

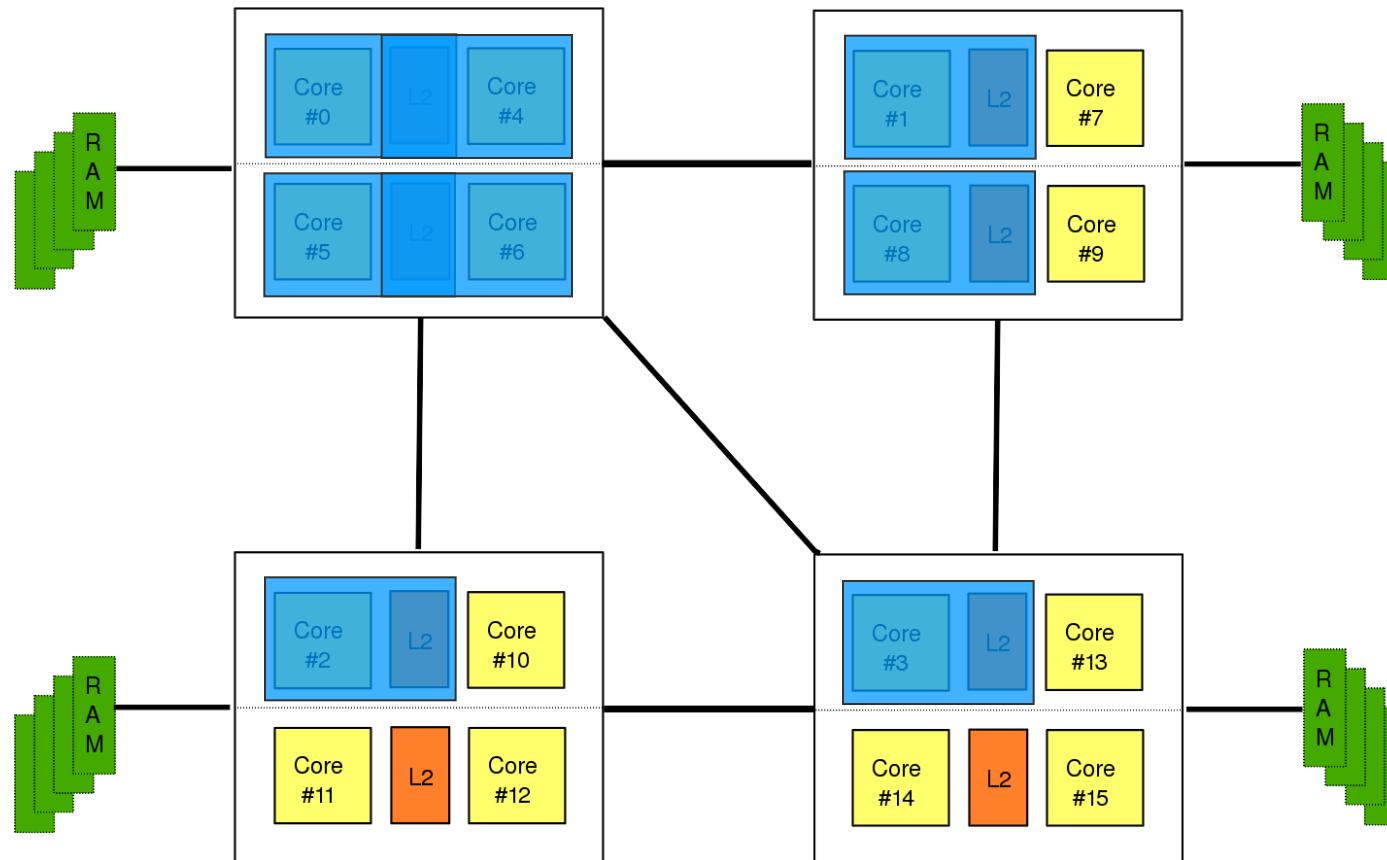
# On the Importance of Thread Placement on Multicore Architectures

**HPCLatAm 2011**  
Keynote  
Cordoba, Argentina  
August 31, 2011

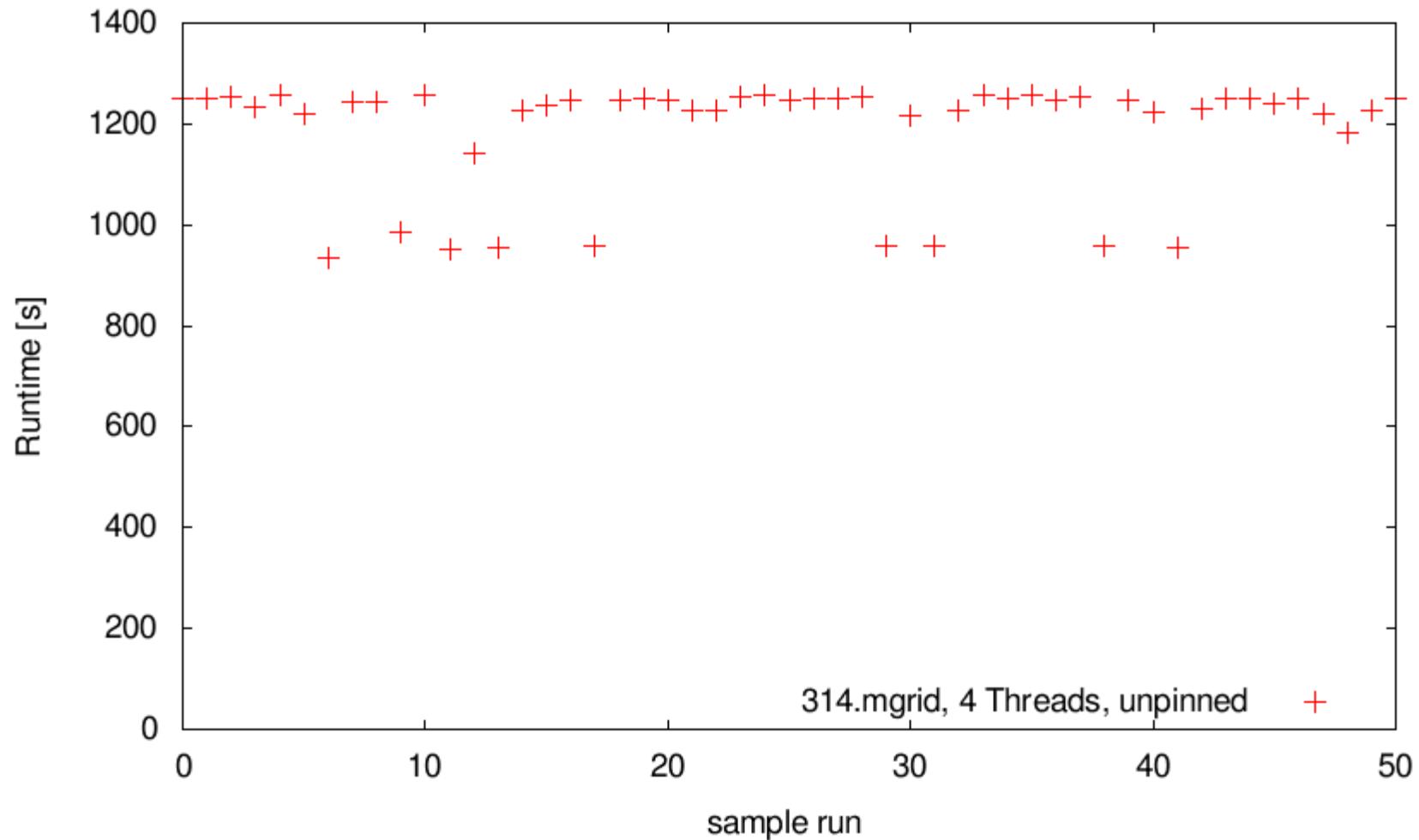
Tobias Klug



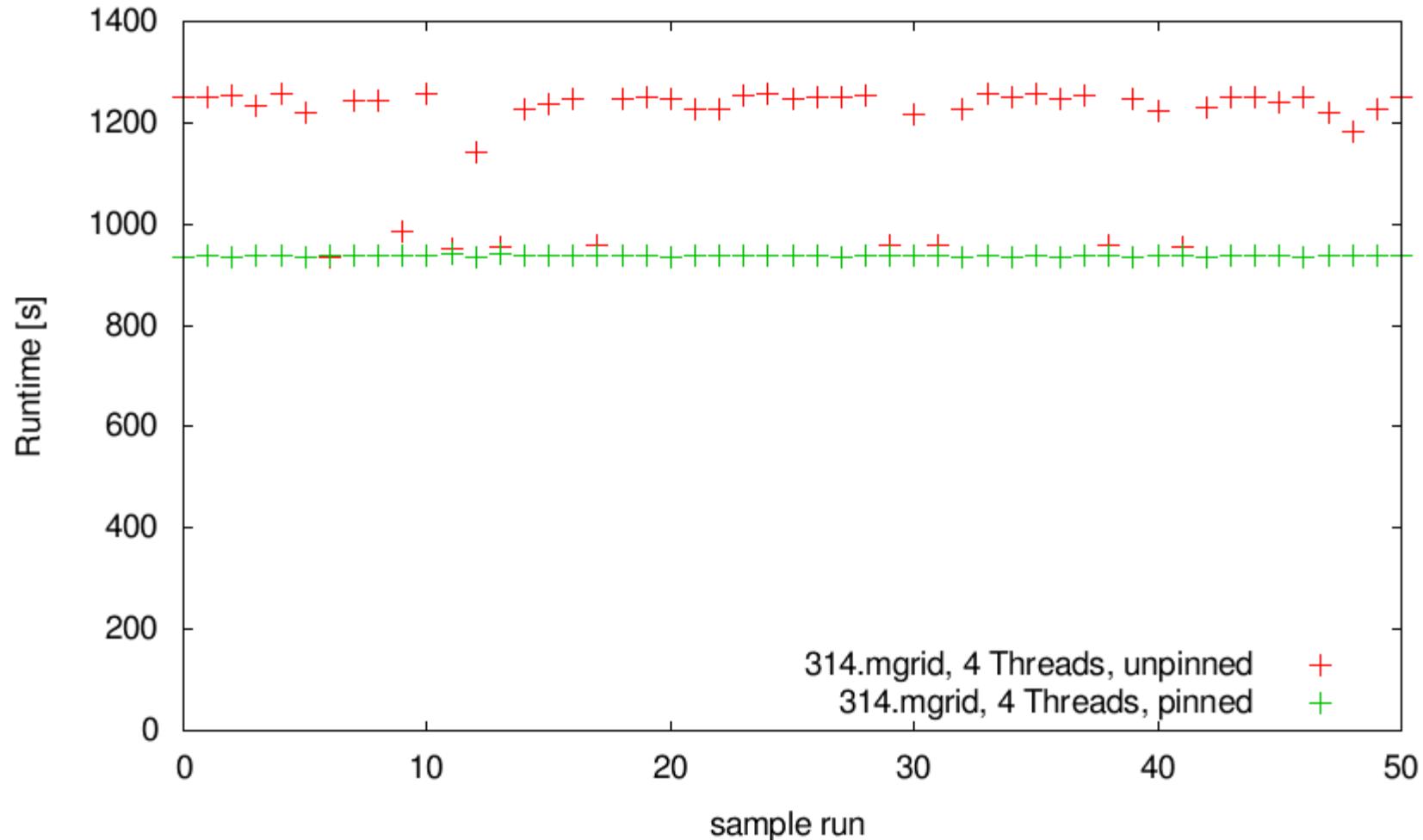
# Motivation: Many possibilities...



