# Advanced Tuning Methodologies

Dr.-Ing. Michael Klemm
Software and Services Group
Intel Corporation
(michael.klemm@intel.com)

# Legal Disclaimer & Optimization Notice

**Optimization Notice**

# Agenda

- Performance Tuning Methodology
- Intel® VTune™ Amplifier XE
  - Fundamental Analysis: Hotspots
  - Finding Issues in Parallel Applications
  - Using the Performance Monitoring Unit
- Optimization Strategies
  - Data Layout and Alignment
  - Vectorization Potential
  - Optimizing Data Transfers
  - Miscellaneous Tuning Hints

# Shameless Plug...

**Authors:** Alexander Supalov, Andrey Semin, Michael Klemm, Chris Dahnken

**Table of Contents:**
Foreword by Bronis de Supinski (CTO LLNL)
Preface
Chapter 1: No Time to Read this Book?
Chapter 2: Overview of Platform Architectures
Chapter 3: Top-Down Software Optimization
Chapter 4: Addressing System Bottlenecks
Chapter 5: Addressing Application Bottlenecks:
Distributed Memory
Chapter 6: Addressing Application Bottlenecks:
Shared Memory
Chapter 7: Addressing Microarchitecture Bottlenecks
Chapter 8: Application Design Implications

ISBN-13 (pbk): 978-1-4302-6496-5
ISBN-13 (electronic): 978-1-4302-6497-2

Order now at http://www.apress.com/9781430264965

# Agenda

- Performance Tuning Methodology
- Intel® VTune™ Amplifier XE
  - Fundamental Analysis: Hotspots
  - Finding Issues in Parallel Applications
  - Using the Performance Monitoring Unit
- Optimization Strategies
  - Data Layout and Alignment
  - Vectorization Potential
  - Optimizing Data Transfers
  - Miscellaneous Tuning Hints

**Software & Services Group, Developer Products Division**

6

# The Software Optimization Process

- The process of improving the software by eliminating bottlenecks so that it operates more efficiently on a given hardware and uses resources optimally
  - Identifying bottlenecks in the target application and eliminating them appropriately is the key to an efficient optimization

- There are many optimization methodologies, which help developers answer the questions of
  - Why to optimize?
  - What to optimize?
  - To what to optimize?
  - How to optimize?

These methods aid developers to reach their performance requirements.

# Performance Analysis Methodology
## Optimization: A Top-down Approach

- Use top down approach
- Understand application and system characteristics
  - Use appropriate tools at each level

| System Config, BIOS, OS Network I/O, Disk I/O, Database Tuning, etc. | **System** |

| Application Design Algorithmic Tuning Driver Turning Parallelization | **Application** |

| Cache/Memory Instructions SIMD others | **Processor** |

# Performance Analysis Methodology
## Optimization: A Top-down Approach



1. Create a Benchmark
2. Collect Data
3. Analyze Data and Identify Performance Problems
4. Fix the problems in your code or system
5. Is Problem Fixed?
6. Are performance requirements met?

- •Repeatable
- •Representative
- •Easy to run
- •Verifiable
- •Measure elapsed time
- •Reasonable coverage
- •Precision

No
No
Yes
Yes

Optimized Code

System
Application
Processor

# When to Stop

- Is architecture at maximum efficiency?
  - What this means: calculating theoretical maximum.
  - Know about best values for CPI or IPC.
  - Know the maximum FLOPS for the data type used.

- Is the performance requirement fulfilled?
  - What are the performance requirements?
  - Incrementally complete optimizations until done.
  - Key question: Are you "happy" with it?

CPI:    Cycles per Instructions
IPC:    Instructions per Cycle
FLOPS: Floating-Point Oper. Per Sec.

# Questions to Ask Yourself

*"We should forget about small efficiencies, say about 97% of the time: premature optimization is the root of all evil."*

— Donald Knuth

Quality code is:

