# 5 Steps to Efficient Vectorization and Memory utilization: Intel Advisor 2017



**1. Compiler diagnostics + Performance Data + SIMD efficiency information**

**2. Guidance: detect problem and recommend how to fix it**

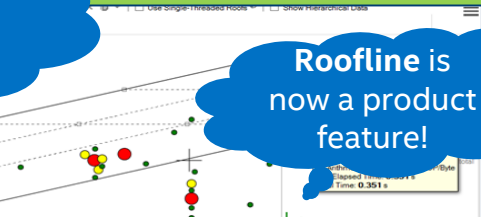**3. "Precise" Trip Counts & FLOPs. Roofline analysis. Characterize your application.**

**4. Loop-Carried Dependency Analysis**

**5. Memory Access Patterns Analysis**

# What's new in "2018" release

# Advisor Survey: **Focus** + **Characterize.**

## **Focus** and order **vectorized loops**



- **Efficiency** – my performance thermometer

- **Recommendations** – get tip on how to improve performance

  - (also apply to scalar loops)

# Data Dependencies – Tough Problem #1
## Is it safe to force the compiler to vectorize?

```
DO I = 1, N
   A(I) = A(I-1) * B(I)
ENDDO
```

or

```
void scale(int *a, int *b)
{
    for (int i = 0; i < 1000; i++)
        b[i] = z * a[i];
}
```

?

**Issue: Assumed dependency present**

The compiler assumed there is an anti-dependency (Write after read – WAR) or true dependency (Read after write – RAW) in the loop. Improve performance by investigating the assumption and handling accordingly.

◇ Enable vectorization
Potential performance gain: Information not available until Beta Update release
Confidence this recommendation applies to your code: Information not available until Beta Update release
The Correctness analysis shows there is no real dependency in the loop for the given workload. Tell the compiler it is safe to vectorize using the `restrict` keyword or a directive.

| ICL/ICC/ICPC Directive | IFORT Directive | Outcome |
|---|---|---|
| #pragma simd or #pragma omp simd | !DIR$ SIMD or !$OMP SIMD | Ignores all dependencies in the loop |
| #pragma ivdep | !DIR$ IVDEP | Ignores only vector dependencies (which is safest) |

**Read More:**

- User and Reference Guide for the Intel C++ Compiler 15.0 > Compiler Reference > Pragmas > Intel-specific Pragma Reference >
  - ivdep
  - omp simd

# Advisor Memory Access Pattern (MAP): know your access pattern



Unit-Stride access

```
for (i=0; i<N; i++)
A[i] = C[i]*D[i]
```

Constant stride access

```
for (i=0; i<N; i++)
point[i].x = x[i]
```

Variable stride access

```
for (i=0; i<N; i++)
A[B[i]] = C[i]*D[i]
```

# Intel® AVX-512 - Comparison

- KNL and future Intel® Xeon® processors share a large set of instructions

- But some sets are not identical

- Subsets are represented by individual feature flags (CPUID)



Intel® microarchitecture code name ...

NHM  SNB  HSW  KNL  Next Intel Xeon

# Advisor 2017: **AVX-512** specific performance insights



- **Native AVX-512 profiling** on KNL

- Precise **FLOPs and Mask Utilization** profiler

- AVX-512 **Advices and "Traits"**

- And more..

  - Performance **Summary for** AVX-512 codes

  - **AVX-512 Gather/Scatter Profiler**

- No access to AVX-512 Hardware yet?

  - Explore AVX-512 code with –axcode flags and new Advisor Survey capability!

# General efficiency (*FLOPS*) vs. VPU-centric efficiency (Vector Efficiency)



High **Vector Efficiency**
Low FLOPS

Low **Vector Efficiency**
High FLOPS

# Survey+FLOPs Report on AVX-512:
## FLOP/s, Bytes and AI, Masks and Efficiency

# Why Mask Utilization Important?

```
for(i = 0; i <= MAX; i++)
    c[i] = a[i] + b[i];
```

**100%**

# Why Mask Utilization Important?

3 elements suppressed

```
for(i = 0; i <= MAX; i++)
    if (cond(i))
        c[i] = a[i] + b[i];
```

SIMD Utilization = 5/8

**62.5%**

# AVX-512 Mask Registers

## 8 Mask registers of size 64-bits

- k1-k7 can be used for predication
  - k0 can be used as a destination or source for mask manipulation operations

## 4 different mask granularities. For instance, at 512b:

- Packed Integer Byte use mask bits [63:0]
  - VPADDB zmm1 {k1}, zmm2, zmm3

- Packed Integer Word use mask bits [31:0]
  - VPADDW zmm1 {k1}, zmm2, zmm3

- Packed IEEE FP32 and Integer Dword use mask bits [15:0]
  - VADDPS zmm1 {k1}, zmm2, zmm3

- Packed IEEE FP64 and Integer Qword use mask bits [7:0]
  - VADDPD zmm1 {k1}, zmm2, zmm3



VADDPD zmm1 {k1}, zmm2, zmm3

| | | Vector Length | | |
|---|---|---|---|---|
| | | 128 | 256 | 512 |
| element size | Byte | 16 | 32 | 64 |
| | Word | 8 | 16 | 32 |
| | Dword/SP | 4 | 8 | 16 |
| | Qword/DP | 2 | 4 | 8 |

# Mask Utilization and FLOPS profiler

- Long-waiting in HPC: accurate HW independent FLOPs measurement tool

- Not just count FLOPs. Has following additions:

  - (AVX-512 only) Mask-aware. Masked-Memory/Unmasked-Compute pattern aware

  - Unique capability to correlate FLOPs with performance data (obtained without instrumentation). Gives FLOPs/s.

- Lightweight instrumentation, PIN-based, benefits from "threadchecker tools" and more generally Advisor framework integration.

# **Characterize and compare**
## AVX-512 against AVX2 versions (on Xeon Phi)

# Highlight "impactful" AVX-512 instructions.
## Survey Static Analysis – AVX-512 *"Traits"*

Presence of  remarkable

performance-impactful

(negative or positive impact)

instructions

| Vectorization Advisor Trait and/or Recommendation | Theoretical Performance Impact Comments | Corresponding AVX-512 Instructions |
|---|---|---|
| Compress / Expand  Trait and Recommendation | >> 4x speedup | `v(p)expand*` `v(p)compress*` |
| Gather / Scatter Trait | Up to 10x slower than contiguous memory access >2x faster than scalar | `v(p)gather*` `v(p)scatter*` |
| Conflict Detection | | `v(p)conflict*` |
| Approximate Reciprocals/Reciprocal SQRT; AVX-512ER | >10x faster than DIV/SQRT | `vrcp*` `vrcsqrt*` `vdiv*` `vsqrt*` |
| Exponent extraction Mantissa extraction Traits | | `vgetexp*` `vgetmant*` |
| L1 (L2) Prefetch L1 (L2) Sparse  prefetch Trait | | `prefetchw*` `vscatterpf*` `vgatherpf*` |

# Highlight "impactful" AVX-512 instructions.
## Survey Static Analysis - AVX-512 *"Traits"*



Summarized Traits in **Survey Report**.