... can lead to non-deterministic runtimes...



... but don't have to



# The `autopin` Approach

- User-level tool
- Start multi-threaded application under `autopin` control
- User can specify pinnings of interest
- Pin threads to cores
- Assess performance of chosen pinning using performance counters
- Try alternative pinnings until optimal pinning is found

# Performance Counters

- Multiple Event Sensors
  - ALU Utilization
  - Branch Prediction
  - Cache Events (L1/L2/TLB)
  - Bus Utilization
- Two Uses:
  - Read: Get Precise Count of Events in Code Regions => Counting
  - Interrupt on Overflow => Statistical Sampling
- Well-known tools:
  - Oprofile
  - Perfctr
  - Intel Vtune
  - Perfmon2

# perfmon2

- Kernel-Patch + library (libpfm)
- Generic interface for PMU access
- Portable: implementations for IA32, x64, IA64, MIPS, Power
- Allows for per-thread and system-wide monitoring
- Support for counting and sampling
- pfmon:
  - attach to running threads
  - fork new processes and attach to them
  - fully exploit performance counters

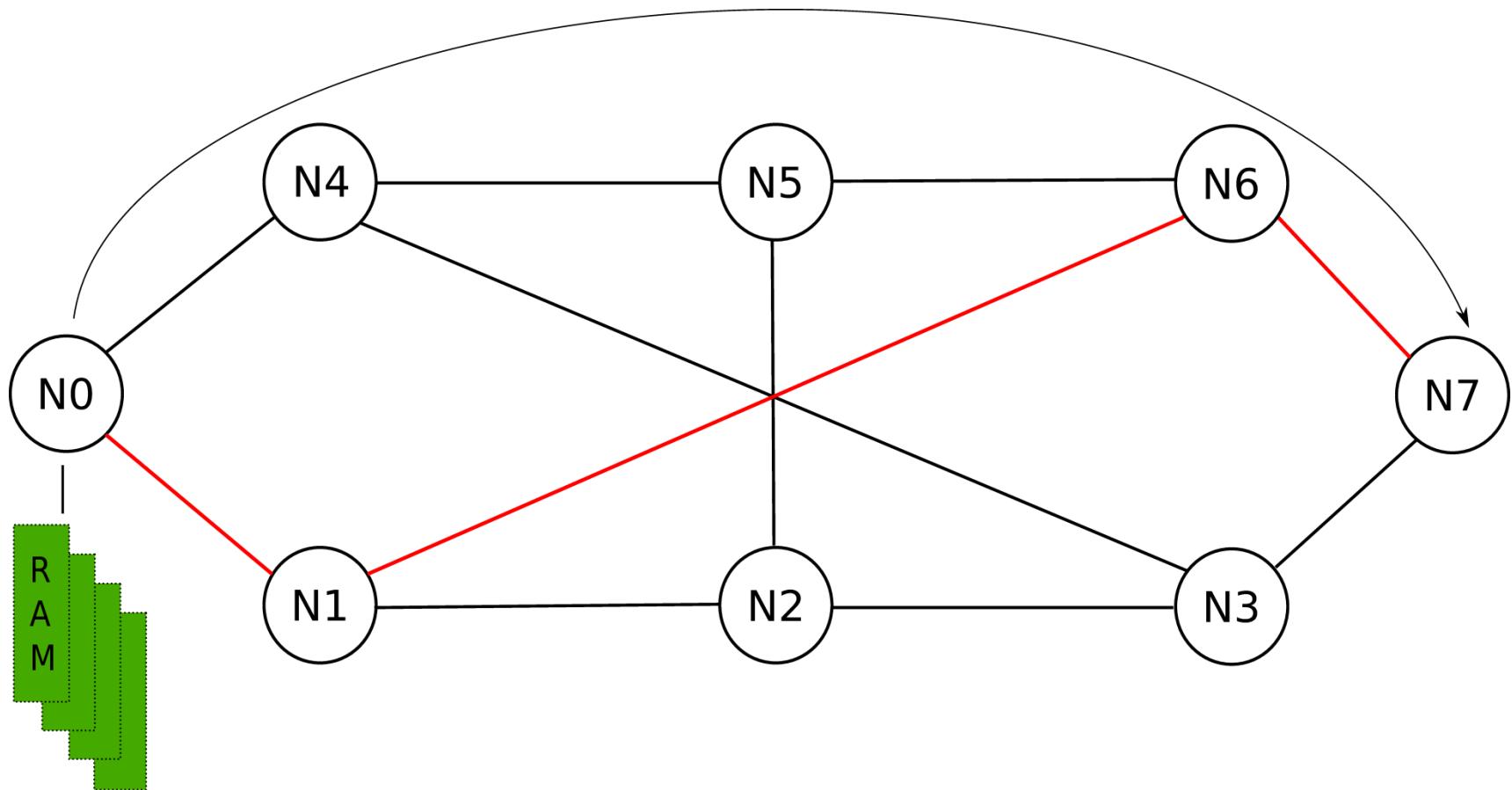
# Algorithm

```
init_autopin (pinningList, initTime, program);

for (i=0; i<numOfPinnings; i++)
{
    pinThreads(pinningList[i]);
    runThreads(warmupTime);
    p1 = readPerformanceCounters();
    runThreads(sampleTime);
    p2 = readPerformanceCounters();
    performanceRate[i] = (p2-p1)/sampleTime;
}

pinThreads(bestPinning);
```

# NUMA automatic page migration



# Automatic Page Migration

- Kernel patch from Lee Schermerhorn
- Thread moves to new NUMA node:  
remove PTE references of “old” NUMA node  
pages are now unmapped
- Next access to page causes page-fault
- Modified kernel routines pull page local  
(migrate on fault)
- Update PTE
- Controlled via cpusets

# Experimental Setup

- Caneland:
  - Intel Tigertown: Quad-Core, 2x4MB L2/socket, 2.93GHz clock rate
  - 4-way, 4x1066MHz FSB, 64MB snoop filter, UMA
- Clovertown:
  - Intel Clovertown: Quad-Core, 2x4MB L2, 2.66GHz clock rate
  - 2-way, 2x1333MHz FSB, UMA
- Barcelona:
  - AMD K10: Quad-Core, 4x512kB L2, 1x2MB L3, 1.9GHz clock rate
  - 2-way, 1000MHz Hypertransport, NUMA
- Linux Kernel 2.6.23 with perfmon2 patches
- Intel Compiler Suite

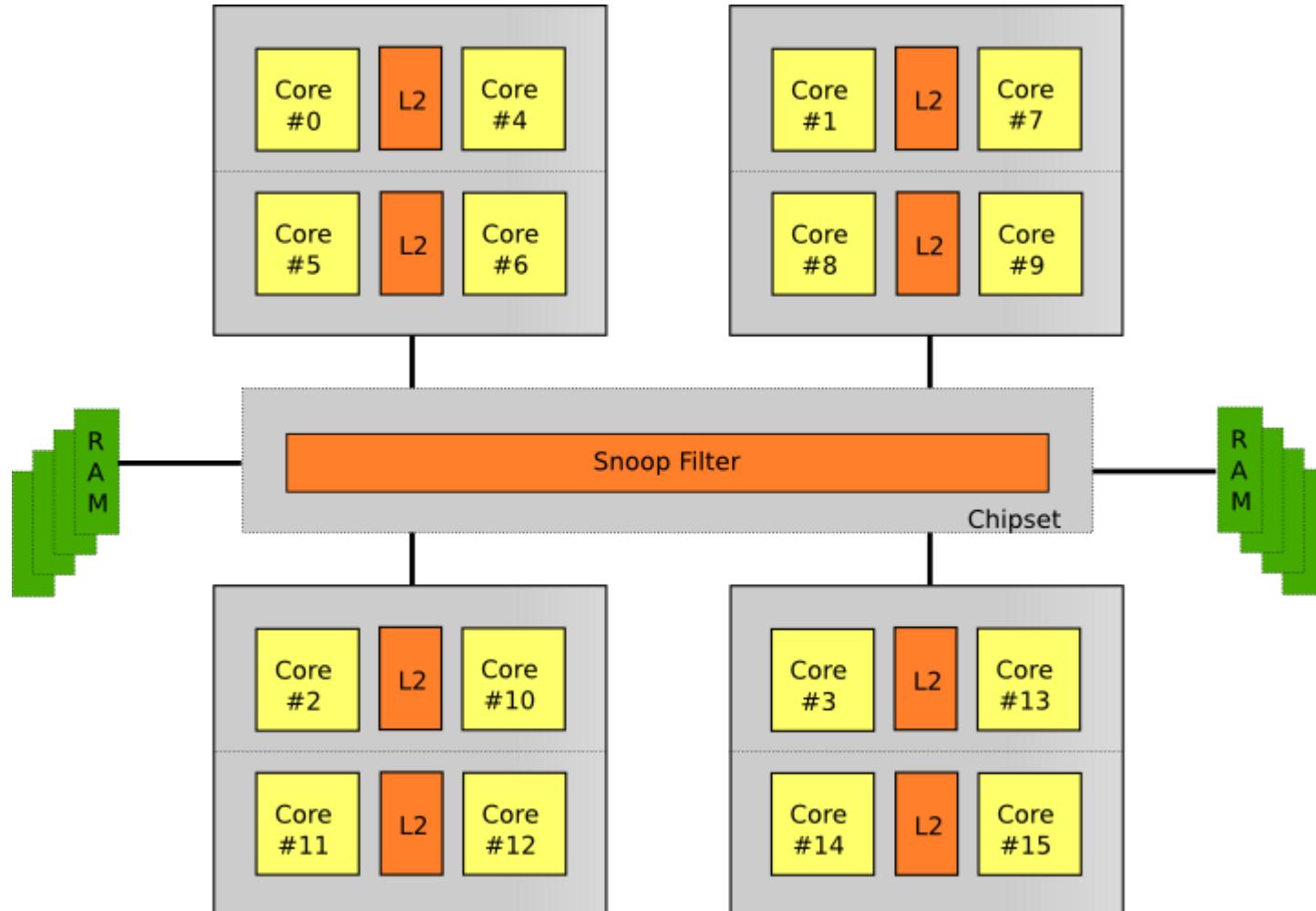
# Evaluation

- **SPEC OMP**
  - Benchmark consists of real scientific applications
  - OpenMP
  - PC: INSTRUCTIONS\_RETIRED
  - Several Multicore-Architekturen examined
- Memory Throughput
- MPI
- Electric power consumption

