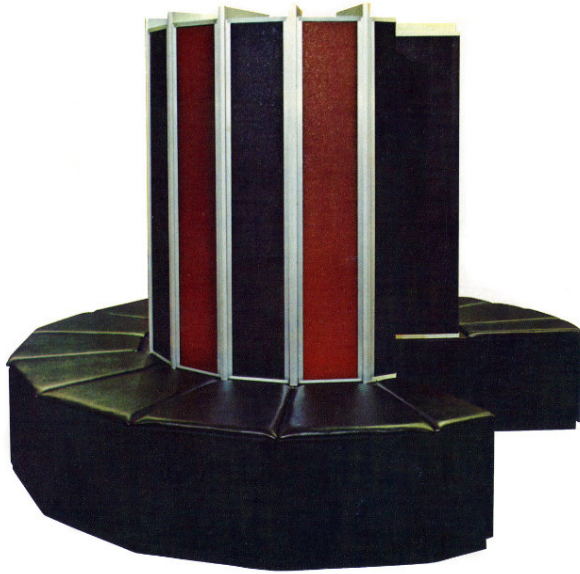# *A Perspective on The Path Forward*
## (Why I'm not *too* worried about Exascale)

**Steve Scott**

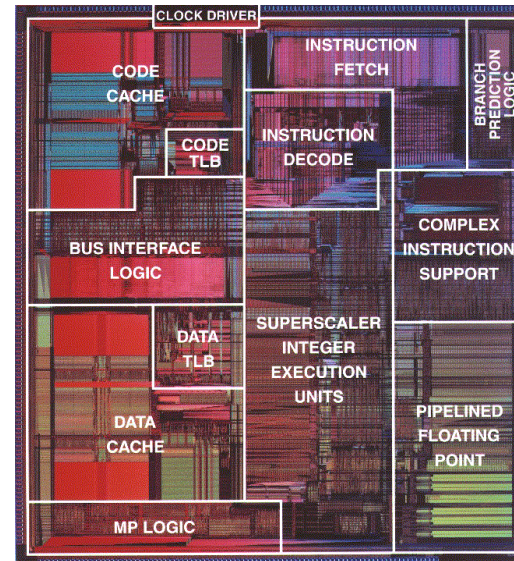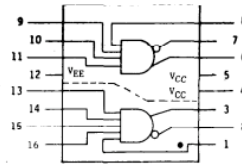Cray CTO

SOS20

March 23, 2016

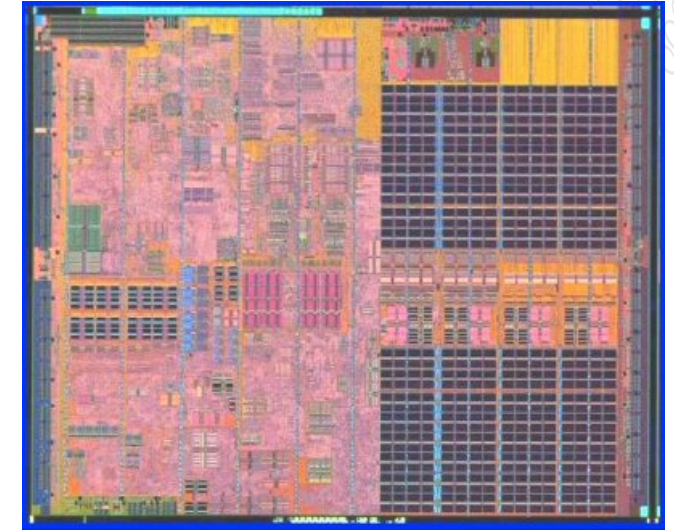# Technology Drives Architecture



**Cray 1, 1976**



**Intel Pentium, 1993**



**Intel Pentium 4 Cedar Mill, 2006**

- ECL 5/4 NAND gate ICs (95%)
- 75K gates.  (3400 PCBs!)
- RISC design
- Vector ISA
- Memory latency 11 clocks

- CMOS VLSI IC
- 3M transistors
- CISC design
- Scalar ISA
- Deep pipelines, complex predictions

- 184M transistors!
- *Very* CISC design
- 31-stage pipeline
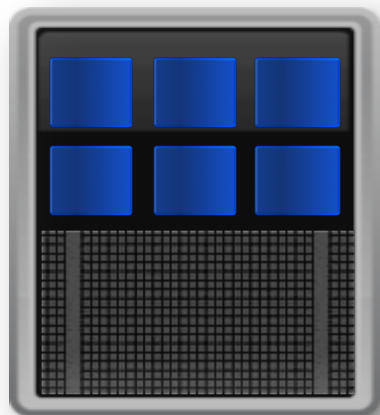- 3.6 GHz in 65nm
- Last of its breed….

