

Fine-Grained Accelerators for Sparse Machine Learning Workloads

ASP-DAC 2017

Asit K. Mishra, Eriko Nurvitadhi, Ganesh Venkatesh,
Jonathan Pearce, Debbie Marr

Accelerator Architecture Lab
Intel Corp.

Disclaimer: The views expressed in this talk are those of the speaker and not his employer.

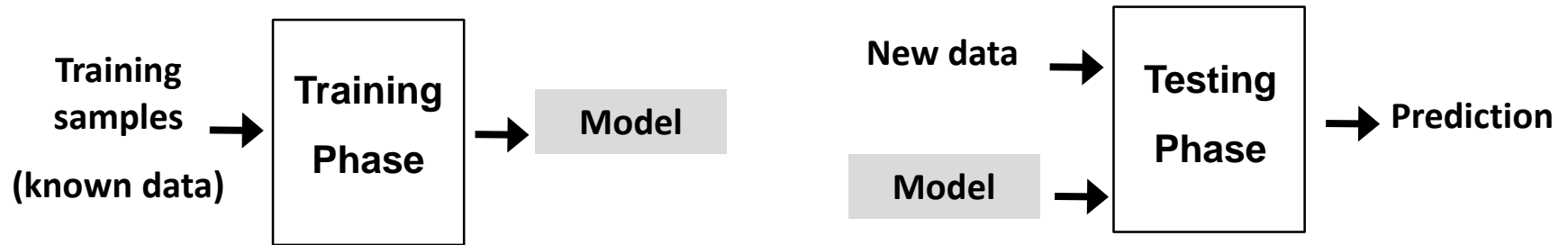
Executive Summary

- Motivation
 - Data analytics growing in importance
 - They rely on machine learning (ML) algorithms
 - Working on datasets that are sparse (texts, ratings)
- This work: accelerate sparse ML workloads
 - Characterized ML workloads → low IPC, mem & branch mispredict stalls, high \$ miss rate
 - Proposed HW accelerator → 4-13x speedups and 9-17x better energy over CPU, with small area

Talk Outline

- **Executive summary**
- **Sparse Machine Learning**
- **Sparse matrix processing**
- **Characterization study**
- **Proposed hardware accelerator**
- **Related work + summary**

Machine Learning for Data Analytics



Data represented as Matrix

Misplaced top-level domain (TLD)

URL	Reputation		Feature 0	Feature 1	...
www.ebay.com.phishy.biz	malicious	sample0	True		
www.ebay.com	normal	sample1	False		
	

Sparse

High-dimensional

Many Real World Datasets Are Large, Sparse, and High-Dimensional

Examples from datasets we studied

Name	Avg. length	n =Max. length	l = # Samples	% Sparsity	Size on disk
E2006	1242	150K	16K	0.83%	485MB
RCV	74	42K	677K	0.15%	1.2GB
Webspam Unigram	86	255	245K	33.38%	268MB
Webspam Trigram	3.7K	399K	245K	0.04%	17GB
Gamevideo	221	1K	97K	22%	225MB
URL	117	2.6M	1.6M	0.003%	1.5GB
CriteoLabs	33	25.2M	32M	0.0001%	25GB
MovieLens	70K users, 10K movies 10M ratings				253MB

Talk Outline

- **Executive summary**
- **Sparse Machine Learning**
- **Sparse matrix processing**
- **Characterization study**
- **Proposed hardware accelerator**
- **Related work + summary**

