

# Optimizing for Cache/Memory

---

Erik Hagersten  
Uppsala University

# Optimizing for the memory system: What is the potential gain?

- Latency difference between L1\$ and mem: ~50x
- Bandwidth difference L1\$ and mem on a MC: ~100x
- At least a factor 2-10x is within reach when optimizing for caches
- Optimizations tactics: Make applications access their data from [L1] cache instead of from memory



# Four tricks for better performance

- Keep the active footprint small
- Use the entire cache line once it has been brought into the cache
- Fetch a cache line prior to its usage
- Avoid putting "ill-suited" data structures into the cache



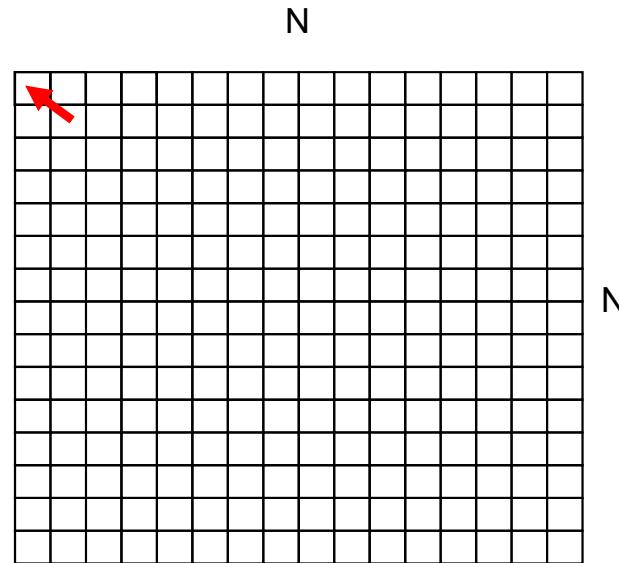
# Final cache lingo slide

- **Miss ratio:** What is the likelihood that a memory access will miss in a cache [%]?
- **Miss rate:** D:o per time unit, e.g. per-second, per-1000-instructions
- **Fetch ratio/rate<sup>\*</sup>:** What is the likelihood that a memory access will cause a fetch to the cache [including HW/SW prefetching]
- **Fetch utilization<sup>\*</sup>:** What fraction of a cacheline was used before it got evicted
- **Writeback utilization<sup>\*</sup>:** What fraction of a cacheline written back to memory contains dirty data



# What can go Wrong? A Simple Example...

Perform a diagonal copy 10 times



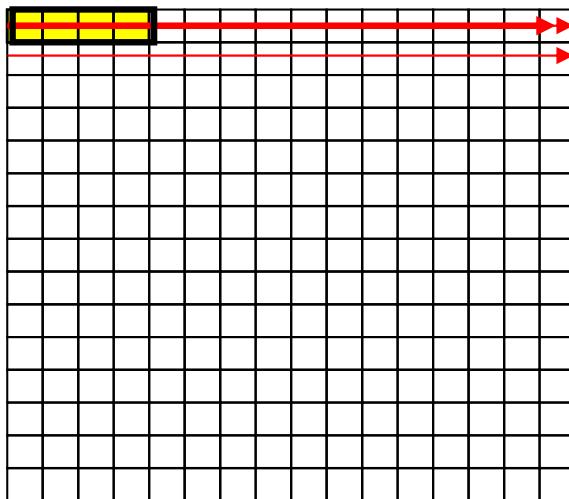


Note that these examples also assume an outer loop that will repeat the code several times...

# Example: Loop order

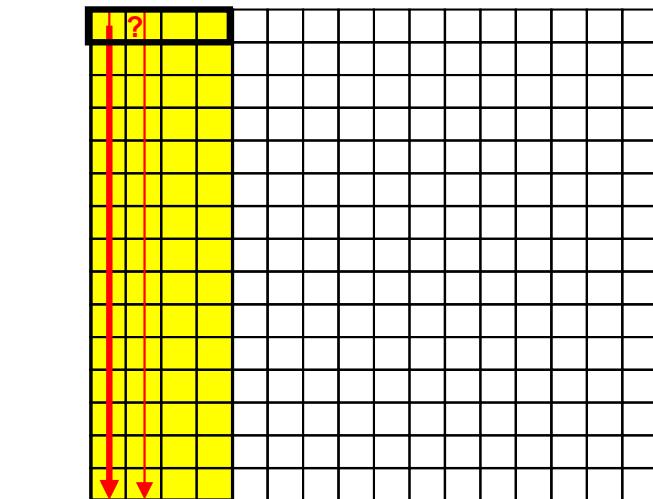
//Optimized Example A

```
for (i=1; i<N; i++) {  
    for (j=1; j<N; j++) {  
        A[i][j]= A[i+1][j+1];  
    }  
}
```

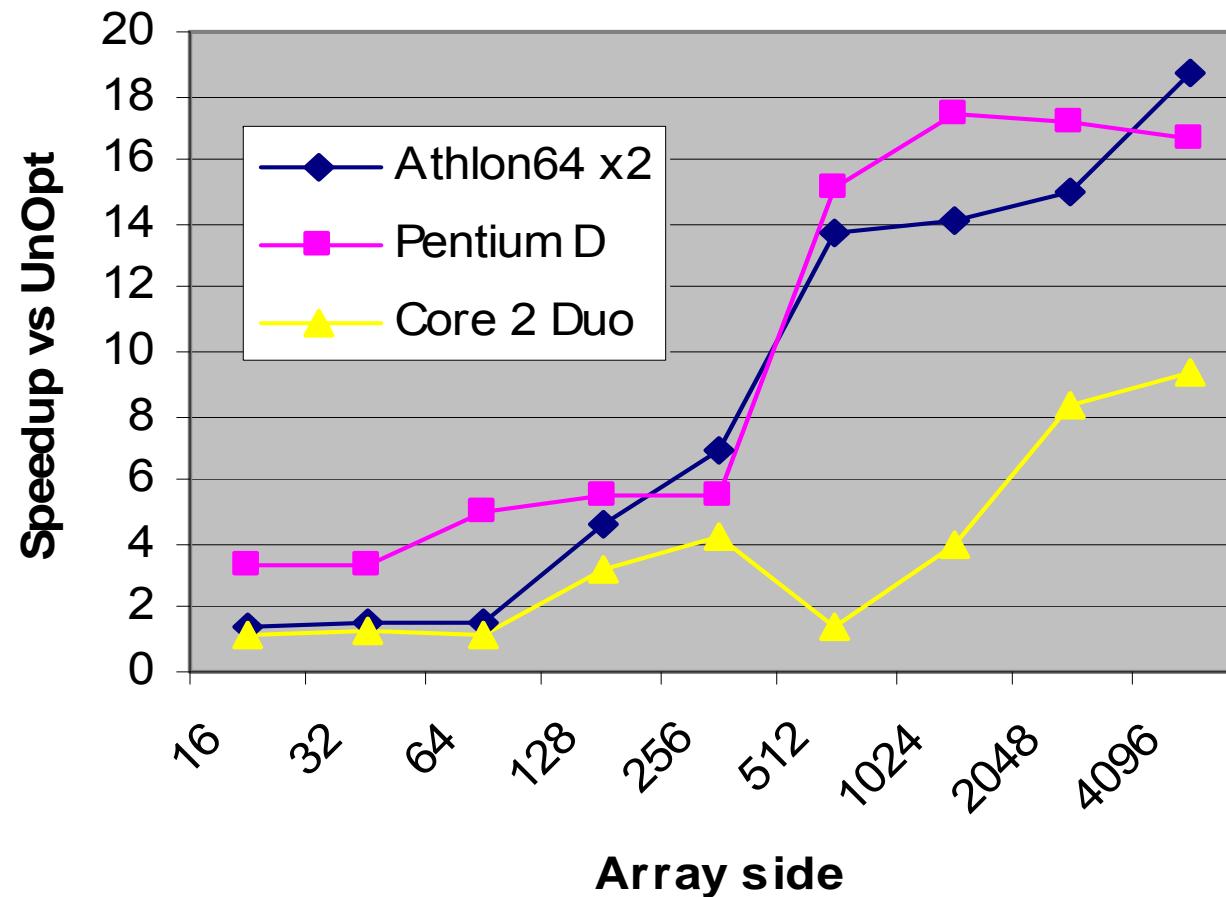


//Unoptimized Example A

```
for (j=1; j<N; j++) {  
    for (i=1; i<N; i++) {  
        A[i][j] = A[i+1][j+1];  
    }  
}
```

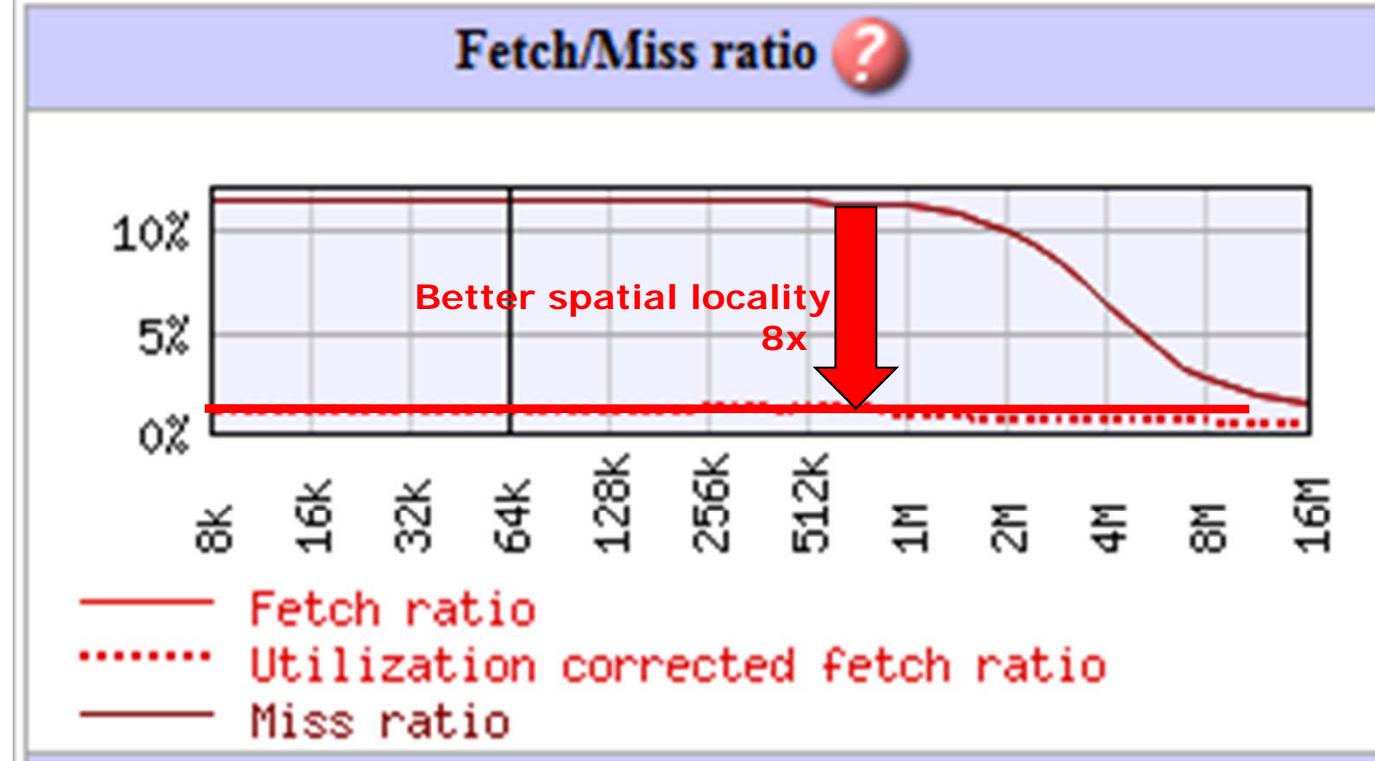


# Performance Difference: Loop order

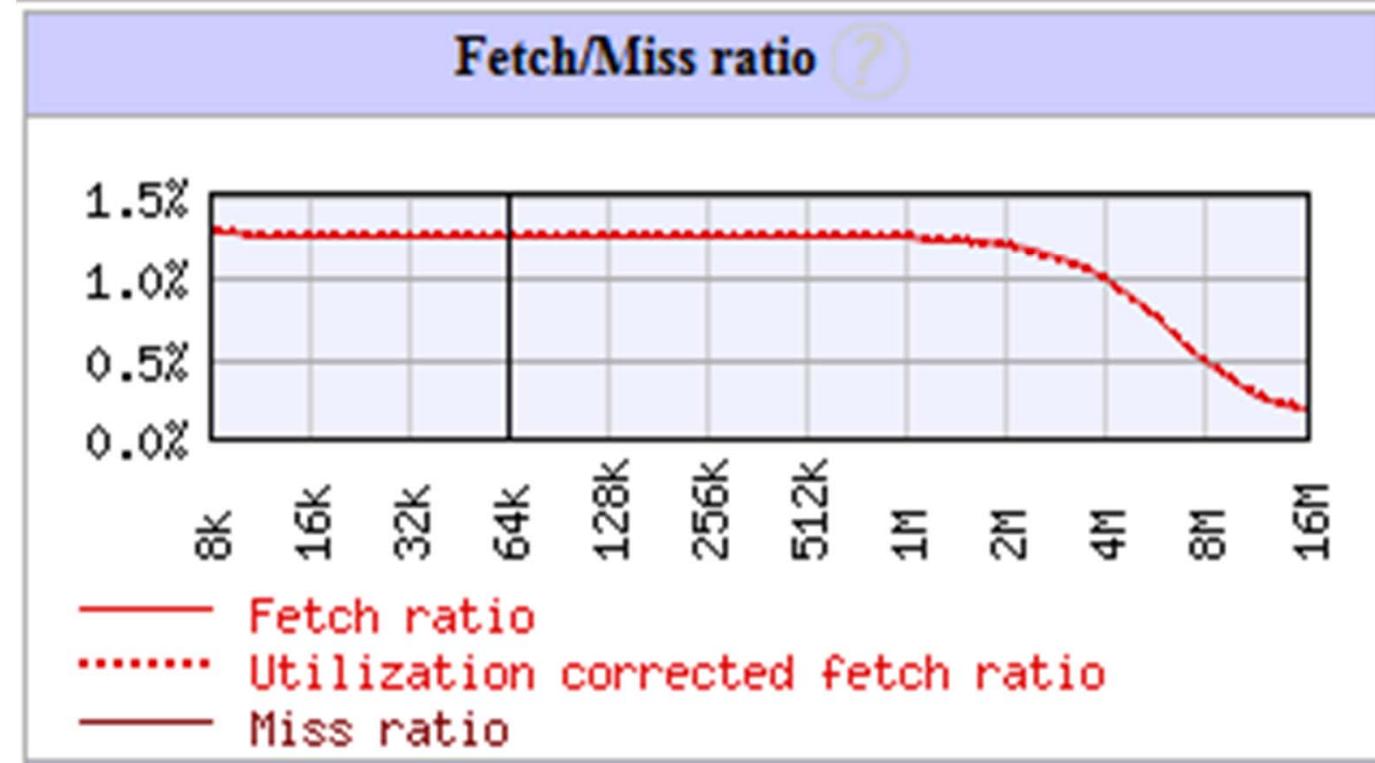




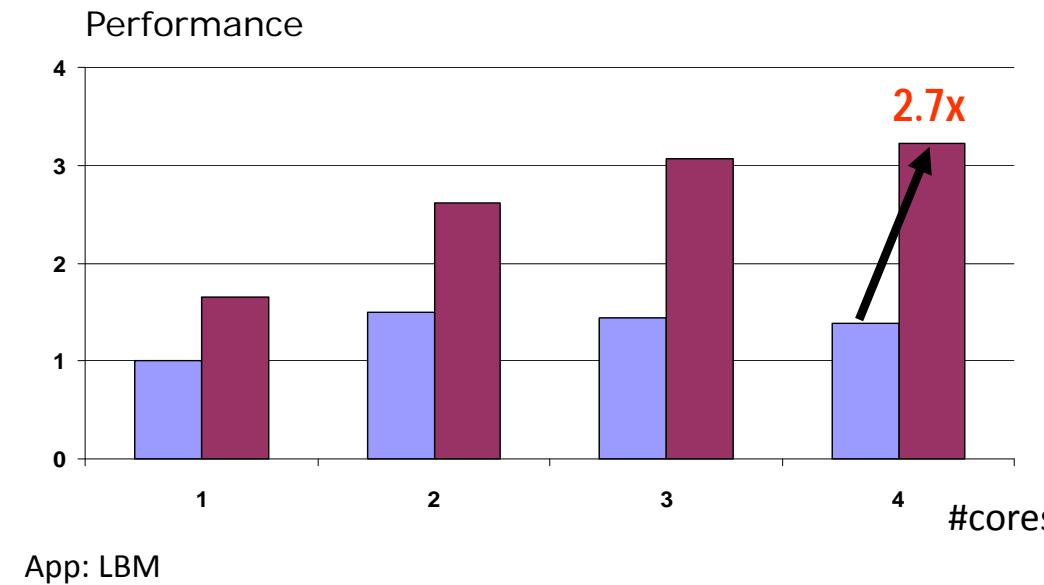
UNOPT:



OPT:



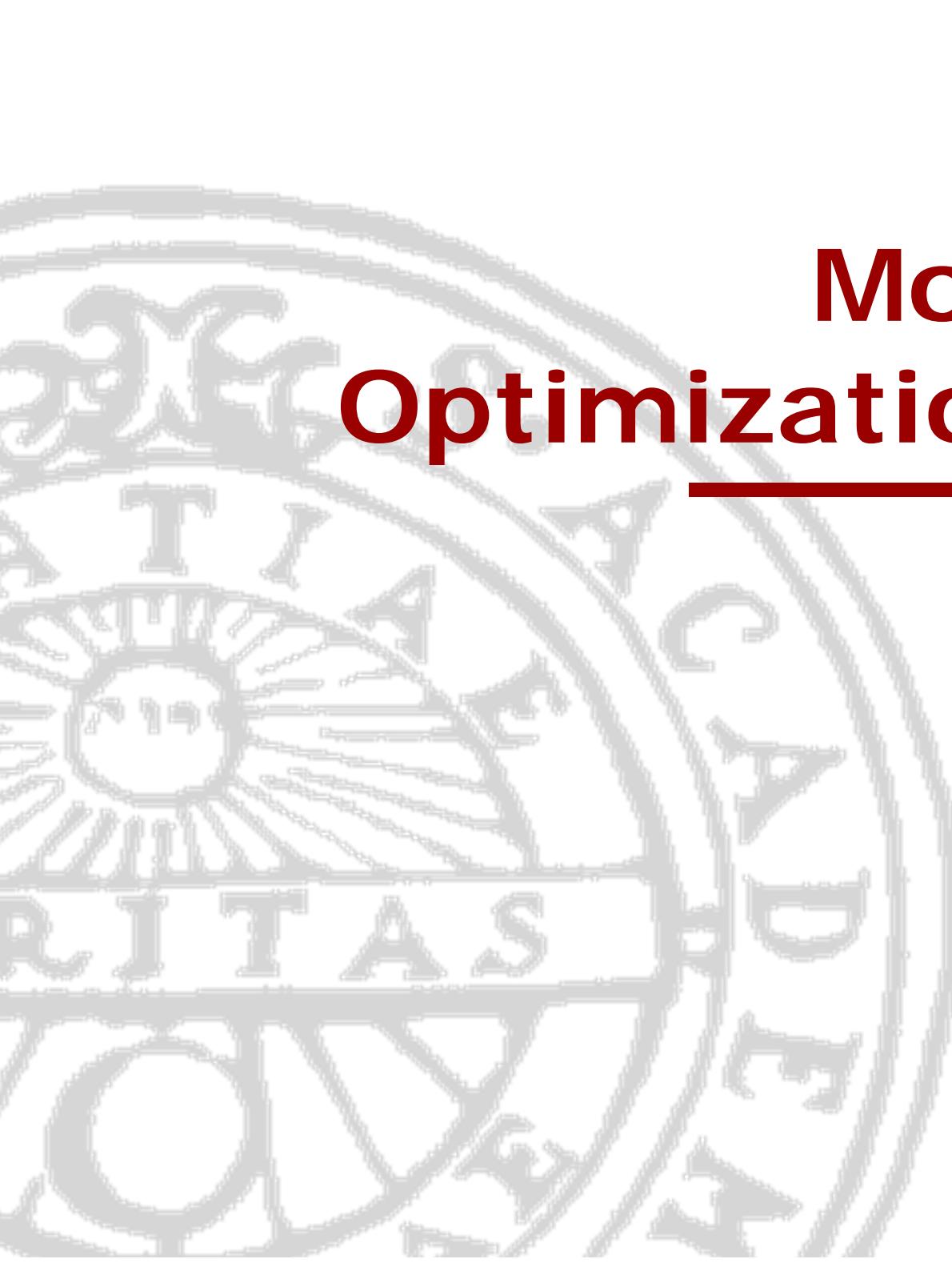
# Example: LBM



Optimization can be rewarding, but costly...

