

# Node-wide Performance (and Power) Introspection Applied to Scheduling

Robert Fowler, Anirban Mandal, Allan Porterfield  
RENCI, UNC-Chapel Hill

renci

RESEARCH \ ENGAGEMENT \ INNOVATION

# Adaptive systems and applications are the future.

- Outline

- Hardware issues:
  - The constrained evolution of processor design.
  - The hardware already adapts itself.
- Off-chip bottlenecks: memory (and I/O).
- Whole node measurement in support of introspection.
- Won't actually get to adaptive, introspective scheduling.

# Dennard Scaling of CMOS Logic.

- Series of papers 1972-1974 by Bob Dennard and others at IBM on scaling properties of CMOS logic circuits (gates and wires!).
- Linear scaling of all transistor parameters.
  - Reduce feature size by a factor of  $S$ . Typically, 0.7/generation.
    - Including gate insulator thickness!
  - Reduce supply voltage ( $V_{dd}$ ) by  $S$  to keep electric field constant.
  - Adjust doping of silicon gate region to compensate.
- Consequences
  - Area shrinks by  $S^2$ ,  $C_{gate}$  and delay ( $1/f$ ) reduced by  $S$ .
  - Power  $\approx CV^2f \rightarrow$  Power per gate goes down by  $S^2$
  - Area and power track each other so power density is unchanged.
  - For a constant die area and design density, power and power density are constant and frequency can be increased.

# Other Aspects of Dennard Scaling.

- Wire resistance/unit length  $\sim S^2$
- Wire capacitance/unit length  $\sim 1$
- RC delay/unit length (unrepeated)  $\sim S^2$
- Die size (D) increases, so “long” wires increase by D
- Unrepeated wire delay  $\sim S^2 D^2$ , repeated  $\sim D \sqrt{S}$   
→ Signals cannot cross the chip in one cycle.

# Moore's law

Empirical observation and self-fulfilling prophecy:

Circuit element count doubles every N months. (N ~18)

- Technological explanation: Features shrink, semiconductor dies grow.
- Dennard scaling: Gate delays decrease. Wires are relatively longer/slower.
  - Dennard scaling has not been perfect in practice and is coming to an end.
- In the past, the focus has been making "conventional" processors faster.
  - Faster clocks
  - Clever architecture and implementation → instruction-level parallelism.
  - Clever architecture (speculation, predication, etc), HW/SW Prefetching, and massive caches ease the "memory wall" problem.
- Problems:
  - Faster clocks --> more power.
  - Power scaling law for CMOS:  $P = \alpha CV^2F$ , but  $F_{\max} \sim V$  so  $P \sim F^3$ 
    - Where  $\alpha$  is proportional to the avg. number of gates active per clock cycle.
  - Smaller transistors + long wires → either slow clock, or pipelined communication.
  - More power goes to overhead: cache, predictors, "Tomasulo", clock, ...
  - Big dies --> fewer dies/wafer, lower yields, higher costs
  - Aggregate effect --> Expensive, power-hog processors on which some signals take 6 cycles to cross.

# The End of Dennard Scaling: Dark and Dim Silicon

- **V<sub>dd</sub> Scaling issues**
    - Initially, designers constrained by standards: 12V, 5V, 3.3V.
    - On-board power regulation now allows V<sub>dd</sub> to be 1V or less.
    - This is getting uncomfortably close to threshold voltages.
    - Decreasing thresholds has rapidly increased leakage current/power.
    - Decreasing f allows operation with higher thresholds.
  - **Gate Insulator issues**
    - Thickness is now ~ 5 atoms
  - **Useful work and duty cycles**
    - Bailey and Snyder (1988) observed that  $\alpha$  was at most a few percent for processors. If  $\alpha$  were much larger, chips would melt.
    - Aggressive architectures have increased  $\alpha$  to do bookkeeping, data movement, ...
  - **“Dark” and “dim” silicon refer to schemes to reduce  $\alpha$  and/or f to reduce power.**
    - Heterogeneous cores and purpose built modules w. power management.
    - Programmable logic and reconfigurable devices.
- ➔ We can now build chips that cannot be run at their full design potential.

# “Dim silicon” adaptation in X64\_64

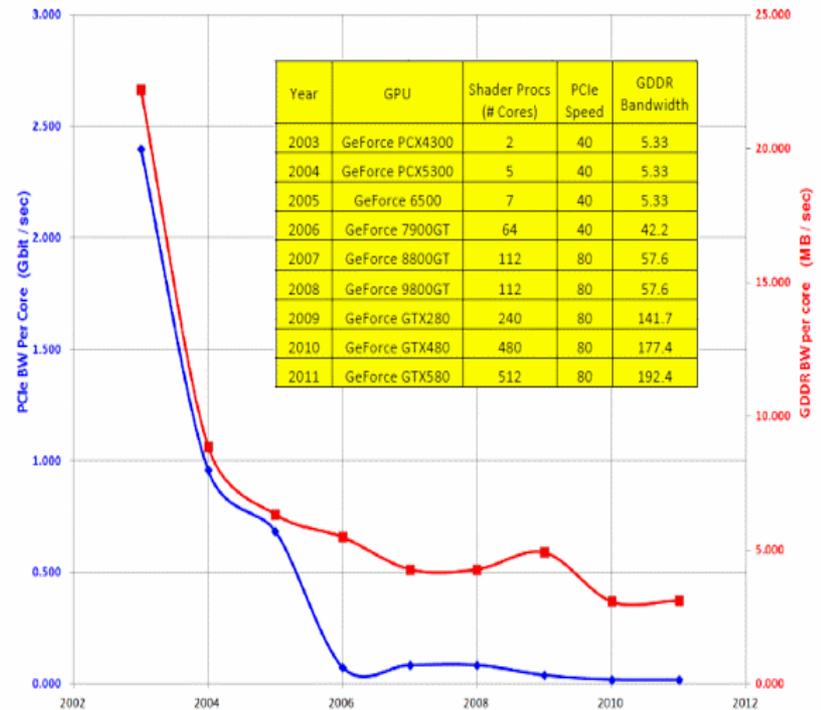
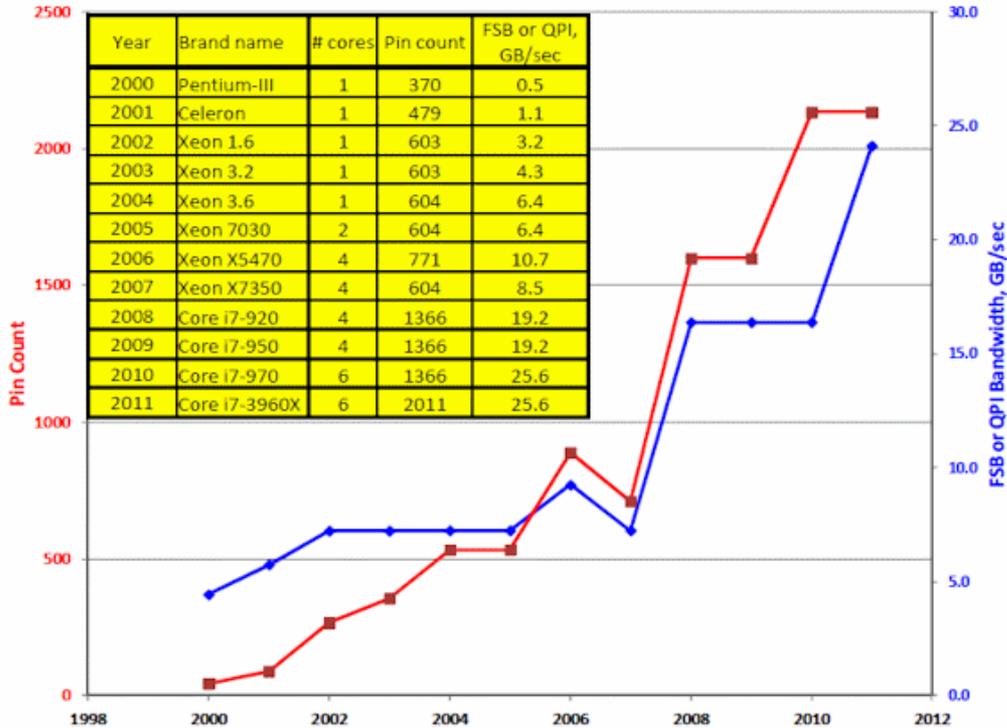
- Intel: On-chip control processor
  - “Turbo” modes throttle f when all cores are active.
    - Run power-efficient, low f, low V in highly parallel code regions.
    - Inefficient high f, high V in sequential regions.
  - Shut down cores
  - Automatic DVFS to keep within thermal budgets.
    - Aggressive turbo mode when one core is active.
    - Brief performance bursts when thermal headroom exists.
  - “Balanced” mode DVFS to trade FLOPS and Watts.
    - Reduce core frequency when code seems to be memory bound.
- AMD:
  - Similar strategies
  - Additionally, A-series can attempt to balance CPU and GPU performance and power budgets.

## Moore's Law/Dennard Scaling for DRAM.

- As more transistors were added to processor chips, they got a lot faster.
  - Dennard scaling for faster transistors.
  - Clever architectures and on-chip concurrency.
- As more transistors were added to memory chips, they got a lot bigger.
  - Cleverness went into reliability, yield, ...
  - Small transistors are fast, but weak (can't drive long wires).
  - Little increase in on chip concurrency.
  - Very low Rent's law (surface/volume ratio) exponent!