Simplify "performance-aware" reading of **Source and Assembly**

# AVX-512 Also Benefit Scalar Code a lot...
# **Survey** Static Analysis - AVX-512 *"Traits"*

# Gather/Scatter Analysis
## Motivation

AVX-512 Gather/Scatter-based vectorization.

Much wider usage than before :

- Makes much more codes (profitably) vectorizable

- Gives good average performance, but often far from optimal.

**Could be 2x faster than scalar mov**

**Could be 10x slower than vmovp***

# Gather/Scatter Analysis
## Advisor MAP detects gather "offset patterns".



| Pattern # | Pattern Name | Horizontal Stride Value | Vertical Stride Value | Example of Corresponding *Fix(es)* |
|---|---|---|---|---|
| 1 | Invariant | 0 | 0 | OpenMP uniform clause, simd pragma/directive, refactoring |
| 2 | Uniform (horizontal invariant) | 0 | Arbitrary | OpenMP uniform clause, simd pragma/directive |
| 3 | Vertical Invariant | Constant | 0 | OpenMP private clause, simd pragma/directive |
| 4 | Unit | 1 or -1 | \|Vertical Stride\| = Vector Length | OpenMP linear clause, simd pragma/directive |
| 5 | Constant | Constant = X | Constant = X*VectorLength | Subject for AoS -> SoA transformation |

# Gather/scatter issue improvements

Compiler may generate gather/scatter instructions despite regular access pattern. In this case, performance can be improved by refactoring the code.

- Detecting regular patterns taking into account masking instructions
- Added new access pattern for gather profiling – Constant (Non-Unit Stride) with adjusted recommendation to transform AOS to SOA

Recommendation: Refactor code with detected regular stride access patterns          Confidence: Low

The Memory Access Patterns Report shows the following regular stride access(es):

| Variable | Pattern |
|---|---|
| block 0x7f049a6ff010 | Constant (non-unit) |

See details in the Memory Access Patterns Report Source Details view.

To improve memory access: Refactor your code to alert the compiler to a regular stride access. Sometimes, it might be beneficial to use the `ipo/Qipo` compiler option to enable interprocedural optimization (IPO) between files.

An array is the most common type of data structure containing a contiguous collection of data items that can be accessed by an ordinal index. You can organize this data as an array of structures (AoS) or as a structure of arrays (SoA). Detected constant stride might be the result of AoS implementation. While this organization is excellent for encapsulation, it can hinder effective vector processing. To fix: Rewrite code to organize data using SoA instead of AoS.
However, the cost of rewriting code to organize data using SoA instead of AoS may outweigh the benefit. To fix: Use Intel SIMD Data Layout Templates (Intel SDLT), introduced in version 16.1 of the Intel compiler, to mitigate the cost. Intel SDLT is a C++11 template library that may reduce code rewrites to just a few lines.

# AVX-512-specific performance trade-offs
## Advisor AVX-512 Recommendations

Increasing Vector Register Size ->



Increase fraction of time spent in Remainders

| Function Call Sites and Loops | 🔥 | Vector Issues | Self Time▼ | Total Time | Type | Vectori... Vector l... |
|---|---|---|---|---|---|---|
| ⊟ ⟳ [loop in fCollisionBGKShanChen$om ... | ☐ | 💡 1 Ineffective peeled/remainder loop(s ... | 0,110s ❶ | 0,110s ❶ | Vectorized (Remainder; [Body]) | AVX512 |
| ⌄ ⟳ [loop in fCollisionBGKShanChen$o ... | ☐ | | 0,110s ❶ | 0,110s ❶ | Vectorized (Remainder) | AVX512 |
| ⌄ ⟳ [loop in fCollisionBGKShanChen$o ... | ☐ | | n/a | n/a | Vectorized (Body) [Not Executed] | AVX512 |
| ⊟ ⟳ [loop in fGetFracSite at lbpGET.cpp:19 ... | ☐ | 💡 1 Ineffective peeled/remainder loop(s ... | 0,060s ❶ | 0,060s ❶ | Vectorized (Peeled; Remainder; [Body]) | AVX512 |
| ⌄ ⟳ [loop in fGetFracSite at lbpGET.cpp ... | ☐ | | 0,040s ❶ | 0,040s ❶ | Vectorized (Peeled) | AVX512 |
| ⌄ ⟳ [loop in fGetFracSite at lbpGET.cpp ... | ☐ | | 0,020s ❶ | 0,020s ❶ | Vectorized (Remainder) | AVX512 |
| ⌄ ⟳ [loop in fGetFracSite at lbpGET.cpp ... | ☐ | | n/a | n/a | Vectorized (Body) [Not Executed] | AVX512 |
| ⊟ ⟳ [loop in fCalcInteraction_ShanChen a ... | ☐ | 💡 1 Ineffective peeled/remainder loop(s ... | 0,060s ❶ | 0,060s ❶ | Vectorized (Remainder; [Body]) | AVX512 |
| ⌄ ⟳ [loop in fCalcInteraction_ShanChe ... | ☐ | | 0,060s ❶ | 0,060s ❶ | Vectorized (Remainder) | AVX512 |
| ⌄ ⟳ [loop in fCalcInteraction_ShanChe ... | ☐ | | n/a | n/a | Vectorized (Body) [Not Executed] | AVX512 |
| ⊞ ⟳ [loop in fGetOneMassSite at lbpGET.c ... | ☐ | 💡 1 Ineffective peeled/remainder loop(s ... | 0,050s ❶ | 0,050s ❶ | Vectorized (Remainder; [Body]) | AVX512 |
| ⊞ ⟳ [loop in fGetTotMomentSite at lbp ... | ☐ | 💡 1 Ineffective peeled/remainder loo ... | 0,040s ❶ | 0,040s ❶ | Vectorized (Remainder) | AVX512 |
| ⊞ ⟳ [loop in fGetOneDirecSpeedSite at lbp ... | ☐ | 💡 1 Ineffective peeled/remainder loop(s ... | 0,030s ❶ | 0,030s ❶ | Vectorized (Remainder) | AVX512 |
| ⊞ ⟳ [loop in fGetOneMassSite at lbpGET.c ... | ☐ | 💡 1 Ineffective peeled/remainder loop(s ... | 0,030s ❶ | 0,030s ❶ | Vectorized (Remainder) | AVX512 |
| ⊞ ⟳ [loop in fGetOneDirecSpeedSite at lbp ... | ☐ | 💡 1 Ineffective peeled/remainder loop(s ... | 0,020s ❶ | 0,020s ❶ | Vectorized (Remainder) | AVX512 |

# Ineffective masked remainder for AVX512 codes

- Compiler generates vector masked remainder due to the number of iterations (trip count) not being divisible by vector length. In case of executing a few iterations, it is ineffective comparing to scalar versions of the loop.

