

INTEL[®] ADVISOR

Jackson Marusarz – Intel[®]

Agenda

- Motivation
- Threading Advisor
 - Threading Advisor Workflow
 - Advisor Interface
 - Survey Report
 - Annotations
 - Suitability Analysis
 - Dependencies Analysis
- Vectorization Advisor & Roofline
 - Vectorization Advisor recap
 - Roofline
 - Memory Access Patterns Analysis
 - Dependencies Analysis
- Summary

Function Call Sites and Loops	Self Time	Type	FLOPS	AI	Why No Vectorization?	Vectorized Loops	Trip Counts	Instruction Set
			GFLOPS			Vect... Efficiency	Gain E... VL	Average Call Count Traits
[loop in main at roofline.cpp:221]	13.562s	Scalar	3.9171	0.179	novector directive...		664	10000000
[loop in main at roofline.cpp:138]	7.563s	Scalar	1.7561	0.045	novector directive...		664	10000000
[loop in main at roofline.cpp:200]	7.328s	Vectorized (Body)	7.2491	0.134		AVX2 31%	1.22x 4	166 10000000 FMA; Inserts; U
[loop in main at roofline.cpp:260]	2.656s	Vectorized (Body)	19.999	0.179		AVX2 100%	5.09x 4	166 10000000 FMA
[loop in main at roofline.cpp:151]	1.828s	Vectorized (Body)	7.2640	0.045		AVX 100%	4.80x 4	166 10000000
[loop in main at roofline.cpp:273]	1.813s	Vectorized (Body)	29.306	0.179		AVX2 100%	4.73x 4	166 10000000 FMA
[loop in main at roofline.cpp:199]	1.781s	Vectorized (Body)	14.913	0.089		AVX2 100%	5.14x 4	166 10000000 FMA

Line	Source	Total Time	%	Loop/F...	%	Traits
247	for (int i = 0; i < ARRAY_SIZE_1; i++) [loop in main at roofline.cpp:247] Vectorized AVX; FMA loop processes Float32; Float64 data type(s) and includes FMA; Inserts; Loop was unrolled by 2	0.406s		7.328s		
248	{					
249	AoS1_X[i] = AoS1_Y[i].a + AoS1_Y[i].a + AoS1_Y[i].b + AoS1_Y[i].b + AoS1_Y[i].b;	6.922s				FMA; Inserts; Unpacks
250	}					

Maximum Program Gain For All Sites: 6.22x

Serial time: 29.338s
Predicted Parallel time: 4.719s

Site Label	Source Location	Impact to Program Gain	Combined Site Metrics, All Instances		Site Instance Metrics, Parallel Time
			Total Serial Time	Total Parallel Time	Site Gain
solve	nqueens_serial.cpp:139	6.22x	29.14s	4.524s	6.44x 4.524s

Site Performance Scalability

Scalability of Maximum Site Gain

Loop Iterations (Tasks) Modeling

Avg. Number of Iterations (Tasks): 70

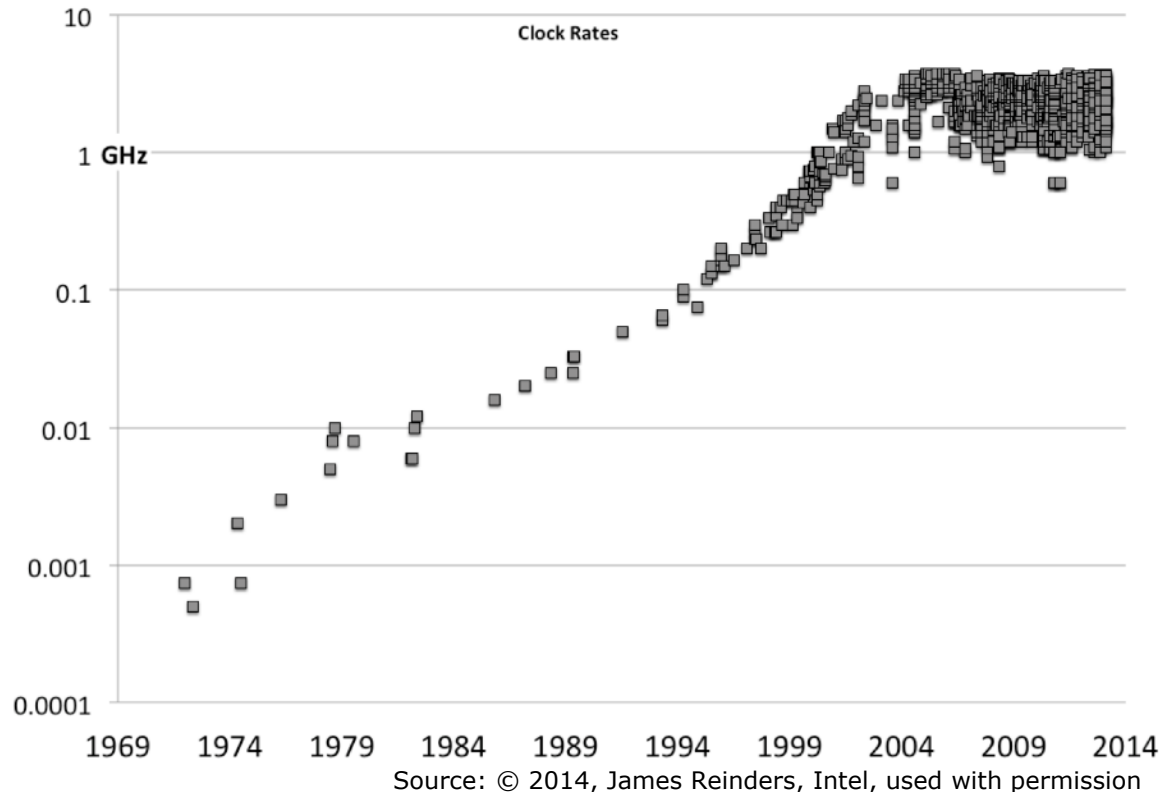
Avg. Iteration (Task) Duration: 0.416s

Options for Runtime Modeling:

- Reduce Site Overhead
- Reduce Task Overhead
- Reduce Lock Overhead
- Reduce Lock Contention
- Enable Task Chunking

The “Free Lunch” is over, really

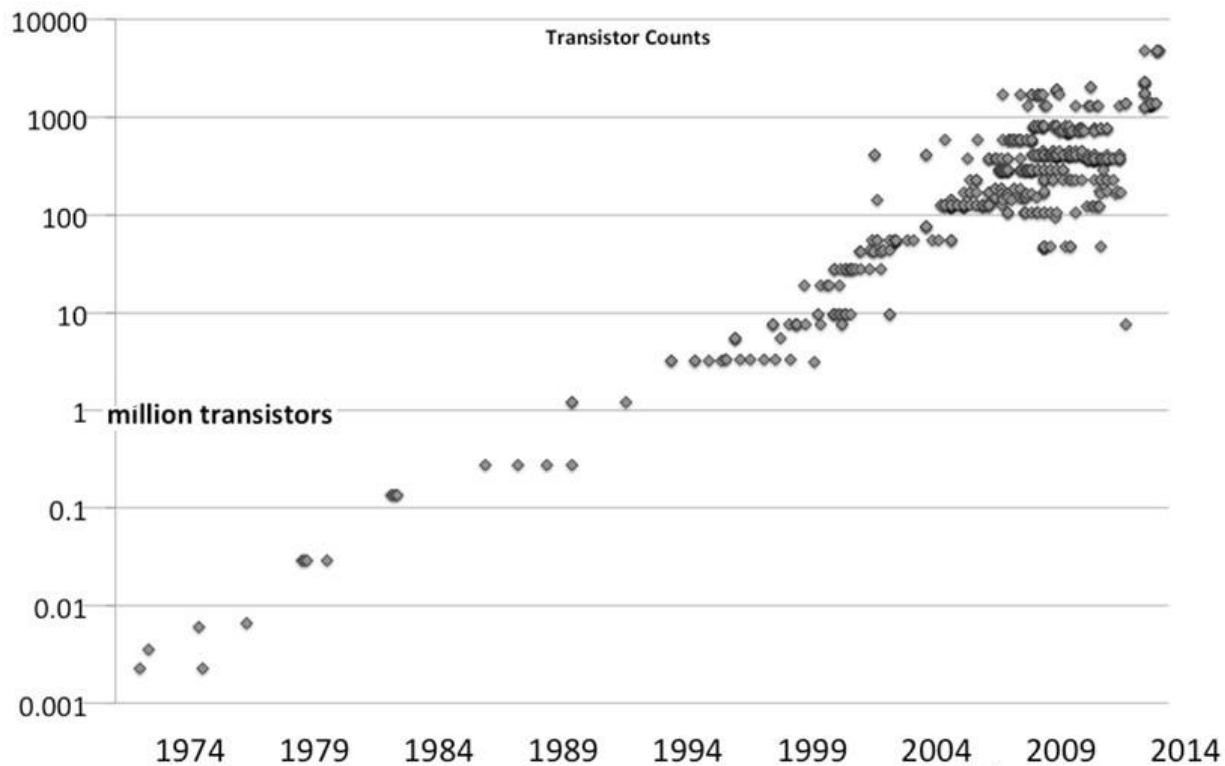
Processor clock rate growth halted around 2005



Software must be parallelized to realize all the potential performance

Moore's Law Is Going Strong

Hardware performance continues to grow exponentially



Source: © 2014, James Reinders, Intel, used with permission

"We think we can continue Moore's Law for at least another 10 years."

Intel Senior Fellow
Mark Bohr
2015

Changing Hardware Impacts Software

More Cores → More Threads → Wider Vectors



Intel® Xeon® Processor

	64-bit	5100 series	5500 series	5600 series	E5-2600	E5-2600 V2	E5-2600 V3	E5-2600 V4	Platinum 8180
Core(s)	1	2	4	6	8	12	18	22	28
Threads	2	2	8	12	16	24	36	44	56
SIMD Width	128	128	128	128	256	256	256	256	512

High performance software must be both

- Parallel (multi-thread, multi-process)
- Vectorized

*Product specification for launched and shipped products available on ark.intel.com.

