# Scalable Resource Composition in a Flat World

Sudhakar Yalamanchili

Computer Architecture and Systems Laboratory
Center for Experimental Research in Computer Systems
School of Electrical and Computer Engineering
Georgia Institute of Technology

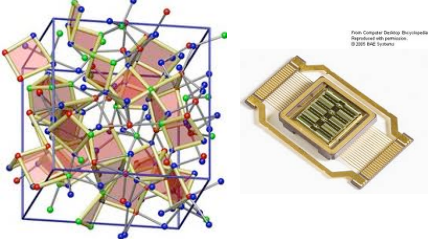*Sponsors: National Science Foundation, NVIDIA, LogicBlox Inc.*

SCHOOL OF ELECTRICAL AND COMPUTER ENGINEERING | GEORGIA INSTITUTE OF TECHNOLOGY

CASL

---

## System Diversity



Phase Change Memory

Amazon EC2 GPU Instances

*Technology Diversity is mainstream*

Cray Blue Waters

Photonics

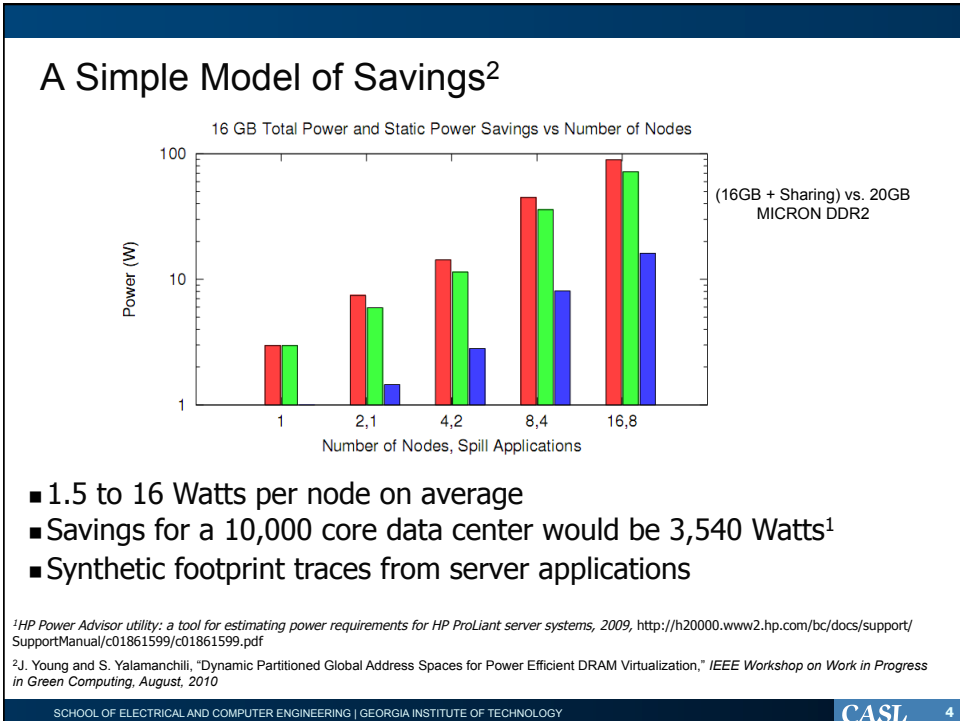SCHOOL OF ELECTRICAL AND COMPUTER ENGINEERING | GEORGIA INSTITUTE OF TECHNOLOGY

CASL   2

## Composing Memory

**Power vs DIMM Size (Kingston DDR3)**



- 32 GB LP
- 16 GB
- 16 GB LP
- 8 GB
- 8 GB LP
- 4 GB
- 4 GB LP

- While 8 GB DIMMs are most cost-effective, larger DIMM sizes are most power-efficient

*Resource Sharing*

**Cost vs. DIMM Size (Kingston DDR3)**



- 32 GB
- 16 GB
- 8 GB
- 4 GB
- 2 GB

*Statistics generated for Supermicro X9DRL-3F Motherboard (8 DIMM slots)*

SCHOOL OF ELECTRICAL AND COMPUTER ENGINEERING | GEORGIA INSTITUTE OF TECHNOLOGY          CASL   3

## A Simple Model of Savings[2]

16 GB Total Power and Static Power Savings vs Number of Nodes



(16GB + Sharing) vs. 20GB MICRON DDR2

- 1.5 to 16 Watts per node on average
- Savings for a 10,000 core data center would be 3,540 Watts[1]
- Synthetic footprint traces from server applications

[1]*HP Power Advisor utility: a tool for estimating power requirements for HP ProLiant server systems, 2009,* http://h20000.www2.hp.com/bc/docs/support/SupportManual/c01861599/c01861599.pdf

[2]*J. Young and S. Yalamanchili, "Dynamic Partitioned Global Address Spaces for Power Efficient DRAM Virtualization," IEEE Workshop on Work in Progress in Green Computing, August, 2010*

SCHOOL OF ELECTRICAL AND COMPUTER ENGINEERING | GEORGIA INSTITUTE OF TECHNOLOGY          CASL   4
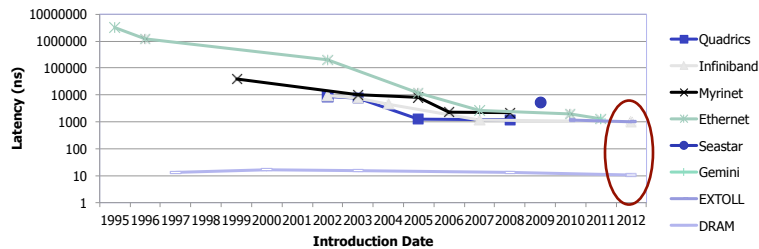
## Bandwidth Trends

**Bandwidth vs. Time for Common Interconnects**



■DRAM to interconnect bandwidth ratio has been steadily dropping

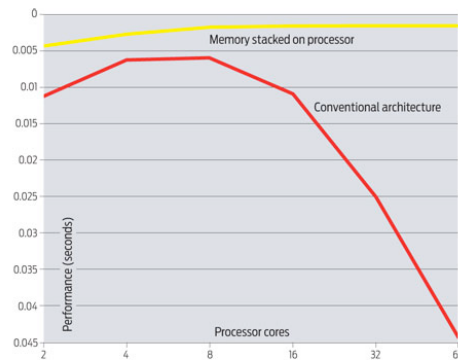**CASL** 5

## Latency Trends

**MPI Ping Latency vs. Time for Common Interconnects**



■MPI latency has steadily approached DRAM read latency
  ■ Hardware switching times in the low hundreds of nanoseconds.
■Note progress in photonics

