

**BERKELEY LAB** 

NATIONAL LABORATORY

# Performance Analysis and Optimization

## Samuel Williams

### <u>SWWilliams@lbl.gov</u>

**Computational Research Division** Lawrence Berkeley National Lab



## Instrumentation

- Application-specific (manual) instrumentation...
  - Most robust
  - Minimal overhead (omp\_get\_time) ٠
  - Insensitive to sampling effects ٠
  - Application-specific knowledge can different based on usage (e.g. different levels of MG) ۲
  - High effort / large reward
- Auto-instrumentation (TAU, Advisor, Vtune)...
  - Minimal effort
  - Integrated visualization
  - Sampling effects can confuse performance analysis
  - Using the same function many different ways can confuse analysis
  - Can have high overhead (Advisor/Vtune) ۲



## **Roofline Model**

- The **Roofline Model** is a throughputoriented performance model...
  - Tracks rates not time •
  - Augmented with Little's Law (concurrency = latency\*bandwidth)
  - Independent of ISA and architecture (applies to CPUs, GPUs, Google TPUs<sup>1</sup>, etc...)
- Informs developers which routines are underperforming the processor's capabilities == which routines to optimize



https://crd.lbl.gov/departments/computer-science/PAR/research/roofline



## **Use by NESAP**

- NESAP is the NERSC KNL application readiness project.
- NESAP used Roofline to drive optimization and analysis on KNL
  - Bound performance expectations (ERT)
  - Use Vtune to quantify DDR and MCDRAM data movement
  - Compare KNL data movement to Haswell (sea of private/coherent L2's vs. unified L3)  $\bullet$
  - Understand importance of vectorization lacksquare
  - Doerfer et al., "Applying the Roofline Performance Model to the Intel Xeon Phi Knights Landing Processor", Intel Xeon Phi User Group Workshop (IXPUG), June 2016.
  - Barnes et al. "Evaluating and Optimizing the NERSC Workload on Knights  $\bullet$ Landing", Performance Modeling, Benchmarking and Simulation of High Performance Computer Systems (PMBS), November 2016.



## **Roofline for NESAP Codes**



- Roofline Model
- •wo/FMA
- Original
- w/Tiling
- w/Tiling+Vect

- Roofline Model
- -wo/FMA
- Original
- w/Tiling
- w/Tiling+Vect



## Intel Advisor

## Intel Advisor is a performance analysis tool (evolved from vector advisor)

### Background

- https://software.intel.com/en-us/intel-advisor-xe
- https://software.intel.com/en-us/articles/getting-started-withintel-advisor-roofline-feature
- https://www.youtube.com/watch?v=h2QEM1HpFgg

### **Running Advisor on NERSC Systems**

<u>http://www.nersc.gov/users/software/performance-and-debugging-tools/advisor/</u>





NoMachine - NERSC

### /global/cscratch1/sd/tkoskela/dram\_roofline/stencil/advi.stencil.aug2.16 - Intel Advisor <@cori05> <2>

File View Help											
Welcome e000 (read-only) 🗶											
📾 Elapsed time: 50.50s 😽 🙆 Vecto	orized	σN	Not Vectorized	Ø						OFF	Sma
FILTER: All Modules - All Sources -	Loop	s Ar	nd Functions	-	A II T	hreads 👻					
🌳 Summary 🛛 💐 Survey & Roofline 🐲	Refine	mei	nt Reports								
징 ⊡	۵	<u>ଡ</u> 	Self Time	T. Ti.	т.	FLOPS GFLOPS <del>▼</del>	AI		W N	Vectorized Vector I	Loop
☐ ☐ [loop in bench stencil ver4\$o			159.595s	15	.v.	23.083	0.117			AVX2	100
E [loop in bench_stencil_ver3\$o		<b>⊚</b> 1	159.953s	15	.ν.	16.274	0.117			AVX2	899
		<b>?</b> 1	160.035s	16	.ν.	15.662	0.117			AVX2	809
<b>⊞</b> ©[loop in bench_stencil_ver1\$o		<b>?</b> 1	159.307s	15	ν.	10.218	0.117			AVX2	809
⊠©[loop in bench_stencil_ver0		01	157.994s 🗆	] 1	s.	9.009	0.117	(	3		
Source Top Down Code Analytics File: stencil_v2.c:29 bench_stencil_v	Assem	bly om	♀ Recomme p\$parallel_fo	nda or@2	tion 25	s 🖻 Why No Vect	orization?				
Lin.	So	urce	2				Total Time	%	L	.oop/Functi	on Tin
25 #pragma omp parallel for							9.890s				
26 <sup>⊞</sup> for(k=1;k <dim+1;k++){< td=""><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td></dim+1;k++){<>											
2/ <sup>™</sup> for(j=1; j <dim+1; j++){<="" td=""><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td></dim+1;>											
20 = for(i=1:i < dim+1:i++)							102 403s			1	157.90
30 int ijk = $i*iStride + 1$	i*iStr	ide	e + k*kStrid	e;			102.1033				137.3
31 new[ijk] = -6.0*old[ijk]	ĸ		]	,			53.651s				
32 + old[ij	k-iStr	ide	e]								
				9	Sele	cted (Total Time):	102.403s				
🧸 🔚 📕 Madvixe-gui			₃ 🕏 emacs	-gtk	@C	ori05-bond0.224	·	cori :			