- Using AVX512 mask profiler and trip-counts data to prove the issue.

Recommendation: Force scalar remainder generation                    Confidence: ⚲Low

The compiler generated a masked vectorized remainder loop that contains too few iterations for efficient vector processing. A scalar loop may be more beneficial. To fix: Force scalar remainder generation using a directive: `#pragma simd novecremainder` or `#pragma vector novecremainder`.

**Example**: Force the compiler to not vectorize the remainder loop

```
void add_floats(float *a, float *b, float *c, float *d, float *e, int n)
{
    int i;
    #pragma simd novecremainder
    for (i=0; i<n; i++)
    {
        a[i] = a[i] + b[i] + c[i] + d[i] + e[i];
    }
}
```

```
#pragma simd reduction(+:mean)
for(int j = 0; j < size; j++) {
        mean += data[order[j]] / N;
        data[order[j]] = 10.f / (j+1);

}
```

*E.g. bad performance if ((size) % (loop_body_vl) == 1), in case of float number it results in 12.5% mask bits utilization only, in addition leads to gathers, scatters...*

**Read More:**

- simd, vector
- Getting Started with Intel Compiler Pragmas and Directives and Vectorization Resources for Intel® Advisor Users

# Start Tuning for AVX-512 without AVX-512 hardware

Intel® Advisor - Vectorization Advisor "axcode feature"

**Multipath Auto-/explicit vectorisation**



CPUID

IA32    SSE4.2    Specialized Path. Set by /Qax option (Linux: –ax)

AVX    SSE3

Default path set by options /arch or /Qx (Linux: -m or –x)

non-intel    intel

Use **–ax** option when compiling to create multiple paths through code

Additional paths can be added by extending the /Qax option e.g. :
/QaxSSE4.2,AVX,SSE3
*(Linux: -axSSE4.2,AVX,SSE3)*

# Viewing non-executed paths

# Start Tuning for AVX-512 without AVX-512 hardware

Intel® Advisor - Vectorization Advisor "axcode feature"

Use –axCOMMON-AVX512 –xAVX compiler flags to generate both code-paths

- AVX(2) code path (executed on Haswell and earlier processors)
- AVX-512 code path for newer hardware

Compare AVX and AVX-512 code characteristics with Intel Advisor



Inserts (AVX2) vs. Gathers (AVX-512)

Speed-up estimate: 13.5x (AVX2) vs. 30.6x (AVX-512)

# From "Old HPC principle" to modern performance model

## "Old" HPC principles:

1. "Balance" principle (e.g. Kung 1986) – hw and software parameters altogether

2. "intensity", "machine balance" -  (FLOP/byte or Byte/FLOP ratio for algorithm or hardware).  E.g. Kennedy, Carr: 1988, 1994: "Improving the Ratio of Memory operations to Floating-Point Operations in Loops ".



## More research catalyzed by memory wall

– 2008, Berkeley: generalized into Roofline Performance Model. Williams, Waterman, Patterson. "Roofline: an insightful visual performance model for multicore"

– 2014: "Cache-aware Roofline model: " Ilic, Pratas, Sousa. INESC-ID/IST, Technical Uni of Lisbon.

# Roofline Performance Model

# Density, Intensity, Machine balance

$$\text{Arithmetic Intensity} = \frac{\text{Total Flops computed}}{\textit{Total Bytes transferred}}$$



0.1-1.0 flops per byte     Typically < 2 flops per byte     O(10) flops per byte

**A r i t h m e t i c   I n t e n s i t y**

SpMV BLAS1,2
Stencils (PDEs)
Lattice Boltzmann Methods
FFTs, Spectral Methods
Dense Linear Algebra (BLAS3)
Particle Methods

O( 1 )     O( log(N) )     O( N )

## OI

$$\text{Operational Intensity} = \frac{\text{Total Flops computed}}{\text{Total Bytes transferred between } \textbf{DRAM (MCDRAM) and LLC}}$$

Implemented in 2017 Update 1

WP

## AI

$$\text{Arithmetic Intensity} = \frac{\text{Total Flops computed}}{\text{Total Bytes transferred between } \textbf{CPU and "memory"}}$$

$$\text{Arithmetic Intensity} = \frac{\text{Total } \textbf{Intops+Flops} \text{ computed}}{\textit{Total Bytes transferred} \text{ between CPU and "memory"}}$$

# Roofline model: Am I bound by VPU/CPU or by Memory?



What makes loops
**A, B, C** different?

# Old approach – pen and paper



Roofline Model of Iso3DFD on Intel Xeon Phi CO-7120P

Run DGEMM

Read the source, count FP ops, loads&stores

Run STREAM

51 adds

4 loads

27 muls

1 store

"3D stencil performance evaluation and auto-tuning on multi and many-core computers", C.Andreolli et.al.

**Cumbersome – but people still did it!**

# Roofline Automation in Intel® Advisor 2017



Each Roof (slope)
Gives peak CPU/Memory throughput
of your **PLATFORM** (benchmarked)

Each Dot
represents loop or function in
**YOUR APPLICATION** (profiled)

**Automatic and integrated – first class citizen in Intel® Advisor**

# Roofline in Intel® Advisor



**Automatic and integrated – first class citizen in Intel® Advisor**

# Find Effective Optimization Strategies

Intel Advisor:  Cache-aware roofline analysis

## Roofs Show Platform Limits

- Memory, cache & compute limits

## Dots Are Loops

- Bigger, red dots take more time so optimization has a bigger impact
- Dots farther from a roof have more room for improvement

## Higher Dot = Higher GFLOPs/sec

- Optimization moves dots up
- Algorithmic changes move dots horizontally



**Which loops should we optimize?**

- A and G have the biggest impact <u>&</u> biggest gap
- B has room to improve, but will have less impact
- E and H are perfectly optimized already

Roofline tutorial video

# Advisor Roofline: under the hood

**Roofline application profile:**

Axis Y: **FLOP/S** = **#FLOP** (mask aware) / **#Seconds**

Axis X: **AI** = **#FLOP** / **#Bytes**

**Seconds**

User-mode **sampling**

*Root access not needed*

**Roofs**

**Microbenchmarks**
Actual peak for the current configuration

Performance = Flops/seconds

AI = Flop/byte



Performance (GFLOPS)

1000 — SP Vector FMA Peak: 449.82 GFLOPS
DP Vector FMA Peak: 221.23 GFLOPS
SP Vector Add Peak: 110.6 GFLOPS
DP Vector Add Peak: 56.21 GFLOPS

Scalar Add Peak: 14.05 GFLOPS

L1 Bandwidth: 1304.97 Gb/sec
L2 Bandwidth: 369.11 Gb/sec
L3 Bandwidth: 180.58 Gb/sec

