

# MagmaDNN: Accelerated Deep Learning Using MAGMA

D. Nichols<sup>1</sup>, K. Wong<sup>1</sup>, S. Tomov<sup>1</sup>, L. Ng<sup>2</sup>, S. Chen<sup>2</sup>,  
A. Gessinger<sup>3</sup>

1 - University of Tennessee, Knoxville

2 - The Chinese University of Hong Kong

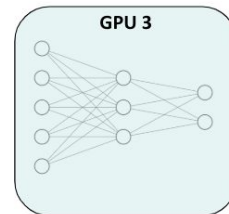
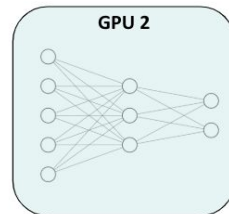
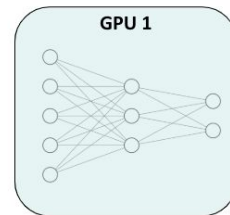
3 - Slippery Rock University

# Organization

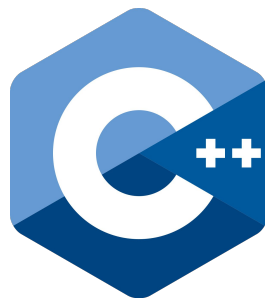
- Motivation
- Magma
- MagmaDNN Overview
  - Framework Overview
  - Compute Graph Optimization
  - Tuning
  - Distributed Training
- Results
- Current and Future Work
- Availability

# Motivation

- Utilize state of the art MAGMA LA framework
- Provide a modular C++ Deep Learning Interface
- Support state of the art distributed training techniques



