

## LA-UR-18-27941

Approved for public release; distribution is unlimited.

Title: Improving Capsaicin Sweep Performance Using Memory Prefetching

Author(s): Halvic, Ian William

Intended for: Presentation to CCS-2

Issued: 2018-08-20

---

**Disclaimer:**

Los Alamos National Laboratory, an affirmative action/equal opportunity employer, is operated by the Los Alamos National Security, LLC for the National Nuclear Security Administration of the U.S. Department of Energy under contract DE-AC52-06NA25396. By approving this article, the publisher recognizes that the U.S. Government retains nonexclusive, royalty-free license to publish or reproduce the published form of this contribution, or to allow others to do so, for U.S. Government purposes. Los Alamos National Laboratory requests that the publisher identify this article as work performed under the auspices of the U.S. Department of Energy. Los Alamos National Laboratory strongly supports academic freedom and a researcher's right to publish; as an institution, however, the Laboratory does not endorse the viewpoint of a publication or guarantee its technical correctness.

# Improving Capsaicin Sweep Performance Using Memory Prefetching



**Ian Halvic**

8/16/18

Texas A&M University

Mentor: Kent Budge



Operated by Los Alamos National Security, LLC for the U.S. Department of Energy's NNSA

# Agenda

8/16/2018



- Prefetching Overview
- Target Routine
- Target Architecture
- Testing Method
- Haswell Results
- KNL Results
- Wrap-Up

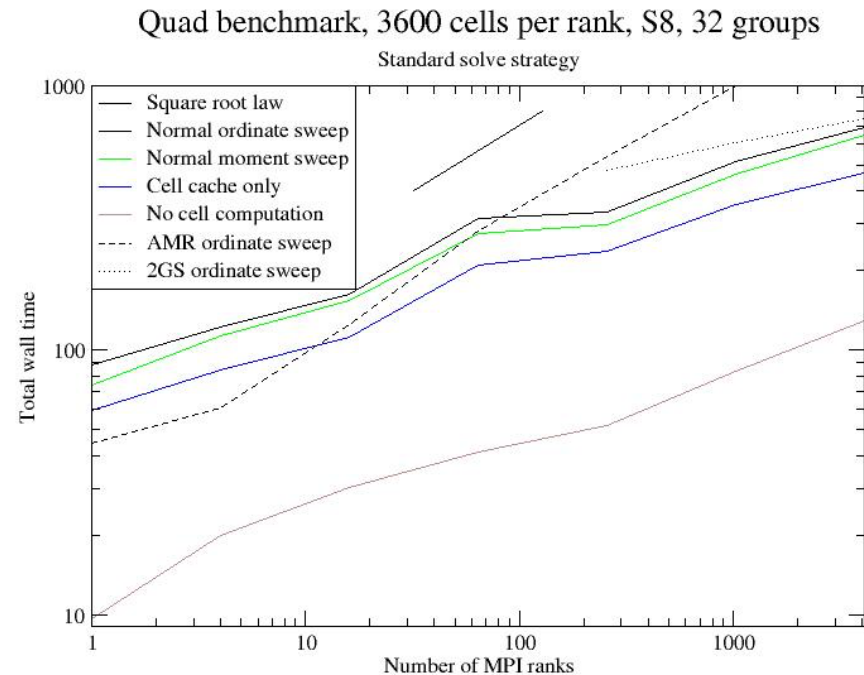
# Prefetching

- **Software Prefetching offers a straightforward method to improve runtime**
  - Fetch data from memory to cache before use (reduce latency)
  - Better utilize memory bandwidth
  - Take advantage of compute heavy sections to have memory working
  - Minimal restructuring of code (no logic changes)
- **Limitations**
  - Prefetch location in code
    - Too late: Wasted overhead
    - Too early: Eviction, Cache pollution
  - Memory bandwidth
  - Cache size
  - Overhead
  - Limited number of fetch “slots”
  - Have to keep hardware (and compiler) prefetching in mind

# Target

- **Attempt to apply technique to Capsaicin Sn solver**

- Using slightly modified code with fixed number of sweeps and no source update (Rocotillo)
- Target is the lumped DFEM Quad grid sweep
- Targeted variables:
  - scattering source
  - $\sigma$  values
  - incident flux



Difference between blue (Cache only) and pink (No cell computation) shows room for memory improvement.

