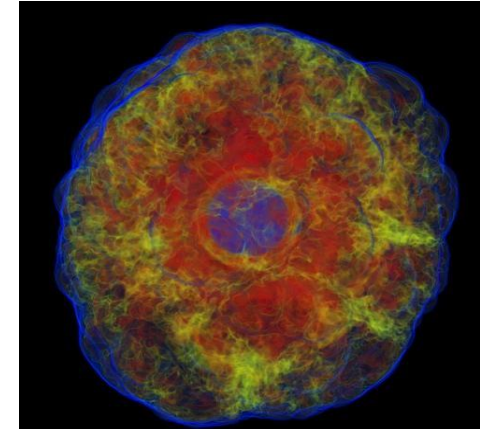
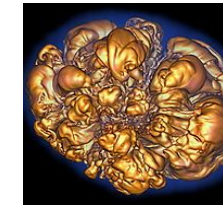
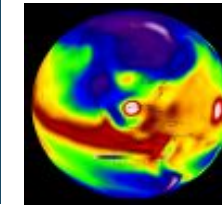
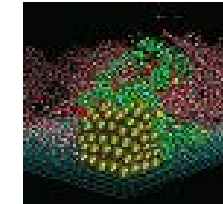
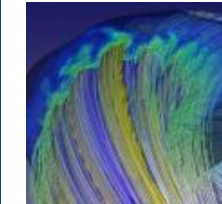
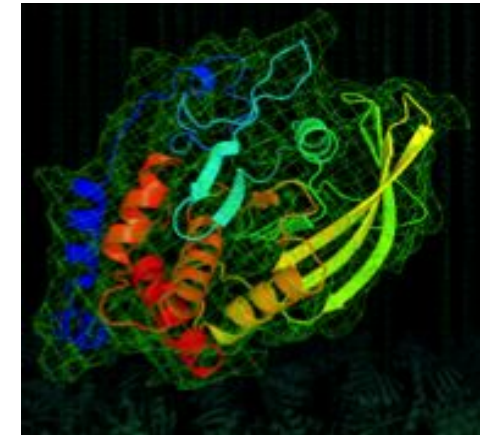
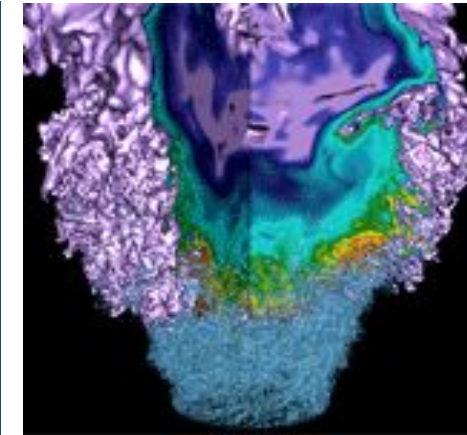


Optimizing Applications: an Ant Farm Approach



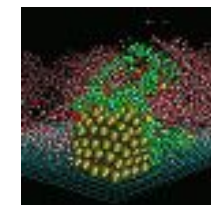
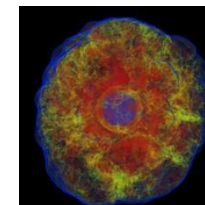
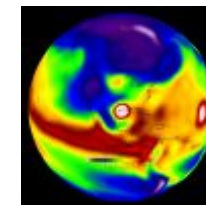
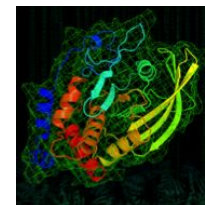
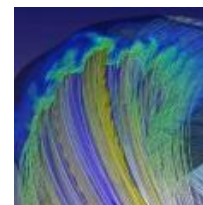
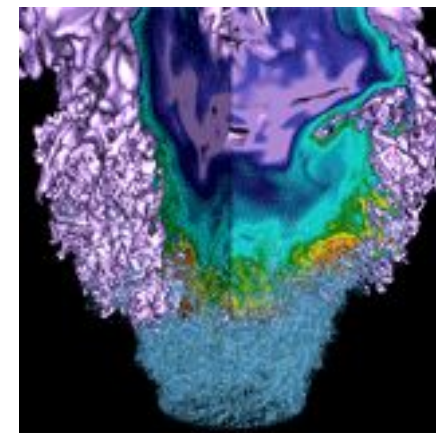
Jack Deslippe
August 2016

Agenda



1. Many Core vs Multi Core
2. Performance Optimization Concepts for Many Core
3. Performance Optimization Strategy for Many Core
4. Example Case Studies

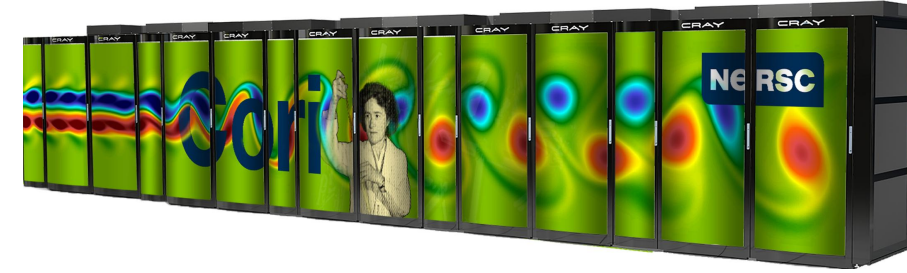
Many Core HPC Systems



Many Core Systems Coming to NERSC, ALCF and More



- **NERSC's Cori will begin to transition the workload to more energy efficient architectures**
- **Cray XC system with over 9300 Intel Knights Landing (Xeon-Phi) compute nodes**
 - Self-hosted, (not an accelerator) manycore processor with 68 cores per node
 - On-package high-bandwidth memory



System named after Gerty Cori, Biochemist and first American woman to receive the Nobel prize in science.

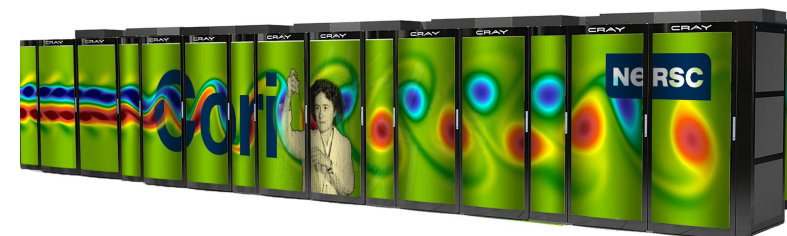
What is different about Many-Core

Edison (Multi-Core):

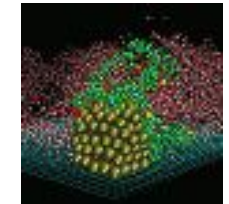
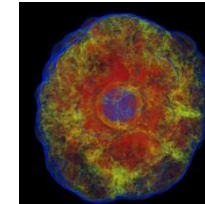
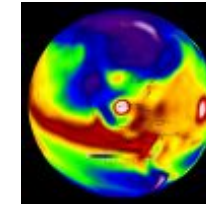
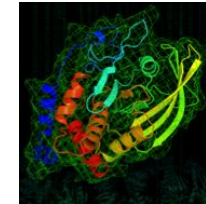
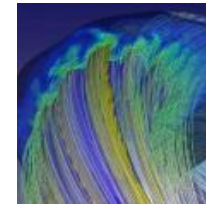
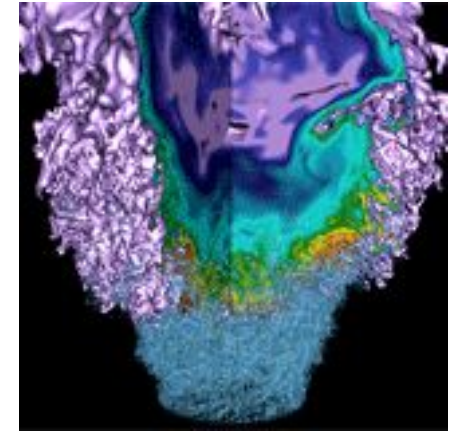
- 5000+ Ivy Bridge Nodes
- **12 Cores Per CPU**
- 24 Virtual Cores Per CPU
- 2.4-3.2 GHz
- Can do **4 Double Precision Operations per Cycle** (+ multiply/add)
- 2.5 GB of Memory Per Core
- **~100 GB/s Memory Bandwidth**

Cori (Many-Core):

- 9000+ Knights Landing Nodes
- **68 Physical Cores Per CPU**
- Up to 272 Virtual Cores Per CPU
- Much slower GHz
- Can do **8 Double Precision Operations per Cycle** (+ multiply/add)
- < 0.3 GB of Fast Memory Per Core
< 2 GB of Slow Memory Per Core
- **Fast Memory has ~ 4-5x DDR4 Bandwidth**



Basic Optimization Concepts



MPI + X?

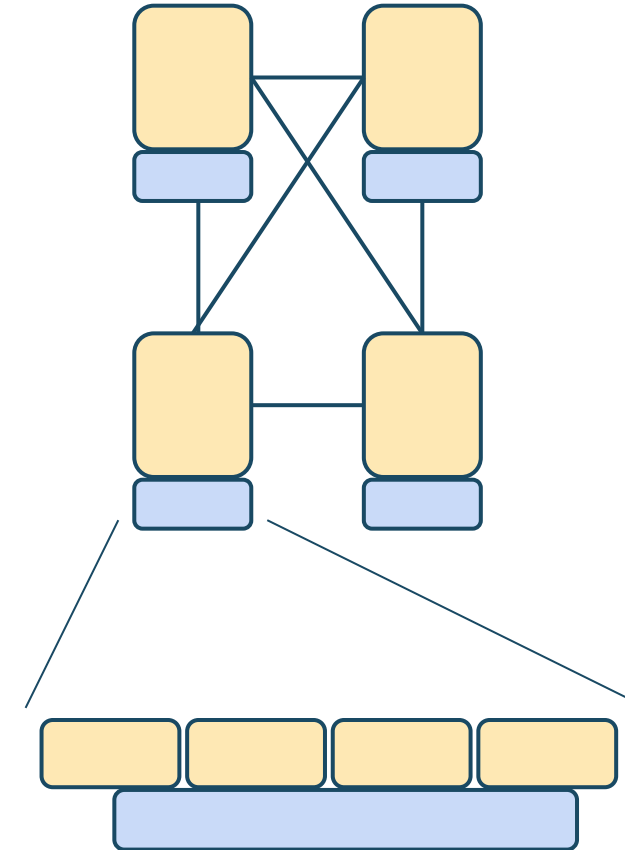
Need to explicitly consider both inter and on-node parallelism in application.

Existing applications may suffer from:

- Memory overhead due to duplicated data in traditional MPI tasks
- Lack of SIMD/Vectorization expressiveness in app.
- Potential MPI latency in all-to-all communication patterns

Possible Solutions:

MPI+MPI, MPI+OpenMP, PGAS (MPI+PGAS), Task Based Programming



PARATEC Use Case For OpenMP

PARATEC computes parallel FFTs across all processors.

Involves MPI all-to-all communication (small messages, latency bound).

Reducing the number of MPI tasks in favor OpenMP threads makes large improvement in overall runtime.

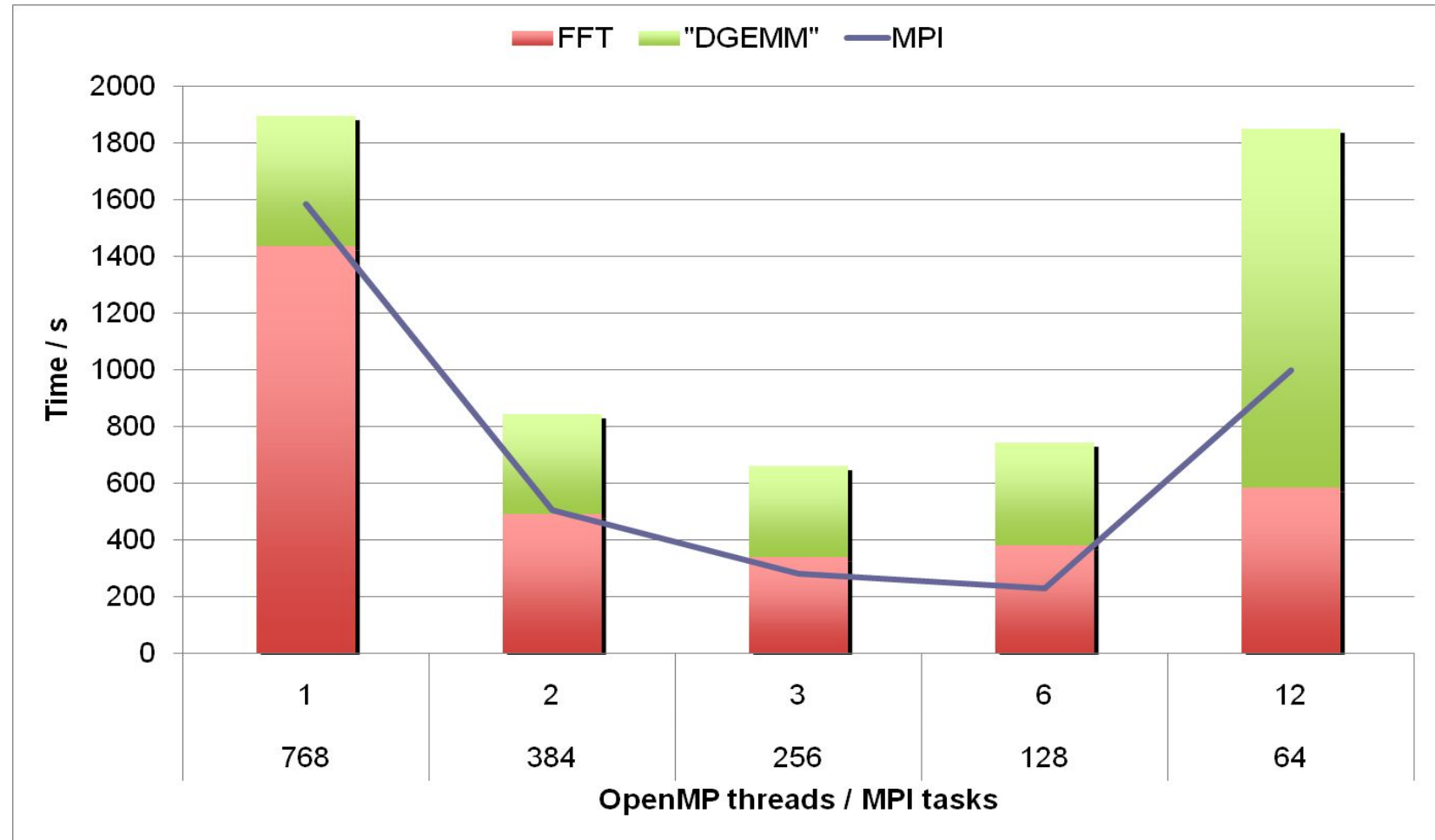



Figure Courtesy of Andrew Canning

There is a another important form of on-node parallelism

```
do i = 1, n  
    a(i) = b(i) + c(i)  
enddo
```


$$\begin{pmatrix} a_1 \\ \dots \\ a_n \end{pmatrix} = \begin{pmatrix} b_1 \\ \dots \\ b_n \end{pmatrix} + \begin{pmatrix} c_1 \\ \dots \\ c_n \end{pmatrix}$$

Vectorization: CPU does identical operations on different data; e.g., multiple iterations of the above loop can be done concurrently. Works best with long/aligned vectors.

