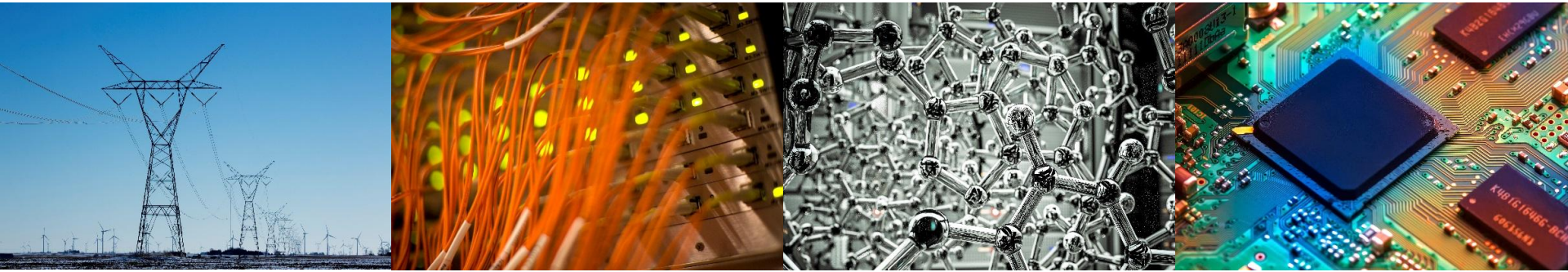


# Bigger GPUs and Bigger Nodes

Carl Pearson ([pearson@illinois.edu](mailto:pearson@illinois.edu))  
PhD Candidate, advised by Professor Wen-Mei Hwu



