

# Programming Systems for Specialized Architectures

## *Interface, Data, Approximation*

***Sarita Adve***

***With: Vikram Adve, Johnathan Alsop, Maria Kotsifakou, Sasa Misailovic,  
Matt Sinclair, Prakaalp Srivastava***

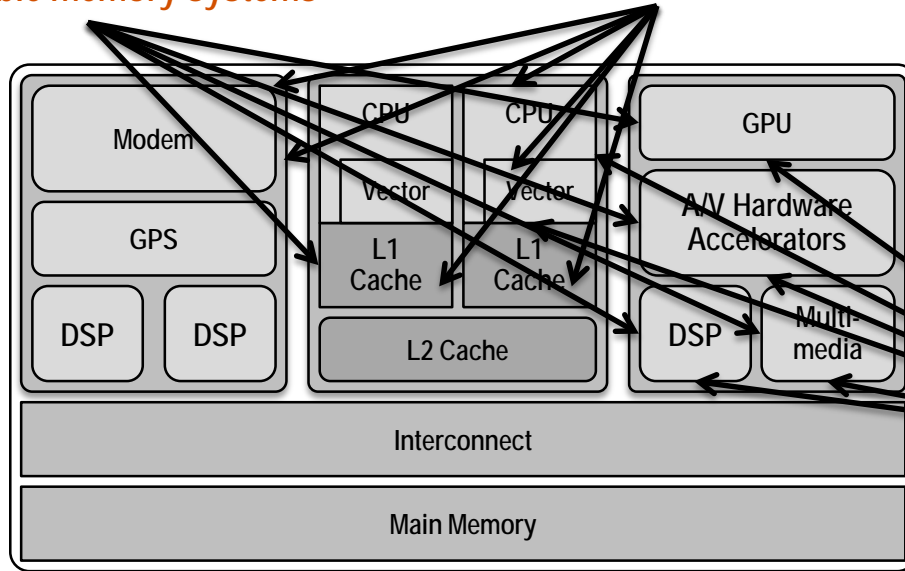
**University of Illinois at Urbana-Champaign  
sadve@illinois.edu**

***Sponsors: NSF, C-FAR, ADA (JUMP center by SRC, DARPA)***

# A Modern Mobile SoC

Incompatible memory systems

Different hardware ISAs



Increasing diversity in & across SoCs  
& supercomputers, data centers, ...

Different parallelism models

Need common **interface** (abstractions): HW-independent SW development, "object code" portability

**Data** movement critical: Memory structures, communication, consistency, synchronization

**Approximation:** Application-driven solution quality trade off to increase efficiency

# Interfaces: Back to the Future

April 7, 1964: IBM announced the 360

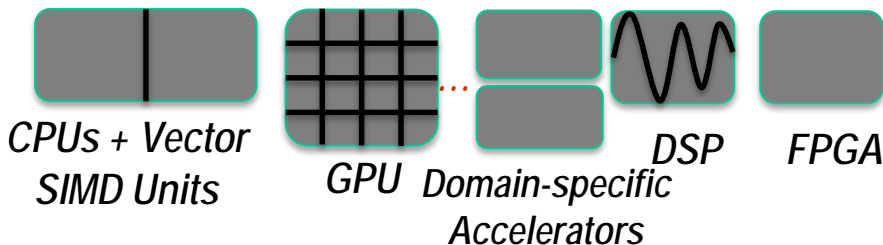
- Family of machines w/ common abstraction/interface/ISA
  - Programmer freedom: no reprogramming
  - Designer freedom: implementation creativity

Not unique

- CPUs : ISAs; Internet : IP; GPUs : CUDA; Databases : SQL; ...

