



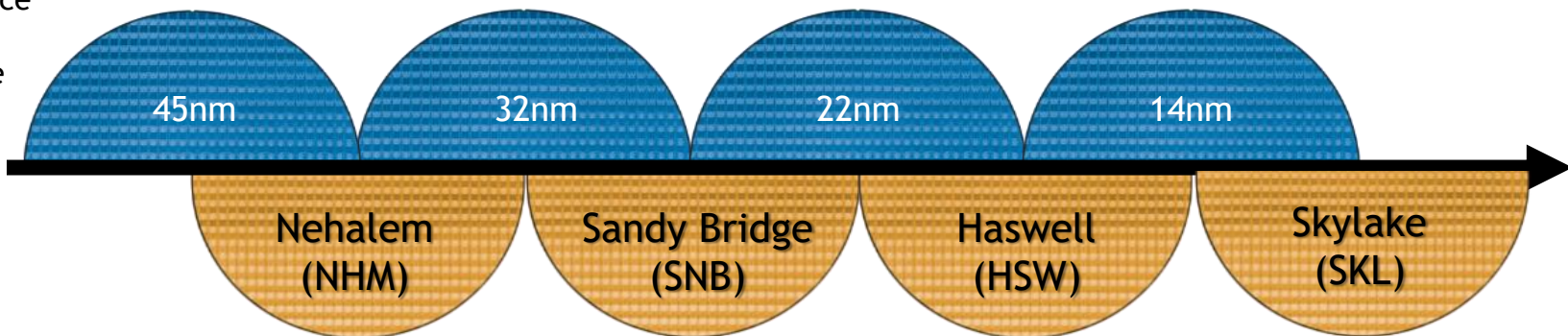
Intel® Software Tools for KNL

甘弛

Developer Products Division/SSG

IA Cores build on a Common Architecture

Scalable Performance
Energy Efficient
Microarchitecture



Highly Parallel
Energy Efficient
Architecture



All dates, product descriptions, availability, and plans are forecasts and subject to change without notice.

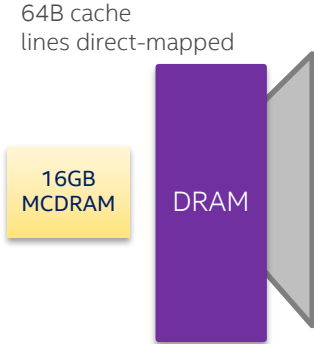
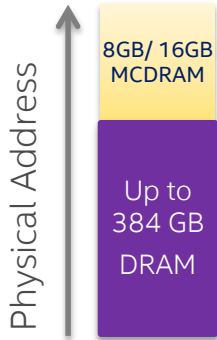
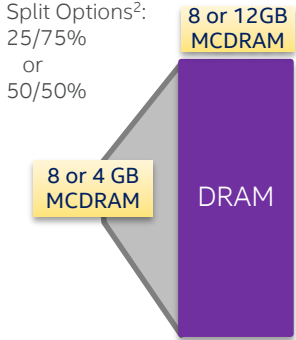
Intel and the Intel logo are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States and other countries. * Other names and brands may be claimed as the property of others. All products, dates, and figures are preliminary and are subject to change without any notice. Copyright © 2016, Intel Corporation.



Integrated On-Package Memory Usage Models

Model configurable at boot time and software exposed through NUMA¹

Platform Memory (DDR4) only available for bootable KNL host processor

	Cache Model	Flat Model	Hybrid Model
			
Description	Hardware automatically manages the MCDRAM as a “L3 cache” between CPU and ext DDR memory	Manually manage how the app uses the integrated on-package memory and external DDR for peak perf	Harness the benefits of both Cache and Flat models by segmenting the integrated on-package memory
Usage Model	<ul style="list-style-type: none">▪ App and/or data set is very large and will not fit into MCDRAM▪ Unknown or unstructured memory access behavior	<ul style="list-style-type: none">▪ App or portion of an app or data set that can be, or is needed to be “locked” into MCDRAM so it doesn’t get flushed out	<ul style="list-style-type: none">▪ Need to “lock” in a relatively small portion of an app or data set via the Flat model▪ Remaining MCDRAM can then be configured as Cache

