# Avoiding hot-spots on two-level direct networks

Abhinav Bhatele, Nikhil Jain, William D. Gropp, Laxmikant V. Kale
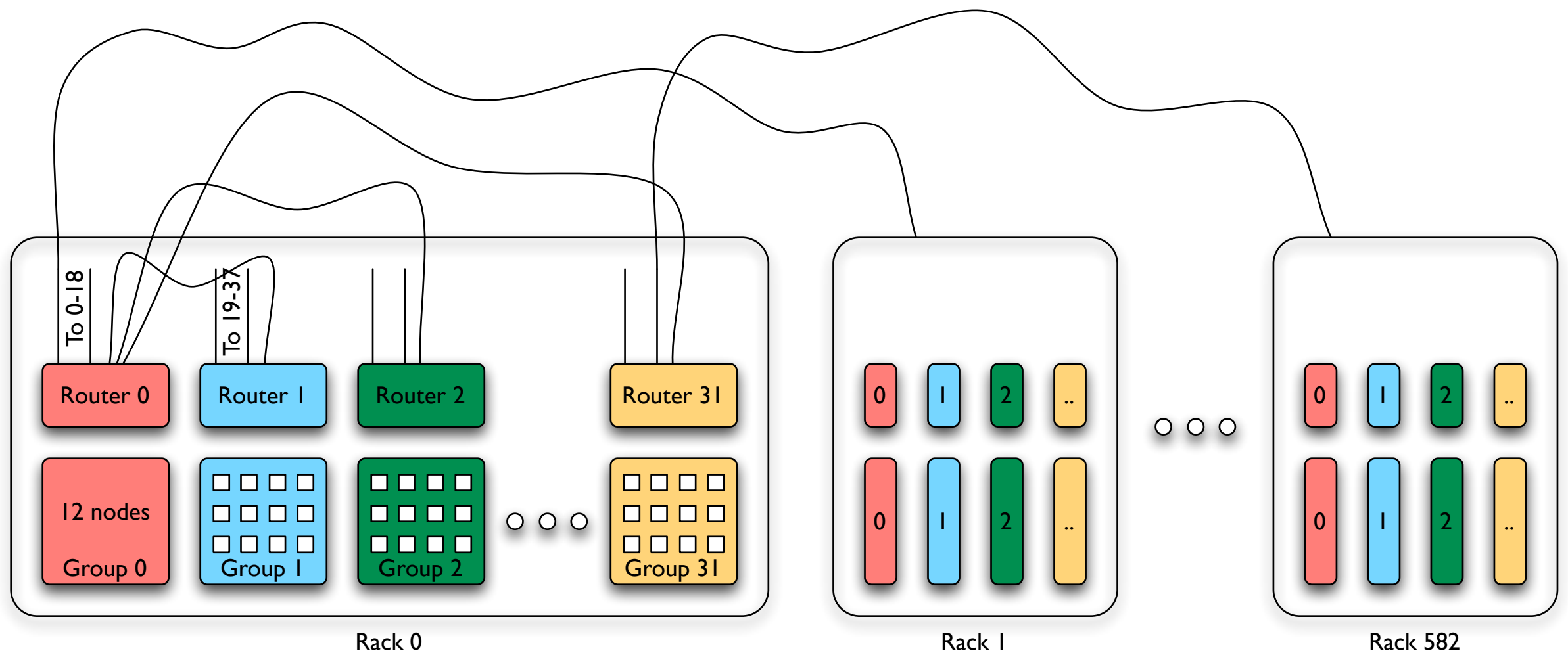
Supercomputing '11 ◆ November 17, 2011

LLNL-PRES-511461

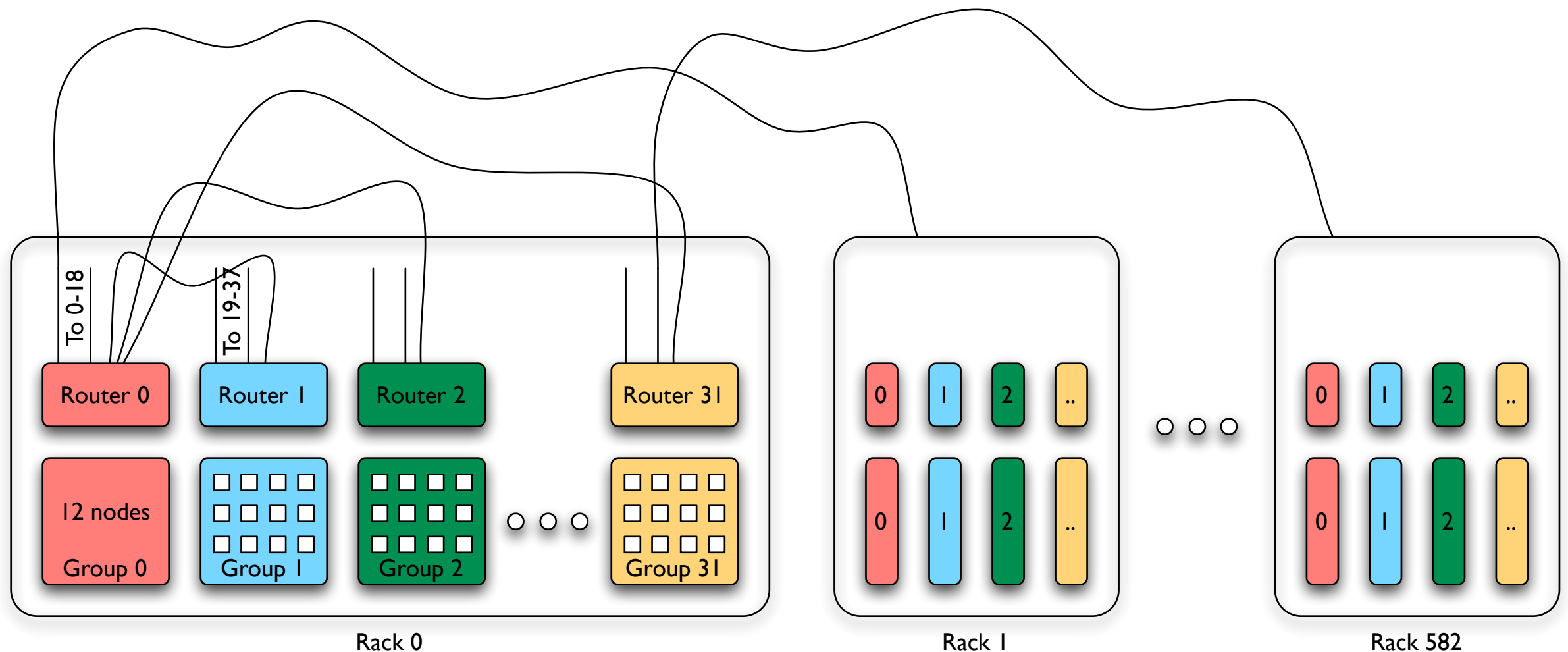Lawrence Livermore National Laboratory, P. O. Box 808, Livermore, CA 94551

# Interconnects for exascale
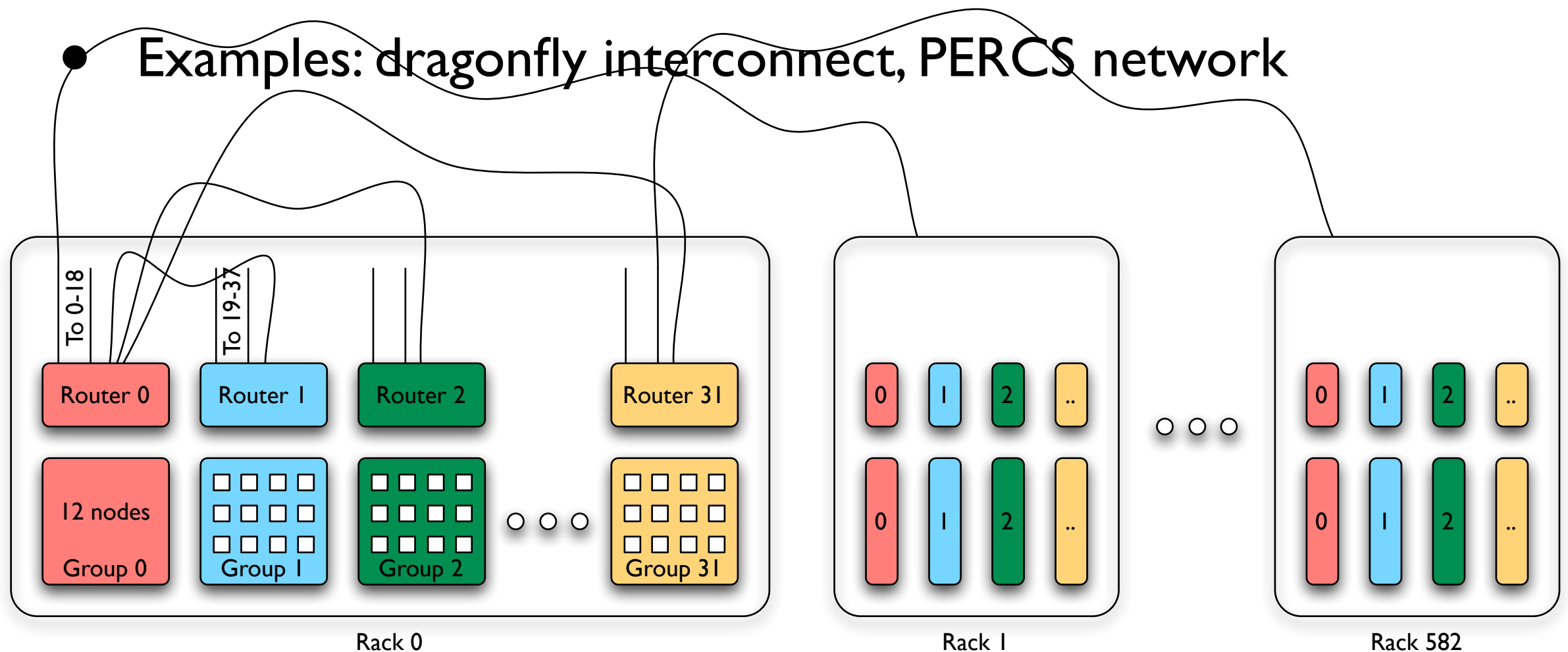
# Interconnects for exascale

- Multi-level direct (all-to-all connection) networks

  - Higher bandwidth links at lower levels

  - Low diameter: few hops on the average
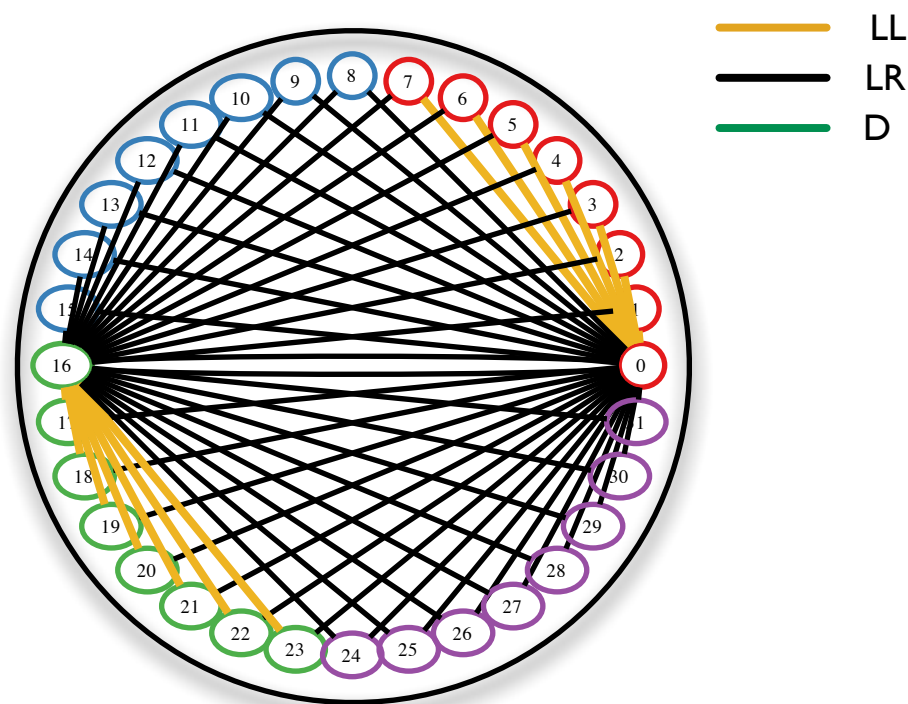
# Interconnects for exascale

- Multi-level direct (all-to-all connection) networks

  - Higher bandwidth links at lower levels

  - Low diameter: few hops on the average

- Examples: dragonfly interconnect, PERCS network

# IBM's PERCS network

- Each QCM has four 8-core POWER7 chips, 8 such QCMs form a drawer, 4 drawers form a supernode

- Two-level network with 512 supernodes

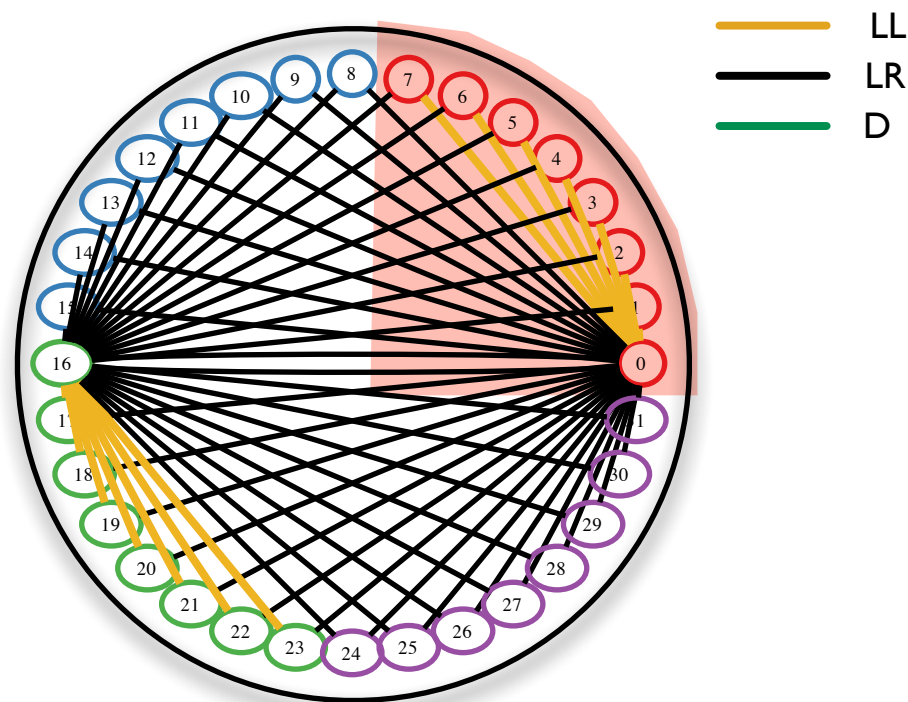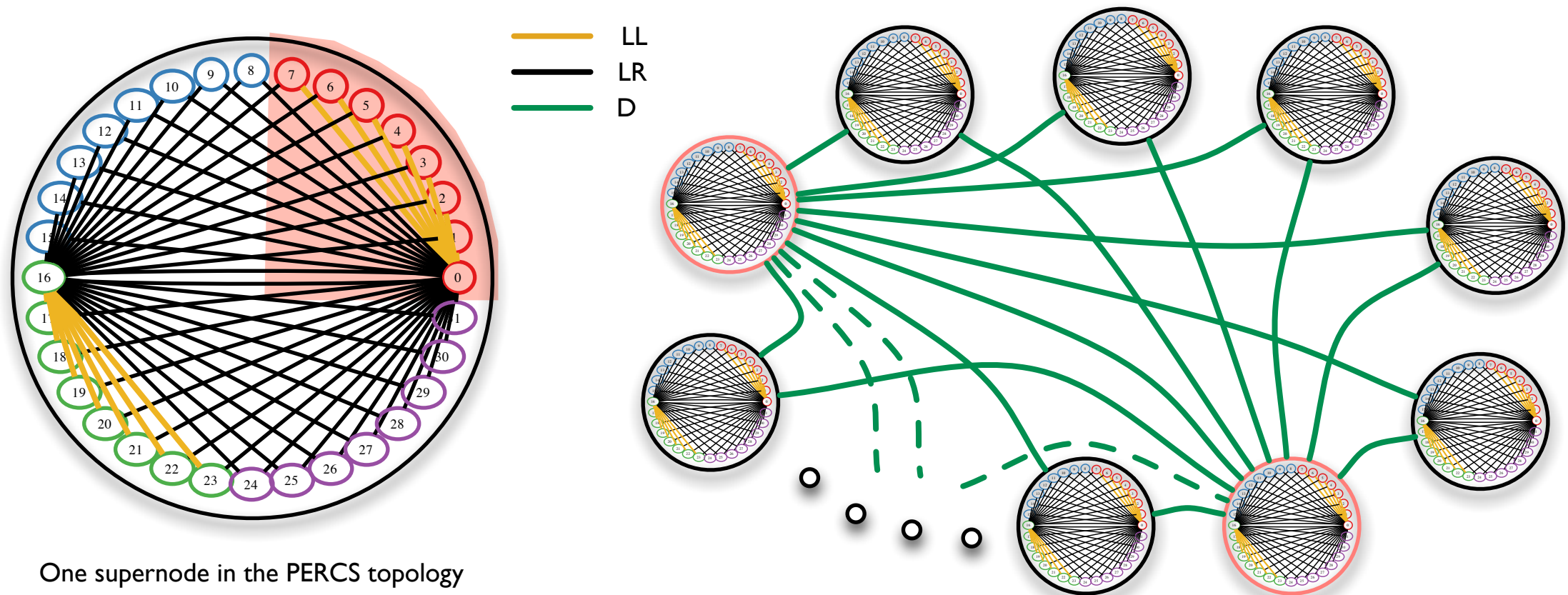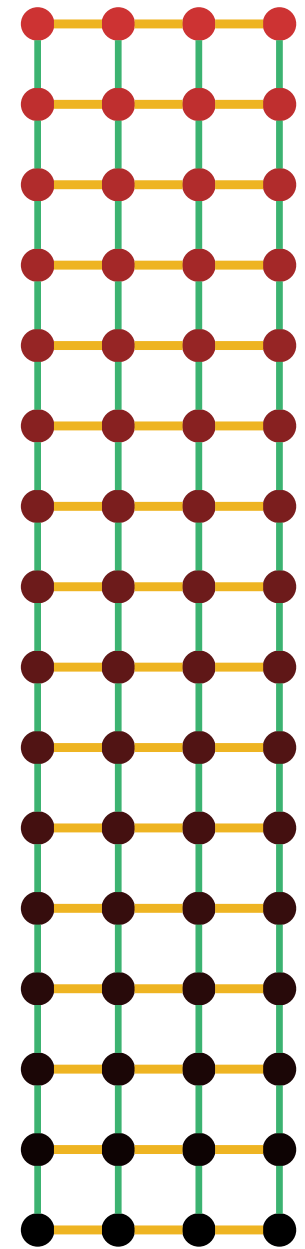- Three types of links: LL (24 GB/s), LR (5 GB/s), D (10 GB/s)