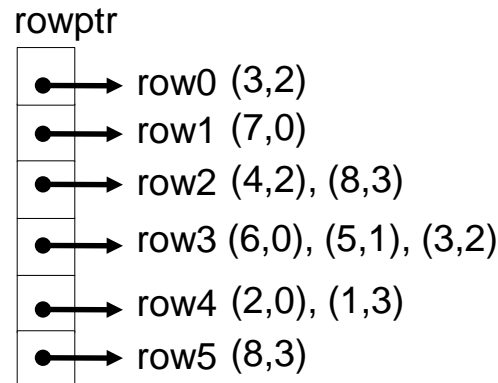
Matrix Formats: CSR vs. CSC

A

0	0	3	0
7	0	0	0
0	0	4	8
6	5	3	0
2	0	0	1
0	0	0	8

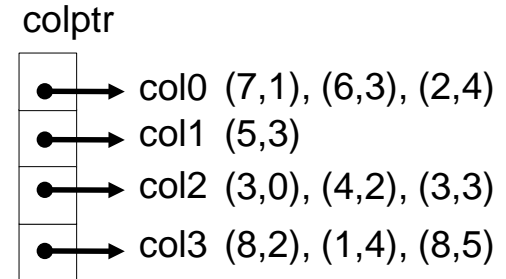
**A matrix
example**

Rows = samples
Columns = features



**Compressed
Sparse Row
(CSR)**

*Good for operation
on samples*



**Compressed
Sparse Column
(CSC)**

*Good for operation
on features*

Example Matrix Operations

spMdV_csr:

spMspV_csc:

spMdV_csc:

Row-oriented sparse
matrix * dense vector

$$\begin{array}{|c|c|c|c|} \hline \text{A} & & \text{x} & \\ \hline 0 & 0 & 3 & 0 \\ \hline 7 & 0 & 0 & 0 \\ \hline 0 & 0 & 4 & 8 \\ \hline 6 & 5 & 3 & 0 \\ \hline 2 & 0 & 0 & 1 \\ \hline 0 & 0 & 0 & 8 \\ \hline \end{array} * \begin{array}{|c|} \hline 1 \\ \hline 2 \\ \hline 4 \\ \hline 3 \\ \hline \end{array} = \begin{array}{|c|} \hline 12 \\ \hline \\ \hline \\ \hline \\ \hline \\ \hline \\ \hline \end{array}$$

...

$$\begin{array}{|c|c|c|c|} \hline 0 & 0 & 3 & 0 \\ \hline 7 & 0 & 0 & 0 \\ \hline 0 & 0 & 4 & 8 \\ \hline 6 & 5 & 3 & 0 \\ \hline 2 & 0 & 0 & 1 \\ \hline 0 & 0 & 0 & 8 \\ \hline \end{array} * \begin{array}{|c|} \hline 1 \\ \hline 2 \\ \hline 4 \\ \hline 3 \\ \hline \end{array} = \begin{array}{|c|} \hline 12 \\ \hline 7 \\ \hline 40 \\ \hline 28 \\ \hline 5 \\ \hline 24 \\ \hline \end{array}$$

Irregular reads on x

Column-oriented sparse
matrix * sparse vector

$$\begin{array}{|c|c|c|c|} \hline \text{A} & & \text{x} & \\ \hline 0 & 0 & 3 & 0 \\ \hline 7 & 0 & 0 & 0 \\ \hline 0 & 0 & 4 & 8 \\ \hline 6 & 5 & 3 & 0 \\ \hline 2 & 0 & 0 & 1 \\ \hline 0 & 0 & 0 & 8 \\ \hline \end{array} * \begin{array}{|c|} \hline 0 \\ \hline 0 \\ \hline 4 \\ \hline 8 \\ \hline \end{array} = \begin{array}{|c|} \hline 12 \\ \hline 0 \\ \hline 16 \\ \hline 12 \\ \hline 0 \\ \hline 0 \\ \hline \end{array}$$

$$\begin{array}{|c|c|c|c|} \hline 0 & 0 & 3 & 0 \\ \hline 7 & 0 & 0 & 0 \\ \hline 0 & 0 & 4 & 8 \\ \hline 6 & 5 & 3 & 0 \\ \hline 2 & 0 & 0 & 1 \\ \hline 0 & 0 & 0 & 8 \\ \hline \end{array} * \begin{array}{|c|} \hline 0 \\ \hline 0 \\ \hline 4 \\ \hline 8 \\ \hline \end{array} = \begin{array}{|c|} \hline 12 \\ \hline 0 \\ \hline 80 \\ \hline 12 \\ \hline 8 \\ \hline 64 \\ \hline \end{array}$$

Irregular reads and writes on y

Scale matrix using scaling factors in x, then update y

$$\begin{array}{|c|c|c|c|} \hline \text{A} & & \text{x (scaling factors)} & \\ \hline 0 & 0 & 3 & 0 \\ \hline 7 & 0 & 0 & 0 \\ \hline 0 & 0 & 4 & 8 \\ \hline 6 & 5 & 3 & 0 \\ \hline 2 & 0 & 0 & 1 \\ \hline 0 & 0 & 0 & 8 \\ \hline \end{array} * \begin{array}{|c|} \hline 1 \\ \hline 2 \\ \hline 1 \\ \hline 2 \\ \hline 1 \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline 0 & 0 & 3 & 0 \\ \hline \end{array} \text{y}$$