*Extend the reach of a socket*

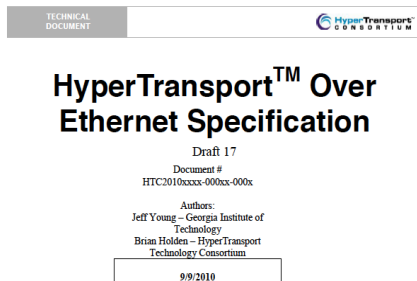**CASL** 6

## The Memory Wall



*"Multicore Is Bad News For Supercomputers"*
IEEE Spectrum 2008

- Data intensive applications

- Memory bandwidth demand is scaling faster than memory interface capacity

*"You can buy bandwidth but you cannot bribe God"*
*- unknown*

*Convert Network Bandwidth into Memory Bandwidth*

SCHOOL OF ELECTRICAL AND COMPUTER ENGINEERING | GEORGIA INSTITUTE OF TECHNOLOGY                CASL    7

---

## Impact on Clustering

TECHNICAL
DOCUMENT                          HyperTransport
                                  CONSORTIUM

**HyperTransport™ Over Ethernet Specification**

Draft 17
Document #
HTC2010xxxx-000xx-000x

Authors:
Jeff Young – Georgia Institute of Technology
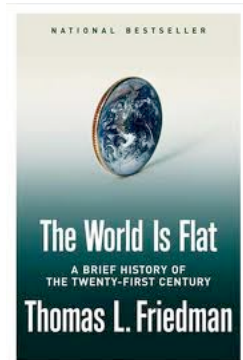Brian Holden – HyperTransport Technology Consortium

9/9/2010

- Combine commodity interconnects and memory systems

- Need flexible hardware level composition of resources

- This is an old idea whose time has come?

Some Examples

- Lim, et al. - Memory Blades for disaggregated memory
- Tolentino, Cameron – Memory Miser OS level support
- Lefurgy, et al. – DRAM server power and DRAM consolidation
- RDMA - Liang '05 low-level implementation for page swapping
- Memscale – UoH, UPC
- Feng et. Al – Green Supercomputing

SCHOOL OF ELECTRICAL AND COMPUTER ENGINEERING | GEORGIA INSTITUTE OF TECHNOLOGY                CASL    8
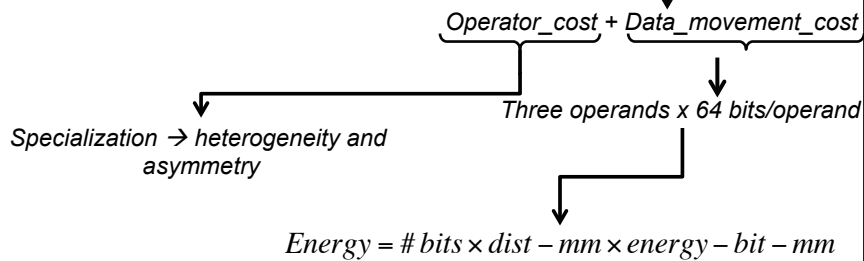
## Flattening Cluster Hierarchies

*Observation I:*
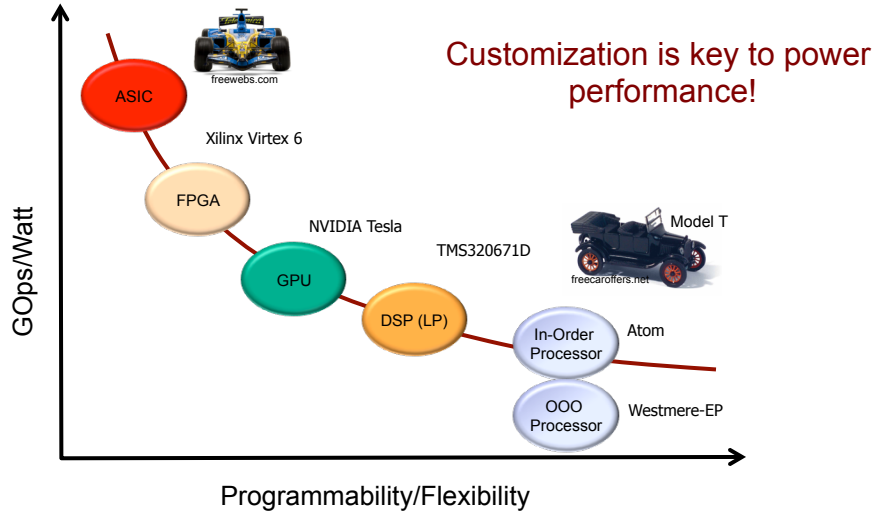*Everyone is getting closer and we*
*need better sharing but…..*

## Post Dennard Performance Scaling

$$Perf\left(\frac{ops}{s}\right) = Power(W) \times Efficiency\left(\frac{ops}{joule}\right)$$

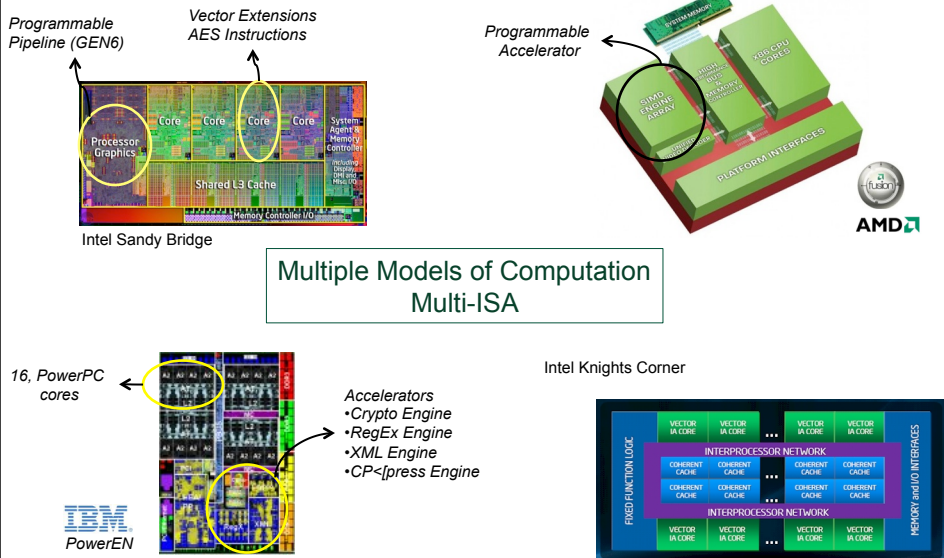Dally, Keynote IITC 2012

$Operator\_cost$ + $Data\_movement\_cost$

*Specialization → heterogeneity and asymmetry*

*Three operands x 64 bits/operand*

$$Energy = \# bits \times dist - mm \times energy - bit - mm$$

## Hardware Power-Performance Tradeoffs



Customization is key to power performance!