DRAM Bandwidth: 22.96 Gb/sec

10

1

0.1

0.01

0.01    0.1    1    10
Arithmetic Intens

**#FLOP**

Binary **Instrumentation**
Does not rely on CPU counters

**Bytes**

Binary **Instrumentation**
Counts operands size (not cachelines)

(intel)

# Getting Roofline data in Intel®Advisor

**Run Roofline** ⊘

▶ Collect  📁  ⌨

**1. Survey Target** ⊘

▶ Collect  ⏸  📁  ⌨

**1.1 Find Trip Counts and FLOPS** ⊘

▶ Collect  📁  ⌨

☑ Trip Counts
☑ FLOPS

| FLOP/S = #FLOP/Seconds | Seconds | #FLOP - Mask Utilization - #Bytes |
|---|---|---|
| Step 1: Survey - Non intrusive. *Representative* - Output: Seconds (+much more) | ✔ | |
| Step 2: Trip counts+FLOPS - Precise, instrumentation based - Physically count Num-Instructions - Output: #FLOP, #Bytes | | ✔ |

# Find Effective Optimization Strategies

Intel Advisor:  Cache-aware roofline analysis

## Roofline Performance Insights

- Highlights poor performing loops
- Shows performance "headroom" for each loop
  - Which can be improved
  - Which are worth improving
- Shows likely causes of bottlenecks
- Suggests next optimization steps

# Classical Roofline vs Cache-Aware Roofline

## Classical Roofline Model

AI = # FLOPS / BYTES  (DRAM →)

Bytes out of a level in memory hierarchy are measured in AI

AI depends on problem size

AI is platform dependent

AI depends on cache reuse

## Cache-Aware Roofline Model

AI = # FLOPS / # BYTES ( → CPU)

Bytes into the cpu from all levels in memory hierarchy are measured in AI

AI is independent of problem size

AI is independent of platform

AI is constant for an algorithm

# CARM vs. ORM Roofline flavors

1. Low Ai, "Stream-like" application. Assume it's well vectorized
   - No cache reuse
     - → DRAM bandwidth bound
     - → DRAM AI = L1 AI
2. Implement L2 cache optimization
   - L2 Cache is fully reused, GFLOPS increase
     - → C-A roofline rises up to the L2 bandwidth limit
     - → CI roofline moves to the right because we are doing less loads from DRAM.
3. Implement L1 cache optimization
   - See 2.
     - → By chance, CI roofline seems bound by the scalar add peak

# Example 2: Compute Bound Application

Legend:
- **Cache-Aware** (red)
- **Classical** (blue)
- **Both Equal** (purple)

X-axis: Arithmetic Intensity (flops/byte)
Y-axis: Attainable Performance (Gflops/s)

1. High AI "particle - like" application.
   - No cache reuse again
   - Compute bound but not using vectorization/FMA/both VPUs
2. Implement vectorization
   - Since we are not touching memory, the AI in both C-A and CI roofline does not change
   - We are fully utilizing VPUs → FLOPS increases
3. Implement FMA use

[1] S. Williams et al. *CACM* (2009), crd.lbl.gov/departments/computer-science/PAR/research/roofline

# Is My Application Bound by a Memory Bandwidth or a Compute Peak?



Often it's a combination of the two

- Applications in area 1 are purely memory bandwidth bound
- Applications in area 3 are purely compute bound
- In area 2 we need more information

# Ask Yourself "Why am I Here?" and "Where am I going?"



Usually, it is more complicated...

You won't be on any ceiling. Or if you are, it is kind of coincidence.

BUT - asking the questions
"why am I not on a higher ceiling?"
and "what should I do to reach it?"
is always productive.

# Perform the right optimization for your region
## Roofline: characterization regions



GFLOPS

Scalar **~2.3** Peak GFLOP/sec

L1 <-> core 155 Gb/sec

DRAM per core 3 Gb/sec

logarithmic scale

**L1/L2/LLC/DRAM-bound**
*Investing into Compute peak could be useless*

**L2/LLC/DRAM/Compute-bound**

AI == 1

**Compute-Bound**
*Investing into Cache/DRAM could be useless*

AI (Flop/Byte)

Optimize memory (cache blocking, etc)

Gray area (need more data to determine right strategy)

Optimize compute (threading, vectorization)

# Interpreting Roofline Data

**Final** Limits

(assuming perfect optimization)

Long-term ROI, optimization strategy

**Current** Limits

(what are my current bottlenecks)

Next step, optimization tactics



Compute bound

**Finally compute-bound**

Invest more into effective CPU/VPU (SIMD) optimization

**Finally memory-bound**

Invest more into effective cache utilization

# PIC (PICADOR) plasma simulation use case



Roofline Model

**Surmin, Meyerov, Gonoskov,**
**NN State University & Institute of Applied Physics, Nizhny Novgorod**

# XGC1 is a PIC Code for Tokamak (Edge) Fusion Plasmas (Koskela et all, LBNL, NERSC)



XGC1 Simulation of edge turbulence in the DIII-D tokamak



Unstructured field-aligned mesh in a poloidal domain

# XGC1: Effect of Optimizations on 1st Order B Interpolation



GFLOPs increase, AI decreases
→ Data alignment should be next optimization target

- **Single KNL quadcache node 1 rank, 64 threads.**
- **Data collected with Advisor survey + tripcounts**
- **Inner loops over blocks of particles added**
  - Scalar function
    → vectorized loops
- **Most time-consuming loops above DRAM bandwidth limit**

Total time: 3.5s → 2.1s
Peak GFLOPS: 4.0 → 16.0

# Roofline Analysis to Tune an MRI Image Reconstruction Benchmark
## The 514.pomriq SPEC ACCEL Benchmark

An MRI image reconstruction kernel described in Stone et al. (2008). MRI image reconstruction is a conversion from sampled radio responses to magnetic field gradients. The sample coordinates are in the space of magnetic field gradients, or K-space.

The algorithm examines a large set of input, representing the intended MRI scanning trajectory and the points that will be sampled.

The input to 514.pomriq consists of one file containing the number of K-space values, the number of X-space values, and then the list of K-space coordinates, X-space coordinates, and Phi-field complex values for the K-space samples.

# Hot loop is vectorized



Intel Advisor summary view

1 vectorized loop that we spend 98.8% of our time in

Need more information to see if we can get more performance

# What is our performance?
## Relative to peak system performance



Our hot loop is below the MCDRAM roof

Potential memory bottleneck

# Get detailed Advice from intel® Advisor

Intel® Advisor code analytics

Loop in ComputeQCPU at computeQ.c:65

**4206.254s**
Vectorized (Body)    Total time

**AVX512F_512    1957.548s**
Instruction Set    Self time

▶ Memory 16% (4)
▶ Compute 33% (8)
● Other 51% (12)
Instruction Mix Summary ⓘ

**Traits**
Gathers, FMA, Mask Manipulations

Average Trip Counts: 12500

Instruction Mix ⓘ
Memory: 4   Compute: 8   Other: 12   Number of Vector Registers: 13

Possible inefficient memory access. Gather stride.