	Introduction	Size	Pins	Cycle Time	Bandwidth
DDR	2000	2 GB	168	5 ns	3.2 MB/sec
DDR2	2003	4 GB	184	3.75 ns	8.5 MB/sec
DDR3	2007(2009)	16 GB	240	5 ns	12.8 MB/sec
DDR4	2012(?)				25.6(?) MB/sec

# Other Trends: Pins and GPU Memory



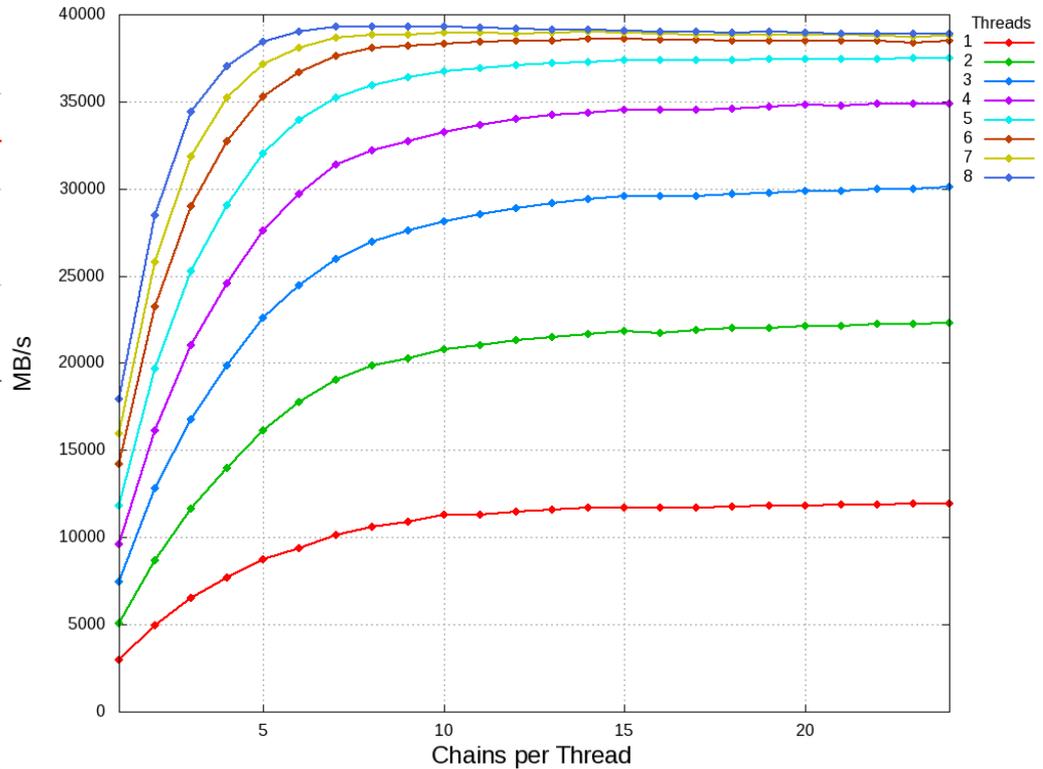
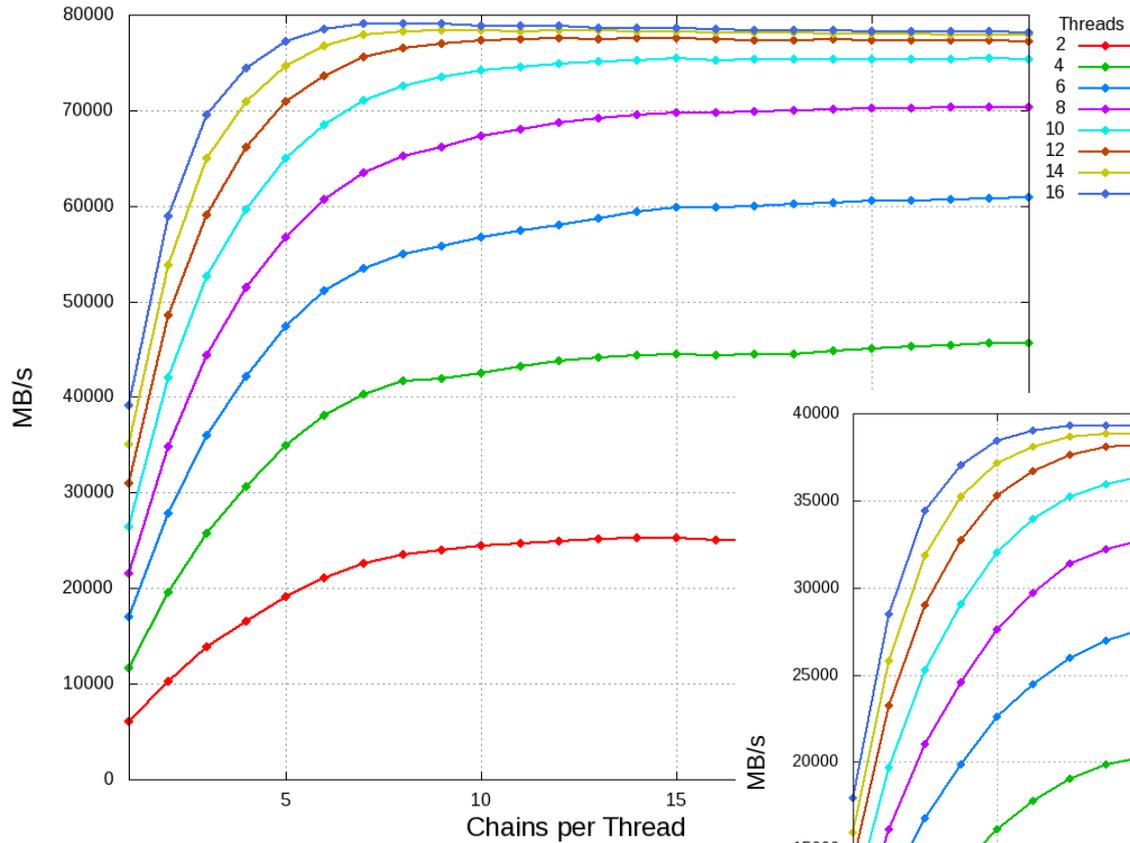
# Little's Law applied to Memory.

- Classic law/lemma in queuing theory
    - (mean # in system/queue) = (arrival rate) (mean residence time)
  - Communication (memory) restatement
    - (concurrency) = (bandwidth) (latency)
- To increase bandwidth without decreasing latency, you have to increase the concurrency of the system
- Wider channels to send more bits per operation.
  - Overlapping, i.e., pipelined, operations.
- Bottleneck → bandwidth plateaus, queuing latency dominates.

# pChase

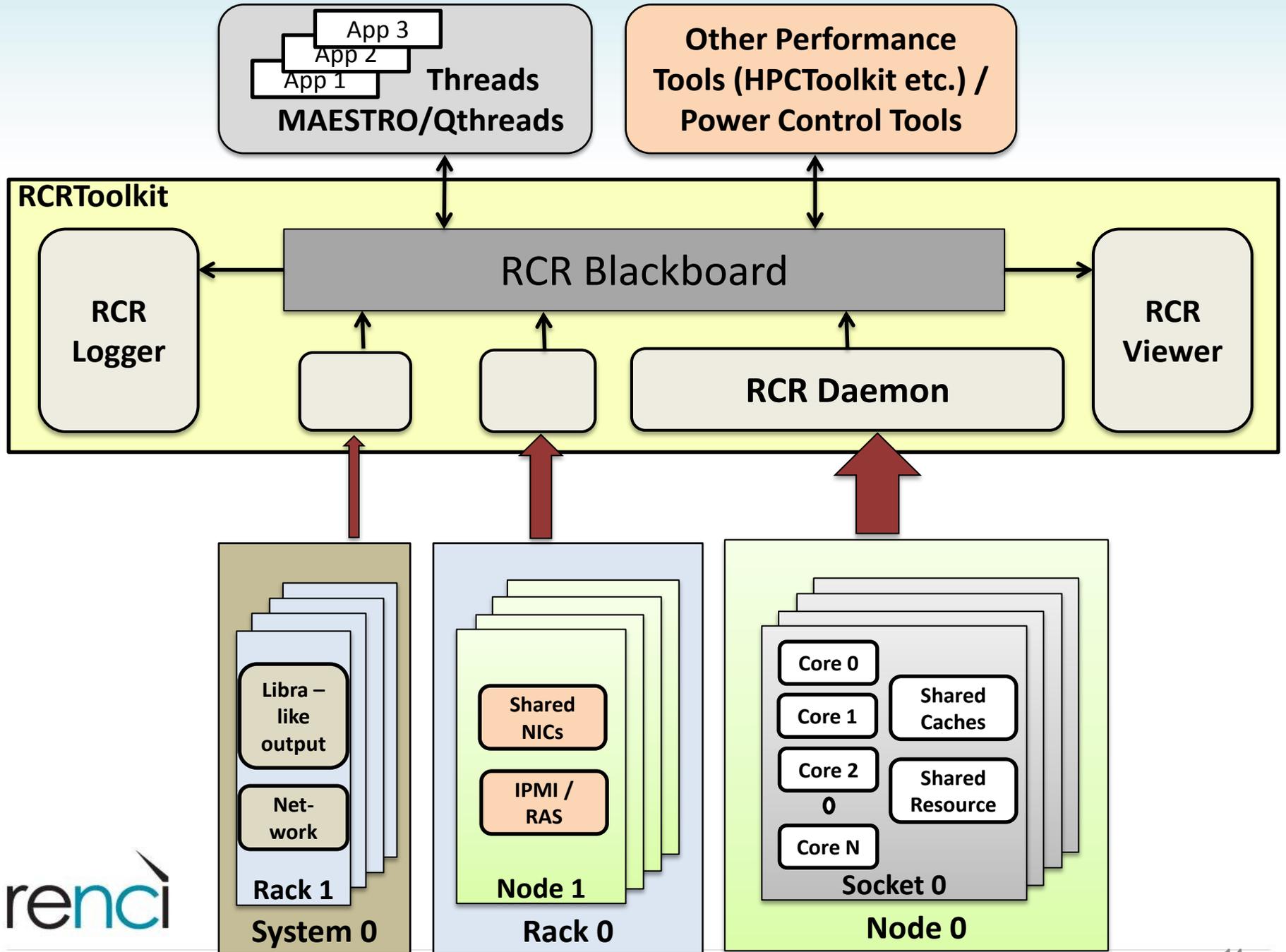
- Developed by Pase and Eckl @IBM
- Multi-threaded benchmark used to test memory throughput under carefully controlled degrees of concurrent accesses
- Each thread executes a controllable number of ‘pointer-chasing’ operations – a memory-reference chain
  - Pointer to the next memory location is stored in the current location. Grow and randomize chain to defeat cache, prefetch.
  - Dereference pointers in k independent chains concurrently, then use them.
- K=1 case measures memory latency.
- Large-k bandwidths are comparable to STREAM measurements at “common” optimization levels.
- Our Modifications
  - Added wrapper scripts around pChase to iterate over different numbers of memory reference chains and threads
  - Added affinity code to control thread and data placement
- Available at <http://pchase.org>

# E5-2680 in “maxperf” mode.



# Resource Centric Performance Reflection: RCRToolkit

- Performance measurement and analysis tool that focuses on shared resources in a system
  - Information and analysis should help applications and system code introspectively, in real-time, to adapt to bottlenecks, power, thermal events.
- RCRToolkit consists of
  - RCRblackboard
  - Several clients for the RCRblackboard
- RCRblackboard
  - Shared memory region (or, currently Google protocol buffers resident in memory) for real time use by producers and consumers of node- and system-wide performance information
  - Information organized in a hierarchy that reflects hardware structure
  - Coordination managed by the RCRblackboard protocol – multiple regions, each owned by a single writer.



