



F U T U R E T E C H N O L O G I E S G R O U P

Throughput Oriented Runtimes for Manycore Clusters

Costin Iancu, Khaled Ibrahim

Tsukuba – March 19



BERKELEY LAB

F U T U R E T E C H N O L O G I E S G R O U P

We know how to build runtimes that scale
with the node count (Latency)

BUT

We have little experience building runtimes
that scale with cores per node (and nodes)
(Throughput)



32 cores



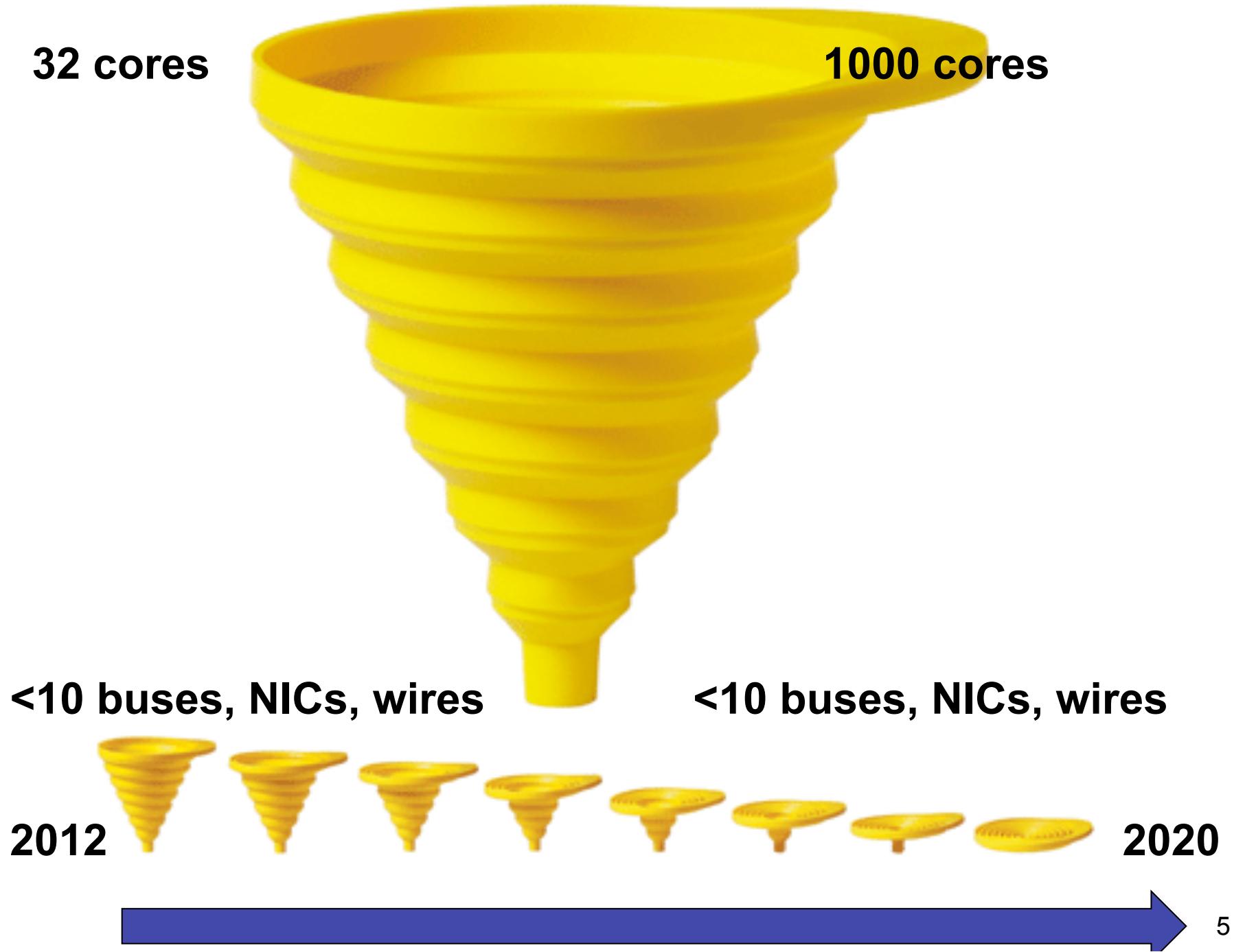
<10 buses, NICs, wires



2012

2020







Outline

F U T U R E T E C H N O L O G I E S G R O U P

- ❖ Network performance in manycore clusters
- ❖ Designing a throughput oriented runtime
- ❖ Future work



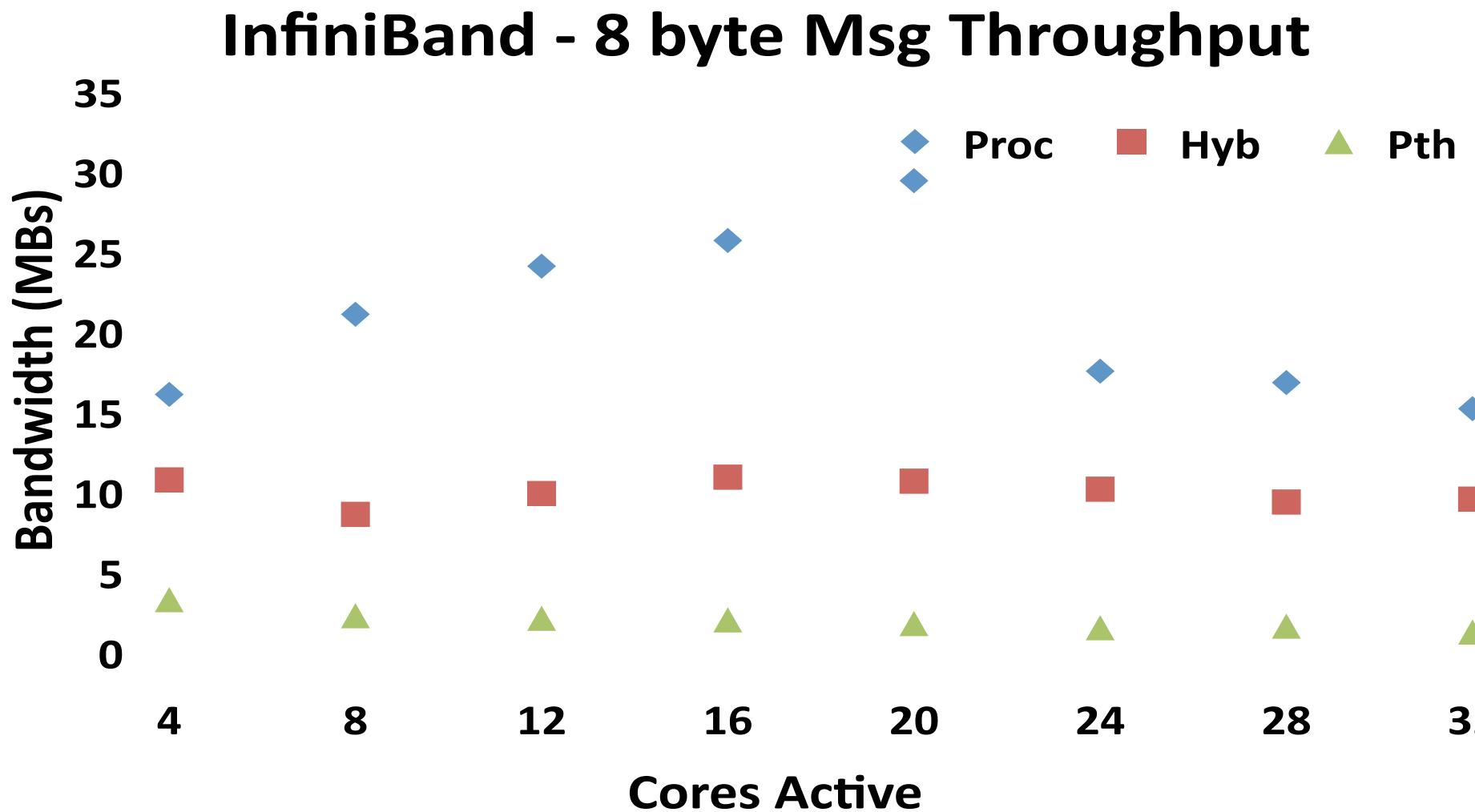
F U T U R E T E C H N O L O G I E S G R O U P