- Portable
- Readable
- Maintainable
- Reliable

**Intelligently sacrifice quality for performance**

# Amdahl's Law Slightly Reinterpreted

1 core, t=140

| |
|---|
| 20 |

↓

| |
|---|
| 50 |

↓

| |
|---|
| 10 |

↓

| |
|---|
| 50 |

↓

| |
|---|
| 10 |

Hotspot

1 core, optimized

| |
|---|
| 20 |

↓

| |
|---|
| 25 |

↓

| |
|---|
| 10 |

↓

| |
|---|
| 25 |

↓

| |
|---|
| 10 |

1 cores, optimized even more

| |
|---|
| 20 |

↓

| |
|---|
| 5 |

↓

| |
|---|
| 10 |

↓

| |
|---|
| 5 |

↓

| |
|---|
| 10 |

A 2x improvement in the hotspots overall leads to 1.5x

A 10x improvement in the hotspots leads to 2.8x

# Agenda

- Performance Tuning Methodology
- Intel® VTune™ Amplifier XE
  - Fundamental Analysis: Hotspots
  - Finding Issues in Parallel Applications
  - Using the Performance Monitoring Unit
- Optimization Strategies
  - Data Layout and Alignment
  - Vectorization Potential
  - Optimizing Data Transfers
  - Miscellaneous Tuning Hints

# Selecting type of data collection



All available analysis types

Different ways to start the analysis

Helps creating new analysis types

Copy correct command line syntax to clipboard

# VTune™ Amplifier XE
## GUI Layout

**Adjust Data Grouping**

Function - Call Stack
Module - Function - Call Stack
Source File - Function - Call Stack
Thread - Function - Call Stack

... (Partial list shown)

**Click [+] for Call Stack**

**Double Click Function to View Source**

**Filter by Timeline Selection (or by Grid Selection)**

Zoom In And Filter On Selection
Filter In by Selection
Remove All Filters

**Filter by Module & Other Controls**

---

**Am Hotspots** - View hotspots colored by CPU usage

| Analysis Target | Analysis Type | Collection Log | Summary | Bottom-up | Top-down Tr |

CPU Time▼

Function - Call Stack

Module

| Idle | Poor | Ok | Ideal | Over |

| | CPU Time | Module |
|---|---|---|
| ⊞ dllStopPlugin | 7.550s | RenderSystem_Direct3D9.DLL |
| ⊞ FireObject::checkColli | 6.389s | SystemProceduralFire.DLL |
| ⊞ FireObject::ProcessFire | 4.592s | SystemProceduralFire.DLL |
| ⊞ BaseThreadInitThunk | 2.566s | kernel32.dll |
| ⊞ Ogre::FileStreamDataS | 2.562s | OgreMain.dll |
| ⊞ TaskManagerTBB::Par | 2.533s | Smoke.exe |
| ⊞ AIScene::GetPOI | 1.710s | SystemAI.DLL |
| ⊞ TaskManagerTBB::Wa | 1.682s | Smoke.exe |
| Selected 189 row(s): | 47.481s | |

10s  20s  30s  40s  50s  60s

Thread
wWinMainCRTStart
Thread (0x1a28)
Thread (0x8e4)
Thread (0x29c8)

CPU Usage

Frames over Time

**Ruler Area**
☑ ▽ Global Mark
☐ ▽▽ Frame
☑ **Thread**
☑ ▇ Running
☑ CPU Time
☐ ▽▽ User Task
☑ **CPU Usage**
CPU Time
☑ **Frames over Ti...**

No filters are applied. ✖ Module: [All] ▼    Call Stack Mode: Only user functions ▼

---

(intel) Software

# VTune™ Amplifier XE
## *GUI Layout*

# VTune™ Amplifier XE
## GUI Layout



**Transitions** — Locks & Waits

**CPU Time** — Hotspots / Lightweight Hotspots

Ruler Area:
- ☑ Frame
- ☑ Thread
  - ☑ Running
  - ☑ Waits
  - ☑ User Task
  - ☑ Transition
- ☑ Thread Concurrency
  - Concurrency
- ☑ Frames over Time
  - Frame Rate

Hovers:

Frame
Start: 29.858s Duration: 0.017s
Frame: 72
Frame Domain: Smoke::Framework::execute()
Frame Type: Good
Frame Rate: 59.8242179

Transition
Transition
wWinMainCRTStartup (0x12d4) to Thread (0x138c) (29.899s to 29.899s)
Sync Object: TBB Scheduler
Object Creation File: taskmanagertbb.cpp
Object Creation Line: 318

User Task
Start: 29.958s Duration: 0.018s
Task Type: Smoke::FrameWork::execute()::Other
Task End Call Stack: Framework::Execute

CPU Time
94.233472%

- Optional: Use API to mark frames and user tasks
- Optional: Add a mark during collection  [Mark Timeline]

# Agenda

- Performance Tuning Methodology
- **Intel® VTune™ Amplifier XE**
  - **Fundamental Analysis: Hotspots**
  - Finding Issues in Parallel Applications
  - Using the Performance Monitoring Unit
- Optimization Strategies
  - Data Layout and Alignment
  - Vectorization Potential
  - Optimizing Data Transfers
  - Miscellaneous Tuning Hints

# Readying Your Application for Intel VTune Amplifier XE

- You should run Amplifier XE on a "Released/Optimized" build.