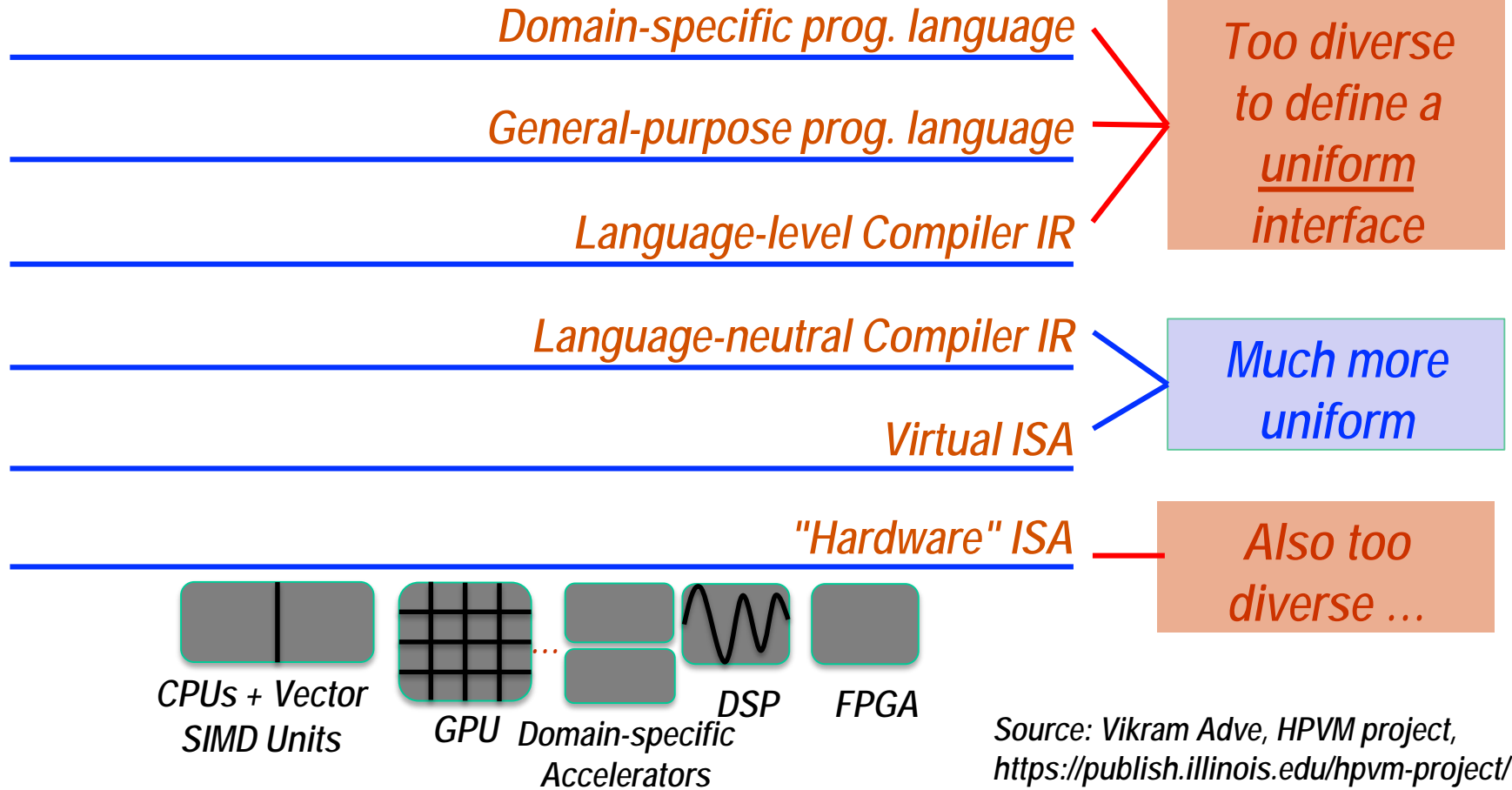
# Current Interface Levels

<u>App. productivity</u>	<u>Domain-specific prog. language</u>	TensorFlow, MXNet, Halide, ...
<u>App. performance</u>	<u>General-purpose prog. language</u>	CUDA, OpenCL, OpenAcc, OpenMP, Python, Julia
<u>Language innovation</u>	<u>Language-level Compiler IR</u>	Delite DSL IR, DLVM, TVM, ...
<u>Compiler investment</u>	<u>Language-neutral Compiler IR</u>	Delite IR, HPVM, OSCAR, Polly
<u>Object-code portability</u>	<u>Virtual ISA</u>	SPIR, HPVM
<u>Hardware innovation</u>	<u>"Hardware" ISA</u>	IBM AS/400, Transmeta, PTX, HSAIL, Codesigned Virtual Machines