Optimization Notice

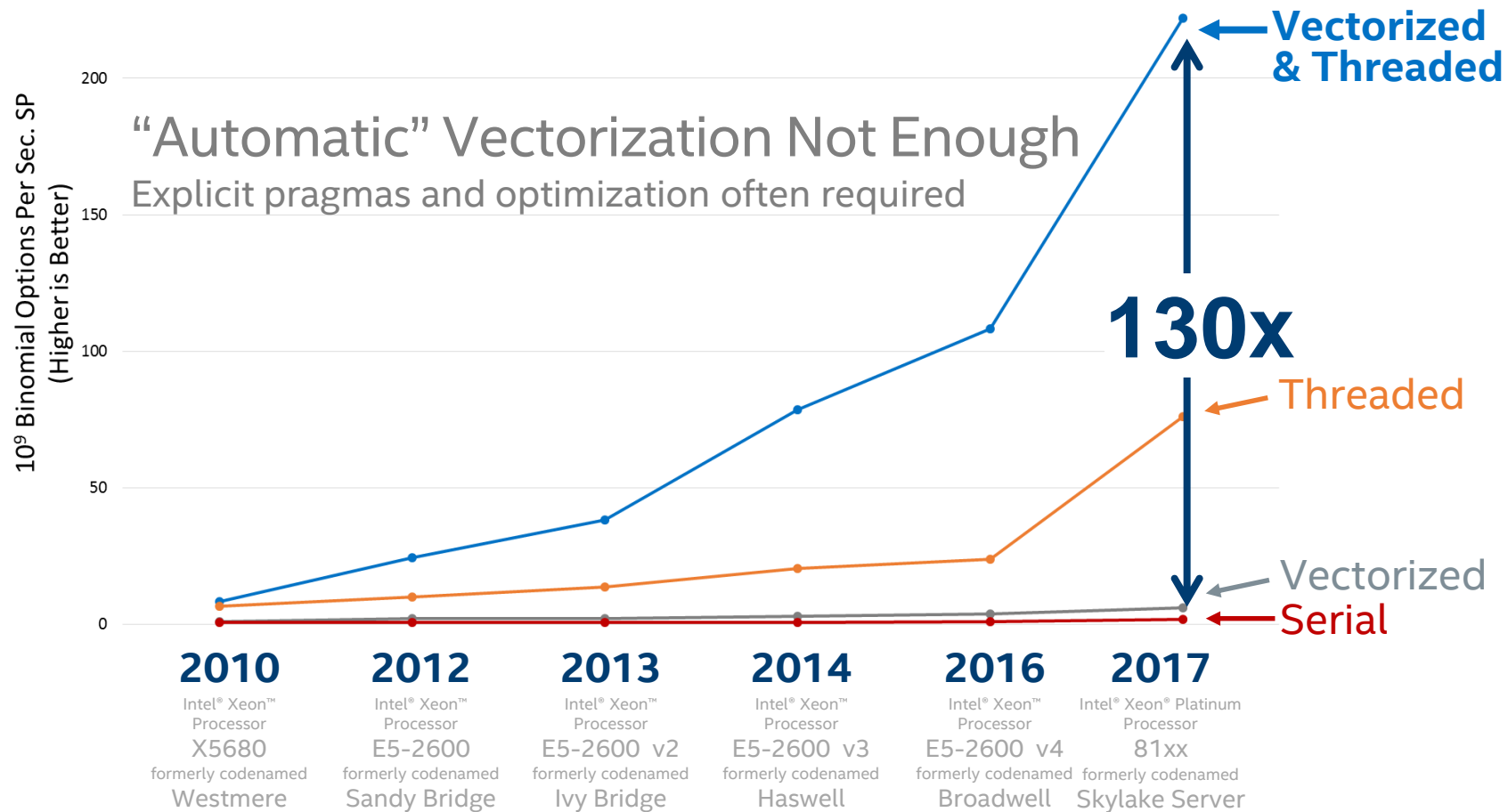
Copyright © 2017, Intel Corporation. All rights reserved.

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



Vectorize & Thread or Performance Dies

Threaded + Vectorized can be much faster than either one alone



The Difference Is Growing With Each New Generation of Hardware

Benchmark results were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown". Implementation of these updates may make these results inapplicable to your device or system. For more complete information about performance and benchmark results, visit www.intel.com/benchmarks. See [Vectorize & Thread or Performance Dies Configurations for 2010-2016 Benchmarks](#) in Backup. Benchmarks source: Intel Corporation.

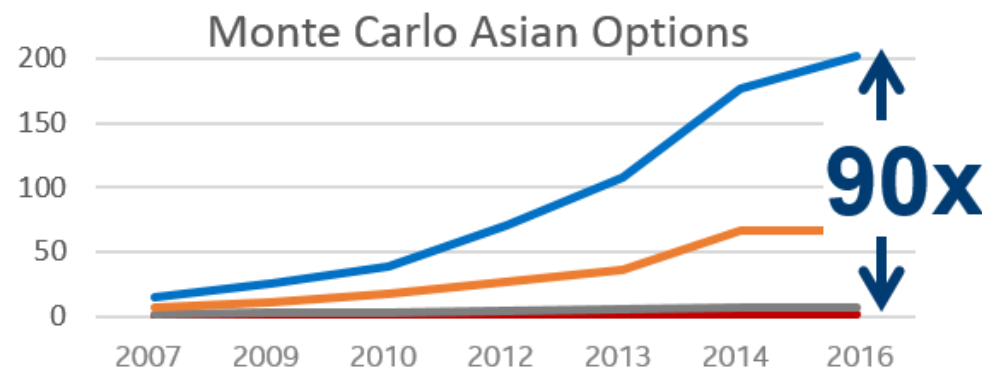
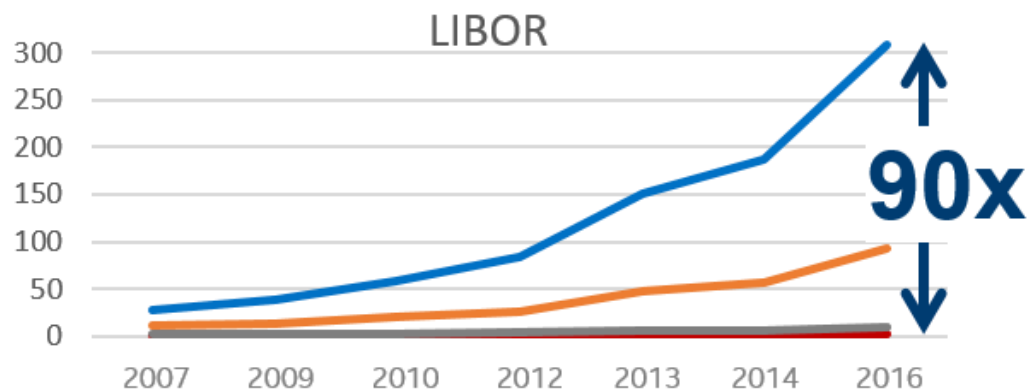
Optimization Notice

Copyright © 2017, Intel Corporation. All rights reserved.
*Other names and brands may be claimed as the property of others.

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, & SSSE3 instruction sets & 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

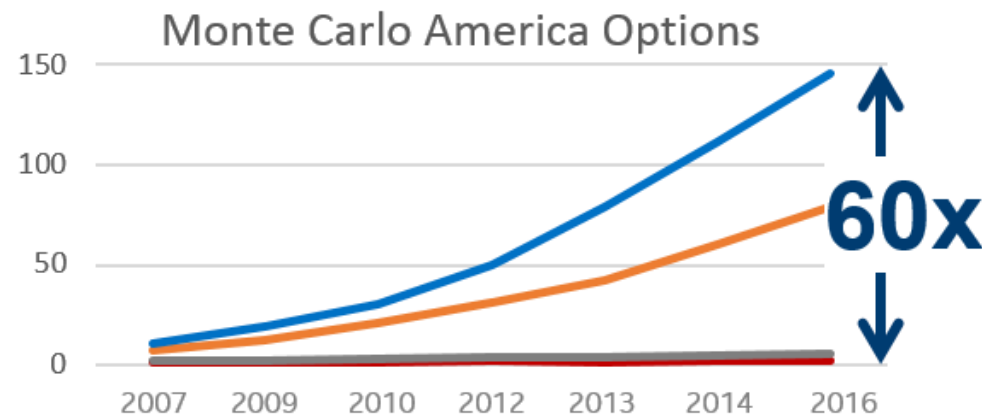
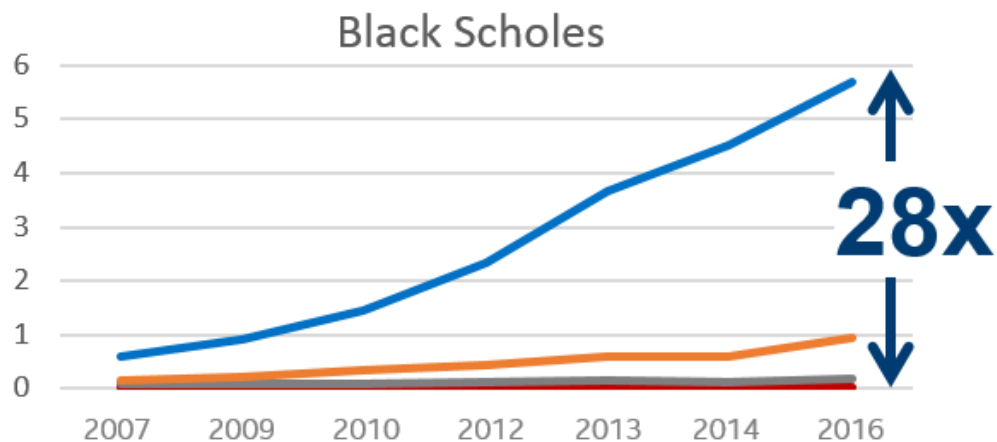


Vectorization & Threading Critical on Modern Hardware



Key:

- Vectorized & Threaded
- Threaded
- Vectorized
- Serial



Benchmark results were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown". Implementation of these updates may make these results inapplicable to your device or system. For more complete information about performance and benchmark results, visit www.intel.com/benchmarks. See [Vectorize & Thread or Performance Dies Configurations for 2010-2016 Benchmarks](#) in Backup. Benchmarks source: Intel Corporation.

Optimization Notice

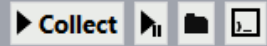
Copyright © 2017, Intel Corporation. All rights reserved.

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

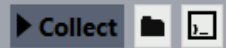
Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, & SSSE3 instruction sets & other optimizations. Intel does not guarantee the availability, functionality, 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



1. Survey Target



1.1 Find Trip Counts and FLOPS

 Trip Counts FLOPS

2. Annotate Sources

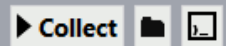
Add Intel Advisor annotations to identify possible parallel tasks and their enclosing parallel sites.