GOps/Watt

ASIC

Xilinx Virtex 6

FPGA

NVIDIA Tesla

GPU

TMS320671D

DSP (LP)

Model T

In-Order Processor — Atom

OOO Processor — Westmere-EP

Programmability/Flexibility

## Consolidation on Chip

Programmable Pipeline (GEN6)

Vector Extensions AES Instructions



Intel Sandy Bridge

Programmable Accelerator



AMD

Multiple Models of Computation
Multi-ISA

16, PowerPC cores



IBM
PowerEN

Accelerators
•Crypto Engine
•RegEx Engine
•XML Engine
•CP<[press Engine

Intel Knights Corner

## Consolidation in a System

perisofparallel.blogspot.com

*So its not just memory that needs to be shared!*

NVIDIA Tesla

**CASL**  13

---

## Post Dennard Performance Scaling

$$Perf\left(\frac{ops}{s}\right) = Power(W) \times Efficiency\left(\frac{ops}{joule}\right)$$
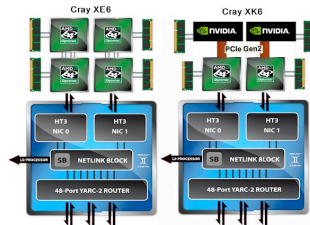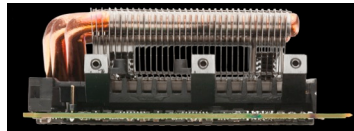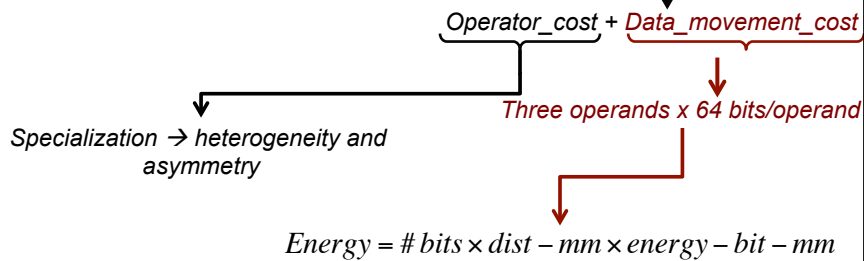
Dally, Keynote IITC 2012

*Operator_cost* + *Data_movement_cost*

*Specialization → heterogeneity and asymmetry*

*Three operands x 64 bits/operand*

$$Energy = \# bits \times dist - mm \times energy - bit - mm$$

**CASL**  14

## Scaling: Key Driver is Energy/Power

Embedded Platforms



Cost of Data Movement

Big Science: To Exascale
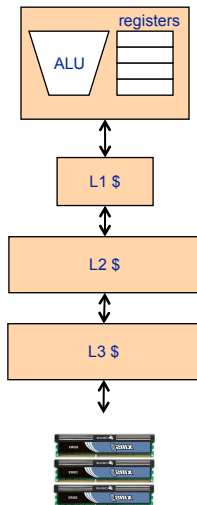


Goal: 1-100 GOps

Goal: 20MW/Exaflop



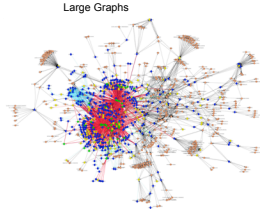*Courtesy: Sandia National Labs :R. Murphy).*

- Sustain performance scaling through massive concurrency
  - New execution models
- Data movement becomes more expensive than computation

## Optimizing Locality



*Observation II:*
*You can hide latency, but you cannot hide energy!*

## A Data Rich World

Large Graphs

*Mixed Modalities and levels of parallelism*
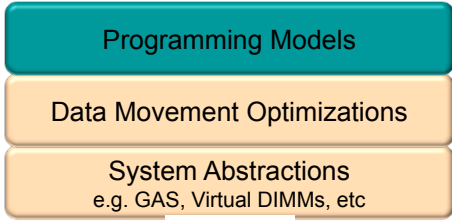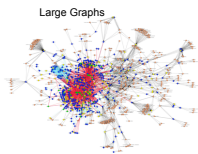
*Irregular, Unstructured Computations and Data*

Pharma

Trend analysis

*Images from math.nist.gov, blog.thefuturescompany.com, mgmtsocpinier.blogspot.com*

Waterexchange.com

conventioninsider.com

## System Model

Large Graphs

| Programming Models | *Domain Specific Languages* |
| --- | --- |
| Data Movement Optimizations | *Compiler and Run-Time Support* |
| System Abstractions e.g. GAS, Virtual DIMMs, etc | *Cluster Wide Hardware Consolidation* |

*Hardware Customization*
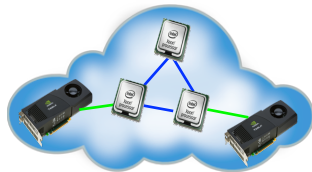
## Application: Data Warehousing

■ On-line and off-line analysis
  ■ Retail analysis
  ■ Forecasting
  ■ Pricing
  ■ ......

■ Combination of data queries and computational kernels

■ Current applications process 1 to 50 TBs of data [1]

■ Potential to change a companies business model!

[1] Independent Oracle Users Group. A New Dimension to Data Warehousing: 2011 IOUG Data Warehousing Survey.

## Databases: *Not* a Traditional Domain of GPUs

LargeQty(p) <-
Qty(q),
q > 1000.
......

**Relational Computations Over Massive Data Sets**

## Database Applications on GPUs

- The good
  - Lots of potential data parallelism
  - If data fits in GPU mem, 2x—27x speedup has been shown
- The bad
  - Very large data set (will not even fit in host memory)
  - I/O bound (GPU has no disk)
  - PCI data transfer takes 15–90% of the total time[*]
- The Ugly
  - Irregular/unstructured accesses to data

| Order | Price | Discount |
|-------|-------|----------|
| 0 | 10 | 10% |
| 1 | 20 | 20% |
| 2 | 10 | 15% |
| 3 | 51 | 14% |
| 4 | 33 | 13% |
| 5 | 22 | 10% |
| ...... | ...... | ...... |

[*] B. He, M. Lu, K. Yang, R. Fang, N. K. Govindaraju, Q. Luo, and P. V. Sander. Relational query co-processing on graphics processors. In TODS, 2009.
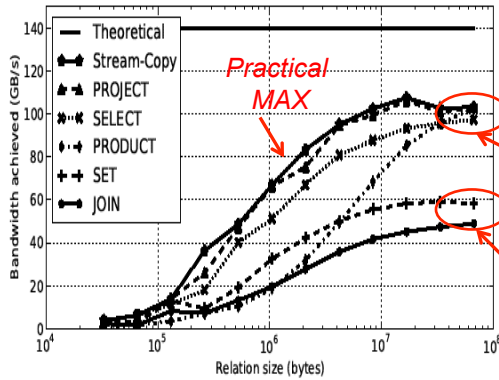
CASL  21

