



Sparsity in ML: Understanding and Optimizing Sparsity in Neural Networks Running on Heterogeneous Systems

Organizers

Wen-mei Hwu, Jinjun Xiong, Mert Hidayetoğlu, Rakesh Nagi,

Jeff Pool, Sitao Huang, Vikram Mailthody

Tutorial | MLSys 2022 | Aug 30, 2022

https://mlsys.org/virtual/2022/tutorial/2199

Slack for discussions: https://bit.ly/3KLQEUN

- 1. Opening remarks and overview of sparsity in ML 1 pm
- 2. Tiled SpMM and its performance model on GPUs 1:50 pm
- 3. Sparse deep neural network inference on FPGAs 2:45 pm
- 4. Break (30 mins)
- 5. 2:4 Sparsity on GPU Tensor Cores 4 pm
- 6. Future work and closing remarks 4:45 pm

Large models become increasingly popular and powerful



Source: Cerebras Systems, HotChips 2022

1000X in 2 years, trillion scale models just a year away

Compression to reduce model complexity and memory capacity

- 1. Down-sizing models: distillation, neural architectural search (NAS)
- 2. Operator factorization: tensor-decomposition $W[n \times m] = AB[n \times r][r \times m]$
- 3. Value quantization: low precision (TF32, FP16, FP8, INT8, Binary...)
- 4. Value compression: compression (gzip, Huffman, delta-encoding..)
- 5. Parameter sharing: reuse parameter across neuron (<u>Shapeshifter</u>)
- 6. Sparsification: input sparsification, prune nodes and weights (layer-wise, channel, random, regularizer,...)

Sparsity in ML



Sparse Attention Layers

Fixed Sparsity Structure

Global Attention

Loca

Attention -Window



Google ai, "Rethinking Attention with Performers,"

Child, et al., "Generating Long Sequences with Sparse Transformers," Arxiv

Convolutional Neural Networks



• Models with massive parameters do memorize input data

- Not all parameters matter (5% efficient connection can give similar accuracy as 100% model)
- Faster training with dense models, but sparse models have less overfitting problems

Extracting Training Data from Large Language Models

Nicholas Carlini	¹ Florian Tramèr ²	Eric Wallace ³	Matthew Jagielski ⁴
Ariel Herbert-Voss	5,6 Katherine Lee ¹	Adam Roberts ¹	Tom Brown ⁵
Dawn Song ³	Úlfar Erlingsson ⁷	Alina Oprea ⁴	Colin Raffel ¹
Google ² Stanford	³ UC Berkeley ⁴ Northeaste	rn University ⁵ Ope	enAI ⁶ Harvard ⁷ Apple

Abstract

It has become common to publish large (billion parameter) language models that have been trained on private datasets. This paper demonstrates that in such settings, an adversary can perform a *training data extraction attack* to recover individual training examples by querying the language model.

We demonstrate our attack on GPT-2, a language model trained on scrapes of the public Internet, and are able to extract hundreds of verbatim text sequences from the model's training data. These extracted examples include (public) personally identifiable information (names, phone numbers, and email addresses), IRC conversations, code, and 128-bit UUIDs. Our attack is possible even though each of the above sequences are included in just *one* document in the training data.

We comprehensively evaluate our extraction attack to understand the factors that contribute to its success. Worryingly, we find that larger models are more vulnerable than smaller models. We conclude by drawing lessons and discussing possible safeguards for training large language models.

1 Introduction



Figure 1: **Our extraction attack.** Given query access to a neural network language model, we extract an individual person's name, email address, phone number, fax number, and physical address. The example in this figure shows information that is all accurate so we redact it to protect privacy.

THE LOTTERY TICKET HYPOTHESIS: FINDING SPARSE, TRAINABLE NEURAL NETWORKS

Jonathan Frankle MIT CSAIL jfrankle@csail.mit.edu Michael Carbin MIT CSAIL mcarbin@csail.mit.edu

Abstract

"A randomly-initialized, dense neural network contains a subnetwork that is initialized such that — when trained in isolation — it can match the test accuracy of the original network after training for at most the same number of iterations"

We present an algorithm to identify winning tickets and a series of experiments that support the lottery ticket hypothesis and the importance of these fortuitous initializations. We consistently find winning tickets that are less than 10-20% of the size of several fully-connected and convolutional feed-forward architectures for MNIST and CIFAR10. Above this size, the winning tickets that we find learn faster than the original network and reach higher test accuracy.













Sparse computation in the form of SpMM

Sparse data access in the form of graph

Three levels of optimization opportunities for SpMM on heterogeneous systems

1. Optimize SpMM data structures and algorithms for generic hardware features



2. Co-optimize SpMM implementation for underlying hardware features







Tiled SpMM acceleration on GPUs – talk by Mert @ 1:50 pm

1. Optimize SpMM data structures and algorithms for generic hardware features



2. Co-optimize SpMM implementation for underlying hardware features





[Hidayetoglu et al: "At-Scale Sparse Deep Neural Network Inference With Efficient GPU Implementation," IEEE GraphChallenge 2020 (Champion)]



2. Co-optimize SpMM implementation for underlying hardware features





[Huang et al: "Accelerating Sparse Deep Neural Network Inference on FPGAs," IEEE GraphChallenge 2019 (Honorable Mention)]



2. Co-optimize SpMM implementation for underlying hardware features





3. Optimize hardware with innovative architecture features



Sparse data access in the form of graph

Feature tensors of big graphs can't fit into GPU's limited memory





Long data access latency and waste of host resources (consuming 3x more memory bandwidth)

[Ying et al., Graph Convolutional Neural Networks for Web-Scale Recommender Systems, KDD 2018]



State-of-the-art graph processing was NOT even able to fully utilize the PCIe bandwidth



[P. Gera, H. Kim, P. Sao, H. Kim, and D. Bader. "Traversing large graphs on GPUs with unified memory". VLDB 2020.]

[Min, Mailthody, Qureshi, Xiong, Ebrahimi, and Hwu., EMOGI: Efficient Memory-access for Out-of-memory Graph-traversal in GPUs, VLDB'21]





[Min, Mailthody, Qureshi, Xiong, Ebrahimi, and Hwu., EMOGI: Efficient Memory-access for Out-of-memory Graph-traversal in GPUs, VLDB'21]

Observation: interleaved and aligned assignment results in coalesced 128-byte PCIe packets

 Warp threads access consecutive 32 floats that align with the 128-byte cacheline granularity



[Min, Mailthody, Qureshi, Xiong, Ebrahimi, and Hwu., EMOGI: Efficient Memory-access for Out-of-memory Graph-traversal in GPUs, VLDB'21]

• Always guaranteeing 128-byte alignment for any data structure can be difficult





A warp may generate two separate PCIe requests to fetch a single 128-byte cacheline.

Existing GNN frameworks (e.g., PyTorch/DGL) do not support zero-copy

- PyTorch always keeps data right next to the processing unit
 - CPU tensor in host memory and GPU tensor in GPU memory
 - No ways for GPU operators to directly access host memory
- Zero-copy tensor is transient
 - Resulting tensors should be GPU tensors in order to be stored in device memory



Native PyTorch tensors are bound to devices.