Vectorization

There is a another important form of on-node parallelism

```
do i = 1, n
  a(i) = b(i) + c(i)
enddo
```



$$\begin{pmatrix} a_1 \\ \dots \end{pmatrix} = \begin{pmatrix} b_1 \\ \dots \end{pmatrix} + \begin{pmatrix} c_1 \\ \dots \end{pmatrix}$$

Intel Xeon Sandy-Bridge/Ivy-Bridge:	4 Double Precision Ops Concurrently
Intel Xeon Phi:	8 Double Precision Ops Concurrently

Vectoriz
above l

of the

Things that prevent vectorization in your code

Compilers want to “vectorize” your loops whenever possible. But sometimes they get stumped. Here are a few things that prevent your code from vectorizing:

Loop dependency:

```
do i = 1, n
  a(i) = a(i-1) + b(i)
enddo
```

Task forking:

```
do i = 1, n
  if (a(i) < x) cycle
  if (a(i) > x) ...
enddo
```

Things that prevent vectorization in your code



Example From NERSC User Group Hackathon - (Astrophysics Transport Code)

```
for (many iterations) {  
    ... many flops ...  
    et = exp(outcome1)  
    tt = pow(outcome2,3)  
    IN = IN * et +tt  
}
```

Things that prevent vectorization in your code

Example From NERSC User Group Hackathon - (Astrophysics Transport Code)

```
for (many iterations) {  
  ... many flops ...  
  et = exp(outcome1)  
  tt = pow(outcome2,3)  
  IN = IN * et +tt  
}
```



```
for (many iterations) {  
  ... many flops ...  
  et(i) = exp(outcome1)  
  tt(i) = pow(outcome2,3)  
}  
for (many iterations) {  
  IN = IN * et(i) + tt(i)  
}
```

Things that prevent vectorization in your code

Example From NERSC User Group Hackathon - (Astrophysics Transport Code)

```
for (many iterations) {  
  ... many flops ...  
  et = exp(outcome1)  
  tt = pow(outcome2,3)  
  IN = IN * et +tt  
}
```



```
for (many iterations) {  
  ... many flops ...  
  et(i) = exp(outcome1)  
  tt(i) = pow(outcome2,3)  
}  
for (many iterations) {  
  IN = IN * et(i) + tt(i)  
}
```

30% speed up for entire application!

Things that prevent vectorization in your code

Original

```
real(8),dimension
  (5,(col_f_nvr-1)*(col_f_nvz-1),
  (col_f_nvr-1)*(col_f_nvz-1)) :: Ms

do index_ip = 1, mesh_Nzml
  do index_jp = 1, mesh_Nrml
    index_2dp = index_jp+mesh_Nrml*(index_ip-1)

    tmp_vol = cs2%local_center_volume(index_jp)
    tmp_f_half_v = f_half(index_jp, index_ip) *
    tmp_vol
    tmp_dfdr_v = dfdr(index_jp, index_ip) *
    tmp_vol
    tmp_dfdz_v = dfdz(index_jp, index_ip) *
    tmp_vol

    tmpr(1:3)= tmpr(1:3)+
    Ms(1:3,index_2dp,index_2D)* tmp_f_half_v
    tmpr(5) = tmpr(5) +
    Ms(4,index_2dp,index_2D)*tmp_dfdr_v +
```

Optimized

```
real (8),dimension
  ((col_f_nvr-1),5,(col_f_nvz-1),
  (col_f_nvr-1)*(col_f_nvz-1)) :: Ms

do index_ip = 1, mesh_Nzml
  do index_jp = 1, mesh_Nrml
    index_2dp = index_jp+mesh_Nrml*(index_ip-1)
    tmp_vol = cs2%local_center_volume(index_jp)
    tmp_f_half_v = f_half(index_jp, index_ip) *
    tmp_vol
    tmp_dfdr_v = dfdr(index_jp, index_ip) * tmp_vol
    tmp_dfdz_v = dfdz(index_jp, index_ip) * tmp_vol

    tmpr(index_jp,1) = tmpr(index_jp,1) +
    Ms(index_jp,1,index_ip,index_2D)*
    tmp_f_half_v
    tmpr(index_jp,2) = tmpr(index_jp,2) +
    Ms(index_jp,2,index_ip,index_2D)*
    tmp_f_half_v
    tmpr(index_jp,3) = tmpr(index_jp,3) +
    Ms(index_jp,3,index_ip,index_2D)*
    tmp_f_half_v
    tmpr(index_jp,5) = tmpr(index_jp,5) +
    Ms(index_jp,4,index_ip,index_2D)*
    tmp_dfdz_v
    + Ms(index_ip,2,index_ip,index_2D)*
    tmp_dfdr_v
```

Example From Cray COE Work on XGC1

Things that prevent vectorization in your code

Original

```
real(8),dimension  
(5,(col_f_nvr-1)*(col_f_nvz-1),  
(col_f_nvr-1)*(col_f_nvz-1)) :: Ms  
  
do index_ip = 1, mesh_Nzml  
  do index_jp = 1, mesh_Nrml  
    index_2dp = index_jp+mesh_Nrml*(index_ip-1)  
  
    tmp_vol = cs2%local_center_volume(index_jp)  
    tmp_f_half_v = f_half(index_jp, index_ip) *  
    tmp_vol  
    tmp_dfdr_v = dfdr(index_jp, index_ip) *  
    tmp_vol  
    tmp_dfdz_v = dfdz(index_jp, index_ip) *  
    tmp_vol  
  
    tmpr(1:3)= tmpr(1:3)+  
    Ms(1:3,index_2dp,index_2D)* tmp_f_half_v  
    tmpr(5) = tmpr(5) +  
    Ms(4,index_2dp,index_2D)*tmp_dfdr_v +
```

Optimized

```
real (8),dimension  
((col_f_nvr-1),5,(col_f_nvz-1),  
(col_f_nvr-1)*(col_f_nvz-1)) :: Ms  
  
do index_ip = 1, mesh_Nzml  
  do index_jp = 1, mesh_Nrml  
    index_2dp = index_jp+mesh_Nrml*(index_ip-1)  
    tmp_vol = cs2%local_center_volume(index_jp)  
    tmp_f_half_v = f_half(index_jp, index_ip) *  
    tmp_vol  
    tmp_dfdr_v = dfdr(index_jp, index_ip) * tmp_vol  
    tmp_dfdz_v = dfdz(index_jp, index_ip) * tmp_vol  
  
    tmpr(index_jp,1) = tmpr(index_jp,1) +  
    Ms(index_jp,1,index_ip,index_2D)*  
    tmp_f_half_v  
    tmpr(index_jp,2) = tmpr(index_jp,2) +  
    Ms(index_jp,2,index_ip,index_2D)*  
    tmp_f_half_v  
    tmpr(index_jp,3) = tmpr(index_jp,3) +  
    Ms(index_jp,3,index_ip,index_2D)*  
    tmp_f_half_v  
    tmpr(index_jp,5) = tmpr(index_jp,5) +  
    Ms(index_jp,4,index_ip,index_2D)*  
    tmp_dfdz_v  
    + Ms(index_ip,2,index_ip,index_2D)*
```

Example From Cray COE Work on XGC1

~40% speed up
for kernel

Memory Bandwidth

Consider the following loop:

```
do i = 1, n
  do j = 1, m
    c = c + a(i) * b(j)
  enddo
enddo
```

Assume, n & m are very large such that a & b don't fit into cache.

Then,

During execution, the **number of loads From DRAM** is

$$n*m + n$$

Memory Bandwidth

Consider the following loop:

```
do i = 1, n
  do j = 1, m
    c = c + a(i) * b(j)
  enddo
enddo
```

Assume, n & m are very large such that a & b don't fit into cache.

Then,

During execution, the **number of loads From DRAM** is

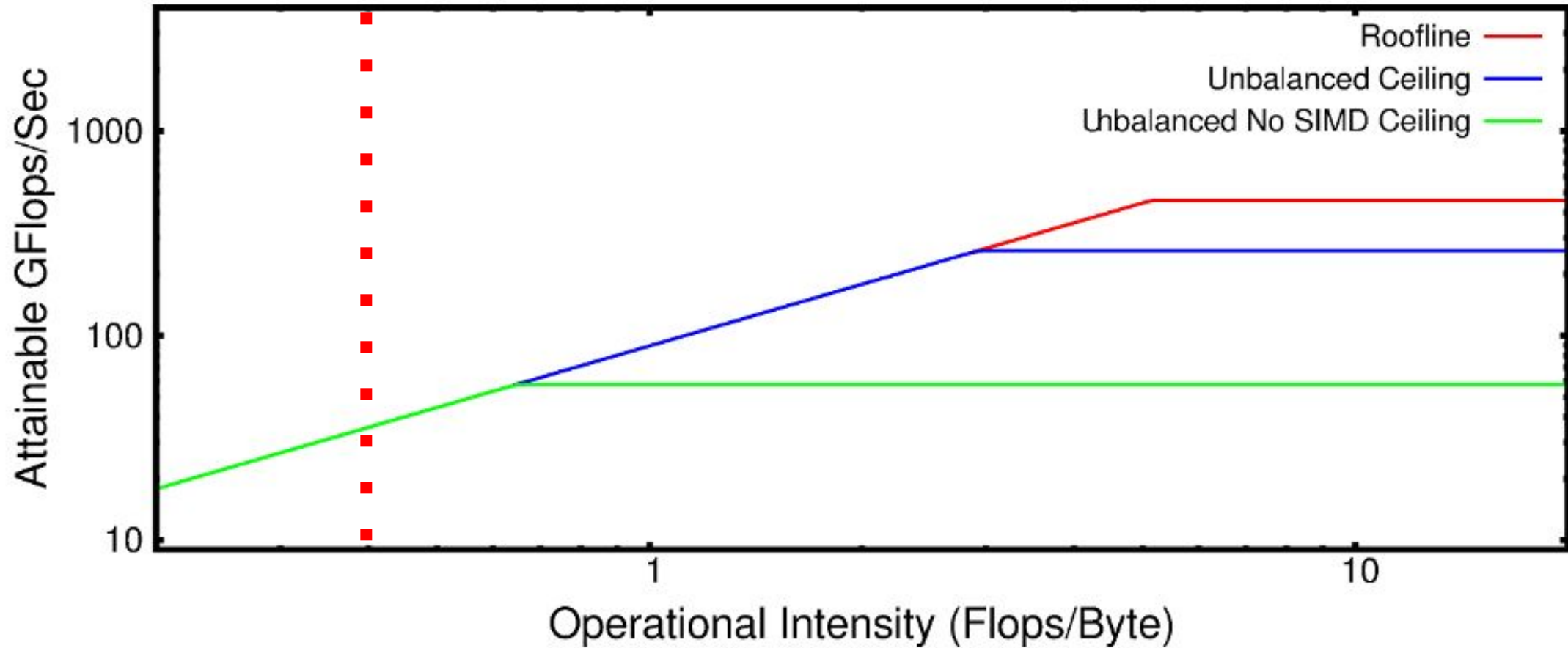
$$n*m + n$$

Requires 8 bytes loaded from DRAM per FMA (if supported). Assuming 100 GB/s bandwidth on Edison, we can **at most achieve 25 GFlops/second** (2 Flops per FMA)

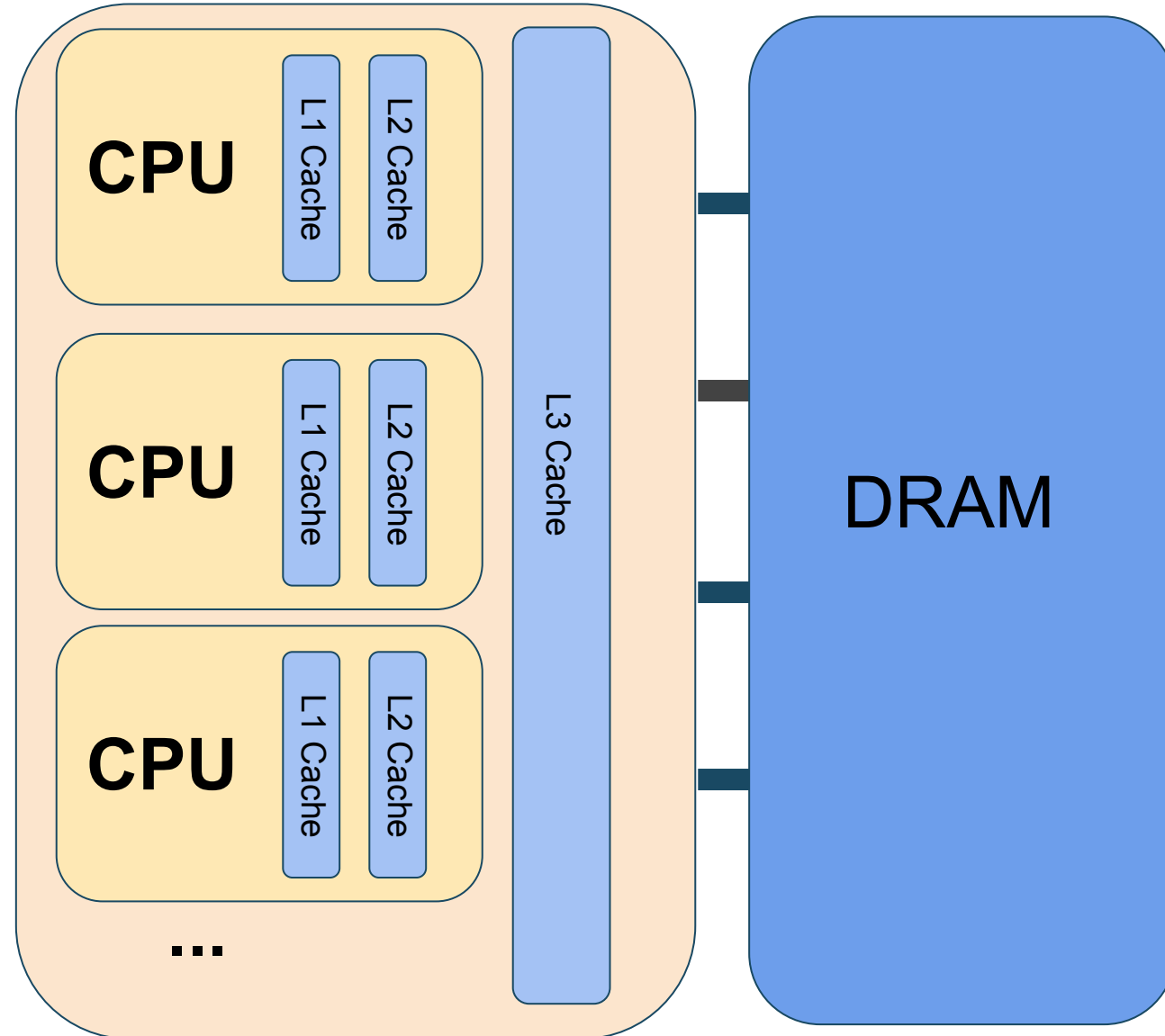
Much lower than 460 GFlops/second peak on Edison node. **Loop is memory bandwidth bound.**

Roofline Model For Edison

Edison Node Roofline Based on Stream of 89GB/s and Peak Flops of 460 GFlop/Sec



Processor Memory Hierarchy on MultiCore



Improving Memory Locality

Improving Memory Locality. Reducing bandwidth required.

```
do i = 1, n
  do j = 1, m
    c = c + a(i) * b(j)
  enddo
enddo
```



```
do jout = 1, m, block
  do i = 1, n
    do j = jout, jout+block
      c = c + a(i) * b(j)
    enddo
  enddo
enddo
```

Loads From DRAM:

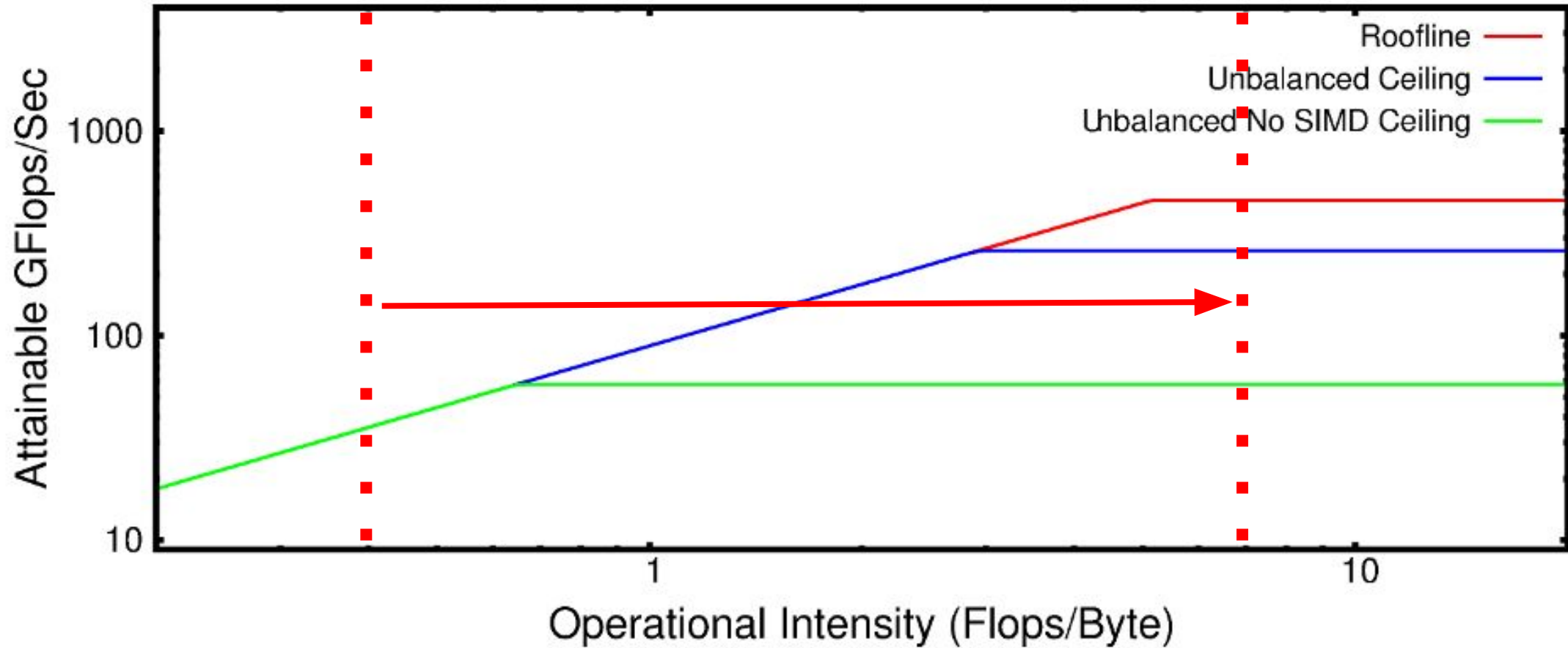
$$n*m + n$$

Loads From DRAM:

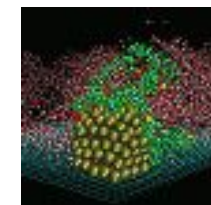
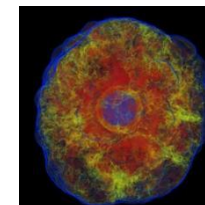
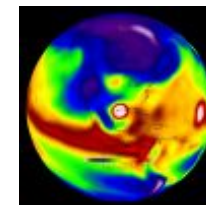
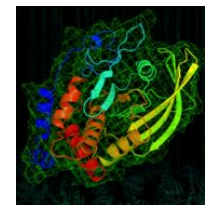
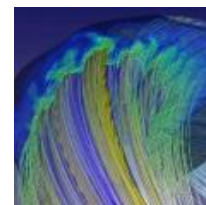
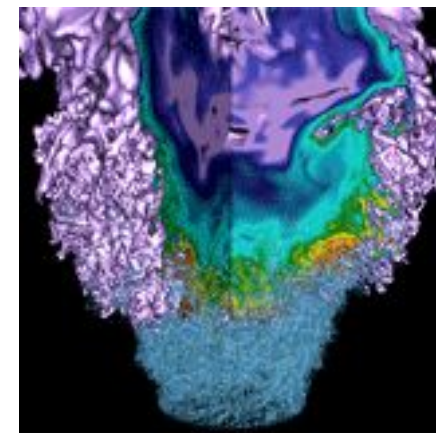
$$m/block * (n+block) = n*m/block + m$$

Improving Memory Locality Moves you to the Right on the Roofline

Edison Node Roofline Based on Stream of 89GB/s and Peak Flops of 460 GFlop/Sec



Optimization Strategy



The Ant Farm!

OpenMP scales only to 4 Threads

large cache miss rate

Code shows no improvements when turning on vectorization

50% Walltime is IO

Communication dominates beyond 100 nodes



Compute intensive doesn't vectorize

Memory bandwidth bound kernel

IO bottlenecks

MPI/OpenMP Scaling Issue

Can you use a library?

Increase Memory Locality

Utilize High-Level IO-Libraries. Consult with NERSC about use of Burst Buffer.

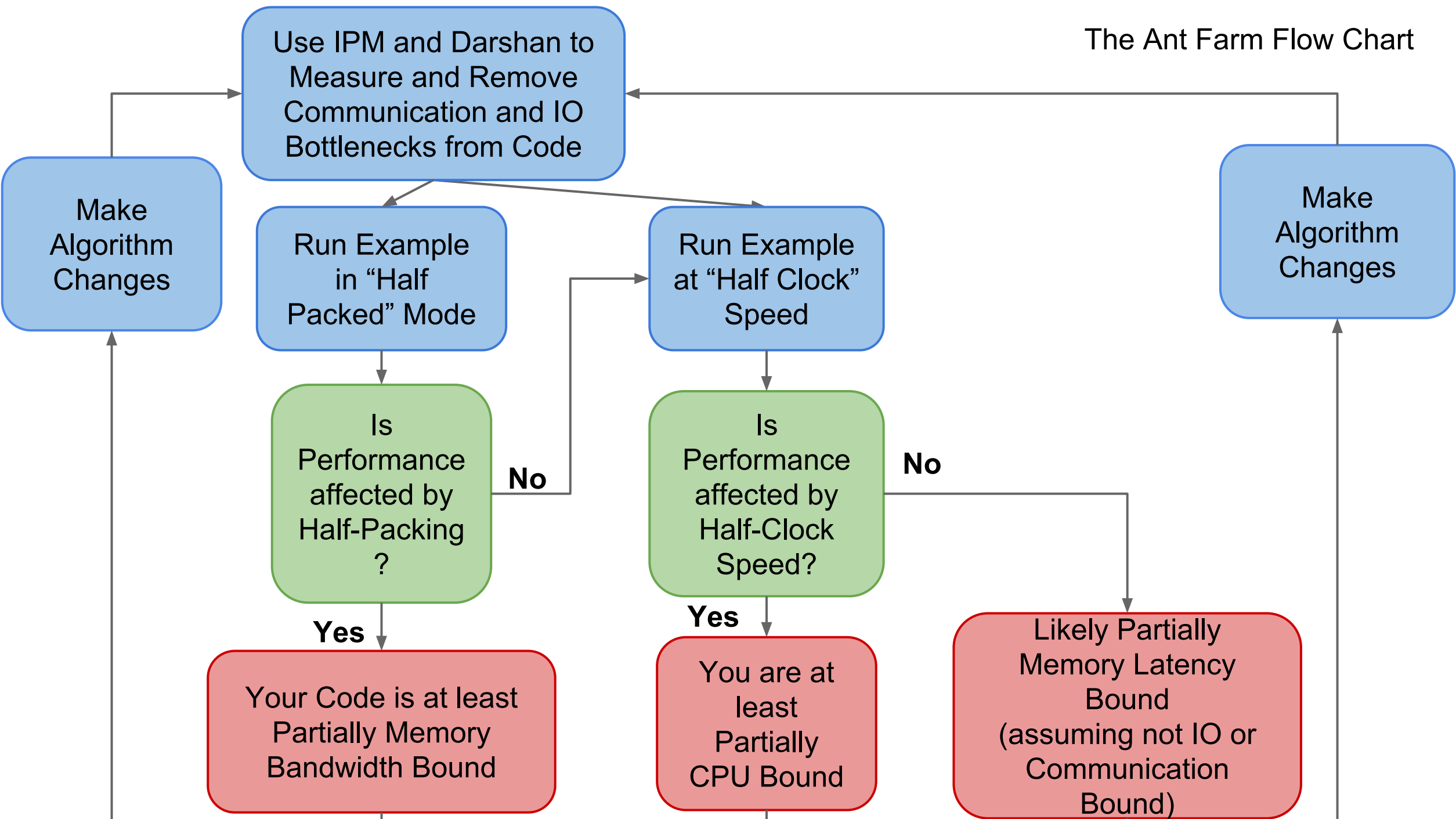
Use Edison to Test/Add OpenMP Improve Scalability. Help from NERSC/Cray COE Available.

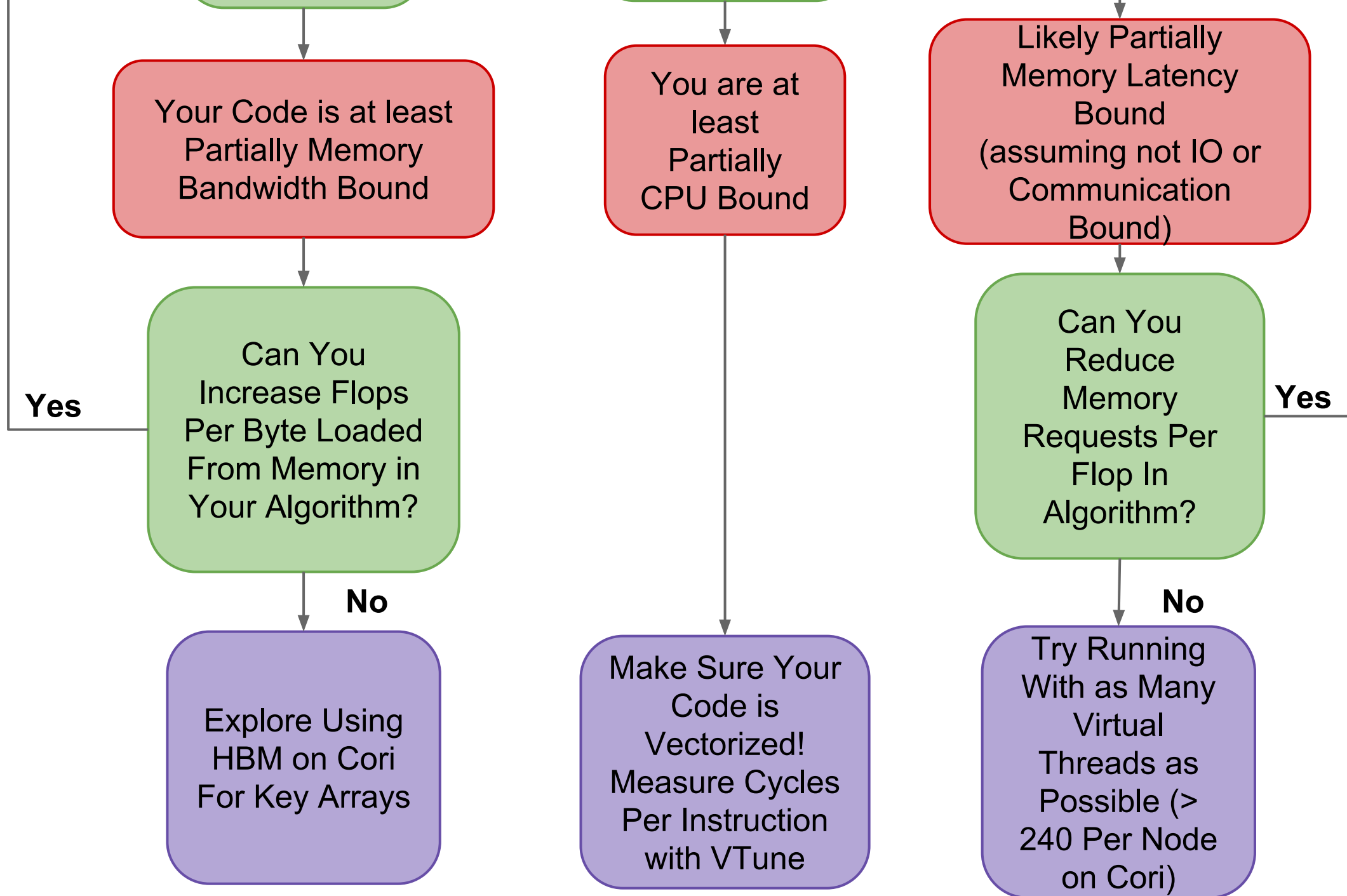
Create micro-kernels or examples to examine thread level performance, vectorization, cache use, locality.

The Dungeon: Simulate kernels on KNL. Plan use of on package memory, vector instructions.

Utilize performant / portable libraries

The Ant Farm Flow Chart





Your Code is at least Partially Memory Bandwidth Bound

Can You Increase Flops Per Byte Loaded From Memory in Your Algorithm?

Yes

No

Explore Using HBM on Cori For Key Arrays

You are at least Partially CPU Bound

Make Sure Your Code is Vectorized! Measure Cycles Per Instruction with VTune

Likely Partially Memory Latency Bound (assuming not IO or Communication Bound)

Can You Reduce Memory Requests Per Flop In Algorithm?

Yes

No

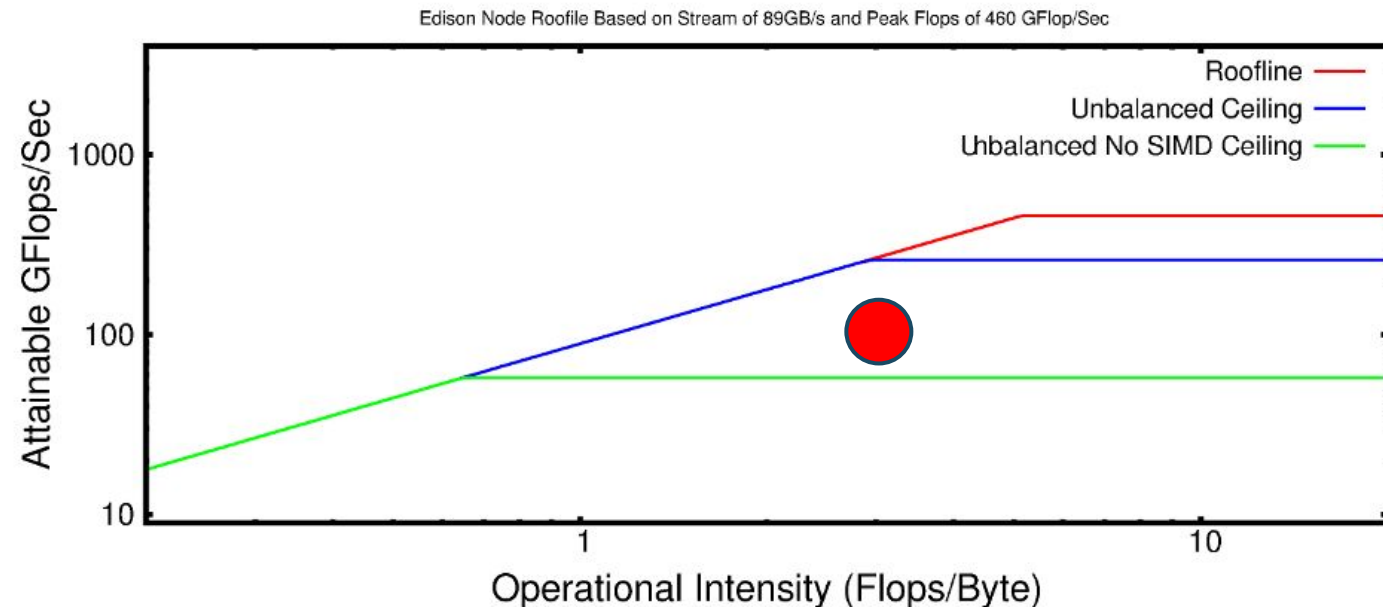
Try Running With as Many Virtual Threads as Possible (> 240 Per Node on Cori)

Are You Memory Bound? Compute Bound? Neither?



