



**ESC 2011 - Bertinoro**

# **Understanding Low-level Performance Tuning**

**Sverre Jarp, Andrzej Nowak**

**CERN openlab**

**October 28<sup>th</sup> 2011**



- 1. Rationale and background**
  - 2. Software performance tuning on a modern PC**
  - 3. Drilling down on performance figures**
  - 4. Practical tips**
  - 5. Tools**
- 
- > In this talk, we focus on processor performance**
  - > We focus on Intel Core processors, but the same techniques can often be applied to other case**

# Rationale and background

# Seven dimensions of performance

## First three dimensions:

Hardware vectors/SIMD

Pipelining

Superscalar

## Next dimension is a “pseudo” dimension:

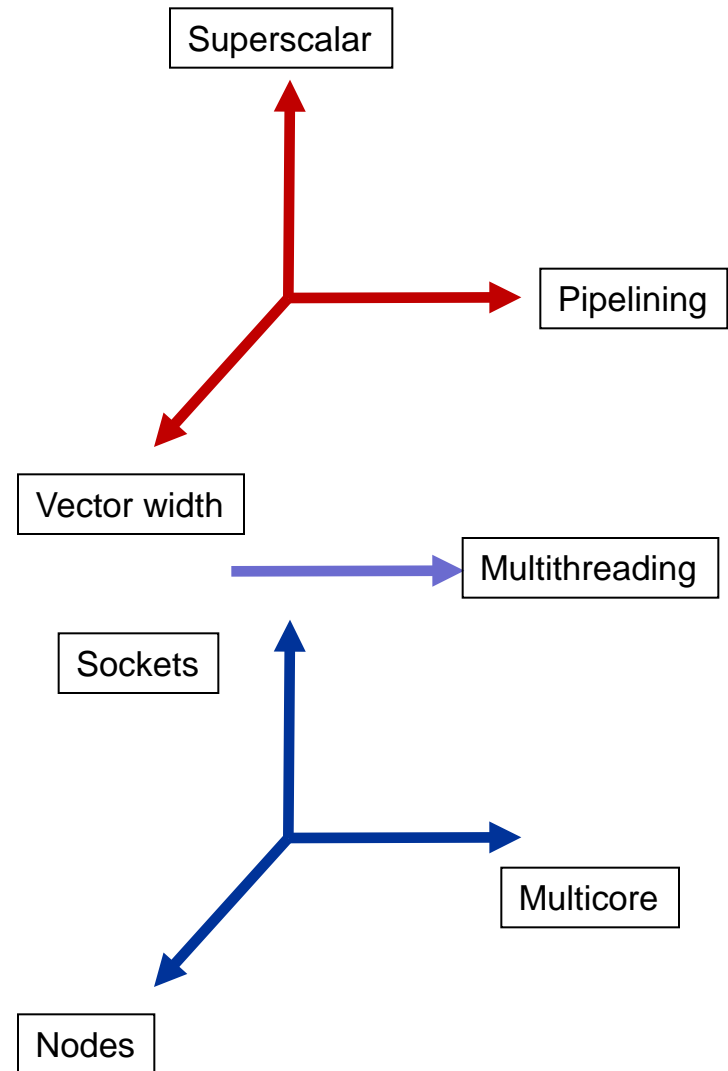
Hardware multithreading

## Last three dimensions:

Multiple cores

Multiple sockets

Multiple compute nodes

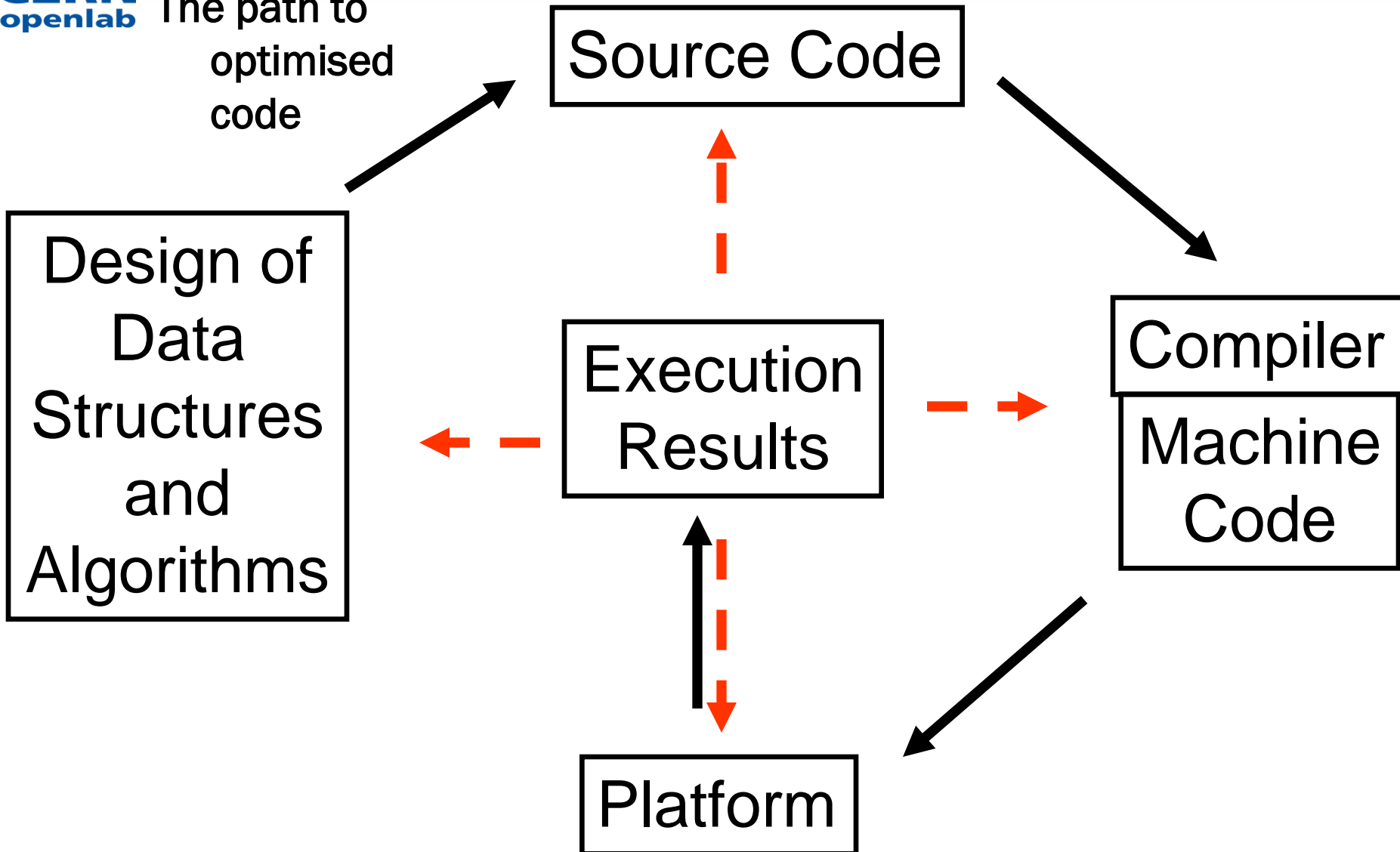




**CERN**  
openlab

# An iterative approach

The path to  
optimised  
code



# Performance tuning levels – reality check

Level	Potential gains	Estimate
Algorithm	Major	~10x-1000x
Source code	Medium	~1x-10x
Compiler level	Medium-Low	~10%-20% (more possible with autovec or parallelization)
Operating system	Low	~5-20%
Hardware	Medium	~10%-30%

There is a facility that can monitor all of the items above  
and their interaction – [hardware counters](#)

# Performance tuning on a modern PC

Key techniques

# Measuring performance

## > **The most common performance measurement unit is time**

- Wall clock time – “how long do I have to wait for it to be done?”
- CPU time – “for how long is the computer busy?”
- Latency – “how long do I have to wait to get an answer?”