- Introduces new "unified tensor" data types in DGL
 - To enable the zero-copy direct data access paradigm for PyTorch
 - Compact and intuitive APIs for seamlessly integration
- Two affinity modes: GPU-affinitive and CPU-affinitive
 - No actual data movements when the affinity is switched

Tensor type	Executor	Storage
CPU	CPU	CPU Memory
CUDA	GPU	GPU Memory
GPU-Affinitive Unified	GPU	CPU Memory
Host-Affinitive Unified	CPU	CPU Memory

Improvements of 16 – 44% for 1 GPU and 65 – 92% for 2 GPUs Freed up CPU for other tasks (less memory access and less power consumption)



Dataset	# Features	GB Size	# Nodes	# Edges
Wikipediaa	315	17.1	13.6 M	437.2 M
Amazon	578	34.0	14.7 M	64.0 M
ogbn-paper	128	57	111 M	1.6 B

System	
CPU	AMD Threadripper 3960X 24C/48T
GPU	2x Ampere RTX 3090 24 GB
Memory	DDR4 3200 MHz 256 GB
CUDA	11.2
PyTorch	1.8.0-nightly
OS	Ubuntu 20.04.1, Linux Kernel 5.8.0



[Min, Wu, Huang, Hidayetoglu, Xiong, Song, and Hwu., Graph Neural Network Training with Data Tiering, KDD'22]



[Min, Wu, Huang, Hidayetoglu, Xiong, Song, and Hwu., Graph Neural Network Training with Data Tiering, KDD'22]











[Min, Wu, Huang, Hidayetoglu, Xiong, Song, and Hwu., Graph Neural Network Training with Data Tiering, KDD'22]





Sparse computation in the form of SpMM

Sparse data access in the form of graph

- EMOGI: zero-copy to allow GPU-centric sparse data access (VLDB'21)
- PyTorch-Direct: unified tensor with zero-copy for GNN training (VLDB'21)
- Data-tiering: model-driven data placement for further performance (KDD'22)

- 1. Opening remarks and overview of sparsity in ML 1 pm
- 2. Tiled SpMM and its performance model on GPUs 1:50 pm
- 3. Sparse deep neural network inference on FPGAs 2:45 pm
- 4. Break (30 mins)
- 5. 2:4 Sparsity on GPU Tensor Cores 4 pm
- 6. Future work and closing remarks 4:45 pm





Tiled SpMM and its Performance Model on GPUs

NVIDIA. University at Buffalo

Mert Hidayetoglu University of Illinois at Urbana-Champaign

Organizers: Mert Hidayetoglu, Wen-mei Hwu, Jinjun Xiong, Rakesh Nagi, Jeff Pool, Sitao Huang, Vikram Mailthody

Tutorial | MLSys 2022 | Aug 30, 2022

https://mlsys.org/virtual/2022/tutorial/2199 Slack for discussions: https://bit.ly/3KLQEUN

- 1. Opening remarks and overview of sparsity in ML 1 pm
- 2. Tiled SpMM and its performance model on GPUs 1:50 pm
- 3. Sparse deep neural network inference on FPGAs 2:45 pm
- 4. Break (30 mins)
- 5. 2:4 Sparsity on GPU Tensor Cores 4 pm
- 6. Future work and closing remarks 4:45 pm

Sparsity in ML



Fundamental Matrix Multiplication Operations


















Fundamental Matrix Multiplication Operations



Sparse Connection Layers Graph Neural Networks Sparse Attention Layers

Sparse Neural Networks



[Lecun et al. NIPS'89] [Han et al. NIPS'15]





Google ai, "Rethinking Attention with Performers,"

Child, et al., "Generating Long Sequences with Sparse Transformers," Arxiv

Graph Neural Networks



Convolutional Neural Networks

Sparse Attention Layers

















The Roofline Model

Slide 47





Arithmetic Intensity

 $I = \frac{\sum \text{Floating Point Operations (GFLOP)}}{\sum \text{Memory Accesses (GB)}}$



Arithmetic Intensity

 $I = \frac{\sum \text{Floating Point Operations (GFLOP)}}{\sum \text{Memory Accesses (GB)}}$











Arithmetic Intensity

I = 0.17 FLOP/Byte



Applies Loop Transformations

Utilizes scratchpad and registers

Improves arithmetic intensity



Applies Loop Transformations

Utilizes scratchpad and registers

Reordering

Tiling Fusion

Improves arithmetic intensity



Applies Loop Transformations

Utilizes scratchpad and registers

Improves arithmetic intensity

Provide on-chip data reuse

Reduce DRAM memory traffic



Applies Loop Transformations

Utilizes scratchpad and registers

Improves arithmetic intensity



Modeling Arithmetic Intensity of Baseline SpMM

Slide 60





CSR SpMM (Baseline)

```
for (int k = 0; k < K; k++)  For K columns
for (int m = 0; m < M; m++)  For M Rows
{
    float acc = 0;
    for (int l = rowPtr[m]; l < rowPtr[m + 1]; l++)
    {
        int idx = index[l];
        float val = value[l];
        acc += B[idx][k] * val;
    }
    C[m][k] = acc;
}</pre>
```





CSR SpMM (Baseline)

```
for (int k = 0; k < K; k++)
for (int m = 0; m < M; m++)
{
    float acc = 0;
    for (int l = rowPtr[m]; l < rowPtr[m + 1]; l++)
    {
        int idx = index[l];
        float val = value[l];
        acc += B[idx][k] * val;
    }
    C[m][k] = acc;
}</pre>
```





CSR SpMM (Baseline)



$$I = \frac{2\mu}{(B_T + B_I)\mu + B_T\mu + B_T + \delta}$$

Algorithmic Intensity



Legend

 μ : # of nonzeroes per row

$$I = \frac{2\mu}{(B_T + B_I)\mu + B_T\mu + B_T + \delta}$$

Algorithmic Intensity





Algorithmic Intensity





Read A







 μ : # of nonzeroes per row B_T : Bytes per value

 B_I : Bytes per index

















second





Baseline Sparse Layer Implementation

- Activation features: dense (col.-major)
- Gather approach
- Each thread computes a single output


Baseline Sparse Layer Implementation

- Activation features: dense (col.-major)
- Gather approach
- Each thread computes a single output
- Sparse weights: CSR
- ReLU is fused with SpMM



Baseline Sparse Layer Implementation

- Activation features: dense (col.-major)
- Gather approach
- Each thread computes a single output
- Sparse weights: CSR
- ReLU is fused with SpMM

Data Access Redundancies

- Weight matrix by all output features
- Input features by different threads



Baseline Sparse Layer Implementation

- Activation features: dense (col.-major)
- Gather approach
- Each thread computes a single output
- Sparse weights: CSR
- ReLU is fused with SpMM

Data Access Redundancies

- Weight matrix by all output features
- Input features by different threads

Data Access Latencies

- Irregular access to input features
- Uncoalesced access to weight matrix



MLSys Tutorial: Sparsity in ML



SpMM Optimization for ML Register Tiling Scratchpad Tiling Fusing Activation

Multi-Level Input Buffering



Hidayetoglu, et al., "At-Scale Deep Neural Network Inference with Efficient GPU Implementation," Arxiv and HPEC'20.