Network Performance



InfiniBand Performance

F U T U R E T E C H N O L O G I E S G R O U P

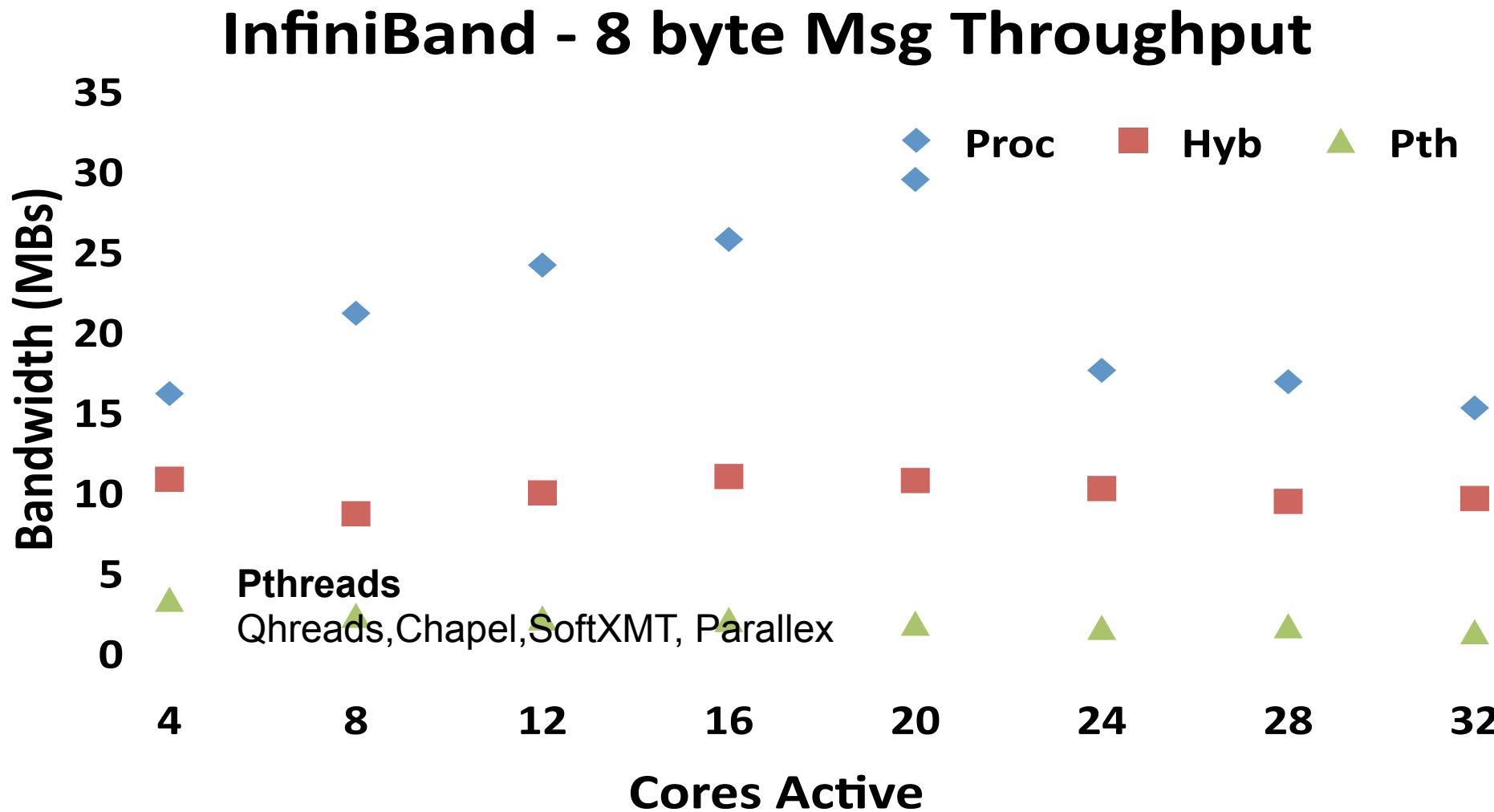


LAWRENCE BERKELEY NATIONAL LABORATORY



InfiniBand Performance

F U T U R E T E C H N O L O G I E S G R O U P

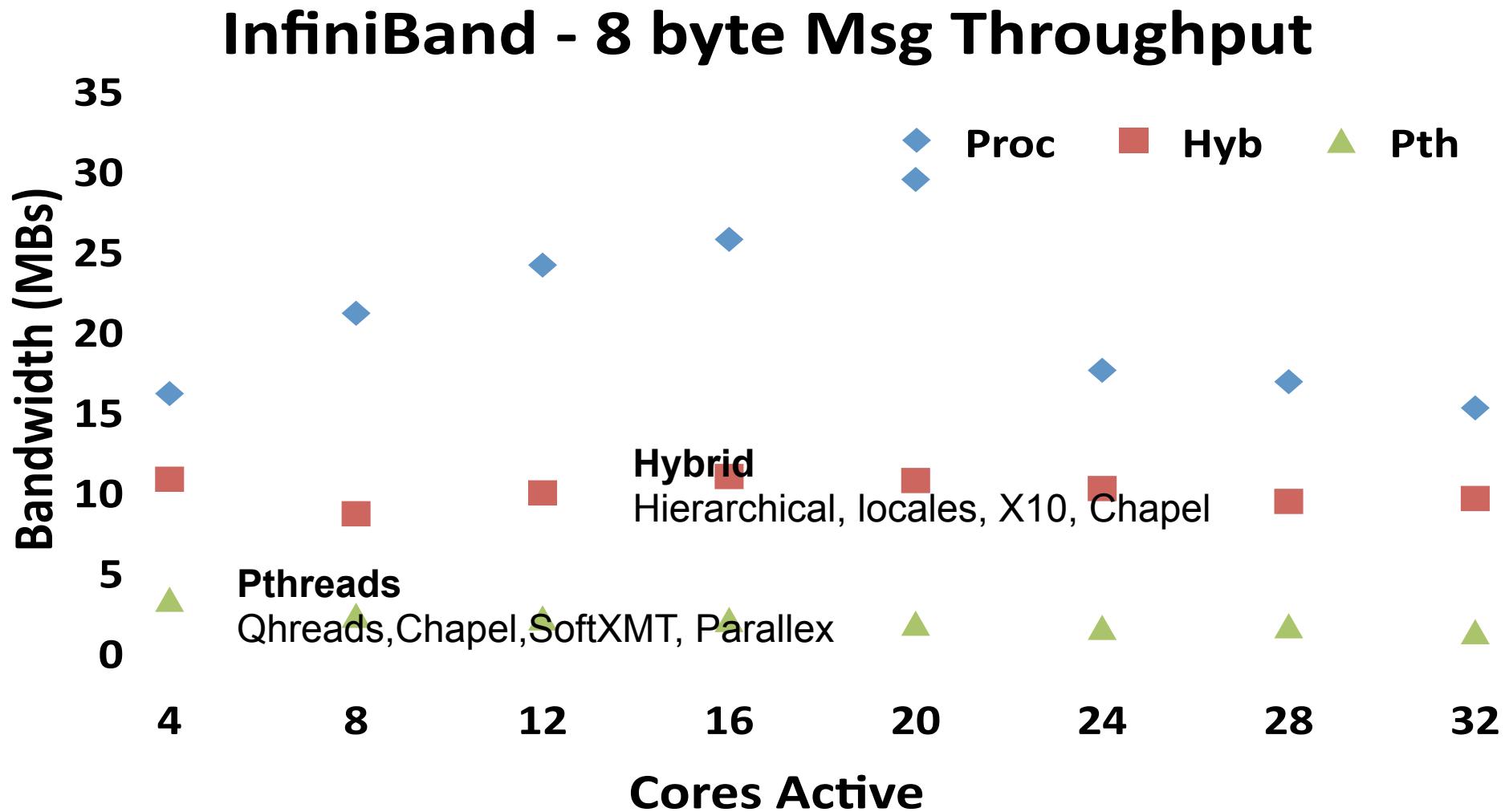


LAWRENCE BERKELEY NATIONAL LABORATORY



InfiniBand Performance

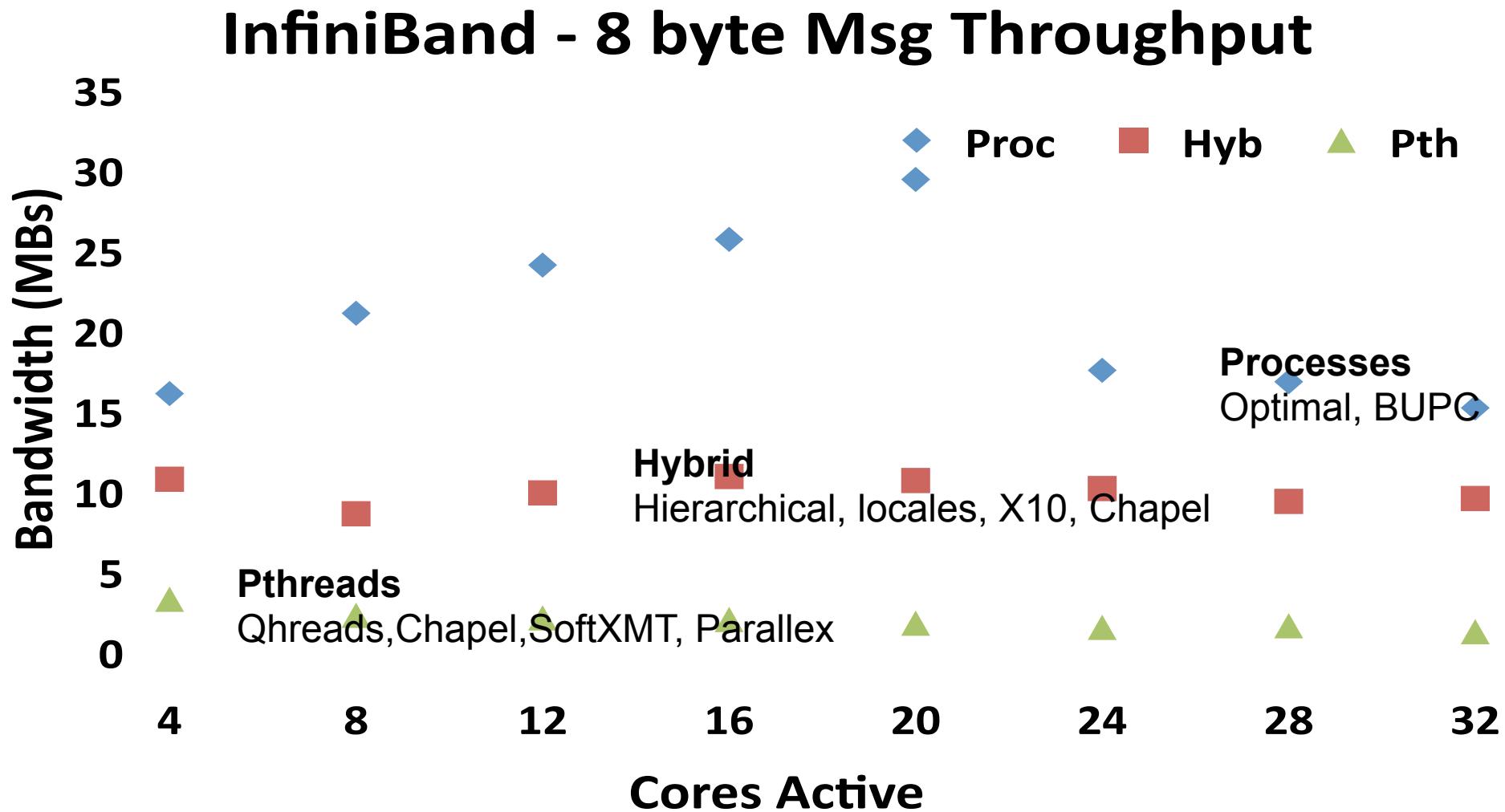
F U T U R E T E C H N O L O G I E S G R O U P





InfiniBand Performance

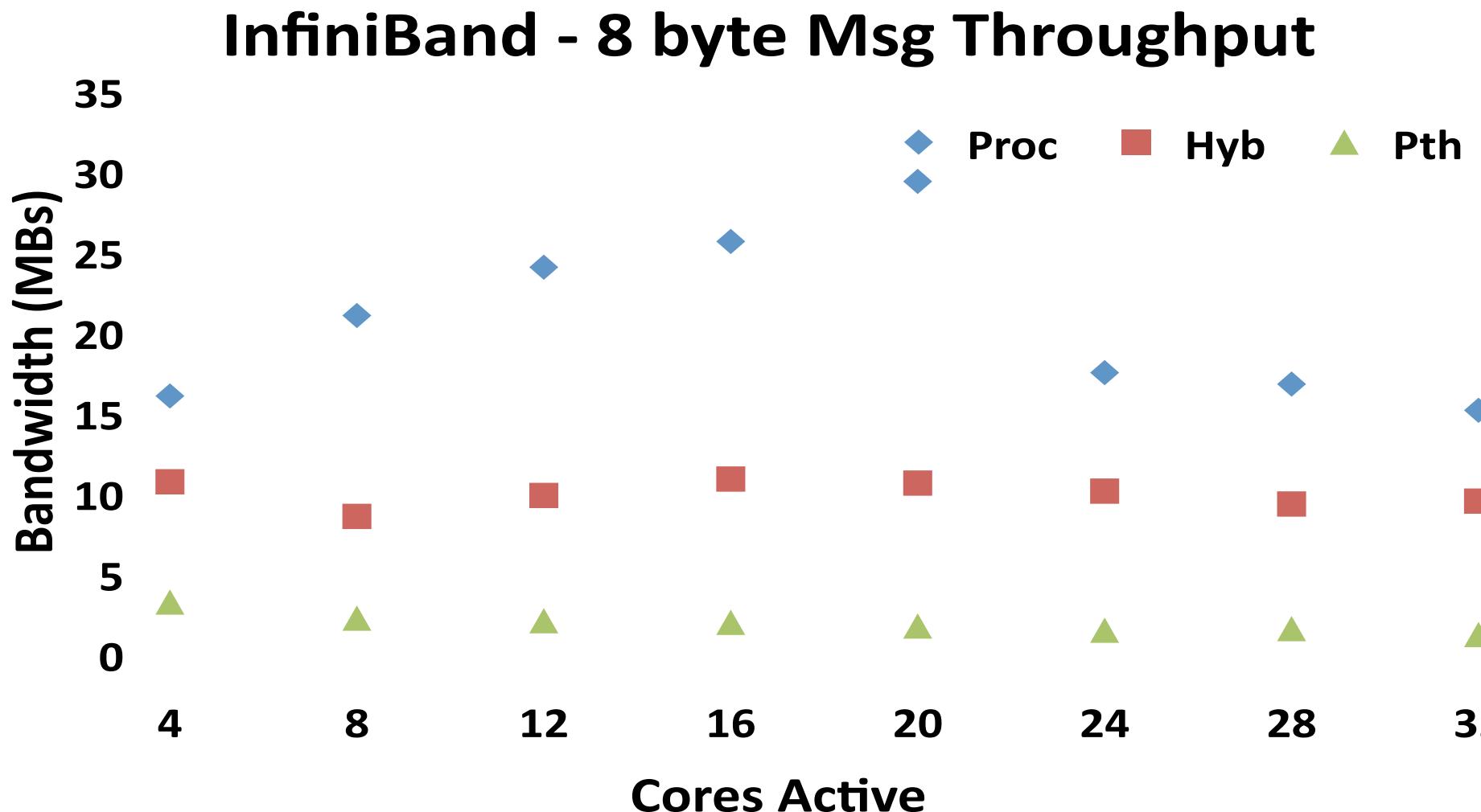
F U T U R E T E C H N O L O G I E S G R O U P