## Statistics for FLOPS And Data Transfers

| | | |
|---|---|---|
| **GFLOPS** | 266.242 | Giga Floating-point Operations Per Second Per-loop GFLOPS = Total FLOP / Elapsed Time. Elapsed time is the exclusive (self-time-based) wall time from the beginning to the end of loop/function execution. For single-threaded applications Elapsed time is equal to Self-Time. |
| **AI** | 0.606 | AI - Arithmetic Intesity - Ratio of Floating-point Operations to L1 Transferred Bytes |
| **Mask Utilization** | 100 | Ratio of Utilized Vector Elements to Total Vector Elements |
| **GFLOP** | 4194.304 | Giga Floating-point Operations |
| **FLOP Per Iteration** | 160 | Floating-point Operations Per Loop Iteration |

Data transfers between CPU and memory sub-system (total traffic, including L1, L2, LLC and DRAM traffic)

**Issue: Possible inefficient memory access patterns present**
Inefficient memory access patterns may result in significant vector code execution slowdown or block automatic vectorization by the compiler. Improve performance by investigating.

◎ Recommendation: Confirm inefficient memory access patterns                                        Confidence: ⓘ Need More Data
There is no confirmation inefficient memory access patterns are present. To confirm: Run a Memory Access Patterns analysis.

Recommendations – need more information, confirm inefficient memory access

(intel)

# Irregular access patterns decreases performance!
## Gather profiling

Run Memory Access Pattern Analysis (MAP)

2.2 Check Memory Access Patterns

▶ Collect

-- Nothing to analyze --

0%: percentage of memory instructions with unit stride or stride 0 accesses

Unit stride (stride 1) = Instruction accesses memory that consistently changes by one element from iteration to iteration

Uniform stride (stride 0) = Instruction accesses the same memory from iteration to iteration

50%: percentage of memory instructions with fixed or constant non-unit stride accesses

Constant stride (stride N) = Instruction accesses memory that consistently changes by N elements from iteration to iteration
Example: for the double floating point type, stride 4 means the memory address accessed by this instruction increased by 32 bytes, (4*sizeof(double)) with each iteration

50%: percentage of memory instructions with irregular (variable or random) stride accesses

Irregular stride = Instruction accesses memory addresses that change by an unpredictable number of elements from iteration to iteration
Typically observed for indirect indexed array accesses, for example, a[index[i]]

- gather (irregular) accesses, detected for v(p)gather* instructions on AVX2 Instruction Set Architecture

# Irregular access patterns
## Bad for vectorization performance

**Details View**

🔴 Gather (irregular) access

Operand Size (bits): 32
Operand Type: bit*16;float32*16
Vector Length: 16
Memory access footprint: 3MB

∨ **Gather/scatter details**

**Pattern: "Constant (non-unit)"**

Instruction accesses values with constant offset from the base:
  - stride within instruction = X
  - stride between iterations = X*vector length

Horizontal stride (bytes): 16
Vertical stride (bytes): 256

Mask is constant
Mask: [1111111111111111]
Active elements in the mask: 100.0%

∨ **Variable references**
Names: block 0x7f0045867010 allocated at main.c:99

Hint: use the Intel Advisor details!

Specific recommendation for your application

**Issue: Inefficient gather/scatter instructions present**

The compiler assumes indirect or irregular stride access to data used for vector operations. Improve memory access by alerting the compiler to detected regular stride access patterns, such as:

| Pattern | Description |
|---|---|
| Invariant | The instruction accesses values in the same memory throughout the loop. |
| Uniform (Horizontal Invariant) | The instruction accesses values in the same memory within the vector iteration. |
| Vertical Invariant | The instruction accesses the memory locations using the same offset across all vector iterations. |
| Unit | The instruction accesses values in contiguous memory throughout the loop, and the stride between vector iterations = vector length. |

◇ Recommendation: Refactor code with detected regular stride access patterns                                    Confidence: ⚲ Low

The Memory Access Patterns Report shows the following regular stride access(es):

(intel)

# Remove gather instructions

step #1 – use newer version of the intel compiler can recognize the access pattern



Gathers replacement is performed by the "Gather to Shuffle/Permutes" compiler transformation

Loop in ComputeQCPU at computeQ.c:65

⟳

**3545.097s**
Vectorized (Body)    Total time

**AVX512F_512    1444.120s**
Instruction Set    Self time

Traits ⌐

2-Source Permutes
Blends ⌐
FMA

Removed gathers

Average Trip Counts: 12500

**Code Optimizations**    ⌃

Compiler: Intel(R) C Intel(R) 64 Compiler for applications running on Intel(R) 64,

Compiler estimated gain: 18.44x

**Code Optimizations Applied By Compiler During Vectorization:**
- Cost Model Was Ignored
- Dependency Analysis Was Ignored
- SIMD

...tistics for FLOPS And Data Transfers

| | | |
|---|---|---|
| GFLOPS | 342.671 | Giga Floating-point Operations Per Second Per-loop GFLOPS = Total FLOP / Elapsed Time. Elapsed time is the exclusive (self-time-based) wall time from the beginning to the end of loop/function execution. For single-threaded applications Elapsed time is equal to Self-Time. |
| AI | 0.08 | AI - Arithmetic Intesity - Ratio of Floating-point Operations to L1 Transferred Bytes |
| Mask Utilization | | Ratio of Utilized Vector Elements to Total Vector Elements |
| GFLOP | 4194 | Giga Floating-point Operations |
| FLOP Per Iteration | | Floating-point Operations Per Loop Iteration |

Data transfers between CPU an...    ...sub-system (total traffic, including L1, L2, LLC and DRAM traffic)

| | | |
|---|---|---|
| in Giga Bytes | 10276. | |
| in Giga Bytes Per Second | 839. | |
| in Bytes Per Loop Iteration | | |