=

$$\begin{array}{|c|c|c|c|} \hline 0 & 0 & 3 & 0 \\ \hline \end{array} \text{y}$$

...

$$\begin{array}{|c|c|c|c|} \hline \text{A} & & \text{x (scaling factors)} & \\ \hline 0 & 0 & 3 & 0 \\ \hline 7 & 0 & 0 & 0 \\ \hline 0 & 0 & 4 & 8 \\ \hline 6 & 5 & 3 & 0 \\ \hline 2 & 0 & 0 & 1 \\ \hline 0 & 0 & 0 & 8 \\ \hline \end{array} * \begin{array}{|c|} \hline 1 \\ \hline 2 \\ \hline 1 \\ \hline 2 \\ \hline 1 \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline 24 & 5 & 12 & 25 \\ \hline \end{array} \text{y}$$

=

$$\begin{array}{|c|c|c|c|} \hline 24 & 5 & 12 & 25 \\ \hline \end{array} \text{y}$$

Irregular reads and writes on y

Talk Outline

- **Executive summary**
- **Sparse Machine Learning**
- **Sparse matrix processing**
- **Characterization study**
- **Proposed hardware accelerator**
- **Related work + summary**

Methodology

- System Under Study
 - 2.7 GHz Intel Ivy Bridge Server (E5-2679 v2)
 - 24 cores, 32KB I-cache, 32KB D-cache,
 - 256 KB private L2 cache, 30 MB shared L3
 - 128 GB DDR3 memory, 60 GB/s max mem bandwidth
- Dataset
 - Real datasets, shown in earlier slide
- Tools
 - Vtune and gprof for hotspot characterizations
 - Sniper simulator to get cache statistics
 - McPat for energy modeling

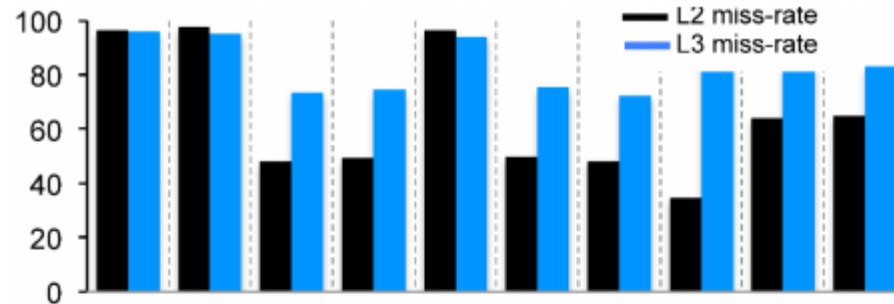
Workloads and Identified Hotspots

Application	Type	Hot code	% Time
Sparse PCA	Dim. reduction	SpVSpV	99%
Kernelized SVM classification	Classification	SpMSpV	96%
Linear SVM classification	Classification	SpMDV, SpVDV	99%
Logistic regression	Classification	SpMDV, SpVDV	98%
Kernelized SVM regression	Regression	SpMSpV	94%
Linear SVM regression	Regression	SpMDV, SpVDV	99%
SLIM	Recom. engine	SpMDV	88%
ALS	Recom. engine	SpMDV	92%
K-means	Clustering	SpVDV	90%

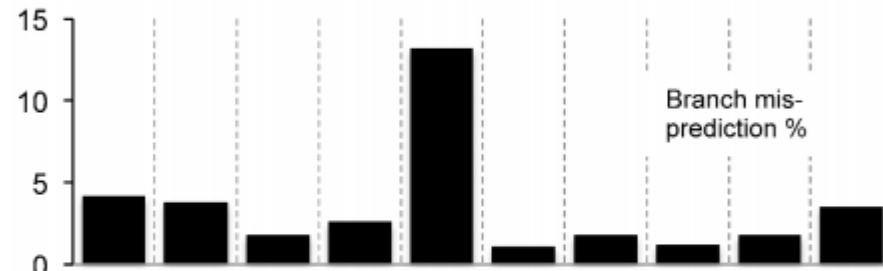
Majority of time spent on sparse matrix/vector ops