- Require expert knowledge about MC and architecture
- Weeks of wading through performance data

→ This fix required one line of code to change



# More Software Optimization Examples

---

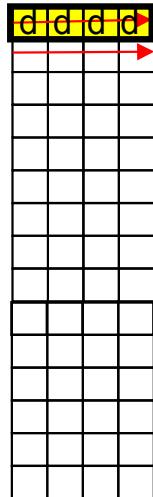
Erik Hagersten  
Uppsala University



# Example: Sparse data utilization

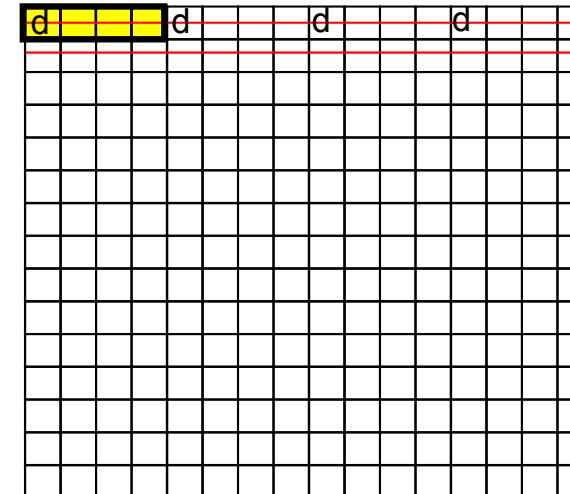
```
//Optimized Example A
for (i=1; i<N; i++) {
    for (j=1; j<N; j++) {
        A_d[i][j]= A_d[i-1][j-1];
    }
}
```

Better spatial locality



↔  
Smaller array

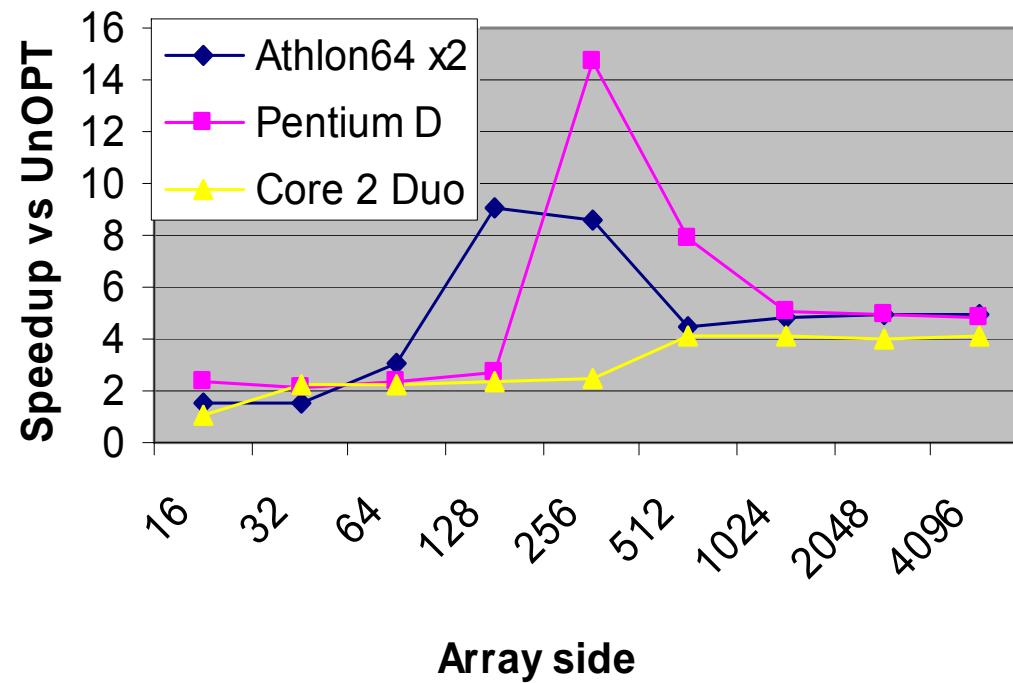
```
//Unoptimized Example A
for (i=1; i<N; i++) {
    for (j=1; j<N; j++) {
        A[i][j].d = A[i-1][j-1].d;
    }
}
```



```
struct vec_type
{
    int a;
    int b;
    int c;
    int d;
};
```

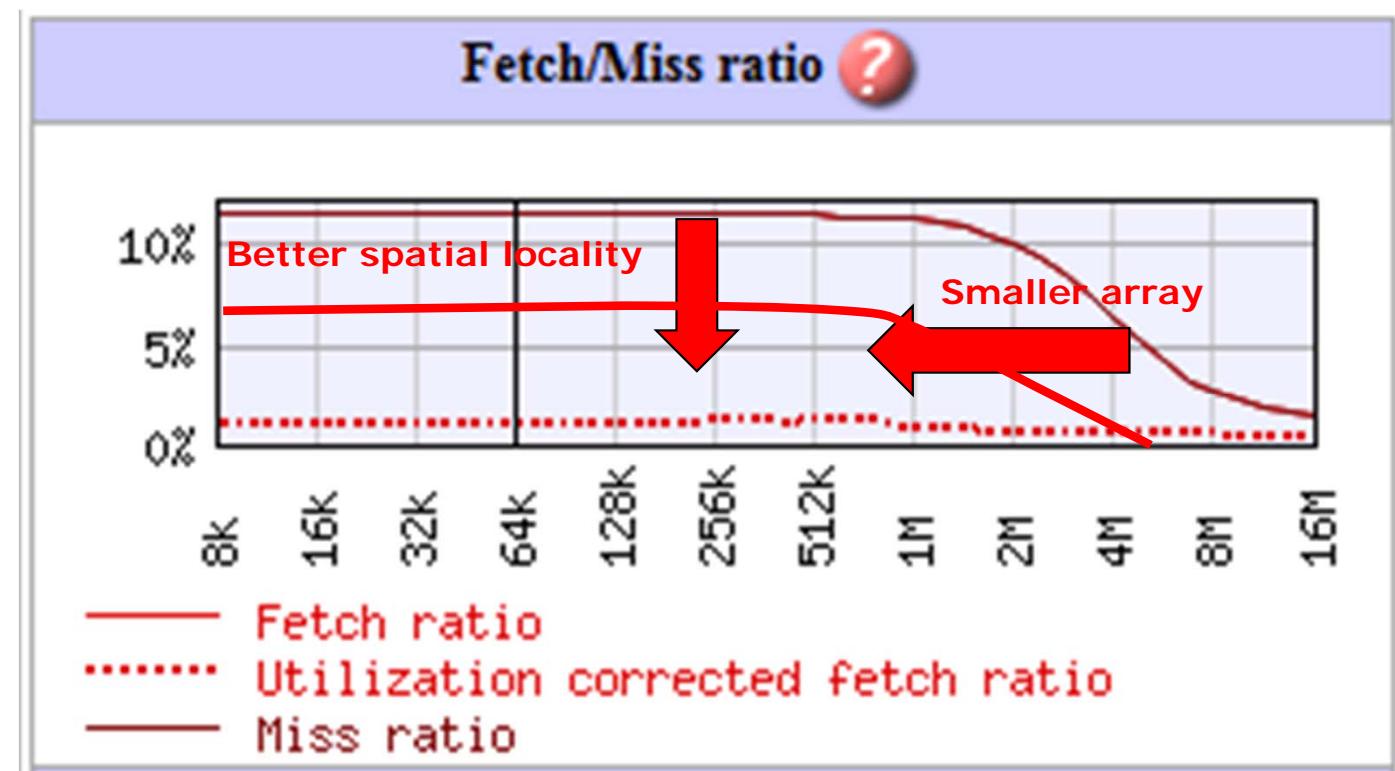


# Performance Difference: Sparse Data

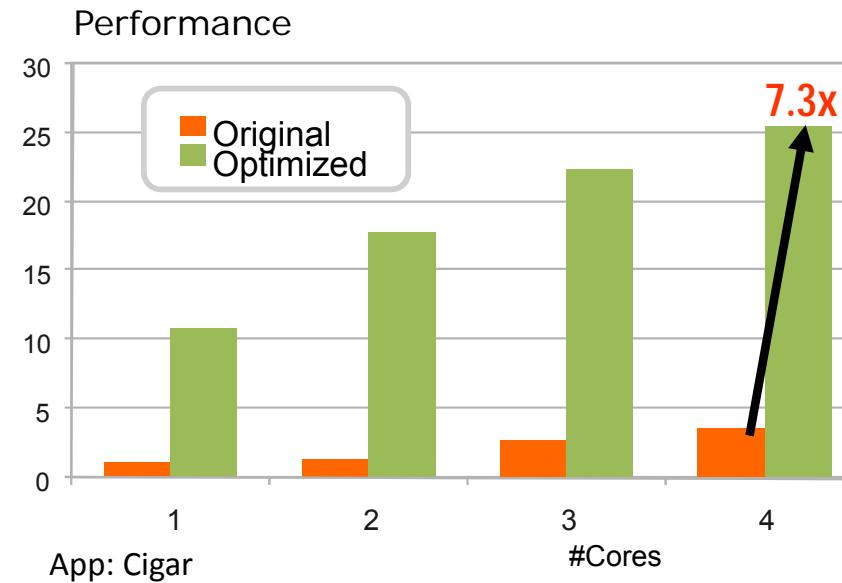




# Example: Sparse data utilization



# Example: Cigar



Looks like a perfect scalable application!

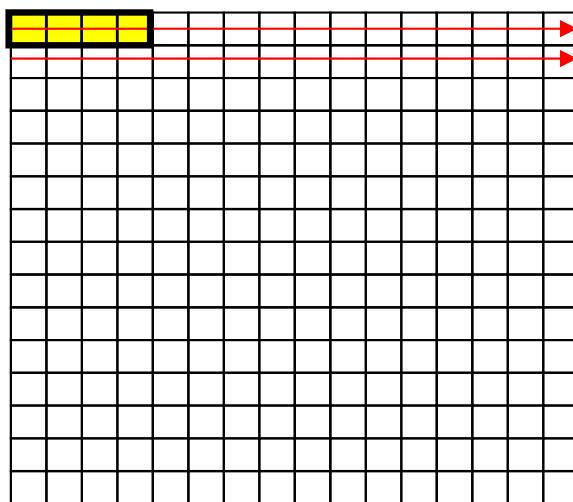
Are we done?

→ Duplicate one data structure



# Example: Sparse data allocation

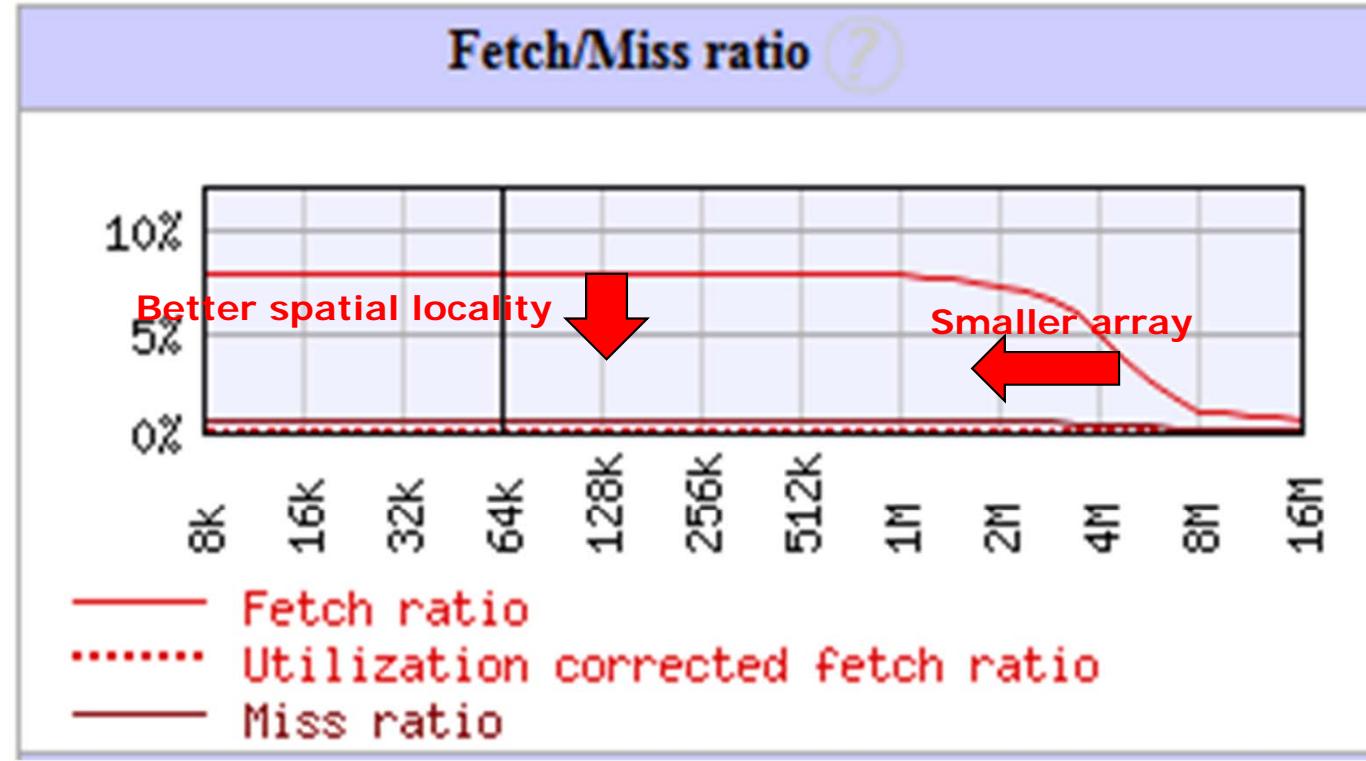
```
sparse_rec sparse [HUGE];  
  
for (int j = 0; j < HUGE; j++)  
{  
    sparse[j].a = 'a'; sparse[j].b = 'b'; sparse[j].c = 'c'; sparse[j].d = 'd'; sparse[j].e = 'e';  
    sparse[j].f1 = 1.0; sparse[j].f2 = 1.0; sparse[j].f3 = 1.0; sparse[j].f4 = 1.0; sparse[j].f5 = 1.0;  
}
```



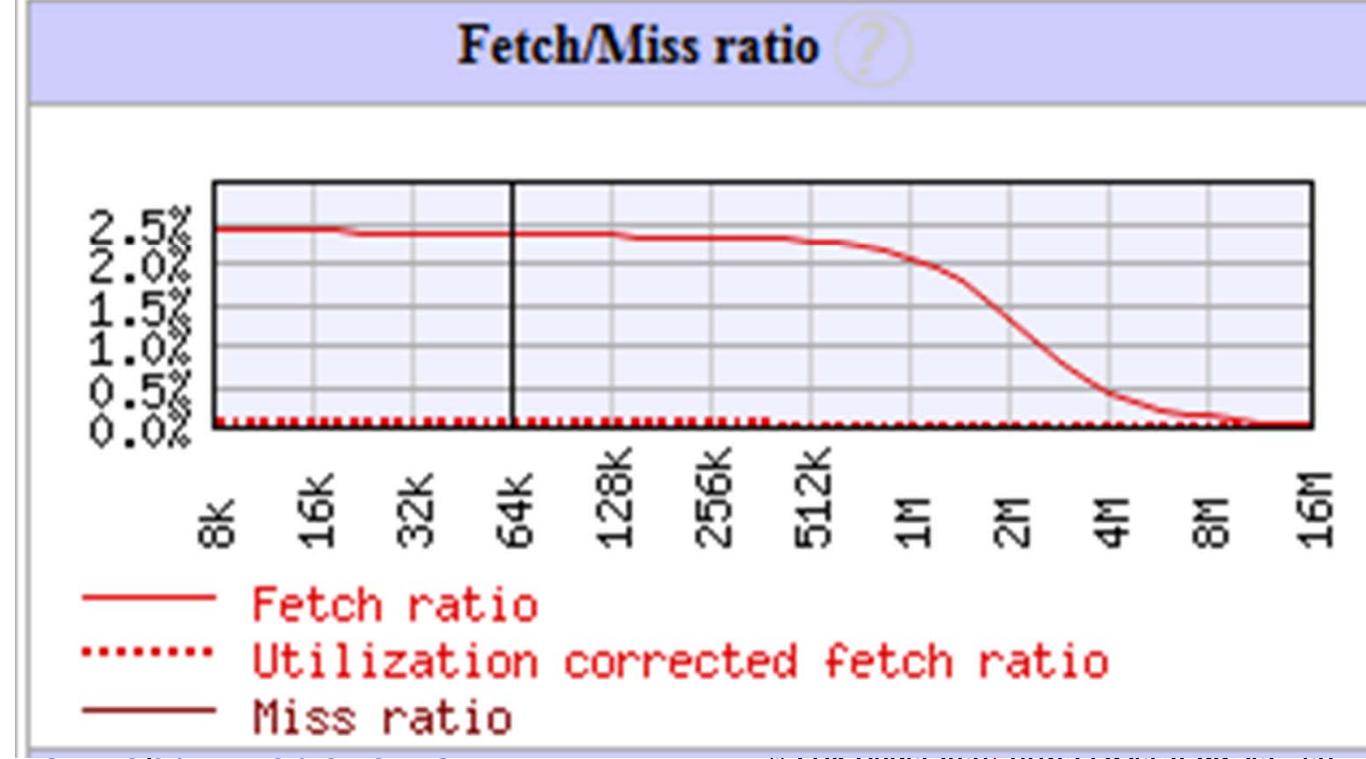
<p>a      f1    b      f2    ...      f1    f2    f3    f4    f5    a b c d e</p> <pre>struct sparse_rec {     // size 80B     char a;     double f1;     char b;     double f2;     char c;     double f3;     char d;     double f4;     char e;     double f5; };</pre>	<pre>struct dense_rec {     //size 48B     double f1;     double f2;     double f3;     double f4;     double f5;     char a;     char b;     char c;     char d;     char e; };</pre>
--	--