- Throughput – “how much of X in a period of time?”

## > **Minimizing time/latency is not the same as maximizing throughput**

- and vice versa – i.e. see Amdahl’s and Gustafson’s laws



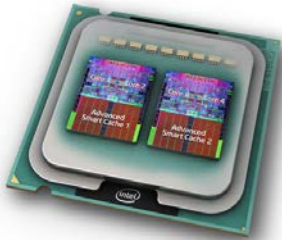
# Performance monitoring in hardware

- > Most modern CPUs are able to provide real-time statistics concerning executed instructions...
- > ...via a **Performance Monitoring Unit** (PMU)
- > The PMU is watching your application in real time! (and everything else that goes on inside the CPU)
- > Limited number of sentries (**counters**) available, but they are versatile
- > Counters monitor **events** as they happen
- > Recorded occurrences are called **samples**
- > **Core i7 and other (Nehalem and later):**
  - 2-4 universal counters: #0, #1, (#2, #3)
  - 3 specialized counters: #16, #17, #18
  - Additional 8 uncore counters: #20-#27



- > **Many events in the CPU can be monitored**
  - A comprehensive list is dependent on the CPU and can be extracted from the manufacturer's manuals or from relevant tools
  - Examples: cache misses, instructions executed, cycles, loads, vector operations
- > **On some CPUs (i.e. Intel Core), some events have bit-masks which limit their range, called “unit masks” or “umasks”**
  - Example: memory instructions retired: “ALL” or “only LOAD” or “only STORE”
- > **Extensive information: Intel Manual 248966-023a**
  - Intel Manual 248966-023a “Intel 64 and IA-32 Architectures Optimization Reference Manual”
- > **AMD CPU-specific manuals**
  - i.e. “BIOS and Kernel Developer’s Guide for AMD Family 10h Processors” or “Software Optimization Guide For AMD Family 10h and 12h Processors”