# And then Dennard scaling ended…

**Power constrained**

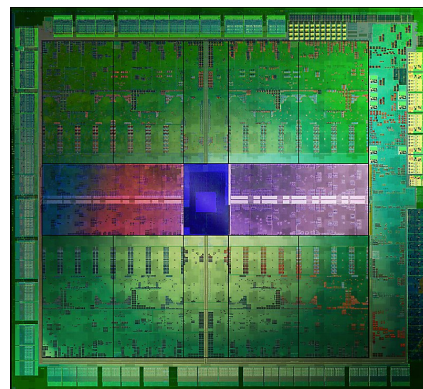**Communication much more expensive than Computation**

# New Processor Landscape
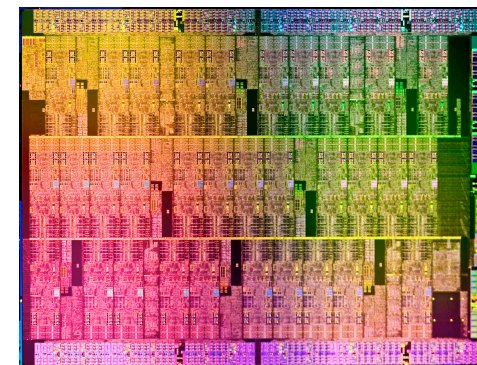## Driven by power efficiency



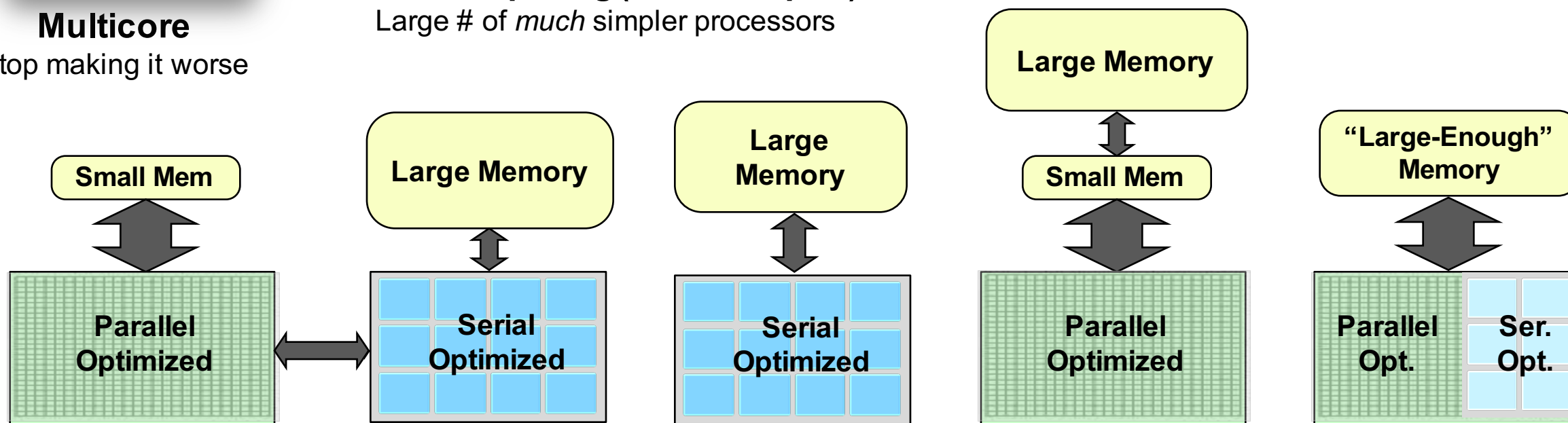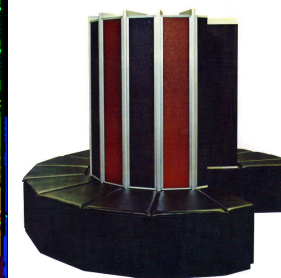**Multicore**
Stop making it worse