1. Determine your roofline position:

<http://www.nersc.gov/users/application-performance/measuring-arithmetic-intensity/>



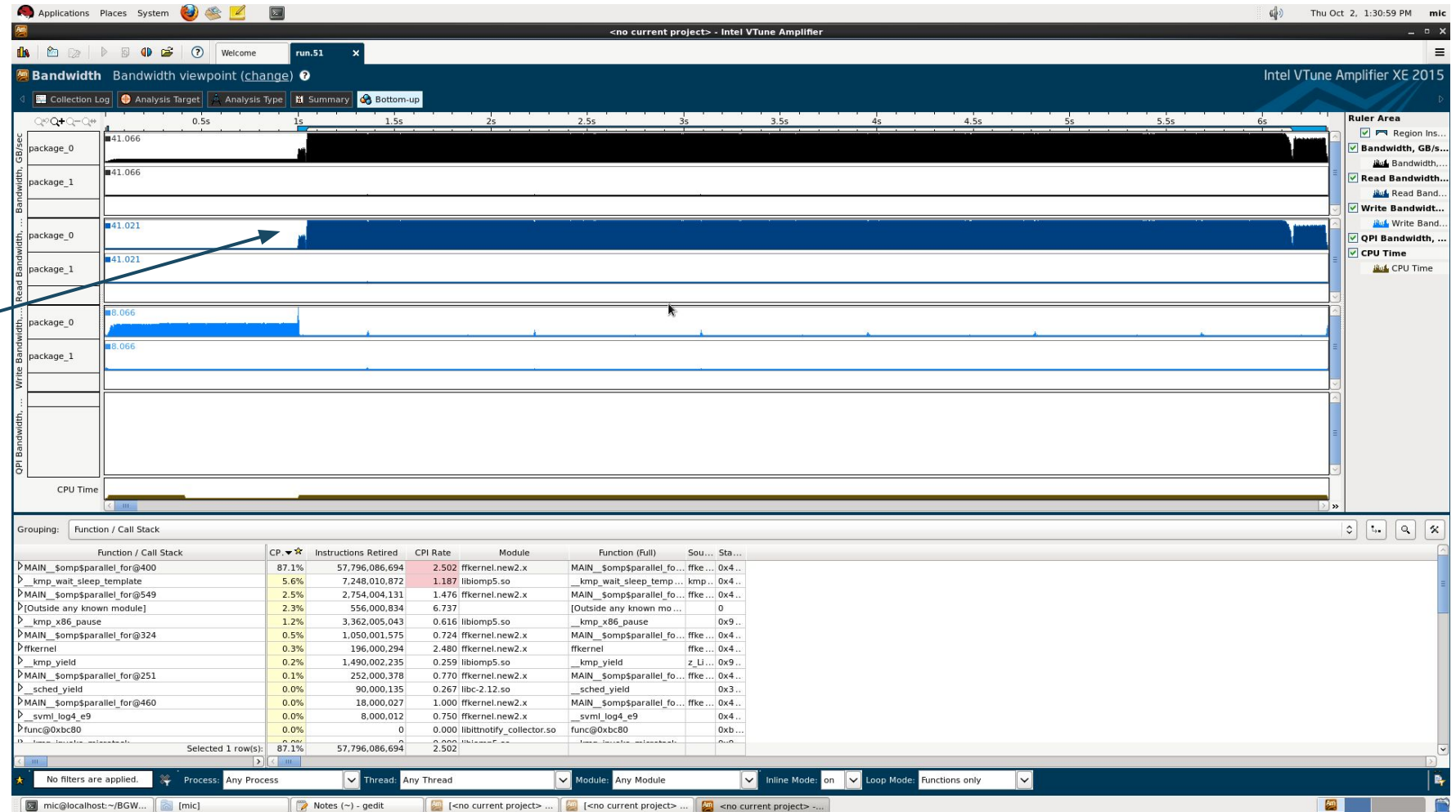
Measuring Your Memory Bandwidth Usage (VTune)

Measure memory bandwidth usage in VTune. (Next Talk)

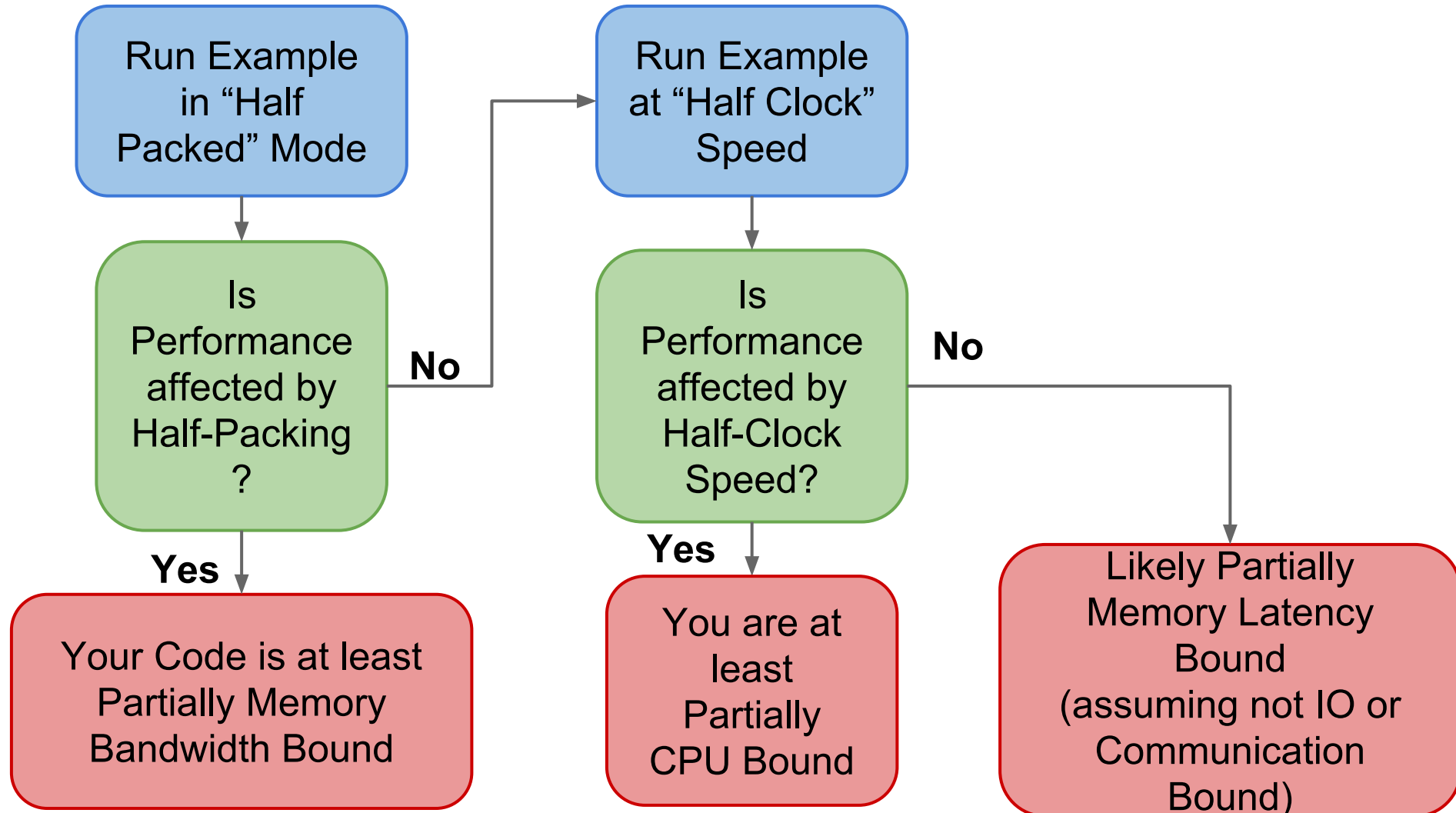
Compare to Stream GB/s.

If 90% of stream, you are memory bandwidth bound.

If less, more tests need to be done.



Are you memory or compute bound? Or both?



Are you memory or compute bound? Or both?

Run Example
in “Half
Packed” Mode

If you run on only half of the cores on a node, each core you do run has access to more bandwidth

```
aprun -n 24 -N 12 -S 6 ...
```

VS

```
aprun -n 24 -N 24 -S 12 ...
```

If your performance changes, you are at least partially memory bandwidth bound

Are you memory or compute bound? Or both?

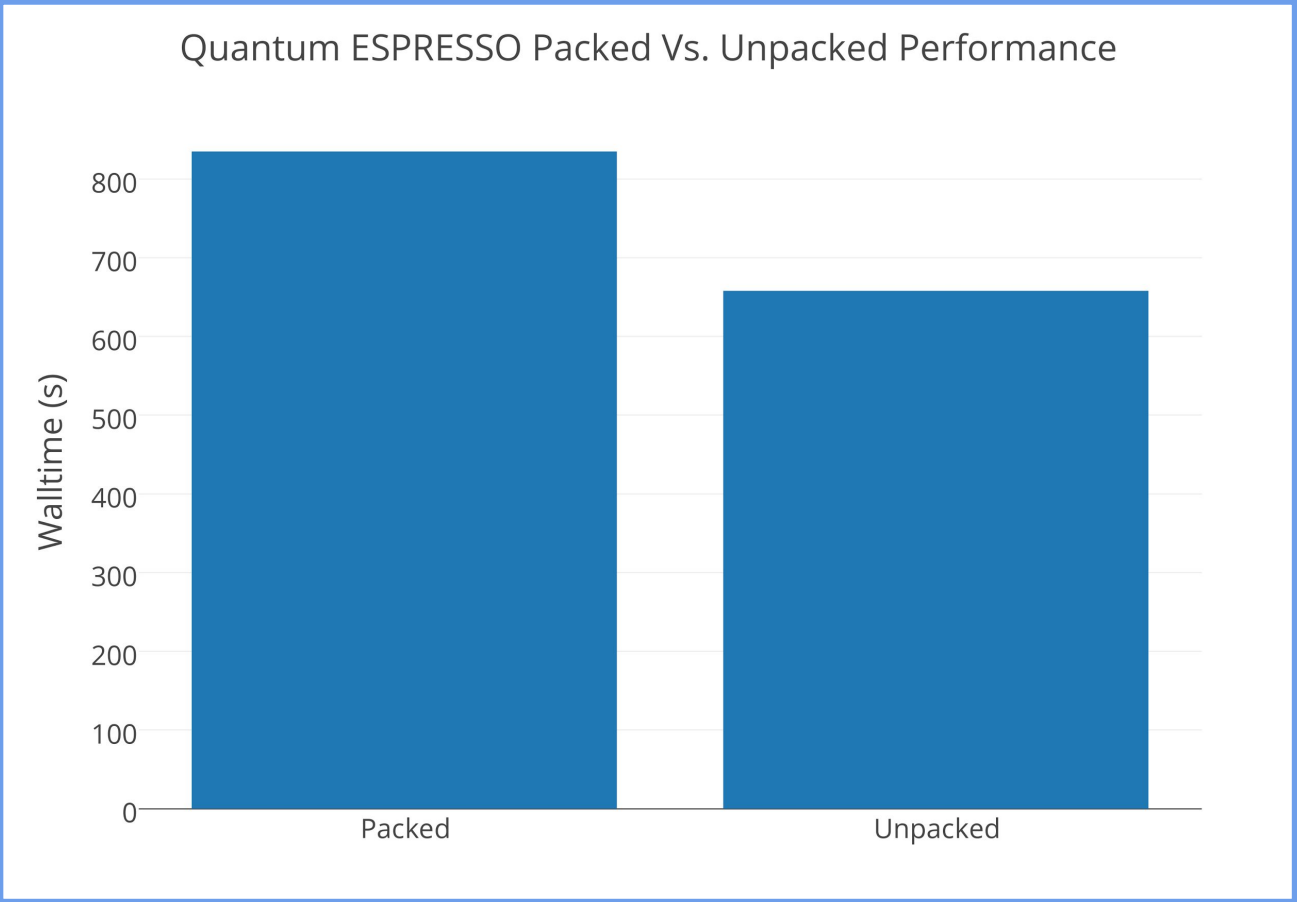
Run Example
in “Half
Packed” Mode

If you run on only half of the cores on a node, each core you do run has access to more bandwidth

```
aprun -n 24 -N
```

2 ...

If your performance



ound

Are you memory or compute bound? Or both?

Run Example
at “Half Clock”
Speed

Reducing the CPU speed slows down computation, but doesn't reduce memory bandwidth available.

```
aprun --p-state=2400000 ...
```

VS

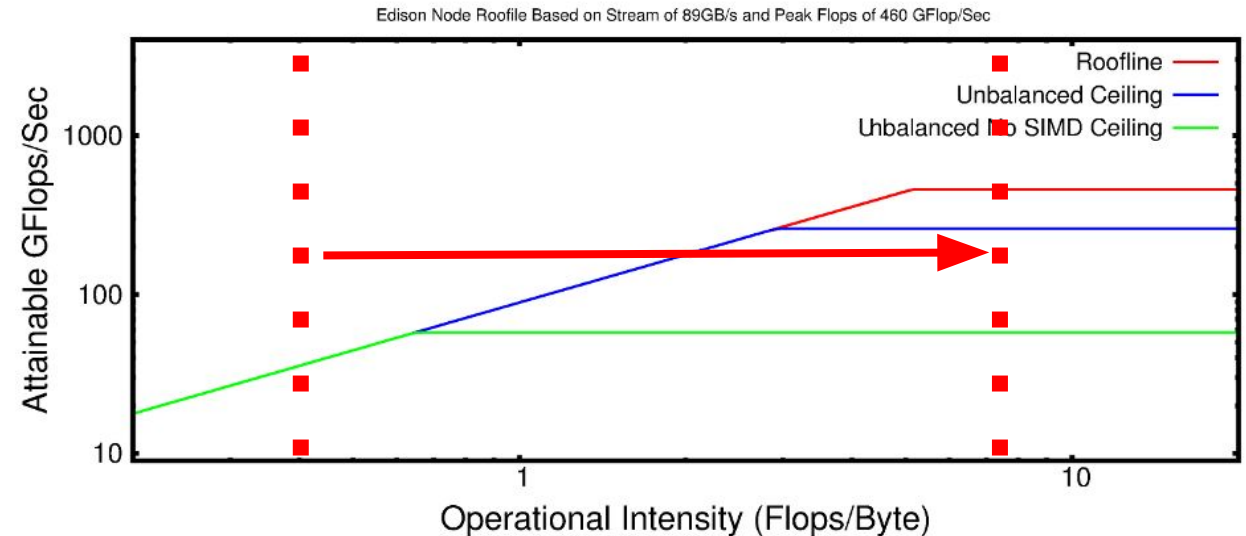
```
aprun --p-state=1900000 ...
```

If your performance changes, you are at least partially compute bound

So, you are Memory Bandwidth Bound?

What to do?

1. Try to improve memory locality, cache reuse



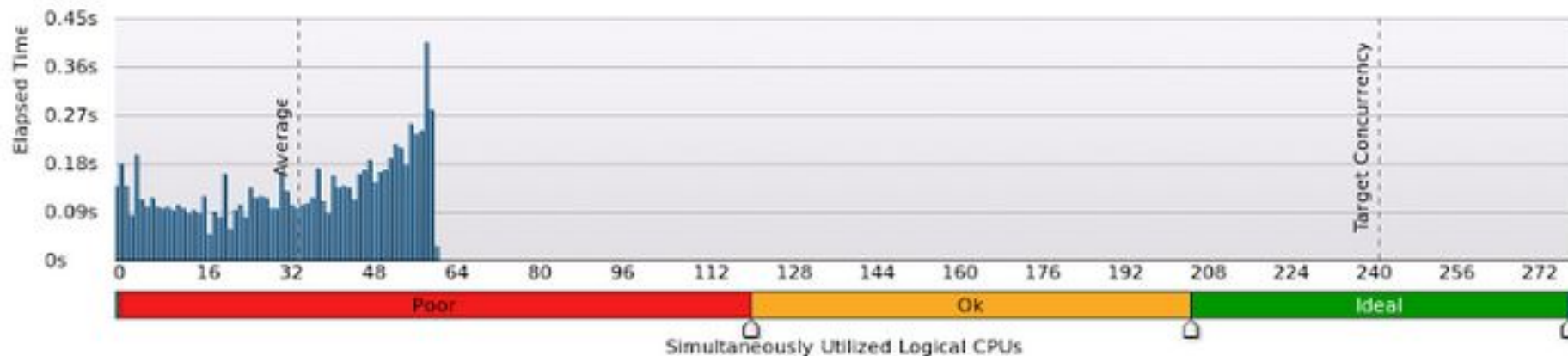
2. Identify the key arrays leading to high memory bandwidth usage and make sure they are/will-be allocated in HBM on Cori.

Profit by getting ~ 5x more bandwidth GB/s.

So, you are Compute Bound?

What to do?

1. Make sure you have good OpenMP scalability. Look at VTune to see thread activity for major OpenMP regions.



2. Make sure your code is vectorizing. Look at Cycles per Instruction (CPI) and VPU utilization in vtune.

See whether intel compiler vectorized loop using compiler flag: `-qopt-report=5`

High latency instructions : Complex-Division (without -fp model fast=2)



Intel VTune Amplifier XE 2015

Hotspots viewpoint (change) ?

Collection Log Analysis Target Analysis Type Summary Bottom-up Caller/Callee Top-down Tree Tasks and Frames gppkernel...