InfiniBand Performance

F U T U R E T E C H N O L O G I E S G R O U P

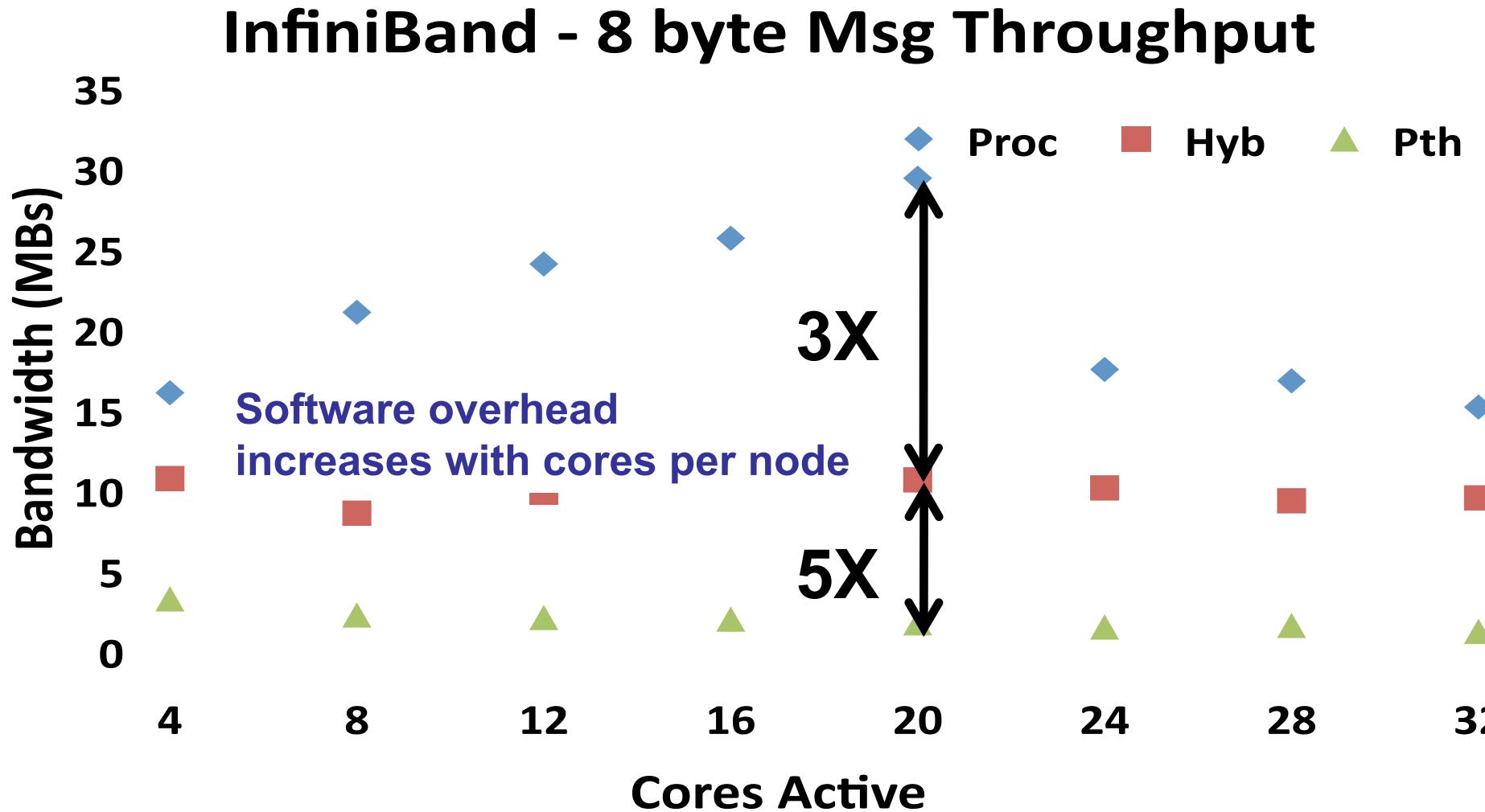


LAWRENCE BERKELEY NATIONAL LABORATORY



InfiniBand Performance

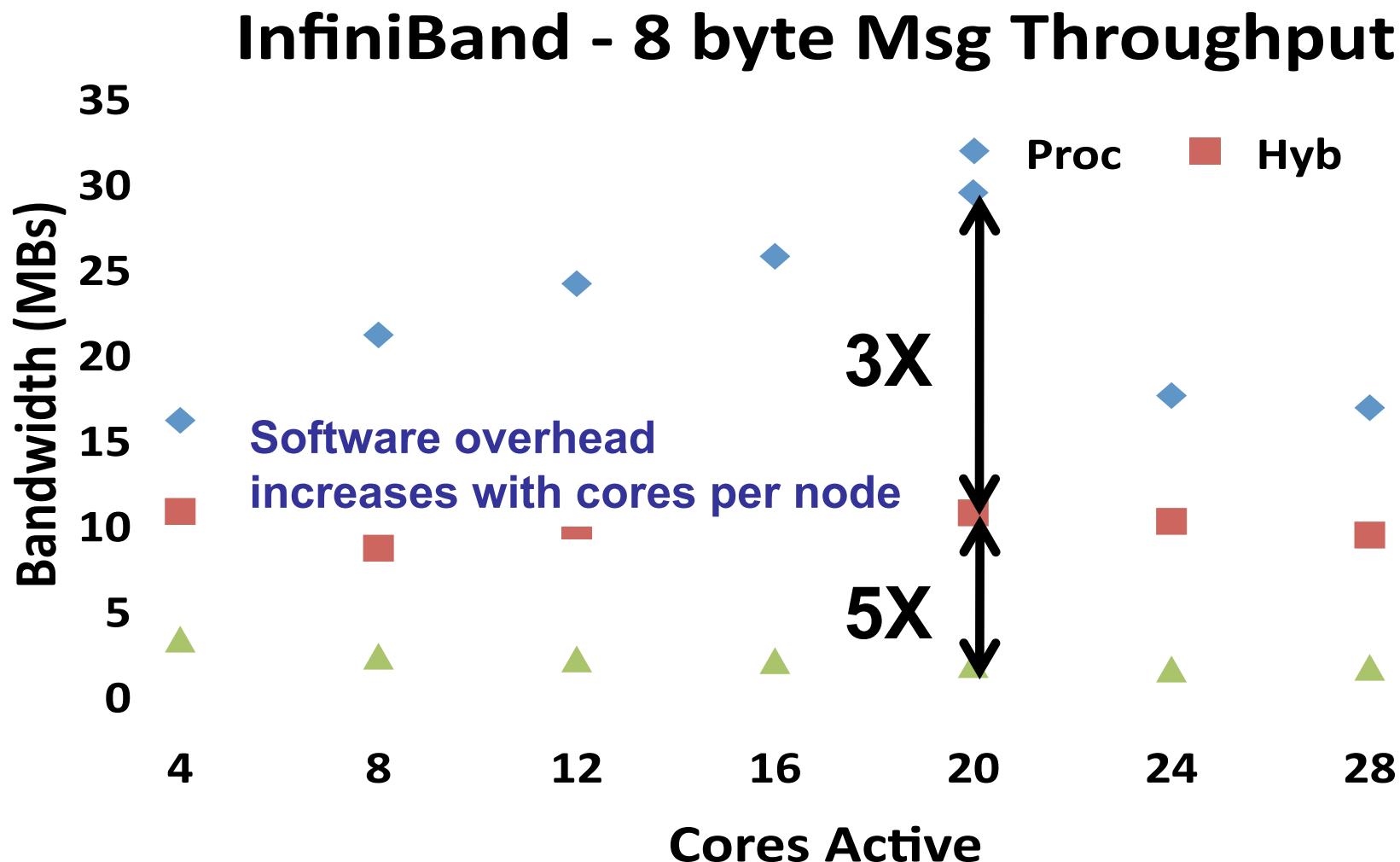
F U T U R E T E C H N O L O G I E S G R O U P





InfiniBand Performance

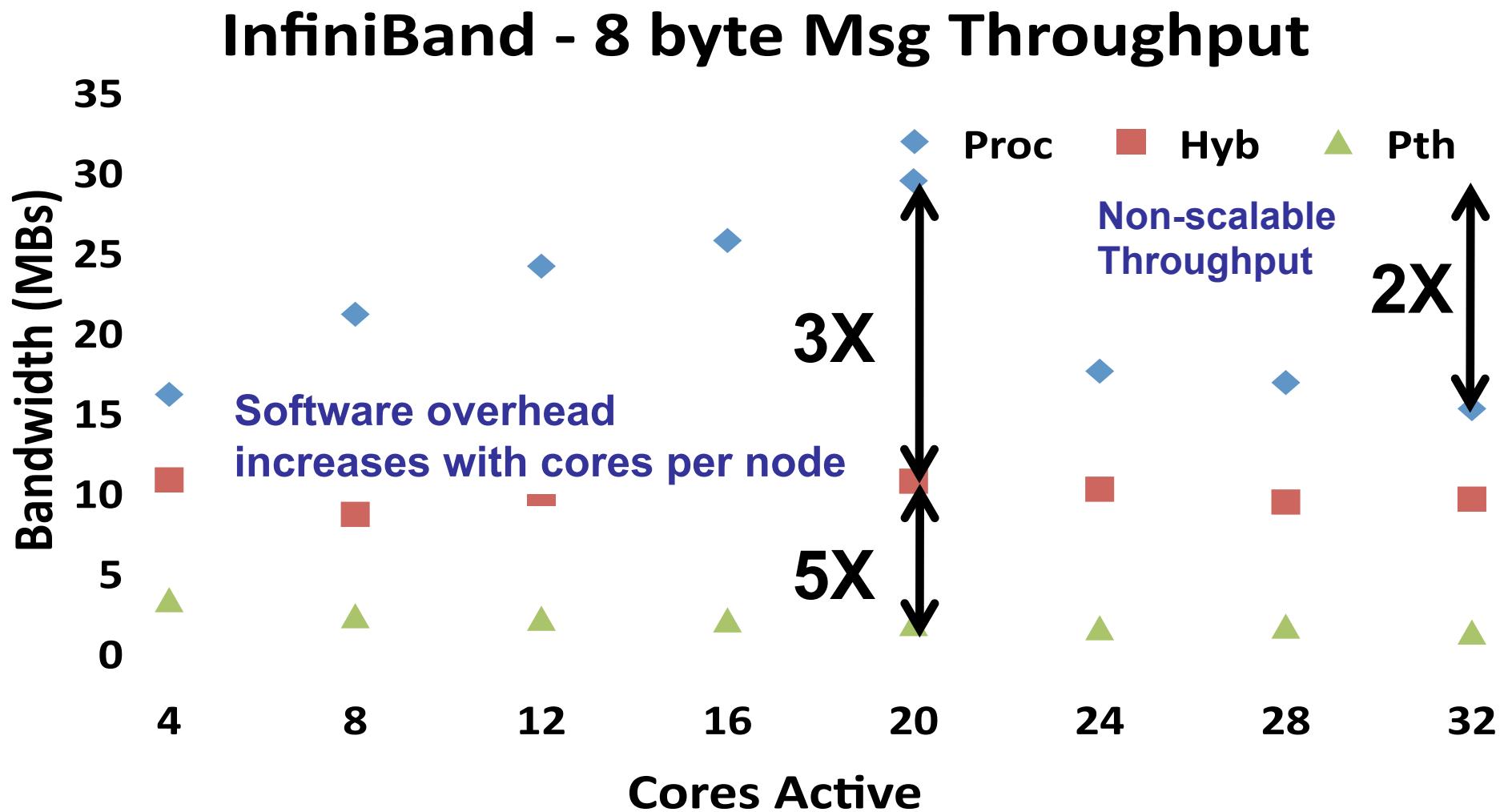
F U T U R E T E C H N O L O G I E S G R O U P