Source: Vikram Adve, HPVM project, <https://publish.illinois.edu/hpvm-project/>

# Which Interface Levels Can Be Uniform?



Source: Vikram Adve, HPVM project, <https://publish.illinois.edu/hpvm-project/>

# One Example

## HPVM: Heterogeneous Parallel Virtual Machine [PPoPP'18]

Parallel program representation for heterogeneous parallel hardware

- **Virtual ISA:** portable virtual object code, simpler translators
- **Compiler IR:** optimizations, map diverse parallel languages
- **Runtime Representation** for flexible scheduling: mapping, load balancing

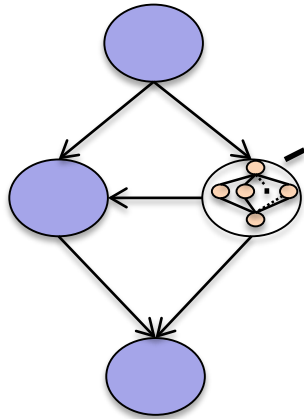
Generalization of LLVM IR for parallel heterogeneous hardware

PPoPP'18: Results on GPU (Nvidia), Vector ISA (AVX), Multicore (Intel Xeon)

Ongoing: FPGA, novel domain-specific SoCs

# HPVM Abstractions

*Hierarchical Dataflow Graph  
with side effects*

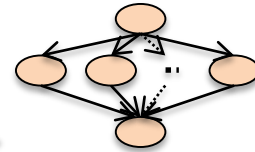


+

*Vector*

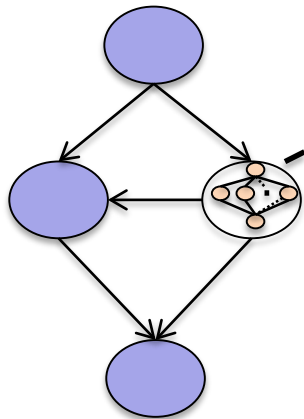
$V_A = \text{load } \langle L4 \times \text{float} \rangle^* A$   
 $V_B = \text{load } \langle L4 \times \text{float} \rangle^* B$   
...  
 $V_C = \text{fmul } \langle L4 \times \text{float} \rangle V_A, V_B$

*or*



# HPVM Abstractions

*Hierarchical Dataflow Graph  
with side effects*

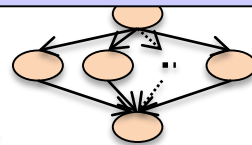


+

*Vector*

$V_A = \text{load } \langle L4 \times \text{float} \rangle^* A$   
 $V_B = \text{load } \langle L4 \times \text{float} \rangle^* B$   
...

- *Task, data, vector parallelism*
- *Streams, pipelines*
- *Shared memory*
- *High-level optimizations*
- *FPGAs (more custom hw?)*