MLSys Tutorial: Sparsity in ML



Multi-Level Input Buffering

Intermediate Data Structures

Multi-Level Input Buffering





Intermediate Data Structures

Multi-Level Input Buffering

Weights Sliced ELLPACK





wdispl windex



53

Tiled SpMM acceleration on GPUs – talk by Mert @ 1:50 pm

1. Optimize SpMM data structures and algorithms for generic hardware features



2. Co-optimize SpMM implementation for underlying hardware features





[Hidayetoglu et al: "At-Scale Sparse Deep Neural Network Inference With Efficient GPU Implementation," IEEE GraphChallenge 2020 (Champion)]

MIT/Amazon/IEEE Graph Challenge

grapchallenge.mit.edu



(a) Fully connected layer;(b) Sparse layer (after pruning)



A sparse DNN with 6 layers, 64 neurons per layer (synthetic, 2 connections per neuron) [1]

[1] J. Kepner et al. Sparse Deep Neural Network Graph Challenge. HPEC 2019.

MIT/Amazon/IEEE Graph Challenge

grapchallenge.mit.edu

Challenges

- Irregular: Frequent random accesses to memory
- Various Scales: from 120 layers, 1024 neurons/layer to 1920 layers 65536 neurons/layer (16.3 GB)

Computation

- Input and parameters are in graph format (sparse)
- CPU performs data preprocessing
- GPU performs DNN inference: ℓ -th layer:

$$\boldsymbol{Y}_{\ell+1} = h(\boldsymbol{W}_{\ell}\boldsymbol{Y}_{\ell} + \boldsymbol{B}_{\ell})$$



Memory Optimizations

- short indices
- half weights (not used for Graph Challenge)
- Compact struct (not used for Graph Challenge)
- Batching for features
- Out-of-core streaming for weights





MIT/Amazon/IEEE Graph Challenge

grapchallenge.mit.edu

HPEC'20 Champion		Benchmark Throughput (# edges per second)										
	This Work		Bisson & Fatica [18]		Davis et al. [20]		Ellis & Rajamanickam [21]		Wang et al. [22]		Wang et al. [23]	
			2019 Champion		2019 Champion		2019 Innovation		2019 Student Innov.		2019 Finalist	
Neurons	Layers	Throughput	Throughput	Speedup	Throughput	Speedup	Throughput	Speedup	Throughput	Speedup	Throughput	Speedup
	120	2.917E+13	4.517E+12	6.46	1.533E+11	190.28	2.760E+11	105.69	1.407E+11	207.32	8.434E+10	345.88
1024	480	2.930E+13	7.703E+12	3.80	2.935E+11	99.83	2.800E+11	104.64	1.781E+11	164.51	9.643E+10	303.84
	1920	2.883E+13	8.878E+12	3.25	2.754E+11	104.68	2.800E+11	102.96	1.896E+11	152.06	9.600E+10	300.30
	120	8.220E+13	6.541E+12	12.57	1.388E+11	592.22	2.120E+11	387.74	1.943E+11	423.06	6.506E+10	1,263.52
4096	480	8.222E+13	1.231E+13	6.68	1.743E+11	471.72	2.160E+11	380.65	2.141E+11	384.03	6.679E+10	1,230.99
	1920	8.232E+13	1.483E+13	5.55	1.863E+11	441.87	2.160E+11	381.11	2.197E+11	374.69	6.617E+10	1,244.02
	120	1.469E+14	1.008E+13	14.57	1.048E+11	1,401.53	1.270E+11	1,156.54	1.966E+11	747.10	3.797E+10	3,867.84
16384	480	1.394E+14	1.500E+13	9.29	1.156E+11	1,206.23	1.280E+11	1,089.38	2.060E+11	676.89	3.747E+10	3,721.66
	1920	1.464E+14	1.670E+13	8.77	1.203E+11	1,216.96	1.310E+11	1,117.56	1.964E+11	745.52	3.750E+10	3,903.72
	120	1.796E+14	9.388E+12	19.13	1.050E+11	1710.29	9.110E+10	1971.24	1.892E11	949.15	-	-
65536	480	1.703E+14	1.638E+13	10.40	1.091E+11	1,560.59	8.580E+10	1,984.38	1.799E+11	946.41	-	-
	1920	1.714E+14	1.787E+13	9.59	1.127E+11	1,520.59	8.430E+10	2,032.86	_	-	-	-

Generalization of Tiled SpMM





2,893 Sparse Matrices Picked the largest 200 samples in

terms of nonzeroes

The matrices in the Collection cover a wide spectrum of domains, including

- 2D or 3D Geometric Domains
 - Structural engineering
 - Computational fluid dynamics
 - Model reduction
 - Electromagnetics
- Non-Geometric Domains
 - Optimization
 - Economic and financial modeling
 - Circuit simulation

- Semiconductor devices
- Thermodynamics
- Material science
- Acoustics

- Computer graphics and vision
- Robotics and kinematics
- Other spatial discretizations

- Theoretical and quantum chemistry
- Chemical process simulation
- Mathematics and statistics

- Power networks
- Other networks and graph-structured data

Tiled SpMM on Roofline Plot





MLSys Tutorial: Sparsity in ML



Throughput on A100 (FP32)





Conclusion

Observation: the bottleneck of sparse computations is data movement over the slowest interconnect, rather than arithmetic operations.



Observation: the bottleneck of sparse computations is data movement over the slowest interconnect, rather than arithmetic operations.

This Work: With the proposed tiling techniques and a careful orchestration of data movement over the memory hierarchy, we can accelerate SpMM on GPUs to a great extent. DRAM->L2\$->L1\$/scratchpad->register **Roofline Model** A100 25,000 19.5 TFLOP/s 20,000 Throughput (GFLOP/s) 12,000 10,000 10,000 \star GBIS DGEMM 55 *****Tiled SpMM 5,000 *****SpMM

10

15 Arithmetic Intensity (FLOP/Byte)

20

25

30

Observation: the bottleneck of sparse computations is data movement over the slowest interconnect, rather than arithmetic operations.

This Work: With the proposed tiling techniques and a careful orchestration of data movement over the memory hierarchy, we can accelerate SpMM on GPUs to a great extent. Work in progress

We developed 'Tiled SpMM' and ...

