

Impact of Write-Allocate Elimination on Fujitsu A64FX

Yan Kang

The Pennsylvania State University
Pennsylvania, USA

Sayan Ghosh

Pacific Northwest National Laboratory
Washington, USA

Mahmut Kandemir

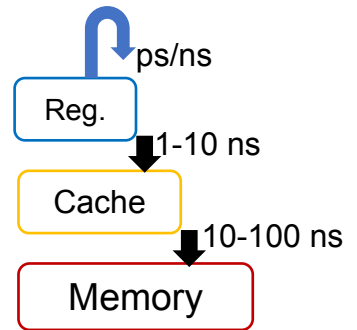
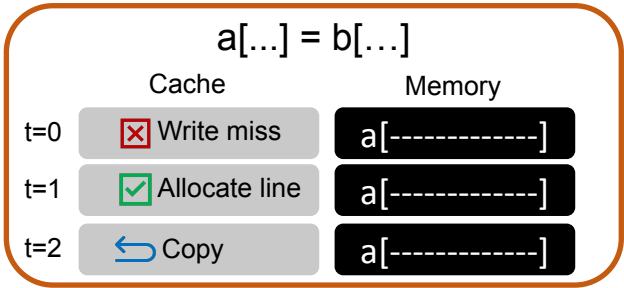
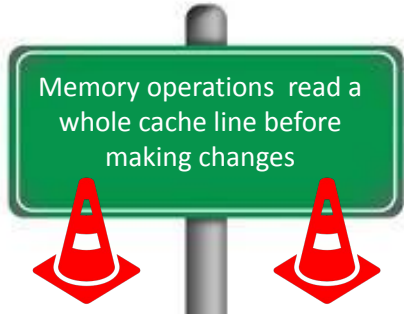
The Pennsylvania State University
Pennsylvania, USA

Andrés Marquez

Pacific Northwest National Laboratory
Washington, USA



Motivation



If the data is not going to be reused soon, what is the point to cache it?

non-temporal write/stores:

Directly write new data into memory instead of read a cache line then modify it

Write-Allocate : Allocate a cache line on a write miss

Code modifications ???

Compiler automatic ???

Performance improvements ???

Write-Allocate Avoidance

Evasion: [Hardware detects if cache line is going to be overwritten] store cache line directly in memory (Intel, non-temporal stores, compiler hints or automatic Spec12M)

Elimination: [Hardware detects if cache line is going to be overwritten] directly write an L2 cache line with zeroes, processor loads cache line avoiding memory read

Fujitsu A64FX: Elimination is available through a special 64-bit instruction (DC ZVA) in the ARMv8-A

Can write-allocate elimination via “zero fill” improve the performance of various applications on Fujitsu A64FX?



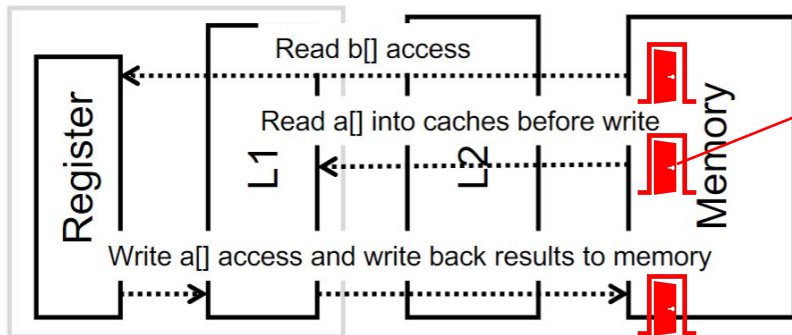
Read Dr. Georg Hager’s blog post and paper:

<https://blogs.fau.de/hager/archives/8997>

<https://onlinelibrary.wiley.com/doi/10.1002/cpe.6512>

Zero Filling in Fujitsu A64FX

Without ZFILL



This is the memory access we are trying to eliminate!

