# Performance Analysis and SW optimization for HPC on Intel® Core™ i7, Xeon™ 5500 and 5600 family  Processors*

Presenter: David Levinthal
Principal Engineer

Business Group, Division: DPD, SSG

Version 1.1.2
July 28, 2010

# Legal Disclaimer

- INFORMATION IN THIS DOCUMENT IS PROVIDED IN CONNECTION WITH INTEL PRODUCTS. NO LICENSE, EXPRESS OR IMPLIED, BY ESTOPPEL OR OTHERWISE, TO ANY INTELLECTUAL PROPERTY RIGHTS IS GRANTED BY THIS DOCUMENT.  EXCEPT AS PROVIDED IN INTEL'S TERMS AND CONDITIONS OF SALE FOR SUCH PRODUCTS, INTEL ASSUMES NO LIABILITY WHATSOEVER AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY, RELATING TO SALE AND/OR USE OF INTEL PRODUCTS INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT OR OTHER INTELLECTUAL PROPERTY RIGHT.

  UNLESS OTHERWISE AGREED IN WRITING BY INTEL, THE INTEL PRODUCTS ARE NOT DESIGNED NOR INTENDED FOR ANY APPLICATION IN WHICH THE FAILURE OF THE INTEL PRODUCT COULD CREATE A SITUATION WHERE PERSONAL INJURY OR DEATH MAY OCCUR.

  Intel may make changes to specifications and product descriptions at any time, without notice. Designers must not rely on the absence or characteristics of any features or instructions marked "reserved" or "undefined." Intel reserves these for future definition and shall have no responsibility whatsoever for conflicts or incompatibilities arising from future changes to them. The information here is subject to change without notice. Do not finalize a design with this information.

  The products described in this document may contain design defects or errors known as errata which may cause the product to deviate from published specifications. Current characterized errata are available on request.

  Contact your local Intel sales office or your distributor to obtain the latest specifications and before placing your product order.

- All products, computer systems, dates, and figures specified are preliminary based on current expectations, and are subject to change without notice.

- Customers, licensees, and other third parties are not authorized by Intel to use Intel code names in advertising, promotion or marketing of any product or service.

- Performance tests and ratings are measured using specific computer systems and/or components and reflect the approximate performance of Intel products as measured by those tests.  Any difference in system hardware or software design or configuration may affect actual performance.  Buyers should consult other sources of information to evaluate the performance of systems or components they are considering purchasing.  For more information on performance tests and on the performance of Intel products, visit Intel Performance Benchmark Limitations

# Risk Factors

The above statements and any others in this document that refer to plans and expectations for the first quarter, the year and the future are forward-looking statements that involve a number of risks and uncertainties. Many factors could affect Intel's actual results, and variances from Intel's current expectations regarding such factors could cause actual results to differ materially from those expressed in these forward-looking statements. Intel presently considers the following to be the important factors that could cause actual results to differ materially from the corporation's expectations. Current uncertainty in global economic conditions pose a risk to the overall economy as consumers and businesses may defer purchases in response to tighter credit and negative financial news, which could negatively affect product demand and other related matters.  Consequently, demand could be different from Intel's expectations due to factors including changes in business and economic conditions, including conditions in the credit market that could affect consumer confidence; customer acceptance of Intel's and competitors' products; changes in customer order patterns including order cancellations; and changes in the level of inventory at customers. Intel operates in intensely competitive industries that are characterized by a high percentage of costs that are fixed or difficult to reduce in the short term and product demand that is highly variable and difficult to forecast. Revenue and the gross margin percentage are affected by the timing of new Intel product introductions and the demand for and market acceptance of Intel's products; actions taken by Intel's competitors, including product offerings and introductions, marketing programs and pricing pressures and Intel's response to such actions; Intel's ability to respond quickly to technological developments and to incorporate new features into its products; and the availability of sufficient supply of  components from suppliers to meet demand. The gross margin percentage could vary significantly from expectations based on changes in revenue levels; capacity utilization; excess or obsolete inventory; product mix and pricing; variations in inventory valuation, including variations related to the timing of qualifying products for sale; manufacturing yields; changes in unit costs; impairments of long-lived assets, including manufacturing, assembly/test and intangible assets; and the timing and execution of the manufacturing ramp and associated costs, including start-up costs. Expenses, particularly certain marketing and compensation expenses, as well as restructuring and asset impairment charges,  vary depending on the level of demand for Intel's products and the level of revenue and profits. The recent financial crisis affecting the banking system and financial markets and the going concern threats to investment banks and other financial institutions have resulted in a tightening in the credit markets, a reduced level of liquidity in many financial markets, and extreme volatility in fixed income, credit and equity markets.  There could be a number of follow-on effects from the credit crisis on Intel's business, including insolvency of key suppliers resulting in product delays; inability of customers to obtain credit to finance purchases of our products and/or customer insolvencies; counterparty failures negatively impacting our treasury operations; increased expense or inability to obtain short-term financing of Intel's operations from the issuance of commercial paper; and increased impairments from the inability of investee companies to obtain financing. Intel's results could be impacted by adverse economic, social, political and physical/infrastructure conditions in the countries in which Intel, its customers or its suppliers operate, including military conflict and other security risks,  natural disasters, infrastructure disruptions, health concerns and fluctuations in currency exchange rates. Intel's results could be affected by adverse effects associated with product defects and errata (deviations from published specifications), and by litigation or regulatory matters involving intellectual property, stockholder, consumer, antitrust and other issues, such as the litigation and regulatory matters described in Intel's SEC reports.

# Performance Analysis Methodology for HPC

- **Measure application performance**
  - **Time or rate of work**
  - **Compare to other platforms**
- **Analyze the contributions to performance bottlenecks methodically**
  - **Top Down**

# Performance Analysis Methodology for HPC

- **Two possible objectives**
  - **Influence future silicon design**
    - **Intel personnel do lots of this**

  - **Modify build and/or source to improve performance**
    - **The sole focus of this presentation**

- **The central objective is to identify performance bottlenecks and <u>estimate the potential gain for fixing them</u>**
  - **Without an accurate estimate of the gain a great deal of effort can be wasted**

(intel®)

# Structure of this presentation

- What would the author do with:
  - A brand new machine
  - A tar ball of 100 million source lines
  - Documented, working build procedure
  - Data set and instructions to run the app
  - And one commandment:

# Make Go Fast

**but get the same answer**

# Presentation Agenda

- Optimization workflow overview
- Event based sampling
  - Why so complicated
  - How the nuts and bolts work
- HPC/Scientific computing overview
- Compiler problems/tuning compiler usage
- Identifying and removing stalls
- Identifying and removing resource saturation
- Identifying and removing non scaling
- PTU features and data interpretation
- Glossary in backup

# Performance Analysis Methodology

- **The steps**
  - **1. make sure the platform is correct**
    - **It should be – some thought went into the specifications**
    - **But don't take this for granted**
  - **2. Use the correct compiler (Intel® Compiler)**
    - **And invoke it correctly**
    - **This should also have already been done…but..**
  - **3. Analyze interaction of SW and micro architecture and tune code/compiler usage**
    - **Intel® VTune™ Analyzer* or better, Intel® Performance Tuning Utility (PTU)**
    - **Iterative process**
  - **4. Parallelize the execution as appropriate**
    - **Batch queue / Intel® MPI Library**
    - **OpenMP** product, Intel® Threading Building Blocks (Intel® TBB), Intel® CILK™ Plus, explicit threading**
- **Iterate on 3 and 4**

*Vtune and CILK are trademarks of Intel Corporation in the U.S. and other countries.
**Other names and brands may be claimed as the property of others.

# Platform Optimization: Step 1

- **1. Make sure the platform is correct**
  - **Enough memory**
    - **Page faults  (Perfmon*, vmstat*)**
      - rates of >100 /sec is cause for investigation
    - **Make sure DIMMs are in identical sets of 6 for DP machines**
      - 3 channel memory controller
      - Best performance with completely uniform dimms
  - **Make sure SATA Bios setting is AHCI, not IDE setting**
    - **Use RAID or SSD if disk speed is critical**
  - **Prefetcher BIOS Settings correct for the app: <u>ON</u>**
    - **Intel® 11.0 compiler can generate SW prefetch**
  - **NUMA BIOS setting correct: <u>ON</u>**
  - **Intel® Hyper-Threading Technology BIOS option set correctly for the application**
    - **HT does not always help HPC**
      - Probably makes little difference

## Disable C states to ensure machine stability when using event based sampling on Corei7/Xeon 5500

**9**

(intel®)

# Compiler Usage Optimization: Step 2

- **2. Optimize the time consuming functions**
  - **Profile functions, and check compiler options**
  - **Intel® VTune™ Analyzer and Intel® PTU have source file granularities**
    - **Data grouped per source file to identify hot files**
  - **Do not assume this has been done**
    - **Build environments are complex**

# Micro architectural Optimization: Step 3

- **3. Identify & Optimize the time-consuming functions**
- **Use performance events methodically to identify performance limitations**
  - **Intel® PTU, Intel® VTune™ Analyzer, etc.**
- **Confirm that compiler really did produce good code (visual inspection of ASM)**
  - **For the components of the code using the cycles**
- **Go after largest, easy things first**
  - **<u>Accurate estimate of potential gain is critical!</u>**

- **Documentation for Intel® Core™ i7 processor Performance Monitoring Unit (PMU) is available**

# Parallelization for HPC : Step 4

- **4. Use as many cores and machines as possible**
  - **Parallel processing by batch queue is OK**
    - **Trivial parallelism**
    - **Hard to beat the throughput**

intel®

# Parallelization for HPC : Step 4

- **4. Use as many cores and machines as possible**
  - **Figure out clean data decomposition**
  - **Intel® MPI Library for process parallel execution**
    - **Minimal shared elements**
    - **Maximal address separation**
  - **OpenMP*, Intel® TBB, CILK, explicit threading for shared memory**
    - **Can reduce all to all MPI API costs**

\* Other names and brands may be claimed as the property of others.

(intel®)

# DP Platform

# Event Based Sampling Analysis

- **Code profiling with performance events can identify where the interaction of the code and data with the microarchitecture is sub optimal**
  - **Ex: What code execution results in load driven cache misses?**
  - **Event_count*Penalty ~ potential gain**
    - **A well defined penalty is <u>essential</u>**
- **Such profiling also provides an execution weighted display of the generated instructions**
  - **Vectorized code was generated but is it being executed?**

**But There are THOUSANDS of Events, Which Ones Matter?**

(intel)

# Which Events you need depends on what problem you wish to study and what you want to accomplish
## Example: Last Level Cache Misses

- What you mean by an LLC miss depends on the exact nature of the question you are asking

- Are you asking about Bandwidth consumption?
  - Due to reads?, RFOs?, HW Prefetch, NT stores? Total?, Code?, SW prefetch?, Cacheable Writebacks?
  - Location of the bandwidth consumption?
  - Source of the data provided?