One supernode in the PERCS topology

# IBM's PERCS network

- Each QCM has four 8-core POWER7 chips, 8 such QCMs form a drawer, 4 drawers form a supernode

- Two-level network with 512 supernodes

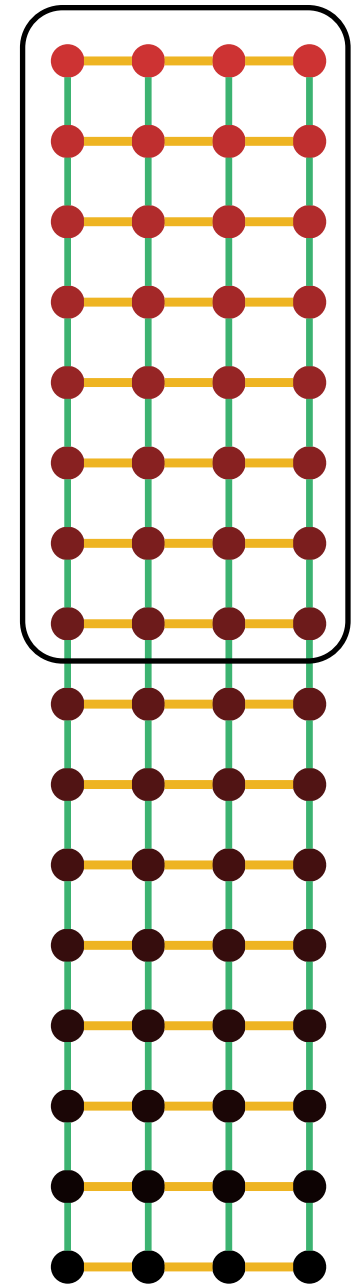- Three types of links: LL (24 GB/s), LR (5 GB/s), D (10 GB/s)



One supernode in the PERCS topology

# IBM's PERCS network

- Each QCM has four 8-core POWER7 chips, 8 such QCMs form a drawer, 4 drawers form a supernode

- Two-level network with 512 supernodes

- Three types of links: LL (24 GB/s), LR (5 GB/s), D (10 GB/s)
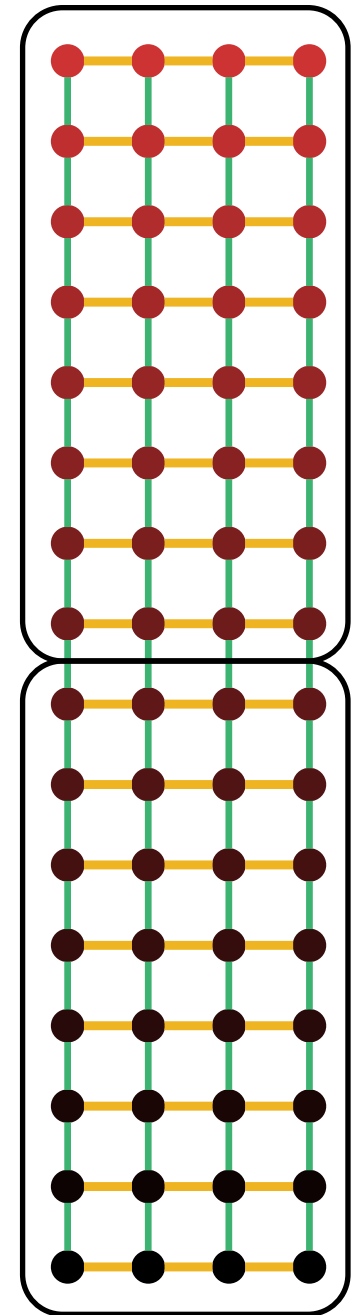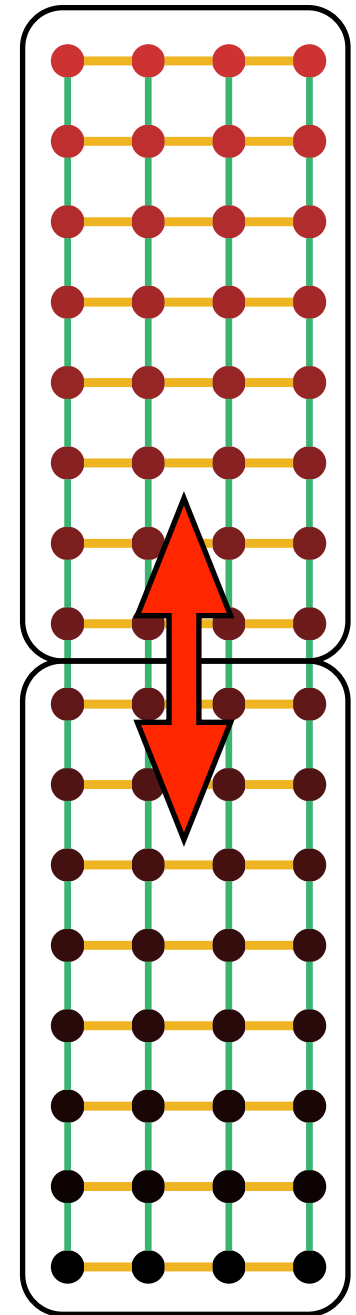


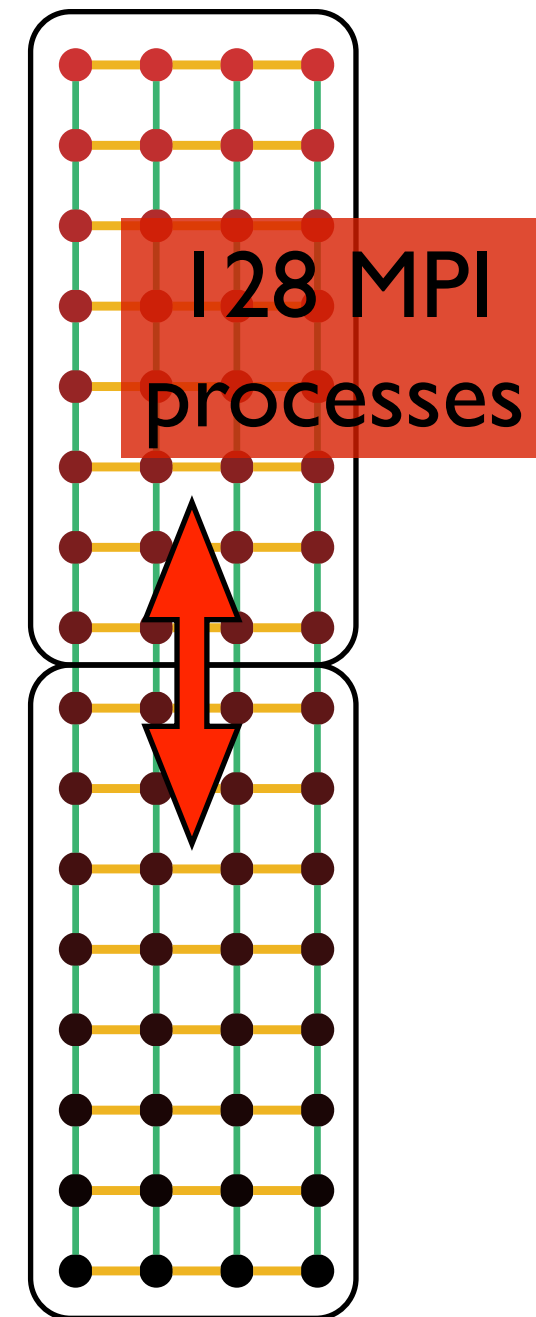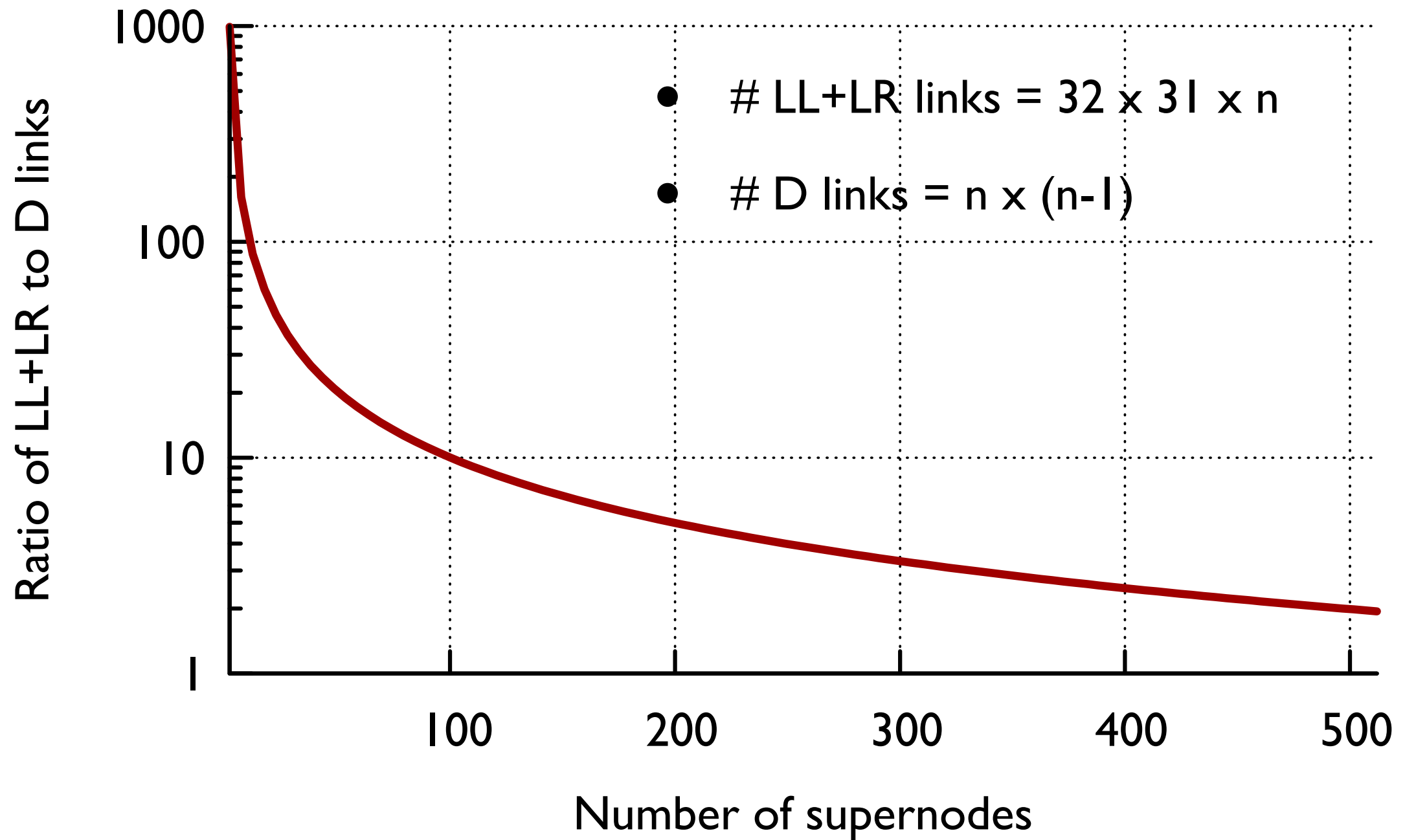One supernode in the PERCS topology

# The 'D' link bottleneck

- Lets say we want to run a 2D Stencil on 2 supernodes (64 QCMs)

- Application communication graph: 16 x 4

- MPI-rank ordered mapping leads to significant contention on the single D link

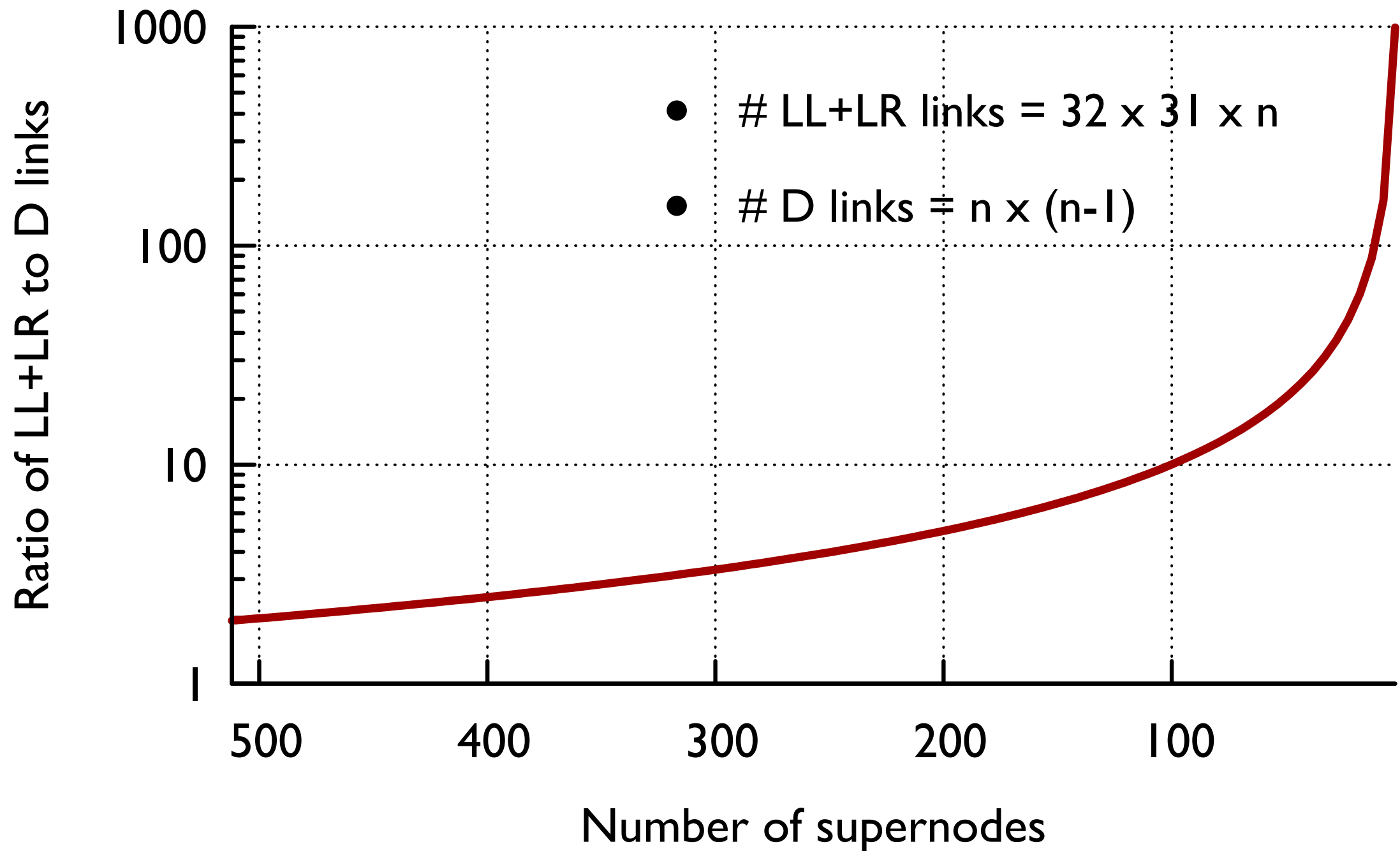- This is true for running any application with O(1) communicating partners per MPI process
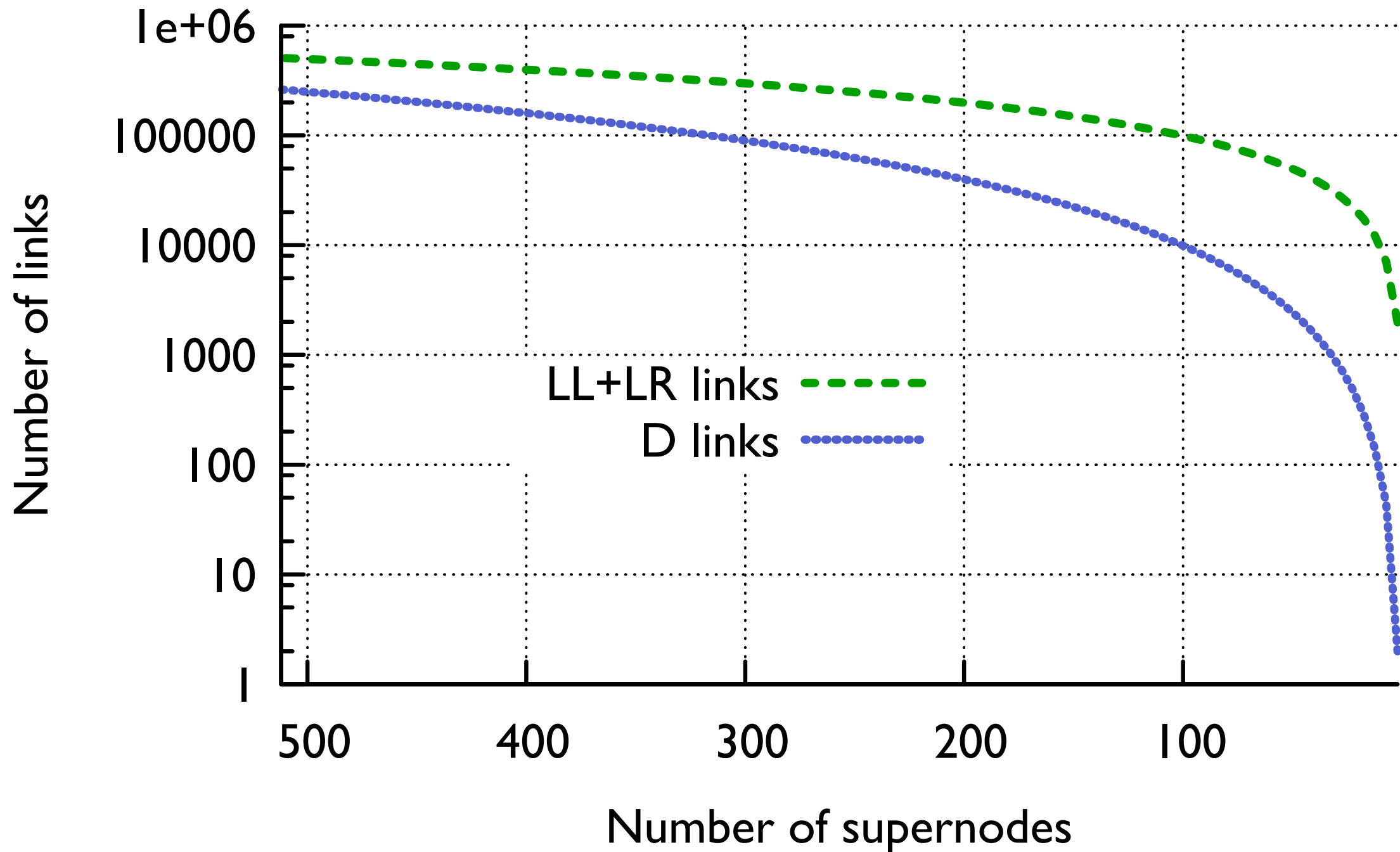
# The 'D' link bottleneck