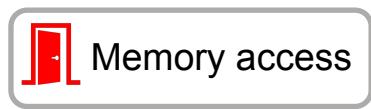
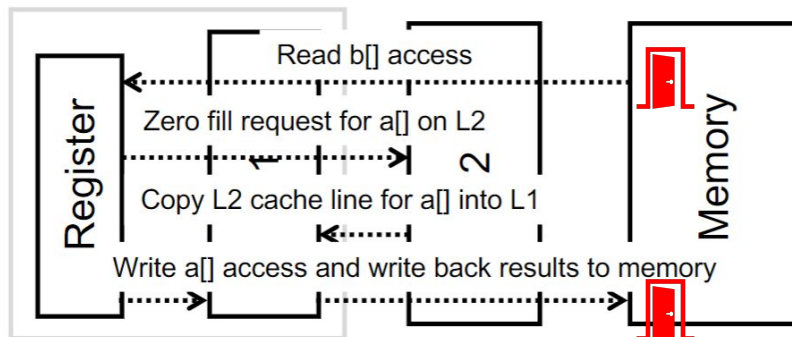
“zero fill” on L2 Cache:

Upon receiving the DC ZVA request, the L2 cache secures the cache line corresponding to the specified virtual address and writes zero data

“zero fill” on L1 Cache:

zero data is written after data in the L1 cache is written back to the L2 cache.

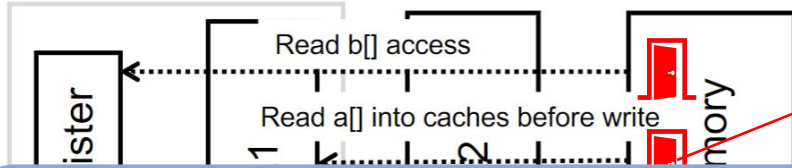
With ZFILL





Zero Filling in Fujitsu A64FX

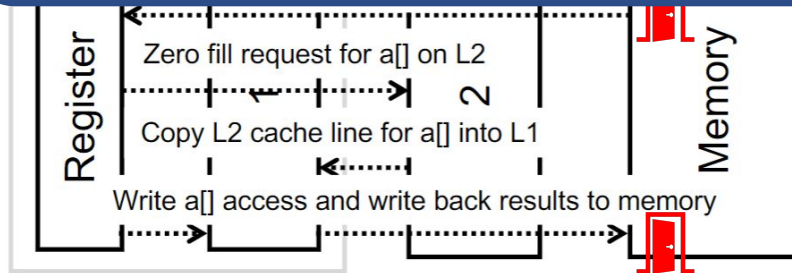
Without ZFILL



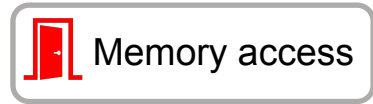
This is the memory access we are trying to eliminate!

Saving memory traffic means improving memory b/w, what's that benchmark to study "sustainable memory b/w"?

the L2



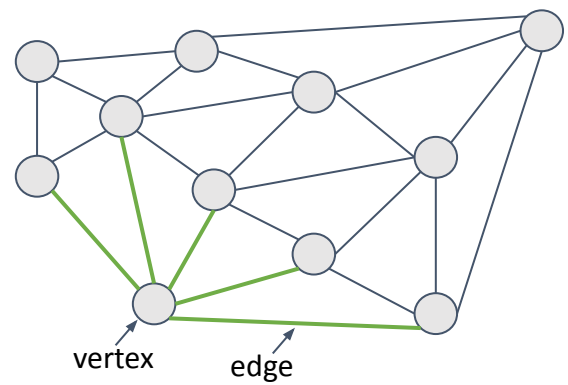
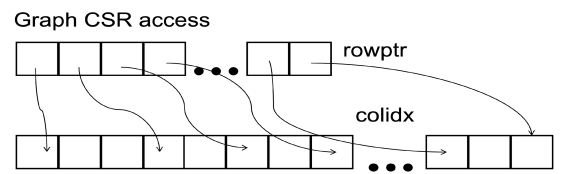
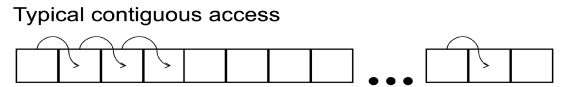
zero data is written after data in the L1 cache is written back to the L2 cache.





