



OPTIMIZE FOR BOTH MEMORY AND COMPUTE ON MODERN HARDWARE USING ROOFLINE MODEL AUTOMATION IN INTEL® ADVISOR

Cedric Andreolli, Sr. Support engineer, Intel IAGS
Zakhar Matveev, Product architect, Intel IAGS
October 10th 2019

Acknowledgments

Cedric Andreoli (Intel),

Sam Williams (LBNL), Aleksandar Ilic (INESC),

Kate Antakova (Intel), Philippe Thierry (Intel)

The Roofline Model

Flop: How many FP operations is your kernel performing

Byte transferred: How many bytes need to be transferred

Performance is expressed in Flop/s

- Number of operations per second
- We usually want to maximize this number

Arithmetic intensity

- Ratio between flops and bytes
- Directly impacts the performance (memory bound/compute bound)

2 Operations

The diagram shows the equation $a_i = b_i + c_i * d_i$. Above the equation, the text "2 Operations" is written in red. Two red arrows point from this text to the addition and multiplication operators. Below the equation, four blue arrows point from a common point to each of the four variables: a_i , b_i , c_i , and d_i .

$$a_i = b_i + c_i * d_i$$

$1W+3R = 4*4\text{bytes} = 16 \text{ bytes}$

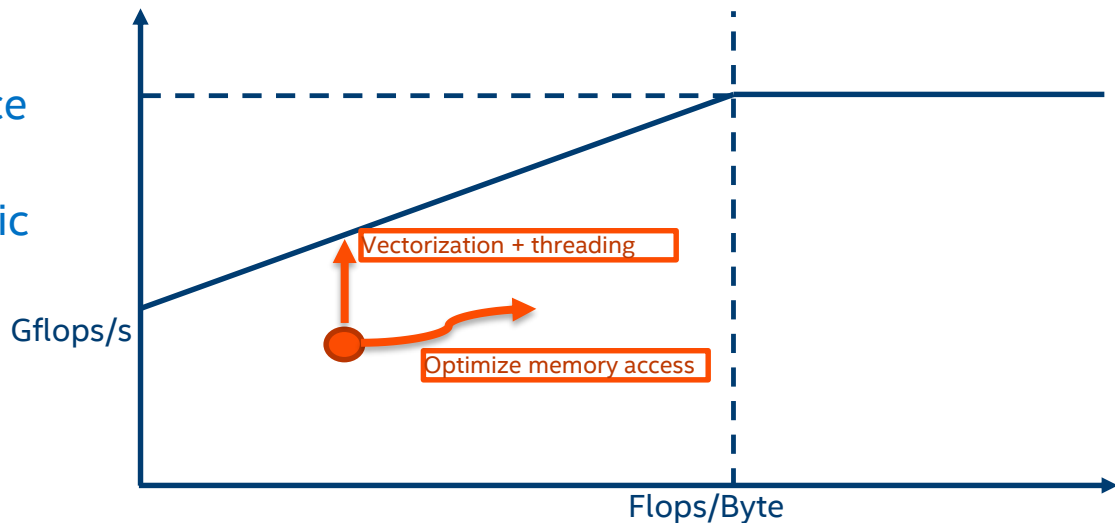
Building a Roofline Model

Roofline model is based on the formula

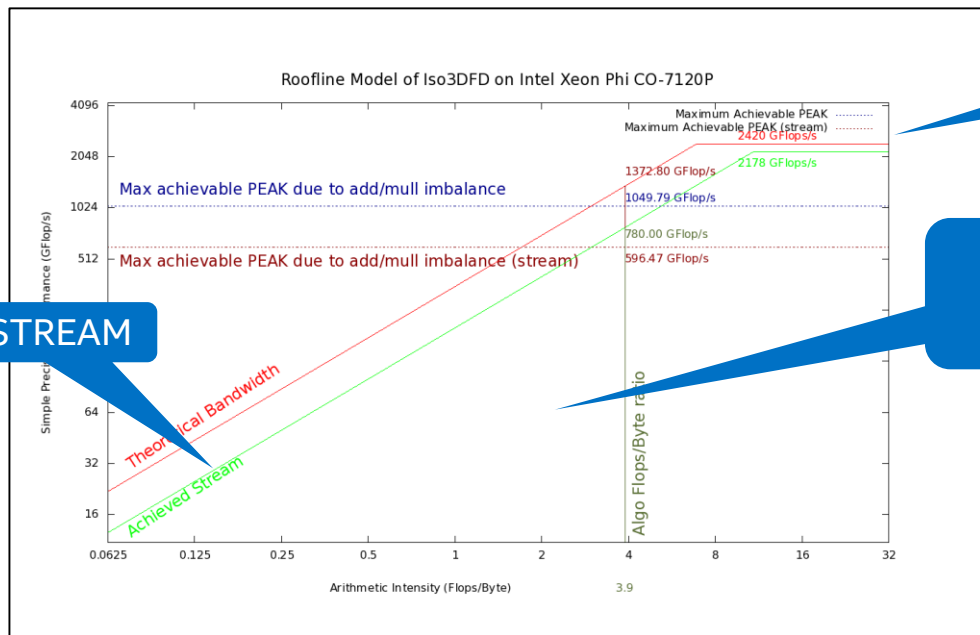
$$\text{Gflop/s} = \min \left\{ \begin{array}{l} \text{Platform PEAK} \\ \text{Platform BW} * \text{AI} \end{array} \right.$$