- Or about Latency/Pipeline stalls
  - Different architectures stall on different things
  - Intel® IA-32/Intel64 Processors' memory access stalls are mostly due to **loads**

**Events needed to measure bandwidth and memory stalls are COMPLETELY different**

(intel®)

# Intel® Xeon™ 5500 load Penalties

| | L1D_HIT | Secondary Miss | L2 Hit | LLC Hit No Snoop | LLC Hit Clean Snoop | LLC Hit Snoop =HITM | Local Dram | Remote Dram | Remote Cache local home Fwd | Remote Cache Remote Home FWD | Remote Cache Local Home HITM | Remote Cache Remote home HITM |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Mem_load_retired .L1d_hit | 0 (By Def) | | | | | | | | | | | |
| Mem_load_retired .Hit_LFB | | 0->Max Val | | | | | | | | | | |
| Mem_load_retired .L2_hit | | | 6 | | Depend on frequency, dimms, bios, etc | | | | | | | |
| Mem_load_retired .LLC_Unshared_hit | | | | ~35 | | | | | | | | |
| Mem_load_retired .other_core_l2 _hit_hitm | | | | | ~60 | ~75 | | | | | | |
| Mem_load_retired .LLC_Miss | | | | | | | ~200 | ~350 | ~180 | ~180 | ~225 -250 | ~370 |
| Mem_uncore_retired .Other_core_l2_hitm | | | | | | ~75 | | | | | | |
| Mem_uncore_retired .Local_Dram | | | | | | | ~200 | | | | ~225 -250 | |
| Mem_uncore_retired .Remote_dram | | | | | | | | ~350 | | | | ~370 |
| Mem_uncore_retired. Remote_cache_ local_home_hit | | | | | | | | | ~180 | | | |

Note: All latencies and memory access penalties shown are merely illustrative. Actual latencies will depend on (among other things) processor model, core and uncore frequencies, type, number and positioning of DIMMS, platform model, bios version and settings. Consult the platform manufacturer for optimal setting for any individual system. Then measure the actual properties of that system by running well established benchmarks.

**The Important Penalties Vary by a Factor of TEN**

(intel)

# Intel® PTU uses profiles to manage complexity

# Intel® PTU predefined collections

- **Cycles and Uops**
  - **Cycle usage and uop flow through the pipeline**
- **Branch Analysis**
  - **Branch execution analysis for loop tripcounts and call counts**
- **General Exploration**
  - **Cycles, instructions, stalls, branches, basic memory access**
- **Memory Access**
  - **Detailed breakdown of off-core memory access (w/wo address profiling)**
- **Working Set**
  - **Precise loads and stores enabling address space analysis**
- **FrontEnd (FE) Investigation**
  - **Detailed instruction starvation analysis**
- **Contested lines**
  - **Precise HITM and Store events**
- **Loop Analysis**
  - **32 events for HPC type codes, w/wo call sites , i.e. including LBR capture**
- **Client Analysis**
  - **54 events for client type codes, w/wo call sites , i.e. including LBR capture**

**Many Possible Issues -> Many Different Events**

(intel)

# Controlling collection

# Performance Monitoring Unit

- **The Performance Monitoring Unit (PMU) consists of a set of counters that can be programmed to count user-selected signals of microprocessor activity**
  - **Cpu_clk_unhalted, inst_retired, LLC_miss, etc..**
- **Counting the number of events that occur in a fixed time period allows workload characterization**
  - **Using a spectrum of events allows a decomposition of the applications activity with respect to the microarchitecture components**
  - **Particularly useful for studying the architecture's strengths and weaknesses running an application**

(intel)

# Performance Monitoring Unit

- **The PMU can be programmed to generate interrupts on counter overflow**
  - **Allows periodic sampling of program counter for any user-chosen event**
    - **Initialize count to (overflow – periodic rate)**
  - **Interrupt Vector Table is programmed with the address of the interrupt handler**
    - **Intel® VTune™ Analyzer driver is invoked by HW on counter overflows and given a program counter where the interrupt (i.e. counter overflow) happened**

- **Identify statistically where events occur in the program**
    - **Application profiling by event**

(intel®)

# SKID:
# IP of causal instruction vs IP of PMI

Event associated with IP +N

| | |
|---|---|
| BPU | |
| Fetch | IP+N |
| Q | |
| Q | |
| decode | |
| decode | |
| Q | |
| Q | |
| Q | |
| RAT | |
| RS | |
| RS | |
| RS | |
| RS | IP+3 |
| exec | IP+2 |
| exec | IP+1 |
| Retire | IP |

Time

L1_ifetch_miss
overflow

delay

LLC_ifetch_miss
overflow

PMI occurs at IP +          IP +m
propagation delays

(intel)

# Analyzing HPC Applications

- **Overview**
- **Loop analysis**
  - **Tripcounts**
  - **Vectorization**
- **Memory access dominated**
  - **Latency dominated**
  - **Bandwidth dominated**
- **Execution dominated**

# Overview

- **Performance Breakdown/cycle accounting can be applied to any scale of a program**

  - Multiple interacting applications-> single apps-> single modules-> source files/functions-> basic blocks

- **Methodology does not change**

  - But can inherit conclusions from higher levels based on importance/cycle cost

- **At all stages in the process look for poorly written, actively executing code that can be improved**

# HPC Applications

- **Dominated by loops**
- **Rarely have pipeline front end problems**
  - **Except for very large binaries (ifetch latency)**
- **Large data sets**
  - **Not cache resident**
  - **Ex: Weather simulation, Oil Reservoir**
  - **Frequently DRAM bandwidth limited**
  - **Or DRAM Latency limited**
- **Occasionally HPC apps are uop flow limited**
  - **Data blocked**
  - **Ex: oil exploration, FFTs**

# What matters when optimizing a loop?

1. **The Trip Count**

2. **The Trip Count**

3. **The TRIP COUNT!**

4. **Variations in the tripcount**

5. **And some other things**

   **BUT..what you do about them depends on**

   **THE TRIP COUNT**

**And of course there are virtually no tools to assist you in determining this..other than printf**

**(you can use PIN..)**

This Will be Discussed Later

intel

# HPC Loops and Memory Access

- **Calculations require data as input and the most severe limitations in a computer are on data access**

  - **CPU speed and efficiency have increased much faster than memory speeds and bandwidth.**

- **Load operations are almost always scheduled almost immediately before consumption (adds, multiplies etc)**

- **Lack of availability will quickly lead to execution stalls**

  - **OOO execution can buy only a few cycles.**

(intel)

# Event Classes: High Level View

1. Execution flow events
   - Cycles, Branches, stalls, uops/inst_retired
   - Guide compiler usage

2. Penalty events
   - Ex: load requiring access to dram
   - Modify code/build to reduce penalties

3. Resource saturation events
   - Bandwidth, load/store buffers, dispatch ports
   - No well defined cost
   - Change data layout/access patterns

4. Architectural characterization
   - Cache accesses, MESI states, snoops
   - Used to improve silicon design, not application performance

5. Instruction mix
   - Do not measure what you think, extremely difficult to validate

(intel)

# Event Classes

1. **Execution flow events: Guide Compiler Usage**

   - **Cycles, Branches, stalls, uops/inst_retired**

2. Penalty events

   - Ex: load requiring access to dram

3. Resource saturation events

   - Bandwidth, load/store buffers, dispatch ports
   - No well defined cost

4. Architectural characterization

   - Cache accesses, MESI states, snoops

5. Instruction mix

# Cycles: Multiple time domains

- **There are actually 4 cycle events on a modern microprocessor**
  - **Core unhalted cycles**
  - **Reference frequency unhalted cycles**
  - **Core halted cycles**
  - **Reference Frequency halted cycles**
- **Core frequency needed for perf issues entirely in the core**
  - **Penalties (ie pipeline stalls) in core cycles**
- **Reference frequency needed for:**
  - **Evaluation of variable frequency effects (Turbo/Power Management)**
  - **Wall clock time utilization**
    - **Ex: Network server applications**
  - **Bandwidth/memory latency**
- **Unhalted events are required for counting modes to work at all**
- **Halted.ref = TSC change – cpu_clk_unhalted.ref**

(intel®)

# Cycle Accounting and Uop Flow

- **Cycles =
  Cycles dispatching to execution units +
  Cycles not dispatching (stalls)**
  - **A trivial truism**
- **Uops dispatched = uops retired +
  speculative uops that are not retired**
  - **Non-retired uops due to mispredicted branches**
    - **Uops_issued.any – uops_retired.slots**
- **Optimization Reduces Total Cycles by**
  - **Reducing stalls**
  - **Reducing retired uops (better code generation)**
  - **Reducing non retired uops (reducing mispredictions)**

(intel)

# (Simplified) Execution in an OOO Engine

- **Two asynchronous components connected by buffering**
  - **Front End provides instructions**
  - **Back End gets data and executes instructions**
  - **Back End trumps Front End**
    - **If BE issues occur, fixing FE issues accomplishes nothing**

Reservation Station

| Inst Fetch Br Predict | → | Decoder | → | Resource Allocation | | RS | dispatch | Execution Units | | ROB: Reorder Buffer |

**Front End FE**

**Back End BE**

Retirement/Writeback

# Identifying Front End Stalls

- **Uop issue**
  - **Uops have been allocated resources**
  - **No downstream blockage (resource_stalls)**
  - **FE Stalls = an instruction delivery problem**
    **= Uops_issued.stall_cycles – Resource_stalls**



```
Inst Fetch    →  Decoder  →  Resource     →→  RS   dispatch  Execution  →  ROB
Br Predict                   Allocation                      Units
```

**Uops_issued**
**Uops_issued.stall_cycles**
**Uops_issued.core_stall_cycles**
**Resource_stalls**

**Retirement/Writeback**

# (Simplified) Execution in an OOO Engine

- **Design optimizes Dispatch to Execution**
  - **Uops wait in RS until inputs are available**
  - **Keeping the Execution Units occupied matters**

# Uop Flow Monitors Execution

- **Uop Execute**
  - **Uops have inputs ?**
  - **No downstream blockage (DIV/SQRT)**
  - **No execution = no progress**

| Inst Fetch Br Predict | → | Decoder | → | Resource Allocation | → | RS | dispatch | Execution Units | → | ROB |

Uops_executed.portX
Uops_executed.core_stall_cycles
Only non HT events in the core

Retirement/Writeback

# Uop Flow Monitors Execution

- **Uop Retire**
  - **All older instructions retired ?**
  - **No retirement = ?  (out of order execution?)**



| Inst Fetch Br Predict | Decoder | Resource Allocation | RS | dispatch | Execution Units | ROB |

**Uops_retired**
**Uops_retired. stall_cycles**
**Precise retirement stall identifies "push out"**

Retirement/Writeback

intel

# Uop Flow



**Fetch / Decode**

To Uncore

**MEU**

| 32 KB Instruction Cache | Next IP | MLC | 32 KB Data Cache |

Inst_written_to_iq

**Execute**

Branch Target Buffer

Instruction Decode (4 issue)

Microcode Sequencer

Register Allocation Table (RAT)

Reservation Stations (RS) 36 entry

Scheduler / Dispatch Ports

Port Port Port Port

Integer Arithmetic — SIMD — FP Add

Integer Arithmetic — SIMD — Integer Shift/Rotate — FP Div/Mul

Integer Arithmetic — SIMD

**Uops_executed**

Load

Store Addr

Store Data

Memory Order Buffer (MOB)

**Retire**

Re-Order Buffer (ROB) – 128 entry

**Uops_Issued**

IA Register Set

**Qualitative: Artistic License employed**

**Uops_retired**

(intel)

# PEBS Basic Events

- **Mechanism:**
  - **counter overflow arms PEBS**
  - **Next event gets captured and raises PMI**
  - **PEBS mechanism captures architectural state information at completion of critical instruction**

| |
|---|
| **instr_retired** |
| **itlb_miss_retired** |
| **uops_retired** |
| **br_instr_retired** |
| **mem_instr_retired.loads** |
| **mem_instr_retired.stores** |

- **Including EIP (+1), even when OS defers PMI**

**For memory events, EIP (+1) is always next instruction**

# Branch Events

- **Measure Control flow through the program**
- **Can be used for**
  - **loop trip counts**
  - **Reconstructing (multi function) execution paths**
  - **Driving inlining, IPO, PGO compilations**
- **Used in conjunction with Last Branch Record (LBR) even more can be done**
  - **Basic block execution counts**
  - **Instruction mix**
  - **Call counts per source**
  - **etc**

(intel)