**MAGMA**

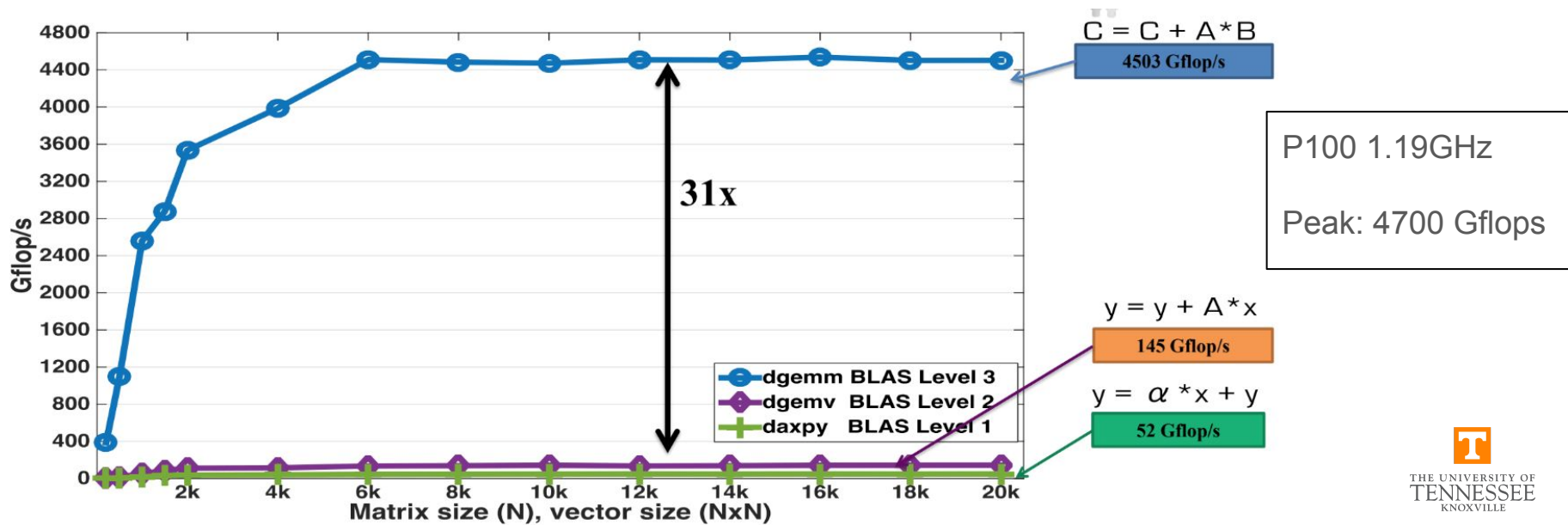


# MAGMA

- Accelerated Linear Algebra on Heterogeneous Architectures

$$g_n (U_{n-1}W_n + b_n)$$

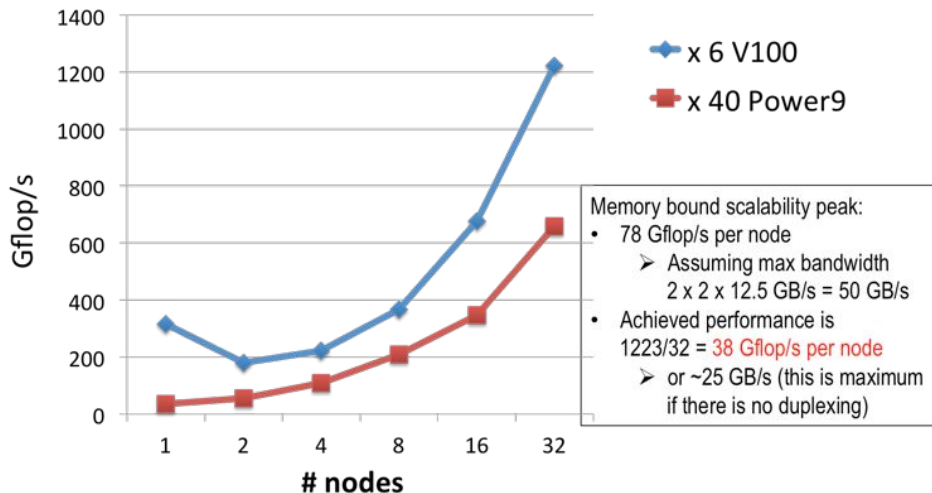
$$Y = A^T \left[ [GgG^T] \odot [B^T dB] \right] A$$



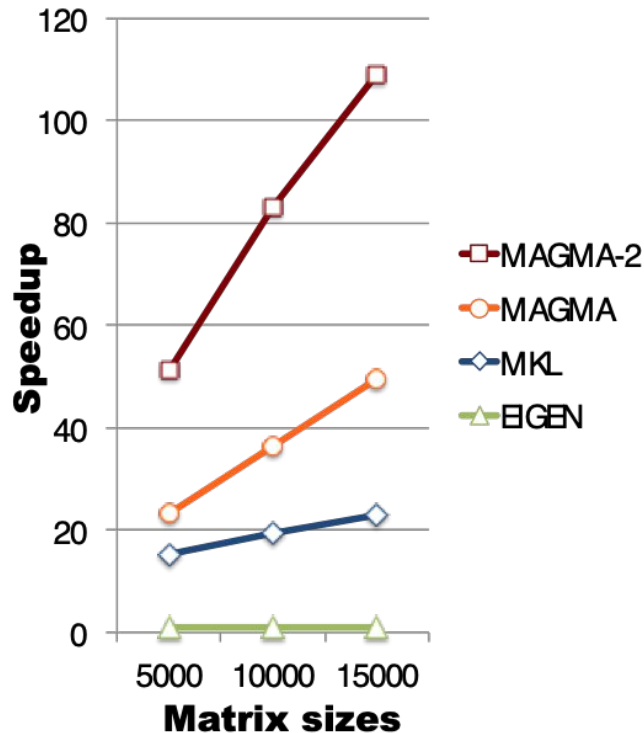
# MAGMA for Data Science

- Better performance than other leading LA packages in SVD
- Scalable FFTs for Convolutions

## Strong scalability of 3D FFT on Summit (N = 1024)

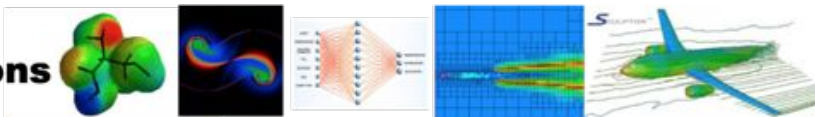


## SVD performance speedup



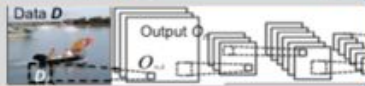
# MagmaDNN Overview

## Applications



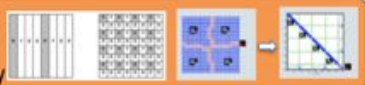
## MagmaDNN

High-performance data analytics and machine learning for many-core CPUs and GPU accelerators



## MAGMA Templates

Scalable LA on new architectures  
Data abstractions and APIs  
Heterogeneous systems portability



## SLATE

Tile algorithms  
LAPACK++  
BLAS++

ScaLAPACK API

MPI

### Single Heterogeneous Node

MAGMA(dense)