# Example Prefetch Call

```
#pragma omp simd
for (unsigned e = 0; e < emax; e+=4)
{
    __builtin_prefetch(&Source[e+(0*offset)],0,3);
    __builtin_prefetch(&Source[e+(1*offset)],0,3);
    __builtin_prefetch(&Source[e+(2*offset)],0,3);
    __builtin_prefetch(&Source[e+(3*offset)],0,3);

    __builtin_prefetch (&sigma[e+(0*offset)],0,3);
    __builtin_prefetch (&sigma[e+(1*offset)],0,3);
    __builtin_prefetch (&sigma[e+(2*offset)],0,3);
    __builtin_prefetch (&sigma[e+(3*offset)],0,3);
}
```

# Target Architectures

- **Tested on Trinitite system (2 architectures)**
- **Intel Haswell E5-2698-v3 processors**
  - 16 cores per socket, 2 sockets per node
  - 40MB L3 per socket (shared).
  - Prefetching into L3 cache
- **Intel Knights Landing co-processors (KNL)**
  - 68 cores per node
  - 34MB L2 (Shared 1MB per 2 cores)
  - High Bandwidth Memory
    - High bandwidth, similar to cache
    - DRAM regime latency
  - Prefetch into L2 cache
  - Lower clock compared to Haswell

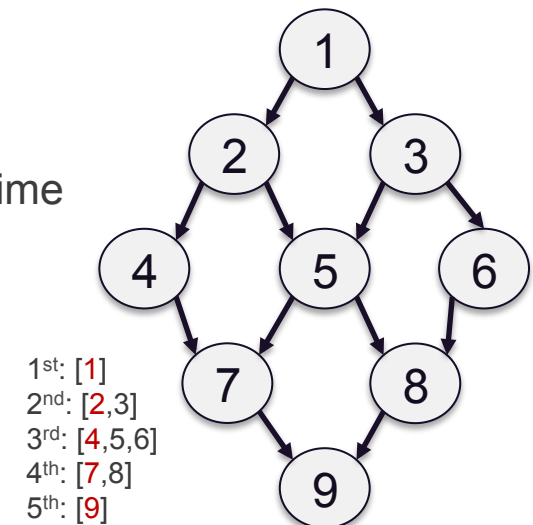


# Testing Method

- **Test problem scales mesh with number of ranks used**
  - 32 groups, N=8
  - Each rank gets an equal number of cells
    - Small problem – 6000 cells per rank
    - Large problem – 89600 cells per rank
      - Memory allocation issues with large problems using this scaling mesh
- **1 set := 5 Baseline & 5 PF runs. Averages used to compute % runtime change**
- **Intel V-tune used to determine cache misses. Typically only run for 1 rank**
- **Scaling of successful prefetch schemes with respect to MPI rank tested out to 4 full nodes**

# Prefetching for next cell (Haswell)

- **Pull data for next cell while computing current cell**
  - Extra overhead to store a link to the next cell
  - Relies on a queue of “computable” cells being populated
  - Can take advantage of current cell matrix-solve compute time
- **Showed promise for single rank**
  - Miss reduction from ~2.2B to 0.5B
  - 5% runtime improvement
- **Poor results for parallel runs**
  - More frequently emptying computable cell stack
  - Potential Bandwidth issues
  - Runtime degradation



Compute stack refreshes for above graph, scheme cannot prefetch for the first element of refresh (shown in red).