1. NUMA = non-uniform memory access

2. As projected based on early product definition

Intel® Parallel Studio XE – Components

	Composer Edition	Professional Edition	Cluster Edition
Intel® C++ compiler	✓	✓	✓
Intel® Fortran compiler	✓	✓	✓
Intel® Math Kernel Library	✓	✓	✓
Intel® Threading Building Blocks library	✓	✓	✓
Intel® Integrated Performance Primitives	✓	✓	✓
Intel® Cilk™ Plus parallel model	✓	✓	✓
OpenMP*	✓	✓	✓
Intel® Advisor XE		✓	✓
Intel® Inspector XE		✓	✓
Intel® VTune™ Amplifier XE		✓	✓
Intel® Data Analytic Acceleration Library (Intel® DAAL)	✓	✓	✓
Intel® MPI library			✓
Intel® Trace Analyzer and Collector			✓
Rogue Wave IMSL* Library	Bundled & Add-on	Add-on	Add-on

Optimization Notice

Copyright © 2015, Intel Corporation. All rights reserved.

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



Generate Code for KNL with Intel Compiler 17.0

Linux* OS:

`-xMIC-AVX512`

Windows* OS:

`/QxMIC-AVX512`

May generate Intel® Advanced Vector Extensions 512 (Intel® AVX-512) Foundation instructions, Intel® AVX-512 Conflict Detection instructions, Intel® AVX-512 Exponential and Reciprocal instructions, Intel® AVX-512 Prefetch instructions for Intel® processors, and the instructions enabled with CORE-AVX2. Optimizes for Intel® processors that support Intel® AVX-512 instructions.

Libraries Optimized for KNL

- Intel® Math Kernel Library 2017
 - Optimized math functions to enable neural networks (CNN and DNN) for deep learning
- Intel® Distribution for Python
 - NumPy/SciPy/Scikit-Learn/Pandas accelerated with Intel® MKL, Intel® MPI, Intel® TBB, Intel® DAAL
- Intel® Data Analytics Acceleration Library 2017
- Intel® MPI Library
 - Usage of specially optimized memcpy for KNL in MPI Library
 - Tuning of shared memory collectives on single KNL Nodes
- Intel® Integrated Performance Primitives (Intel® IPP)
- Intel® TBB Library

KNL Profiling Available Today

VTune™ Amplifier XE 2016 Update 2

Memory Access analysis

- Shows performance problems by memory hierarchy
- Measure DRAM and MCDRAM bandwidth
- Helps define data structures to allocate to MCDRAM

General Exploration analysis

- Efficiency of code passing through the core pipeline

Scalability analysis with Advanced Hotspots

- Serial vs Parallel time
- MPI and OpenMP imbalance, overhead cost, parallel loop parameters

Custom PMU event collection

Grouping: Bandwidth Domain / Bandwidth Utilization Type / Memory Object / Allocation Stack

Bandwidth Domain / Bandwidth Utilization Type / Memory Object / Allocation Stack	Memory Bound	Loads	Stores	LLC Miss Count	Average Latency (cycles)
DRAM, GB/sec	0.657	125,874,377,622	16,061,040...	130,507,830	40
High	0.750	28,236,084,708	5,014,875,...	75,304,518	91
stream.c:180 (76 MB)		900,002,700	654,009,810	18,301,098	495
stream.c:179 (76 MB)		1,050,003,150	667,210,008	33,301,998	487
stream.c:181 (76 MB)		1,434,004,302	907,213,608	20,101,206	412
Selected 1 row(s):	1.000	126,000,378	21,600,324	300,018	61

OpenMP Analysis. Collection Time[?]: 28.061

Serial Time (outside any parallel region)[?]: 12.203s (43.5%)

Serial Time of your application is high. It directly impacts application Elapsed Time and scalability. Explore options for parallelization, algorithm or microarchitecture tuning of the serial part of the application.

Parallel Region Time[?]: 15.858s (56.5%)

Estimated Ideal Time[?]: 5.005s (17.8%)

OpenMP Potential Gain[?]: 10.853s (38.7%)

The time wasted on load imbalance or parallel work arrangement is significant and negatively impacts the application performance and scalability. Explore OpenMP regions with the highest metric values. Make sure the workload of the regions is enough and the loop schedule is optimal.

