

# Tuning sparse matrix kernels for multicore

*U.C. Berkeley / LBNL*

→ Sam Williams, Lenny Oliker, John Shalf, James Demmel, Katherine Yelick

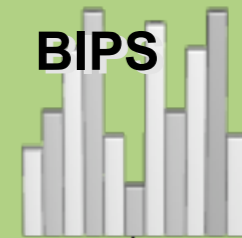
*Georgia Tech*

→ Richard Vuduc

CSE/CS 8803 PNA: Parallel Numerical Algorithms

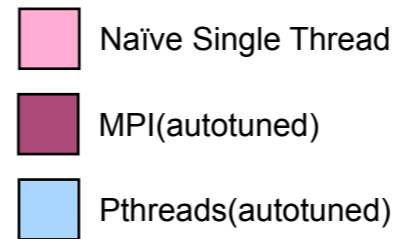
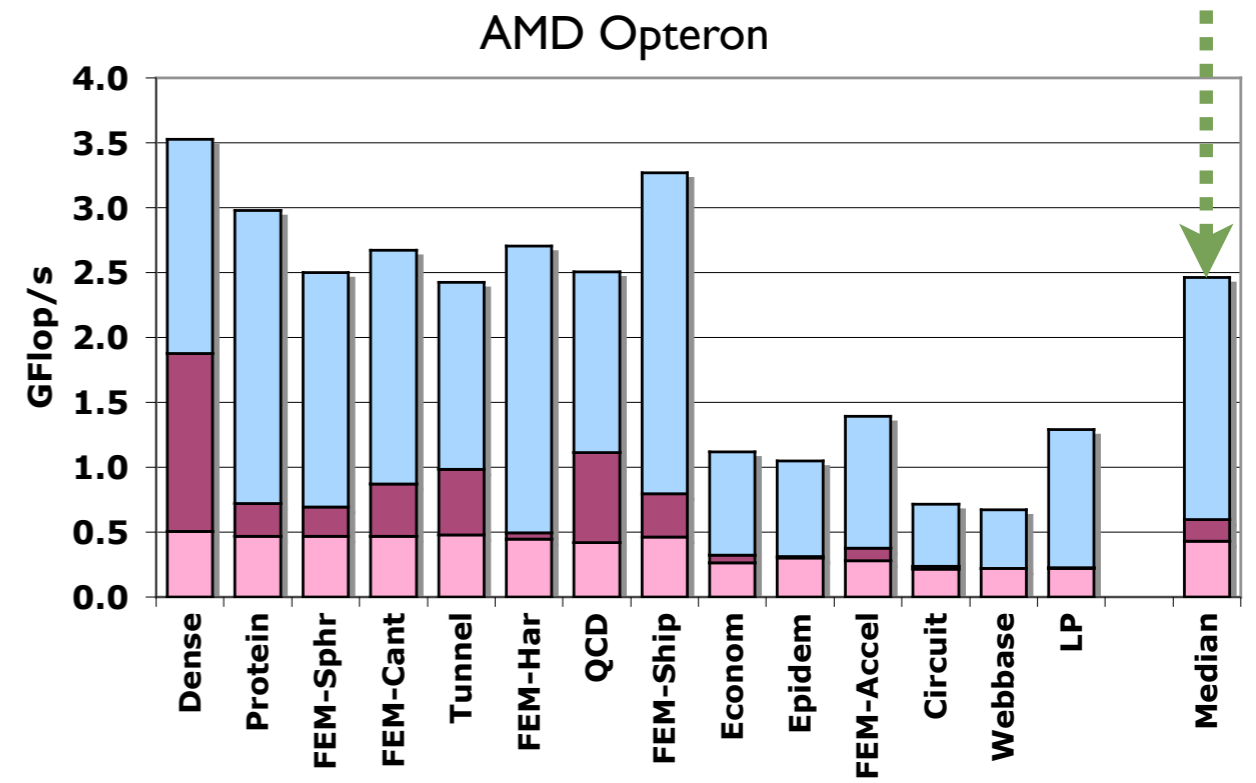
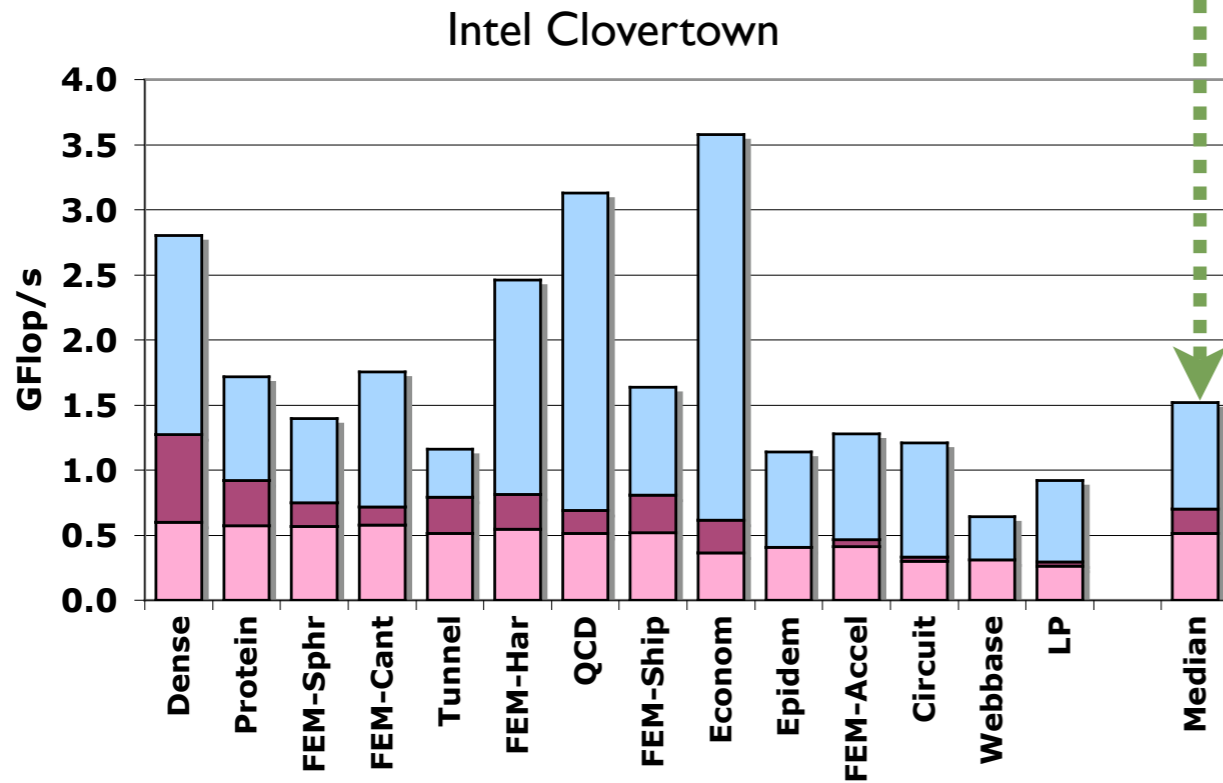
[L.25] Tuesday, April 15, 2008

**bebop**  
Berkeley Benchmarking and OPTimization Group  
bebop.cs.berkeley.edu



~2x on 2x4 cores

~4x on 2x2 cores



Multicore-specific  
vs. "Off-the-shelf"

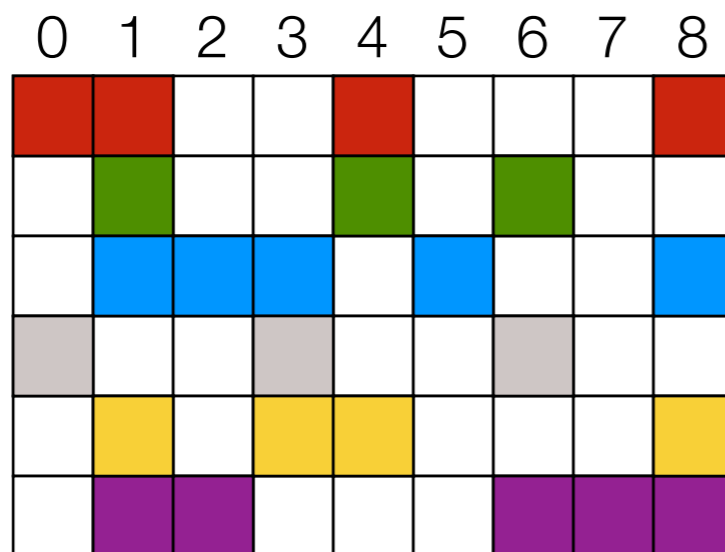
Sparse matrix-vector multiply (SpMV)  
PETSc 2.3.0+MPICH 1.2.7+OSKI 1.0.1h

# Preview

- ▶ Comparisons
  - ▶ Intel 4-core, AMD 2-core, Sun Niagara2, STI Cell
  - ▶ Multicore-specific vs. “Off-the-shelf” MPI
- ▶ Key conclusions
  - ▶ Machine with the most complex architecture needs most tuning & achieves lowest performance
  - ▶ 50–60% DRAM bandwidth vs. 90% on Cell
  - ▶ Place, Prefetch, Compress, Block—and tune!

- ▶ **The need for automatic tuning  
(A 20-year retrospective)**

# Compressed sparse row (CSR) format



value



index



row\_pointer



# SpMV performance: Rules of thumb

value



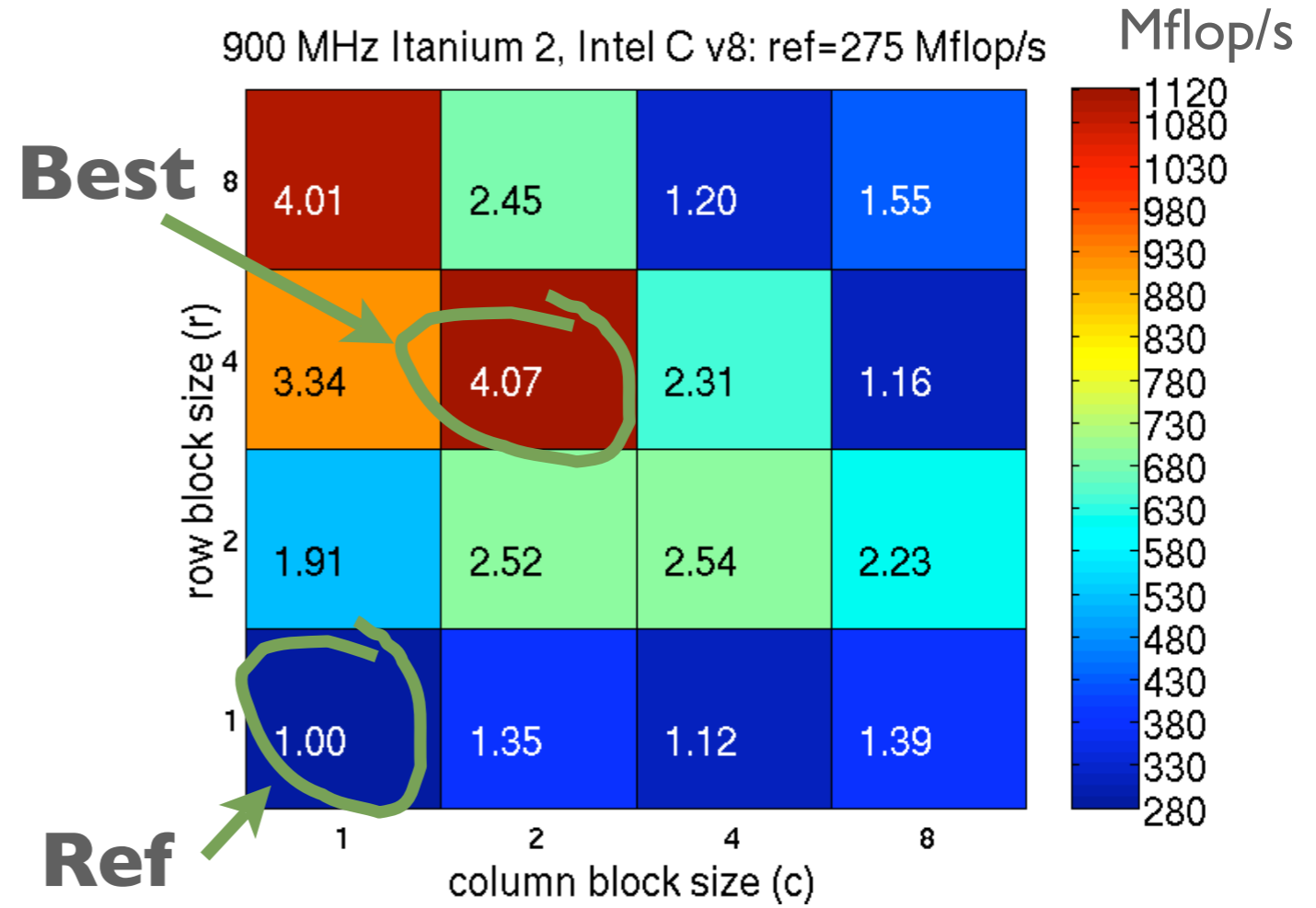
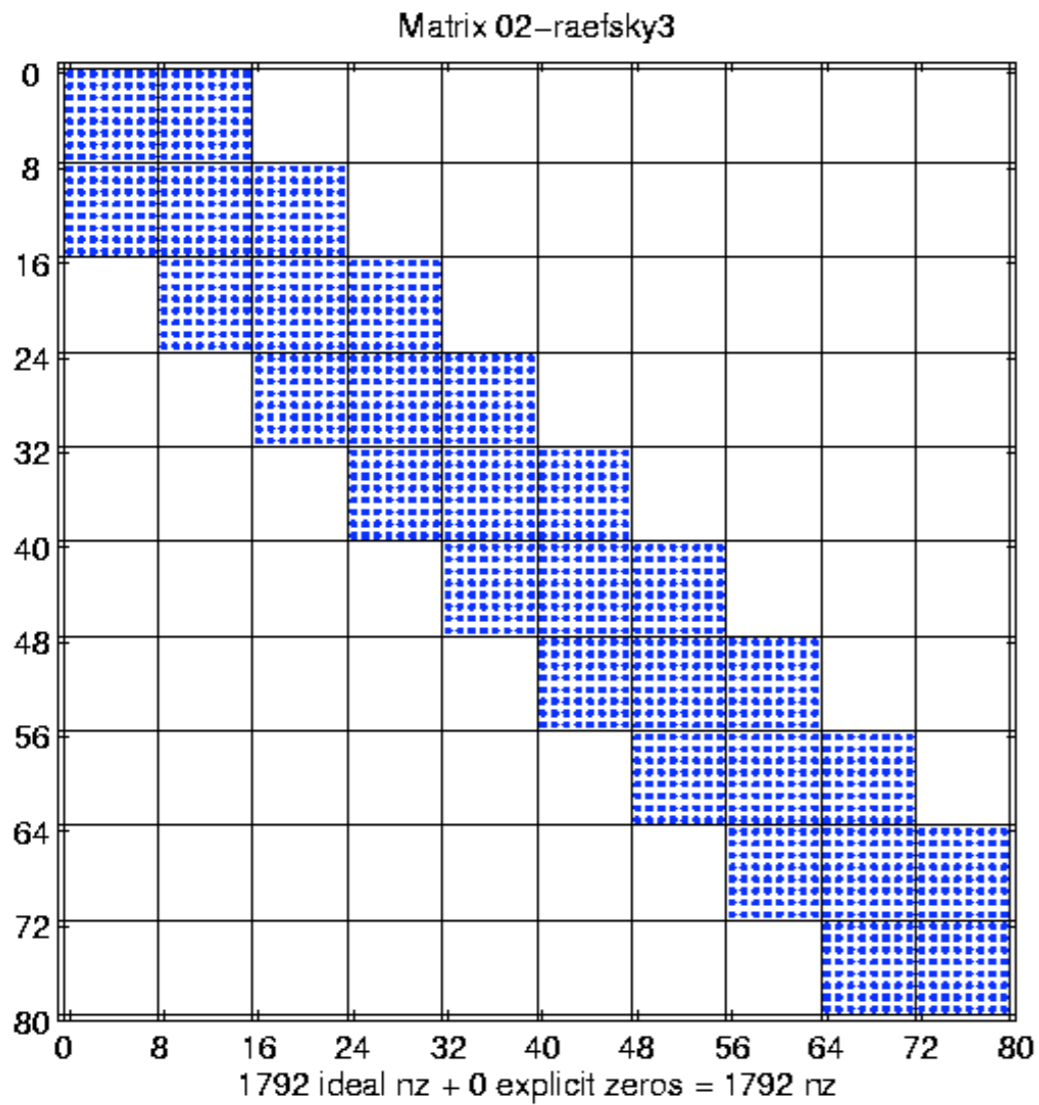
index



row\_pointer

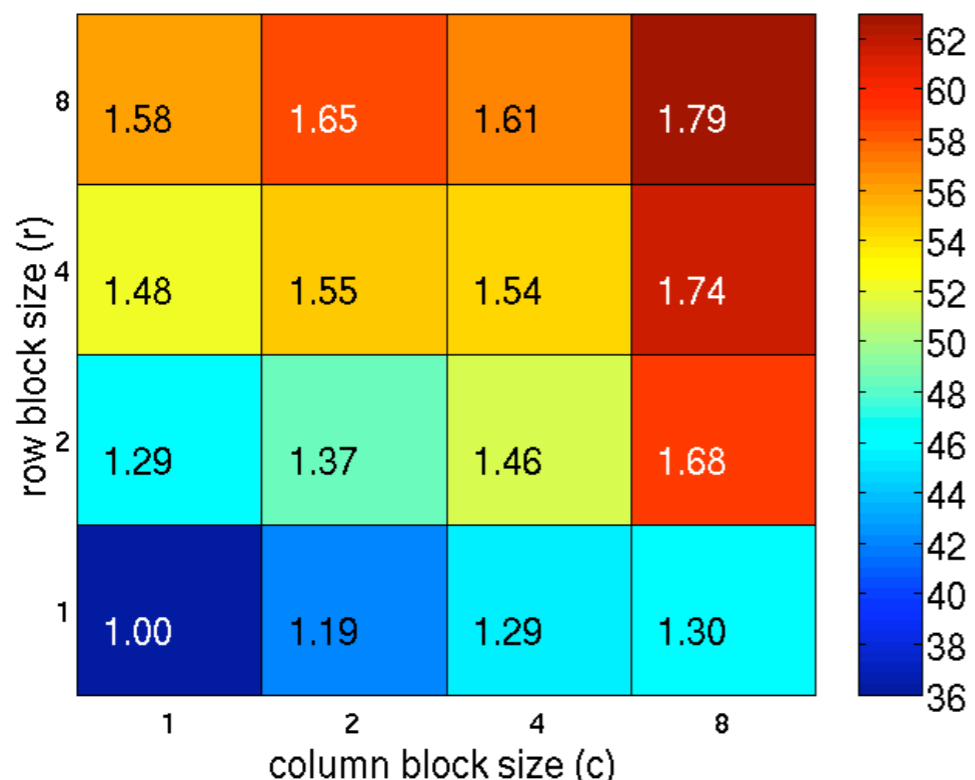


- ▶ Low computational intensity, vs. dense case
- ▶ Serial performance ~ 10% peak or less
- ▶ Bandwidth limited → compress
- ▶ Expect 1.5x speedup max, 32b ints + 64b vals

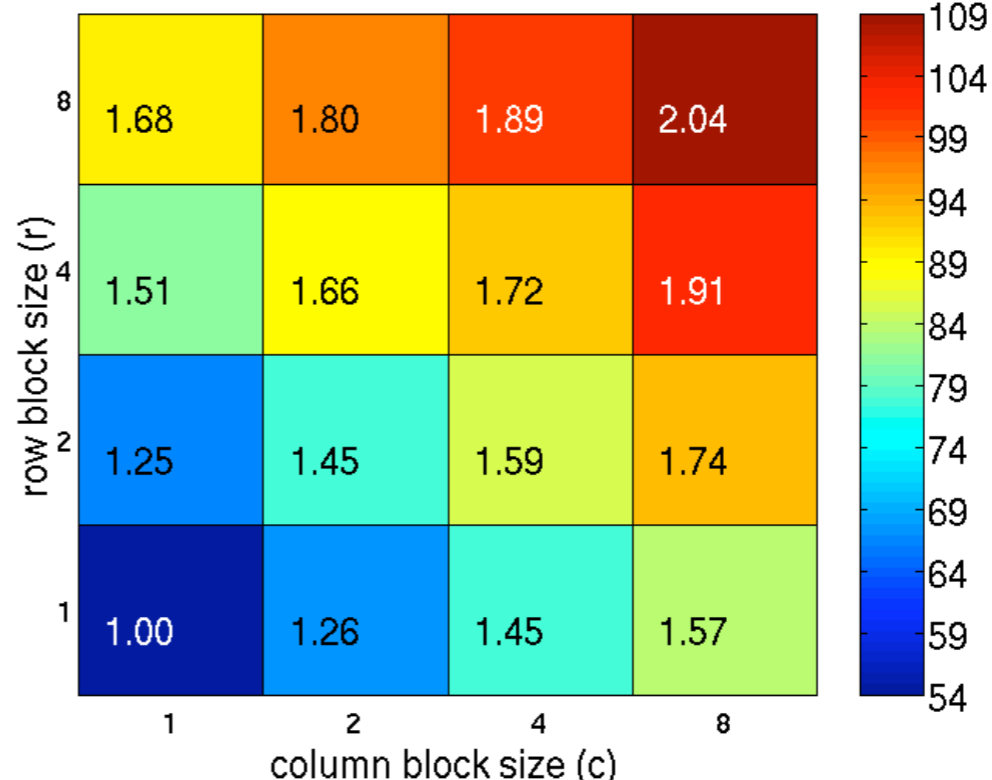


Optimal block size not obvious

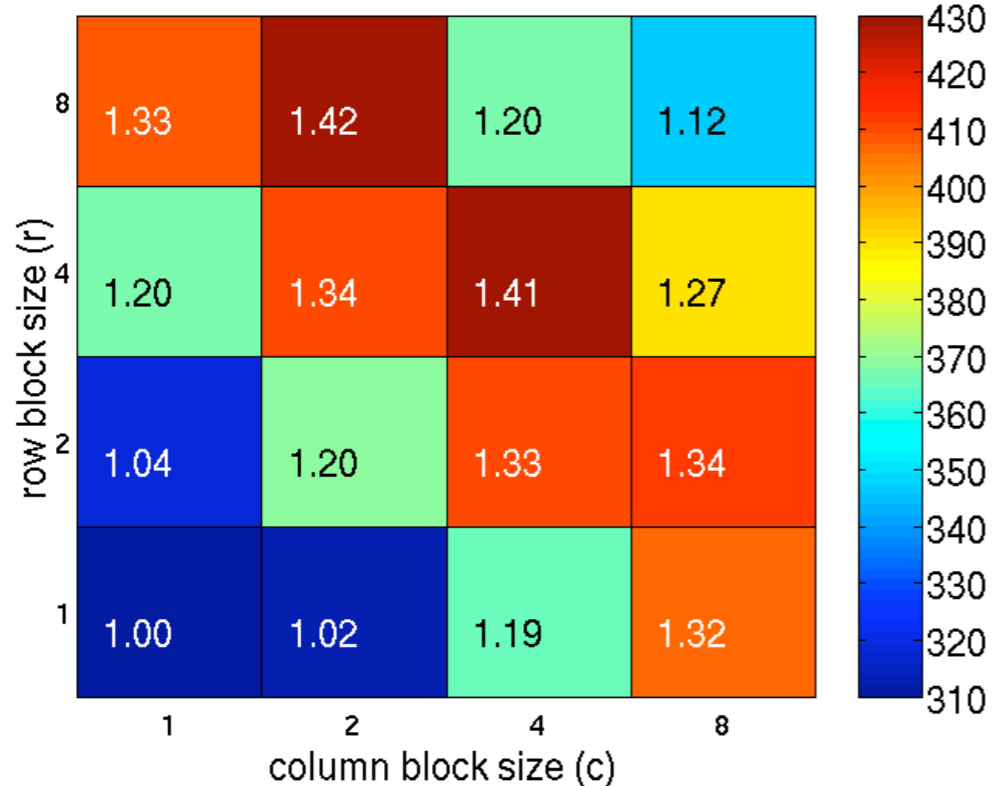
333 MHz Sun Ultra 2i, Sun C v6.0: ref=35 Mflop/s



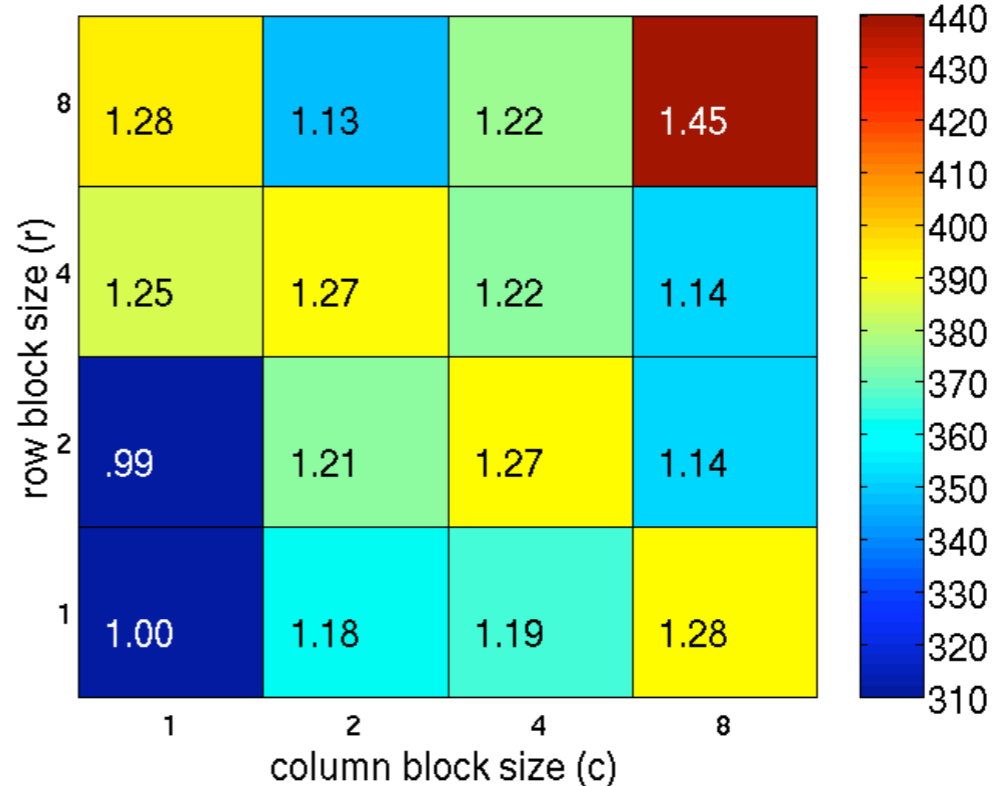
900 MHz Ultra 3, Sun CC v6: ref=54 Mflop/s