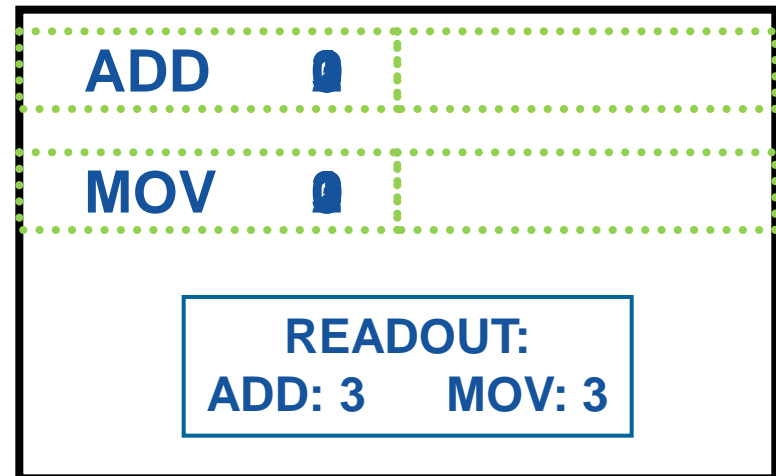
# The Performance Monitoring Unit



Let's monitor  
ADD and MOV  
instructions



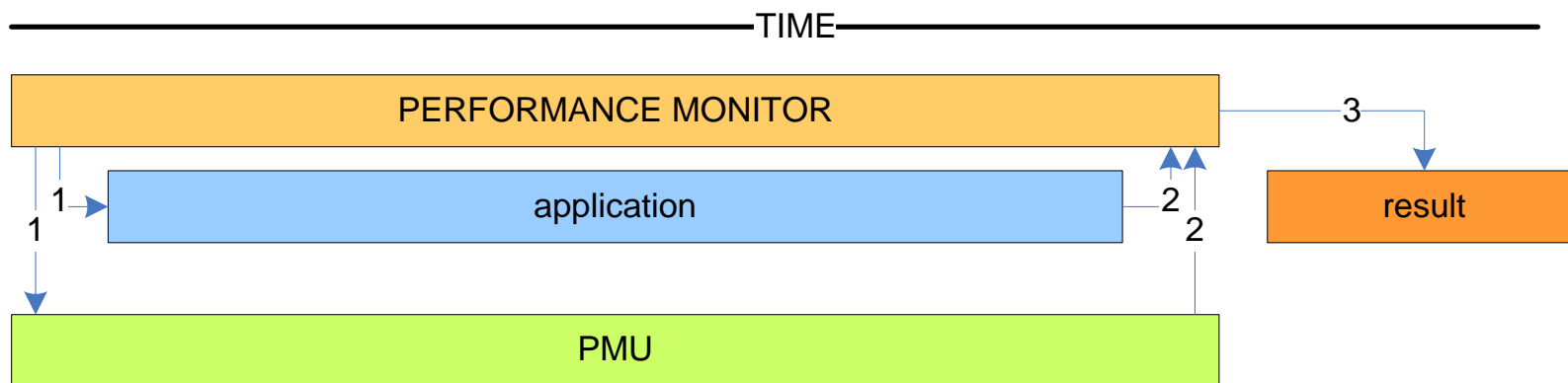
**RETIRED INSTRUCTIONS**  
(=successful & useful execution)



**PERFORMANCE MONITORING UNIT**  
(PMU)

# Popular methods for performance monitoring: Counting

1. Programming the PMU with the specified event(s)
2. Reading the elapsed counts
3. Producing the result

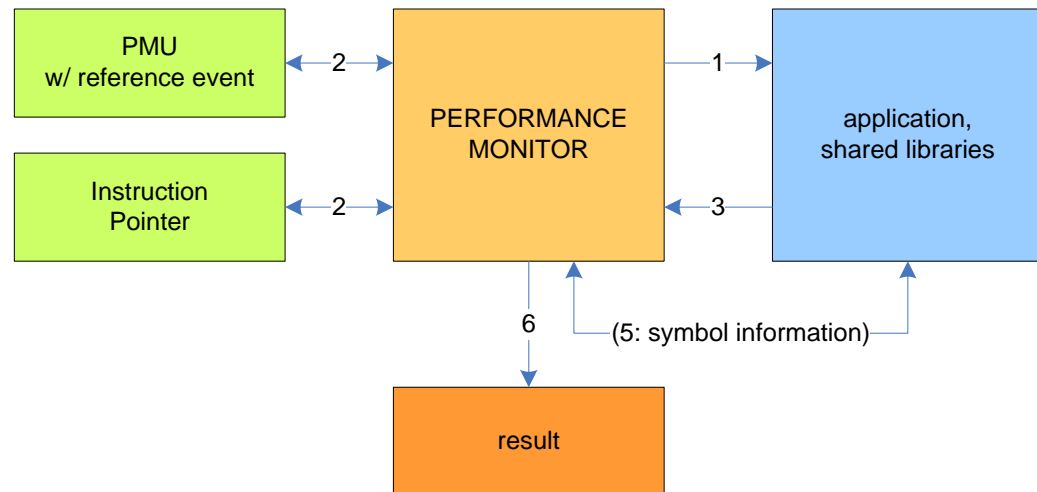


# Popular methods for performance monitoring: Sampling (a.k.a. “Profiling”)

1. The PMU is programmed with the event(s) of interest
2. The application is started
3. Hardware performance counters count in the background
  - > When a pre-programmed value is reached, a performance monitoring interrupt is sent
  - > The interrupt handler can perform a variety of operations, but most of the time it will try to register the state of the machine, especially the instruction pointer – this is called a sample

## 4. Finalization

- Once the application finishes, the performance monitoring software consolidates and renders the results. There are different ways to present the captured data, and they also depend on the type of counters used.



# Popular methods for performance monitoring: Other useful mechanisms

## > Event multiplexing

- Since the amount of simultaneously monitored events is limited, a technique called “multiplexing” can be used. It consists of frequent, periodical time-based switching of the monitored events on the counters. Final results are produced through extrapolation, and the choice of events in each group is important for accuracy.

## > System wide profiling vs. application profiling

- Some tools monitor the whole operating system and later allow you to restrict the results just to your application or library. Some tools (like pfmon) have built in context switching logic which allows for direct monitoring of applications.