Roofline Axis

- Vertical axis is the performance in Gflop/s
- Horizontal axis is the arithmetic intensity (Flop/Byte)
- Using log scales



Old approach – pen and paper to get Roofline data



Run STREAM

Run DGEMM

Read the source,
count FP ops,
loads&stores

4 loads

51 adds

27 muls

1 store

```
for(int bz=HALF_LENGTH; bz<n3; bz+=n3_Tblock)
  for(int by=0; by<n2; by+=n2_Tblock)
    for(int bx=0; bx<n1; bx+=n1_Tblock)
      MIN(bz+n3_Tblock, n3);
      MIN(by+n2_Tblock, n2);
      MIN(n1_Tblock, n1-bx);
      for(int iz=0; iz<izEnd; iz++) {
        ptr_prev_base = ptr_prev_base + iz*nln2 + iy*n1 + bx;
        ptr_next_base = ptr_next_base + iz*nln2 + iy*n1 + bx;
        float* vel = ptr_vel_base + iz*nln2 + iy*n1 + bx;
        for(int ix=0; ix<ixEnd; ix++) {
          value = 0.0;
          value += coeff[ix]*coeff[0];
          for(int ir=1; ir<=HALF_LENGTH; ir++) {
            value += coeff[ir] * (prev[ix + ir] + prev[ix - ir]);
          }
          value += coeff[ir] * (prev[ix + ir*n1] + prev[ix - ir*n1]);
          value += coeff[ir] * (prev[ix + ir*nln2] + prev[ix - ir*nln2]);
          next[ix] = 2.0f* prev[ix] - next[ix] + value*vel[ix];
        }
      }
    }
```

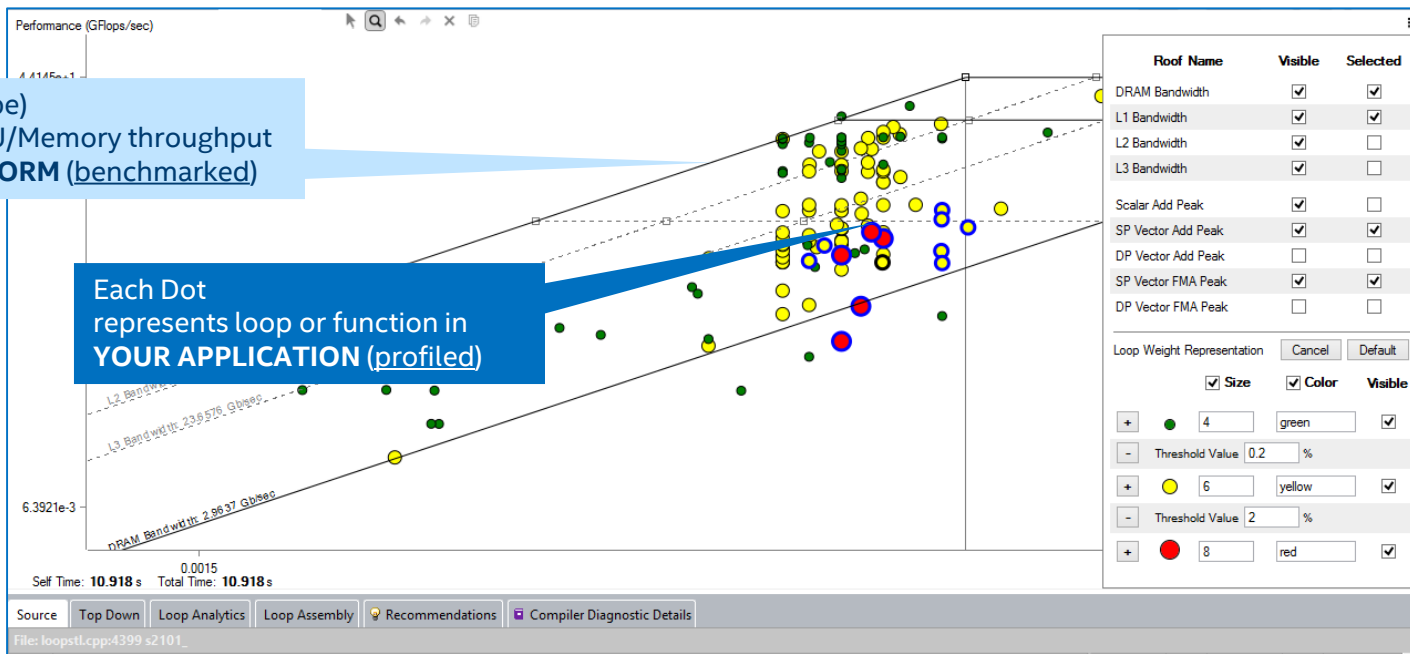
“3D stencil performance evaluation and auto-tuning on multi and many-core computers”, C.Andreolli et.al.