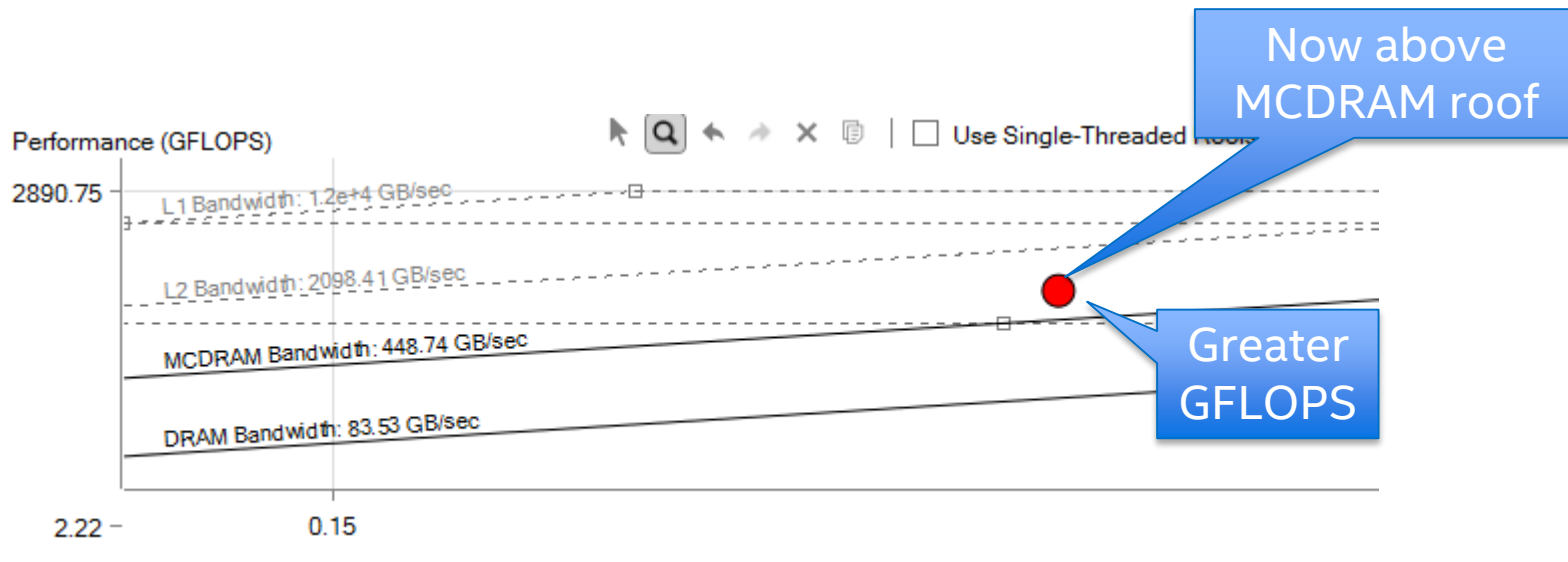Increased GFLOPS (from 266.42 to 342.67)

(intel)

# Remove gather instructions
step #1 – newer version of the intel compiler can recognize the access pattern

# Remove gather instructions

## step #2 – Use structure of arrays instead of array of structures T

```cpp
struct kValues {
  float Kx;
  float Ky;
  float Kz;
  float PhiMag;
};

SDLT_PRIMITIVE(kValues, Kx, Ky, Kz, PhiMag)


sdlt::soa1d_container<kValues> inputKValues(numK);
auto kValues = inputKValues.access();

  for (k = 0; k < numK; k++) {
    kValues [k].Kx() = kx[k];
    kValues [k].Ky() = ky[k];
    kValues [k].Kz() = kz[k];
    kValues [k].PhiMag() = phiMag[k];
  }

auto kVals = inputKValues.const_access();
#pragma omp simd private(expArg, cosArg, sinArg) reduction(+:QrSum, QiSum)
    for (indexK = 0; indexK < numK; indexK++) {
      expArg = PIx2 * (kVals[indexK].Kx() * x[indexX] +
      kVals[indexK].Ky() * y[indexX] +
      kVals[indexK].Kz() * z[indexX]);

      cosArg = cosf(expArg);
      sinArg = sinf(expArg);

      float phi = kVals[indexK].PhiMag();
      QrSum += phi * cosArg;
      QiSum += phi * sinArg;
    }
```

This is a classic vectorization efficiency strategy

But it can yield poorly designed code
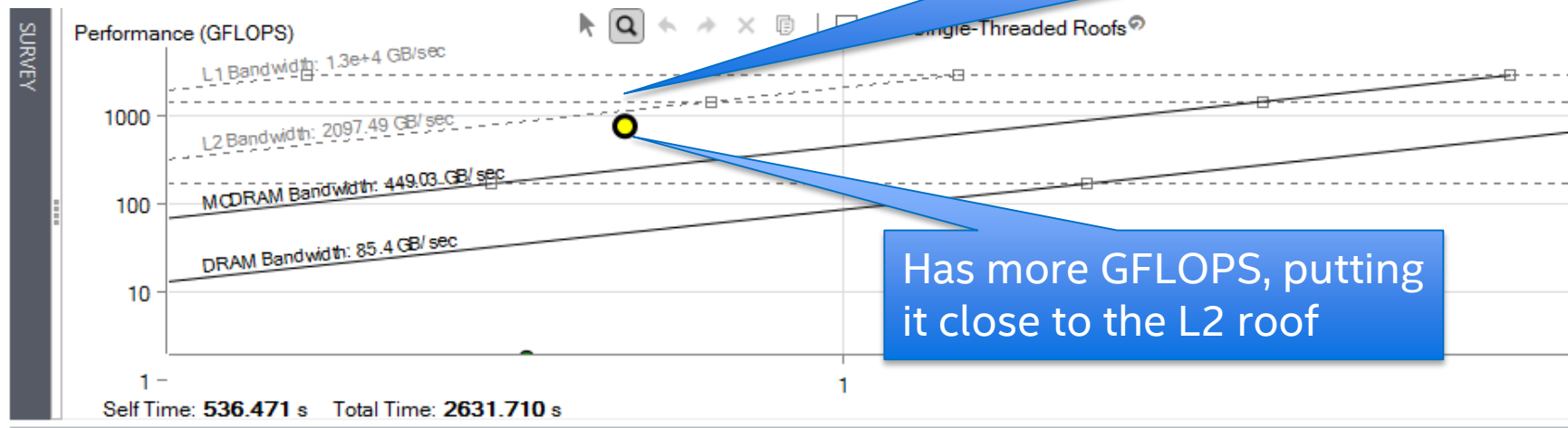
Intel® SIMD Data Layout Templates makes this transformation easy and painless!

# Remove gather instructions
## step #2 – Transform code using the Intel® SIMD Data Layout Templates



The loop is no longer red. This means it takes less time now

Has more GFLOPS, putting it close to the L2 roof

The total performance improvement is almost 3x for the kernel and 50% for the entire application.