MAGMABatched

MAGMASparse

Shared memory

BLASAPI

LAPACKAPI

Batched BLASAPI

OpenMP

MKL

ESSL

cuBLAS

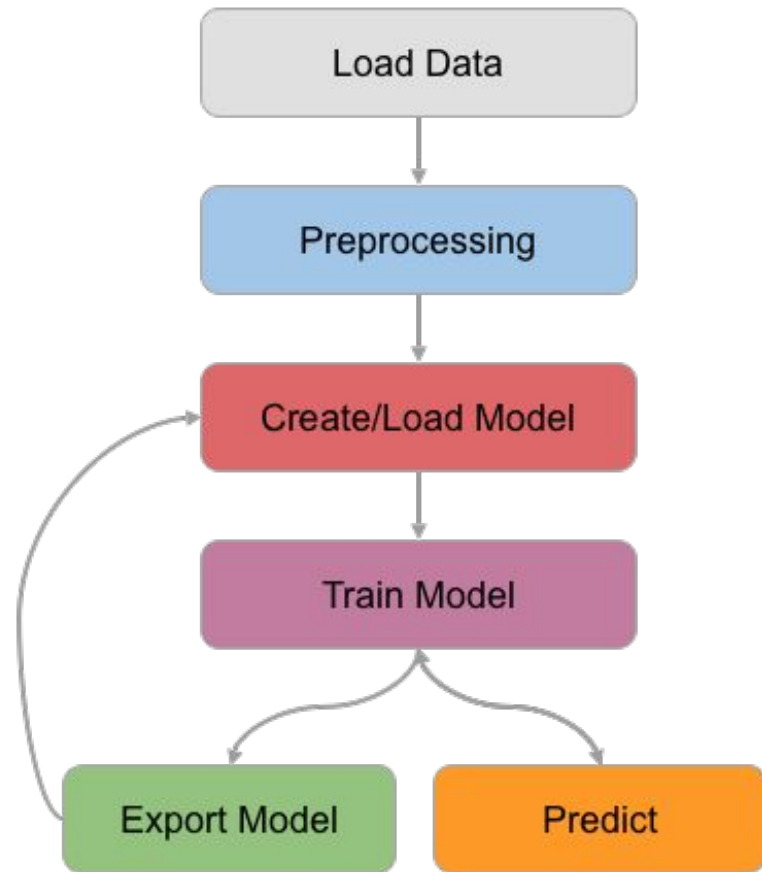
ACML

LA libraries

Standard LA APIs

Run-time/comm. APIs

Vendor Libraries



# Framework Overview

Memory  
Manager

Tensor

Operation &  
Graph

Optimizer

Model

# Framework Overview

Memory  
Manager

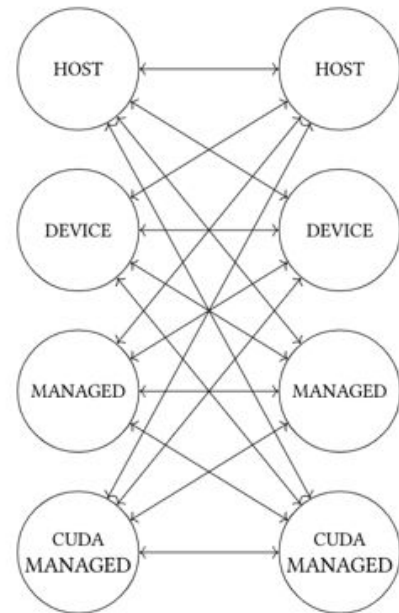
Tensor

Operation &  
Graph

Optimizer

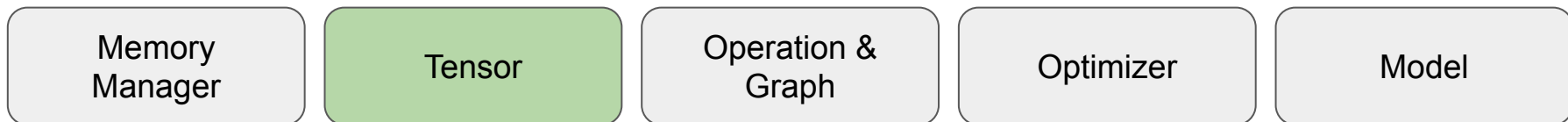
Model

```
size_t n_vals = 20;  
// HOST, DEVICE, MANAGED, CUDA_MANAGED  
memory_t mem_type = MANAGED;  
device_t device_id = 0;  
MemoryManager<float> mm(n_vals, mem_type, device_id);
```