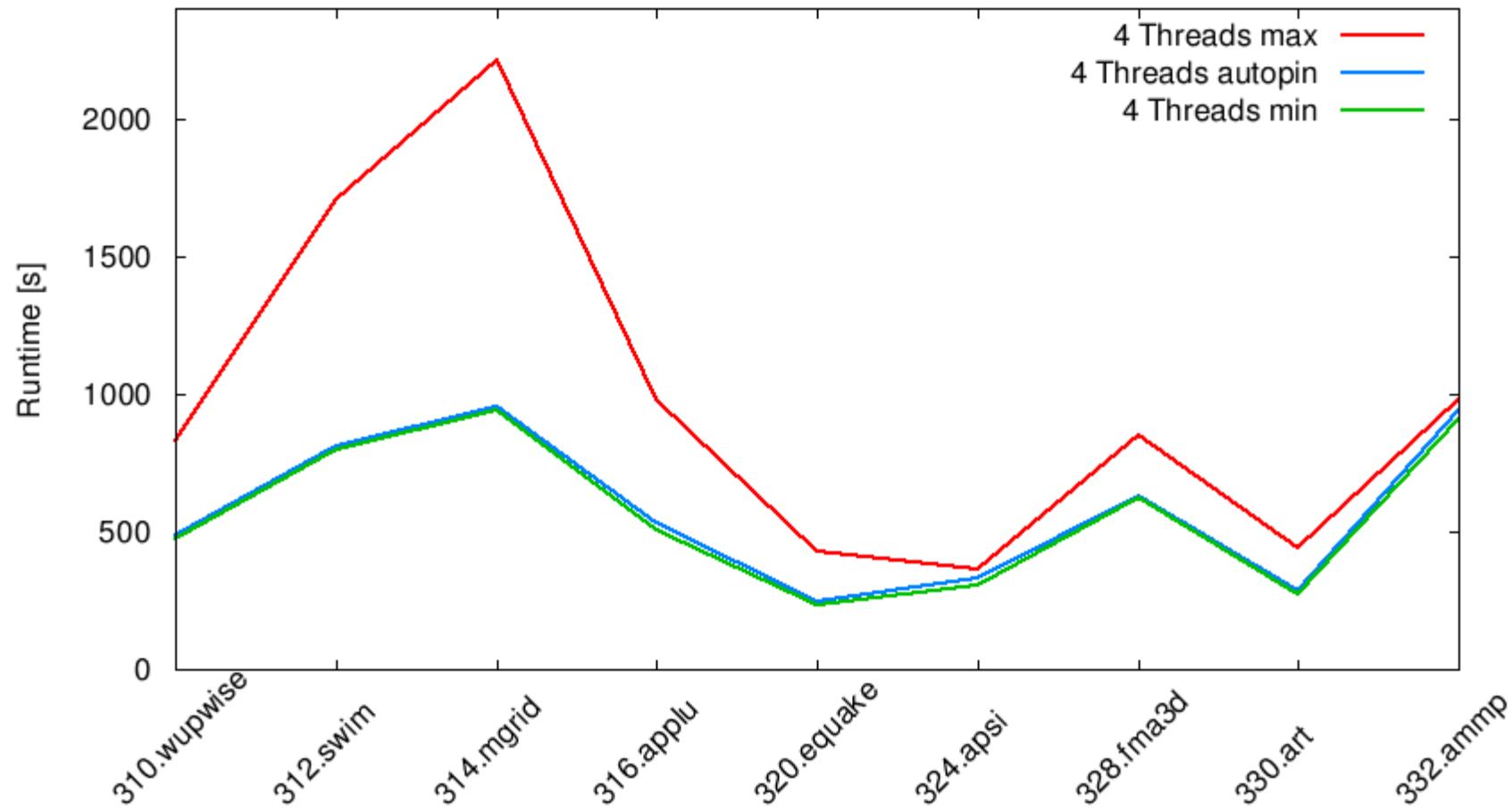
# SPEC OMP

Benchmark	Description
310.wupwise	Quantum chromodynamics
312.swim	shallow water modeling
314.mgrid	multi-grid solver in 3D potential field
316.applu	parabolic/elliptic partial differential equations
318.galgel	fluid dynamics analysis of oscillatory instability
320.equake	finite element simulation of earthquake modeling
324.apsi	weather prediction
326.gafort	genetic algorithm code
328.fma3d	finite-element crash simulation
330.art	neural network simulation of adaptive resonance theory
332.Ammp	computational chemistry

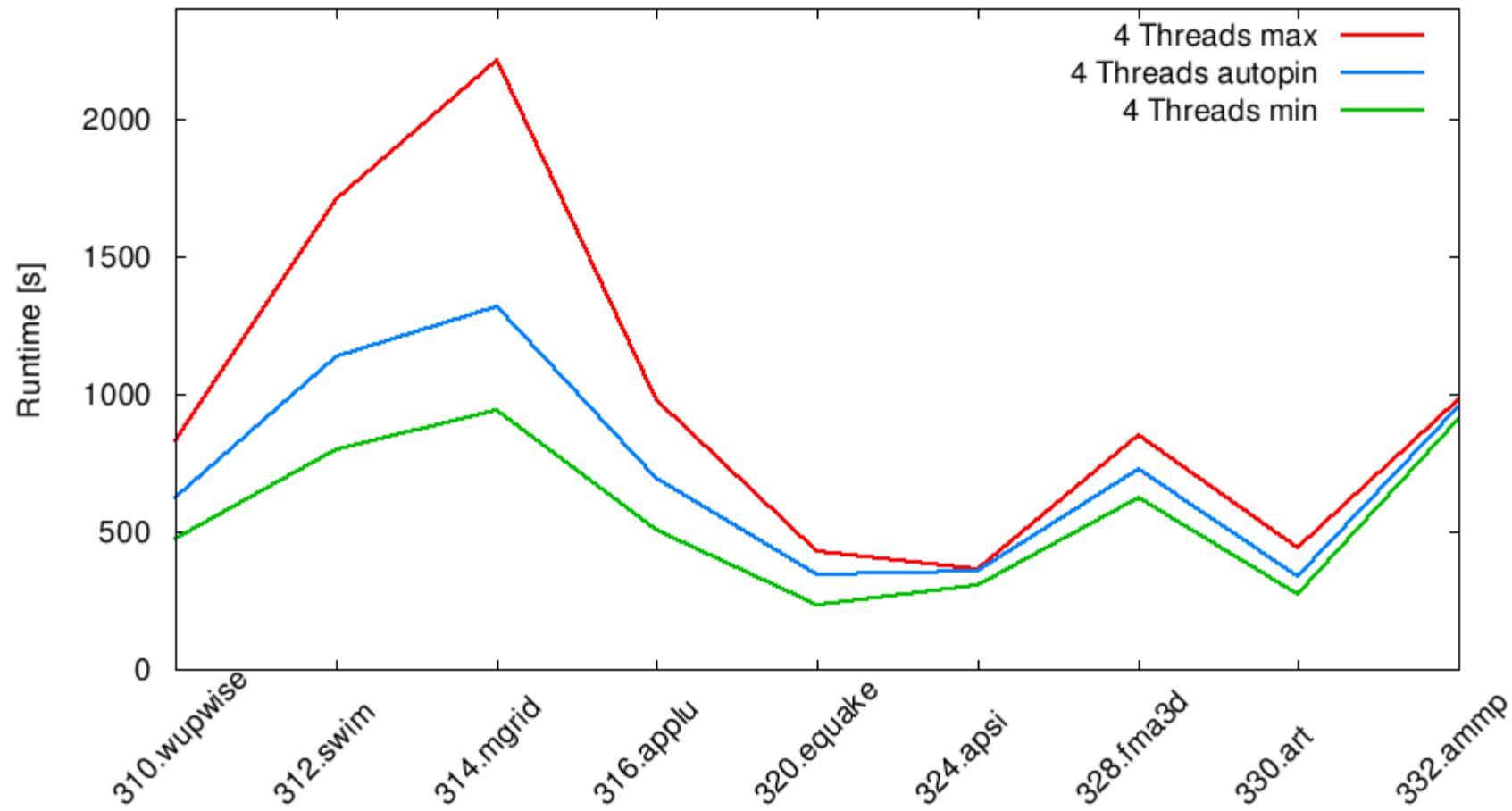
# Caneland



## Caneland: Runtimes (10s sample time)



## Caneland: Runtimes (30s sample time)



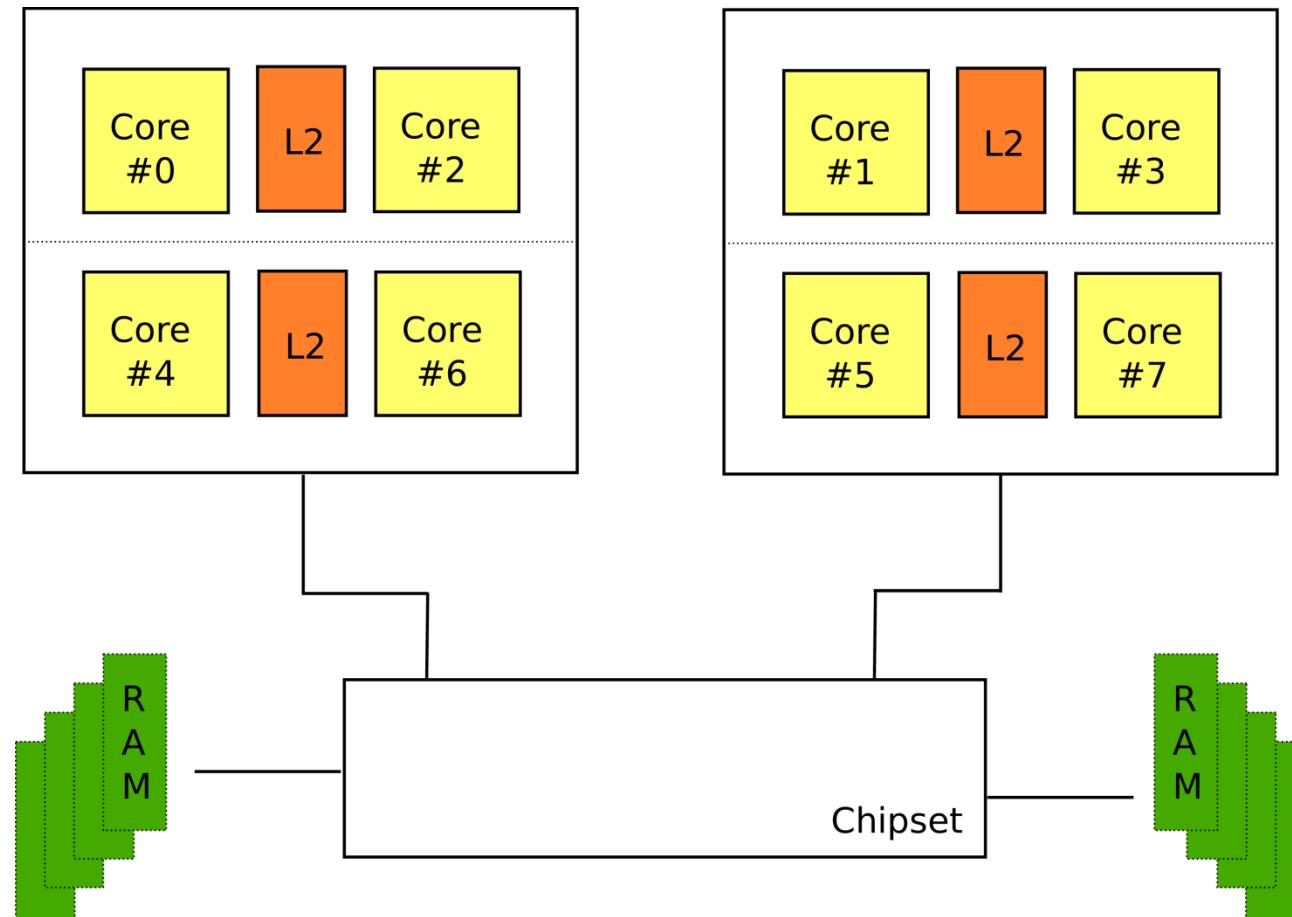
# Results

	Caneland			Clovertown		Barcelona	
	2	4	8	2	4	2	4
310.wupwise							
312.swim							
314.mgrid							
316.aplu							
320.equake							
324.apsi							
328.fma3d							
330.art							
332.ammp							

