

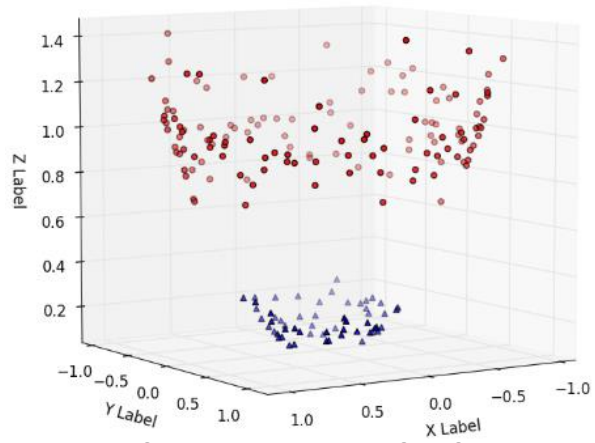
Enhancing Programmable Accelerators for Sparsity

Vidushi Dadu*, Jian Weng*,
Sihao Liu*, Tony Nowatzki*

*UCLA

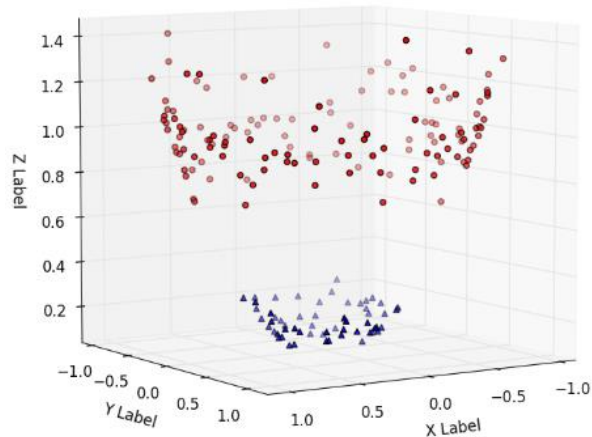
Yarch, HPCA 2019, Feb 17

Irregular workloads are ubiquitous

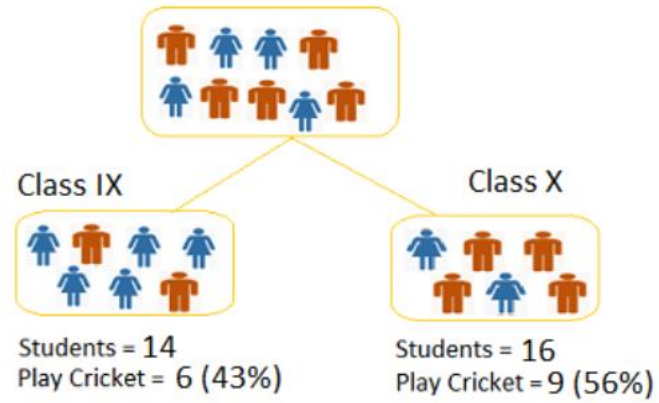


Kernel-SVM on High-dim Data

Irregular workloads are ubiquitous

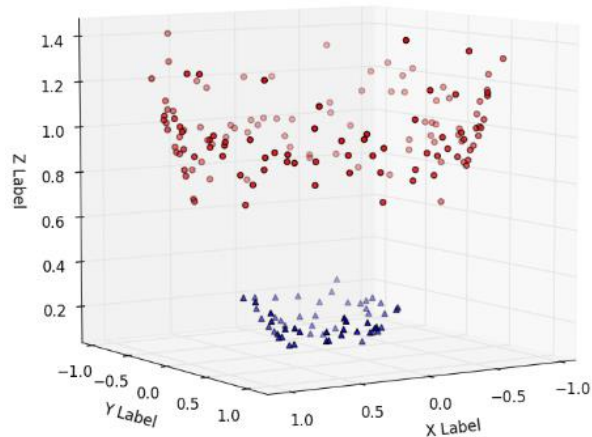


Kernel-SVM on High-dim Data

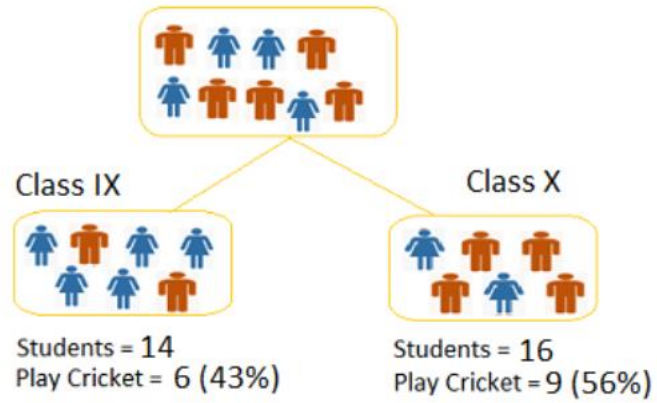


Decision Tree training

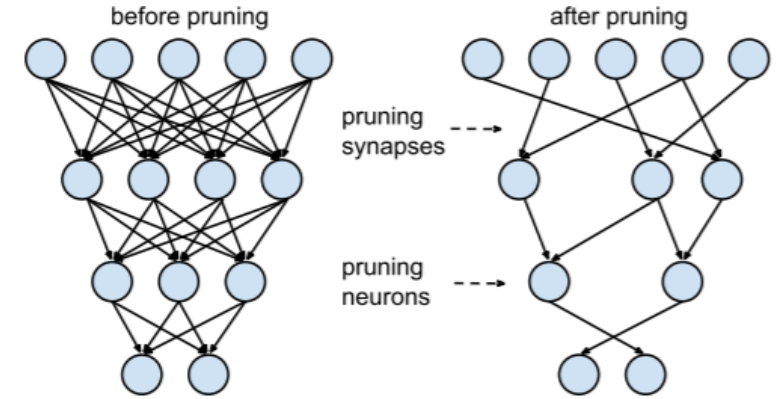
Irregular workloads are ubiquitous



Kernel-SVM on High-dim Data

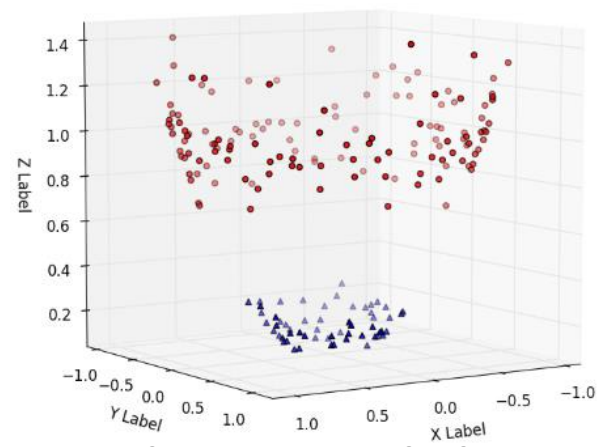


Decision Tree training

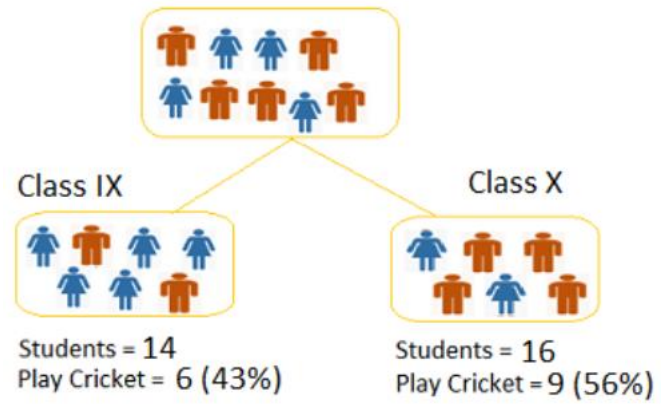


Pruned Deep Neural Networks

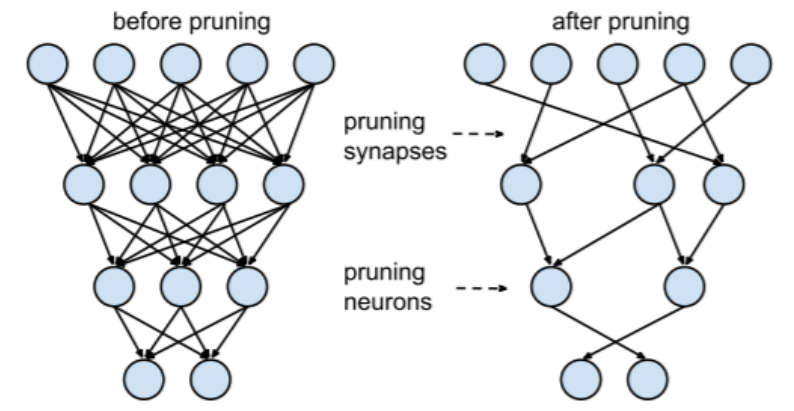
Irregular workloads are ubiquitous



Kernel-SVM on High-dim Data