## UNOPT:



## OPT:





# Loop Fusion

```
/* Unoptimized */

for (i = 0; i < N; i = i + 1)
    for (j = 0; j < N; j = j + 1)
        a[i][j] = 2 * b[i][j];

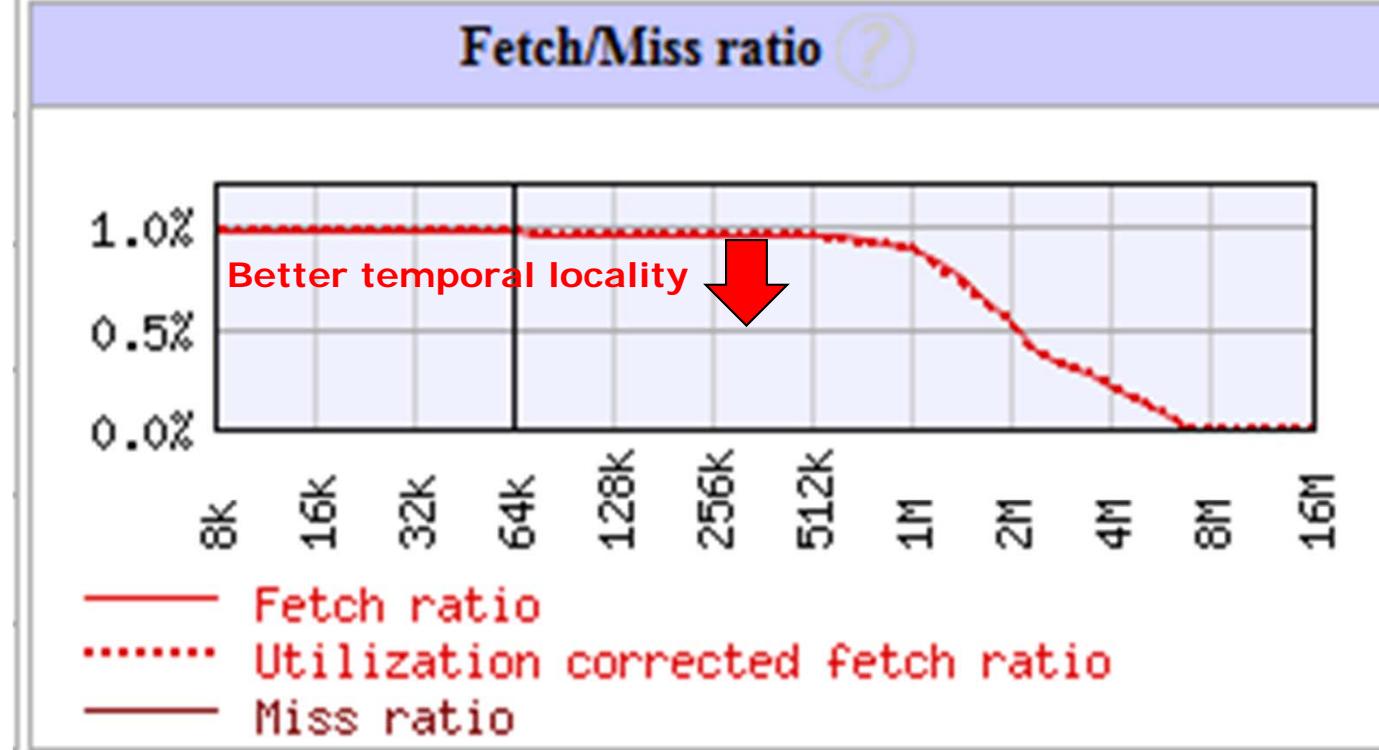
for (i = 0; i < N; i = i + 1)
    for (j = 0; j < N; j = j + 1)
        c[i][j] = K * b[i][j] + d[i][j]/2

/* Optimized */

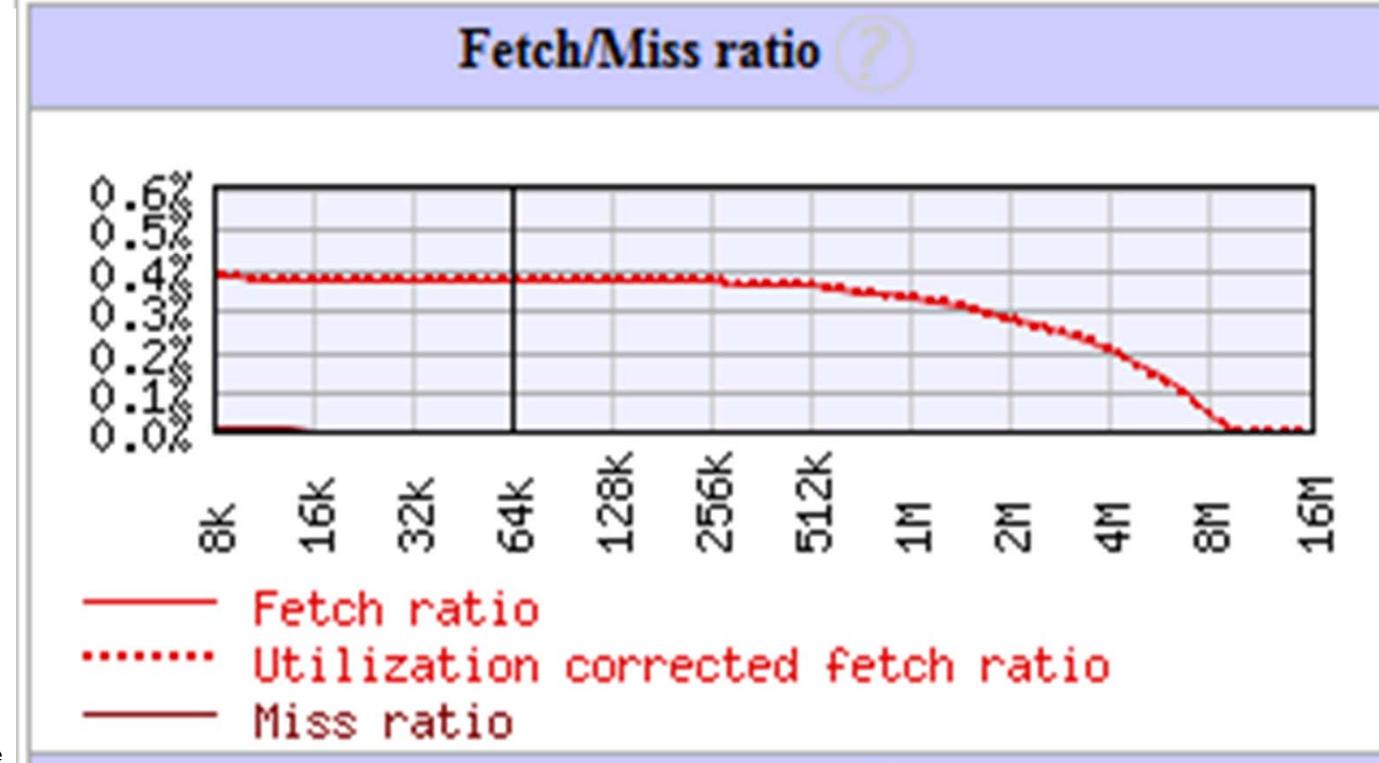
for (i = 0; i < N; i = i + 1)
    for (j = 0; j < N; j = j + 1) {
        a[i][j] = 2 * b[i][j];
        c[i][j] = K * b[i][j] + d[i][j]/2;
    }
```



## UNOPT:

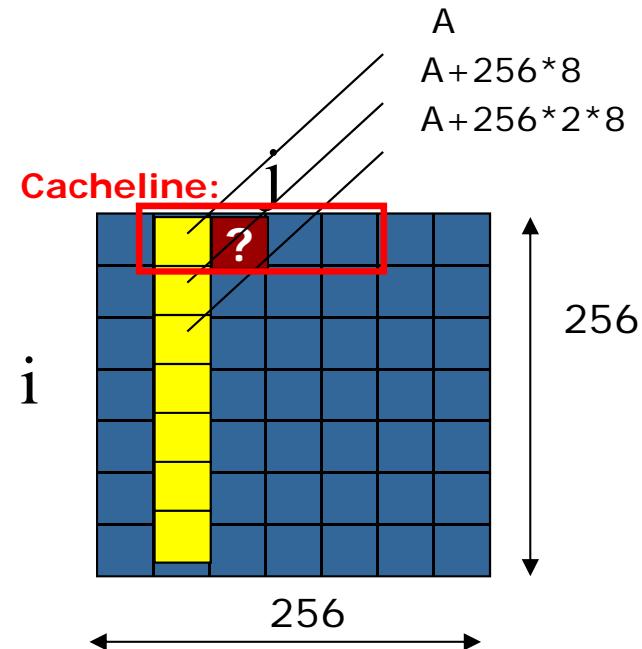


## OPT:

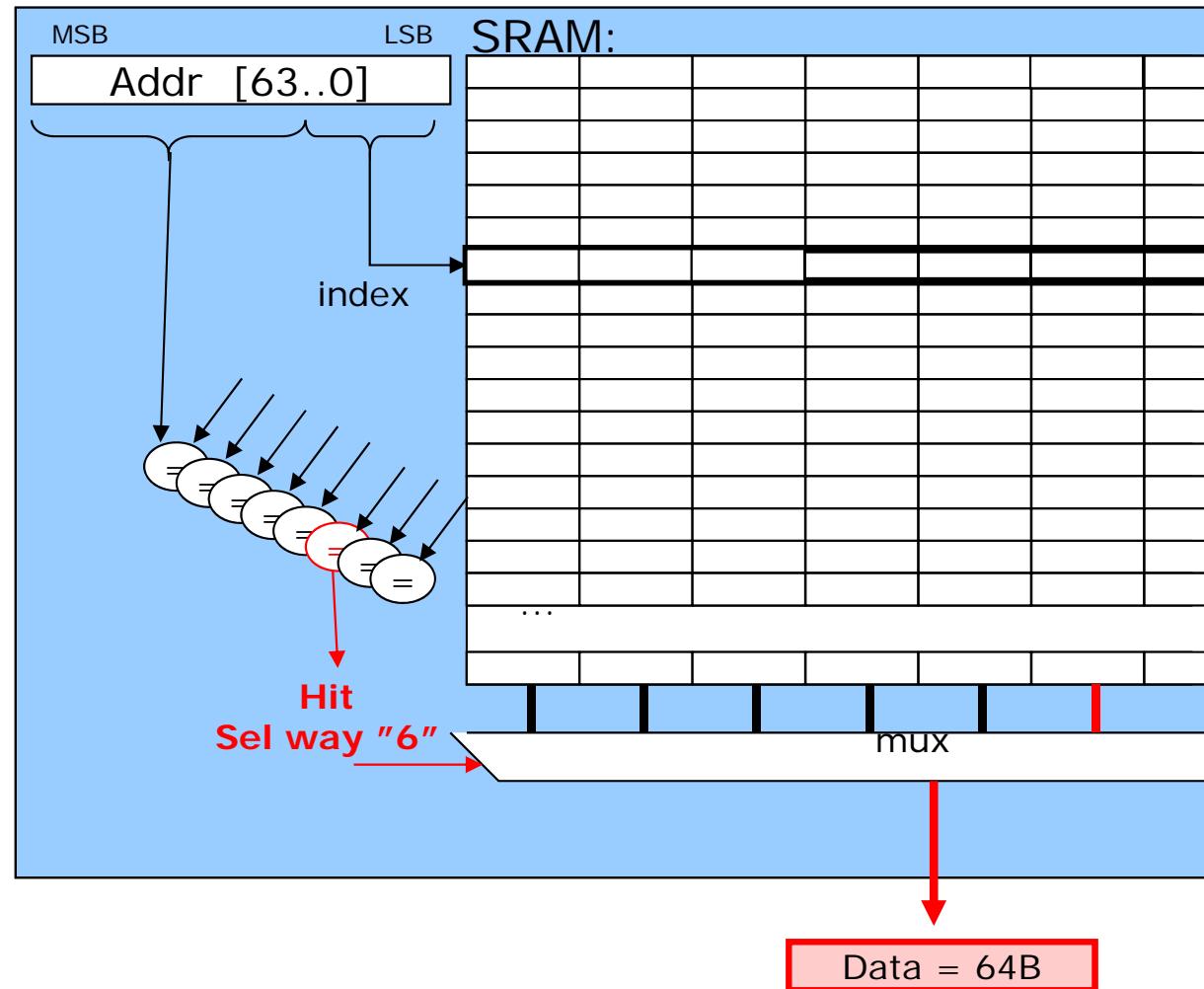




# Padding of data structures

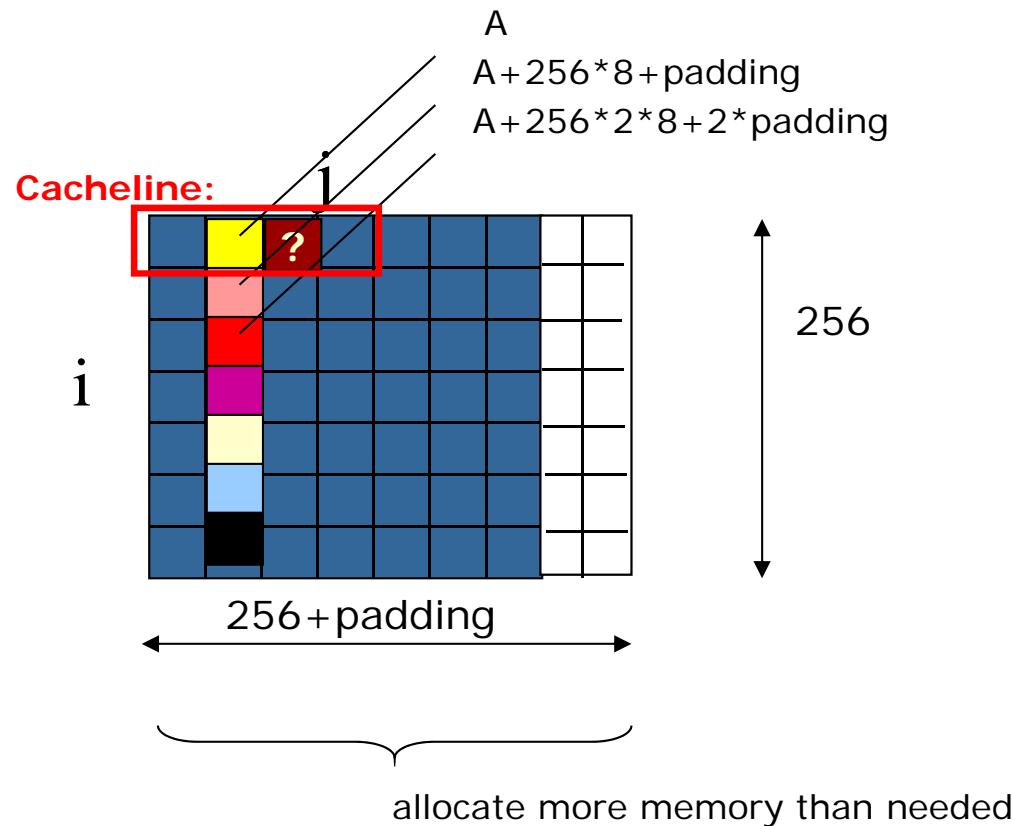


Generic Cache:

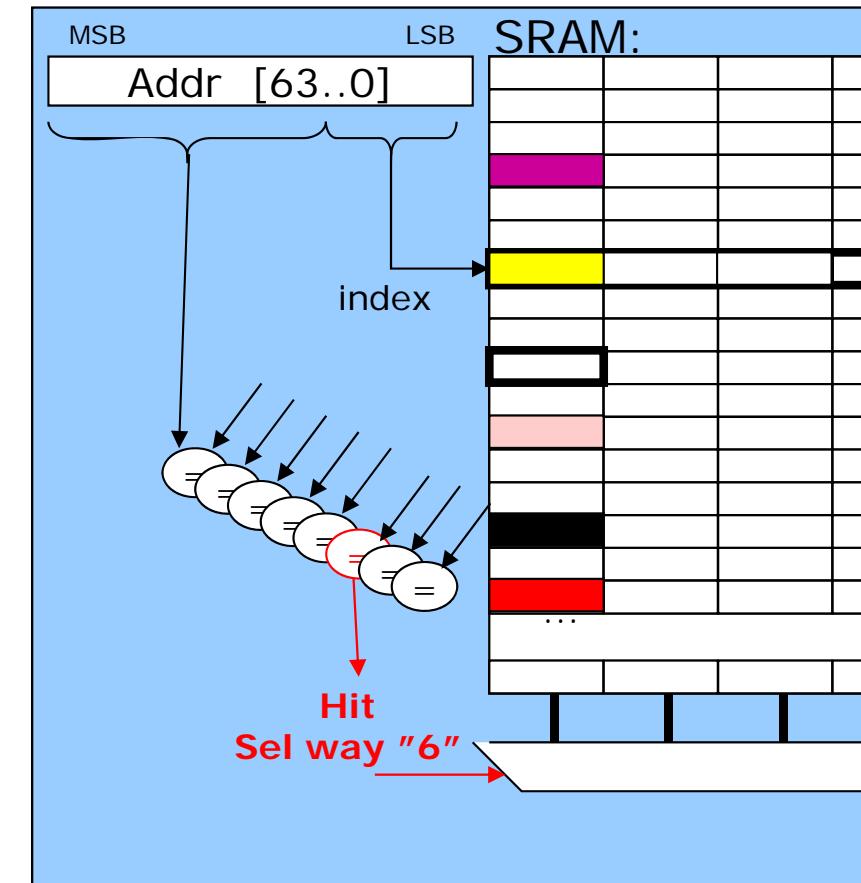




# Padding of data structures



Generic Cache:

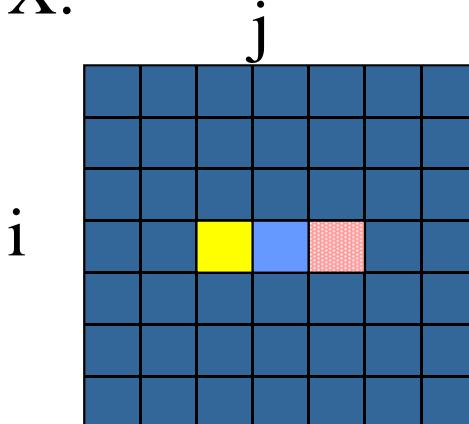




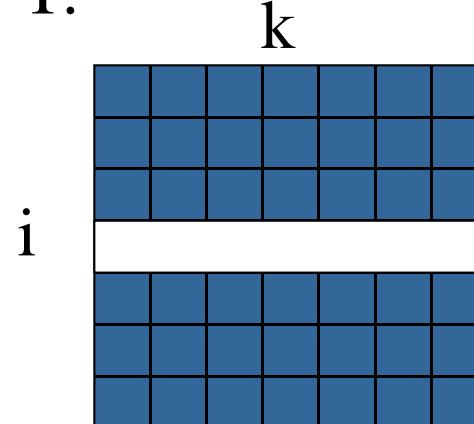
# Blocking

```
/* Unoptimized ARRAY: X = Y * Z */  
  
for (i = 0; i < N; i = i + 1)  
    for (j = 0; j < N; j = j + 1)  
        {r = 0;  
         for (k = 0; k < N; k = k + 1)  
             r = r + y[i][k] * z[k][j];  
         x[i][j] = r;  
     };
```

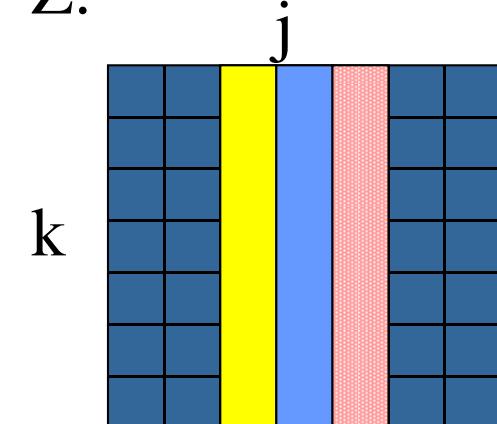
X:



Y:



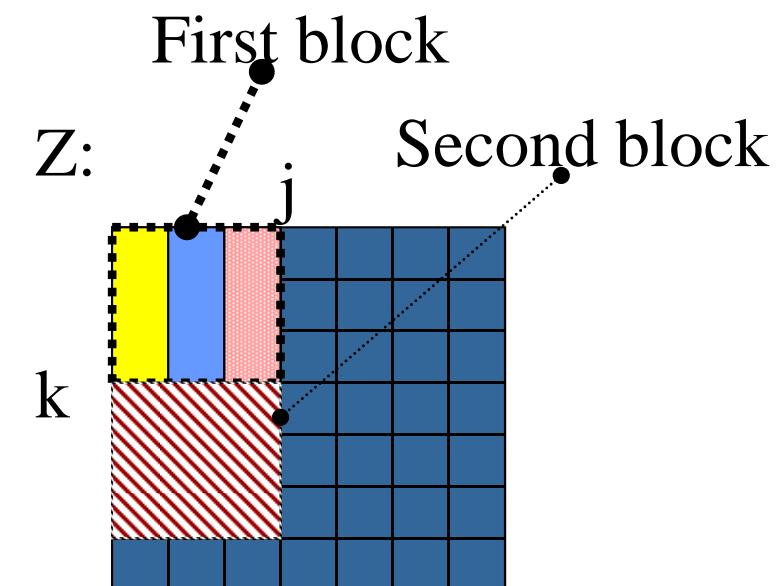
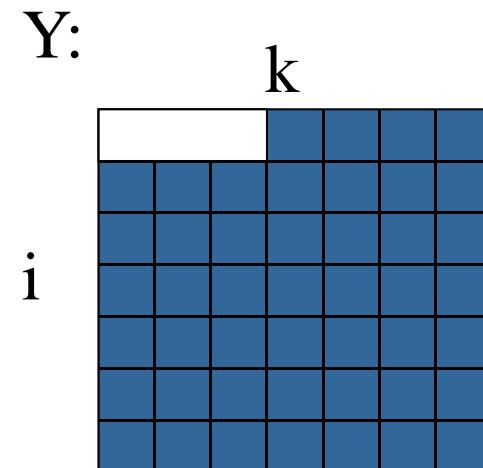
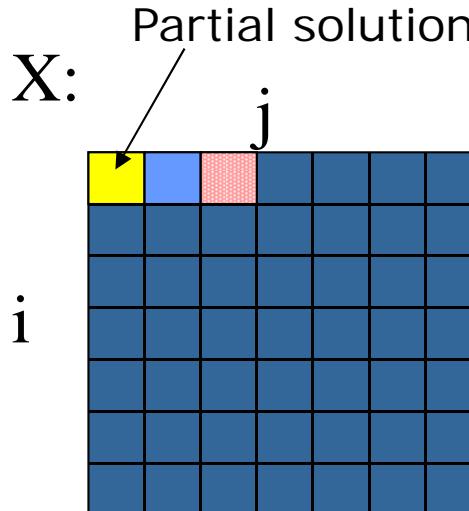
Z:





# Blocking

```
/* Optimized ARRAY: X = Y * Z */  
for (jj = 0; jj < N; jj = jj + B)  
for (kk = 0; kk < N; kk = kk + B)  
for (i = 0; i < N; i = i + 1)  
    for (j = jj; j < min(jj+B,N); j = j + 1)  
        {r = 0;  
         for (k = kk; k < min(kk+B,N); k = k + 1)  
             r = r + y[i][k] * z[k][j];  
         x[i][j] += r;  
        };
```