- Benchmarked over 200 sparse matrices from a wide-spectrum of applications
- Verified performance bounds with the **proposed** analytical model
- Obtained 16x speedup over cuSPARSE

Product: Tiled SpMM is open source for other application developers.

https://github.com/merthidayetoglu/SpDNN_Challenge2020



Accelerating Sparse Deep Neural Network Inference on FPGAs

NVIDIA. University at Buffalo

Sitao Huang

Assistant Professor, EECS, University of California Irvine

Organizers: Mert Hidayetoğlu, Wen-mei Hwu, Jinjun Xiong, Rakesh Nagi,

Jeff Pool, Sitao Huang, Vikram Mailthody

Tutorial | MLSys 2022 | Aug 30, 2022

https://mlsys.org/virtual/2022/tutorial/2199

- 1. Opening remarks and overview of sparsity in ML 1 pm
- 2. Tiled SpMM and its performance model on GPUs 1:50 pm
- 3. Sparse deep neural network inference on FPGAs 2:45 pm
- 4. Break (30 mins)
- 5. 2:4 Sparsity on GPU Tensor Cores 4 pm
- 6. Future work and closing remarks 4:45 pm

Three levels of optimization opportunities for SpMM on heterogeneous systems

1. Optimize SpMM data structures and algorithms for generic hardware features



2. Co-optimize SpMM implementation for underlying hardware features

3. Optimize Hardware with innovative architecture features





Three levels of optimization opportunities for SpMM on heterogeneous systems

1. Optimize SpMM data structures and algorithms for generic hardware features



2. Co-optimize SpMM implementation for underlying hardware features





In This Talk

- Why FPGA?
- A Sparse DNN Inference Accelerator on FPGA
- PyLog: Python-based FPGA Design Flow
- Custom Sparse DNN Accelerators
- Future Topics

Why FPGAs?

FPGA: Field Programmable Gate Array

- Customizable, reconfigurable logic
- Low-latency, energy-efficient solutions
- Available in many different sizes



- Cloud computing, HPC, smart phone, edge devices, etc.
- Wide applications: medical, video processing, telecom, space, etc.



A Sparse DNN Inference Accelerator on FPGA

• Accelerating *sparse* large-scale DNNs on FPGAs



(a) Fully connected layer;(b) Sparse layer (after pruning)



A sparse DNN with 6 layers, 64 neurons per layer (synthetic, 2 connections per neuron) [1]

[1] J. Kepner et al. Sparse Deep Neural Network Graph Challenge. HPEC 2019.

A Sparse DNN Inference Accelerator on FPGA

Challenges

- Irregular: Frequent random accesses to memory
- Large-scale: from 120 layers, 1024 neurons/layer to 1920 layers 65536 neurons/layer (16.3 GB)

Computation

- Input and parameters are in graph format (sparse)
- Host CPU performs data preprocessing
- FPGA performs DNN inference: ℓ -th layer: $\mathbf{Y}_{\ell+1} = h(\mathbf{Y}_{\ell}\mathbf{W}_{\ell} + \mathbf{B}_{\ell})$

 $N_a \rightarrow N_b$: edge in graph

	Na	N_b	weight		Na	N _b	weight
input.tsv:	1323	10	1	param-L1.tsv:	1	1	6.250000e-02
	55292	12	1		2	1	6.250000e-02
	44323	14	1		65	1	6.250000e-02
	1323	41	1		66	1	6.250000e-02
	1323	42	1		129	1	6.250000e-02
	• • •	• • •	• • •			• • •	• • •

FPGA Accelerator Design

FPGA Design Challenges

- Limited on-chip resources (memory size, ports, DSP, etc.)
- Limited memory bandwidth
- Programmability

Our Solution

- Coarse-grained asynchronous accelerator pool
- Sparse dot product units with buffers
- Fully pipelined datapath
- Design created with C/C++ high-level synthesis (HLS)



Sitao Huang, Carl Pearson, Rakesh Nagi, Jinjun Xiong, Deming Chen, and Wen-mei Hwu. *Accelerating Sparse Deep Neural Network on FPGA*. HPEC 2019.

Grid Representation

- Rows: input samples (images), independent
- Columns: DNN layers, dependent



Computation illustrated with a 2D grid

Ping-Pong Buffering (Double Buffering)

- For both *data loading* and *computation*
- Alternate load and use buffers (Buffer A and B)





Other buffering optimizations in DNN accelerators:

- CNNs: fine-grained crosslayer pipelining
- Memory block partitioning to increase # of ports

Computation illustrated with a 2D grid

Feature vectors and parameters: use sparse or dense format?

- Treating both as sparse: SpMV becomes set intersection
- Access pattern in parameters is regular while pattern in feature vectors is less regular
- Design decision: Feature vectors: dense format; parameters: sparse format



Multi-level tiling

- 1 across input samples (input batch)
- Pused layers (inter-layer)
- ③ within a layer, across neurons (intra-layer)





Computation illustrated with a 2D grid

Dynamic workload balancing

- Multiple accelerators instantiated (at block design level)
- CPU dynamically assigns jobs to idle accelerators and collect returned results



Algorithm 1 Dynamic Workload Assignment
Input: Number of input samples N, pack size S, accelerator
pool P = {P[0],, P[m-1]}.
1: curr_img $\leftarrow 0$, acc_ptr $\leftarrow 0$
2: size $\leftarrow MIN(S, N - curr_img)$
3: while curr_img $< N$ do
4: if P[acc_ptr].ISIDLE() then
5: ASSIGN(curr_img, size, P[acc_ptr])
6: $curr_img \leftarrow curr_img + size$
7: else if P[acc_ptr].IsDone() then
8: COLLECTRESULTS(P[acc_ptr])
9: ASSIGN(curr_img, size, P[acc_ptr])
10: $curr_img \leftarrow curr_img+size$
11: end if
12: $acc_ptr \leftarrow (acc_ptr+1) \%$
13: end while

FPGA Accelerator Evaluation

- FPGA: Xilinx Virtex-7 VC709 board (250MHz)
- CPU baseline: AMD Operon 6272 Processors (64 cores), MATLAB Sparse Matrix library
- Dataset: 1024 neurons per layer x 120 layers (32x32 image classification)



Sitao Huang, EECS, UC Irvine

High-Level Synthesis

- **High-level synthesis (HLS)**, also referred to as C synthesis, electronic system-level (ESL) synthesis, algorithmic synthesis, or behavior synthesis
- HLS takes an abstract behavioral specification of a digital system, and finds a register-transfer level structure that realizes the given behavior

INPUT:

- A high-level, algorithmic description
 - Control structures (if/else, loop, subroutines)
 - Concurrent and sequential semantics
 - Abstract data types
 - Logical and arithmetic operators
- A set of constraints
 - Speed, power, area, interconnect style
 - A library of pre-specified components

Reference: Register-Transfer Level (RTL), Abstract Data Type (ADT)

OUTPUT:

• A register-transfer level description for further synthesis and optimization
High-Level Synthesis Process

Resource Allocation:

 Allocating resources (library components) to each of the operations, buses, MUXes, and registers for storage

Scheduling:

 Scheduling the operations in the control dataflow graph (CDFG) to minimize area, time and/or power

Binding:

- Determining the time of use of each component
- e.g., which registers to use and when



PyLog: Python-based FPGA Design Flow



PyLog Language and Compiler Highlights

High-level operators, design space exploration



given the HW resource of the target platform

Python statements, compute customization

```
@pylog
def pl_matmul(a, b, c, d):
    buf = np.empty([16, 16], pl_fixed(8, 3))
    pragma("HLS array_partition variable=buf")
    def matmul(a, b, c):
        for i in range(32):
            for j in range(32).unroll(4):
               tmp = 0.0
               for k in range(32).pipeline():
                   tmp += a[i][k] * b[k][j]
               c[i][j] = tmp
```

Unified FPGA and host programming

```
import numpy as np
from pylog import *
```