2 GHz Pentium M, Intel C v8.1: ref=308 Mflop/s

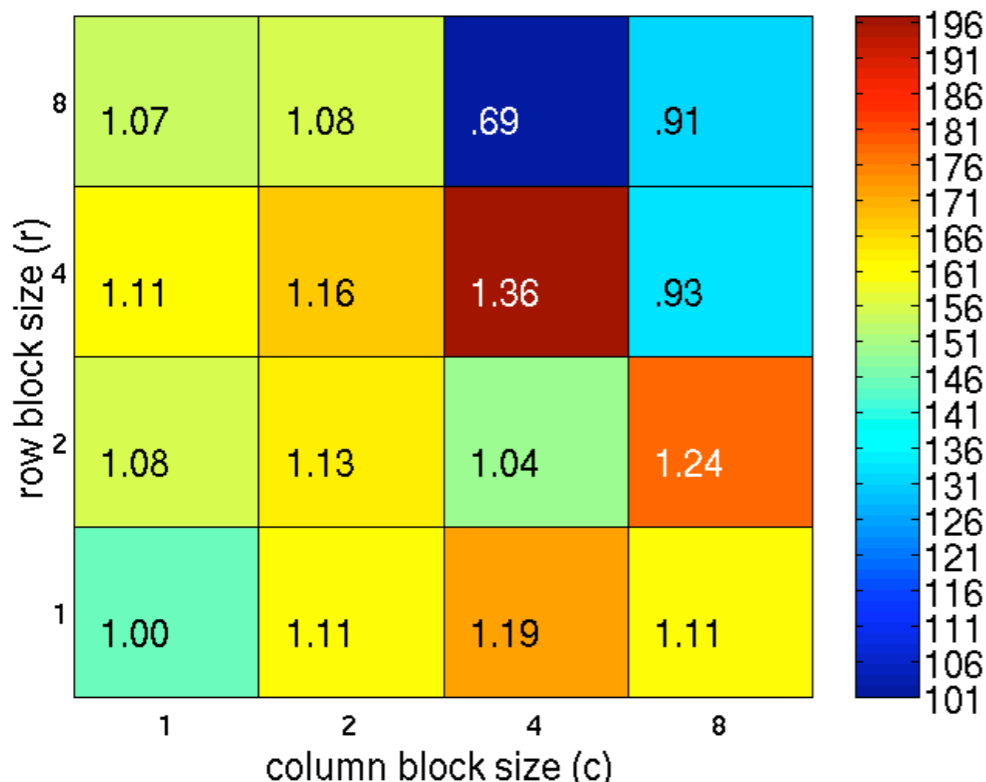


1.4 GHz Opteron, gcc 3.4.2: ref=308 Mflop/s

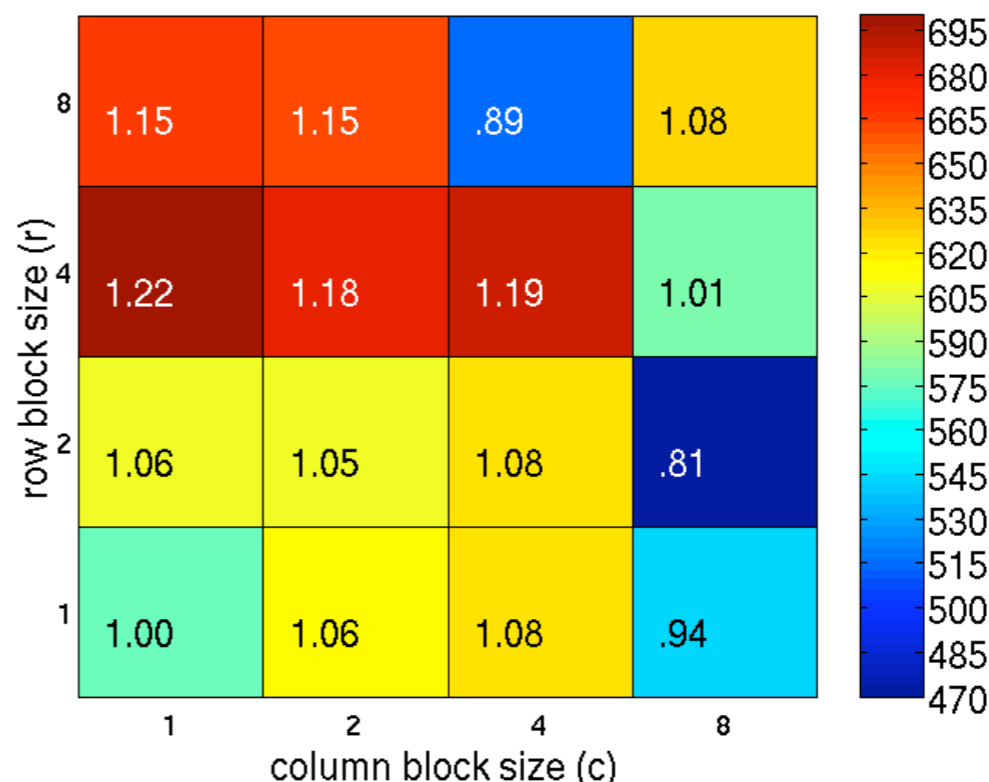




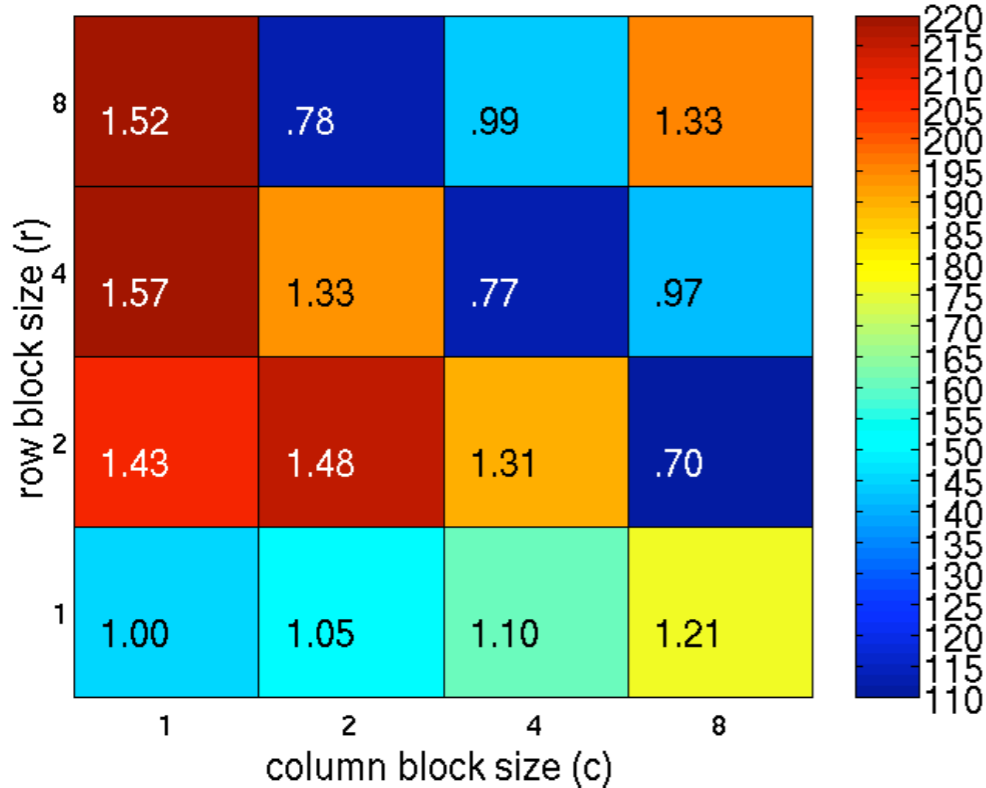
375 MHz Power3, IBM xlc v6: ref=145 Mflop/s



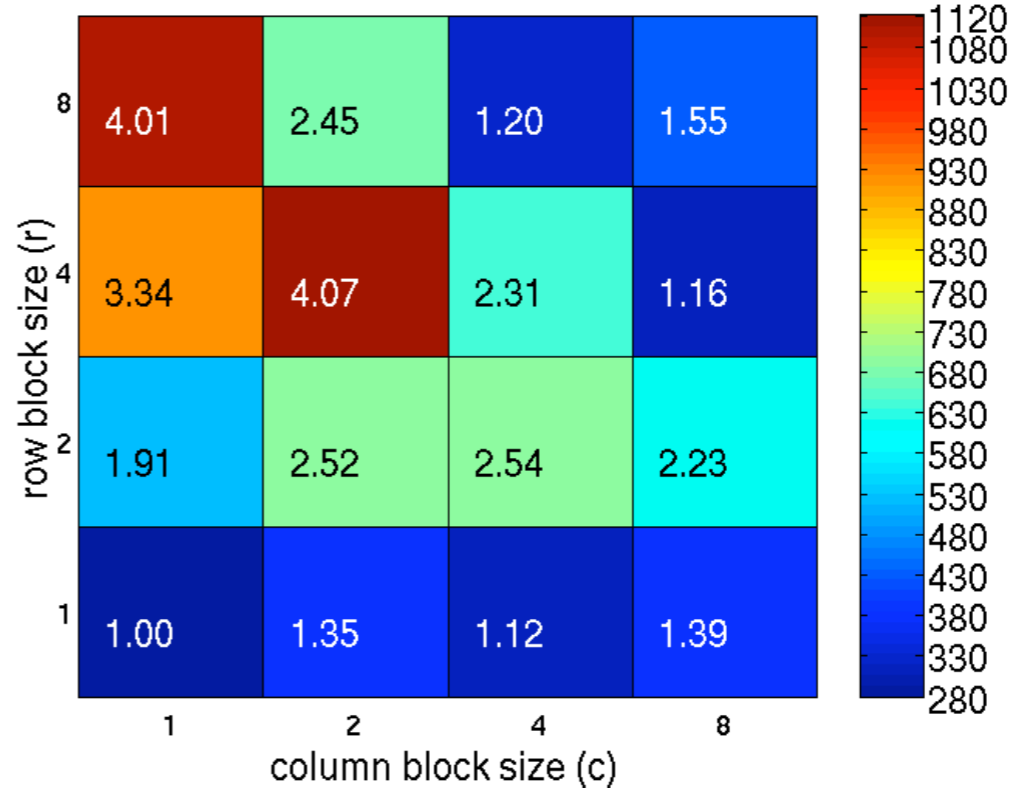
1.3 GHz Power4, IBM xlc v6: ref=577 Mflop/s



800 MHz Itanium, Intel C v7: ref=146 Mflop/s



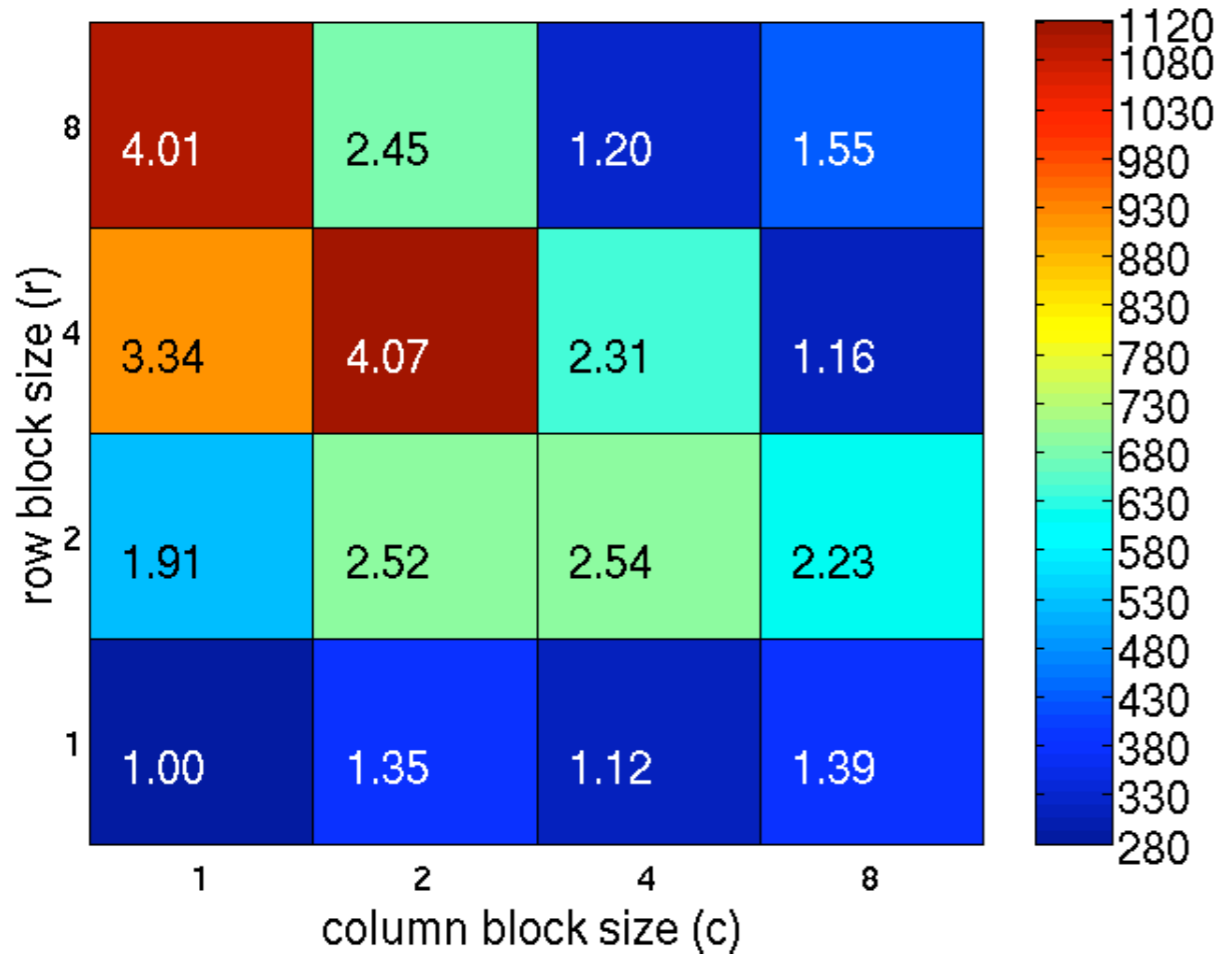
900 MHz Itanium 2, Intel C v8: ref=275 Mflop/s



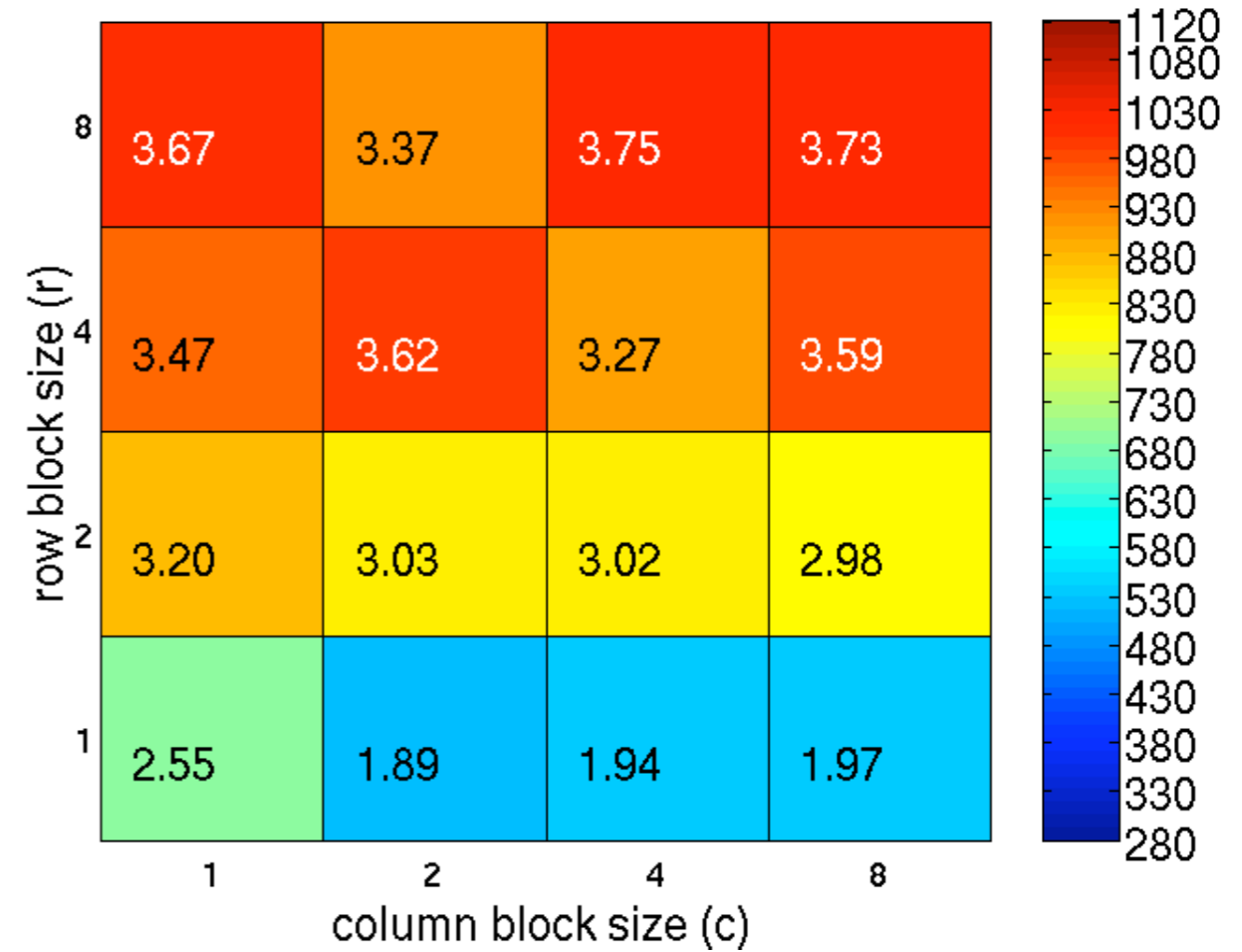
Better, worse,  
or about the  
same?



900 MHz Itanium 2, Intel C v8: ref=274 Mflop/s

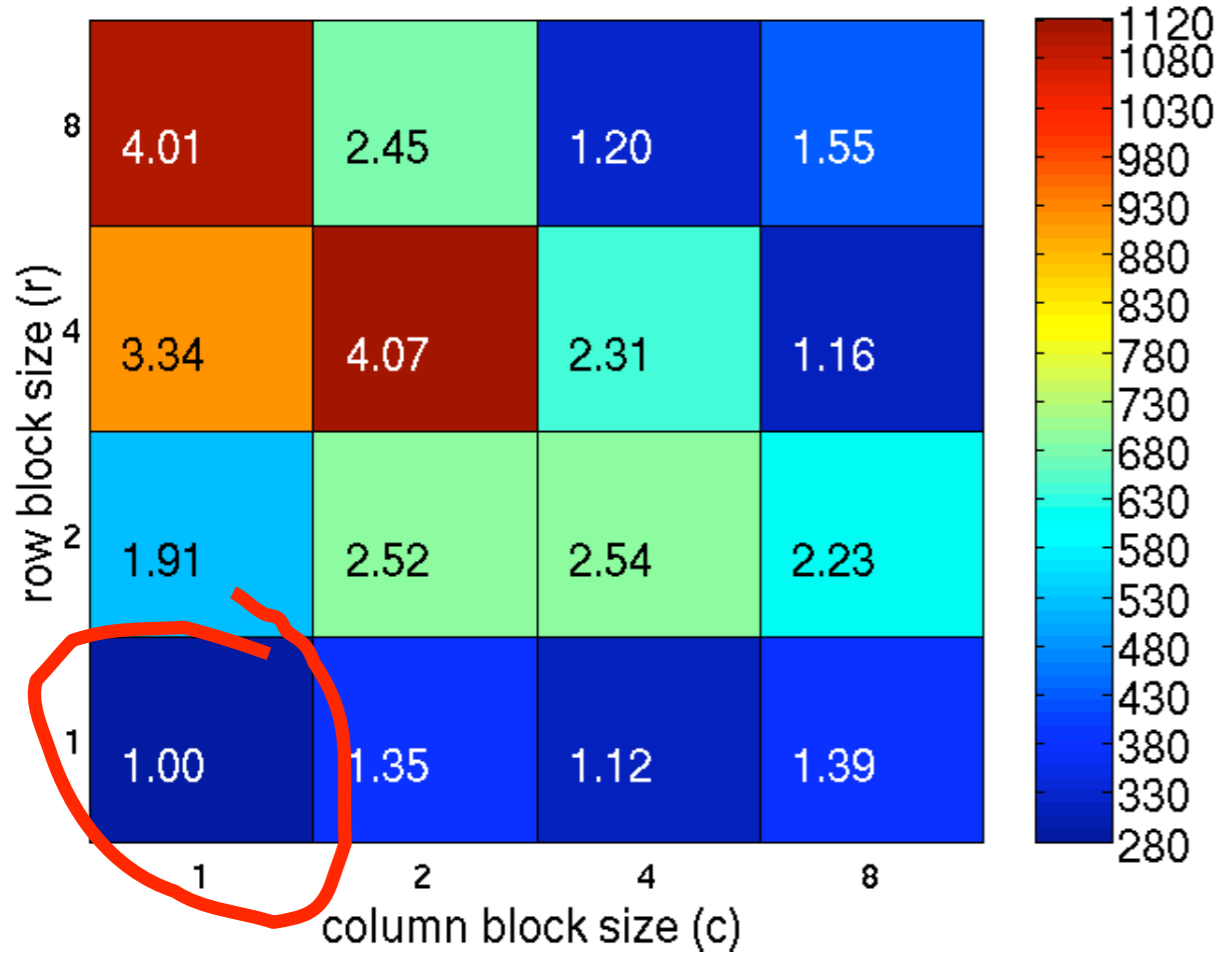


1.3 GHz Itanium 2, Intel C v8: ref=699 Mflop/s

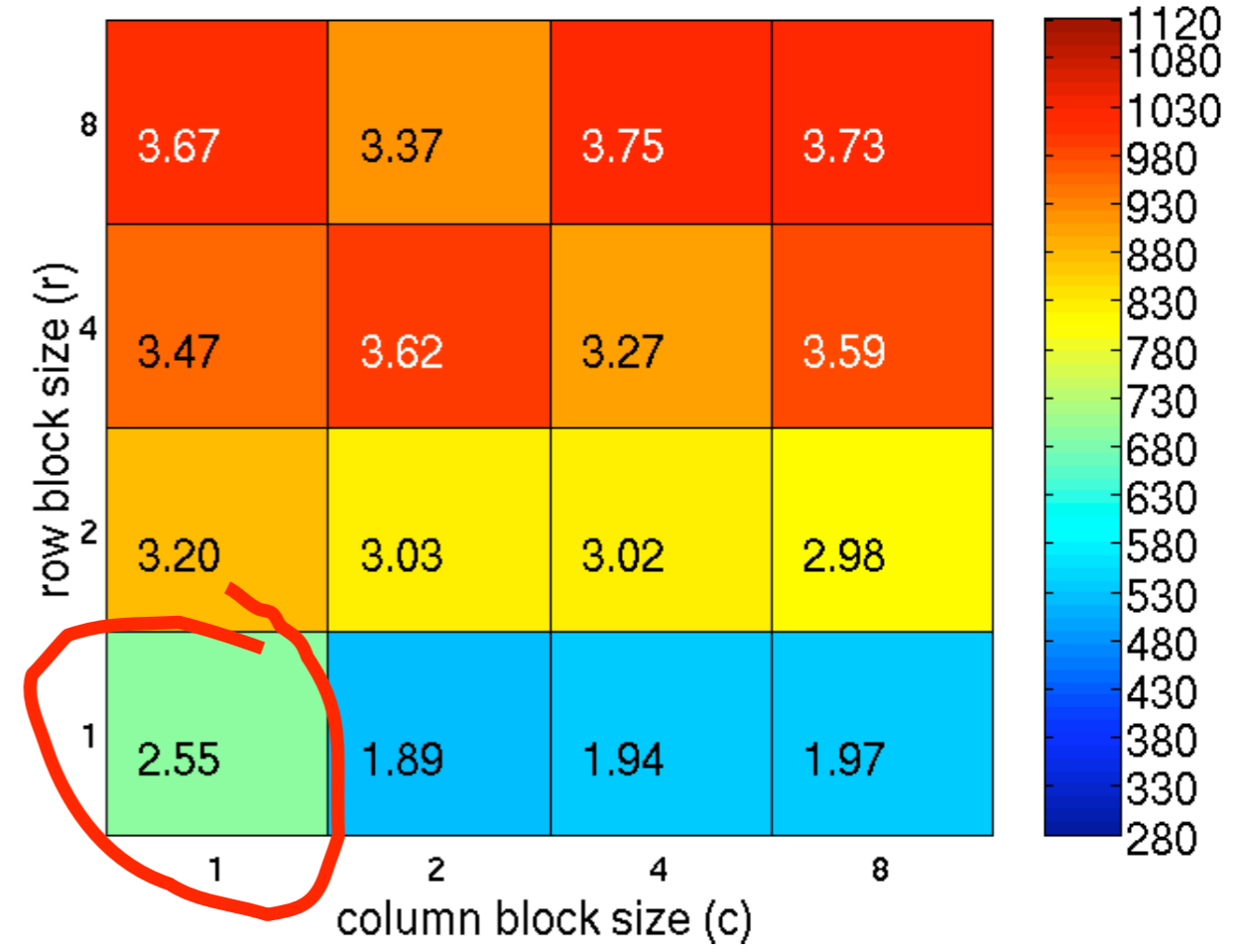


Example 1: Processor upgrade.

900 MHz Itanium 2, Intel C v8: ref=274 Mflop/s

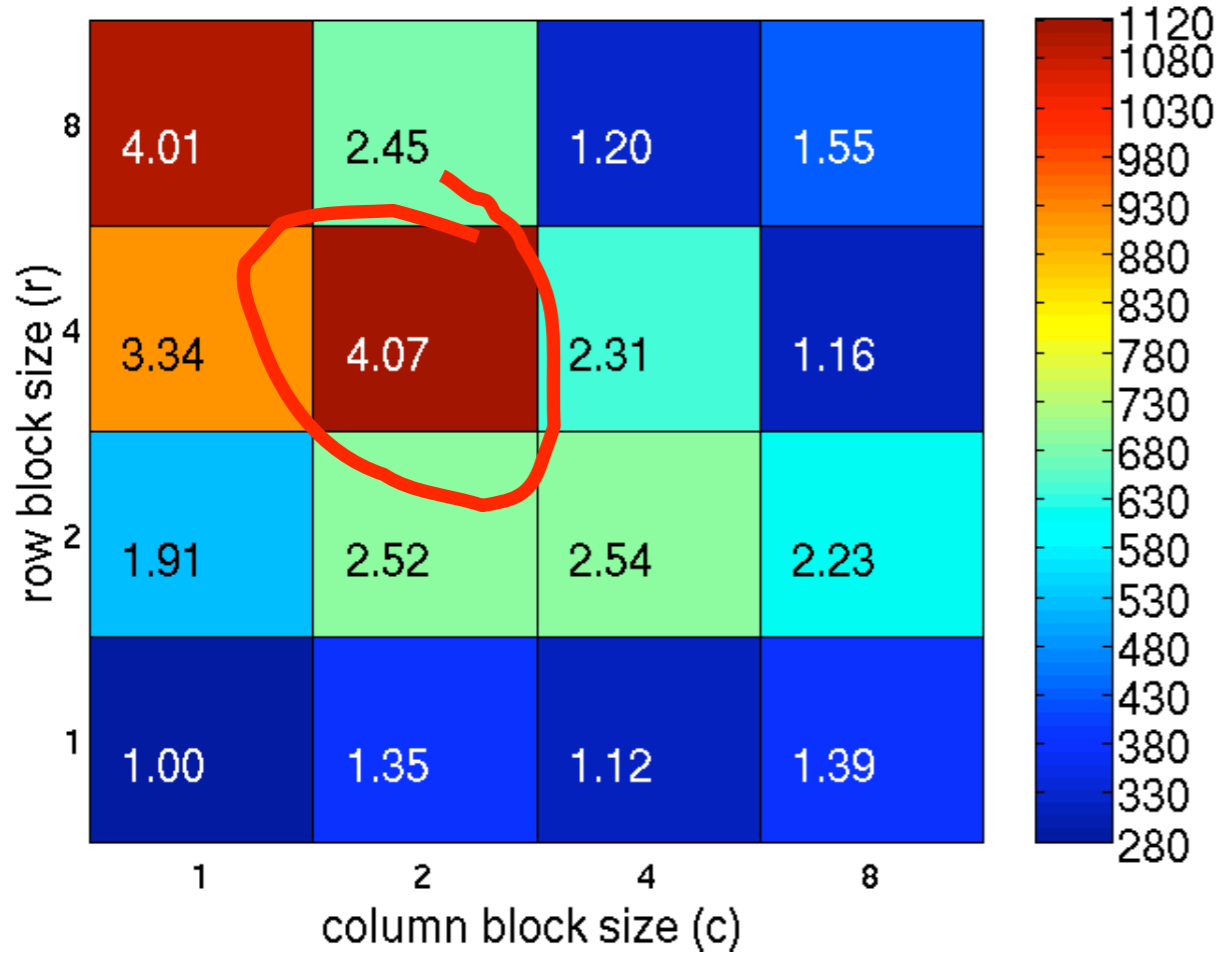


1.3 GHz Itanium 2, Intel C v8: ref=699 Mflop/s

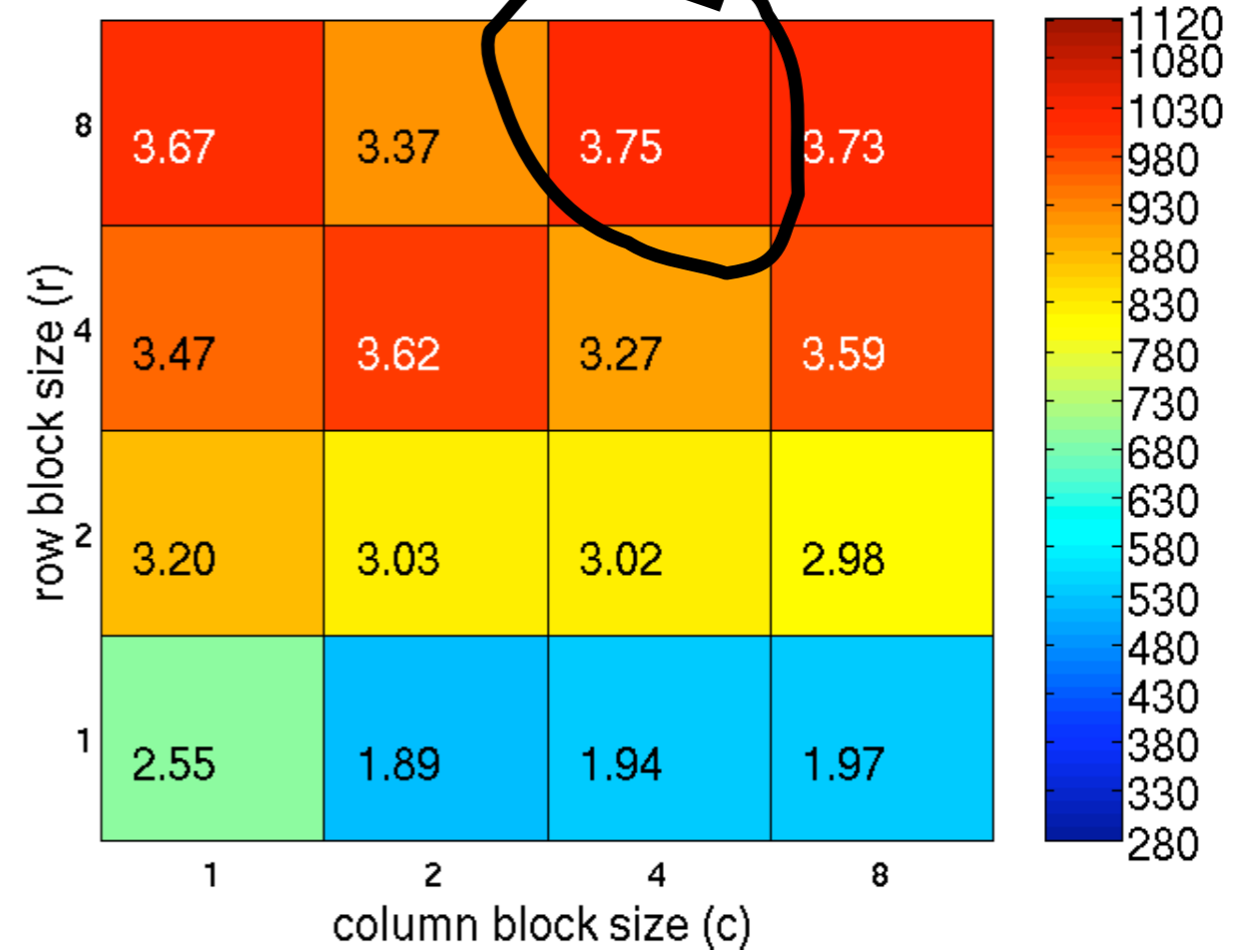


Reference improves...