**GPU computing (Nvidia Kepler)**
Large # of *much* simpler processors

**Vector Computing (Intel Xeon Phi)**
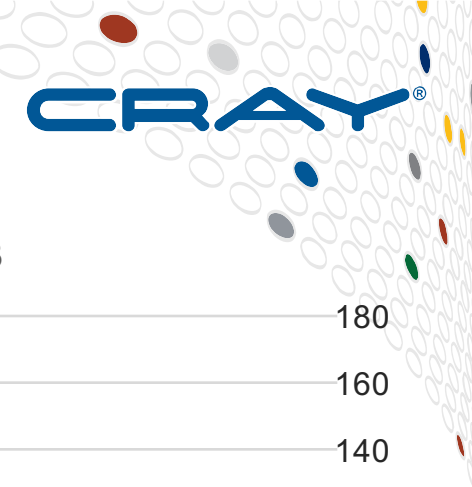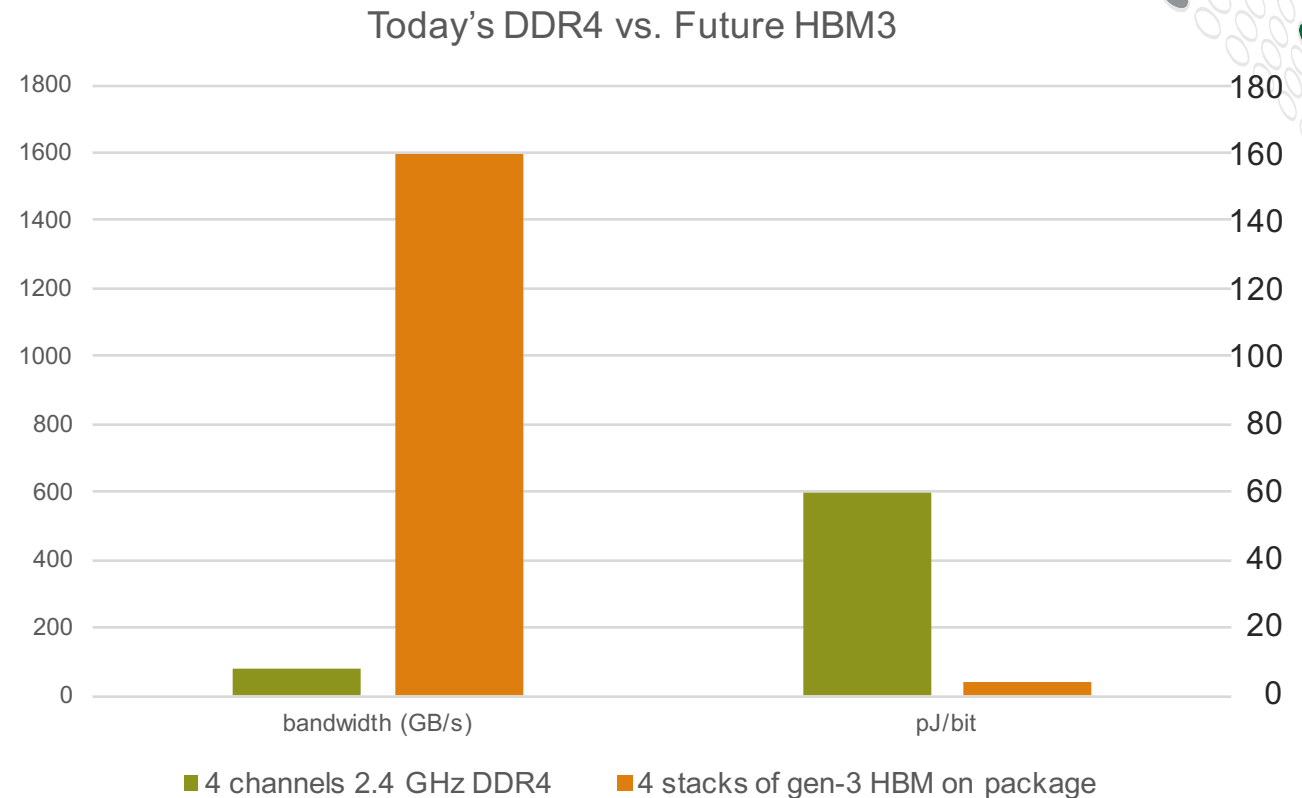Parallelism with low complexity & control overhead

| Small Mem | Large Memory | Large Memory | Large Memory / Small Mem | "Large-Enough" Memory |

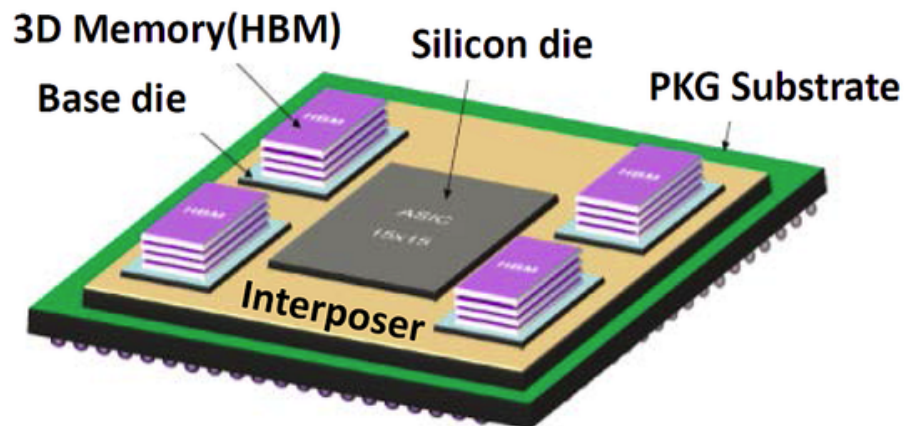| Parallel Optimized | Serial Optimized | Serial Optimized | Parallel Optimized | Parallel Opt. / Ser. Opt. |