# Core RCRblackboard Clients

- **RCRdaemon**
  - Uses Linux MSR driver in conjunction with configuration information accessible through libpfm4 (and other sources) to access off-core HPM counters.
  - Off-core counters in PCI CONFIG space coming soon.
  - Can co-exist with tools using Linux PerfEvents using on-core counters.
  - Configuration file specifies which HPM events to monitor, as well as a set of derived measures (“meters”) computable with simple arithmetic.
- **RCRlogger and RCRviewer**
  - Read RCRblackboard information and output a log for post-execution analysis
  - On-line monitoring.
- **MAESTRO thread scheduler**
  - Adaptive, locality-aware scheduling of over-partitioned applications, e.g., OpenMP tasks, loops with guided self-scheduling.
  - Monitors shared resource usage and adapts scheduling during periods of high utilization
  - Writes scheduling decision information to RCRblackboard, including some source attribution at the level of OpenMP loops and tasks.

# Hot-wiring HPCToolkit to leverage RCRToolkit

- Experimental extension to HPCToolkit to act as an RCRblackboard client
- Goal: combine third-person, system-wide performance measures with “first-person” call stack profiling.
- Accomplished by
  - Command line extension to specify events to split based on RCR predicates.
  - Modifying HPCToolkit’s event sampling code to examine RCRblackboard for existence of shared resource bottlenecks.
  - Monitoring whether the sampled event occurs when shared resource utilization exceeds threshold set in RCRdaemon configuration
- Original HPCToolkit event is split into two sub-events that can be viewed using an unmodified HPCViewer

# Hot-wiring HPCToolkit (2)

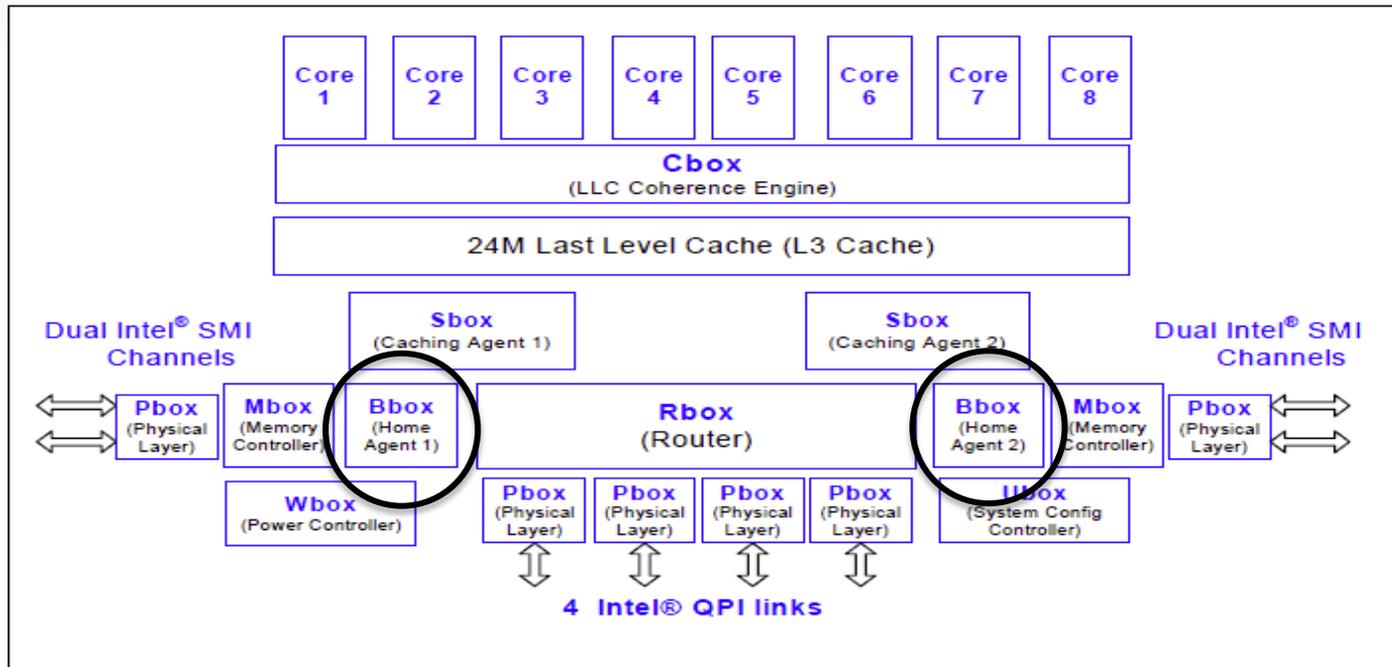
- RCR augmented HPCToolkit metrics
  - Extended command line argument to “hpcrun” command
  - “hpcrun” event specification string was extended by adding a suffix string, delimited by # and ^, which corresponds to a RCRToolkit derived system-wide metric with a boolean value
  - Example: If PAPI\_L2\_TCM is the base event and is passed as an extended specification with hpcrun, and if the RCRToolkit derived metric is a threshold variable corresponding to full utilization of a memory channel, we will have two HPCToolkit metrics – one for normal PAPI\_L2\_TCM events and other for PAPI\_L2\_TCM events that occurred during memory contention
  - Currently, user has a dictionary of mapping between contentious events and shared memory locations. (Symbolic specification via GPB “real soon now”)
  - Augmented metric appears with a “RCR-” prefix in HPCViewer, for example RCR-PAPI\_L2\_TCM
  - Could go beyond simple predicates. Multiple sub-events? Integrate RCR metrics rather than just a simple histogram?

# Examples

- Analyzed memory performance of three applications/benchmarks
  - Memory performance: significant source of bottlenecks on multi-core systems
  - Lattice QCD “chroma”, Lattice Boltzmann Magneto HydroDynamics (LBMHD), and Barcelona OpenMP FFT
- Each run is simultaneously measured by
  - RCRToolkit
  - Modified HPCToolkit hot-wired with RCRToolkit
- Two systems
  - Dell PowerEdge M910 test system (MMQ)
  - Two socket E5-2680 system

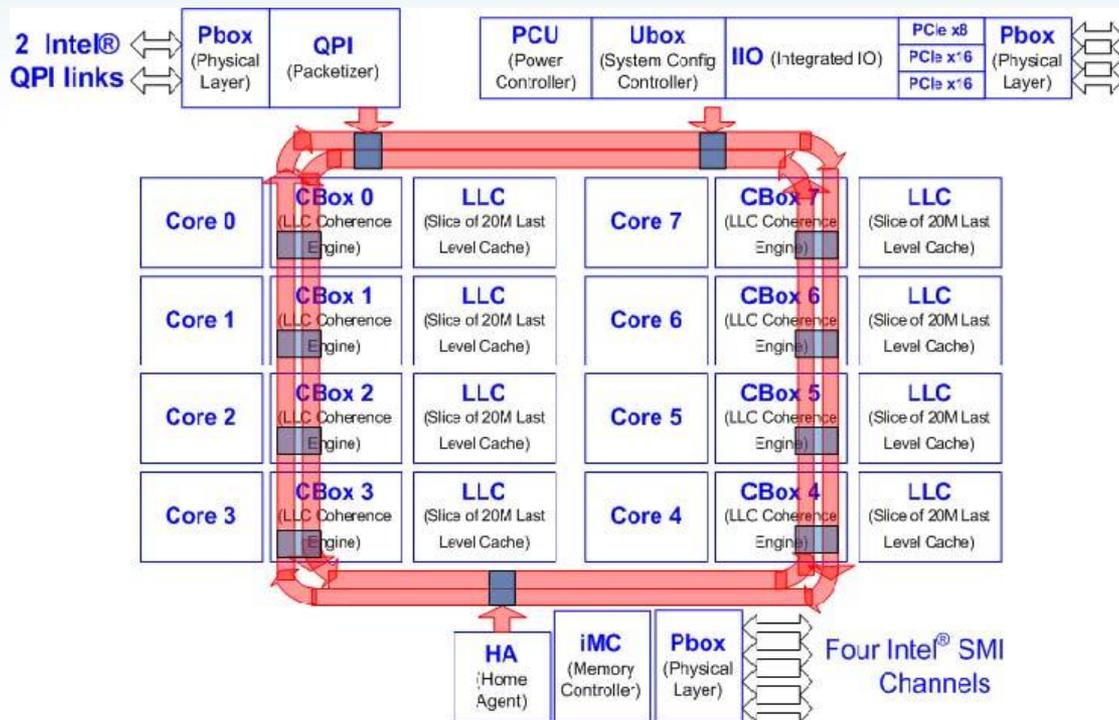
# Nehalem EX Uncore.

Figure 2-2. Intel® Xeon® Processor 7500 Series Block Diagram



- B-box has an In Memory Table (IMT) that tracks all in-flight memory block operations and ensures that they are all unique
- Uncore counter, `IMT_VALID_OCCUPANCY` tracks number of unique entries in the IMT == number of concurrent memory operations in flight
- RCRdaemon calculates IMT Average Occupancy =  $(\text{count} * 32 / \text{cycles})$