# Basic Branch Analysis

- **Vastly improved precise branch monitoring capabilities**
  - **Branches retired**
  - **16 deep LBR**
    - LBR can be filtered by branch type and privilege level
  - **One per SMT**
    - Not merged when SMT disabled
  - **Only taken branches are captured**
- **Precise BR retired by branch type**
  - **Calls, conditional and all branches**
  - **Coupled with LBR capture yields**
    - Call counts
    - "HW call graph"
    - Basic block execution counts

(intel)

# Branch Analysis

- **Precise branch events on NHM enable**
  - **Function call counts**
  - **Function arguments (em64T only)**
  - **Taken fraction/branch**
- **Mispredicted Branches must be counted with Non-PEBS events BR_MISP_EXEC.\* and BR_INST_EXEC.\* on Corei7/Xeon 5500**
- **Br_misp_retired.\* on Xeon 5600 (PEBS)**

# Branch Analysis: Call Counts

- **Call counts require sampling on calls**
  - Sampling on anything else introduces a "trigger bias" that cannot be corrected for
- **Requires PEBS buffer to identify which branch caused the event**
  - EIP+1 results in capturing call target
- **Requires LBR to identify source and target**
  - Matching PEBS EIP with LBR target

# Precise Conditional Branch Retired

- **Counted loops that actually use the induction variable will frequently keep the tripcount in a register for the termination test**

  - **E.g. heavily optimized triad with the Intel compiler has**
    **Addq $0x8, %rcx**
    **Cmpq %rax, %rcx**
    **Jnge    triad+0x27**

- **Average value of RAX is the tripcount**

# Branch Analysis: Function Arguments (Intel64 only)

- **Functions with "few" (<6?) arguments use registers for argument values**

- **Capturing full PEBS buffer + LBR on calls_retired event allows measurement of distribution of argument values per calling site**
  - **E.g. length of memcpy,memset**

# Processing LBRs

| Branch_0 | Target_0 |
|----------|----------|

| Branch_1 | Target_1 |
|----------|----------|

- **All instructions between Target_0 and Branch_1 are retired 1 time**

- **All Basic Blocks between Target_0 and Branch_1 are executed 1 time**

- **All Branch Instructions between Target_0 and Branch_1 are not taken**

**So it would all Seem Very Straight Forward**

(intel)

# Shadowing and Precise Data Collection

- **The time between the counter overflow and the PEBS arming creates a "shadow", during which events cannot be collected**
      **~8 cycles?**

- **Ex: conditional branches retired**
  - **Sequence of short BBs (< 3 cycles in duration)**
  - **If branch into first overflows counter, Pebs event cannot occur until branch at end of 4th BB**
  - **Intervening branches will never be sampled**

# Shadowing

Assume 10 cycle shadow for this example

| | |
|---|---|
| 20 | |
| 20 | |
| 2 | |
| 2 | |
| 2 | |
| 2 | |
| | |
| 20 | |
| | |
| 20 | |

**O**

**P**

**C**   **O**

     **O**

       **O**

         **O**

           **O**

**P**    **P**    **P**    **P**

            **P**

**C**   **C**   **C**   **C**   **C**

| | |
|---|---|
| | |
| N | |
| N | |
| 0 | |
| 0 | |
| 0 | |
| | |
| 0 | |
| | |
| 5N | |

O means counter overflow
P means PEBS enabled
C means interupt occurs

# Reducing Shadowing Impact

- Some "events" will never occur!
  - Falling into shadowed window
- Use LBR to extend range of the single sample
- Count the number of objects in LBR and increment count for all of them by 1/15
  - Since you have only one sample

# Minimizing Shadowing Impact on BB Execution Count

Cycles/branch taken

Pebs Samples taken

Number of LBR entries

| Cycles/branch taken |
|---|
| 20 |
| 20 |
| 2 |
| 2 |
| 2 |
| 2 |
| 20 |
| 20 |

O
P
C  O
O
O
O
O
P
P
P
P
P
C  C  C  C  C

| Pebs Samples |
|---|
| |
| N |
| N |
| 0 |
| 0 |
| 0 |
| 0 |
| 5N |

| LBR entries |
|---|
| 15N |
| 15N |
| 15N |
| 15N |
| 16N |
| 17N |
| 18N |
| 19N |

In this example there are always 15 BB's covered in the LBR.

Incrementing the BB execution count for each BB detected in the LBR, by 1/15 seen in the LBR path will greatly reduce the effect of shadowing

Many more with 20 Cycles/branch taken

Many more with N samples taken

Many more with 15 N LBR Entries

(intel)

# Branch Filtering

| LBR Filter Bit Name | Bit Description | bit |
|---|---|---|
| CPL_EQ_0 | Exclude ring 0 | 0 |
| CPL_NEQ_0 | Exclude ring3 | 1 |
| JCC | Exclude taken conditional branches | 2 |
| NEAR_REL_CALL | Exclude near relative calls | 3 |
| NEAR_INDIRECT_CALL | Exclude near indirect calls | 4 |
| NEAR_RET | Exclude near returns | 5 |
| NEAR_INDIRECT_JMP | Exclude near unconditional near branches | 6 |
| NEAR_REL_JMP | Exclude near unconditional relative branches | 7 |
| FAR_BRANCH | Exclude far branches | 8 |

# Branch Filtering

- ## User near calls only
  - ### Tracking back from OS critical sections to user function that caused the problem
  - ### Lack of returns may be an issue in some cases
    - #### But not for HPC ☺
  - ### Use static call analysis to clean up chains

- ## User and OS near calls only
  - ### Profiling OS call stacks
  - ### Eliminating leaf functions may be complicated by lack of returns
    - #### Don't remove returns if this is a problem
    - #### Use BTS to capture deeper stack
  - ### Issue: cannot exclude unconditional jumps without excluding calls

(intel®)

# Precise cycles can be constructed from any PEBS event

- **Allow profiling code sections screened with STI/CLI semantics**
  - **Ring 0 OS critical sections**
- **PEBS sampling mechanism may loose interrupts during halted state**
  - **Instruction retirement required to generate performance monitoring interrupts (PMI)**

    **Counts will not occur without PEBS being invoked**

(intel)

# Using cycles to optimize the optimizations

- **Profile the application for cycle usage and uop flow.**
  - **Identify hot functions**
  - **Check asm of FP intensive code for correct instruction mix**
    - **X87 is slower than SSE**
    - **Intel® Compiler has FP-model flags and many pragmas**
- **Vectorize long tripcount loops**
  - **-SSE4.2 uses unaligned loads more aggressively**
    - **Align data whenever possible**
  - **Check loop tripcounts with br events and register values (described later)**
    - **Interchange loop orders to get long loops as inner loop**
      - **Change multi dimensional array layout as needed**
    - **Completely unroll short tripcount (<~7) inner loops**
    - **Split/merge loops depending on code size**
    - **Predicate hoist constant condition if's out of loops**
    - **Etc, etc , etc...I could write a book**

(intel®)

# Using cycles to optimize the optimizations

- **C++ and large binaries: Only optimize what uses cycles**
  - **Use call counts to drive compiler inlining**
    - **Compiler needs to evaluate a large enough scope to do its best work**
    - **Particularly functions/methods invoked inside loops**
  - **Size vs Speed**
    - **Extremely large binaries need to minimize size**
      - **-Os (linux) –O1 (windows)**
  - **Conditional Branch Mispredictions**
    - **HW prediction is shockingly good**
      - **Cost is unretired uop flow (uops_issued.any – uops_retired.slots)**
      - **Optimize case statement order, lowers uops_retired**
- **Use Intel Compiler LIBM,MKL, tbbmalloc, tbbmalloc_proxy**
  - **Intel linker with LD_PRELOAD env variable**
  - **-L/path/to/intel/libs –limf etc**
  - **http://software.intel.com/en-us/articles/optimizing-without-breaking-a-sweat/**

# Thoughts on optimizing large OOP code bases

- **Classic OOP will result in code bases of small functions integrated together to invoke the algorithm**

- **Signatures**
  - **Low instruction_retired/call_retired**
  - **High call_retired/branch_retired**
  - **High indirect_call/call_retired**
  - **High uops_issued.core_stall_cycles – resource_stalls.any**
  - **High $\sum$latency(source)*ifetch_miss(source)**

(intel®)

# How big are the CERN programs



Cacheline access frequency evaluated by sorting cachelines by their accesses

Thus a binary working set size measurement

# Optimizing large Object Oriented Code

- **Inlining is the advice of choice but things are more complicated.**
- **Inlining increases binary size and can make ifetch misses more costly and code slows down**
  - **Even if fewer in overall number**
- **Ifetch miss events have among the largest IP skids of all events**
  - **They can show up in the wrong function**
- **Large codes built of many small methods can result in flat cycle profiles**
  - **It can take thousands of functions to account for 80% of the clock cycle samples**
  - **Thus thousands of functions must be optimized to achieve a significant performance improvement**

# Optimizing large Object Oriented Code

- **The author knows of no proven methodology to correct the cost of excessive taken branches and the resulting flat cycle profile.**
  - **Need fewer calls,**
    - **instructions required for calling conventions**
  - **Larger functions to allow the compiler to see the whole calculation and do a better job**
  - **Larger shared objects to allow greater effect from IPO**
    - **Create shared objects using just the hot methods to avoid excessive inlining**

- **This has to be applied to enough methods to account for 80->95% of the cycles**

  **Mostly this is about reducing the total instruction count**

# Thoughts on optimizing large OOP code bases

- **Function calls result in added instructions**
  - **Call and return**
  - **Trampolines required for position independent code/ shared object cross invocations**
    - **Indirect branches can be more costly**
  - **Freeing & restoring registers for local use**
    - **Mostly an ia32 issue**
  - **Setting and reading function arguments**
    - **Larger on ia32 due to required use of stack**
- **Virtual function calls (function pointers) increase indirect call instructions and associated pointer loads**

(intel)

# Thoughts on optimizing large OOP code bases

- **Does a call graph help?**
  - **Unlikely**
  - **Provides the direct path back to main**
    - **Usually sampled on time**
    - **Does not provide call counts in most cases**
  - **Does not identify clusters of active (excessive) call activity**

(intel®)

# Thoughts on optimizing large OOP code bases

- **A modest proposal:**
- **Use LBRs and static analysis to evaluate frequency and cost of function calls**
  - **the call count**
  - **count taken branches between call and arrival in function**
  - **Get count of indirect branches invoked**
  - **Add cost for function arguments**
  - **Add a cost for push/pop of registers**

(intel®)

# Thoughts on optimizing large OOP code bases

- **A modest proposal:**
- **Use social network analysis/network theory to identify clusters of active, costly function call activity**
  - **Web search on Social networking/social networking analysis**
- **Order clusters by total time and/or total "cost"**
  - **Split time of functions shared between clusters by call counts**
  - **Calls have a direction**
    - **Utility functions must not be viewed as bridges**

(intel)

# Thoughts on optimizing large OOP code bases

- **A modest proposal:**

- **Manually reduce function count in hot clusters by explicit code inlining**

  - **Prioritize work by call overhead cost to be gained**

  - **Duplicate code as needed**

  - **Reduce cross shared object call counts**

(intel)

# Using cycles to optimize the optimizations

- **PEBS near call event + LBRs to get call counts/source**
  - **Selecting source files to compile with enhanced inlining**
    - **IPO can be enahnced when used with PGO**

- **PEBS near call event + registers (em64T) to get function arguments**
  - **Fix memset/memcpy calls with short lengths**
  - **Excessive calls to malloc/free due to constructor/destructor?**
    - **Identify small malloc's/free's**
    - **Let the compiler allocate small structures statically rather than malloc and free them excessively**

# Using cycles to optimize the optimizations

- **Optimize only functions that use significant cycles**
  - **Reduces build time**
  - **Minimize fighting the compiler**
    - **Changing optimizations or compilers in large builds can be problematic**