# Results

	Caneland			Clovertown		Barcelona	
	2	4	8	2	4	2	4
310.wupwise							
312.swim							
314.mgrid							
316.aplu							
320.equake							
324.apsi							
328.fma3d							
330.art							
332.ammp							

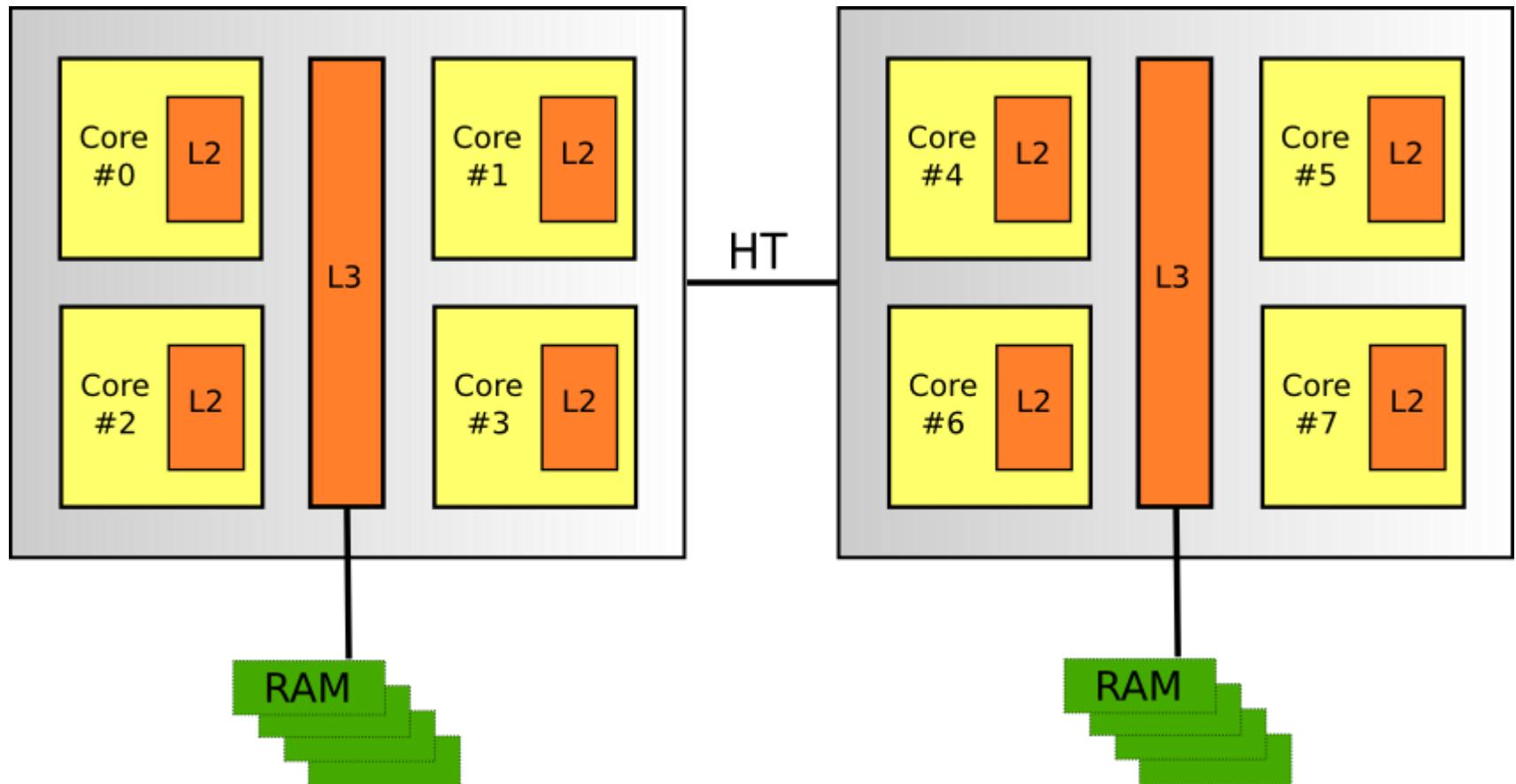
# Clovertown



# Results

	Caneland			Clovertown		Barcelona	
	2	4	8	2	4	2	4
310.wupwise							
312.swim							
314.mgrid							
316.aplu							
320.equake							
324.apsi							
328.fma3d							
330.art							
332.ammp							

# Barcelona



## Results (w/o NUMA patch)

	Caneland			Clovertown		Barcelona	
	2	4	8	2	4	2	4
310.wupwise							
312.swim							
314.mgrid							
316.aplu							
320.equake							
324.apsi							
328.fma3d							
330.art							
332.ammp							

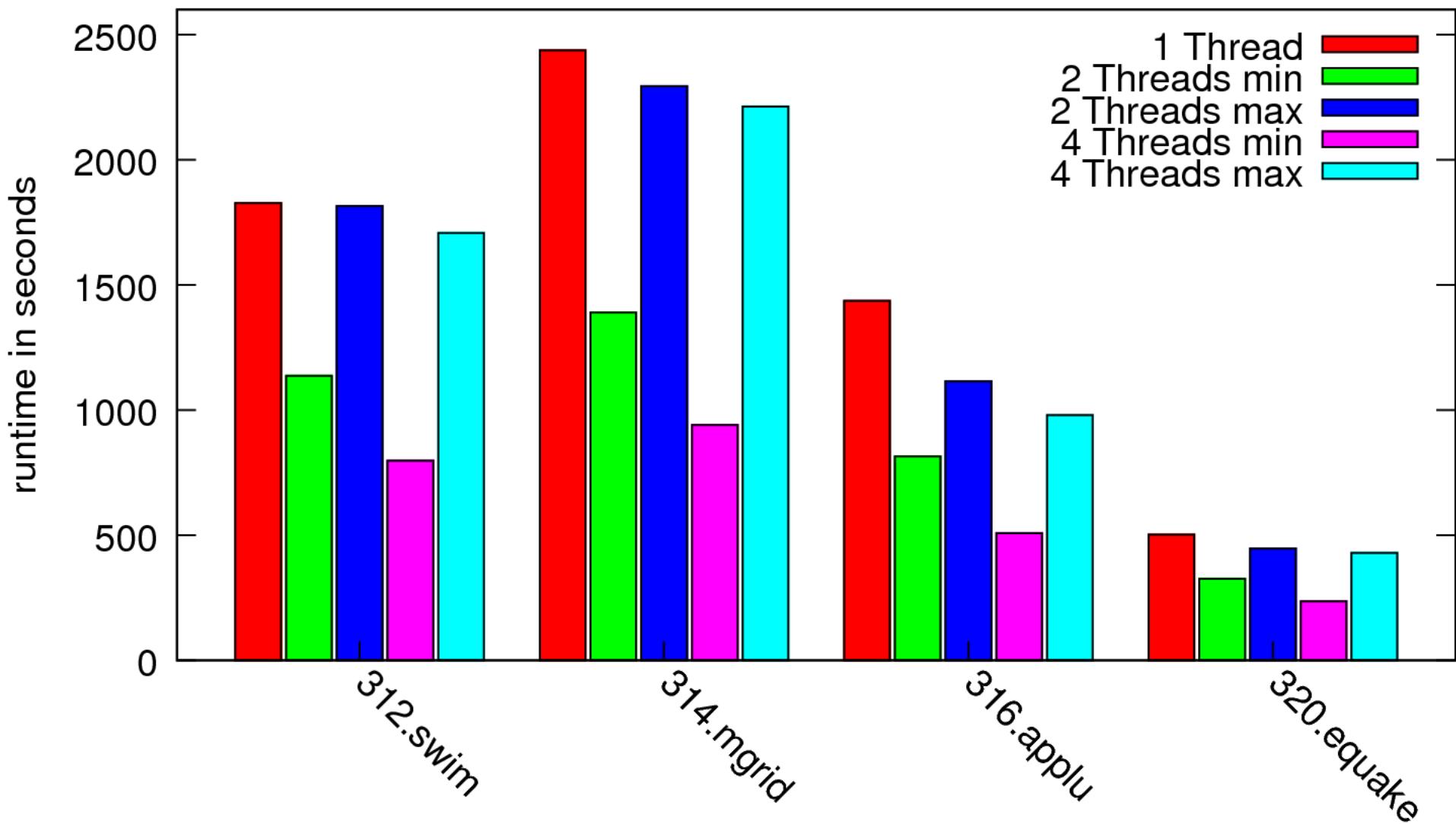
## Results (with NUMA patch)

	Caneland			Clovertown		Barcelona	
	2	4	8	2	4	2	4
310.wupwise							
312.swim							
314.mgrid							
316.aplu							
320.equake							
324.apsi							
328.fma3d							
330.art							
332.ammp							