- Lets say we want to run a 2D Stencil on 2 supernodes (64 QCMs)

- Application communication graph: 16 x 4

- MPI-rank ordered mapping leads to significant contention on the single D link

- This is true for running any application with O(1) communicating partners per MPI process

# The 'D' link bottleneck

- Lets say we want to run a 2D Stencil on 2 supernodes (64 QCMs)

- Application communication graph: 16 x 4

- MPI-rank ordered mapping leads to significant contention on the single D link

- This is true for running any application with O(1) communicating partners per MPI process

# The 'D' link bottleneck

- Lets say we want to run a 2D Stencil on 2 supernodes (64 QCMs)

- Application communication graph: 16 x 4

- MPI-rank ordered mapping leads to significant contention on the single D link

- This is true for running any application with O(1) communicating partners per MPI process

# The 'D' link bottleneck

- Lets say we want to run a 2D Stencil on 2 supernodes (64 QCMs)

- Application communication graph: 16 x 4

- MPI-rank ordered mapping leads to significant contention on the single D link

- This is true for running any application with O(1) communicating partners per MPI process

128 MPI processes

# The 'D' link bottleneck



- # LL+LR links = 32 x 31 x n
- # D links = n x (n-1)

Ratio of LL+LR to D links

Number of supernodes

# The 'D' link bottleneck



- # LL+LR links = 32 x 31 x n
- # D links = n x (n-1)

# The 'D' link bottleneck

# Software stack/runtime choices

- ## Job scheduler:

    - Granularity of allocation: QCM (node), drawer, supernode

    - Contiguous allocation, random allocation or careful topology-aware allocation

- ## Routing:

    - Direct versus random indirect

- ## Mapping

    - Is it important for optimal performance?

# Simulation study – BigSim

# Simulation study − BigSim

- Use BigSim for emulation and simulation of a future machine - detailed packet-level network simulation

  - Compute time prediction, trace collection, simulation

# Simulation study − BigSim

- Use BigSim for emulation and simulation of a future machine - detailed packet-level network simulation

  - Compute time prediction, trace collection, simulation

- Three different benchmarks

  - 2-dimensional five-point stencil

  - 4-dimensional nine-point stencil

  - n-targets multicast pattern

# Simulation study – BigSim

- Use BigSim for emulation and simulation of a future machine - detailed packet-level network simulation

  - Compute time prediction, trace collection, simulation

- Three different benchmarks

  - 2-dimensional five-point stencil

  - 4-dimensional nine-point stencil

  - n-targets multicast pattern

- Two job allocation sizes

  - 64 supernodes

  - 300 supernodes

# Simulation study − BigSim

- Use BigSim for emulation and simulation of a future machine - detailed packet-level network simulation

  - Compute time prediction, trace collection, simulation

- Three different benchmarks

  - 2-dimensional five-point stencil     WRF

  - 4-dimensional nine-point stencil

  - n-targets multicast pattern

- Two job allocation sizes

  - 64 supernodes

  - 300 supernodes

# Simulation study – BigSim

- Use BigSim for emulation and simulation of a future machine - detailed packet-level network simulation

  - Compute time prediction, trace collection, simulation

- Three different benchmarks

  - 2-dimensional five-point stencil    WRF

  - 4-dimensional nine-point stencil    MILC

  - n-targets multicast pattern

- Two job allocation sizes

  - 64 supernodes

  - 300 supernodes

# Simulation study – BigSim

- Use BigSim for emulation and simulation of a future machine - detailed packet-level network simulation

  - Compute time prediction, trace collection, simulation

- Three different benchmarks

  - 2-dimensional five-point stencil      WRF

  - 4-dimensional nine-point stencil      MILC

  - n-targets multicast pattern      NAMD

- Two job allocation sizes

  - 64 supernodes

  - 300 supernodes

# Prediction Methodology

- Compute time prediction: run on a 3.8 GHz Power7 processor to get timings for sequential computation

- Emulation: obtain traces by running on 512-1360 cores of a 1.9 GHz Power5 cluster

- Simulations on one node of a SGI Altix 1000 shared memory machine

# Mapping and Routing

- Default mapping: MPI rank-ordered mapping

- Blocking: at the level of nodes, then drawers and supernodes

- Random mapping on nodes

- Routing choices:

  - Indirect routing w/ default mapping

  - Indirect routing w/ random drawers mapping

# 4-dimensional stencil

- Representative of MILC, a Lattice QCD code

- Each MPI task has 64 x 64 x 64 x 64 elements

- Size of messages exchanged = 2 MB

# Experiments

- Direct routing:

  - Default MPI rank-ordered mapping (DEF)

  - Blocking MPI tasks at the level of nodes (BNM)

  - Blocking at the level of drawers (BDM)

  - Blocking at the level of supernodes (BSM)

  - Random mapping at the level of nodes (RNM)

  - Random mapping at the level of drawers (RDM)

- Indirect routing

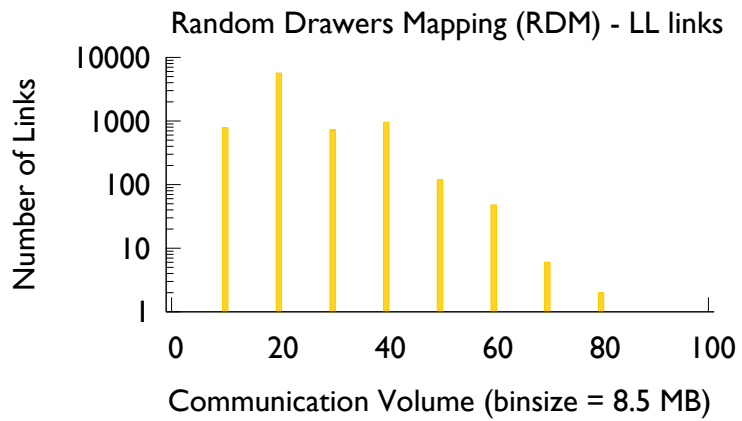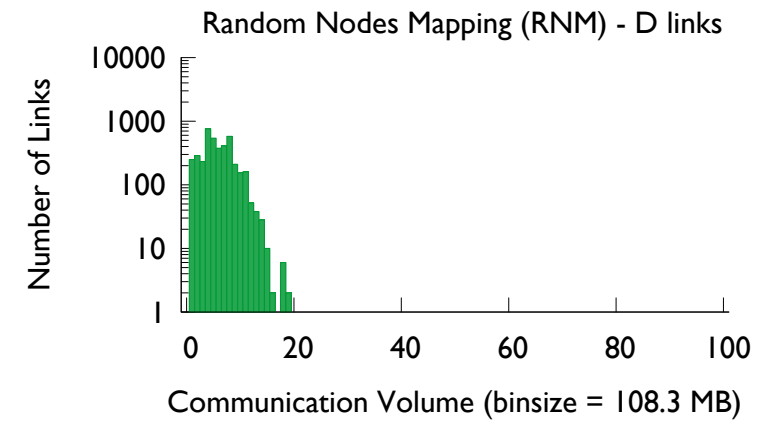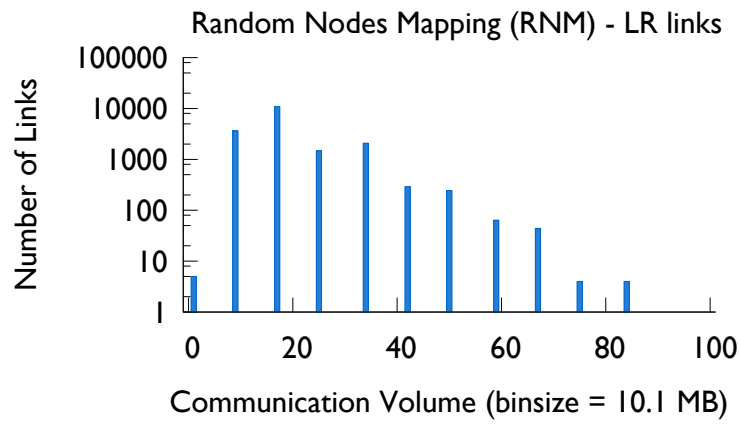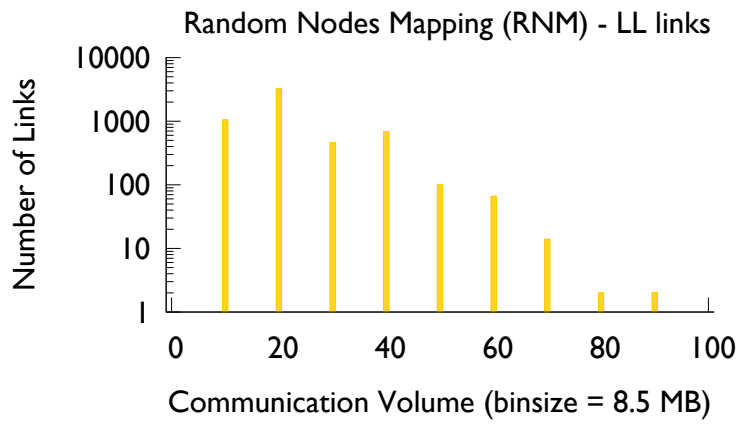  - Default MPI rank-ordered mapping (DFI)

  - Random mapping at the level of drawers (RDI)

**4D Stencil on 64 supernodes**

4D Stencil on 64 supernodes

4D Stencil on 64 supernodes

Abhinav Bhatele @ Supercomputing '11

**4D Stencil on 64 supernodes**

Default Mapping (DEF) - LL links
Number of Links / Communication Volume (binsize = 8.5 MB)

Default Mapping (DEF) - LR links
Number of Links / Communication Volume (binsize = 10.1 MB)

Default Mapping (DEF) - D links
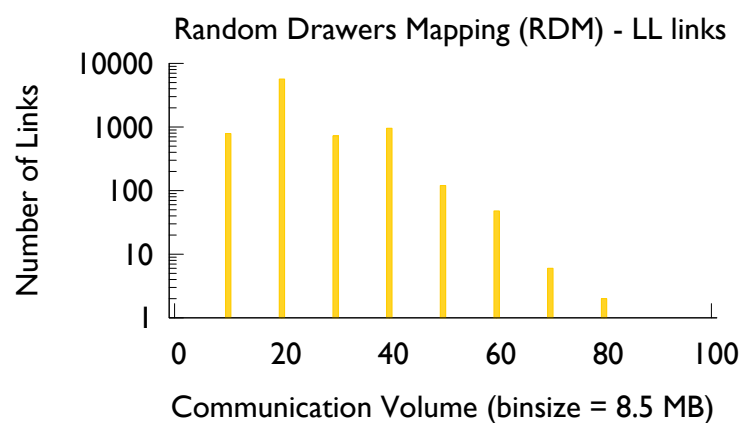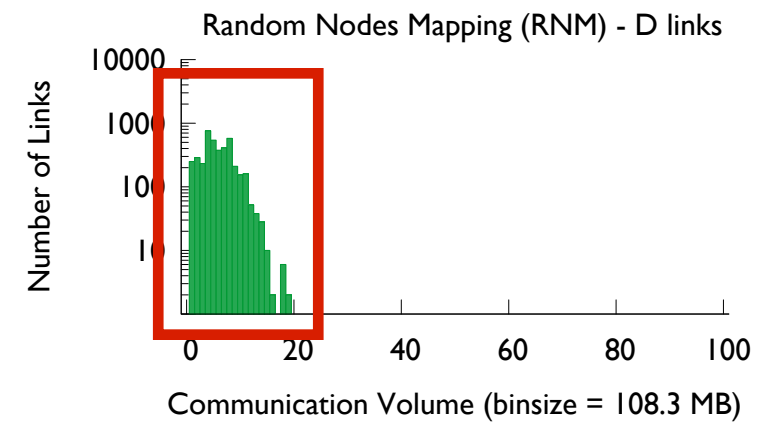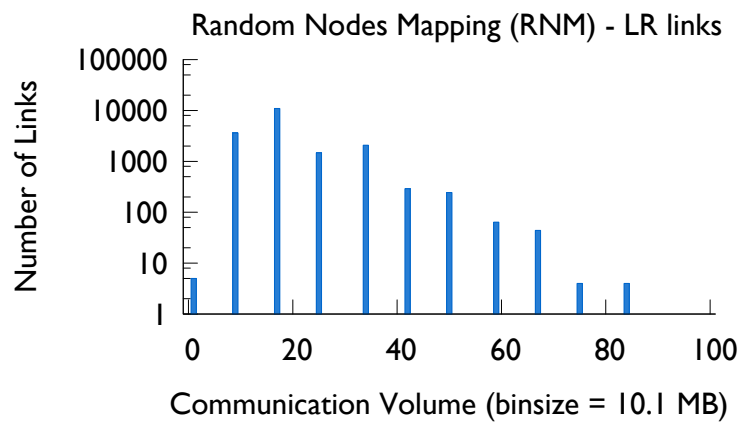Number of Links / Communication Volume (binsize = 108.3 MB)

Blocked Nodes Mapping (BNM) - LL links
Number of Links / Communication Volume (binsize = 8.5 MB)

Blocked Nodes Mapping (BNM) - LR links
Number of Links / Communication Volume (binsize = 10.1 MB)

Blocked Nodes Mapping (BNM) - D links
Number of Links / Communication Volume (binsize = 108.3 MB)

Blocked Drawers Mapping (BDM) - LL links
Number of Links / Communication Volume (binsize = 8.5 MB)

Blocked Drawers Mapping (BDM) - LR links
Number of Links / Communication Volume (binsize = 10.1 MB)

Blocked Drawers Mapping (BDM) - D links
Number of Links / Communication Volume (binsize = 108.3 MB)

Blocked Supernodes Mapping (BSM) - LL links
Number of Links / Communication Volume (binsize = 8.5 MB)