# Popular methods for performance monitoring:

## Other useful mechanisms

### > **Instrumented vs. non-instrumented monitoring**

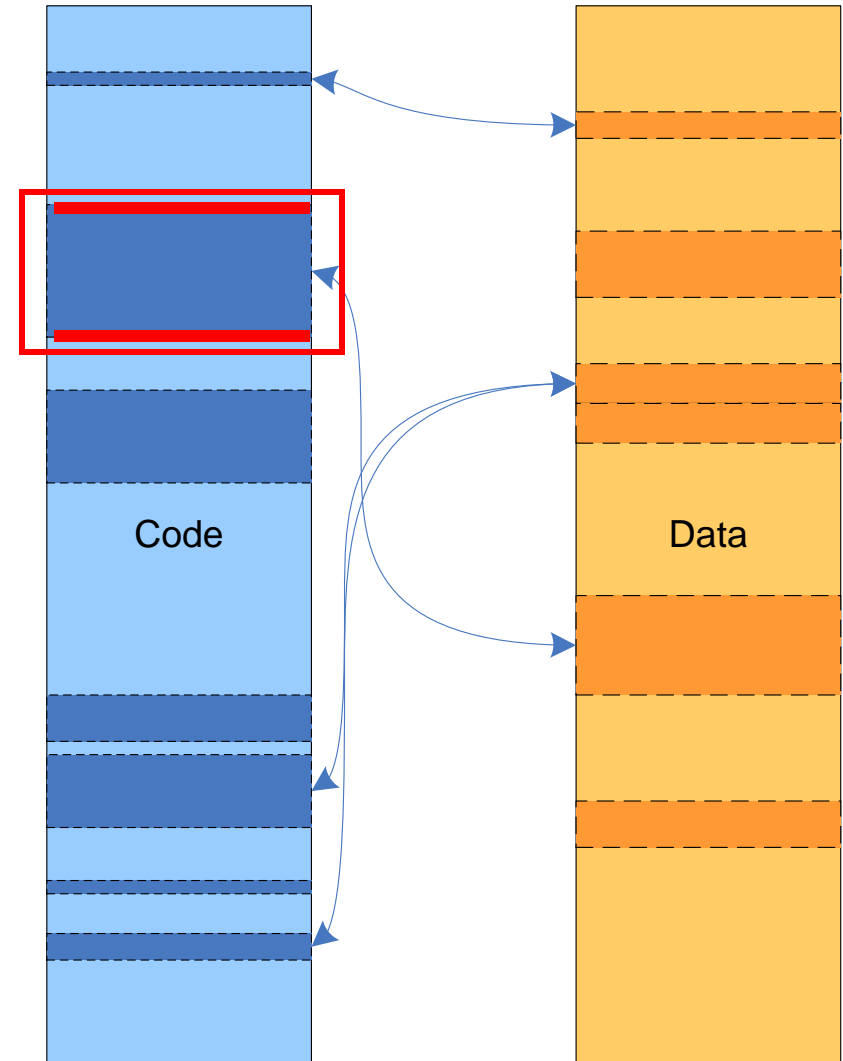
- Some analysis types require binary instrumentation. That means that the monitoring tool inserts probes into the binary code and changes the way the monitored application is running. Such activities might slow down the application by orders of magnitude!
- Other analysis types don't require binary instrumentations and can be performed in the background without disturbing the application. The performance penalty is usually minor, in the order of 1-2%.

### > **User level vs. kernel level**

- Most tools will allow you to choose whether you wish to monitor on the user level or the kernel level (e.g. pfmon). Kernel level monitoring is very handy to debug drivers or system call abuse. You can get a profile and counts just as you would with a userland application.

# Popular methods for performance monitoring: Triggers and triggering

- > **Automatically start or stop monitoring**
- > **Trigger types:**
  - Code
  - Data
- > **A symbol name...**
  - i.e. “foobar”
- > **...or an address**
  - i.e. 0x8103b91e





# Common performance figures

And how to interpret them

# Basic information about your program:

## Recap

### > The amount of:

- instructions executed
- processor cycles spent on the program
- transactions on the bus

### > The amount/percentage of:

- memory loads and stores
- floating point operations
- vector operations (SIMD)
- branch instructions
- cache misses

# Advanced information about your program

## > **The amount and type of:**

- micro-ops executed
- SIMD instructions executed (and the kind)
- resource stalls within the CPU

## > **Cache access characteristics**

- A rich set on Intel Core CPUs
- Requests (missed / hit / total / exclusive or shared / store or read)
- Lines modified / evicted / prefetched

## > Too much information available?

- Low level and fine grained events can be combined to produce ratios (so called “derived events”)

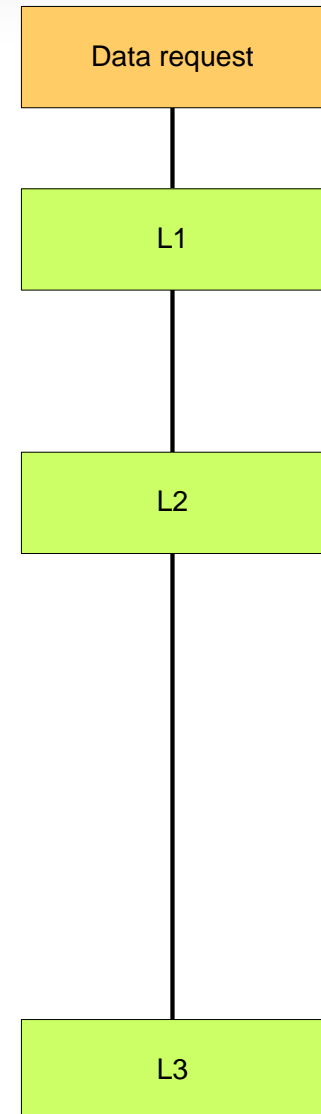
## > Examples (discussion follows):

- Cycles per instructions [CPI]
- Cache miss ratio or impact
- Branch misprediction ratio
- Modified data sharing ratio
- Percentage of wasted cycles
- Bus occupancy

**Mapping performance monitoring data onto your  
source code and environment requires care and  
experience**

# Cache misses

- > If the requested item is not in the polled cache, a higher level has to be consulted (**cache miss**)
- > **Significant impact on performance**
  - Memory access issues are very common, yet hard to fix
- > **Ratio:**  
**LAST LEVEL CACHE MISSES / LAST LEVEL CACHE REFERENCES**
- > **Tips:**
  - A last level cache hit ratio below 95% is considered to be catastrophic!
  - Usually the figure should be above 99%
  - The overall cache miss rate might be low (misses / total instructions), but the resource stalls figure might be high; always check the cache miss percentage