# On-Package Memory Can Restore Balance

- Standard DDR memory BW has not kept pace with CPUs

- Expect processors to adopt stacked, on-package memory

- HBM:
  - 10x higher BW, 10x less energy/bit
  - Much lower latency
  - Costs ~2x DDR4 per bit
  - JDEC standard with multiple sources

Today's DDR4 vs. Future HBM3

Chart axis labels (left): 1800, 1600, 1400, 1200, 1000, 800, 600, 400, 200, 0
Chart axis labels (right): 180, 160, 140, 120, 100, 80, 60, 40, 20, 0
X-axis: bandwidth (GB/s) | pJ/bit

Legend:
- ■ 4 channels 2.4 GHz DDR4
- ■ 4 stacks of gen-3 HBM on package



3D Memory(HBM)
Base die
Silicon die
PKG Substrate
Interposer

*May drive us to smaller, simpler nodes that are balanced with on-package memory*

# Deeper Memory and Storage Hierarchy



**Node memory moving on package**

On Node
- CPU
- Memory (DRAM)

Off Node
- Storage (HDD)

On Node
- CPU
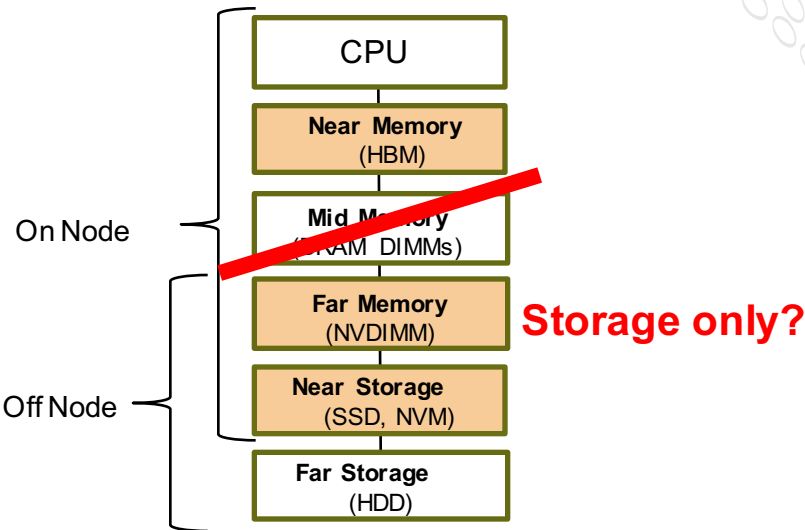- Near Memory (HBM)
- Mid Memory (DRAM DIMMs)

Off Node
- Far Memory (NVDIMM) — **Storage only?**
- Near Storage (SSD, NVM)
- Far Storage (HDD)

**Cold storage moving to disk**

**Primary storage moving to Flash**

3D Xpoint    PCM    ReRAM    STT-MRAM

**New technologies coming to bridge memory-Flash gap**

# Storage Will Scale

- **APEX requirement: Time to checkpoint 80% memory < (0.005)*JMTTI**
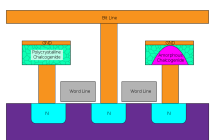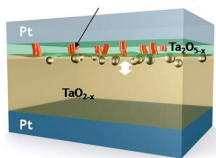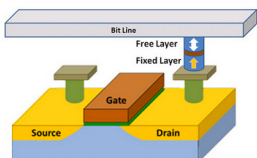    - Extrapolate to Exascale sytem:
    - Assume saving 80% of 32 PB of memory and a JMTTI of 10 hours
      $\Rightarrow$ requires checkpoint bandwidth of ~150 TB/s   (doable with distributed Flash)
    - Primary resiliency issue is dealing with *undetected* errors…

- **Storage latencies dropping faster than compute increasing**
    - Flash O(100) faster than  disk
    - NVRAM is O(100) faster than Flash

- **But there's lots of work to do on storage *architecture*..**
    - Reducing software overheads for Flash and NVRAM timescale
    - Metadata scaling and resiliency (relax Posix consistency?)
    - Namespace flexibility
    - Support for non-POSIX file systems (KVS, NoSQL, Spark RDDs, etc.)