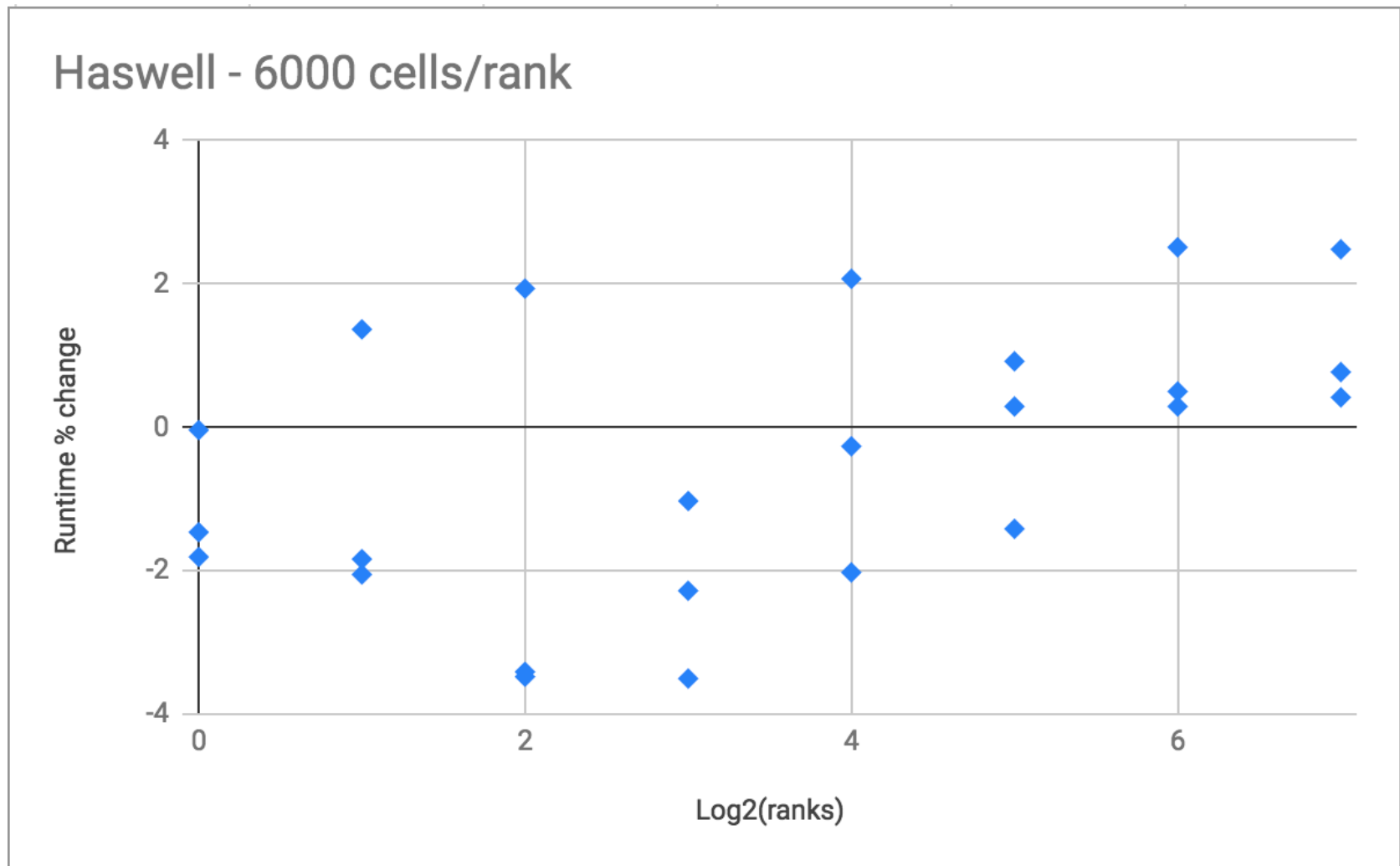
***This method has issues with parallel scaling, graph dependence.***

# Prefetching within cell (Haswell)

- **Focus primarily on improving performance within a single cell**
- **Numerous configurations tested**
  - Many which reduce misses but offer no runtime improvement
- **Best performing scheme (1 rank)**
  - Miss reduction from ~2.2B to 0.5B (similar to next cell PF method)
  - Runtime improvement of ~3%
- **Expanding to parallel reduces runtime impact**
  - Issues set in very quickly, ~8 to 16 ranks
  - Runtimes became more varied with higher ranks, difficult to analyze