		$\odot$	$\otimes$
			•
rt Mode IN1	e° Tel Ai	DVISOR 20	a, 017
5			
ency	Gai	n Es	V
0%	5.2	7x	4
<mark>%</mark>	3.5	4	
%	3.2	4	
%	3.2	4	
	Ve	Traite	
	/0	ITALS	
94s 📃			
		FMA	



	NoMachine - NERSC									
🞽 🕐 /global/cscratch1/sd/tkoskela/dram_roofline/stencil/advi.dram.stencil.aug2.16 - Intel Advisor <@cori05> 👘 🐼 📀 📀									$\odot$ $\otimes$ $\otimes$	
File View	Help									
Welcome	Welcome e000 x									
<ul> <li>Elapsed time: 50.40s</li> <li>Vectorized</li> <li>Not Vectorized</li> <li>MKL</li> <li>FILTER: All Modules</li> <li>All Sources</li> <li>Loops And Functions</li> <li>All Threads</li> </ul>							Q VISOR 2018			
	Summary Survey & Roome Reports									
Perfor	mance (G	FLOPS)			, ose single-rin				ata	
1688.1	18- 1 1 Band	width: 1.1e+4 GB/	sec				5 = = = = = = = =	DP Vector FMA Peak: 16	88.18	FLOPS
<b>→</b>	LI Dans	width: 3542.01-GE	3/sec γ						22.07	FLOPS
	B	width: 1003.46-GE	3/ <u>5eC</u>			<b>_</b>		Scalar Add Peak: 1	15 16 0	FLOPS?
	L3 Daris	00.50	CRISEC ?	·····					10,10 0	
	DRAM E	andwidth: 128.58	GD/300		0					d it langue for fail
	,			0	[loop in bencl	n_stencil_ve	r4\$omp	\$parallel_for@193 at s	tencil_	v2.c:193]
0.6	Total Performance: 18.98 GFLOPS									
	Total L1 Arithmetic Intensity: 0.41 FLOP/Byte									·
		0.05		· <u> </u> ·	Total Flansed	Time: <b>0.000</b> S				
	Self Elaps	ed Time: <b>0.000</b> s	Total Elap	sed Time: <b>9.956</b> s		Time. <b>3.330</b> 3	,		_	
Source 7	lop Down	Code Analytics	Assembly	Recommendations	🖣 Why No Vecto	orization?				
File: sten	cil_v2.c:1	93 bench_stenc	il_ver4\$om	p\$parallel_for@193	I	I				
Lin.			Source			Total Time	%	Loop/Function Time	%	Traits
191 StartTime = omp_get_wtime();										
192 <sup>™</sup> while(ElapsedTime < TIME){										
193 # #pragma omp parallel for schedule(static,1)					8.004ms 160161.000ms					
194 for(tile=0;tile <jtiles*ktiles;tile++){< th=""><th>8.004ms</th><th></th><th></th><th></th><th></th></jtiles*ktiles;tile++){<>					8.004ms					
Selected (Total Time).					0.0041115					
	Advixe-ou	ıi		⇒ Semacs-atk@cori	05-bond0 224		cori ·			
	nu inc gu			3 Chinaco dettactori	00 00100.224		2011.			



### .....



## <u>likwid</u>



• Used to characterize AMReX ECP apps...







## **Threading vs. Processes**

- Threads provide no inherent compute advantage over processes
- Threads incur additional overhead (omp parallel, single, barrier, ...) == slower in the perfectly parallel world
- Threads provide easy access to shared memory...
  - Some codes are easier to parallelize with threads than MPI
  - Easier to avoid data duplication (memory requirements) with threads than processes
  - Threads can access shared data in cache rather than copying data between processes  $\bullet$
  - Using threads simplifies topology-aware MPI process mapping (MPICH\_RANK\_REORDER) is often insufficient)... topology-aware is still important on Aries/Dragonfly and IB/FatTrees
  - Using threads provides on-ramp to GPUs and other emerging architectures



## **Threading Experiments/Analysis**

- Fix process concurrency, increase threading ...
  - Ideally, function runtime should scales as O(1/omp\_num\_threads)
  - Identify functions that plateau (saturation)
  - Identify functions that are flat (sequential bottlenecks) lacksquare
  - Identify functions that increase (duplicated work... e.g. use of Fortran's sum(a(:))) ullet
- Fix hardware concurrency (cores), increase threads while reducing processes.... (32x1, 16x2, 8x4, 4x8, 2x16, 1x32)...
  - Distinguish true flat MPI (no –fopenmp) from hybrid w/1 thread (-fopenmp + OMP NUM THREADS=1) == threading overhead
  - Identify functions that are flat, better with threads, or better with processes ۲
  - Especially useful for communication routines





## **KNL-specific issues**

- GNU and Intel runtimes treat OMP\_PLACES/PROC\_BIND differently...
  - use same compiler for all threaded routines / know which settings to use
  - (I use KMP\_AFFINITY with Intel)
- If you can fit in 16GB, run in quadflat...
  - avoids cache aliasing issues where some nodes are slower than others
  - Can manifest as abnormally high MPI\_Wait times and degraded scalability
  - Use srun ... numactl –m1 ./a.out ... (or use –p1 or no numactl but allocate key data in HBM)
- use large contiguous blocks with 2M pages for data and MPI buffers...
  - Better performance for complex (less streaming) memory access patterns
  - minimizes NIC TLB pressure == higher MPI bandwidth
- Examine vectorization reports...
  - KNL is very slow if not vectorized •





# Questions