Type inference and type checking

- No explicit type annotation required
- Top function has NumPy arrays as arguments, which carries input type and shape information
- Type engine infers type and shape of each object in the code

PyLog High-Level Operations

- map operation: map(f, x0, x1,...)
 - Repeatedly apply function f to each element in *x0, x1,...*
 - Allow access neighbor elements inside function
- Dot operation: dot product of two arrays

Example: vector addition (a and b have same shape) $z = map(Lambda \ a, \ b: \ a + b, \ x, \ y)$

(a) 2D map PyLog
$$y = map(Lambda a_0, a_1, ...; op(a_0[-1, 0], a_1[0, 2], ...), x_0, x_1, ...)$$

L1: for(int i1 = 0; i1 < x_0 .dim(0); i1++) {
HLS C L2: for(int i2 = 0; i2 < x_0 .dim(1); i2++) {
 $y[i1][i2] = op(x_0[i1-1][i2], x_1[i1][i2+2], ...); }}$

(b) 1D conv PyLog
$$y = map(Lambda a: a[-1]+a[0]+a[1], x[1:-1])$$

HLS C $L1: for(int i1 = 1; i1 < x.dim(0)-1; i1++) \{ y[i1] = x[i1-1] + x[i1] + x[i1+1]; \}$

PyLog Supported Platforms

Platforms supported by PyLog

FPGA Platform Type	Platforms	Synthesis Flow	Runtime Library
SoCs and MPSoCs	ZedBoard, PYNQ, Ultra96	Vivado HLS + Vivado	PYNQ
PCIe-Based High-End FPGAs	Amazon EC2 F1 instance (XCVU9p), Alveo series (U200, U250, U280)	Vivado HLS + Vitis	PYNQ

- Design and deploy for different platforms easily
 - Simply change @pylog mode, no extra coding needed

Synthesize for AWS F1: @pylog(mode='hwgen', board='aws_f1')
Deploy on AWS F1: @pylog(mode='deploy', board='aws_f1')

Synthesize for PYNQ: @pylog(mode='hwgen', board='pynq')
Deploy on PYNQ: @pylog(mode='deploy', board='pynq')





PyLog Evaluation: Expressiveness

- Compare LoC (Lines of Code) of HLS C code vs PyLog code
 - Normalized to HLS C LOC
 - 70% reduction compared to HLS C code



*Yi-Hsiang Lai, Yuze Chi, Yuwei Hu, Jie Wang, Cody Hao Yu, Yuan Zhou, Jason Cong, and Zhiru Zhang. HeteroCL: A Multi-Paradigm Programming Infrastructure for Software-Defined Reconfigurable Computing. FPGA 2019.

PyLog Evaluation: Accelerator Performance

Compare accelerator performance vs CPU (C++) performance and HeteroCL accelerator performance

- Platform: Amazon AWS F1 instance (Xilinx Virtex UltraScale+ XCVU9P)
- 3.17x and 1.24x faster than optimized CPU and FPGA accelerators



CPU: Optimized CPU code running on single thread on 8-core Intel Xeon E5-2686 v4 CPU

SpMV in PyLog

Baseline version

```
@pylog
def spmv(row_ptr, column_idx, values, y, x):
    for i in range(num_neuron):
        y0 = 0
        for k in range(row_ptr[i], row_ptr[i+1]):
            y0 += values[k] * x[column_idx[k]]
        y[i] = y0
```

SpMV in PyLog

• (lightly) optimized version

```
@pylog
def spmv(row_ptr, column_idx, values, y, x):
    for i in range(64).pipeline():
        v0 = 0
        for k in range(row ptr[i], row ptr[i+1], 7):
            pragma('HLS pipeline II=7')
            yt = values[k] * x[column idx[k]]
            for j in range(7):
                if ((k+j) < row ptr[i+1]):</pre>
                    yt += values[k+j] *x[column_idx[k+j]]
            y0 += yt
        y[i] = y0
```

PyLog is open-source. Give it a try!

PyLog usage

To use PyLog, import pylog and simply add PyLog decorator <code>@pylog</code> to the function that you'd like to synthesize into an FPGA accelerator. Pass NumPy arrays to the decorated function and call the decorated function. Then run the whole Python program. In the following example, <code>vecadd</code> functior



https://github.com/hst10/pylog

PyLog 2.0: New Features

PyLog x Heterogeneous Computing



New heterogeneous platforms



- Returns a new function that combines kernels
- Connects inputs and outputs with kernels

a b kernel 1 d c kernel 2

Xilinx Versal devices

Collaborative computing:

- Across multiple adaptable logic kernels
- Across scalar engines, adaptable engines, and AIEs
- Across multiple ACAP / FPGA devices Data placement and movement:
- Data sources: PCIe, DDR, HBM, Ethernet, other engines, etc.

Generic programmable logic: Vitis connectivity configuration

[connectivity] Creating Multiple Instances of a Kernel #nk=<kernel name>:<number>:<cu_name>... nk=vadd:3:vadd_X.vadd_Y.vadd_Z

Mapping Kernel Ports to Global Memory #sp=<compute unit name>.<interface name>:<bank name>

sp=ccompute_unit_name>.<internate_name>.
sp=cnn_1.m_axi_gmem:DDR[0]

Specify streaming connections

#stream_connect=<cu_name>.<output_port>:<cu_name>.<input_port>
stream_connect=vadd_1.stream_out:vadd_2.stream_in

AI Engines: Dataflow model (on-going)



// create nets to connect kernels and IO ports
connect<window<128> > net0 (in, k1.in[0]);
connect<window<128> > net1 (k1.out[0], k2.in[0]);
connect<window<128> > net2 (k2.out[0], out);
};

Other Optimizations for Sparse Accelerators

- Optimizations for HBM [1, 2]
 - HBM bandwidth utilization
 - On-chip memory utilization
 - Compute occupancy
- New compute patterns [3]
 - Outer product based GEMM
 - Huffman tree scheduler for balanced reduction





Yixiao Du, Yuwei Hu, Zhongchun Zhou, Zhiru Zhang. High-Performance Sparse Linear Algebra on HBM-Equipped FPGAs Using HLS: A Case Study on SpMV. FPGA 2022.
 Linghao Song, Yuze Chi, Licheng Guo, and Jason Cong. Serpens: A High Bandwidth Memory Based Accelerator for General-Purpose Sparse Matrix-Vector Multiplication. DAC 2022.
 Zhekai Zhang, Hanrui Wang, Song Han, William J. Dally. SpArch: Efficient Architecture for Sparse Matrix Multiplication. HPCA 2020.