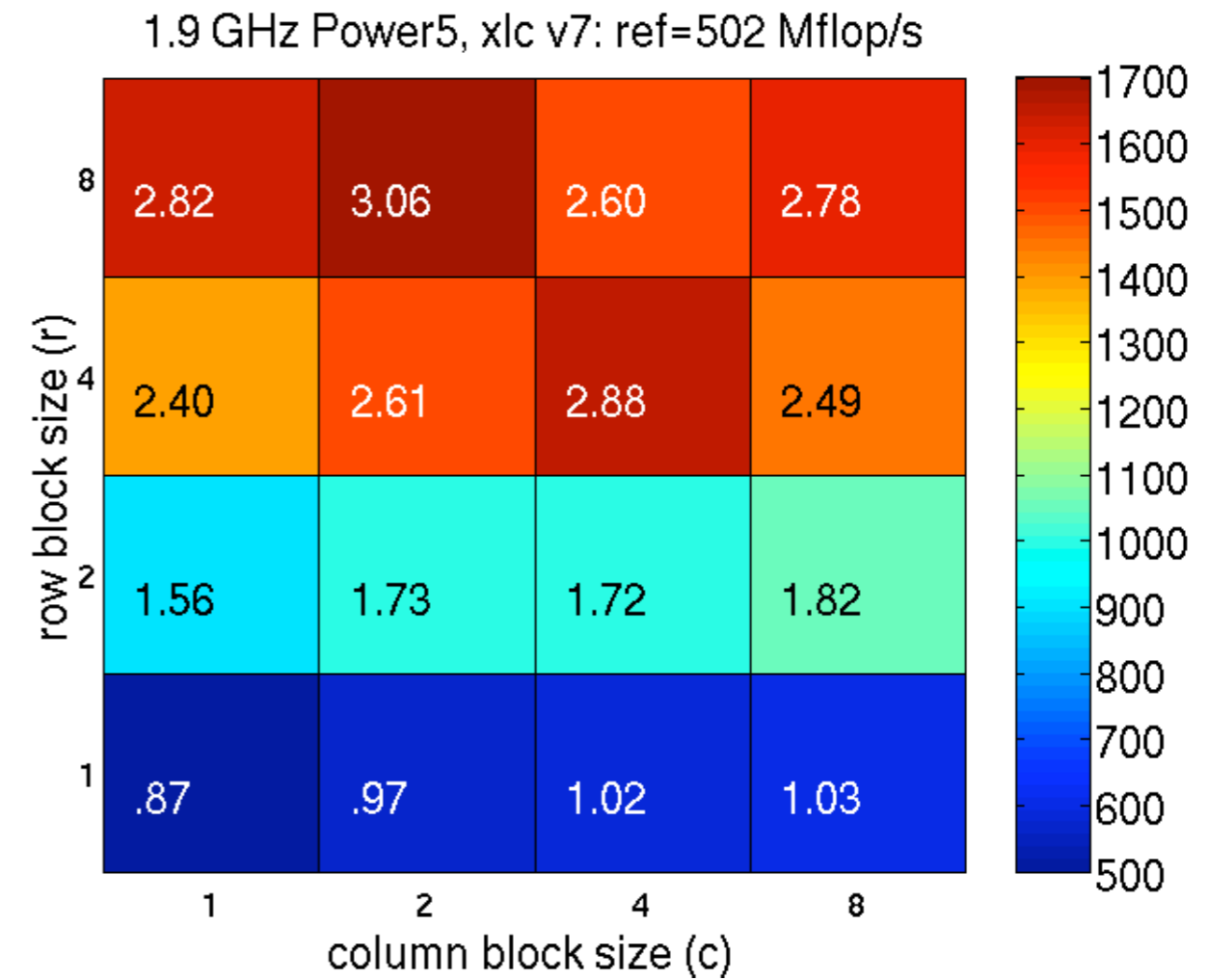
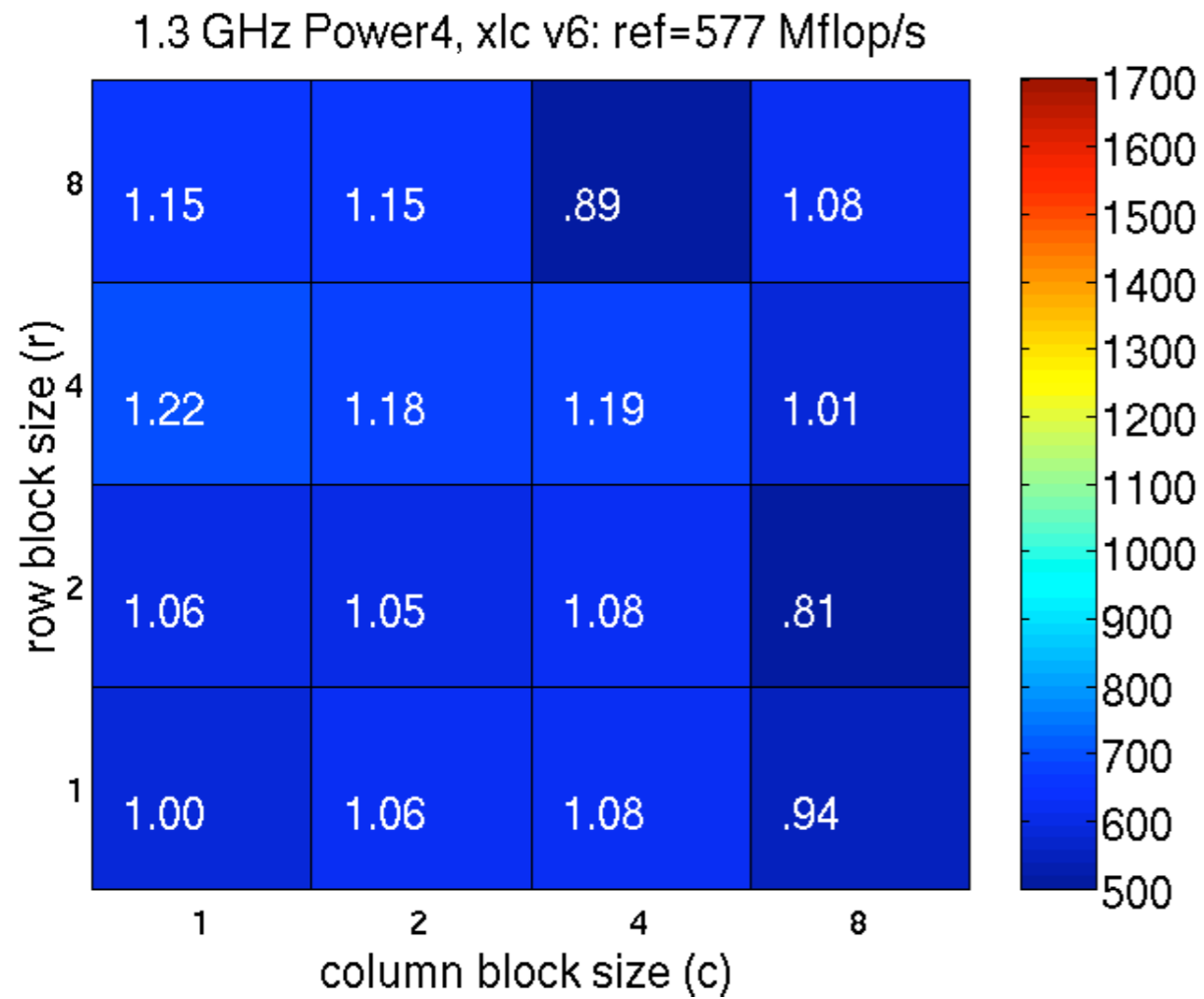
900 MHz Itanium 2, Intel C v8: ref=274 Mflop/s



1.3 GHz Itanium 2, Intel C v8: ref=699 Mflop/s

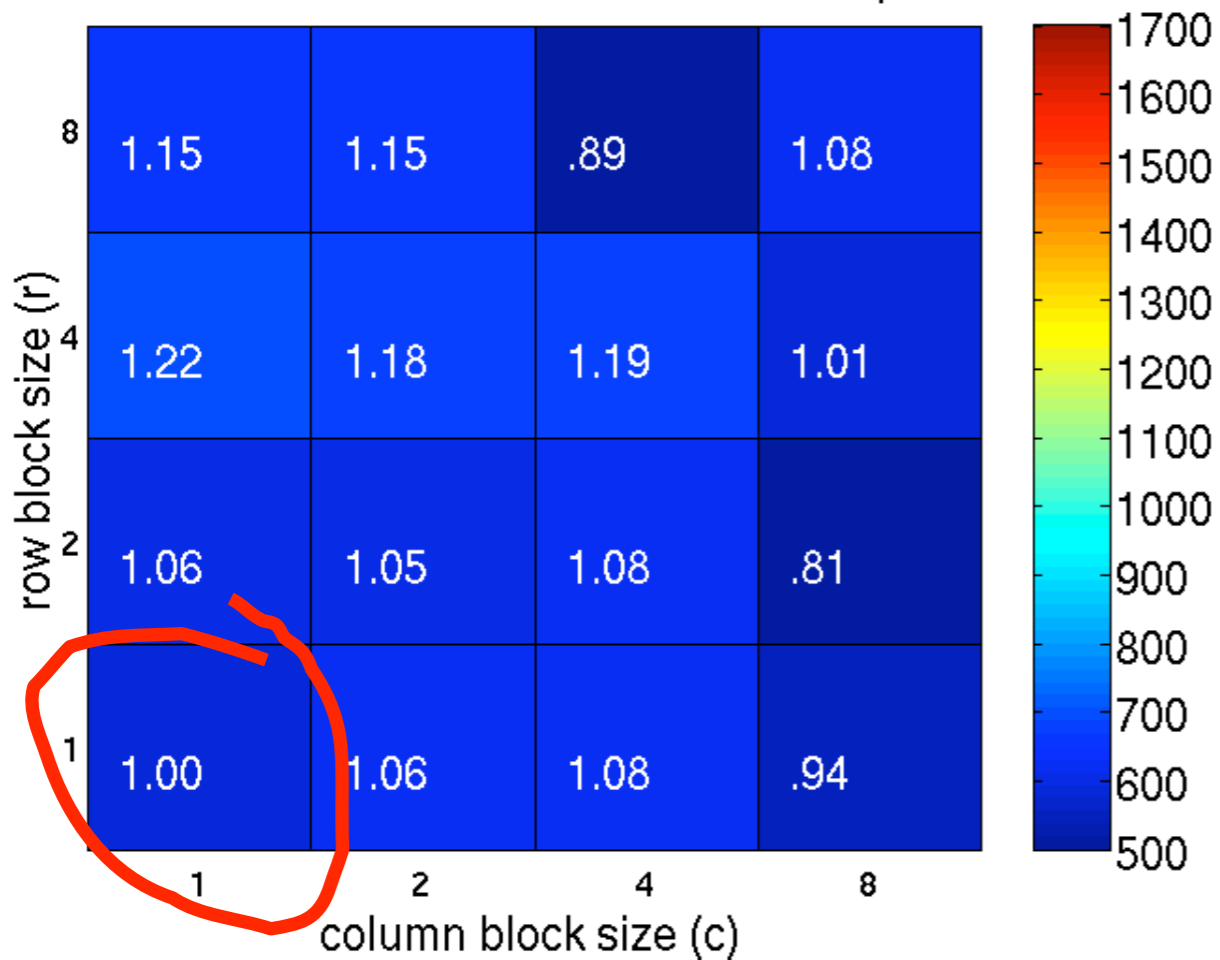


... and best possible worsens?

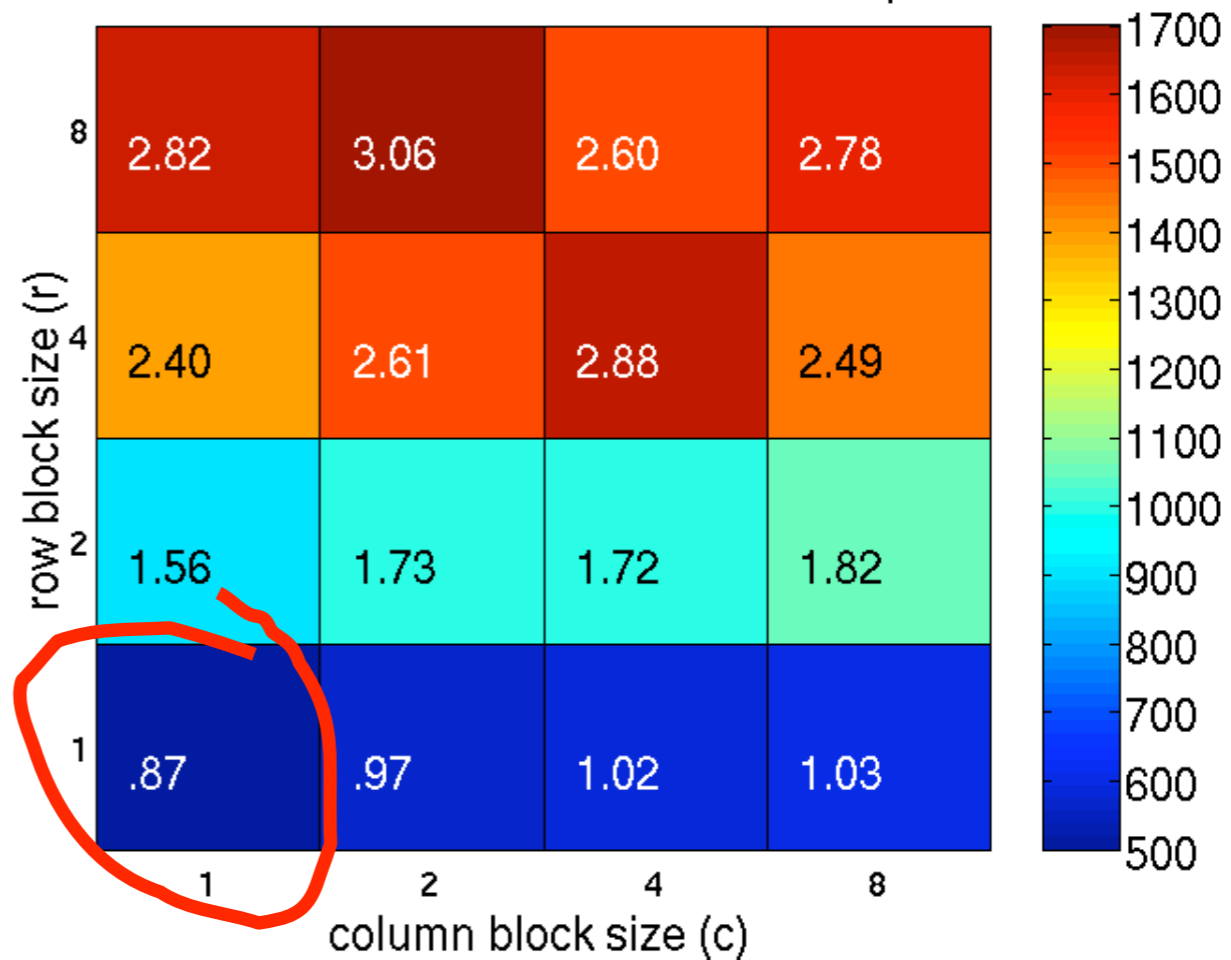


Example 2: Power4  $\Rightarrow$  Power5.

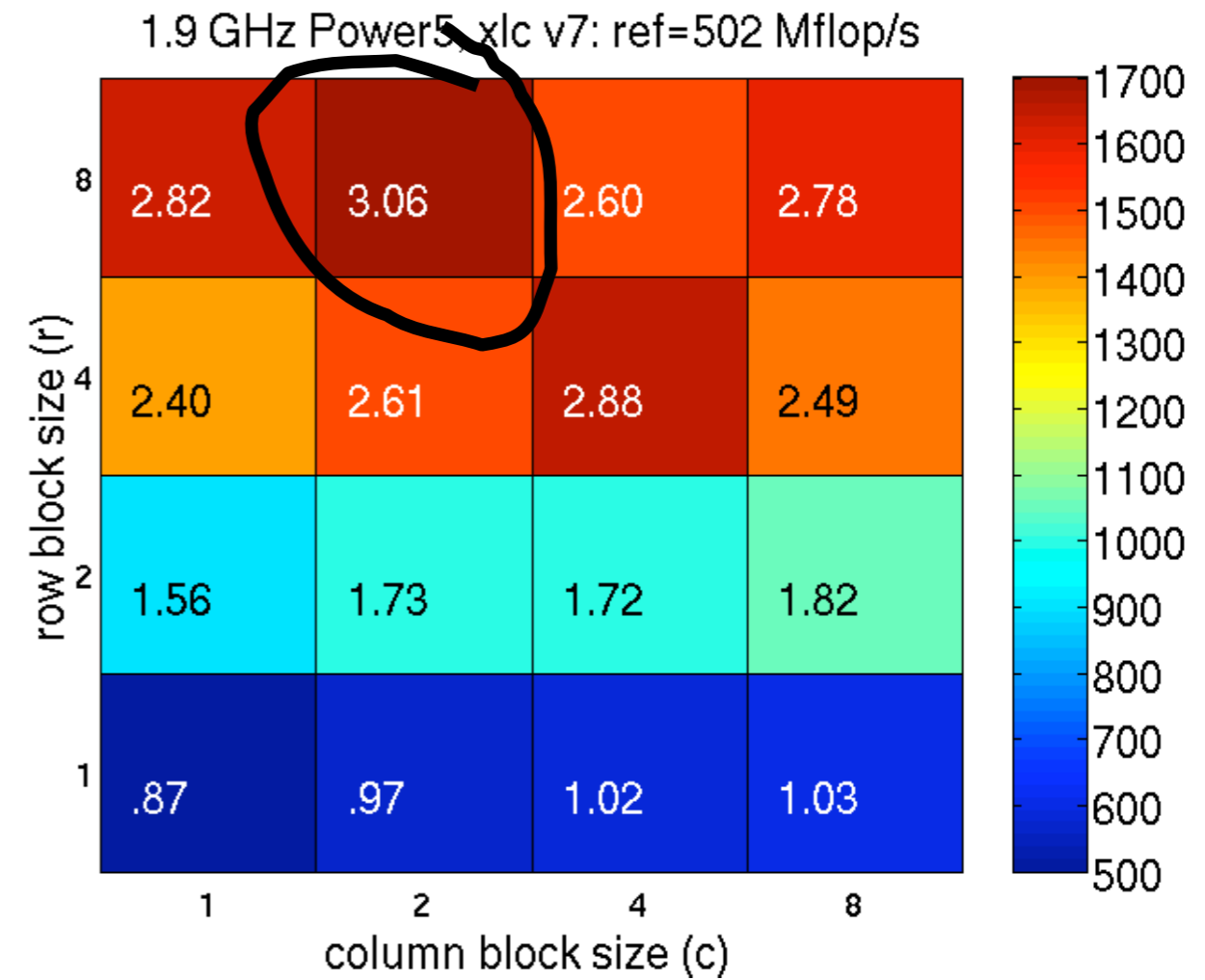
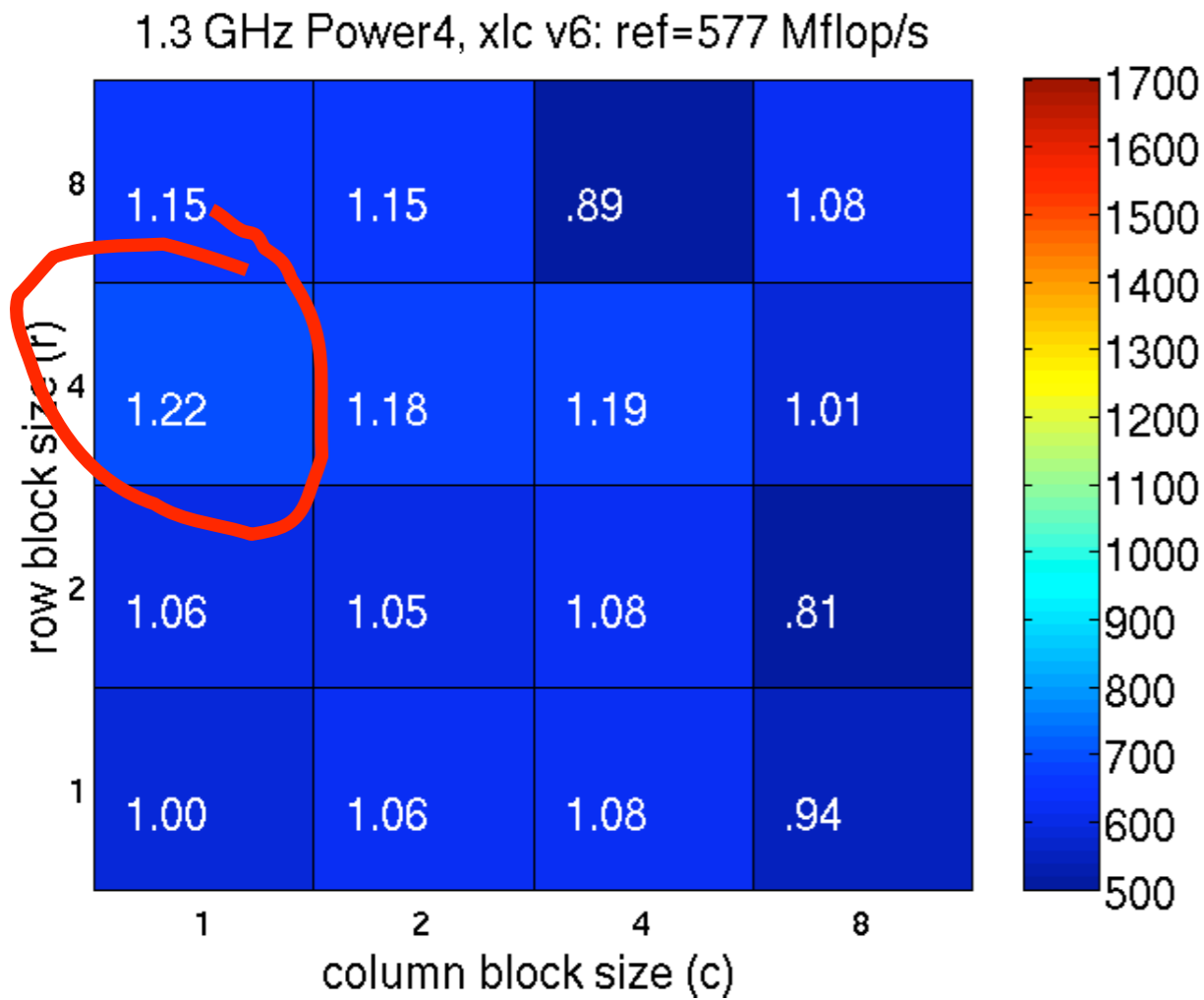
1.3 GHz Power4, xlc v6: ref=577 Mflop/s



1.9 GHz Power5, xlc v7: ref=502 Mflop/s



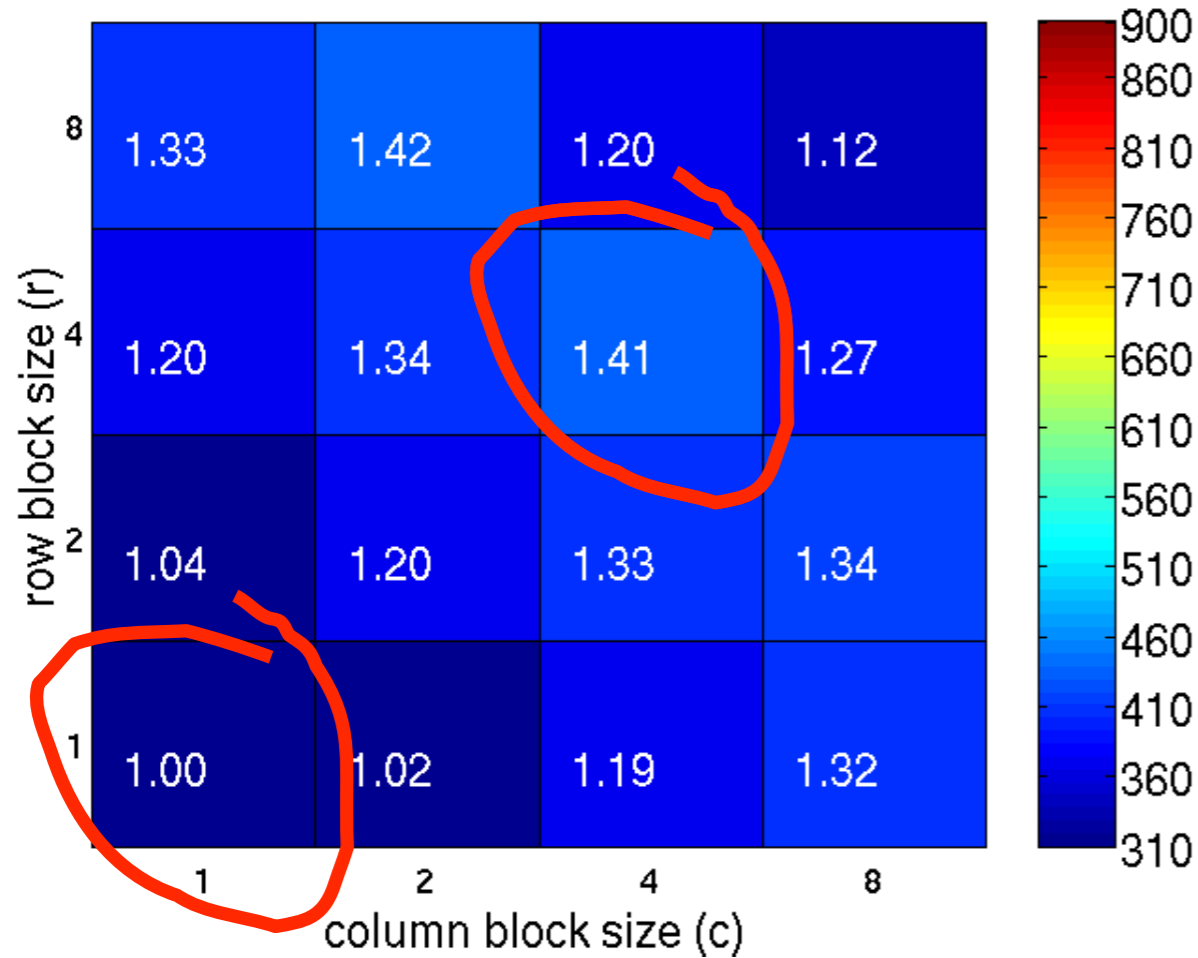
Reference gets worse!



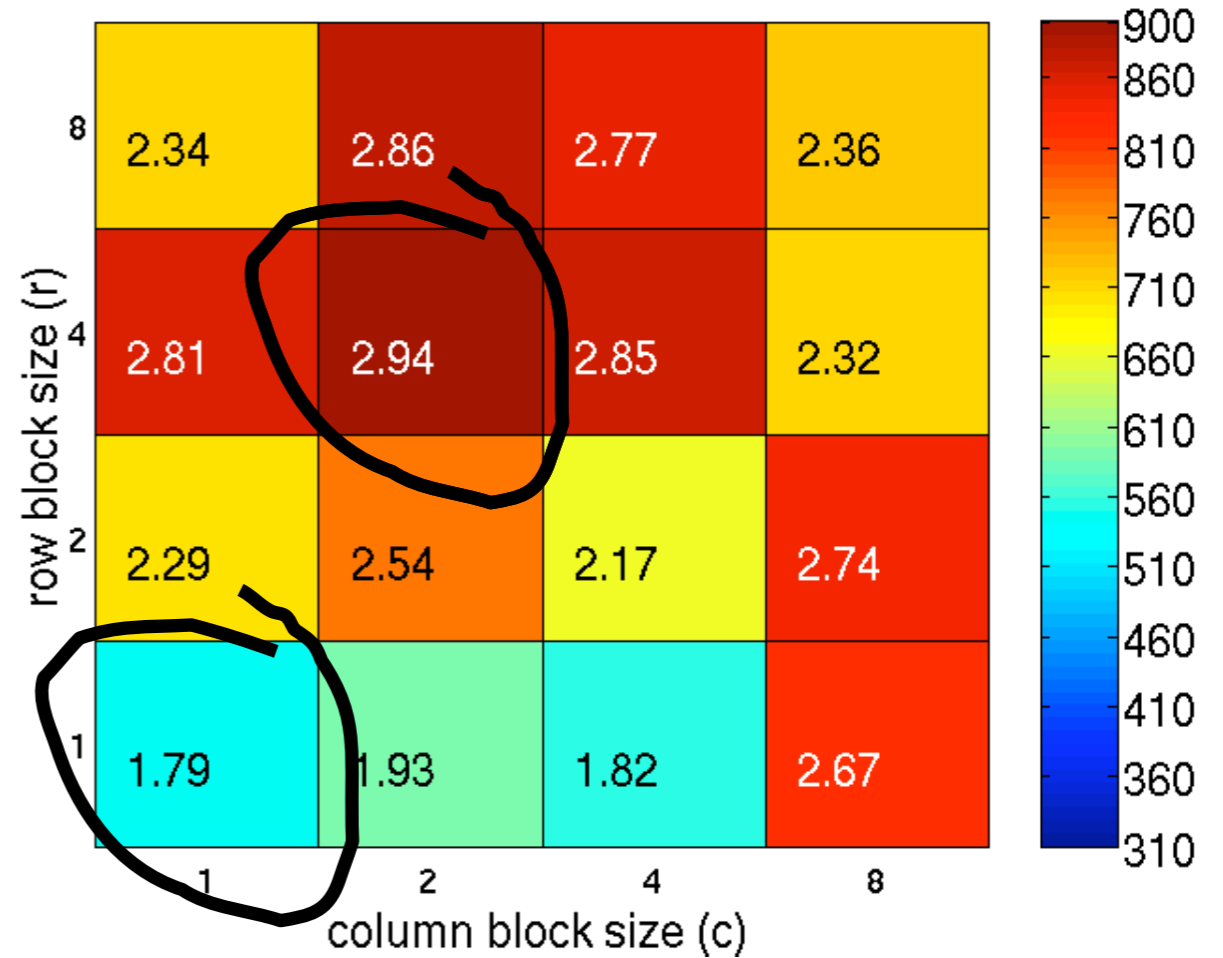
Importance of tuning increases.