- Symbols Allow you to view the source (not just the assembly)
  - Windows: /Zi
  - Linux: –g

- Intel Threading Runtimes need instrumented runtimes
  - TBB: Define TBB_USE_THREADING_TOOLS
  - OpenMP: Use Intel Dynamic Version of OpenMP

- Call Stack Mode – Requires use of the dynamic version of the C Runtime library to properly attribute System Calls
  - Windows use:/MD(d)
  - Linux do not use: -static

# Analysis Types
## *Basic Hotspots*

- For each sample, capture execution context, time passed since previous sample and thread CPU time
- Allows time spent in system calls to be attributed to user functions making the call
- Provides additional knobs:

  - The defaults for Hotspot analysis are configurable and can be done so by creating a custom analysis type inherited from Hotspots



**User-mode Sampling and Tracing Settings**

| | |
|---|---|
| CPU sampling interval, ms: | 10 |
| Collect CPU sampling data: | With stacks |
| Collect signalling API data: | No |
| Collect synchronization API data: | No |
| Collect I/O API data: | No |

☑ Collect timeline data

# Analysis Types
## *Advanced Hotspots*

- Similar to Hotspot Analysis
  - Sampling is performed with the SEP collector
  - Driver is required
- Stack walking can be performed
  - Only hotspots are reported by default
  - More complex data collection is possible, incl. trip counts
- Samples are taken more frequently, but may have less accurate timing information
- Analysis may be performed for a single application or for the entire system

# Agenda

- Performance Tuning Methodology
- **Intel® VTune™ Amplifier XE**
  - Fundamental Analysis: Hotspots
  - **Finding Issues in Parallel Applications**
  - Using the Performance Monitoring Unit
- Optimization Strategies
  - Data Layout and Alignment
  - Vectorization Potential
  - Optimizing Data Transfers
  - Miscellaneous Tuning Hints

# Issues in Parallel Applications

- Load imbalance
  - Work distribution is not optimal
  - Some threads are heavily loaded, while others idle
  - Slowest thread determines total speed-up

- Locking issues
  - Locks prohibit threads to concurrently enter code regions
  - Effectively serialize execution

- Parallelization overhead
  - With large no. of threads, data partition get smaller
  - Overhead might get significant (e.g. OpenMP startup time)

# Issues in Parallel Applications

- Load imbalance
  - Work distribution is not optimal
  - Some threads are heavily loaded, while others idle
  - Slowest thread determines total speed-up

- Locking issues
  - Locks prohibit threads to concurrently enter code regions
  - Effectively serialize execution

- Parallelization overhead
  - With large no. of threads, data partition get smaller
  - Overhead might get significant (e.g. OpenMP startup time)

# Threading Analysis Terminology



- **Elapsed Time**: 6 seconds
- **CPU Time**: T1 (4s) + T2 (3s) + T3 (3s) = 10 seconds
- **Wait Time**: T1(2s) + T2(2s) + T3 (2s) = 6 seconds

# Analysis Types
## *Concurrency*

# Hotspots Analysis vs. Concurrency Analysis

- Hotspot Analysis and Concurrency Analysis are similar:

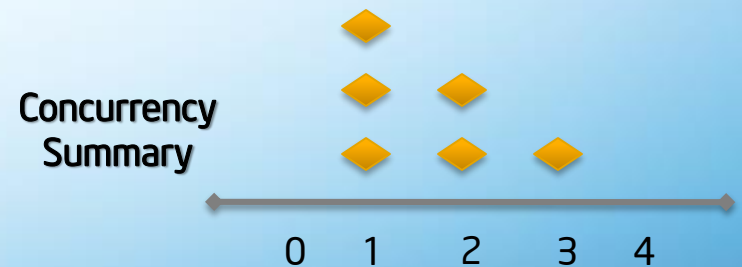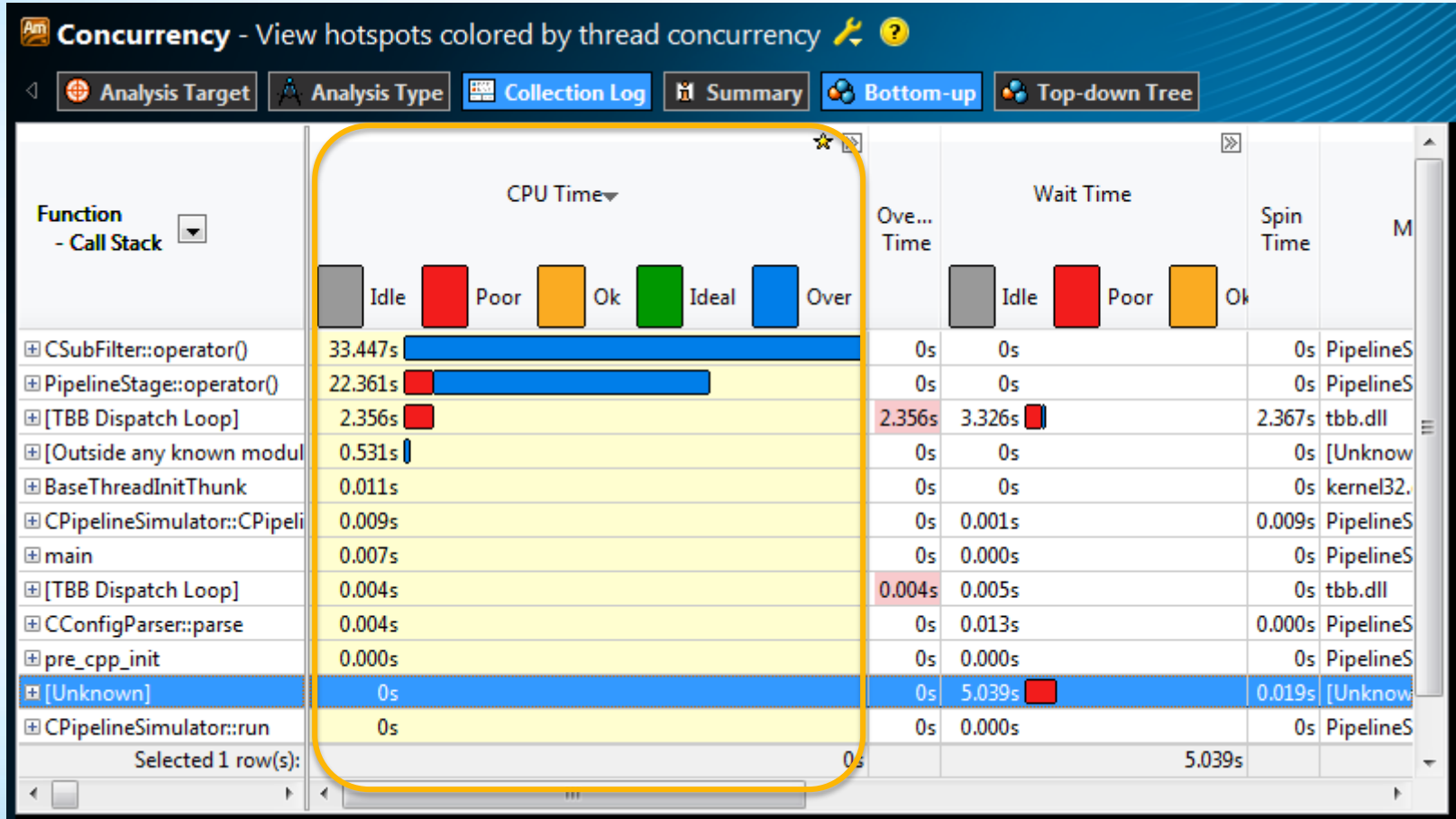# Thread Concurrency Histogram

# Issues in Parallel Applications

- Load imbalance
  - Work distribution is not optimal
  - Some threads are heavily loaded, while others idle
  - Slowest thread determines total speed-up

- Locking issues
  - Locks prohibit threads to concurrently enter code regions
  - Effectively serialize execution

- Parallelization overhead
  - With large no. of threads, data partition get smaller
  - Overhead might get significant (e.g. OpenMP startup time)

# Analysis Types
## *Lock and Waits*

# Timeline Visualizes Thread Behavior

- Retrieve additional information about waiting threads

# Timeline Visualizes Thread Behavior

- Retrieve additional information on thread transitions

# Drilling down into Thread Behavior

- Reveal source code that causes thread transitions

# Agenda

- Performance Tuning Methodology
- Intel® VTune™ Amplifier XE
  - Fundamental Analysis: Hotspots
  - Finding Issues in Parallel Applications
  - **Using the Performance Monitoring Unit**
- Optimization Strategies
  - Data Layout and Alignment
  - Vectorization Potential
  - Optimizing Data Transfers
  - Miscellaneous Tuning Hints

# Performance Monitoring Unit (PMU)

- **Per-core PMU**:
  - Each core provides 2 programmable counters and 1 fixed counter.
  - The programmable per-core counters can be configured to investigate front-end/micro-op flow issues, stalls inside a processor core.