Running VTune Amplifier from the command-line

On self-boot KNL machines ensure the `amplxe-cl` command is installed. See the “`amplxe-cl –help`” command for complete details.

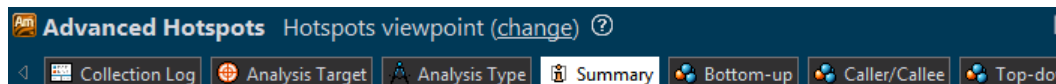
To collect:

- Hotspots:
 - `amplxe-cl –collect advanced-hotspots -- myapp.out`
- General Exploration:
 - `amplxe-cl –collect general-exploration -- myapp.out`
- Memory Access:
 - `amplxe-cl –collect memory-access -- myapp.out`

Advanced Hotspots Analysis

Supports OpenMP* analysis

Stack-sampling is enabled. However, call counts and trip counts are not supported.



Elapsed Time[?]: 4.506s

CPU Time[?]: 566.499s

Effective Time[?]: 195.181s

Spin Time[?]: 352.865s

A significant portion of CPU time is spent waiting. Use this metric to discover which synchronizations are parameters, changing the lock implementation (for example, by backing off then descheduling), or adjust

Overhead Time[?]: 18.453s

Instructions Retired: 320,824,000,000

CPI Rate[?]: 2.642

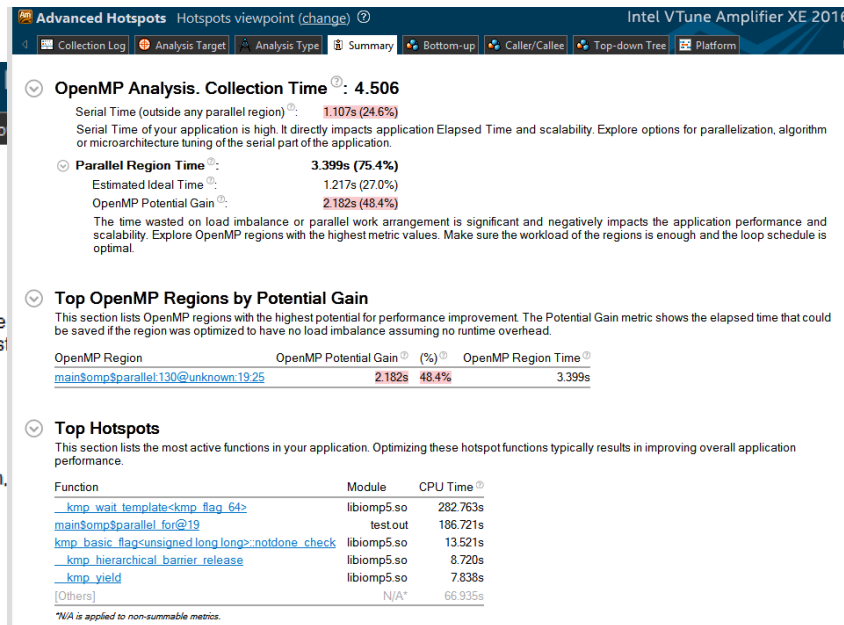
The CPI may be too high. This could be caused by issues such as memory stalls, instruction starvation, instructions. Explore the other hardware-related metrics to identify what is causing high CPI.