+ Steps to annotate

3. Check Suitability



4. Check Dependencies



THREADING ADVISOR

Serial Modeling Has Multiple Benefits

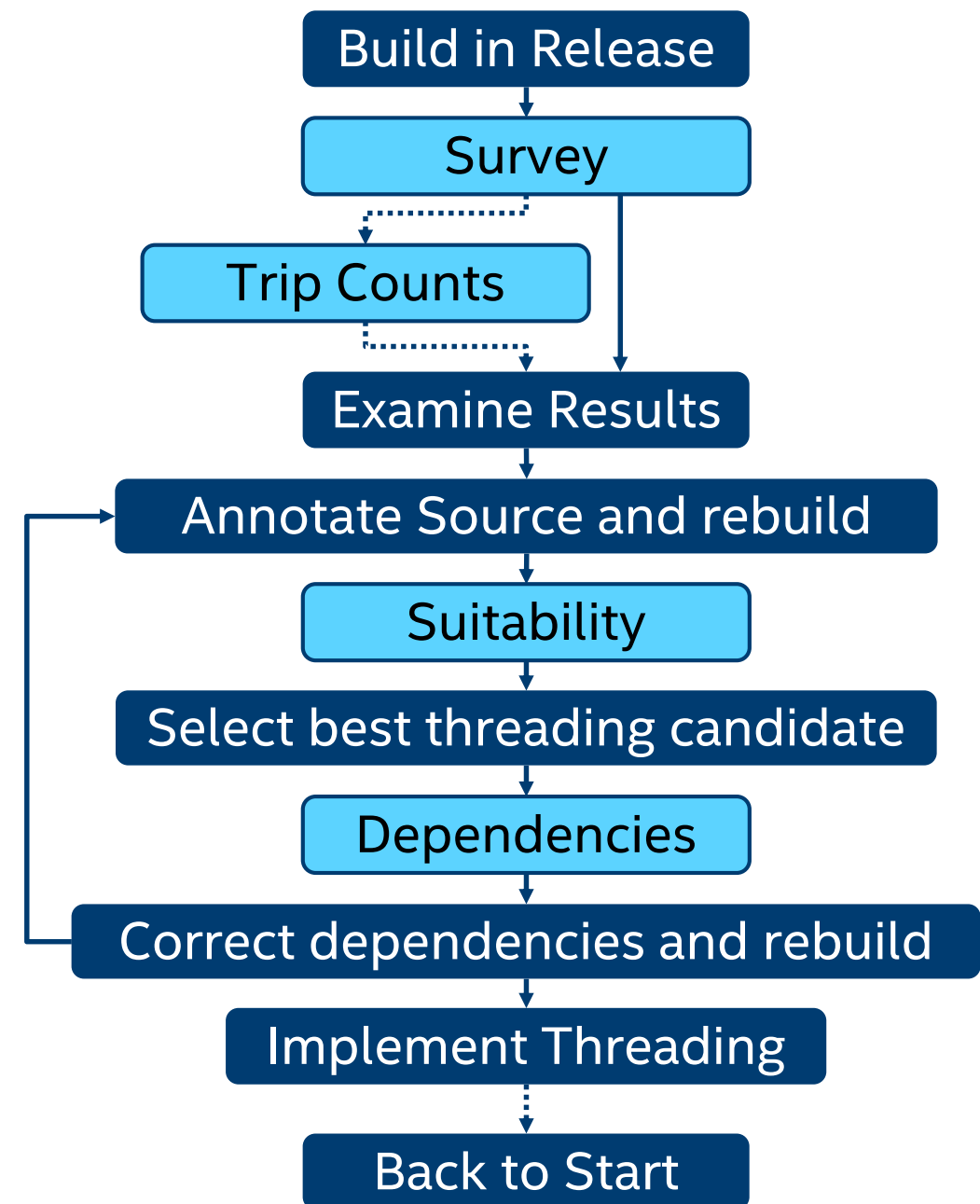
Intel® Advisor

- 1) Your application can't fail due to bugs caused by incorrect parallel execution.
(It's running serially.)
- 2) You can easily experiment with several different proposals before committing to the expense of implementation.
 - a) Measure performance - focus on where it will pay off.
 - b) Predict scalability, load balancing and overheads.
 - c) Predict (and avoid) data races
- 3) All of your test suites should still pass.
Validate the correctness of your transformations.
- 4) You can use Advisor on partially or completely parallelized code.

Design, measure and test before implementation

Threading Advisor Workflow

- Use the **Survey** to find good potential threading sites.
 - Optionally, follow up with **Trip Counts** to find information about iteration and call counts.
- Annotate your code.
- Use **Suitability** to predict how much performance improvement the proposed threading model will create under specific, editable conditions.
- Use **Dependencies** to determine whether the proposed model is safe, and what needs to be done to correct it.



Survey Report

Threading Advisor

Tip:

Survey sorts by Self Time by default. This is good for Vector Advisor, but for Threading Advisor, you may want to sort by Total Time.

- The Survey Report has lots of information, but most of it is more relevant to Vector Advisor.
- Look for outer loops or functions with high Total Time.
- In this example, setQueen has a high Total Time. It's recursive, but is originally called from a loop in Solve. That makes the loop in Solve a good potential candidate.

Function Call Sites and Loops	Vector Issues	Self Time	Total Time	Type
f setQueen		2.202s	38.480s	Function
[loop in setQueen at nqueens_serial.cpp:116]		0.297s	33.520s	Scalar
f _sclr_common_main_seh		0.000s	5.250s	Function

Function Call Sites and Loops	Total Time %	Self Time	Total Time	Type
Total	100.0%	0.000s	5.250s	
RtlUserThreadStart	100.0%	0.000s	5.250s	Function
RtlUserThreadStart	100.0%	0.000s	5.250s	Function
BaseThreadInitThunk	100.0%	0.000s	5.250s	Function
_sclr_common_main_seh	100.0%	0.000s	5.250s	Function
main	100.0%	0.000s	5.250s	Function
solve	100.0%	0.000s	5.250s	Function
[loop in solve at nqueens_serial.cpp:1	100.0%	0.000s	5.250s	Scalar
setQueen	100.0%	0.000s	5.250s	Function
[loop in setQueen at nqueens_se	100.0%	0.000s	5.250s	Scalar
setQueen	100.0%	0.000s	5.250s	Function
[loop in setQueen at nque	100.0%	0.000s	5.250s	Scalar
setQueen	100.0%	0.047s	5.250s	Function
[loop in setQueen at	98.5%	0.031s	5.172s	Scalar
setQueen	97.9%	0.031s	5.141s	Function
[loop in setQue	96.1%	0.016s	5.047s	Scalar
setQueen	95.8%	0.219s	5.031s	Function
[loop in se	88.1%	0.031s	4.625s	Scalar
setQueen	87.5%	0.282s	4.594s	Function

Annotating Your Code

- Annotations are notes to Advisor. They are *not* parallelization commands. They do not affect the way the program itself runs.
- They mark places Advisor should treat as locks or parallel sites.
- To use annotations, you must include the appropriate header/module.

C/C++

- In source files where annotations are used, add:
#include <advisor-annotate.h>
- Add `<install_dir>/include` to your include directories.

FORTRAN

- In source files where annotations are used, add:
use advisor_annotate
- Add
`<install_dir>/include`
to your include directories.

C#

- In source files where annotations are used, add:
using AdvisorAnnotate;
- Add the C# annotations definition file to your project.

- The Advisor User's Guide contains a section on Annotations with full documentation, examples, and instructions on the above if you forget.

Common Annotations

- **Site Annotations**

- Indicate locations which contain a number of tasks which can run in parallel with each other. Code within the site that is not part of a task is executed only by the entering thread, but can still run in parallel with the tasks.
- Consist of a Site Begin and a Site End annotation.

- **Task Annotations**

- Indicate chunks of code which can run in parallel with other tasks in the same site, including any additional copies of themselves which may be launched.
- Consist of a Task Begin and a Task End annotation.

- **Iteration Task Annotations**

- Indicate that the entirety of a loop's contents should be considered as a task.
- Consist only of an Iteration Task annotation. The end is implied.

- **Lock Annotations**

- Indicate where Advisor should simulate locking.
- Consist of a Lock Acquire annotation and a Lock Release annotation.

Suitability Analysis

- Using your annotations, Advisor models how the program would behave in parallel, and predicts performance in specified hypothetical circumstances.

Summary Survey & Roofline Refinement Reports Annotation Report Suitability Report INTEL ADVISOR 2018

Maximum Program Gain For All Sites: 5.74x

Serial time: 5.391s
Predicted Parallel time: 0.938s

Target System: CPU Threading Model: Intel TBB CPU Count: 8

Site Label	Source Location	Impact to Program Gain	Combined Site Metrics, All Instances			Site Instance Metrics,
			Total Serial Time	Total Parallel Time	Site Gain	Parallel Time
solve	nqueens_serial.cpp:...	5.74x	5.229s	0.777s	6.73x	0.777s

Site Performance Scalability Site Details

Scalability of Maximum Site Gain

Loop Iterations (Tasks) Modeling

Avg. Number of Iterations (Tasks): 14

Avg. Iteration (Task) Duration: 0.374s

Runtime Modeling

Type of Change

- Reduce Site Overhead
- Reduce Task Overhead
- Reduce Lock Overhead
- Reduce Lock Contention
- Enable Task Chunking

Gain Benefit if Enabled

Apply

Indicate how many CPUs, what kind of system, and what threading model to make predictions on.

Select a site to view site-specific info in the bottom pane.

Calculate on the assumption you're using framework constructs that address these issues.

See how things would change if you altered the duration and/or number of iterations/tasks.

Suitability Analysis

- The white dots are the estimated number of times speedup at the given number of CPUs.
- The bullseye is placed on the one corresponding to the number of CPUs you have selected in the controls.
- If your bullseye is in:

Red:

No benefit, possible slowdown.

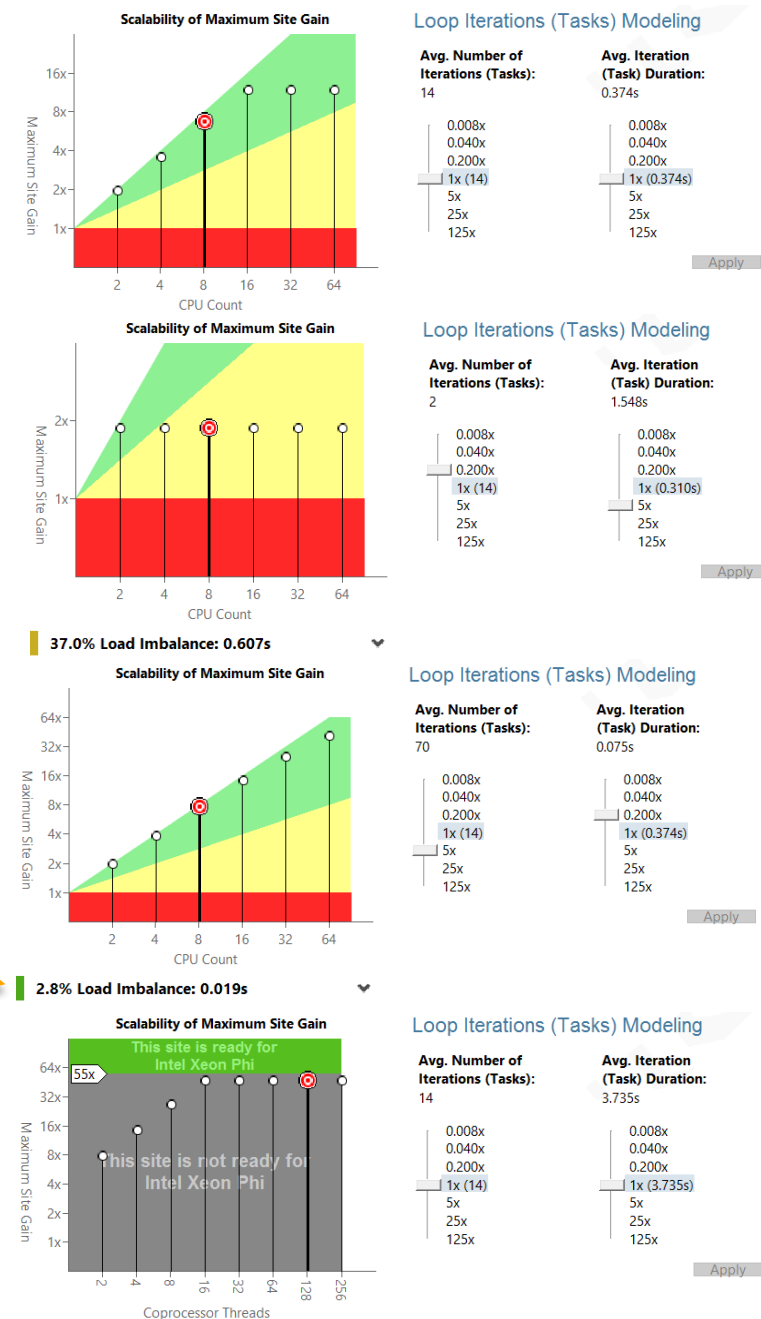
Yellow:

Poor scaling, needs improvement.