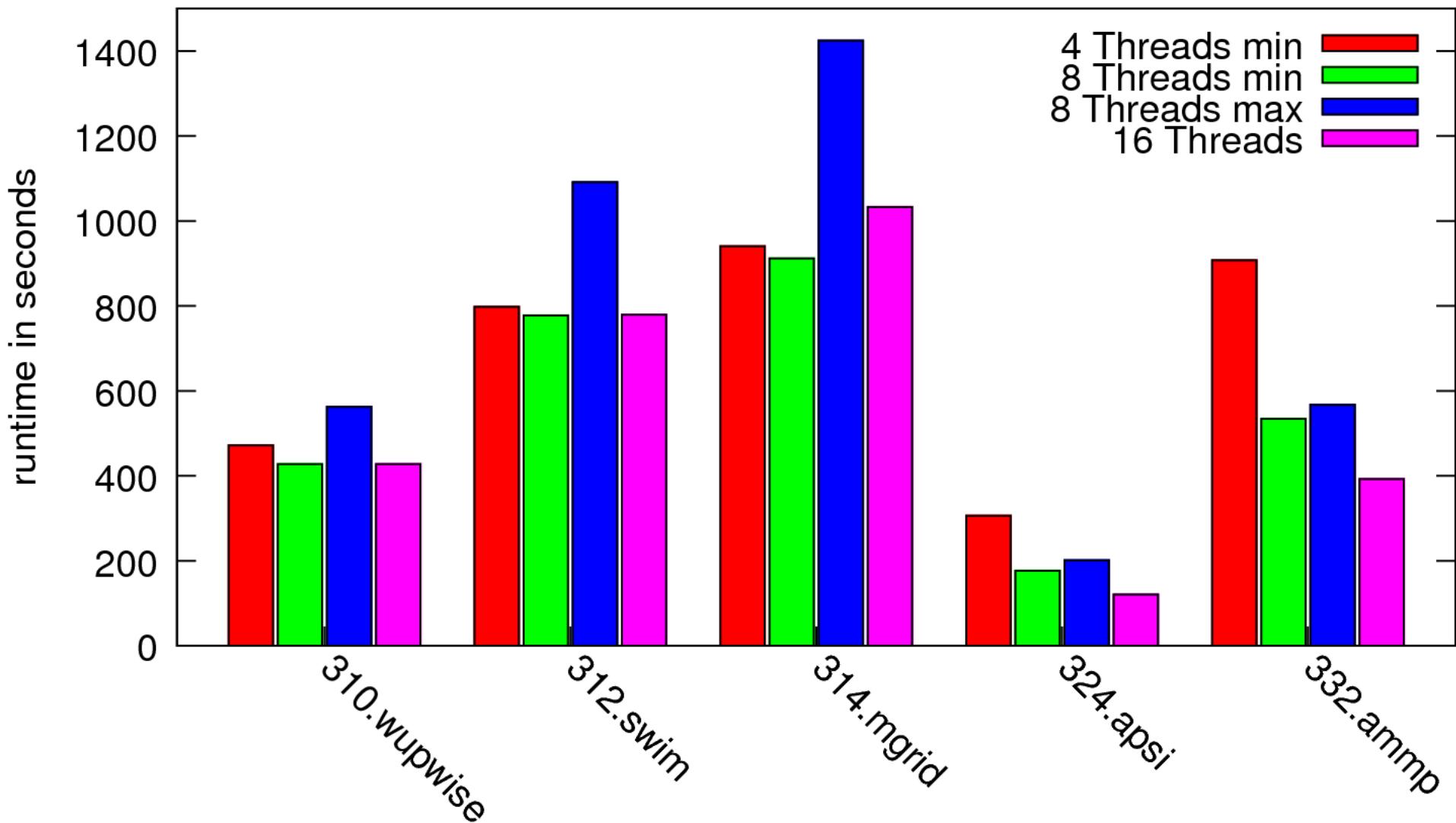
## Results SPEC OMP

- Optimal pinning found in all but 2 cases  
autopin's alternative less than 5% slower
- Overhead less than 3% on UMA platform
- Overhead less than 7,5% on NUMA platform  
(Kernel level page migration)