CPU Frequency Ratio[?]: 1.071

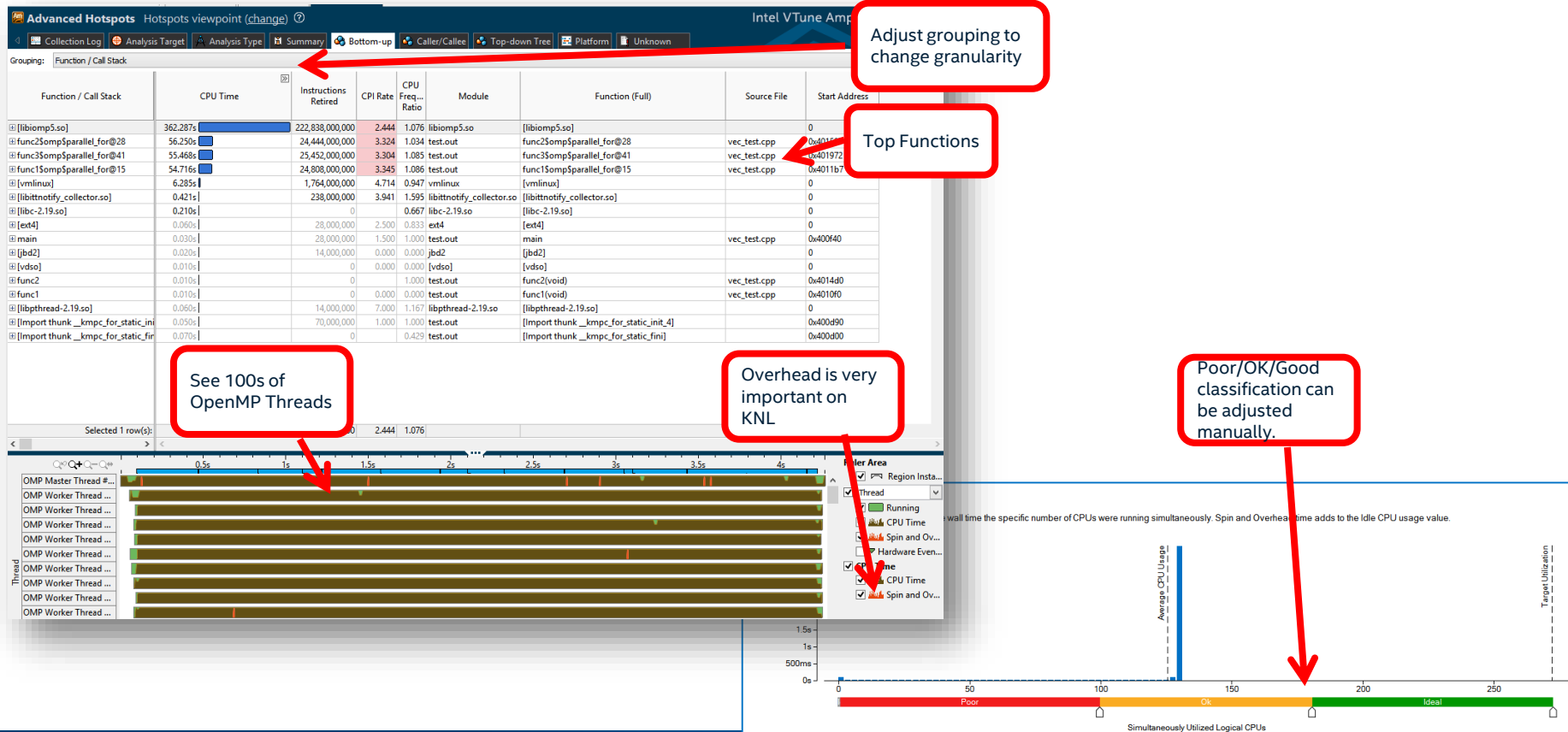
Total Thread Count: 130

Paused Time[?]: 0s

Total Thread Count much higher on Phi. Threading is vital for performance.



Advanced Hotspots Analysis



VTune Amplifier Tips

VTune Finalization:

- Finalization might be slow on KNL. Finalize on Xeon.
- Disable auto finalization with: `-no-auto-finalize`

Large amount of raw data collected:

- Appropriately select the app run duration using: `-target-duration-type=<veryshort/short/medium/long>`
- Change the default data limit as required.

Power throttling:

- Keep an eye on the CPU frequency ratio. If this ratio changes significantly during the run then you might be seeing throttling or turbo effects.

VTune Amplifier Tips (cont.)

Event multiplexing:

- Similar to KNC, KNL has only 2 general purpose counters. Hence, when collecting a large number of events the data might be statistically invalid.
- Try changing the target duration type or allow multiple runs.