## Research Thrusts

- I: Optimized implementations of primitives
  - Relational algebra
  - Data management within the GPU memory hierarchy

- II: In-core processing
  - Cluster wide memory aggregation techniques
  - Change the ratio of host memory size to accelerator memory size

- III: Data movement optimizations
  - Between hosts and (local or remote) accelerators
  - Within an accelerator

CASL  22

## I. Relational Algebra Primitives on GPUs
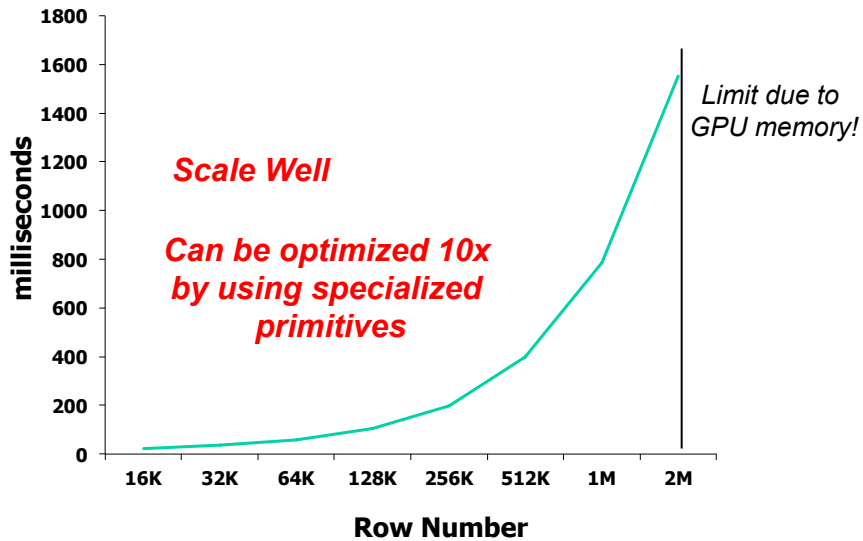
**Raw Performance (C2050)**

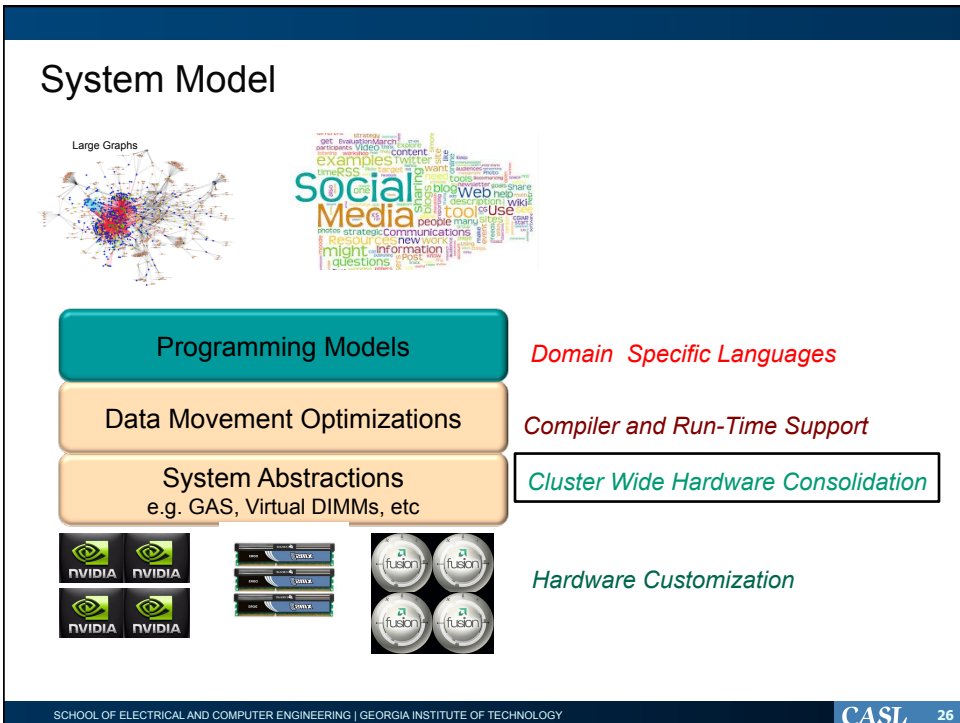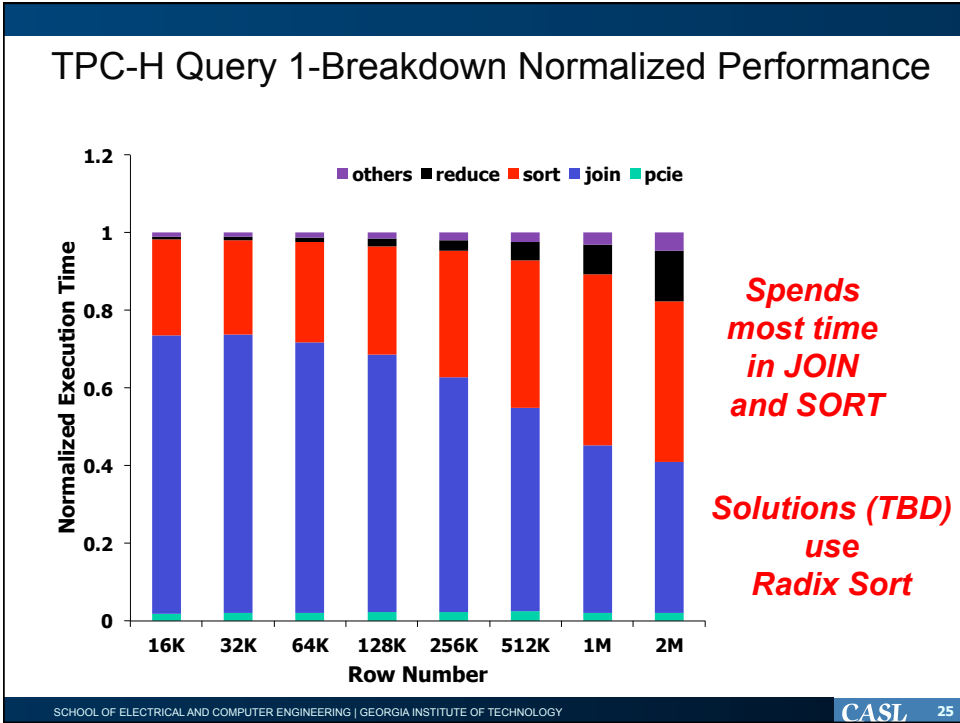*Fastest in GPU*



- Multi-stage algorithm (under review)