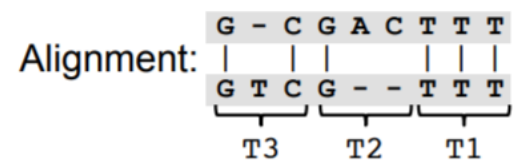
Decision Tree training



Pruned Deep Neural Networks

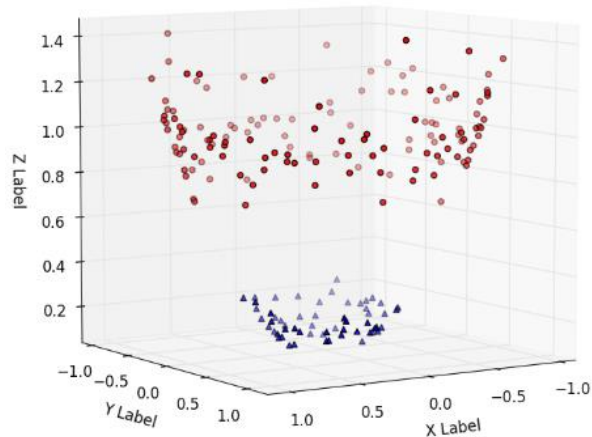
	T1				T2				T3		
	C	T	T	T	C	G	A	C	G	C	
G	0	0	0	0	0	2	0	0	2	0	
T	0	2	2	2	0	0	0	0	1	0	
T	0	2	4	4	2	0	0	2	0	3	
T	0	2	4	6	1	4	3	2			

$(i_{off}=3, j_{off}=3)$ $(i_{off}=3, j_{off}=1)$ $(i_{off}=2, j_{off}=3)$

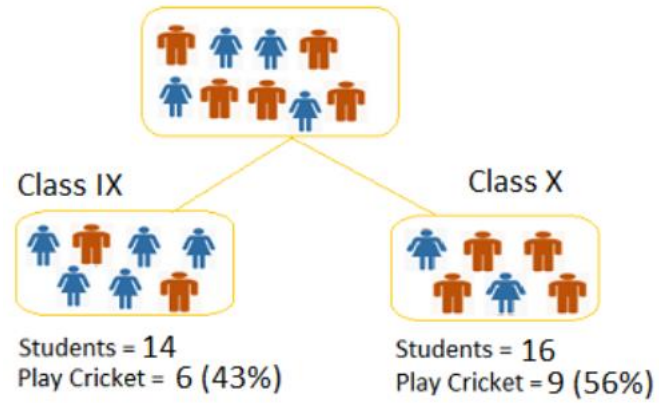


Genomics

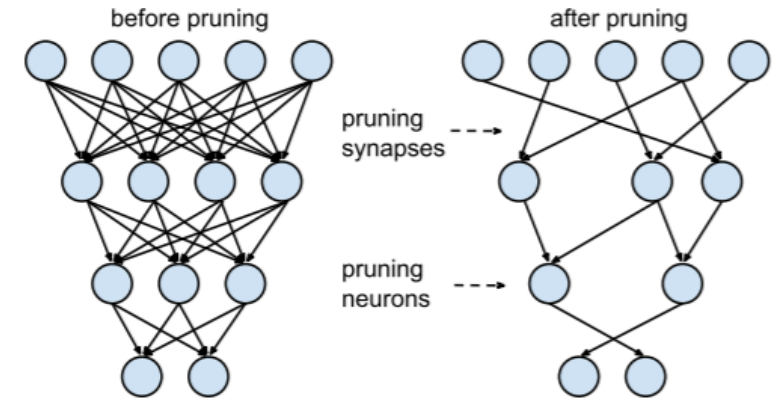
Irregular workloads are ubiquitous



Kernel-SVM on High-dim Data



Decision Tree training



Pruned Deep Neural Networks

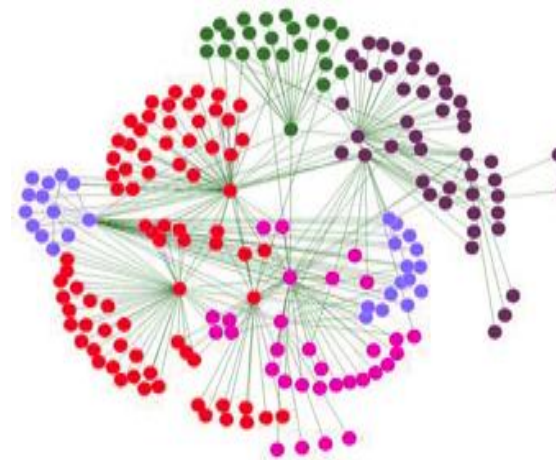
	T1				T2				T3			
	C	T	T	T	C	G	A	C	G	C		
G	0	0	0	0	G	0	2	0	0	G	2	0
T	0	2	2	2	T	0	0	0	0	T	1	0
T	0	2	4	4	C	2	0	0	2	C	0	3
T	0	2	4	6	G	1	4	3	2			

$(i_{off}=3, j_{off}=3)$ $(i_{off}=3, j_{off}=1)$ $(i_{off}=2, j_{off}=3)$

Alignment:

	G	-	C	G	A	C	T	T	T
G	T	C	G	-	-	T	T	T	
	T3			T2		T1			

Genomics



Graph Processing

Sparsity in Workloads Requires Complex Architecture Mechanisms

Workloads/Kernels

Sparse Factorization

GBDT Training

Dynamic Sparsification

Graph Traversal

Shortest Path

Challenging Properties

Indirect Memory Access

Control-dependent Memory

Atomic Updates

Dynamic Parallelism

Load Balancing

Conditional Computation

Heterogeneous Datatypes

Memory
Control
type

Data-9

Sparsity in Workloads Requires Complex Architecture Mechanisms

Workloads/Kernels

Sparse Factorization

GBDT Training

Dynamic Sparsification

Graph Traversal

Shortest Path

Challenging Properties

Indirect Memory Access

Control-dependent Memory

Atomic Updates

Dynamic Parallelism

Load Balancing

Conditional Computation

Heterogeneous Datatypes

Memory
Control
type

Data-7

Sparsity in Workloads Requires Complex Architecture Mechanisms

Workloads/Kernels

Sparse Factorization

GBDT Training

Dynamic Sparsification

Graph Traversal

Shortest Path

Challenging Properties

Indirect Memory Access

Control-dependent Memory

Atomic Updates

Dynamic Parallelism

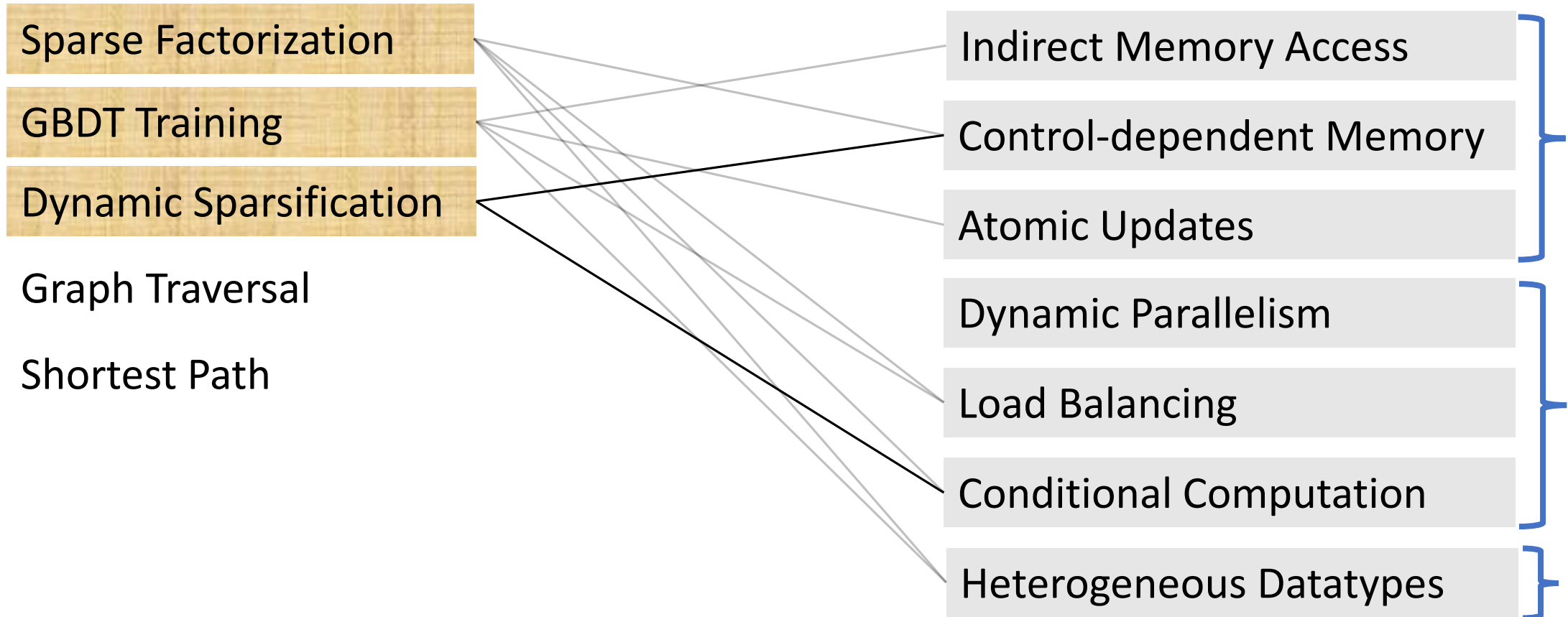
Load Balancing

Conditional Computation

Heterogeneous Datatypes

Memory
Control
type

Data-∞



Sparsity in Workloads Requires Complex Architecture Mechanisms

Workloads/Kernels

Sparse Factorization

GBDT Training

Dynamic Sparsification

Graph Traversal

Shortest Path

Challenging Properties

Indirect Memory Access

Control-dependent Memory