Benchmarking decisions

- ❑ STREAM is “best case” memory b/w benchmark
 - ❑ Does not represent irregular cases, most applications
- ❑ Graphs – irregular memory accesses
 - ❑ Applications perform repetitive *neighborhood accesses*
- ❑ NEVE is a benchmark, like STREAM for graphs (has COPY, SUM and MAX) - $|V| * |E| * \#ops / t$



STREAM		Graph Neighborhood Access Kernels	
name	kernel	name	kernel
COPY:	$a(i) = b(i)$	Neighbor Copy:	$a(i, j) \leftarrow COPY(\forall i \in V \forall j \in neighbor(i) weight(j))$
SCALE:	$a(i) = q * b(i)$	Neighbor Add:	$b(i) \leftarrow SUM(\forall j \in V \forall j \in neighbor(i) weight(j))$
SUM:	$a(i) = b(i) + c(i)$	Neighbor Max:	$c(i) \leftarrow MAX(\forall i \in V \forall j \in neighbor(i) weight(j))$
TRIAD:	$a(i) = b(i) + q * c(i)$		

Can return MB/s!

Explicit “Zero Fill” formulation for graph neighborhood accesses

```

1 static const int DISTANCE = 100;
2 static const int ELEMS_CACHE_LINE = 256 / sizeof(double);
3 static const int OFFSET = DISTANCE * ELEM_CACHE_LINE;
4
5 static inline void zfill(double * a) {
6     asm volatile("dc zva, %0" : : "r"(a));
7 }
8
9 #pragma omp parallel
10 {
11     int const tid = omp_get_thread_num();
12     int const nthreads = omp_get_num_threads();
13     int chunk = nvertices / nthreads;
14     double* const zfill_limit = c + (tid+1)*chunk - OFFSET;
15
16     #pragma omp for schedule(static)
17     for (int j=0; j<nvertices; j+=ELEMS_CACHE_LINE) {
18         int const * __restrict__ const jrowptr = rowptr + j;
19         double * __restrict__ const jbuf = buf + j;
20         double sum = 0.0;
21
22         if (jbuf+OFFSET < zfill_limit)
23             zfill(jbuf+OFFSET);
24
25         for (int i=0; i<ELEMS_CACHE_LINE; ++i) {
26             for (int e=jrowptr[i]; e<jrowptr[i+1]; ++e) {
27                 sum += colidx[e].weight;
28             }
29             jbuf[i] = sum;
30         }
31     } // loop over vertices
32 } // openmp

```

Explicit assembly to
invoke DC ZVA

Each thread works on fixed
chunk of iterations over $|V|$
(work is variable)

Block outermost loop
over vertices

Invoke zero fill in strides larger
than L2 prefetch distance

Inner loop, where the zfill virtual
address will be invoked several
times (trip count unknown)

Explicit “Zero Fill” formulation for graph neighborhood accesses

```

1 static const int DISTANCE = 100;
2 static const int ELEMS_CACHE_LINE = 256 / sizeof(double);
3 static const int OFFSET = DISTANCE * ELEM_CACHE_LINE;
4
5 static inline void zfill(double * a) {
6     asm volatile("dc zva, %0" : "r"(a));
7 }
8
9 #pragma omp parallel
10 {

```

Explicit assembly to
invoke DC ZVA



"Zero filling" may not yield performance benefits for irregular graph workloads unless the number of vertices in a graph is significantly larger than the cache line size and the standard deviation of the vertex degrees is relatively low.

```

20 double sum = 0.0;
21
22 if (jbuf+OFFSET < zfill_limit)
23     zfill(jbuf+OFFSET);
24
25 for (int i=0; i<ELEMS_CACHE_LINE; ++i) {
26     for (int e=jrowptr[i]; e<jrowptr[i+1]; ++e) {
27         sum += colidx[e].weight;
28     }
29     jbuf[i] = sum;
30 }
31 } // loop over vertices
32 } // openmp