Future Works

- New optimizations for new architecture features
 - HBM, CXL, etc.
 - Hard accelerator cores
- Higher level programming abstractions and synthesis flow
 - Python-based programming to improve programmability and DSE efficiency
- System-level optimizations for heterogeneous platforms
 - Extending optimizations to system level
 - Support heterogeneous environments with diverse accelerators

Summary

- Optimizing sparse computation on custom accelerators
 - Understand platform strength and limitations
 - Balance memory vs compute
 - Exploring design space of multiple dimensions
- High-level programming flows improves productivity
 - PyLog: Python-based FPGA programming flow
 - PyLog 2.0: supporting heterogenous systems
 - Opens up new dimensions of design space
- Future directions







Jeff Pool and Rakesh Nagi

jpool@nvidia.com

nagi@illinois.edu

OVIDIA. University at Buffalo

Tutorial | MLSys 2022 | Aug 30, 2022

https://mlsys.org/virtual/2022/tutorial/2199

Outline



- 1. 2:4 Structured Sparsity Why and How
- 2. Performance
- 3. Maintaining Accuracy
- 4. Channel Permutations
- 5. Q&A



Structure:

- Unstructured: irregular, no pattern of zeros
- Structured: regular, fixed set of patterns





Granularity:

- Finest: prune individual values
- Coarser: prune blocks of values
- Coarsest: prune entire layers

Ц		_				
\vdash		_				
⊢		_				
⊢		-				
⊢					-	
┝╋		-			H	Η
H						
H						





Lots of research in two areas:

- High amounts (80-95%) unstructured, fine-grained sparsity
- Coarse-grained sparsity for simpler acceleration

Challenges not resolved for these approaches:

Accuracy loss

High sparsity often leads to accuracy loss of a few percentage points, even after advanced training techniques
Absence of a training approach that works across different tasks and networks (Recipe)

- •Training approaches to recover accuracy vary from network to network, often require hyper-parameter searches
- Lack of **speedup**
 - •Math: unstructured data struggles to take advantage of modern vector/matrix math instructions
 - •Memory access: unstructured data tends to poorly utilize memory buses, increases latency due to dependent sequences of reads
 - •Storage overheads: metadata can consume 2x more storage than non-zero weights, undoing some of compression benefits

Fine-grained structured sparsity for Tensor Cores

- 50% fine-grained sparsity
- •2:4 pattern: 2 values out of 4 must be zero

Addresses the three challenges:

•Accuracy: maintains accuracy of the original, unpruned network

- •Medium sparsity level (50%), fine-grained
- •Recipe: shown to work across tasks and networks

•Speedup

- Math: Specialized Tensor Core support for sparse math
- •Structured nature allows aligned, regular memory accesses and minimal metadata overheads









Compressed Matrix:

- Data: ½ size
- •Metadata:

2b per nonzero element (index into its location in the 4-element region) 16b data: 12.5% overhead 8b data: 25% overhead



Applicable for:

Convolutions

•GEMMs of all types (linear layers, MLPs, recurrent cells, transformer blocks, ...)

Inputs: sparse (weights), dense (activations) Outputs: dense (activations)





- 1. 2:4 Structured Sparsity Why and How
- 2. Performance
 - 1. Tensor Core
 - 2. **GEMM Operation**
 - 3. Full Network



Sparse Tensor Cores are 2x Faster

INPUT OPERANDS	ACCUMULATOR	A100 TOPS	Dense vs. FFMA	Sparse Vs. FFMA
FP32	FP32	19.5	-	-
TF32	FP32	156	8X	16X
FP16	FP32	312	16X	32X
BF16	FP32	312	16X	32X
FP16	FP16	312	16X	32X
INT8	INT32	624	32X	64X



Sparse GEMMs are up to 1.95x Faster

- Snapshot of cuSPARSELt and cuBLAS performance
- A100, 1410MHz



INT8 (TN) cuSPARSELt vs. cuBLAS Performance GEMM-M = GEMM-N = 10240



A realtime detection network is 1.56x faster

- RetinaNet with a ResNet-34 backbone
- 1920x1280 input, batch size = 1
- NVIDIA Orin SOC, 1275MHz

Dense Latency (ms)	Sparse Latency (ms)	Speedup
15.66	10.07	1.56x

Outline



- 1. 2:4 Structured Sparsity Why and How
- 2. Performance
- 3. Maintaining Accuracy
 - 1. Goals
 - 2. Recipe
 - 3. Results



- **1. Maintain accuracy**
- 2. Apply to various tasks, network architectures, and optimizers
- 3. Do not require hyperparameter searches

We're sharing a simple recipe we've found that satisfies these criteria – others exist!

MAccelerating Sparse Deep Neural Networks, Mishra et al., https://arxiv.org/abs/2104.08378



1. Train (or obtain) a dense network

2. Prune for 2:4 sparsity

- By magnitude, set 2 out of every 4 weights to zero
- Each group of 4 weights is independent
- Store a mask of which values were set to zero

3. Repeat the original training procedure

- Same hyperparameters as step (1)
- Initialize weights from step (2)
- Maintain the 0 pattern from step (2), no need to recompute the mask



ASP: Automatic SParsity (PyTorch)

- Part of the APEX library
- Three lines of code add 2:4 sparsity to your final network

```
import torch
from apex.contrib.sparsity import ASP
device = torch.device('cuda')
model = TheModelClass(*args, **kwargs) # Define model structure
model.load_state_dict(torch.load('dense_model.pth'))
optimizer = optim.SGD(model.parameters(), lr=0.01, momentum=0.9) # Define optimizer
ASP.prune_trained_model(model, optimizer)
x, y = DataLoader(...) #load data samples and labels to train the model
for t in range(500):
    y_pred = model(x)
    loss = loss_fn(y_pred, y)
    optimizer.zero_grad()
    loss.backward()
    optimizer.step()
```

torch.save(model.state_dict(), 'pruned_model.pth') # checkpoint has weights and masks





Network		Accuracy				
Network	Dense FP16	Sparse FP16	Sparse INT8	Sparse INT8		
ResNet-34	73.7	73.9 0.2	73.7 -	-		
ResNet-50	76.6	76.8 0.2	76.8 0.	.2		
ResNet-101	77.7	78.0 0.3	77.9 -	-		
ResNeXt-50-32x4d	77.6	77.7 0.1	77.7 -	-		
ResNeXt-101-32x16d	79.7	79.9 0.2	79.9 0.	.2		
DenseNet-121	75.5	75.3 -0.2	75.3 -0	.2		
DenseNet-161	78.8	78.8 -	78.9 0.	.1		
Wide ResNet-50	78.5	78.6 0.1	78.5 -	-		
Wide ResNet-101	78.9	79.2 0.3	79.1 0.	.2		
Inception v3	77.1	77.1 -	77.1 -	-		
Xception	79.2	79.2 -	79.2 -	-		
VGG-16	74.0	74.1 0.1	74.1 0.	.1		
VGG-19	75.0	75.0 -	75.0 -	-		

Maintaining Accuracy – Our Recipe's Results (COCO 2017, bbox AP)



Natwork		Accuracy		
Network	Dense FP16	Sparse FP16	Sparse INT8	
MaskRCNN-RN50	37.9	37.9 -	37.8 -0.1	
SSD-RN50	24.8	24.8 -	24.9 0.1	
FasterRCNN-RN50-FPN-1x	37.6	38.6 1.0	38.4 0.8	
FasterRCNN-RN50-FPN-3x	39.8	39.9 -0.1	39.4 -0.4	
FasterRCNN-RN101-FPN-3x	41.9	42.0 0.1	41.8 -0.1	
MaskRCNN-RN50-FPN-1x	39.9	40.3 0.4	40.0 0.1	
MaskRCNN-RN50-FPN-3x	40.6	40.7 0.1	40.4 -0.2	
MaskRCNN-RN101-FPN-3x	42.9	43.2 0.3	42.8 -0.1	
RetinaNet-RN50-FPN-1x	36.4	37.4 1.0	37.2 0.8	
RPN-RN50-FPN-1x	45.8	45.6 -0.2	45.5 -0.3	

Many more in our whitepaper...