- **Uncore PMU:**
  - Uncore of the coprocessor has four counters to monitor uncore events
  - Can be used to investigate memory behavior and global on-chip issues

# Event Based Performance Analysis

## Event Based Sampling(EBS)

- Both architectural and non-architectural processor events can be monitored using sampling and counting technologies

  **Sampling:** Allows to profile all active software on the system, including operating system, device driver, and application software.

  - Event-based samples are collected periodically after a specific number of processor events have occurred while the program is running

  - The program is interrupted, allowing the interrupt handling driver to collect the Instruction Pointer (IP), load module, thread and process ID's

  - Instruction pointer is then used to derive the function and source line number from the debug information created at compile time

# How Event Based Sampling (EBS) Works

**Select Event Signal** → **Count Down** → **"Sample After" Value**

**Underflow to Zero**

**Internal Interrupt Controller**

**Interrupt CPU to take a sample Performance Monitoring Interrupt (PMI)**

- A performance counter increments on the CPU every time an event occurs
- A sample of the execution context is recorded every time a performance counter overflows

**Events = samples * sample after value**

# Native Launch configuration



- Application settings:
  - Application: ssh
  - Parameters:          mic0 "<app startup>"
  - Working directory:   Usually does not matter
  - Don't forget to set search directories under "All files"

# Application Configuration

# Configuring a User-defined Analysis

Software

# Some useful events and metrics

| Scenario | Event name(s) |
|---|---|
| Wall-clock profiling | `CPU_CLK_UNHALTED,INSTRUCTIONS_EXECUTED` (or `EXEC_STAGE_CYCLES`) |
| Main memory bandwidth | `L2_DATA_READ_MISS_MEM_FILL,` `L2_DATA_WRITE_MISS_MEM_FILL` |
| L1 Cache misses | `DATA_READ_MISS_OR_WRITE_MISS` |
| TLB misses and page faults | `DATA_PAGE_WALK, LONG_DATA_PAGE_WALK,` `DATA_PAGE_FAULT` |
| Vectorized code execution | `VPU_INSTRUCTIONS_EXECUTED,` `VPU_ELEMENTS_ACTIVE` |
| Various hazards | `BRANCHES_MISPREDICTED` |
| Cycles per instruction | `CPU_CLK_UNHALTED /` `INSTRUCTIONS_EXECUTED` |
| Memory Bandwidth (used by all cores at once) | `(L2_DATA_READ_MISS_MEM_FILL +` `L2_DATA_WRITE_MISS_MEM_FILL) * 64 /` `CPU_CLK_UNHALTED / Frequency` |

# Example: Hotspots of OpenFOAM

# Agenda

- Performance Tuning Methodology
- Intel® VTune™ Amplifier XE
  - Fundamental Analysis: Hotspots
  - Finding Issues in Parallel Applications
  - Using the Performance Monitoring Unit
- **Optimization Strategies**
  - **Data Layout and Alignment**
  - Vectorization Potential
  - Optimizing Data Transfers
  - Miscellaneous Tuning Hints

# Data Layout – Common Layouts

## Array-of-Structs (AoS)



- Pros:
  Good locality of
  {x, y, z}.
  1 memory stream.

- Cons:
  Potential for
  gather/scatter.

## Struct-of-Arrays (SoA)



- Pros:
  Contiguous
  load/store.

- Cons:
  Poor locality of
  {x, y, z}.
  3 memory streams.

## Hybrid (AoSoA)



- Pros:
  Contiguous
  load/store.
  1 memory stream.

- Cons:
  Not a "normal"
  layout.

# Data Layout – Why It's Important

- ## Instruction-Level
  - Hardware is optimized for contiguous loads/stores.
  - Support for non-contiguous accesses differs with hardware. (e.g. AVX2/KNC gather)

- ## Memory-Level
  - Contiguous memory accesses are cache-friendly.
  - Number of memory streams can place pressure on prefetchers.