InfiniBand Performance

F U T U R E T E C H N O L O G I E S G R O U P

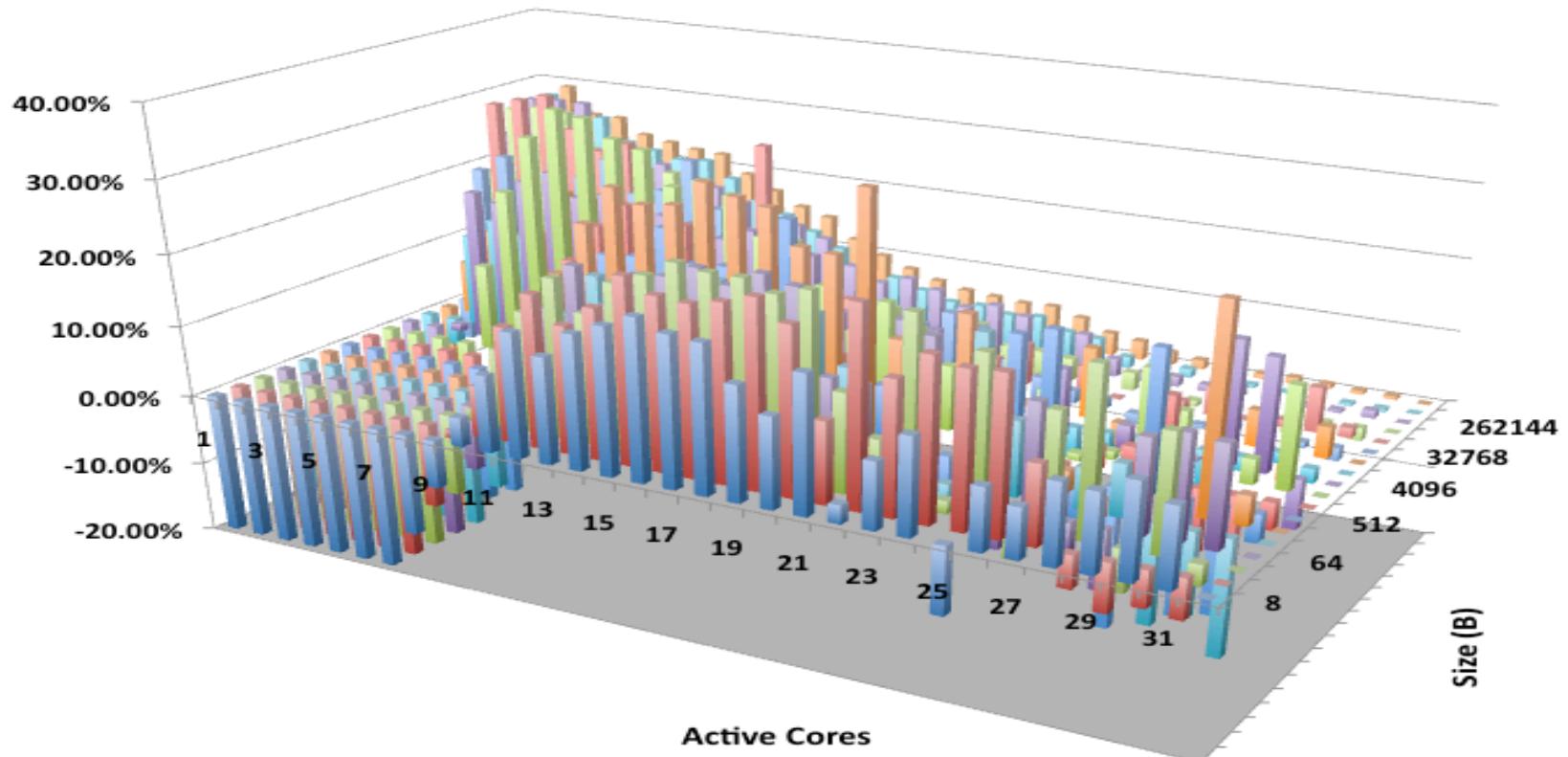




BUPC/GASNet on InfiniBand

F U T U R E T E C H N O L O G I E S G R O U P

Throughput Improvement when Restricting Active Cores: InfiniBand

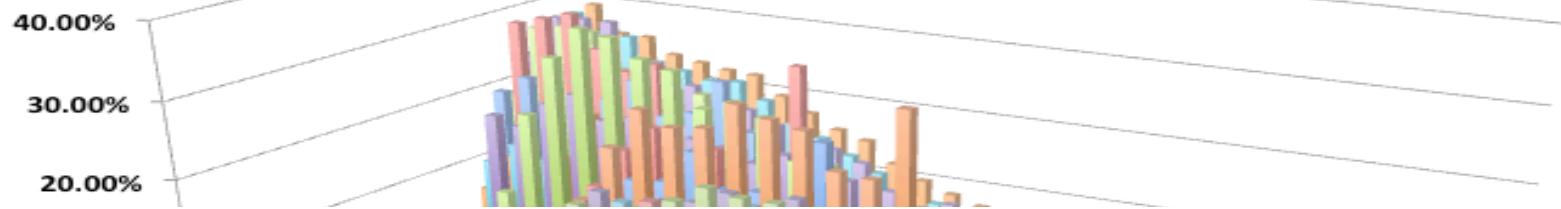




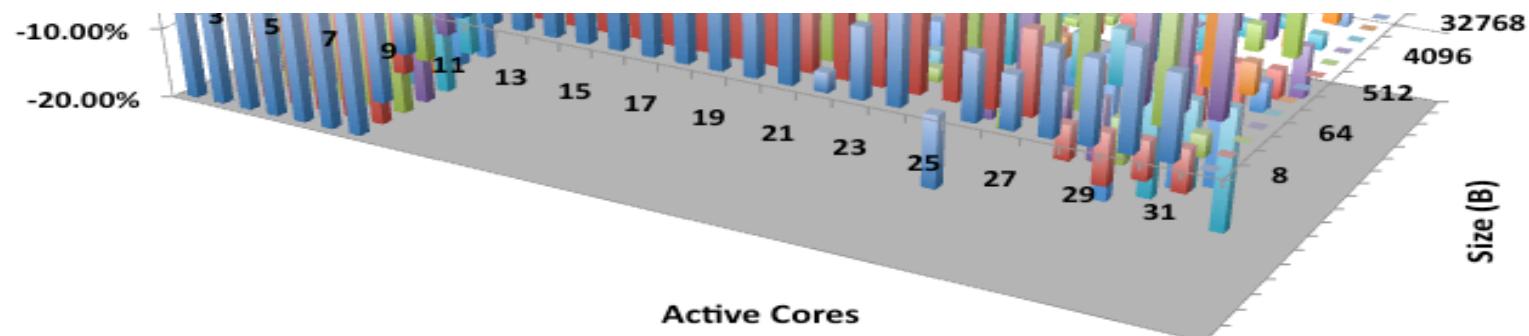
BUPC/GASNet on InfiniBand

F U T U R E T E C H N O L O G I E S G R O U P

Throughput Improvement when Restricting Active Cores: InfiniBand



- 10. Serializing communication using 16 cores 40% faster
 - than using 32 cores

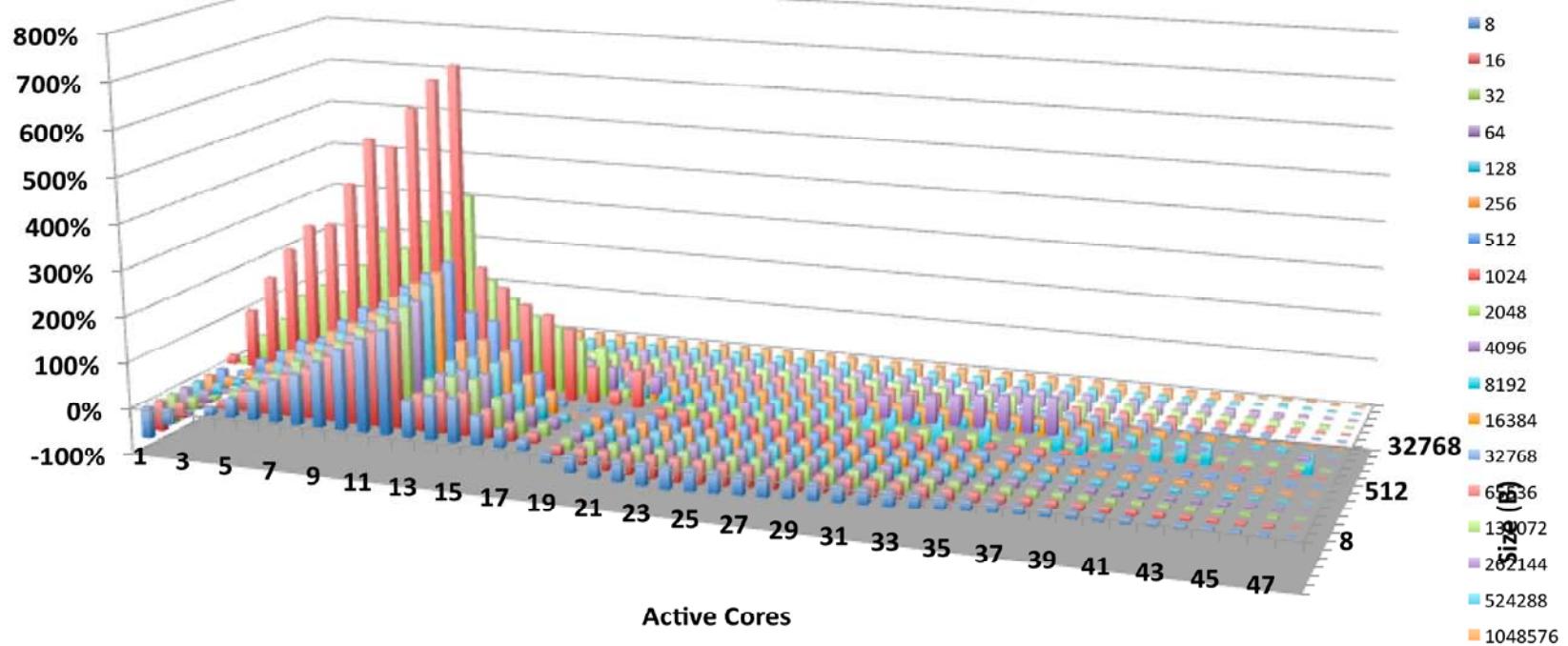




Cray MPI on Gemini

F U T U R E T E C H N O L O G I E S G R O U P

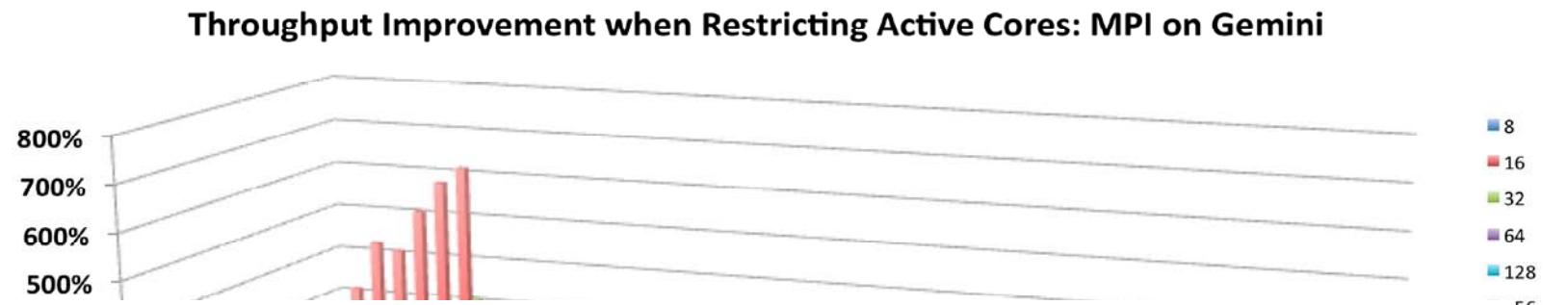
Throughput Improvement when Restricting Active Cores: MPI on Gemini



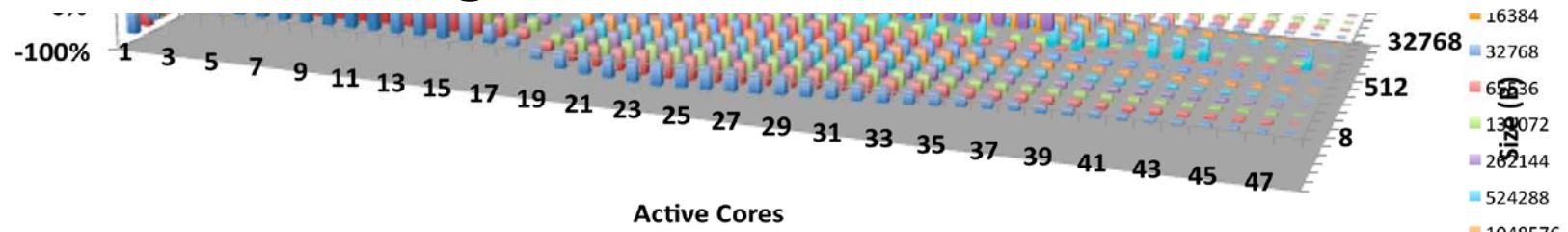


Cray MPI on Gemini

F U T U R E T E C H N O L O G I E S G R O U P



Serializing communication using 12 cores is 6X faster than using 48 cores





Network Performance

F U T U R E T E C H N O L O G I E S G R O U P

- ❖ **Modern NICs are saturated by a small number of messages (cores*messages = K)**
- ❖ **Low latency is not enough**
 - OpenIB Verbs and Cray DMAPP throughput worse than GASNet
- ❖ **Common behavior across software and hardware**
 - InfiniBand (Mellanox), Cray Gemini have similar behavior
 - MPI, OpenIB and DMAPP have similar behavior
- ❖ **Behavior is likely to get worse**



