AutoTM: Automatic Tensor Movement in Heterogeneous Memory Systems using Integer Linear Programming

> **Mark Hildebrand**¹, Jawad Khan², Sanjeev Trika², Jason Lowe-Power¹, Venkatesh Akella¹

> > ¹ University of California, Davis ² Intel Corporation

https://github.com/darchr/AutoTM March 12, 2020

Executive Summary

Problem

Automatic two-level memory management for Deep Neural Networks

Idea

- Profile Guided Optimization
- Model as an Integer Linear Program (ILP)

Results

- Replace 50-80% DRAM with NVDIMMs with geometric mean 27.1% performance loss.
- **3x** better performance than real hardware cache.

Outline

Background

AutoTM

Profiling ILP Modeling

Results

Wrap Up

Why Deep Neural Networks

Why Deep Neural Networks



Image: https://www.microsoft.com/en-us/research/blog/turing-nlg-a-17-billion-parameter-language-model-by-microsoft/ 4/29

Why Deep Neural Networks

T-NLG

Can we use multiple levels of memory to train large models on a single machine?



Image: https://www.microsoft.com/en-us/research/blog/turing-nlg-a-17-billion-parameter-language-model-by-microsoft/ 4/29

Heterogeneous Memory Systems

- Two types of memory.
- Same memory controller.
- Both are byte addressable.
- NVDIMMs for high capacity and low cost



NVDIMM Style

Heterogeneous Memory Systems

- Two types of memory.
- Same memory controller.
- Both are byte addressable.
- NVDIMMs for high capacity and low cost

Challenges

- All tensors in NVDIMMs memory is **too slow**.
- DRAM as a cache for NVDIMMs also too slow.
- Intelligent memory management required.



NVDIMM Style

Outline

Background

$\mathsf{Auto}\mathsf{T}\mathsf{M}$

Profiling ILP Modeling

Results

Wrap Up

AutoTM

Goal

Minimize execution time

- Arbitrary computation graph
- Size constraint on fast memory



AutoTM

Goal

Minimize execution time

- Arbitrary computation graph
- Size constraint on **fast** memory

How

- Place tensors in **fast** or **slow** memory.
- Optimal tensor movement



AutoTM

Goal

Minimize execution time

- Arbitrary computation graph
- Size constraint on **fast** memory

How

- Place tensors in **fast** or **slow** memory.
- Optimal tensor movement

Strategy

- Profile kernel performance.
- Model tensor assignment as ILP.



Kernel Profiling

Profile performance of kernels for all tensor IO locations.

Kernel	IO Tensor Locations		
	T1	T2	Т3
	DRAM	DRAM	DRAM
	DRAM	DRAM	PMM
	DRAM	PMM	DRAM
K2	DRAM	PMM	PMM
	PMM	DRAM	DRAM
	PMM	DRAM	PMM
	PMM	PMM	DRAM
	PMM	PMM	PMM

Table: Profile space for kernel K2.



Kernel Profiling

Profile performance of kernels for all tensor IO locations.

Kernel	IO Tensor Locations		
	T1	T2	Т3
	DRAM	DRAM	DRAM
	DRAM	DRAM	PMM
	DRAM	PMM	DRAM
K2	DRAM	PMM	PMM
	PMM	DRAM	DRAM
	PMM	DRAM	PMM
	PMM	PMM	DRAM
	PMM	PMM	PMM

Table: Profile space for kernel K2.







Tensor Lifetime Flow Graph for tensor T1





Tensor Lifetime Flow Graph for tensor T1



Tensor Lifetime Flow Graph for tensor T1



Tensor Lifetime Flow Graph for tensor T1



Tensor Lifetime Flow Graph for tensor T1









17/29





Time moving tensor t_1

18/29



 $t \in \mathcal{L}(k)$ 19/29

Variations of AutoTM

Name	Description	PMM System	GPU System
Static	Tensor's can't move	1	×
Synchronous	Tensor's move but block computation	1	√
Asynchronous	Tensor movement con- current with computa- tion	•	√

Outline

Background

AutoTM

Profiling ILP Modeling

Results

Wrap Up

Experiments!

Software

- Modified the *ngraph*¹ compiler.
- Julia's *JuMP*² package for ILP modeling.
- Gurobi³ as the ILP solver.

Hardware

- 1.5 TB OptaneTM DC PMM
- 384 GiB DRAM

¹https://github.com/NervanaSystems/ngraph ²https://github.com/JuliaOpt/JuMP.jl

³gurobi.com

Experiments!

Software

- Modified the *ngraph*¹ compiler.
- Julia's *JuMP*² package for ILP modeling.
- Gurobi³ as the ILP solver.

Hardware

- 1.5 TB OptaneTM DC PMM
- 384 GiB DRAM

¹https://github.com/NervanaSystems/ngraph
²https://github.com/JuliaOpt/JuMP.jl
³gurobi.com

Conventional	Batchsize	Memory (GB)
Inception v4	1024	111
Vgg 19	2048	143
Resnet 200	512	132
DenseNet 264	512	115

Experiments!

Software

- Modified the *ngraph*¹ compiler.
- Julia's *JuMP*² package for ILP modeling.
- Gurobi³ as the ILP solver.

Hardware

- 1.5 TB OptaneTM DC PMM
- 384 GiB DRAM

Conventional	Batchsize	Memory (GB)
Inception v4	1024	111
Vgg 19	2048	143
Resnet 200	512	132
DenseNet 264	512	115

Large	Batchsize	Memory (GB)
Inception v4	6144	659
Vgg 416	128	658
Resnet 200	2560	651
DenseNet 264	3072	688

¹https://github.com/NervanaSystems/ngraph

²https://github.com/JuliaOpt/JuMP.jl

³gurobi.com







 \circ Just using PMM is **too slow**.





 \circ Best performance when working-set fits in memory. $_{^{25/29}}$

2LM DRAM Cache







26/29





- Avoid Dirty Writebacks
- Lower Memory Contention





Software management outperforms hardware management by up to **3**x.

Outline

Background

AutoTM

Profiling ILP Modeling

Results

Wrap Up

Limitations

- Static computation graphs.
- Kernel profiling overhead.
- ILP solution times.
- ILP solution may be hard to interpret.

Limitations

- Static computation graphs.
- Kernel profiling overhead.
- ILP solution times.
- ILP solution may be hard to interpret.



Conclusion

AutoTM: A technique for managing tensors in heterogeneous memory systems.

- Profiling for Kernel Performance.
- Use ILP to optimally assign tensor location and movement.
- Three formulations: *Static, Synchronous, Asynchronous.*

We show

- Reduce **DRAM** requirement.
- Significant performance improvement over hardware solutions.

Code Available: https://github.com/darchr/AutoTM

Common Questions

Asynchronous Movement on PMMs

- Interference between DRAM and PMM.
- Low bandwidth and difficulty of DMA.
- Performance of kernels greatly impacted due to copy kernels.

GPU

synchronous synchronous oracle



RNNs - more complex models

- AutoTM is limited to static computation graphs.
- RNNs have dynamic behavior (i.e. unrolling based on sequence length).
- RNNs can be implemented statically.
- Key ideas from AutoTM can be used for dynamic workloads.



Concluding Conclusion

AutoTM: A technique for managing tensors in heterogeneous memory systems.

- Profiling for Kernel Performance.
- Use ILP to optimally assign tensor location and movement.
- Three formulations: Static, Synchronous, Asynchronous.

We show

- Reduce **DRAM** requirement.
- Significant performance improvement over hardware solutions.

Code Available: https://github.com/darchr/AutoTM