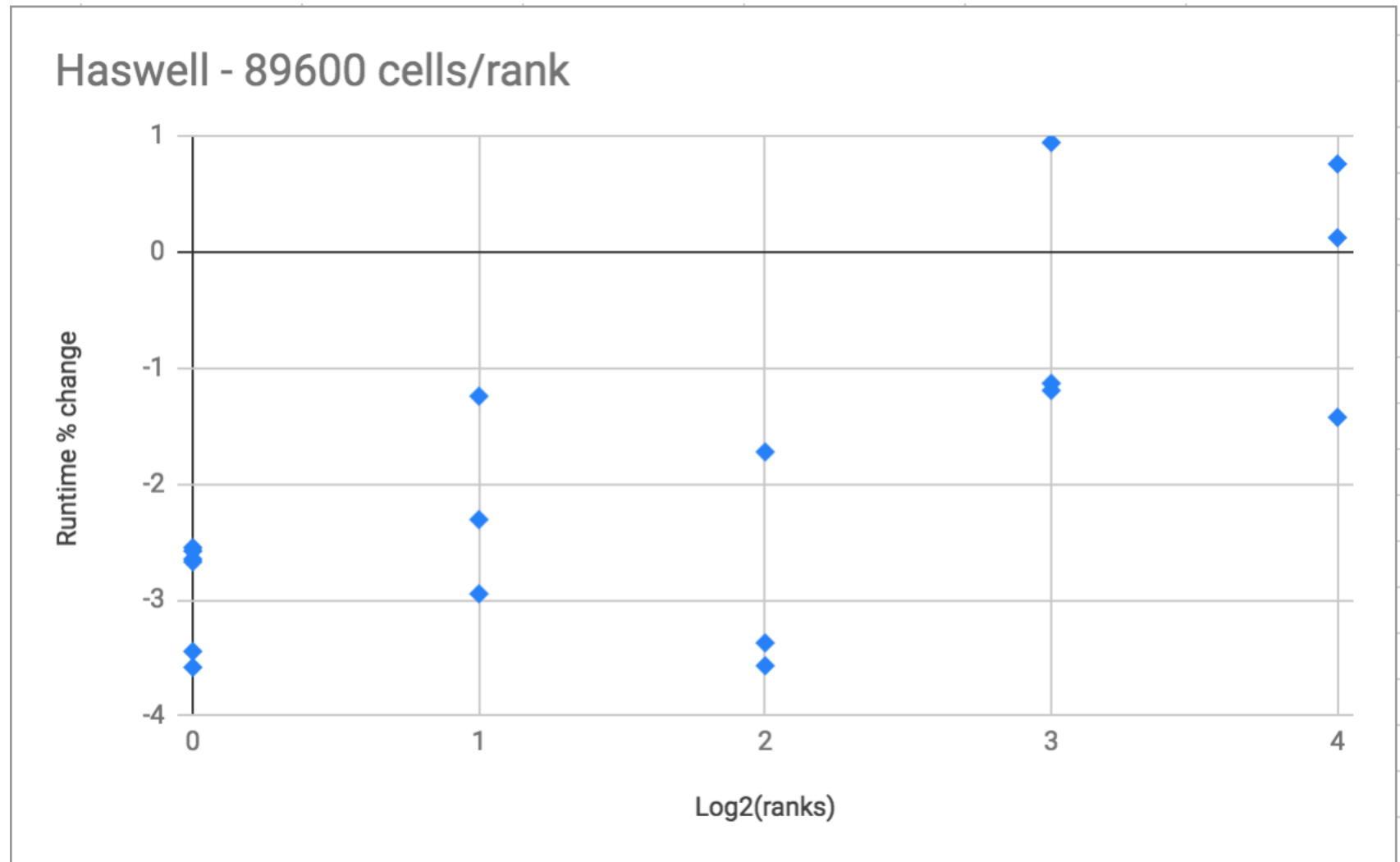
# Haswell: Runtime Change Scaling

Small problem tested from 1 to 128 ranks (4 nodes)



# Haswell: Runtime Change Scaling

Larger problem tested from 1 to 16 ranks



# Haswell: V-tune and Runtime Comparison

## Elapsed Time<sup>?</sup>: 127.653s

CPU Time <sup>?</sup> :	115.753s	
Memory Bound <sup>?</sup> :	61.3%	of Pipeline Slots
L1 Bound <sup>?</sup> :	7.3%	of Clockticks
DRAM Bound <sup>?</sup> :		
DRAM Bandwidth Bound <sup>?</sup> :	0.0%	of Elapsed Time
NUMA: % of Remote Accesses <sup>?</sup> :	0.0%	
Loads:	77,299,118,904	
Stores:	48,601,057,988	
LLC Miss Count <sup>?</sup> :	2,242,934,568	
Local DRAM Access Count <sup>?</sup> :	2,165,864,974	
Remote DRAM Access Count <sup>?</sup> :	0	
Remote Cache Access Count <sup>?</sup> :	700,021	
Average Latency (cycles) <sup>?</sup> :	40	
Total Thread Count:	1	
Paused Time <sup>?</sup> :	0s	