**I** ILLINOIS

Electrical & Computer Engineering

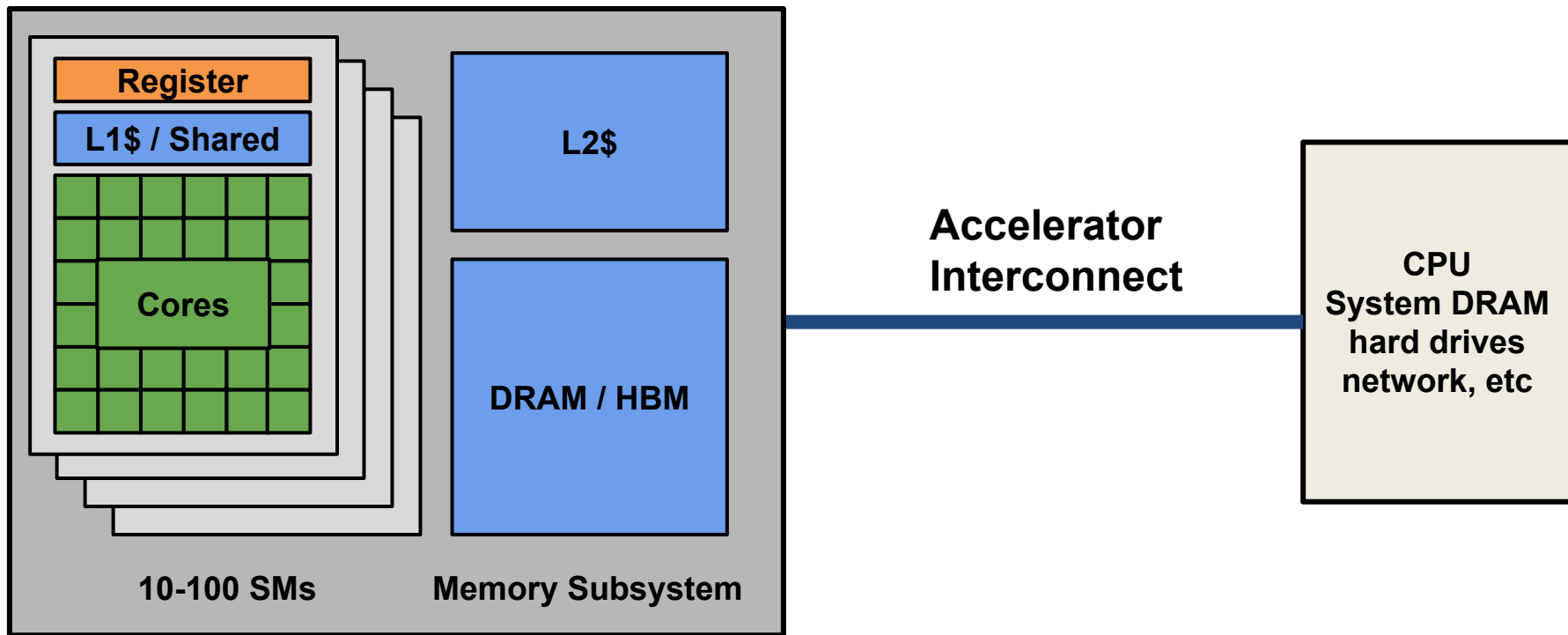
COLLEGE OF ENGINEERING

# Outline

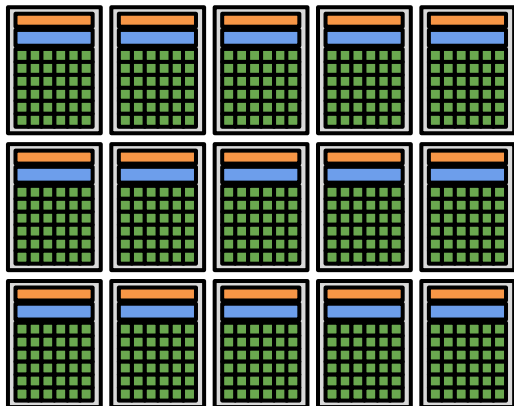
Experiences from working with domain experts to develop GPU codes on Blue Waters

- Kepler and Volta GPUs
- HPC Kepler to Volta Speedup
- Blue Waters, Summit, Sierra
- Intra-Node Communication Performance

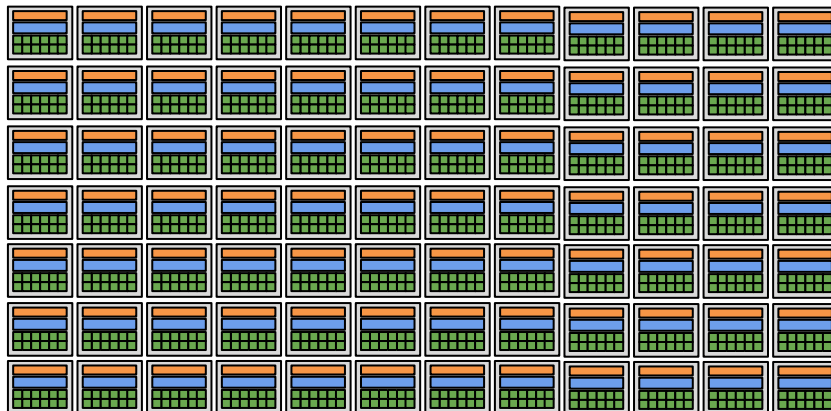
# GPU Architecture Bird's Eye View (not to scale)