- Push to memory-bound

- Simple primitives are close to maximum performance

- Improved primitives under development

## TPC-H Query 1-Overall Performance



*Limit due to GPU memory!*

*Scale Well*

*Can be optimized 10x by using specialized primitives*

## TPC-H Query 1-Breakdown Normalized Performance



**Spends most time in JOIN and SORT**

**Solutions (TBD) use Radix Sort**

## System Model



Large Graphs

**Programming Models** — *Domain Specific Languages*

**Data Movement Optimizations** — *Compiler and Run-Time Support*

**System Abstractions** e.g. GAS, Virtual DIMMs, etc — *Cluster Wide Hardware Consolidation*
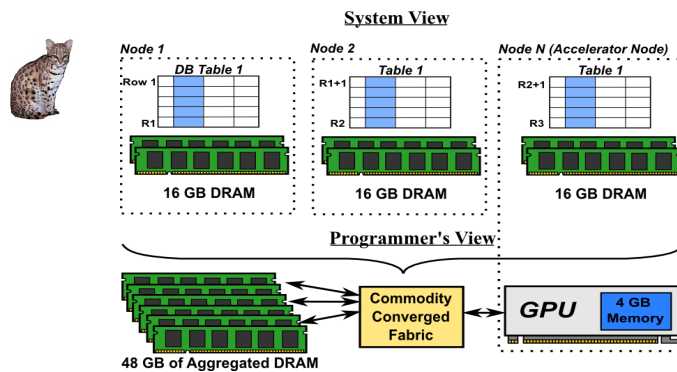
*Hardware Customization*

## II. In-Core Processing



- Cluster-based memory aggregation
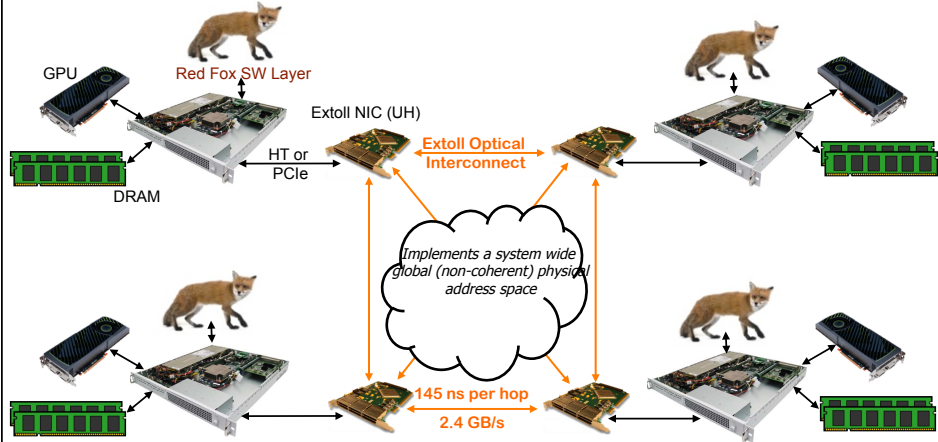- Hardware support for global non-coherent, physical address space system
- Change the ratio of host-memory : GPU-memory
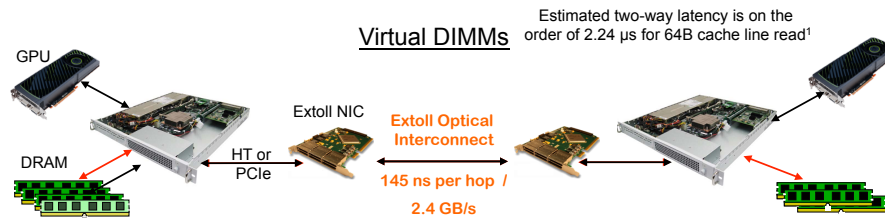
## Oncilla: Fabrics for Accelerator Clouds



- Goal: Efficient memory aggregation for accelerators in data centers
- Solution: Use Global Address Spaces (GAS) and commodity fabrics (HT, QPI, PCIe, 10GE, IB)
  - Support in-core databases using software from Red Fox project

# Oncilla Infrastructure

GPU

Red Fox SW Layer

Extoll NIC (UH)

HT or PCIe

DRAM

Extoll Optical Interconnect

*Implements a system wide global (non-coherent) physical address space*
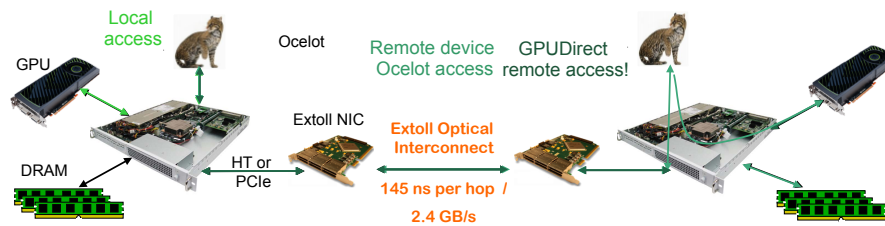
145 ns per hop

2.4 GB/s

- Low-latency, commodity hardware (Extoll) for efficient memory and GPU aggregation and Red Fox SW layer supports DB queries on remote nodes
- Collaboration with University of Heidelberg (UH), Polytechnic University of Valencia, AIC Inc., LogicBlox Inc.

CASL  29

---

# Some Candidate Systems Concepts

Estimated two-way latency is on the order of 2.24 µs for 64B cache line read[1]