Atomic Updates

Dynamic Parallelism

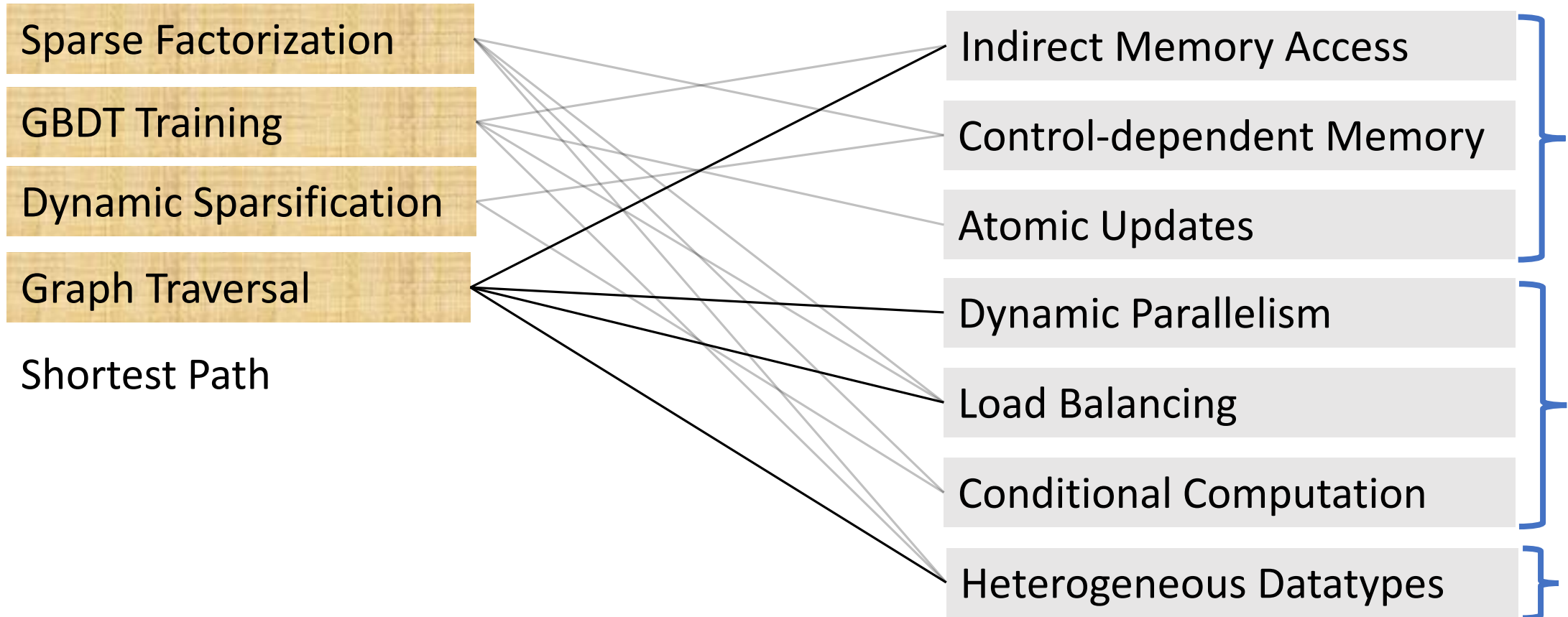
Load Balancing

Conditional Computation

Heterogeneous Datatypes

Memory
Control
type

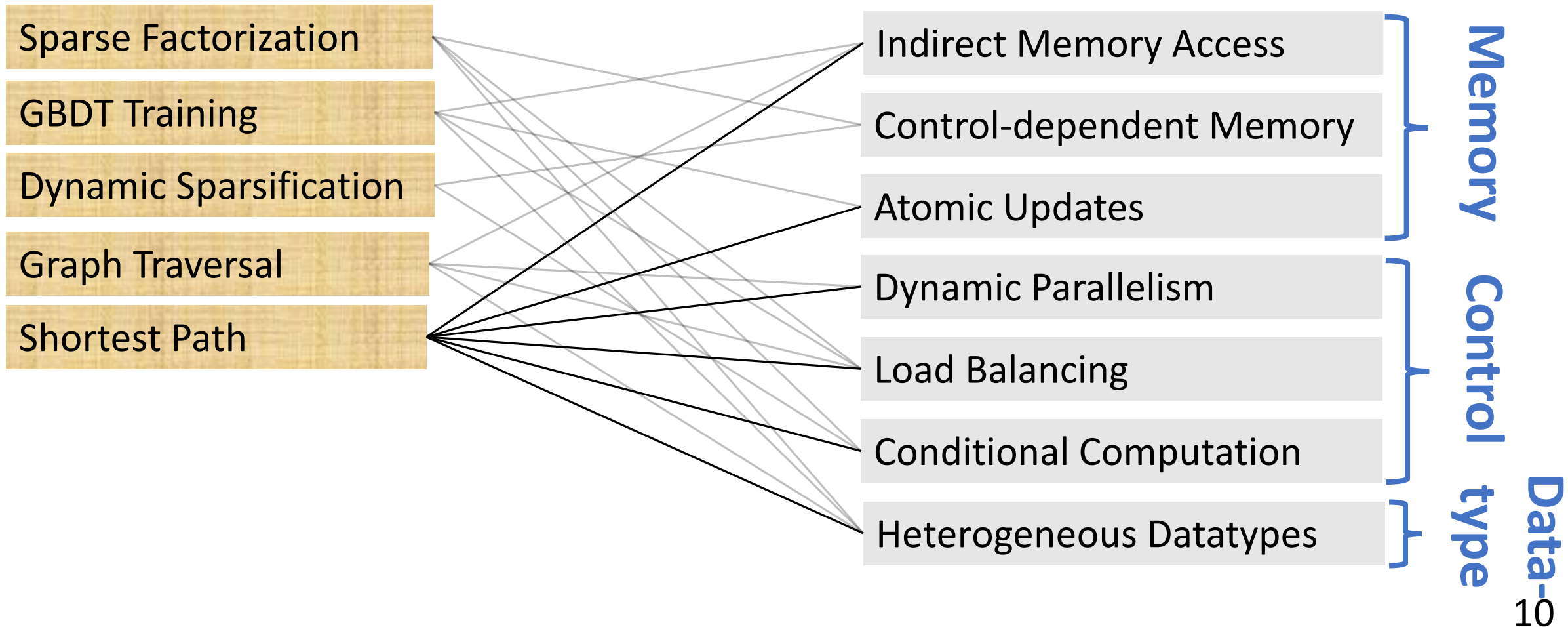
Data-6



Sparsity in Workloads Requires Complex Architecture Mechanisms

Workloads/Kernels

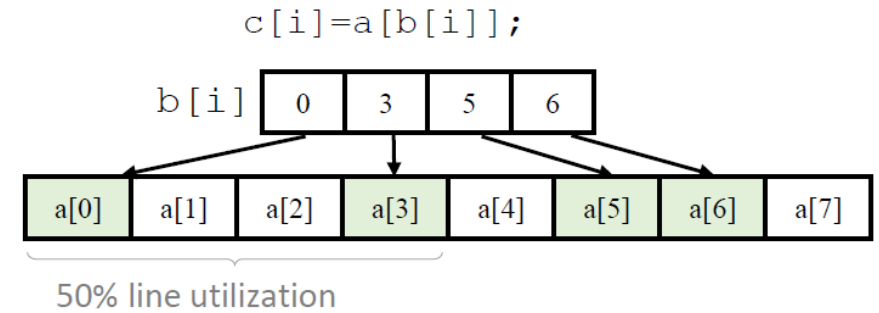
Challenging Properties



GPUs are insufficient for irregularity

Memory Irregularity (Wide mem)

Accesses to random memory location causes inefficient cache line utilization.



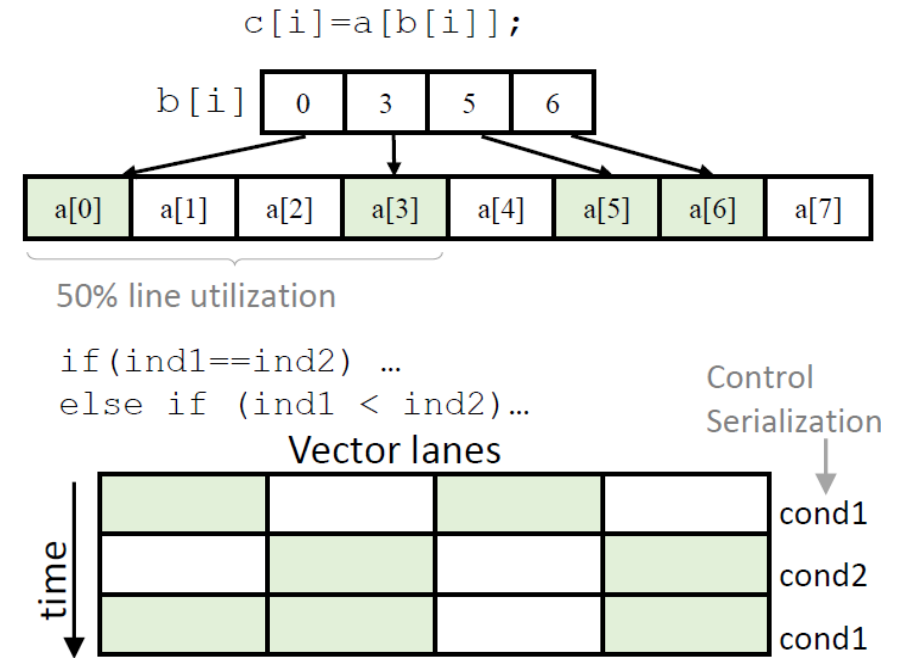
GPUs are insufficient for irregularity

Memory Irregularity (Wide mem)

Accesses to random memory location causes inefficient cache line utilization.

Control Irregularity (SIMD)

Data-dependent compute leads to masking in vector architectures.



GPUs are insufficient for irregularity

Memory Irregularity (Wide mem)

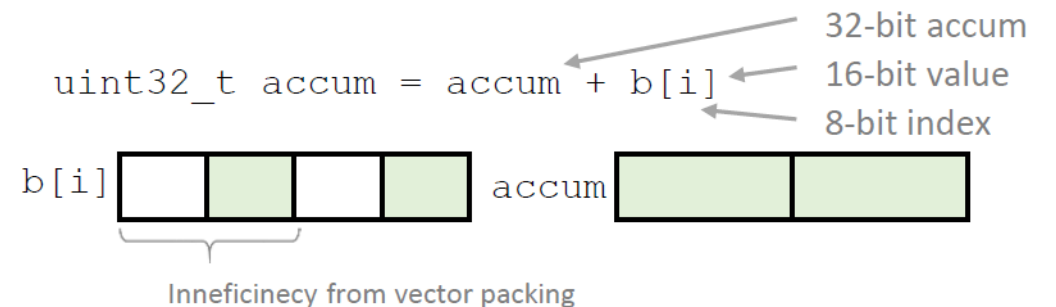
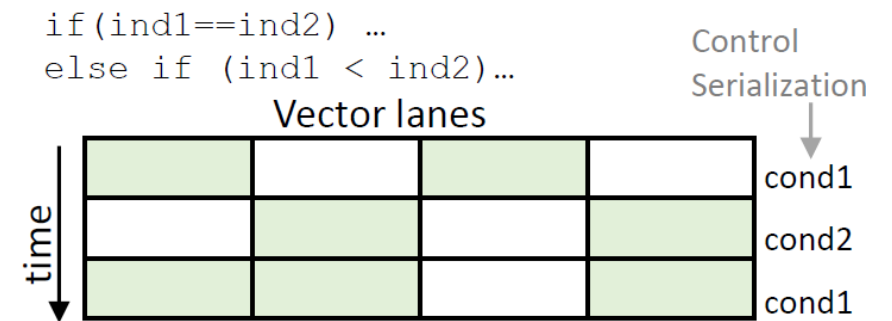
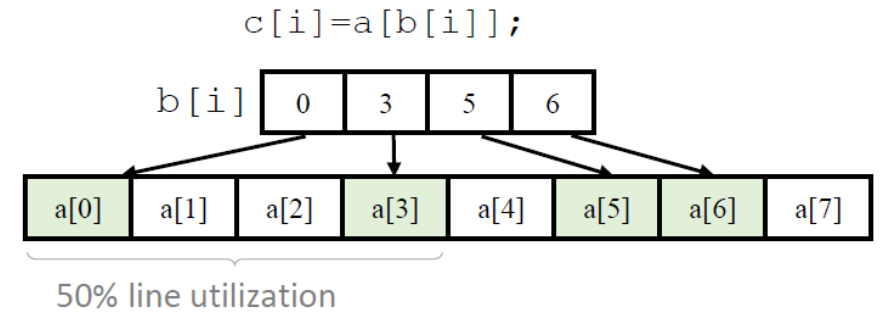
Accesses to random memory location causes inefficient cache line utilization.

Control Irregularity (SIMD)

Data-dependent compute leads to masking in vector architectures.

Datatype Irregularity (Fixed vector width)

Irregular algorithms doesn't allow efficient packing of lower datatypes.



Is Massive Scalar Processor Sufficient?

Is Massive Scalar Processor Sufficient?

There are a couple of problems:

- **General purpose overheads:** Maintaining the program counter and precise state limits performance.

