CIM Research

Memory Centric High Performance Computing Panel

Jonathan Beard

11 November 2018



Specialization drives performance...



More cores



Year

Energy of Data Movement

Cost of obtaining 8B from DRAM for a single DP FLOP



5



Locality is, from a certain point of view.

Observation: Manipulating distance by moving processor can expand locality and decrease latency for some workloads

L1

te X L2 LOng distance

L3

Core

*Storage Class Memory - NVM

PE

Short

distance

arm Research

SCM*

Data



Locality is, from a certain point of view.



*Storage Class Memory - NVM

Relativity

Bandwidth Scaling

Benchmark



Memory Access Latency

Processor to Memory speed (clock cycles) ratio

Latency relative to core





Scalable parallelism



There's no such thing as processing in memory

Challenge: breaking the semantic barrier





Interface standards

Challenge

• Locality-based targeting





Data Movement

We solved one problem, but created another...

Challenge

- Data layout transformation is critical
- Where to transform, how to program, and how best to virtualize are the biggest questions...



L1D Cache Used Line Utilized versus **Cache Line** Wasted Wasted GUPS: 80% *CoMD*: 50% *mcb*: 40% LULESH: 20% **DGEMM: 10%**

arm Research

L3

Keeping dark bandwidth coherent

Hmmm....maybe there's a better way.

- *Dark Bandwidth* blows up transfers, moving one cache line actually moves far more than that!
- With every move comes coherence traffic, often lots.
- Synchronization also takes time...



<u>Challenge</u>

- Improve performance / transparency of communications between all processing elements.
- Do we really need coherence? (likely not always)





Grand Challenges





1: Efficiency of data movement, logic is cheap, movement is expensive – future systems must capitalize on both data with reuse and streaming data, *Dark Bandwidth* must be avoided

2: Multiple drivers / compilers / software stacks are multi-million-dollar efforts for each vendor – will developers even adopt? – Reducing cost is a huge disruptor!

3: Communications / scalability of cores is not good with current coherence methods – but specialization and more cores are the future, Post-Moore.

4: Virtualization and translation for accelerators is an afterthought at the moment, extending the virtual memory model eases programming – but can we do more?

Thank You! Danke! Merci! 谢谢! ありがとう! **Gracias!** Kiitos!

CIM Research