Cumbersome – but people still did it!

("Cache-aware") Roofline Automation , since 2016

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

Questions to answer with Roofline: for your loops / functions

1

Am I doing well? How far am I from the pick?

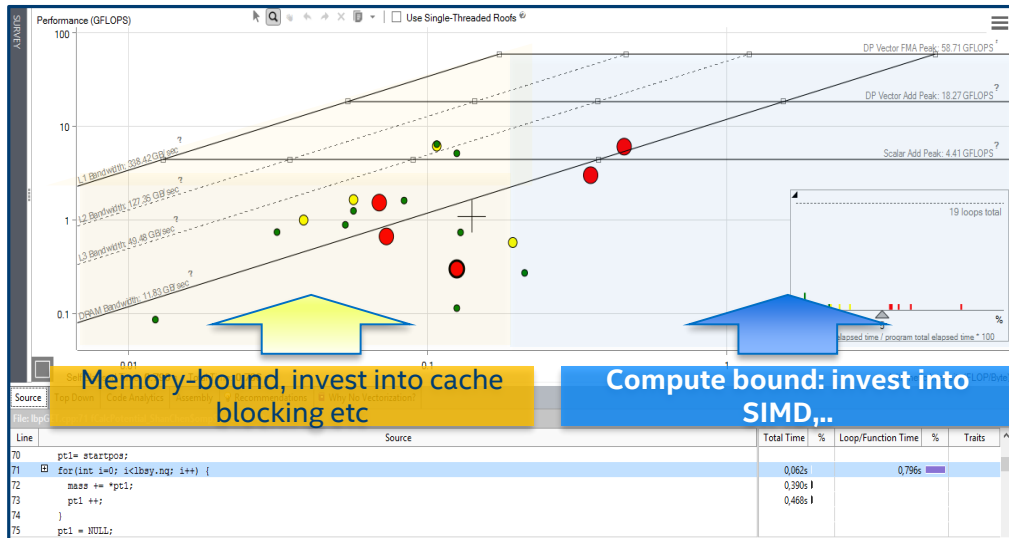
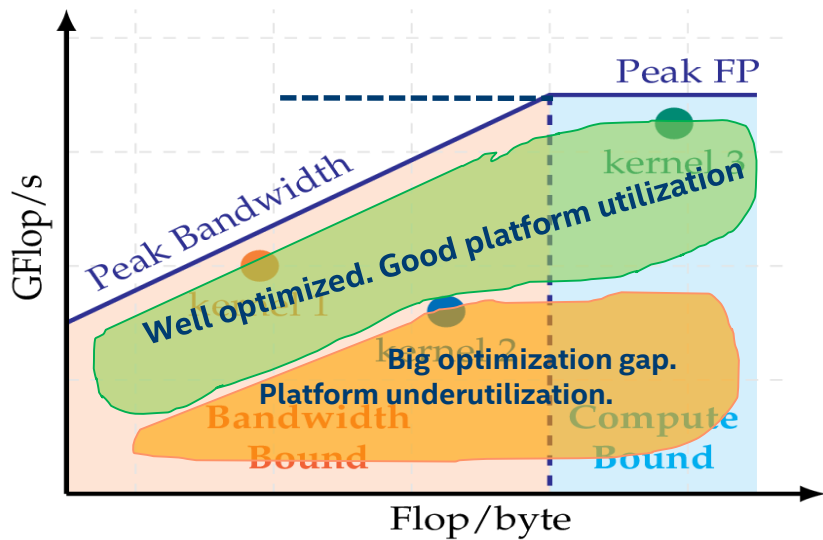
(do I utilize hardware well or not?)

2

Final Bottleneck?

(where will be my limit after I done all optimizations?)

Long-term ROI, optimization strategy



The Roofline Model in Intel® Advisor

First Implementation: Cache-Aware Roofline Model (CARM)

- Based on instrumentation
- 2 runs: one for sampling and timing loops and functions (low overhead); second one for instrumentation
- Algorithmic version of the Roofline Model, optimization usually does not impact AI
- 😊 **Really powerful to characterize an algorithm**
- ☹️ **Not always easy to find your current #1 bottleneck**

Currently exposed as a technical
preview feature, need:
export
ADVIXE_EXPERIMENTAL=int_roofline

Additional Implementation: Integrated or « mutli-level » Roofline Model (IRM)

- Based on cache simulation, evaluate the traffic between each memory subsystem (L1/L2/LLC/DRAM)
- 😊 **Incorporates the « original Roofline model », provides meaningful information for the #1 bottleneck improvement**
- ☹️ **Requires more time to run and more visual aids to comprehend**