- > **Branch prediction** is a process inside the CPU which determines whether a conditional branch in the program is anticipated by the hardware to be taken or not
- > Typically: prediction based on history
- > The effectiveness of this hardware mechanism heavily depends on the way the software is written
- > The penalty for a mispredicted branch is usually severe (the pipelines inside the CPU get flushed and execution stalls for a while)

# Branch prediction ratios

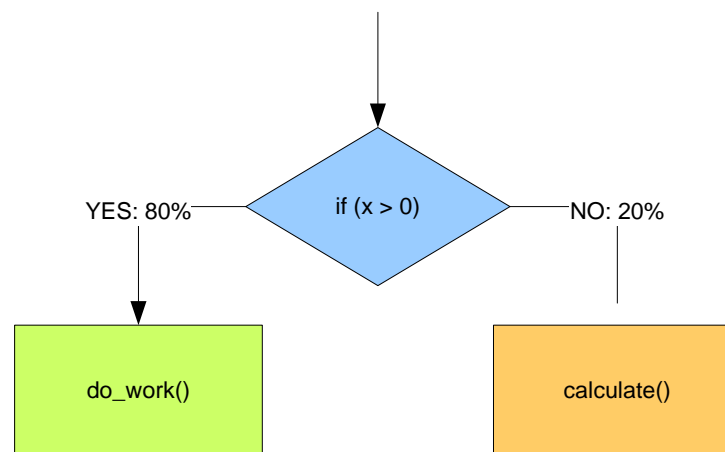
## > The percentage of branch instructions

$\text{BRANCH INSTRUCTIONS} / \text{ALL INSTRUCTIONS}$

## > The percentage of mispredicted branches

$\text{MISPREDICTED BRANCHES} / \text{BRANCH INSTRUCTIONS}$

- The number of incorrectly predicted branches is typically rather low: 20%, 10%, 5%, .... (?)





# Floating point operations

- > Often a significant portion of work of an application
- > May be accelerated using AVX or SSE (SIMD)
- > Related events on the Intel Core microarchitecture:
  - “traditional” x87 FP ops
  - Packed/Scalar single computational SIMD
  - Packed/Scalar double computational SIMD
  - SIMD micro-ops in total
- > Non-computational SIMD instructions can also be counted

# Performance tuning tools

An update on popular performance monitoring facilities

# Popular performance tuning software (1)

## > Linux “perf” – the new performance monitoring subsystem in the Linux kernel

- Pros:
  - Low level access to some counters
  - No patching needed for the application nor the kernel
  - Growing community and toolset
- Cons:
  - A new (and somewhat crude) implementation that still needs time
  - Initial version: dangerously oversimplified

## > perfmon2 – a powerful legacy performance monitoring subsystem for Linux

- Pros:
  - Low level access to all counters
  - No recompilation needed for the application
  - Well established toolset
- Cons:
  - Recompiled kernel or kernel patch needed
  - Development slowed down because of “perf”, shrinking community

# Popular performance tuning software (2)

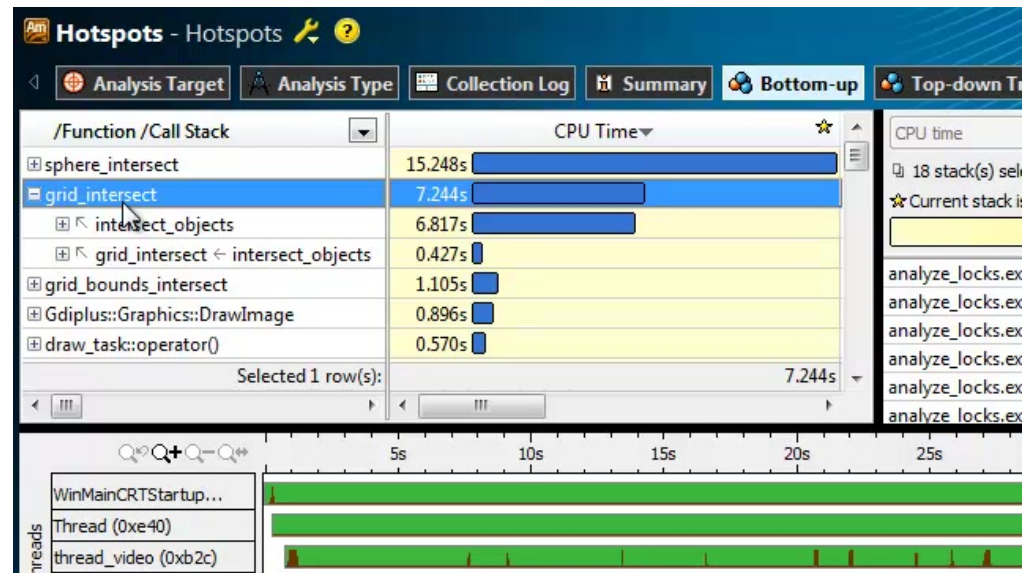
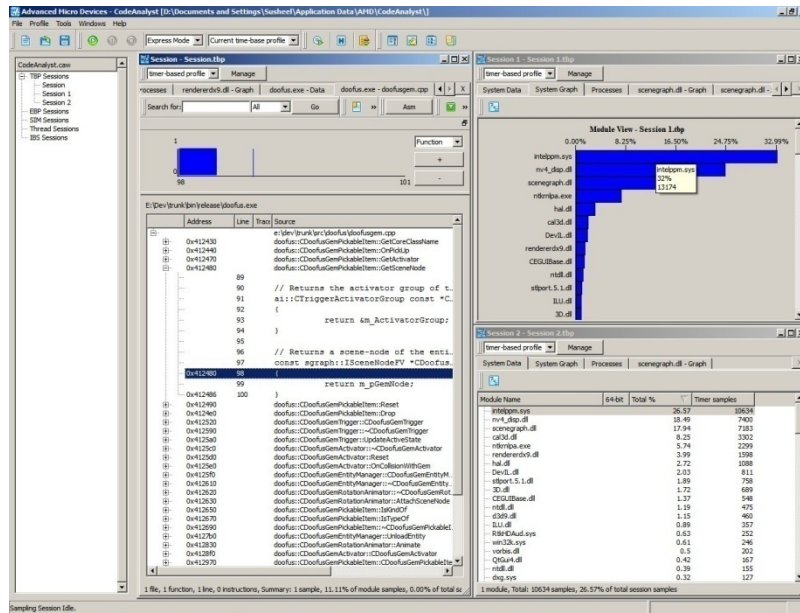
- > **Igprof – Covered by Lassi earlier in the week**
- > **gprof – flat profiler**
  - Recompilation needed
- > **oprofile – flat profiler and event based sampling**
  - Flat profiles
  - Antique; Kernel driver needed
- > **Valgrind**
  - Synthetic software CPU
  - Simulates cache misses and branch mispredictions, memory space profiler, function call relationships
- > **Intel tools - new redesigned tools, native to Linux**
  - “Inspector” - Memory and threading inspector (extended functionality of Thread Checker)
  - “VTune Amplifier” - Performance inspector (extended functionality of VTune and Thread Profiler)