Large (Exa) Scale = Throughput

F U T U R E T E C H N O L O G I E S G R O U P

- ❖ **Future large scale systems are manycore**
 - Homogeneous or heterogeneous hardware
 - Asymmetric software architecture or performance
- ❖ **Multithreading proven to work for irregular applications**
- ❖ **Many research projects advocate “fine grained” multithreading** (SoftXMT, Qthreads, ParalleX, SWARM)
 - Energy efficiency
 - Data movement
 - Latency hiding
 - Load balancing
 - Performance and productivity
- ❖ **Fine grained multithreading requires network throughput**



BERKELEY LAB

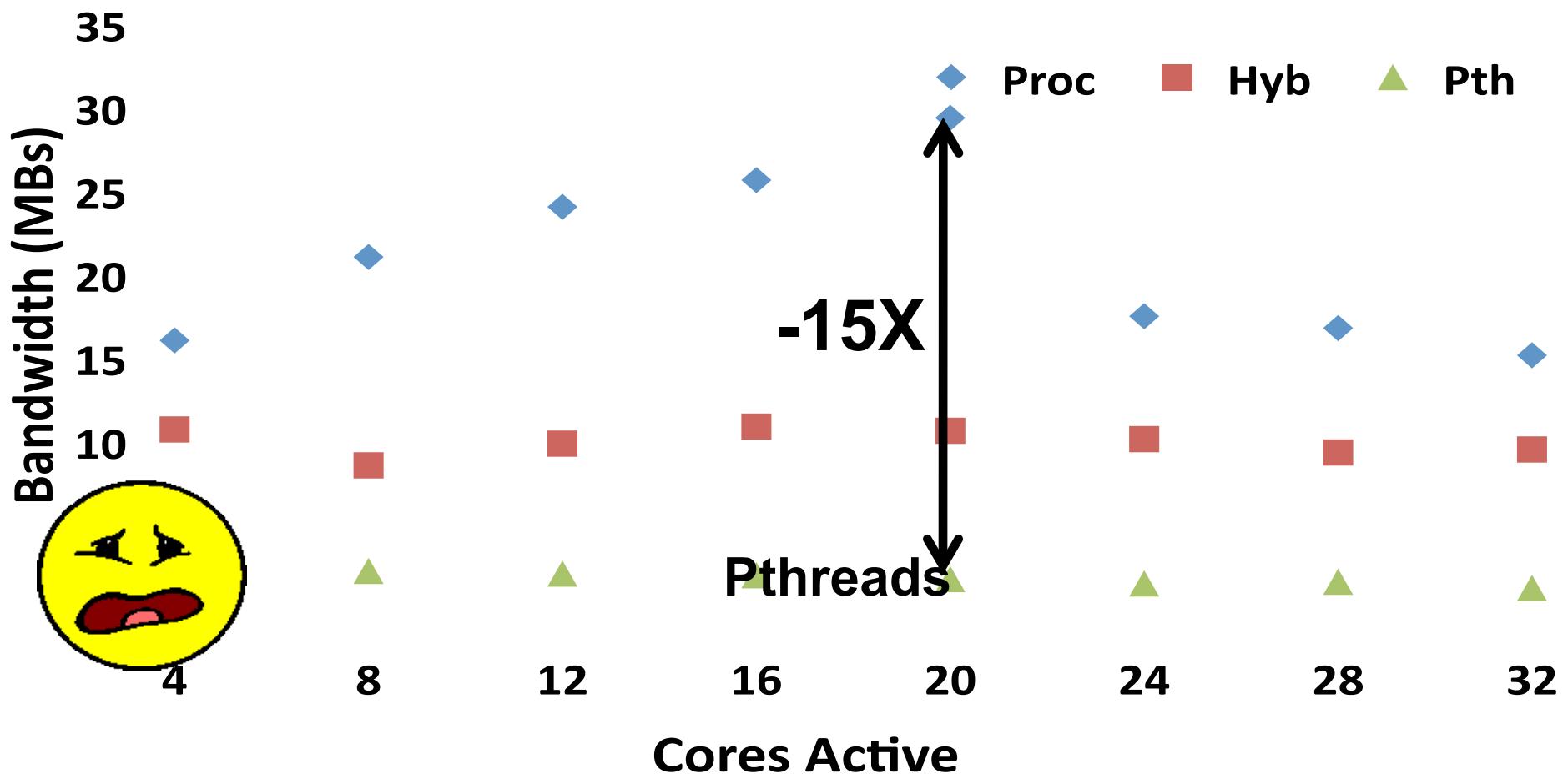
F U T U R E T E C H N O L O G I E S G R O U P

**We do not know that we do not know how to
build integrated “threading” and networking
runtimes that provide good throughput**

(building user level threading mapped on pthreads)



F U T U R E T E C H N O L O G I E S G R O U P



LAWRENCE BERKELEY NATIONAL LABORATORY



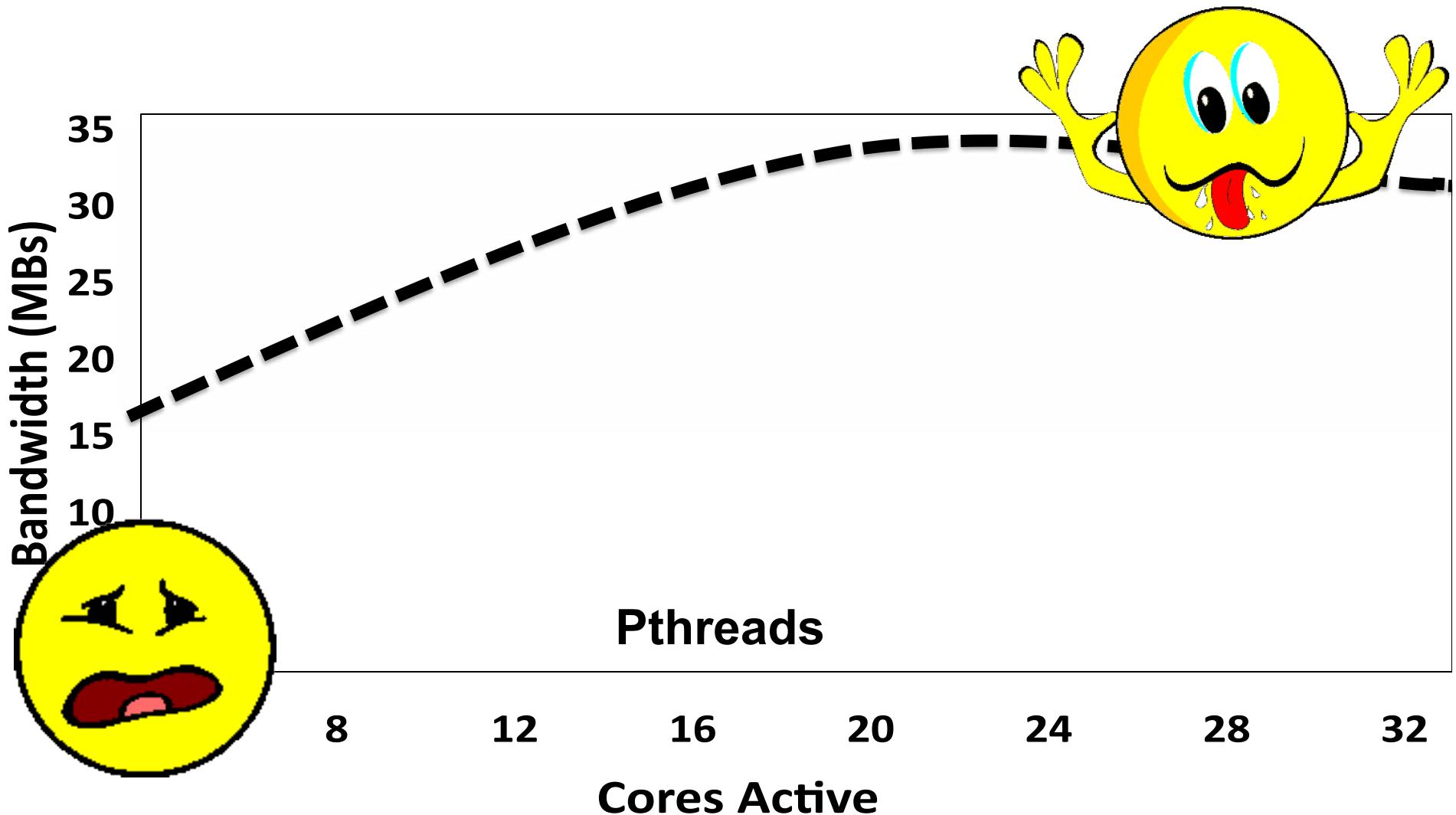
F U T U R E T E C H N O L O G I E S G R O U P

Throughput Oriented Runtime for Large Scale Manycore Systems (THOR)