Blocked Supernodes Mapping (BSM) - LR links
Number of Links / Communication Volume (binsize = 10.1 MB)
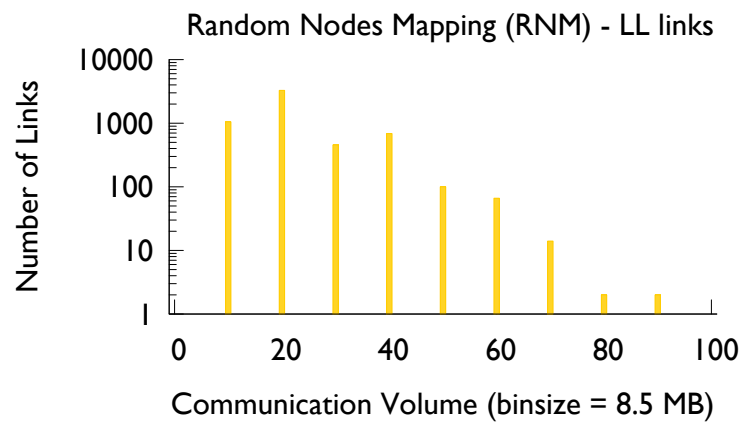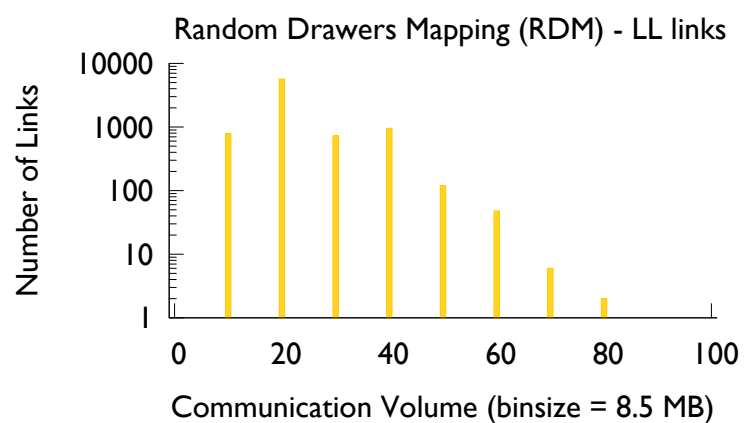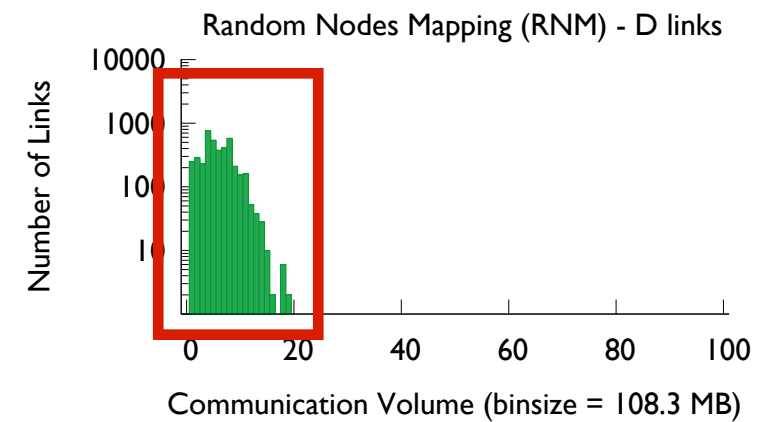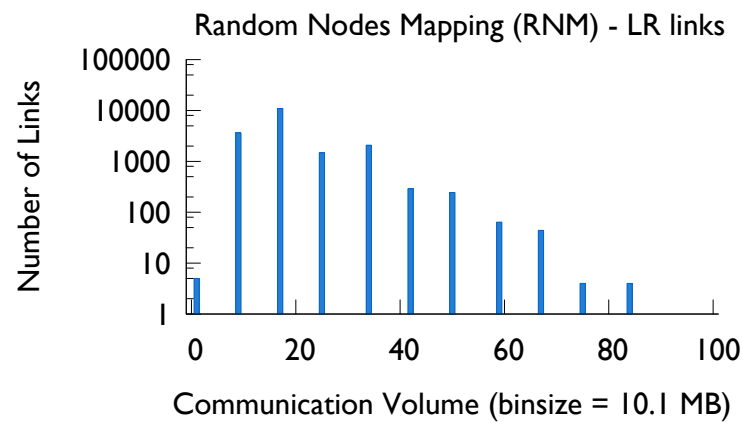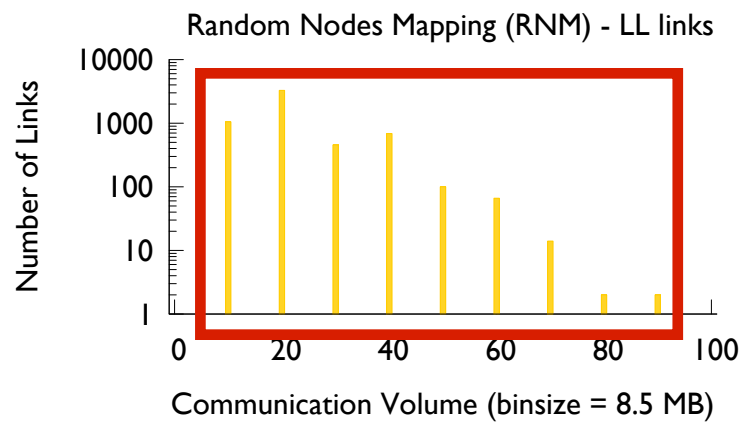
Blocked Supernodes Mapping (BSM) - D links
Number of Links / Communication Volume (binsize = 108.3 MB)

4D Stencil on 64 supernodes

Random Nodes Mapping (RNM) - LL links
Random Nodes Mapping (RNM) - LR links
Random Nodes Mapping (RNM) - D links

Random Drawers Mapping (RDM) - LL links
Random Drawers Mapping (RDM) - LR links
Random Drawers Mapping (RDM) - D links

Default with Indirect (DFI) - LL links
Default with Indirect (DFI) - LR links
Default with Indirect (DFI) - D links

Random Drawers with Indirect (RDI) - LL links
Random Drawers with Indirect (RDI) - LR links
Random Drawers with Indirect (RDI) - D links

4D Stencil on 64 supernodes

Random Nodes Mapping (RNM) - LL links
Random Nodes Mapping (RNM) - LR links
Random Nodes Mapping (RNM) - D links

Random Drawers Mapping (RDM) - LL links
Random Drawers Mapping (RDM) - LR links
Random Drawers Mapping (RDM) - D links

Default with Indirect (DFI) - LL links
Default with Indirect (DFI) - LR links
Default with Indirect (DFI) - D links

Random Drawers with Indirect (RDI) - LL links
Random Drawers with Indirect (RDI) - LR links
Random Drawers with Indirect (RDI) - D links

4D Stencil on 64 supernodes

Random Nodes Mapping (RNM) - LL links
Random Nodes Mapping (RNM) - LR links
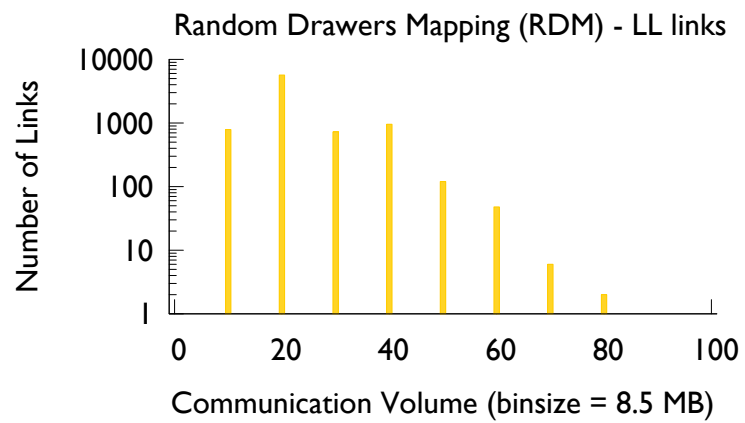Random Nodes Mapping (RNM) - D links

Random Drawers Mapping (RDM) - LL links
Random Drawers Mapping (RDM) - LR links
Random Drawers Mapping (RDM) - D links

Default with Indirect (DFI) - LL links
Default with Indirect (DFI) - LR links
Default with Indirect (DFI) - D links

Random Drawers with Indirect (RDI) - LL links
Random Drawers with Indirect (RDI) - LR links
Random Drawers with Indirect (RDI) - D links

4D Stencil on 64 supernodes

LLNL-PRES-511461     Abhinav Bhatele @ Supercomputing '11     13

4D Stencil on 64 supernodes

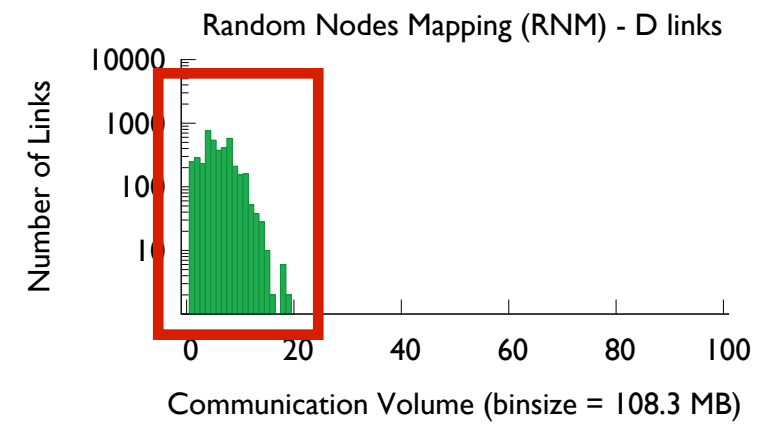LLNL-PRES-511461                    Abhinav Bhatele @ Supercomputing '11                    13
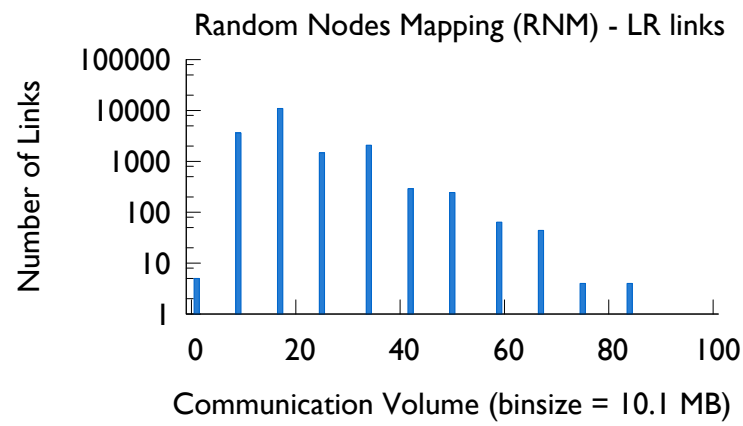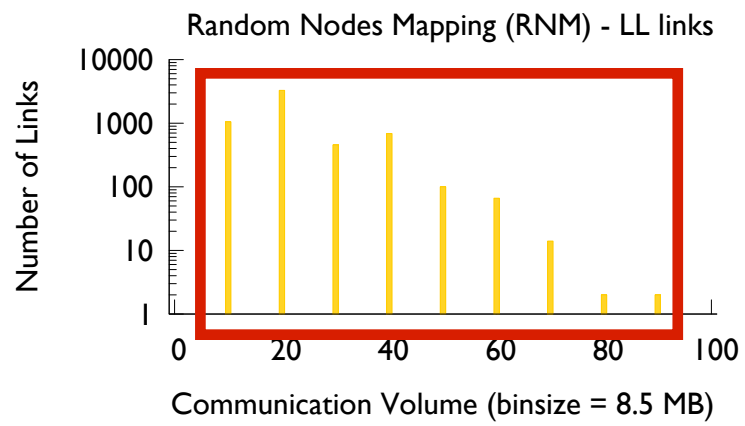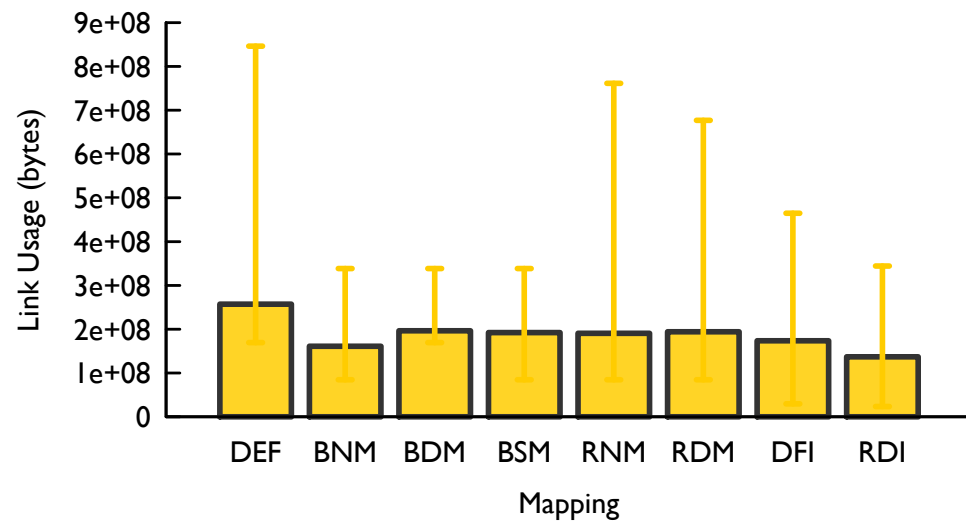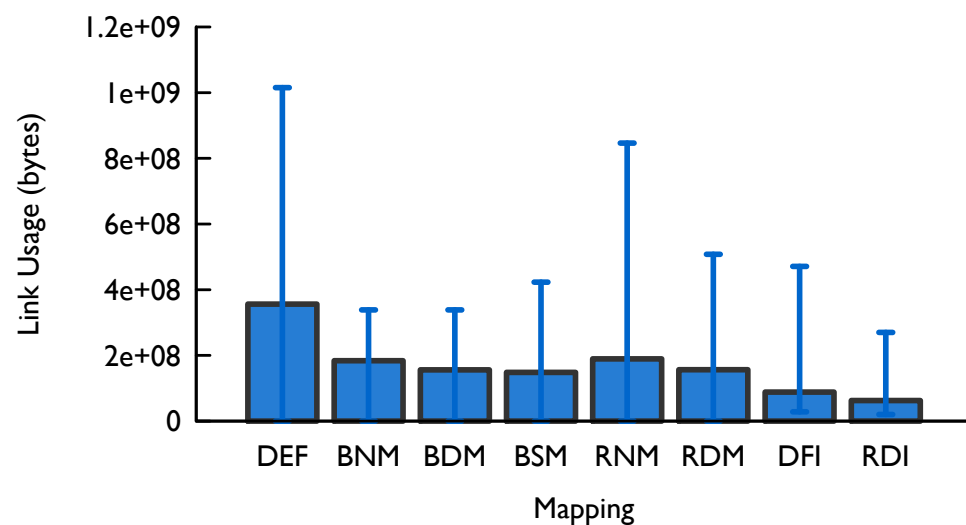
4D Stencil on 64 supernodes
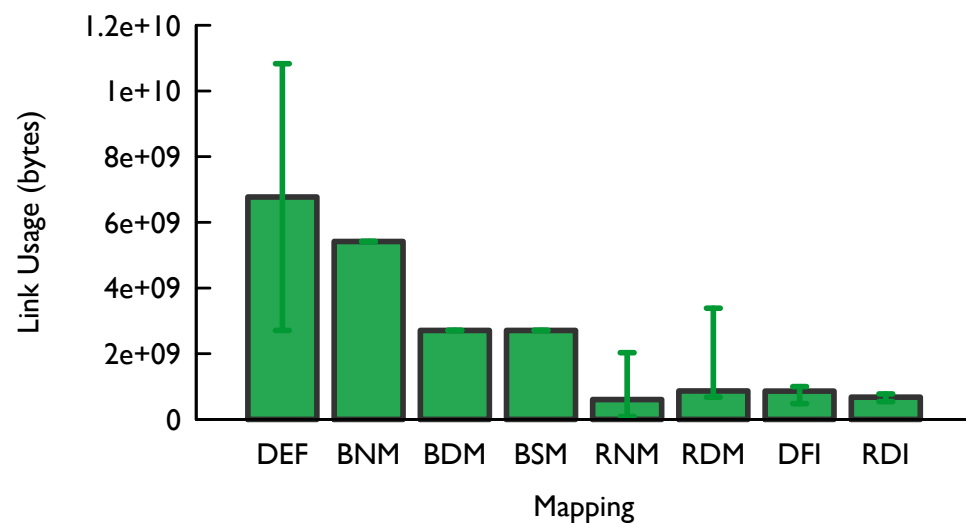
4D Stencil on 64 supernodes

4D Stencil on 64 supernodes (LL links)