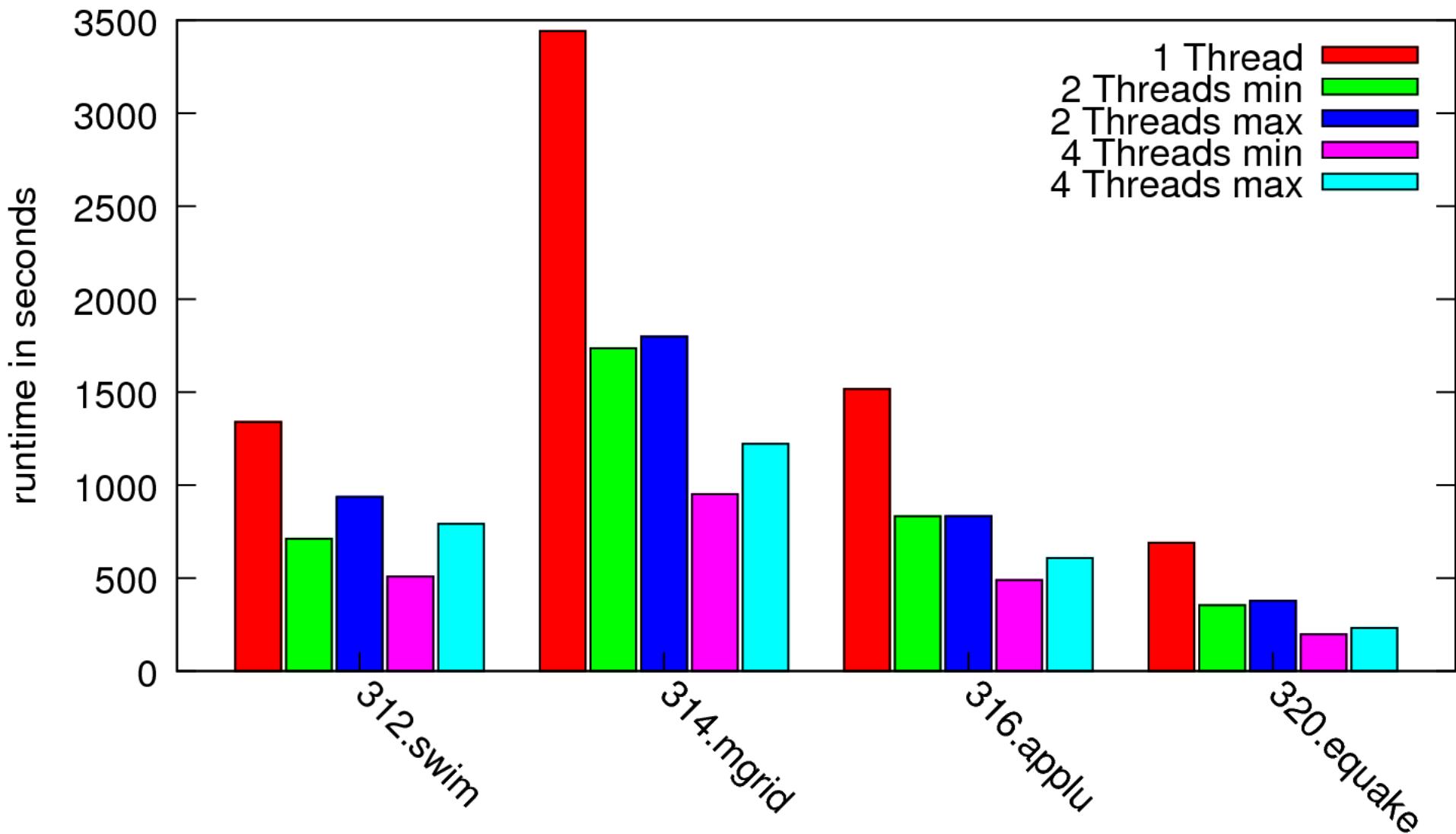
## SPEC OMP 1–4 Threads Caneland



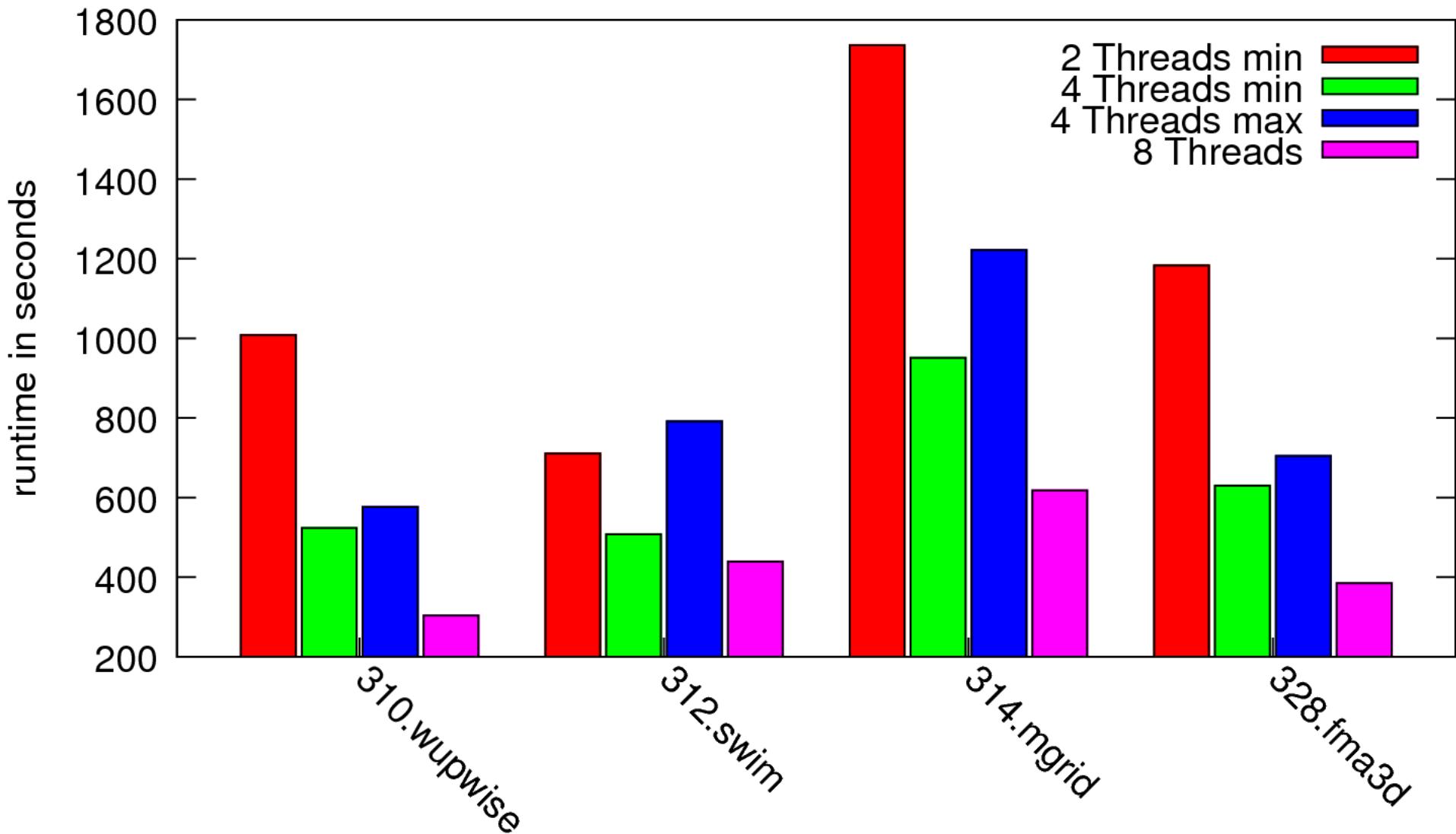
## SPEC OMP 4–16 Threads Caneland



## SPEC OMP 1–4 Threads Barcelona



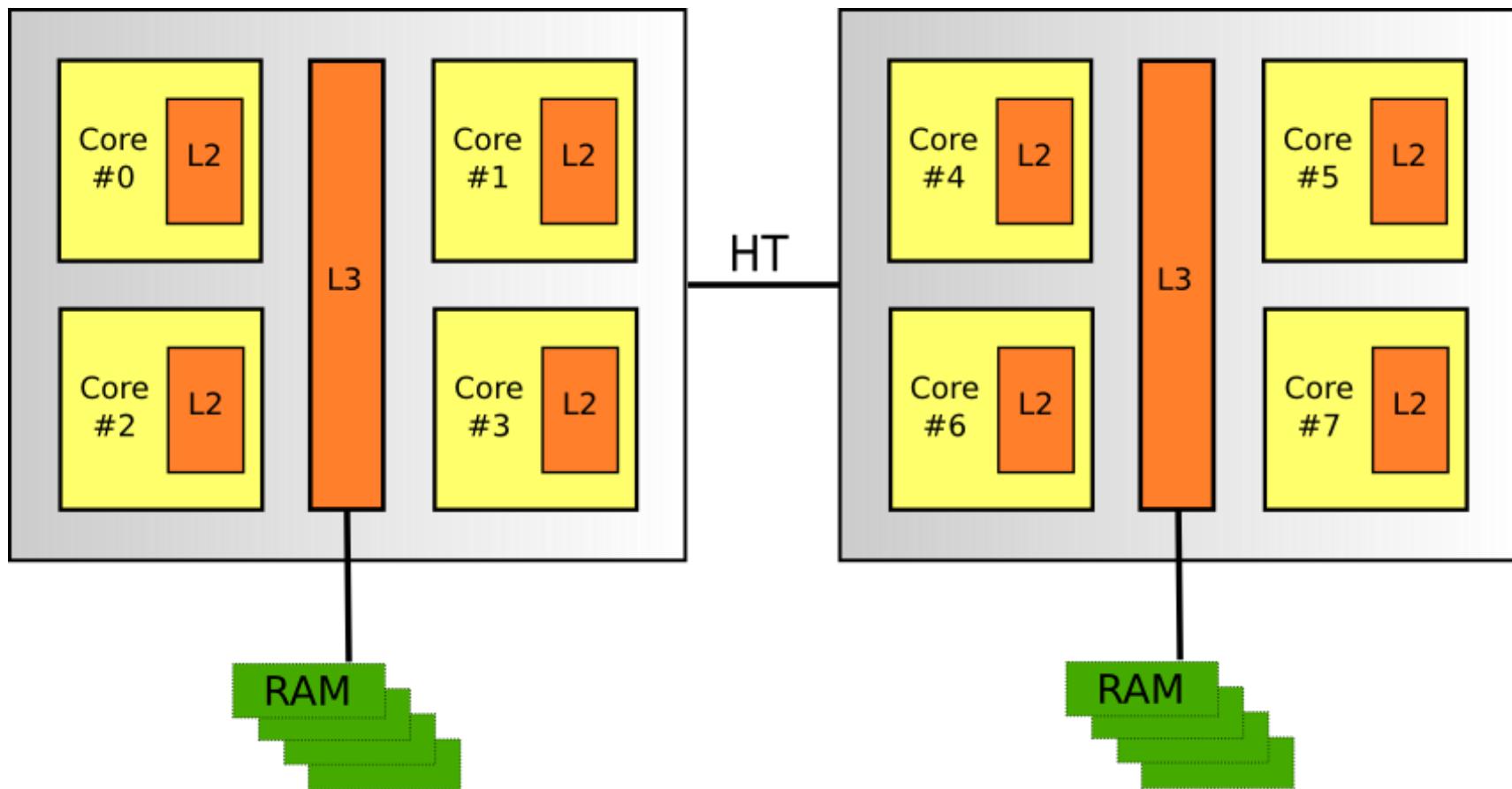
## SPEC OMP 4–8 Threads Barcelona



# Evaluation

- SPEC OMP
  - Benchmark consists of real scientific applications
  - OpenMP
  - PC: INSTRUCTIONS\_RETIRED
  - Several Multicore-Architekturen examined
- Memory Throughput
- MPI
- Electric power consumption

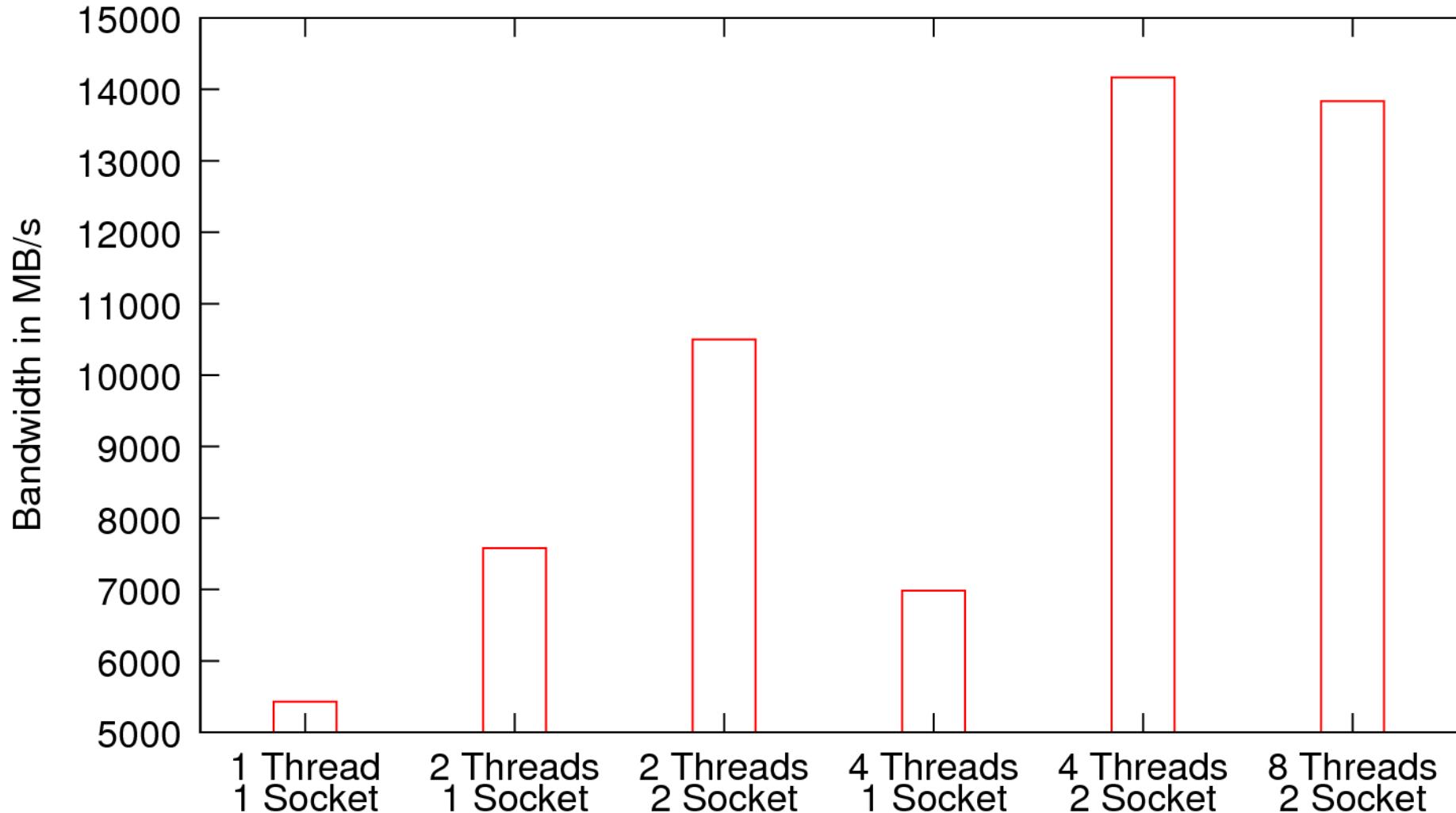
# Barcelona



# Memory Bandwidth

- STREAM John McCalpin
- synthetic benchmark
- copy and computation operations on large FP arrays
- Reusage of data avoided
- **copy:**  $a[i] = b[i]$   
**scale:**  $a[i] = q * b[i]$   
**sum:**  $a[i] = b[i] + c[i]$   
**triad:**  $a[i] = b[i] + q * c[i]$