Virtual DIMMs

GPU

Extoll NIC

Extoll Optical Interconnect

DRAM

HT or PCIe

145 ns per hop /

2.4 GB/s
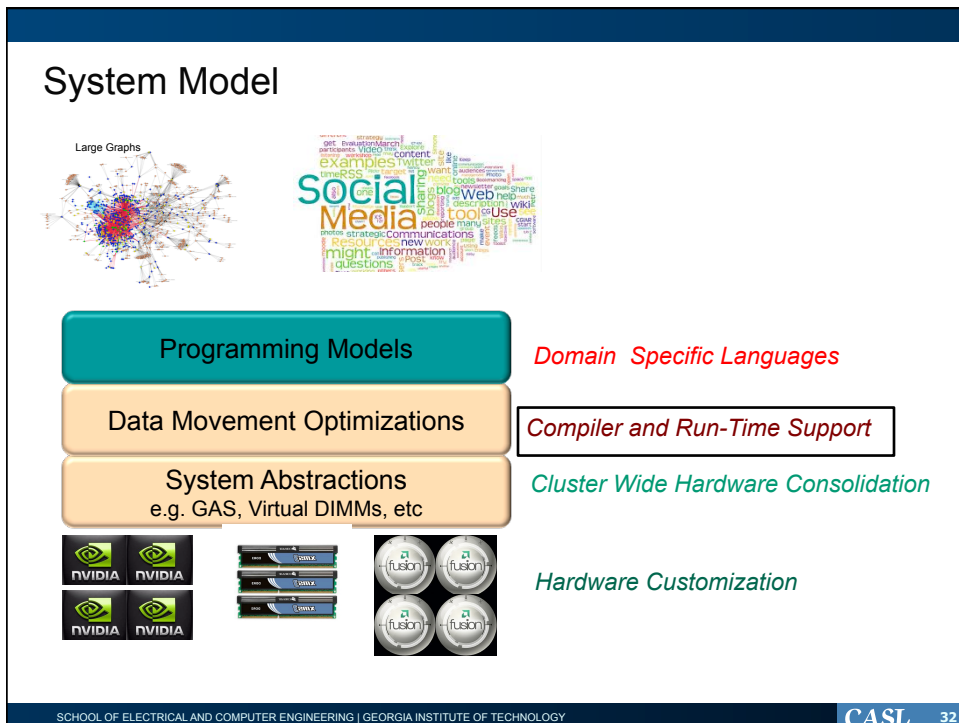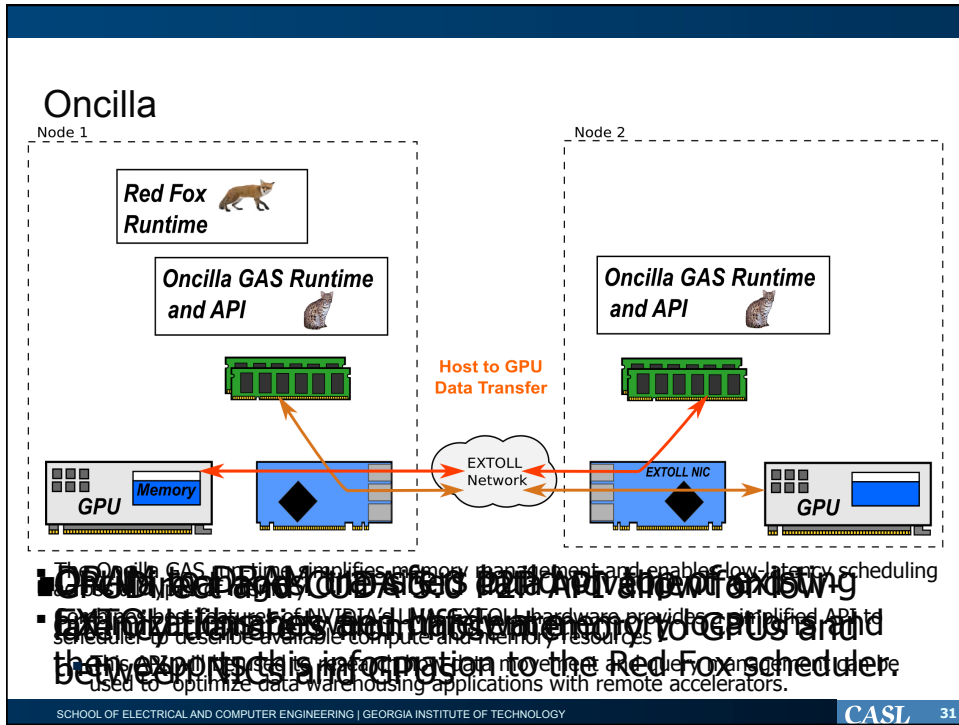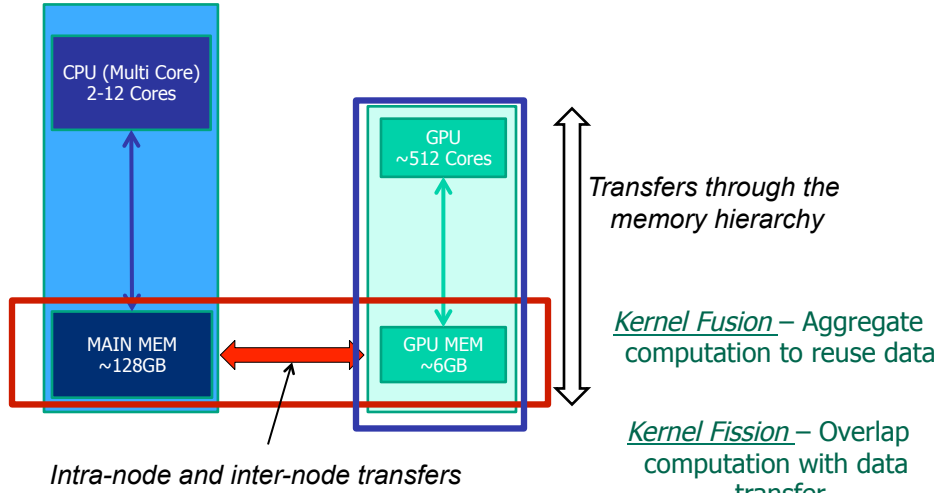
1) Young, J., Yalamanchili, S., Dynamic Partitioned Global Address Spaces for Power-Efficient DRAM Virtualization, WIPGC at IGCC, 2010
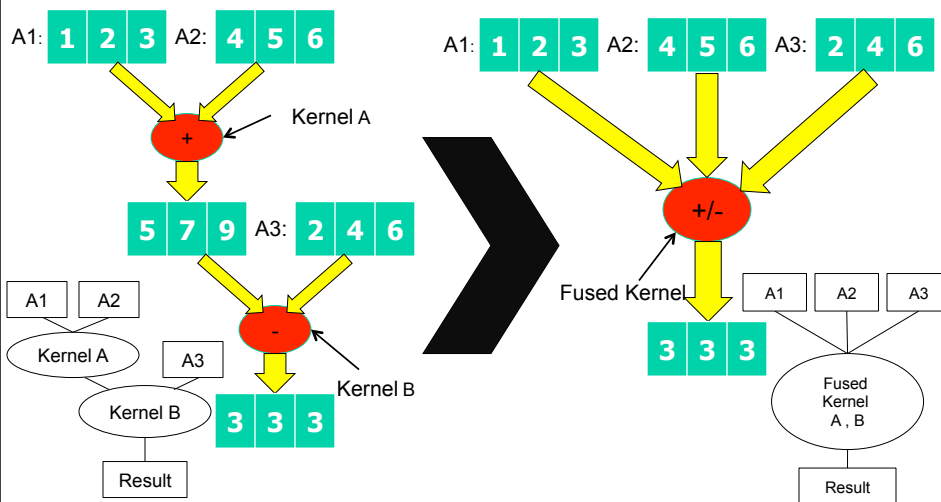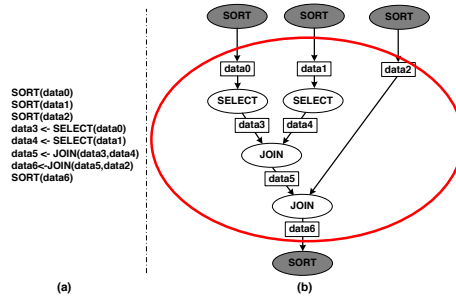
Remote GPU Access

Local access

Ocelot

Remote device Ocelot access

GPUDirect remote access!