Application Characteristics

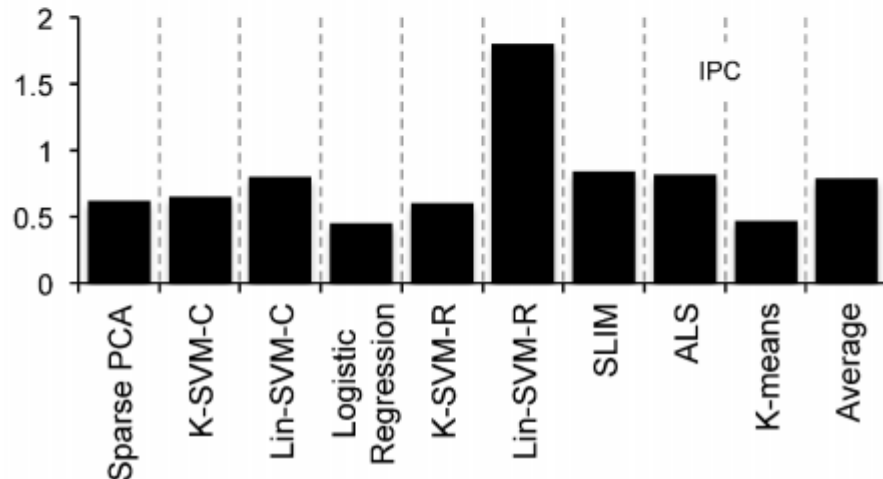
High cache miss rate



High Branch Misprediction



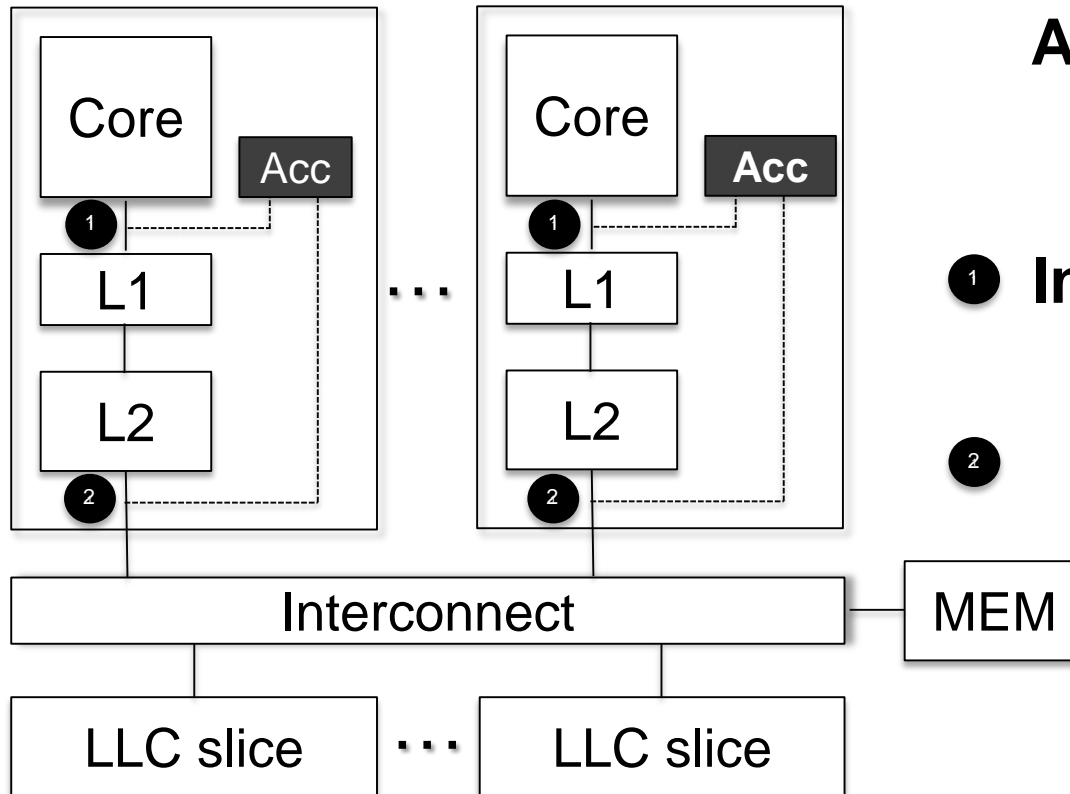
Low IPC



Talk Outline

- **Executive summary**
- **Sparse Machine Learning**
- **Sparse matrix processing**
- **Characterization study**
- **Proposed hardware accelerator**
- **Related work + summary**