## Elapsed Time<sup>?</sup>: 126.159s

CPU Time <sup>?</sup> :	114.328s	
Memory Bound <sup>?</sup> :	50.7%	of Pipeline Slots
L1 Bound <sup>?</sup> :	12.6%	of Clockticks
DRAM Bound <sup>?</sup> :		
DRAM Bandwidth Bound <sup>?</sup> :	0.0%	of Elapsed Time
NUMA: % of Remote Accesses <sup>?</sup> :	0.0%	
Loads:	83,645,509,290	
Stores:	50,582,817,439	
LLC Miss Count <sup>?</sup> :	473,228,392	
Local DRAM Access Count <sup>?</sup> :	459,213,776	
Remote DRAM Access Count <sup>?</sup> :	0	
Remote Cache Access Count <sup>?</sup> :	0	
Average Latency (cycles) <sup>?</sup> :	44	
Total Thread Count:	1	
Paused Time <sup>?</sup> :	0s	

- Approximately 75% reduction in LLC misses
- Scaling issues may be caused by shared L3 cache and related resources
  - L2 prefetching tested but resulted in runtime degradation

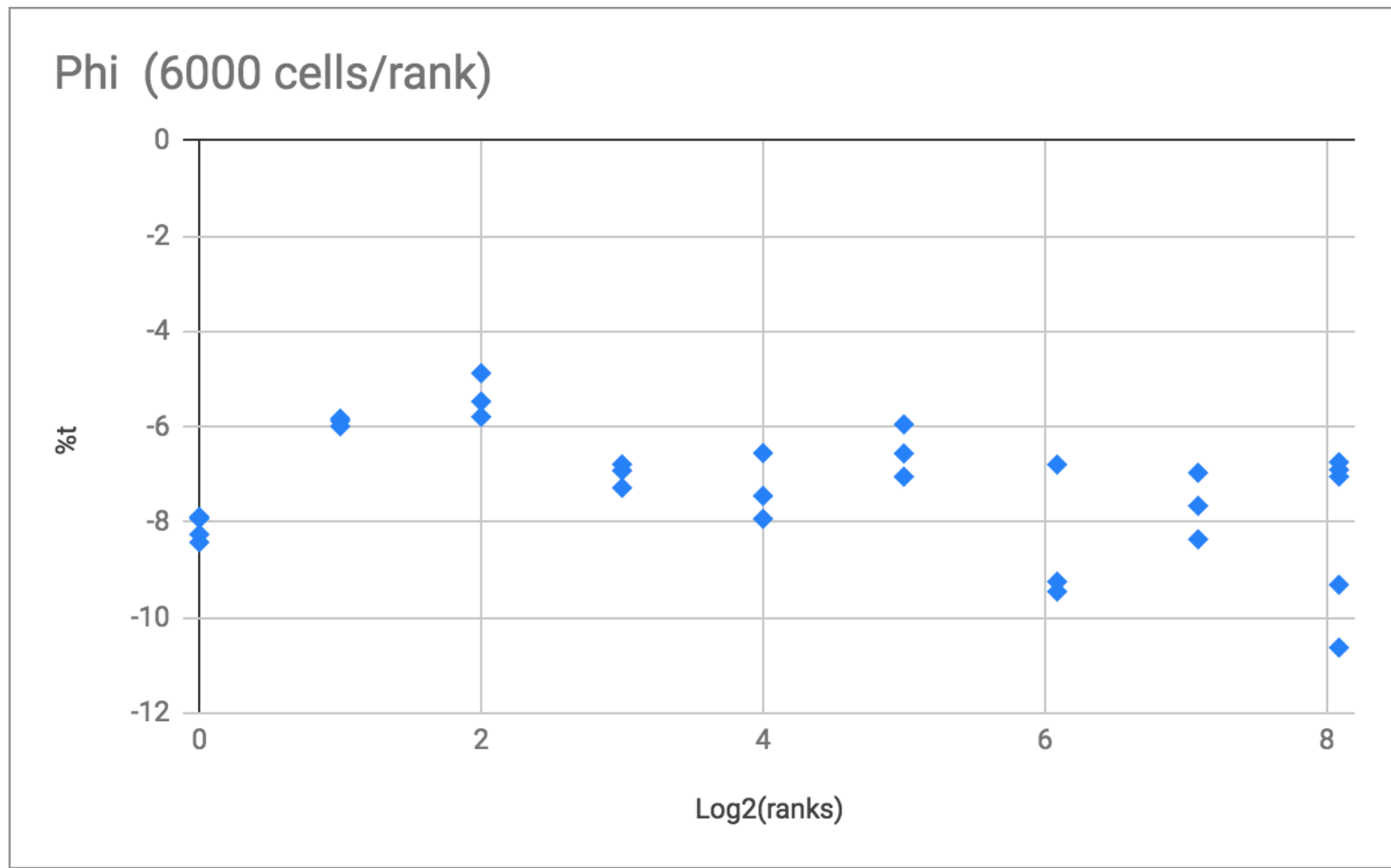
***Even single cell PF schemes can have issues with scaling***