# Sandy Bridge E5-26xx

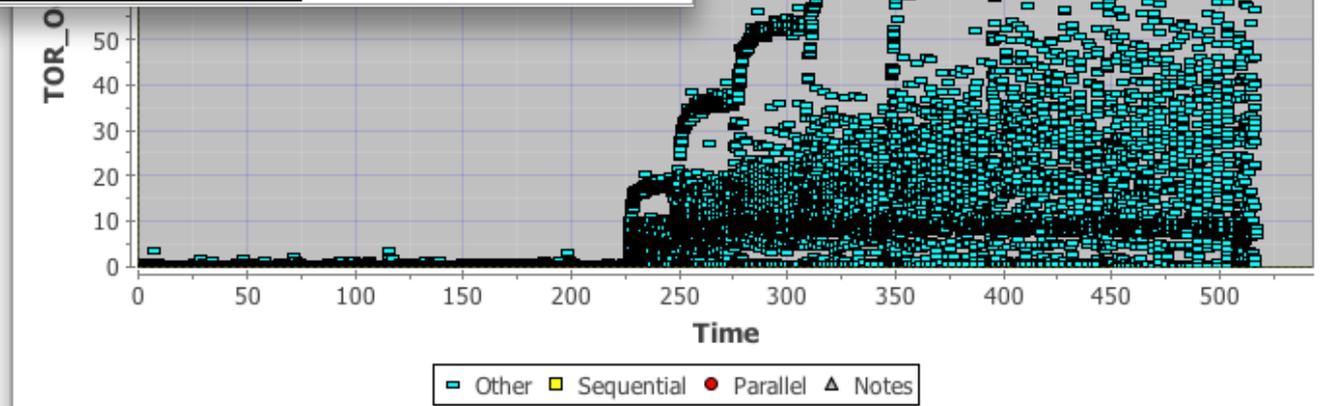
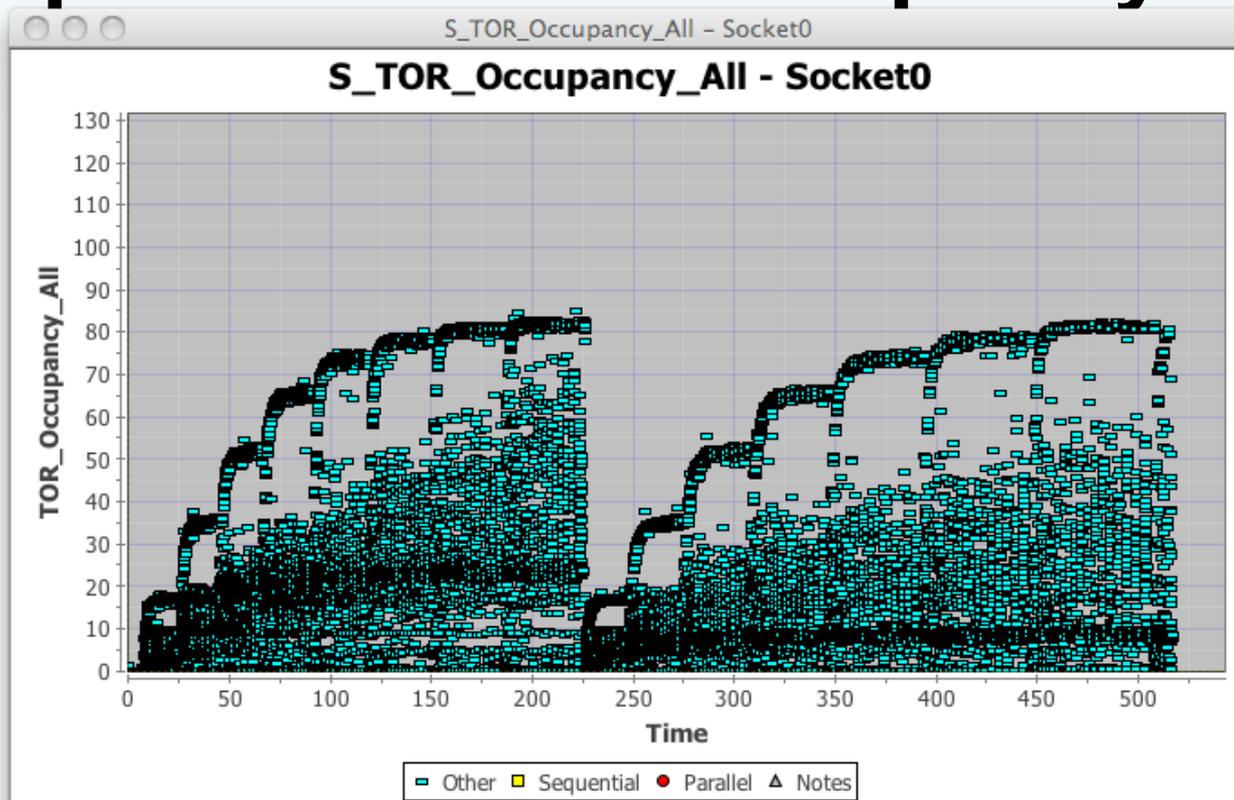


- Rich collection of uncore boxes.
- LLC is split into 8 pieces, each behind a CBox
- PMUs for CBoxes, Ubox, and PCU are accessible through MSR space.
- All other PMUs are in the PCI CONFIG space.
- LLC/Memory requests are queued in a per CBox TOR.
- We measure energy (power) consumption and TOR occupancy.

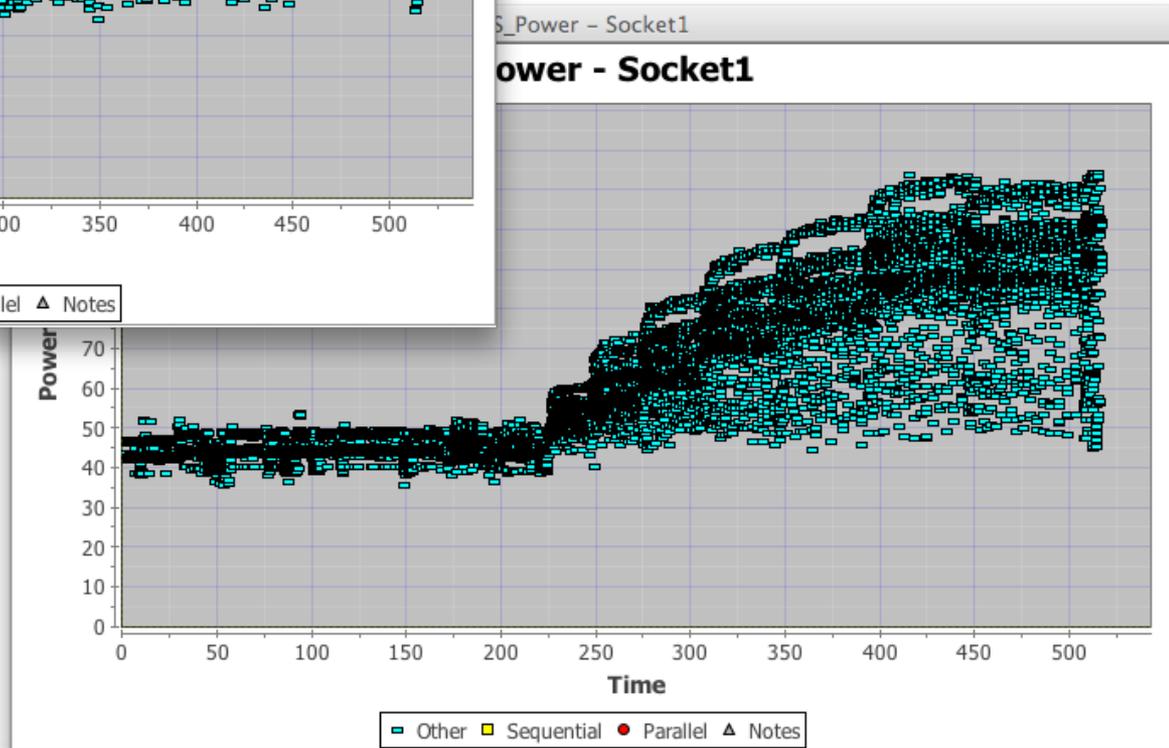
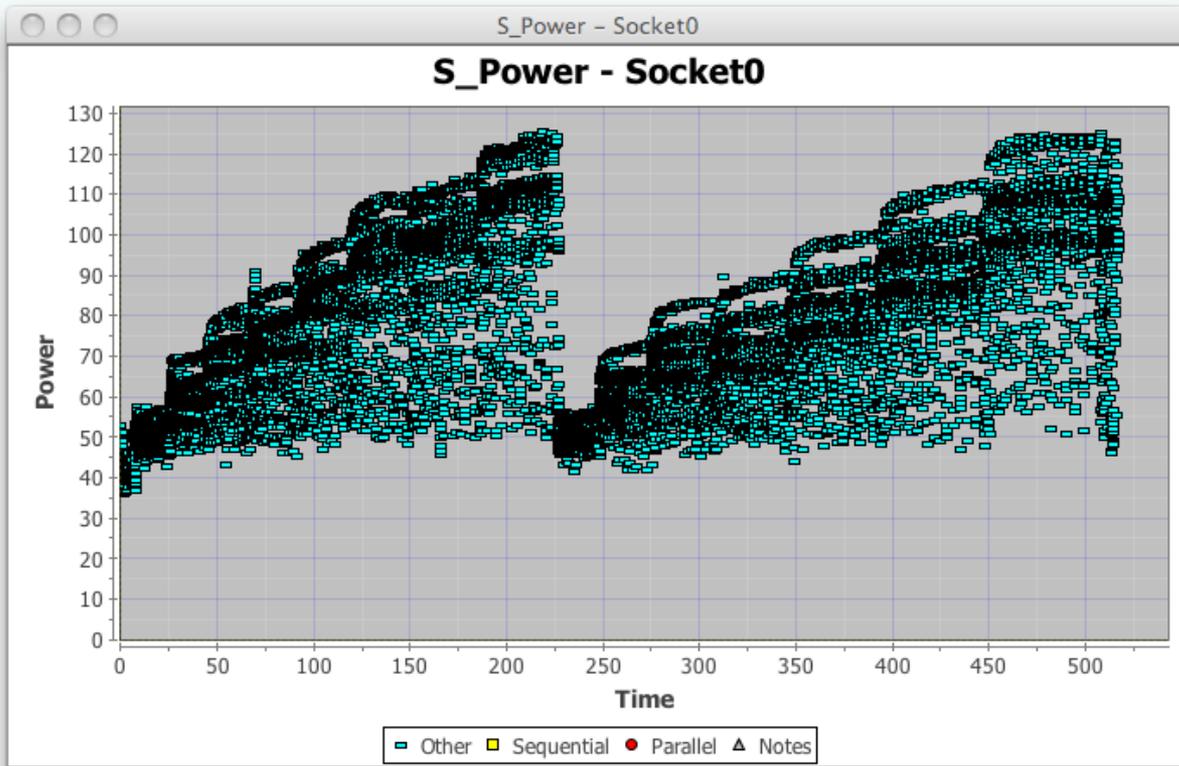
# Validation and threshold selection.

- Used the pChase benchmark code to generate controlled level of concurrent memory operations
- Bandwidth is observed to increase with load up to a practical threshold after which requests are queued
- Using pChase, determined that memory contention becomes severe when IMT Average Occupancy reaches 23, and saturates at about 27
- Chose 23 as a threshold defining a memory controller bottleneck.