- Generative tasks (GANs)
- Translation
- Language modeling

Outline



- 1. 2:4 Structured Sparsity Why and How
- 2. Performance
- 3. Maintaining Accuracy
- 4. Channel Permutations
 - 1. Problem/Solution
 - 2. Applying Permutations
 - 3. Finding Permutations
 - 4. Results



Small, parameter-efficient models can lose accuracy after pruning

ILSVRC'12	Default 2:4		
Network	Δ Top 1	Δ Top 5	
MobileNet v2	-1.98	-1.13	
MobileNet v3 (Small)	-2.73	-1.67	
SqueezeNet v1.0	-4.01	-2.96	
SqueezeNet v1.1	-1.25	-0.93	
ShuffleNet v2	-1.35	-0.87	
EfficientNet B0	-1.27	-0.52	



Column (channel) permutations allow better pruning choices





C is part of the inner dimension of a GEMM, so permuting that dimension of both operands does not change the result

- Permuting the *C* dimension of the **weights** is done offline
- Permuting the *C* dimension of layer *i*'s **activations** can be done by permuting the *K* dimension (output channels) of layer *i*-1, which is accomplished by permuting the weights' rows (offline)
- 1. Maintain magnitude with a desired weight permutation
- 2. Match the new layout in the corresponding activations
- 3. Cause the permutation in (2)





Applying permutations to a whole network requires care



M Channel Chermatations for N:M Sparsity, Pool and Yu, NeurIPS 2021


There are many, many unique permutations with N:M sparsity for even modest channel counts.

We adapted and improved past work¹, resulting in a heuristic approach that performs well:

- Gives good quality permutations (minimizes lost weight magnitude)
- Converges quickly in practice
- Scalable: trade off solution quality for search time

¹TETRIS: TilE-matching the TRemendous Irregular Sparsity, Ji et al., NeurIPS 2018



Some terminology:

- Stripe: vertical slice of a matrix consisting of one M-wide group (4 values in 2:4 sparsity)
- Stripe Group: a group of stripes (e.g., {0,1}, {0,3}, {2,3,6}, etc.) concatenated together





• Brute-force: a stripe group's best permutation is found exhaustively

2-stripe groups =8 columns =35 unique permutations3-stripe groups =12 columns =5,775 unique permutations4-stripe groups =16 columns =2,627,625 unique permutations



- Brute-force: a stripe group's best permutation is found exhaustively
- Divide-and-conquer: all potential stripe groups are brute-forced independently Improvement in post-pruning magnitude is collected for all groups: {0,1} = +1.4 {0,2} = +0.3
 - $\{1,2\} = +1.5$



- Brute-force: a stripe group's best permutation is found exhaustively
- **Divide-and-conquer**: all potential stripe groups are brute-forced independently
- **Greedy**: the stripe group with the largest improvement is selected That group's permutation is applied to that group Columns outside of the winning group are unchanged



- Brute-force: a stripe group's best permutation is found exhaustively
- **Divide-and-conquer**: all potential stripe groups are brute-forced independently
- Greedy: the stripe group with the largest improvement is selected

The greedy step changes the columns that constitute many groups, so we iterate:

- 1. Generate all stripe groups
- 2. Divide-and-conquer the groups
- 3. If there is a stripe group with a nonzero improvement:
 - 1. Take the greedy step with the winner
 - 2. GOTO(1)
- 4. Converged!



This heuristic converges, but it's probably to a local optimum.

Bounded regressions can find a better solution:

- After convergence, choose two columns at random and swap them
- This permutation will be worse but that's okay!
- One iteration of the greedy heuristic will either:
 - Find the previous solution and swap the columns back
 - Find a permutation in a stripe group with a better improvement, and repeat until converging at a better solution

We have two knobs to trade solution quality with search time:

- 1. The size of the stripe groups to optimize (default: 2 groups = 8 columns)
- 2. The number of bounded regression attempts to allow (default: 100)



Optimal Found: how many times the strategy found the optimal permutation





Optimal Found: how many times the strategy found the optimal permutation



Quality of Permutation Search Strategies



Optimal Found: how many times the strategy found the optimal permutation





Optimal Found: how many times the strategy found the optimal permutation





Optimal Found: how many times the strategy found the optimal permutation



Increase escape attempts **and** the size of the stripe groups for even better solutions.

Quality of Permutation Search Strategies



	Baseline (ILSVRC2012)		Default 2:4		Permuted 2:4 (ours)	
Network	Top 1	Тор 5	∆ Тор 1	∆ Тор 5	Δ Тор 1	∆ Тор 5
MobileNet v2	71.55	90.28	-1.98	-1.13	0.01	0.02
MobileNet v3 (Small)	67.67	87.40	-2.73	-1.67	0.10	0.15
MobileNet v3 (Large)	74.04	91.34	-0.91	-0.40	0.10	0.04
SqueezeNet v1.0	58.09	80.42	-4.01	-2.96	0.60	0.54
SqueezeNet v1.1	58.21	80.62	-1.25	-0.93	0.03	0.05
MNASNet 1.0	73.24	91.36	-1.25	-0.58	0.02	0.00
ShuffleNet v2	68.32	88.36	-1.35	-0.87	0.10	0.01
EfficientNet B0	77.25	93.59	-1.27	-0.52	0.04	0.07
ResNet-50	76.16	92.88	0.05	0.13	0.13	0.26
ResNeXt-50	77.62	93.70	0.01	0.05	0.13	0.07
DenseNet-161	77.14	93.56	0.82	0.41	0.92	0.52



Part 2: 2:4 (N:M) Sparsity

Rakesh Nagi

Department of Industrial and Enterprise Systems Engineering University of Illinois at Urbana-Champaign



- Sequence Approach
 - Form a sequence of *C* elements = *C*!
 - But the order of these partitions does not matter (divide by (C/M)!)
 - Order within each partition does not matter (divide by $(M!)^{(C/M)}$)



- Partition Approach

 - For first partition choose *M* from $C = \begin{pmatrix} C \\ M \end{pmatrix} = \frac{C(C-1)(C-2)...(C-M+1)}{M!}$ For 2nd partition choose *M* from $C-M = \begin{pmatrix} C M \\ M \end{pmatrix} = \frac{(C-M)(C-M-1)(C-M-2)...(C-2M+1)}{M!}$
 - But different orderings of (C/M) parts don't matter





Let's review performance of various mathematical programming optimization approaches.