System Architecture



**Accelerator Tightly
Connected to CPU**

1 Interface to L1 cache

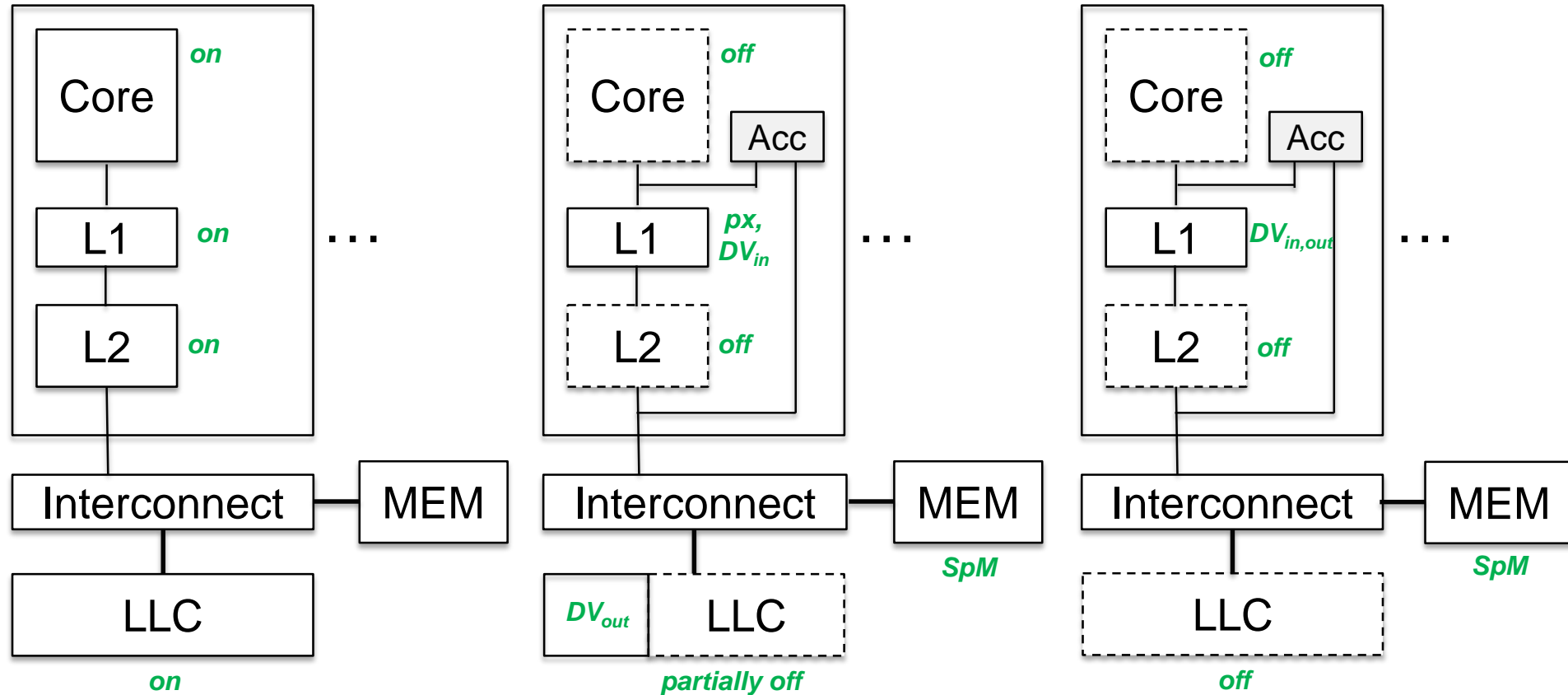
**2 Direct channel to
external memory**