# Kepler

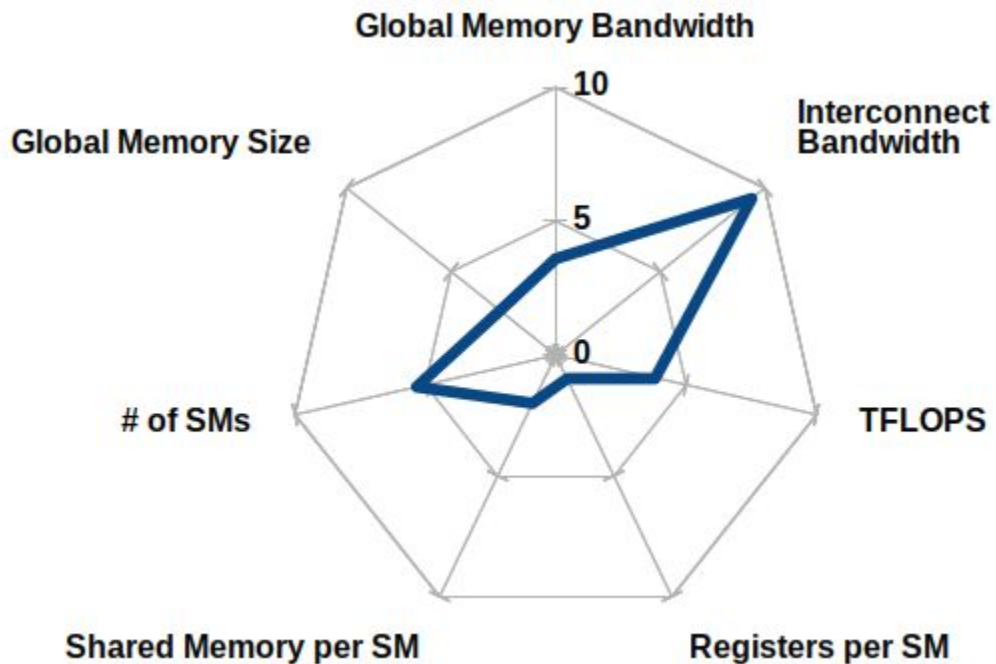


# Volta



	Number of SMs	Maximum Blocks / SM	Shared Memory / SM	Registers / SM	Single Precision Rate	Global Memory Bandwidth
<b>K20X (Kepler)</b>	15	16	48 KB	64 K	3.94 TFLOPS	250 GB/s
<b>V100 (Volta)</b>	80	32	96 KB	64 K	15 TFLOPS	900 GB/s

# K20x to V100: Architectural Parameters



# HPC Case Studies

## AWP-ODC

Tom Jordan, Yifeng Cui  
Southern California Earthquake Center  
University of Southern California

Anelastic Wave propagation

Solves a velocity-stress formulation of  
the 3D wave equation

## ChaNGa

Tom Quinn  
University of Washington

Charm N-body Gravity Solver

Collisionless N-body simulations

# AWP and ChaNGa V100 Speedup

	<b>Vs. P100</b>	<b>Vs. K20x (Blue Waters)</b>
<b>ChaNGa</b>	3.28	4.73
<b>AWP</b>	1.71	5.19

# AWP Detail

<b>SP over p100</b>	<b>SP over K20X</b>
1.711	5.188

	<b>K20x</b>		<b>V100</b>	
	<b>Kernel 1</b>	<b>Kernel 2</b>	<b>Kernel 1</b>	<b>Kernel 2</b>
<b>GPU Time</b>	72.4 %	27.5 %	70.1 %	29.3 %
<b>Mem BW</b>	145.7 GB/s	136.1 GB/s	726.7 GB/s	600.2 GB/s
	Latency-Limited		Bandwidth-Limited	



# AWP Optimizations

## Large Blocks to Capture Reuse

Reuse in fast memory

Blocks / SM limited by registers and SMs



## Uneven Architectural Change

Many more SMs

More memory per SM

Same registers per SM



## Unclear Tradeoff

Fine-grained parallelism: more work for GPU, less reuse

# Takeaways

## *Laissez-faire* Approach:

- 3-5x kernel speedup over optimized Kepler

- 3-5x interconnect speedup over optimized Kepler

- Larger problem to fill GPU

## Redesign/Rewrite Approach:

- Finer-grained parallelism to fill GPU

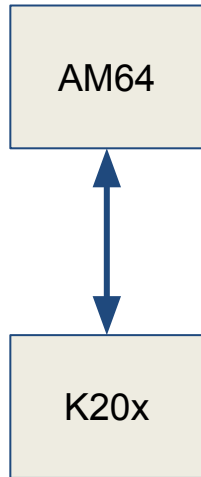
- Harder to capture reuse (key to performance)