GPU

Extoll NIC

Extoll Optical Interconnect

DRAM

HT or PCIe

145 ns per hop /

2.4 GB/s

CASL  30

# Oncilla

Node 1 | Node 2

**Red Fox Runtime**

**Oncilla GAS Runtime and API**

**Oncilla GAS Runtime and API**

**Host to GPU Data Transfer**

GPU | *Memory*

EXTOLL Network

*EXTOLL NIC*

GPU

- The Oncilla GAS runtime simplifies memory management and enables low-latency scheduling
- EXTOLL hardware provides a simplified API to describe available compute and memory resources
- them exporting this information to the Red Fox scheduler.

*(overlapping text: CPU workloads can use both API and workflow flexibility between host and remote memory together with CUDA and CPU host features of NVIDIA's UNAGE EXTOLL hardware provides a simplified API to... used to optimize data warehousing applications with remote accelerators.)*

---

# System Model

Large Graphs

| Programming Models | *Domain Specific Languages* |
| Data Movement Optimizations | *Compiler and Run-Time Support* |
| System Abstractions<br>e.g. GAS, Virtual DIMMs, etc | *Cluster Wide Hardware Consolidation* |
| | *Hardware Customization* |

III. Data Movement Optimizations

CPU (Multi Core)
2-12 Cores

GPU
~512 Cores

*Transfers through the memory hierarchy*

MAIN MEM
~128GB

GPU MEM
~6GB

*Kernel Fusion* – Aggregate computation to reuse data

*Kernel Fission* – Overlap computation with data transfer

*Intra-node and inter-node transfers*

SCHOOL OF ELECTRICAL AND COMPUTER ENGINEERING | GEORGIA INSTITUTE OF TECHNOLOGY

CASL    33



Kernel Fusion

A1: 1 2 3    A2: 4 5 6

Kernel A

+

5 7 9    A3: 2 4 6

A1    A2

Kernel A    A3

Kernel B

- 

3 3 3

Kernel B

Result

A1: 1 2 3    A2: 4 5 6    A3: 2 4 6

+/-

Fused Kernel

3 3 3

A1    A2    A3

Fused Kernel A , B

Result

SCHOOL OF ELECTRICAL AND COMPUTER ENGINEERING | GEORGIA INSTITUTE OF TECHNOLOGY
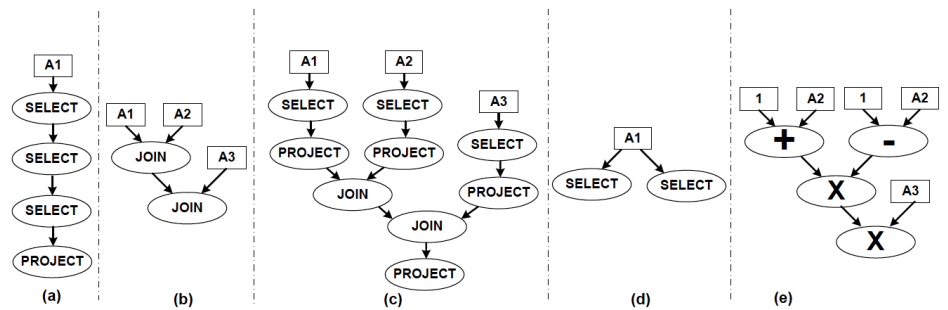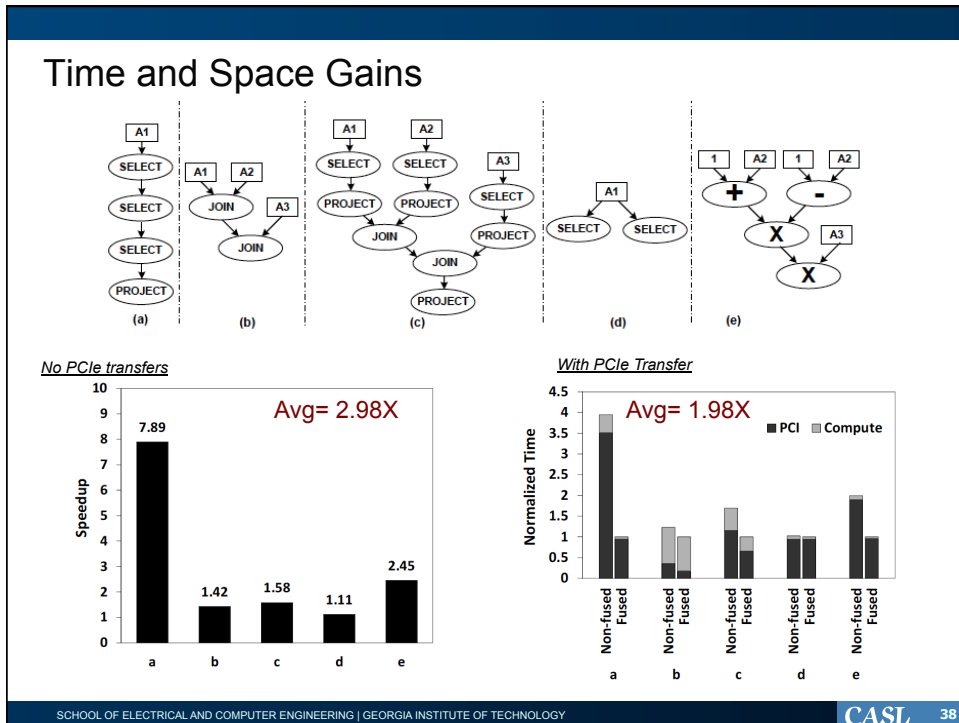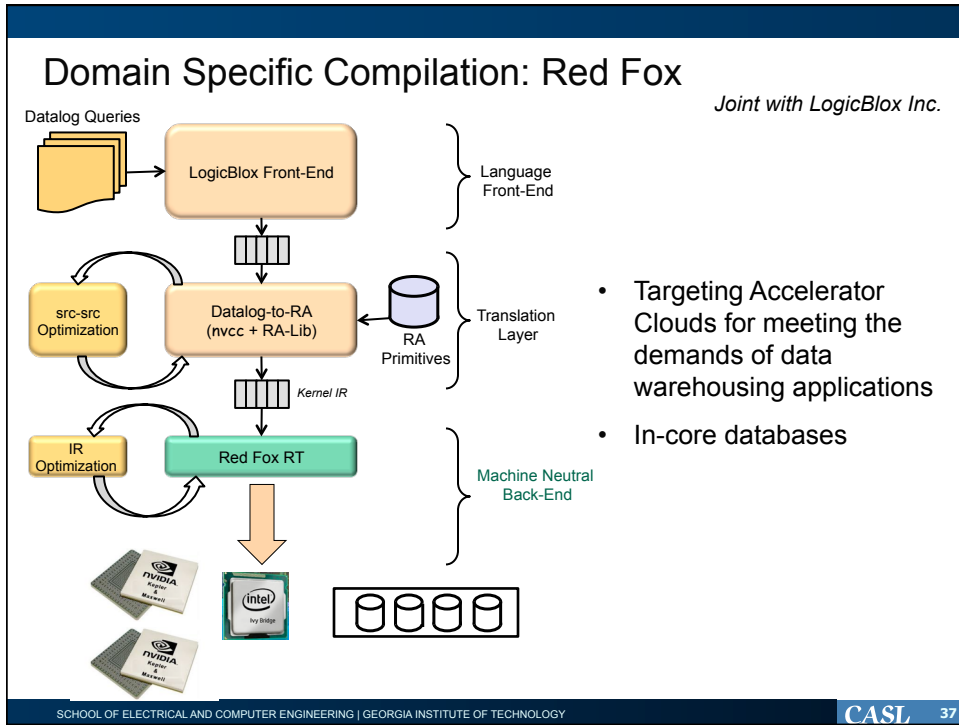
CASL    34