*N different parallelism models* → *single unified model*

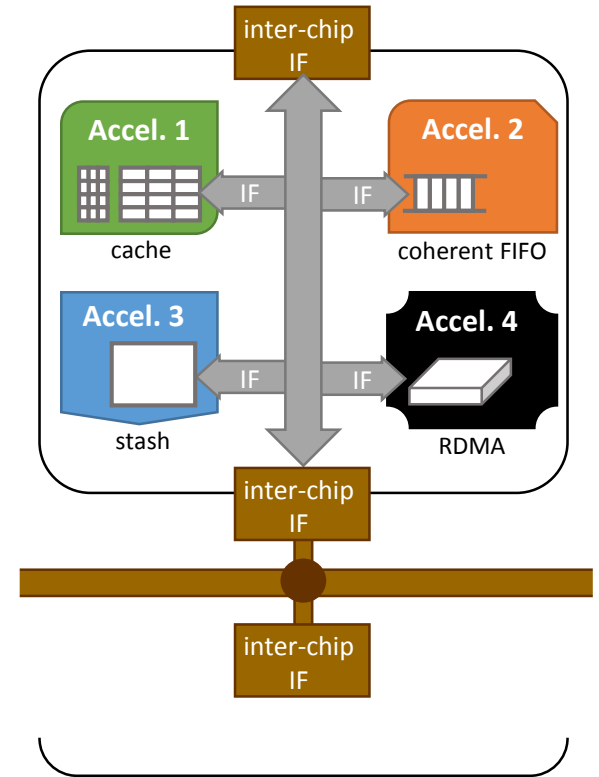


# Data

Data movement critical to efficiency

- Memory structures
- Communication
- Coherence
- Consistency
- Synchronization

Uniform communication interface for hardware  
Abstract to software interface



# Application-Customized Accelerator Communication Arch

Problem: Design + Integrate

Multiple accelerator memory systems + Communication

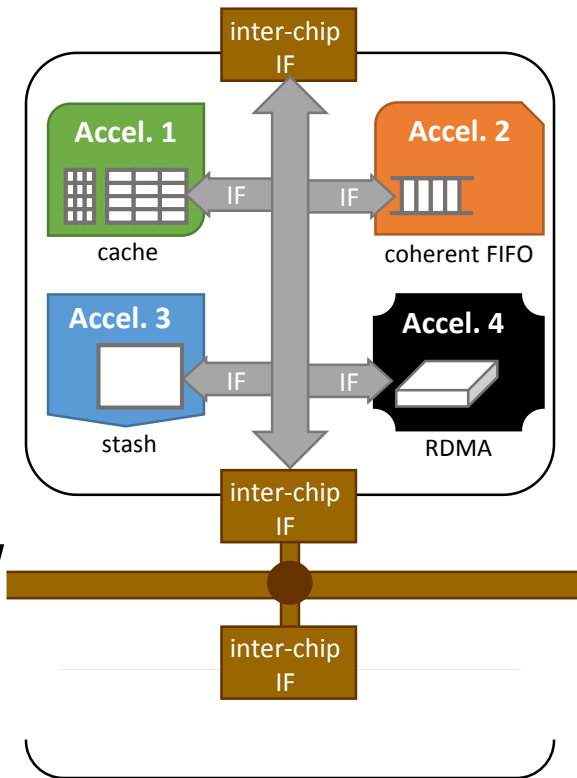
Challenges:

- Friction between different app-specific specializations
- Inefficiencies due to deep memory hierarchy
- Multiple scales: on-chip to cloud

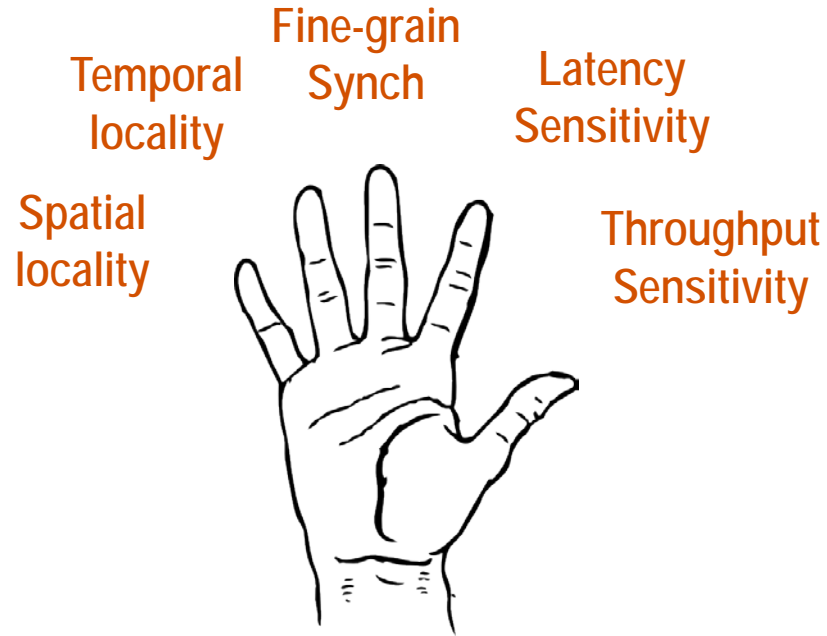
New accelerator communication architecture

- Coherent, global address space
- App-specialized coherence, comm, storage, soln quality