High Bandwidth Memory Analysis on KNL

Memory Bandwidth often is a limiting factor of compute intensive applications on multi-core systems

MCDRAM – High Bandwidth Memory with much greater bandwidth speedup to alleviate this problem

Limited MCDRAM size might require selective data object placement to HBM (for flat and hybrid MCDRAM modes)

Memory Access analysis helps to identify memory objects
for HBM placement to benefit the most

HBW Analysis Steps – Case Study

Use Case:

- miniFE* benchmark from Mantevo Suite

Platform:

- KNL with MCDRAM in flat mode

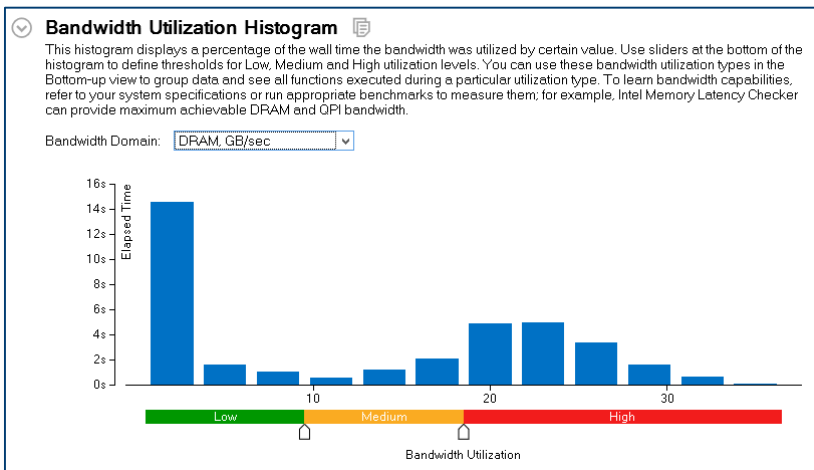
VTune Amplifier XE 2017 Beta with NDA package

HBW Analysis Steps – Case Study. Step 1

Run the original application under Memory Access analysis

Explore DRAM Bandwidth histogram to see if the app is bandwidth bound

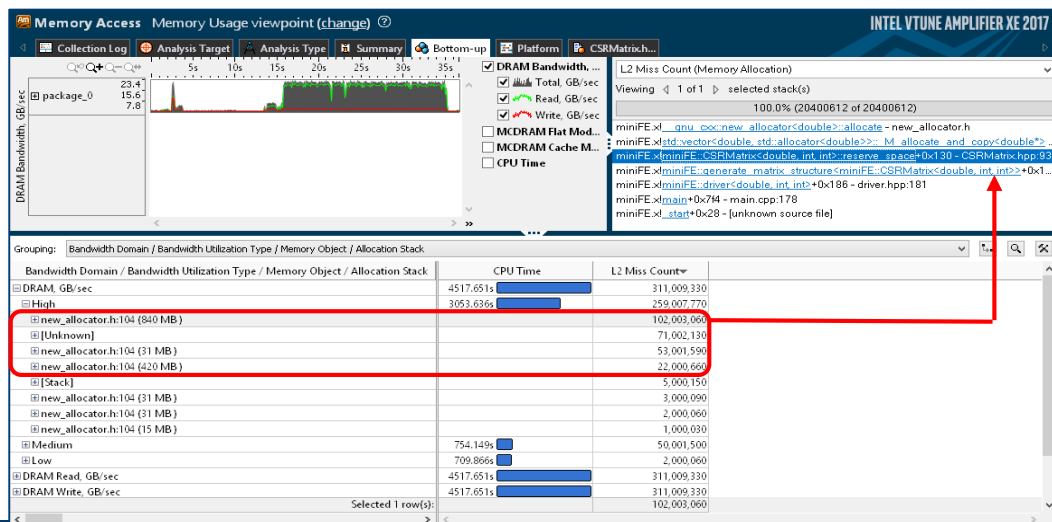
Significant portion
of application time
spent in high memory
bandwidth utilization
The app may benefit
from MCDRAM



HBW Analysis Steps – Case Study. Step 2

Investigate the memory allocations inducing bandwidth

- “Bandwidth Domain/Bandwidth Utilization Type/Memory Object/Allocation Stack” grouping with expansion by “DRAM/High” and sorting by L2 Miss Count



```
void reserve_space(unsigned nrows, unsigned ncols_per_row)
{
    rows.resize(nrows);
    row_offsets.resize(nrows+1);
    packed_cols.reserve(nrows * ncols_per_row);
    packed_coefs.reserve(nrows * ncols_per_row);

    #pragma omp parallel for
    for(MINIFE_GLOBAL_ORDINAL i = 0; i < nrows; ++i) {
        rows[i] = 0;
        row_offsets[i] = 0;
    }
}
```

Focus on allocations inducing L2 misses
Allocation stack shows the allocation place in user's code