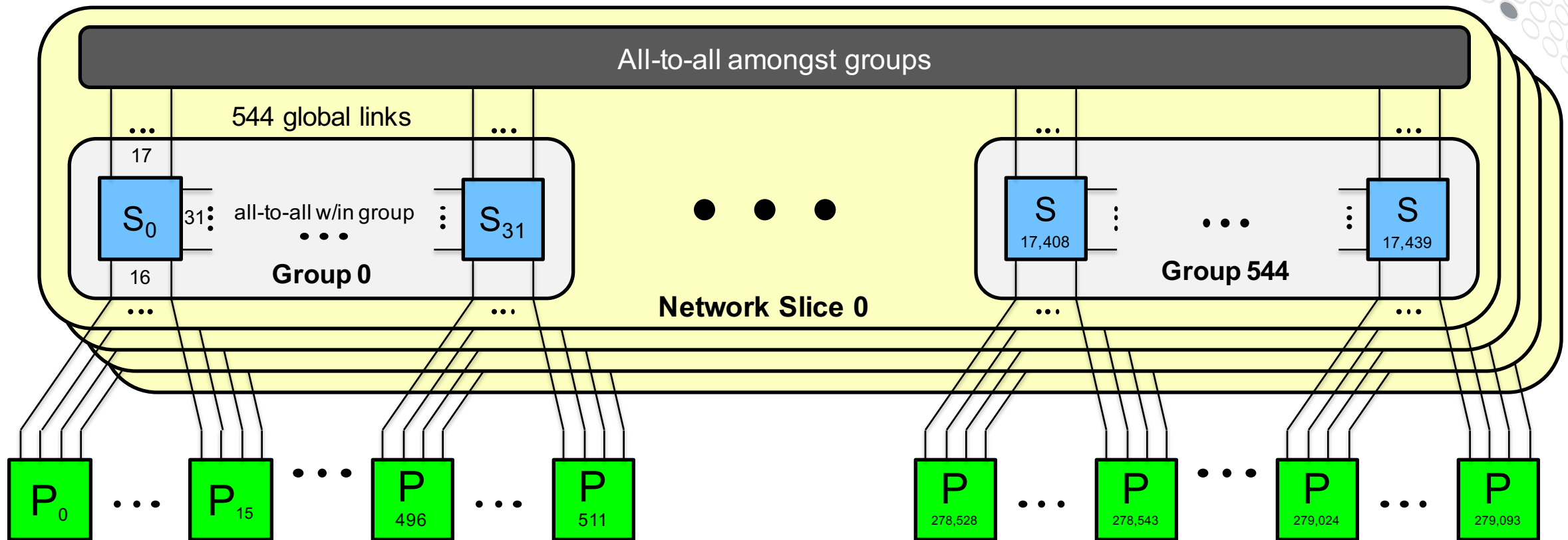# Cost- and Power-Efficient Networks

- **Cray pioneered the use of high radix routers in HPC**
  - Became optimal due to **technology shift**
    - Faster signaling permits narrower links
  - Reduced network diameter (number of hops)
    - $\Rightarrow$ Lower latency and cost
  - But… higher radix network require longer cable lengths

- **Optics enables longer cable lengths**
  - Now cost-effective above a few meters (and dropping)
  - Cost, bandwidth and power are insensitive to cable length

- **Future systems will based on hybrid, electrical-optical networks**
  - Cost-effective, scalable global bandwidth
  - Very low network diameter (small number of hops) $\Rightarrow$ very energy efficient