# Popular performance tuning software (3)

## > Intel products (soon available at CERN):

- VTune Amplifier (image on right), Inspector
- PTU (Performance Tuning Utility; next slide)
- Thread Profiler (legacy)

## > AMD CodeAnalyst (image on left)



# Popular performance tuning software (4)

Intel(R) Performance Tuning Utility - loop.c - Eclipse Platform

File Edit Navigate Project Run Window Help

2007-06-27-19-31-45 loop.c X

Source Assembly Control Graph Event of Interest: Clocksicks

Line	Source	Co...	Ins...
597	xr34_l[i] = 0.0;	654	122
598	}		
599	#if !ORG_HIGHEST_SFB		
600	}		
601	else /* cut off the (high...		
602	{		
603	for (i = start; i < en...		
604	xr34_l[i] = 0.0;		
605	}		
606	#endif		
607	total_energy += energy_l[sf...		
608			
609	if (sfb < max_used_sfb_l)	2	
610	xmin_l[sfb] = ratio->l[...	2	2
611	}		
612			
613	for (sfb = min_used_sfb_s; sfb...		
614	{		
615	start = scalefac_band_short...		

Total Selected: 654 122

Address	L...	Assembly	C...	I...
▼ Block 26 5... iteration_loop+03efh:				
0x635F	597	mov eax, DWORD PTR [esp...	9	1
0x6363	597	fild QWORD PTR [eax+edx]	2	
0x6366	597	fabs	1	1
0x6368	597	fild st(0)	8	
0x636A	597	fcomp QWORD PTR [_getch+0...	1	
0x6370	597	instsw ax	2	
0x6372	597	test ah, 0x40h	6	1
0x6375	597	jnz iteration_loop+0439h	4	2
▼ Block 27 5... iteration_loop+0407h:				
0x6377	597	fild st(0)		1
0x6379	597	fsqrt		
0x637B	597	fauil st(0), st(1)	220	61
0x637D	597	fsqrt	11	
0x637F	597	fstp QWORD PTR [edx]	369	55
0x6381	597	fild st(0)	5	
0x6383	597	fauil st(0), st(1)		
0x6385	597	faddp st(2), st(0)		
0x6387	597	fcom QWORD PTR [edi*8+_g...	7	

Total Selected (32 instructions): 654 122

Block 27 Line 597

Block 30 Line 597

Block 31 [iteration\_loop+044ch] (Clockticks: 13)

Address	L...	Assembly	C...
0x63BC	581	inc esi	2
0x63BD	581	add edx, 0x8h	
0x63C0	581	cap esi, ebx	9
0x63C2	581	jnge iteration_loop+03efh	2