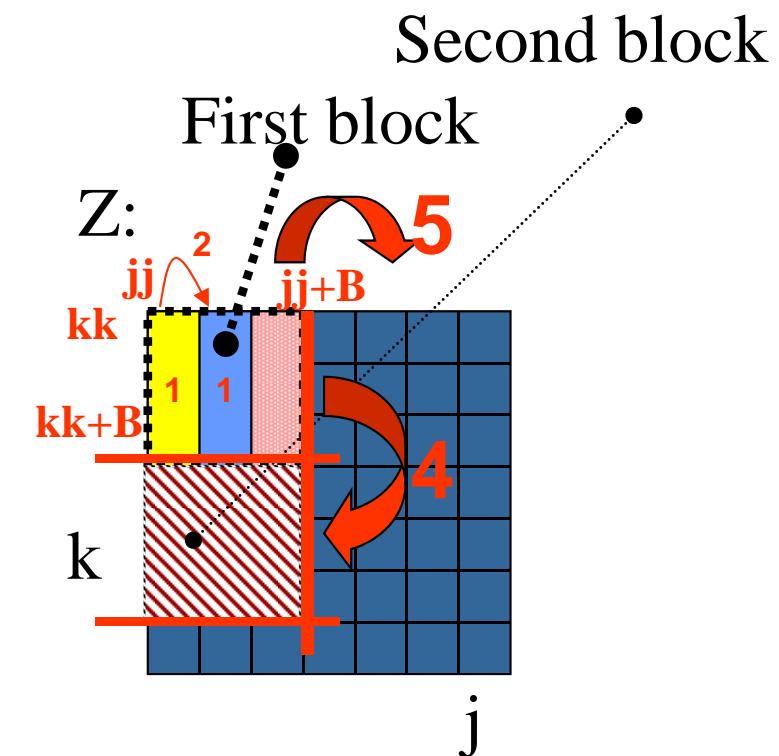
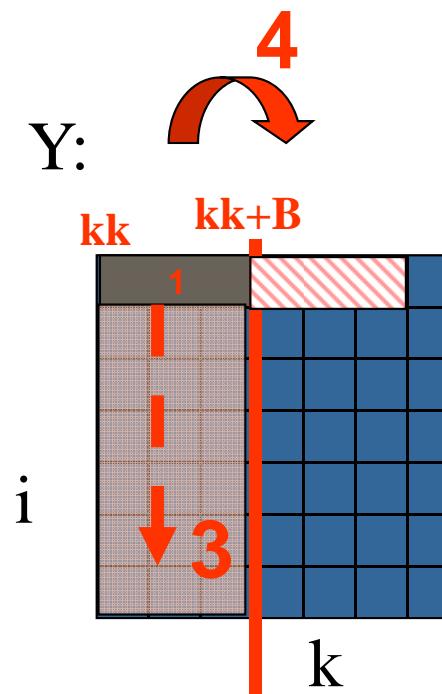
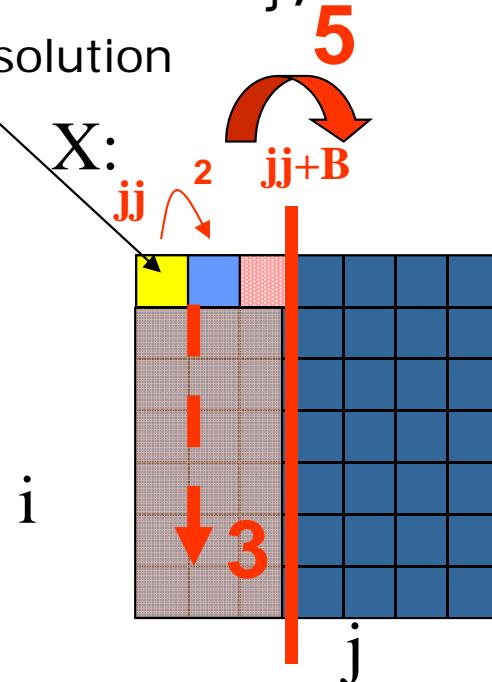


# Blocking: the Movie!

```
/* Optimized ARRAY: X = Y * Z */  
for (jj = 0; jj < N; jj = jj + B)  
for (kk = 0; kk < N; kk = kk + B)  
for (i = 0; i < N; i = i + 1)  
    for (j = jj; j < min(jj+B,N); j = j + 1)  
        {r = 0;  
         for (k = kk; k < min(kk+B,N); k = k + 1) /* Loop 1 */  
             r = r + y[i][k] * z[k][j];  
         x[i][j] += r;  
     };
```

/\* Loop 5 \*/  
/\* Loop 4 \*/  
/\* Loop 3 \*/  
/\* Loop 2 \*/  
/\* Loop 1 \*/

Partial solution



# Software Prefetch Optimizations

---

Erik Hagersten

Uppsala University, Sweden

[eh@it.uu.se](mailto:eh@it.uu.se)



# SW Prefetching

```
/* Unoptimized */  
for (j = 0; j < N; j++)  
    for (i = 0; i < N; i++)  
        x[j][i] = 2 * x[j][i];
```

```
/* Optimized */  
for (j = 0; j < N; j++)  
    for (i = 0; i < N; i++) {  
        PREFETCH x[j+1][i]  
        x[j][i] = 2 * x[j][i];  
    }
```

Typically, the HW prefetcher will successfully prefetch sequential streams

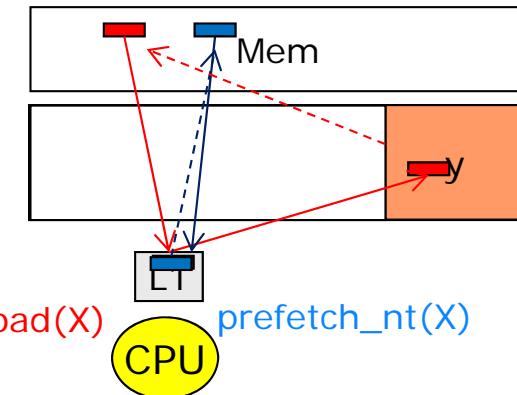
Are there examples where SW prefetching can do any better?



# Cache Waste

```
/* Unoptimized */
for (s = 0; s < ITERATIONS; s++) {
    for (j = 0; j < HUGE; j++)
        x[j] = x[j+1]; /* will hog the cache but not benefit*/
    for (i = 0; i < SMALLER_THAN_L2_CACHE; i++)
        y[i] = y[i+1]; /* will be evicted between usages */
}

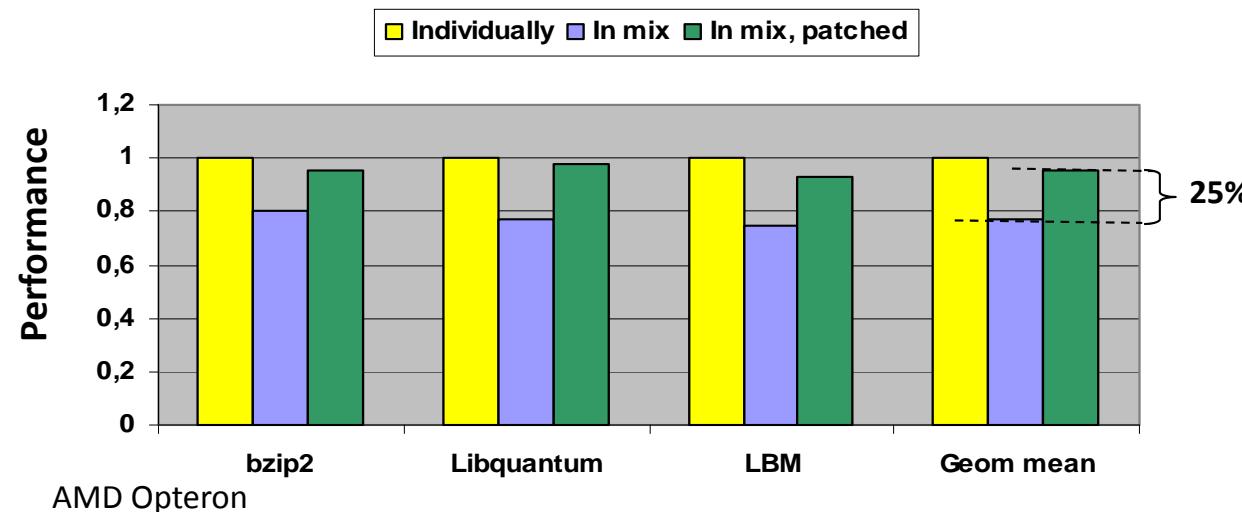
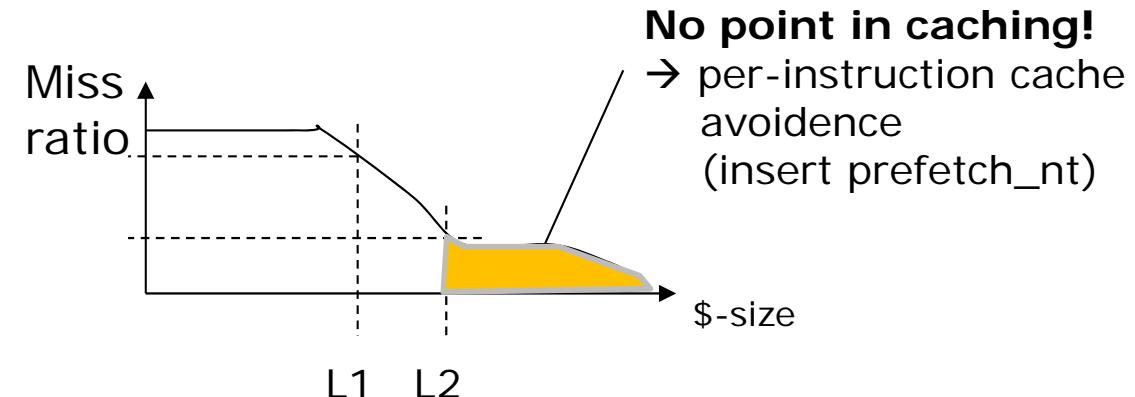
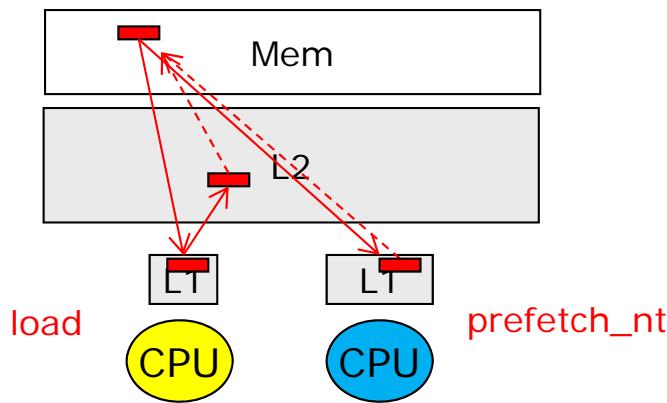
/* Optimized */
for (s = 0; s < ITERATIONS; s++) {
    for (j = 0; j < HUGE; j++) {
        PREFETCH_NT x[j+1] /* will be installed in L1, bypass L2 */
        x[j] = x[j+1];
    for (i = 0; I < SMALLER_THAN_L2_CACHE; i++)
        y[i] = y[i+1]; /* will always hit in the cache*/
    }
}
```



- Can also be beneficial if applications are co-scheduled on MC and share cache with other applications.



# UART: Avoiding cache waste (A. Sandberg)

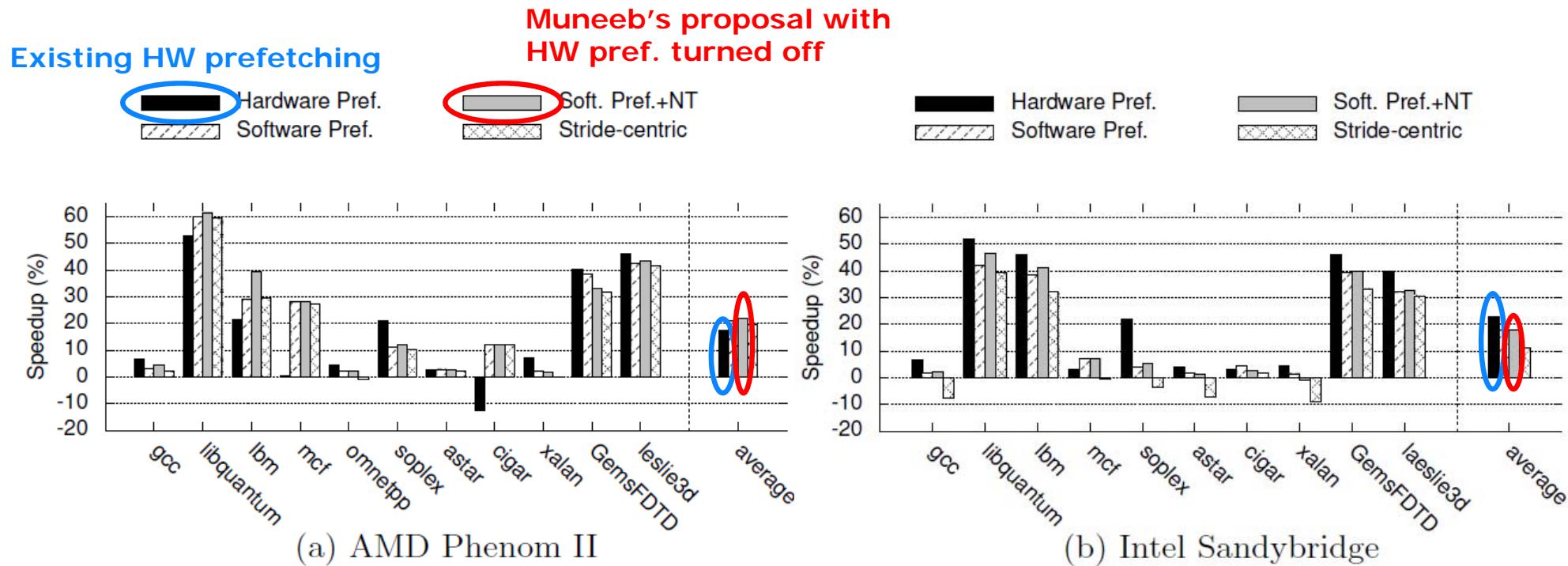




## UART: Resource-lean SW prefetching (Muneeb)

- Use the “UART modeling” (StatStack) to “understand” the application (run once with a 20% overhead)
- Find instructions that miss in the cache
- Capture their popular strides
- Figure out the necessary prefetch distances
- Figure out the best non-temporal strategy
- Do cost/benefit analysis
- Insert prefetches into binary automatically (actually: into ASM)
  - ✿ Use the right stride
  - ✿ Use the right prefetch distance

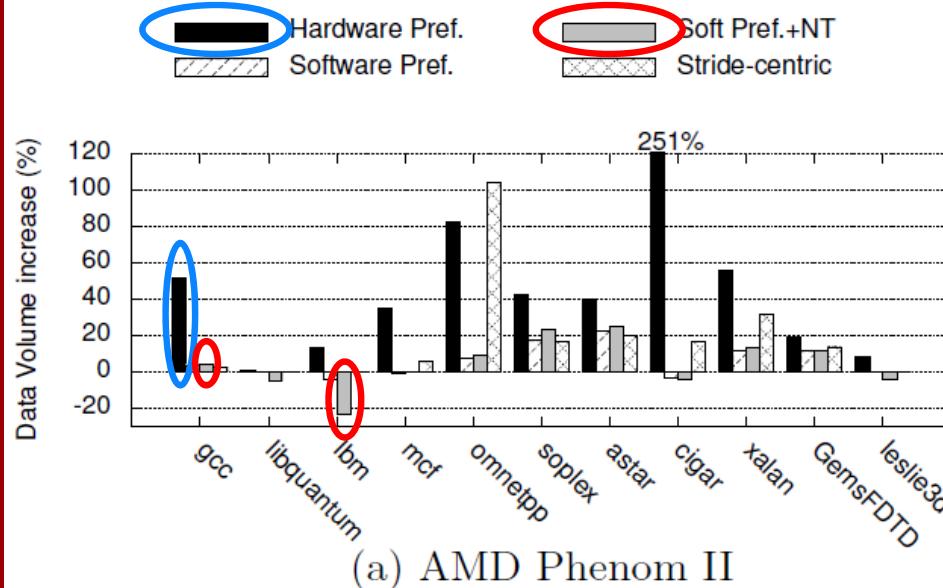
# Speedup comparison



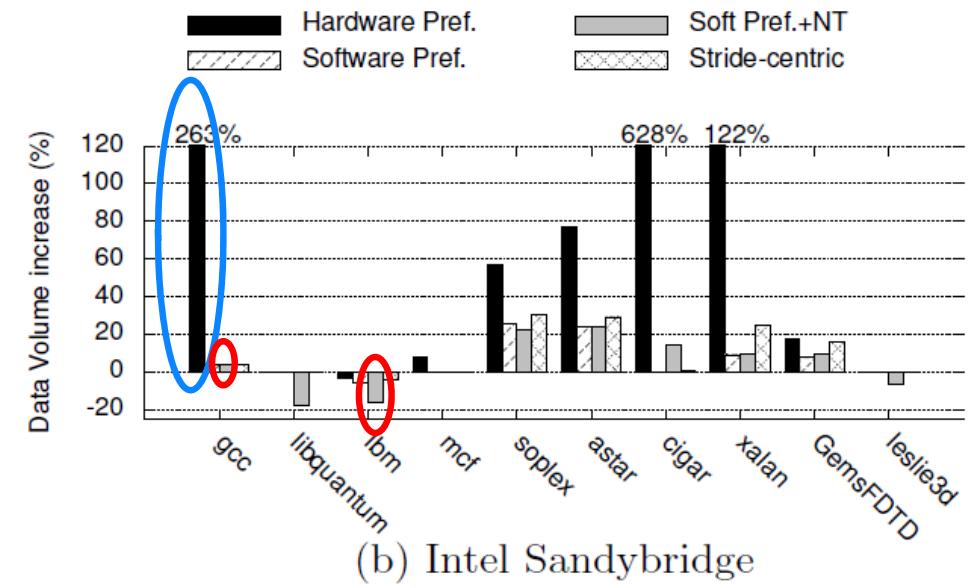


# Data Volume Increase Caused by Prefetching

Existing HW prefetching



Muneeb's proposal with  
HW pref. turned off

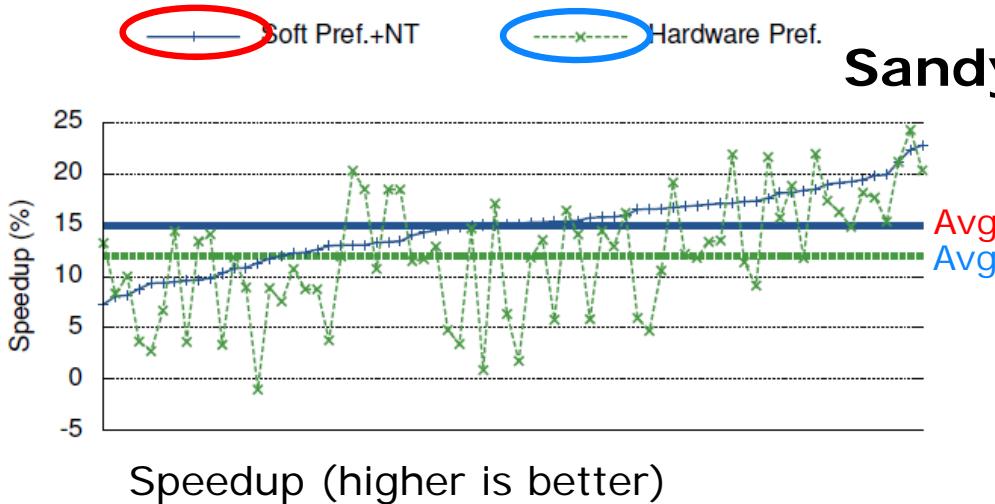




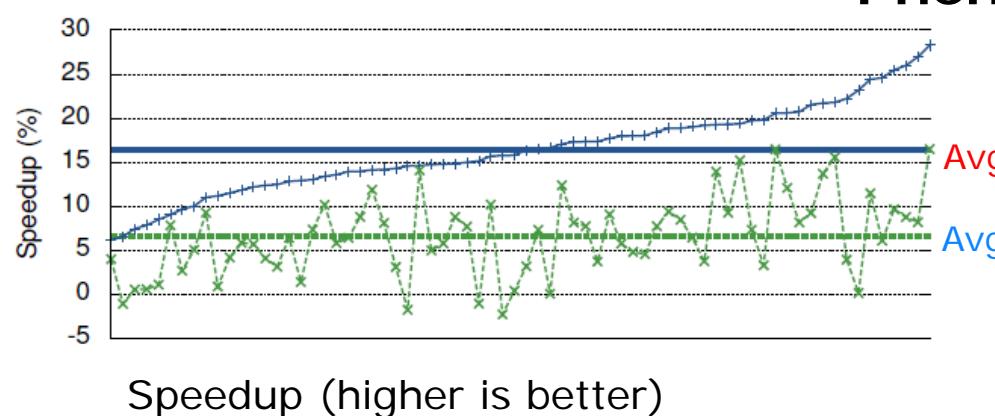
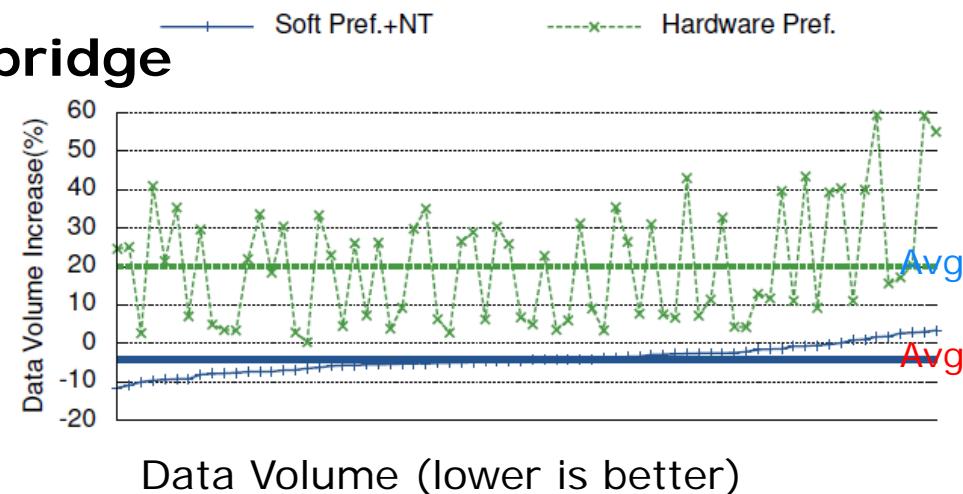
# Running many apps at the same time on a MC