4D Stencil on 64 supernodes (LR links)
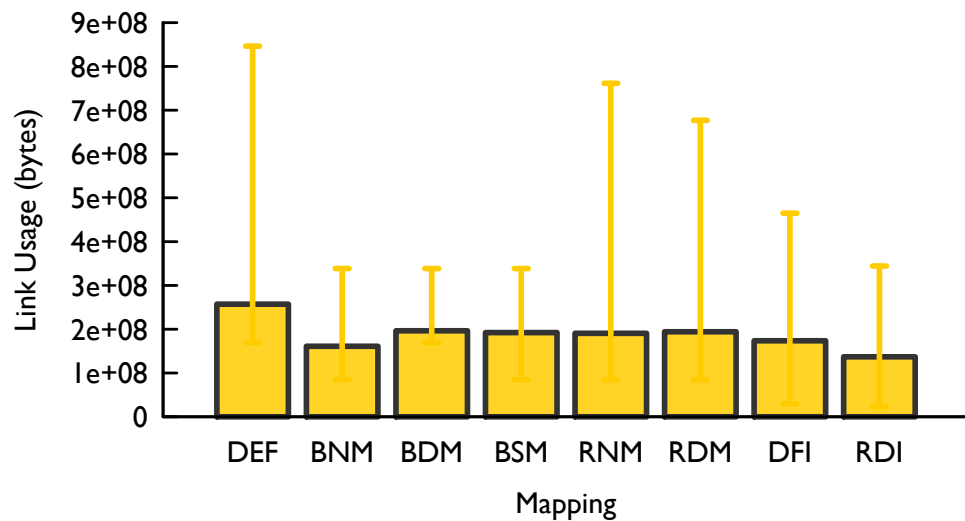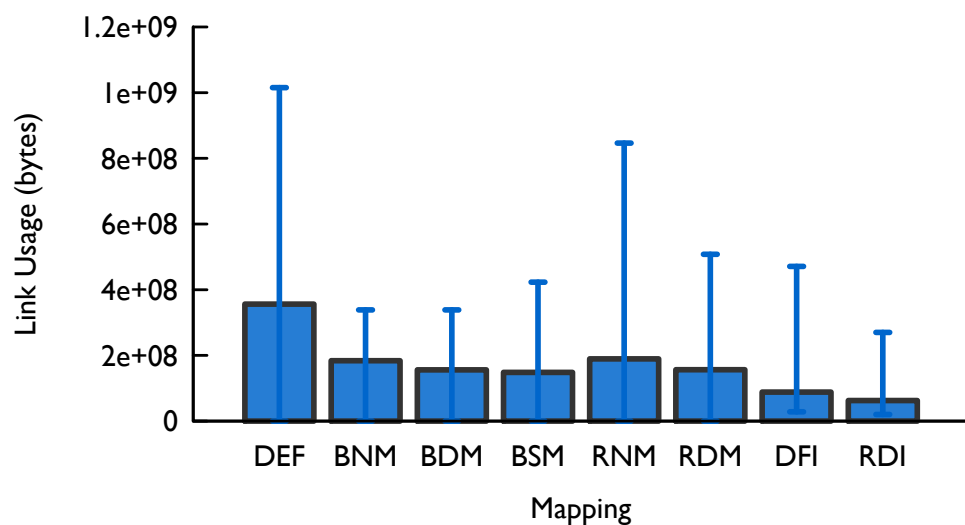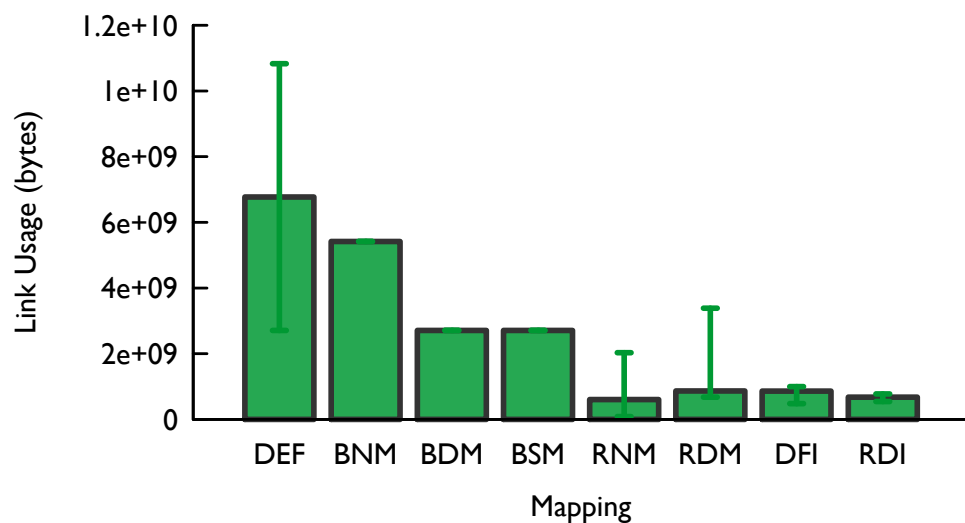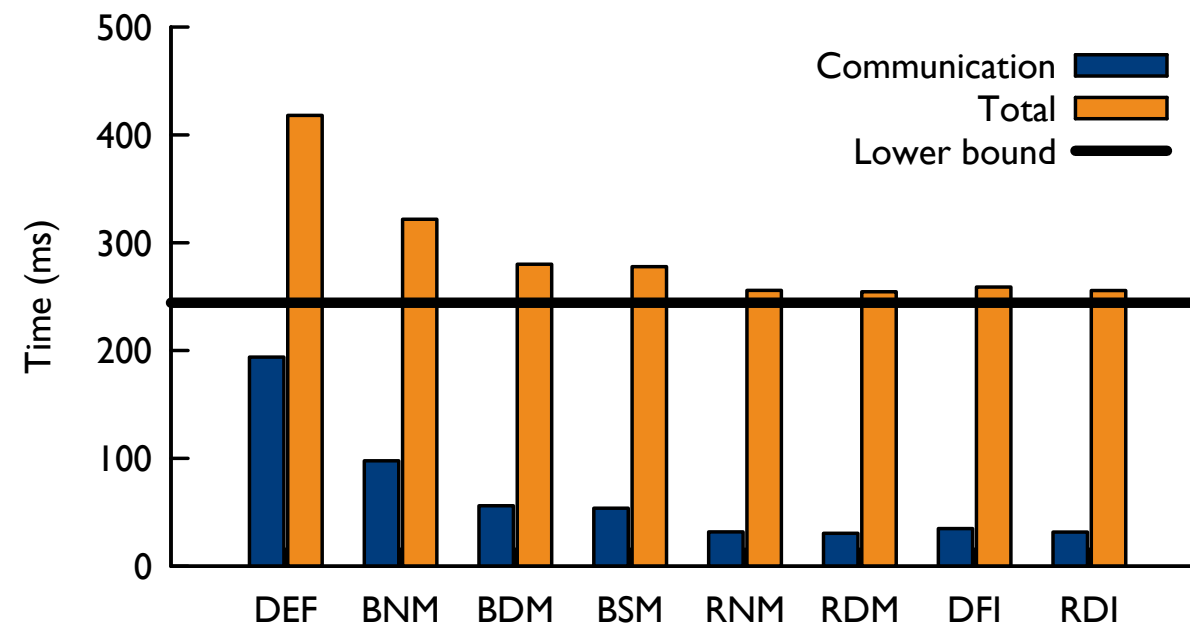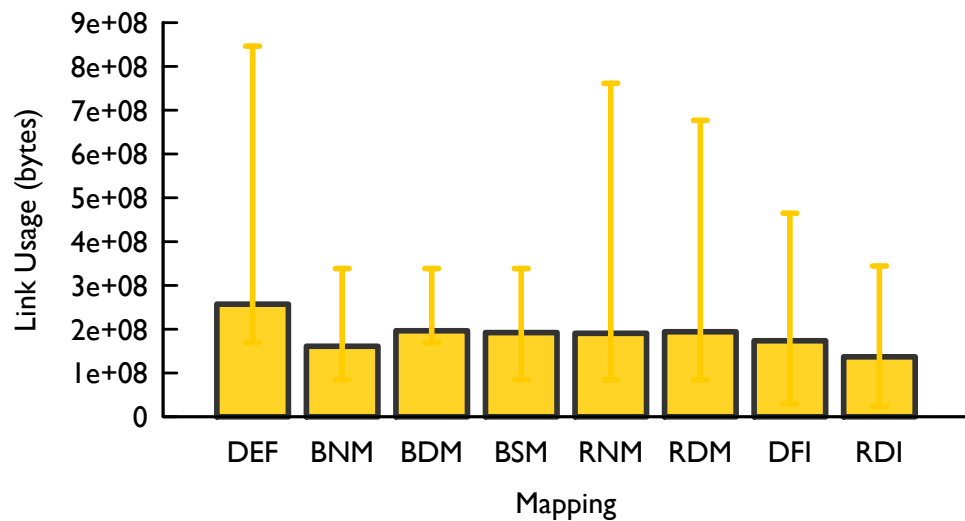
4D Stencil on 64 supernodes (D links)

Performance of 4D Stencil on 64 supernodes

4D Stencil on 64 supernodes

4D Stencil on 64 supernodes (LL links)

4D Stencil on 64 supernodes (LR links)

4D Stencil on 64 supernodes (D links)

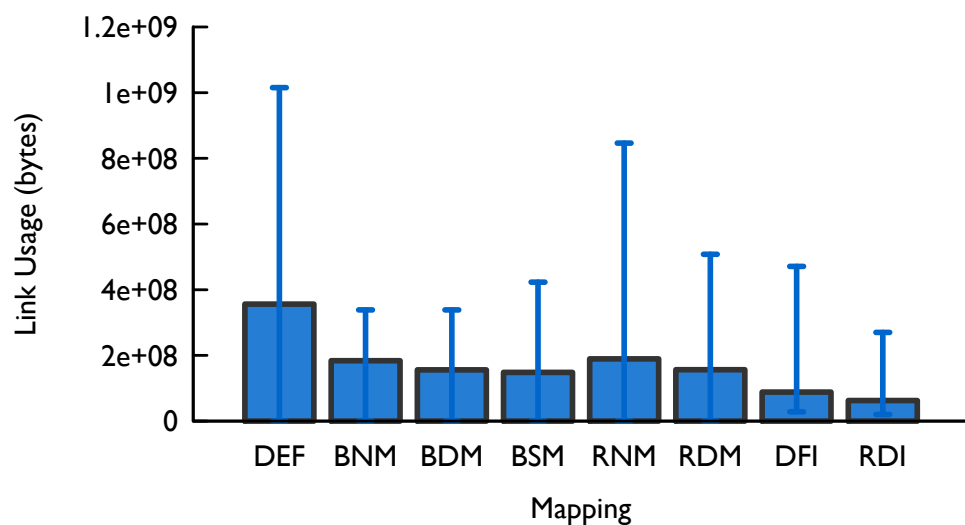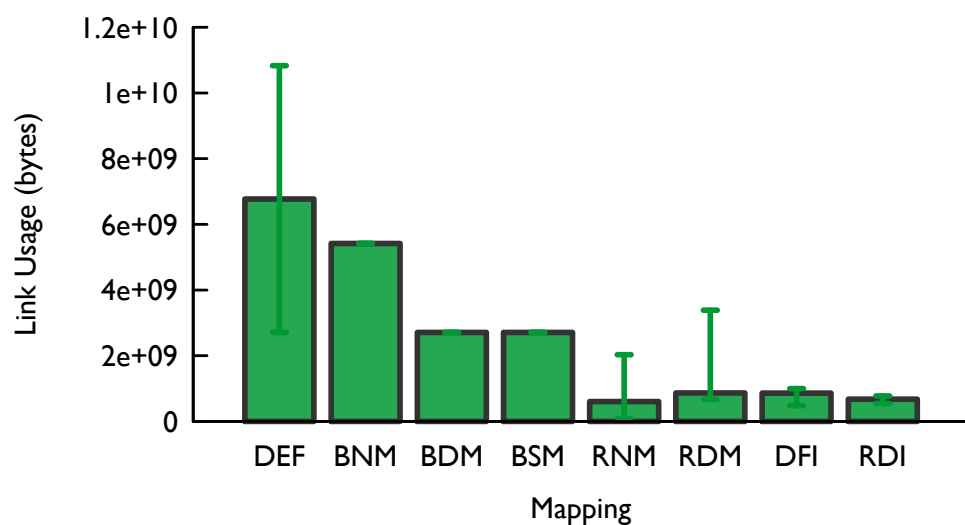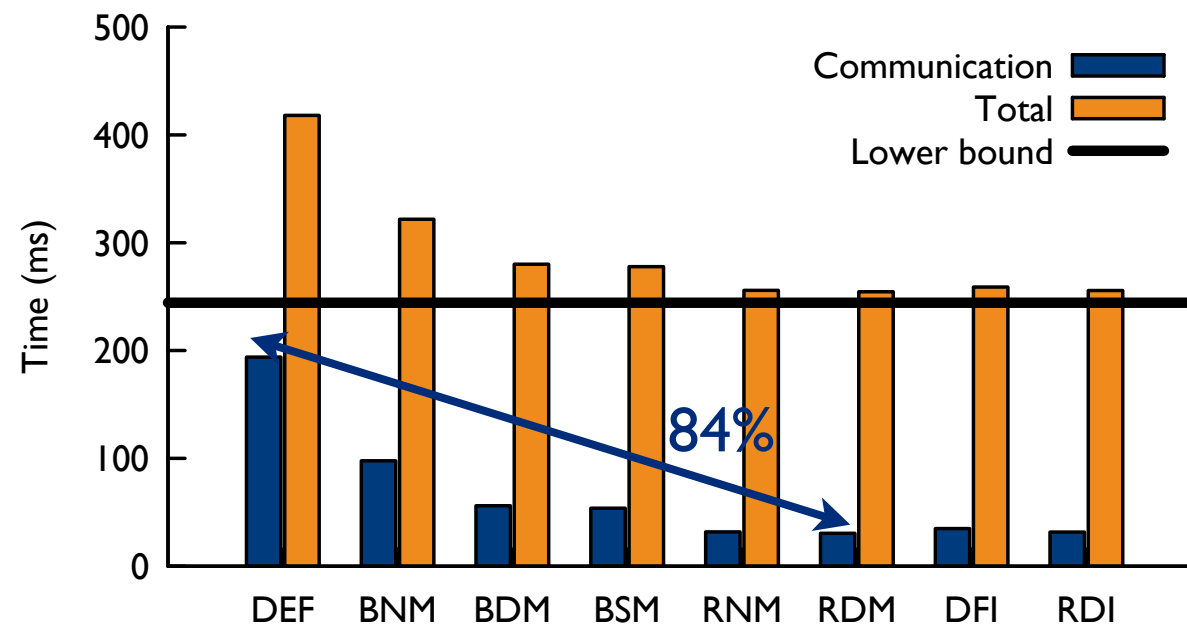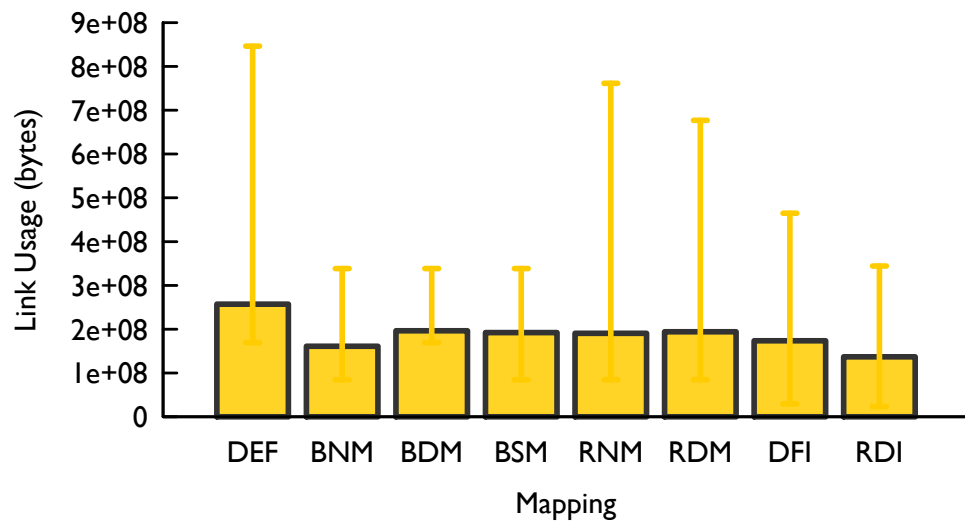Performance of 4D Stencil on 64 supernodes