## Kernel Fusion Benefits

- Smaller Data Footprint
  - Reduction in Memory Accesses
  - Temporal Data Locality
  - Reduction in Traffic
  - Larger Input Data
- Larger Optimization Scope
  - Common Computation Elimination
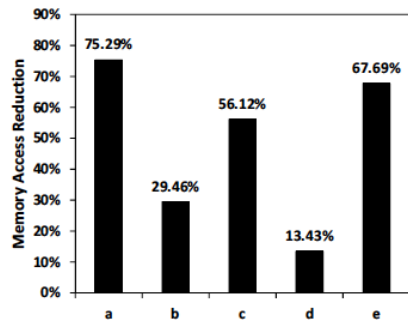  - Improved Compiler Optimization Benefits

SORT(data0)
SORT(data1)
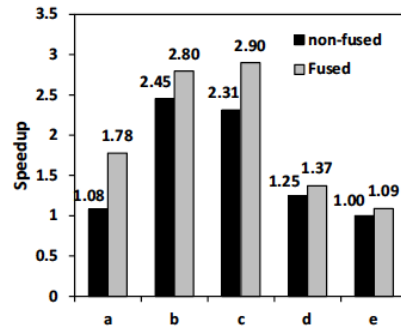SORT(data2)
data3 <- SELECT(data0)
data4 <- SELECT(data1)
data5 <- JOIN(data3,data4)
data6<-JOIN(data5,data2)
SORT(data6)

(a)    (b)

## Common RA Combinations of TPC-H

(a)    (b)    (c)    (d)    (e)

# Domain Specific Compilation: Red Fox

*Joint with LogicBlox Inc.*

Datalog Queries

LogicBlox Front-End

Language Front-End

src-src Optimization

Datalog-to-RA (nvcc + RA-Lib)

RA Primitives

Translation Layer

*Kernel IR*

IR Optimization

Red Fox RT

Machine Neutral Back-End

- Targeting Accelerator Clouds for meeting the demands of data warehousing applications
- In-core databases

CASL 37

# Time and Space Gains



No PCIe transfers

Avg= 2.98X

| | a | b | c | d | e |
|---|---|---|---|---|---|
| Speedup | 7.89 | 1.42 | 1.58 | 1.11 | 2.45 |

With PCIe Transfer

Avg= 1.98X

PCI ■ Compute

CASL 38

19

## Memory and Optimization Scope

**Memory Accesses**

**Impact of Optimization Scope – O3 vs. O0**
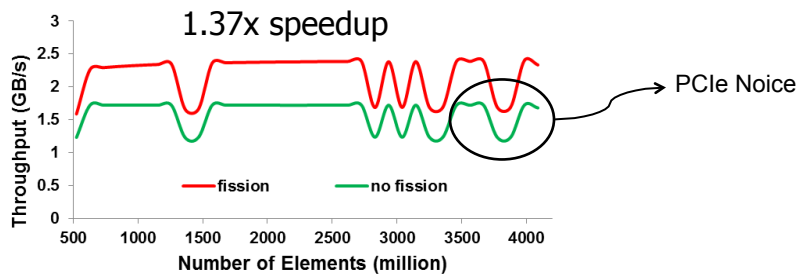
## TPC-H Queries

**Query 1**

**Query 21**



Avg= 1.25X (3.18X w/o SORT and PCIe )

Avg= 1.22X

## Example of Kernel Fission

**GPU MEM**

| | | |
|---|---|---|
| CTA0 | GPU->CPU | CPU->GPU |
| CTA1 | GPU Computation | GPU->CPU |
| CTA2 | CPU->GPU | GPU Computation |
| | **Cycle 0** | **Cycle 1** |

Reminiscent of Software Pipelining

**1.37x speedup**



PCIe Noice

— fission  — no fission

Throughput (GB/s) vs Number of Elements (million)

## People



**Gregory Diamos**
Dynamic optimizations (Harmony, Ocelot, LLVM Bridge)



**Andrew Kerr**
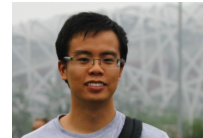Program Transformations & Optimizations for Data Parallel Computation (Ocelot, LLVM Bridge, VSIPL)



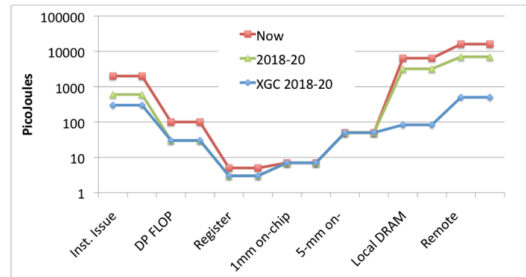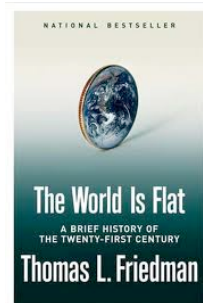**Haicheng Wu**
Dynamic Opytimizations (Ocelot)



**Jeff Young**
Integrated Networks-Memory, Oncilla



**Si Li**
Correctness & Emulation Tools, GP architectures

## Summary



- Refactor cluster architectures for
- Flexible hardware composition of resources and
- Migrate to a communication-centric model of algorithms, systems, and optimizations

CASL   43