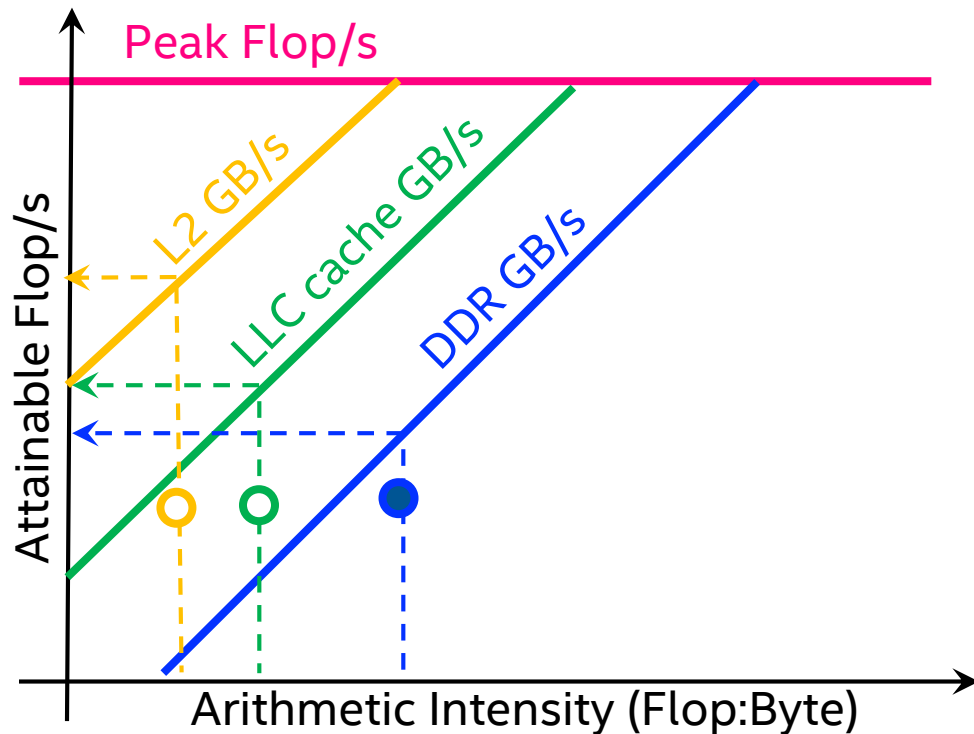
3 Integrated Roofline. What is my current limit?

Performance is limited by minimum of intercepts:

- L1,
- L2,
- LLC,
- DRAM,
- CPU

In this case: *by DRAM*

Hierarchical Roofline for one kernel

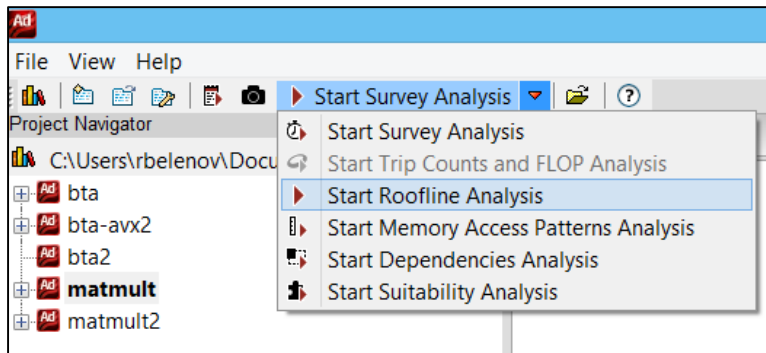




GETTING ROOFLINE PROFILE

Running roofline collection

GUI



Command line

```
> source advise-vars.sh  
> advise-cl -collect roofline -project-dir <advisor project>  
-- <your app>
```

- Simple click or single command
- **Two collections** are launched one after another
 - Low-overhead sampling to get execution times
 - Binary instrumentation for FLOPs and traffic
- Roofs are obtained via benchmarks

Roofline collection – advanced

Launch two collections explicitly – may be required for MPI

```
> advixe-cl -project-dir <your project> -collect survey -- <your app>  
> advixe-cl -project-dir <your project> -collect tripcounts -- <your app>
```

Customize cache configuration

```
> advixe-cl -project-dir <your project> -collect roofline  
--cache-config 4:1024k:32w -- <your app>
```

- “What if” analysis/projection
- More info in backup slides

How to generate **IRM** Roofline profile?*

```
$ source advixe-vars.sh
```

```
$ export ADVIXE_EXPERIMENTAL=int_roofline
```

1st method. Not compatible with MPI applications :

```
$ advixe-cl -collect roofline -enable-cache-simulation --project-dir  
./your_project -- <your-executable-with-parameters>
```

2nd method (compatible with MPI applications and more flexible):

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

```
$ advixe-cl -collect tripcounts -enable-cache-simulation -flop --project-dir  
./your_project -- <your-executable-with-parameters>
```