# Framework Overview



Scalar   Vector   Matrix   Tensor

1

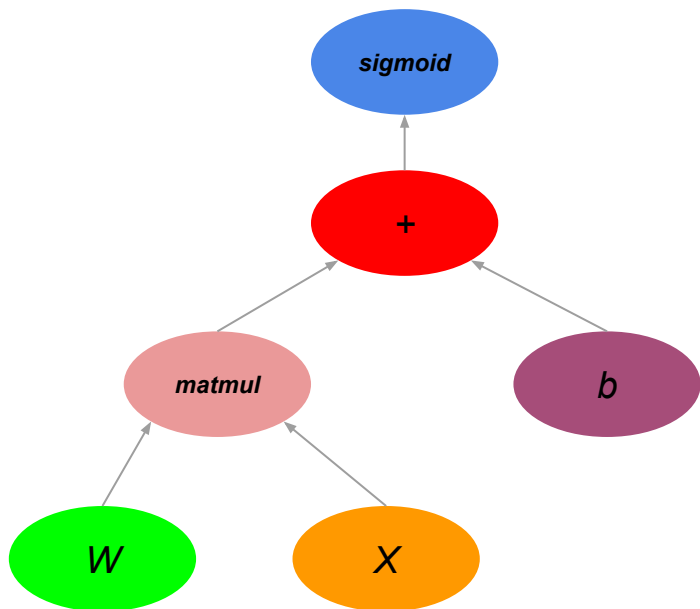
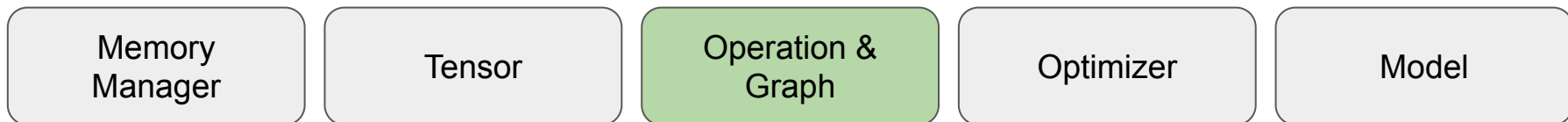
$\begin{bmatrix} 1 \\ 2 \end{bmatrix}$

$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$

$\begin{bmatrix} \begin{bmatrix} 1 & 2 \end{bmatrix} & \begin{bmatrix} 3 & 2 \end{bmatrix} \\ \begin{bmatrix} 1 & 7 \end{bmatrix} & \begin{bmatrix} 5 & 4 \end{bmatrix} \end{bmatrix}$

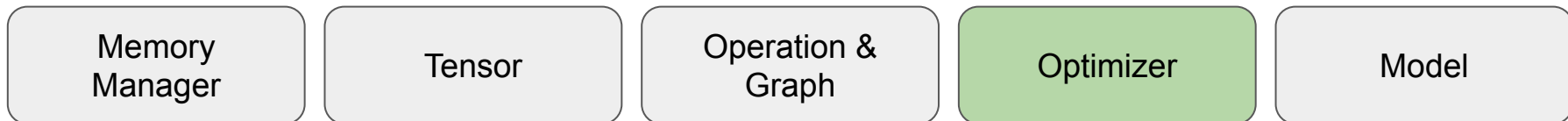
```
Tensor<float> x({32,28,28}, {UNIFORM, {-1.0f,1.0f}}, DEVICE);  
std::cout << x.get({1,3,0});  
x.set({1,3,0}, 8.0f);
```

# Framework Overview



```
auto X = op::var<float>("X", {5,3},{UNIFORM});  
auto W = op::var<float>("W", {6,5},{UNIFORM});  
auto b = op::var<float>("b", {6,3},{UNIFORM});  
  
auto transform = op::add( op::matmul(W, X), b );  
  
Tensor<float> *output = transform->eval();
```

# Framework Overview



```
auto x = op::var<float>("x", NONE);
auto c = op::var<float>("c", {CONSTANT, -2.0f});

optimizer::GradientDescent opt(0.05);

opt.minimize( op::add(op::pow(x, 2), c), {x});
```