Source Assembly Assembly grouping: Address

S. Li.	Source	Address	Sou.. Line	Assembly	Effective Time by Utilization
					Idle Poor Ok Ideal Over
466	scht = scht + scha(ig)	0x408745	481	vunpckhpd %ymm3, %ymm3, %ymm3	0.001s
467	endif	0x408749	480	vmovapd %xmm5, %xmm15	
468		0x40874d	480	vmovsdq %xmm15, -0x28(%rbp)	0.202s
469	else	0x408752	480	fldq -0x28(%rbp), %st0	0.456s
470	! !dir\$ no unroll	0x408755	480	vunpckhpd %xmm5, %xmm5, %xmm11	0.001s
471	do ig = igbeg, min(igend,igmax)	0x408759	480	fld %st0, %st0	
472	! do ig = 1, igmax	0x40875b	480	vmovsdq %xmm11, -0x28(%rbp)	0.184s
473		0x408760	480	fmul %st1, %st0	0.444s
474	wdiff = wxt - wtilde_array(ig,my_igp)	0x408762	480	vextractf128 \$0x1, %ymm5, %xmm9	0.006s
475		0x408768	480	fldq -0x28(%rbp), %st0	
476	cden = wdiff	0x40876b	480	fld %st0, %st0	0.183s
477	!rden = cden * CONJG(cden)	0x40876d	480	fmul %st1, %st0	0.418s
478	!rden = 1D0 / rden	0x40876f	480	vmovsdq %xmm12, -0x28(%rbp)	0.006s
479	!delw = wtilde_array(ig,my_igp) * CONJG(cden) * rden	0x408774	480	faddp %st0, %st2	0.001s
480	cden = 1 / cden	0x408776	480	fxch %st1, %st0	0.196s
481	delw = wtilde_array(ig,my_igp) * cden	0x408778	480	fdivr %st3, %st0	0.462s
482	delwr = delw*CONJG(delw)	0x40877a	480	fldq -0x28(%rbp), %st0	0.113s
483	wdiffr = wdiff*CONJG(wdiff)	0x40877d	480	vmovsdq %xmm7, -0x28(%rbp)	0.192s
484		0x408782	480	fld %st0, %st0	0.418s
485	! JRD: Complex division is hard to vectorize. So, we help the compiler.	0x408784	480	fmul %st4, %st0	0.001s
486	scha(ig) = mygpvar1 * aqsntemp(ig,n1) * delw * I_eps_array(ig,n1)	0x408786	480	fxch %st1, %st0	0.025s
487	! scha_temp = mygpvar1 * aqsntemp(ig,n1) * delw * I_eps_array(ig,n1)	0x408788	480	fmul %st3, %st0	0.602s
488		0x40878a	480	fldq -0x28(%rbp), %st0	0.002s
489	! JRD: This if is OK for vectorization	0x40878d	480	fld %st0, %st0	0.026s
490	if (wdiffr.gt.limittwo .and. delwr.lt.limitone) then	0x40878f	480	fmulp %st0, %st5	0.185s
491	scht = scht + scha(ig)	0x408791	480	vunpckhpd %xmm9, %xmm9, %xmm4	0.404s
492	endif	0x408796	480	fxch %st4, %st0	0s

Selected 1 row(s): Highlighted 217 row(s):



Are you latency bound?

You may be memory latency bound (or you may be spending all your time in IO and Communication).

If running with hyper-threading on Edison improves performance, you **might** be latency bound:

```
aprun -j 2 -n 48 ....
```

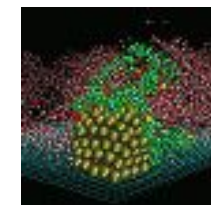
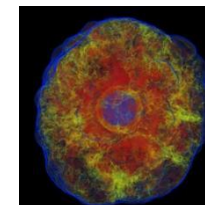
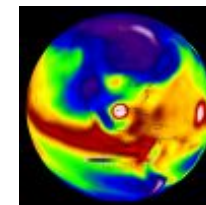
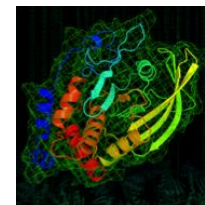
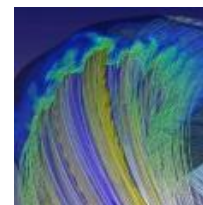
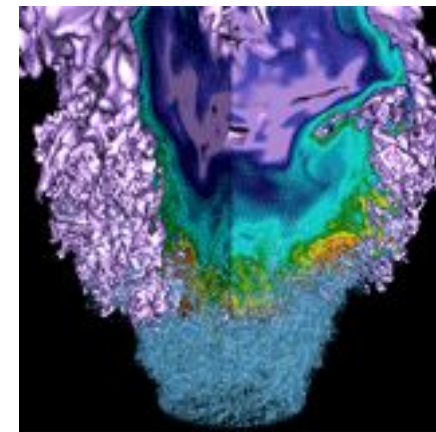
VS

```
aprun -n 24 ....
```

If you can, try to reduce the number of memory requests per flop by accessing contiguous and predictable segments of memory and reusing variables in cache as much as possible.

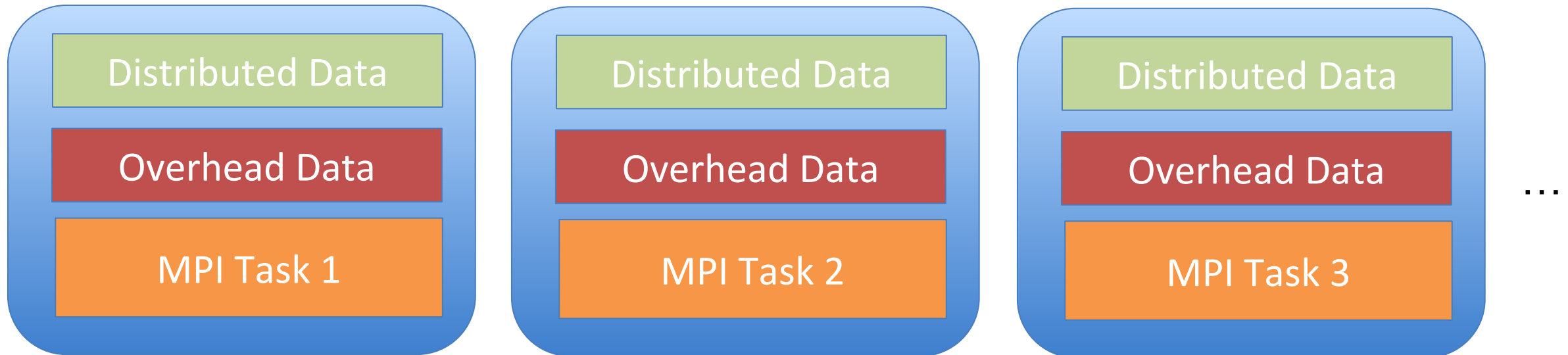
On Cori, each core will support up to 4 threads. Use them all.

NESAP Case Study



BerkeleyGW Use Case

- ★ Big systems require more memory. Cost scales as N_{atoms}^2 to store the data.
- ★ In an MPI GW implementation, in practice, to avoid communication, data is duplicated and **each MPI task has a memory overhead.**
- ★ Users sometimes forced to use 1 of 24 available cores, in order to provide MPI tasks with enough memory. **90% of the computing capability is lost.**



In house code (I'm one of main developers). Use as “prototype” for App Readiness.

Significant Bottleneck is large matrix reduction like operations. Turning arrays into numbers.

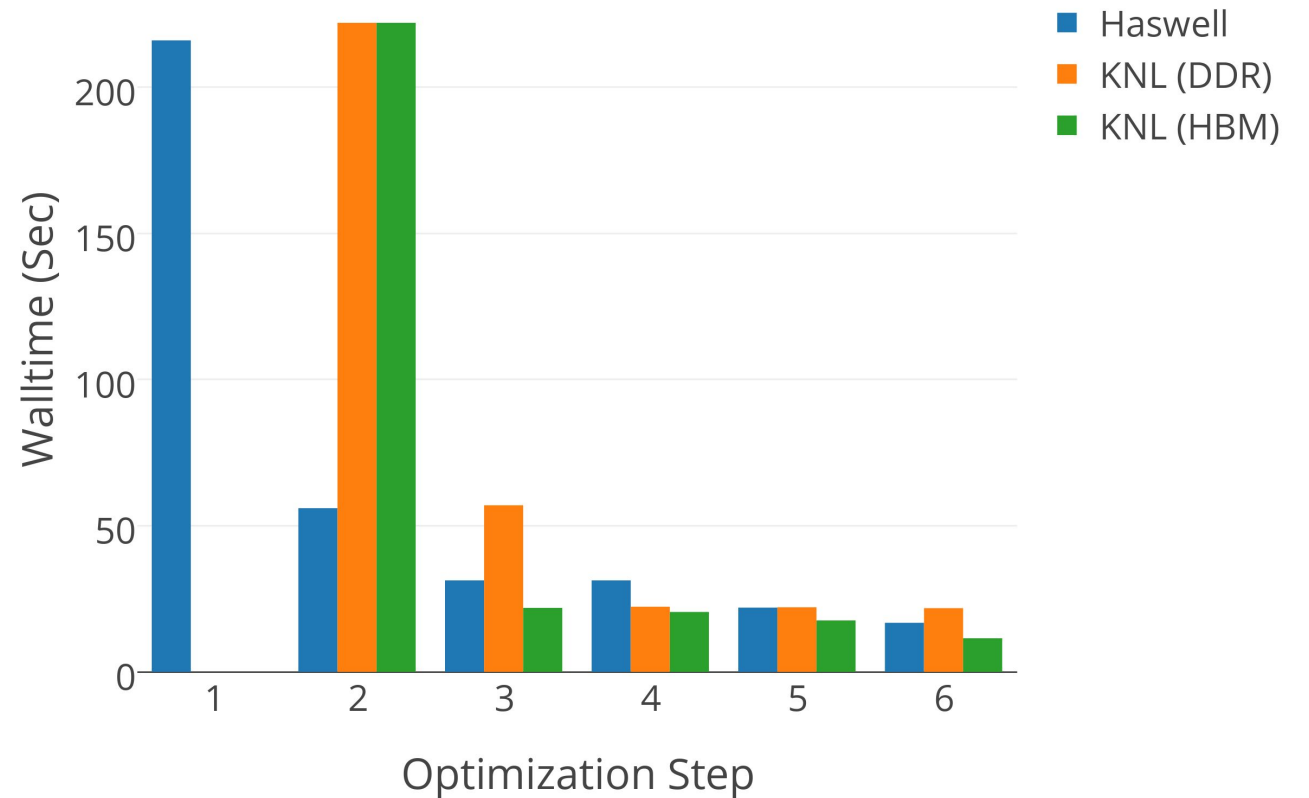
$$\langle n\mathbf{k} | \Sigma_{\text{CH}}(E) | n'\mathbf{k} \rangle = \frac{1}{2} \sum_{n''} \sum_{\mathbf{q}\mathbf{G}\mathbf{G}'} M_{n''n}^*(\mathbf{k}, -\mathbf{q}, -\mathbf{G}) M_{n''n'}(\mathbf{k}, -\mathbf{q}, -\mathbf{G}') \\ \times \frac{\Omega_{\mathbf{G}\mathbf{G}'}^2(\mathbf{q}) (1 - i \tan \phi_{\mathbf{G}\mathbf{G}'}(\mathbf{q}))}{\tilde{\omega}_{\mathbf{G}\mathbf{G}'}(\mathbf{q}) (E - E_{n''\mathbf{k}-\mathbf{q}} - \tilde{\omega}_{\mathbf{G}\mathbf{G}'}(\mathbf{q}))} v(\mathbf{q} + \mathbf{G}')$$

Optimization Path

Optimization process for Kernel-C (Sigma code):

1. Refactor (3 Loops for MPI, OpenMP, Vectors)
2. Add OpenMP
3. Initial Vectorization (loop reordering, conditional removal)
4. Cache-Blocking
5. Improved Vectorization
6. Hyper-threading

Optimization Process



Vectorization

```
!$OMP DO reduction(+:achtemp)
do my_igp = 1, ngpown
  ...
  do iw=1,nfreq ! nfreq is 3

    scht=0D0
    wxt = wx_array(iw)

    do ig = 1, ncouls

      !if (abs(wtilde_array(ig,my_igp) * eps(ig,my_igp)) .lt. TOL) cycle

      wdifff = wxt - wtilde_array(ig,my_igp)
      delw = wtilde_array(ig,my_igp) / wdifff
      ...
      scha(ig) = mygpvar1 * aqsntemp(ig) * delw * eps(ig,my_igp)
      scht = scht + scha(ig)

    enddo ! loop over g
    sch_array(iw) = sch_array(iw) + 0.5D0*scht

  enddo

  achtemp(:) = achtemp(:) + sch_array(:) * vcoul(my_igp)

enddo
```

ngpown typically in 100's to 1000s. Good for many threads.

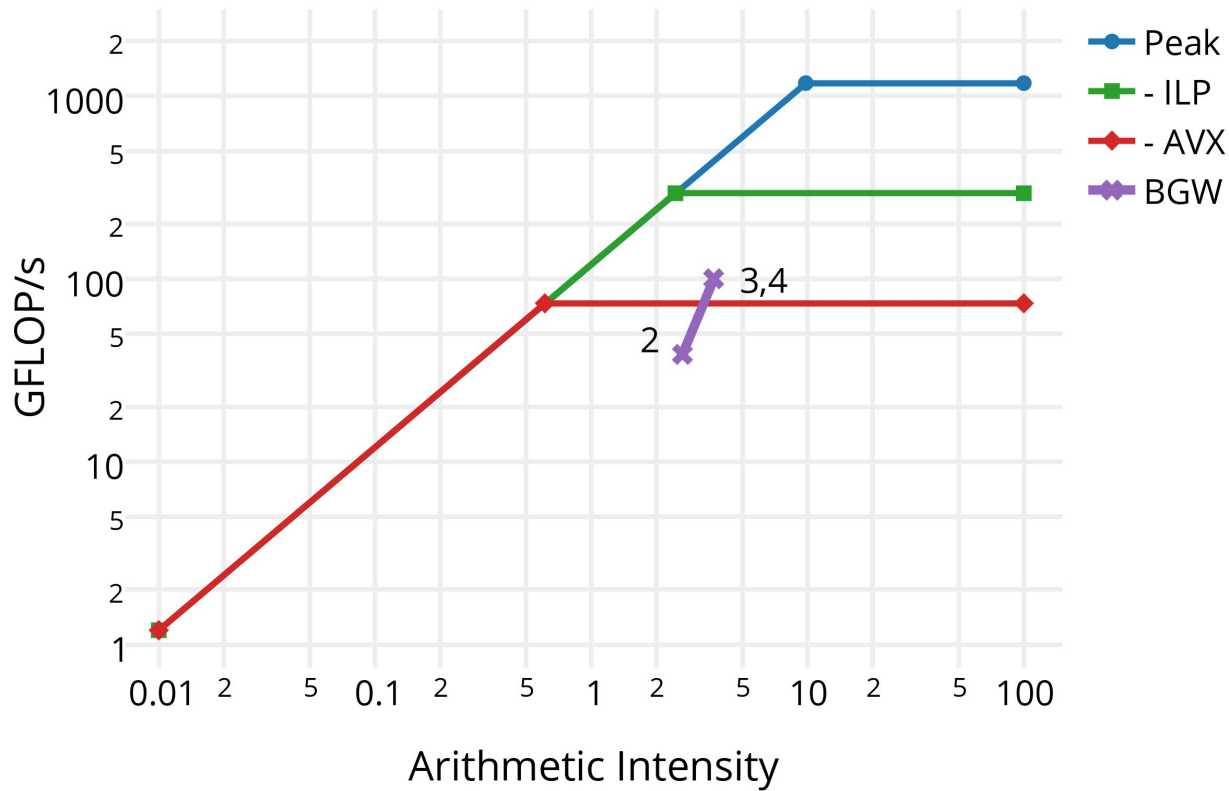
Original inner loop. Too small to vectorize!

ncouls typically in 1000s - 10,000s. Good for vectorization.