Green:

Ideal scaling!

- Advisor provides warnings when it predicts that your specifications would lead to issues like load imbalance.



Optimization Notice

Copyright © 2017, Intel Corporation. All rights reserved.
 *Other names and brands may be claimed as the property of others.



Dependencies Analysis

Threading Advisor

- This is the same analysis as in Vectorization Advisor. It works with annotations as well as selections in the survey report.
- Add lock annotations or reorganize code to resolve reported dependencies, then re-run the analysis to confirm the problem has been resolved.
- Run suitability again to check that you still get good improvement.
- Once you're happy with Advisor's predictions, replace the annotations with actual parallelism and locks.

The screenshot displays the Intel Threading Advisor interface. At the top, there are tabs for Summary, Survey & Roofline, Refinement Reports, Annotation Report, and Suitability Report. Below the tabs, a table shows analysis results for a loop in solve at nqueens_serial.cpp:111. The table has columns for Site Location, Loop-Carried Dependencies, Strides Distribution, Access Pattern, and Max. Site Footprint. The Loop-Carried Dependencies column shows RAW:1, WAR:1, and WAW:1. Below the table, a code snippet is shown with annotations: `137 //ADVISOR SUITABILITY EDIT: Uncomment the three annotations below to model`, `138 // parallelizing the body of this for() loop.`, `139 ANNOTATE_SITE_BEGIN(solve);`, and `140 for (int i=0; i < boardsize; i++) {`. Below the code, there are tabs for Dependencies Report and Recommendations. The Dependencies Report tab is active, showing a table of Problems and Messages. The table has columns for ID, Type, Site Name, Sources, Modules, and State. The table lists five problems: P1 (Parallel site information), P3 (Read after write dependency), P4 (Write after write dependency), and P5 (Write after read dependency). Below the Problems and Messages table, there is a section for Read after write dependency: Code Locations, which has columns for ID, Instruction..., Description, Source, Function, Variable refer..., Module, and State. This section lists three code locations: X3 (Read), X4 (Write), and X5 (Parallel site).

Site Location	Loop-Carried Dependencies	Strides Distribution	Access Pattern	Max. Site Footprint
[loop in solve at nqueens_serial.cpp:111]	RAW:1 WAR:1 WAW:1	No information available	No information available	No information available

```
137 //ADVISOR SUITABILITY EDIT: Uncomment the three annotations below to model
138 // parallelizing the body of this for() loop.
139 ANNOTATE_SITE_BEGIN(solve);
140 for (int i=0; i < boardsize; i++) {
141     ...
142 }
```

ID	Type	Site Name	Sources	Modules	State
P1	Parallel site information	solve	nqueens_serial.cpp	1_nqueens_serial.exe	✓ Not a problem
P3	Read after write dependency	solve	nqueens_serial.cpp	1_nqueens_serial.exe	New
P4	Write after write dependency	solve	nqueens_serial.cpp	1_nqueens_serial.exe	New
P5	Write after read dependency	solve	nqueens_serial.cpp	1_nqueens_serial.exe	New

ID	Instruction...	Description	Source	Function	Variable refer...	Module	State
X3	0x401c04	Read	nqueens_serial.cpp:111	setQueen	nrOfSolutions	1_nqueens_serial.exe	New
X4	0x401c04	Write	nqueens_serial.cpp:111	setQueen	nrOfSolutions	1_nqueens_serial.exe	New
X5	0x401b5f	Parallel site	nqueens_serial.cpp:139	solve		1_nqueens_serial.exe	New

Add Parallel Framework

Summary | Survey Report | Annotation Report | Suitability Report | Correctness Report

Maximum program gain[®]: 5.20x (8 CPUs, Intel TBB Threading Model)

These annotated parallel sites were detected:

Parallel Site	Maximum Site Gain [®]	Correctness Problems
solve (nqueens_annotated.cpp:113)	6.51x	2 ⚠ 0

Consider adding parallel site and task annotations around these time-consuming loops found during Survey an

Loop	Source Location	CPU Total Time [®]
setQueen	nqueens_annotated.cpp:96	1.8252s
solve	nqueens_annotated.cpp:117	1.8252s
setQueen	nqueens_annotated.cpp:69	0.1976s
setQueen	nqueens_annotated.cpp:69	0.1877s
setQueen	nqueens_annotated.cpp:69	0.1346s

Here is the list of source locations

Here are templates for popular parallel frameworks

Serial Code with Intel Advisor Annotations	Parallel Code using Intel TBB
<pre>// Locking ANNOTATE_LOCK_ACQUIRE(); Body(); ANNOTATE_LOCK_RELEASE();</pre>	<pre>// Locking can use various mutex types provided // by Intel TBB. For example: #include <tbb/tbb.h> ... tbb::mutex g_Mutex; ... { tbb::mutex::scoped_lock lock(g_Mutex); Body(); }</pre>
<pre>// Do-All Counted loops, one task ANNOTATE_SITE_BEGIN(site); For (I = 0; I < N; ++I) { ANNOTATE_ITERATION_TASK(task); {statement;} } ANNOTATE_SITE_END();</pre>	<pre>// Do-All Counted loops, using lambda // expressions #include <tbb/tbb.h> ... tbb::parallel_for(0,N,[&](int I) { statement; });</pre>
<pre>// Create Multiple Tasks ANNOTATE_SITE_BEGIN(site); ANNOTATE_TASK_BEGIN(task1); statement-or-task1; ANNOTATE_TASK_END(); ANNOTATE_TASK_BEGIN(task2); statement-or-task2; ANNOTATE_TASK_END(); ANNOTATE_SITE_END();</pre>	<pre>// Create Multiple tasks, using lambda // expressions #include <tbb/tbb.h> ... tbb::parallel_invoke([&]{statement-or-task1;}, [&]{statement-or-task2;});</pre>

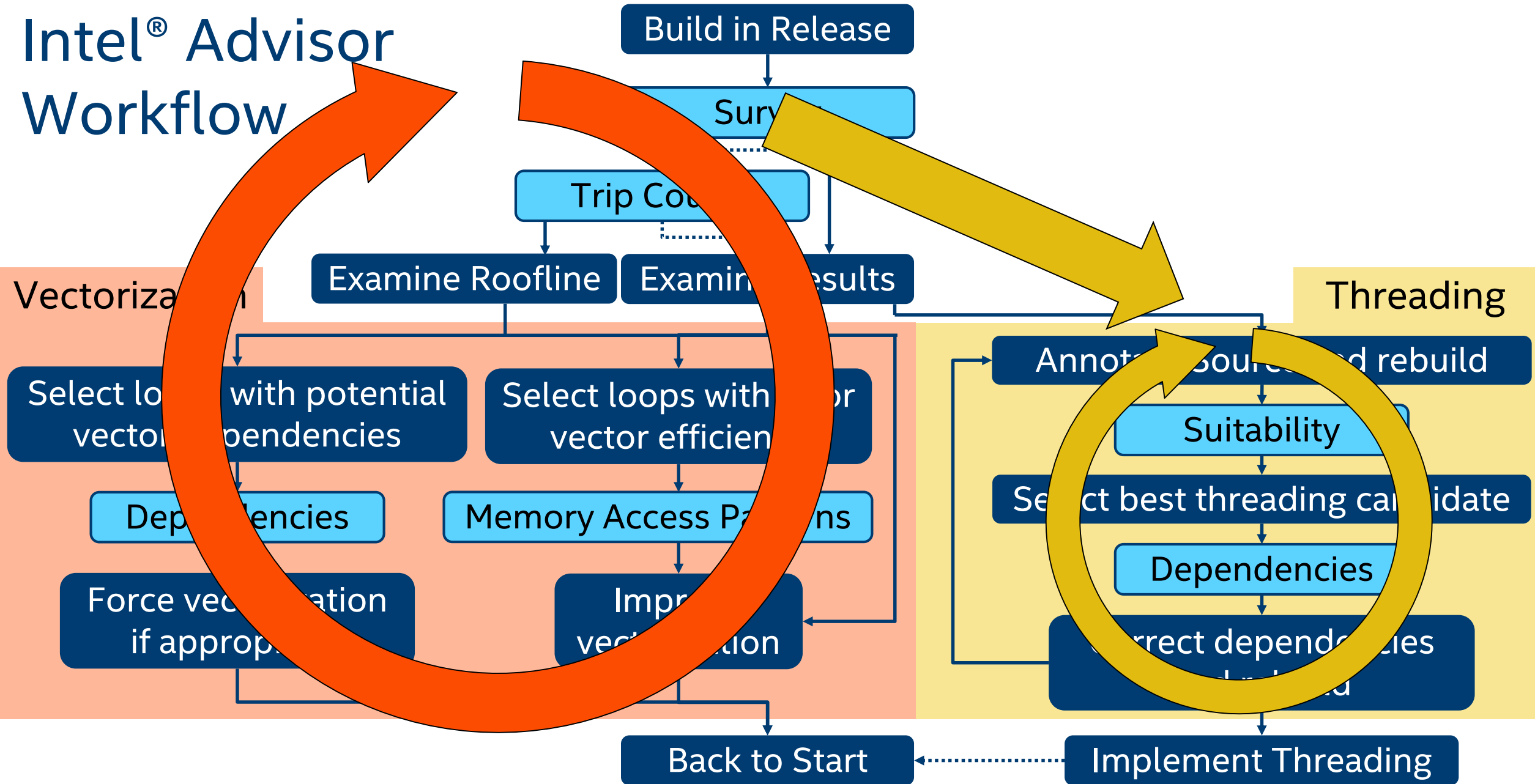
Threading Model:

- Intel TBB
- Other
- Intel Cilk Plus
- OpenMP
- Microsoft TPL

Intel® Advisor

- Contains overhead metrics for popular parallel frameworks
- Quickly prototype and evaluate alternatives
- Detailed help pages for popular parallel frameworks

Intel® Advisor Workflow



Optimization Notice

Copyright © 2017, Intel Corporation. All rights reserved.
 *Other names and brands may be claimed as the property of others.





Vectorization
Workflow



Threading
Workflow

OFF Batch mode



Run Roofline

▶ Collect  

1. Survey Target

▶ Collect   

1.1 Find Trip Counts and FL...

▶ Collect  

Trip Counts



FLOPS

Mark Loops for Deeper Anal...

Select loops in the Survey Report for Dependencies and/or Memory Access Patterns analysis.