2 GHz Pentium M, Intel C v8.1: ref=306 Mflop/s

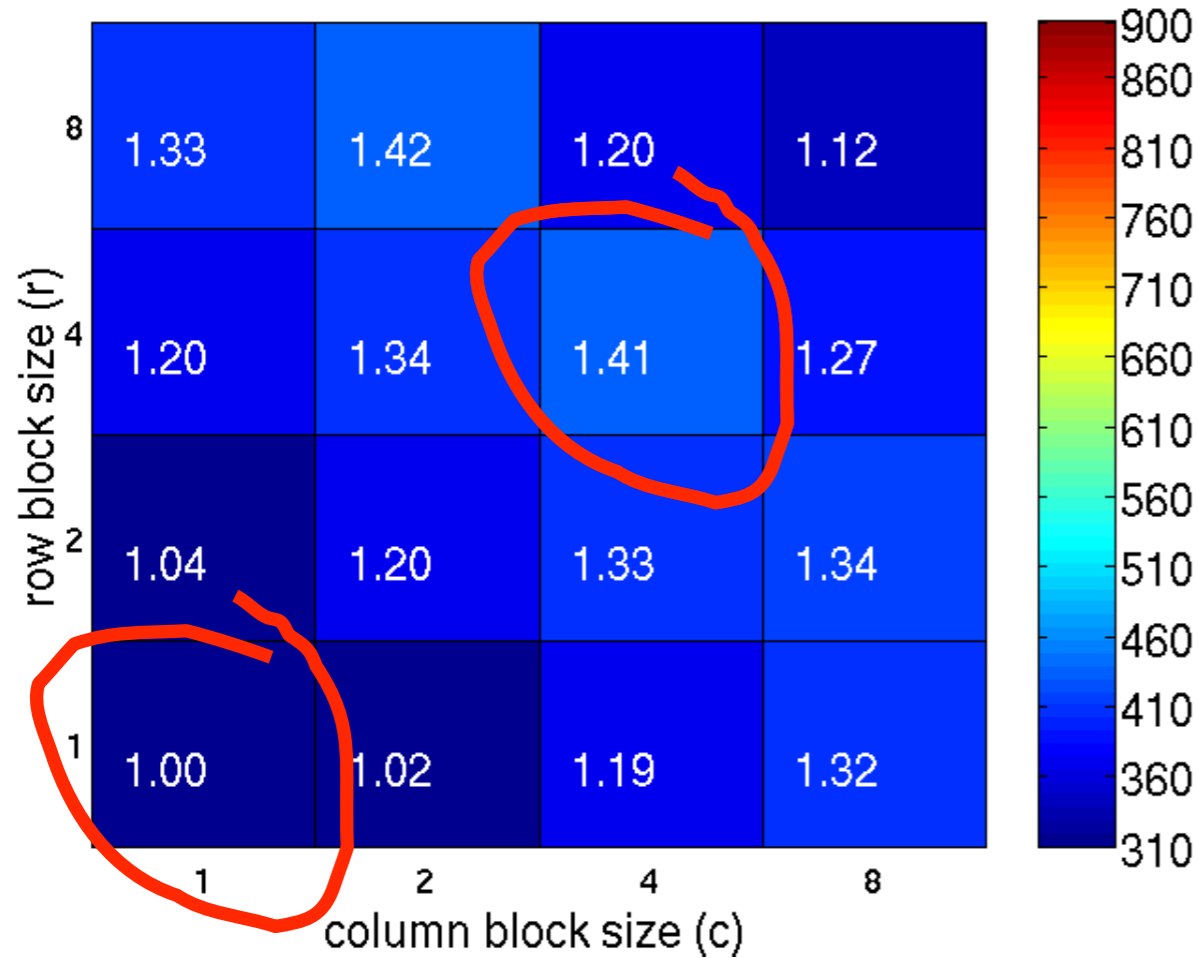


2.33 GHz Core 2 Duo, Intel C v9.1: ref=549 Mflop/s

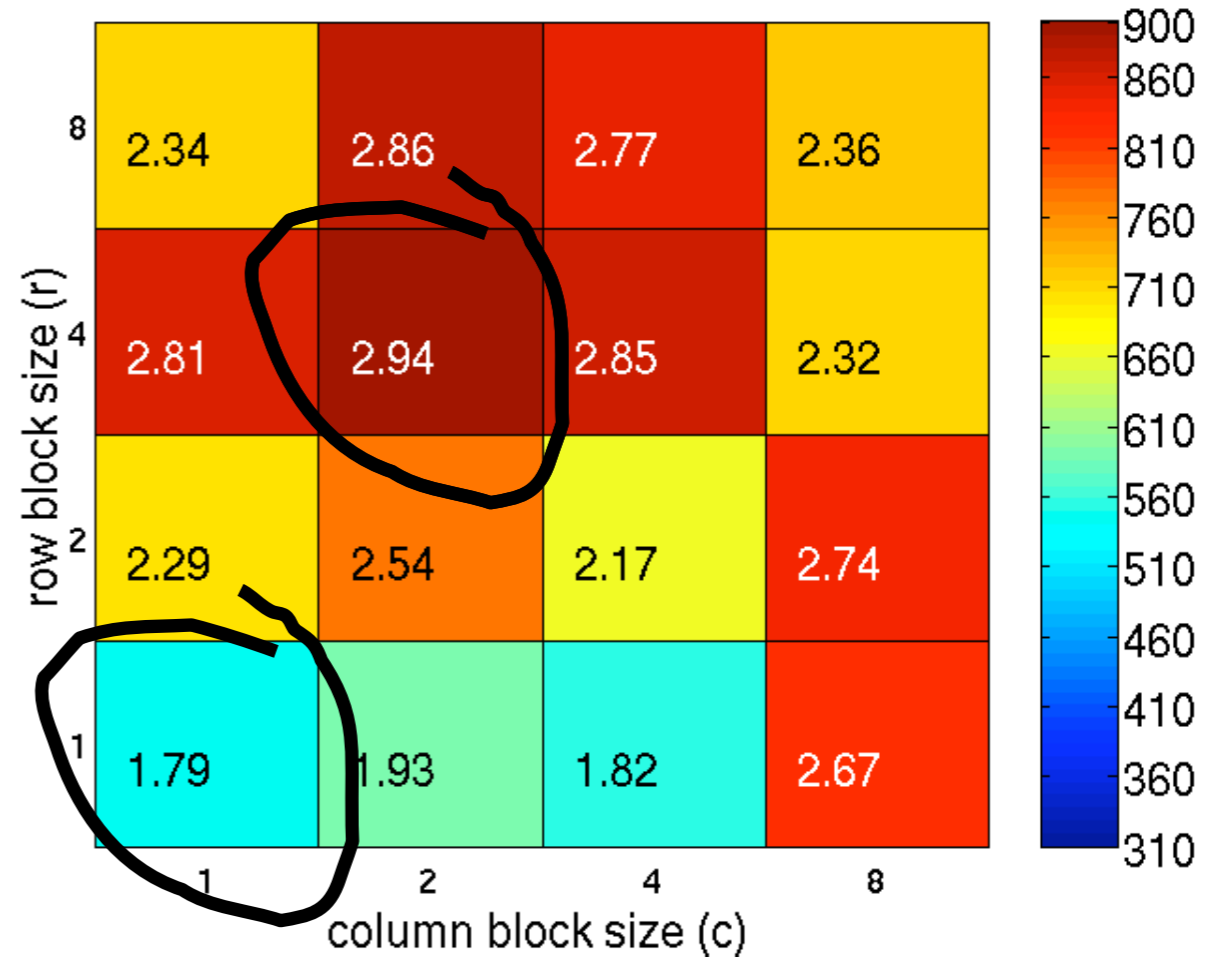


Example 3: Single- to dual-core.

2 GHz Pentium M, Intel C v8.1: ref=306 Mflop/s

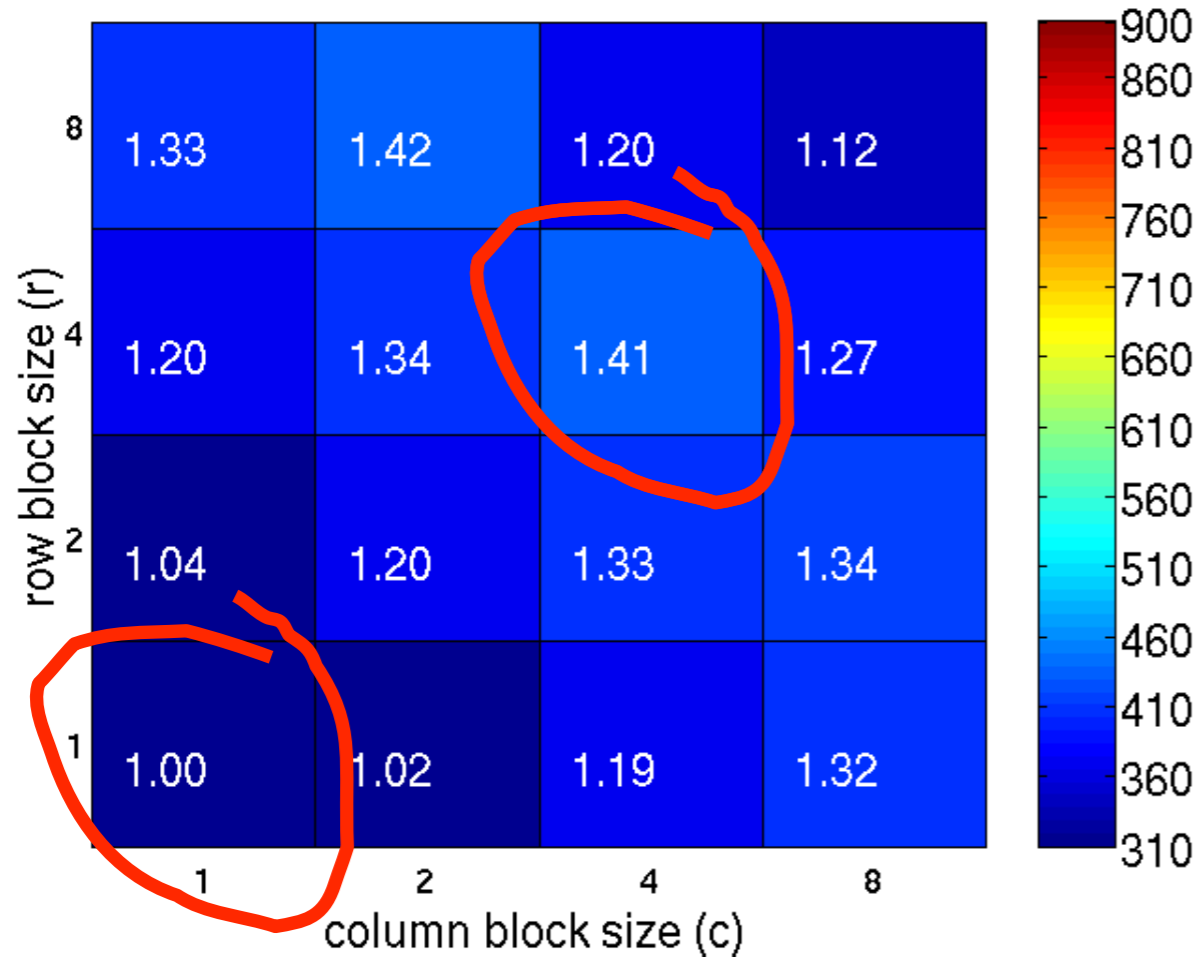


2.33 GHz Core 2 Duo, Intel C v9.1: ref=549 Mflop/s

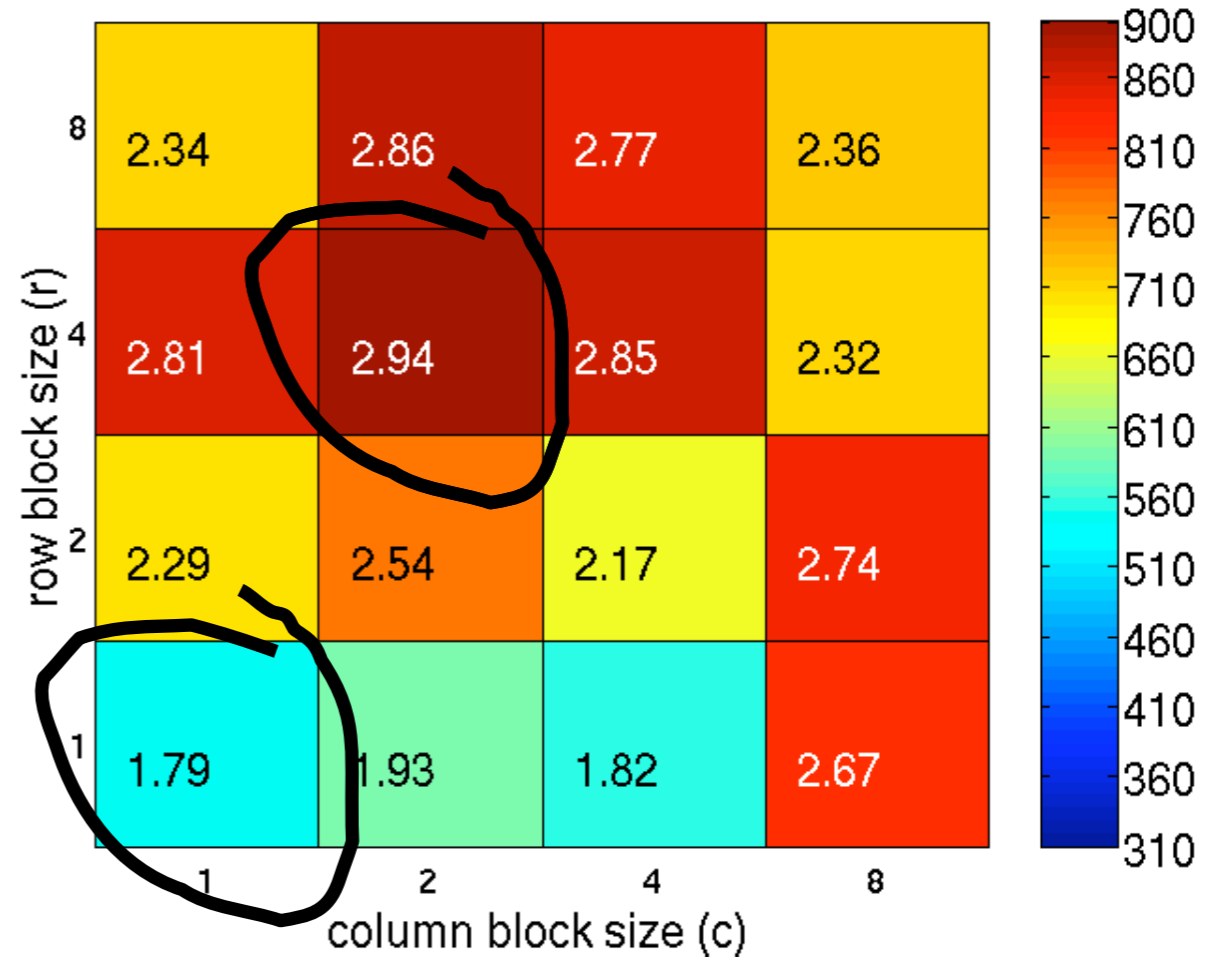


Reference & best improve...

2 GHz Pentium M, Intel C v8.1: ref=306 Mflop/s

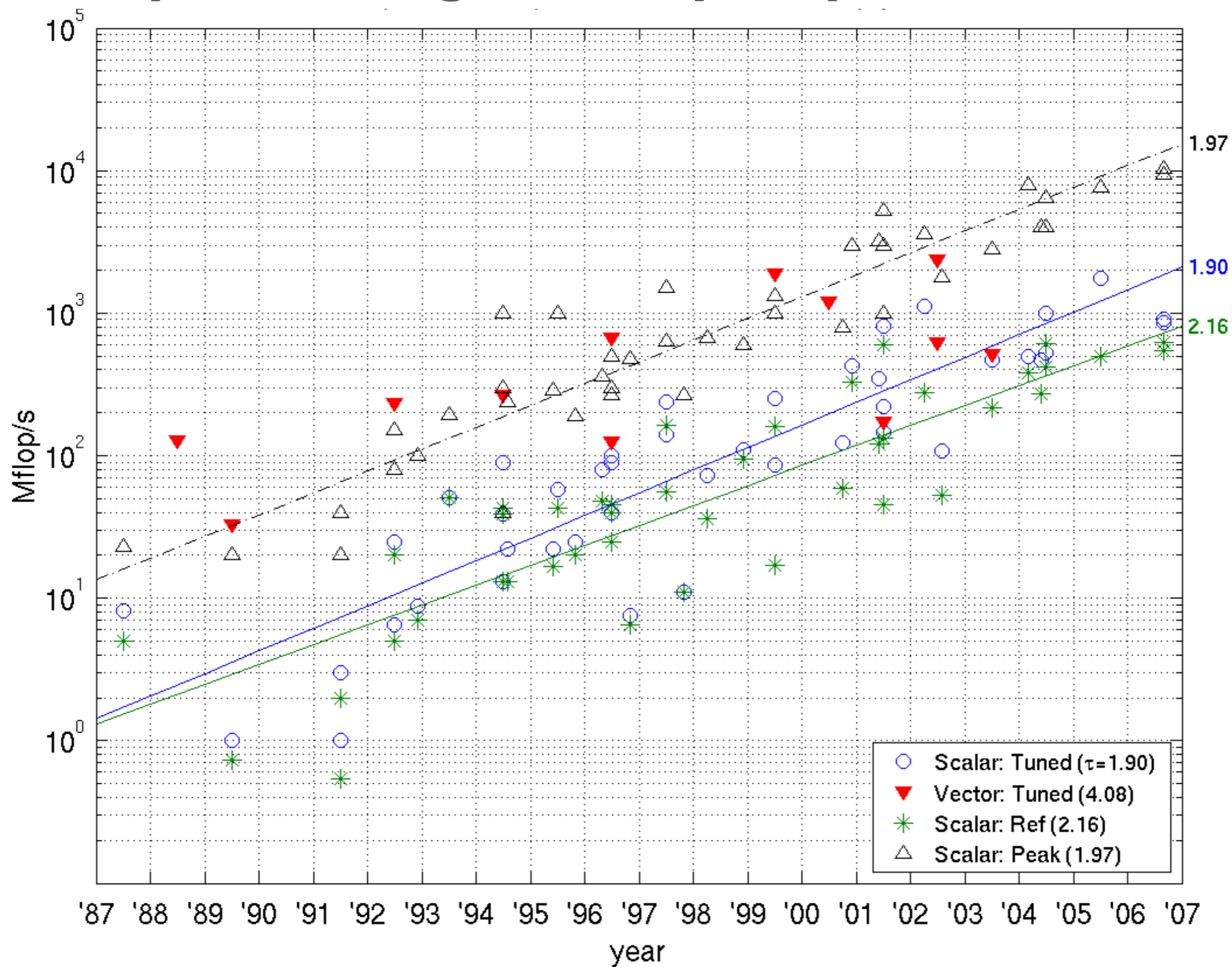


2.33 GHz Core 2 Duo, Intel C v9.1: ref=549 Mflop/s

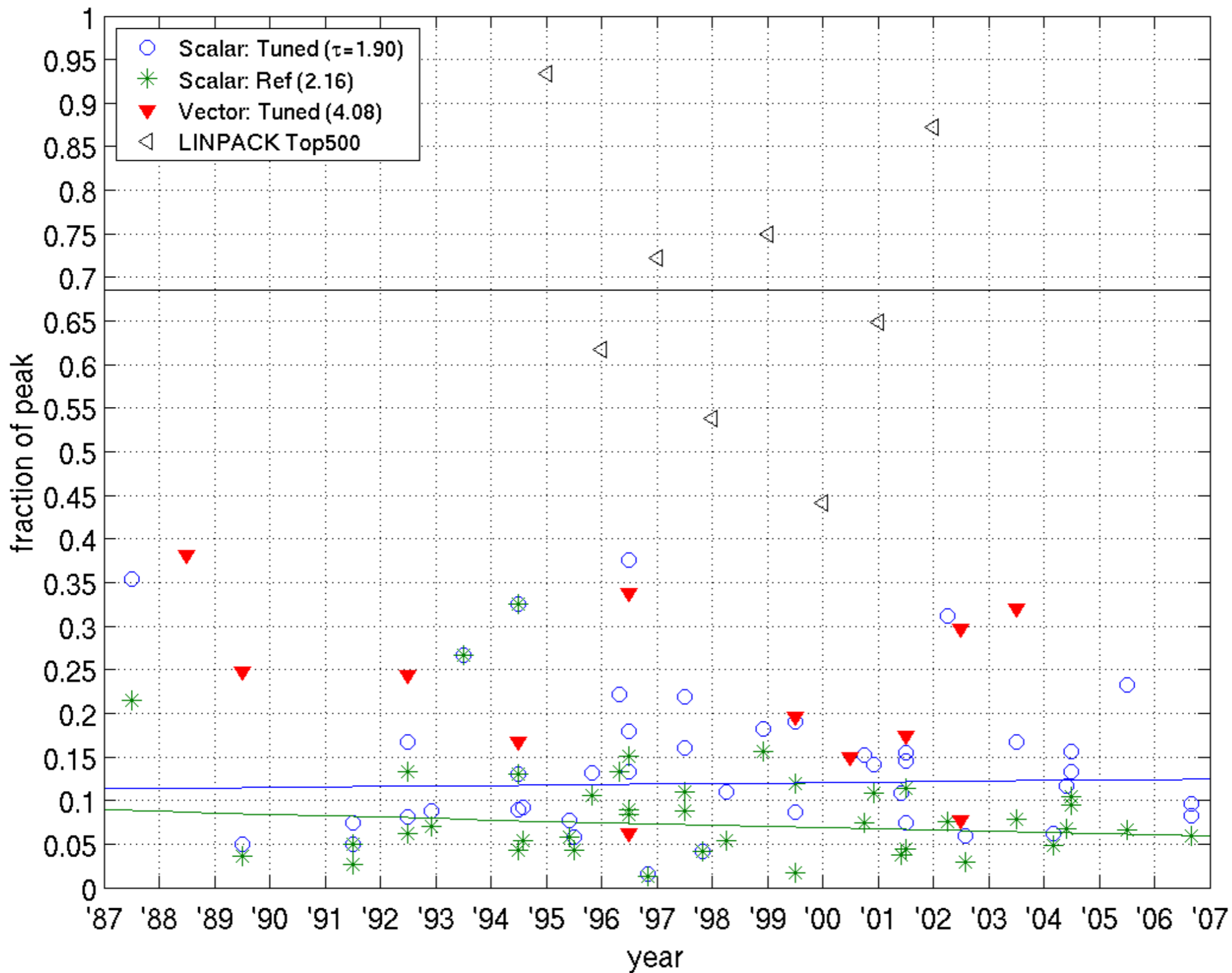


Fraction of peak: 11% → 9.5%

# 20 years of single-core SpMV performance



# 20 years of single-core SpMV performance



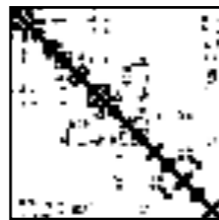
- ▶ **Tuning for multicore:  
Matrices and machines**

**2K x 2K Dense matrix  
stored in sparse format**

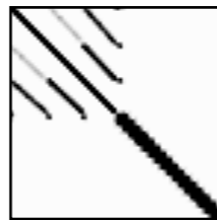


Dense

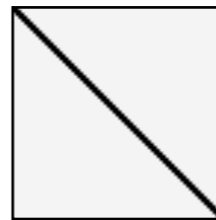
**Well Structured**  
(sorted by nonzeros/row)



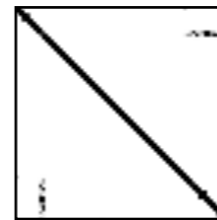
Protein



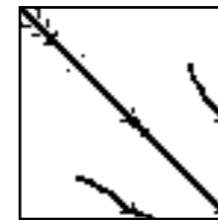
FEM /  
Spheres



FEM /  
Cantilever



Wind  
Tunnel



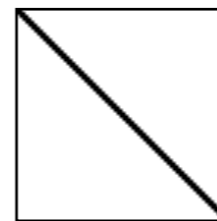
FEM /  
Harbor



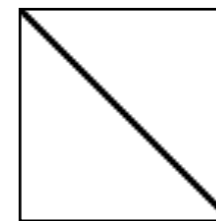
QCD



FEM /  
Ship



Economics



Epidemiology

**Poorly Structured  
Mix**



FEM /  
Accelerator



Circuit

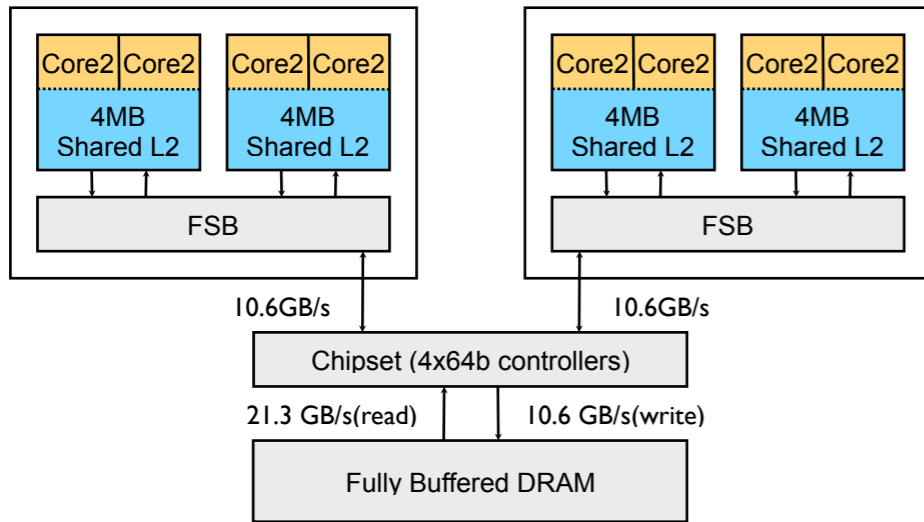


webbase

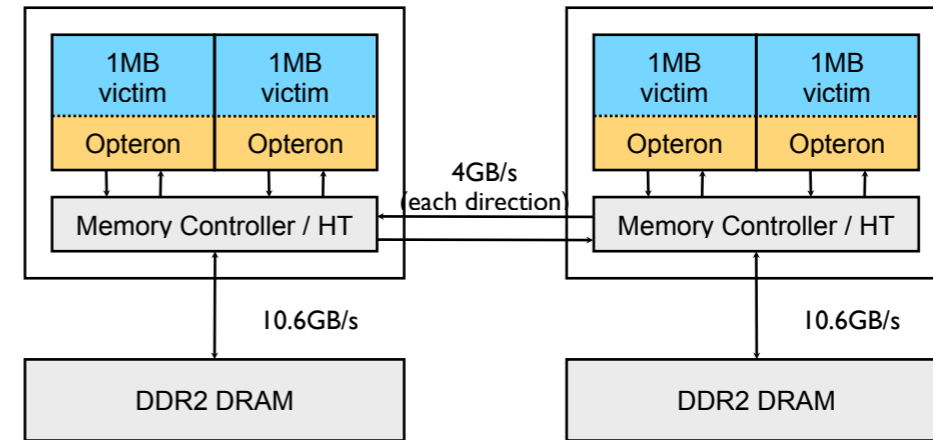
**Extreme Aspect Ratio  
(linear programming)**



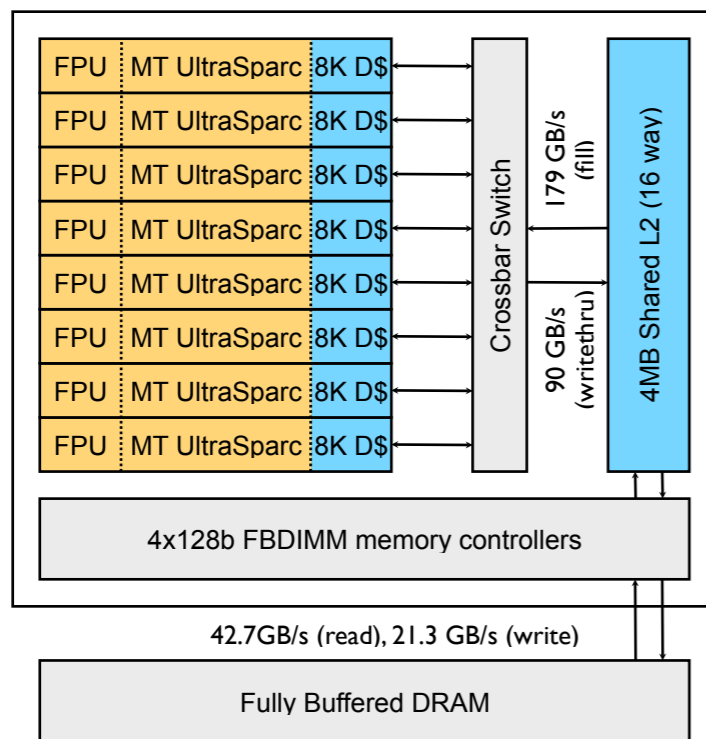
$n \sim 2k$  to  $1M$   
 $nnz \sim 2$  to  $12M$   
 $nnz/row \sim 3$  to  $3k$



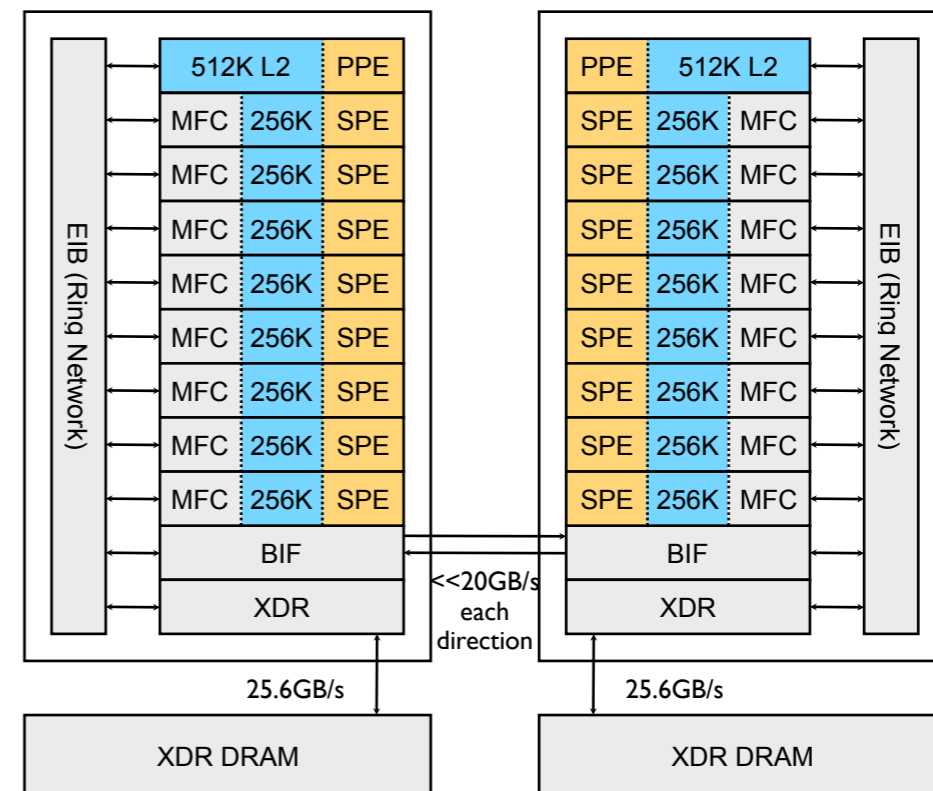
Intel Clovertown (2 x 4)



AMD Opteron X2 (2 x 2)