Communication
Total
Lower bound

84%

# 4D Stencil on 64 supernodes

### 4D Stencil on 64 supernodes (LL links)



### 4D Stencil on 64 supernodes (LR links)



### 4D Stencil on 64 supernodes (D links)



### Performance of 4D Stencil on 64 supernodes



Communication
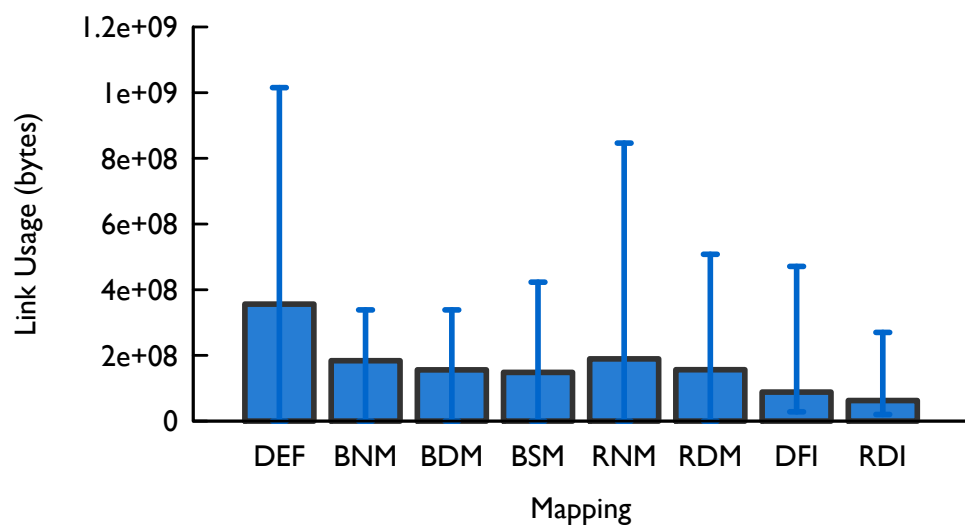Total
Lower bound

39%

84%

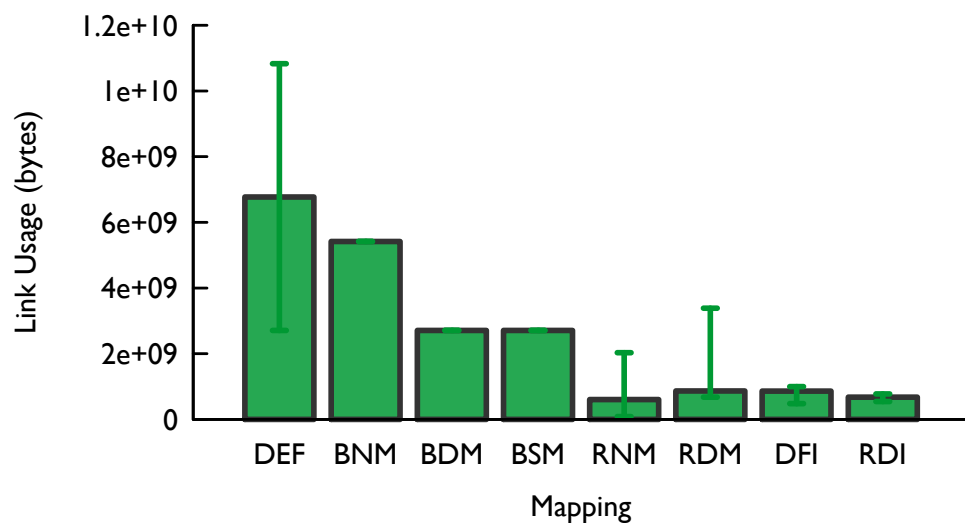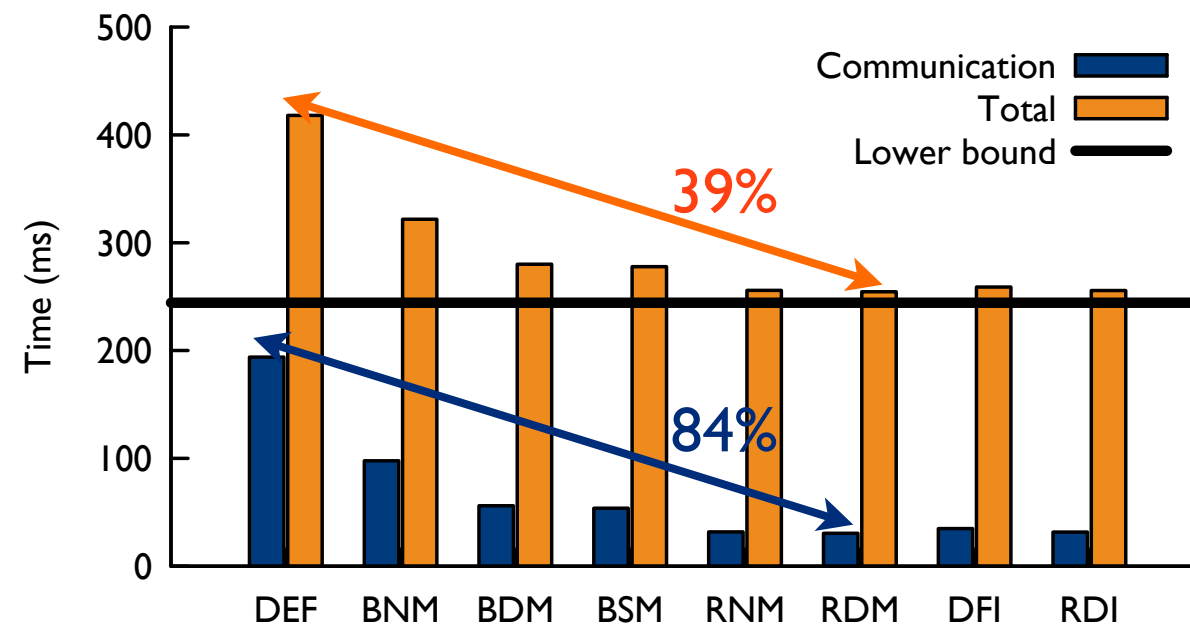# 4D Stencil on 64 supernodes



4D Stencil on 64 supernodes (LL links)



4D Stencil on 64 supernodes (LR links)
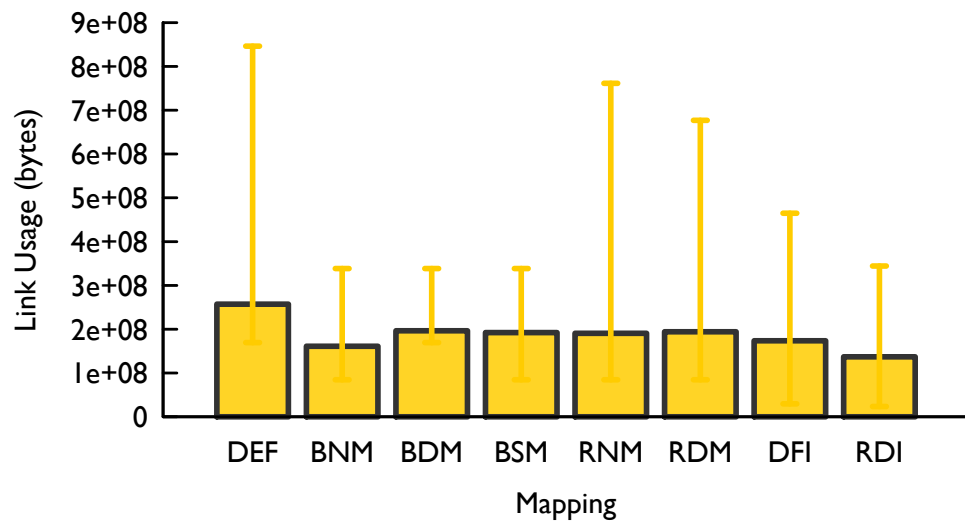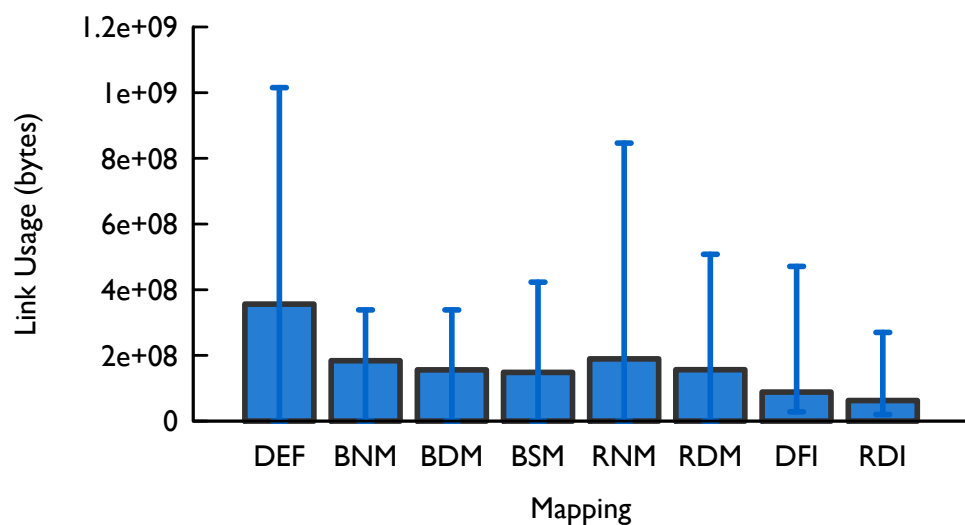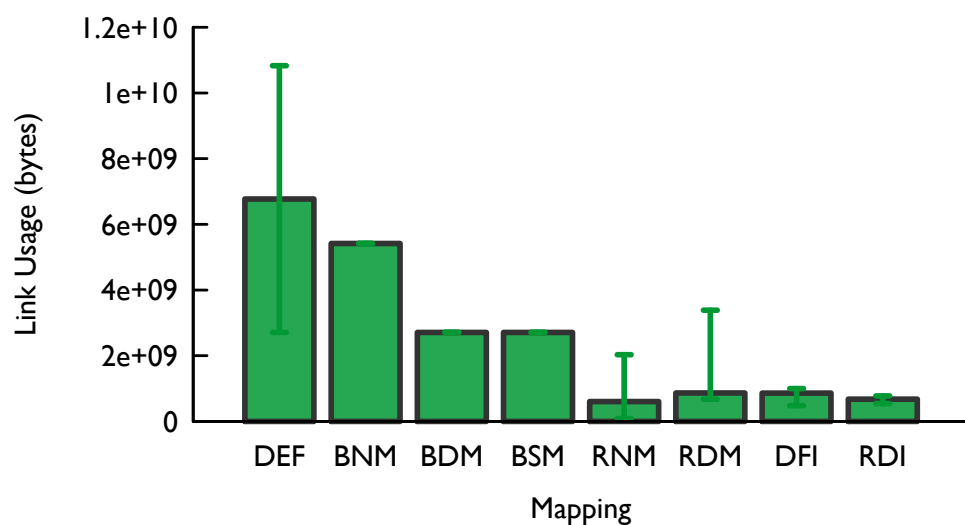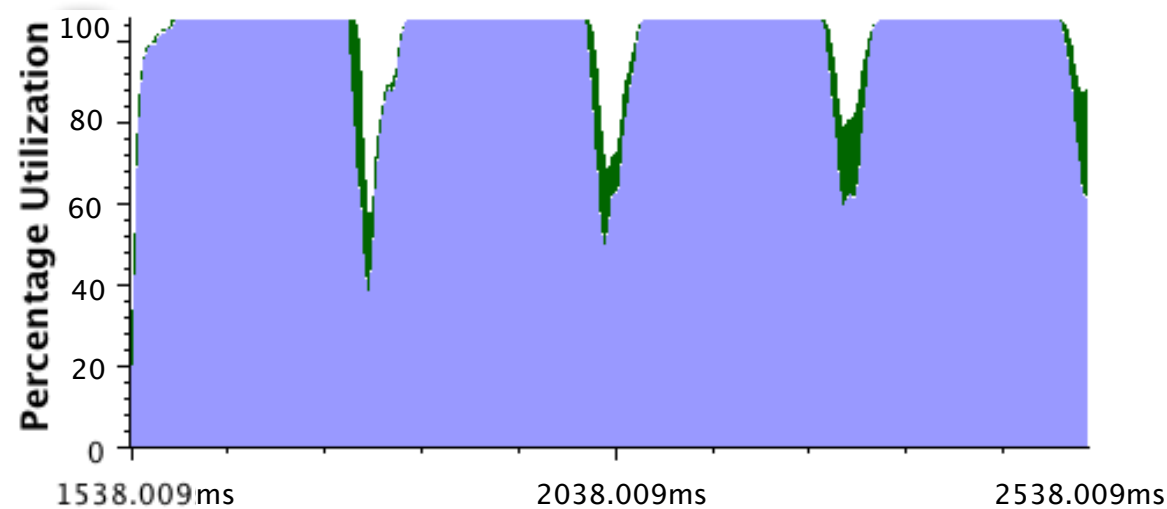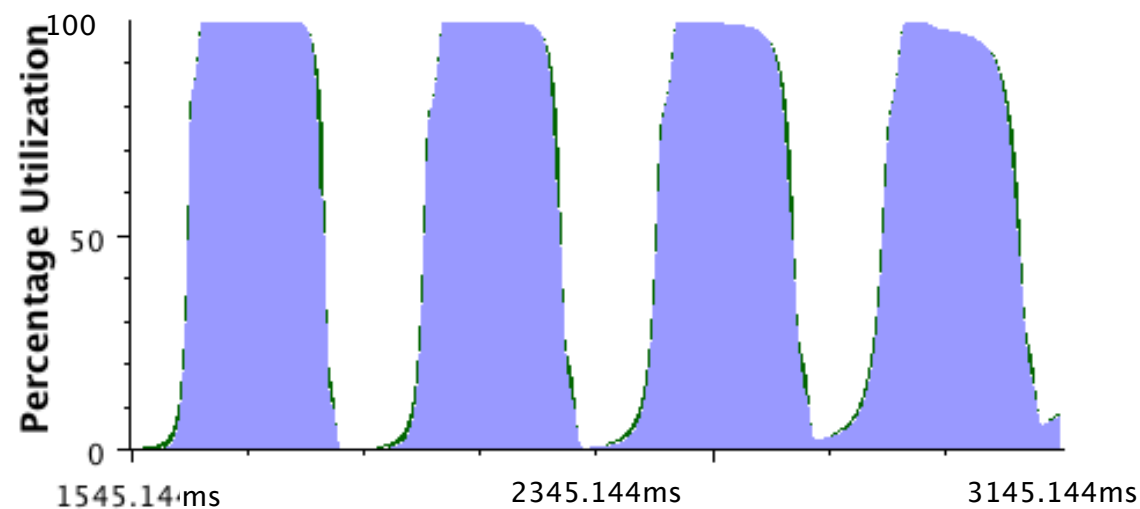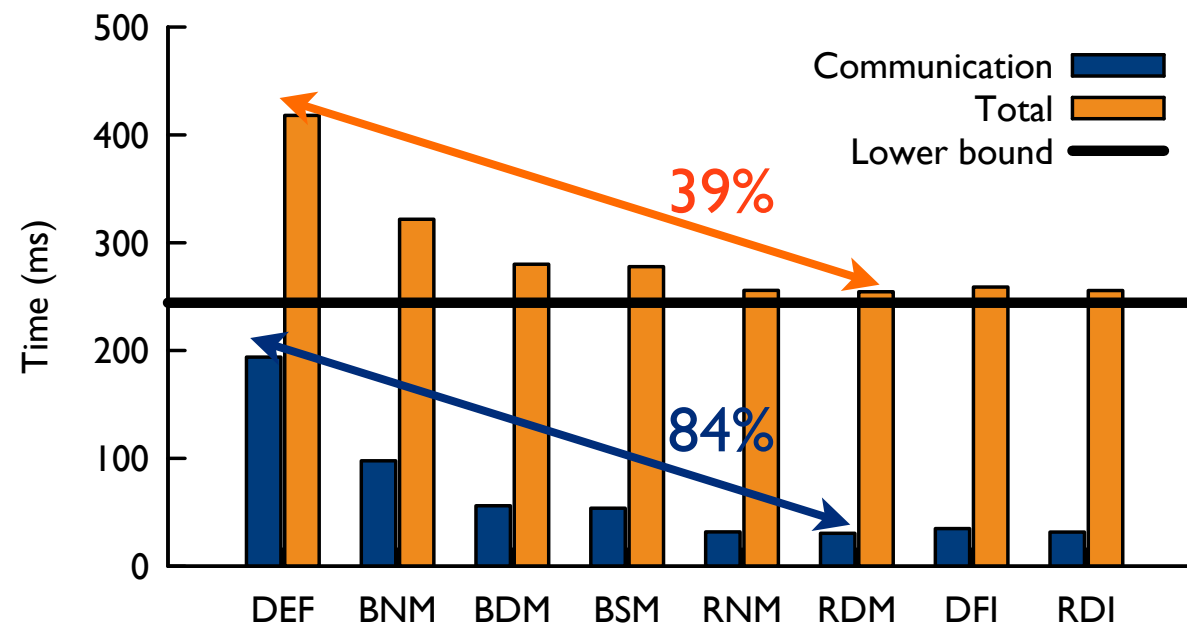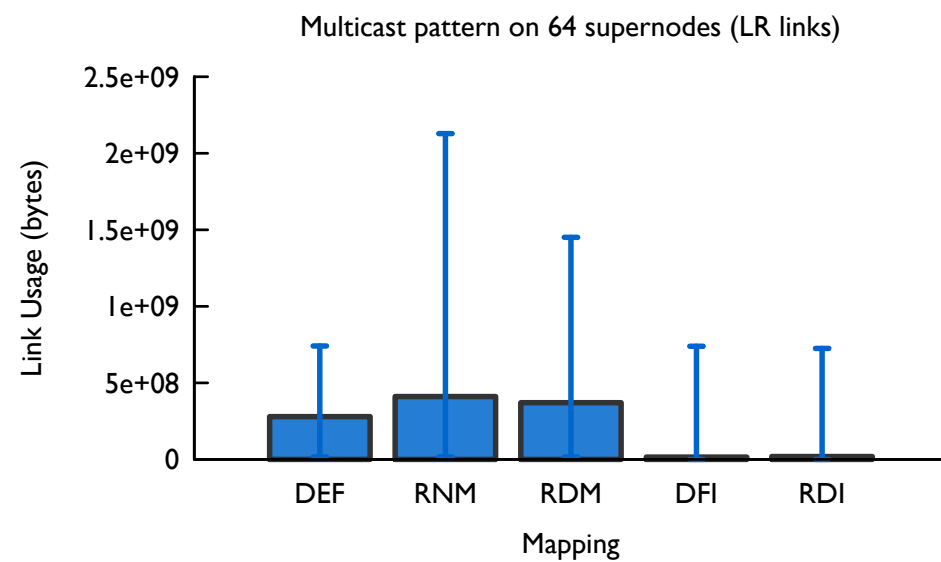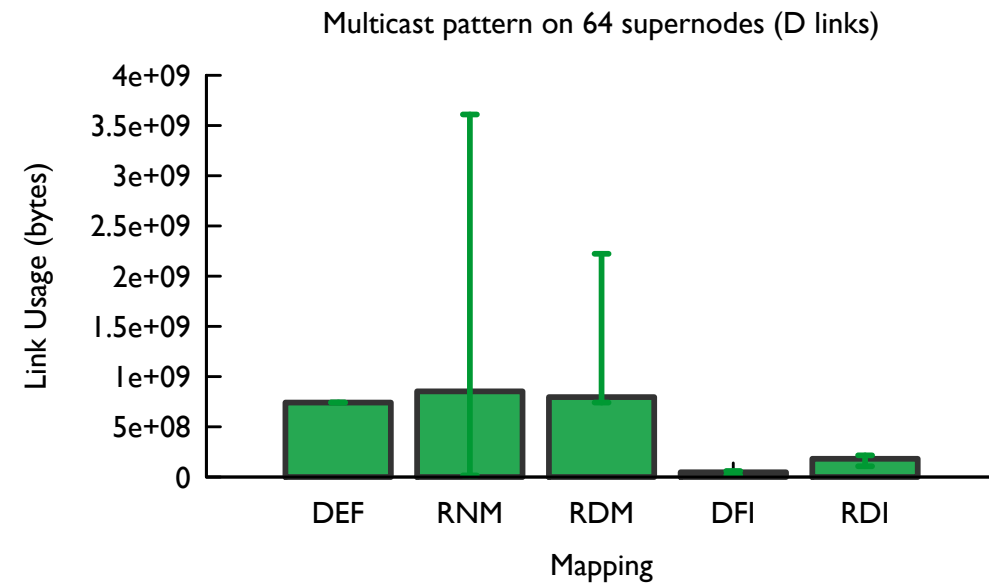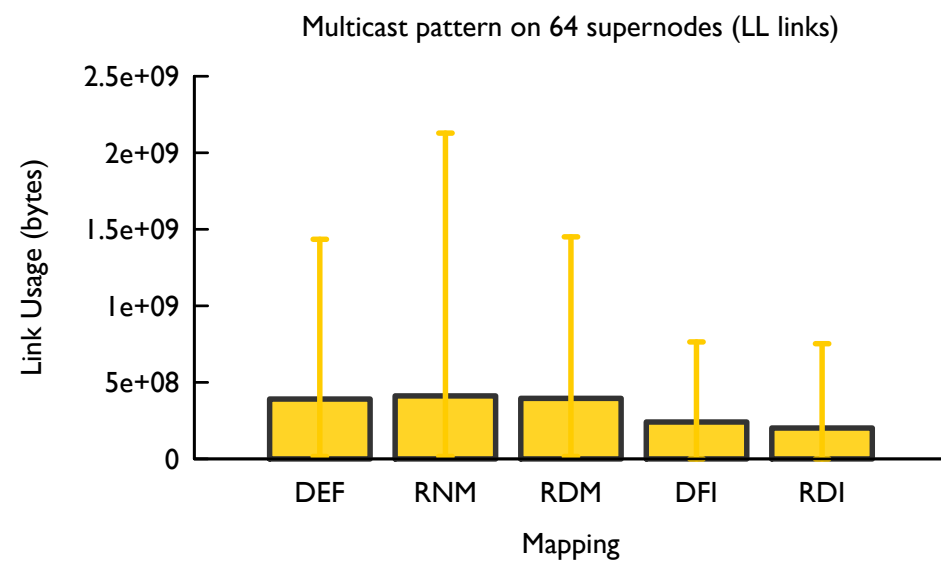