Block 32

Block 33 Line 609

Block 35

Line 566

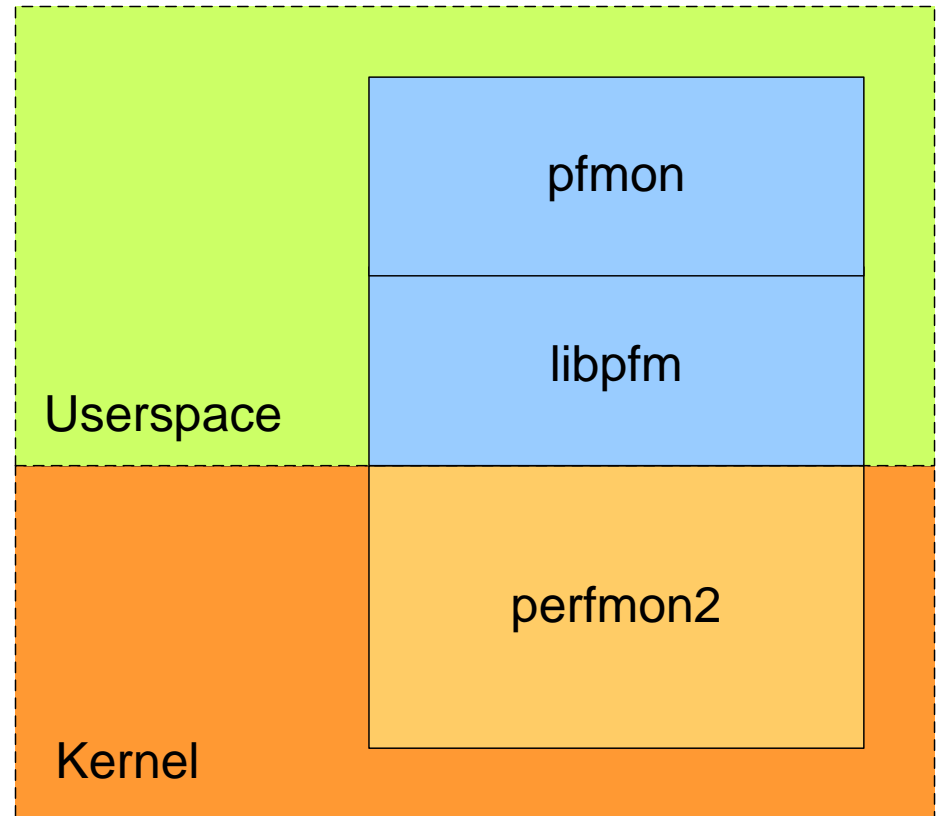
Block 24

# Perfmon2 & pfmon

A real-world performance monitoring framework example  
... and a demo

# Perfmon2 architecture

- > We use it as an example of a robust performance monitoring framework for Linux
- > perfmon2 – kernel part
- > libpfm – userspace interface for perfmon
- > pfmon – “example” userspace application, perfmon2 client





# Perfmon2 – description and rationale

## > **Resides in the kernel**

- Available as a kernel patch
- Very basic functionality

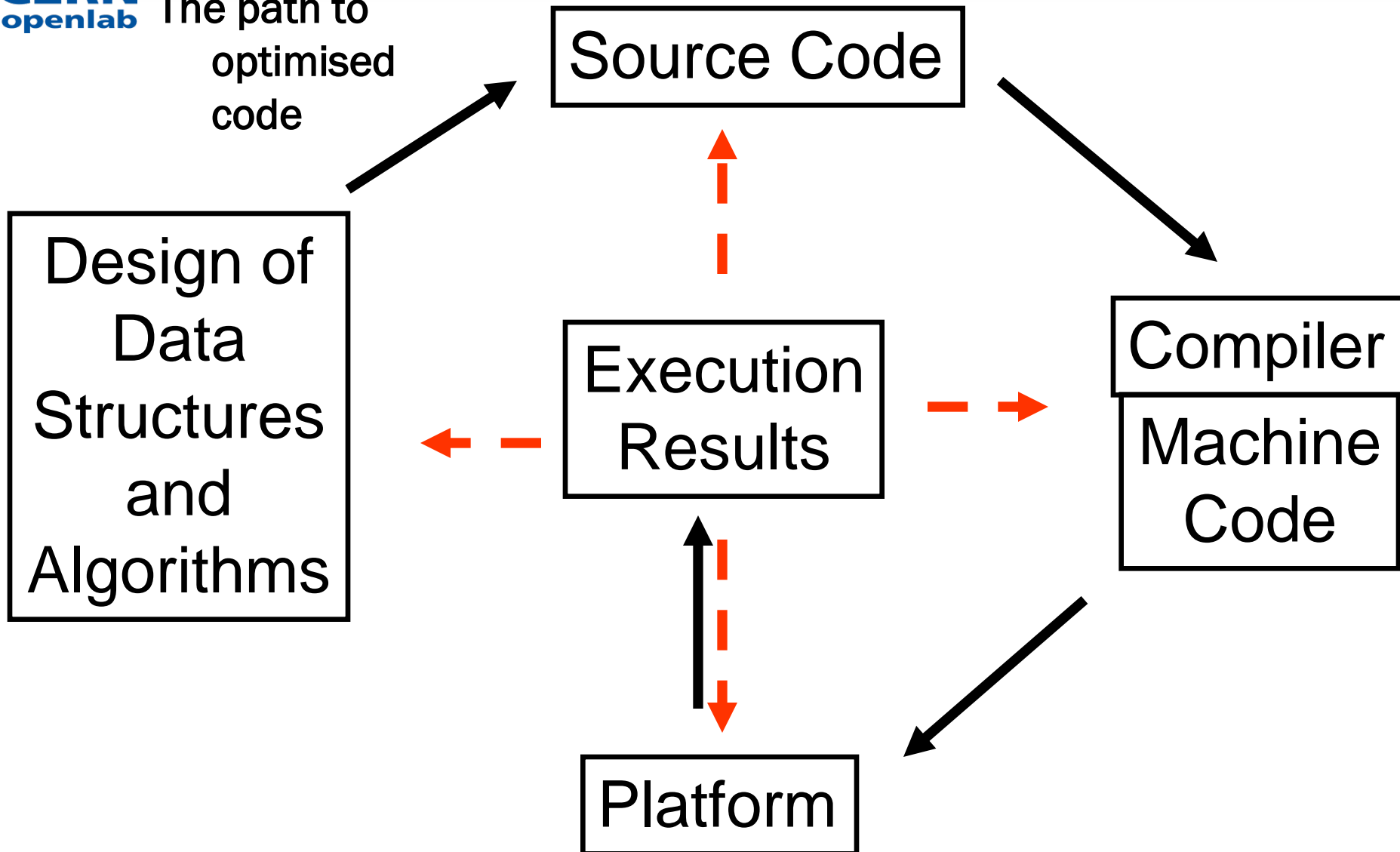
## > **Why was it chosen?**

- Simple to use, lightweight, yet robust
- Support for numerous architectures:
  - x86, x86-64, ia64, powerpc, cell / ps3, mips, sparc
- Supported by numerous hardware vendors:
  - HP Labs, AMD, IBM, Intel, Sony, Toshiba, Cray, SiCortex, Broadcom
  - Also received support from Google and RedHat
- Built on well established experience
- Long history of good support