-- There are no marked loops --


2.1 Check Dependencies

▶ Collect  

-- Nothing to analyze --

2.2 Check Memory Access P...

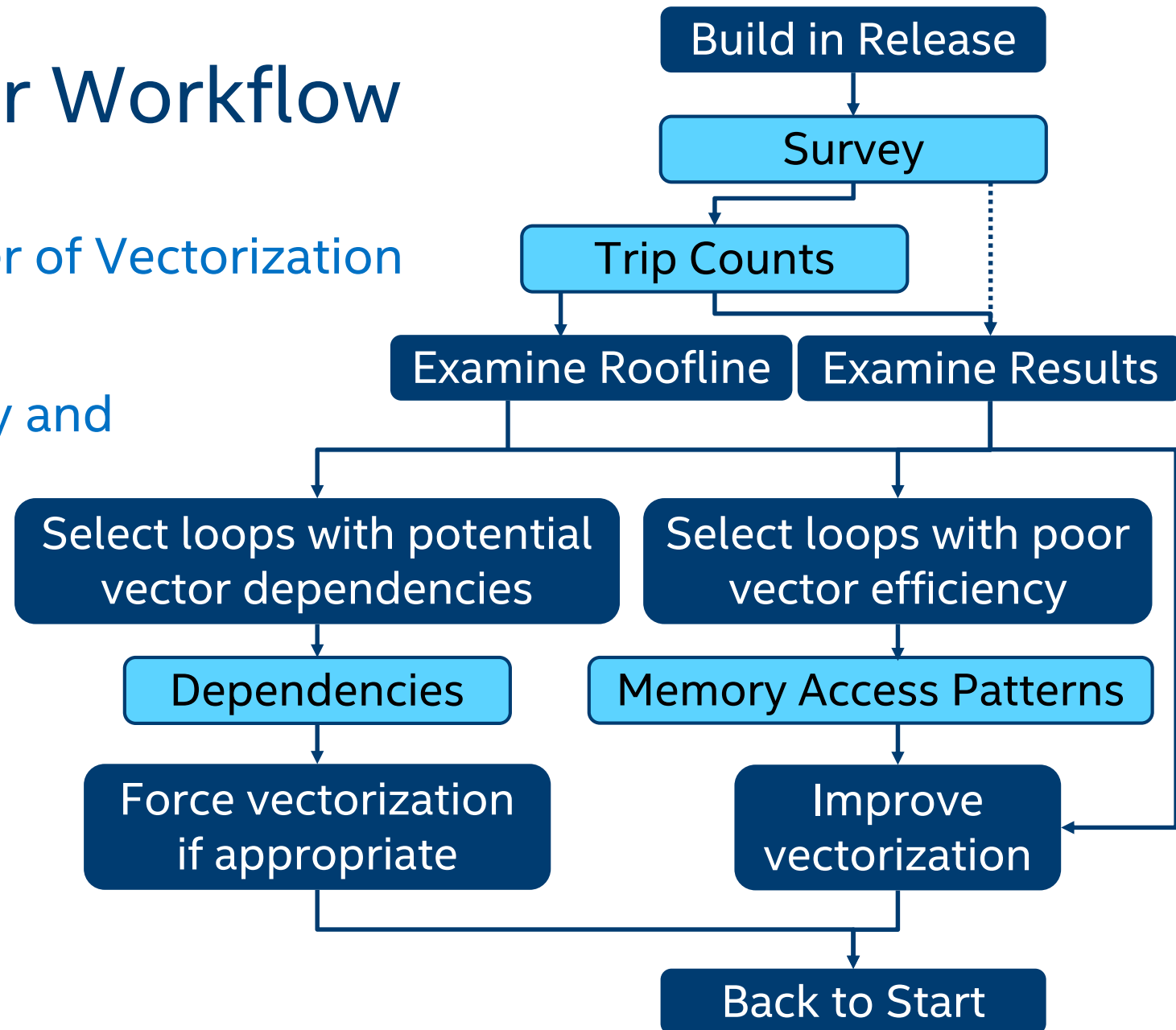
▶ Collect  

 Re-finalize Survey

VECTORIZATION ADVISOR & ROOFLINE

Vectorization Advisor Workflow

- **Survey** is the bread and butter of Vectorization Advisor! All else builds on it!
- **Trip Counts** adds onto Survey and enables the **Roofline**.
- **Dependencies** determines whether it's safe to force a scalar loop to vectorize.
- **Memory Access Patterns** diagnoses vectorization inefficiency caused by poor memory striding.



Survey

Vectorization Advisor

Tip:
For vectorization, you generally only care about loops. Set the type dropdown to "Loops".

Efficiency is important!

$$\text{Efficiency} = 100\% \frac{\text{Speedup}}{\text{Vec. Length}}$$
 The black arrow is 1x. Gray means you got less than that. Gold means you got more. You want to get this value as high as possible!

- Function/Loop Icons**
- Scalar Function
 - Vector Function
 - Scalar Loop
 - Vector Loop

Vectorizing a loop is usually best done on innermost loops. Since it effectively divides duration by vector length, you want to target loops with high self time.

Function Call Sites and Loops	Vector Issues	Self Time	Total Time	Type	Why No Vectorization?	Vectorized Loops			
						Vect...	Efficiency	Gain...	VL .
[loop in main at example.cpp:38]	1 Assumed depend...	0.391s	0.391s	Scalar	vector dependen...				
[loop in main at example.cpp:64]	1 Possible inefficien...	0.297s	0.297s	Vector...		AVX2	2%	0.37x	16
[loop in main at example.cpp:51]	1 Possible inefficien...	0.094s	0.094s	Vector...	1 vectorizatio...	AVX2	8%	1.23x	16
[loop in main at example.cpp:26]		0.030s	0.030s	Vector...		AVX2	100%	7.98x	8
[loop in main at example.cpp:14]	3 Assumed depend...	0.000s	0.000s	Scalar	vector dependen...				
[loop in main at example.cpp:23]		0.000s	0.030s	Scalar	inner loop w...				

Expand a vectorized loop to see it split into body, peel, and remainder (if applicable).

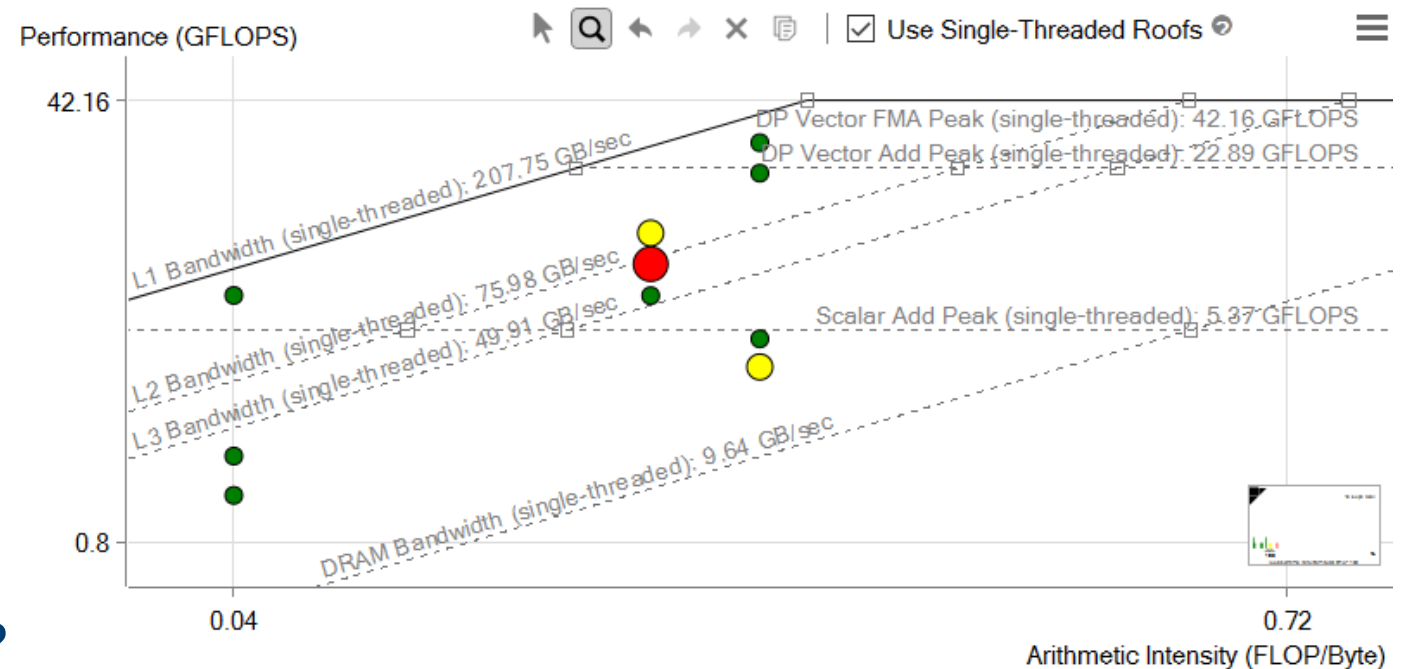
Advisor *advises* you on potential vector issues. This is often your cue to run MAP or Dependencies. Click the icon to see an explanation in the bottom pane.

The Intel Compiler embeds extra information that Advisor can report in addition to its sampled data, such as why loops failed to vectorize.

What is a Roofline Chart?

A Roofline Chart plots application performance against hardware limitations.

- Where are the bottlenecks?
- How much performance is being left on the table?
- Which bottlenecks can be addressed, and which *should* be addressed?
- What's the most likely cause?
- What are the next steps?

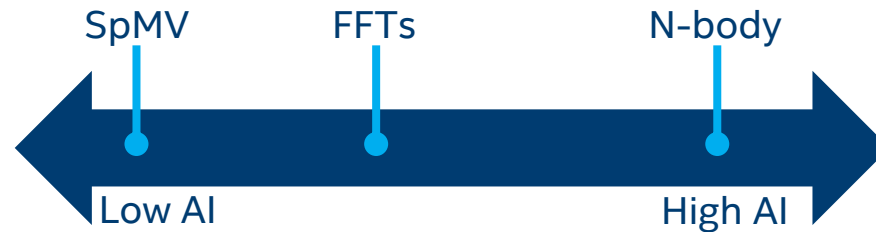


Roofline first proposed by University of California at Berkeley:
[Roofline: An Insightful Visual Performance Model for Multicore Architectures](#), 2009
Cache-aware variant proposed by University of Lisbon:
[Cache-Aware Roofline Model: Upgrading the Loft](#), 2013