Muneeb's proposal with  
HW pref. turned off

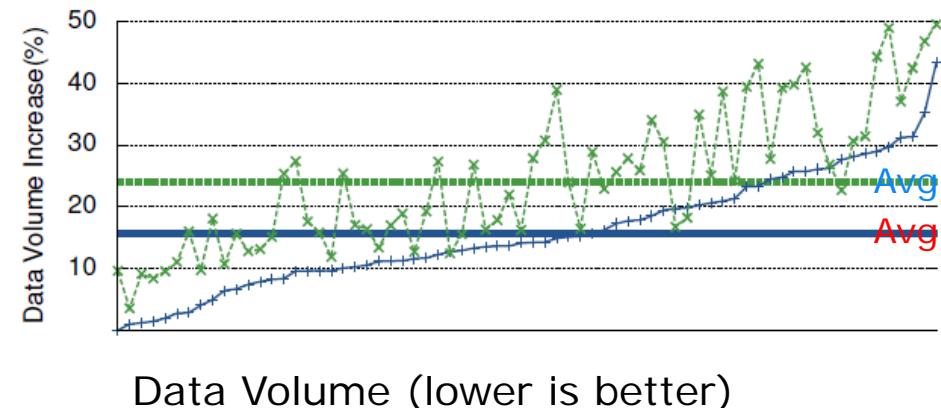
Existing HW prefetching



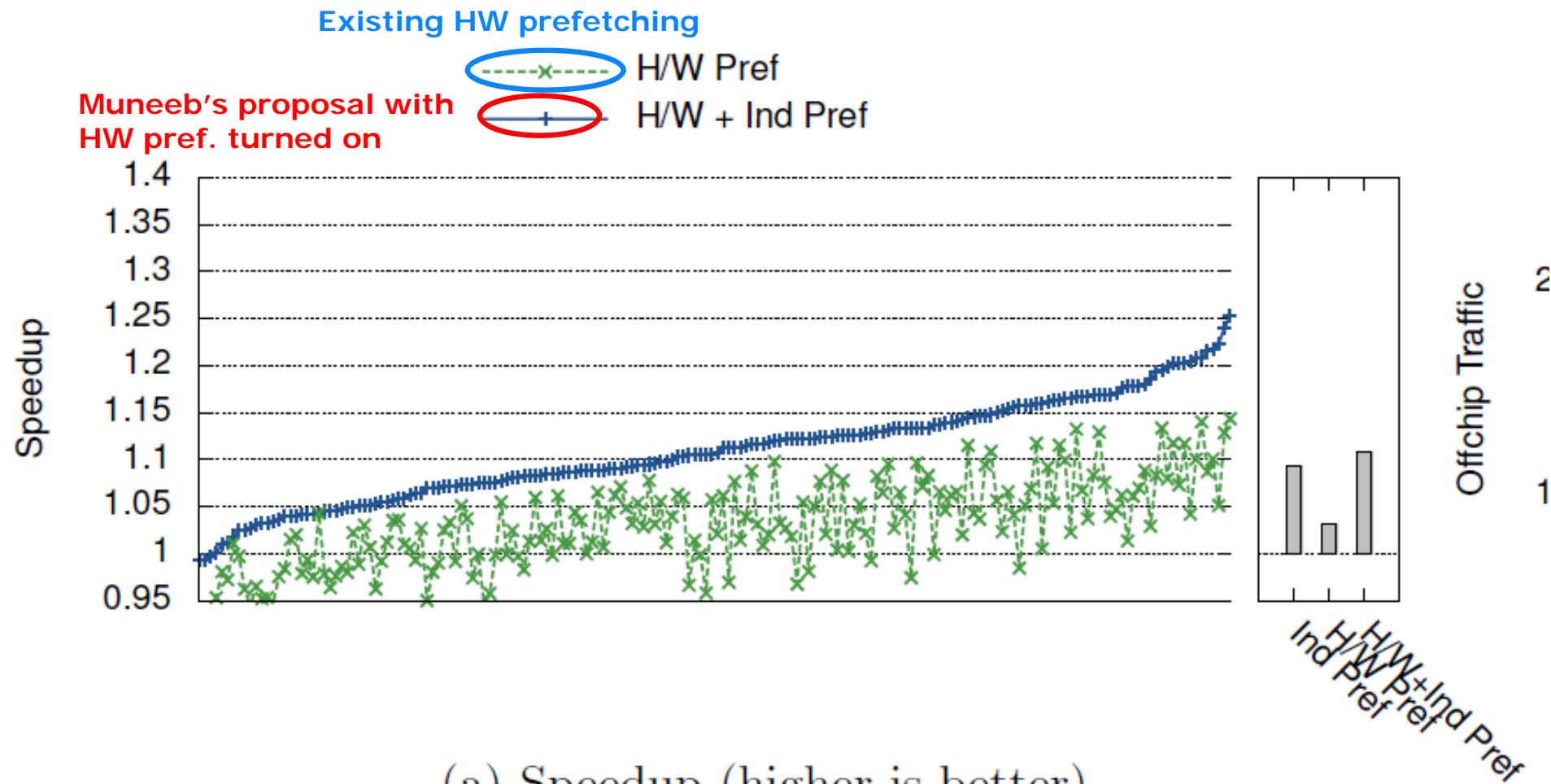
## Sandybridge



## Phenom II



# Now also non-strided prefetching



# Example of Performance Tools: Vtune/ThreadSpotter

---

Erik Hagersten

Uppsala University, Sweden

[eh@it.uu.se](mailto:eh@it.uu.se)

# Some performance tools

## Free tools

- Oprofile
- GNU: gprof
- AMD: code analyst
- Google performance tools
- Virtual Inst: High Productivity Supercomputing  
(<http://www.vi-hps.org/tools/>)
- ...

## Commercial tools

- Intel: Vtune and many more
- HP: Multicore toolkit (some free, some not)
- ThreadSpotter (of course😊)
- ...

Covered here

Covered here

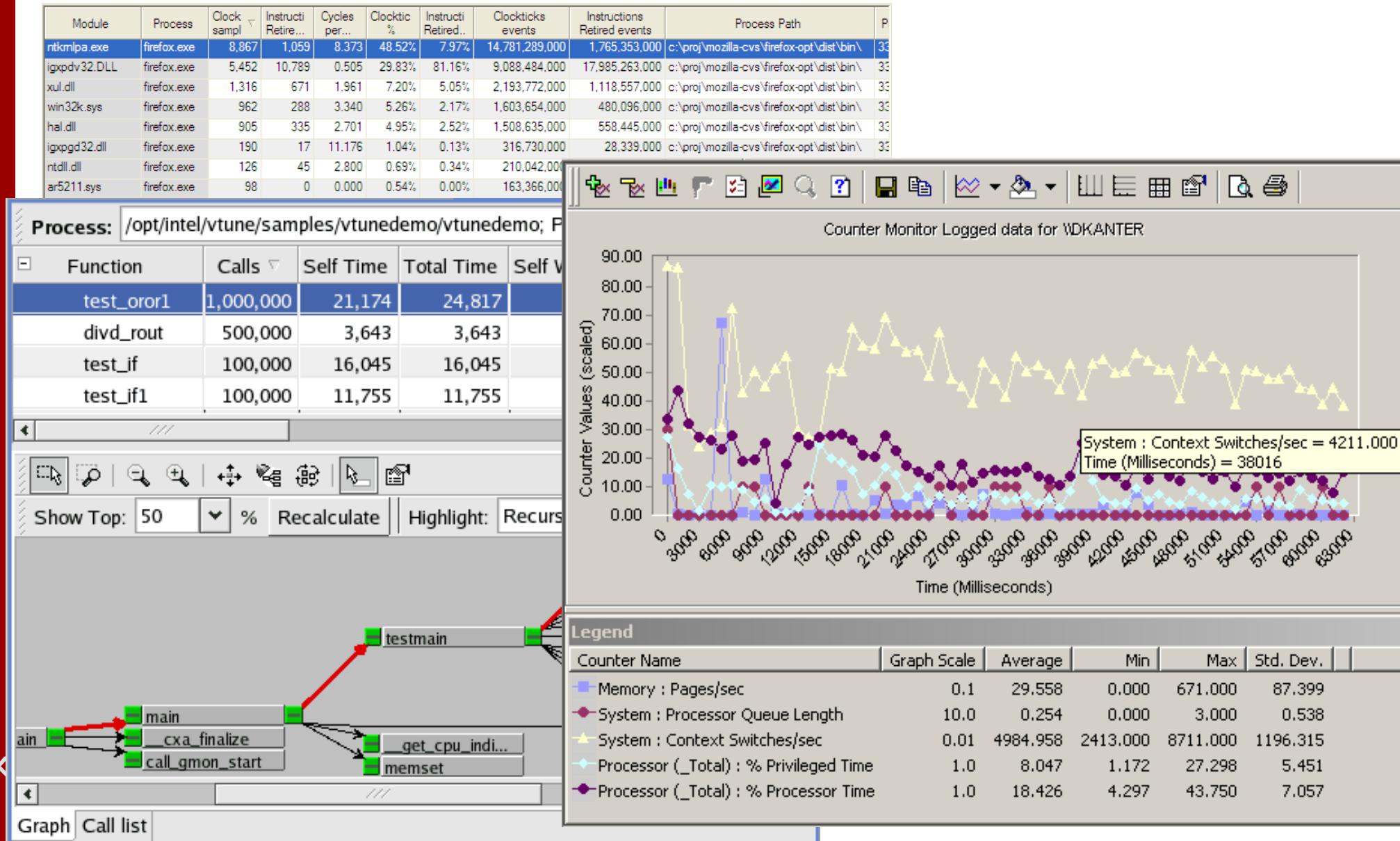


# VTune

- Broad-spectrum, low-level tool
  - ✿ Profiling (where do you spend your time)
  - ✿ HW ctr-based information
    - CPI...
    - I1, D1, L2, L3 cache misses...
    - Branch prediction
    - Per: application, function, instruction
- Often requires architecture knowledge to understand the stats (so you should be fine😊)



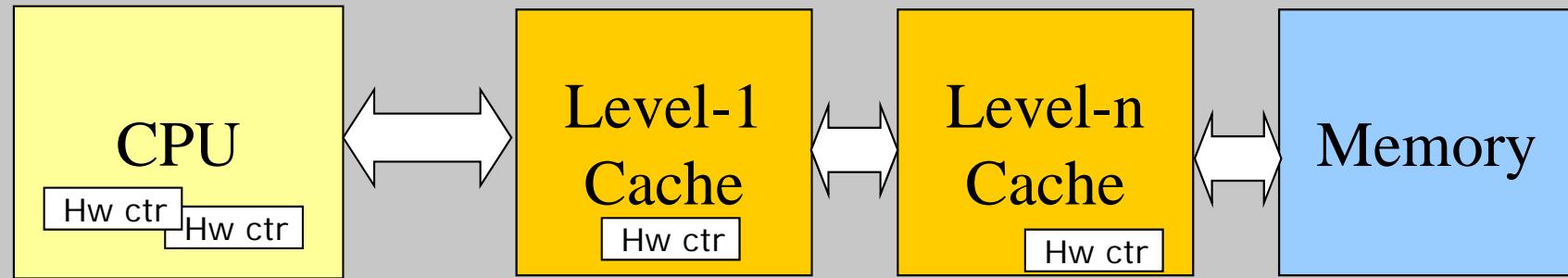
# Example of Vtune info





# Using Programmable HW Counters

Real HW



Runs at near native speed!

# ThreadSpotter: "Cache Tool for Dummies"



UPPSALA  
UNIVERSITET

Source:  
C, C++, Fortran

What?

```
/* Unoptimized Array Multiplication: x = y * z N = 1024 */
for (i = 0; i < N; i += 1)
    for (j = 0; j < N; j += 1)
        {r = 0;
        for (k = 0; k < N; k = k + 1)
            r = r + y[i][k] * z[k][j];
        x[i][j] = r;
        }
/* Unoptimized Array Multiplication: x = y * z N = 1024 */
for (i = 0; i < N; i += 1)
    for (j = 0; j < N; j += 1)
        {r = 0;
        for (k = 0; k < N; k = k + 1)
            r = r + y[i][k] * z[k][j];
        x[i][j] = r;
        }
```

How?

Any Compiler



Host System

The screenshot shows the ThreadSpotter application interface. At the top, there's a navigation bar with links like 'Release Notes', 'Fedora Project', 'Acumem Report', etc. Below the navigation is a table titled 'Bandwidth Issues' with columns: Loop / Issue, Summary, % of fetches, Utilization, HW-Prefetch, and Rand. Several rows are listed, such as 'Inefficient loop nesting' and 'Loop fusion'. To the right of the table is a code editor window showing C code with annotations. A red arrow points from the 'What?' text to this code editor. Below the table is a detailed analysis pane for 'Issue #2: Cache line utilization', which includes sections for 'Statistics for instructions of this issue', 'Instructions involved in this issue', and 'Instructions previously writing to related data'. A red arrow points from the 'Where?' text to this analysis pane. At the bottom of the interface, there's a status bar with browser-like tabs and a message 'file:/// - 8.1. Poor Cache Line Utilization ...'.

Where?

"Application locality"

Finger  
Print  
(~4MB)

UART technology

Analysis

Advice

Target System  
Parameters





Loop / Issue	Summary	% of fetches	Utilization	HW-Prefetch	Randomness
1 / 3	Poor utilization	29.4%	12.4%	100.0%	Low
1 / 4	Loop fusion	29.4%	12.4%	97.6%	Low
1 / 1	Inefficient loop nesting	29.2%	12.6%	0.0%	Low
3 / 9	Loop fusion	4.4%	11.8%	97.3%	Low
3 / 8	Poor utilization	4.4%	23.7%	100.0%	Low
4 / 13	Loop fusion	4.2%	12.7%	96.7%	Low
4 / 12	Loop fusion	4.2%	12.7%	96.7%	Low
7 / 18	Poor utilization	4.2%	25.1%	100.0%	Low
4 / 10	Poor utilization	4.2%			Low

List of Bandwidth SlowSpots

## Issue #1: Inefficient loop nesting

This instruction group also show symptoms of: Cache line utilization, Hot-spot.

### + Statistics for instructions of this issue

### + Instructions involved in this issue

### + Loop statistics

### + Loop instructions

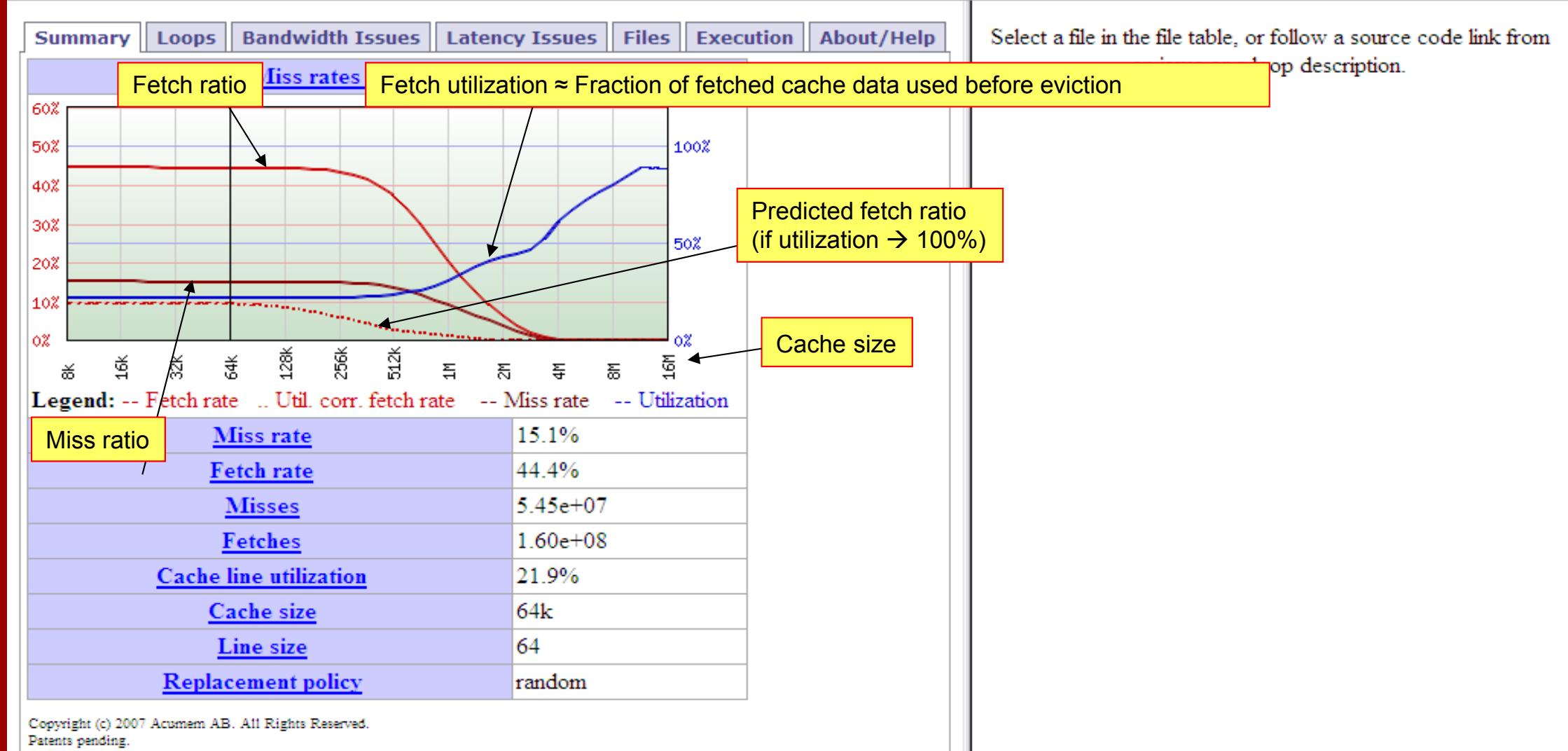
Copyright (c) 2007 Acumem AB. All Rights Reserved.  
Patent pending.

Explaining what to do

```
600 tnorm += f1_layer[ti].P * f1_layer[ti].P;
601
602 if (ttemp != f1_layer[ti].P)
603     tresult=0;
604 }
605 f1res = tresult;
606
607 /* Compute F1 - Q values */
608
609 tnorm = sqrt((double)tnorm);
610 for (tj=0;tj<numf1s;tj++)
611     f1_layer[tj].Q = f1_layer[tj].P;
612
613 /* Compute F2 - y values */
614 for (tj=0;tj<numf2s;tj++)
615 {
616     Y[tj].y = 0;
617     if ( !Y[tj].reset )
618         for (ti=0;ti<numf1s;ti++)
619             Y[tj].y += f1_layer[ti].P * bus[ti][tj];
620
621
622     /* Find match */
623     winner = 0;
624     for (ti=0;ti<numf2s;ti++)
625     {
626         if (Y[ti].y > Y[winner].y)
627             winner = ti;
628     }
629
630
631     #ifdef DEBUG
632         if (DB1) print_f12();
633         if (DB1) printf("\n num iterations for p to stabilize = %i \n",
634 #endif
635         match_confidence=simtest2();
636         if ((match_confidence) > rho)
637         {
638             /* If the winner is not the default F2 neuron (the
639             * winner is always the first one in the array)
```

Spotting the crime

Klar





# Resource Sharing Example

## Libquantum

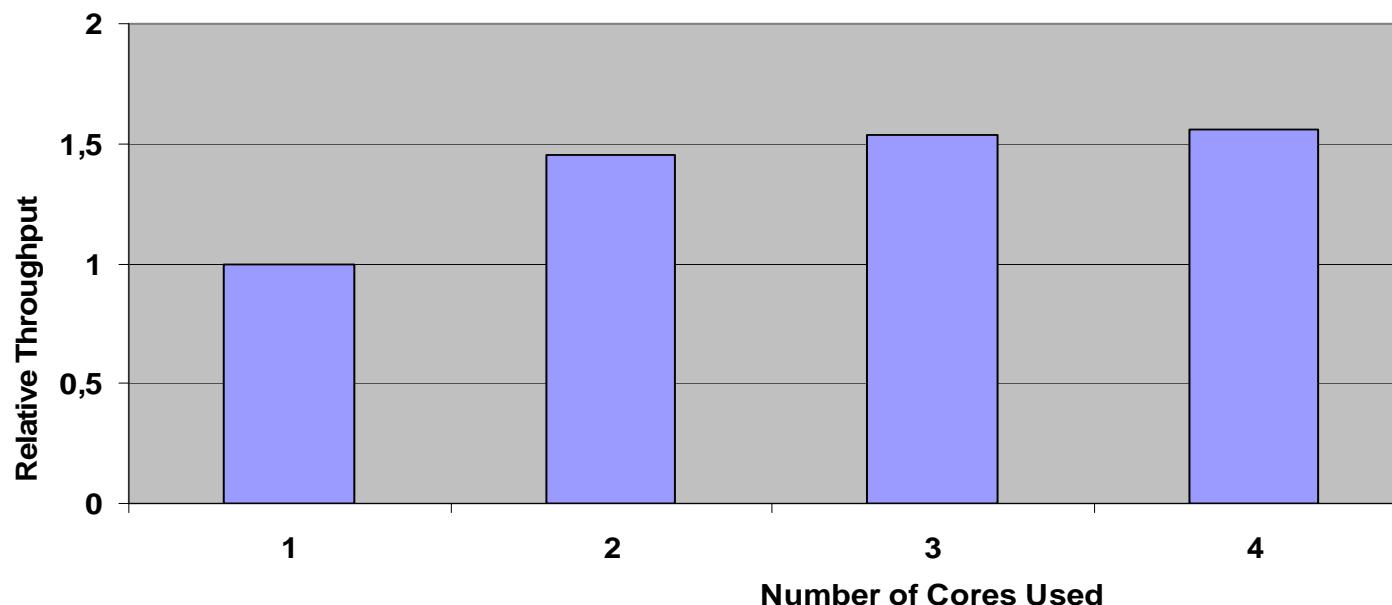
A quantum computer simulation

Widely used in research (download from: <http://www.libquantum.de/>)

4000+ lines of C, fairly complex code.