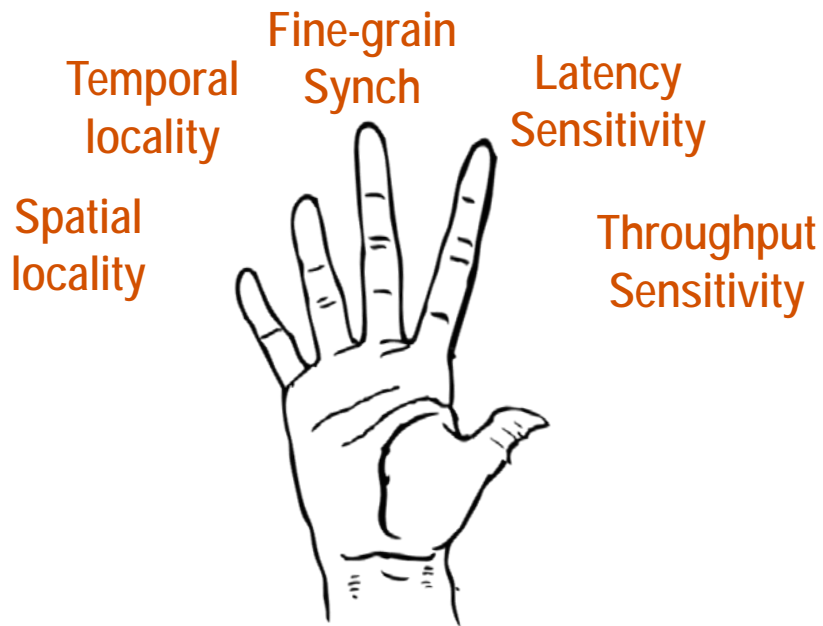
One example next focused on coherence: Spandex [ISCA'18]



# Heterogeneous devices have diverse memory demands

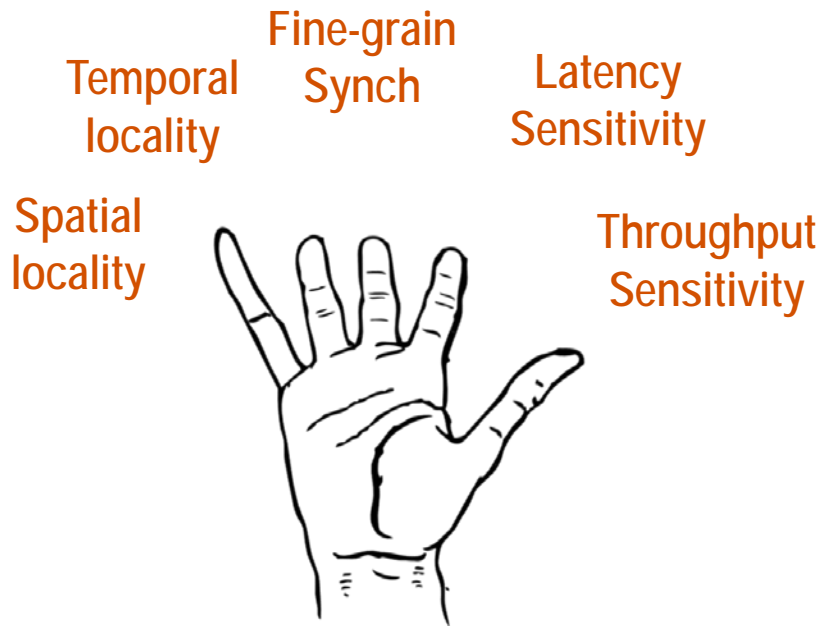


# Heterogeneous devices have diverse memory demands



Typical **CPU** workloads:  
fine-grain synch, latency sensitive

# Heterogeneous devices have diverse memory demands



Typical **GPU** workloads:  
spatial locality, throughput sensitive

# MESI coherence targets CPU workloads

Protocol properties	MESI
Granularity	Line
Stale data invalidation	Writer-invalidate
Write propagation	Ownership

Good for:



CPU

## MESI

- Coarse-grain state
  - ✓ Spatial locality
  - ✗ False sharing
- Writer-initiated invalidation
  - ✓ Temporal locality for reads
  - ✗ Overheads limit throughput, scalability
- Ownership-based updates
  - ✓ Temporal locality for writes
  - ✗ Indirection if low locality

# GPU coherence fits GPU workloads

Protocol properties	MESI	GPU coherence
Granularity	Line	Reads: line writes: word
Stale data invalidation	Writer-invalidate	Self-invalidate
Write propagation	Ownership	Write-through

Good for:



CPU



GPU

## GPU Coherence

- Fine-grain writes
  - ✓ No false sharing
  - ✗ Reduced spatial locality
- Self invalidation
  - ✓ Simple, scalable
  - ✗ Synch limits read reuse
- Write-through caches
  - ✓ Simple, low overhead
  - ✗ Synch limits write reuse