Roofline Metrics

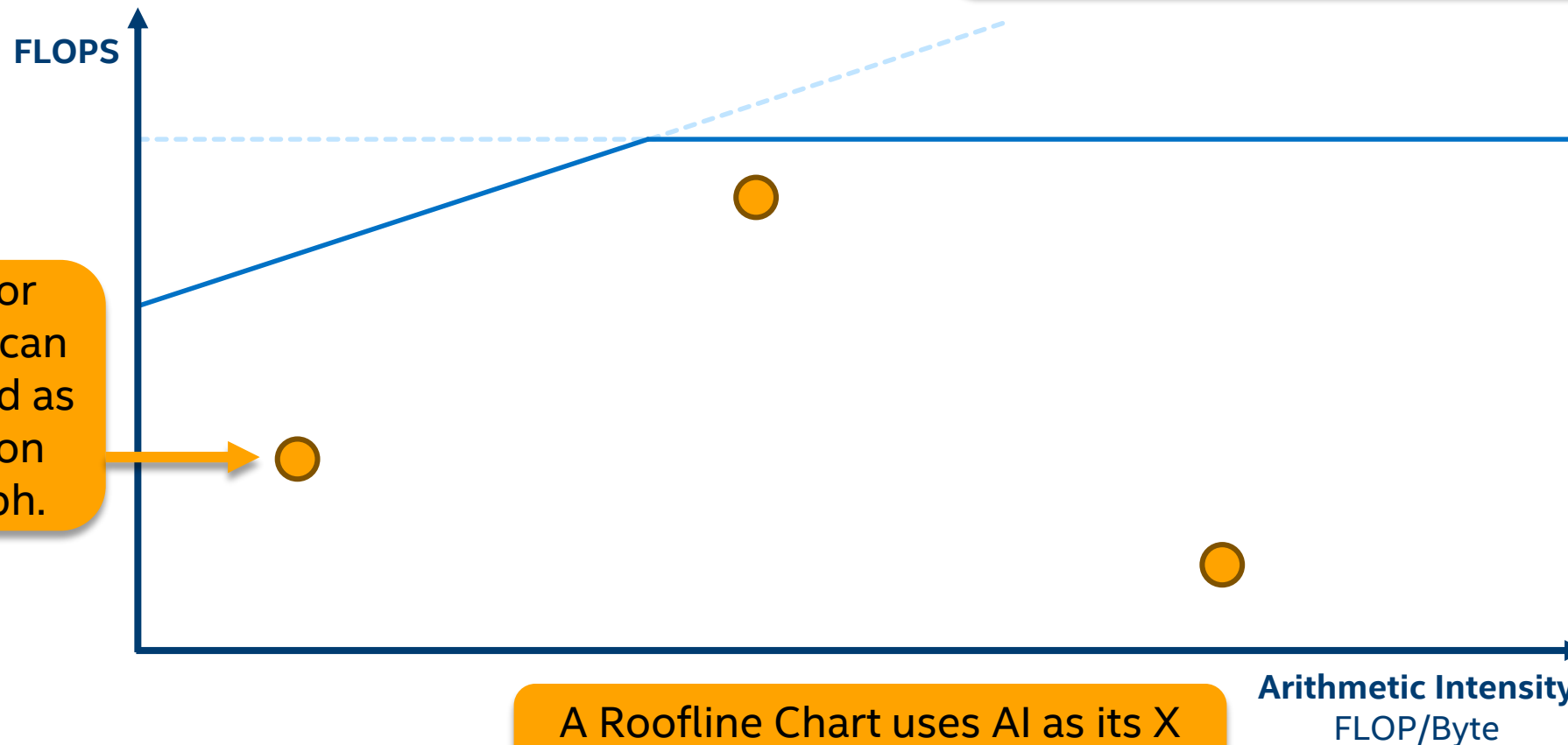
Roofline is based on Arithmetic Intensity (AI) and FLOPS.

- **Arithmetic Intensity:** FLOP / Byte Accessed
 - This is a characteristic of your algorithm



- **FLOPS:** Floating-Point Operations / Second
 - Is a measure of an implementation (it achieves a certain FLOPS)
 - **And** there is a maximum that a platform can provide

Plotting a Roofline Chart



The maximum FLOPS as a product of ops/byte (AI) and maximum bytes supplied per second is a diagonal line.

A loop or function can be plotted as a point on the graph.

The CPU's maximum FLOPS can be plotted as a horizontal line.

A Roofline Chart uses AI as its X axis and FLOPS as its Y axis.

Classic vs. Cache-Aware Roofline

Intel® Advisor uses the Cache-Aware Roofline model, which has a different definition of Arithmetic Intensity than the original (“Classic”) model.

Classical Roofline

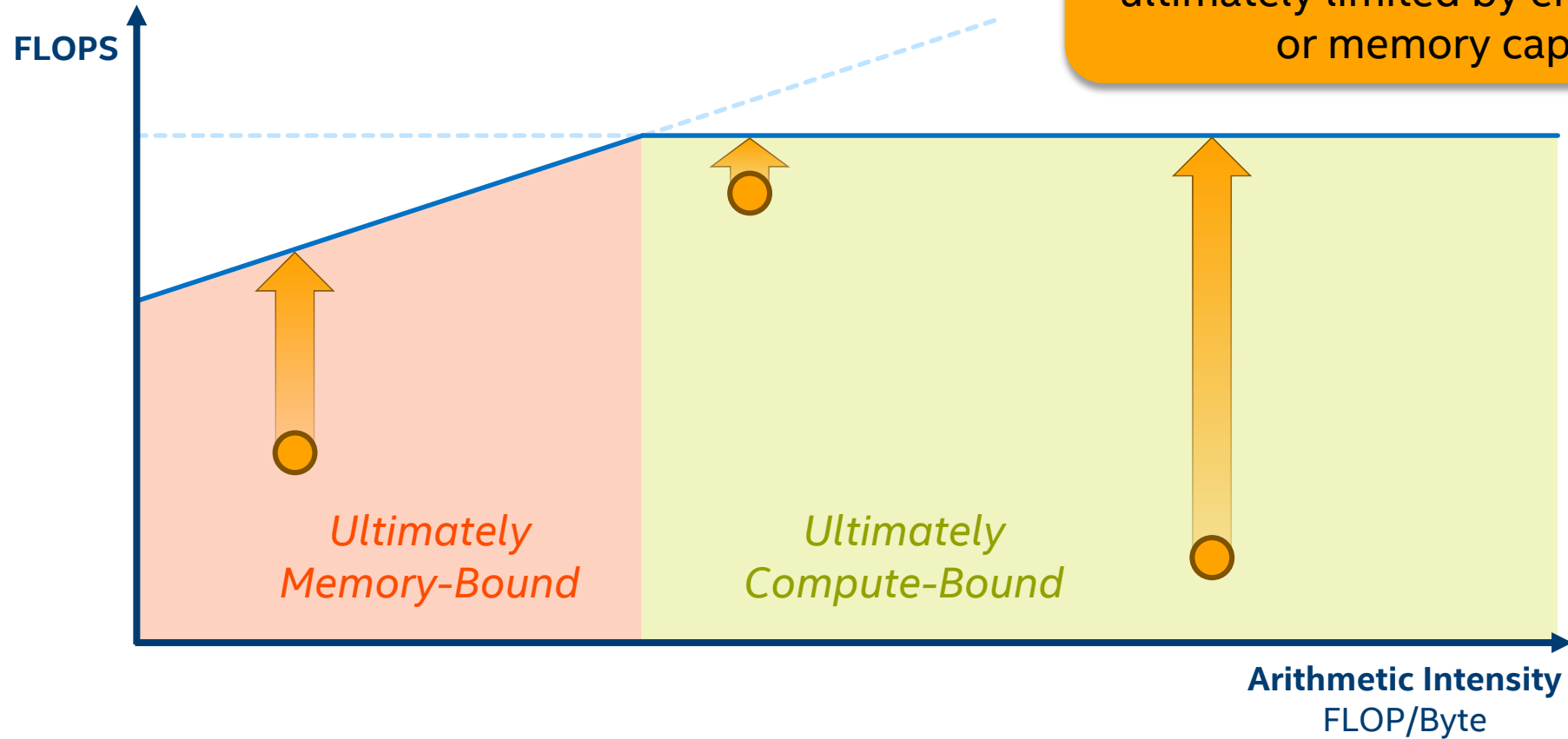
- Traffic measured from one level of memory (usually DRAM)
- AI may change with data set size
- AI changes as a result of memory optimizations

Cache-Aware Roofline

- Traffic measured from all levels of memory
- AI is tied to the algorithm and will not change with data set size
- Optimization does not change AI*, only the performance

**Compiler optimizations may modify the algorithm, which may change the AI.*

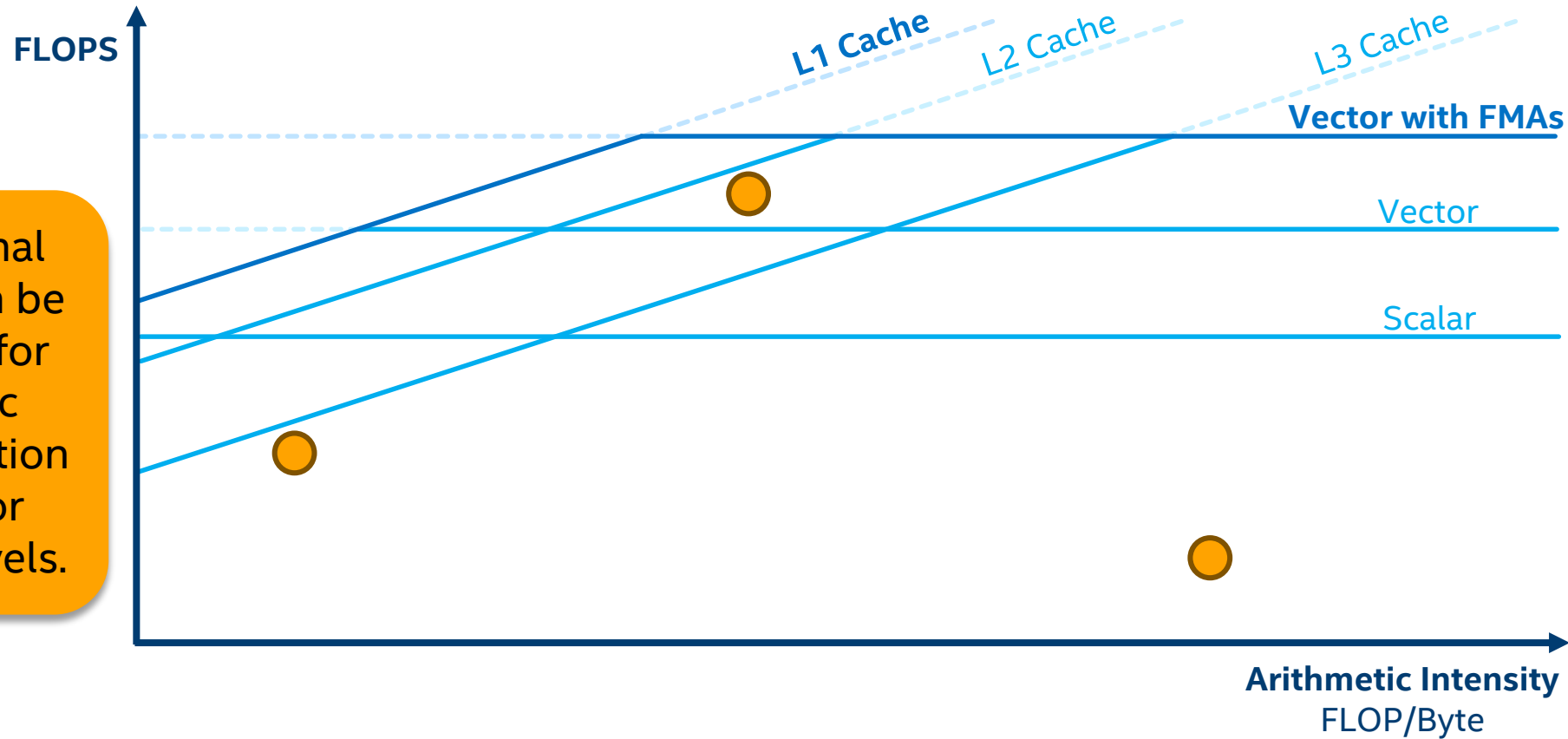
Ultimate Performance Limits



Optimization Notice

Copyright © 2017, Intel Corporation. All rights reserved.
*Other names and brands may be claimed as the property of others.