# Nodes are Getting Bigger

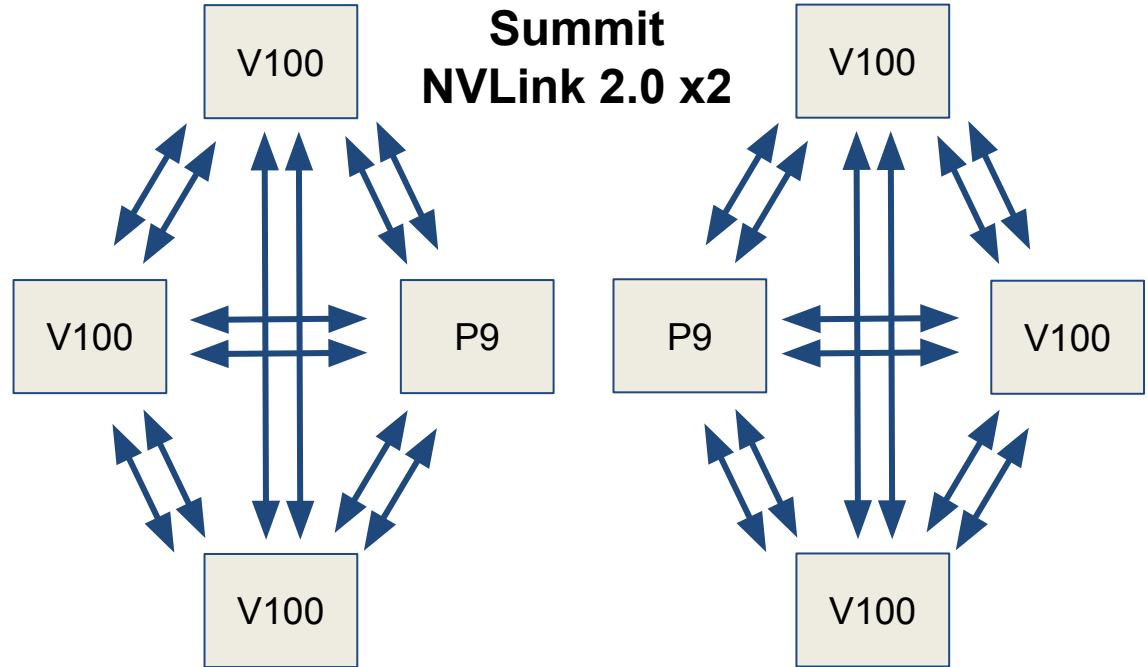
	BW	Summit <sup>1</sup> (ORNL)					
<b>CPU</b>	1x AMD64 32 threads 16 FP	POWER9 88 threads 22 FP			POWER9 88 threads 22 FP		
<b>GPU</b>	K20X 6 GB 4 TF	V100 16 GB 15 TF	V100 16 GB 15 TF	V100 16 GB 15 TF	V100 16 GB 15 TF	V100 16 GB 15 TF	V100 16 GB 15 TF
<b>Accelerator Interconnect (unidirectional)</b>	PCIe 2x16 8 GB/s	NVLink 2.0 x2 50 GB/s					
<b>Memory</b>	32GB	512 GB					

# Blue Waters XK and Summit Intra-Node Interconnects

**Blue Waters  
PCIe 2.0 x16**



**Summit  
NVLink 2.0 x2**



# System Performance Research

## CUDA

Microbench: <https://github.com/rai-project/microbench>

## Neural Networks

MLModelScope: <http://ml-arc-minsky.netlify.com/>

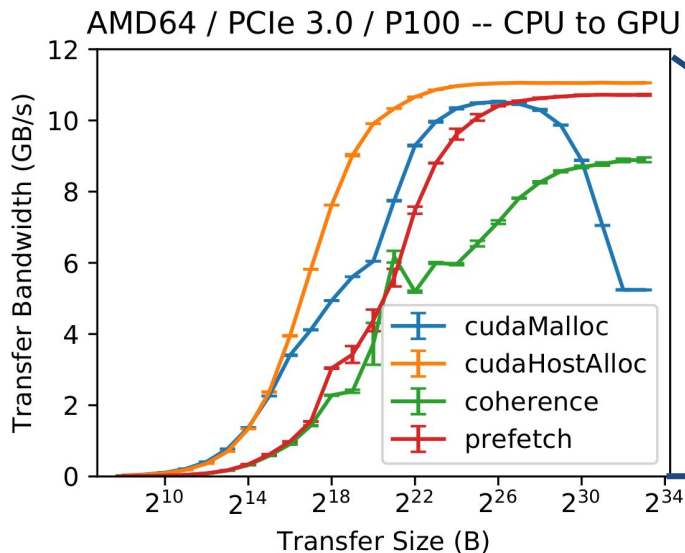
## Future Directions:

Quick application-driven architecture design

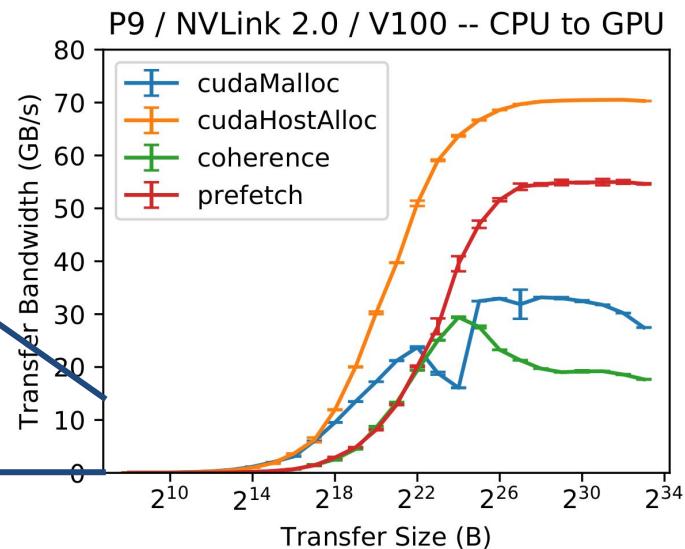
Performance modeling of neural networks

# Faster Interconnects

PCIe 3.0 x16 (2x BW)  
15.8 GB/s



NVLink 2.0 x3 (1.5x Summit)  
75 GB/s

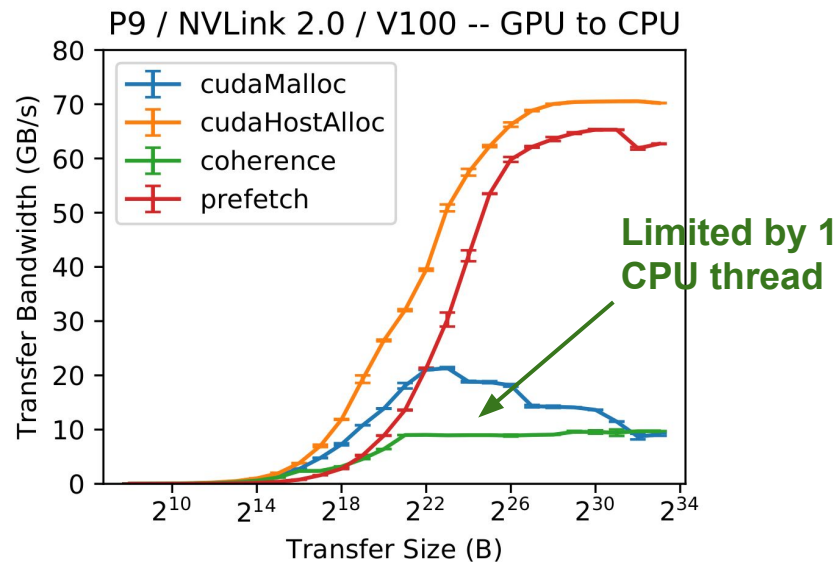
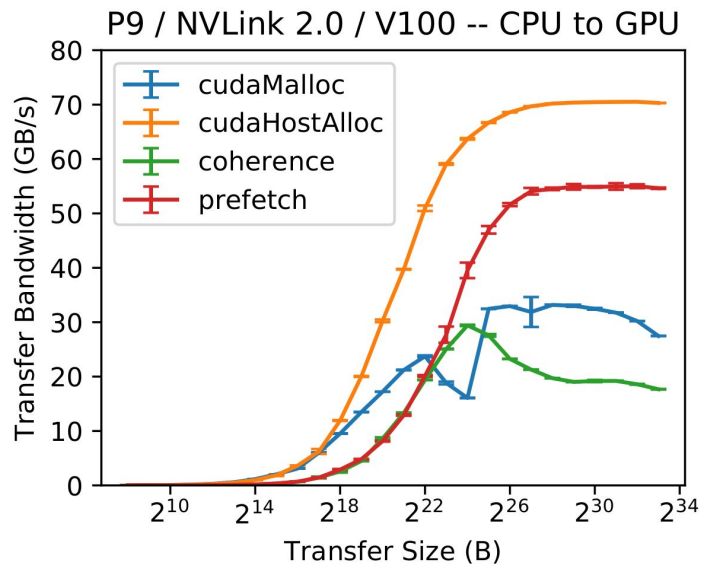