**First 64 port router
Cray X2 (2005)**

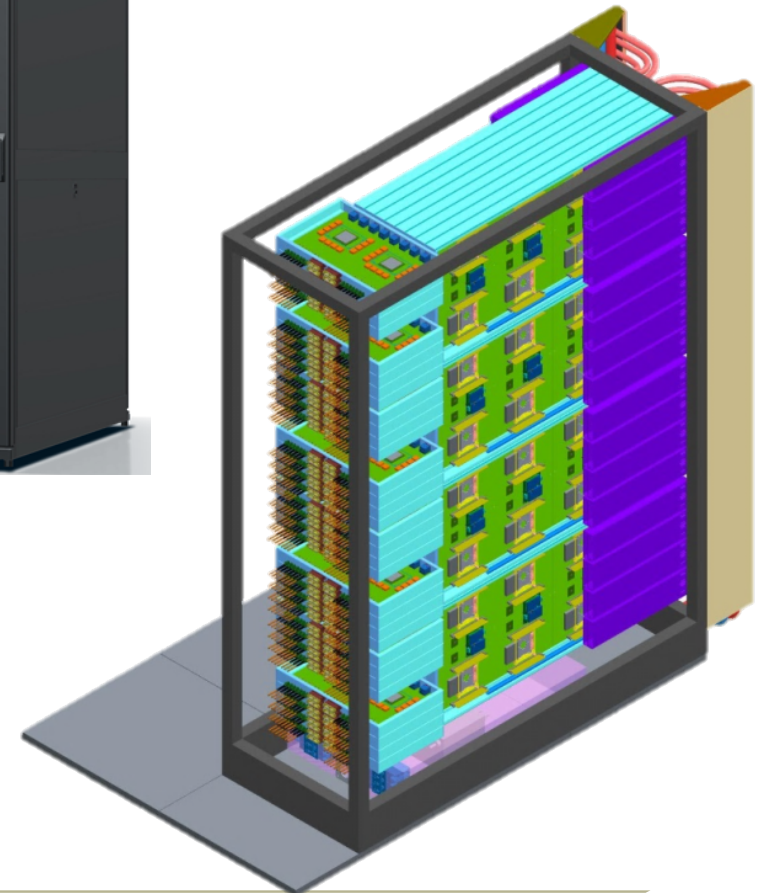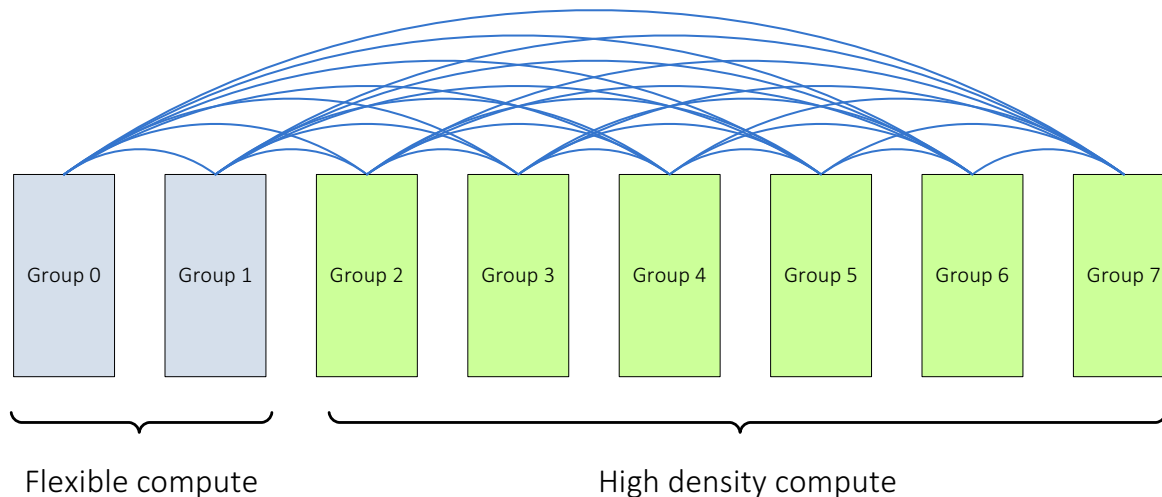# Example Dragonfly Network with a 64-port Switch



- Scales to 279K endpoints, with a network diameter of 3 hops!
- Only a *single* hop over a long (optical) link
- Narrow links allow sliced network for configurable bandwidth

# Next-Gen Shasta System Infrastructure

- **Single system with choice of:**
  - Cabinet type and cooling infrastructure
  - Processor type
  - Software stack
  - Interconnect
- **Extensible to Exascale and Beyond**
  - Power & cooling headroom
  - Network and processor configurability

| Group 0 | Group 1 | Group 2 | Group 3 | Group 4 | Group 5 | Group 6 | Group 7 |

Flexible compute                High density compute

# Summary of Future Machines

- **Computers are not getting faster…  just wider**
  - O(EF) with O(GHz) clocks → O(B) way parallelism!

- *Vertical* **locality much more important than** *horizontal* **locality**

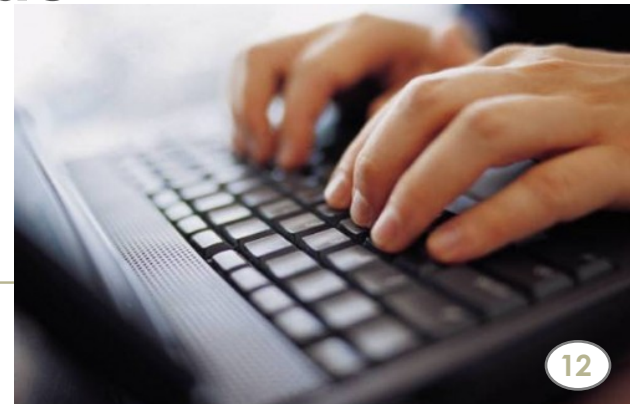| Dimension | Latency Hit | Bandwidth Hit | Energy Hit |
|---|---|---|---|
| Within node | ~200x | ~200x | > 500x |
| Across nodes | ~25x | ~8x | ~5x |

*\* If include local NVM, within node grows, across nodes shrinks*

- **Parallelism is multi-dimensional (and heterogeneous?)**
  - Vectorization + threading + multi-node
  - Processors optimized for serial performance *or* power efficiency  (not both)