Runs an experiment in ~30 min

Throughput improvement:



**Demo Time!**

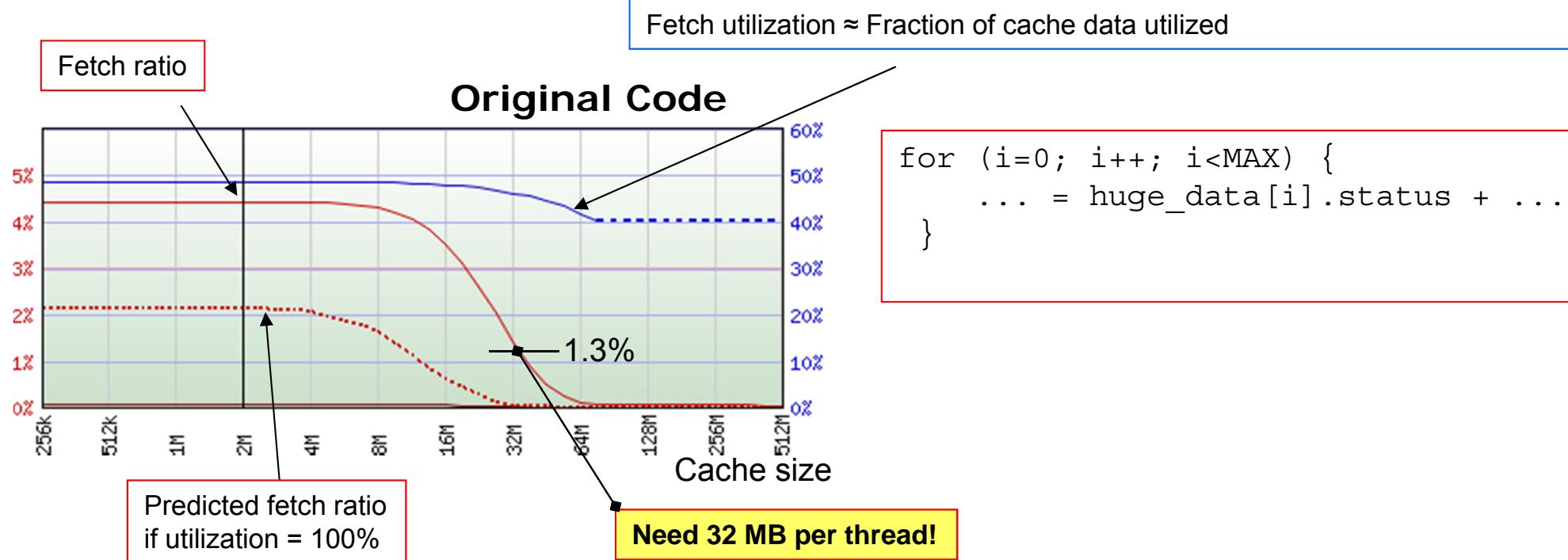
**Libquantum**

[Libquantum](http://user.it.uu.se/~eh)

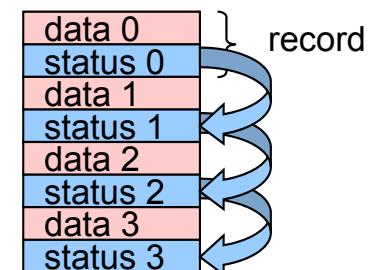


# Fetch Utilization Analysis

## Libquantum



Main data structure  
(vector of structs)

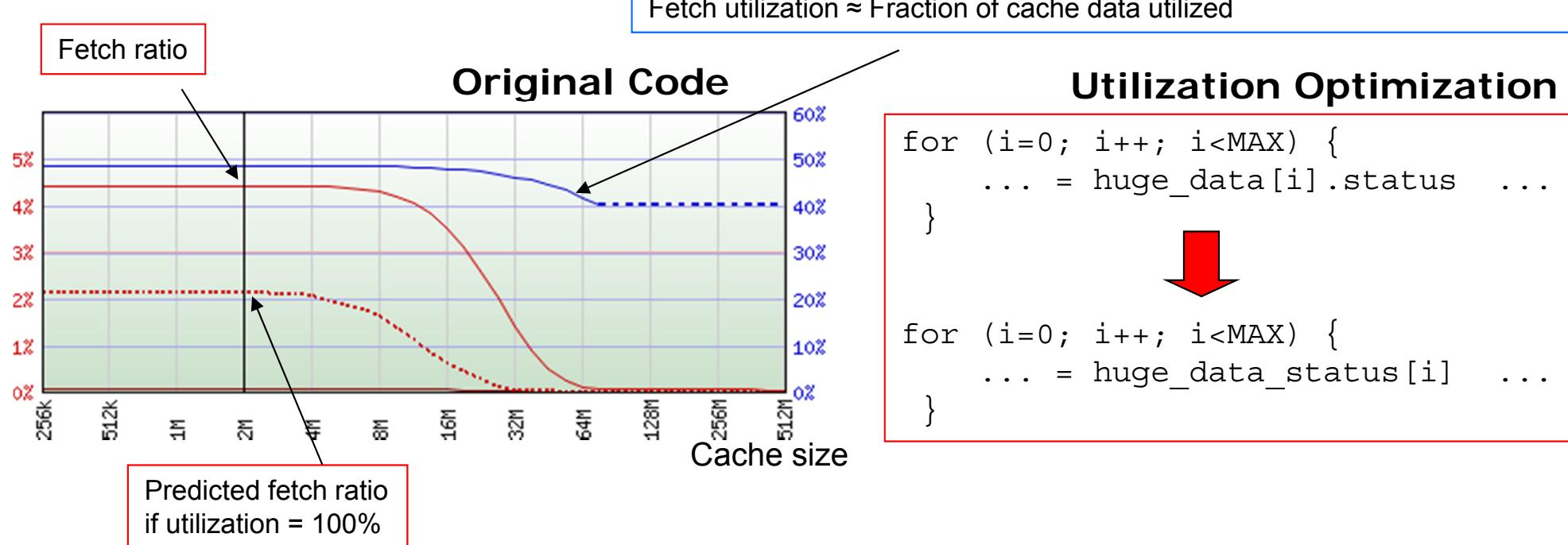


Only accessing status  
data in main loop



# Utilization Analysis

## Libquantum



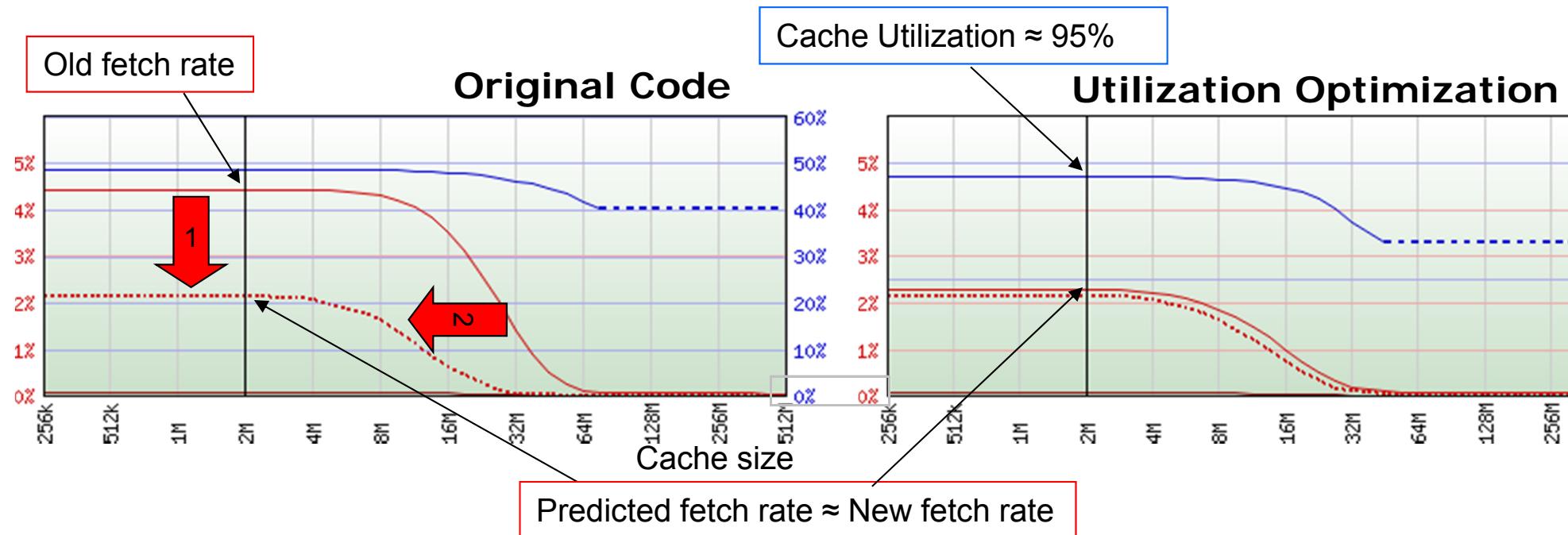
## ThreadSpotter's First Advice: Improve Utilization

→ Change to “struct of vectors”

- Involves ~20 lines of code
- Takes a non-expert 30 min



# Utilization Optimization

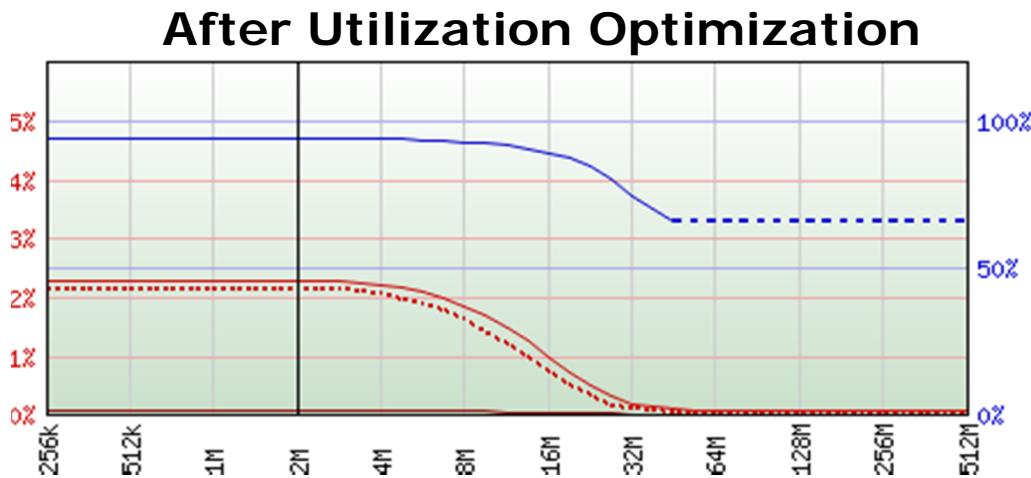


## Two positive effects from better utilization

1. Each fetch brings in more useful data → lower fetch ratio
2. The same amount of useful data can fit in a smaller cache → shift left



# Reuse Analysis



### Utilization + Fusion Optimization

```
...  
toffoli(huge_data_status, ...)  
toffoli(huge_data_status, ...)  
...  
...  
fused_toffoli(huge_data_status, ...)  
...
```

## Second-Fifth ThreadSpotter Advice: Improve reuse of data

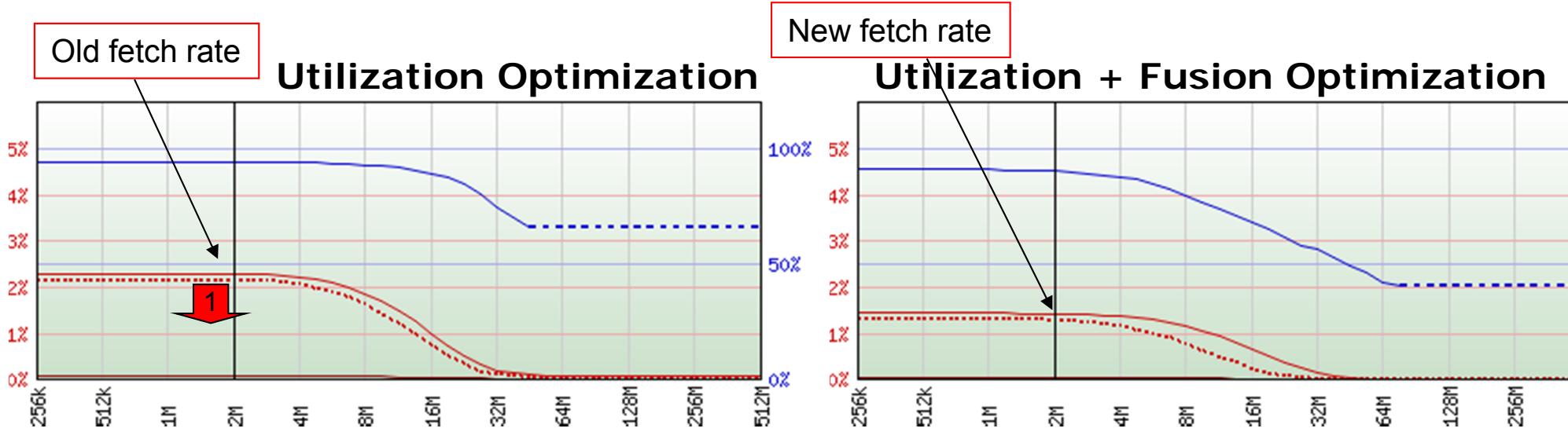
→ Fuse functions traversing the same data

- Here: four new fused functions created
- Takes a non-expert < 2h



# Effect: Loop Fusion

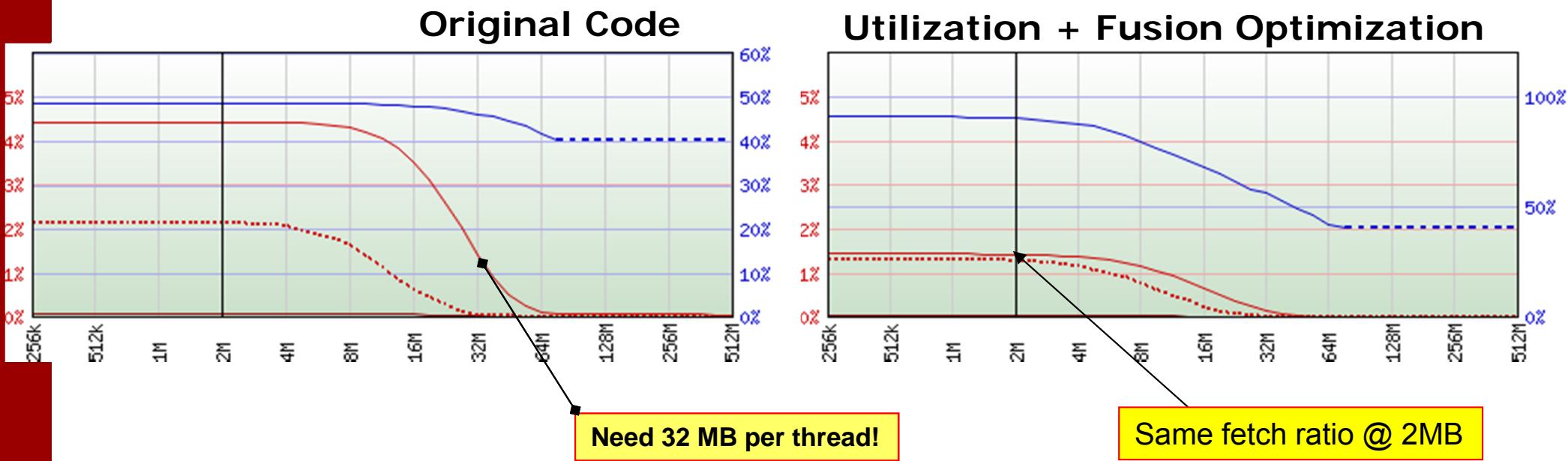
SPEC CPU2006-462.libquantum



- The miss in the second loop goes away
- Still need the same amount of cache to fit “all data”



# Utilization + Reuse Optimization

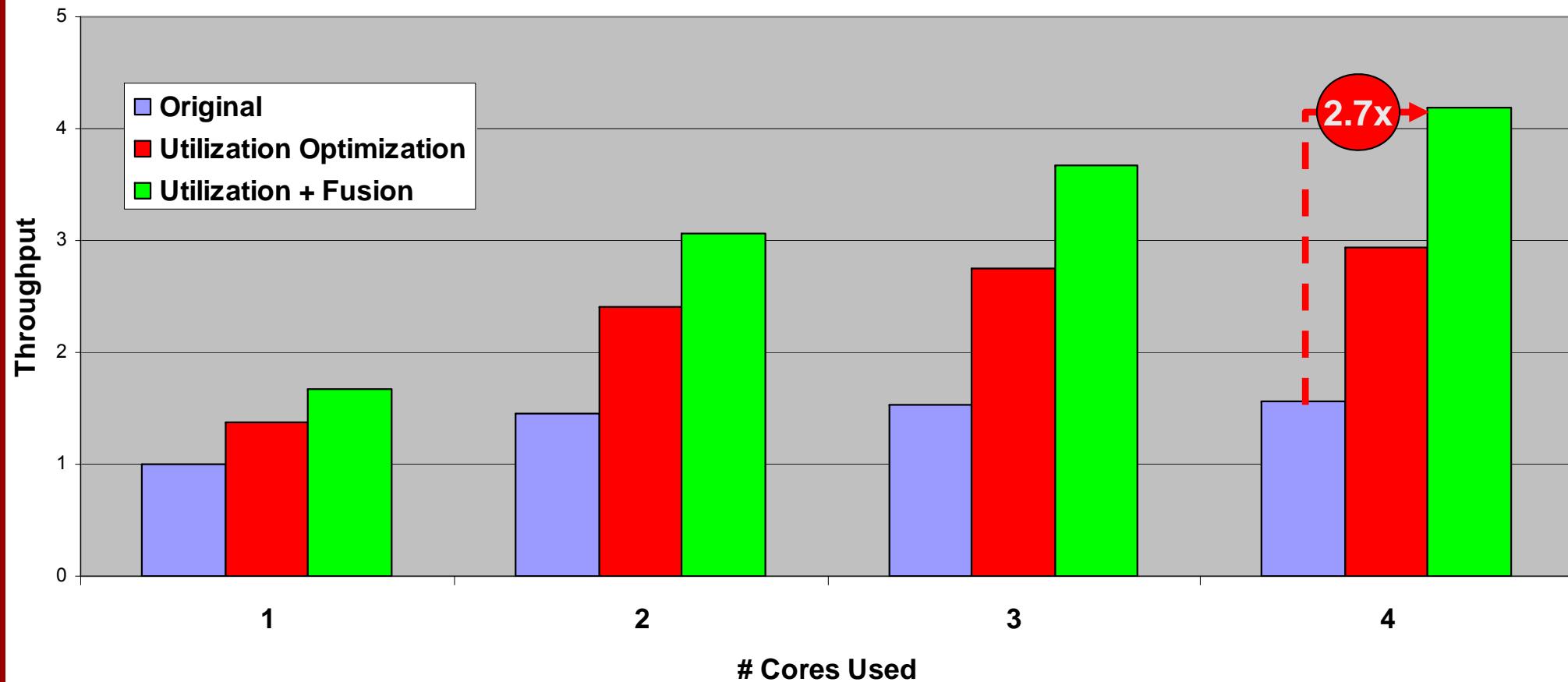


- Fetch rate down to 1.3% for 2MB caches
- Same as a 32 MB cache originally



# Summary

## Libquantum





# Can Compilers Find These Optimizations?

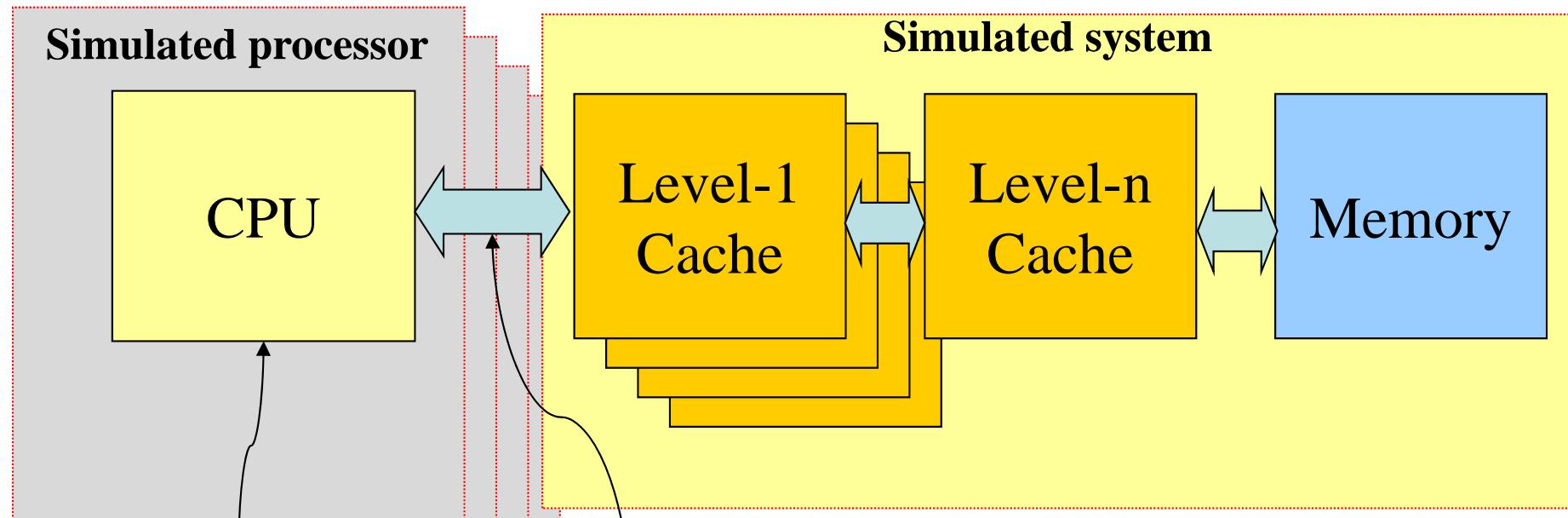
- That's what they are trying to do...
- They do a good job at small simple examples
- Sometimes fail for large complex applications