- **Move gcc/icc and create script called gcc/icc**

```
#!/bin/sh
if echo $@ | grep -f /tmp/sourcefilelist.txt > /dev/null ;
  then /opt/intel/Compiler/11.0/083/bin/intel64/icc.ori –g -fast $@;
  else gcc.ori -g -O2 $@;
fi
```

# Using cycles to optimize the optimizations

- **PTU sometimes shows *.h files as source**

- **Generate a list of c/cpp files as follows:**
  - **Export list of functions from Intel® PTU**
  - **Create script grepf.sh to grep for defined symbols:**
    **#!/bin/sh**
    **if nm --defined-only --demangle $1 | grep -f $2 > /dev/null ; then echo `basename $1 .o`.cpp; fi**
  - **Find hot object files and remember cpp files:**
    **find -name "*.o" -exec grepf.sh '{}' /tmp/functionlist.txt \; > /tmp/sourcefilelist.txt**

- **This will produce sourcefilelist that only includes targets of compiler**

# Event Classes

1. Execution flow events

   – Cycles, Branches, stalls, uops/inst_retired

2. **Penalty events
   Change code to remove the penalty**

   – **Ex: load requiring access to dram**

3. Resource saturation events

   – Bandwidth, load/store buffers, dispatch ports

   – No well defined cost

4. Architectural characterization

   – Cache accesses, MESI states, snoops

5. Instruction mix

(intel)

# Memory Access

- **Load instruction uses virtual address to access memory space**
- **HW translates that to physical address to access caches**
  - **DTLB does this**
- **Access is hierarchical**
  - **Check L1D first**
  - **If (miss) check if Line Fill Buffer (LFB) allocated**
  - **If(LFB miss) allocate LFB, escalate miss to L2**
  - **If(miss L2) get Super Queue (SQ) slot, escalate to uncore**

# Memory Access Penalties

- **Load misses cause execution stalls**
  - **In most cases store misses will not stall execution**
    - **Data to be stored is held in store buffer until desired line is in L1d, thus execution continues**

- **Loads that hit LFBs overlap in time with original line request**
  - **If the original request was a load, the original miss accounts for the entire penalty**
  - **If there are multiple load request to the LFB the least costly would be the penalty**
  - **Not all load misses are equally costly**

# Stall Decomposition on Intel® Core™ i7 Processors

- **Same basic methodology as on Intel® Core™2 processors***

- **Basic strategy is to identify the largest penalty event contributions first**
  - **Work your way down to smaller contributors**

- **FE starvation can now be measured**
  - **And no branch misprediction flush penalty**

- **Only both_threads_stalled can be measured at execution**

  - **SMT will make $\Sigma\text{events}_i\text{*penalties}_i >$ both_thread_stalled**

  - **ALU_only stalls can be measured per thread**
    - **Ports 0,1 and 5**

* Intel, the Intel logo, Intel Core and Core Inside are trademarks of Intel Corporation in the U.S. and other countries.

# Stall Decomposition: $\Sigma events_i * penalties_i$
## The Elephants

- **LLC, L2, and DTLB misses are the large penalty, common events**
- **LLC activity must be measured at L2 for it to have core, PID, TID context**
    - **Uncore has no ability to track core, PID or ThreadID**
    - **Uncore event collection not yet supported**
- **Figure of merit: Events*Penalty/cycles**
    - **Samples_ev*SAV(ev)*Penalty(ev)/ ( Samples_cyc*SAV(cyc) )**
    - **If SAV(ev) = SAV(cyc)/Penalty(ev)**
    - **FOM = Samples_ev/Samples_cyc**
    - **This is ~ how the default SAVs are set**
    - **Minimizes required screen area in the data display**

# Stall Decomposition: $\Sigma events_i * penalties_i$ The Elephants

- **Figure of merit: Events*Penalty/cycles**
  - **Overcounts when there are temporally overlapping penalties**
  - **Compilers can hoist loads. So make sure there are stalls as well**
    - **PEBS event uops_retired.stall_cycles should pile up very close to instructions suffering large penalties**
  - **The combination provides the answer to the critical question:**

**Is the fix worth the effort?**

(intel)

# Penalty Events: Memory Access

- **Intel® Core™ i7 processor memory access events are "per source"**
  - **How many times cacheline came from "here"**

- **Unique sources have unique Penalties**
  - **DP system has ~10 sources outside a core**
  - **Large number of performance events**

- **Memory access events are precise**
  - **HW captures IP and register values**
  - **Sample + Disassembly => Reconstruct Address**

- **Latency Event captures IP, load latency, data source and address**
  - **Similar to Itanium® Processor Family* Data Ear**

\* Itanium is a trademark or registered trademark of Intel Corporation or its subsidiaries in the United States and other countries.

# Offcore Response Latencies

- **LLC Hit that does not need snooping**
  - **LLC latency ~ 35-40 cycles**
- **LLC Hit requiring snoop, clean response ~65**
- **LLC Hit requiring snoop, dirty response ~75**
- **LLC Miss from remote LLC ~ 200 cycles**
- **LLC Miss from local Dram ~60 ns**
- **LLC Miss from remote Dram ~100 ns**

Note: All latencies and memory access penalties shown are merely illustrative. Actual latencies will depend on (among other things) processor model, core and uncore frequencies, type, number and positioning of DIMMS, platform model, bios version and settings. Consult the platform manufacturer for optimal setting for any individual system. Then measure the actual properties of that system by running well established benchmarks.

(intel)

# Memory Access PEBS Events

## Identify LLC and DTLB load miss

– **Precise load events do not include DCU prefetch/ L2 prefetch**

| Name | Penalty | Umask | Umask_name |
|---|---|---|---|
| mem_load_retired | 0 | 0x1 | L1D_HIT |
| | 6 | 0x2 | L2_HIT |
| | ~35 | 0x4 | LLC_HIT_UNSHARED* |
| | ~75 | 0x8 | OTHER_CORE_L2_HIT_HITM* |
| | depends | 0x10 | LLC_MISS |
| | depends | 0x40 | HIT_LFB |
| | | 0x80 | DTLB_MISS* |

**LLC_HIT_UNSHARED should be LLC_HIT_NO_SNOOP**
**OTHER_CORE_L2_HIT_HITM should be LLC_HIT_SNOOP**
**DTLB_MISS counts primary and secondary DTLB misses on CoreI7**
**Only counts primary on Xeon™ 5600 Family Processors**
**Penalty for DTLB miss is not a constant**
**Also use Dtlb_load_misses.walk_cycles on Xeon™ 5600 Family Processors**

Note: All latencies and memory access penalties shown are merely illustrative. Actual latencies will depend on (among other things) processor model, core and uncore frequencies, type, number and positioning of DIMMS, platform model, bios version and settings. Consult the platform manufacturer for optimal setting for any individual system. Then measure the actual properties of that system by running well established benchmarks.

(intel)

# Precise Uncore Response Xeon™ 5500 Family Processors

- **Load response from LLC, another core, local DRAM, remote socket, remote DRAM and IO**

| Name | Penalty | Umask | Umask_name |
| --- | --- | --- | --- |
| mem_uncore_retired | ~85 | 0x2 | OTHER_CORE_L2_HITM |
| | ~185 | 0x8 | REMOTE_CACHE_ LOCAL_HOME_HIT |
| | ~200 | 0x20 | LOCAL_DRAM |
| | ~350 | 0x40 | REMOTE_DRAM |
| | | 0x80 | IO |

**OTHER_CORE_L2_HITM should be LOCAL_HITM**

Note: All latencies and memory access penalties shown are merely illustrative. Actual latencies will depend on (among other things) processor model, core and uncore frequencies, type, number and positioning of DIMMS, platform model, bios version and settings. Consult the platform manufacturer for optimal setting for any individual system. Then measure the actual properties of that system by running well established benchmarks.

# Precise Uncore Response
# Xeon™ 5600 Family Processors

- **Load response from LLC, another core, local DRAM, remote socket, remote DRAM and IO**

| Name | Penalty | Umask | Umask_name |
|---|---|---|---|
| mem_uncore_retired | ~85 | 0x2 | LOCAL_HITM |
| | ~375 | 0x4 | REMOTE_HITM |
| | ~220 | 0x8 | LOCAL_DRAM_AND_ REMOTE_CACHE_HIT |
| | ~375 | 0x10 | REMOTE_DRAM |
| | | 0x80 | UNCACHEABLE |

Note: All latencies and memory access penalties shown are merely illustrative. Actual latencies will depend on (among other things) processor model, core and uncore frequencies, type, number and positioning of DIMMS, platform model, bios version and settings. Consult the platform manufacturer for optimal setting for any individual system. Then measure the actual properties of that system by running well established benchmarks.

(intel®)

# Precise Store DTLB miss

| Name | Event | Umask | Umask_name |
|------|-------|-------|------------|
| mem_store_retired | 0x0c | 0x1 | DTLB_MISS* |
|  |  | 0x2 | dropped events |

**DTLB_MISS counts primary and secondary DTLB misses on Corel7**
**Only counts primary on Xeon™ 5600 Family Processors**

# Overlapping Memory access penalties Xeon 5600 family: Offcore_request_outstanding

| Event Name | umask | cmask, inv |
|---|---|---|
| OFFCORE_REQUESTS_OUTSTANDING.ANY.READ | 0x8 | |
| OFFCORE_REQUESTS_OUTSTANDING.ANY.READ_NOT_EMPTY | 0x8 | 1,0 |
| OFFCORE_REQUESTS_OUTSTANDING.DEMAND.READ_CODE | 0x2 | |
| OFFCORE_REQUESTS_OUTSTANDING.DEMAND.READ_CODE_NOT_EMPTY | 0x2 | 1,0 |
| OFFCORE_REQUESTS_OUTSTANDING.DEMAND.READ_DATA | 0x1 | |
| OFFCORE_REQUESTS_OUTSTANDING.DEMAND.READ_DATA_NOT_EMPTY | 0x1 | 1,0 |
| OFFCORE_REQUESTS_OUTSTANDING.DEMAND.RFO | 0x4 | |
| OFFCORE_REQUESTS_OUTSTANDING.DEMAND.RFO_NOT_EMPTY | 0x4 | 1,0 |

Offcore_requests_outstanding.demand.read_data_not_empty = cycles there is at least one request from L1d that had to be satisfied by escalation to uncore
Includes L1d HW prefetch, loads and SW_prefetch
**Defines upper limit of memory access penalties due to L2 miss**

(intel)

# So what do you do?

- Load driven misses resulting in pipeline stalls can be fixed by
  - Use longest tripcount loop to drive strategy
  - Change loop order/data layout to give HW prefetcher a chance
    - Divide large structures by usage (See MILC)
    - Structures of arrays rather than arrays of structures
  - Make sure buffer initialization is consistent with usage
    - Make remote_dram misses local dram misses & cut latency in half

- DTLB misses: use large pages

(intel®)

# So what do you do?

- Load driven misses resulting in pipeline stalls can be fixed by

- SW prefetch _mm_prefetch(addr, hint)
  <ia32intrin.h>

  - Use LOAD_HIT_PRE to identify when prefetch distance is too small

    - Min prefetch dist (iter) ~ 200/(uops_per_iteration/3)

      - For local dram
      - Will change as latency changes

  - long inner loop-> prefetch ahead in inner loop

  - Short inner loop-> prefetch 1,2 iterations ahead on outer

  - Reused linked list -> create indirect address array

  - #pragma openmp for (guided)
                    will cause havoc

  - Volume 2 of that book

  - SW prefetches will not help a BW limited application

# Other Penalties

- **Divides and SQRT (Arith.Cycles_div_active)**
  - **Vectorize**
  - **Save reciprocals that are reused**
  - **Merge with bandwidth limited loops**
- **Store Forwarding (Load_Block.overlap_store)**
  - **Event only on Xeon™ 5600**
  - **Use Intel Compiler**
  - **Be careful with data type sizes (keep consistent)**
- **FP exceptions (uops_decoded.ms)**
  - **Use Intel compiler (no x87, FTZ)**
  - **Uninitialized values in simd registers**
- **No ability to measure stalls associated with chained long latency instructions**
  - **Sum = a+b+c+d+e…evaluated left to right**

# Instruction Starvation

- Lots of calls to small functions can lead to starving the pipeline of instructions
  - Only L2 prefetchers prefetch instructions
- Uops_issued.core_stall_cycles – resource_stalls.any = cycles BE wants instructions, but does not get them
  - This is more accurate with HT off
- Can be cross checked on Xeon™ 5600 processor with offcore_requests_outstanding.demand.read _code_not_empty (for L2 miss)

(intel)

# Decomposing instruction starvation

| Event | Penalty |
|---|---|
| l2_rqsts.ifetch_hit | ~6 |
| offcore_response_0.demand_ifetch.local_cache | ~35 |
| offcore_response_0.demand_ifetch.local_dram | ~200 |
| offcore_response_0.demand_ifetch.remote_dram | ~350 |

**Ifetch miss events have among the largest IP skids of all performance events. The IP can easily have been on in a previously executing function at the time the ifetch miss occurred. See slide 23**
**Uncertainties are also larger, due to the many buffers in the pipeline**
**Instruction starvation does not occur unless the buffers drain**

Note: All latencies and memory access penalties shown are merely illustrative. Actual latencies will depend on (among other things) processor model, core and uncore frequencies, type, number and positioning of DIMMS, platform model, bios version and settings. Consult the platform manufacturer for optimal setting for any individual system. Then measure the actual properties of that system by running well established benchmarks.