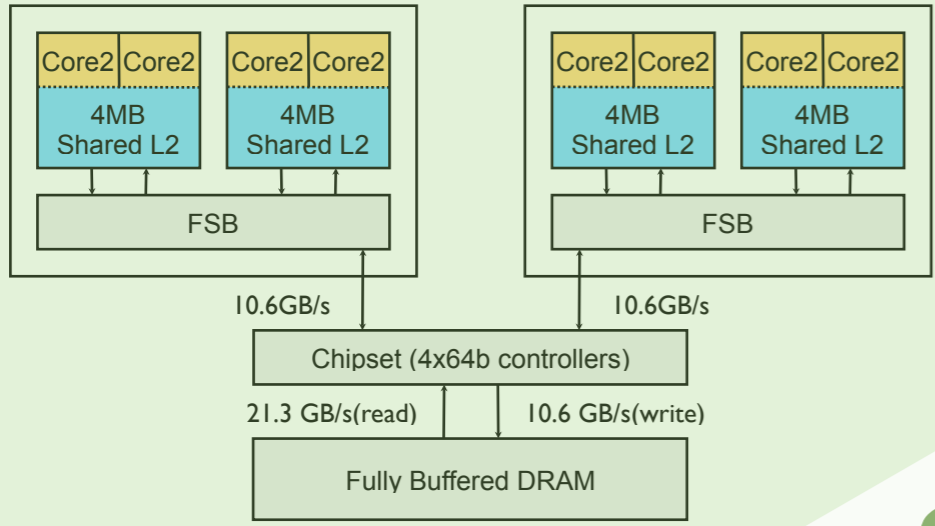


Sun Niagara2 (1 x 8 x 8)

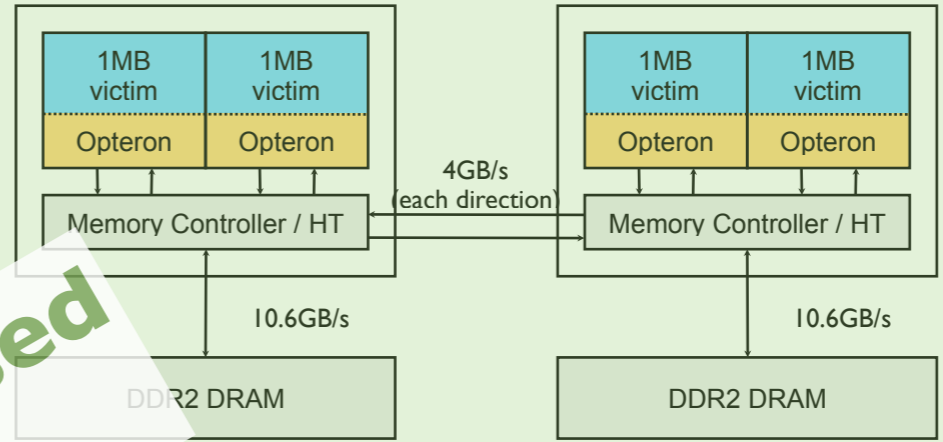


STI Cell Blade (2 x 8)

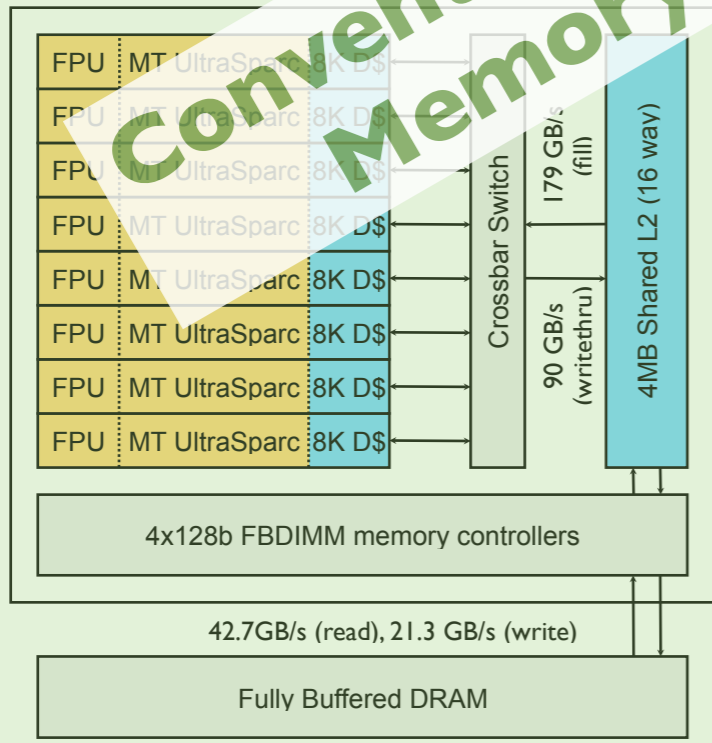




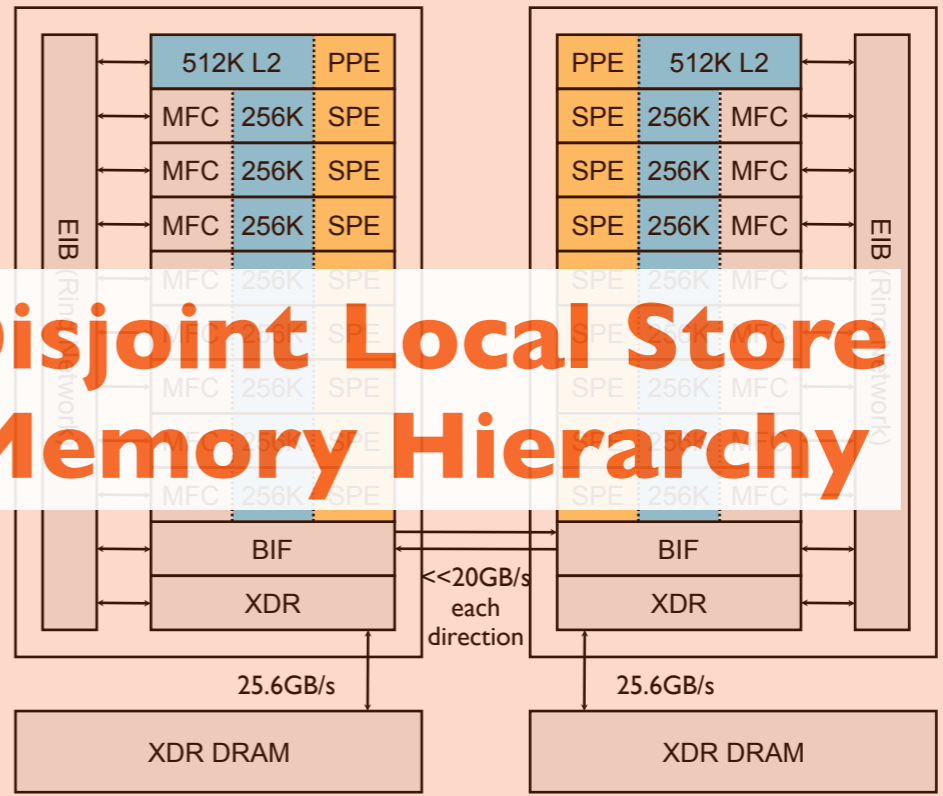
Intel Clovertown



AMD Opteron



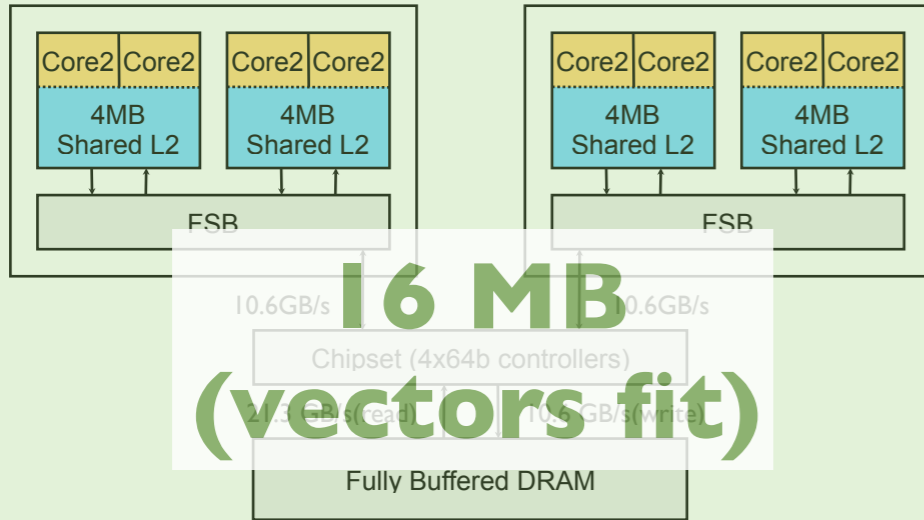
Sun Niagara2



IBM Cell Blade

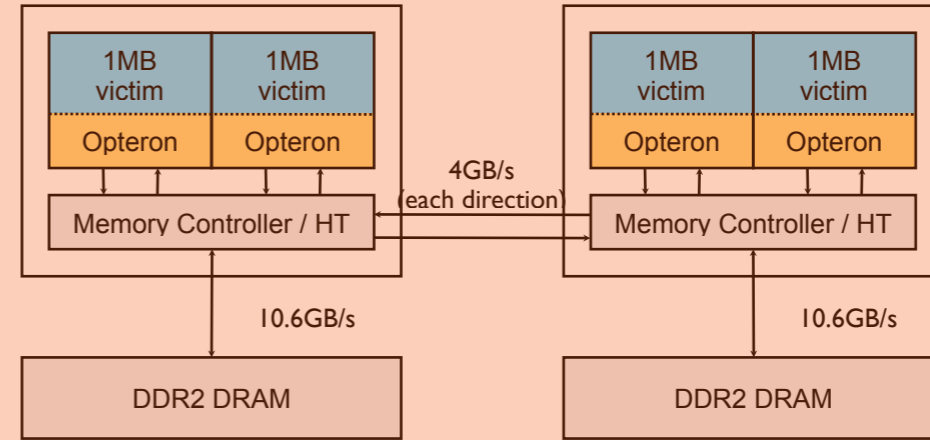
Conventional Cache-based Memory Hierarchy

Disjoint Local Store Memory Hierarchy

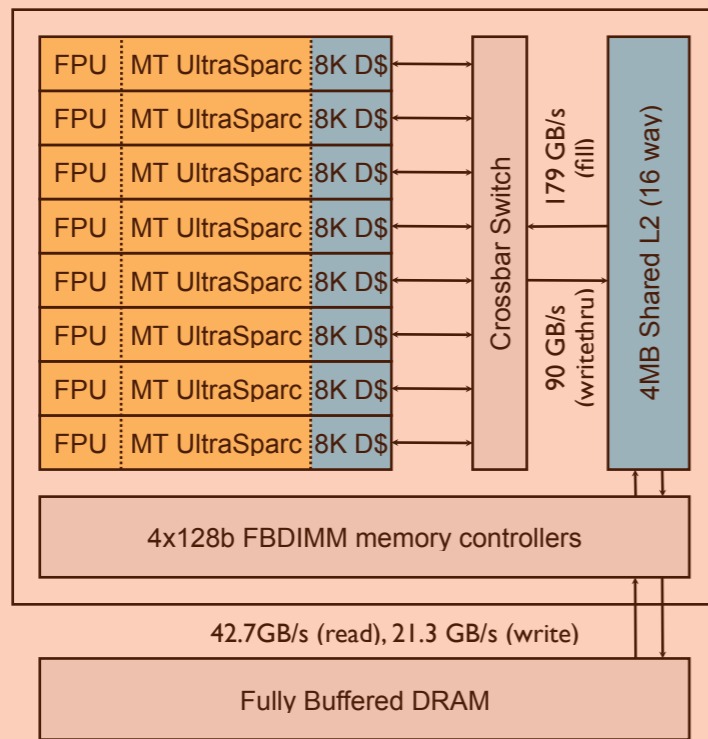


**16 MB**  
(vectors fit)

Intel Clovertown

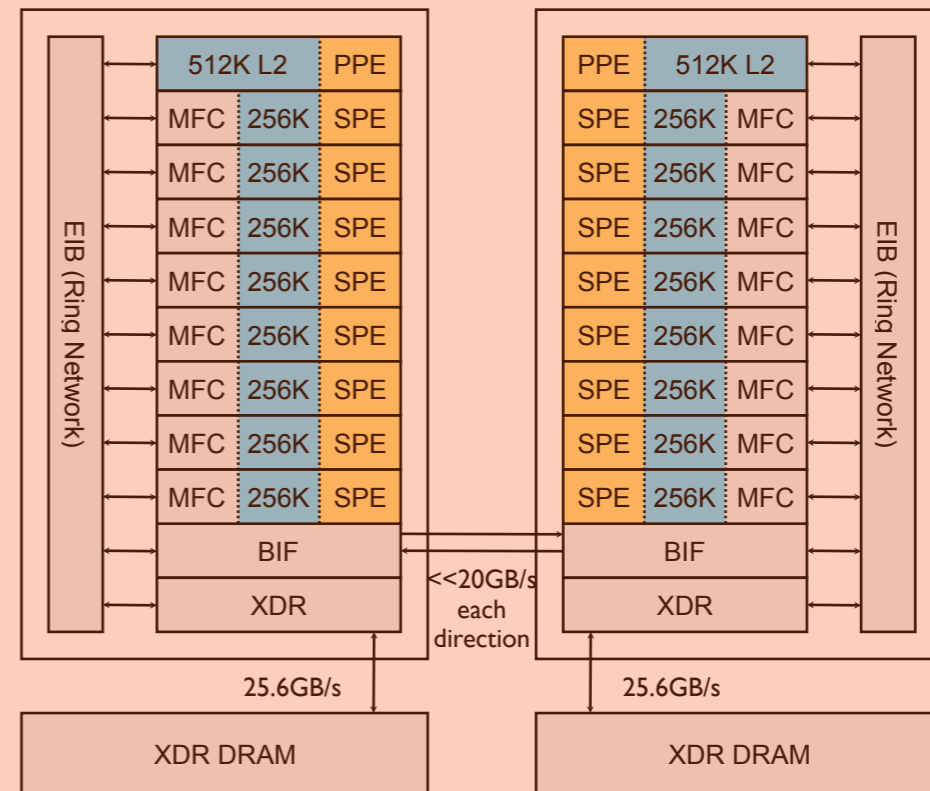


AMD Opteron

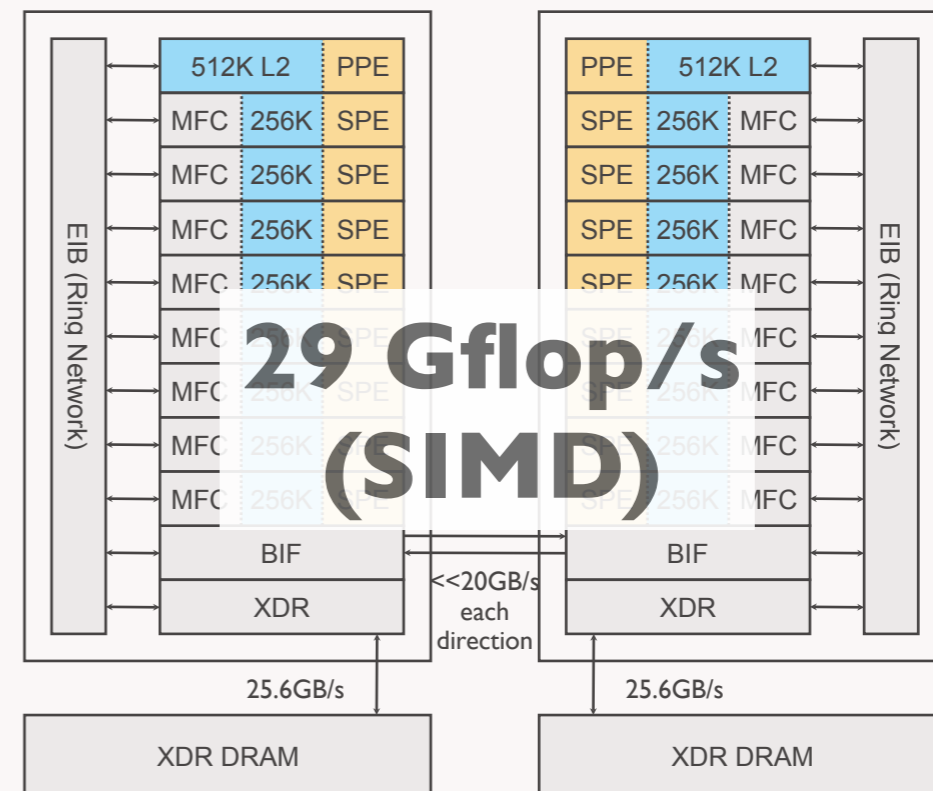
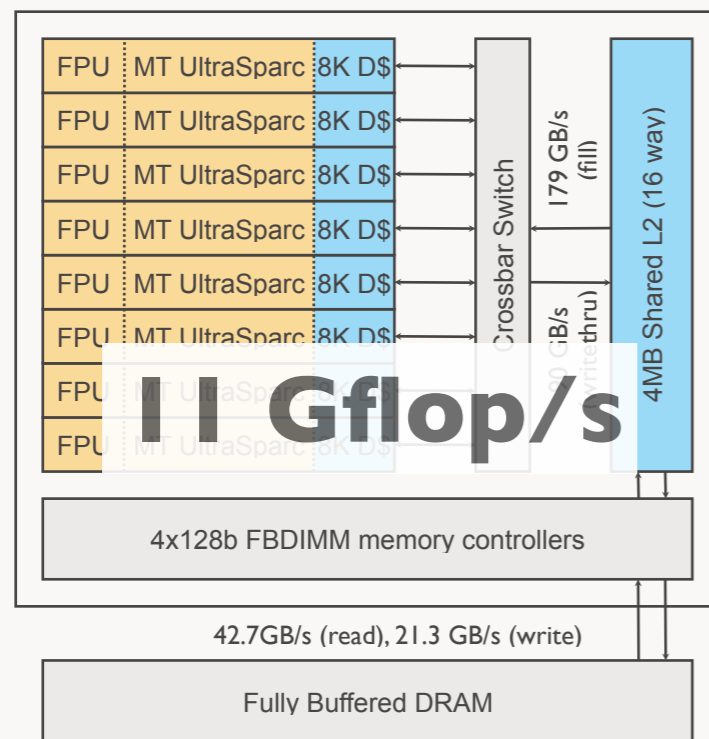
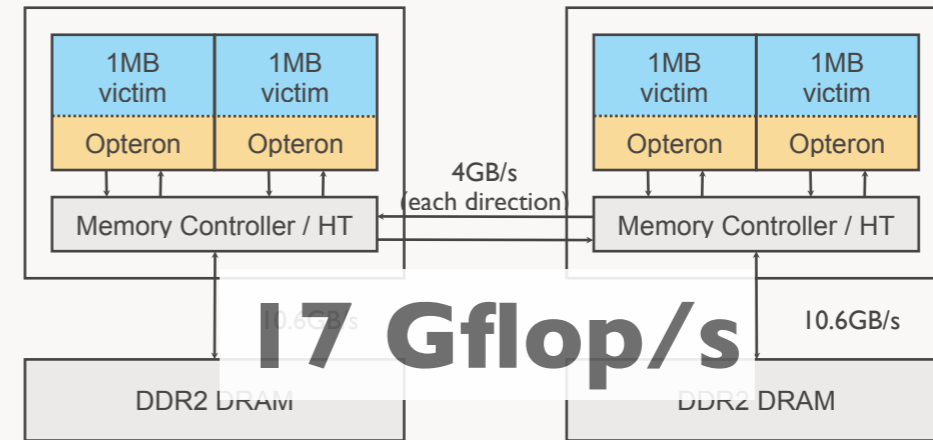
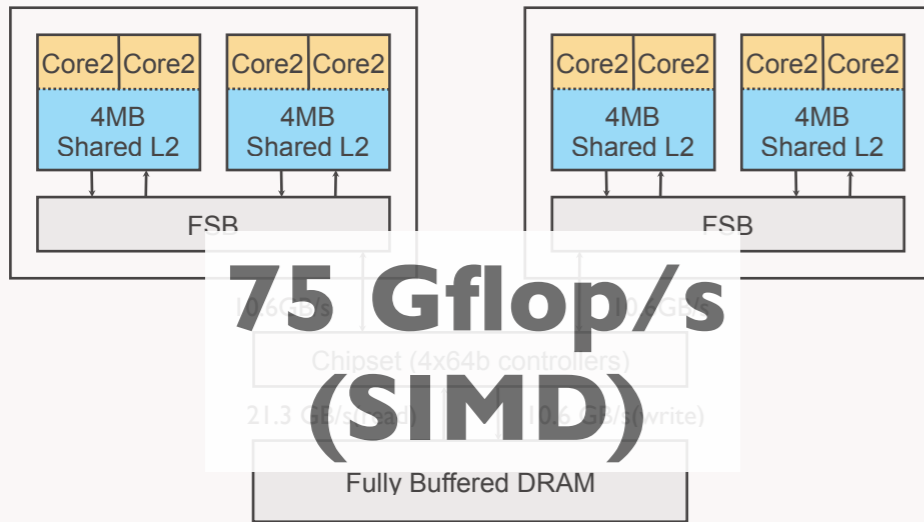


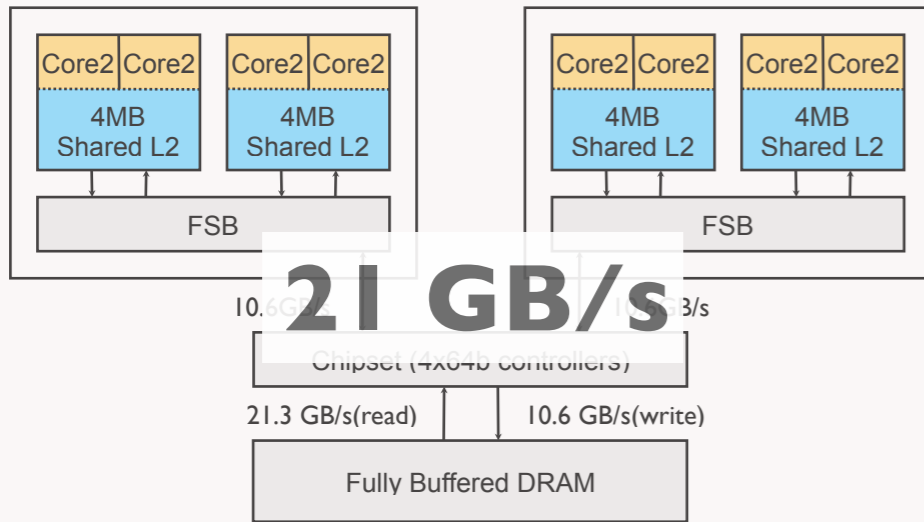
**4 MB**

Sun Niagara2

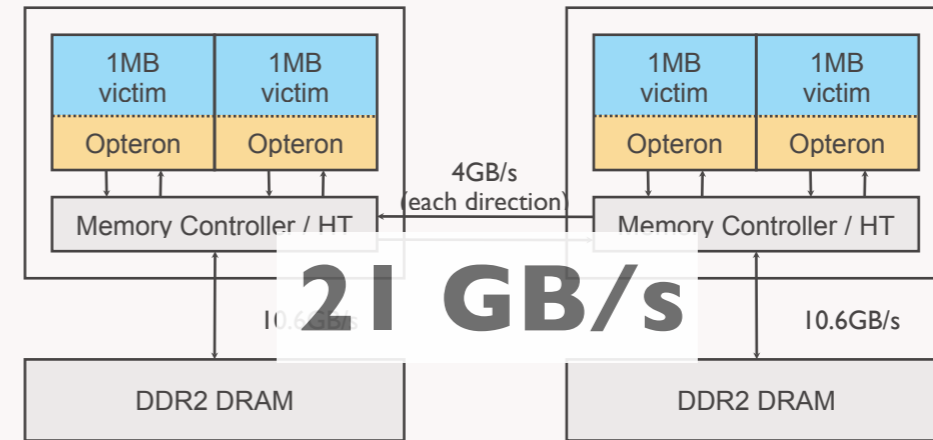


IBM Cell Blade

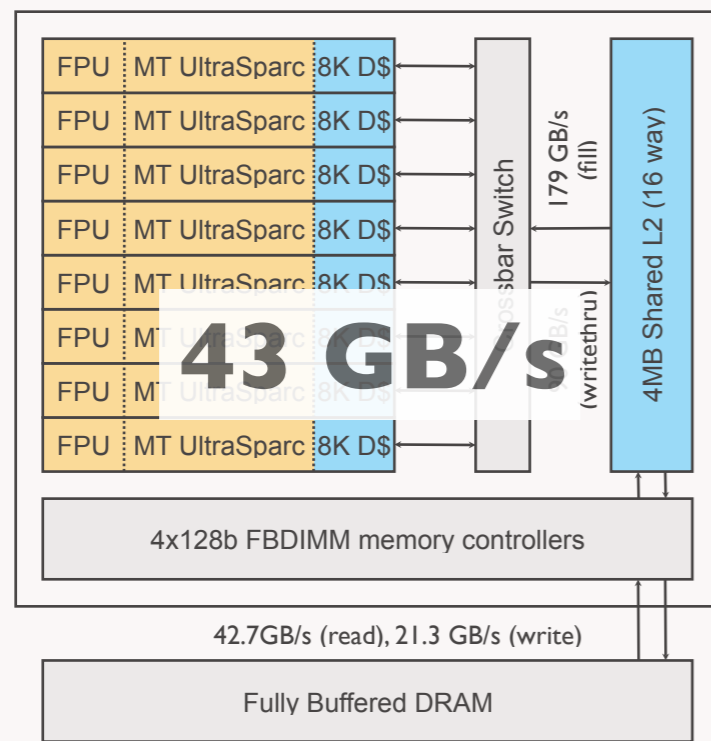




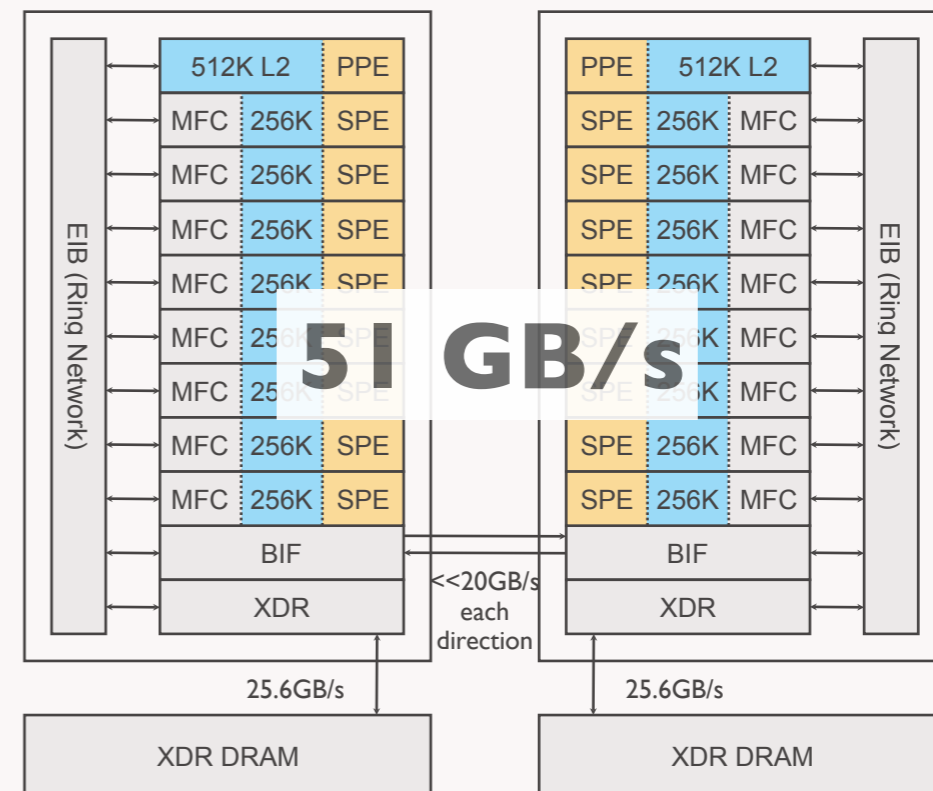
Intel Clovertown



AMD Opteron

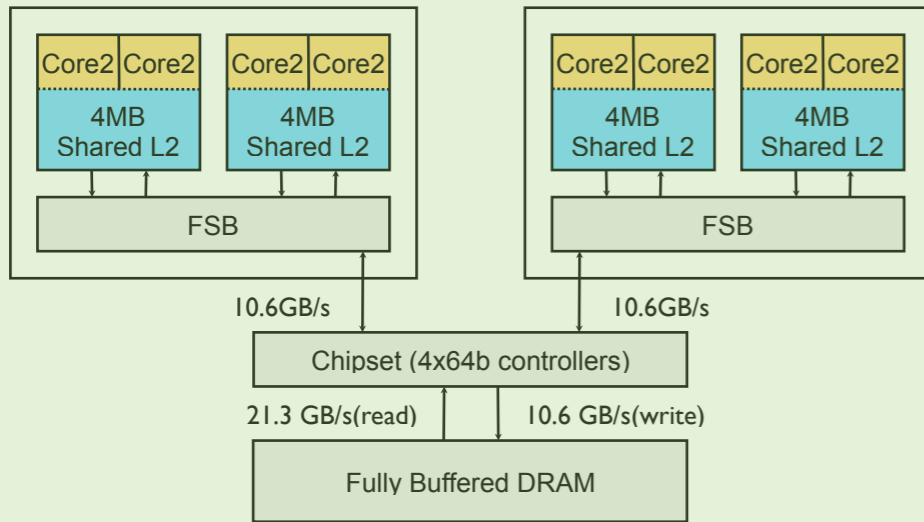


Sun Niagara2



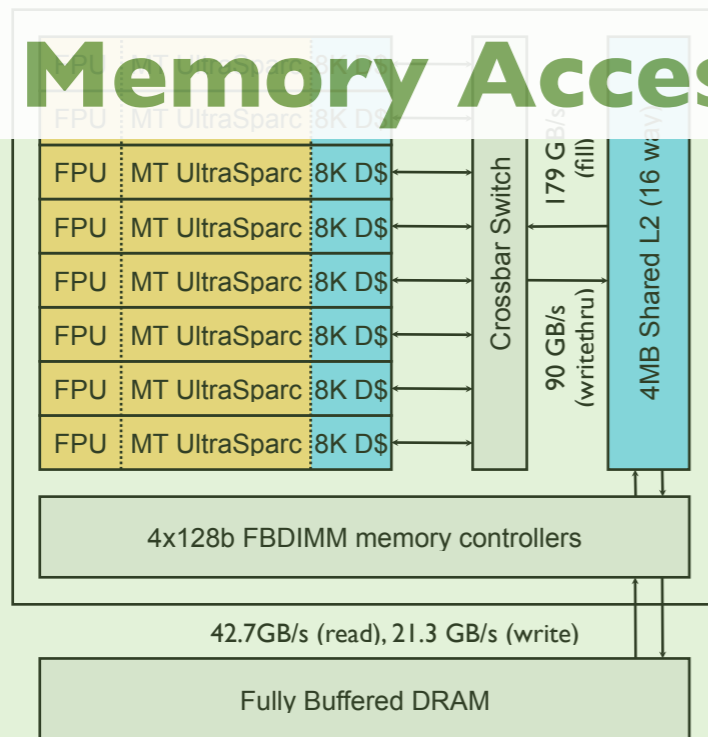
IBM Cell Blade

(Peak read bandwidth)