# KNL: Prefetching Within Cell

- **Using same prefetch scheme as best performing Haswell**
- **Single rank performance gain**
  - Miss reduction from 160M to 40M
  - Frequent ~7%-12% Runtime improvement (up to ~14% observed)
- **Appears to scale to multirank runs well**
  - Similar runtime improvement observed on 272 ranks (4 nodes)

# KNL: Runtime Change Scaling

Small problem tested from 1 to 272 ranks (4 nodes)





# KNL: V-tune and Runtime Comparison

**Elapsed Time<sup>②</sup>: 23.401s**

CPU Time<sup>⑦</sup>: 21.480s

☑ **Memory Bound:**

L2 Hit Rate<sup>⑦</sup>: 56.5%

L2 Hit Bound<sup>③</sup>: 11.2% of Clockticks

➤ **L2 Miss Bound<sup>⑦</sup>: 100.0% 🚩 of Clockticks**

MCDRAM Bandwidth Bound<sup>③</sup>: 0.0%

DRAM Bandwidth Bound<sup>⑦</sup>: 0.0% of Elapsed Time

L2 Miss Count<sup>③</sup>: 162,004,860

MCDRAM Hit Rate: 99.3%

MCDRAM HitM Rate: 89.6%

Total Thread Count: 1

Paused Time<sup>③</sup>: 0s

**Elapsed Time<sup>②</sup>: 21.392s** 📄

CPU Time<sup>⑦</sup>: 19.560s

☑ **Memory Bound:**

L2 Hit Rate<sup>⑦</sup>: 84.9%

L2 Hit Bound<sup>③</sup>: 12.7% of Clockticks

➤ **L2 Miss Bound<sup>⑦</sup>: 30.5% 🚩 of Clockticks**

MCDRAM Bandwidth Bound<sup>③</sup>: 0.0%

DRAM Bandwidth Bound<sup>⑦</sup>: 0.0% of Elapsed Time

L2 Miss Count<sup>③</sup>: 39,001,170

MCDRAM Hit Rate: 99.4%

MCDRAM HitM Rate: 91.7%

Total Thread Count: 1

Paused Time<sup>③</sup>: 0s

- Once again approximately 75% reduction in LLC misses
- Better scaling may be explained by the “less shared” nature of the LLC for the KNL compared to Haswell.

***KNLs show promise for this prefetching scheme***

# Less Misses $\neq$ Reduced Runtime

- **The vast majority of schemes tested on Haswell reduced misses (often by a significant amount) but had little effect on runtime**
  - Potentially overloading fetch “slots”?
  - Prefetch call overhead?
    - Thought: Put in pointless Prefetch calls and see the effect of overhead
- **Poor scaling despite similar cache miss reduction**
  - At 16 ranks, Haswell still showed ~75% miss reduction
  - More confusing: At 16 ranks KNL only showed ~55% miss reduction

# Wrap-Up

- **Reduction in cache misses will not always correspond to decreased runtimes**
  - Could be overtaxing fetching resources
  - Some overhead to consider
- **Prefetching for next graph node results in graph dependence, which complicates scaling to higher ranks.**
- **More performance improvements due to prefetching seem to be available for KNL compared to Haswell**
  - Possibly due to fetching into shared L3 cache for Haswell compared to less-shared L2 for KNL (16 fast cores in 40MB versus 2 slow cores in 1 MB)
- **Implementation of prefetching for complex code seems to require a lot of guess and check**

# Future Directions

- **Test prefetching scheme on larger “real” problems with full transport solve**
- **Adapt scheme to other mesh geometries beyond what was tested here**
- **Deeper dive into why a 75% reduction in cache misses does little for Haswell, as well as its poor scaling**
- **Gather more scaling data points (runtime and miss changes) out to higher number of ranks**
  - Can probably script this