# **Fast Cache Modeling Technology**

---

Erik Hagersten  
Uppsala University

# Understanding Performance and Power: Traditional simulation model



Code:

```
set A,%r1  
ld [%r1],%r0  
st %r0,[%r1+8]  
add %r1,1,%r1  
ld [%r1+16],%r0  
add %r0,%r5,%r5
```

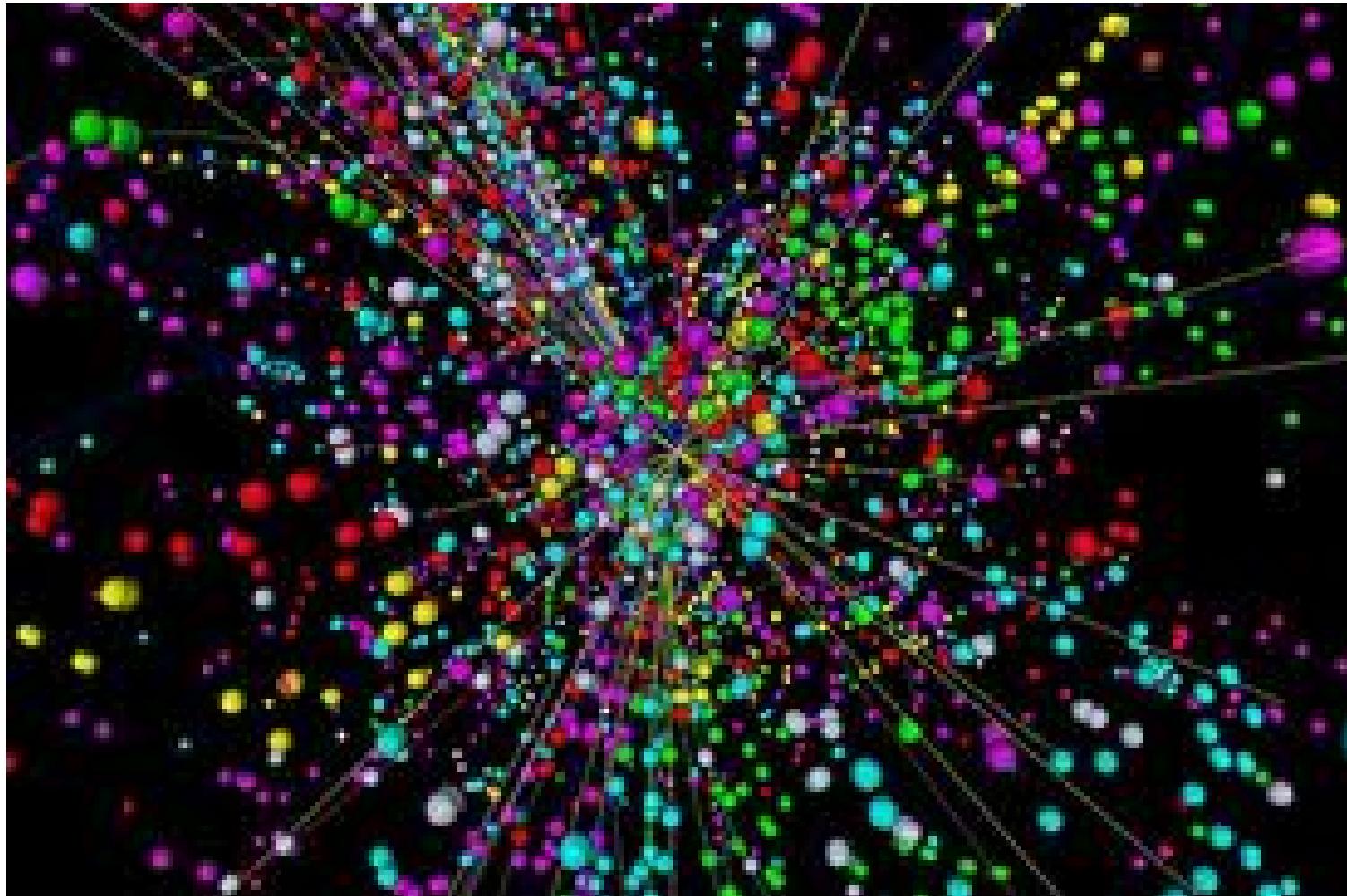
Memory ref:  
1:read A  
2:write B  
3:read C  
4:write B  
[...]

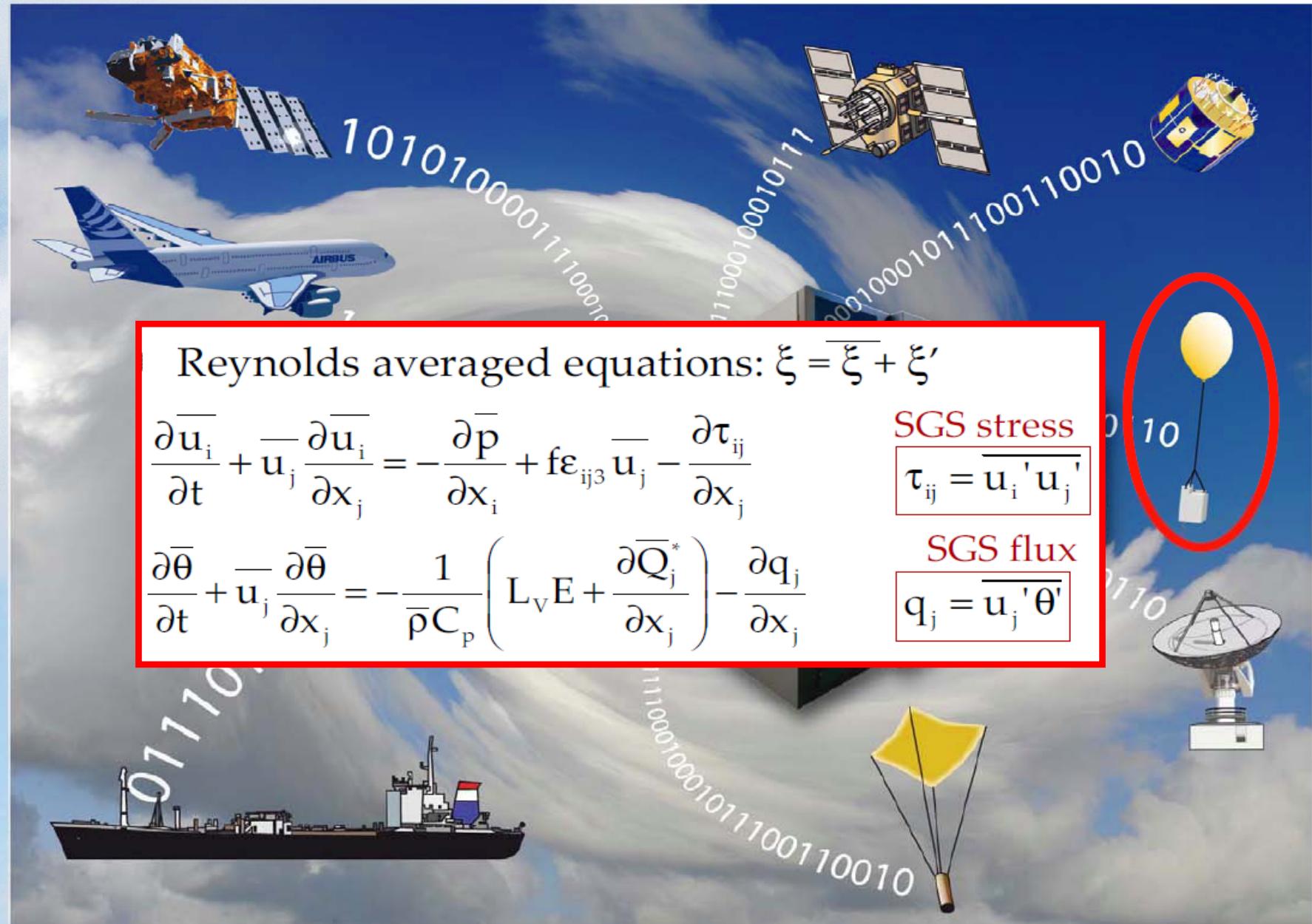
Overhead ~100 - 100,000x  
(UART is also trying to  
make traditional simulation  
faster...)

# Our Goal

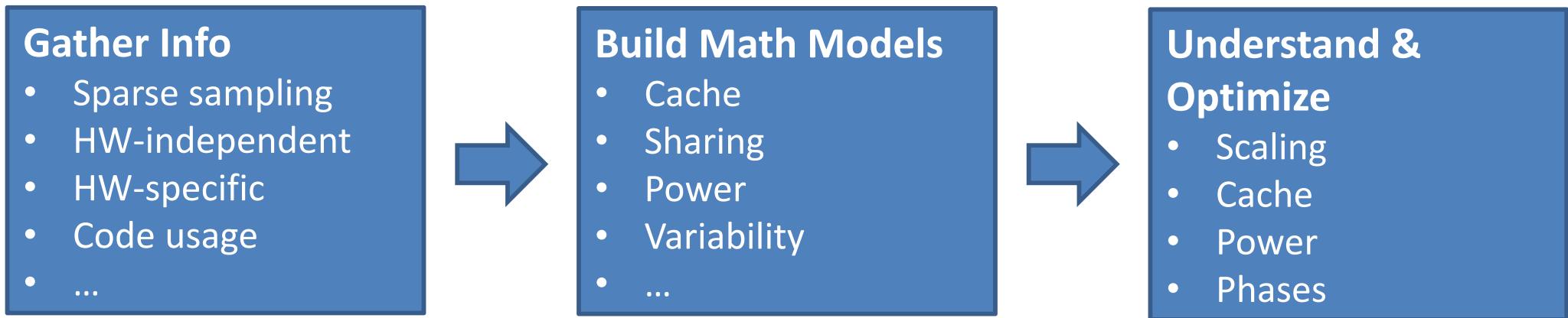
- Estimate performance and power effects for
  - Programmers
  - Compilers / compiler writers
  - JIT
  - Runtime systems
  - (... and Architects)
- Runtime Overhead <2x (i.e., <100% over native)
- Architecturally-independent in-data
- Flexible models (answer what-if questions)

# Simulating Particles in Physics





# Architecture Research at Uppsala: Overall approach for modeling



## Gather Info

- Sparse sampling
- HW-independent

## Build Models

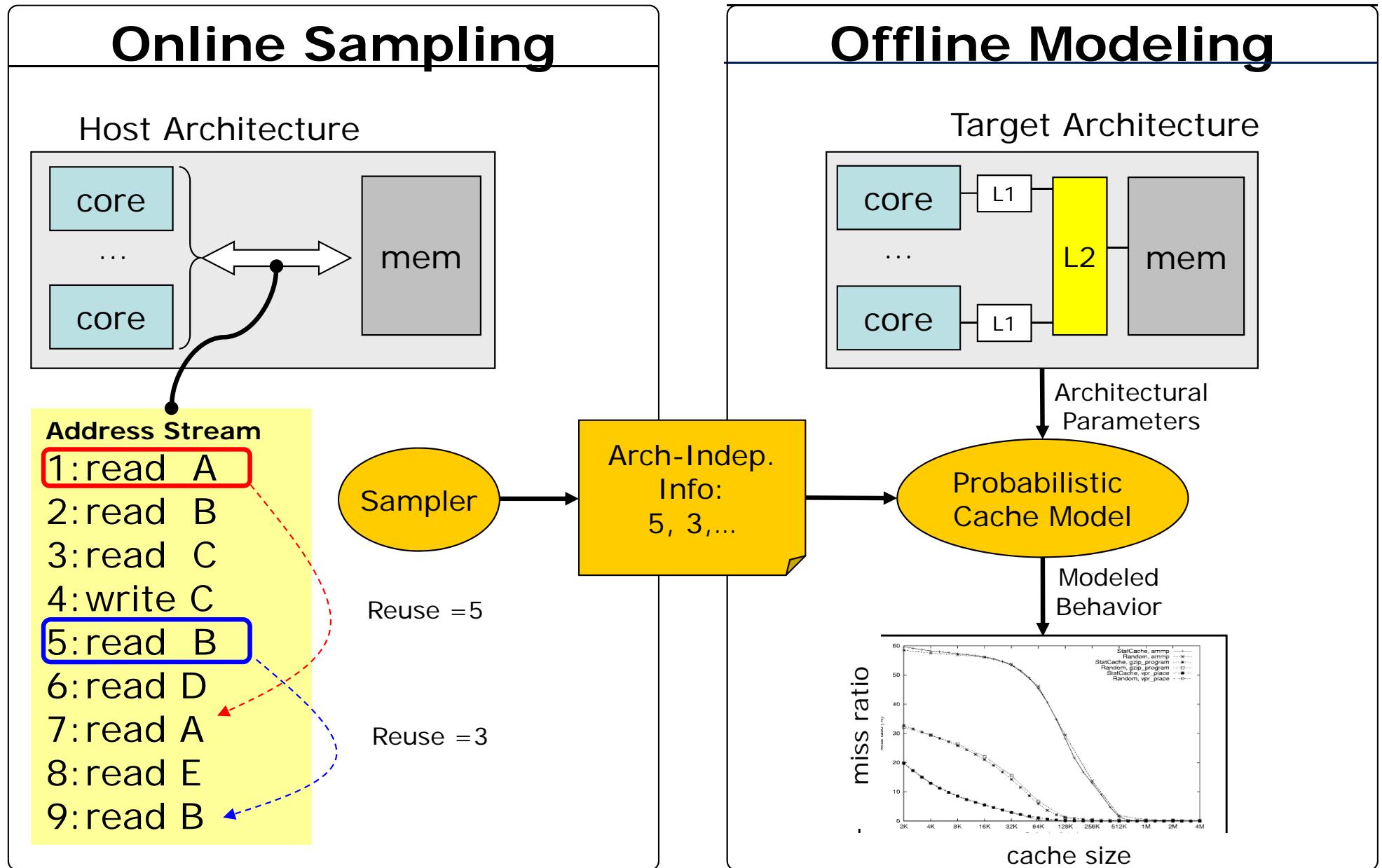
- Cache
- Sharing

## Understand & Optimize

- Cache

**HERE: “SIMULATORS” ON STEROIDS  
FOR UNDERSTANDING APPLICATION  
BEHAVIOUR**

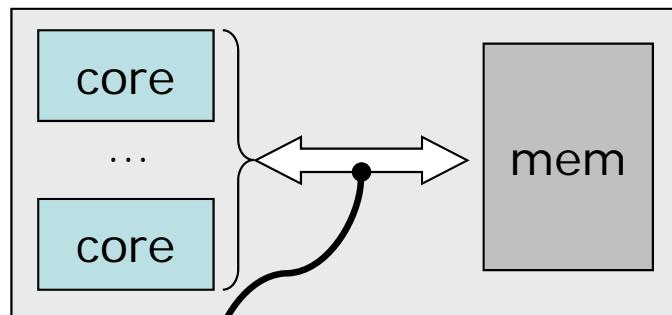
# 1. Example: Fast Cache Modeling



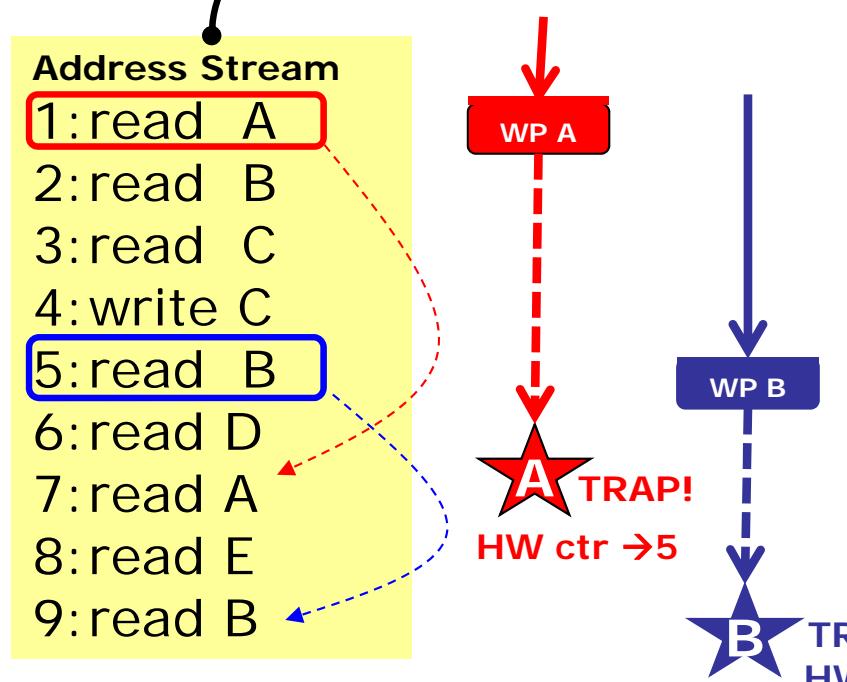
# Efficient Reuse Sampling

## Online Sampling

Host Architecture



Typically, 50.000 samples per application  
 Typical sample rate 1:100.000  
 Typical overhead: 20% over native execution



**Init HW ctr to overflow on next instruction to sample**

**Set watchpoint for the data address it accesses**

**Resume execution**

**The execution traps on next access to that address  
 HW ctrs has been counting intervening mem ops**

# Modeling caches (random replacement)

$P_{hit}$  after one cache replacement:  $(1 - 1/\text{Size})$

$P_{hit}(R)$ , after  $R$  cache replacements:  $(1 - 1/\text{Size})^R$

$P_{miss}(R) = 1 - P_{hit}(R)$ :  $(1 - (1 - 1/\text{Size})^R)$

Assume that we know the average miss ratio  $M$  for the application

$P_{miss}(A)$ , after  $A$  accesses, assuming miss ratio  $M$ :  $(1 - (1 - 1/\text{Size})^{M \cdot A})$

For  $N$  samples, each with a reuse  $r(i)$ , total misses are:

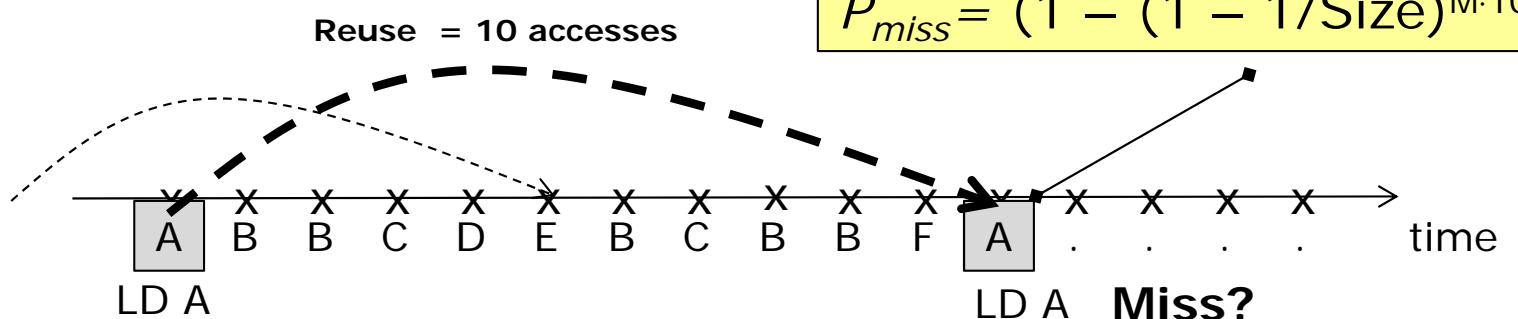
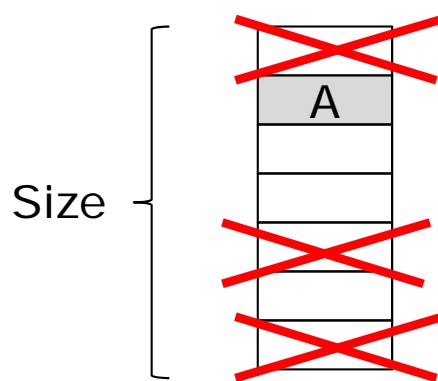
$$\sum_{k=1}^N (1 - (1 - 1/\text{Size})^{M r(k)})$$

...which should be equal to  $N$  time miss ratio ( $M$ ):

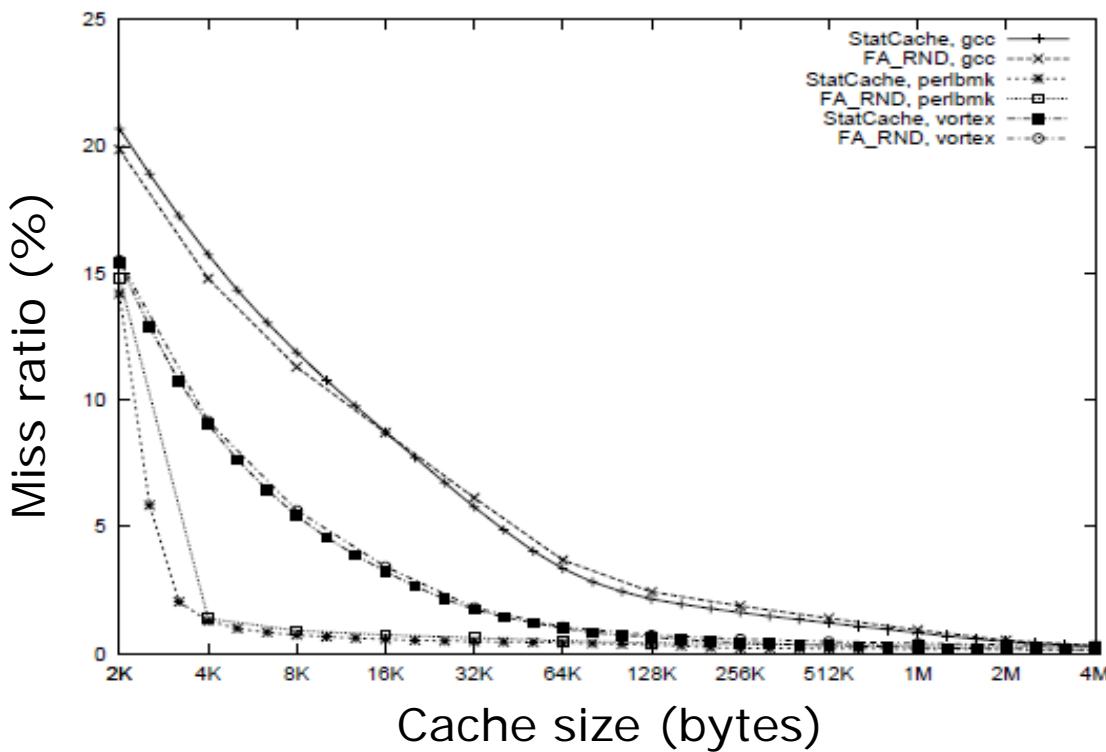
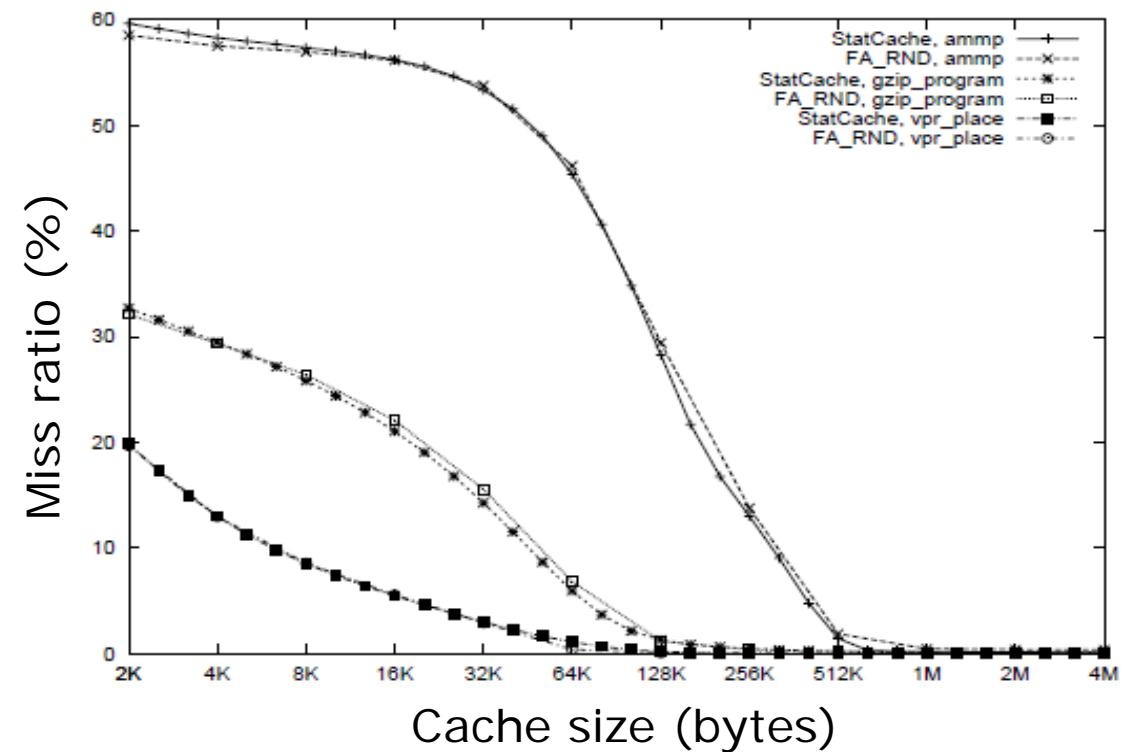
$$N M = \sum_{k=1}^N (1 - (1 - 1/\text{Size})^{M r(k)})$$

...where  $M$  can be solved numerically for any value of Size

CACHE:



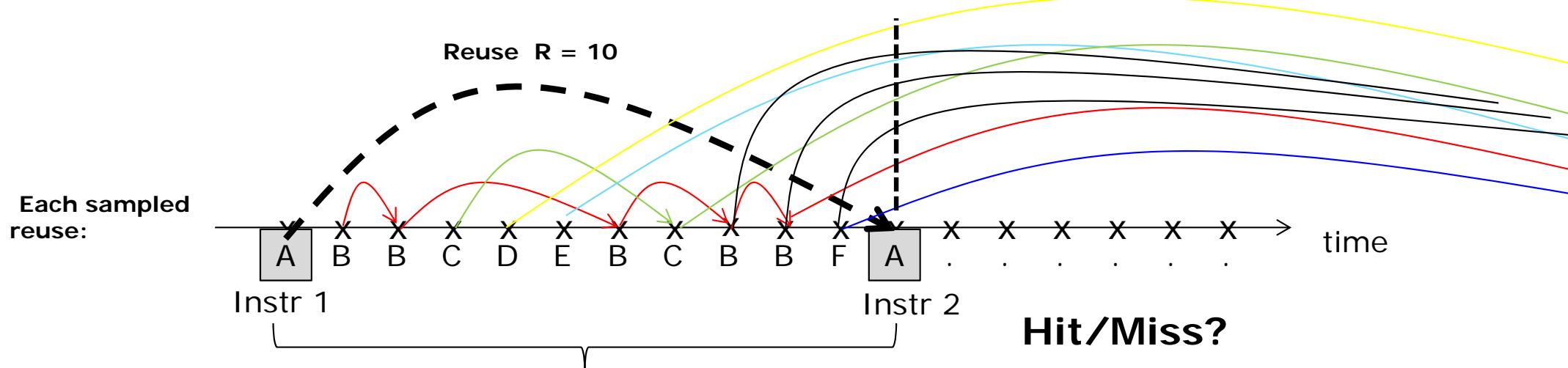
# Accuracy for modeling caches with random replacement



[Berg, ISPASS 2004], [Berg, SIGMETRICS 2005]

- Randomly halting instructions is a more tricky
  - Application phase behavior requires “windowing”
- [Sembrant, CGO 2010]
- How to get down to 20% overhead

# StatStack: Modeling LRU Caches

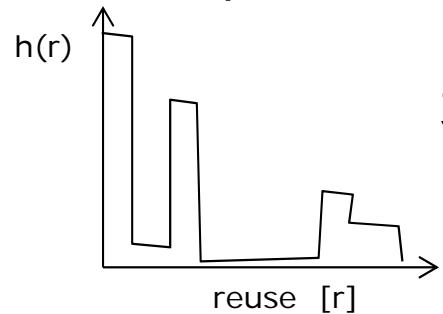


**Stack distance: #Unique cache lines (=5) → Hit for cache size > 5**  
**(...but capturing stack distances is expensive ☹)**

**If we knew all the other reuse distances: #Jump-outs = Stack distance**  
**(...but we do not know all the reuses ☹)**

**Can be estimated using the reuses [from the vicinity]**

Reuse histogram  
for all samples

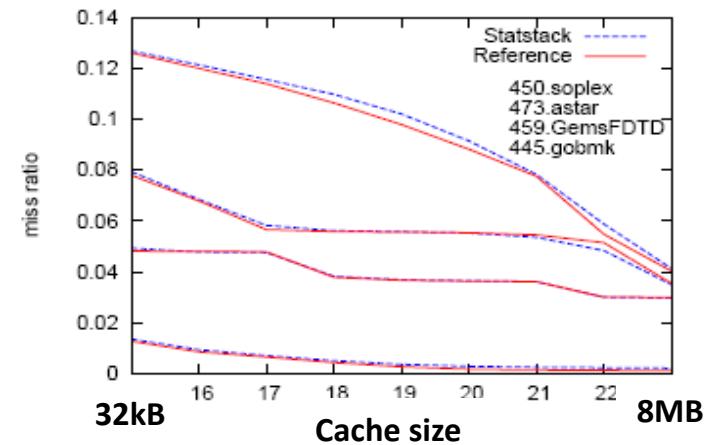
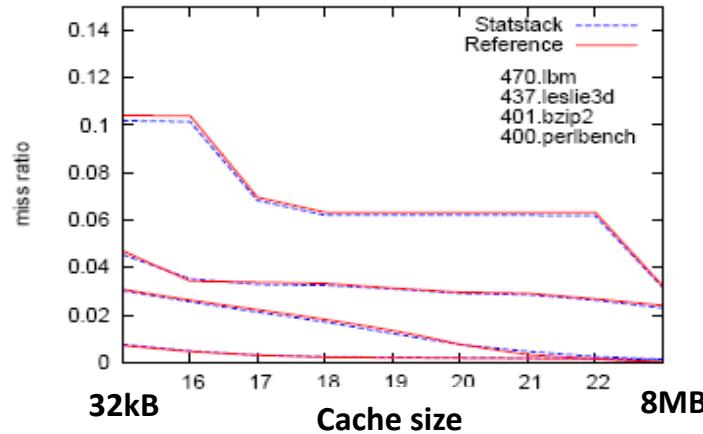
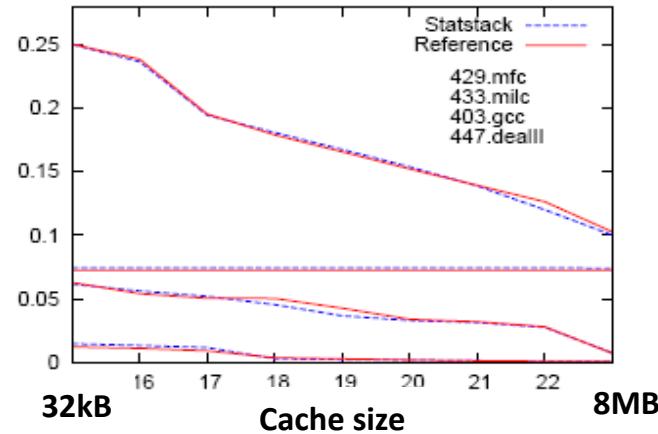


$$\text{Stack-distance}(R) = \sum_{i=1}^{R-1} (P(\text{reuse} > i))$$

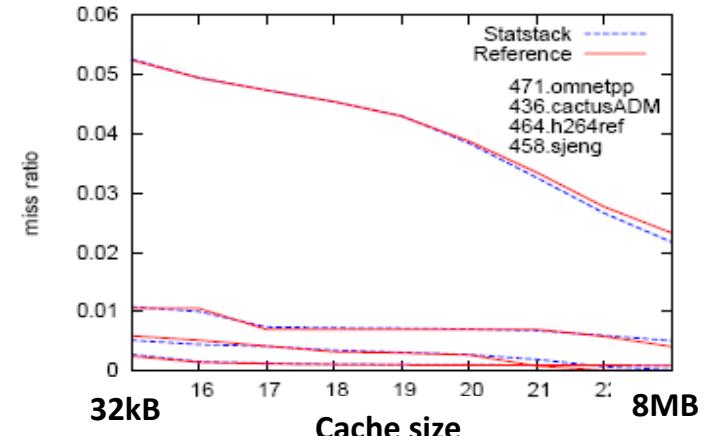
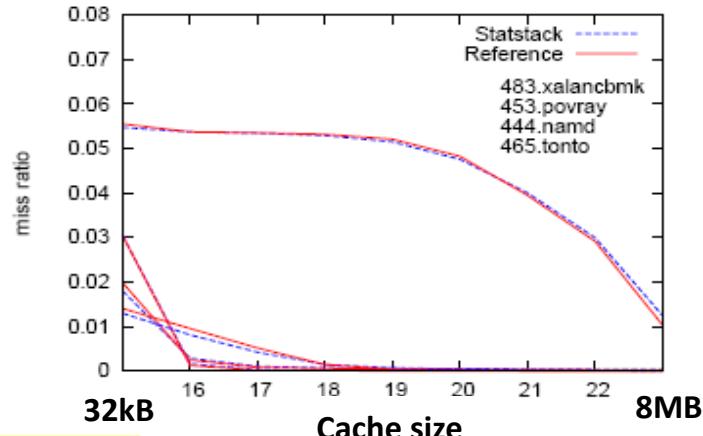
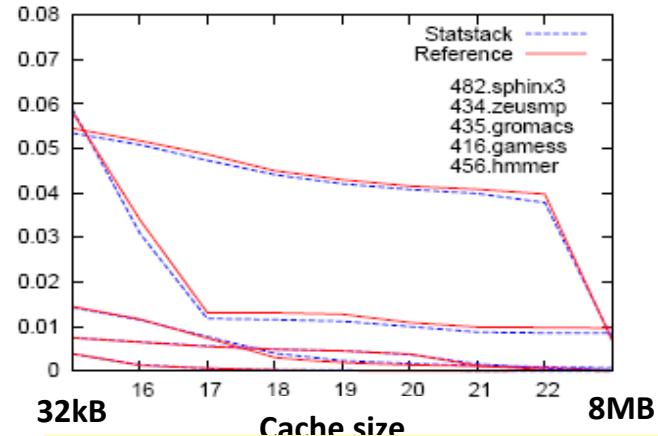
**Hit if (cache size > Stack-distance(R))**

# Accuracy for modeling LRU caches of all sizes

Miss ratio

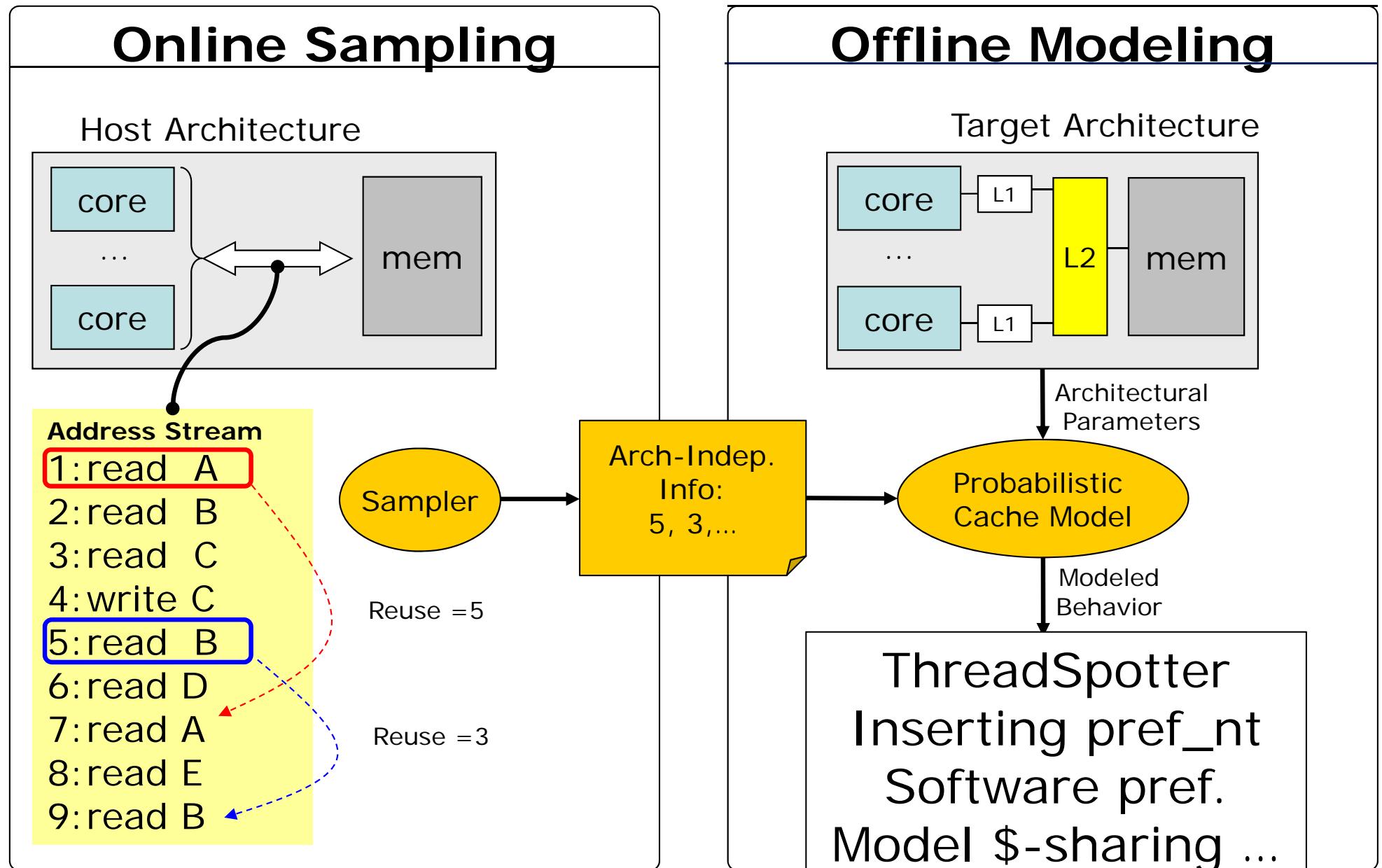


Miss ratio



[Eklov, ISPASS 2008]

# 1. Example: Fast Cache Modeling





# UPMARC

Uppsala Programming for  
Multicore Architectures  
Research Center

- Uppsala Programming for Multicore Architecture Center
- 62 MSEK grant / 10 years [\$9M/10y]  
+ related additional grants at UU = 130MSEK
- Research areas:

Erik:

- Performance modeling
- New parallel algorithms
- Scheduling of threads and resources
- Testing & verification
- Language technology
- MC in wireless and sensors