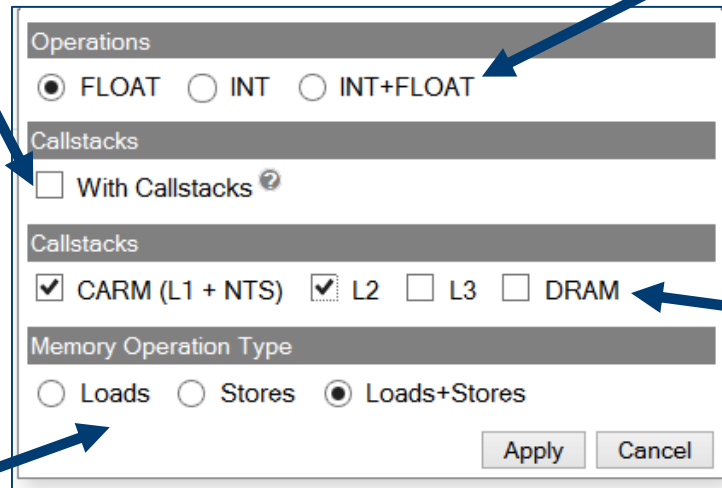
(optional) copy data to your UI desktop system

```
$ advixe-gui ./your_project
```

IRM How-to: <https://software.intel.com/en-us/articles/integrated-roofline-model-with-intel-advisor>

Chart configuration

Aggregate data over calltree
(currently works only for CARM, not integrated roofline!)



The screenshot shows a configuration window with the following sections and options:

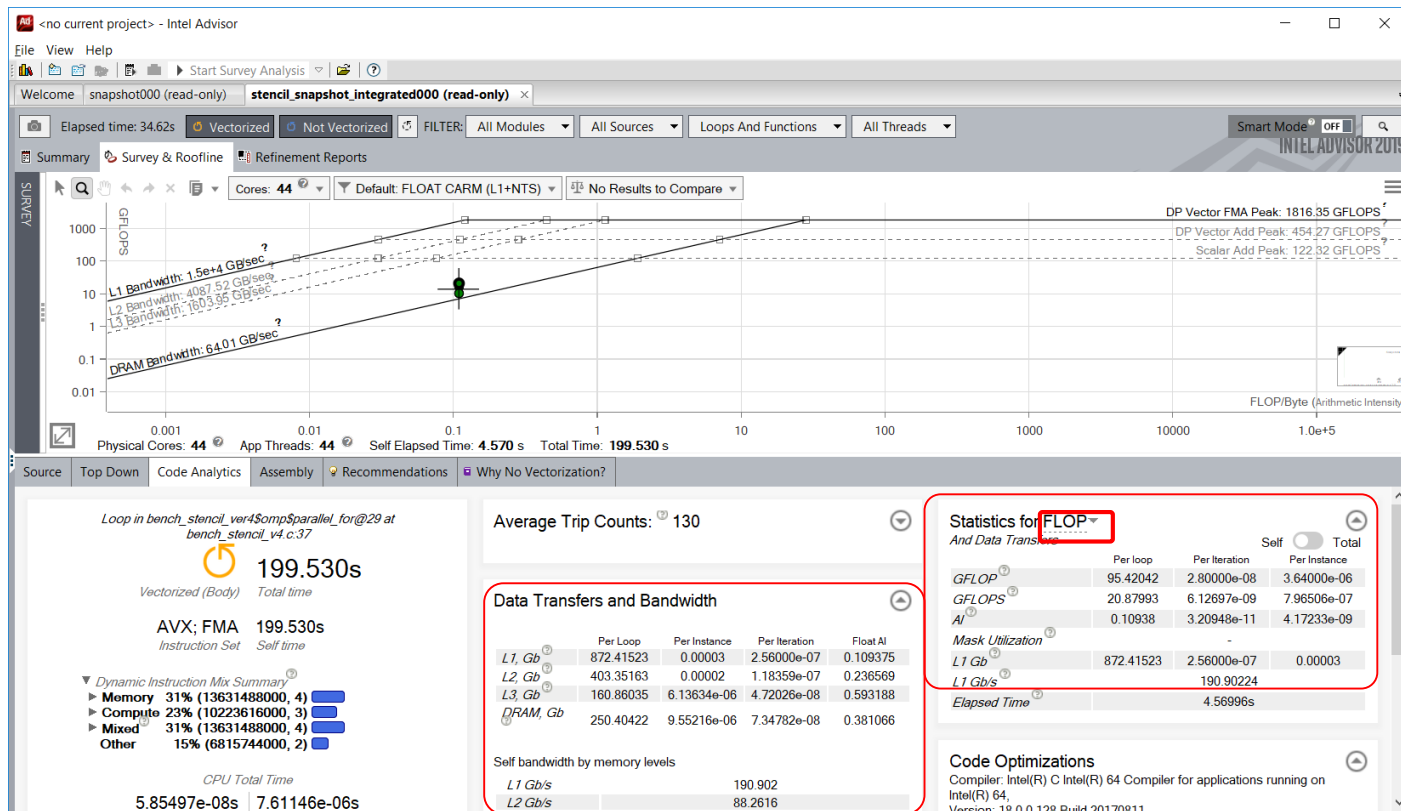
- Operations:** Radio buttons for ☒ FLOAT, ☐ INT, and ☐ INT+FLOAT. An arrow points to this section from the text "Select which operations are counted".
- Callstacks:** A checkbox for ☐ With Callstacks. An arrow points to this section from the text "Aggregate data over calltree (currently works only for CARM, not integrated roofline!)".
- Callstacks:** Checkboxes for ☒ CARM (L1 + NTS), ☒ L2, ☐ L3, and ☐ DRAM. An arrow points to this section from the text "Select memory levels".
- Memory Operation Type:** Radio buttons for ☐ Loads, ☐ Stores, and ☒ Loads+Stores. An arrow points to this section from the text "Select only loads or stores".
- Buttons:** "Apply" and "Cancel" buttons at the bottom right.