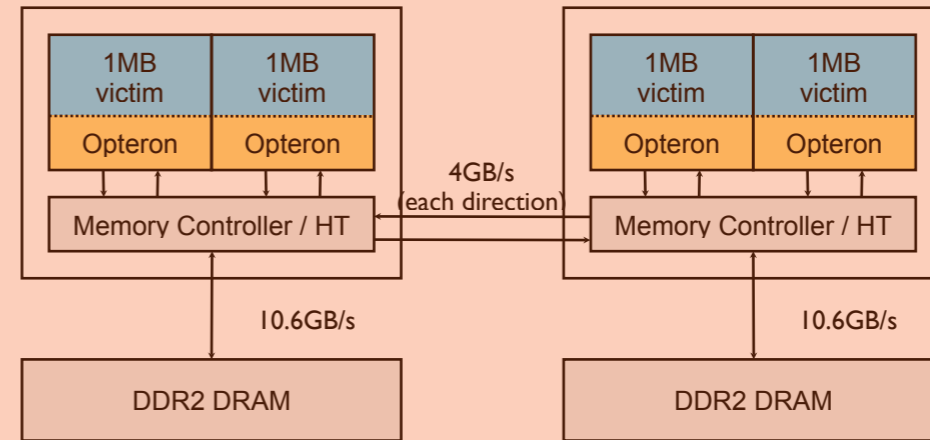


Intel Clovertown

**Uniform  
Memory Access**

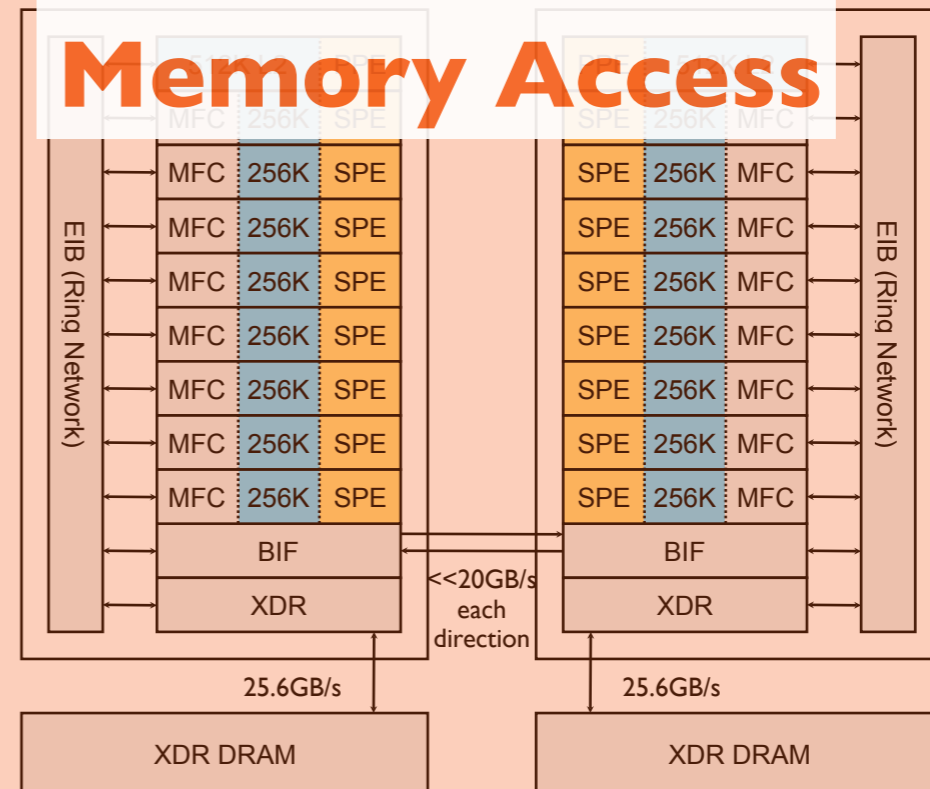


Sun Niagara2



AMD Opteron

**Non-uniform  
Memory Access**



IBM Cell Blade

(Peak read bandwidth)

## ▶ **Naïve codes on cache-based machines**

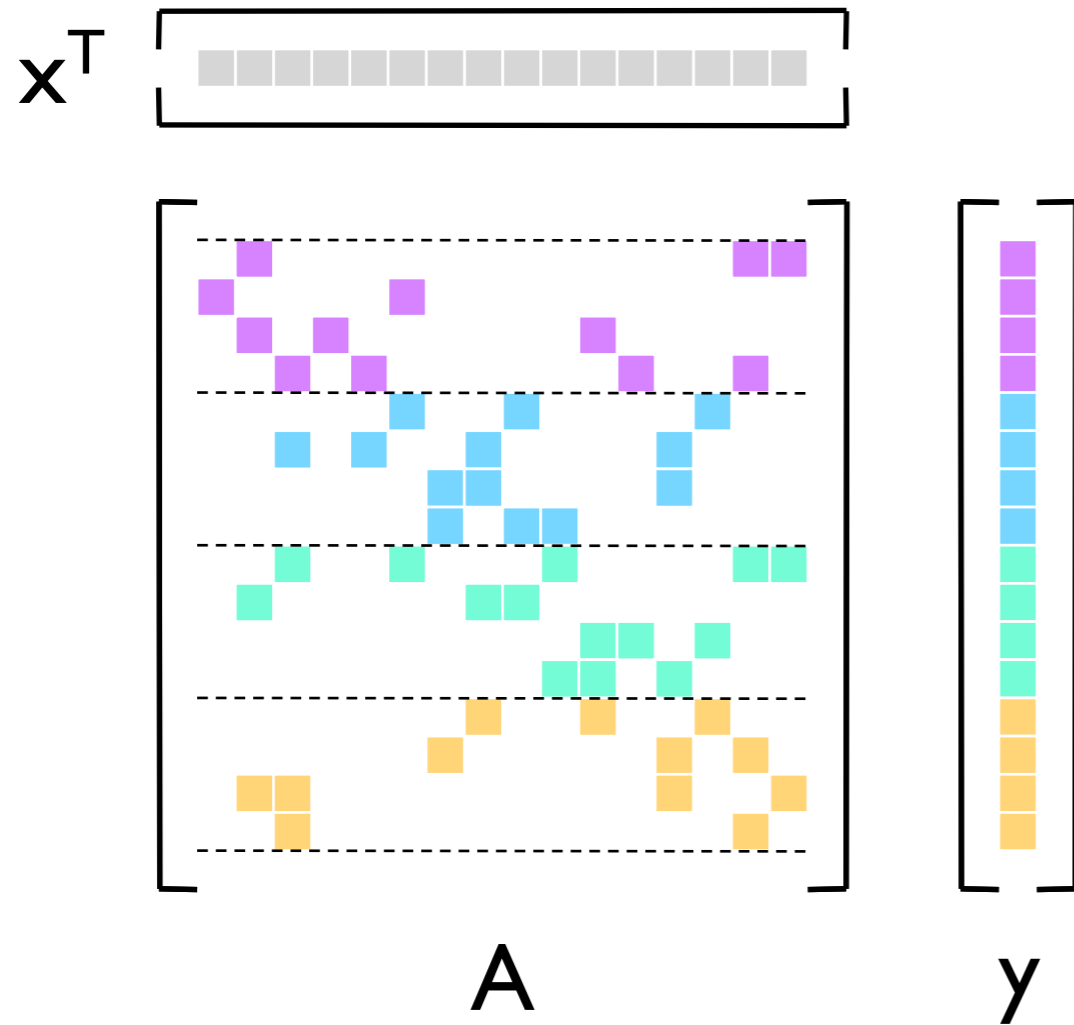
### ▶ Serial

- ▶ Compressed sparse row (CSR) storage
- ▶ “Vanilla” C
- ▶ “Best” compiler options, with some search

### ▶ PThreads

- ▶ I-D block rows, in CSR
- ▶ Roughly equal nnz / block
- ▶ Barrier before & after SpMV

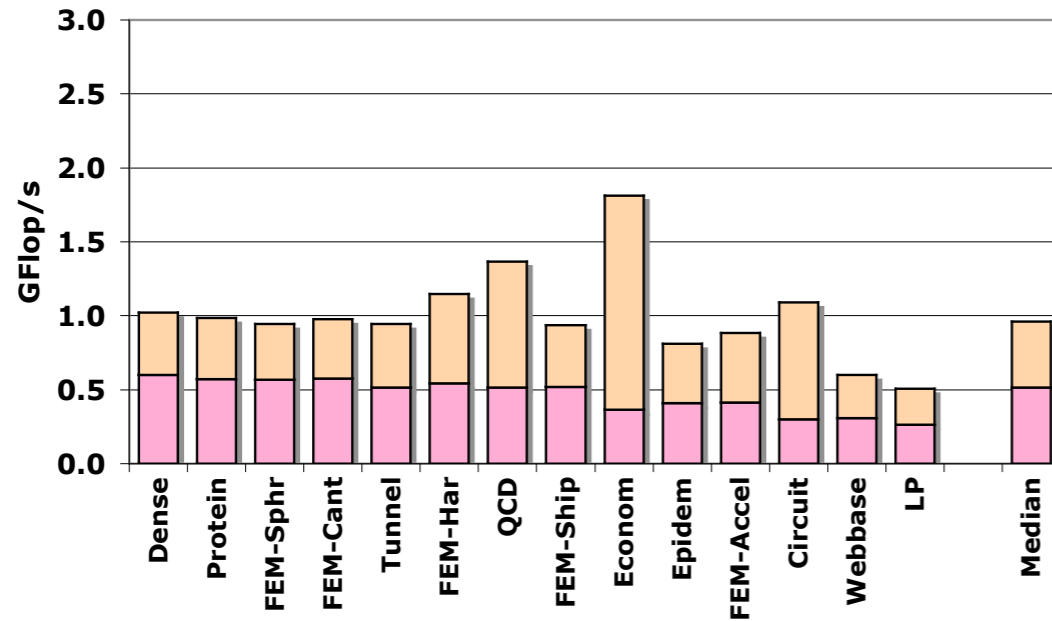
# Parallelization



- ▶ I-D block rows, in CSR
- ▶ Roughly equal nnz / block
- ▶ Barrier before & after SpMV

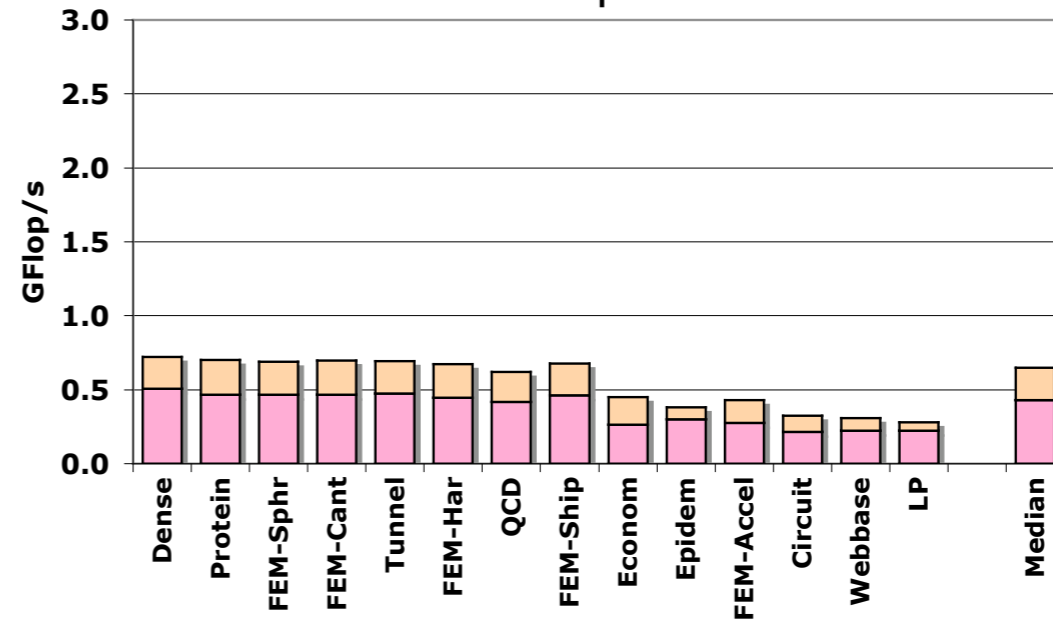
**8 cores  $\Rightarrow$  1.9x**

Intel Clovertown

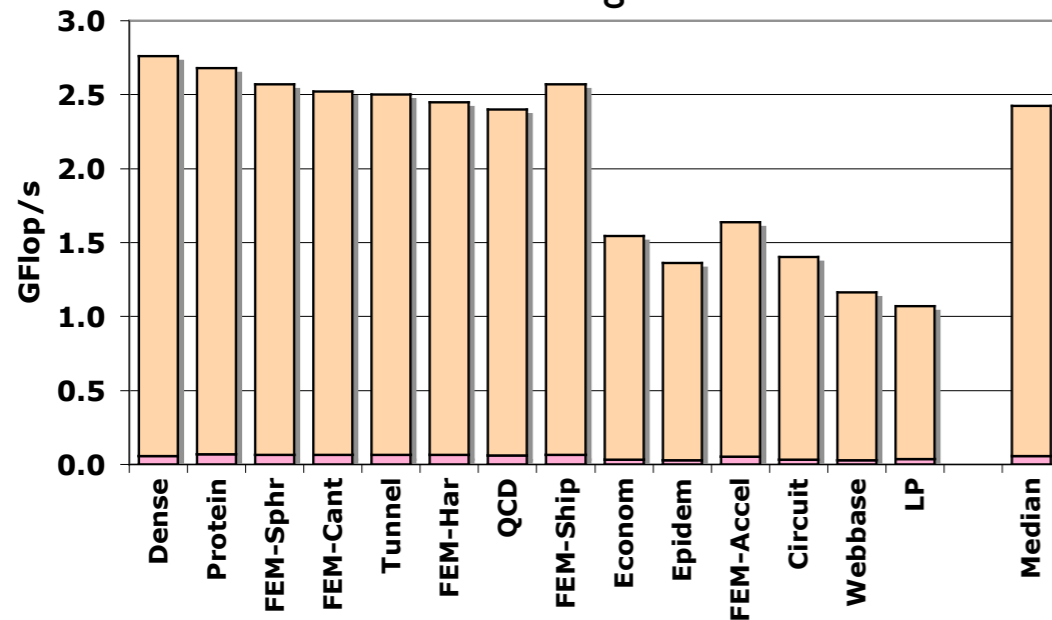


**4 cores  $\Rightarrow$  1.5x**

AMD Opteron



Sun Niagara2



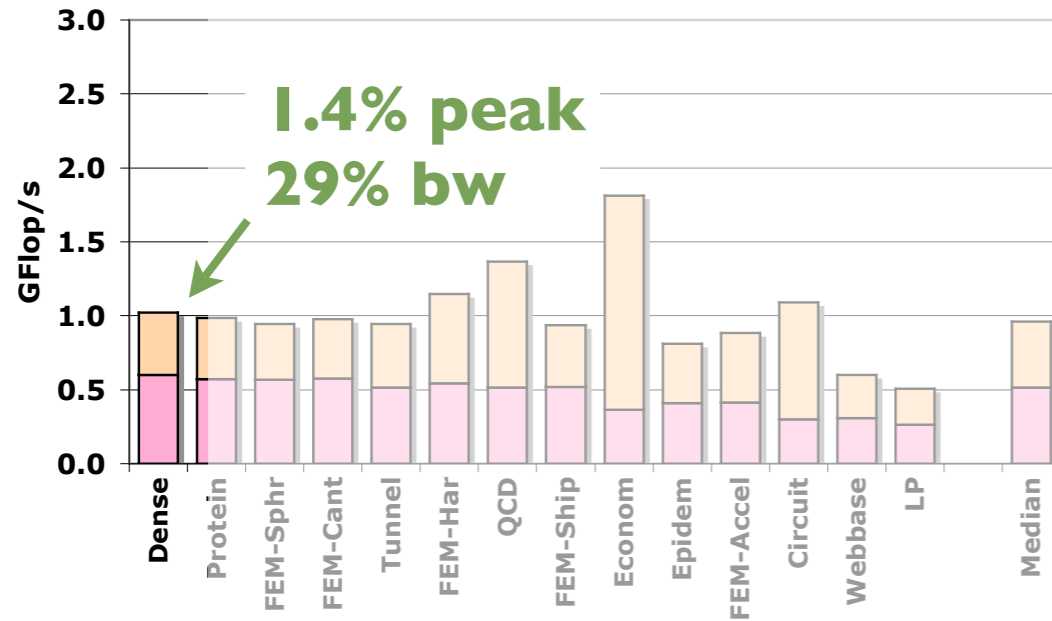
**64 threads  $\Rightarrow$  41x**

Naive Pthreads  
Naive Single Thread



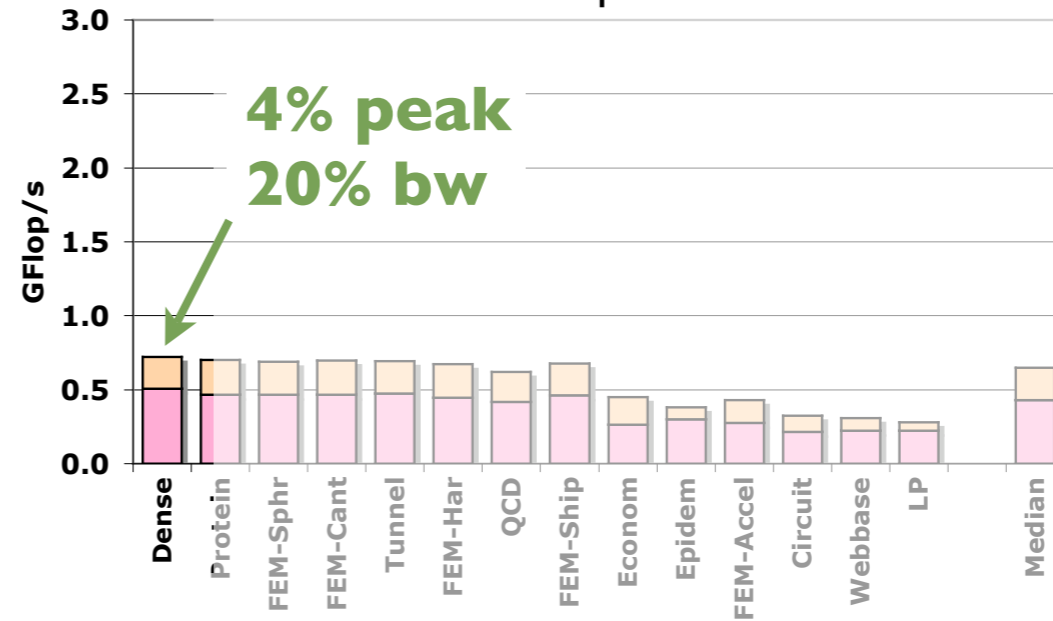
8 cores ⇒ 1.9x

Intel Clovertown

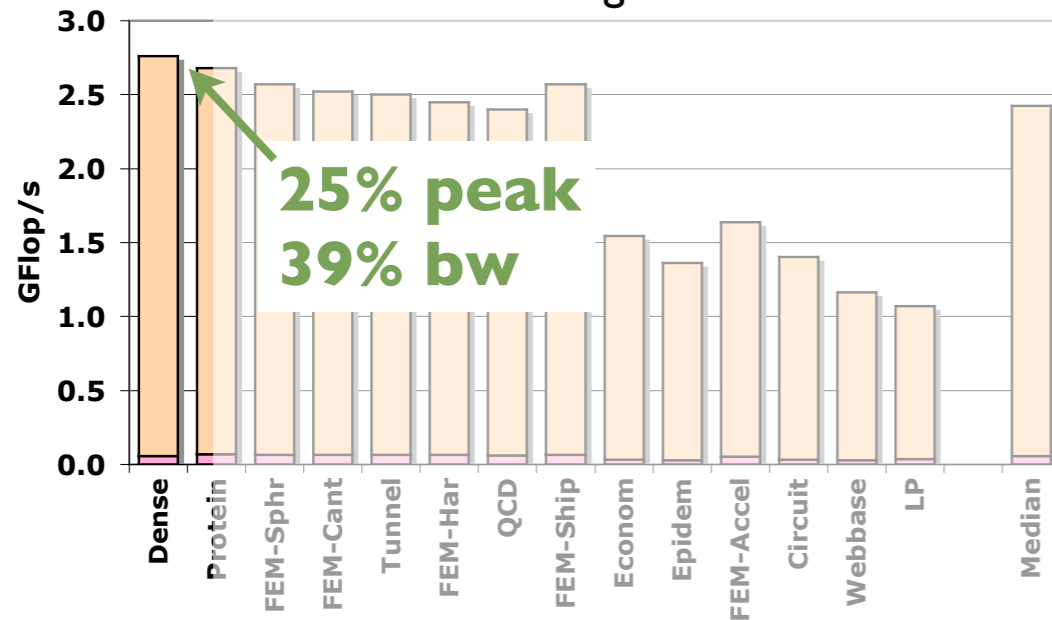


4 cores ⇒ 1.5x

AMD Opteron



Sun Niagara2



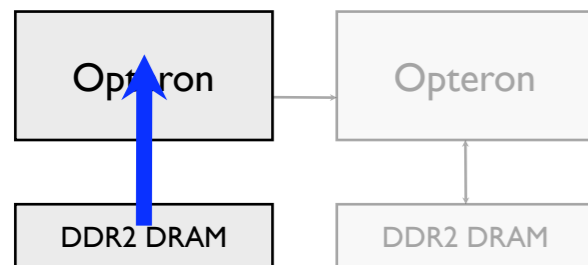
64 threads ⇒ 41x

Naive Pthreads  
Naive Single Thread

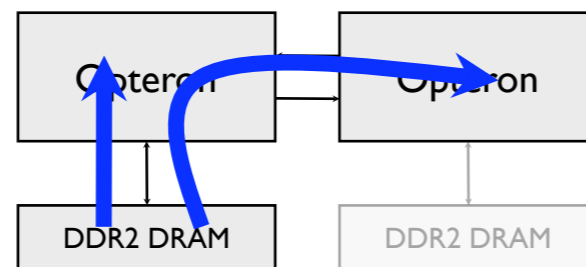
## ▶ **Multicore optimizations**

- ▶ Bind
- ▶ Prefetch
- ▶ Compress
- ▶ Block

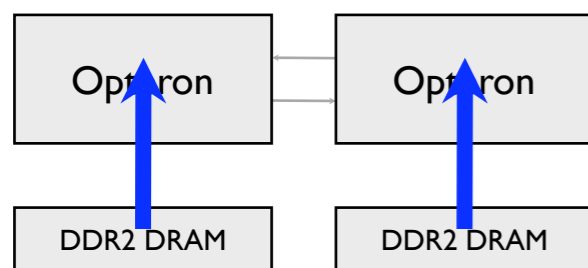
# Bind



**Single Thread**



**Multiple Threads,  
One memory controller**



**Multiple Threads,  
Both memory controllers**

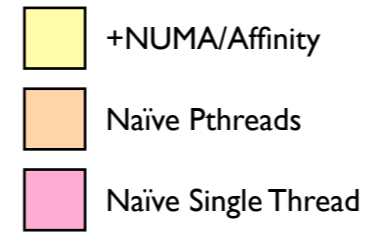
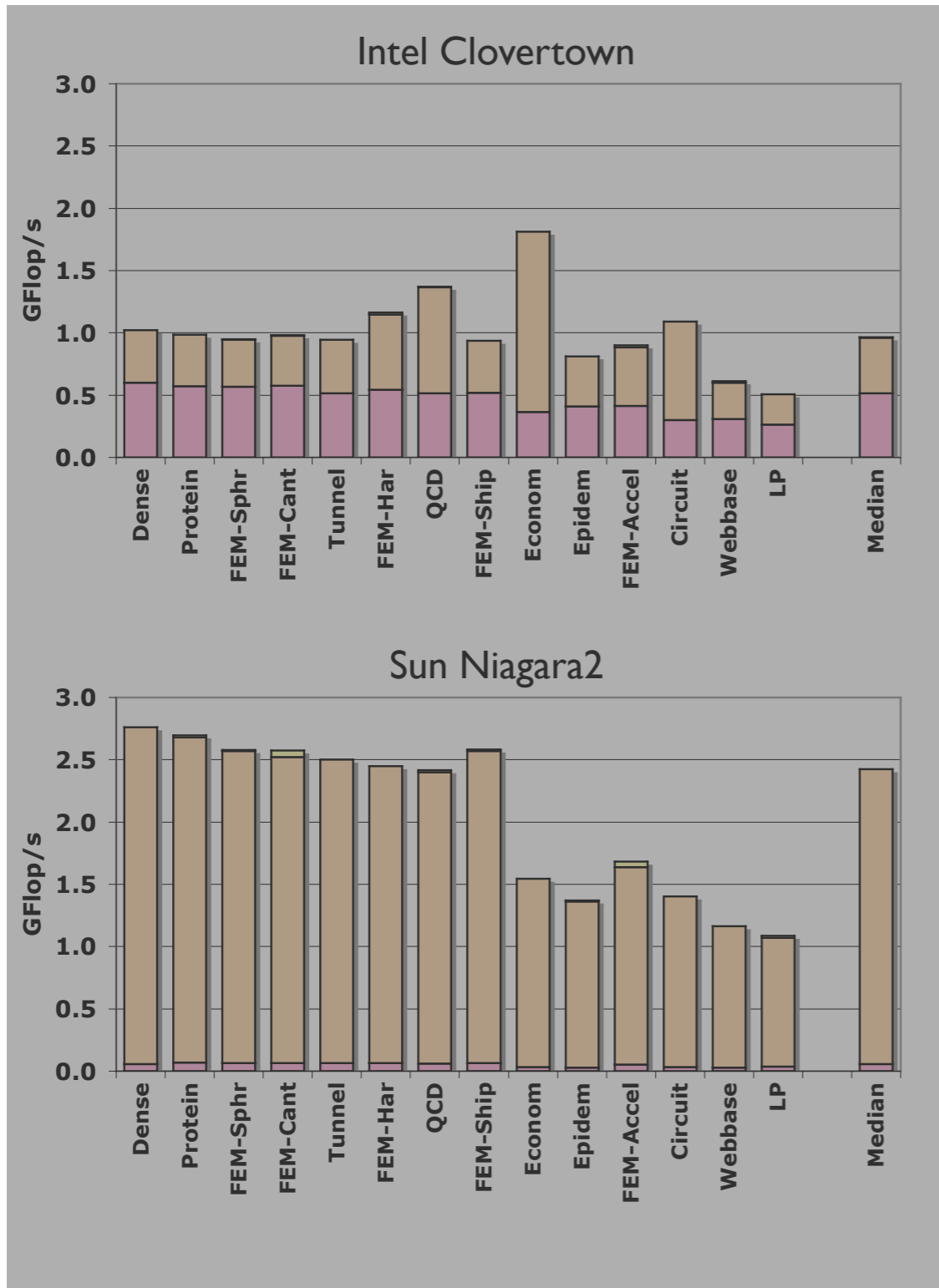
**Idea:**  
**Exploit NUMA, affinity**