# Unified Memory

Allocations accessible from CPU and GPU  
Implicit data transfer (no cudaMemcpy)

	GPU 0	GPU 1	CPU	
<code>cudaSetDevice(0); cudaMallocManaged(&amp;a,...);</code>				
<code>a[page0] = 0; // gpu0</code>				
<code>a[page1] = 1; // gpu1</code>				Page fault and migration
<code>a[page2] = 2; // cpu</code>				Page fault and migration
<code>cudaMemAdvise(a, gpu1, cudaMemAdviseSetPreferredLocation); a[page1] = 1; // cpu</code>				Write served over NVLink
<code>cudaMemPrefetchAsync(a, gpu1);</code>				Bulk page migration

# P9 Unified Memory Performance

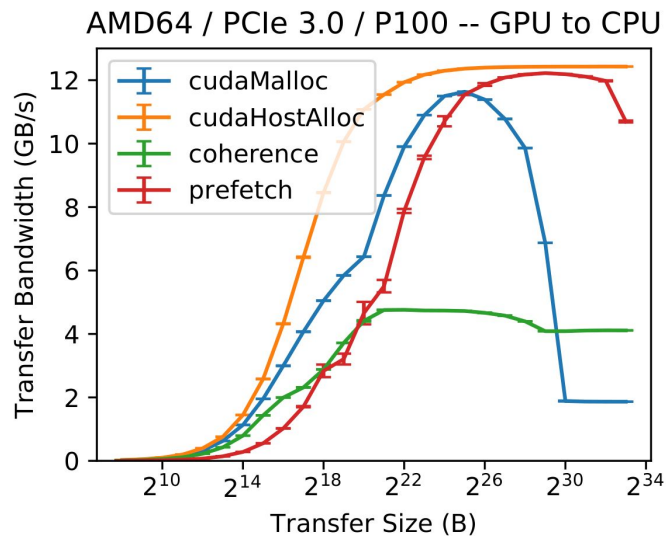
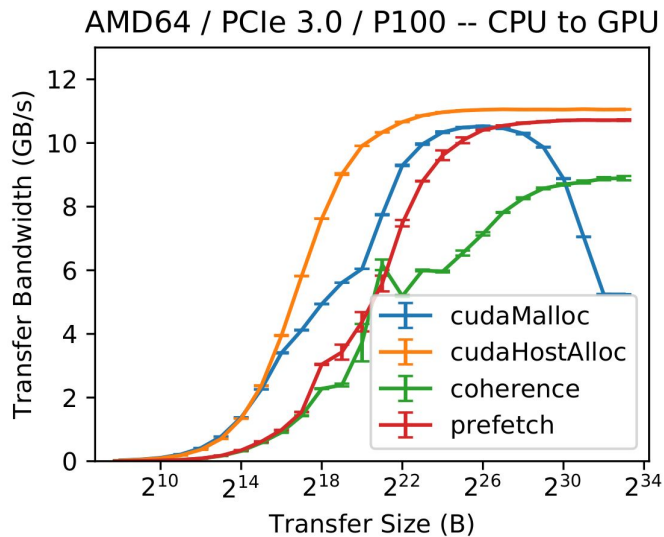