# Roofline access and how-to **command line** example

> **source advixe-vars.sh**

> advixe-cl **--collect survey** --project-dir ./your_project --

<your-executable-with-parameters>

> advixe-cl **--collect tripcounts -flops-and-masks** --project-dir
./your_project -- <your-executable-with-parameters>

> **advixe-gui** ./your_project

FLOP/S =
#FLOP/Seconds

1st pass
Obtain "Seconds"
1.1x overhead

2nd pass
Obtain #FLOP count:
3x-5x overhead

Launch GUI

# MPI example (slurm)

1st step:

**srun** –n <num-of-ranks> -c <num_of_cores_per_rank> **advixe-cl** –v -**collect survey** -project-dir=<same_dir_name> -data-limit=0 <your_executable>

2nd step:

**srun** –n <num-of-ranks> -c <num_of_cores_per_rank> **advixe-cl** –v -**collect tripcounts –flops-and-masks** –project-dir=<same_dir_name> -data-limit=0 <your_executable>

# Observe slower Survey analysis or "finalization"?
(1.5x analysis slow-down or more)

Change default call stacks processing mode (***especially for Fortran***)

advixe-cl –collect survey –stackwalk-mode=online –no-stack-stitching

Disable system modules and non-interesting modules processing:

advixe-cl –collect survey -module-filter-mode=include -module-filter=**foo.so**

# Observe slow tripcounts/FLOP analysis ??

( > 8x slower than native and more )

Consider combinations:

1. FLOPS **only**, no TripCounts:

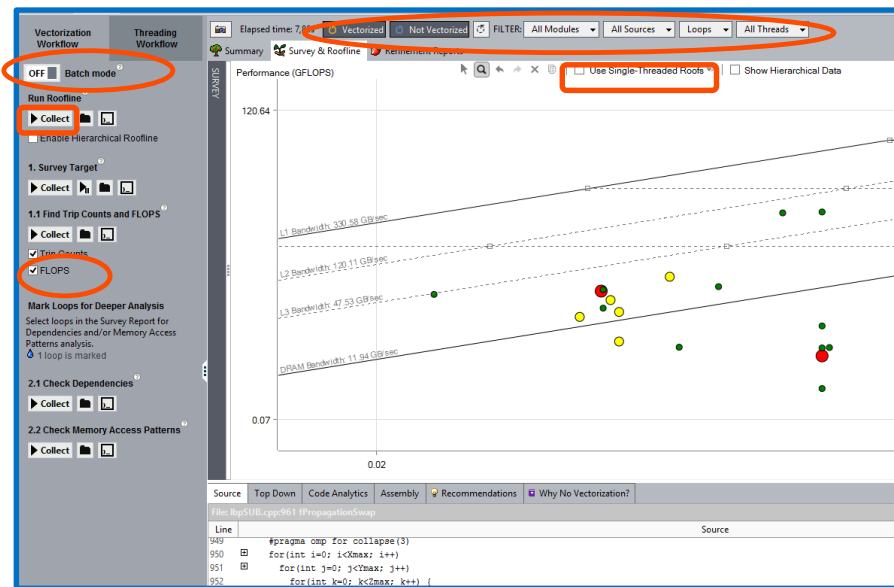    advixe-cl -collect tripcounts –flops-and-masks –no-trip-counts

2. no FLOPS , TripCounts **only**, (->No Roofline):

    advixe-cl -collect tripcounts

3. FLOPS **and** TripCounts :

    advixe-cl -collect tripcounts –flops-and-masks
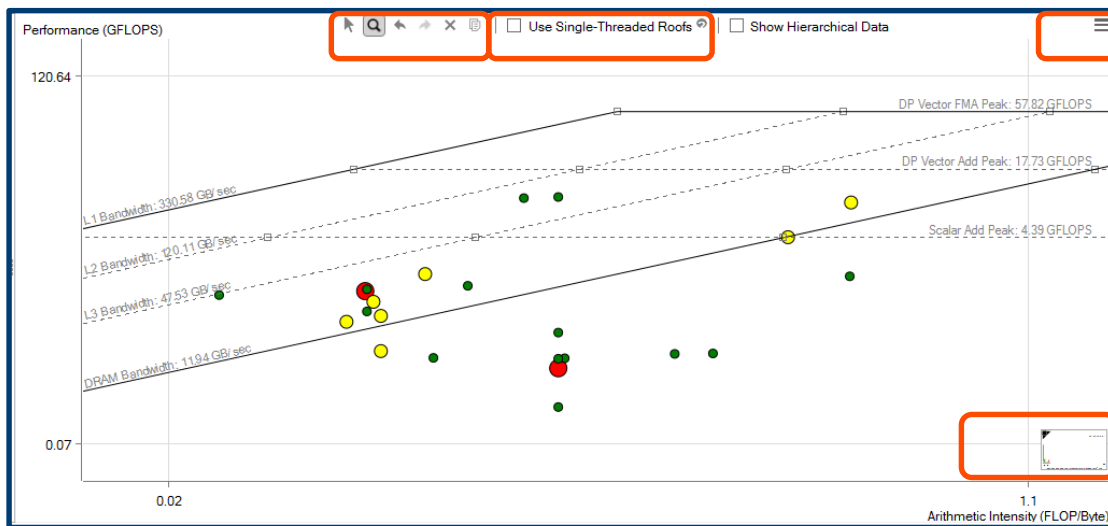
# Roofline GUI access and how-to: **GUI**



1) "**Run Roofline**": most automated way.

2) You can also use **two separate runs**:

1. Survey

2. TripCounts (remember to switch **FLOPs** ON)
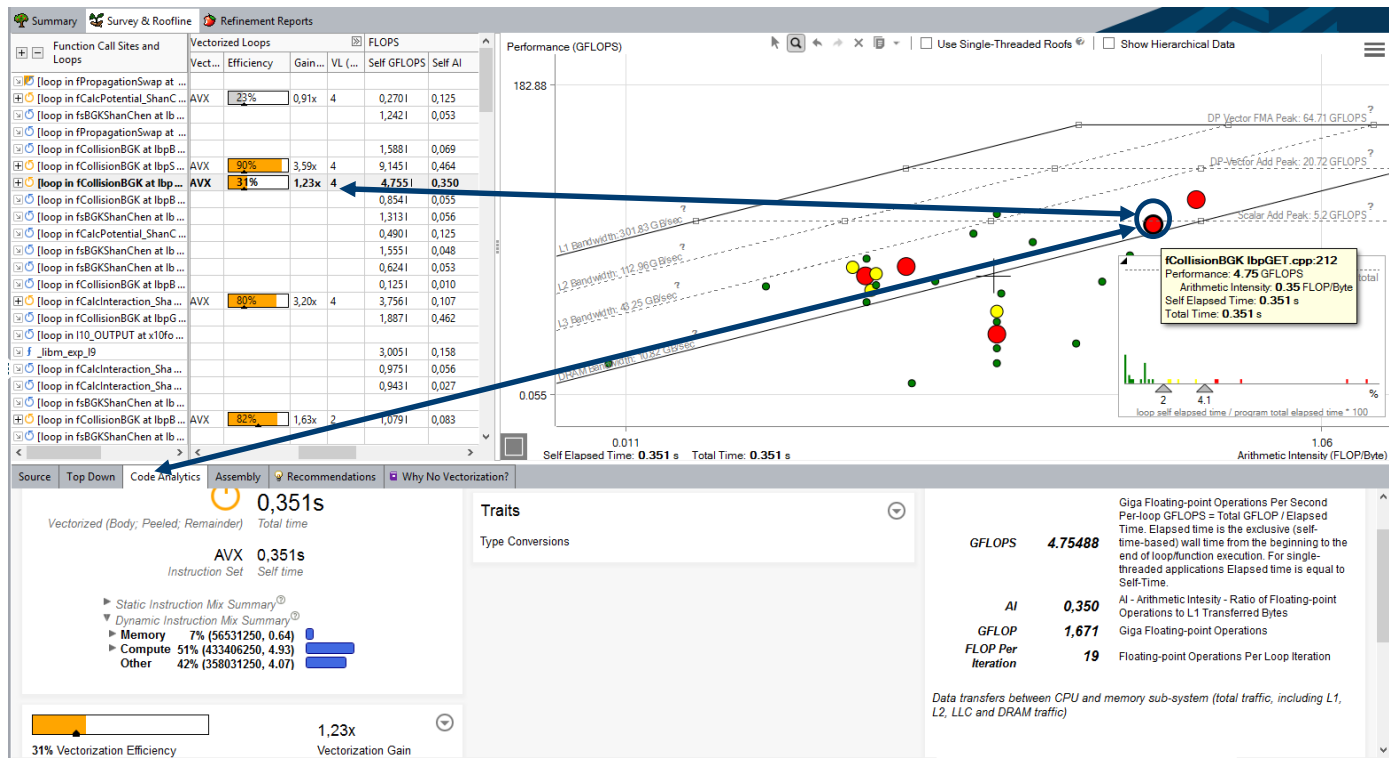
3) **Batch** Mode
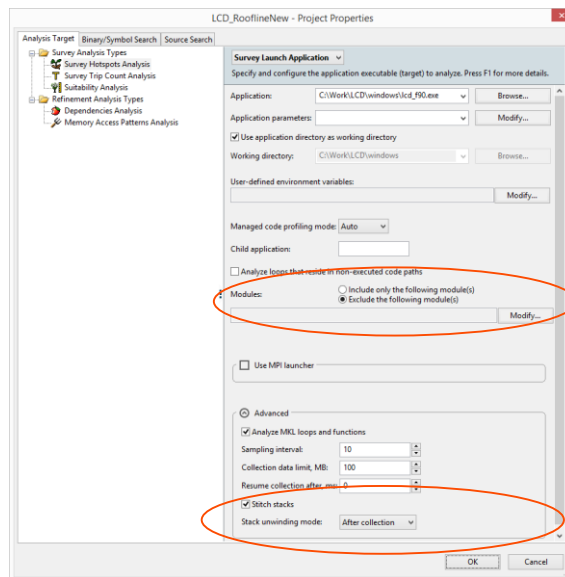
# Roofline Chart
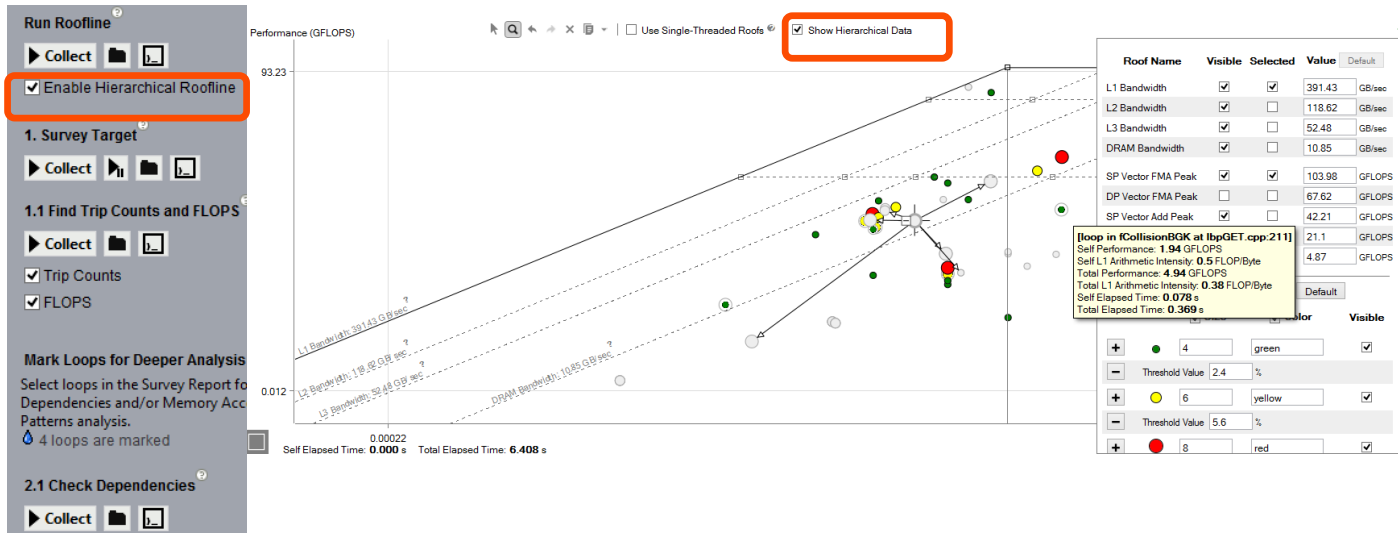
# Use Vectorization and Roofline views together

# Observe slower Survey analysis or "finalization"?

(1.5x slower than native run and more )

Configuration via GUI:

# Hierarchical (top-down) Roofline: new in 2018 release



```
export ADVIXE_EXPERIMENTAL=roofline_ex
```

# Hierarchical Roofline (based on stacks w/ FLOPS )

```
> source advixe-vars.sh
```

```
> export ADVIXE_EXPERIMENTAL=roofline_ex
```

```
> advixe-cl --collect survey --project-dir ./your_project -- <your-executable-with-
parameters>
```

```
> advixe-cl --collect tripcounts -flops-and-masks -callstack-flops --project-dir
./your_project -- <your-executable-with-parameters>
```

```
> export ADVIXE_EXPERIMENTAL=roofline_ex
```

```
> advixe-gui ./your_project
```

# OpenLAB location of Advisor 2017 Update 3

Update 3:

 **/oplashare/sw/Intel/advisor_2017_update**
                              accessible from openlab machines


Also consider installing advisor on your local laptop. Just copy advisor*.tar.gz
from  **/oplashare/sw/Intel/advisor_2017_update/** to your laptop, unpack, run advixe-genvars.sh


You'll need to point $INTEL_LICENSE_FILE to license server in openlab

# BACK-UP

# Why Do We Need the Roofline Model?

## Need a sense of absolute performance when optimizing applications

- How do I know if my performance is good?

- Why am I not getting peak performance of the platform?

## Many potential optimization directions

- How do I know which one to apply?

- What is the limiting factor in my app's performance?

- How do I know when to stop?