Bind adjacent blocks to adjacent cores using

Cell: libnuma

Opteron: Custom, based on Linux scheduler

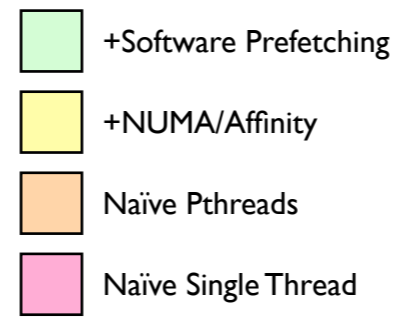
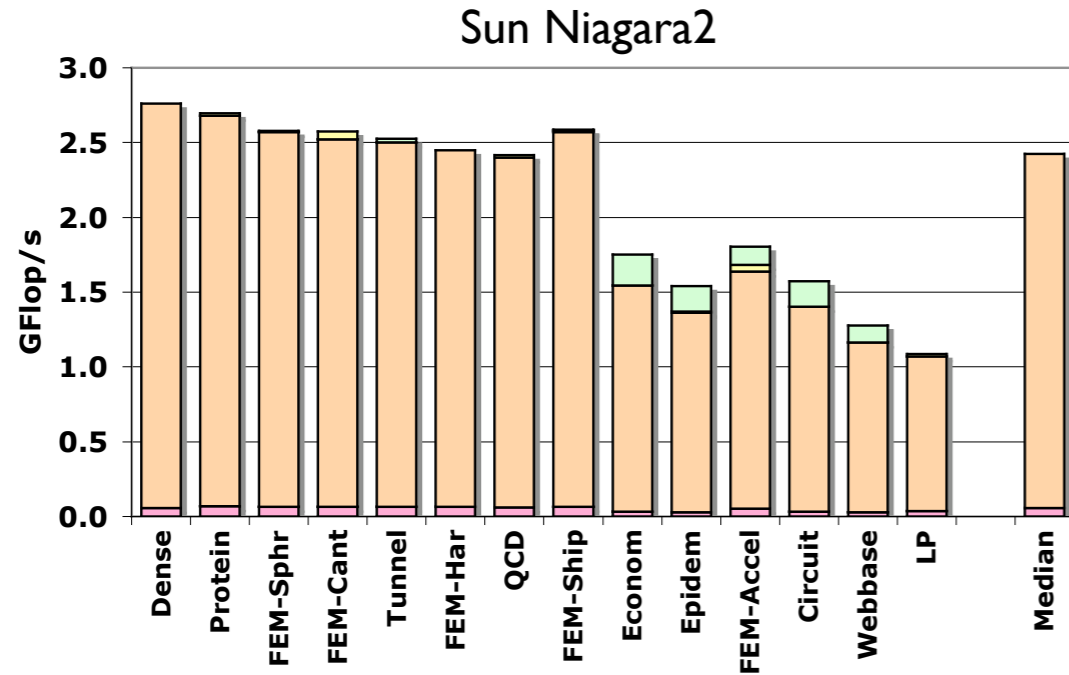
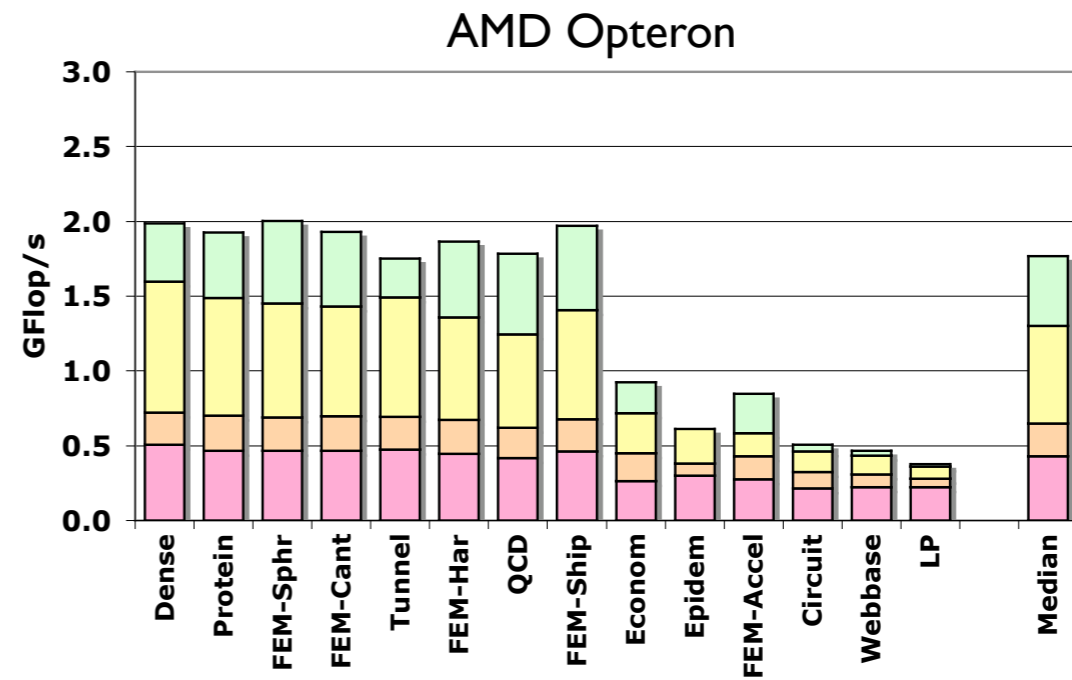
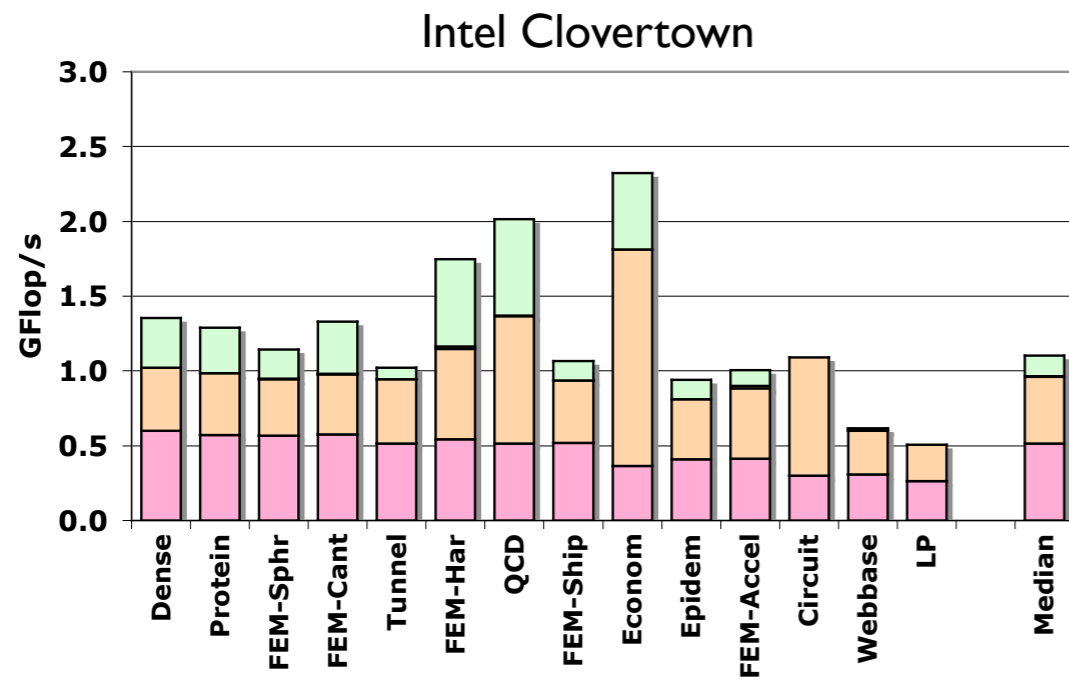
# Bind



# Prefetch

- ▶ Search over
  - ▶ all distances from 0 (none) to 1 KB
  - ▶ all combinations of target
- ▶ x86: Prefetch into L1, indicating temporal mode
- ▶ Niagara2: L2 only, but still works

# Prefetch

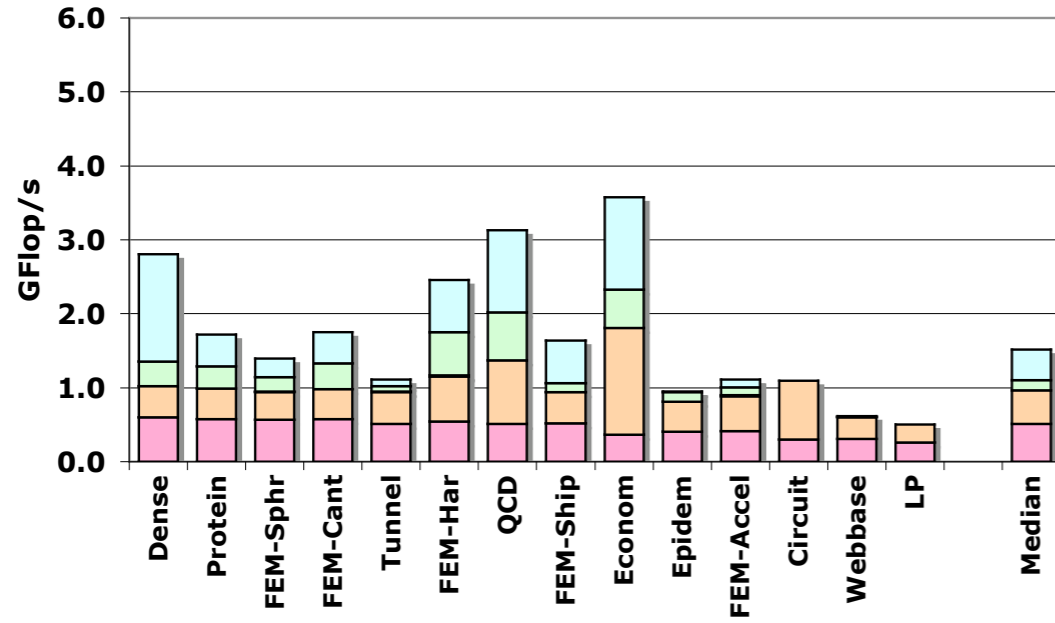


# Compress

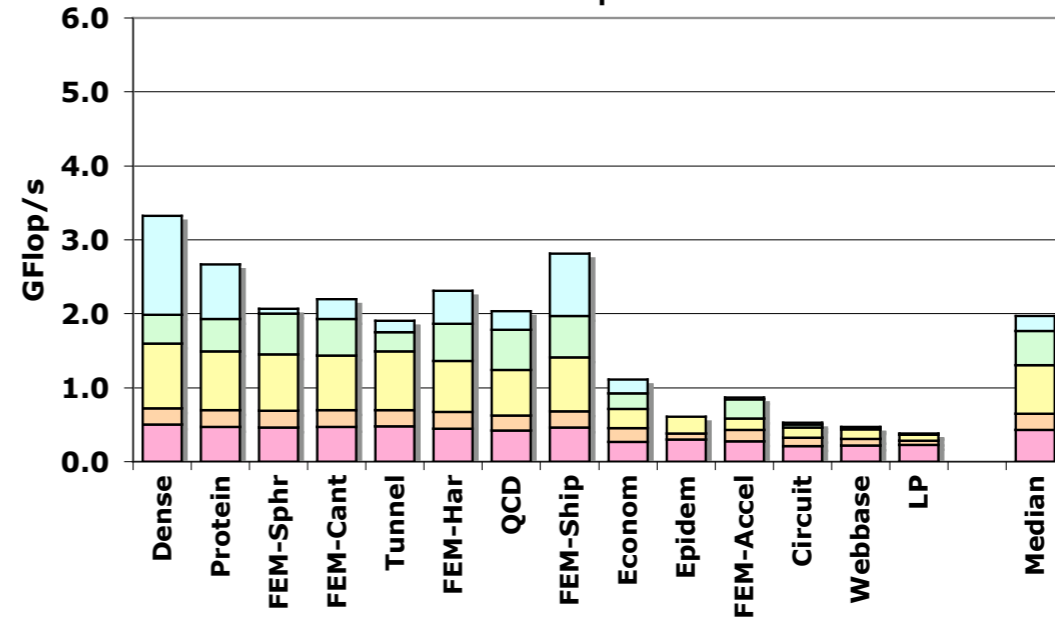
- ▶ Intuition: Reading matrix dominates run-time
- ▶ Reduce integer index volume and size
  - ▶ Power of 2 register blocking
  - ▶ BCSR / BCOO formats
  - ▶ 16-bit vs. 32-bit indices

# Compress

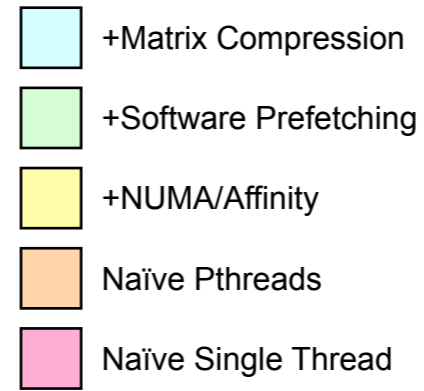
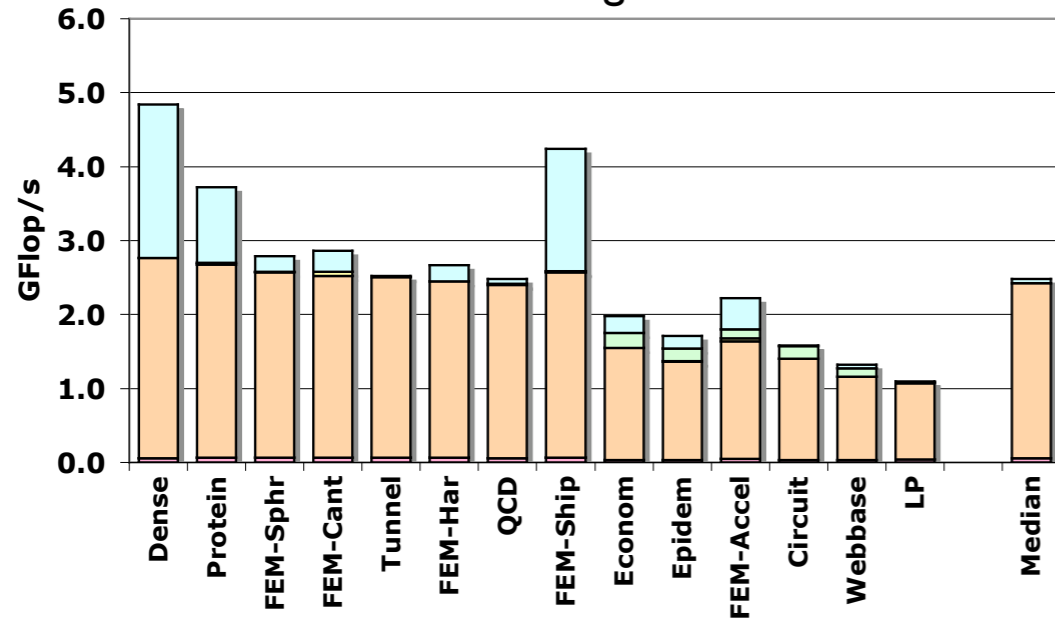
## Intel Clovertown



## AMD Opteron

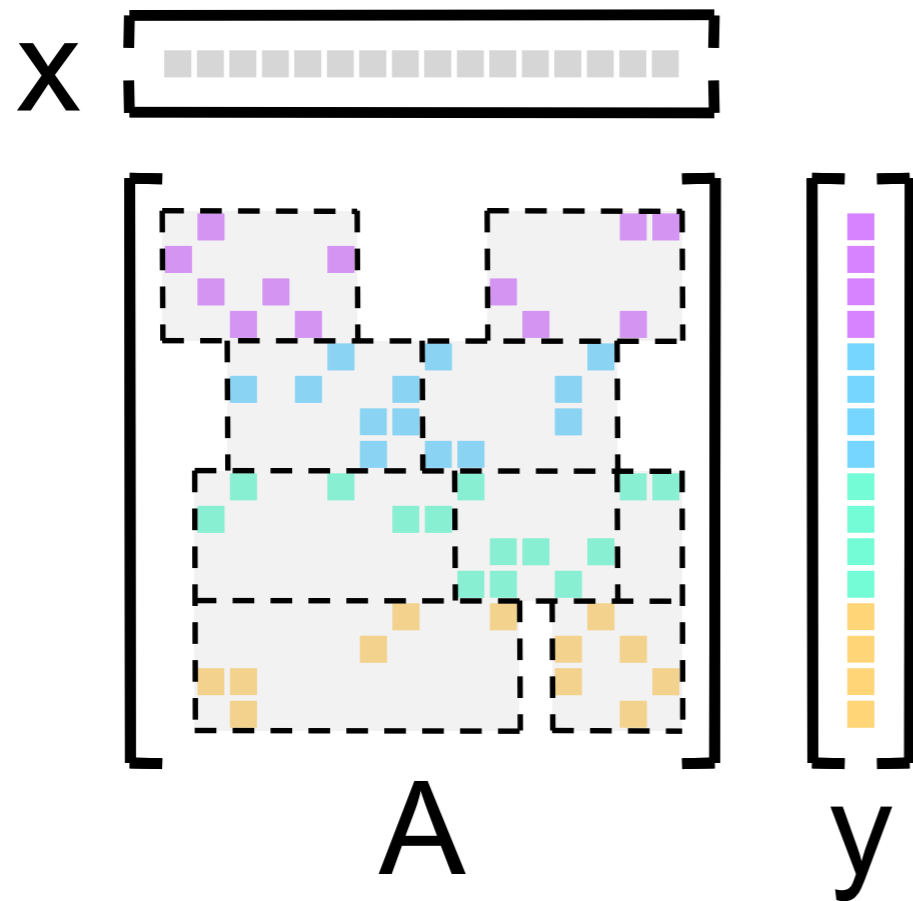


## Sun Niagara2





# Block (Cache & TLB)

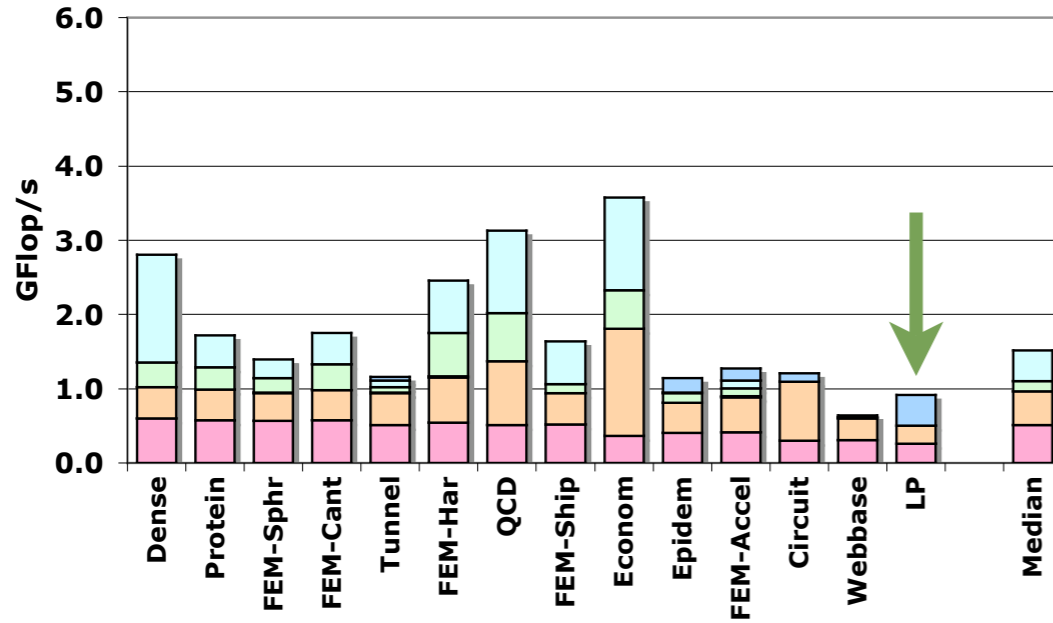


Like dense case,  
localizes 'x' loads for reuse.

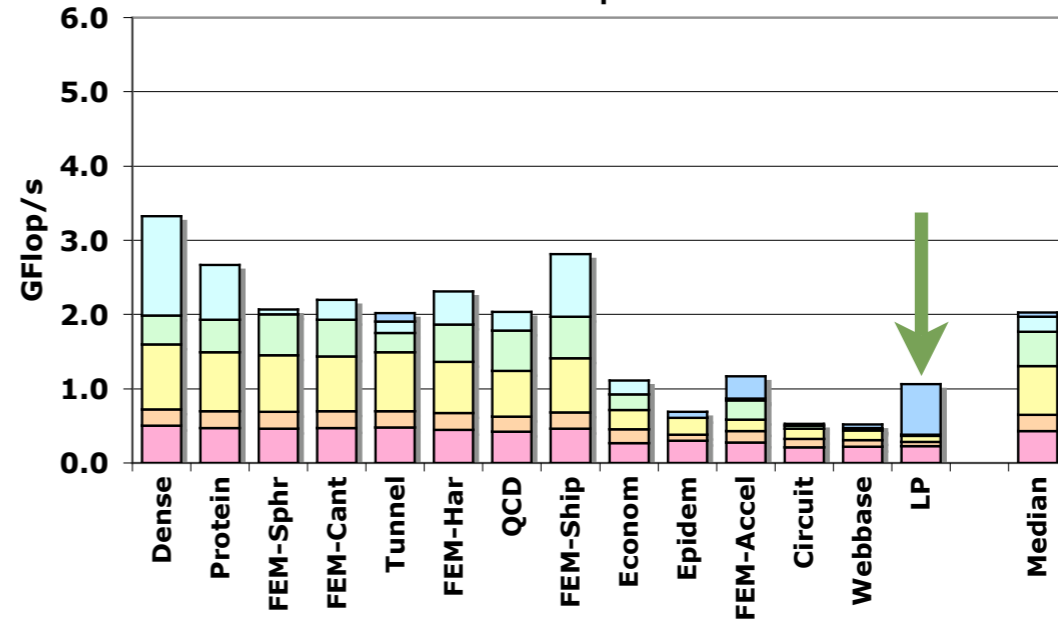
Partition into cache blocks.  
Tune each block.

# Block

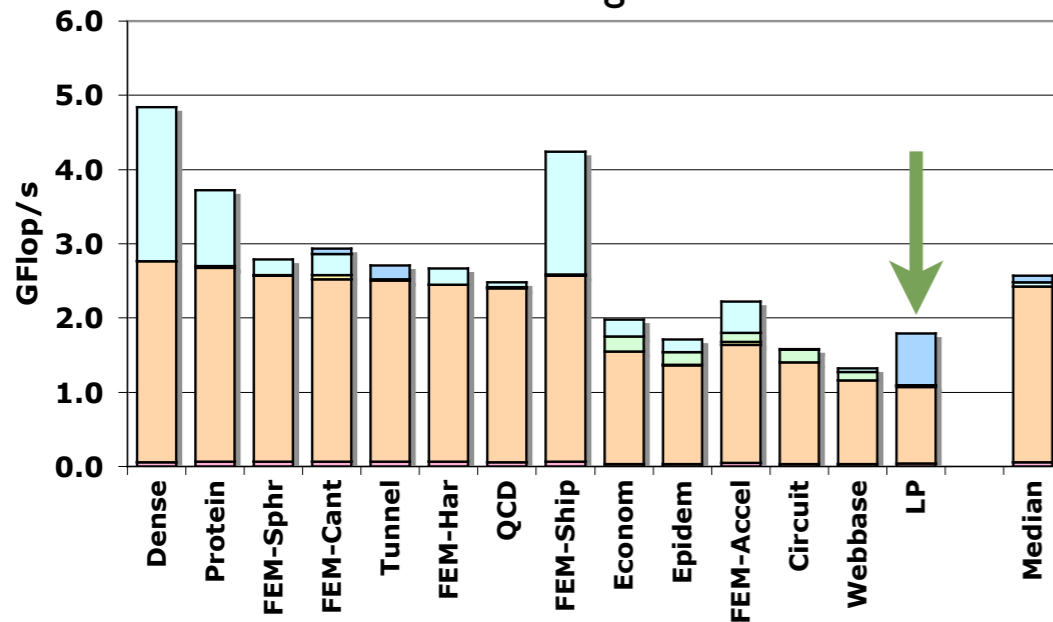
## Intel Clovertown



## AMD Opteron



## Sun Niagara2



- +Cache/TLB Blocking
- +Matrix Compression
- +Software Prefetching
- +NUMA/Affinity
- Naïve Pthreads
- Naïve Single Thread

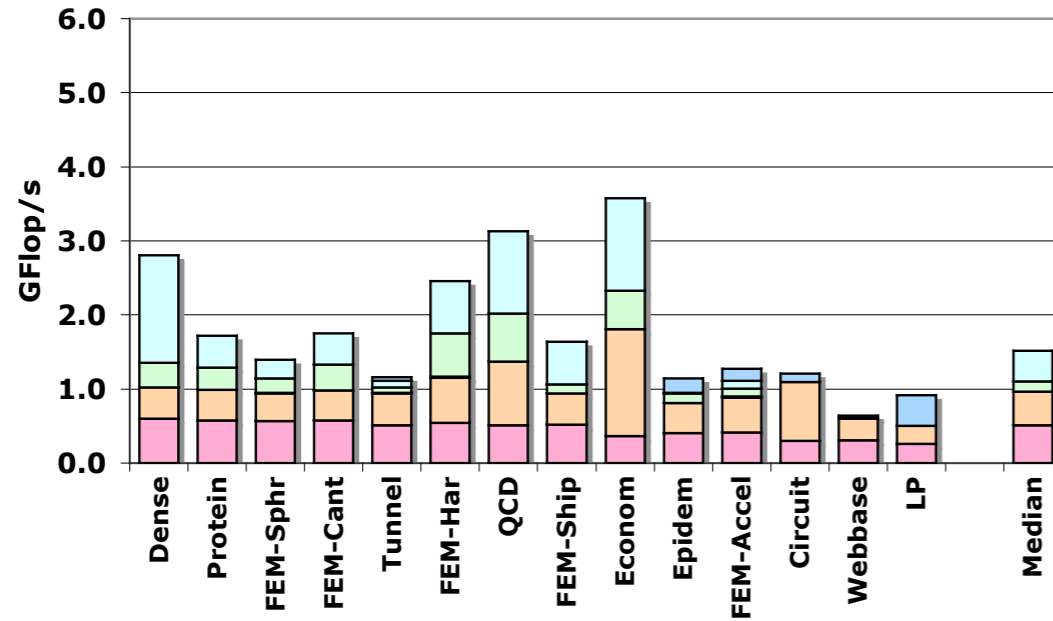
- ▶ Multicore optimizations
  - ▶ Bind
  - ▶ Prefetch
  - ▶ Compress
  - ▶ Block
  - ▶ **Buy (more memory)**

# Banks, ranks, and DIMMs

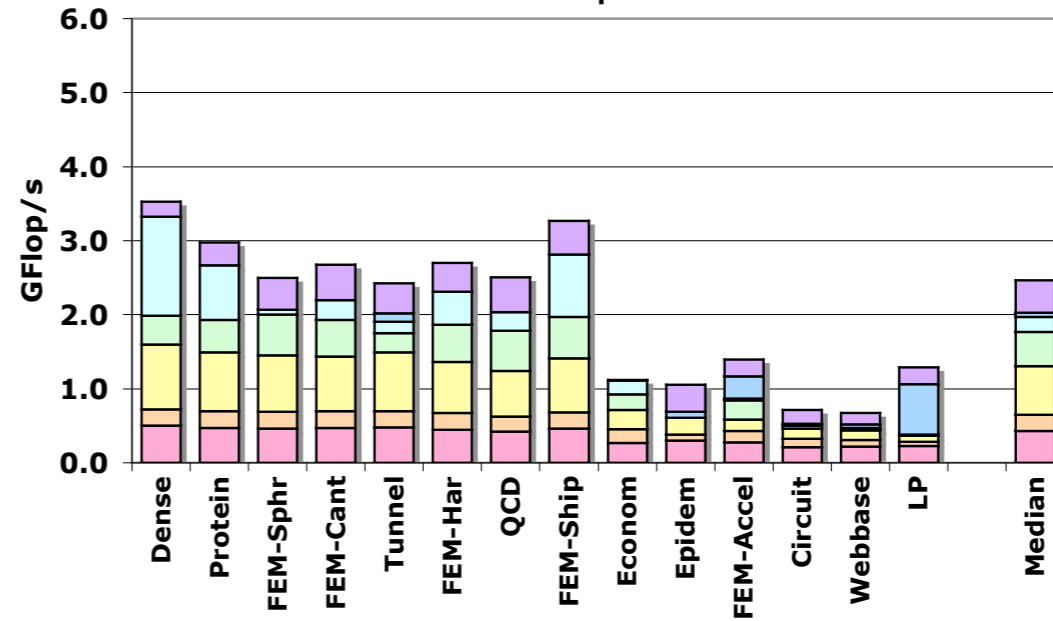
- ▶ More threads  $\Rightarrow$  more streams to memory
  - ▶ Memory controller reorders some no. of requests
  - ▶ Addressing/bank conflicts more likely
- ▶ Idea: Add more DIMMs, configure ranks
- ▶ (Our Clovertown was fully populated)

# Buy (DIMMs)

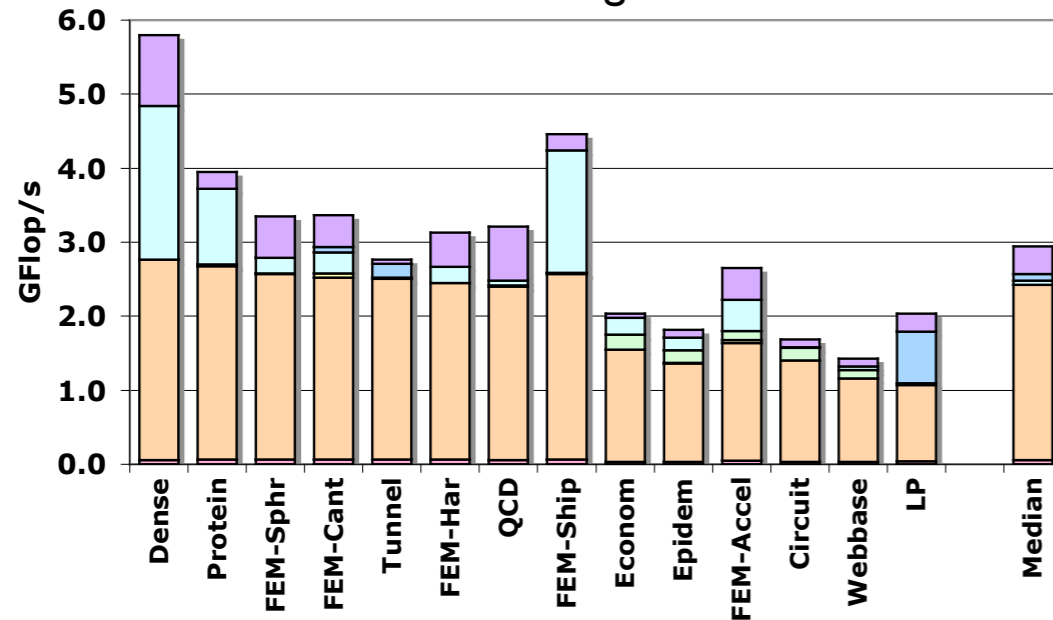
## Intel Clovertown



## AMD Opteron



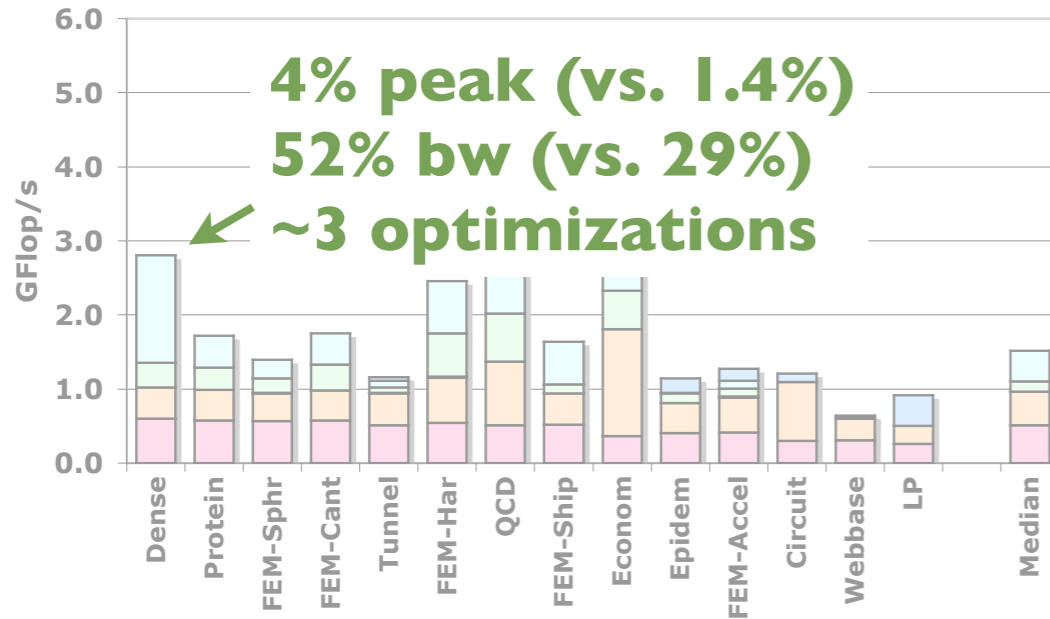
## Sun Niagara2



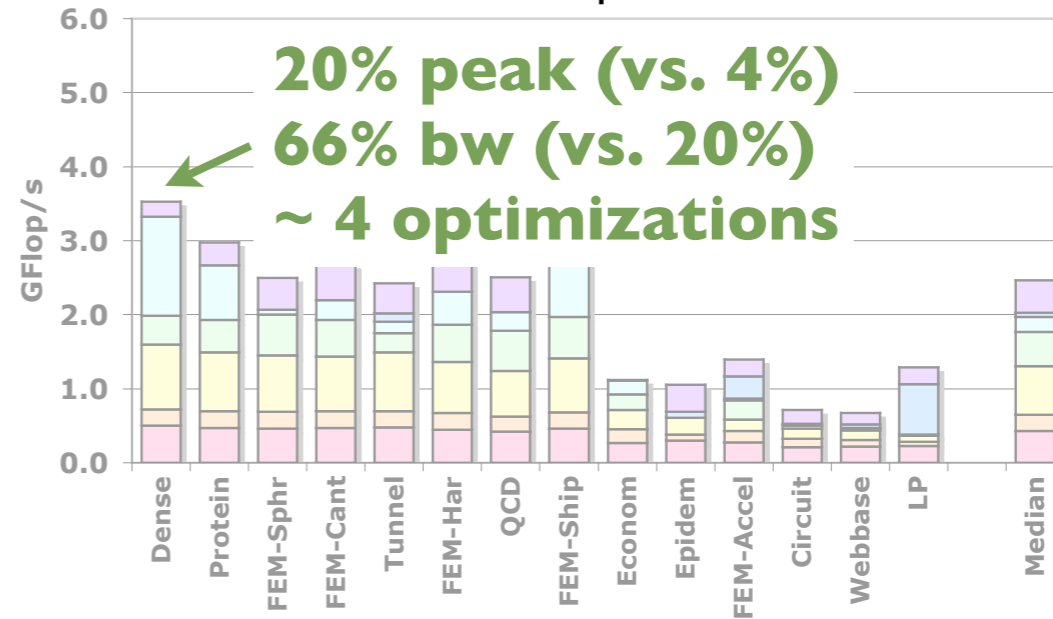
- +More DIMMs, Rank configuration, etc...
- +Cache/TLB Blocking
- +Matrix Compression
- +Software Prefetching
- +NUMA/Affinity
- Naive Pthreads
- Naive Single Thread