**Coherence: 30% of explicit management**

**Prefetch: 50-80% of explicit**



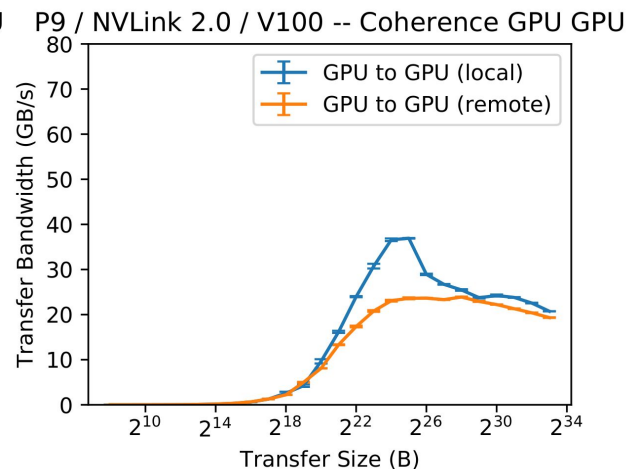
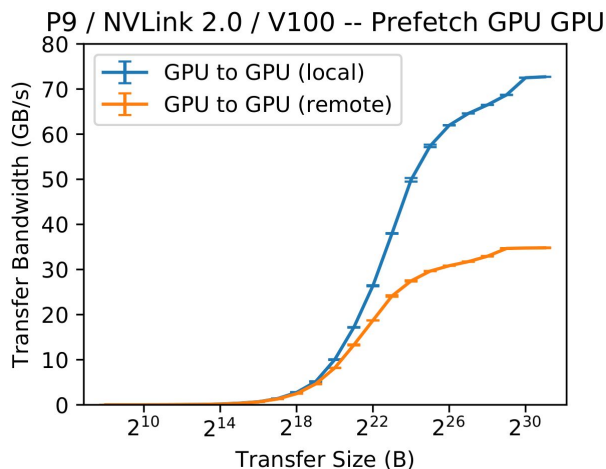
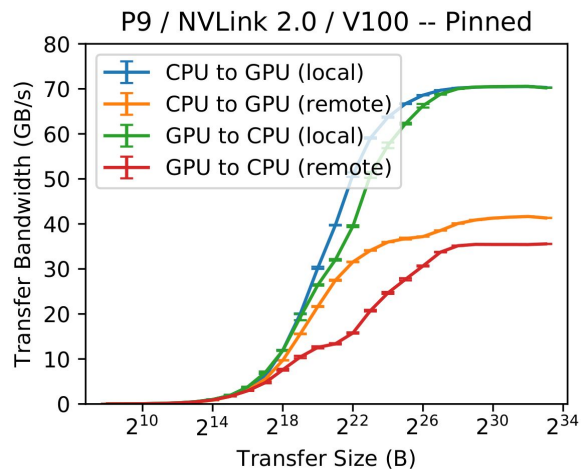
# AMD64 Unified Memory Performance



**Coherence: 30-70% of explicit management**

**Prefetch: 50-95% of explicit**

# Device Affinity



**Data placement on big nodes can have a dramatic communication impact**

# MLModelScope: Neural Network Performance Data

<http://ml-arc-minsky.netlify.com>

(model -- machine -- framework) triples

- ( AlexNet -- Jetson TX-1 -- Tensorflow )
- ( VGG19 -- AWS P2 X-large -- MxNet )

Neural-network performance primitive benchmarks

# Thank You

<https://cwpearson.github.io>  
[pearson@illinois.edu](mailto:pearson@illinois.edu)

## Special thanks to

- Professor Wen-Mei Hwu
- John Larson, Simon Garcia de Gonzalo, Zaid Qureshi, Mert Hidayetoglu, Abdul Dakkak and Cheng Li (University of Illinois)
- Isaac Gelado (NVIDIA)
- Jinjun Xiong and I-Hsin Chung (IBM)
- The IBM-ILLINOIS Center for Cognitive Computing Systems Research (C3SR) - a research collaboration as part of the IBM Cognitive Horizon Network.