# pChase TOR Occupancy on SB



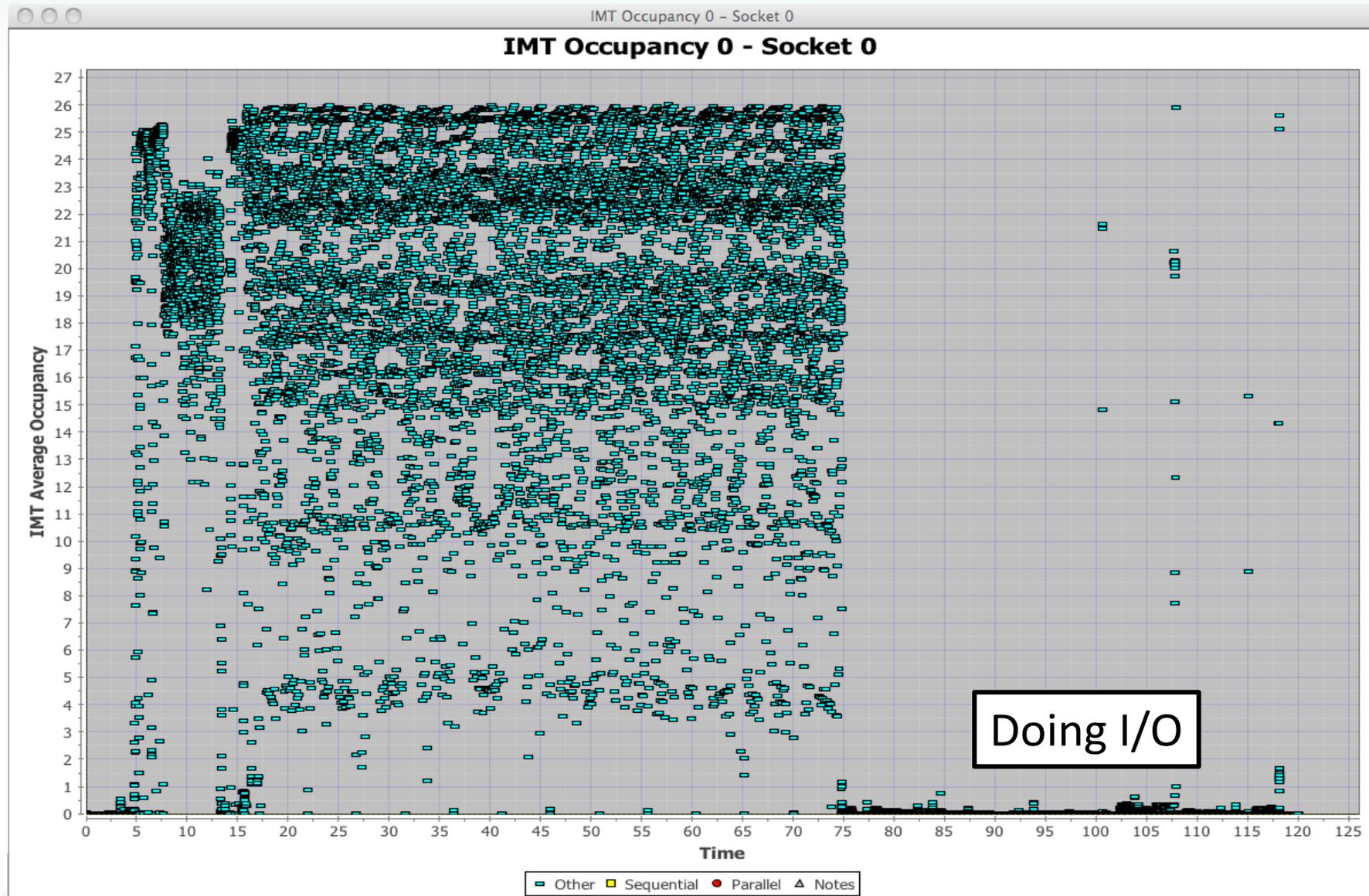
# pChase Power on Sandy Bridge 2600



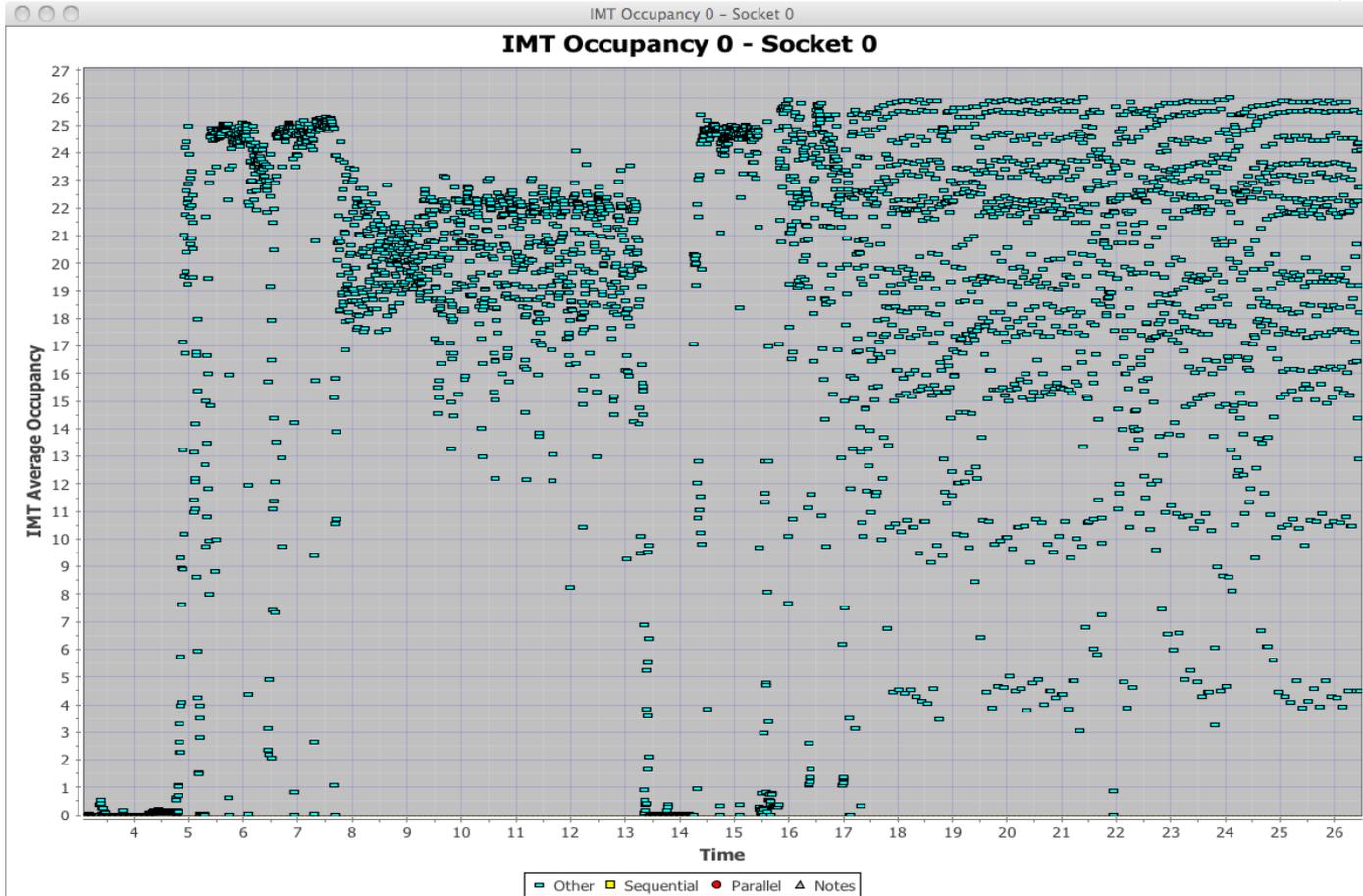
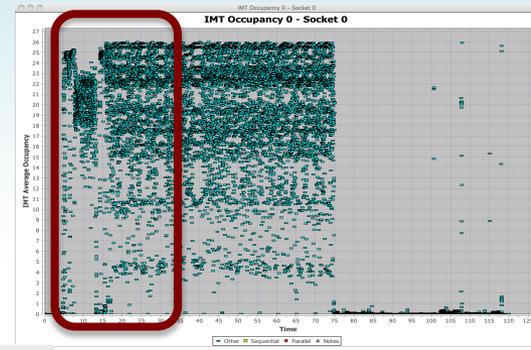
# Evaluation: Lattice QCD “chroma” Application

- Open-source C++ based software system from US SciDAC QCD initiative
- “chroma” is based on the QDP++ library that implements data-parallel programming constructs for lattice field theory
- Application runs used unmodified MPI build of Chroma version 3.7.3 and QDP++ version 1.35.1
  - Used ‘clover’ sample input from the FermiLab QCD benchmark suite
- Ran on all 32 cores of MMQ
  - For runs using HPCToolkit, RCRdaemon was running with a threshold for Average IMT Occupancy set to 23