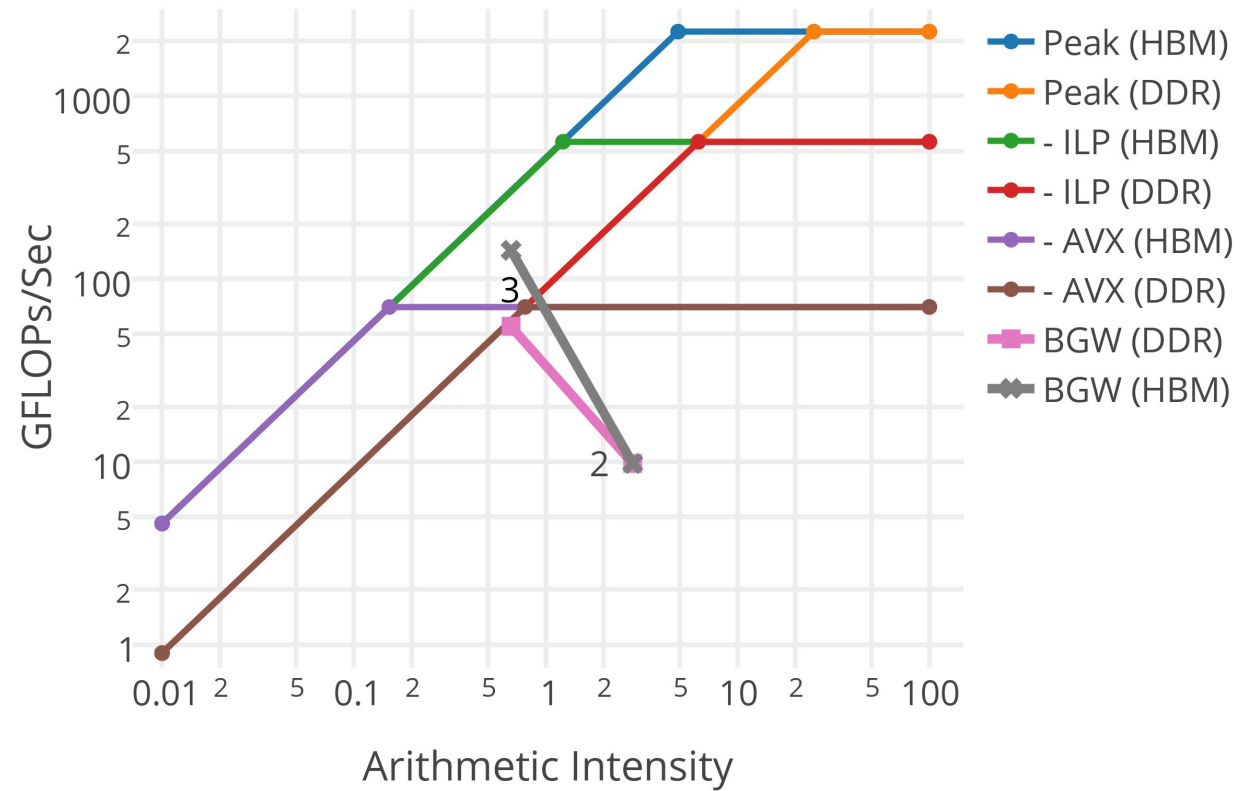
Attempt to save work breaks vectorization and makes code slower.

Change in Roofline

Haswell Roofline Optimization Path



KNL Roofline Optimization Path



The loss of L3 on MIC makes locality more important.

Why KNC worse than Haswell for GPP Kernel?

- 2S Haswell 27.9s KNC 39.9s (Bandwidth bound on KNC but not on Haswell)

```
!$OMP DO
do my_igp = 1, ngpown
  do iw = 1 , 3
    do ig = 1, igmax
      load wtilde_array(ig,my_igp) 819 MB, 512KB per row
      load aqsntemp(ig,n1) 256 MB, 512KB per row
      load l_eps_array(ig,my_igp) 819 MB, 512KB per row
      do work (including divide)
```

Required Cache size to reuse 3 times:

1536 KB

L2 on KNL is 512 KB per core

L2 on Has. is 256 KB per core

L3 on Has. is 3800 KB per core

Without blocking we spill out of L2 on KNC and Haswell. But, Haswell has L3 to catch us.

Why KNC worse than Haswell for GPP Kernel?

- 2S Haswell 27.9s KNC 39.9s (Bandwidth bound on KNC but not on Haswell)

```
!$OMP DO
```

```
do my_igp = 1, ngpown
```

```
  do igbeg = 1, igmax, igblk
```

```
    do iw = 1 , 3
```

```
      do ig = igbeg, min(igbeg + igblk,igmax)
```

```
        load wtilde_array(ig,my_igp) 819 MB, 512KB per row
```

```
        load aqsntemp(ig,n1) 256 MB, 512KB per row
```

```
        load l_eps_array(ig,my_igp) 819 MB, 512KB per row
```

```
        do work (including divide)
```

Required Cache size to reuse 3 times:

1536 KB

L2 on KNL is 512 KB per core

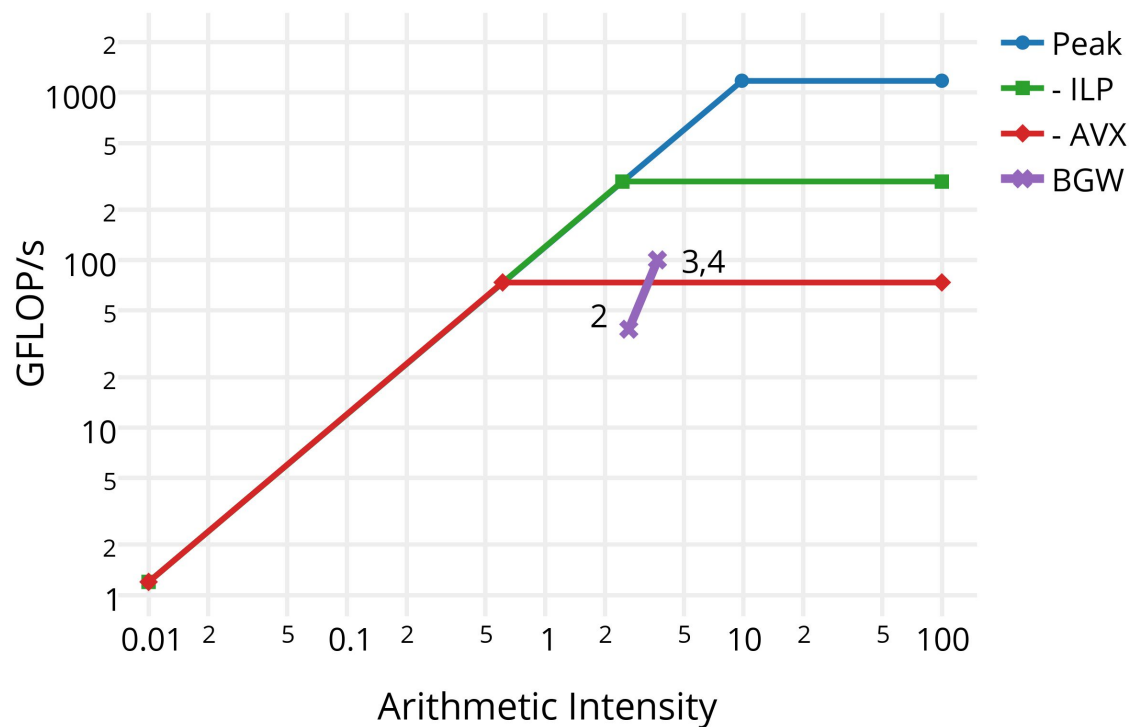
L2 on Has. is 256 KB per core

L3 on Has. is 3800 KB per core

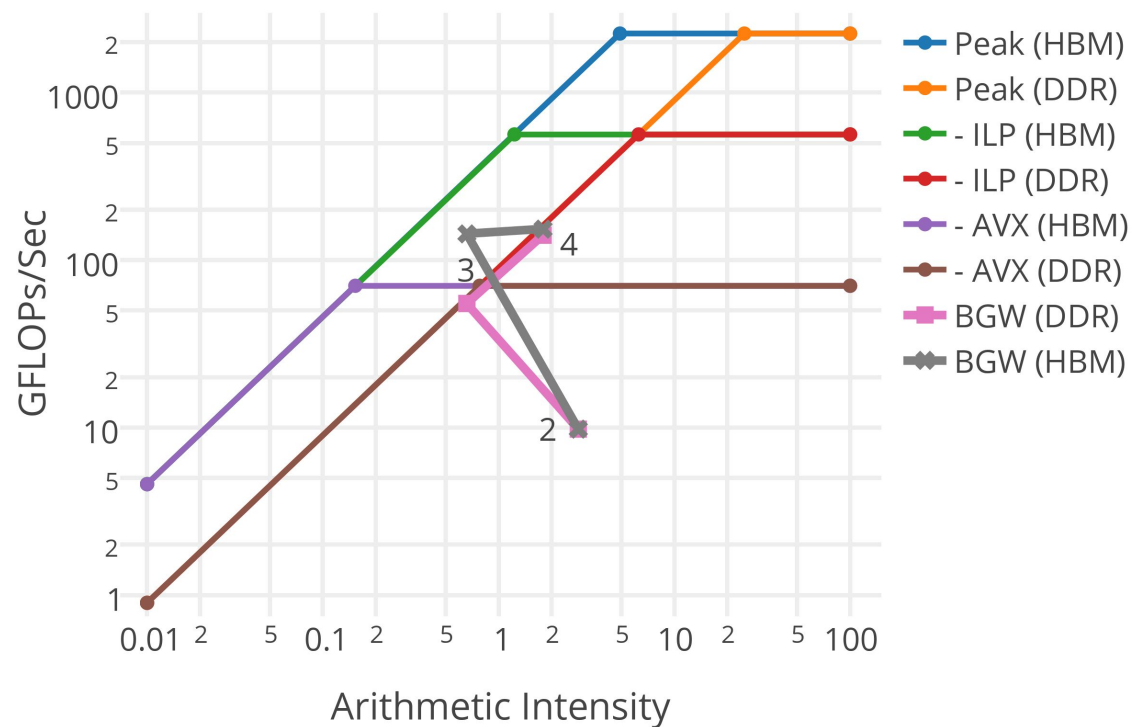
Without blocking we spill out of L2 on KNC and Haswell. But, Haswell has L3 to catch us.

Cache Blocking Optimization

Haswell Roofline Optimization Path



KNL Roofline Optimization Path



Why Complex Divides so Slow?



Found significant x87 instructions from 1/complex_number instead of AVX/AVX-512

Source	Address	Sou. Line	Assembly	Effective Time by Utilization
scht = scht + scha(ig)	0x408745	481	vunpckhpd %ymm3, %ymm3, %ymm3	0.001s
endif	0x408749	480	vmovapd %xmm5, %xmm15	
	0x40874d	480	vmovsdq %xmm15, -0x28(%rbp)	0.202s
else	0x408752	480	fldq -0x28(%rbp), %st0	0.456s
! dir\$ no unroll	0x408755	480	vunpckhpd %xmm5, %xmm5, %xmm11	0.001s
do ig = igbeg, min(igend, igmax)	0x408759	480	fld %st0, %st0	
! do ig = 1, igmax	0x40875b	480	vmovsdq %xmm11, -0x28(%rbp)	0.184s
wdiff = wxt - wtilde_array(ig, my_igp)	0x408760	480	fmul %st1, %st0	0.444s
	0x408762	480	vextractf128 \$0x1, %ymm5, %xmm9	0.006s
	0x408768	480	fldq -0x28(%rbp), %st0	
cden = wdiff	0x40876b	480	fld %st0, %st0	0.183s
!rden = cden * CONJG(cden)	0x40876d	480	fmul %st1, %st0	0.418s
!rden = 1D0 / rden	0x40876f	480	vmovsdq %xmm12, -0x28(%rbp)	0.006s
!delw = wtilde_array(ig, my_igp) * CONJG(cden) * rden	0x408774	480	faddp %st0, %st2	0.001s
cden = 1 / cden	0x408776	480	fxch %st1, %st0	0.196s
delw = wtilde_array(ig, my_igp) * cden	0x408778	480	fdivr %st3, %st0	0.462s
delwr = delw * CONJG(delw)	0x40877a	480	fldq -0x28(%rbp), %st0	0.113s
wdiff = wdiff * CONJG(wdiff)	0x40877d	480	vmovsdq %xmm7, -0x28(%rbp)	0.192s
	0x408782	480	fld %st0, %st0	0.418s
! JRD: Complex division is hard to vectorize. So, we help the compiler.	0x408784	480	fmul %st4, %st0	0.001s
scha(ig) = mygpvar1 * aqsntemp(ig, n1) * delw * I_eps_array(ig, n1)	0x408786	480	fxch %st1, %st0	0.025s
! scha_temp = mygpvar1 * aqsntemp(ig, n1) * delw * I_eps_array(ig, n1)	0x408788	480	fmul %st3, %st0	0.602s
	0x40878a	480	fldq -0x28(%rbp), %st0	0.002s
! JRD: This if is OK for vectorization	0x40878d	480	fld %st0, %st0	0.026s
if (wdiff.gt.limittwo .and. delwr.lt.limitone) then	0x40878f	480	fmlp %st0, %st5	0.185s
scht = scht + scha(ig)	0x408791	480	vunpckhpd %xmm9, %xmm9, %xmm4	0.404s
endif	0x408796	480	fxch %st4, %st0	0s

