

# Reducing Dark Bandwidth Through Data Reduction Near Memory

Jonathan Beard

Staff Research Engineer, Arm Inc. UCAR SEA, Boulder, CO

4 April 2018

Twitter: @jonathan\_beard

Copyright © 2018 Arm Limited

## **Disclaimer**

This is a research project, not product.

Any changes to architecture must go through the standard architecture review board process

Nothing within this presentation should be taken as a direction of future architecture or product plans for Arm Inc/Ltd.

Please ask questions related to the presentation!

## The lunch line is closed

- > Put simply, frequency scaling is not a solution
- Performance curve must continue, we need FLOPS/IOPS/etc.

## > As a consequence:

- > Core counts on chip and accelerators on chip will become more numerous
- Integration will be closer, move onto system-on-chip (and multi-part chip)
- > What used to move off-chip before now moves on-chip, further stressing on chip interconnects
- Increased compute density of the on-chip demands even more off-chip data movement bandwidth
- > As a result, less memory bandwidth/cache/interconnect capacity per core on average....

## > Core problem:

- > We can build machines with tons of FLOPS, but we can't feed them
- > Data movement is a key problem...guess what, we move a lot of data we don't need to move







## **Cores are reuse optimized**





# Dark Bandwidth:

# Data moved throughout the memory hierarchy that is not used after it is moved.





## How we get data

move -16(%r1), %r2



## How we get data

move -16(%r1), %r2

























## **Reuse Distance**

If data within a cache line is used more than once, then the cache line is said to be reused

We need to answer: how many other cache lines are used in between accesses?







## **Reuse Distance**

If data within a cache line is used more than once, then the cache line is said to be reused

We need to answer: how many other cache lines are used in between accesses?



## **Reuse Distance**

If data within a cache line is used more than once, then the cache line is said to be reused

We need to answer: how many other cache lines are used in between accesses?





How much of a cache line is used before it is evicted?

We moved all the data labeled waste through that entire maze of wires only to evict it again!!







## **Some Real Applications – Linked List**





**Some Real Applications – Others** 

### Cache Line Utilization

24 © Arm 2018

arm

## **Dance of reuse and utilization**



## **Dark Bandwidth**

- Here's what we're actually getting (pointer-chasing application)
- > Avg. 12.5% line utilization



30 GB/s - 60 GB/s

## **Dark Bandwidth**

- Here's what we're actually getting (pointer-chasing application)
- > Avg. 12.5% line utilization



# Dark Bandwidth:

Unused data moved throughout the memory hierarchy not at request of programmer but because of hardware design.

## What this means for power



## What this means for power



## What this means for time to solution



arm

# Extant options to solve the problem



## **Conventional Gather/Scatter**

3 cache lines



## **Improving on scatter/gather**



- Reduce cache misses for sparse workloads
- Improve Bandwidth Utilization

## **Programmable DMA**

- Generally limited to page gather/move/scatter
- Usually coherent (lots of transactions), wastes bus utilization (i.e., all traffic goes through coherent network)
- Often limited when it comes to finer granularity
- > Typically no persistent state
- > Not solving quite the same problem!

### **Toy Example:**



## What we need to do



## **Sparse Data Reduction (mechanisms)**

### Standard Gather-Scatter (In Core)

**In-Memory Gather-Scatter** 



• Smaller binary (e.g., single gather vs. multiple loads)

- Smaller binary (e.g., single gather vs. multiple loads)
- Improved Cache Line Utilization
- Reduced Overall Data Movement

## **SPiDRE in words**

Decoupling access and execute

## The idea

- Programmers know what data they want (usually), "SQL for Memory"
- Finding Compressed Sparse Row manipulations for new problems often takes a lot of time (bad for productivity)
- Why not build an interface that allows gathering (potentially) in any device of only the data you (the programmer) needs
- "Bulk" byte-addressed memory
- Doesn't break coherence...with some caveats
- Can create multiple windows of original segment



# The SPiDRE



## **Programming SPiDRE (one approach)**

User inserted allocate (new virtual and physical space)

```
User inserted rearrange function
```

```
template < class DST, class SRC >
static std::size_t rearrange(
   DST * const dst,
   SRC * const src,
   const std::size_t nitems,
   rearrange_func_t< DST, SRC > src_dst );
```

```
template < class S, class S_Prime >
 static void sync( S * const s,
     S_Prime * const s_p,
     const std::size_t nitems );
```

```
Synchronize S' -> S
```

template < typename TYPE >
 static void free( TYPE \*ptr );















## The big picture





## **Sparse translation through IOMMU**





# Sim Setup & Results



# **Emulation Environment (SPiDRE)**



## Fixed Stride (1GB data set) - Gather SPiDRE



## Fixed Stride (1GB data set) - Gather SPiDRE



## Fixed Stride (1GB data set) - Gather SPiDRE



arm

## **Random Gather - SPiDRE**



## **HPC Mini-apps**

All single threaded executions for initial study

### CoMD

- Used Lennard-Jones potentials
- Simple port to gather within the main loop

### LULESH

- Multiple ports with varying rearrange to use distances
- Varied size and iteration count
- Demonstrated sparse data reduction combined with programmer placed pre-fetch hints







**Emulator Based - Latency** 



© Arm 2018

## **Speed-up by average stride**



orm

## Early Results Speedup

- Emulator on actual hardware:
  - 2x-5x speed-up for random access
  - -~1.4 2.6x speed-up for Page Rank
  - -~2x speed-up for LULESH

- Full system gem5 (early results)
  - -~10x speed-up HPCG



## Conclusions

- SPiDRE is an interface and hardware acceleration infrastructure to gather data near memory/storage and make it dense (reducing bandwidth utilization, enabling more vectorization)
- We've shown a 2-10x speedup and significant data movement reduction on several applications, definitely more room (some cases greater than 2x)
- Not shown, increase in vectorization....but pretty obvious.
  Future work
  - Translation for near-memory compute (inc. gather/scatter)
  - > Programming models, compilers (incorporate cost models with llvm, etc.)
  - More full system simulation

# Questions...

Twitter: @jonathan\_beard Email: jonathan.beard@arm.com