How to improve efficiency: custom config for each matrix ops

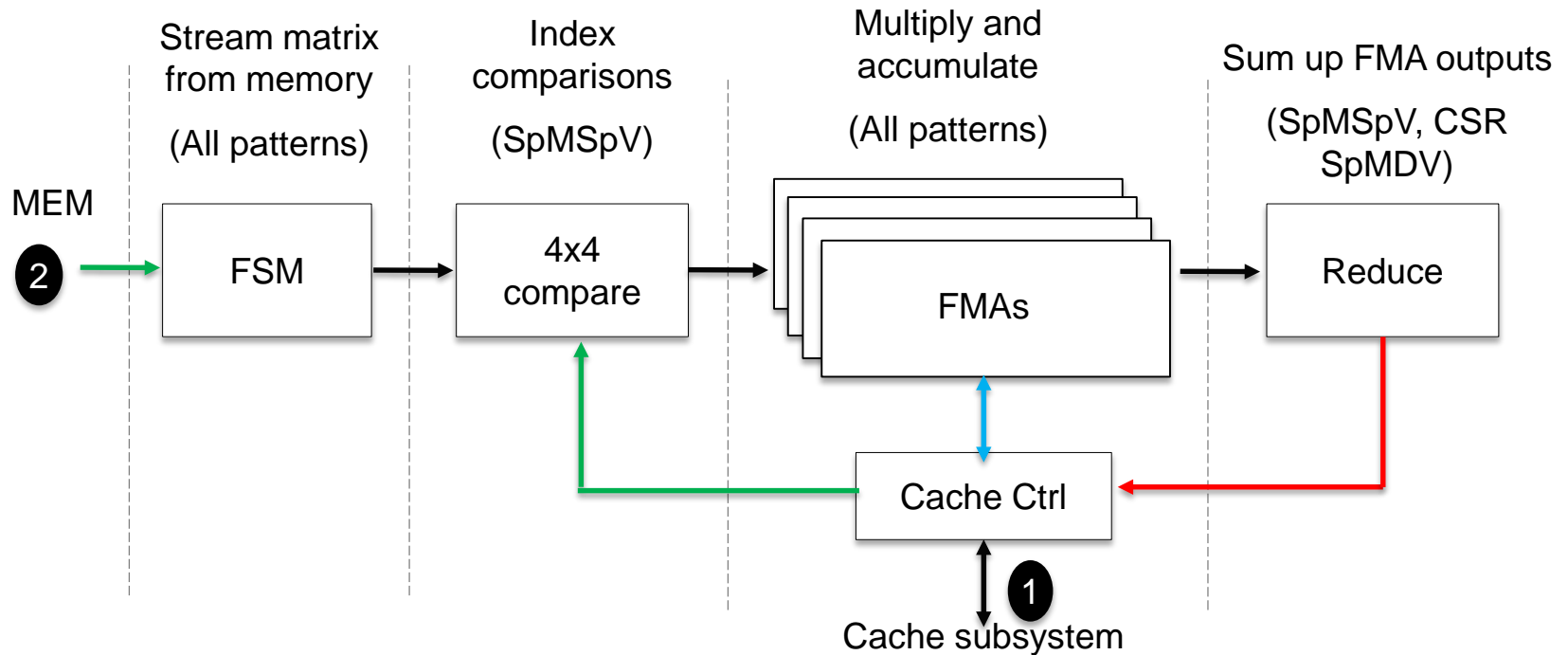
Baseline

spMspV, spMDV

reduce(spM,DV)