Select which operations are counted

Select memory levels

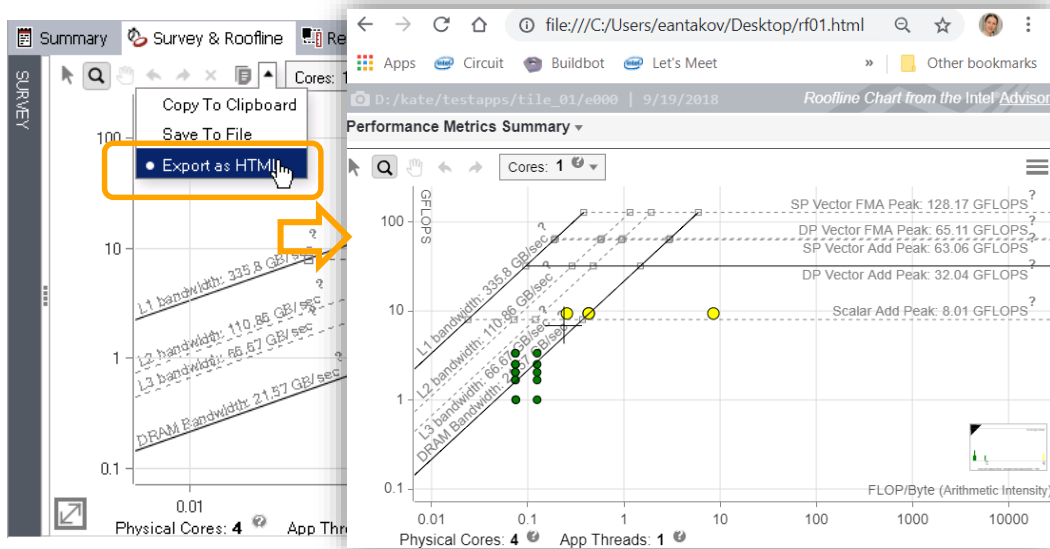
Select only loads or stores

Use “Code Analytics” in conjunction with Roofline



Exporting Integer and Integrated Roofline as HTML

GUI: Use Export as HTML button



Command line:

export
ADVIXE_EXPERIMENTAL=int_roofline

advixe-cl --report roofline
-data-type=float
-memory-level=L2
-memory-operation-type=load
-project-dir /path/to/project/dir

Possible

data types: float, int, mixed
memory levels: L1, L2, L3, DRAM
memory operation types: load, store, all

- Export Roofline from command line does not need GUI sub-system on clusters
- Useful for rooflines quick exchange

Roofs configuration (and more...)

Customize roof calculation:

- Fully loaded system assumed by default
- Tick if most cores are idle

Change hotspot classification

Roofs Settings

☐ Use single-threaded benchmark results to build roofs

| Roof Name | Visible | Selected | Value | Default |
|--------------------|-------------------------------------|--------------------------|---------|---------|
| L1 Bandwidth | <input checked="" type="checkbox"/> | <input type="checkbox"/> | 1283.25 | GB/sec |
| L2 Bandwidth | <input checked="" type="checkbox"/> | <input type="checkbox"/> | 375.94 | GB/sec |
| L3 Bandwidth | <input checked="" type="checkbox"/> | <input type="checkbox"/> | 177.2 | GB/sec |
| DRAM Bandwidth | <input checked="" type="checkbox"/> | <input type="checkbox"/> | 21.89 | GB/sec |
| SP Vector FMA Peak | <input checked="" type="checkbox"/> | <input type="checkbox"/> | 444.88 | GFLOPS |
| DP Vector FMA Peak | <input checked="" type="checkbox"/> | <input type="checkbox"/> | 219.15 | GFLOPS |
| SP Vector Add Peak | <input checked="" type="checkbox"/> | <input type="checkbox"/> | 112.3 | GFLOPS |
| DP Vector Add Peak | <input checked="" type="checkbox"/> | <input type="checkbox"/> | 56.14 | GFLOPS |
| Scalar Add Peak | <input checked="" type="checkbox"/> | <input type="checkbox"/> | 13.82 | GFLOPS |

Load... Save...

Loop Weight Representation Cancel Default

☒ Size Color

Time

+ 4 green

- Threshold Value 2 %

+ 6 yellow

- Threshold Value 12 %

+ 8 red

No self time loop color ☐ #eee

Load... Save...