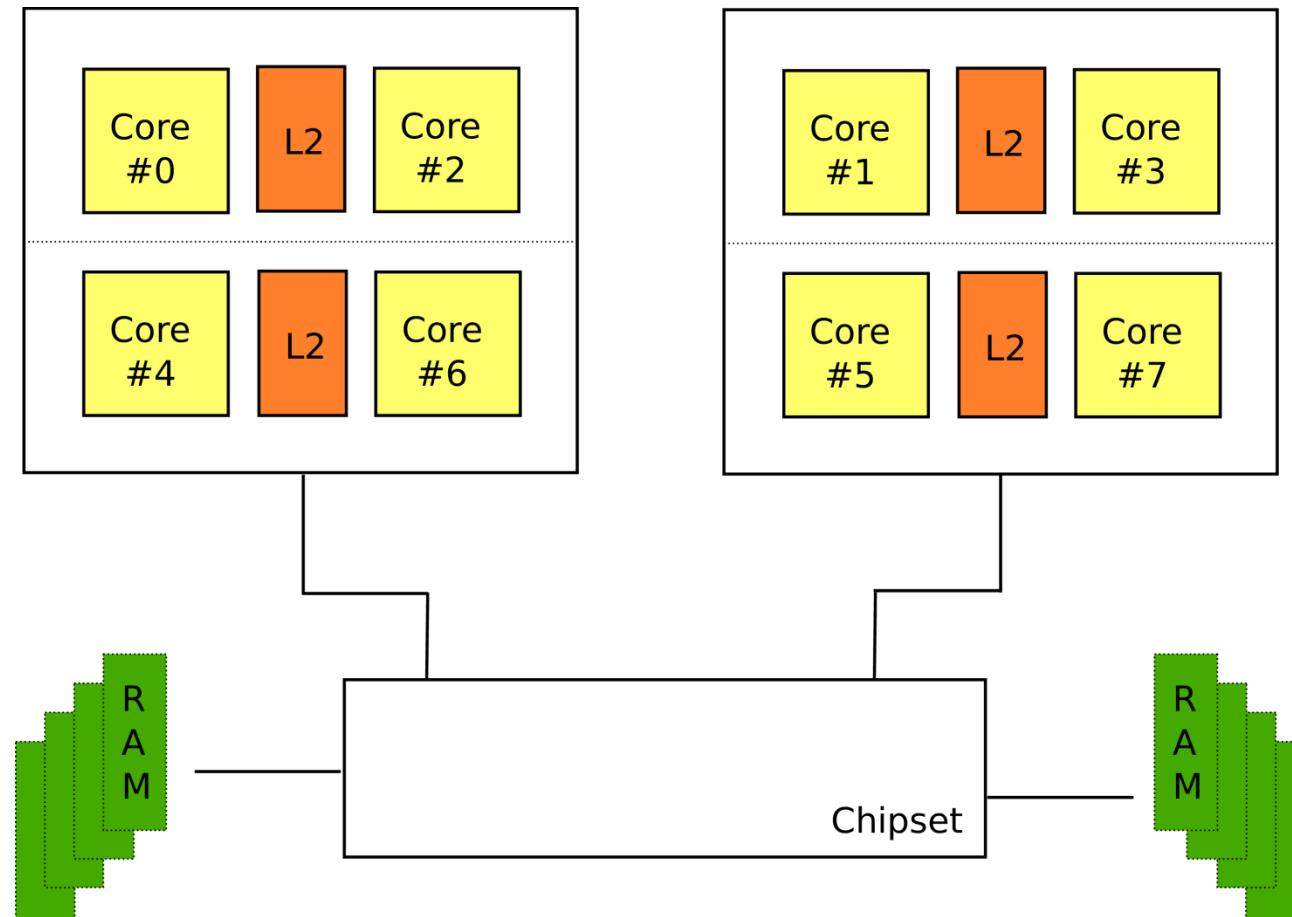
## Memory Bandwidth Barcelona



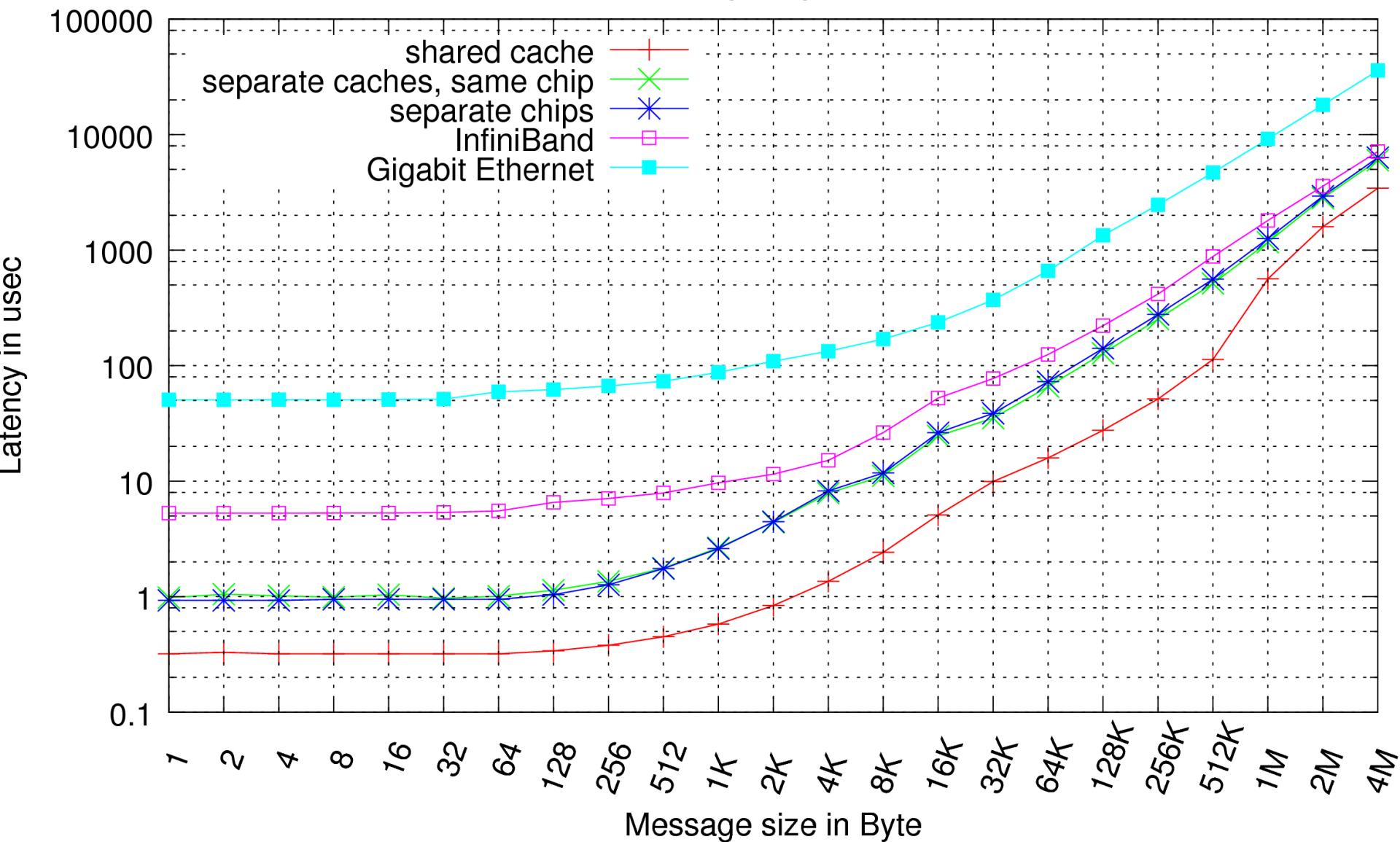
# Evaluation

- SPEC OMP
  - Benchmark consists of real scientific applications
  - OpenMP
  - PC: INSTRUCTIONS\_RETIRED
  - Several Multicore-Architekturen examined
- Memory Throughput
- MPI
- Electric power consumption