- > **Console based interface to libpfm/perfmon2**
- > **Provides convenient access to performance counters**
- > **Wide range of functionality:**
  - Counting events
  - Sampling in regular intervals
  - Flat profile
  - System wide mode
  - Triggers
  - Different data readout “plug-ins” (modules) available

# Remember: An iterative approach

The path to  
optimised  
code



# An example: Matrix Multiplication

**Why matmul?**

- 1) Simple, analytically understandable**
- 2) Computationally intensive**
  - a) but memory accesses get in your way, if you are not careful**
- 3) FLP-OPS:  $N^3$  with SSE (DP)**

# Matmul (1): Source code

Two separate source files: matrix.c, multiply.c

```
#define NUM 1024
#define DIM 1024
static double  a[DIM][DIM], b[DIM][DIM], c[DIM][DIM];
....
start = _rdtsc();    // Shown on Monday
multiply_d(a,b,c);
stop = _rdtsc();
```

```
unsigned long i,j,k;
for(i=0;i<NUM;i++) {
    for(k=0;k<NUM;k++) {
        for(j=0;j<NUM;j++) {
            c[i][j] = c[i][j] + a[i][k] * b[k][j];
        }
    }
}
```

```
Minimum FP-OPS = 1.073.741.824
Minimum CYCLES = 536.870.912
```

# Needed pfmon commands

> Find the counters of interest

```
pfmon -l
```

> Decide if “umasks” are required

> Use multiplexing

```
pfmon -iMEM_INST_RETIRED
```

```
Name      : MEM_INST_RETIRED
```

```
Desc       : Memory instructions retired
```

```
Umask-01   : 0x01 : [LOADS] : Instructions retired which cont
```

```
Umask-02   : 0x02 : [STORES] : Instructions retired which con
```

```
pfmon -ecpu_clk_unhalted:thread_p,inst_retired:any_p,
```

```
fp_comp_ops_exe:sse_fp,last_level_cache_references
```

```
    -emem_inst_retired:loads,mem_inst_retired:stores,
```

```
last_level_cache_misses --eu-c --switch-timeout=2
```

```
./matrix
```

# Matmul (2): Run 1

Compile with icc and O2

Fraction of theoretical throughput limit = 13.25%

4.428.551.833	CPU_CLK_UNHALTED:THREAD_P
7.527.231.268	INST_RETIRED:ANY_P
2.150.475.685	FP_COMP_OPS_EXE:SSE_FP
7.234.570	LAST_LEVEL_CACHE_REFERENCES
3.221.299.472	MEM_INST_RETIRED:LOADS
1.074.821.795	MEM_INST_RETIRED:STORES
555.421	LAST_LEVEL_CACHE_MISSES

# Matmul (2): Run 3

Compile with `icc` and `O3, ipo`

Vector-friendly change

Fraction of th	Fraction of throughput limit = <u>27.02%</u>	
4.428.551.833	2.153.265.512	CPU_CLK_UNHALTED:THREAD_P
7.527.231.268	2.683.340.732	INST_RETIRED:ANY_P
2.150.475.685	1.083.687.270	FP_COMP_OPS_EXE:SSE_FP
7.234.570	70.547.893	LAST_LEVEL_CACHE_REFERENCES
3.221.299.472	1.082.271.127	MEM_INST_RETIRED:LOADS
1.074.821.795	538.017.571	MEM_INST_RETIRED:STORES
555.421	170.399	LAST_LEVEL_CACHE_MISSES

Compared to previous run



# Matmul (3): Run 6b

Compile with icc and O3, ipo

Transpose one matrix

Unroll by 4

Cache-friendly changes

Fraction of the	Fraction of throughput limit = <u>66.35%</u>	
2.153.265.512	875.098.637	CPU_CLK_UNHALTED:THREAD_P
2.683.340.732	2.202.587.599	INST_RETIRED:ANY_P
1.083.687.270	1.078.769.120	FP_COMP_OPS_EXE:SSE_FP
70.547.893	4.455.041	LAST_LEVEL_CACHE_REFERENCES
1.082.271.127	512.036.721	MEM_INST_RETIRED:LOADS
538.017.571	240.648.551	MEM_INST_RETIRED:STORES
170.399	254.890	LAST_LEVEL_CACHE_MISSES

Compared to previous run

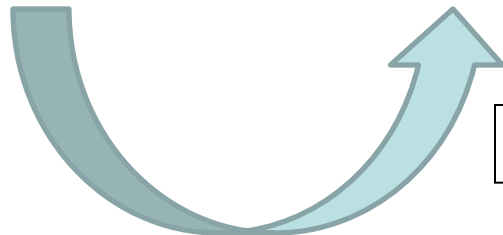
# Matmul (4): Overall gain is 5x

Compile with icc and O3, ipo

Transpose b-matrix

Unroll by 4

Fraction of the	Fraction of throughput limit = <u>66.35%</u>
4.428.551.833	875.098.637 CPU_CLK_UNHALTED:THREAD_P
7.527.231.268	2.202.587.599 INST_RETIRED:ANY_P
2.150.475.685	1.078.769.120 FP_COMP_OPS_EXE:SSE_FP
7.234.570	4.455.041 LAST_LEVEL_CACHE_REFERENCES
3.221.299.472	512.036.721 MEM_INST_RETIRED:LOADS
1.074.821.795	240.648.551 MEM_INST_RETIRED:STORES
555.421	254.890 LAST_LEVEL_CACHE_MISSES



Compared to first run

## > Low-level performance monitoring works best inside-the-core

- On a restricted piece of code:
  - one central algorithm
- Source code fully available
- Good knowledge of compiler and platform
- The theoretically best result is known

# Q & A