```

over vertices

Invoke zero fill in strides larger
than L2 prefetch distance



Inner loop, where the zfill virtual
address will be invoked several
times (trip count unknown)





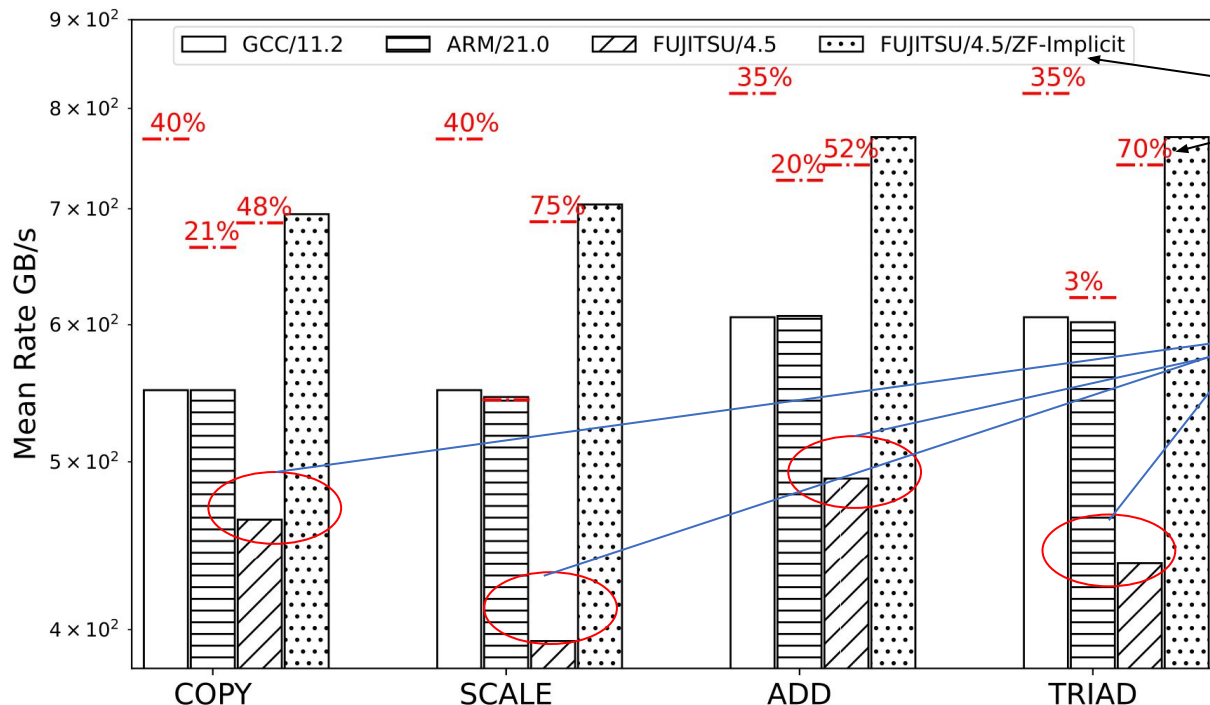
Benchmarks and applications for evaluations

Benchmark Scenarios	Tested Kernels
STREAM	Copy
	Scale
	Add
	Triad
Graph Neighborhood Kernels(NEVE)	Add
	Copy
	Max

Expecting STREAM to be the best case!