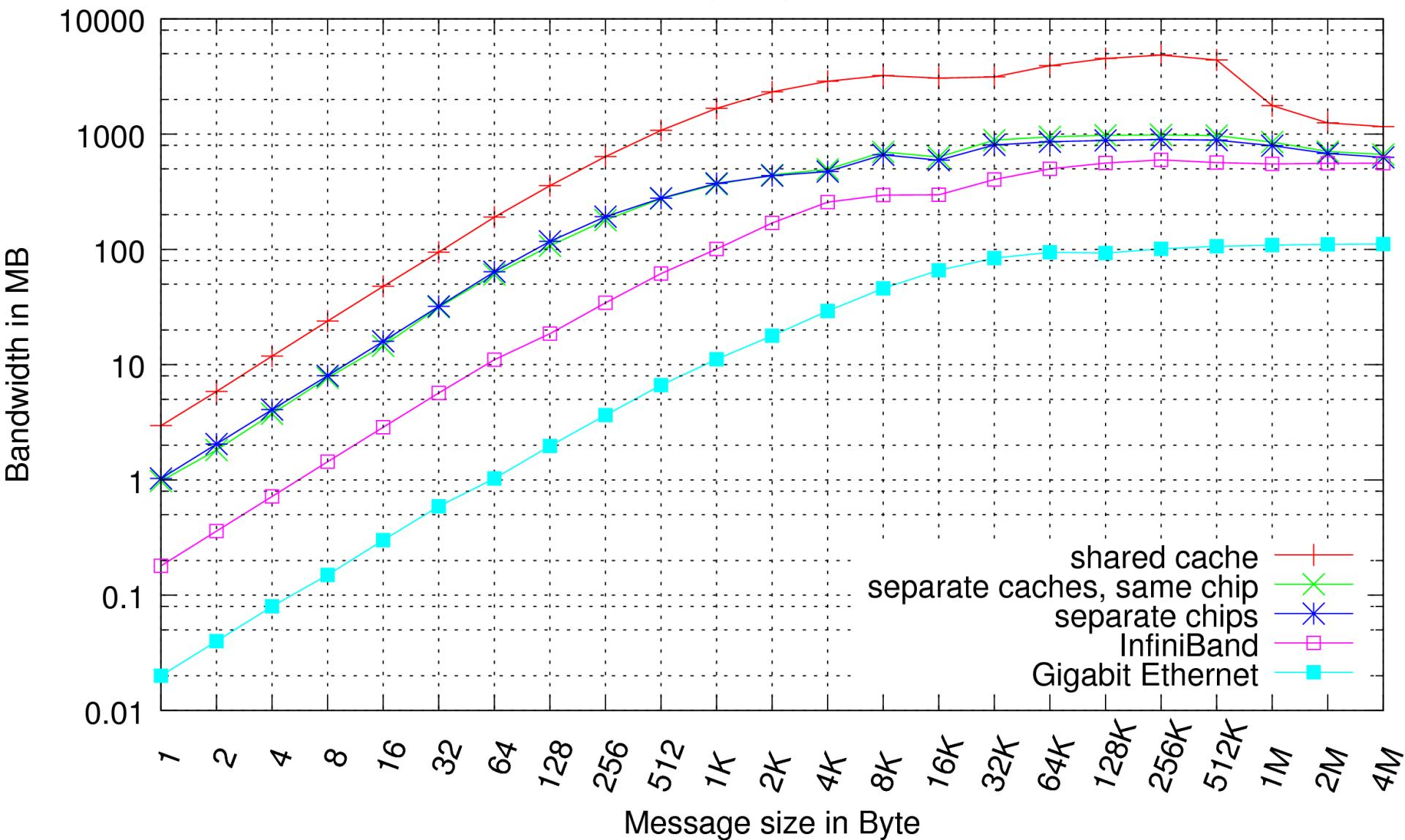
# Clovertown



# IMB PingPong Latency



# IMB PingPong Bandwidth

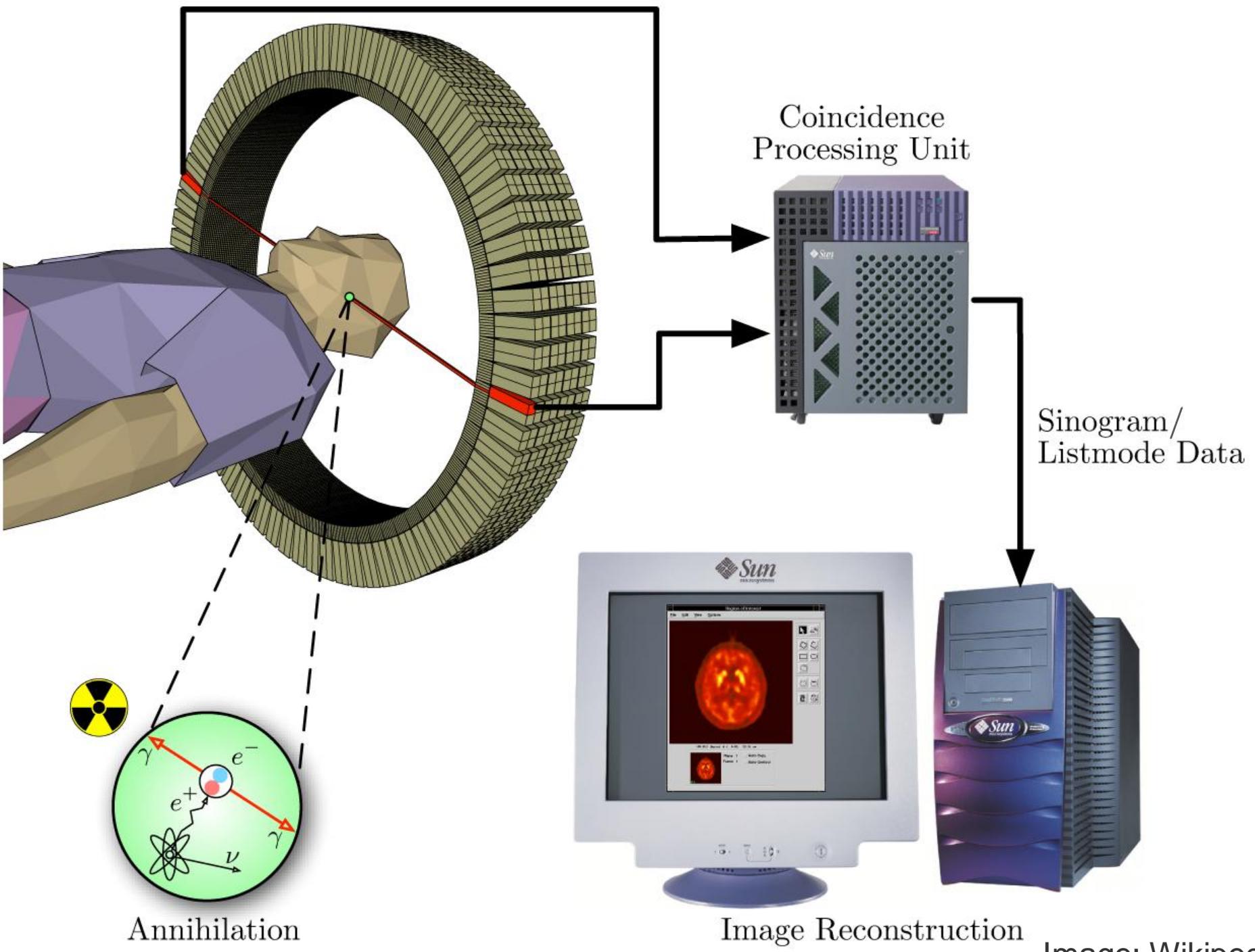


# Evaluation

- SPEC OMP
  - Benchmark consists of real scientific applications
  - OpenMP
  - PC: INSTRUCTIONS\_RETIRED
  - Several Multicore-Architekturen examined
- Memory Throughput
- MPI
- Electric power consumption

# PET (Positron Emission Tomography)

- Nuclear medicine imaging
- Visualizes functional processes  
(e.g. tumor diagnostics)
- fixed detector ring around patient
- radioisotopes injected into body
- Positron vs. electron → 2 photons 180 degree
- coincidence circuit



# Image Reconstruction

$$\underline{g} = \underline{\mathbf{A}} \underline{f}$$

$\underline{g}$  known measurement vector

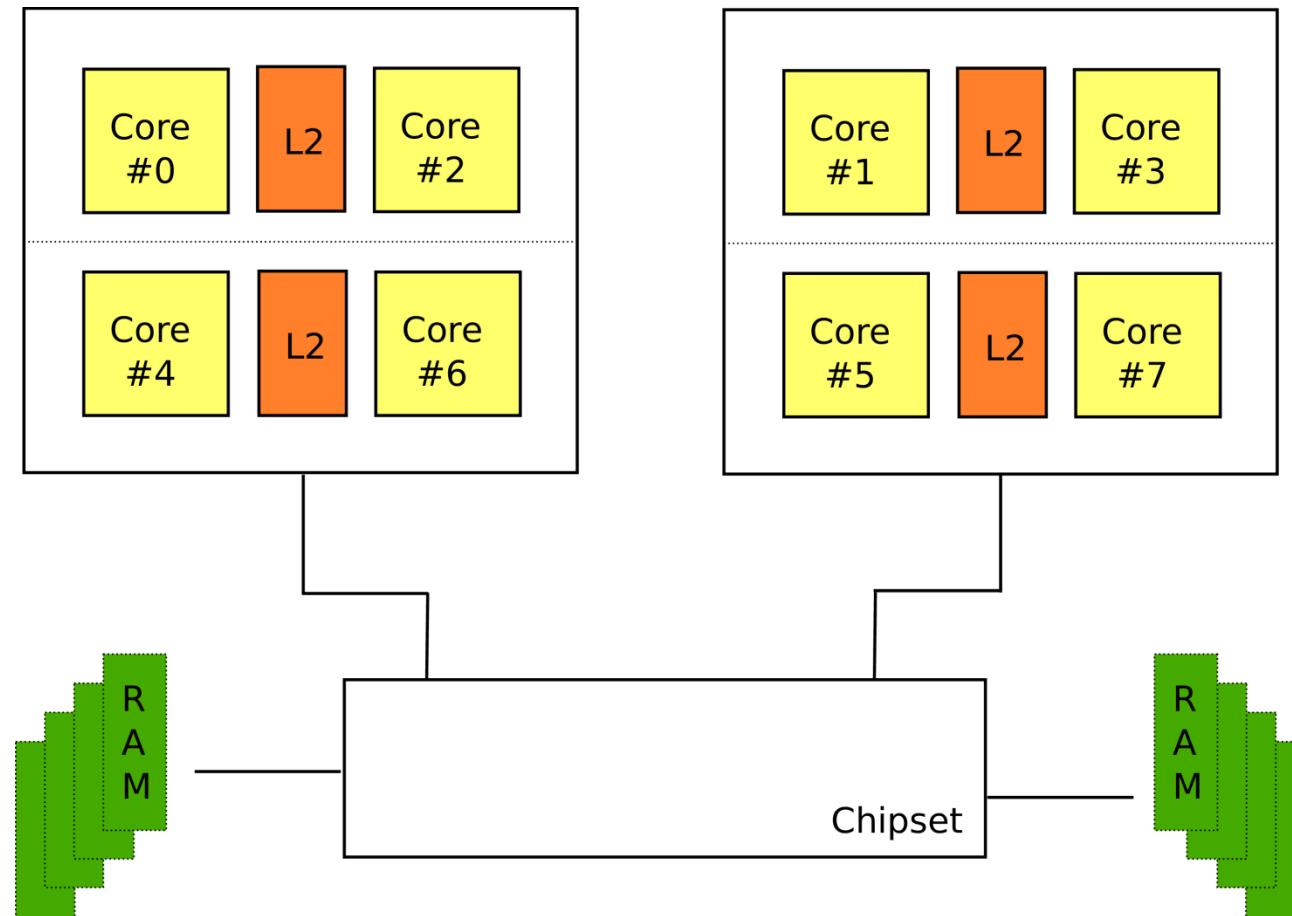
$\underline{f}$  unknown image vector

$\underline{\mathbf{A}}$  system matrix

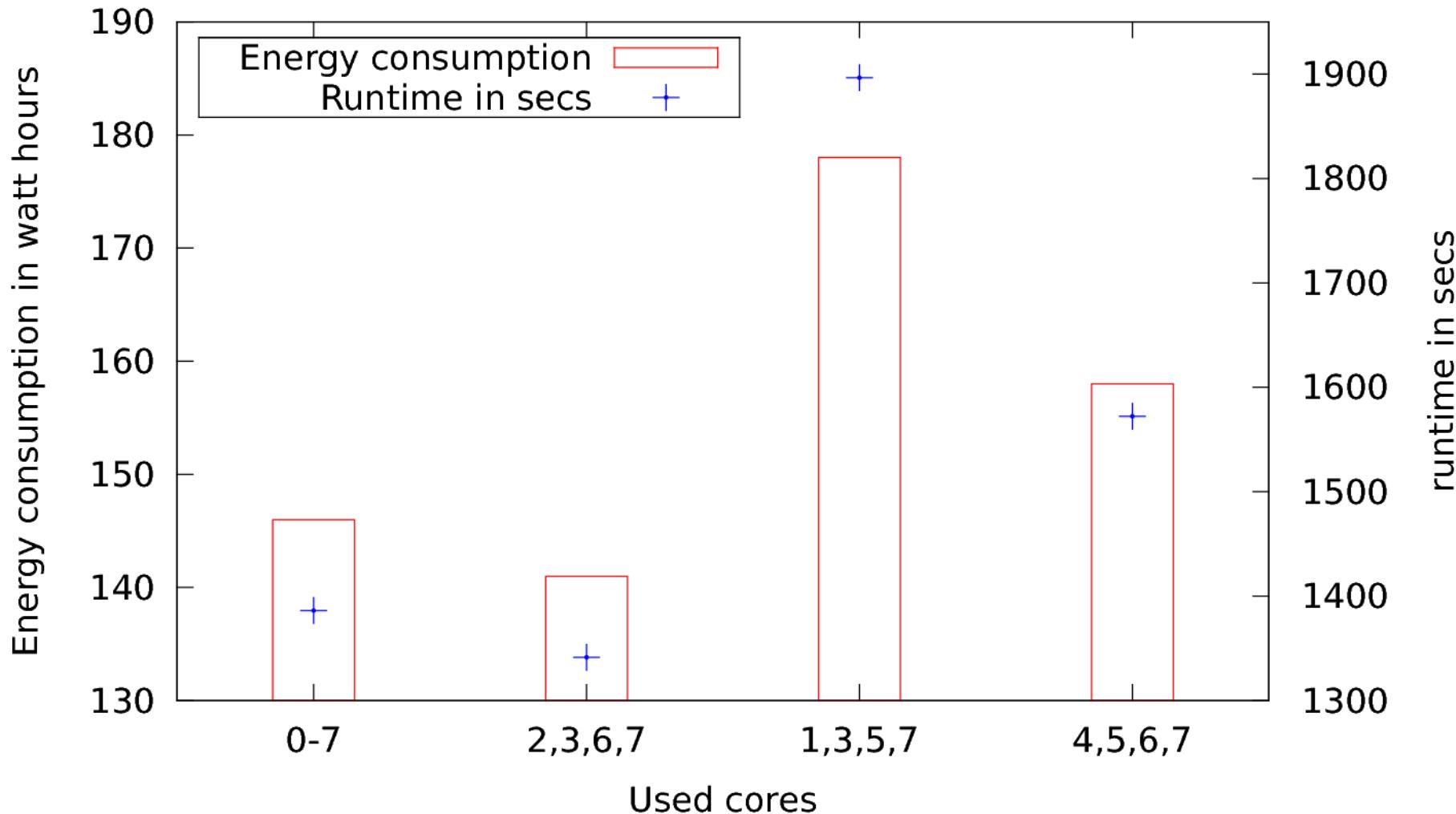
(describes characteristics of detector ring)

MLEM approximates linear system

# Clovertown



## Energy consumption petMLEM



# Conclusion and Outlook

- Pinning is essential on multicore systems
- Will become even more important on many core architectures
- Tools can reliably find optimal pinnings on UMA and NUMA architectures
- Outlook – autopin2
  - new design: perf performance counters subsystem
  - Flexible and modular Design:
    - perfmon, perf
    - Energy, runtime, user defined objective functions
    - Back channel from application to autopin2

# Questions?