4D Stencil on 64 supernodes (D links)

Performance of 4D Stencil on 64 supernodes
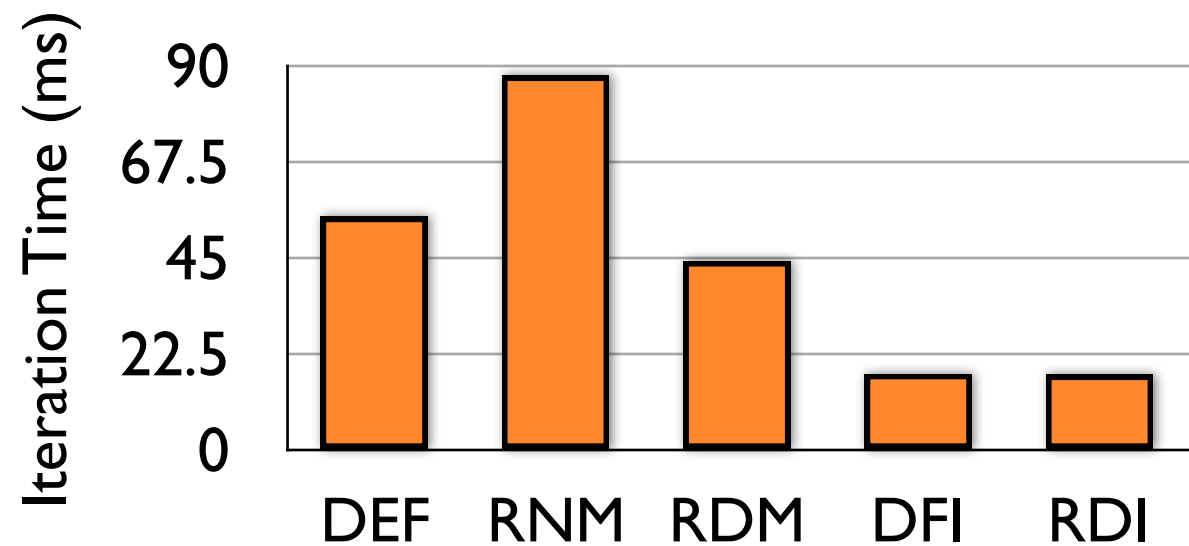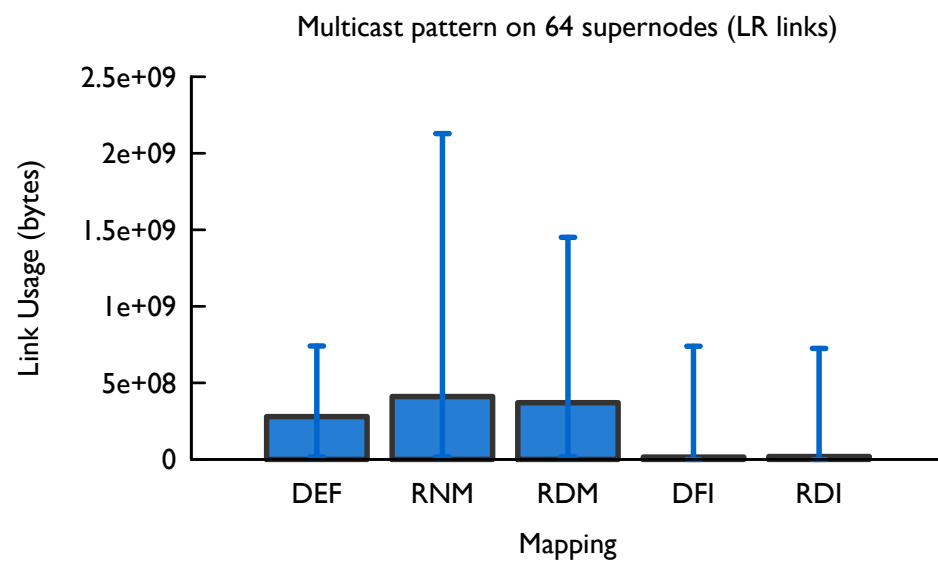
# n-targets multicast pattern

- Similar to the communication pattern in NAMD

- Each MPI task sends messages to 14 others

- Message size = 1 MB

# n-targets multicast pattern

# n-targets multicast pattern



Multicast pattern on 64 supernodes (LL links)

Multicast pattern on 64 supernodes (D links)

Multicast pattern on 64 supernodes (LR links)

# n-targets multicast pattern



Multicast pattern on 64 supernodes (LL links)

Multicast pattern on 64 supernodes (D links)

Multicast pattern on 64 supernodes (LR links)

68%

4D Stencil on 300 supernodes (LL links)

Link Usage (bytes)

Mapping: DEF BNM BDM BSM RNM RDM DFI RDI

4D Stencil on 300 supernodes (LR links)

Link Usage (bytes)

Mapping: DEF BNM BDM BSM RNM RDM DFI RDI

4D Stencil on 300 supernodes (D links)

Link Usage (bytes)

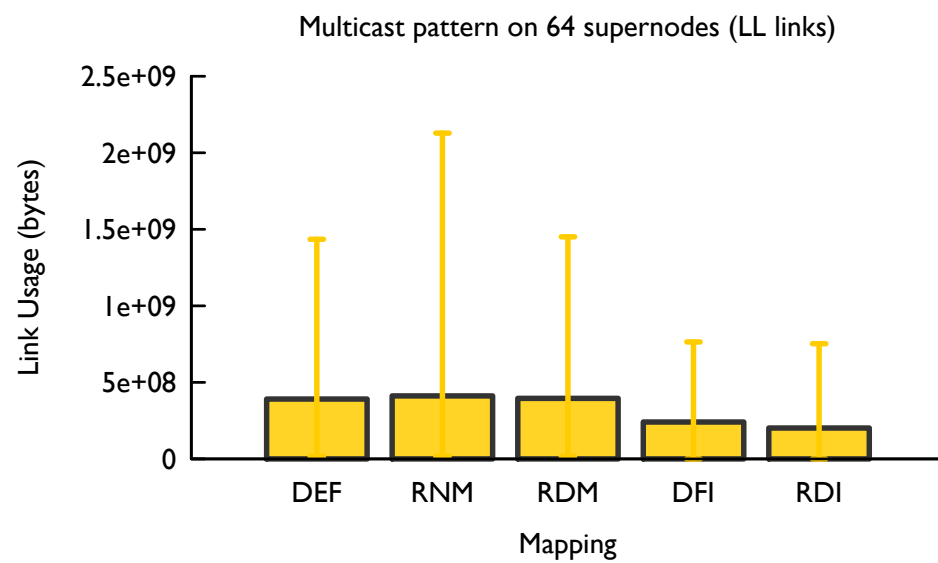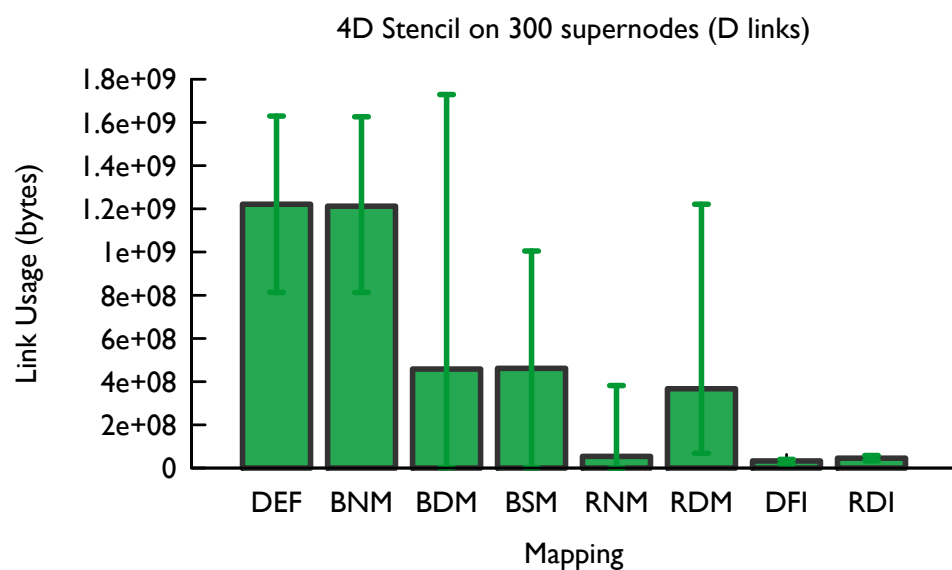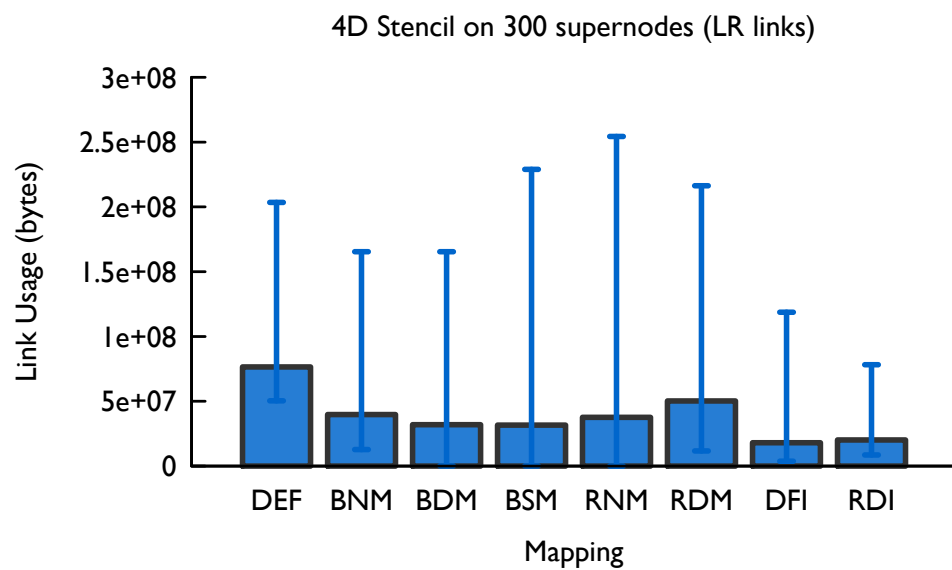Mapping: DEF BNM BDM BSM RNM RDM DFI RDI

- Largest detailed network simulation so far = 307,200 MPI tasks

- Non-power-of-2 leads to more complex mapping

- Message size = 1 MB

Text

4D Stencil on 300 supernodes

- Largest detailed network simulation so far = 307,200 MPI tasks

- Non-power-of-2 leads to more complex mapping

- Message size = 1 MB

Text
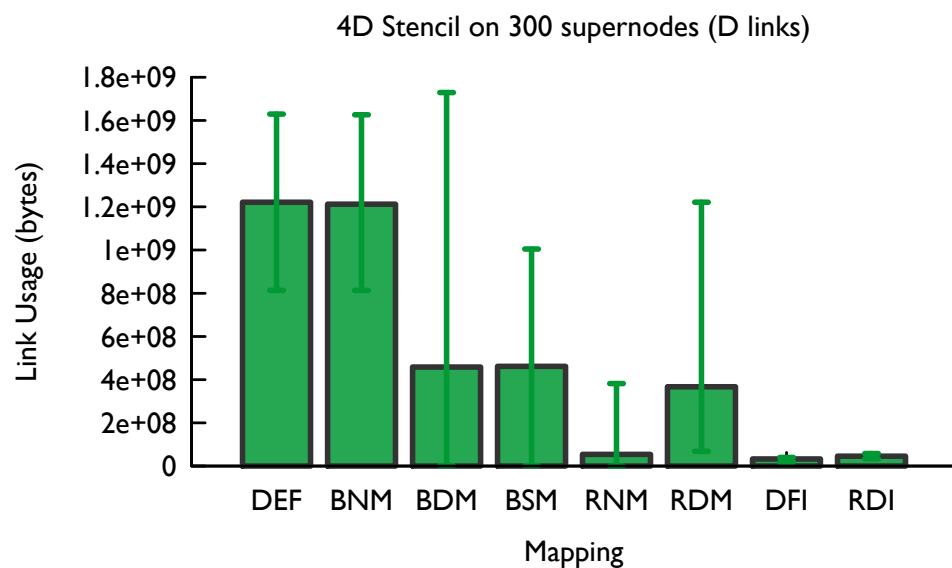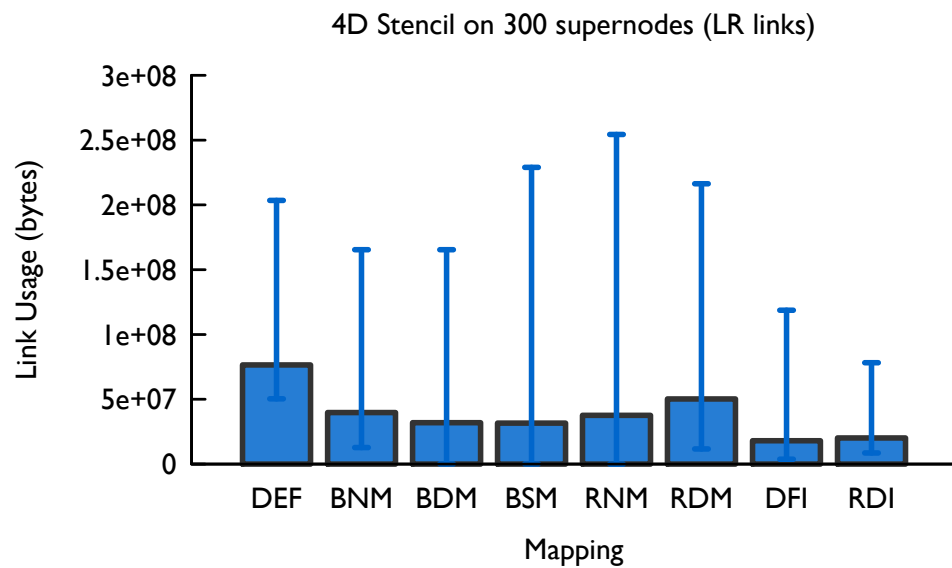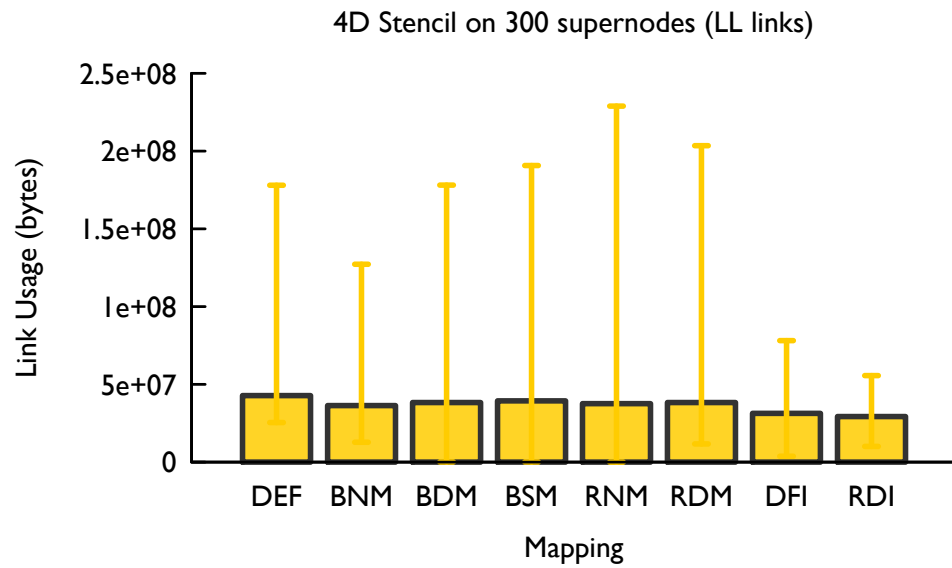
4D Stencil on 300 supernodes