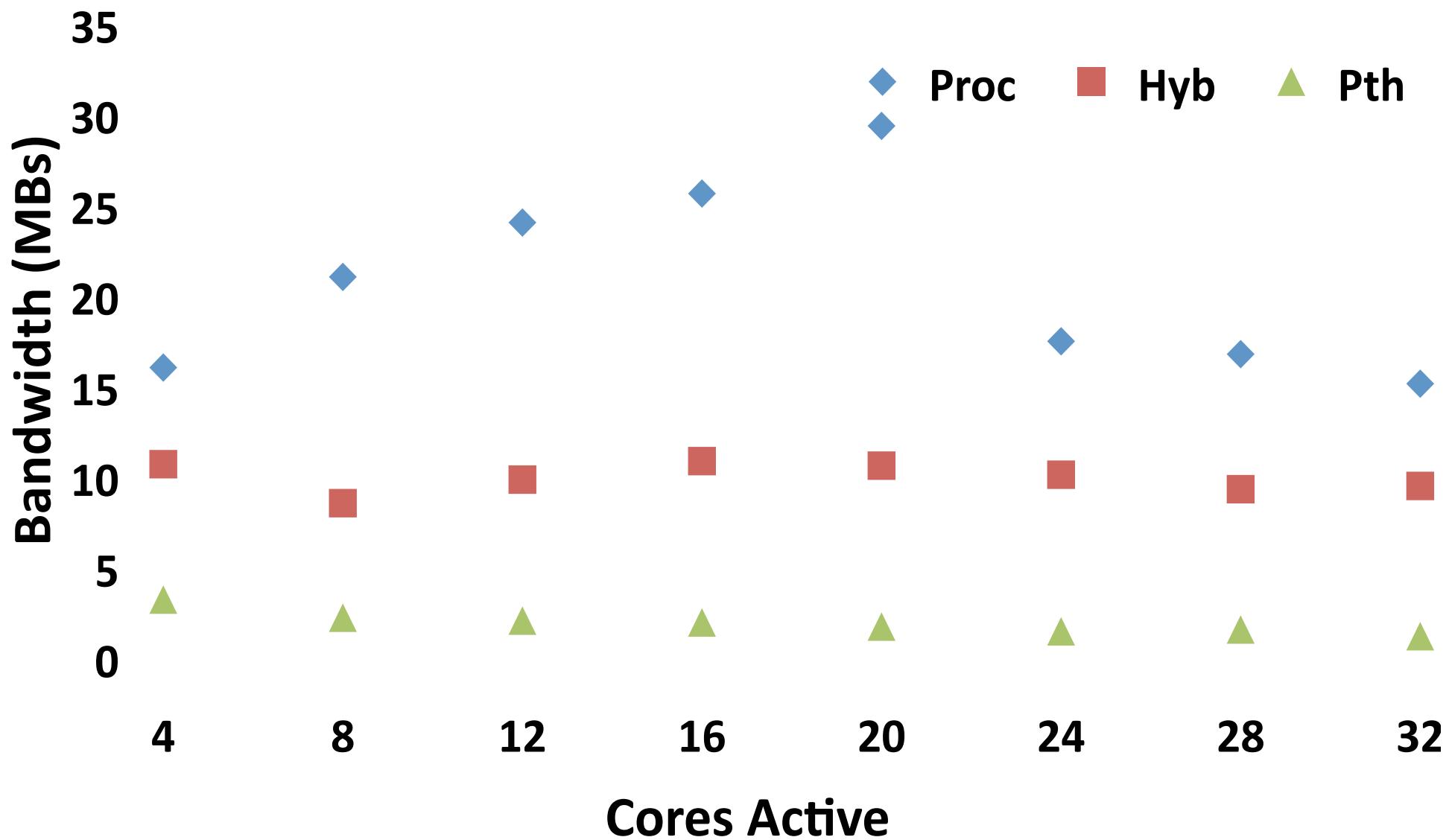


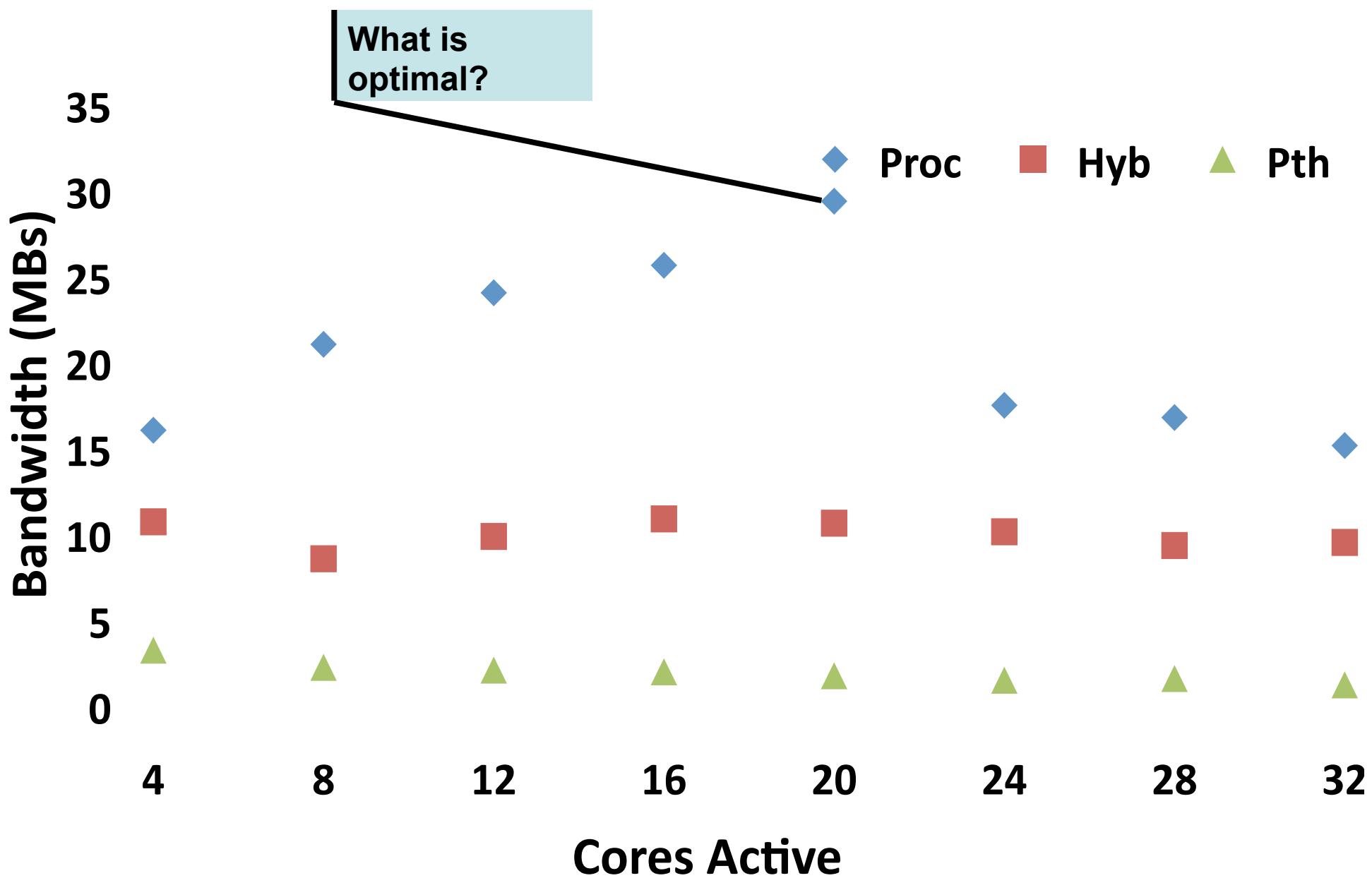
THOR

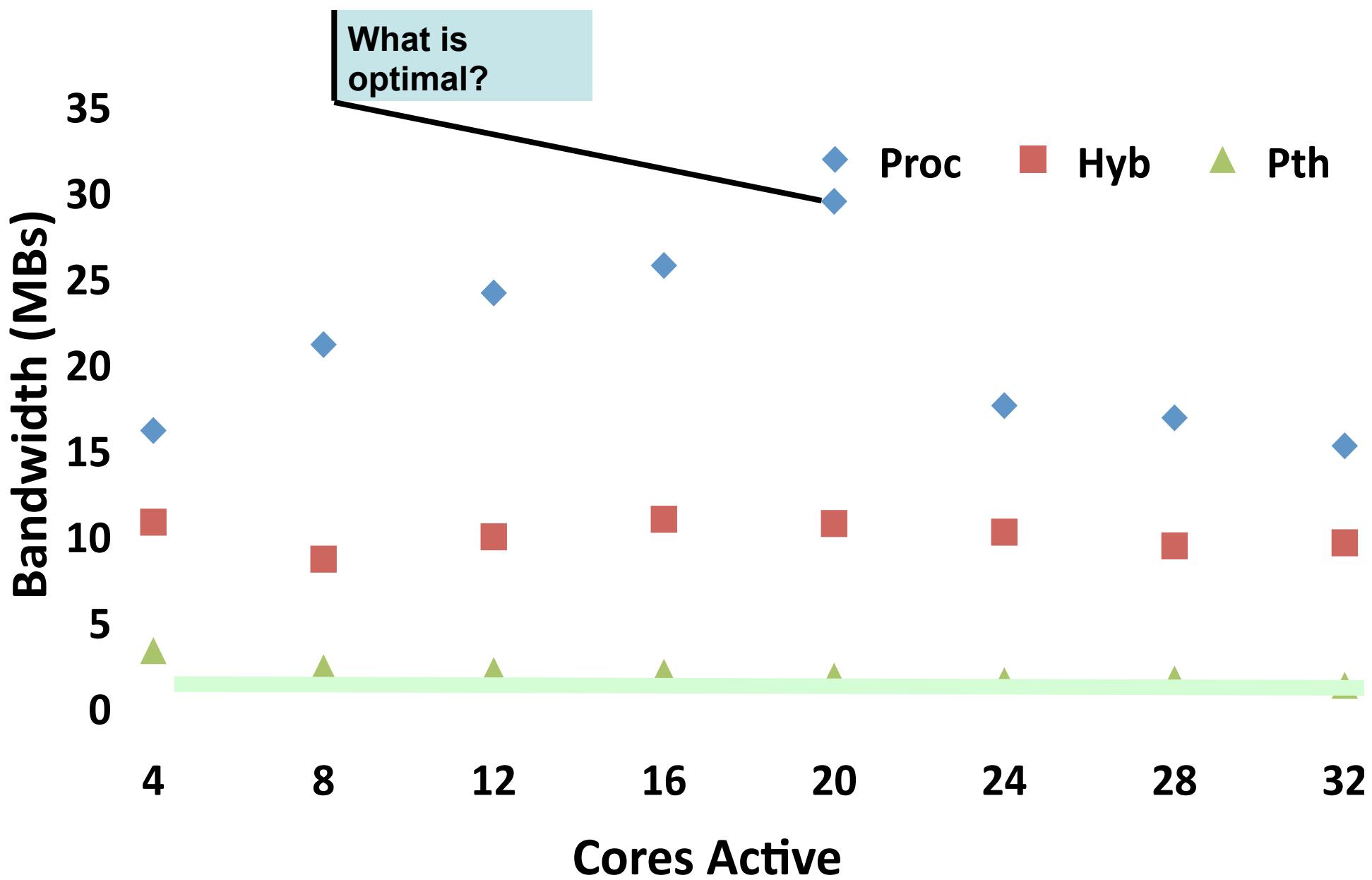
F U T U R E T E C H N O L O G I E S G R O U P