# Where do we stand?

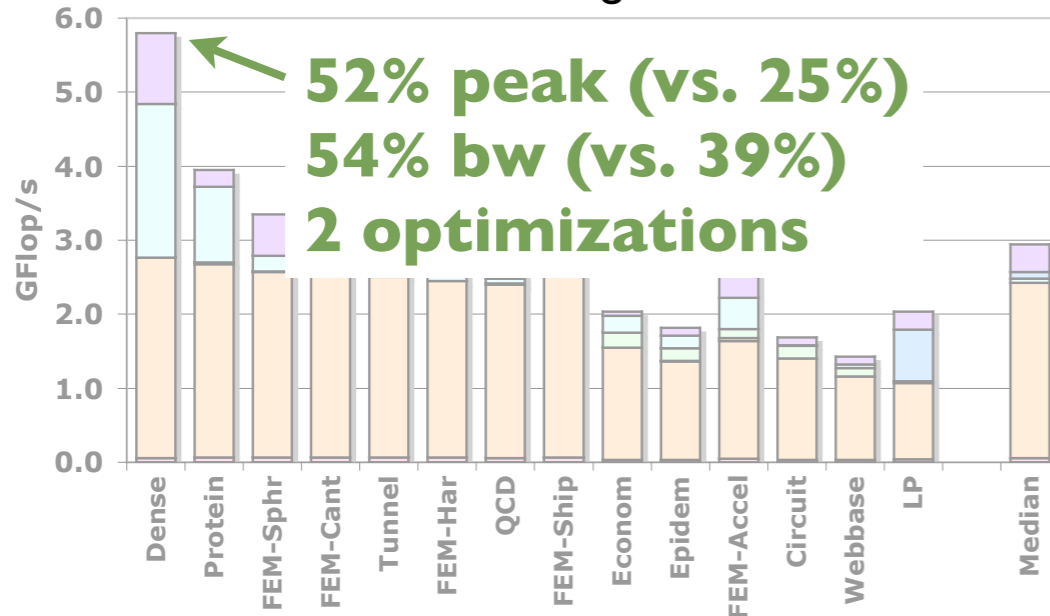
Intel Clovertown



AMD Opteron



Sun Niagara2

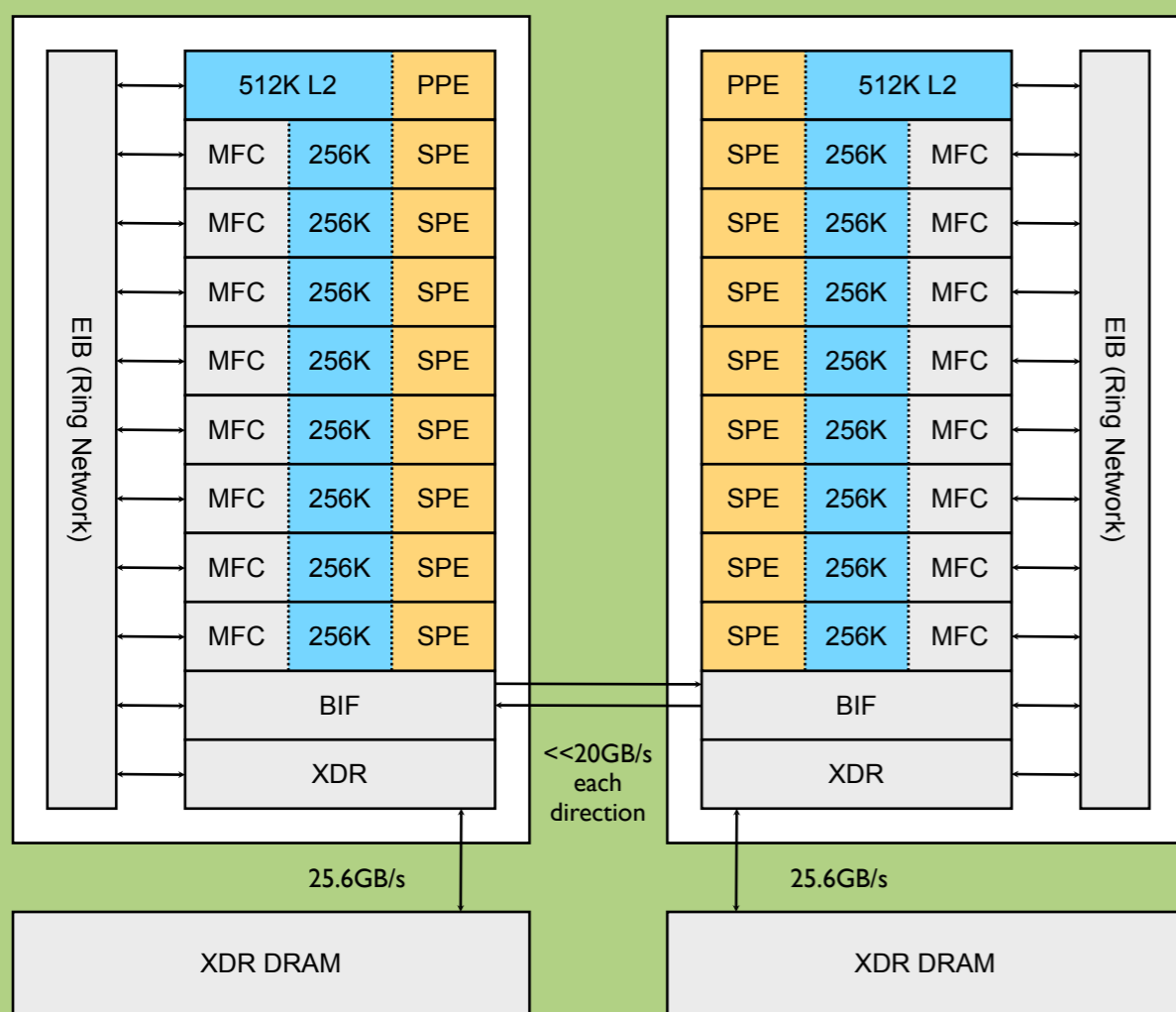


- +More DIMMs, Rank configuration, etc...
- +Cache/TLB Blocking
- +Matrix Compression
- +Software Prefetching
- +NUMA/Affinity
- Naïve Pthreads
- Naïve Single Thread

▶ Cell

- ▶ Consider double-precision, even if slow (14x)
- ▶ Cache blocking  $\Rightarrow$  Local-store blocking
- ▶ Forced 2x1 blocking  $\Rightarrow$  extra traffic

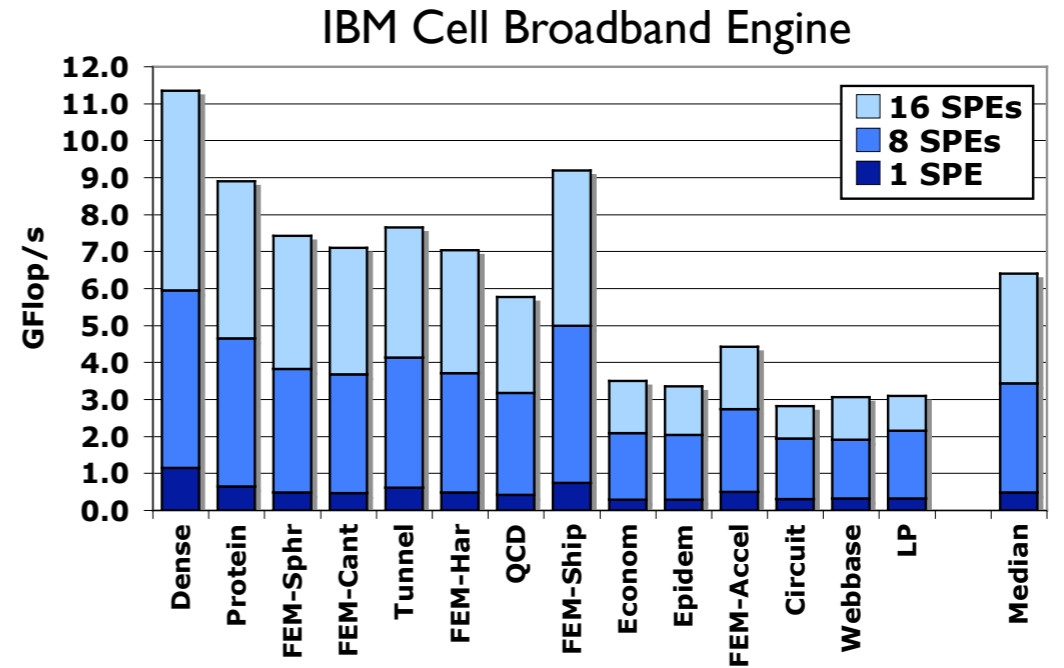
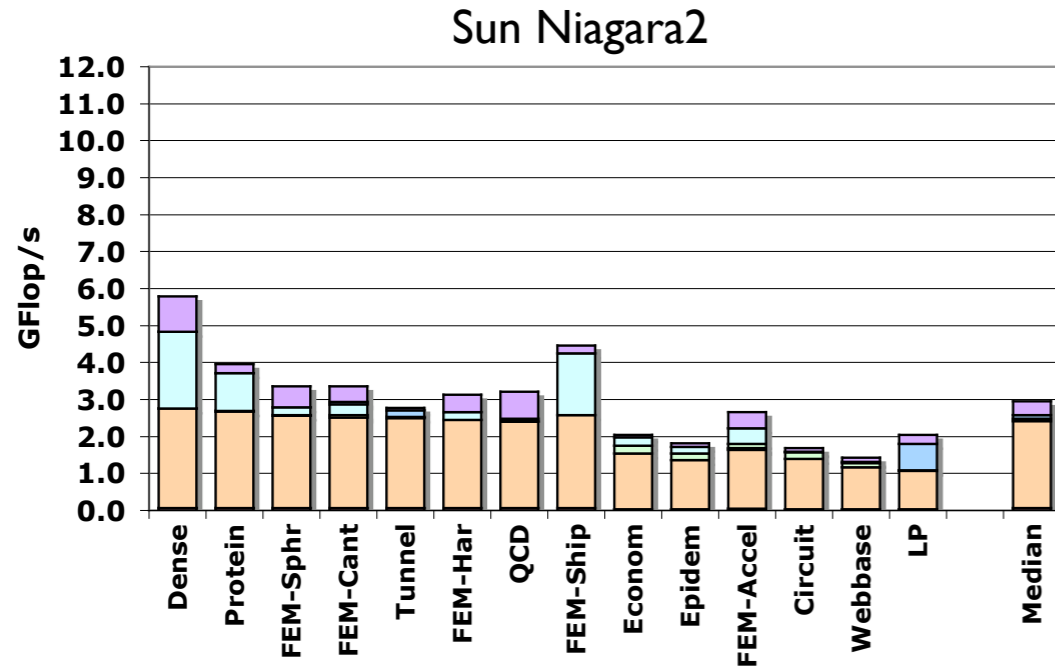
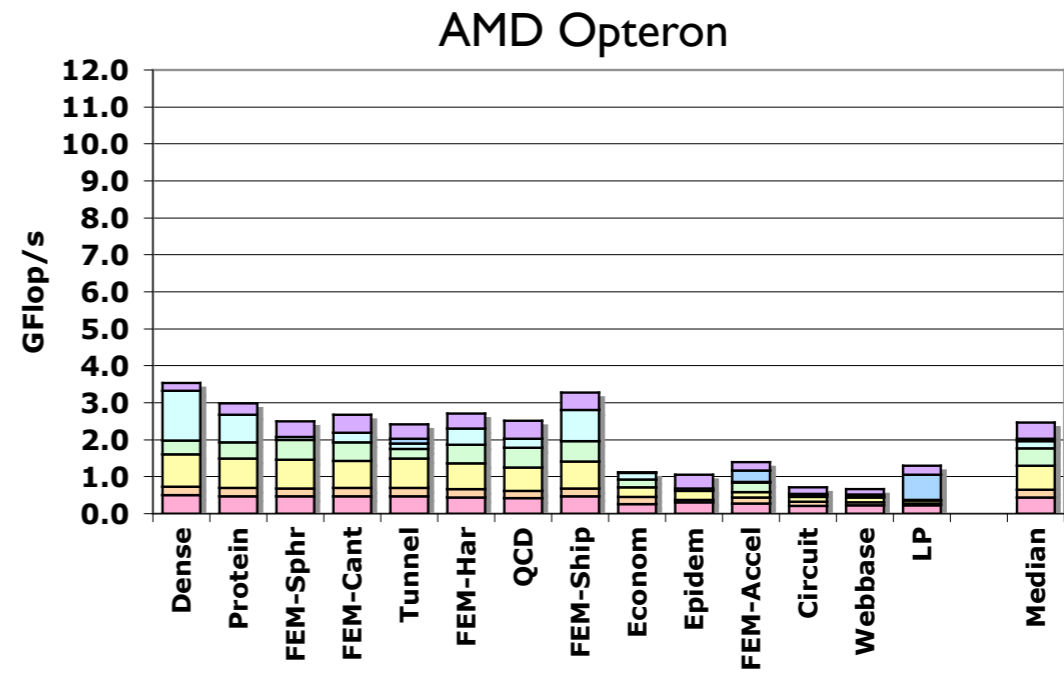
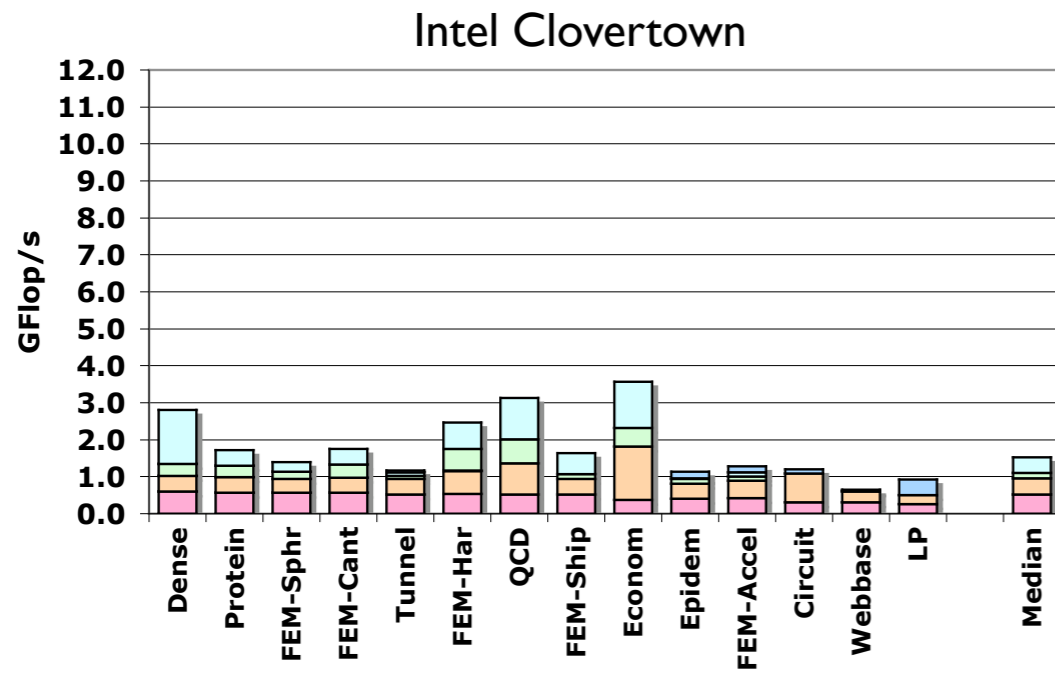
# Cell

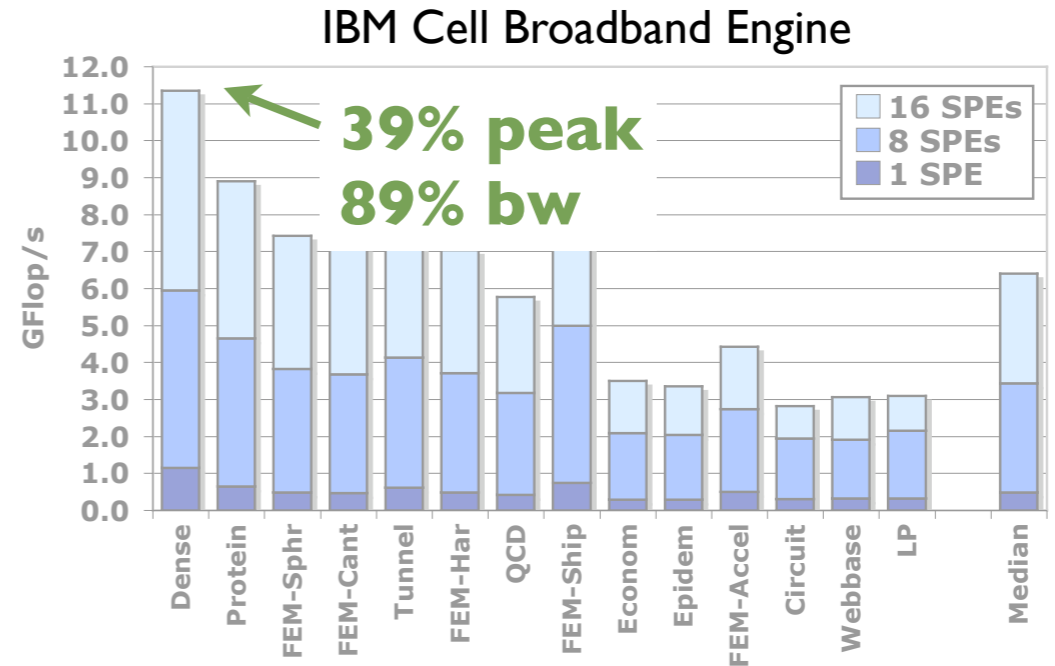
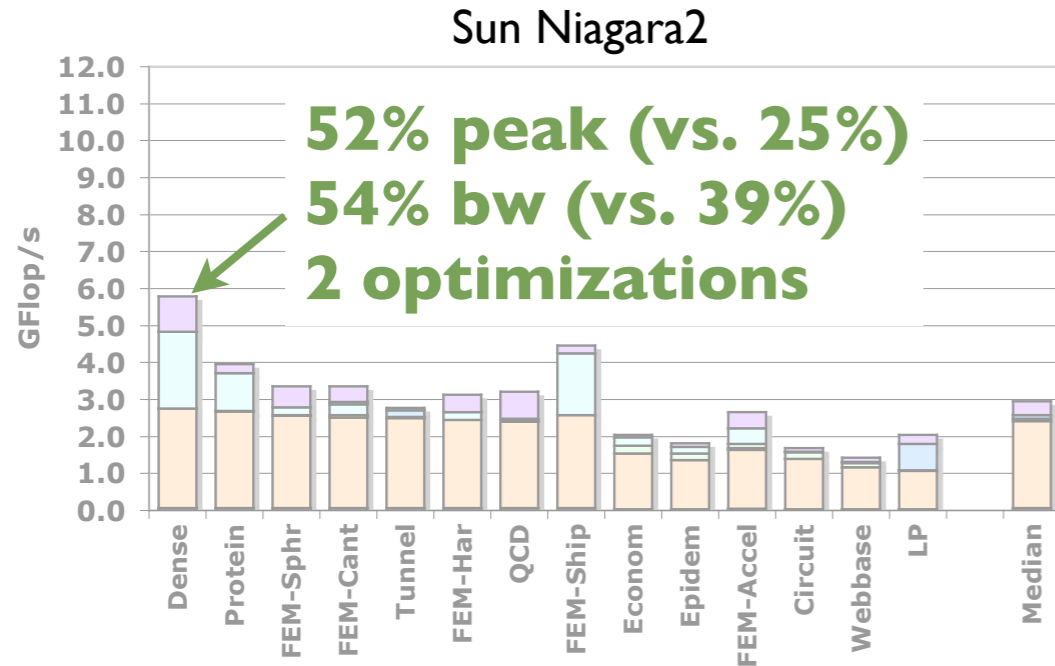
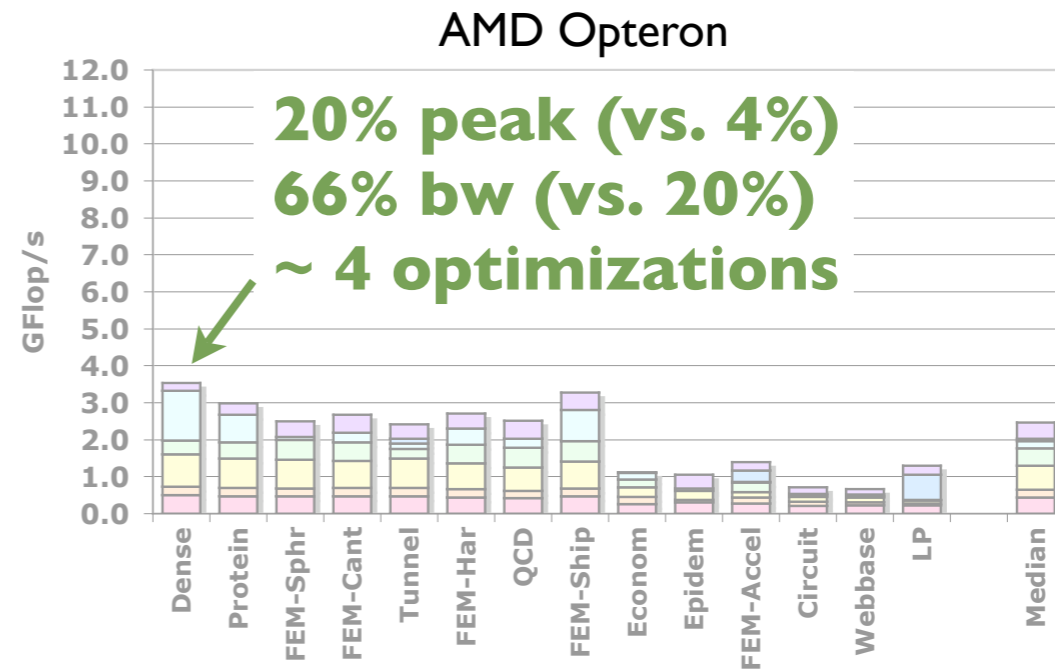
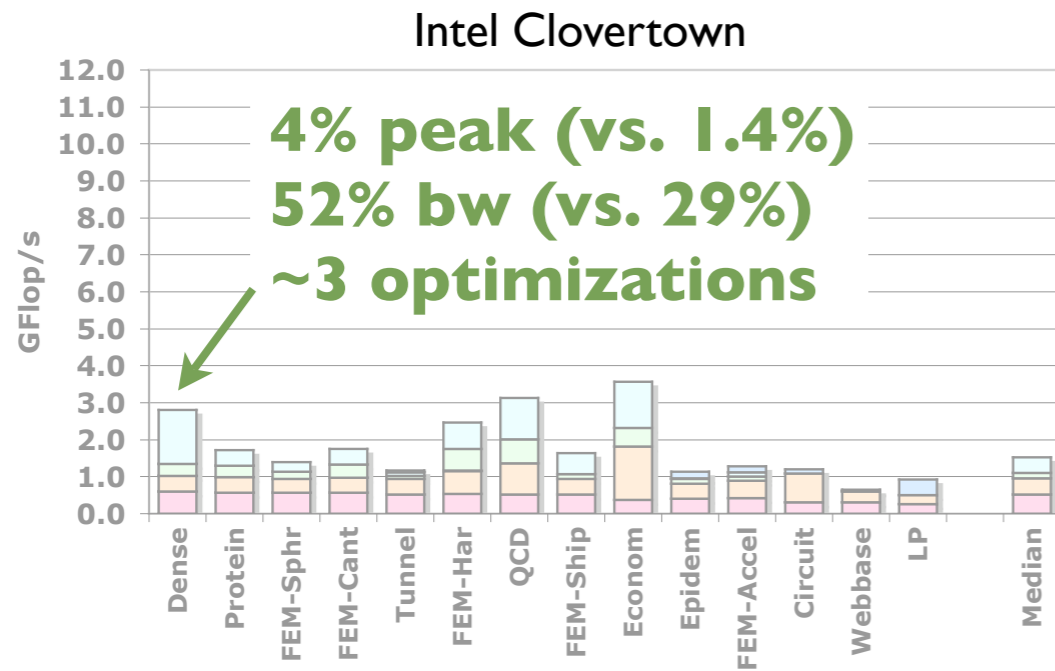


STI Cell Blade (2 x 8)

- ▶ Single vs. double: 14x
- ▶ Cache-block  $\Rightarrow$  Local-store
- ▶ SIMD forces 2x1 blocking

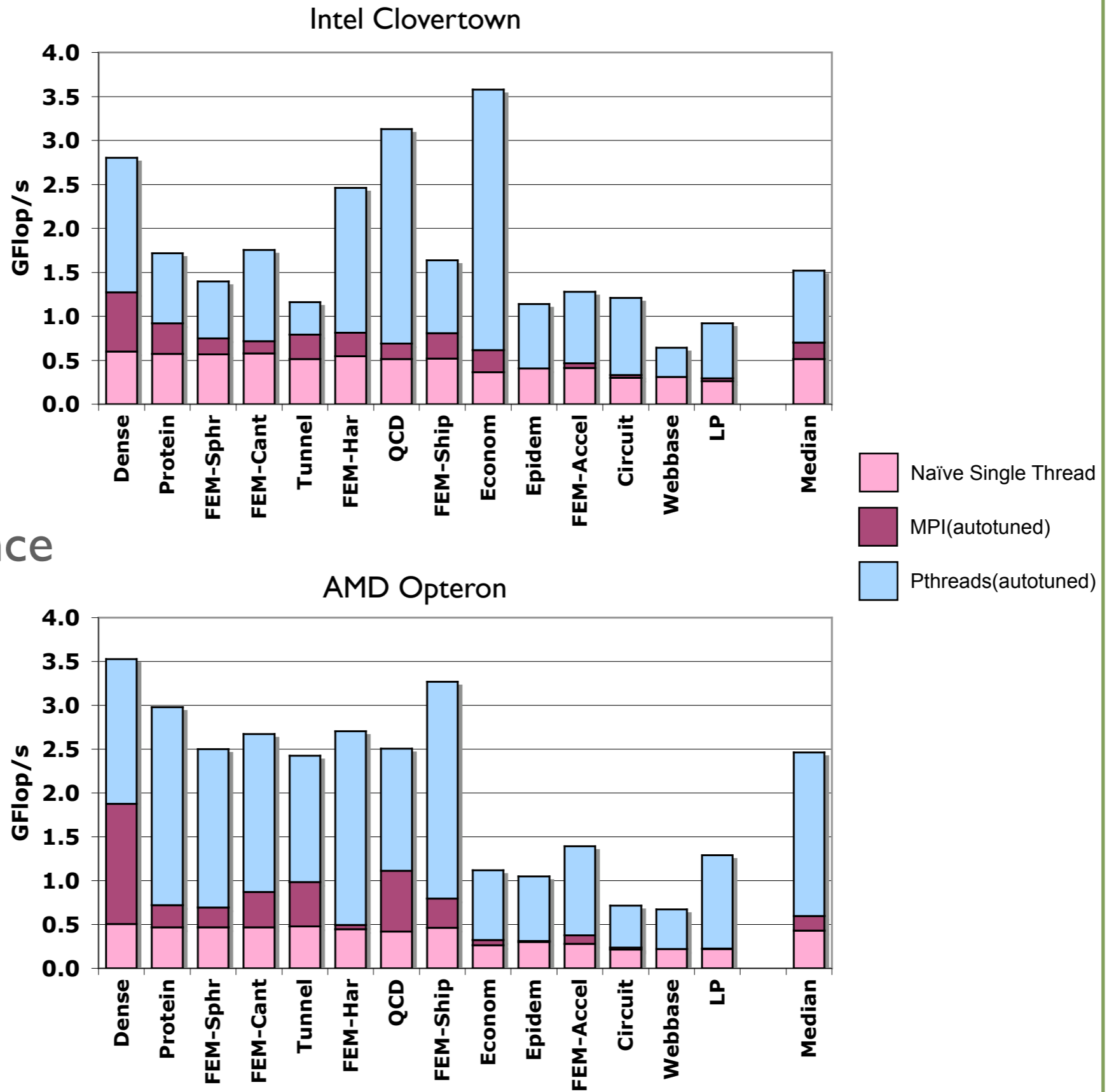






- ▶ An “off-the-shelf” approach
  - ▶ OSKI (tuned serial)
  - ▶ + PETSc (parallel MPI)
  - ▶ + MPICH (shmem)

**Issues:**  
Copies  
Load imbalance



# Review

- ▶ Key conclusions
  - ▶ The most complex architecture needed most tuning yet achieved lowest performance
  - ▶ 50–60% DRAM bandwidth vs. 90% on Cell
  - ▶ Place, Prefetch, Compress, Block (& Buy, Tune)
- ▶ For power/performance issues, see SC'07 paper
- ▶ Status and future work
  - ▶ Integration into OSKI (heuristics)
  - ▶ Improved “off-the-shelf” comparisons (MPI+OMP)