# Data Alignment – Why It's Important

Cache Line 0

| 0 | 1 | 2 | 3 | ... | ... | 6 | 7 |
|---|---|---|---|-----|-----|---|---|

Cache Line 1

| 8 | 9 | ... | ... | ... | ... | ... | ... |
|---|---|-----|-----|-----|-----|-----|-----|

| 0 | 1 | 2 | 3 |
|---|---|---|---|

| 6 | 7 | 8 | 9 |
|---|---|---|---|

**Aligned Load**

- Address is aligned.
- One cache line.
- One instruction.

**Unaligned Load**

- Address is not aligned.
- Potentially multiple cache lines.
- Potentially multiple instructions.

# Data Alignment – Why It's Important

- Cache Associativity
  - L1 and L2 on Knights Corner are 8-way associative.
    - L1: 8 cache lines that are 32KB/8 = 4KB apart.
    - L2: 8 cache lines that are 512KB/8 = 64KB apart.

- Set Conflicts
  - Occur when references are a multiple of 4K (L1) or 64K (L2) apart.
  - Look for high cache miss rate, even though working set < cache capacity.
  - Solution is to pad arrays appropriately.

# Data Alignment – Sample Applications

## 1) Align memory

- _mm_malloc(bytes, 64)                    /          !dir$ attributes align:64

## 2) Access Memory in an aligned way

- for (i = 0; i < N; i++) { array[i] … }

## 3) Tell the compiler!

- #pragma vector aligned        /        !dir$ vector aligned
- __assume_aligned(p, 16)      /        !dir$ assume_aligned (p, 16)
- __assume(i % 16 == 0)        /        !dir$ assume (mod(i, 16) .eq. 0)

# Data Alignment – Real-life Applications



| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| 8 | 9 | … | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |

Data

# Data Alignment – Real-life Applications



| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| 8 | 9 | … |   |   |   |   |   |

Data

# Data Alignment – Real-life Applications

# Data Alignment – Real-life Applications



Data

Halo

# Data Alignment – Real-life Applications

# Data Alignment – Real-life Applications



Data

Halo

Padding

Not strictly necessary…

# Example: saxpy

```c
void saxpy(float *y, float *x, float a, int n) {
    int i;
    for (i = 0; i < n; ++i){
        y[i] = a*x[i] + y[i];
    }
}
```

# The Effect of Data Alignment on saxpy

## Unaligned

```
vloadunpackld (%rsi,%r14,4),%zmm1
vprefetch1 0x400(%rsi,%r14,4)
vloadunpackld 0x40(%rsi,%r14,4),%zmm2
vprefetch0 0x200(%rsi,%r14,4)
vloadunpackhd 0x40(%rsi,%r14,4),%zmm1
vprefetche1 0x400(%rdi,%r14,4)
vloadunpackhd 0x80(%rsi,%r14,4),%zmm2
vprefetch0 0x200(%rdi,%r14,4)
vfmadd213ps (%rdi,%r14,4),%zmm0,%zmm1
vprefetch1 0x440(%rsi,%r14,4)
vfmadd213ps 0x40(%rdi,%r14,4),%zmm0,%zmm2
vprefetch0 0x240(%rsi,%r14,4)
vmovaps %zmm1,(%rdi,%r14,4)
vprefetche1 0x440(%rdi,%r14,4)
vmovaps %zmm2,0x40(%rdi,%r14,4)
vprefetch0 0x240(%rdi,%r14,4)
```

## Aligned

```
vmovaps (%rsi,%r9,4),%zmm2
vprefetch1 0x400(%rsi,%r9,4)
vmovaps 0x40(%rsi,%r9,4),%zmm3
vprefetch0 0x200(%rsi,%r9,4)
vfmadd213ps (%rdi,%r9,4),%zmm0,%zmm2
vprefetche1 0x400(%rdi,%r9,4)
vfmadd213ps 0x40(%rdi,%r9,4),%zmm0,%zmm3
vprefetch0 0x200(%rdi,%r9,4)
vmovaps %zmm2,(%rdi,%r9,4)
vprefetch1 0x440(%rsi,%r9,4)
vmovaps %zmm3,0x40(%rdi,%r9,4)
vprefetch0 0x240(%rsi,%r9,4)
vprefetche1 0x440(%rdi,%r9,4)
mov %al,%al
vprefetch0 0x240(%rdi,%r9,4)
```

# Data Alignment – Host

- AVX instructions are less sensitive to unaligned accesses.