Can significantly speed up by

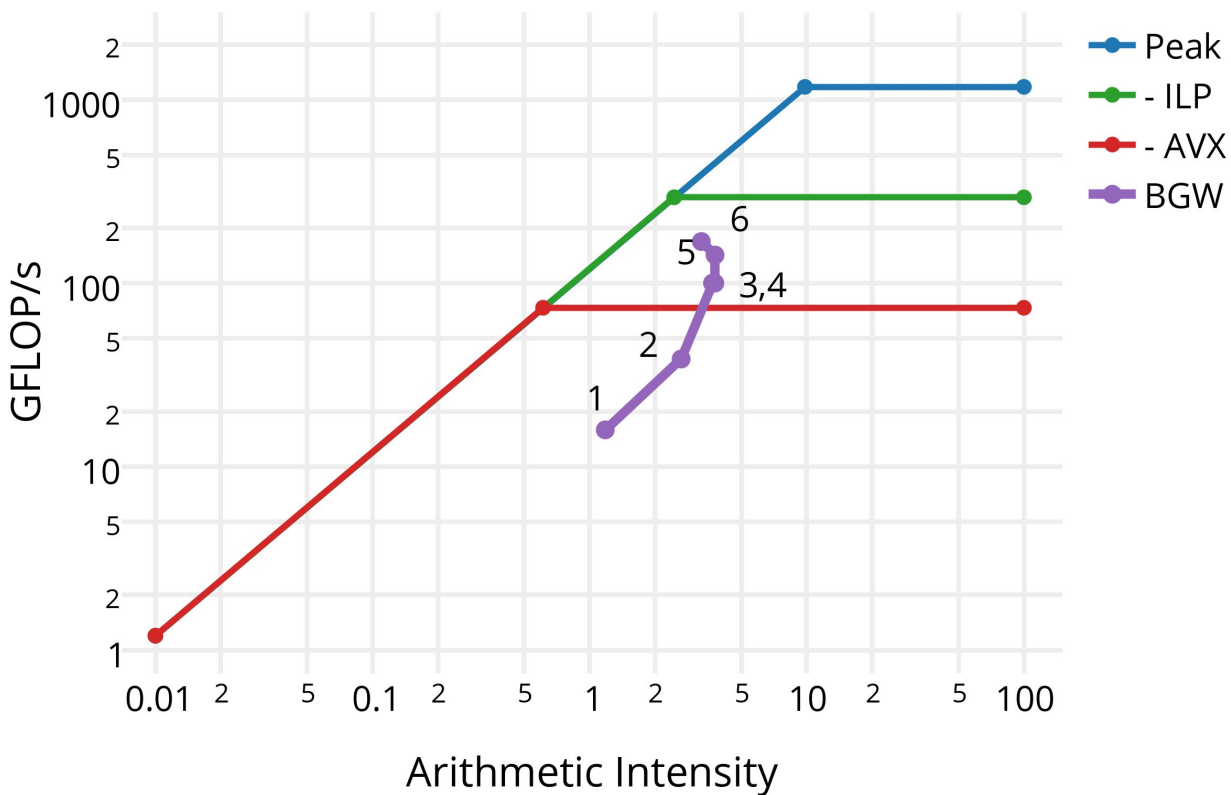
a) Doing complex divide manually

Or

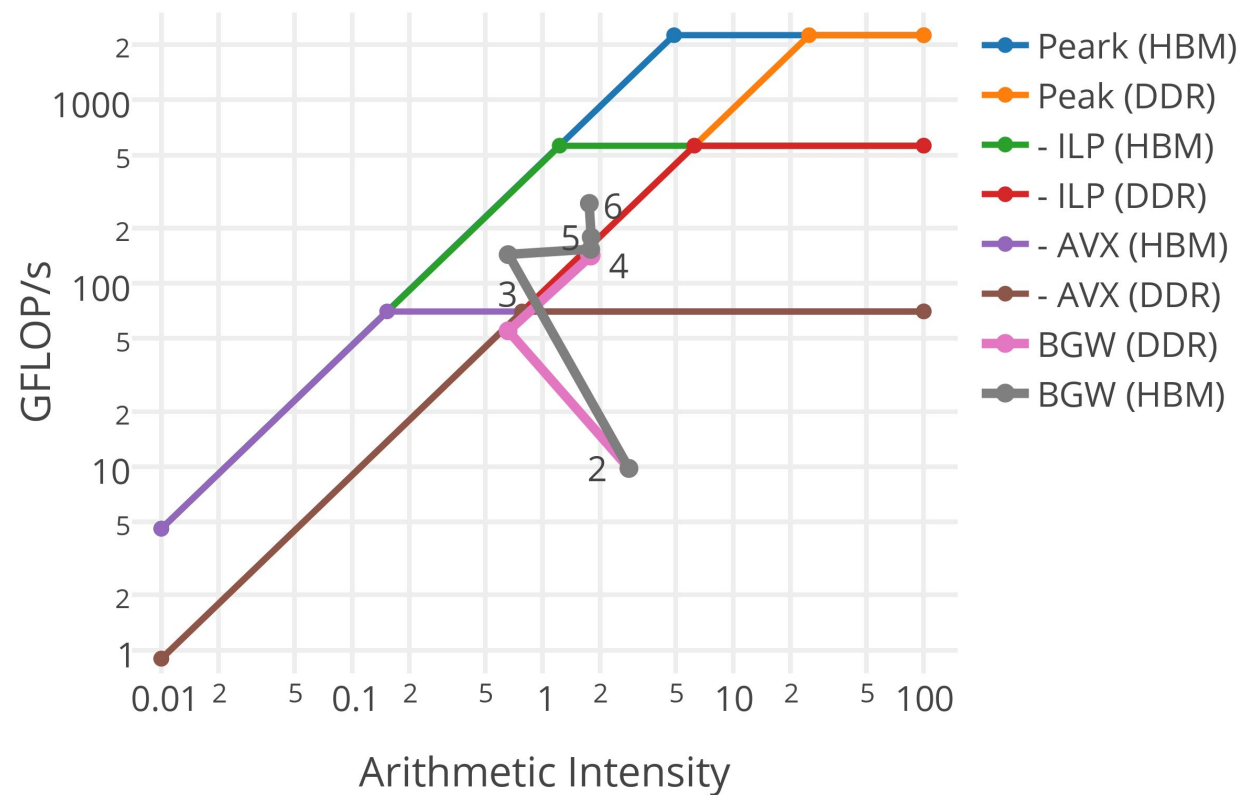
b) Using -fp-model fast=2

Additional Speedups from Hyperthreading

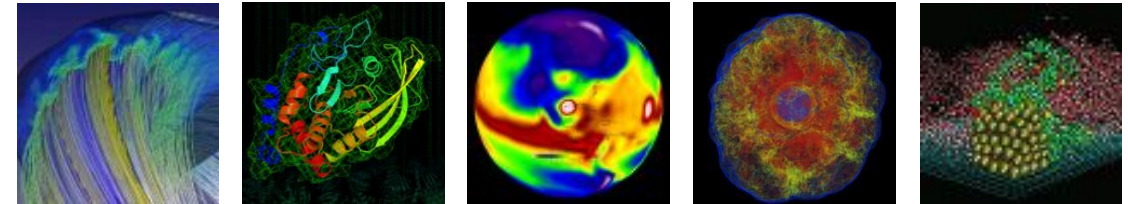
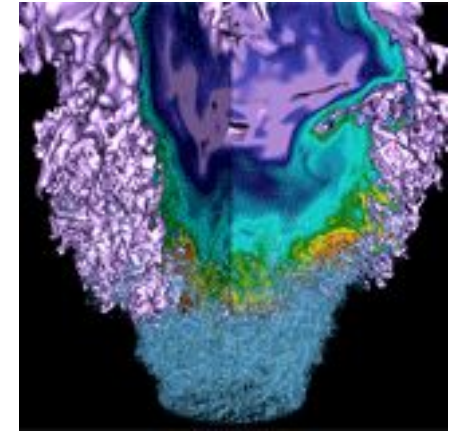
Haswell Roofline Optimization Path



KNL Roofline Optimization Path

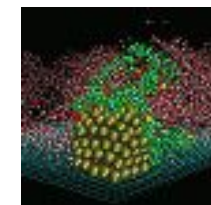
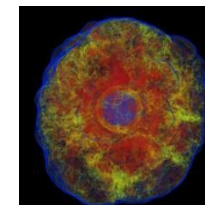
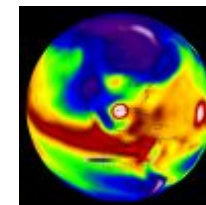
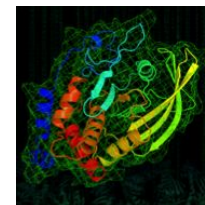
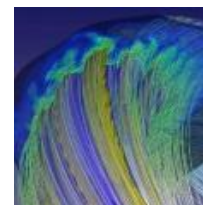
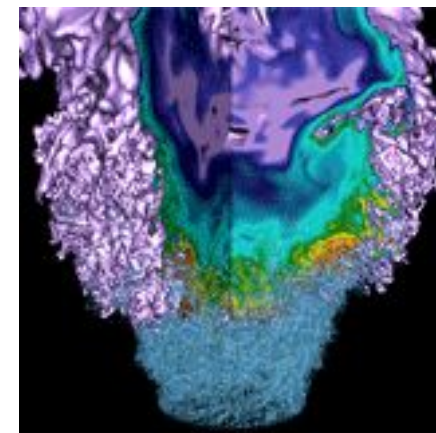


Conclusions



1. Optimizing code is not always straightforward. It is a continual discovery process that involves many sequential and coupled changes.
2. Use profiling tools to find and characterize hotspots.
3. Understanding bandwidth and compute limitations of hotspots are key to deciding how to improve code.

The End (Extra Slides)

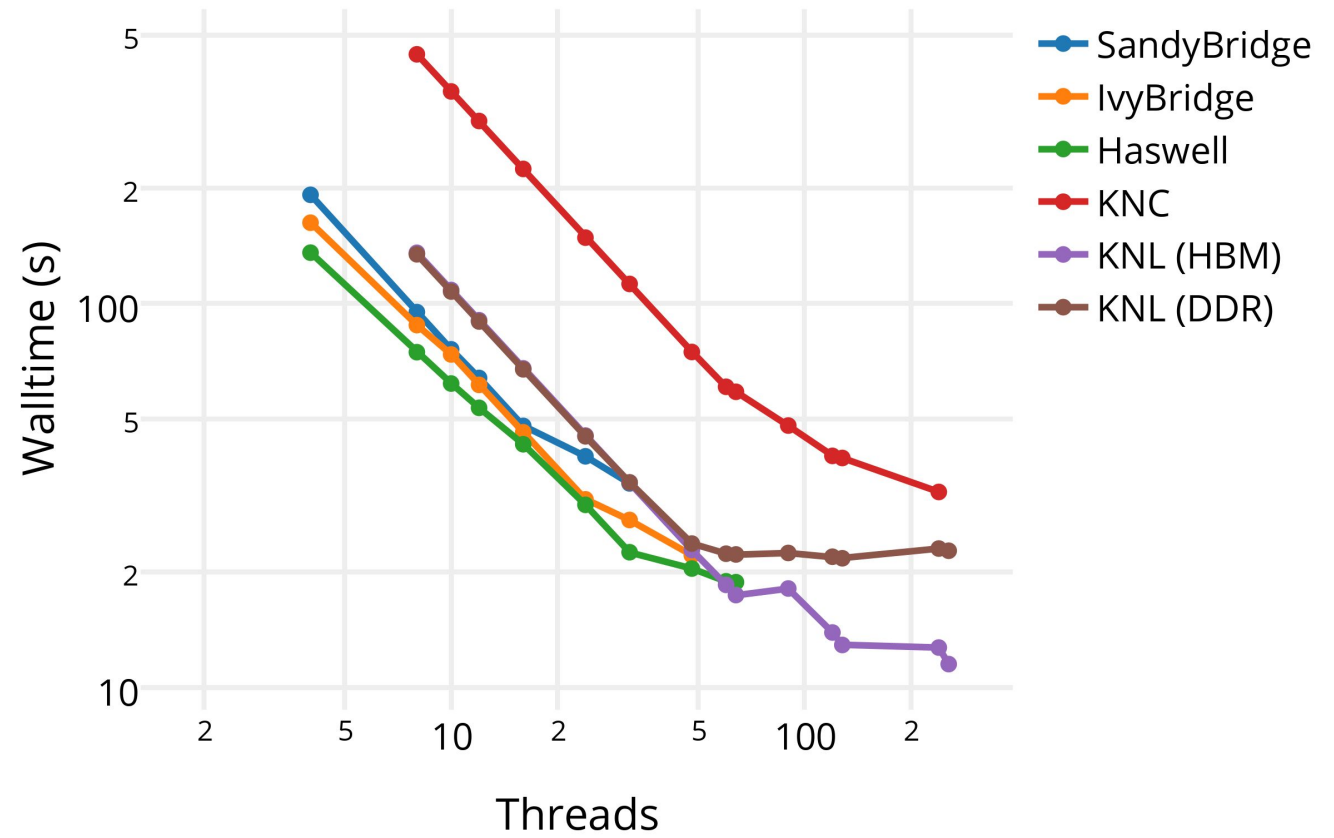


Thread Scaling

KNL DDR performance saturates at around 50 threads, becomes memory bandwidth limited.

KNL MCDRAM performance beats dual socket Haswell by 63%.

Kernel C Thread Scaling



Why Complex Divides so Slow?

Code performance now limited by complex divides

why??

For complex division in performance critical loop, I had already removed the explicit complex divide but what is faster?

a) $c = 1 / c$ vs. b) $r = c * \text{conjg}(c)$
 $r = 1 / r$
 $c = \text{conjg}(c) * r$

c/d) Compiling with/without -fp-model fast=2

Real-Division (with or without -fp model fast=2)



Intel VTune Amplifier XE 2015

Advanced Hotspots Hotspots viewpoint (change) ?

Collection Log Analysis Target Analysis Type Summary Bottom-up Caller/Callee Top-down Tree Tasks and Frames gppkernel... x

Source Assembly Assembly grouping: Address

S. Li.	Source	Address	Sou.. Line	Assembly	Effective Time by Utili
					Idle Poor Ok Ideal
469	else	0x4087b0	477	vmovddup %ymm3, %ymm2	0.455s
470	! !dir\$ no unroll	0x4087b4	477	vmovddup %ymm1, %ymm3	0.214s
471	do ig = igbeg, min(igend,igmax)	0x4087b8	477	vfmaddsub231pd %ymm2, %ymm7, %ymm14	0.349s
472	! do ig = 1, igmax	0x4087bd	477	vfmaddsub231pd %ymm3, %ymm5, %ymm8	0.045s
473		0x4087c2	477	vperm2f128 \$0x20, %ymm8, %ymm14, %ymm7	0.494s
474	wdiff = wxt - wtilde_array(ig,my_igp)	0x4087c8	477	vperm2f128 \$0x31, %ymm8, %ymm14, %ymm14	0.565s
475		0x4087ce	477	vunpcklpd %ymm14, %ymm7, %ymm1	0.423s
476	cden = wdiff	0x4087d3	478	vdivpd %ymm1, %ymm13, %ymm7	0.141s
477	rden = cden * CONJG(cden)	0x4087d7	480	vunpckhpd %xmm7, %xmm7, %xmm0	14.800s
478	rden = 1D0 / rden	0x4087db	480	vinsertrf128 \$0x1, %xmm0, %ymm7, %ymm8	0.727s
479	!rden = 1D0 + rden	0x4087e1	480	vmovddup %ymm8, %ymm14	1.869s
480	delw = wtilde_array(ig,my_igp) * CONJG(cden) * rden	0x4087e6	480	vextractf128 \$0x1, %ymm7, %xmm7	1.249s
481	delwr = delw*CONJG(delw)	0x4087ec	480	vmulpy (%r14,%rdi,1), %ymm14, %ymm8	0.005s
482	wdiffrr = wdiff*CONJG(wdiff)	0x4087f2	480	vunpckhpd %xmm7, %xmm7, %xmm0	3.126s
483		0x4087f6	480	vinsertrf128 \$0x1, %xmm0, %ymm7, %ymm7	0.015s
484	! JRD: Complex division is hard to vectorize. So, we help the compiler.	0x4087fc	480	vmovddup %ymm7, %ymm14	
485	scha(ig) = mygpvar1 * aqsntemp(ig,n1) * delw * I_eps_array(ig,n1)	0x408800	480	vshufpd \$0x5, %ymm8, %ymm8, %ymm7	0.032s
486	! scha_temp = mygpvar1 * aqsntemp(ig,n1) * delw * I_eps_array(ig,n1)	0x408806	480	vmulpd %ymm7, %ymm6, %ymm6	0.619s
487		0x40880a	485	vmovupdy (%r14,%r15,1), %ymm7	3.080s
488	! JRD: This if is OK for vectorization	0x408810	480	vmulpy 0x20(%r14,%rdi,1), %ymm14, %ymm0	0.019s
489	if (wdiffrr.gt.limittwo .and. delwr.lt.limitone) then	0x408817	480	vfmaddsub213pd %ymm6, %ymm8, %ymm2	0.399s
490	scht = scht + scha(ig)	0x40881c	485	vunpckhpd %ymm7, %ymm7, %ymm14	3.034s
491	endif	0x408820	485	vmovupdy -0x50(%rbp), %ymm7	0.017s
492		0x408825	480	vshufpd \$0x5, %ymm0, %ymm0, %ymm8	0.025s
493	! scha_mult = merge(1.0,0.0,wdiffrr.gt.limittwo .and. delwr.lt.limitone)	0x40882a	480	vmulpd %ymm8, %ymm15, %ymm15	0.014s
494	! scht = scht + scha(ig)*scha_mult	0x40882f	485	vmovupdy 0x20(%r14,%r15,1), %ymm8	0.637s
495		0x408836	485	vmulpd %ymm7, %ymm14, %ymm6	0.333s

Selected 1 row(s):

Highlighted 3 row(s):



Complex-Division (with -fp model fast=2)

Intel VTune Amplifier XE 2015

Hotspots viewpoint (change) ?

Collection Log Analysis Target Analysis Type Summary Bottom-up Caller/Callee Top-down Tree Tasks and Frames gppkernel...