minimize  $x^2 + c$

with respect to  $x$

# Framework Overview

Memory  
Manager

Tensor

Operation &  
Graph

Optimizer

Model

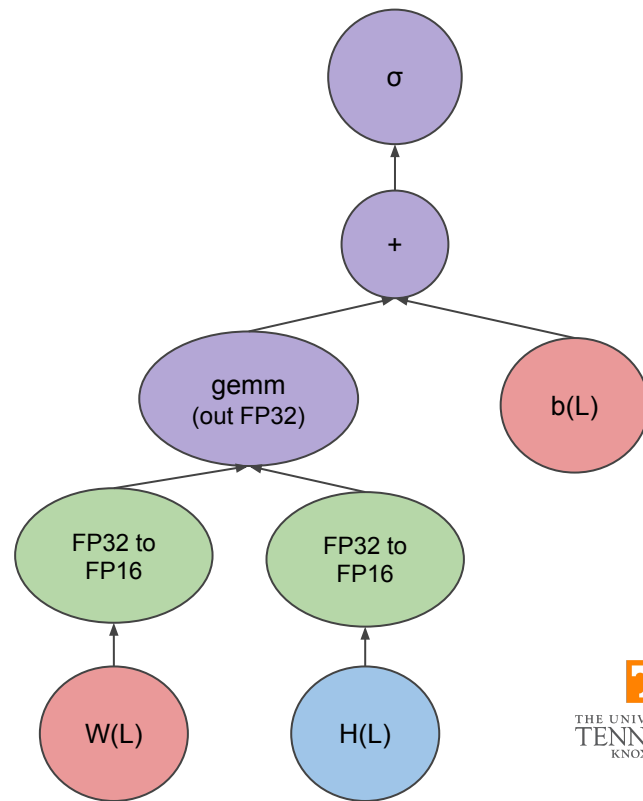
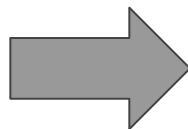
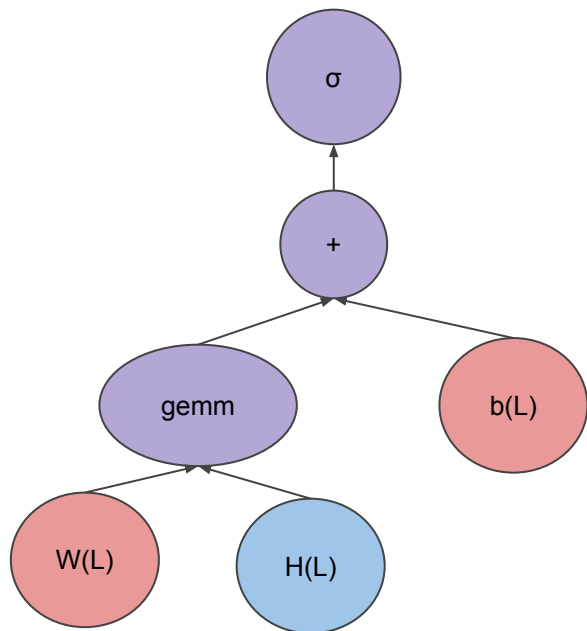
```
Tensor<float> data({60000, 28*28}, HOST);
io::read_csv_to_tensor(data, "mnist_data_set.csv");
Tensor<float> labels({60000, 10}, HOST);
io::read_csv_to_tensor(labels, "mnist_labels_set.csv");

/* create a vector of layers ... */

model::NeuralNetwork<float> model(layers_vector, optimizer::CROSS_ENTROPY, optimizer::ADAM,
{batch_size, n_epochs, learning_rate});
model.fit(data, labels, params_out, verbose);
```

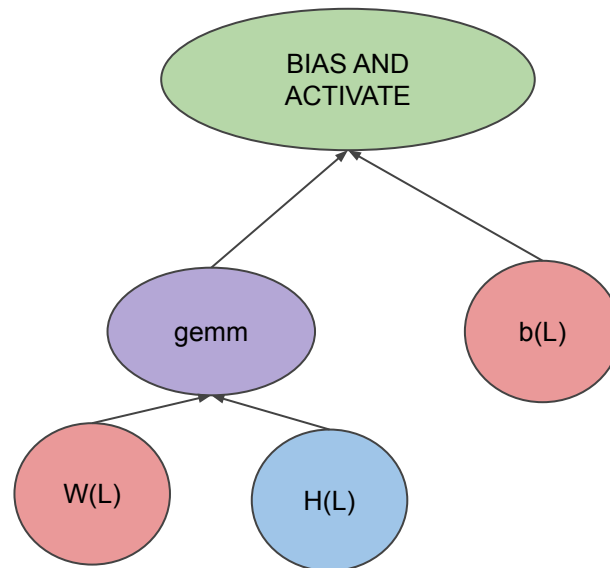
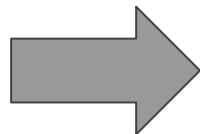
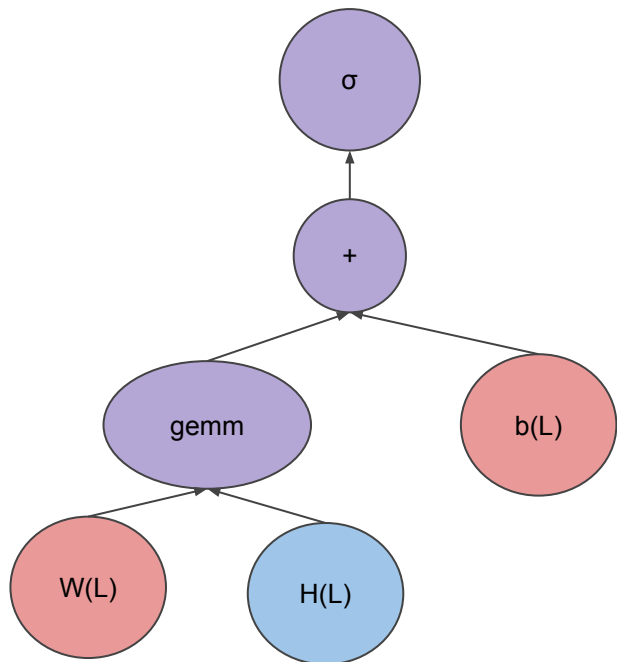
# Compute Graph Optimization