(intel)

# Instruction Access Penalties

- **Demand Ifetch: offcore_response.demand_ifetch.\***
  - **Usually associated with function calls followed by taken branches in LARGE binaries**
  - **IPO, force inlining**
  - **PGO to reduce taken branches**
  - **shrink sizes of other functions**
  - **Change order of link command**
  - **Offcore_response.demand_ifetch.local_dram**
    - **Sw_prefetch(&foo(),1);   ?????**
  - **Offcore_response.demand_ifetch.remote_dram**
    - **Run 1 copy of binary per socket**
      - **Must have two complete copies on the disk**
  - **Offcore_response.demand_ifetch.llc_hit_no_other_core**
    - **Sw prefetch?, PGO, IPO**
- **ITLB misses: use large Itlb pages**

(intel®)

# Reducing calls and *.so

- Use linker and a control list to identify internal and external functions in *.so to reduce the use of trampolines
  - icpc -Wl,-z,defs -L/External -L/Linker -Wl,-version-script,export.tmp

```
$ cat export.tmp
{
  global:
        _Foo1;
        _Foo2;
  local:
        _Bar1;
        _Bar2;
   };
```

# Reducing calls and *.so

- Identifying the internal functions is not simple

- Use LBRs, and sfdump5 (see backup) to identify call chains between *.so

- Merge source files into fewer *.so
  - This will improve effectiveness of PGO/IPO

- Use global/local file of previous slide to reduce trampolines

NOTE: Author has never personally done this, so he does not know if it really works, or if the syntax is really correct.

(intel)

# Event Classes

1.  Execution flow events
    - Cycles, Branches, stalls, uops/inst_retired
2.  Penalty events
    - Ex: load requiring access to dram
3.  **Resource saturation events**
    - **Bandwidth, ld/st buffers, dispatch ports**
    - **No well defined cost**
4.  Architectural characterization
    - Cache accesses, MESI states, snoops
5.  Instruction mix

# Resource Limitation Events

- **Resource limitation is usually only a problem when the resource is saturated**
  - There is ~no cost* for bandwidth until the bandwidth is close to saturated
    - *Latency depends weakly on BW on Corei7
- **Lost cycles due to resource saturation can be hard to measure**
- **Only way to determine bandwidth limit is to measure it**
  - Count cachelines transferred/cycle for triad
    - (w/wo SSE NT stores)
  - Depends on the number of triad threads
- **Resource saturation results in no gain from HT**

# Resource Limitation: Memory Bandwidth

- **Usually needs HW (or SW) prefetch**
  - **Load latencies will restrict execution otherwise**
    - **Exception: for(i=0;i<len;i++)a[i] = b[addr[i]];**

- **Limit depends on**
  - **number and location of concurrent threads consuming large numbers of lines**
    - **For asynchronous execution this becomes ~impossible to know**
  - **core and uncore frequencies**
  - **type, number, size, location of dimms**
  - **bios version and settings**
  - **Motherboard**

- **Measured in cycles/cacheline transferred**
  - **Triad with/wo RFO result in ~ same limit!**
  - **All "BW" events discussed here count cachelines transferred**

(intel)

# Triad bandwidth vs thread count



Socket 0 only

Socket 0 full, adding threads on socket 1

Open MP triad on machine with unmatched dimms

## NUMA = ON

Note: All latencies and memory access penalties shown are merely illustrative. Actual latencies will depend on (among other things) processor model, core and uncore frequencies, type, number and positioning of DIMMS, platform model, bios version and settings. Consult the platform manufacturer for optimal setting for any individual system. Then measure the actual properties of that system by running well established benchmarks.

(intel)

# Latency stalls vs Bandwith saturation

- A latency stalled program has a small number of outstanding data cachelines in flight simultaneously

```
i=0;
While(mystruc->next !=0){
    mystruc=mystruc->next;
    a[i] = mystruc->b_val;
    i++;
}
```

Only one (possibly 2) loads in flight at a time

- Clearly a triad with prefetchers enabled in BW limited

# Gather, OOO execution and Bandwidth saturation

**Consider:**
   **For(i=0;i<len;i++)A[i] = B[ADDR[i]];**

**A data collection might show something like  1000 cycle samples, 200 instruction retired samples and 5000 mem_uncore_retired.local_dram samples**

**The mem_uncore SAV is 10K, the cycle SAV is 2 million**
**This absorbs the 200 cycle penalty..so the ratio of the samples is the ratio of the cycles…**

**Clearly, there are more cycles in dram access than cycles executed.**

# Gather, OOO execution and Bandwidth saturation

In a gather loop the RS acts as a prefetcher.
There are 6 uops/iteration  -> ~5 iterations in the RS?
except the loads go out immediately..
there is no dependency so the 2 loads can be executed,
the incr, cmp and branch can execute, again as there are no dependencies
so only the stores pile up
This would suggest ~30 iterations in flight at a time

the number of load buffers might be what blocks FE uop issue
there are 48 and 2/iteration are needed

The loads of ADDR[i] are sequential and thus HW prefetched.
All the stalls are on the load of B[ADDR[I]]
Thus the events fall on the next instruction.

The mem_uncore_retired.local_dram events are all overlapped..
Thus events*penalties overcounts by a huge factor

# Bandwidth per core

- **Much more complicated than on Intel® Core™2 processors**
  - **Bandwidth limit depends on number of threads using maximum BW and core position of those threads**
    - **CAN ONLY BE MEASURED**
  - **No single event counts total cachelines in+out to memory /core**
    - **Cacheable writebacks are written to LLC and written to memory at a later time**
    - **Offcore_response.data_ifetch.all_dram**
      - **However, WB ->dram makes no sense**
    - **Local vs remote memory**
    - **NT SSE Stored cachelines are problematic**

# Offcore_Response: Breaking Down Off-core Memory Access

- **Matrix type event**
  - **Request type X Response type**
    - **65025 possible real combinations  (65535 – 2 X 255)**
  - **Request and Response programmed in MSRs**
  - **OR(Request bits true) .AND.  OR(Response bits true)**
  - **Ex: all LLC misses = set bits 0,1,2,3,4,5,6,11,12,13,14**
    - **787F**

- **Solves problem of averaging over widely differing penalties**

- **Only one version of the event (b7/msr 1a6)**
  - **offcore_response_0**

(intel®)

# Memory Access: Off-core Access

- Offcore_Response_0
  - "umasks" set with MSRs 1a6
  - Two versions on XEON 5600 processor family
    - Programming a little different

| | Bit position | Description |
|---|---|---|
| **Request** | **0** | **Demand Data Rd = DCU reads (includes partials, DCU Prefetch)** |
| **Type** | **1** | **Demand RFO = DCU RFOs** |
| | **2** | **Demand Ifetch = IFU Fetches** |
| | **3** | **Writeback = MLC_EVICT/DCUWB** |
| | **4** | **PF Data Rd = MPL Reads** |
| | **5** | **PF RFO = MPL RFOs** |
| | **6** | **PF Ifetch = MPL Fetches** |
| | **7** | **OTHER** |
| **Response** | **8** | **LLC_HIT_UNCORE_HIT** |
| **Type** | **9** | **LLC_HIT_OTHER_CORE_HIT_SNP** |
| | **10** | **LLC_HIT_OTHER_CORE_HITM** |
| | **11** | **LLC_MISS_REMOTE_HIT_SCRUB** |
| | **12** | **LLC_MISS_REMOTE_FWD** |
| | **13** | **LLC_MISS_REMOTE_DRAM** |
| | **14** | **LLC_MISS_LOCAL_DRAM** |
| | **15** | **IO_CSR_MMIO** |

# Offcore_response Reasonable Combinations

| Request Type | MSR Encoding |
|---|---|
| ANY_DATA | xx11 |
| ANY_IFETCH | xx44 |
| ANY_REQUEST | xxFF |
| ANY_RFO | xx22 |
| COREWB | xx08 |
| DATA_IFETCH | xx77 |
| **DATA_IN** | **xx33** |
| DEMAND_DATA | xx03 |
| DEMAND_DATA_RD | xx01 |
| DEMAND_IFETCH | xx04 |
| DEMAND_RFO | xx02 |
| OTHER | xx80 |
| PF_DATA | xx30 |
| PF_DATA_RD | xx10 |
| PF_IFETCH | xx40 |
| PF_RFO | xx20 |
| PREFETCH | xx70 |

| Response Type | MSR Encoding |
|---|---|
| ANY_CACHE_DRAM | 7Fxx |
| ANY_DRAM | 60xx |
| ANY_LLC_MISS | F8xx |
| ANY_LOCATION | FFxx |
| IO_CSR_MMIO | 80xx |
| LLC_HIT_NO_OTHER_CORE | 01xx |
| LLC_OTHER_CORE_HIT | 02xx |
| LLC_OTHER_CORE_HITM | 04xx |
| LCOAL_CACHE | 07xx |
| LOCAL_CACHE_DRAM | 47xx |
| LOCAL_DRAM | 40xx |
| REMOTE_CACHE | 18xx |
| REMOTE_CACHE_DRAM | 38xx |
| REMOTE_CACHE_HIT | 10xx |
| REMOTE_CACHE_HITM | 08xx |
| REMOTE_DRAM | 20xx |

**NT local stores counted by 0200 not 4000**

A bit different on Xeon 5600 Processor Family

# Total Memory Bandwidth

- **Delivered + Speculative Traffic to local memory**

  – **Reads and Writes Per Source**
    - UNC_QHL_REQUESTS.IOH_READS
    - UNC_QHL_REQUESTS.IOH_WRITES
    - UNC_QHL_REQUESTS.REMOTE_READS (includes RFO and NT store)
    - UNC_QHL_REQUESTS.REMOTE_WRITES (includes NT Stores)
    - UNC_QHL_REQUESTS.LOCAL_READS (includes RFO and NT Store)
    - UNC_QHL_REQUESTS.LOCAL_WRITES (no NT stores)

- **Precise totals can be measured in IMC**
  – **But cannot be broken down per source**
    - UNC_IMC_NORMAL_READS.ANY  (or by channel, includes RFO)
    - UNC_IMC_WRITES.FULL.ANY (or by channel, includes NT stores)

# A few particularly useful events for measuring BW

- Offcore_response.data_in.local_dram
  - Read BW (per core) from local dram
- Offcore_response.data_in.remote_dram
  - Read BW (per core) from remote dram
    - Indicates NUMA locality problem
- Uncore events get totals but only in counting mode with no data/core
  - Unc_imc_normal_reads.any
    - Total read cachelines from this mem controller
  - Unc_imc_writes.full.any
    - Total written cachelines to this mem controller

(intel®)

# Latency vs Bandwidth

- On Xeon™ 5600 processors the average occupancy of the super queue can be evaluated as offcore_requests_outstanding.any.reads/ cpu_clk_unhalted.thread

- If this is large then the loop is likely BW limited

- If it is small and the event counts indicate a memory access problem due to loads then it is likely to be a latency issue

(intel)

# Triad bandwidth vs thread count



Socket 0 only

Socket 0 full,
adding threads
on socket 1

Open MP triad on machine with matched dimms

## NUMA = ON

Note: All latencies and memory access penalties shown are merely illustrative. Actual latencies will depend on (among other things) processor model, core and uncore frequencies, type, number and positioning of DIMMS, platform model, bios version and settings. Consult the platform manufacturer for optimal setting for any individual system. Then measure the actual properties of that system by running well established benchmarks.