Sub-Roofs and Current Limits



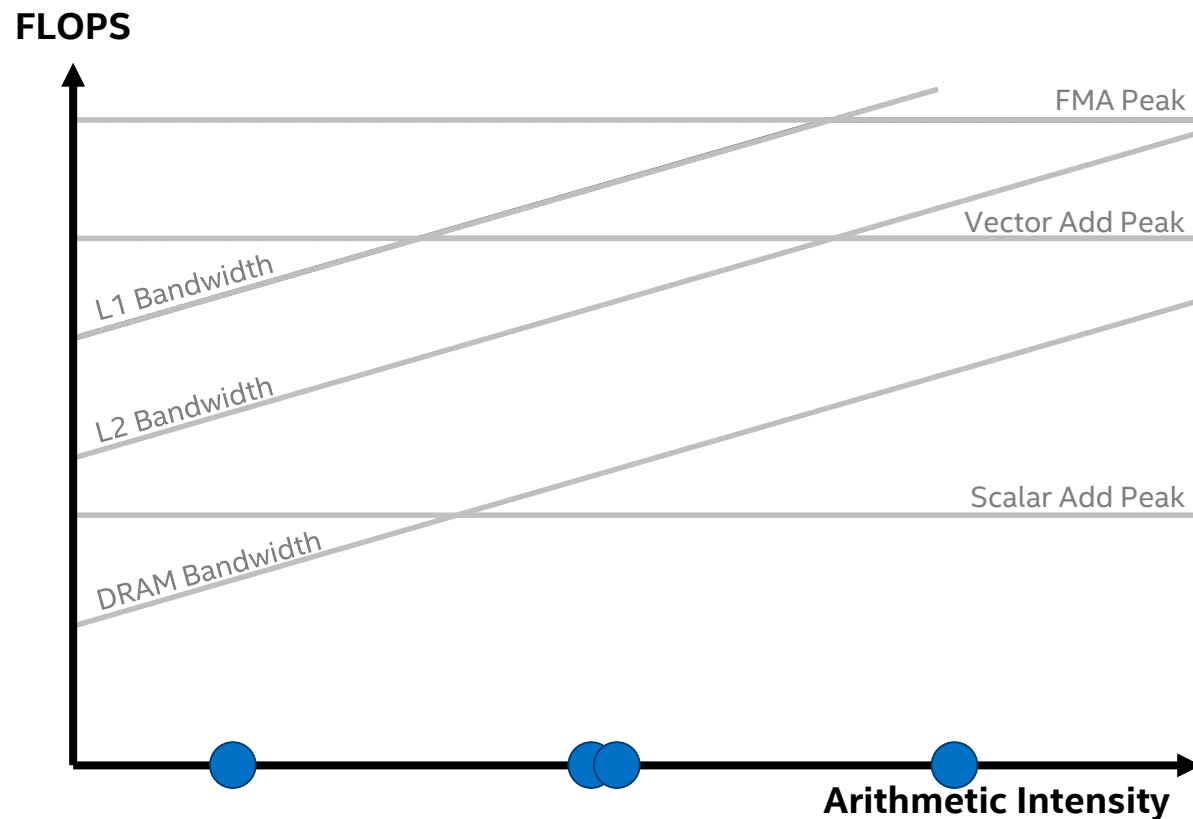
Additional roofs can be plotted for specific computation types or cache levels.

These sub-roofs can be used to help diagnose bottlenecks.

Cache-Aware Roofline Concept

- Prior to collecting data, Advisor runs quick benchmarks to measure hardware limitations.
 - Computational limitations
 - Memory Bandwidth limitations
- These form the performance “roofs”.
- Loops and functions have algorithms and therefore a specific AI.
- Their performance in FLOPS is also measured.
- Optimization changes performance. The goal is to go as far up as possible.

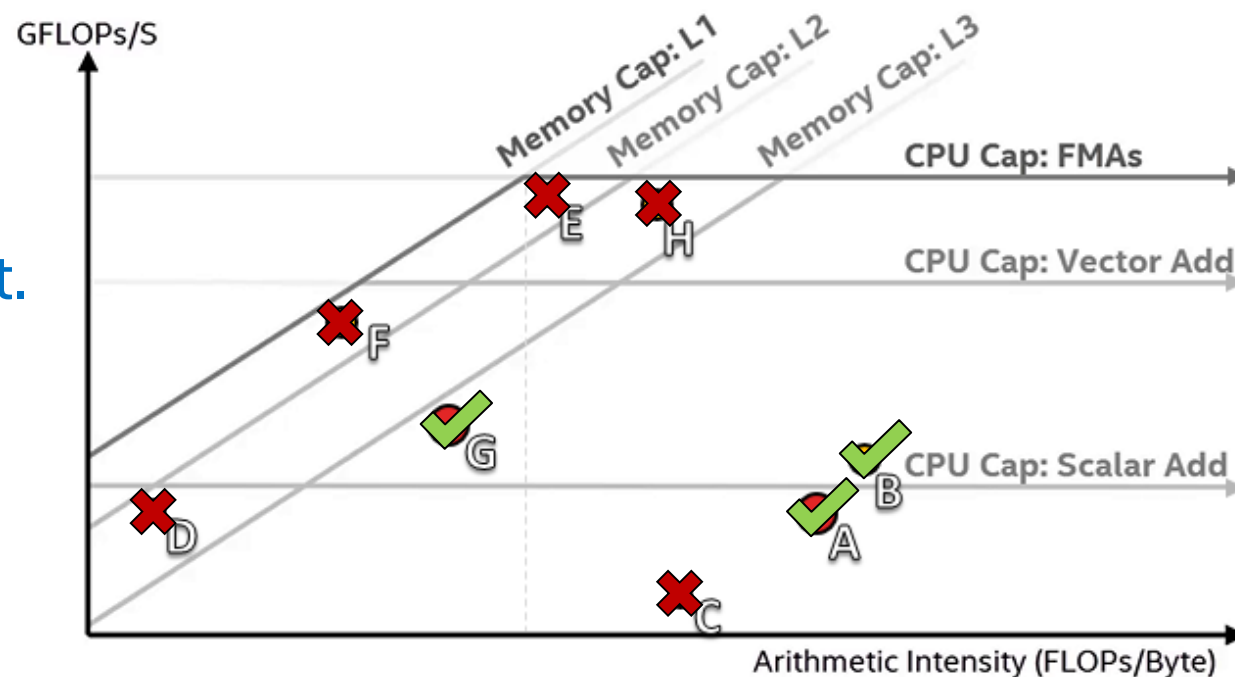
Video Available: *Roofline Analysis in Intel® Advisor 2017*



Roofline first proposed by University of California at Berkeley:
[Roofline: An Insightful Visual Performance Model for Multicore Architectures](#), 2009
Cache aware variant proposed by Technical University of Lisbon:
[Cache-Aware Roofline Model: Upgrading the Loft](#), 2013

Cache-Aware Roofline Usage

- Target the bigger, redder dots that are farthest below the top roofs first.
 - Big and red = takes up a lot of time.
 - Far below roofs = lots of room for improvement.
- Roofs below: unlikely bottlenecks, unless you're *just* above them.
- Roofs above: very likely bottlenecks! The closest roofs above are the most likely causes of bottlenecking.



Note:

Not every algorithm can break every roof!
Some algorithms inherently can't avoid certain bottlenecks.

Cache-Aware Roofline

Next Steps

If under or near a memory roof...

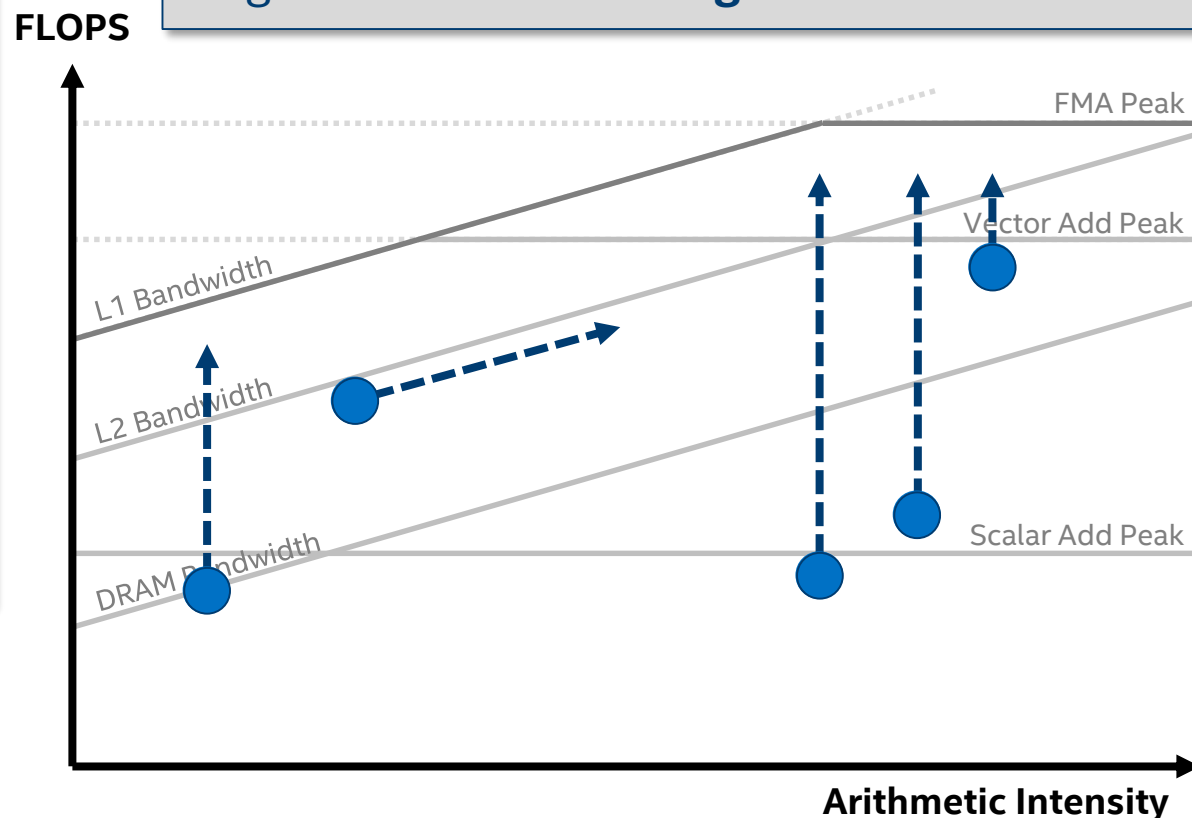
- Try a MAP analysis. Make any appropriate **cache optimizations**.
- If cache optimization is impossible, try **reworking the algorithm to have a higher AI**.

If Under the Vector Add Peak

Check “Traits” in the Survey to see if FMAs are used. If not, try altering your code or compiler flags to **induce FMA usage**.

If just above the Scalar Add Peak

Check **vectorization efficiency** in the Survey. Follow the recommendations to improve it if it's low.




If under the Scalar Add Peak...

Check the Survey Report to see if the loop vectorized. If not, try to **get it to vectorize** if possible. This may involve running Dependencies to see if it's safe to force it.