Application scenarios	Targeted kernels
Graph500 Breadth First Search	Next frontier list update is similar to graph neighborhood Copy
Louvain graph clustering	Modularity computation requires summing data, similar to graph neighborhood Add.
GAP benchmark suite	
Breadth First Search (BFS)	Next frontier list update is similar to graph neighborhood Copy.
PageRank (PR and PR(SPMV))	Score update is similar to STREAM Copy.
Connected Components (CC and CC(SV))	Singleton partition assignment is similar to STREAM Copy.
Betweenness Centrality (BC)	Aggregation of betweenness scores similar to graph neighborhood Add.
Rodinia benchmark suite	
HotSpot & HotSpot3D	Grid cell update of tiled 2D/3D transient iterative solver similar to STREAM Copy.
LavaMD	Particle updates in the home box similar to STREAM Add.
SRAD(V2)	Image data update is similar to STREAM Add.
Needleman-Wunsch (NW)	Global copy from local reference is similar to STREAM Copy.

STREAM benchmark evaluations (GCC, ARM and FCC)

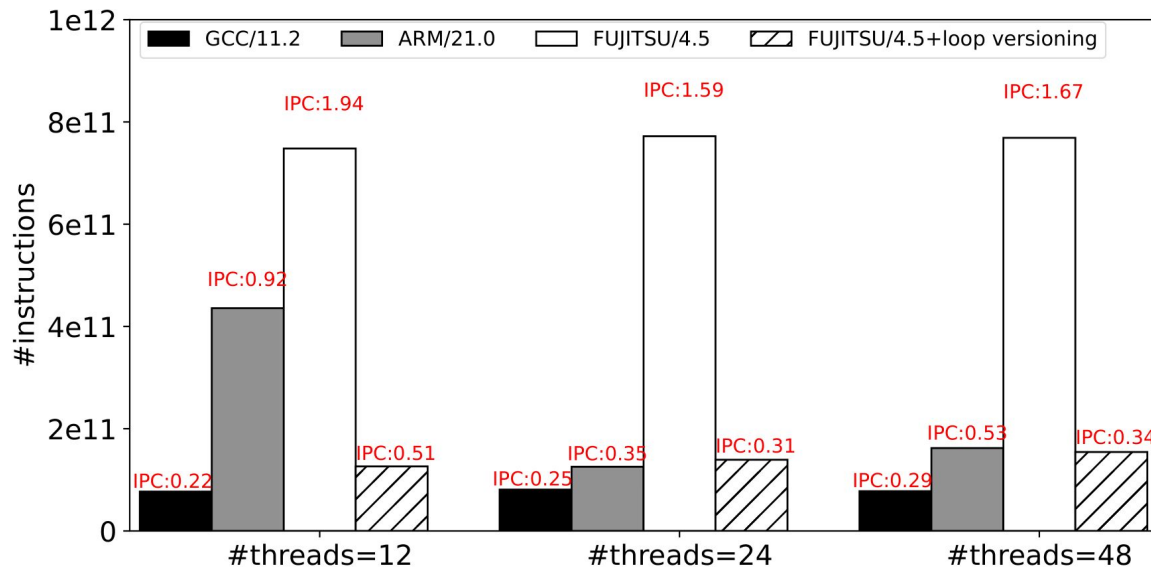


Fujitsu has a compiler option (-Kzfill), referred as implicit version [does not work for C++ compiler]

Fujitsu is showing a huge performance gap comparing with GCC/ARM ???

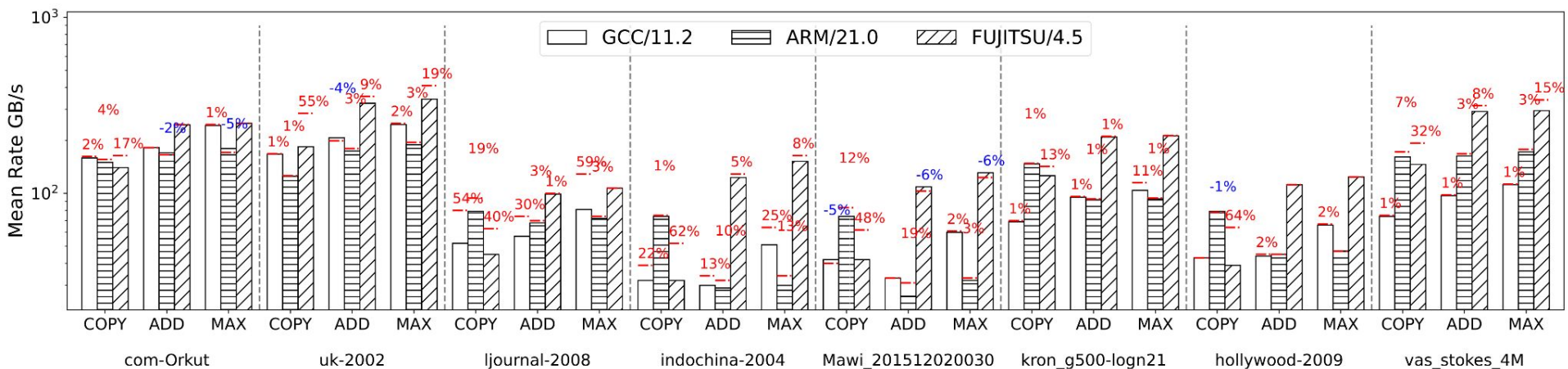
All compilers demonstrate improvements, FCC up to 70%!