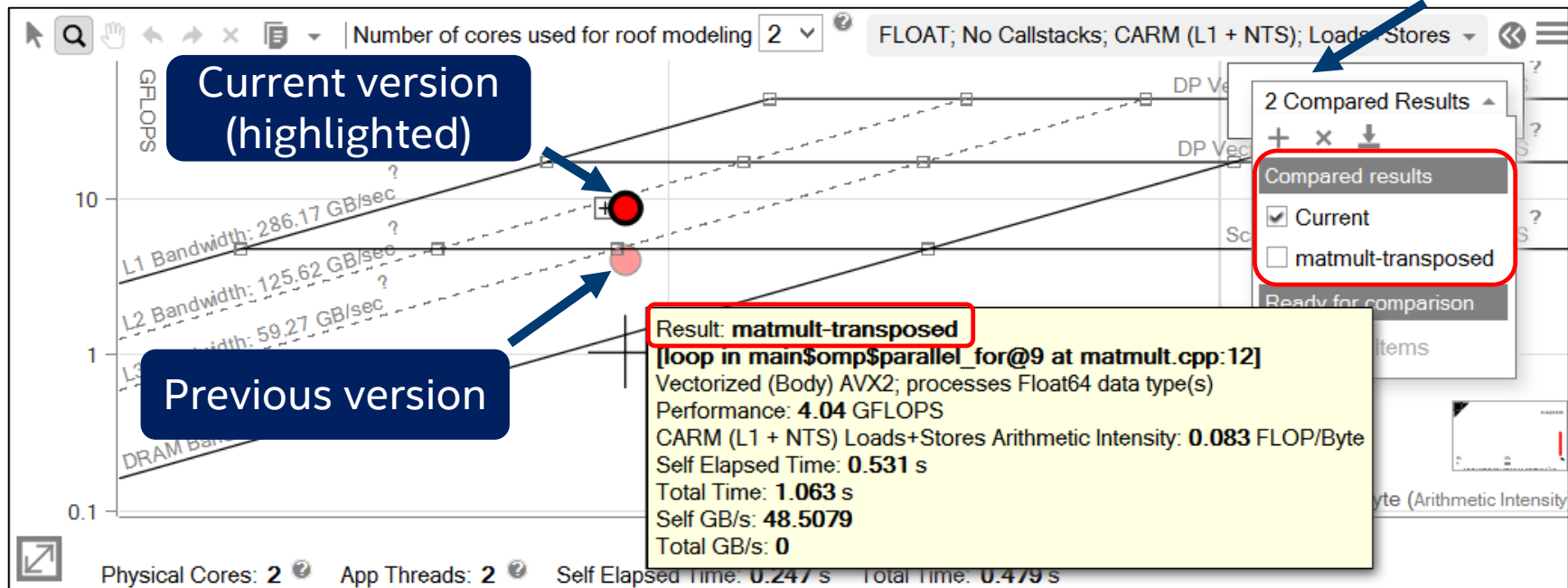
Hide or highlight the roofs

Edit values manually (if needed)

Choose dots coloring strategy (vector-vs-scalar or slow-vs-fast)

Compare results

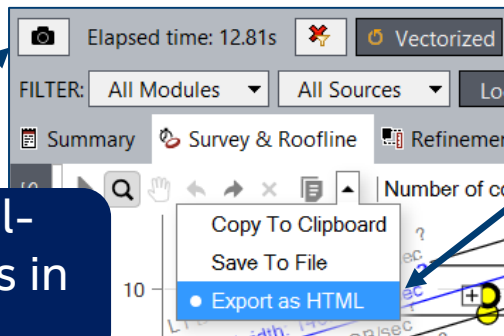
Loaded results for two versions



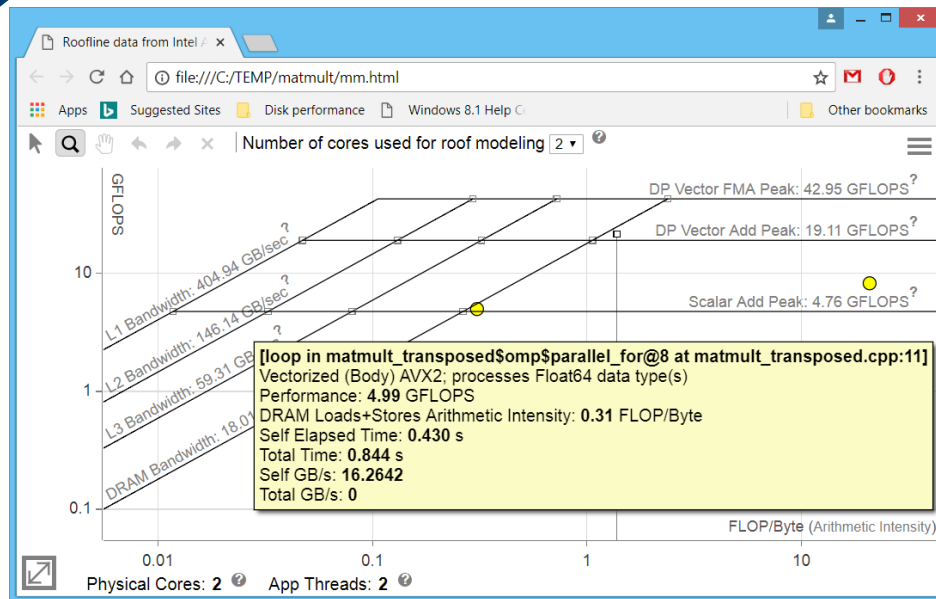
Easy to check optimization progress

Share with others

Snapshot (full-featured, opens in Advisor)