Memory Access Patterns Analysis

Collecting a MAP

- If you have low vector efficiency, or see that a loop did not vectorize because it was deemed “possible but inefficient”, you may want to run a MAP analysis.
- Advisor will also recommend a MAP analysis if it detects a possible inefficient access pattern.
- Memory access patterns affect vectorization efficiency because they affect how data is loaded into and stored from the vector registers.
- Select the loops you want to run the MAP on using the checkboxes. It may be helpful to reduce the problem size, as MAP only needs to detect patterns, and has high overhead.
 - Note that if changing the problem size requires recompiling, you will need to re-collect the survey before running MAP.

	Vector Issues
<input checked="" type="checkbox"/>	1 Possible inefficient memory access patterns present

Memory Access Patterns Analysis

Reading a MAP

- MAP is color coded by stride type. From best to worst:
 - **Blue** is unit/uniform (stepping by 1 or 0)
 - **Yellow** is constant (stepping a set distance)
 - **Red** is variable (a changing step distance)
- Click a loop in the top pane to see a detailed report below.
 - The strides that contribute to the loop are broken down in this table.
 - Important information includes the size of the stride, the variable being accessed, and the source.
 - Not all strides will come from your code!

Videos Available:
Stride and Memory Access Patterns
and
Memory Access 101



Site Location	Strides Distrib ...	Access Pa ...	Max. Site Footprint	Recommendations
[loop in main ...	76% / 0% / 24 ...	Mixed stri ...	64KB	
[loop in main ...	76% / 0% / 24 ...	Mixed stri ...	64KB	
[loop in main ...	70% / 6% / 24 ...	Mixed stri ...	564MB	
[loop in main ...	100% / 0% / 0 ...	All unit str ...	70KB	
[loop in main ...	33% / 67% / 0 ...	Mixed stri ...	616MB	1 Inefficient memory access pa

ID	Stride	Type	Source	Modules	Nested Func ...	Variable references
P1	36000	Constant stride	stride.cpp:49	stride.exe		tableA, tableB
P2	36000	Constant stride	stride.cpp:49	stride.exe		results
P7		Parallel site info ...	stride.cpp:47	stride.exe		
P19	0	Uniform stride	stride.exe:0x...	stride.exe	_svml_atan4	
P23	0	Uniform stride	svml_dispmd...	svml_dispmd.dll	_svml_atan2	
P1.	-12; -8; ...	Variable stride	svml_dispmd...	svml_dispmd.dll	_svml_atan2...	
P1.	-12; -8; ...	Variable stride	svml_dispmd...	svml_dispmd.dll	_svml_atan2...	

Dependencies Analysis

Vectorization Advisor

- Generally, you don't need to run Dependencies analysis unless Advisor tells you to. It produces recommendations to do so if it detects:
 - Loops that remained unvectorized because the compiler was playing it safe with autovectorization.
 - Outer loop vectorization opportunities
- Use the survey checkboxes to select which loops to analyze.
- If no dependencies are found, it's safe to force vectorization.
- Otherwise, use the reported variable read/write information to see if you can rework the code to eliminate the dependency.

☑ Recommendation: Confirm Confidence: 🟡 Need More Data
dependency is real
There is no confirmation that a real (proven) dependency is present in the loop. To confirm: Run a Dependencies analysis.

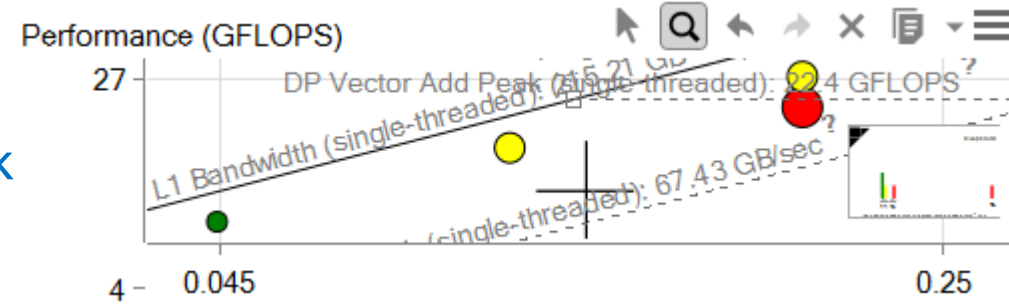
☑ Recommendation: Check Confidence: 🟠 Low
dependencies for outer loop
It is not safe to force vectorization without knowing that there are no dependencies. **Disable inner vectorization before check Dependency.** To check: Run a [Dependencies analysis](#).

Summary		Survey & Roofline		Refinement Reports	
Site Location		Loop-Carried Dependencies			
⊕	🔗 [loop in main at example.c...	🟢 No dependencies found			
⊕	🔗 [loop in main at example.c...	🔴 RAW:1			

Summary

- V** **Survey** – Find the most promising sites for threading, see the meat of the vectorization information, and get recommendations from Advisor.
- T**
- V** **Trip Counts & FLOPS** – Add to your Survey report to help fine-tune vector efficiency and capability, as well as unlock the powerful **Roofline** to visualize your bottlenecks and help direct your efforts.
- T**
- T** **Suitability** – Predict how well your proposed threading model will scale under certain conditions quickly and easily.
- V** **Dependencies** – Prove or disprove the existence of parallel dependencies and learn how to fix them.
- T**
- V** **Memory Access Patterns** – See how you traverse your data and how it affects your vector efficiency and cache bandwidth usage.

Function Call Sites and Loops	Self Time	Vectorized Loops			
		Vect..	Efficiency	Gai...	VL.
[loop in main at exa...	0.188s				
[loop in main at exa...	0.031s	AVX2	26%	2.09x	8
[loop in main at exa...	0.000s	AVX2	100%	8.00x	8



Site Performance Scalability Site Details

Scalability of Maximum Site Gain

Loop Iterations (Tasks) Modeling

Runtime Modeling

Apply

Site Location	Loop-Carried Dependencies	Strides Distribution
[loop in main..	No dependencies found	0% / 0% / 100%
[loop in main..	No information available	0% / 100% / 0%
[loop in main..	RAW:1	100% / 0% / 0%

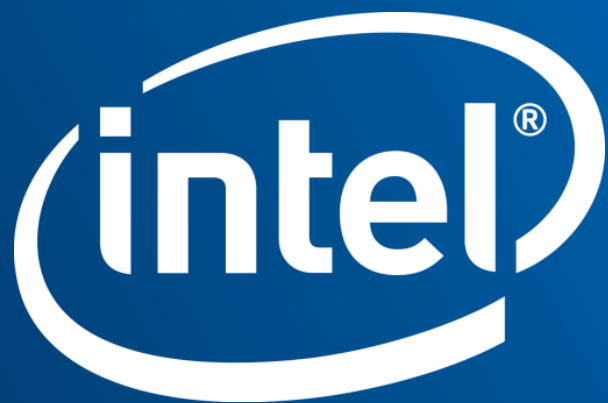
Threading Advisor Demo

Optimization Notice

Copyright © 2017, Intel Corporation. All rights reserved.

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



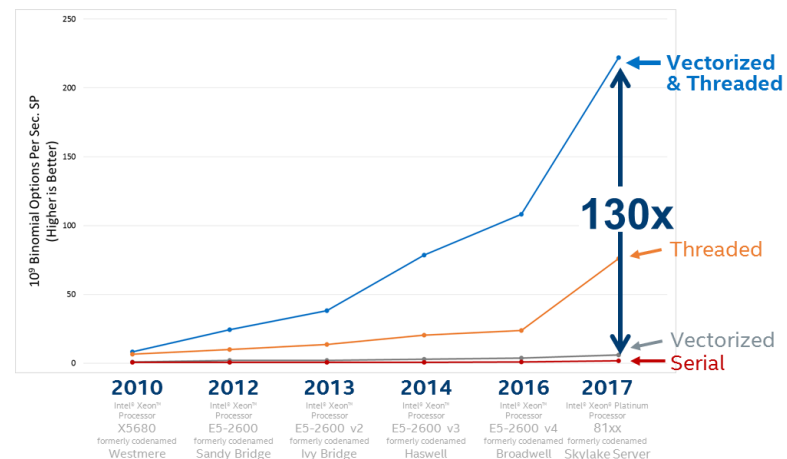


Configurations for 2010-2017 Benchmarks

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

Performance measured in Intel Labs by Intel employees



Platform Hardware and Software Configuration

	Platform	Unscaled Core Frequency	Cores/Socket	Num Sockets	L1 Data Cache	L2 Cache	L3 Cache	Memory	Memory Frequency	Memory Access	H/W Prefetchers Enabled	HT Enabled	Turbo Enabled	C States	O/S Name	Operating System	Compiler Version
WSM	Intel® Xeon™ X5680 Processor	3.33 GHz	6	2	32K	256K	12 MB	48 MB	1333 MHz	NUMA	Y	Y	Y	Disabled	Fedora 20	3.11.10-301.fc20	icc version 17.0.2
SNB	Intel® Xeon™ E5 2690 Processor	2.9 GHz	8	2	32K	256K	20 MB	64 GB	1600 MHz	NUMA	Y	Y	Y	Disabled	Fedora 20	3.11.10-301.fc20	icc version 17.0.2
IVB	Intel® Xeon™ E5 2697v2 Processor	2.7 GHz	12	2	32K	256K	30 MB	64 GB	1867 MHz	NUMA	Y	Y	Y	Disabled	RHEL 7.1	3.10.0-229.el7.x86_64	icc version 17.0.2
HSW	Intel® Xeon™ E5 2600v3 Processor	2.2 GHz	18	2	32K	256K	46 MB	128 GB	2133 MHz	NUMA	Y	Y	Y	Disabled	Fedora 20	3.15.10-200.fc20.x86_64	icc version 17.0.2
BDW	Intel® Xeon™ E5 2600v4 Processor	2.3 GHz	18	2	32K	256K	46 MB	256 GB	2400 MHz	NUMA	Y	Y	Y	Disabled	RHEL 7.0	3.10.0-123.el7.x86_64	icc version 17.0.2
BDW	Intel® Xeon™ E5 2600v4 Processor	2.2 GHz	22	2	32K	256K	56 MB	128 GB	2133 MHz	NUMA	Y	Y	Y	Disabled	CentOS 7.2	3.10.0-327.el7.x86_64	icc version 17.0.2
SKX	Intel® Xeon® Platinum 81xx Processor	2.5 GHz	28	2	32K	1024K	40 MB	192 GB	2666 MHz	NUMA	Y	Y	Y	Disabled	CentOS 7.3	3.10.0-514.10.2.el7.x86_64	icc version 17.0.2

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. For more information go to <http://www.intel.com/performance>

Optimization Notice

Copyright © 2017, Intel Corporation. All rights reserved.

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

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, & SSSE3 instruction sets & 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



Legal Disclaimer & Optimization Notice

Benchmark results were obtained prior to implementation of recent software patches and firmware updates intended to address exploits referred to as "Spectre" and "Meltdown". Implementation of these updates may make these results inapplicable to your device or system.

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. For more complete information visit www.intel.com/benchmarks.

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.

Copyright © 2018, Intel Corporation. All rights reserved. Intel, Xeon, Xeon Phi, Core, VTune, and the Intel logo 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