Loop versioning for STREAM



Enabling loop versioning option allows the compiler to perform software pipelining optimizations, bringing the overall performance close to ARM/GCC.

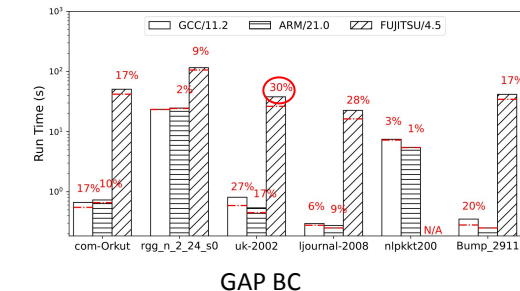
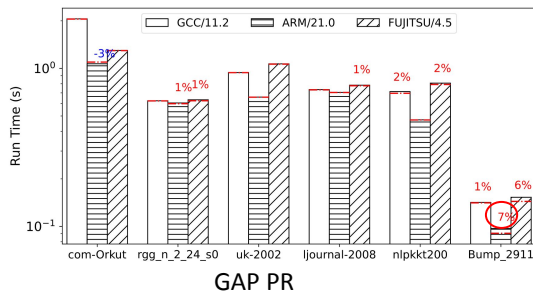
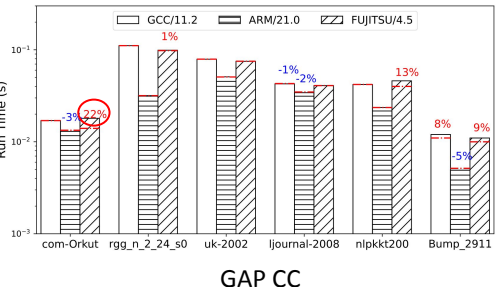
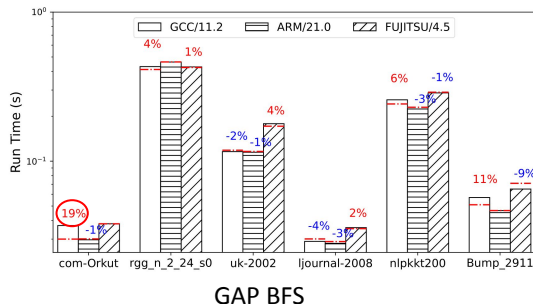
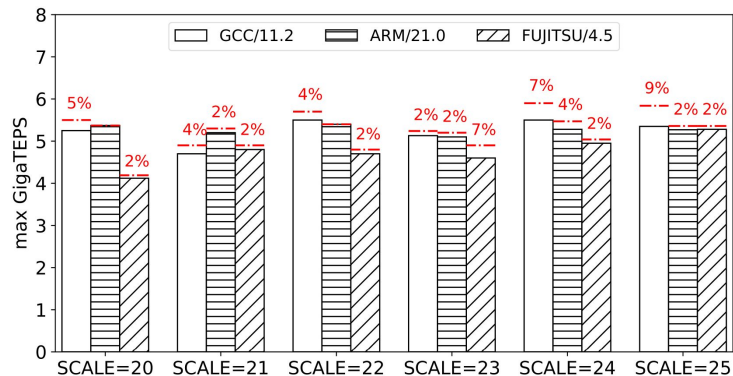
Graph benchmark evaluations (GCC, ARM and FCC)



