4 A100 or A40 NVIDIA GPUs)
  - Slingshot network (mostly – NIC update coming)
  - Performance in B/s (higher is better)
- Default mapping provides poor performance
  - Cart is MPI\_Cart\_create – also MPI\_COMM\_WORLD
  - Nodec uses node-awareness, inspired by max-rate model
  - Nodech extends to socket (3-level)



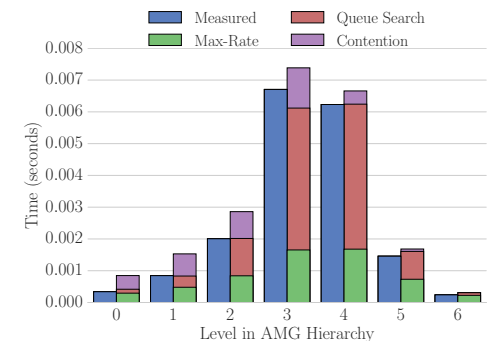
# How Important is Network Topology?

- No answer yet, but...
- 432 nodes, 3D halo exchange on Blue Waters
  - Requested a cube of nodes, used non-standard routines to implement mapping for network topology
- Part of study into scalable Krylov methods (looking to avoid the blocking MPI\_Allreduce)
- Nodecart version provides most of the benefit with no need for network topology information
- Some (nontrivial) further benefit possible by taking network topology into account
- But the largest contribution comes from node-awareness
- Thanks to Paul Eller for these results



# Node Aware Algorithms

- Can use max rate model to design better algorithms for SpMV, related operations
  - May need to add additional terms – as in queue search overhead
  - May have unexpected results – e.g., MPI Direct from GPU not always best on multicore nodes (upcoming paper with Lockhart, Bienz, and Olson)
- Work with Amanda Bienz, Shelby Lockhart, Luke Olson
  - Amanda Bienz, William D. Gropp, and Luke N. Olson. Node aware sparse matrix–vector multiplication. *Journal of Parallel and Distributed Computing*, 130:166–178, 2019.
  - A. Bienz, L. Olson, and W. Gropp. Node-aware improvements to allreduce. In 2019 IEEE/ACM Workshop on Exascale MPI (ExaMPI), pages 19–28, Nov 2019.
  - Amanda Bienz, Luke N. Olson, William D. Gropp, and Shelby Lockhart. Modeling data movement performance on heterogeneous architectures. In 2021 IEEE High Performance Extreme Computing Conference (HPEC), pages 1–7, 2021.



---

# Summary

- Need to enable algorithm development
  - Great to see so many talks at this meeting embracing the need to match algorithms to real hardware – and take advantage of specialization
  - As others have stated, need the “right” performance model to drive algorithm design
  - Max-rate model can guide some node-aware algorithms
  - <something about extensions to multiGPU nodes>