Accelerator Internal

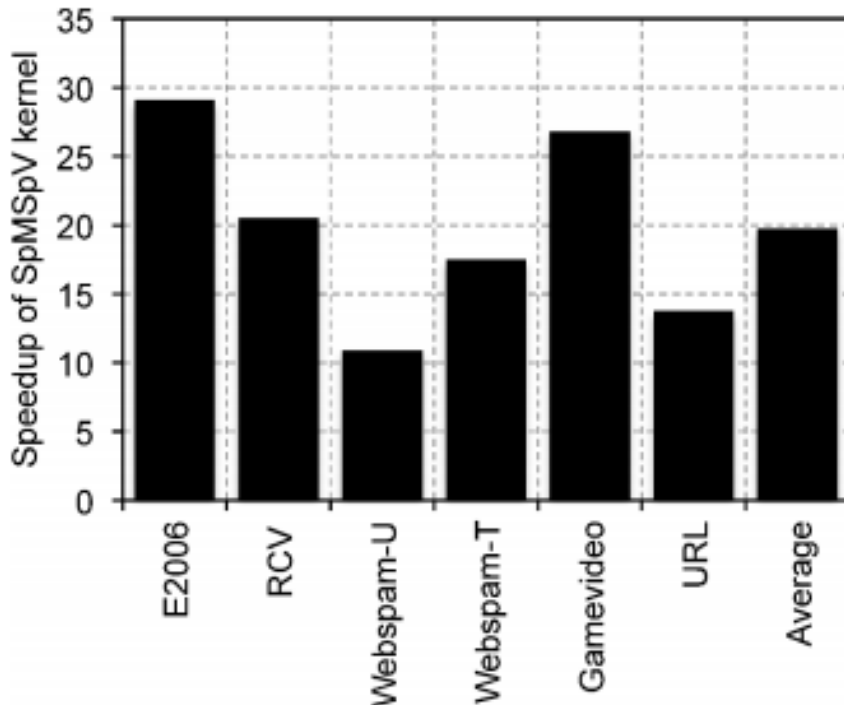


Support matrix operations used in ML workloads under study

(See details in the paper)

Results

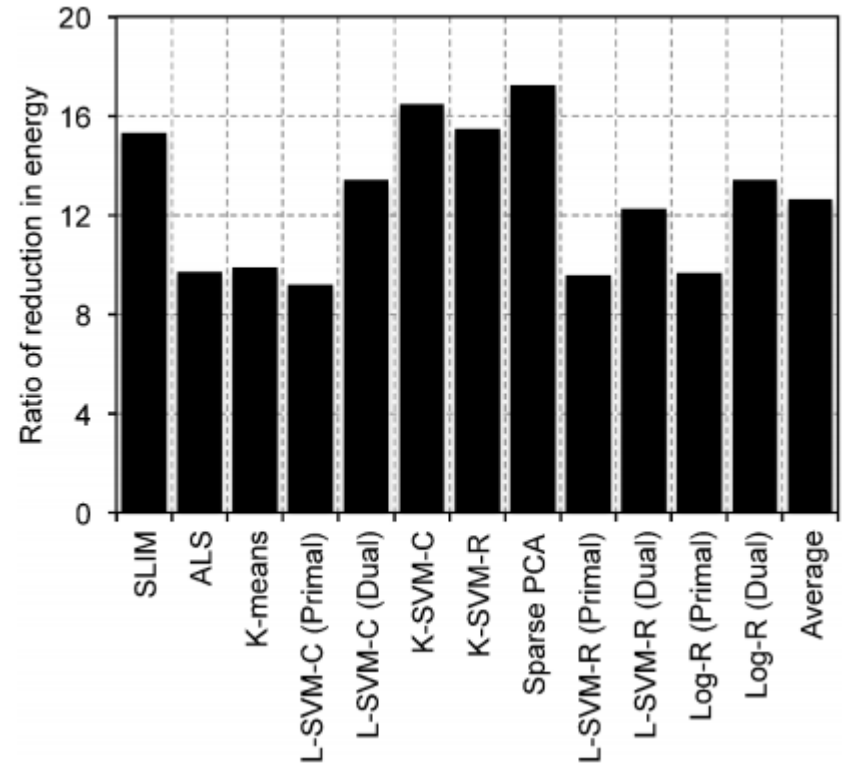
**~20x speedup for SpMSPv
(K-SVM, sparse PCA)**



**Due to parallel index matching &
custom datapath**

**SpMDV is bandwidth bound,
perf limited by mem system**

**Energy efficiency
improvements**



~12x avg better energy

**From turning off core and
parts of mem subsystem**

Talk Outline

- **Executive summary**
- **Sparse Machine Learning**
- **Sparse matrix processing**
- **Characterization study**
- **Proposed hardware accelerator**
- **Related work + summary**

Related Work

- Sparse matrix & sparse ML accelerator
 - Many proposals target only 1 sparse matrix format/op
 - Our previous work on sparse ML accelerator did not tightly integrate accelerator blocks with CPU

- Other ML accelerators
 - Many proposals for individual workloads
 - Many proposals for neural networks and/or dense data

Summary

- Sparse ML growing in importance
 - Sparsity from unstructured data (e.g., texts, ratings)
- We characterized various sparse ML workloads
 - Most runtime spent on sparse matrix op hotspots
- We proposed HW accelerator for these matrix ops
 - Tightly coupled with CPU and mem system
 - Improve efficiency dramatically (performance, energy)