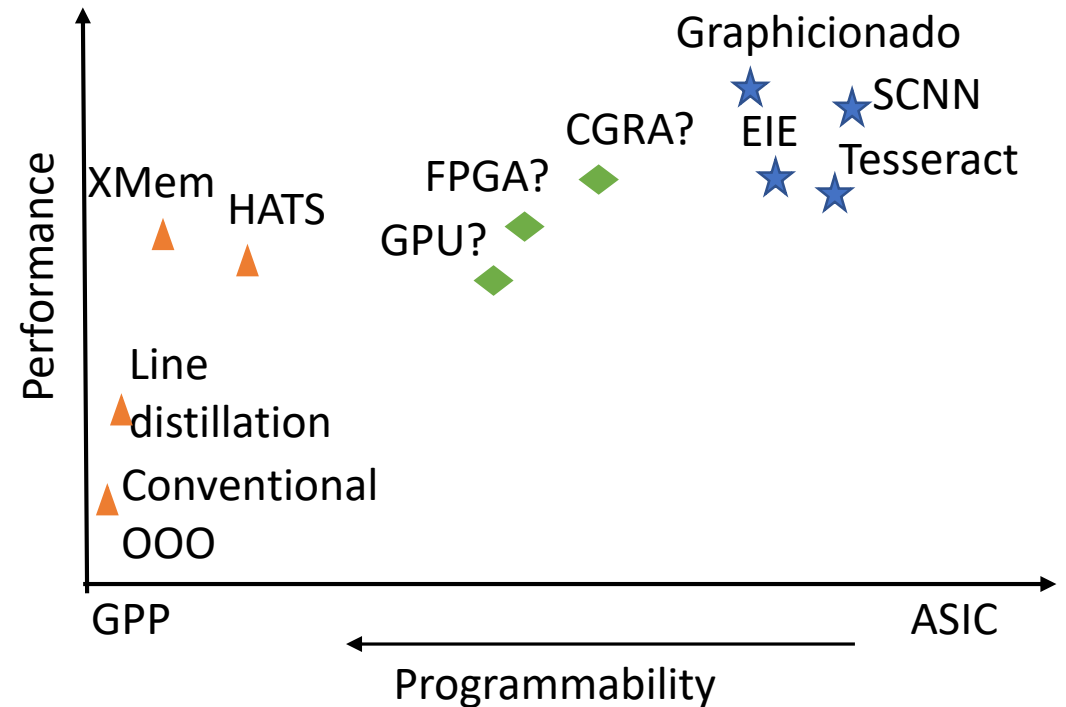
Is Massive Scalar Processor Sufficient?

There are a couple of problems:

- **General purpose overheads:** Maintaining the program counter and precise state limits performance.
- **Programmability:** Such architecture hurts the performance of regular algorithms which have high acceleration potential.

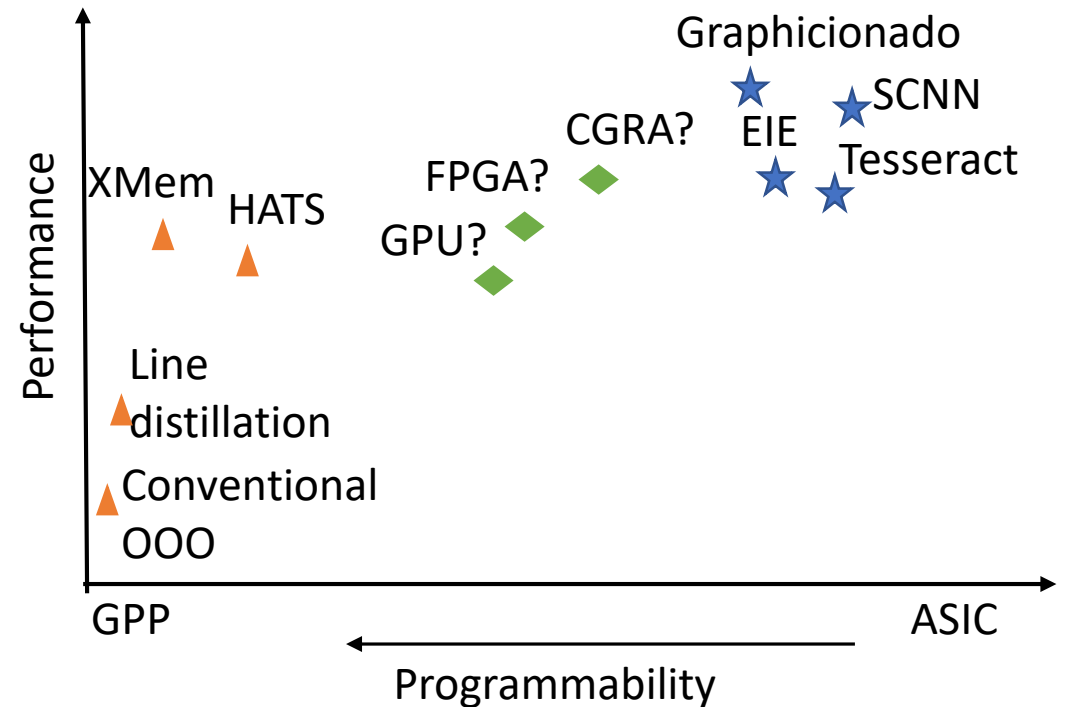
Hope: Sparse Accelerators Have Been Successful

1. **SCNN**: for sparse conv. layer
2. **EIE**: for sparse FC layers
3. **Graphicionado**: graph processor
4. **HATS**: Locality-aware scheduling for graph processing
5. **XMem**: Programmer hints to prefetchers, caching policies



Hope: Sparse Accelerators Have Been Successful

1. **SCNN**: for sparse conv. layer
2. **EIE**: for sparse FC layers
3. **Graphicionado**: graph processor
4. **HATS**: Locality-aware scheduling for graph processing
5. **XMem**: Programmer hints to prefetchers, caching policies



It will be useful if we have an architecture which gives high performance on the set of workloads we care about most.

Goal: Is there an accelerator paradigm which performs well on most irregular workloads?

Goal: Is there an accelerator paradigm which performs well on most irregular workloads?

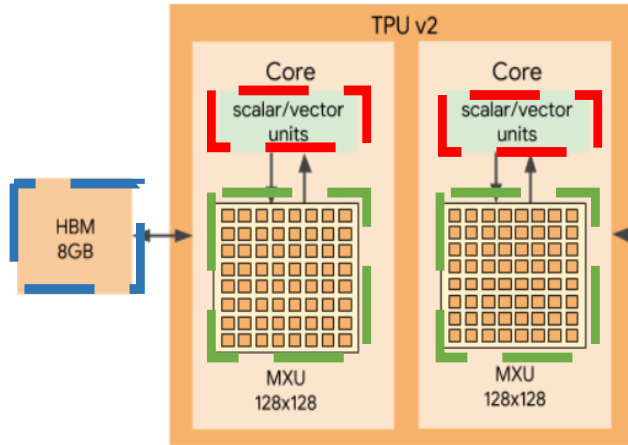
In this talk, I will focus on **Irregular Machine Learning**.

In machine learning domain, both dense and sparse scenarios are pretty common.

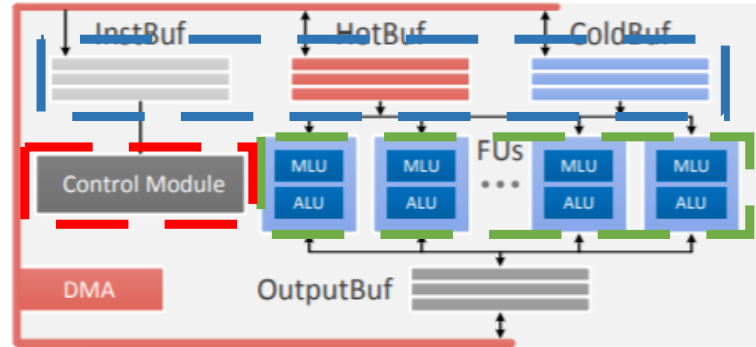
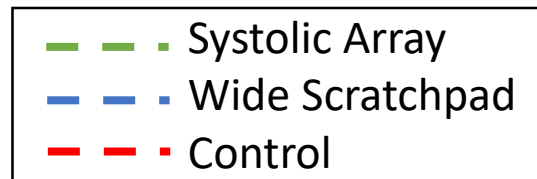
Sub-goal: Design a programmable sparse accelerator while maintaining efficiency for the dense workloads.