- Mixed Precision Training
- Make use of Volta Tensor Cores



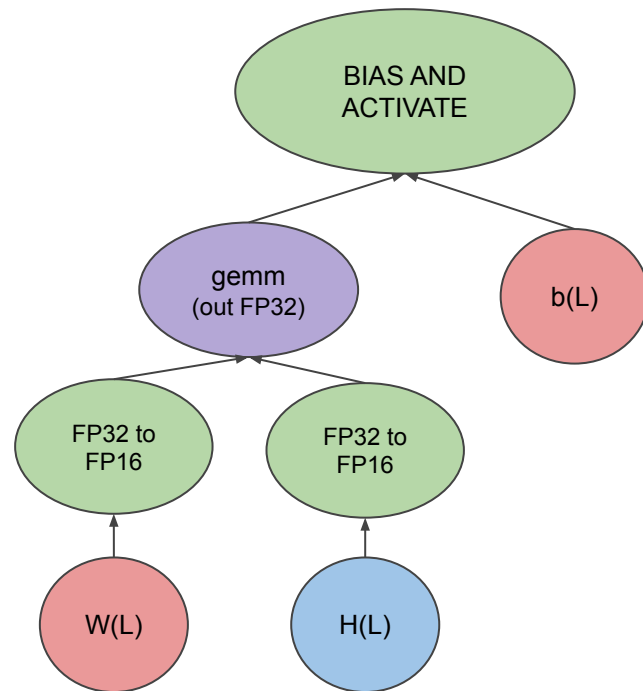
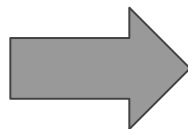
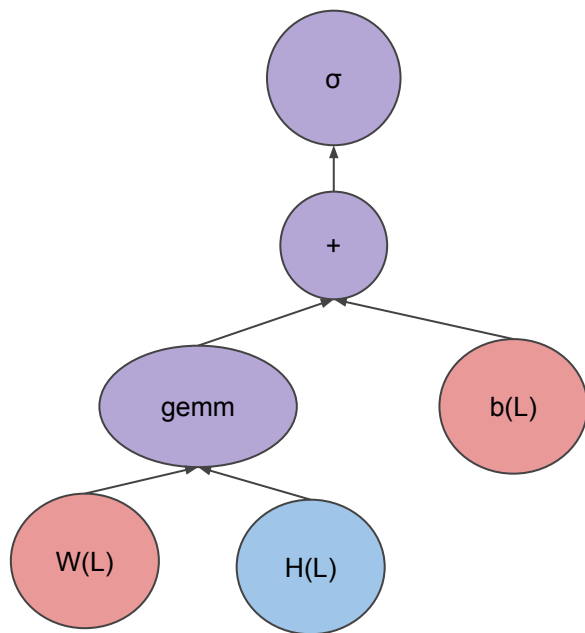
# Compute Graph Optimization (cont.)

- Fused Operations



# Compute Graph Optimization

- Combined



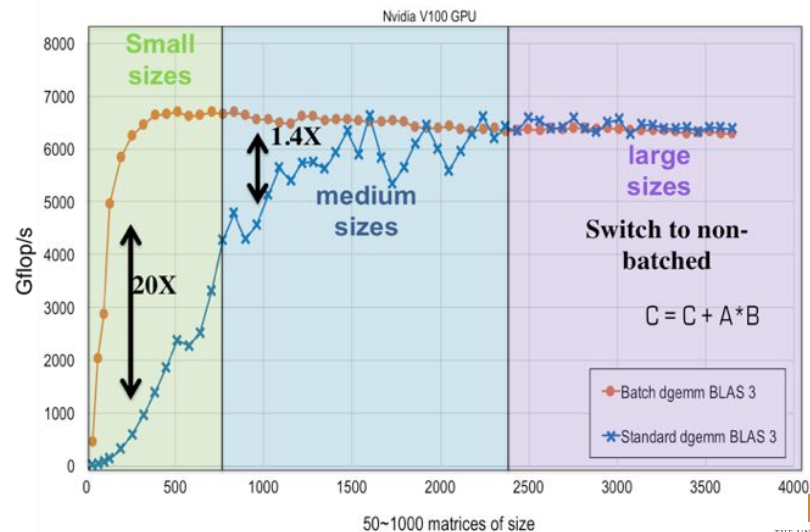
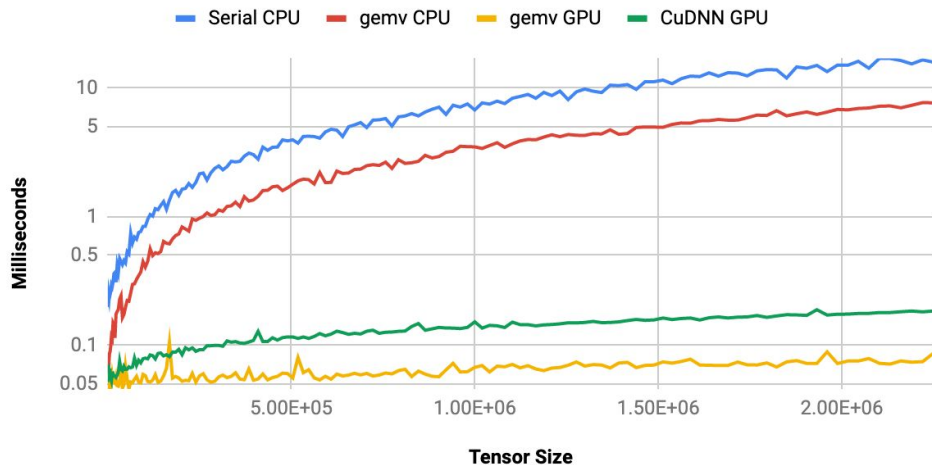
# Tuning