Source Assembly Address

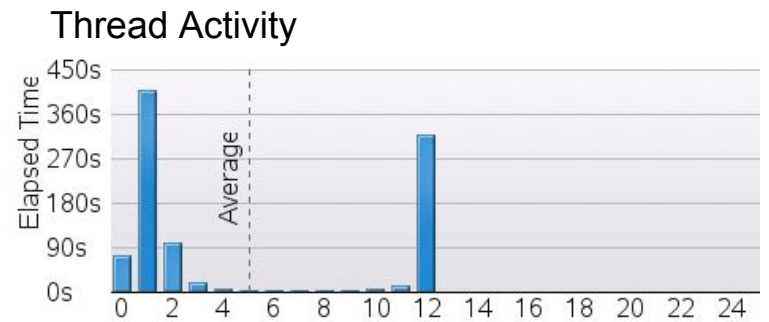
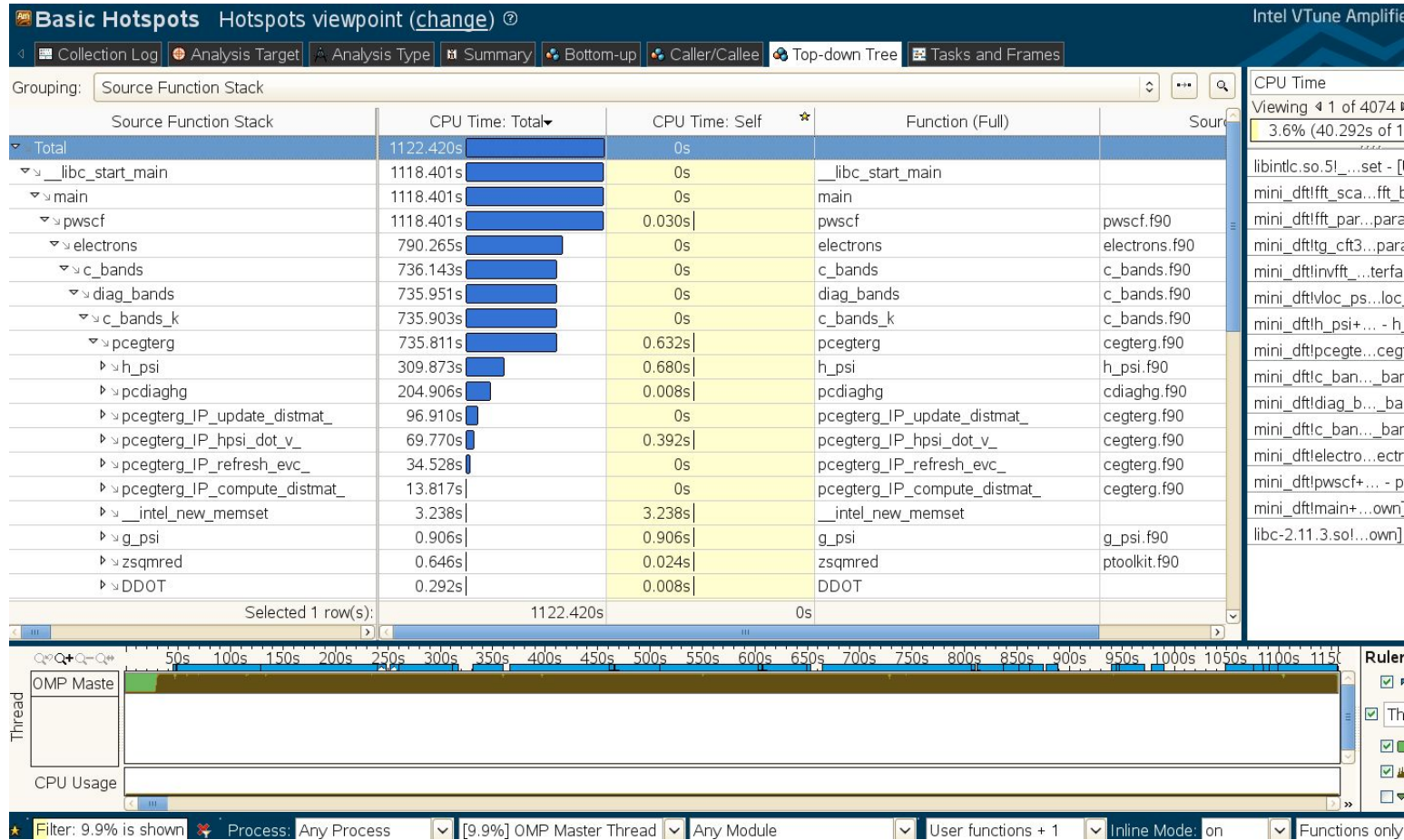
S. Li.	Source	Address	Sou.. Line	Assembly	Effect
472	! do ig = 1, igmax	0x4085d9	471	add \$0x4, %rax	0.013s
473		0x4085dd	474	vmovupdy (%r14,%rdi,1), %ymm10	0.094s
474	wdiff = wxt - wtilde_array(ig,my_igp)	0x4085e3	480	vmovupsy -0x70(%rbp), %ymm2	0.390s
475		0x4085e8	474	vmovupdy 0x20(%r14,%rdi,1), %ymm1	0.429s
476	cden = wdiff	0x4085ef	474	vsubpd %ymm10, %ymm0, %ymm5	3.454s
477	!rden = cden * CONJG(cden)	0x4085f4	474	vsubpd %ymm1, %ymm0, %ymm3	0.083s
478	!rden = 1D0 / rden	0x4085f8	480	vmulpd %ymm5, %ymm5, %ymm7	0.505s
479	!delw = wtilde_array(ig,my_igp) * CONJG(cden) * rden	0x4085fc	480	vunpckhpd %ymm5, %ymm5, %ymm9	0.342s
480	cden = 1 / cden	0x408600	480	vmulpd %ymm2, %ymm9, %ymm15	0.019s
481	delw = wtilde_array(ig,my_igp) * cden	0x408604	480	vmovddup %ymm5, %ymm8	0.080s
482	delwr = delw*CONJG(delw)	0x408608	480	vfmaddsub213pd %ymm15, %ymm13, %ymm8	0.213s
483	wdiffr = wdiff*CONJG(wdiff)	0x40860d	480	vshufpd \$0x5, %ymm7, %ymm7, %ymm6	0.380s
484		0x408612	480	vaddpd %ymm6, %ymm7, %ymm7	0.001s
485	! JRD: Complex division is hard to vectorize. So, we help the compiler.	0x408616	480	vshufpd \$0x5, %ymm8, %ymm8, %ymm9	0.075s
486	scha(ig) = mygpvar1 * aqsntemp(ig,n1) * delw * I_eps_array(ig,n1)	0x40861c	480	vdivpd %ymm7, %ymm9, %ymm6	0.232s
487	! scha_temp = mygpvar1 * aqsntemp(ig,n1) * delw * I_eps_array(ig,n1)	0x408620	480	vmulpd %ymm3, %ymm3, %ymm8	2.619s
488		0x408624	480	vunpckhpd %ymm3, %ymm3, %ymm9	0.267s
489	! JRD: This if is OK for vectorization	0x408628	480	vmulpd %ymm2, %ymm9, %ymm2	0.114s
490	if (wdiffr.gt.limittwo .and. delwr.lt.limitone) then	0x40862c	480	vmovddup %ymm3, %ymm15	0.062s
491	scht = scht + scha(ig)	0x408630	480	vfmaddsub213pd %ymm2, %ymm13, %ymm15	0.425s
492	endif	0x408635	480	vshufpd \$0x5, %ymm8, %ymm8, %ymm9	0.525s
493		0x40863b	480	vaddpd %ymm9, %ymm8, %ymm8	0.107s
494	! scha_mult = merge(1.0,0.0,wdiffr.gt.limittwo .and. delwr.lt.limitone)	0x408640	480	vshufpd \$0x5, %ymm15, %ymm15, %ymm7	0.031s
495	! scht = scht + scha(ig)*scha_mult	0x408646	480	vdivpd %ymm8, %ymm7, %ymm9	0.481s
496		0x40864b	481	vunpckhpd %ymm10, %ymm10, %ymm10	15.927s
497	enddo ! loop over g	0x408650	481	vshufpd \$0x5, %ymm6, %ymm6, %ymm2	0.107s
498		0x408655	481	vmulpd %ymm2, %ymm10, %ymm10	0.003s

Selected 1 row(s):

Highlighted 40 row(s):



Profile Your Application (VTune / CrayPat)



Approximation:

a. Real Division

b. Complex Division

c. Complex Division
+ -fp-model fast=2

?

Wall Time:

6.37 seconds

4.99 seconds

5.30 seconds

Approximation:

Wall Time:

a. Real Division	→	6.37 seconds
b. Complex Division	→	4.99 seconds
c. Complex Division + -fp-model=fast	→	5.30 seconds

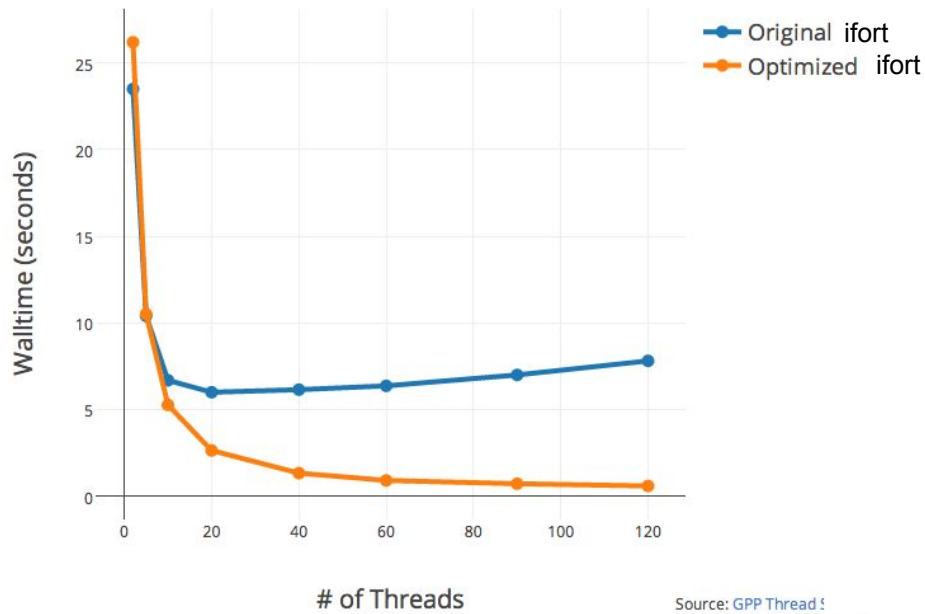
Approximation:

Wall Time:

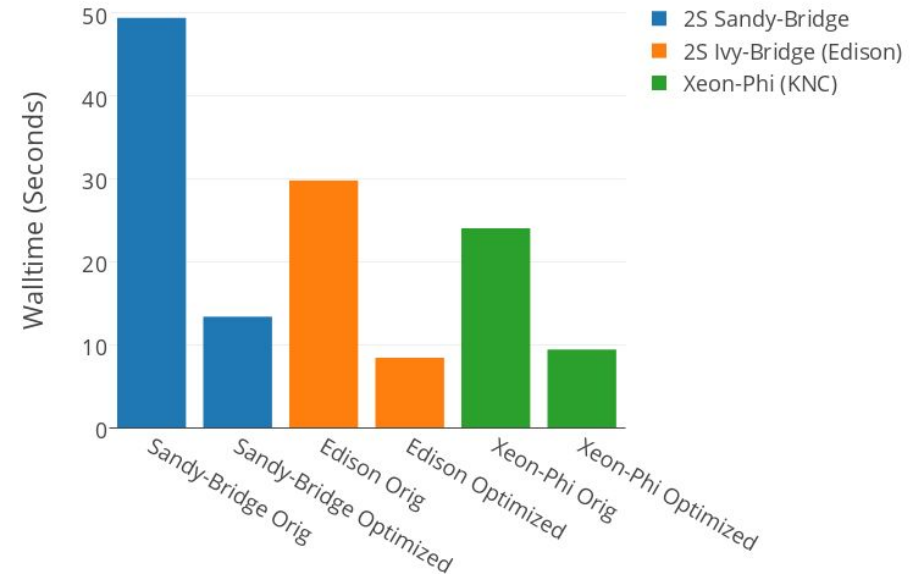
a. Real Division	6.37 seconds
b. Complex Division	4.99 seconds
c. Complex Division + -fp-model fast=2	5.30 seconds
d. Complex Division + -fp-model=fast=2 + !dir\$ nounroll	4.89 seconds

Early NESAP (Advances with Cray and Intel) Advances

Thread Scaling in BerkeleyGW GPP Kernel on Xeon-Phi

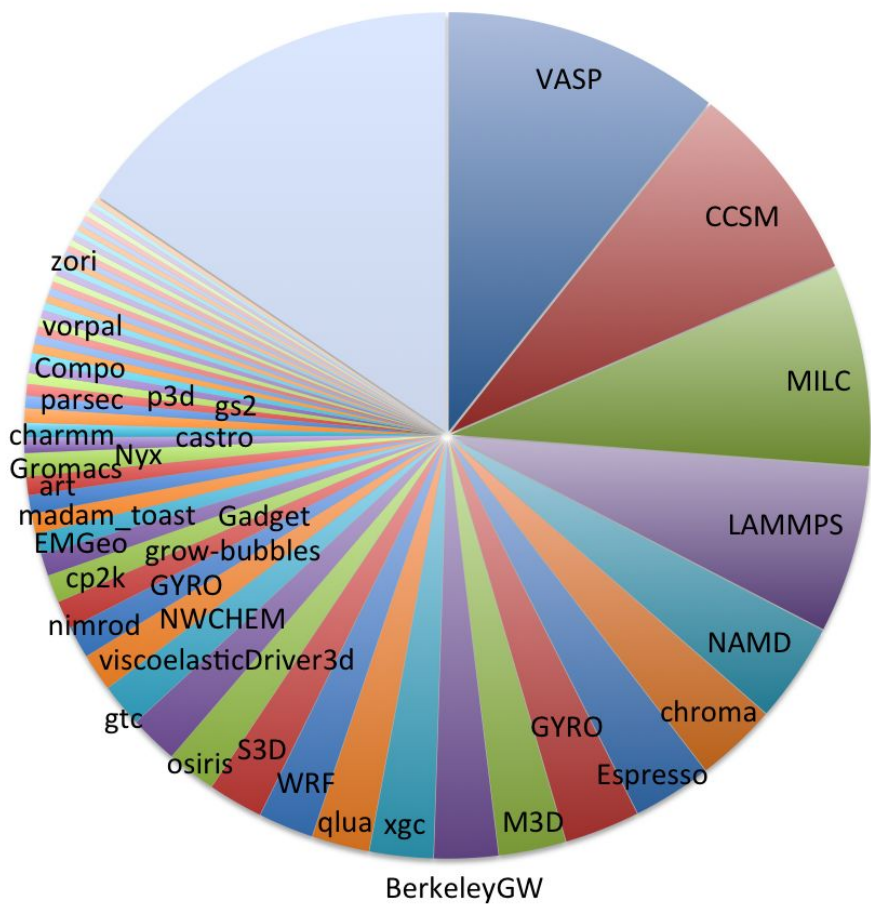


BerkeleyGW FF Kernel Runtimes on Xeon and Xeon-Phi (Nathan)

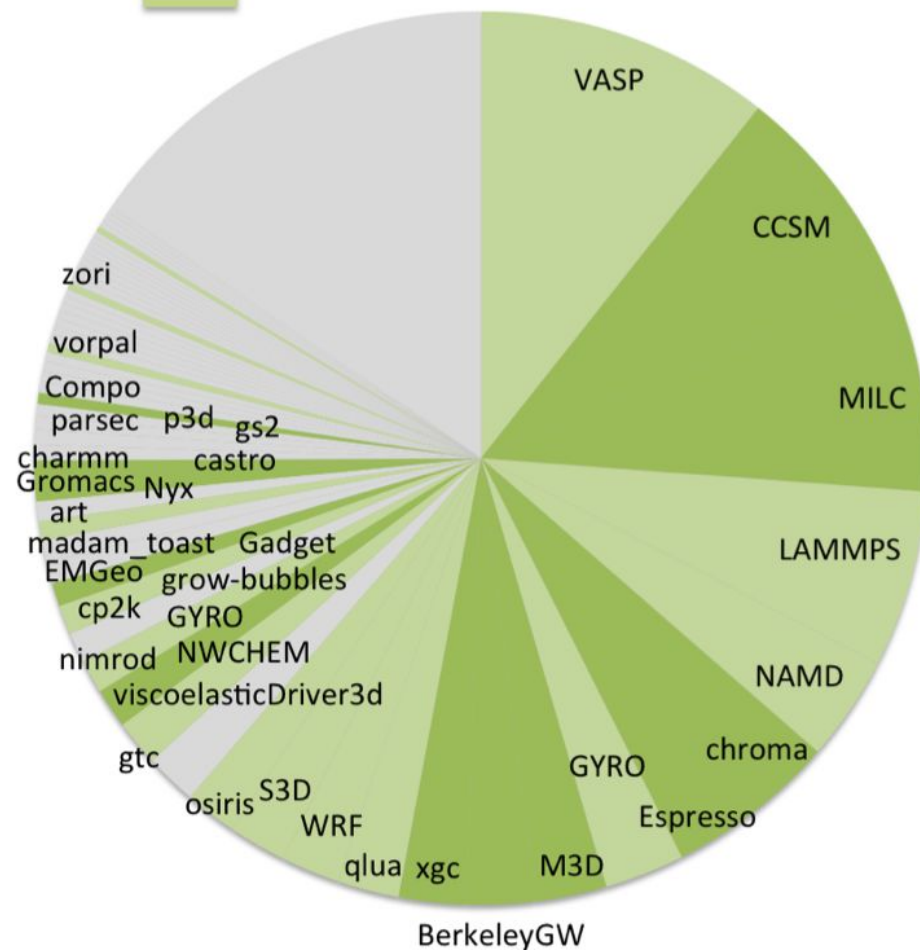


	Overall Improvement	Notes
BGW GPP Kernel	0-10%	Pretty optimized to begin with. Thread scalability improved by fixing ifort allocation performance. Unoptimized to begin with. Cache reuse improvements Moved threaded region outward in code Created custom vector matmuls
BGW FF Kernel	2x-4x	
BGW Chi Kernel	10-30%	
BGW BSE Kernel	10-50%	

Breakdown of Application Hours on Hopper and Edison 2013



NESAP Tier-1, 2 Code
 NESAP Proxy Code or Tier-3 Code



Early Lessons Learned



Cray and Intel very helpful in profiling/optimizing the code. See following slides for using Intel resources effectively

Generating small tangible kernels is important for success

Targeting Many-Core greatly helps performance back on Xeon.

Complex division is slow on (particularly on KNC)

BGW 1.0 vs 1.1 Sigma Performance