Our Approach: Start with a Dense Programmable Accelerator

Our Approach: Start with a Dense Programmable Accelerator



Google TPU v2



PuDianNao (ASPLOS'15)

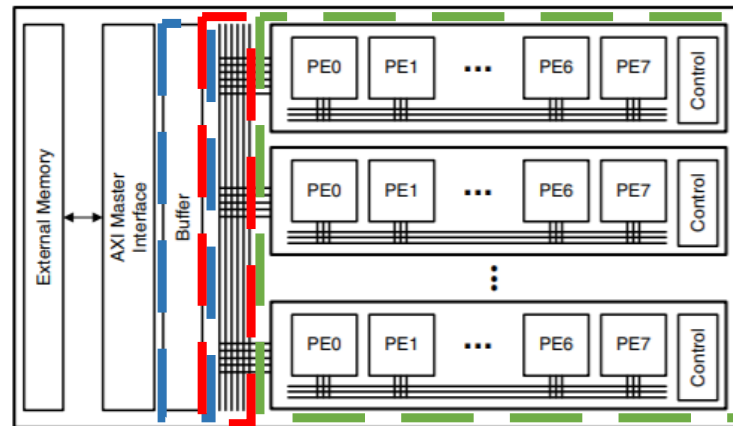
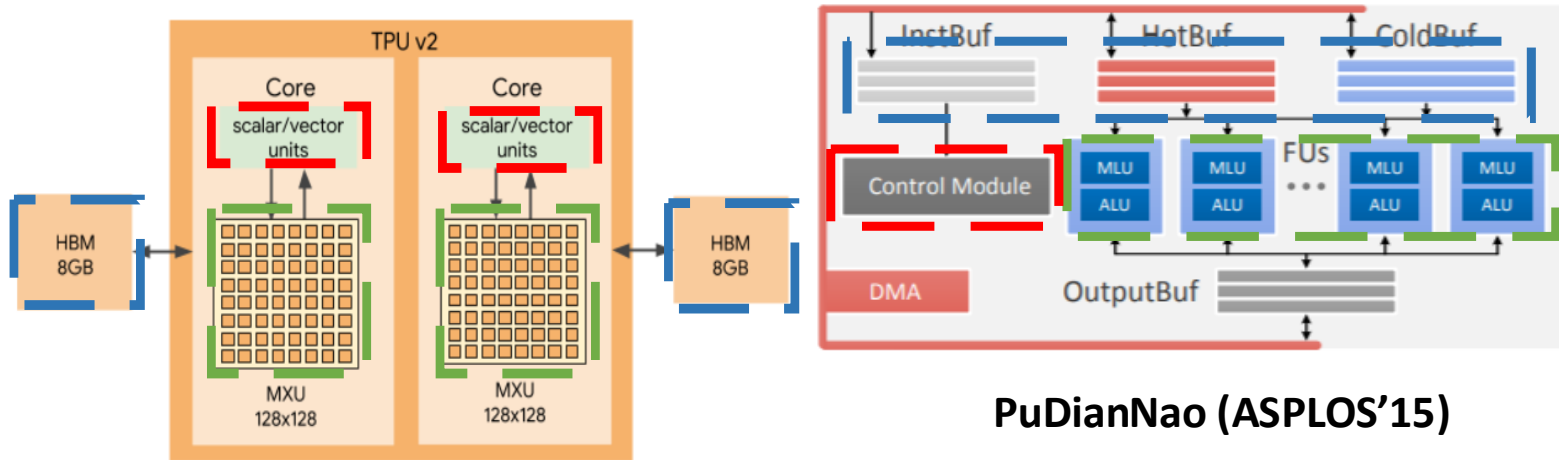


Tabla (HPCA'16)

Our Approach: Start with a Dense Programmable Accelerator



Google TPU v2

PuDianNao (ASPLOS'15)

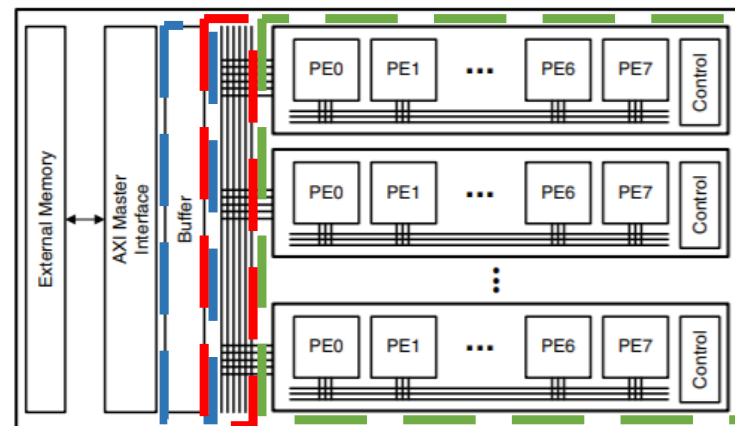
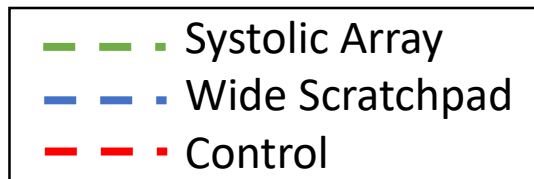
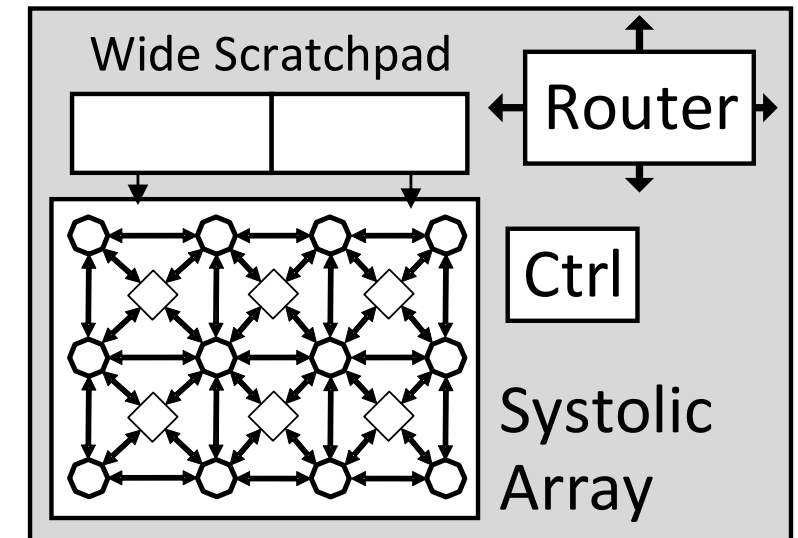
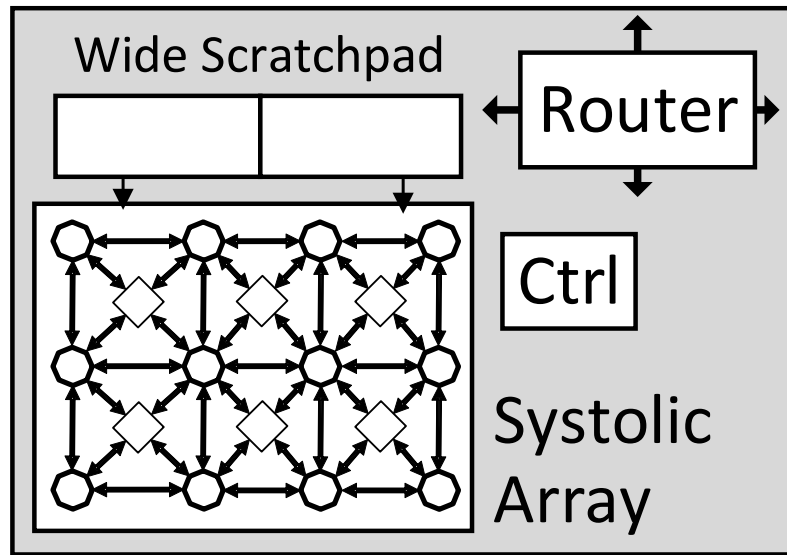