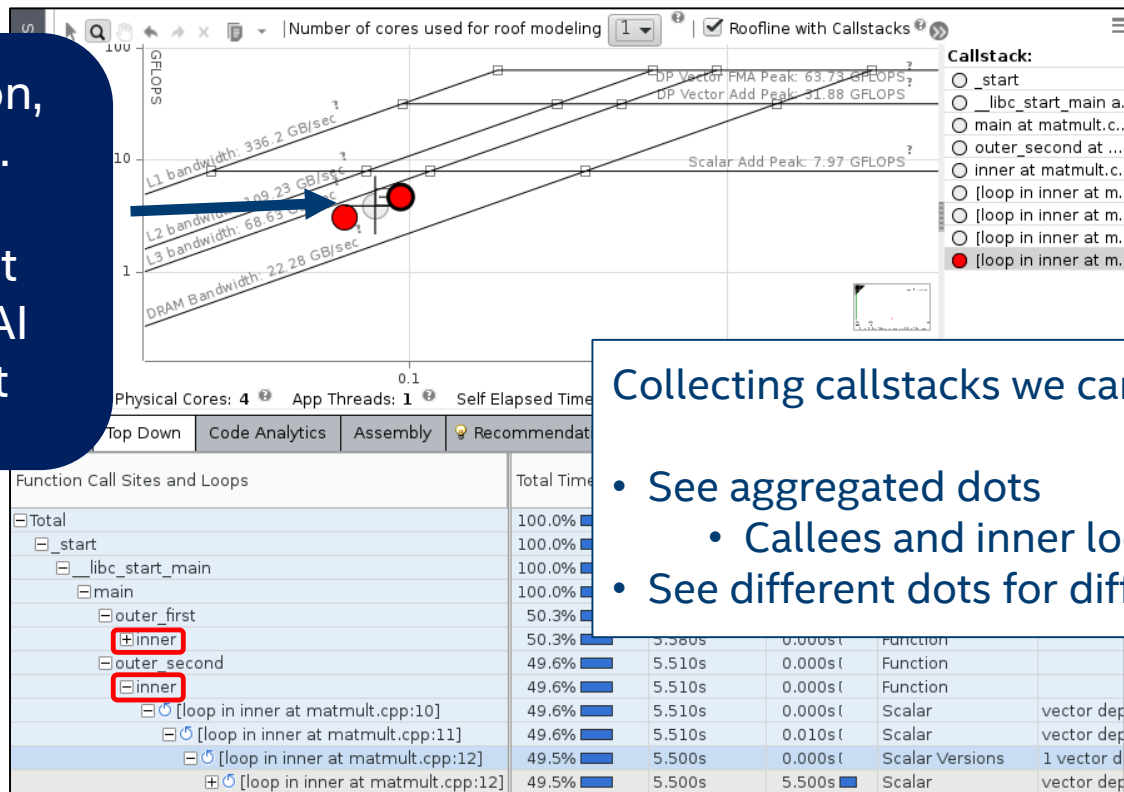
Standalone *interactive* HTML
(limited functionality)
Share roofline by email! - with
colleagues or your manager



A few words about callstacks

Same function,
same loop...

But different
FLOPS and AI
on different
callpaths!



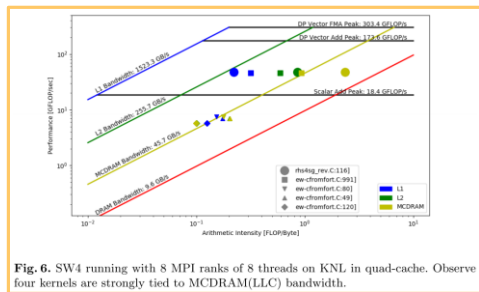
Collecting callstacks we can

- See aggregated dots
 - Callees and inner loops included
- See different dots for different callchains

Available in current product version, does not work with hierarchical roofline yet

Python API

- Fairly new Advisor Extensibility/customization mechanism. Actively used internally in Intel
- Up to 500 metrics for each loop/function. Really easy to use:
 - <advisor_install_dir>/pythonapi/examples
 - \$python **survey_bottomup.py** <project_dir>
 - Generate your own customized roofline charts



From ISC'18 paper
(cudos Tuomas)

```
import sys

try:
    import advisor

    project = advisor.open_project(sys.argv[1])
    survey = project.load(advisor.SURVEY)

    for row in survey.bottomup:
        # row as string
        print(row)
        # row as iterator
        for key in row:
            # row as dictionary
            print('{}: {}'.format(key, row[key]))
```