- MagmaDNN tunes tensor and deep learning kernels

- Magma tunes matrix algebra kernels

## Tensor Reductions in MagmaDNN

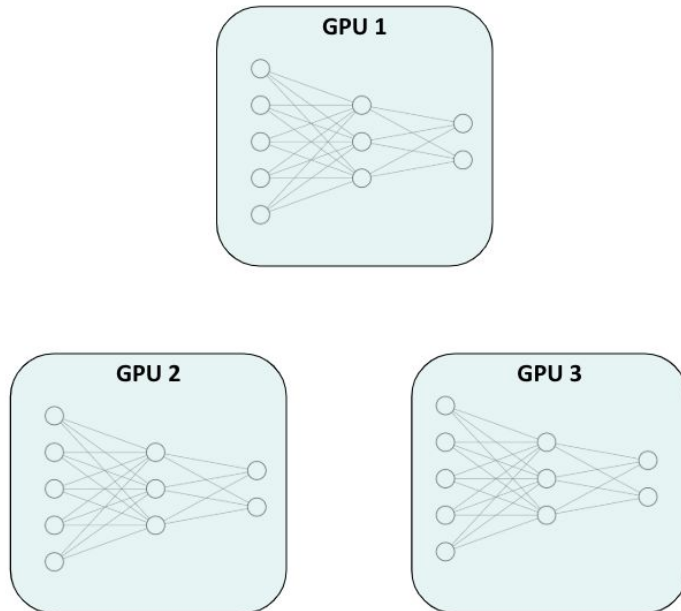
Data Collected on P100 GPU



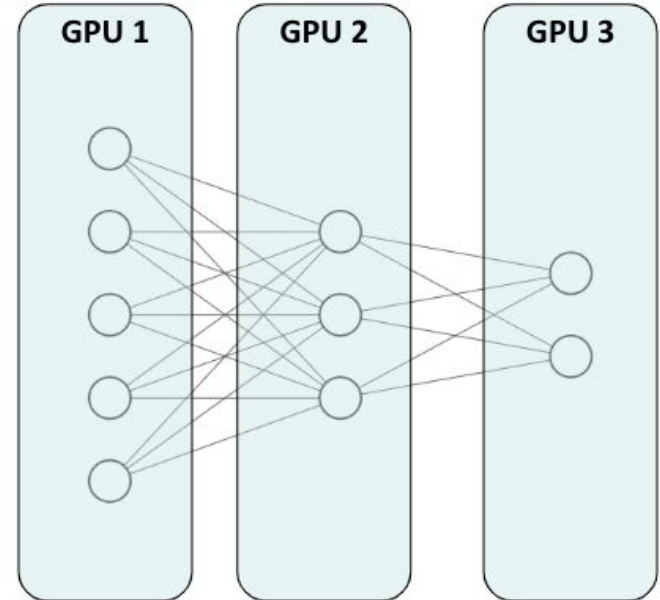


# Distributed Training

## Data Parallelism

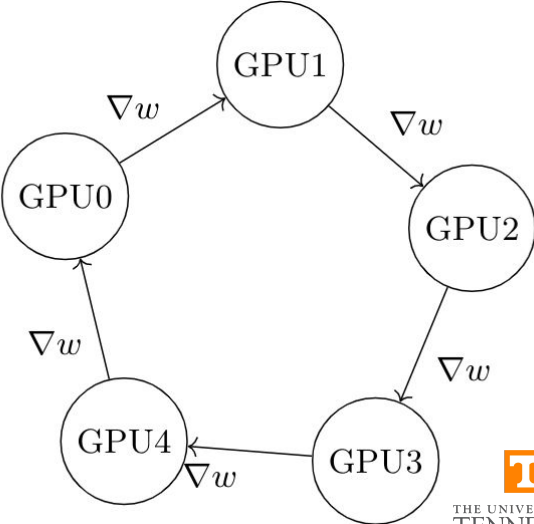
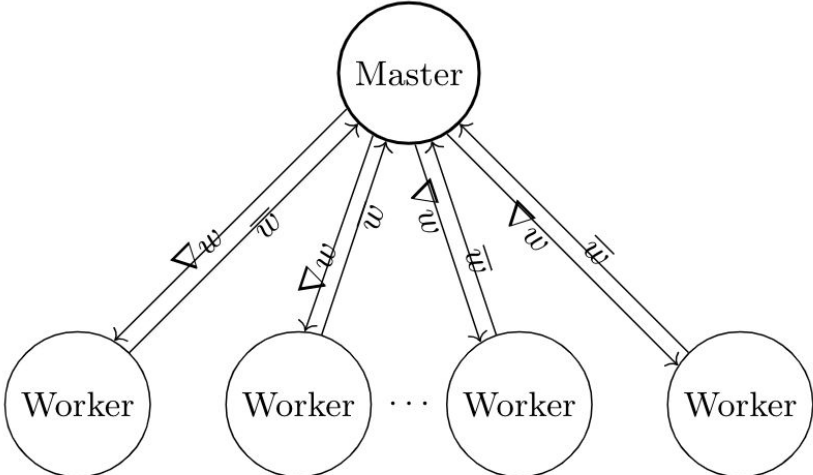


## Model Parallelism



# Distributed Training

- ASGD
- AllReduce
- Ring Reduce

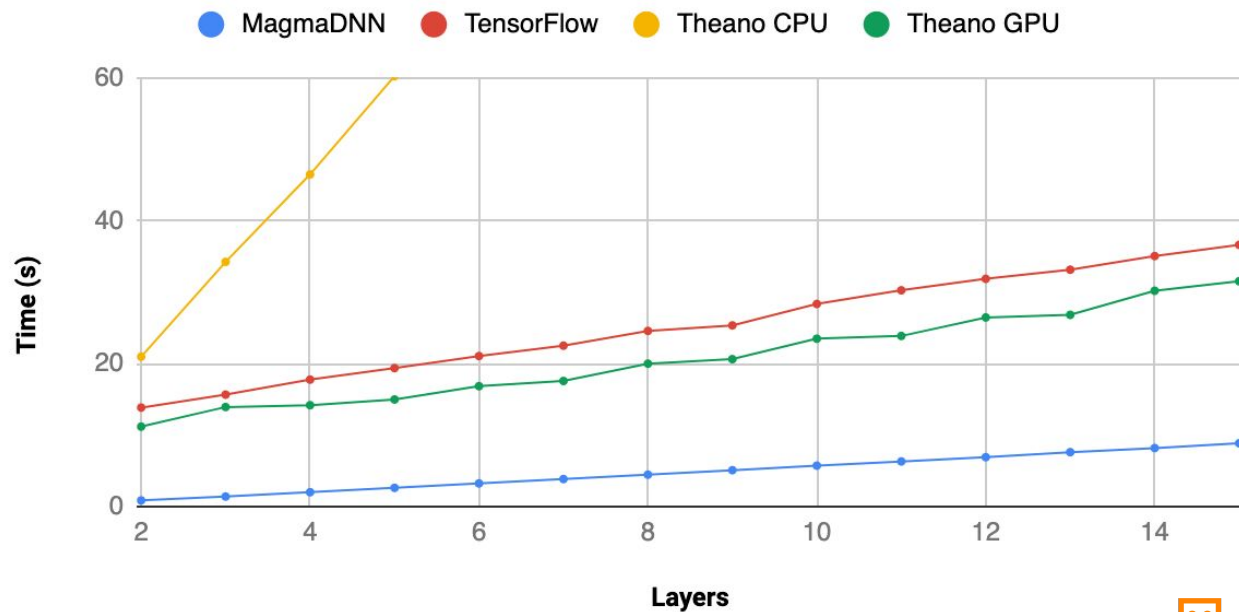


# Results

- Best performance across a single node
- Best scaling with network size

## MLP Time Comparison

Profiled on Nvidia 1050 Ti



# Summary

- Accelerated single node training times
- Modern C++ interface
- Support for distributed training

# Current and Future Work

- Competitive distributed ResNet-50 training time
- Move to new C++ standard
- More “Bells and Whistles”
- Development
- Python Interface

# Acknowledgements

- Joint Institute for Computer Science (JICS)
- National Institute for Computational Sciences (NICS)
- University of Tennessee, Knoxville (UTK)
- NSF Award #1659502
- Extreme Science and Engineering Discovery Environment (XSEDE)
- BP High Performance Computing

# Availability

Development: <https://github.com/MagmaDNN/magmadnn>

Releases: <https://bitbucket.org/icl/magmadnn>

Tutorials: <https://github.com/MagmaDNN/magmadnn/tree/master/docs/tutorials>

Contact: Daniel Nichols (danielnichols@utk.edu)