Tabla (HPCA'16)

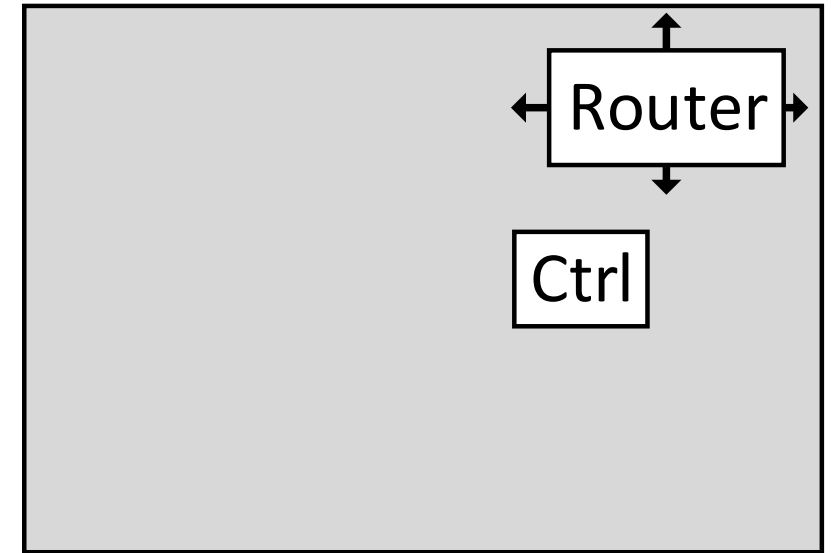


Stereotypical Dense Accelerator Core

Overview of Sparsity-Enabled SPU core

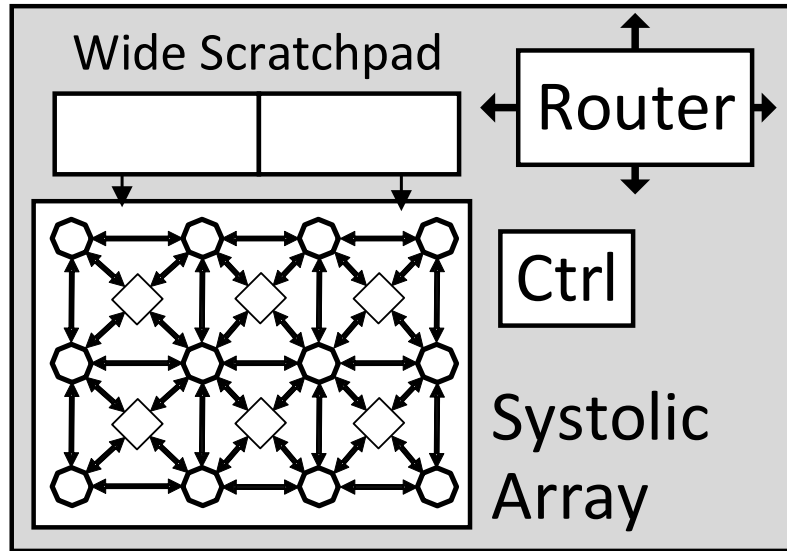


(a) Stereotypical Dense Accelerator Core



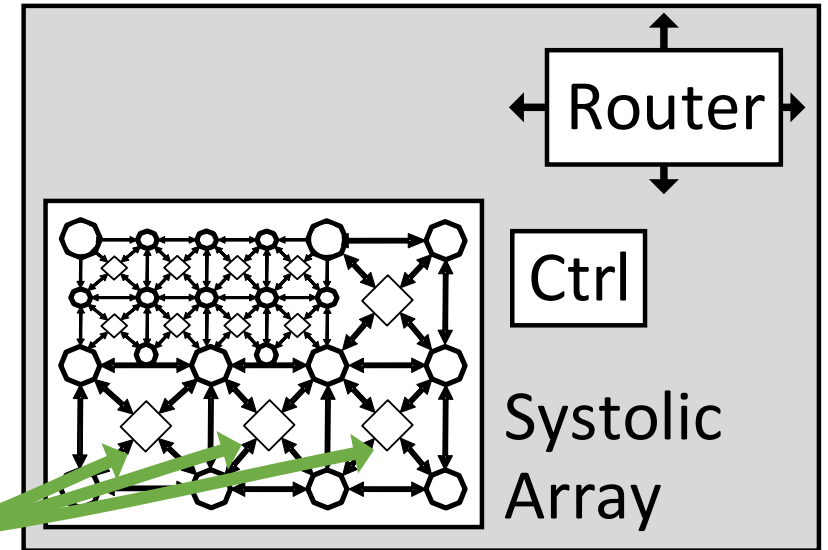
(b) Sparsity Enabled Accelerator (SPU Core)

Overview of Sparsity-Enabled SPU core



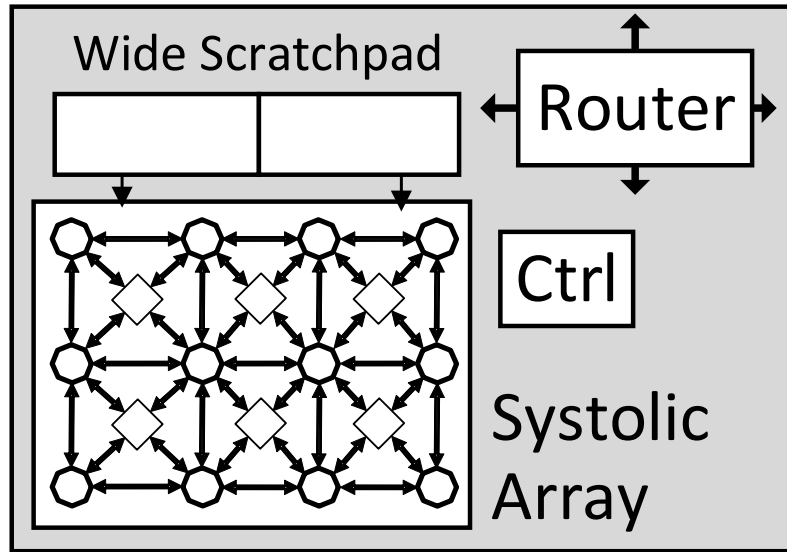
(a) Stereotypical Dense Accelerator Core

Systolic array with novel meta-reuse control flow



(b) Sparsity Enabled Accelerator (SPU Core)