# “chroma” Results: RCRToolkit

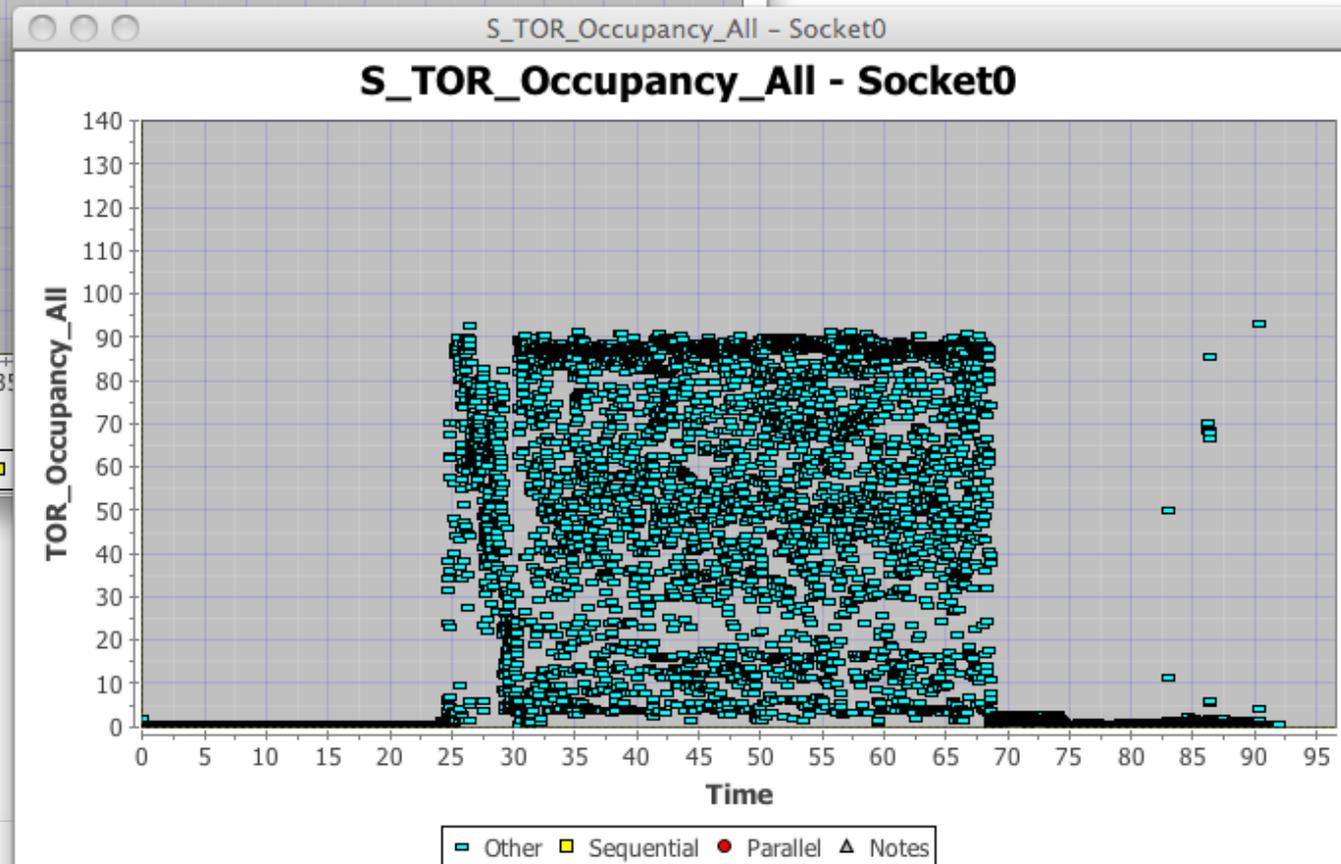
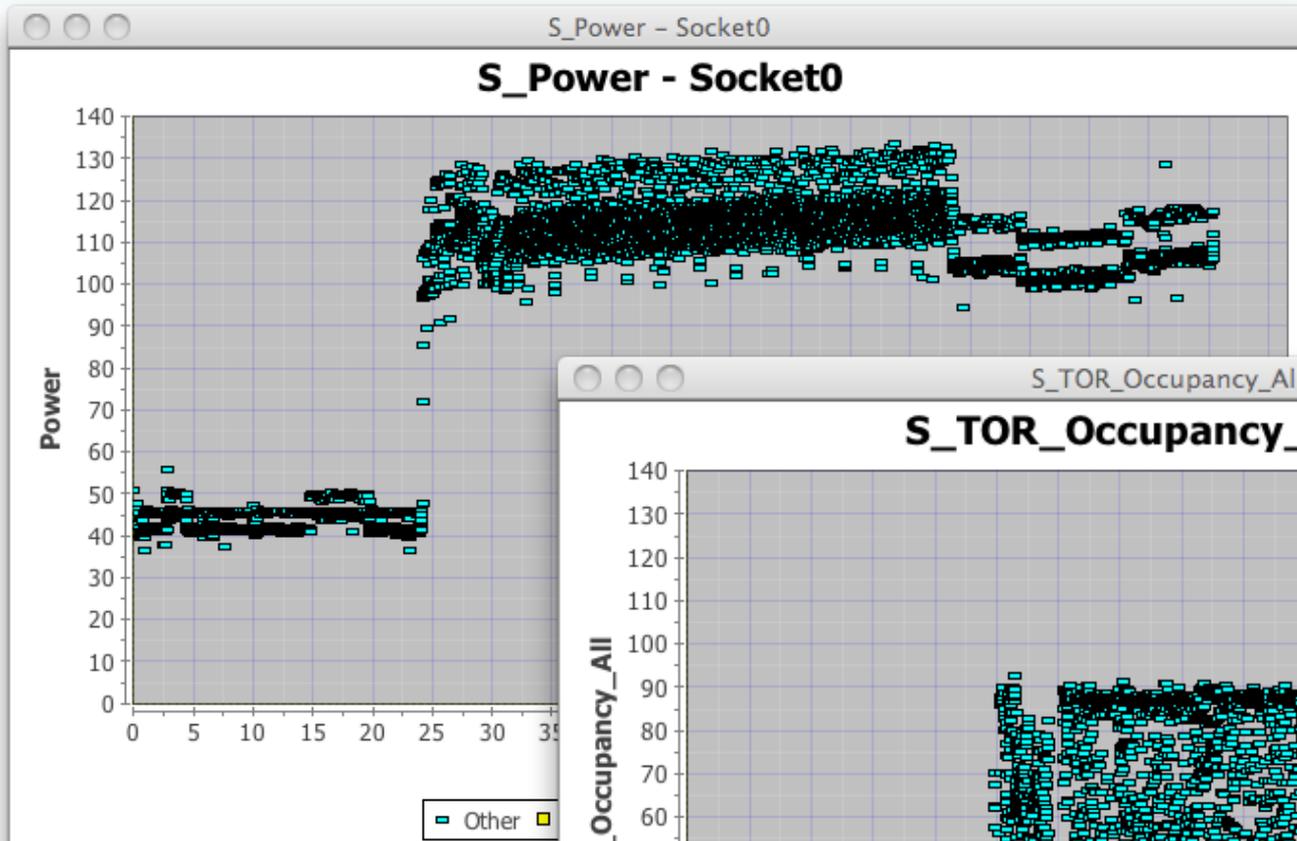


# “chroma” Results: RCRToolkit





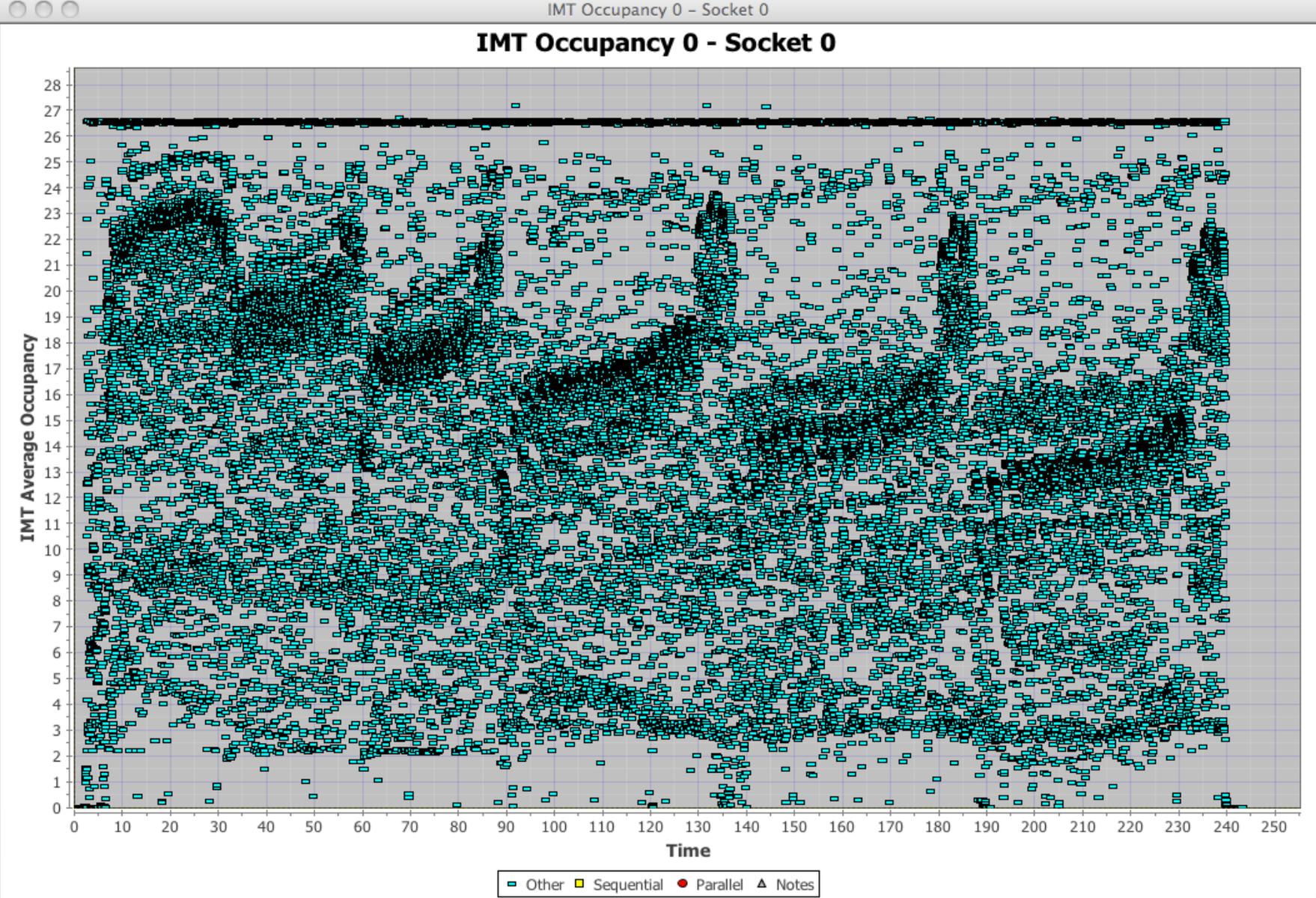
# Chroma Power and L3 queue on SB.



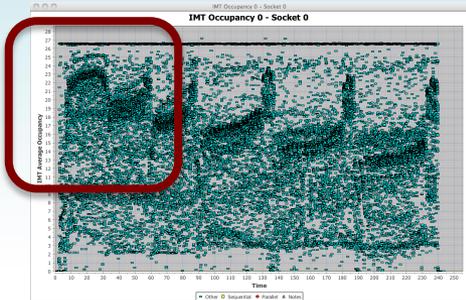
# Evaluation: LBMHD Benchmark

- The LBMHD benchmark models homogeneous isotropic turbulence in dissipative magneto-hydrodynamics
- LBMHD code was obtained from S. Williams at Lawrence Berkeley National Lab
- Ran on all 32 cores of MMQ
  - For runs using HPCToolkit, RCRdaemon was running with a threshold for Average IMT Occupancy set to 23