(intel)

# Average super queue occupancy



Socket 0 only

socket0
socket1

Socket 0 full,
adding threads
on socket 1

**Variation per core is ~ 20%**
**Many temporally overlapping requests identifies**
**bandwidth saturation**

Note: All latencies and memory access penalties shown are merely illustrative. Actual latencies will depend on (among other things) processor model, core and uncore frequencies, type, number and positioning of DIMMS, platform model, bios version and settings. Consult the platform manufacturer for optimal setting for any individual system. Then measure the actual properties of that system by running well established benchmarks.

# Average super queue occupancy



Evaluated with no knowledge of thread count

Note: All latencies and memory access penalties shown are merely illustrative. Actual latencies will depend on (among other things) processor model, core and uncore frequencies, type, number and positioning of DIMMS, platform model, bios version and settings. Consult the platform manufacturer for optimal setting for any individual system. Then measure the actual properties of that system by running well established benchmarks.

# Identifying bandwidth saturation

- Identifying BW saturation by measuring bytes/time is complicated by the BW limit changing with the number of threads consuming BW (slide 90)

- Non concurrent execution, with some threads consuming large BW, while others consume little, can make identifying saturation extremely difficult

# Identifying bandwidth saturation

- Average SQ occupancy limit varies less with thread count/concurrency

- It does not distinguish between LLC hits and LLC misses

- Recipe:
  - Identify problematic functions with <SQ occup>
  - Use offcore_response events to determine the fraction associated with LLC hits vs misses

# But what is the potential gain?

- **<u>None of this measures what is needed!</u>**
  - **It does not tell us if the fix is worth the effort!**
- **The fix is to reduce the number of lines transferred**
  - **Consume more data per line transferred**
- **Gain**
  - **BW_time = total_lines/BW_limit**
  - **Exec_time = time to execute instructions**
    - **Memory latency of ~0**
  - **Time = MAX(BW_time, Exec_time)**
  - **Completely BW limited ~ change_in_total_lines/BW_limit**

Problem: cannot measure exec time,
BW limit is absurdly complex in general
(must assume synchronous execution)

(intel)

# An example

**Double \*a, \*b;**
**For(i=0; i<len; i+=8)a[i] = sqrt(b[i]);**

**We might be able to compress a and b to transfer fewer lines**

**Double \*ap, \*bp;**
**For(i=0; i<len/8; i++)ap[i] = sqrt(bp[i]);**

**But would it actually go any faster?**

**No, The SQRT latency ~ matches the BW limit**

# Estimating the gain

- Exec time ~ uops_retired.slots/`*3'*+ arith.cycles_div_active
  - Undercounts cycles associated with chained long latency uops
- Optimized BW time = Adjusted_lines/Max_bw
- Gain ~ Cpu_clk_unhalted.thread – MAX(Optimized BW time, Exec Time)
- Many Uncertainties, but better than nothing
  - Assumptions about concurrency of high BW usage
  - Assumptions about cycles associated with chained long latency uops
  - Is uops/3 realistic?

(intel)

# What do you do about Bandwidth?

- Data layout change is usually best
  - Fix buffer initialization to make remote_dram small
  - Fix order of structure elements (big to small)
  - Eliminate unused structure elements
    - Divide structures into parallel structures by use
  - Measure data consumed/cacheline in
    - Sum load/store in loops (ignore stack pointer, +=)
    - Multiply by total tripcount & divide by 64*offcore_response.data_in.local_dram
  - Fix nested loop order

- Measure data_in with prefetchers on & off
  - If difference is large
    - Change data layout to help HW prefetcher or
    - Consider sw prefetching everything and disabling HW prefetchers

# OOO resource Saturation

- **Load buffer saturation (resource_stalls.ld)**
  - **In HPC, frequently due to bandwidth saturation**
- **Store buffer saturation (resource_stalls.st)**
  - **This will cause stores to stop the pipeline**
  - **Usually associated with stores missing l1d/l2 etc**
    - **SW prefetch, change layout to help HW prefetch**
- **Port saturation (uops_executed.portX/cycles)**
  - **Most common for load port (2)**
    - **Avoid loop distribution (F90)**
    - **Merge loops to reuse data while available**
    - **Align data and vectorize**

(intel)

# Less than ideal multi core scaling

- Perfect scaling results in the number of perf events (summed over cores) being constant

- Difference of event counts can identify locality using cycles and some reasons for non scaling behavior
  - Cacheline access contention can cause non scaling
    - Load-hitm and store address analysis identifies this

- Most non scaling due to resource saturation and evaluated as a ratio: events/wall_cycles
  - Wall_cycles ~ cycles/active cores
          or Cpu_clk_unhalted.thread max(ICPU)
  - **<u>Cannot be seen in difference display</u>**

# Sources/signatures of non scaling

- Turbo
  - Having this on results in large drop from 1->2
- Smaller share of LLC
  - Decrease in LLC hits, increase in LLC miss
- Increase in page faults
  - More threads require more memory
- Asymmetry associated with core 0
  - OS induced imbalance
- Context switching
  - OS's love to move things around, being the boss!
  - Don't know about logical cores & double up on one physical core, while other phys cores are idle

(intel)

# Sources/signatures of non scaling

- Saturating a resource
  - Ex: Bandwidth
  - Code optimization increases resource saturation
- Shared memory application specific
  - Serial execution
  - Overly contested lock access
  - False sharing (non overlapping access to a line)
- NUMA based non scaling
  - Increase in *.remote_dram
- HT can be viewed as a way to recover scaling

(intel)

# More sources of non scaling

- Load imbalance
  - Increase in halted cycles
- MPI global operations
  - increase in time associates with MPI global APIs
    - Ex: allreduce
- Synchronous message passing
  - "Intrinsically" non scaling

(intel)

# Resolving non scaling issues

- Disable turbo while doing measurements
- Disable HT while doing measurements
- Pin all affinities
  - OS's love to move things
  - Old OS's will schedule 2 threads on a physical core while leaving other physical cores idle. This increases with thread count
- Make sure there is enough memory
  - /proc/meminfo->Active (?)
- Do 1 thread baseline on a core other than 0
- Increased LLC miss
  - Usual approaches to fixing these, see previous

# Resolving non scaling issues

- Bandwidth issues
  - Check data decomposition for sepparation
  - Improve data layout to reduce cacheline usage
  - See previous section on BW issues
- Excessive lock contention
  - Use finer grained locking
  - Use faster locking APIs
  - Make sure the global update is really needed
    - Can you continue working with local copy
- False sharing
  - Put 64 bytes between data elements

# Resolving non scaling issues

- NUMA related non scaling
  - Remote dram data access
    - Improve buffer initialization for local access
    - Make multiple copies for each socket
  - Remote dram ifetch access
    - Make two binaries on the disk and affinity pin per socket

- MPI global operations
  - Use openMP within a box to reduce MPI nodes
  - Use good MPI library

(intel)

# Resolving non scaling issues

- Load imbalance
  - Seen as halted cycles
    - TSC difference for successive cpu_clk_unhalted.ref != SAV
  - Work queue approach dynamically restores balance
    - At a cost
      - NUMA locality can be lost
      - SW prefetching can become unpredictable within a thread
  - Estimate work during data decomposition to create balanced work rather than balanced iteration count
  - Save some iterations for final work queue balancing

# Graphical tool needed to organize data viewing

- **Workflow of event based performance analysis is extremely complicated**
  - **Requires an enormous number of features/options to enable all possible tasks**
  - **Automation is very difficult**

- **To do a lot of things requires a lot of options**
  - **Many docking windows, menus, buttons**
  - **Easier to make a tool for a knowledgable user**

- **The data collection is the easy part**
  ## Interpreting the data and determining the correct action is the hard part

(intel)

# Tool Requirements

- Maximize data density

  - Required quantity of data is enormous

- Integrated source/asm display

- Ability to restart sessions later

- Difference utility to monitor changes

- Minimize mouse clicks

- Predefined event lists

- Predefined penalty file

  - Cycle accounting

  - dynamic column layout

# Primary display shows offending events and even call counts

# Set the Granularity to LOOPS

# Get Tuning Advice for the Selected Event/Ratio: Highlighting the Event Row Enables Explanation

# Get Tuning Advice for the Selected Event/Ratio: Highlighting the Event Row Enables Explanation

# Differences of EBS Measurments

- Intel® PTU supports an analysis of differences of experiments
- This requires
  - Event names must be the ~same
  - Load Modules have the same names
    - They can be the same, with data taken on different machines
    - They can be different but built from the same source
      - Allowing differences to be analyzed down to source view
    - They can be completely different (sources and binaries)
      - PTU will compare functions with the same names for modules with the same names

- Identify compiler differences/regressions
- Multi core scaling

**For perfect scaling and identical work,
total event counts, summed over cores,
will be equal**

(intel)

# Data blocked 2X2 unrolled Matrix Multiply compiled at -O2 (Binary = o2\matrix_blk2.exe) Cycle_Usage Profile

# Data blocked 2X2 unrolled Matrix Multiply compiled at -O3 –QxT   (Binary = xt\matrix_blk2.exe) Cycle_Usage Profile

# Only Significant Difference is Cycle Count Create Difference Display

- Control click to select 2 experiments
- Right click to select "Compare Experiments"

# Differences of Samples
# Differences in Cycles Shown in msec to Correct for Comparison of Machines at Different Frequencies



**Scaling Analysis: Sort by Time and see what causes non scaling**

# Drill down by Double Click on Function to Source in difference view

- It is likely to ask where to find the source file

# Same Source can Display Difference per Source Line

# Shift Right click to Highlight a Region and Display Subtotal at the Bottom

# Select "Assembly (1st Exp.)"
# Only Contributing Basic Blocks are Displayed

# Select "Assembly (2nd Exp.)"
## Only Contributing Basic Blocks are Displayed
## Now for BOTH Binaries

# Export Selected Source and the Contributing Basic Blocks from Both Binaries to a Single CSV Spread Sheet
## Instant Compiler Regression Bug Report

# Measuring non parallel execution

- With turbo enabled, non parallel execution will result in a frequency boost to the core executing the serial code

- The serial functions can be identified using the filtering capability of the over time display

# Single threaded execution with turbo boost enabled

# Zoom in on frequency multiplier select range and filter up

# Source View Shows what is Executed



**This is Vectorized**

# Cmp in Blk 15 Controls Loop, Comparing R8 and R11. R8 increments by 48 (30H)

# Register Values Collected with Precise Event Br_inst_retired.all_branches in Blk 11 Yield Values for R11 (14 samples)

# Select the Asm Line, Right Click and Show Register Statistics

# Tripcount is constant (min=max=avg, rms=0) and Equals 786432/48 = 16384 Which is the 4-Dim Lattice size for this Problem