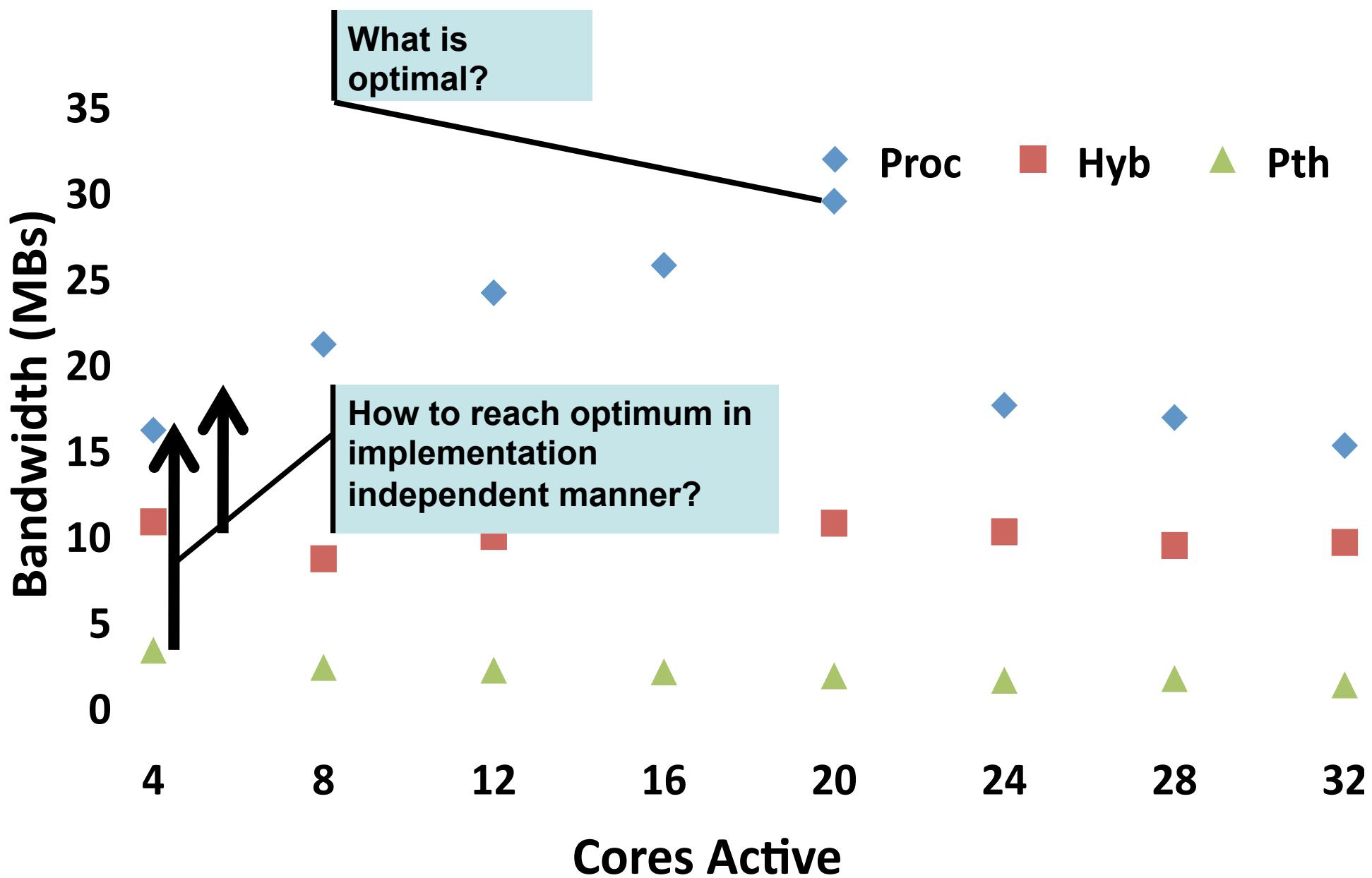


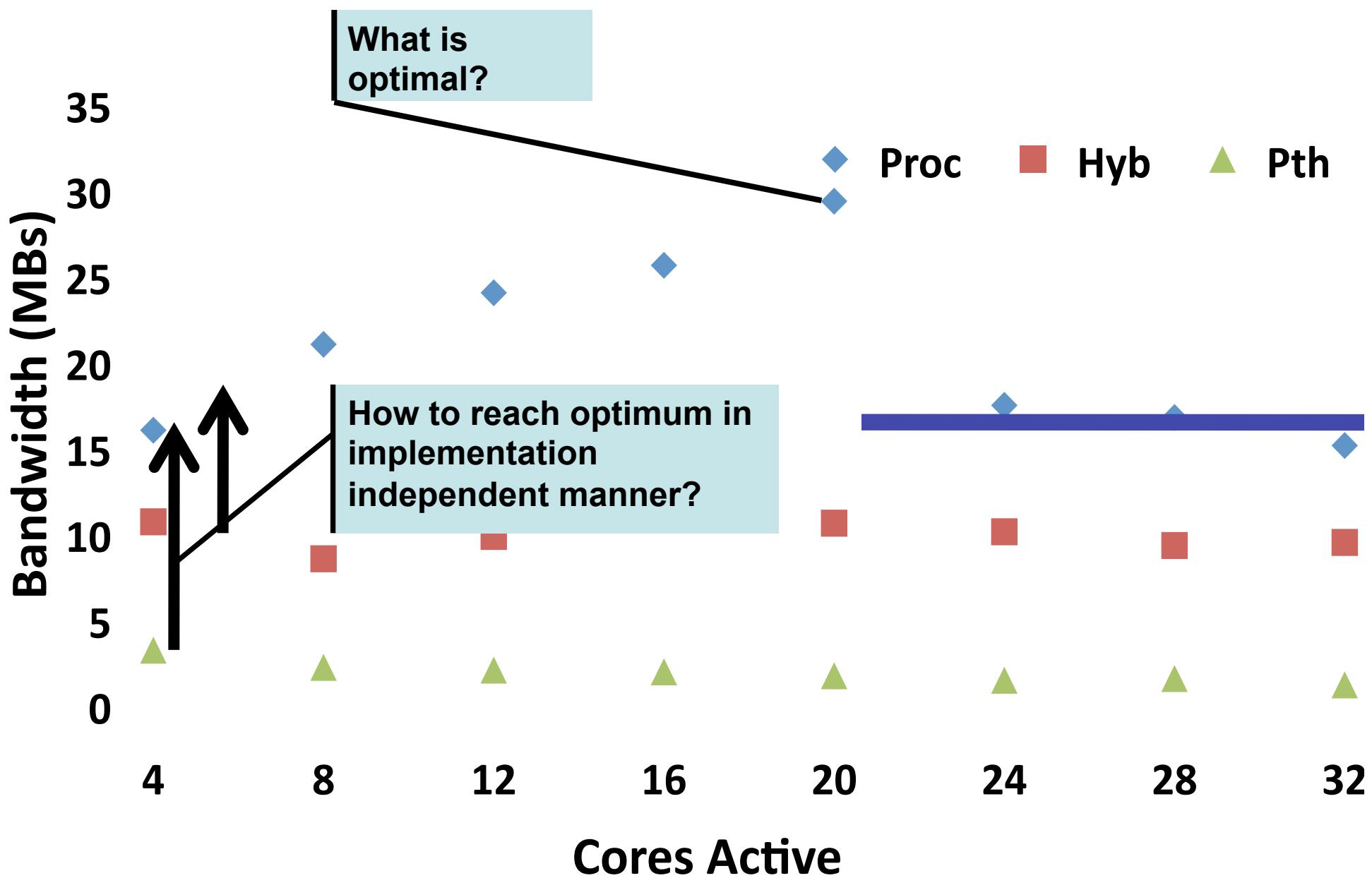
LAWRENCE BERKELEY NATIONAL LABORATORY

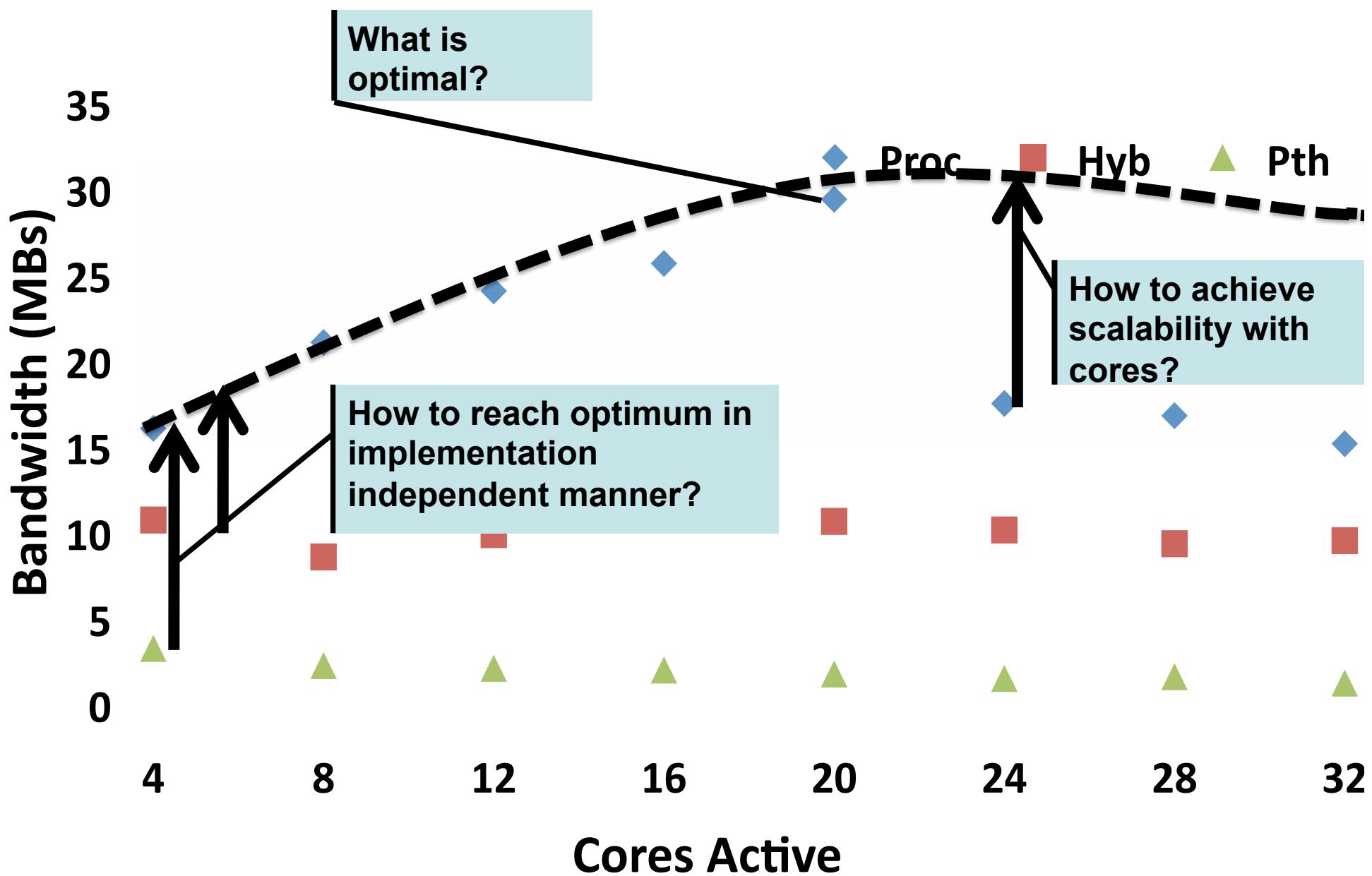


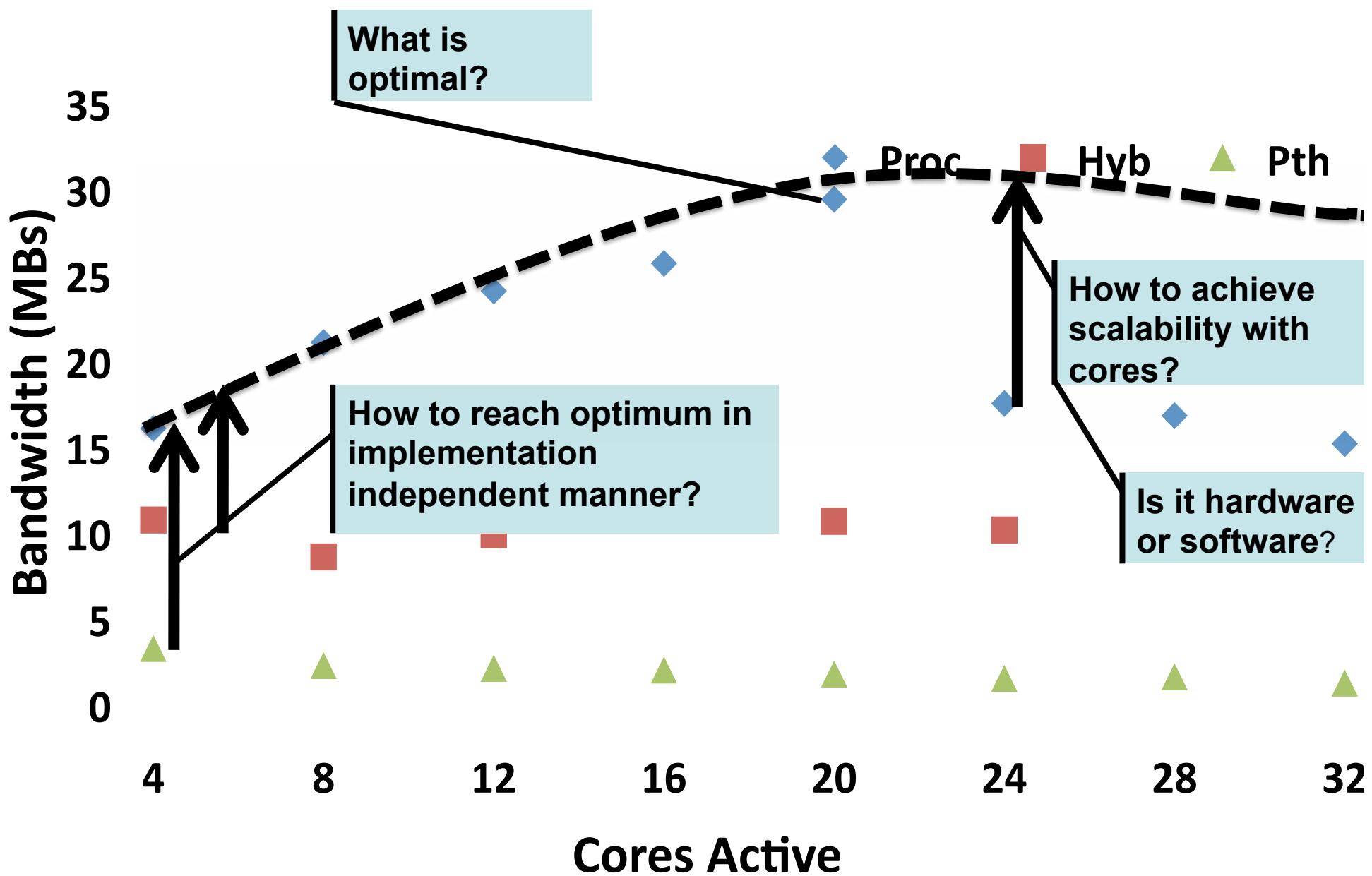














Goals

F U T U R E T E C H N O L O G I E S G R O U P

- ❖ **Efficiency and productivity layer for heavily threaded/asynchronous runtimes**
 - Decouple application/runtime level concurrency from runtime concurrency
 - Manage asynchrony for clients
 - Provide optimal throughput for
 - Any implementation (pthreads, procs ...)
 - Any hardware architecture (asymmetric, heterogeneous)
 - Any message mix
 - Any source, target



THOR Architecture

F U T U R E T E C H N O L O G I E S G R O U P

**Programming Models
(SPMD, task and data parallel) (UPC, Chapel)**

Runtimes: SoftXMT, BUPC, Qthreads

Admission Control Layer

Optimization Layer

Scheduling Layer

GASNet/MPI

**T
H
O
R**

Driven by runtime analysis and performance models



Thor Layers

F U T U R E T E C H N O L O G I E S G R O U P

❖ Admission Control Layer

- Congestion Avoidance
- Flow Control
- Memory Consistency/Ordering
- Dispatch to Optimization Services

❖ Optimization Layer

- Coalescing
- Aggregation

❖ Scheduling Layer

- Reordering
- Instantiate and Retire Communication to Network



Scalable/Portable Design

F U T U R E T E C H N O L O G I E S G R O U P

❖ Multiple implementations for asymmetric/heterogeneous hardware

- Inline: mechanisms implemented in a distributed manner, e.g. GASNet/MPI
- Proxy: servers acting on behalf of clients

❖ Open loop control

- With as little “global” state as possible

❖ Declarative behavior

- Intuitive (human descriptions) – e.g. train of messages
- Annotated by compilers/optimizers