# DeNovo is good fit for CPU and GPU

Protocol properties	MESI	GPU coherence	DeNovo
<b>Granularity</b>	Line	Reads: line writes: word	Reads: flexible Writes: word
<b>Stale data invalidation</b>	Writer-invalidate	Self-invalidate	Self-invalidate
<b>Write propagation</b>	Ownership	Write-through	Ownership

Good for:



CPU



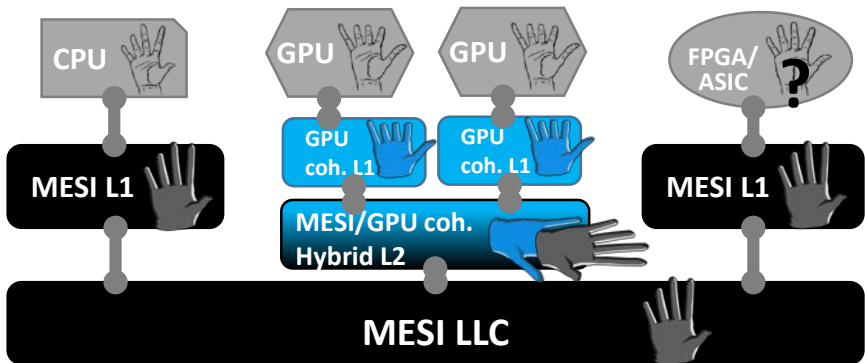
GPU



CPU or GPU

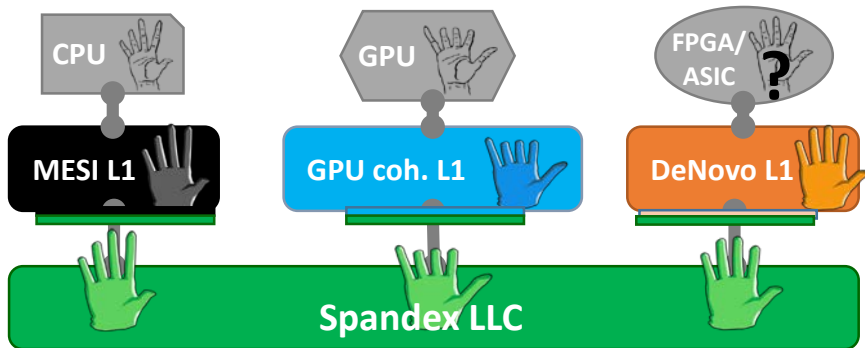


# Integrating Diverse Coherence Strategies



## Existing Solutions: MESI-based LLC

- Accelerator Requests forced to use MESI
- Added latency for inter-device communication
- MESI is complex: extensions are difficult



## Spandex: DeNovo-based interface [ISCA'18]

- Supports write-through and write-back
- Supports self-invalidate and writer-invalidate
- Supports requests of variable granularity
- Directly interfaces MESI, GPU coherence, hybrid (e.g. DeNovo) caches

# Example: Collaborative Graph Applications

Vertex-centric algorithms: distribute vertices among CPU, GPU threads

Application	Access Pattern	Important Dimension	Results
Pull-based PageRank	Read neighbor vertices, Update local vertex	<b>Flat LLC</b> avoids indirection for read misses	Spandex LLC ⇒ <b>37%</b> better exec. time <b>9%</b> better NW traffic
Push-based Betweenness Centrality	Read local vertex, Update (RMW) neighbor vertices	<b>Ownership-based write propagation</b> exploits locality in updates	DeNovo at GPU ⇒ <b>18%</b> better exec. time <b>61%</b> better NW traffic

# Looking Forward...

Software Innovations



Synchronization  
locality

Data locality,  
visibility

Coarse-grain  
operations

Producer/consumer  
relationships

**HPVM + DRF Consistency + ???**

Hardware Innovations

hLRC adaptive  
laziness

Spandex  
dynamic caches

Coherent  
scratchpads  
Stash, ISCA'15

HBM caches

Hardware  
queues

NVRAM

# Approximation

How to express quality of solution from the application to the hardware?

Integrate approximation (quality) into the interface

# Summary

- Interfaces
- Data
- Approximation