- Largest detailed network simulation so far = 307,200 MPI tasks

- Non-power-of-2 leads to more complex mapping

- Message size = 1 MB

Text

Abhinav Bhatele @ Supercomputing '11

4D Stencil on 300 supernodes

- Largest detailed network simulation so far = 307,200 MPI tasks

- Non-power-of-2 leads to more complex mapping
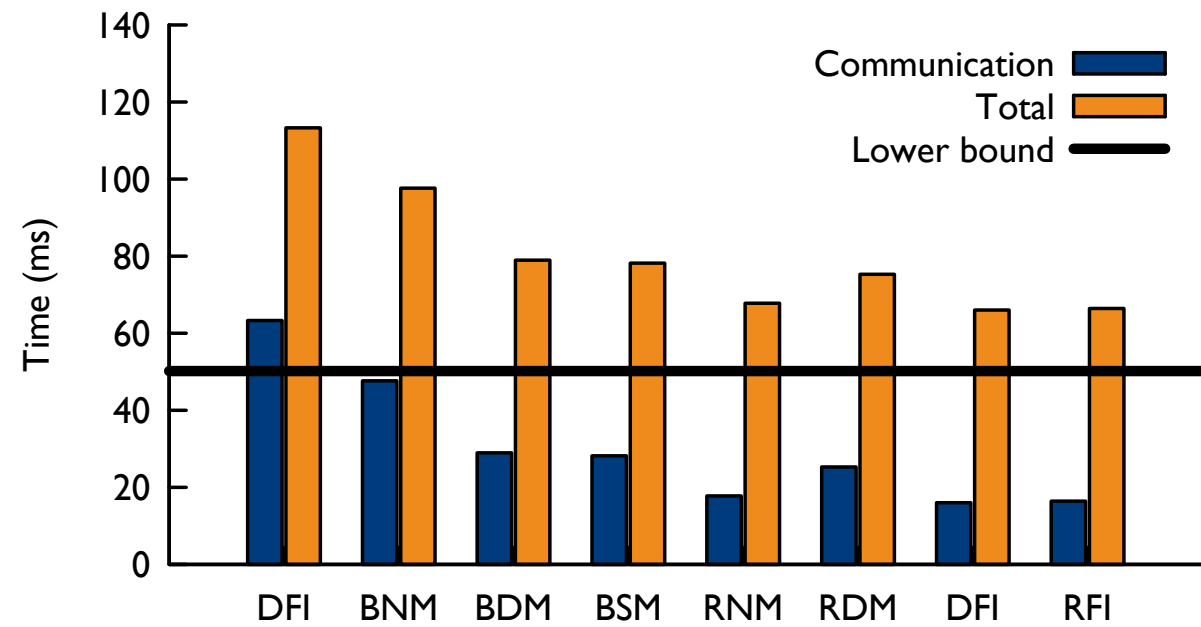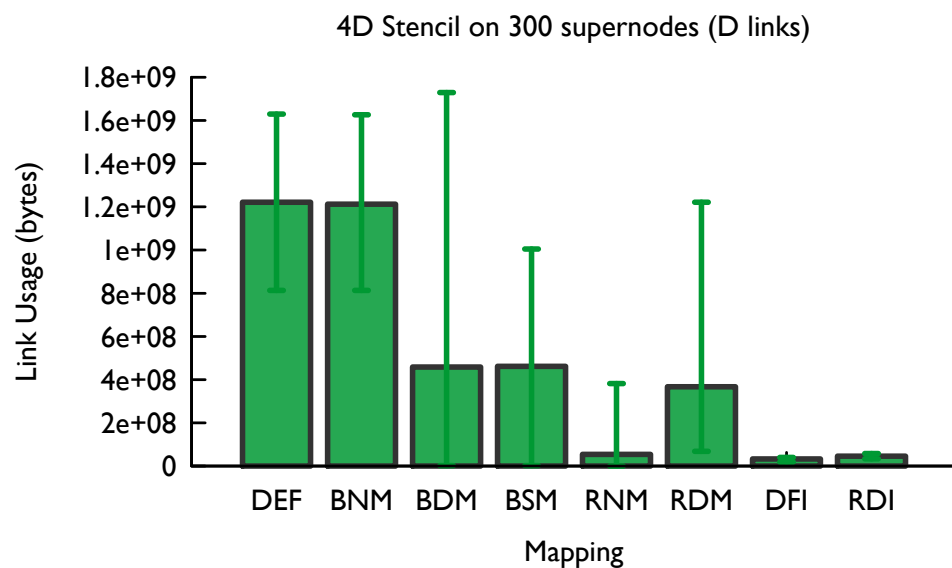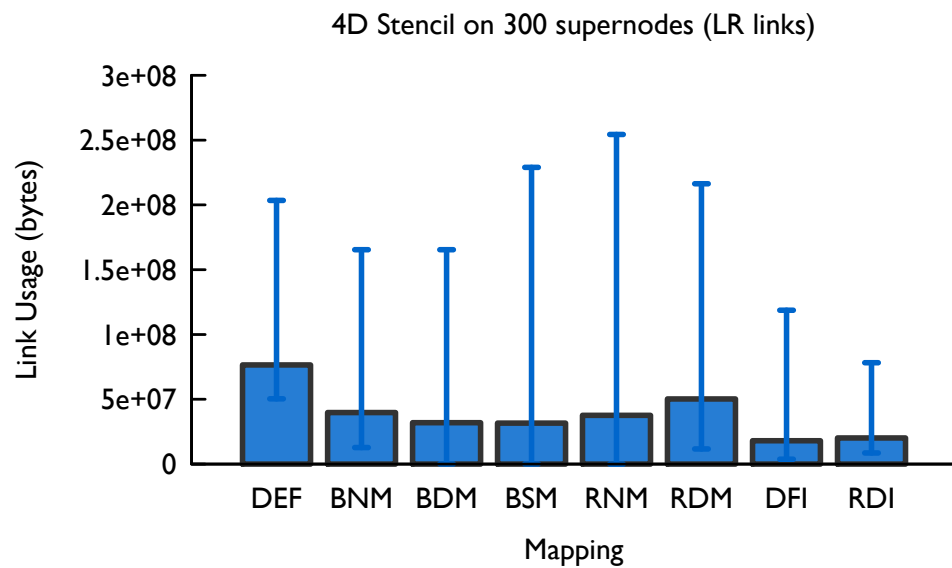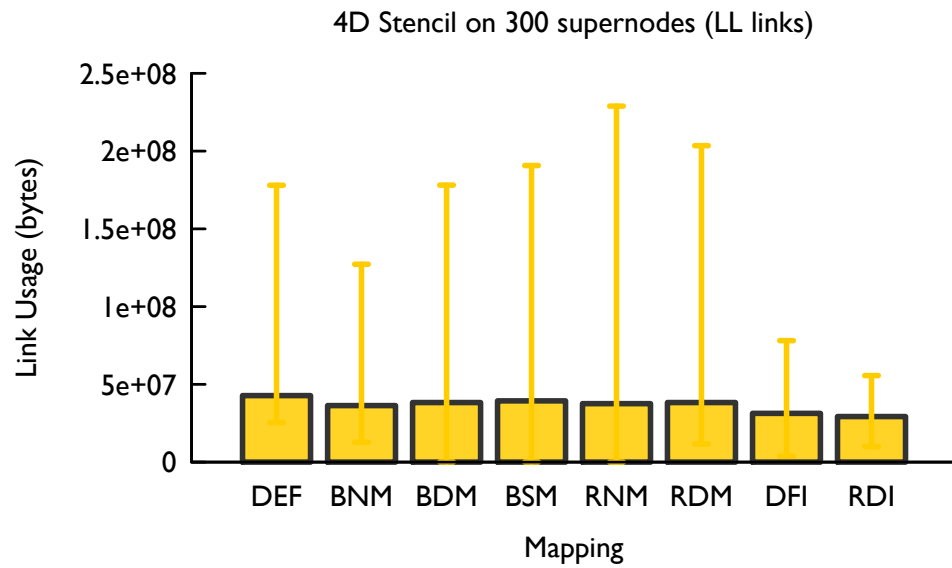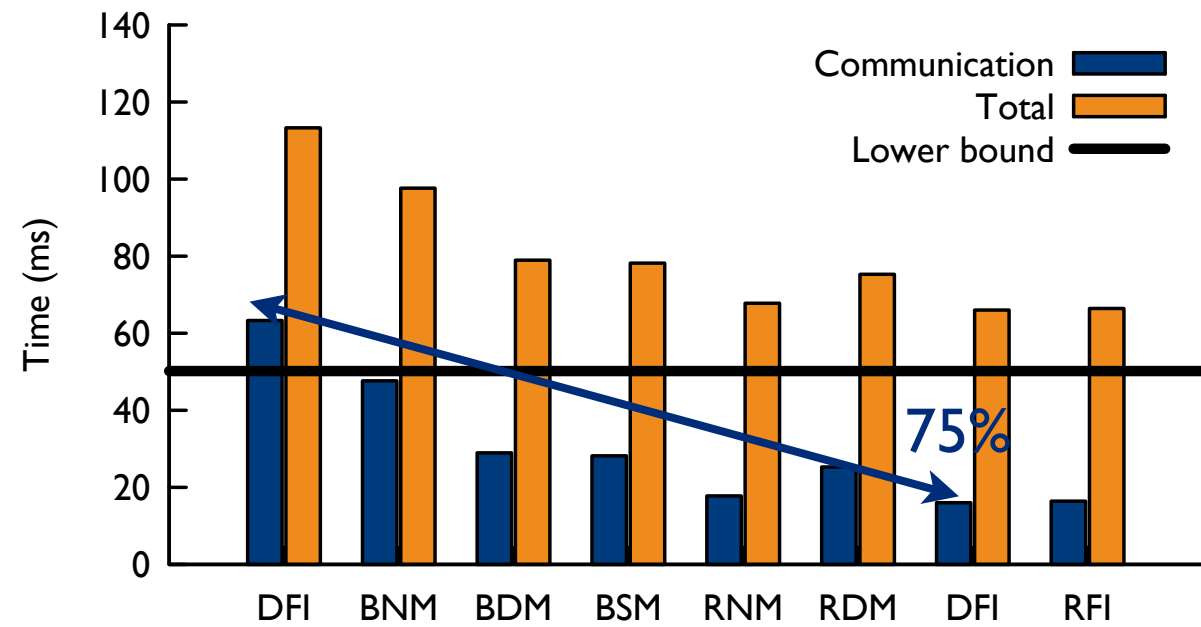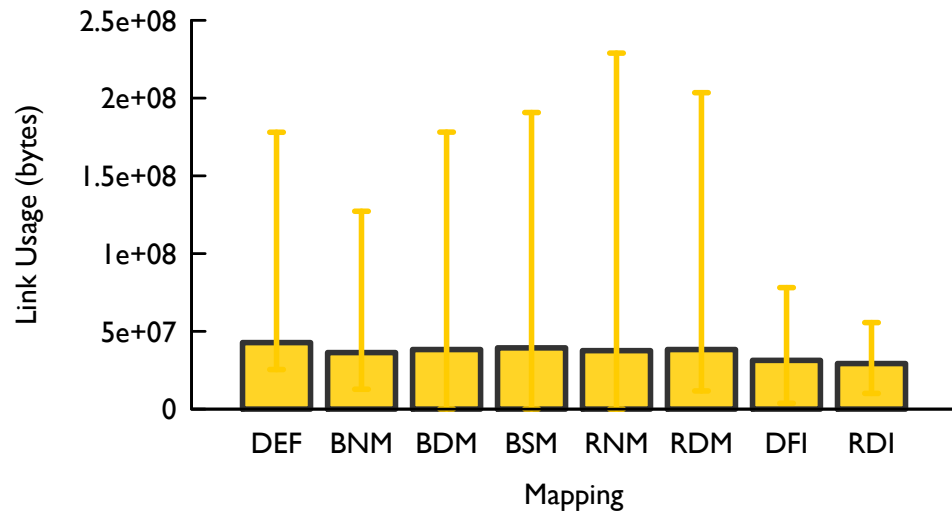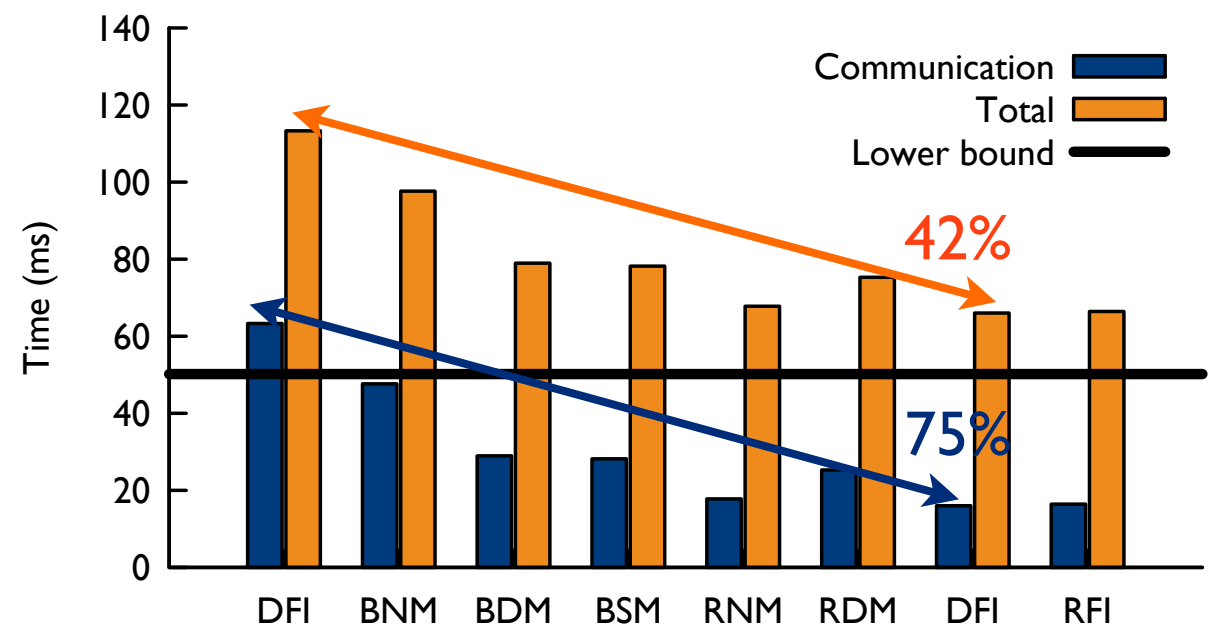
- Message size = 1 MB

Text

# Summary

- Default MPI rank-ordered mapping on multi-level direct networks can lead to hot-spots

- Packet-level simulation to assist machine architects and application developers in making routing and mapping choices

- Conclusions:

  - With direct routing, random mapping at node granularity is best

  - With indirect routing, default mapping is good enough

- Utility of simulation-based analysis to analyze algorithms and design choices for future machines

# Questions?

More information at: http://charm.cs.illinois.edu/research/topology

Thanks to Ryan Mokos (PPL) for implementing the PERCS network model.
Funding support: NSF grant OCI 07-25070, DOE grant DE-SC0001845

Supercomputing '11 ◆ November 17, 2011

LLNL-PRES-511461

Lawrence Livermore National Laboratory, P. O. Box 808, Livermore, CA 94551