- Still makes sense to align data on the host
  - Cache aligned
  - 4K page aligned

- Data alignment ensures that DMA transfers work best when offloading data to the target
  - Ideal alignment: 4K pages

# Agenda

- Performance Tuning Methodology
- Intel® VTune™ Amplifier XE
  - Fundamental Analysis: Hotspots
  - Finding Issues in Parallel Applications
  - Using the Performance Monitoring Unit
- **Optimization Strategies**
  - Data Layout and Alignment
  - **Vectorization Potential**
  - Optimizing Data Transfers
  - Miscellaneous Tuning Hints

# Loop Optimizations

- SIMD works best, if
  - data is properly aligned
  - data is stored in large portions of stride 1
  - loops have a high trip count

- Vectorization potential
  - Average utilization of SIMD register for a loop
  - VP = (N/vp) / ceil(N/vl)

- Vectorization potential can be measured
  - % vector instructions: VPU_INSTRUCTIONS_EXECUTED/INSTRUCTIONS_EXECUTED
  - Avg. % elements used per vector: VPU_ELEMENTS_ACTIVE/INSTRUCTIONS_EXECUTED

# Example Loop: NWChem

```
do p4=1,p4d
do p5=1,p5d
do p6=1,p6d
do h1=1,h1d
do h7=1,h7d
do h2=1,h2d
do h3=1,h3d
triplesx(h3,h2,h1,p6,p5,p4) =
  triplesx(h3,h2,h1,p6,p5,p4) –
    t2sub(h7,p4,p5,h1)*
      v2sub(h3,h2,p6,h7)
end do
end do
end do
end do
end do
end do
end do
```

- Each loop runs on the order of 10 to 25 iterations

- Example: trip count 19

- Vectorization potential
  - SIMD length: 8
  - 2 registers to cover 16 iterations
  - remainder loop: 3 trips

  80%

# Example Loop: NWChem - Optimized

```
do p4=1,p4d
do p5=1,p5d
do p6=1,p6d
do h1=1,h1d
do h7=1,h7d
do h3h2=1,h2d*h3d
triplesx(h3h2,h1,p6,p5,p4) =
  triplesx(h3h2,h1,p6,p5,p4) -
    t2sub(h7,p4,p5,h1)*
      v2sub(h3h2,p6,h7)
end do
end do
end do
end do
end do
end do
```

- Fused loop has about 1 or 2 orders of magnitudes more trips

- Example: trip count 19

- Vectorization potential
  - SIMD length: 8
  - 45 registers to cover 360 iterations
  - remainder loop: 1 trips

  98%

# Agenda

- Performance Tuning Methodology
- Intel® VTune™ Amplifier XE
  - Fundamental Analysis: Hotspots
  - Finding Issues in Parallel Applications
  - Using the Performance Monitoring Unit
- **Optimization Strategies**
  - Data Layout and Alignment
  - Vectorization Potential
  - **Optimizing Data Transfers**
  - Miscellaneous Tuning Hints

# Optimizing Data Transfers

- PCIe is a bottleneck compare to memory
  - High latency (on the order of microseconds)
  - Low bandwidth (1 to 2 orders of magnitude less than GDDR)

- Avoid data transfers as much as possible

- Overlap data transfers and computation where possible

- Stay on the target device as long as possible

# Double Buffering Example

- Overlap computation and communication
- Generalizes to data domain decomposition

# Double Buffering Example - Main

```
int main(int argc, char* argv[])  {
   int i;
   double st_time, end_time;
   double sync_tm, async_in_tm;
   //  Allocate & initialize in1, res1,
   // in2, res2 on the host
```

Allocate arrays on target: alloc_if(1)

Retain for duration of sample: free_if(0)

in1 and in2 represent 2 separate buffers.

```
#pragma offload_transfer target(mic:0) in(cnt) \
      nocopy(in1, res1, in2, res2 : length(cnt) \
   alloc_if(1) free_if(0) )


do_async_in();


// Validate results and print timings
```

Free target allocations: free_if(1)

```
#pragma offload_transfer target(mic:0) \
      nocopy(in1, res1, in2, res2 : length(cnt) \
   alloc_if(0) free_if(1) )
```