Scalability with Cores

F U T U R E T E C H N O L O G I E S G R O U P

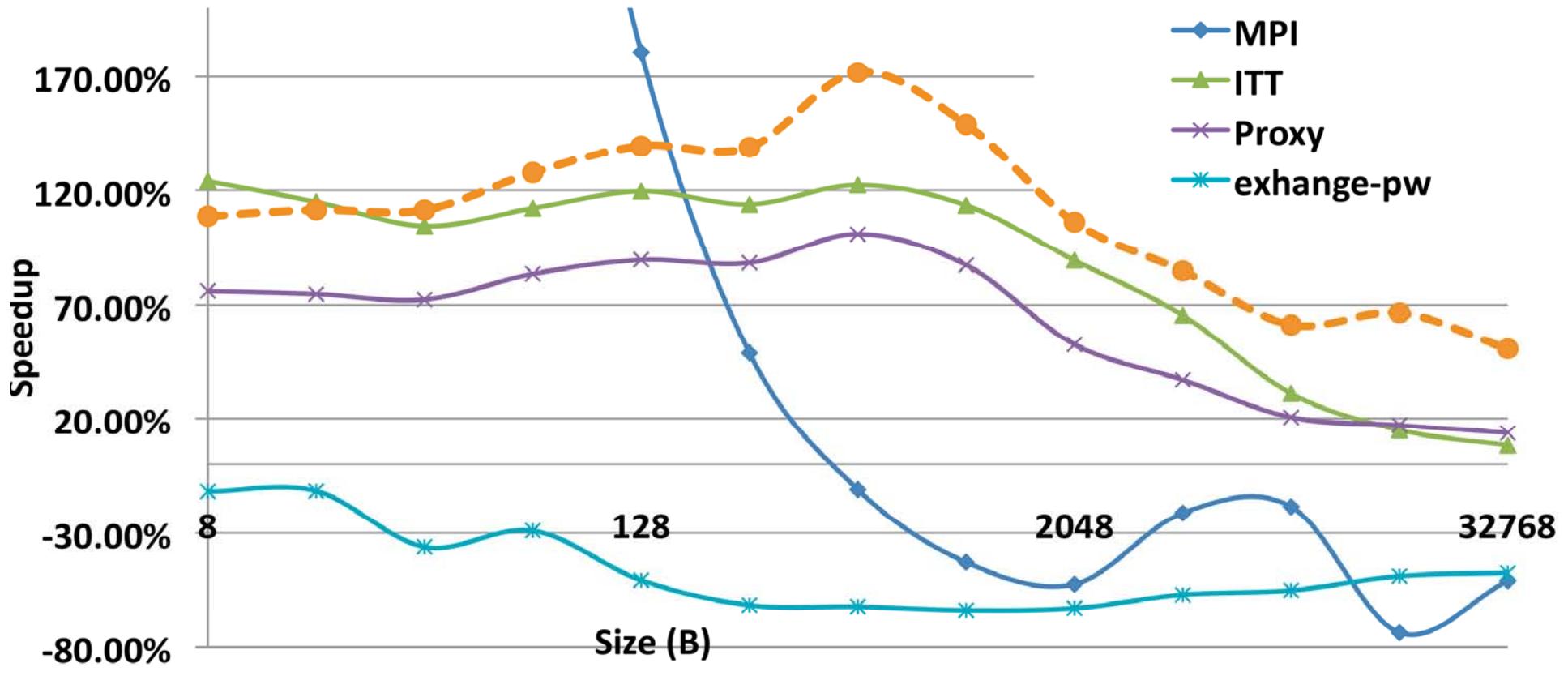
- ❖ **Congestion avoidance runtime prototyped**
 - BUPC/GASNet/InfiniBand
 - Cray UPC/DMAPP/Gemini
- ❖ **Admission Control + Scheduling Layer**
 - Not well tuned yet
- ❖ **Results:**
 - **4X** performance improvement for all-to-all
 - **70%** improvement on GUPS/HPCC RA
 - **17%** on NAS Parallel Benchmarks

To appear as “Congestion Avoidance on Manycore Clusters” in ICS 2012



All-to-all InfiniBand 1024 Cores

F U T U R E T E C H N O L O G I E S G R O U P

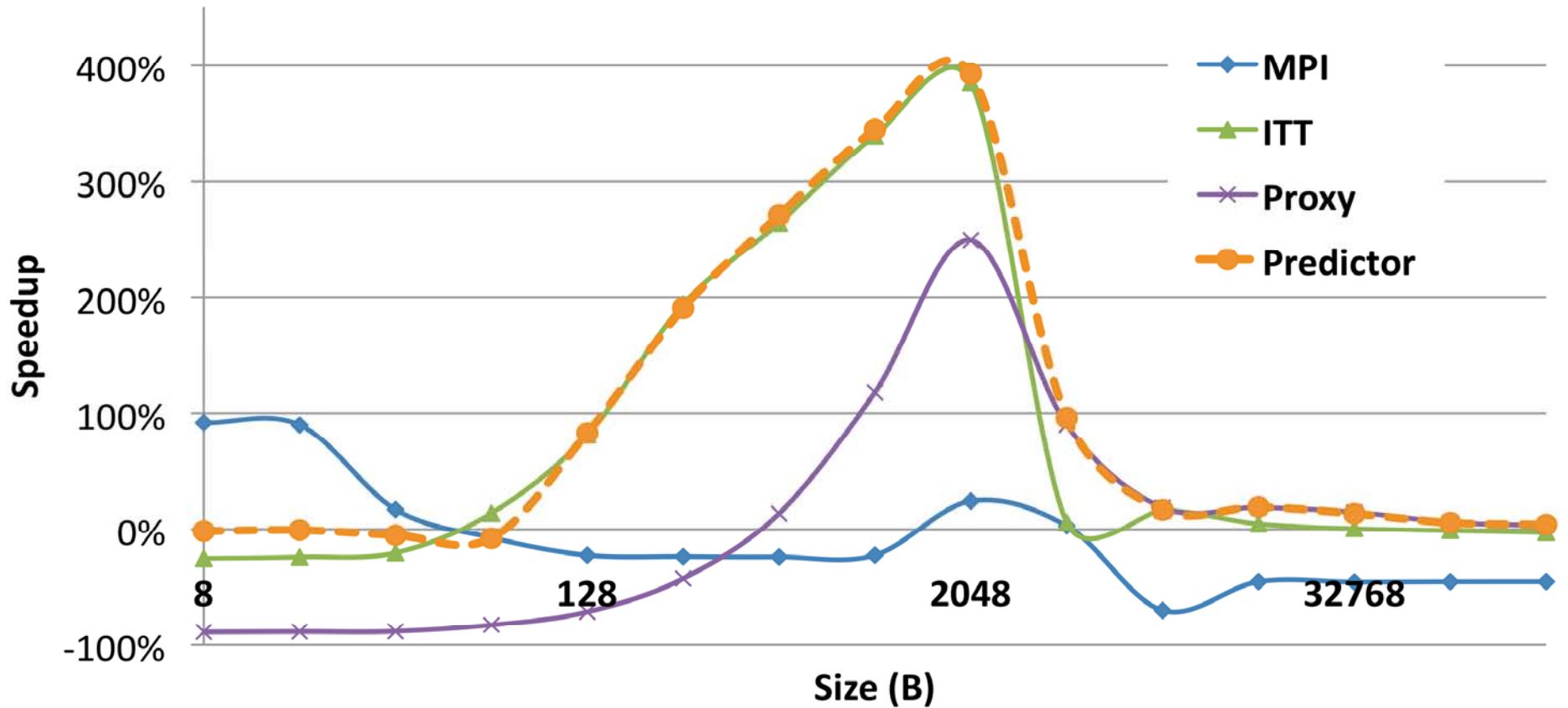


**Speedup over GASNet tuned all-to-all - 2x
Performance Portable – single implementation**



All-to-all Gemini 768 Cores

F U T U R E T E C H N O L O G I E S G R O U P



**Speedup over Cray UPC all-to-all - 4x
Performance Portable – single implementation**



F U T U R E T E C H N O L O G I E S G R O U P

Future Work



Principles of HW/SW Design

F U T U R E T E C H N O L O G I E S G R O U P

- ❖ Why don't we get maximal throughput?
- ❖ How do we get maximal throughput?



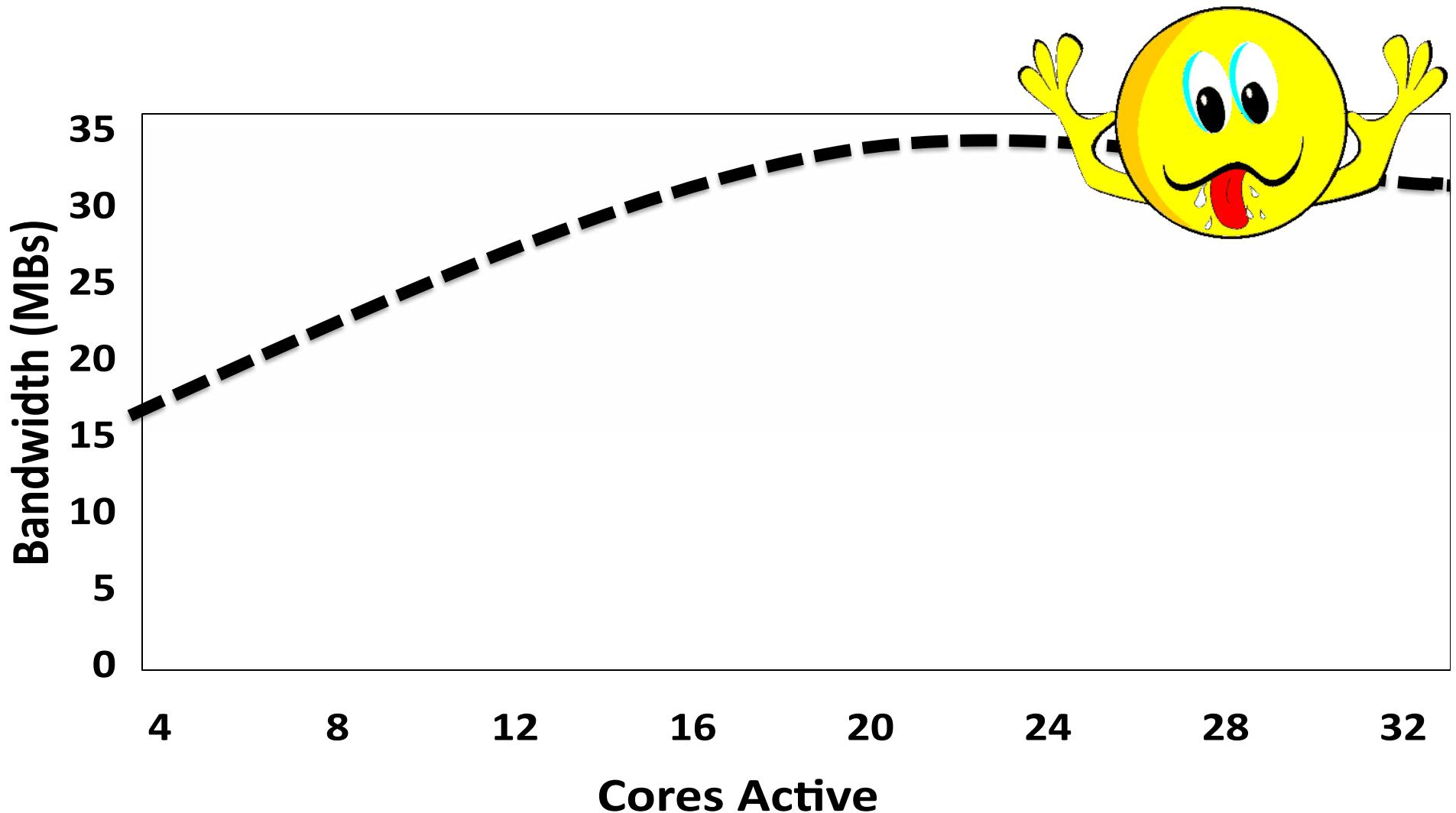
F U T U R E T E C H N O L O G I E S G R O U P

Thank You!



THOR

F U T U R E T E C H N O L O G I E S G R O U P



LAWRENCE BERKELEY NATIONAL LABORATORY



BERKELEY LAB

THOR Timeline

F U T U R E T E C H N O L O G I E S G R O U P

- ❖ **Year 1:** Performance study and prototype implementation
- ❖ **Year 2:** Admission Control Layer + Scheduling Layer
- ❖ **Year 3-4:** Optimization Layer and tuning on target systems

- ❖ **Demonstration:**
 - Languages: UPC, Chapel
 - Runtimes: BUPC, SoftXMT (Qthreads?)
 - Networks: InfiniBand, Cray
 - CPUs: Intel/AMD, Knights Corner/Landing (other suggestions)