MLSys Tutorial: Sparsity in ML

Binary Quadratic Program (BQP) - Naïve Version

- Indexes: i for column; j for row; k for partition (k=C/4)
- Constant: $A_{ij} \ge 0$ Matrix (weight) value in column i and row j
- Binary variables:
 - $x_{ik} \in \{0,1\}$ 1 if column i is assigned to partition k; 0 otherwise.
 - $y_{ij} \in \{0,1\}$ 1 if entry column i and row j is selected for deletion; 0 otherwise.
- Constraints:
 - $\sum_k x_{ik} = 1 \forall i$ (every column is assigned to exactly one partition)
 - $\sum_{i} x_{ik} = 4 \forall k$ (every partition has exactly 4 columns)
 - $\sum_{i} x_{ik} y_{ij} = 2 \forall j, k$ (delete 2 entries from every partition) <= Quadratic!
- Objective Function:
 - Min $\sum_{ij} A_{ij} y_{ij}$ (Minimize smallest entries to obtain 2: 4 sparsity)



Number of

 C^2/M

and

Binary Variables

C x R (# Rows)

Binary Linaear Program (BLP)

- Indexes: i for column; j for row; k for partition
- Constant: $A_{ij} \ge 0$ Matrix (weight) value in column i and row j
- Binary variables:
 - $x_{ik} \in \{0,1\}$ 1 if column i is assigned to partition k; 0 otherwise.
 - $1 \ge z_{ijk} \ge 0$; 1 if entry column i and row j in partition k is selected for deletion; 0 o/w.
- Constraints:
 - $\sum_k x_{ik} = 1 \forall i$ (every column is assigned to exactly one partition)
 - $\sum_{i} x_{ik} = 4 \forall k$ (every partition has exactly 4 columns)
 - $z_{ijk} \le x_{ik} \forall i, j, k$ (can only delete entry if column i is in partition k)
 - $\sum_{i} z_{ijk} = 2 \forall j, k$ (delete 2 entries from every partition)
- Objective Function:
 - Min $\sum_{ijk} A_{ij} z_{ijk}$ (Minimize smallest entries to obtain 2:4 sparsity)





Further Analysis: 4-Dimensional Assignment Formulation



- Symmetry in BLP and BQP:
- 4-Dimensional Network Structure:
 - C × 4 nodes
 - Nodes represent columns,

edges represent columns of same partition $Z_{ijkl} \in \{0, 1\}$, represents the 4D assignment

 $(i-j, j-k, k-l) \rightarrow (i,j,k,l)$



k = 2 k = 1 Solution 1 2 3 4 5 6 7 8 1 Solution 2 5 6 2 3 7 8 4 1

min $\sum_{i=1}^{N} \sum_{j=1}^{N} \sum_{k=1}^{N} \sum_{l=1}^{N} C_{ijkl} Z_{ijkl}$
s.t. $\sum_{j=1}^{N} X_{ij}^{s} \leq 1 \forall i \in \{1, 2, \cdots, N\} \forall s \in \{1, 2, 3\}$
$\sum_{i=1}^{N} X_{ij}^{s} \leq 1 \hspace{0.2cm} \forall j \in \{1,2,\cdots,N\} \hspace{0.2cm} \forall s \in \{1,2,3\}$
$\sum_{i=1}^N\sum_{j=1}^N X_{ij}^s=rac{N}{4}~~orall s\in\{1,2,3\}$
$\sum_{j=1}^{N} \sum_{s=1}^{3} X_{ij}^{s} = 1 \forall i \in \{1, 2, \cdots, N\}$
$\sum_{k=1}^{N}\sum_{l=1}^{N}Z_{ijkl}=X_{ij}^{1} \forall i \in \{1, 2, \cdots, N\} \ \forall j \in \{1, 2, \cdots, N\}$
$\sum_{l=1}^{N}\sum_{i=1}^{N}Z_{ijkl}=X_{jk}^{2} \forall j \in \{1, 2, \cdots, N\} \ \forall k \in \{1, 2, \cdots, N\}$
$\sum_{j=1}^{N}\sum_{k=1}^{N}Z_{ijkl}=X_{kl}^{3} \forall k \in \{1, 2, \cdots, N\} \ \forall l \in \{1, 2, \cdots, N\}$
$X_{ik} \in \{0,1\}, \hspace{0.2cm} 0 \leq Z_{ijkl} \leq 1$



- Given, Q = {q : q \subset {1, 2, 3, \cdots , N}, |q| = 4}
- For some $q \in Q : q \rightarrow (i, j, k, l)$ define deletion cost of partition $D_{ijkl} = D_q$

• Define a constant matrix H:

$$H_{iq} = \begin{cases} 1 & \text{if column-}i \in q \\ 0 & \text{otherwise} \end{cases}$$

• $X_q \in \{0, 1\}$ represents the selection of q among the optimal column partitions. min $\sum_{q \in Q} D_q X_q$

s.t.
$$\sum_{q \in Q} H_{iq}X_q = 1 \quad \forall i \in \{1, 2, \cdots, N\}$$

 $X_q \in \{0, 1\}$

• Note: We need O(N⁴) cost terms. Set-Covering is NP-hard.



Comparison of Math Programming Formulations



• Note: We need O(N⁴) cost terms. Set-Covering is NP-hard.

but timeout at

120

- Given the Optimal N:M Channel Sparsity problem is NP-hard, we need scalable heuristics for larger problem instances.
- Recall hill-climbing heuristics discussed earlier.
- Simulated Annealing is a probabilistic method for achieving near global optimality.
- (We are currently testing this approach.)
- Other metaheuristics can also be applied.











OVIDIA. University at Buffalo

Organizers

Wen-mei Hwu, Jinjun Xiong, Mert Hidayetoğlu, Rakesh Nagi, Jeff Pool, Sitao Huang, Vikram Mailthody

Tutorial | MLSys 2022 | Aug 30, 2022

https://mlsys.org/virtual/2022/tutorial/2199

Slack for discussions: https://bit.ly/3KLQEUN

Summary

1. Opening remarks and overview of sparsity in ML



- 1. EMOGI, PyTorch-Direct are available in DGLv0.8+ with UnifiedTensor representation
- 2. <u>https://github.com/illinois-impact/EMOGI</u>
- 2. Tiled SpMM and its performance model on GPUs
 - 1. PyTorch binding is available on request. Code:

https://github.com/merthidayetoglu/SpDNN_Challenge2020

- 3. Sparse deep neural network inference on FPGAs
- 4. 2:4 Sparsity on GPU Tensor Cores
- 5. Future work and closing remarks

Existing Models Larger Than CPU Memory



Source: Cerebras Systems, HotChips 2022

Lower compute utilization with scale-out!

MLSys Tutorial: Sparsity in ML

Slide 172

1. Very large data objects with sparse access patterns



- What if the data object does not fit into the host memory of one compute node?
- E.g., efficient communication for SpMM across **multiple compute nodes**
- E.g., efficient access to node feature tensors that has to be accessed from **storage**

• Training spare models rather than sparsifying dense models to enable training of models that are orders of magnitude larger than what we can train today

• What else can we do to the compute devices and memory interfaces to support sparsity much better?



Problems With CPU-Centric Model



Complex tiling strategies

I/O and memory management overhead

Gross I/O amplification

Interested in collaborating with us, please reach out to us in slack:

https://bit.ly/3KLQEUN

Email: vsm2@Illinois.edu or hidayet2@illinois.edu

Thank you!

MLSys Tutorial: Sparsity in ML