HBW Analysis Steps – Case Study. Step 3

Allocate object using High Bandwidth Memory

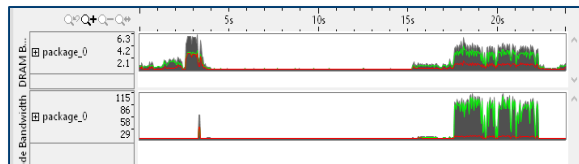
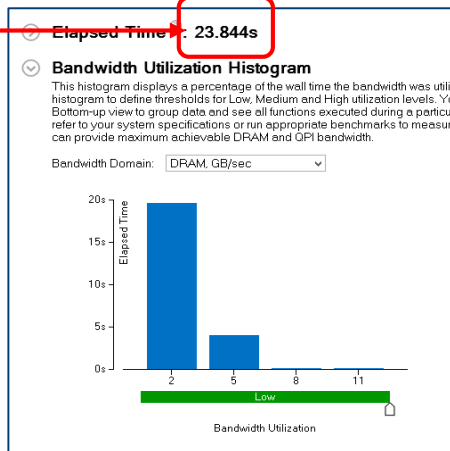
- Specifying a custom memory allocator class for stored vector elements

```
template<typename Scalar,  
        typename LocalOrdinal,  
        typename GlobalOrdinal,  
        typename ComputeNode>  
struct CSRMatrix {  
...  
    typedef Scalar    ScalarType;  
    typedef LocalOrdinal LocalOrdinalType;  
    typedef GlobalOrdinal GlobalOrdinalType;  
    typedef ComputeNode ComputeNodeType;  
  
    bool has_local_indices;  
    std::vector<GlobalOrdinal> rows;  
    std::vector<LocalOrdinal> row_offsets;  
    std::vector<LocalOrdinal> row_offsets_external;  
    std::vector<GlobalOrdinal> packed_cols;  
    std::vector<Scalar> packed_coefs;  
    LocalOrdinal num_cols;  
    ComputeNode& compute_node;  
  
...  
}
```



```
        typename LocalOrdinal,  
        typename GlobalOrdinal,  
        typename ComputeNode>  
struct CSRMatrix {  
...  
    typedef Scalar    ScalarType;  
    typedef LocalOrdinal LocalOrdinalType;  
    typedef GlobalOrdinal GlobalOrdinalType;  
    typedef ComputeNode ComputeNodeType;  
  
    bool has_local_indices;  
    std::vector<GlobalOrdinal, hbwmalloc::hbwmalloc_allocator<GlobalOrdinal> > rows;  
    std::vector<LocalOrdinal, hbwmalloc::hbwmalloc_allocator<GlobalOrdinal> > row_offsets;  
    std::vector<LocalOrdinal, hbwmalloc::hbwmalloc_allocator<GlobalOrdinal> > row_offsets_external;  
    std::vector<GlobalOrdinal, hbwmalloc::hbwmalloc_allocator<GlobalOrdinal> > packed_cols;  
    std::vector<Scalar, hbwmalloc::hbwmalloc_allocator<GlobalOrdinal> > packed_coefs;  
    LocalOrdinal num_cols;  
    ComputeNode& compute_node;  
  
...  
}
```

Rerun the benchmark



Intel and the Intel logo are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States and other countries.

All products, dates, and figures are preliminary and are subject to change without any notice. Copyright © 2016, Intel Corporation.