# Source/Asm View Text Search Utility

# Data Address Profiling and False Sharing Detection

**Data Mining in 2 Dimensional Model**

- **Sorting** – **repositioning segments of the axes**
- **Applying granularity** – **changing scale of the axis**
- **Filtering** - **projecting slices onto another dimension**

Data Address

Data HotSpots

Each **Sample** now described by **IP** and **Data address** (plus other characteristics)

IP / Code

Code Hotspots

**Filtering by cachelines marked as "falsely-shared" isolate the causing instructions And the data objects**

(intel)

# Data Address Profiling and False Sharing Detection

**Sampling during app execution** → **Symbolization & Data Address reconstruction** → **Aggregation**

**Precise Event Sampling**: events associated with memory operations, e.g. MEM_INST_RETIRED.LOADS, MEM_INST_RETIRED.STORES…

**Pin threads affinity**

**Iterate over Samples and PEBS records in ebs.tb5**

**Using the binary, identify the instruction that overflowed event counter -> IP-1**

**Sample: IP, data address, threadID..**

**To aggregate addresses into cachelines:**

### Binary

```
         . . .
0x7F5D  mulpd     xmm1,xmm2
0x7F61  movaps    xmm2,XMMWORD PTR [rsp+0230h]
0x7F69  mulpd     xmm7,xmm2
0x7F6D  subpd     xmm3,xmm1
0x7F71  movsd     xmm1,MMWORD PTR [rcx+rbx+rho_i.0+018h]
0x7F7A  movhpd    xmm1,MMWORD PTR [rcx+rbx+rho_i.0+020h]
0x7F83  mulpd     xmm1,xmm8
         . . .
```

IP-1

ebs.tb5

**Sample record:**
**IP, process, module, threadID..**

**PEBS record:**
**IP, rax, rbx, rcx**

**Same cacheline accessed by different threads at different offsets True and False Sharing Next foils Illustrate GUI Navigation**

```
00000000055CC9900
00000000055CC9908
00000000055CC9910
00000000055CC9916
00000000055CC9928
00000000055CC9938
```

&&...FFFFC0

| Cacheline Address / Offset / Thread ... | Contributors | MEM...L1D_MISS |
|---|---|---|
| ▽ 0x00000000055cc900 | Offsets: 5 Threads: 3 | 31 (0.0%) |
| ▽ Offset:0x00(0) | Threads: 1 | 21 (0.0%) |
| ▷ Thread:00003fbb(0014) | Functions: 3 | 21 (0.0%) |
| ▽ Offset:0x38(56) | Threads: 1 | 5 (0.0%) |
| ▷ Thread:00003fbd(0009) | Functions: 3 | 5 (0.0%) |
| ▽ Offset:0x08(8) | Threads: 2 | 1 (0.0%) |
| ▷ Thread:00003fbb(0014) | Functions: 3 | 0 (0.0%) |
| ▷ Thread:00003fbc(0015) | Functions: 1 | 1 (0.0%) |
| ▽ Offset:0x10(16) | Threads: 1 | 3 (0.0%) |
| ▷ Thread:00003fbc(0015) | Functions: 2 | 3 (0.0%) |
| ▽ Offset:0x28(40) | Threads: 1 | 1 (0.0%) |
| ▷ Thread:00003fbc(0015) | Functions: 1 | 1 (0.0%) |

(intel)

# Use Cacheline Access Count to Measure Working Set Size



Performance comparison difference may be due to Cache Size

# NEW – Exact latency / Latency Histogram

- Exact latency in CPU cycles for loads collected with Latency events
- Intel® PTU offers a latency histogram
  - Can be filtered by selected hotspots
  - IP and address spreadsheets, and memory histogram can be filtered by latency region (shown below)

# Array of Structures
# (address-base)% struct_size
# Most structure elements never accessed

# Filtering to a Single Thread Displays the Data Decomposition

# A Different Thread

# Example: False Sharing
# What is it and why is it a Problem

- Cache coherency protocols require that all cores use the most current version of every cacheline
- Shared lines can be modified by any thread
  - Causing lines to be renewed regularly, if any thread writes to any byte in the line
    - (replace an invalid state copy with new valid copy)
  - Line renewal can cause a cache miss by other threads
  - and a 40-300 cycle execution stall
    - Depending on cacheline location
- False sharing is when different threads access non-overlapping regions of a cacheline

**False Sharing Causes Avoidable 40-300 Cycle Stalls
For Every Read Following a Write by Another Thread**

(intel)

# Synthetic Example: Heavy Contention on this Line -- Multiple Threads Accessing Different Offsets Indicate False Sharing (Identified by Rose Highlighting)

# Expanding the "arrow" we see the 2 threads access the line at Different Offsets...This is False Sharing

# Select the falsely shared cacheline (now blue) and Filter the Hotspot view to only Display Accesses to that Line (multiple lines also work)

# Only Events Referencing the Selected Line(s) are now in the Hotspot View Double Click to reach source/ASM view

# The Pointer "sum" is Causing the False Sharing

# NUMA cacheline access

# A NHM Socket is a Caching Agent and a Home Agent



Socket 1

Cores

Uncore

GQ

L L C

Caching Agent

Q PI

Home Agent

IMC          QHL

Socket 2

Cores

Uncore

GQ

Caching Agent

L L C

Q P I

Home Agent

QHL          IMC

(intel)

# Simple Data Read

# RdData request after LLC Miss to Local Home (Clean Rsp)



**DRd**

Cores

Cores

Uncore

**Cache Lookup**

GQ

**SnpData**

L L C

**Cache Miss**

**[Send Snoop to LLC]**

**RspI**

Q P I

**SnpData**

**RspI**

Uncore

**[ Broadcast snoops to all other caching agents) ]**

GQ

**Cache Lookup**

**Cache Miss**

**SnpData**

**RspI**

Q P I

**RdData**

**DataC_E_CMP**

**Allocate in E state [I→ E]**

L L C

**[ Sending Req to Local Home (socket2 owns this address) ]**

**[Fill complete to Socket2]**

**Speculative memory Rd**

QMC

QHL

QHL

QMC

**Data**

Socket 1

Socket 2

(intel)

# RdData request after LLC Miss to Local Home (Hitm Response)

# Uncore Opcode Match events

- **Match address, opcode using an MSR**
  - **37 bit address match**
  - **8 bit opcode match**

| Event | Event code | Umask |
|-------|------------|-------|
| **UNC_ADDR_OPCODE_MATCH.IOH_REQUEST_TRACKER** | **35** | **01** |
| **UNC_ADDR_OPCODE_MATCH.REMOTE_CORES_REQUEST_TRACKER** | **35** | **02** |
| **UNC_ADDR_OPCODE_MATCH.LOCAL_CORES_REQUEST_TRACKER** | **35** | **04** |

- **Local Home data read, remote LLC hit**
  - **Ev=35, umask = 2, opcode = RspFwdS = 0001 1010, opcode only**
- **Local Home data read, remote LLC hitm**
  - **Ev=35, umask = 2, opcode = RspIWb = 0001 1101, opcode only**
- **RFO and perhaps other cases also (E->E problematic)**

# Summary

- Event based sampling performance analysis is extremely powerful on Intel® Core™ i7, XEON™ 5500 and 5600 Processor Families

- Correct methodology is essential

- Correct usage of events is essential

- Intel® PTU simplifies task

# backup

# Low level utilities

- PTU low level utilities can be invoked from the command line by adding the PTU bin directory to the path

- Low level PMU collector is SEP
  - Invoked by vtsarun
  - Data is stored in file called tbsXXXYYYY.tb5
  - sep –start –ex 16 –ec "CPU_CLK_UNHALTED.THREAD:sa=2000000,UOPS_RETIRED.ANY,UOPS_RETIRED.STALL_CYCLES" –app ./myapp –args " arg1 arg2"
    - :sa=VAL explicitly sets SAV value for the event preceding it
    - -ex 16 causes sep to add PEBS buffer to event record
      - Selecting data profile does the same thing

(intel)

# Low level utilities

- sep –start –ex 16 –ec "CPU_CLK_UNHALTED.THREAD:sa=2000000,UOPS_RETIRED.ANY,UOPS_RETIRED.STALL_CYCLES,BR_INST_RETIRED.NEAR_CALL:lbr=2" –app ./myapp –args " arg1 arg2"
  - Event names must be upper case
  - :lbr=VAL turns on LBR capture with filter value determined by VAL
    - Filter values can be determined with profile editor and show command button

| LBR Value | Filter Result |
|-----------|---------------|
| 1 | All Branches |
| 2 | All Calls |
| 3 | User Calls |
| 4 | All Calls & Ret |
| 5 | User Calls & Ret |

# Low level utilities

- sfdump5 creates test output based on data in tb5 file

- sfdump5 tbsXXXZZZ.tb5 –modules > modules.txt
  - Summary of data
    - Total number of samples and events=samples*SAV
      - Events ordered by "event number"
    - Total number of samples/module/event_type

(intel)

# Example sfdump5 output

```
Event Summary
CPU_CLK_UNHALTED.THREAD
   2396           = Samples collected due to this event
   2000000        = Sample after value used during collection
   4792000000     = Total events (samples*SAV)
INST_RETIRED.ANY
   1327           = Samples collected due to this event
   2000000        = Sample after value used during collection
   2654000000     = Total events (samples*SAV)

Module View (all values in decimal)


Module                   Process
  Event                              Events%  Samples  Events      Module Path
---------------------------------------------------------------------------------------------------------------------------
triad                    triad
  CPU_CLK_UNHALTED.THREAD            90.40%   2166    4332000000    /home/vtune/snb3/triad_src/triad
  INST_RETIRED.ANY                  89.98%   1194    2388000000
vmlinux                  triad
  CPU_CLK_UNHALTED.THREAD            4.47%    107     214000000     vmlinux
  INST_RETIRED.ANY                  4.97%    66      132000000
```

- **Thus CPU_CLK_UNHALTED.THREAD is event 0    "ei-00"**
- **Thus Inst_RETIRED.ANY is event 1    "ei-01"**

# Low level utilities

- Sfdump5 tbsXXXZZZ.tb5 /dumpsamples > samples.txt
  - Text dump of all samples
  - All sample records in a given file are same length
  - Length = SUM of all required fields for all events
    - If PEBS record is collected for PEBS events, the corresponding fields exist for non PEBS event but are zero filled
    - Events with LBR collection are only collected with other events that have SAME LBR filter value
      - 33 X 64 bits are added

# /dumpsamples example output

```
00000208  64--0033:0x0000000000400DF9-0  p-0x0000231C    c-00 t-0x0000231C    sgno-
   0x00000001  ei-00  tsc-0x0003C06F0CF15DD4   triad
```

- 00000208 is the record number
- 64--0033:0x0000000000400DF9-0 tells you this is a 64 bit binary and the IP of the interupt was 0x0000000000400DF9
- p-0x0000231C   gives the process ID
- c-00   the core number of the interupt   in this case 0
- t-0x0000231C  the thread ID
- ei-00   the event number
  - thus this is an record triggered by CPU_CLK_UNHALTED.THREAD
  - See –modules output to determine event numbers for a particular collection
- tsc-0x0003C06F0CF15DD4   the Time Stamp Counter
- Triad  the load module name