# LBMHD on Nehalem EX: RCRToolkit

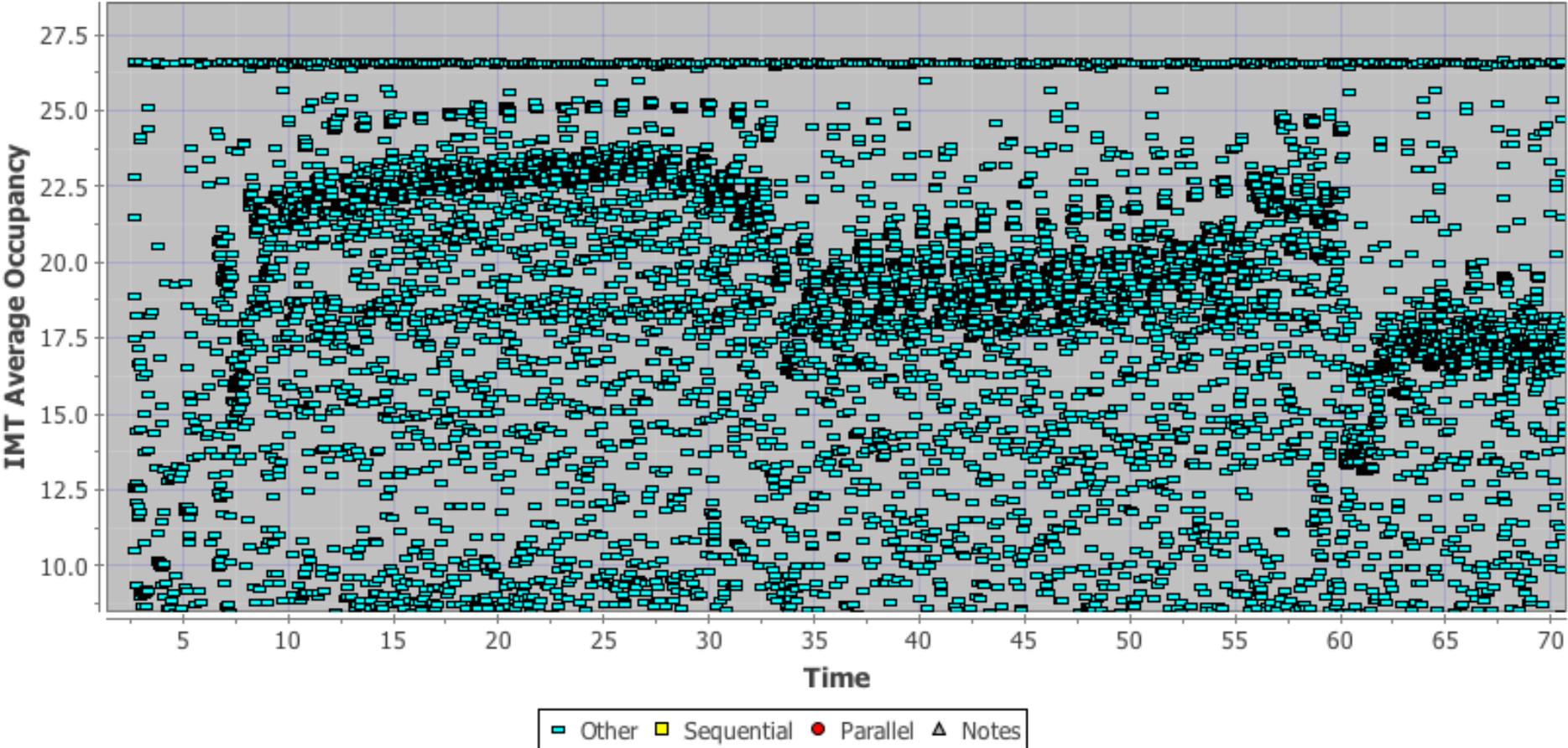


# LBMHD Results: RCRToolkit



IMT Occupancy 0 - Socket 0

## IMT Occupancy 0 - Socket 0



# LBMHD: HPCToolkit hot-wired with RCRToolkit

hpcviewer: lbmhd

```
init.c 28 bench.c
46 int64_t EndX = guide->collision.StartX+Parameters.XDimPerThread;
47 int64_t EndY = guide->collision.StartY+Parameters.YDimPerThread;
48 int64_t EndZ = guide->collision.StartZ+Parameters.ZDimPerThread;
49 for(i=0;i<=11;i++){
50   Init_MyGrid(L->F[i] , StartX, EndX, StartY, EndY, StartZ, EndZ, Up);
51 }
52 for(i=12;i<=26;i++){
53   Init_MyGrid(L->F[i] , StartX, EndX, StartY, EndY, StartZ, EndZ, Up);
54   Init_MyGrid(L->G[0][i], StartX, EndX, StartY, EndY, StartZ, EndZ, Up);
55   Init_MyGrid(L->G[1][i], StartX, EndX, StartY, EndY, StartZ, EndZ, Up);
56   Init_MyGrid(L->G[2][i], StartX, EndX, StartY, EndY, StartZ, EndZ, Up);
57 }
58   Init_MyGrid(L->V[0] , StartX, EndX, StartY, EndY, StartZ, EndZ, Up);
59   Init_MyGrid(L->V[1] , StartX, EndX, StartY, EndY, StartZ, EndZ, Up);
60   Init_MyGrid(L->V[2] , StartX, EndX, StartY, EndY, StartZ, EndZ, Up);
61   Init_MyGrid(L->B[0] , StartX, EndX, StartY, EndY, StartZ, EndZ, Up);
62   Init_MyGrid(L->B[1] , StartX, EndX, StartY, EndY, StartZ, EndZ, Up);
63   Init_MyGrid(L->B[2] , StartX, EndX, StartY, EndY, StartZ, EndZ, Up);
64   Init_MyGrid(L->R , StartX, EndX, StartY, EndY, StartZ, EndZ, Up);
65 }
66
67 void
68 p
69 i
70
71 C
```

PAPI_TOT_CYC:Sum (l)	PAPI_L2_TCM:Sum (l)	RCR-PAPI_L2_TCM:Sum (l) ▼	Contention Percentage
1.64e+13 100 %	1.25e+11 100 %	2.96e+10 100 %	2.36e+01

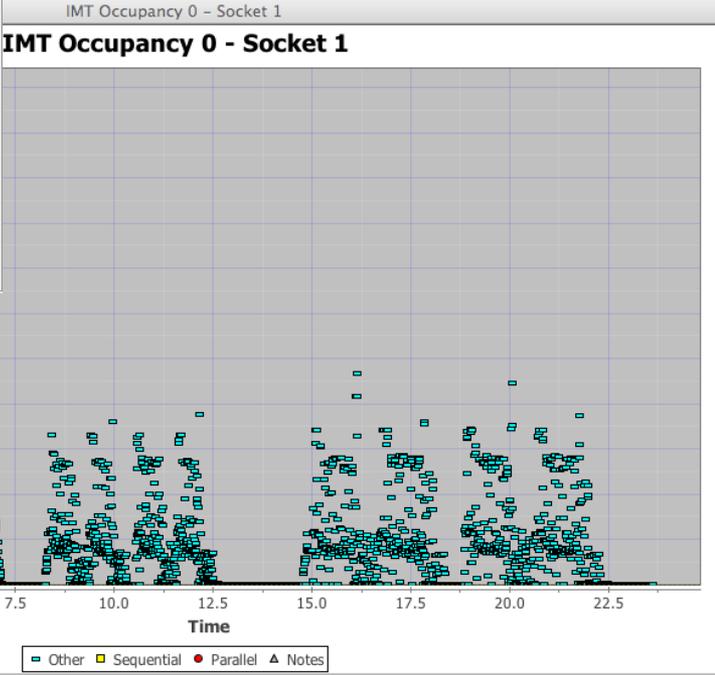
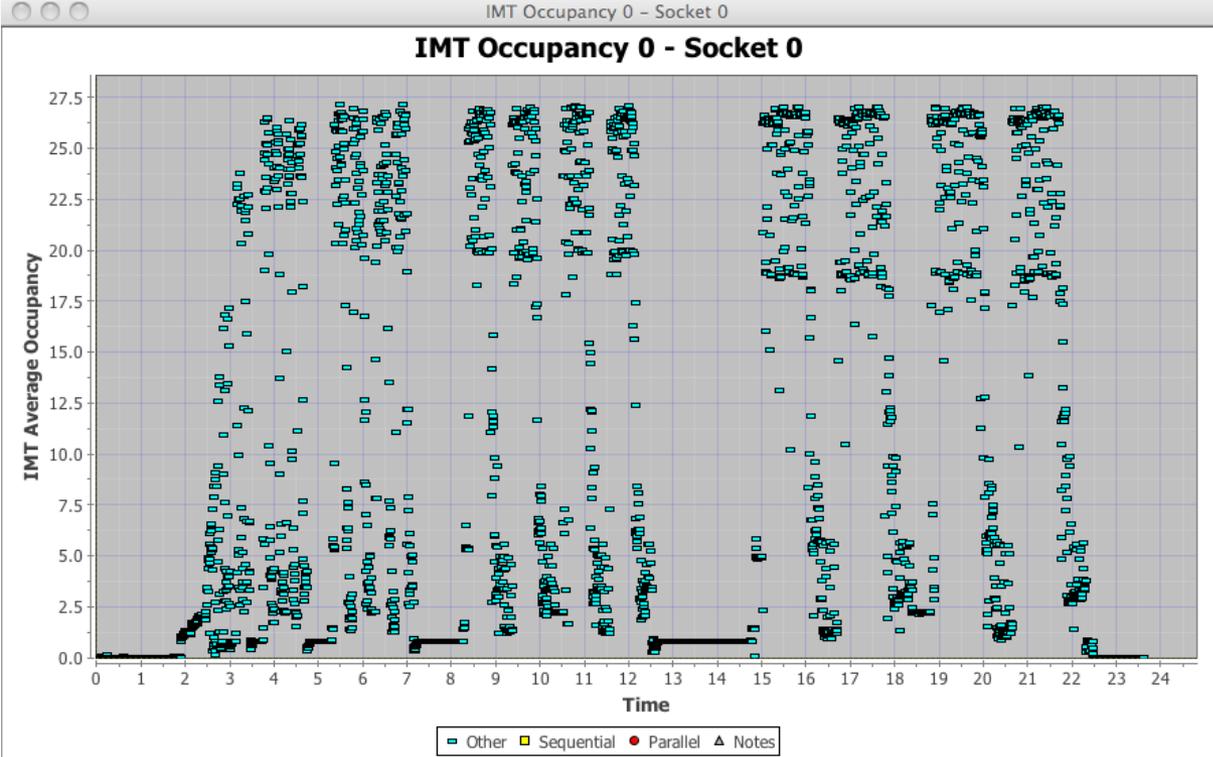
Calling Context View Callers View Flat View

Scope	PAPI_TOT_CYC:Sum (l)	PAPI_L2_TCM:Sum (l)	RCR-PAPI_L2_TCM:Sum (l) ▼	Contention Percentage
Experiment Aggregate Metrics	1.64e+13 100 %	1.25e+11 100 %	2.96e+10 100 %	2.36e+01
main	1.64e+13 100 %	1.25e+11 100 %	2.96e+10 100 %	2.36e+01
bench pthreads	1.64e+13 100.0	1.25e+11 100 %	2.96e+10 100 %	2.36e+01
loop at bench.c: 363	1.64e+13 100.0	1.25e+11 100 %	2.96e+10 100 %	2.36e+01
pthread create	1.60e+13 97.4%	1.21e+11 96.8%	2.86e+10 96.6%	2.36e+01
bench pthreads each	1.60e+13 97.4%	1.21e+11 96.8%	2.86e+10 96.6%	2.36e+01
loop at bench.c: 273	1.59e+13 97.2%	1.21e+11 96.8%	2.86e+10 96.6%	2.36e+01
loop at bench.c: 274	1.59e+13 97.2%	1.21e+11 96.8%	2.86e+10 96.6%	2.36e+01
loop at bench.c: 280	1.59e+13 97.2%	1.21e+11 96.8%	2.86e+10 96.6%	2.36e+01
loop at bench.c: 284	1.59e+13 97.2%	1.21e+11 96.8%	2.86e+10 96.6%	2.36e+01
loop at bench.c: 285	1.59e+13 97.2%	1.21e+11 96.8%	2.86e+10 96.6%	2.36e+01
loop at bench.c: 295	1.35e+13 82.1%	1.03e+11 82.3%	1.57e+10 53.0%	1.52e+01
Init Lattice pthreads each	9.95e+11 6.1%	9.20e+09 7.3%	9.10e+09 30.7%	9.89e+01
loop at init.c: 52	7.58e+11 4.6%	7.03e+09 5.6%	6.94e+09 23.4%	9.87e+01
loop at init.c: 49	1.51e+11 0.9%	1.46e+09 1.2%	1.45e+09 4.9%	9.93e+01
inlined from init.c: 3	8.62e+10 0.5%	7.10e+08 0.6%	7.10e+08 2.4%	1.00e+02
Init Lattice pthreads each	9.80e+11 6.0%	8.90e+09 7.1%	3.78e+09 12.8%	4.25e+01
inlined from barrier.c: 28	4.99e+11 3.0%	4.00e+07 0.0%	4.00e+07 0.1%	1.00e+02
inlined from bench.c: 267	4.37e+08 0.0%			
Init Lattice pthreads each	1.33e+10 0.1%			
Init Lattice pthreads each	1.30e+10 0.1%			
loop at bench.c: 295	1.35e+13 82.1%	1.03e+11 82.3%	1.57e+10 53.0%	1.52e+01
Init Lattice pthreads each	9.95e+11 6.1%	9.20e+09 7.3%	9.10e+09 30.7%	9.89e+01
loop at init.c: 52	7.58e+11 4.6%	7.03e+09 5.6%	6.94e+09 23.4%	9.87e+01