- Used different graphs – implies different structure/work-per-loop
- ZFILL: degradation of up to 6% but also up to 64% improvement (FCC)
- Since the “zero fill” stride length can be greater than the median #edges for certain graphs, it can have a limited impact and, in some cases, may incur overheads
- No compiler automatic DC ZVA guaranteed



Graph Application Evaluations

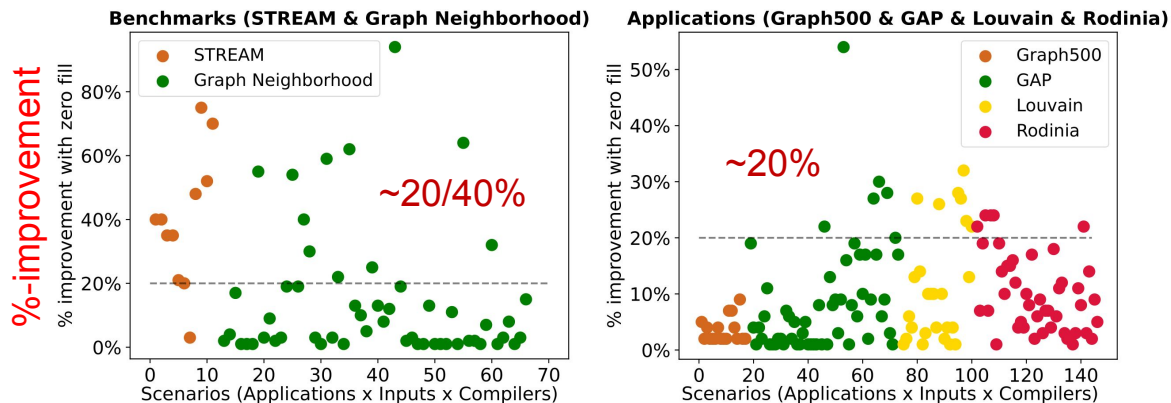


- Does not improve performance where there is limited work in the ZFILL section
- ~15% improvement when there is sufficient work
- No compiler automatic DC ZVA guaranteed

Non-temporal store patterns may not positively impact overall application execution times if they are not on the critical performance path.

Observations

- NEVE exhibit about 2–5x performance degradation compared to STREAM
- End-to-end improvements between 5-20% for benchmarks and diverse application scenarios due to "zero fill" adaptations
- Performance improvements of up to 32% in Louvain clustering and median improvements of 5–17% in GAP PR, CC, and BC benchmarks



Limitations of Zero Filling

- Number of vertices in a graph may not larger than the cache line size.
- Regular prefetching may benefit more for short buffer streaming write.
- Non-temporal store may not be on the critical performance path.
- Applications may not be written to exploit non-temporal stores.

Summary & Future Steps

Demonstrated the impact of write allocate elimination on A64fx for various applications and show cased the improvement brought by “Zero fill”

Identified possible causes might lead to our observations on the varied performance improvements, such as software pipelining optimization, SIMD extensions, etc.

What exactly reasons behind each or all applications variations on A64FX ???

In what condition it maximize improvement?

Where to apply?

In what condition it requires minimal modifications ?

Compiler automatic DC ZVA generation?

Acknowledgements

- PNNL LDRD Data Model-Convergence (PI: Sayan Ghosh, PNNL)
- DOE ASCR Advanced Memory to Support Artificial Intelligence for Science (AIAMS, PI: Andrés Márquez, PNNL)
- Penn State HPCL (Prof. Mahmut Kandemir)
- Ookami testbed support (Dr. Eva Siegmann and team, SBU)
- IWAHPCE'24 paper reviewers