# Double Buffering – do_asynch_in, evens

```
void do_async_in() {
    float lsum;
    int i;
    lsum = 0.0f;

#pragma offload_transfer target(mic:0)   \
        in(in1 : length(cnt) alloc_if(0) free_if(0) ) signal(in1)


    for (i=0; i < iter; i++) {
        if (i%2 == 0) {

            #pragma offload_transfer target(mic:0) if(i!=iter-1)  \
                    in(in2 : length(cnt) alloc_if(0) free_if(0) ) \
                    signal(in2)


            #pragma offload target(mic:0) nocopy(in1) wait(in1)  \
                    out(res1 : length(cnt) alloc_if(0) free_if(0) )
            {
                compute(in1, res1);
            }

            lsum = lsum + sum_array(res1);
        }
```

> Begin an initial transfer of the first dataset for the target to work on.  Transfer begins immediately, is non-blocking, and will signal when complete.

> For even loop iterations (except the final), first start another non-blocking transfer of the next dataset.

> While that transfer progresses, process the previous dataset through the compute() function, first waiting for its transfer to complete, then return a result.  Execution on the host waits for the function to return.

# Double Buffering – do_asynch_in, odds

```
else {

    #pragma offload_transfer target(mic:0) if(i!=iter-1) \
        in(in1 : length(cnt) alloc_if(0) free_if(0) ) \
        signal(in1)




    #pragma offload target(mic:0) nocopy(in2) wait(in2) \
        out(res2 : length(cnt) alloc_if(0) free_if(0) )
    {
        compute13(in2, res2);
    }

    lsum = lsum + sum_array(res2);
  }
}
async_in_sum = lsum / (float) iter;
}
```

For odd iterations (except the final), work on the other buffer.  Start another non-blocking transfer.

Offload the compute function, and host waits for the result. Repeat, alternating between the buffers (in1 & in2).

# Agenda

- Performance Tuning Methodology
- Intel® VTune™ Amplifier XE
  - Fundamental Analysis: Hotspots
  - Finding Issues in Parallel Applications
  - Using the Performance Monitoring Unit
- **Optimization Strategies**
  - Data Layout and Alignment
  - Vectorization Potential
  - Optimizing Data Transfers
  - **Miscellaneous Tuning Hints**

# Miscellaneous Tuning Hints

- Check if data transfer is necessary
  - Keep data as long on the device as possible
  - Avoid transferring scratch pad memory

- Use 2MB pages for offload data transfers
  - (May) result in less page faults
  - (May) improve utilization of the TLB cache
  - (May) increase time needed to initialize pages

- Controlled by environment variable:
  ```
  export MIC_USE_2MB_BUFFERS=16k
  ```

intel
Software

# Miscellaneous Tuning Hints

- Use 2MB pages for dynamic memory
- Reserve 2MB pages in the kernel:
  - Remember this leaves less memory available for other programs!
  - Example (reserves 128 2MB pages):
    
    echo "128" > /proc/sys/vm/nr_hugepages
- Use mmap to request allocation of 2MB pages:

```
#include <sys/mman.h>
void *p;
size_t s = 64*2*1024*1024; // 128MB
…
p = mmap(0, s, PROT_READ|PROT_WRITE,
        MAP_ANONYMOUS|MAP_PRIVATE|MAP_HUGETLB, -1, 0);
if (p == MAP_FAILED)
    perror("mmap failed");
```

intel
Software

# Miscellaneous Tuning Hints

- Always use thread affinity on the coprocessor
  - Avoid expensive thread migration between cores
  - Example:
    ```
    export MIC_ENV_PREFIX=MIC
    export MIC_KMP_AFFINITY=compact / scatter / balanced
    ```

  - Affinity policies:
    - compact:        keep threads close
    - scatter:        distribute threads as far apart as poissible
    - balanced:       create *compact* groups and *scatter* these

  - Simpler approach:
    ```
    export KMP_PLACE_THREADS=60c,3t,0o
    ```

# Miscellaneous Tuning Hints

- Do not overload the coprocessor
  - Multiple offloads from different host MPI ranks
  - Multiple offloads from different host threads


- Partition the device:
  - Rank 0: export MIC_PLACE_THREADS=30c,3t,0o
  - Rank 1: export MIC_PLACE_THREADS=30c,3t,30o


- May help to improve relative performance
  - Two offloads will half the compute load each might be more beneficial
  - Less synchronization / threading overhead

# Summary

- Use profiling tools to find hotspots and offload candidates

- Use a performance analysis tool to find threading issues

- Optimize offloaded code
  - Data layout and alignment
  - Vectorization potential