- **Interconnects won't look that different than today**
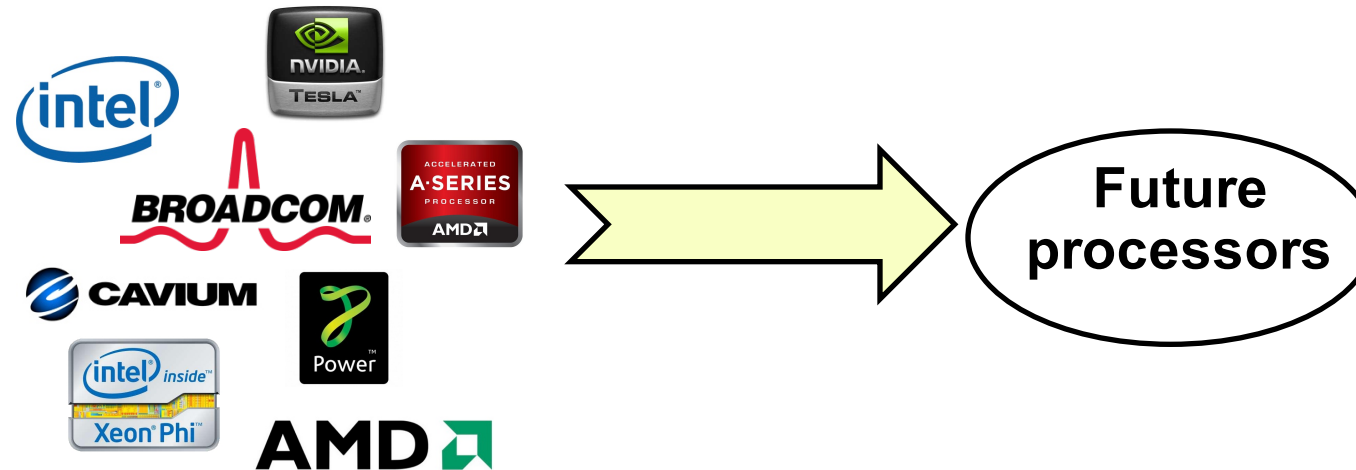
# Implications for Programmers

- **May need to move to more threading on the node**
  - All-MPI often won't deliver maximum performance
- **Must vectorize low-level loops**
  - 8-30x performance improvement on array operations
- **Must avoid serial scalar code**
  - Inherently slower and less power-efficient
  - On "accelerated" nodes, either
    - creates traffic between accelerator and host, or
    - runs 3-4x slower than on a serial-optimized core
- **Must pay a *lot* more attention to locality within node**
  - Think about data placement and movement
  - Consider "sub-optimal" algorithms that limit data motion

# Would like to code for future machines in a portable way
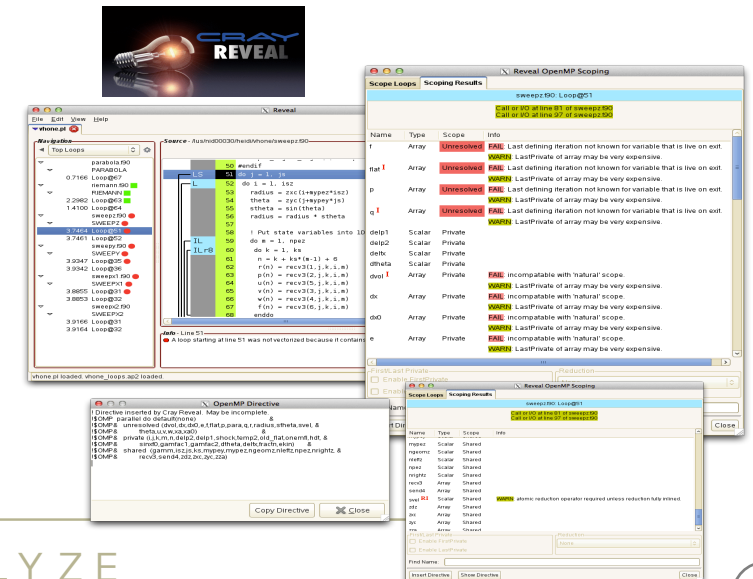
- **Spatial and Temporal Portability**



- **Separation of labor**
  - Programmer *exposes* parallelism and locality
  - Compiler, tools, and runtime map onto specific hardware
  - Optimized libraries for various platforms

# Bold Prediction:

- **Future HPC Programming Model:   MPI + OpenMP**

- **Can we make this easier?**
  - Threading, vectorization, data placement

- **Recent poll at NERSC found 80% of apps use single level of parallelism**

- **Why & when to convert to hybrid programming model?**
  - When code becomes network bound
  - Load balancing and synchronization overheads become large
  - Excessive memory used by straight MPI
  - To take advantage of hybrid compute nodes