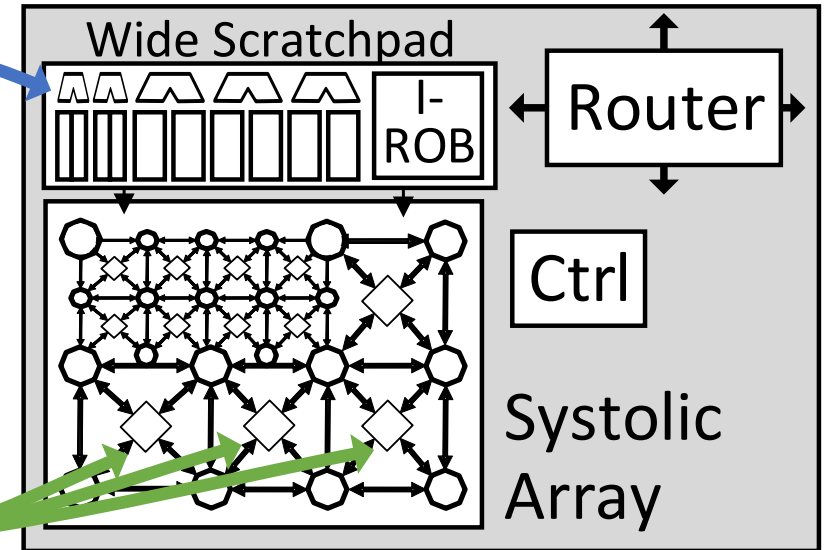
Overview of Sparsity-Enabled SPU core



(a) Stereotypical Dense Accelerator Core

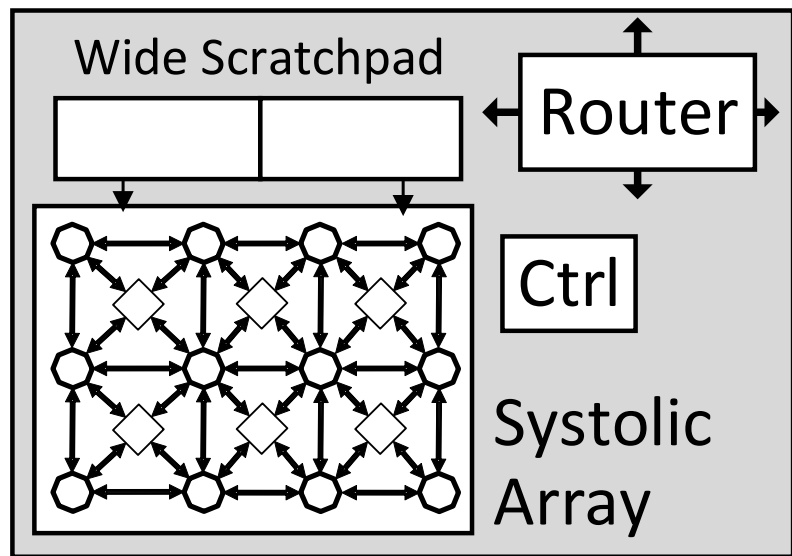
Compute-Enabled
High-bandwidth
Indirect Scratchpad

Systolic array with
novel meta-reuse
control flow



(b) Sparsity Enabled Accelerator (SPU Core)

Overview of Sparsity-Enabled SPU core

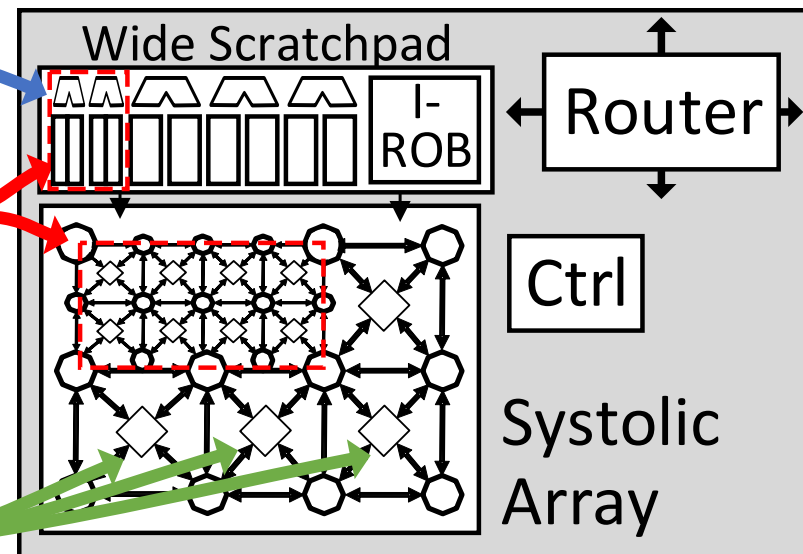


(a) Stereotypical Dense Accelerator Core

Compute-Enabled
High-bandwidth
Indirect Scratchpad

Decomposable
Memory/Network/
Compute

Systolic array with
novel meta-reuse
control flow



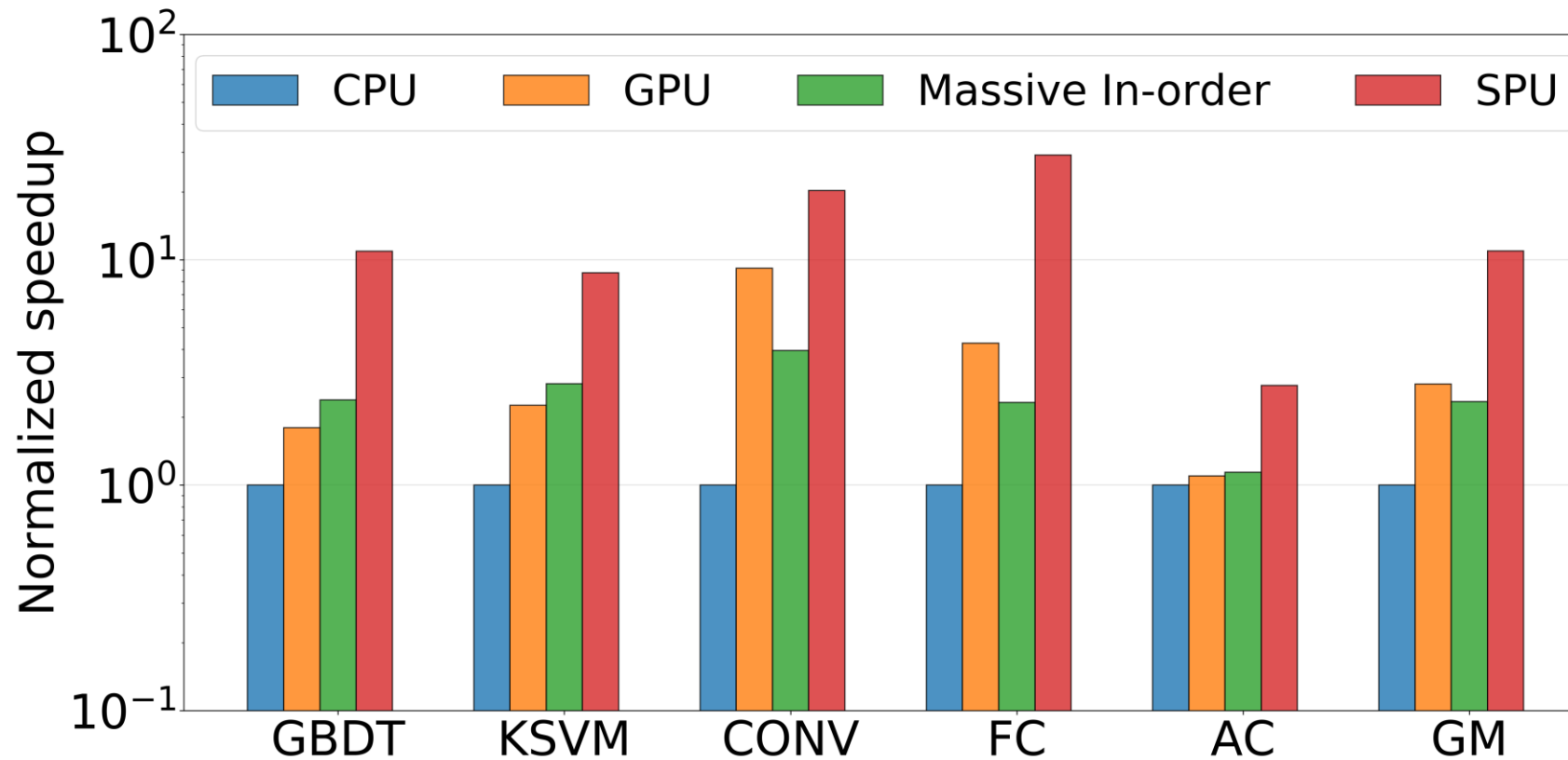
(b) Sparsity Enabled Accelerator (SPU Core)

Evaluation Methodology

- gem5 cycle level simulator
- Dataflow compiler
- Benchmarks (top-5 ML algorithms using by Facebook in 2018*):
 - KSVM, Conv layer, FC layer, GBDT, Arithmetic Circuits
- Datasets: Open-source practical datasets
- Baselines:
 - CPU: 4-core CPU (reference only)
 - GPU: Nvidia P4000
 - Massive In-order: SPU but with 512 inorder cores
 - SPU: Sparse Proc. Unit (64 Core)

*Kim Hazelwood, et al. Applied machine learning at facebook: A datacenter infrastructure perspective. In High Performance Computer Architecture (HPCA), 2018

SPU achieves 2.2-6.8x performance over GPU



Ongoing/Future Work

- Compiler for streaming-dataflow architecture.
- Specialized cache hierarchy
- Hardware-software (ISA) primitives for other irregular domains:
 - Graph processing
 - Genomics
 - Compression/Decompression