# /dumpsamples example output LBRs

```
00000091  64--0033:0x0000000000400694-0  p-0x00000A0A     c-00 t-0x00000A0A     sgno-
    0x00000001 ei-00  tsc-0x000000C43DECAFB1   extra_00-0x0000000000000006 extra_01-
    0x0000000000400A2C extra_02-0x00000000004009C4 extra_03-0x000000000040095C extra_04-
    0x00000000004008E6 extra_05-0x000000000040086E extra_06-0x0000000000400806 extra_07-
    0x000000000040074A extra_08-0x00000000004006E2 extra_09-0x0000000000401061 extra_10-
    0x0000000000400D7F extra_11-0x0000000000400D97 extra_12-0x0000000000400C52 extra_13-
    0x0000000000400BEC extra_14-0x0000000000400B84 extra_15-0x0000000000400AFC extra_16-
    0x0000000000400A94 extra_17-0x0000000000400976 extra_18-0x000000000040090E extra_19-
    0x0000000000400888 extra_20-0x0000000000400820 extra_21-0x00000000004007B8 extra_22-
    0x00000000004006FC extra_23-0x0000000000400694 extra_24-0x0000000000400648 extra_25-
    0x0000000000400D38 extra_26-0x0000000000400CC2 extra_27-0x0000000000400C06 extra_28-
    0x0000000000400B9E extra_29-0x0000000000400B36 extra_30-0x0000000000400AAE extra_31-
    0x0000000000400A46 extra_32-0x00000000004009DE call_chain
```

- record number is 91
- Event number (ei) is 0
- Extra_01 -> extra_16 are the branch source addresses
- Extra_17 -> extra_32 are the branch target addresses
- extra_00 points to the most recent LBR source entry
    - In this case extra_06
- Most recent target is extra_(extra_00+17)
    - Thus last target is extra_23 = extra_23-0x0000000000400694
    - And PEBS IP field is        = 64--0033:0x0000000000400694-0

# /dumpsamples example output PEBS

00000445  64--0033:0x0000000000401665-0  p-0x00000978     c-00 t-0x00000978     sgno-0x00000001 ei-00  tsc-0x0000011CF7198F6F   extra_00-0x0000000000000202 extra_01-0x0000000000401665 extra_02-0x00000123F1DE149A extra_03-0x0000000000000001 extra_04-0x0000000000000000 extra_05-0x00000123F1DE149A extra_06-0x000000001B4E4355 extra_07-0x000000004ABCE4E1 extra_08-0x00007FFFA989B710 extra_09-0x00007FFFA989B6A0 extra_10-0x0000000000000000 extra_11-0x0000000000000001 extra_12-0x00007FFFA989B400 extra_13-0x0000003731E97DD0 extra_14-0x0000000000400720 extra_15-0x00007FFFA989B860 extra_16-0x0000000000000000 extra_17-0x0000000000000000 extra_18-0x00007FFFA989B6F8 extra_19-0x0000000000000041 extra_20-0x0000000000000038 extra_21-0x000000000000FFFF extra_22-0x0000000000000000 store_fwd_lnx2

- Event number (ei) is 0 (in this case the latency event)
- Extra_01 is Event IP
  - IP of instruction after the instruction that caused the interupt  ("IP+1")
- Extra_02-> extra_17 are the register values at the completion of the offending instruction

(intel)

# PEBS Buffer field definitions

| | |
|---|---|
| (x)->r_flags | //extra_00 |
| (x)->linear_ip | //extra_01 |
| (x)->rax | //extra_02 |
| (x)->rbx | //extra_03 |
| (x)->rcx | //extra_04 |
| (x)->rdx | //extea_05 |
| (x)->rsi | //extra_06 |
| (x)->rdi | //extra_07 |
| (x)->rbp | //extra_08 |
| (x)->rsp | //extra_09 |
| (x)->r8 | //extra_10 |
| (x)->r9 | //extra_11 |
| (x)->r10 | //extra_12 |
| (x)->r11 | //extra_13 |
| (x)->r12 | //extra_14 |
| (x)->r13 | //extra_15 |
| (x)->r14 | //extra_16 |
| (x)->r15 | //extra_17 |
| (x)->data_linear_address | //extra_18 |
| (x)->data_source | //extra_19 |
| (x)->latency | //extra_20 |

# Precise Events

- **Significant expansion of PEBS capability on Intel® Core™ i7 Processors**
    - **4 events simultaneously**
    - **Latency event = IPF data ear + bit pattern for data source**
    - **Branches retired by type**
    - **Calls retired + LBR gives call counts**
    - **Calls_retired + full PEBS gives function arguments on Intel64**

# Data Access Analysis and PEBS

- **Data address profiling for loads and stores can be done as it is on Intel® Core™2 Processor Family**
  - **Full PEBS buffer + disassembly to identify registers with valid addresses at time of capture**
  - **Mem_inst_retired.load**
    - Cannot deal with mov rax,[rax] type instruction
  - **Mem_inst_retired.store**
    - Not subject to constraint of loads
  - **Inst_retired.any**
    - Cannot deal with EIP+1 = first instr of Basic Block

## PEBS Buffer

| 63 | BTS Buffer Base | 0 |
|----|:---:|----|
| | BTS Index | |
| | BTS Absolute Maximum | |
| | BTS Interrupt Threshold | |
| | PEBS Buffer Base | |
| | PEBS Index | |
| | PEBS Absolute Maximum | |
| | PEBS Interrupt Threshold | |
| | PEBS Counter Reset 0 | |
| | PEBS Counter Reset 1 | |
| | PEBS Counter Reset 2 | |
| | PEBS Counter Reset 3 | |

| Merom/Penryn - Format 0000b |
|:---:|
| Nehalem - Format 0001b |

| 63 | RFLAGS | 0 |
|----|:---:|----|
| | RIP | |
| | RAX | |
| | RBX | |
| | RCX | |
| | RDX | |
| | RSI | |
| | RDI | |
| | RBP | |
| | RSP | |
| | R8 | |
| | R15 | |
| | Global Perf Overflow MSR | |
| | Data Linear Address | |
| | Data Source (encodings) | |
| | Latency (core cycles) | |

(intel)

# Load Latency Threshold Event:

- **Ability to trigger count on minimum latency**
  - **Core cycles from load execute->data availability**
- **Linear address in PEBS buffer**
  - **<u>Allows driver to collect physical address</u>**
  - **Only total measurement of local/remote home access**
- **Data source captured in bit pattern**
  - **Actual NUMA source revealed**
- **Only ONE latency event/min thresh can be taken per run**
  - **Minimum latency programmed with MSR**
  - **Global per core**
    - 0x3F6 MS_PEBS_LD_LAT_THRESHOLD bits 15:0
  - **HW samples loads**
    - **EX: Sampling fraction for local dram= mem_inst_retired.latency_gt_128(DS= A or C) /mem_uncore_retired.local_dram**

(intel®)

# Front End/Decode Analysis

- **Instruction decode BW has lower maximum**
- **Instruction flow interruption at RAT output**
  - **UOPS_ISSUED.STALL_CYCLES – RESOURCE_STALLS.ANY**
  - **HT ON**
    - **subtract half the cycles as well**
    - **Or UOPS_ISSUED.CORE_STALL_CYCLES- RESOURCE_STALLS.ANY**
- **ILD_STALL.LCP_STALL**

# NUMA, Intel® QuickPath Interconnect, and Intel ® Xeon 5500/5600 Processor DP systems

- **Intel® QuickPath Interconnect (Intel® QPI) will greatly increase memory bandwidth of our platforms**

- **Integrated memory controllers on each socket access DIMMs**
  - **Intel® QPI provides cache coherency**
  - **Bandwidth improves by a lot**

- **Bandwidth improvement comes at a price**
  - **Non-Uniform Memory Access (NUMA)**
  - **Latency to DIMMs on remote sockets is ~2X larger**

**Pealing away the Bandwidth layer reveals the NUMA Latency layer**

(intel®)

# NUMA Modes on DP Systems Controlled in BIOS

- **Non-NUMA**
  - **Even/Odd lines assigned to sockets 0/1**
    - **Line interleaving**

- **NUMA mode**
  - **First Half of memory space on socket 0**
  - **Second half of memory space on socket 1**

(intel)

# Non-Uniform Memory Access and Parallel Execution

- **Parallel processing is intrinsically NUMA friendly**
  - **Affinity pinning maximizes local memory access**
  - **Message Passing Interface (MPI)**
  - **Parallel submission to batch queues**
  - **Standard for HPC**
- **Shared memory threading is more problematic**
  - **Explicit threading, OpenMP\* product, Intel® Threading Building Blocks (Intel® TBB)**
  - **NUMA friendly data decomposition (page-based) has not been required**
  - **OS-scheduled thread migration can aggravate situation**

*Other names and brands may be claimed as the property of others.

# HPC Applications will see Large Performance Gains due to Bandwidth Improvements

- **A remaining performance bottleneck may be due to Non-Uniform Memory Access latency**

- **This next level in the performance onion was not really addressed**
  - **Other performance tools offered little insight**
  - **Default usage of Non-NUMA BIOS settings**
    - **Except for some HPC accounts**

- **Intel® PTU data access profiling feature was designed to address NUMA**
  - **NHM events were designed to provide the required data**

# Gather and OOO execution

| | no prefetch | pref = 8 | pref = 16 | pref = 32 | pref = 64 | pref = 96 |
|---|---|---|---|---|---|---|
| 2 fp ops | 34.5 | 34.9 | 34.2 | 37.2 | 38.7 | 38.9 |
| 4 fp ops | 44.5 | 34.5 | 33.6 | 38 | 42.2 | 41.4 |
| 8 fp ops | 74.8 | 34.8 | 34.1 | 38.7 | 42.7 | 41.7 |
| 16 fp ops | 108.9 | 34.6 | 34 | 42.2 | 50.9 | 45.6 |

Data collected on Core™ 2 processor, prefetchers on

# Glossary

- PMU: Performance Monitoring Unit
  - Assembly of counters and programmable crossbars that allow counting and profiling using user selectable events

- FE: core pipeline Front End
  - Responsible for branch prediction, instruction fetch, decode to uops, allocation of OOO backend resources

- BE: core pipeline Backend
  - Stage uops waiting for inputs, execute upon availability, retire in order

# Glossary

- RS: reservation station
  - Where uops are staged for execution waiting for availability of their inputs
- ROB: Reorder Buffer
  - Where uops wait prior to retirement until all older uops have retired and execution path is confirmed. Second point corrects when uops are executed on a mispredicted path.
- RAT: Resource Allocation Table
  - Allocates BE resources for uops prior to issuing them from front end of pipeline to the backend

*intel*

# Glossary

- Cachelines are 64 bytes

- LLC: Last level Cache
  - L3 on these processors

- LFB: line fill buffer
  - Buffers used for transfering cachelines into and out of L1D

- WB: writeback
  - Modified data is written back to higher level in memory subsystem on eviction

- RFO: Read for Ownership
  - Stores require cachelines are in exclusive ownership state so they can be modified

# Glossary

- Prefetch, by hardware (HW) or by explicit instruction (SW)
  - Request cacheline prior to execution of consuming instruction (load/store) with intention of hiding latency

- BW: bandwidth
  - Data moved/unit time. I prefer cachelines/cycle as that is what is measured

- Latency: time required to transfer a single line from source to usage.

(intel®)

# Glossary

- SIMD: Single instruction multiple data
  - SSE parallel execution mode
  - AKA vectorization
- X87: legacy floating point computation mode. In contrast to SSE FP instructions
- NT: Non Temporal
  - Data store mode that writebacks data in 64 byte aligned contiguous 64 byte chunks directly to dram without RFO
- HITM: Hit Modified
  - Snoop response when line is found in modified state in another cache

# Glossary

- HT: Intel® Hyper-threading Technology
  - Execution mode allowing uops from two threads to be executed in an intermingled flow, without an OS context switch, through a single core pipeline.

- Turbo: Intel® Turbo Boost Technology
  - Adjusting core frequency upwards on active core when other cores are under utilized, while staying within required power envelope. Enhances performance of single threaded execution