- **Programming tools are going to be critical**
  - Exposing parallelism (especially higher in call chain)
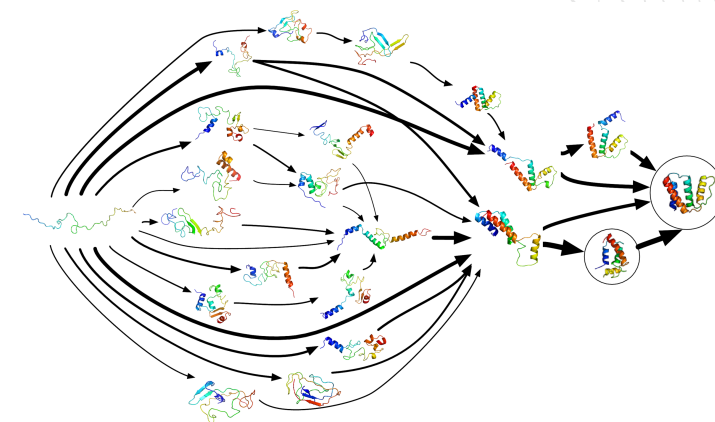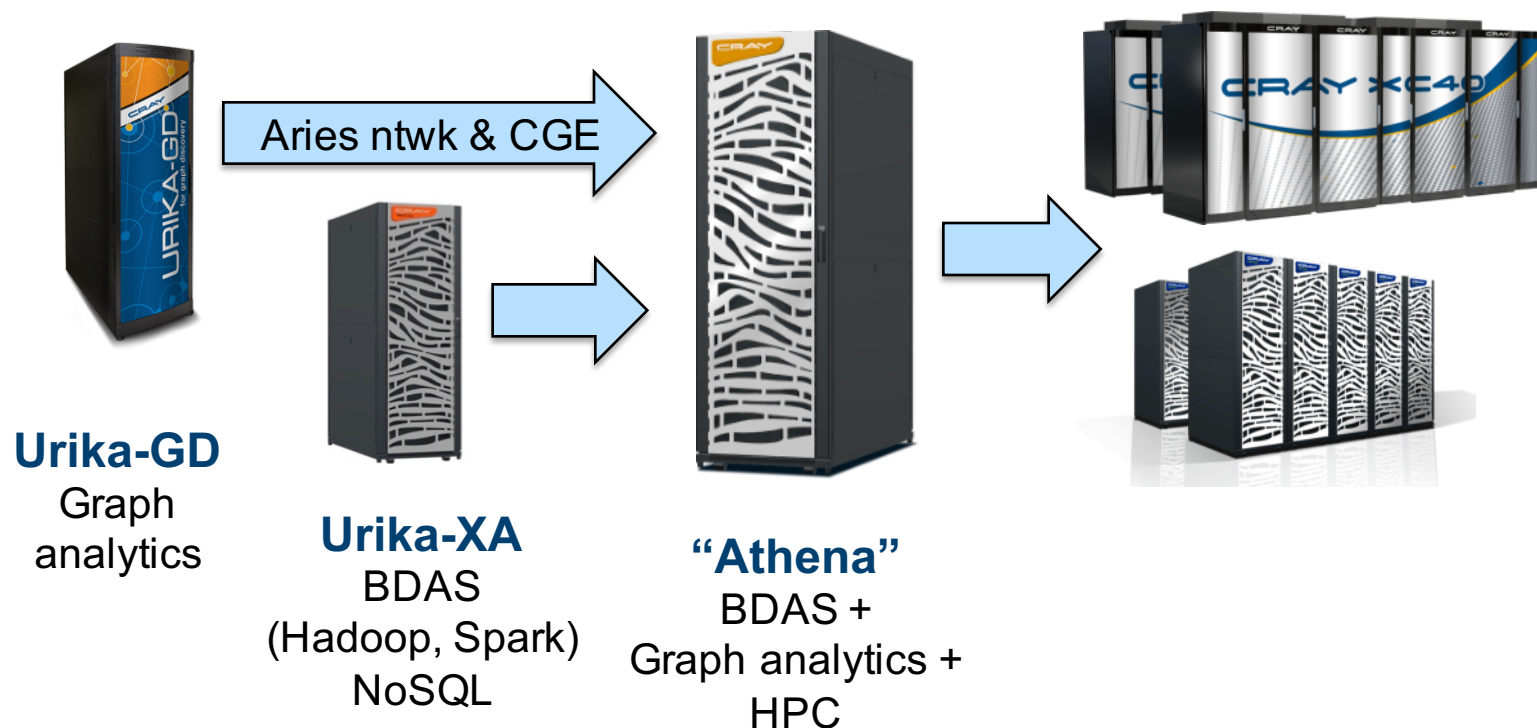  - Data placement and movement in the memory hierarchy

# Beyond Classic HPC

# Merging of HPC and Data Analytics



Aries ntwk & CGE

**Urika-GD**
Graph analytics

**Urika-XA**
BDAS
(Hadoop, Spark)
NoSQL

**"Athena"**
BDAS +
Graph analytics +
HPC

HPC + Analytics workflows

*Why combine HPC and Analytics solutions in a single box?*

HPC underneath the covers

# *Thank You*

# *Questions?*