# Evaluation: FFT Benchmark

- Benchmark from the Barcelona OpenMP Tasks Suite (BOTS)
- FFT computes 1-dimensional FFT of a vector of  $n$  complex values using the Cooley-Turkey algorithm
  - Recursive algorithm that divides a DFT into smaller DFTs
  - Each division generates multiple OpenMP tasks
  - Version ran was compiled with ROSE source-to-source compiler and executed with the Qthreads runtime
- Ran on all 32 cores of MMQ
  - For runs using RCRToolkit, RCRdaemon was running with a threshold for Average IMT Occupancy set to 23

# FFT on Nehalem EX: RCRToolkit



# FFT: HPCToolkit hot-wired with RCRToolkit

hpcviewer: fft\_open\_mp

```

483  _p_jc = (p_index_ * mj2);
484  _p_jd = _p_jc;
485  _p_wjw[0] = ( *w)[[_p_jw * 2] + 0];
486  _p_wjw[1] = ( *w)[[_p_jw * 2] + 1];
487  if ( *sgn < 0.0) {
488  _p_wjw[1] = -_p_wjw[1];
489  }
490  for (_p_k = 0; _p_k < mj; _p_k++) {
491  ( *c)[[_p_jc + _p_k] * 2] + 0] = (( *a)[[_p_ja + _p_k] * 2] + 0] + ( *b)[[_p_jb + _p_k] * 2] + 0]);
492  ( *c)[[_p_jc + _p_k] * 2] + 1] = (( *a)[[_p_ja + _p_k] * 2] + 1] + ( *b)[[_p_jb + _p_k] * 2] + 1]);
493  _p_ambr = (( *a)[[_p_ja + _p_k] * 2] + 0] - ( *b)[[_p_jb + _p_k] * 2] + 0]);
494  _p_ambu = (( *a)[[_p_ja + _p_k] * 2] + 1] - ( *b)[[_p_jb + _p_k] * 2] + 1]);
495  ( *d)[[_p_jd + _p_k] * 2] + 0] = (_p_wjw[0] * _p_ambr) - (_p_wjw[1] * _p_ambu);
496  ( *d)[[_p_jd + _p_k] * 2] + 1] = (_p_wjw[1] * _p_ambr) + (_p_wjw[0] * _p_ambu);
497  }
498  }
499  }
500  }
501  stat
502  {
503  do
504  do
505  int *n2 = (int *)(((struct OUT__2_1527__data *)__out_argv) -> OUT__2_1527__data::n2_p);

```

PAPI_TOT_CYC:Sum (I)	PAPI_L2_TCM:Sum (I)	RCR-PAPI_L2_TCM:Sum (I)	Contention Percentage
1.37e+12 100 %	2.69e+09 100 %	1.46e+09 100 %	5.43e+01

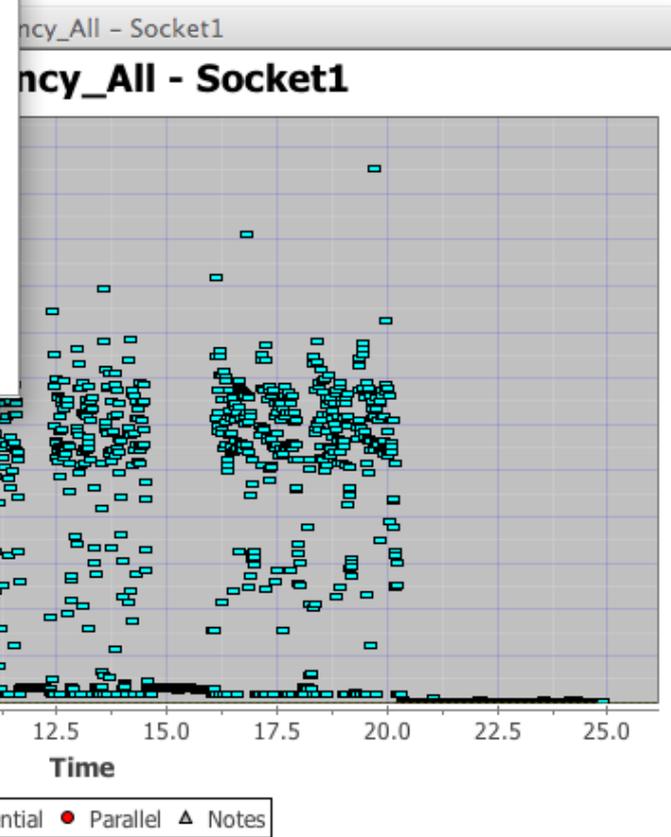
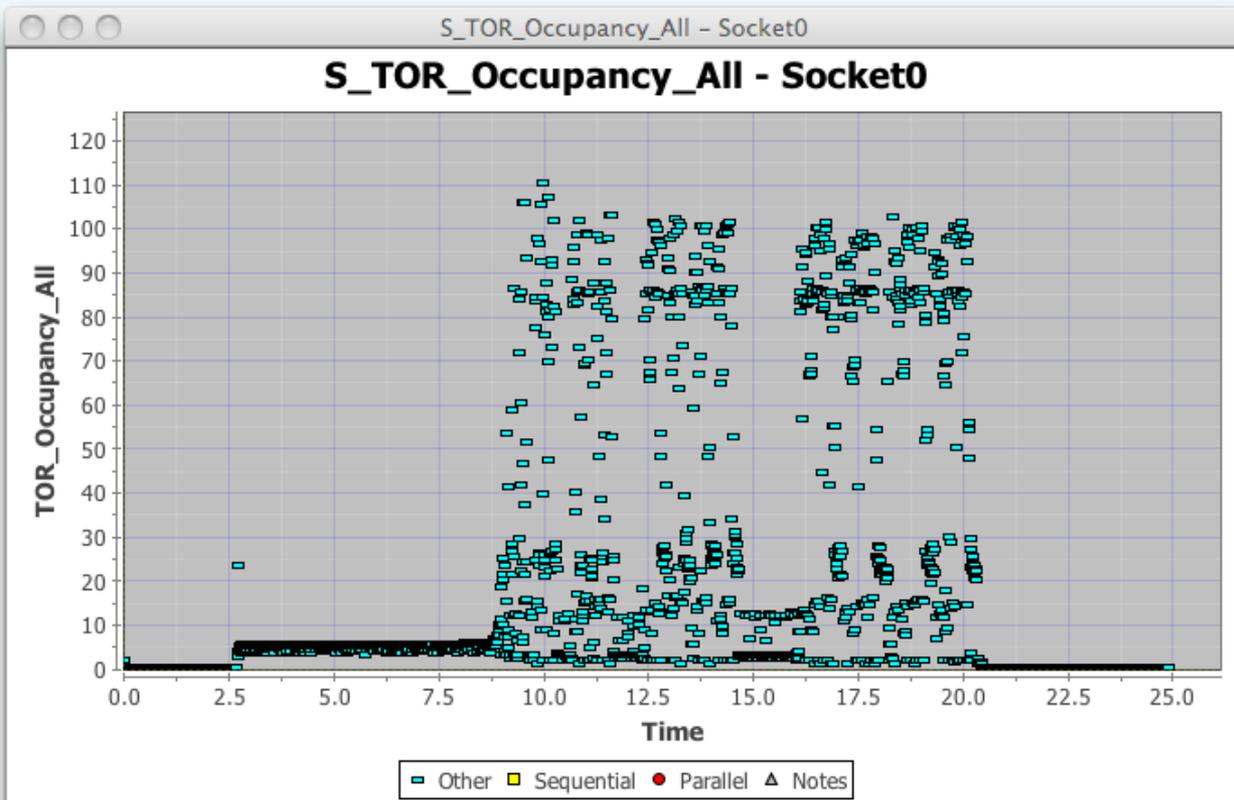
Calling Context View | Callers View | Flat View

Scope	PAPI_TOT_CYC:Sum (I)	PAPI_L2_TCM:Sum (I)	RCR-PAPI_L2_TCM:Sum (I)	Contention Percentage
Experiment Aggregate Metrics	1.37e+12 100 %	2.69e+09 100 %	1.46e+09 100 %	5.43e+01
monitor main	8.40e+11 61.2%	1.26e+09 46.8%	3.50e+08 24.0%	2.78e+01
atthread wrapper	4.91e+11 35.8%	1.25e+09 46.5%	1.06e+09 72.6%	8.48e+01
qloop step wrapper	4.90e+11 35.7%	1.23e+09 45.7%	1.04e+09 71.2%	8.46e+01
OUT_1_1527	4.71e+11 34.3%	1.22e+09 45.4%	1.04e+09 71.2%	8.52e+01
loop at fft_open_mp.rose: 479	4.66e+11 33.9%	1.22e+09 45.4%	1.04e+09 71.2%	8.52e+01
loop at fft open mp.rose: 490	4.57e+11 33.3%	1.20e+09 44.6%	1.02e+09 69.9%	8.50e+01
fft open mp.rose: 486	3.77e+09 0.3%			
inlined from fft open mp.rose: 454	3.15e+09 0.2%	2.00e+07 0.7%	2.00e+07 1.4%	1.00e+02
fft open mp.rose: 487	7.74e+08 0.1%			
fft open mp.rose: 479	3.58e+08 0.0%			
fft open mp.rose: 485	3.13e+08 0.0%			
fft open mp.rose: 488	3.60e+07 0.0%			
XOMP loop default	5.37e+09 0.4%			
loop at fft open mp.rose: 490	4.57e+11 33.3%	1.20e+09 44.6%	1.02e+09 69.9%	8.50e+01
OUT_3_1527	4.83e+09 0.4%			

321M of 513M

# FFT TOR Occupancy on SB

We know, from the MMQ results, that data is poorly balanced between memories. Why is there an imbalance at the TOR?



# FFT Power on SB