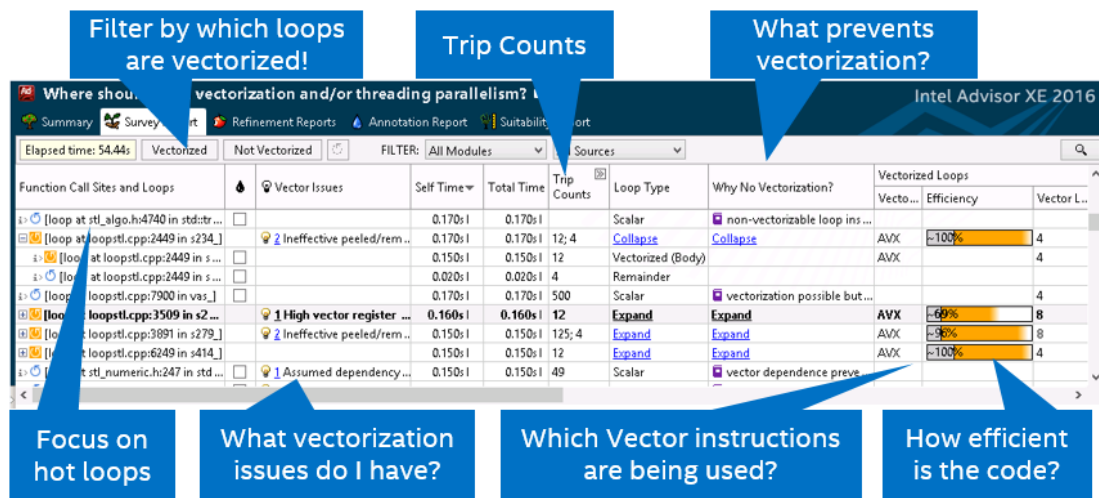
be claimed as the property of others.

Easier Vectorization – Faster Code Faster

Intel® Advisor – Vectorization Advisor – New for 2016

The Right Data At Your Fingertips

- Sorts loops by potential performance gain
- Easy to read compiler reports on your source
- Vectorization tips
- Trip counts
- Dependencies
- Memory access pattern data

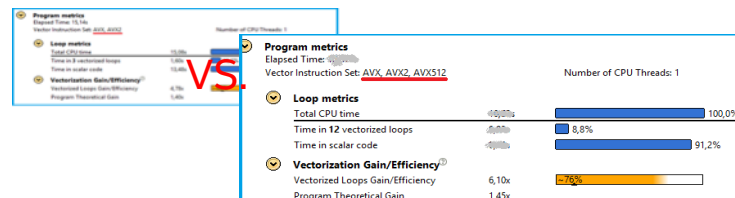


Optimize for AVX-512 Even Without Hardware

Intel® Advisor 2016

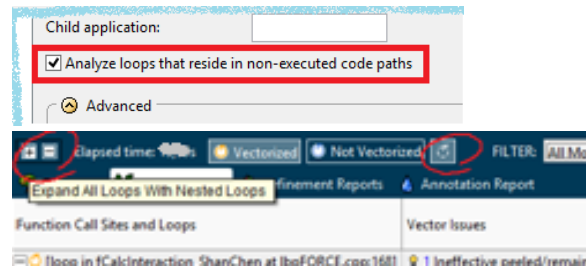
Native Advisor profiling for KNL and future platforms

- Vectorization Summary quickly compares AVX-512 vs. other architecture



Optimize for AVX-512 even without AVX-512 hardware

- Explores AVX-512 code paths characteristics
- Create speed-up estimates for AVX512 (traits for different ISA, optimizations applied by compiler, vector length and etc.)



Survey Report provides AVX-512 Traits

- Advisor highlights AVX-512 instructions which could significantly affect vector code performance (e.g. Compress / Expand, Scatter, Conflict Detection)

The screenshot shows the 'Survey Report' window in Intel Advisor. The 'Function Call Sites and Loops' section is expanded, showing the 'Loop in [source code]' and 'Vectorized (Body)' sections. The 'Compress' section is highlighted, showing the following code snippet:

```
#pragma ivdep
for (i=0; i<BUFF_SIZE; i++)
{
    if ( source[i] > 0 )
    {
        dest[i++] = source[i];
    }
}
```

The 'Traits' column shows 'Mask Manipulations' and 'Compresses'.

Legal Disclaimer & Optimization Notice

INFORMATION IN THIS DOCUMENT IS PROVIDED “AS IS”. NO LICENSE, EXPRESS OR IMPLIED, BY ESTOPPEL OR OTHERWISE, TO ANY INTELLECTUAL PROPERTY RIGHTS IS GRANTED BY THIS DOCUMENT. INTEL ASSUMES NO LIABILITY WHATSOEVER AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY, RELATING TO THIS INFORMATION INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT OR OTHER INTELLECTUAL PROPERTY RIGHT.

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products.

Copyright© 2016, Intel Corporation. All rights reserved. Intel, the Intel logo, Atom, Xeon, Xeon Phi, Core, VTune, and Cilk are trademarks of Intel Corporation in the U.S. and other countries.

Optimization Notice

Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice.

Notice revision #20110804

