

# Neural Network Hyperparameter Optimization

Chris Ouyang (CUHK Mathematics)

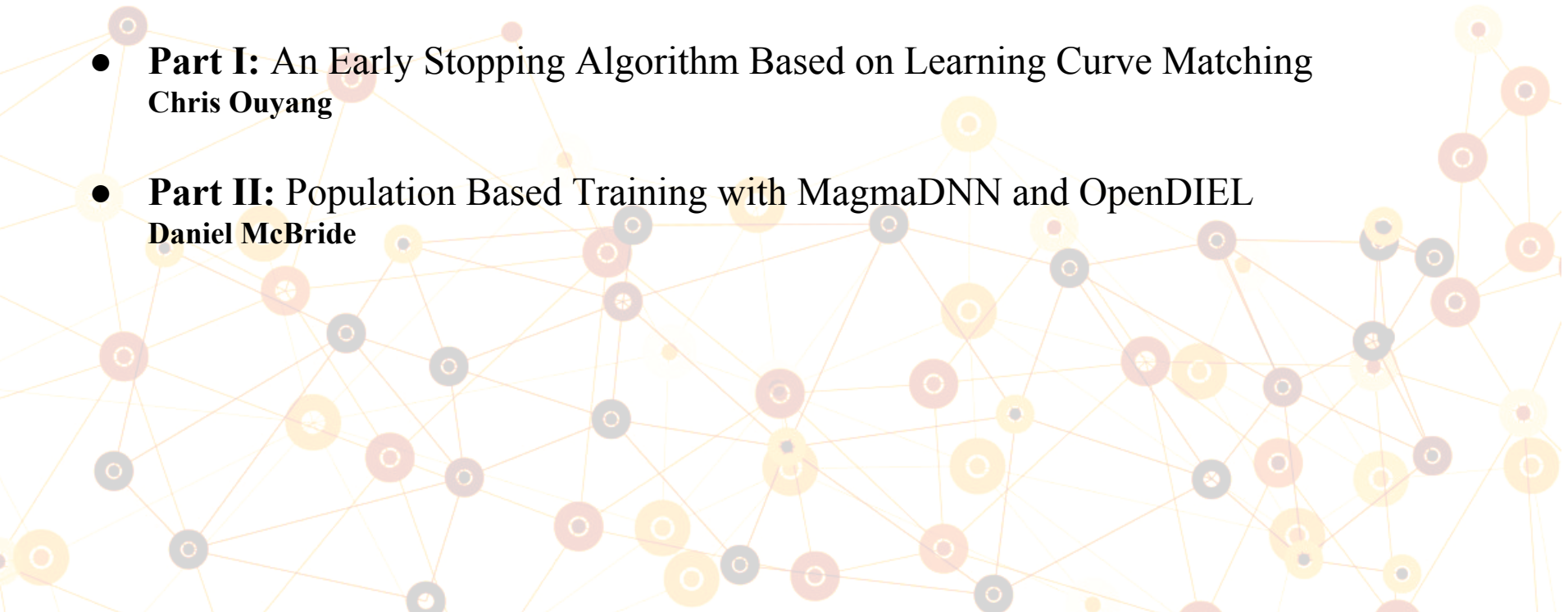


Daniel McBride (UTK Mathematics)



# Presentation Outline

- **Introduction**
- **Part I: An Early Stopping Algorithm Based on Learning Curve Matching**  
Chris Ouyang
- **Part II: Population Based Training with MagmaDNN and OpenDIEL**  
Daniel McBride



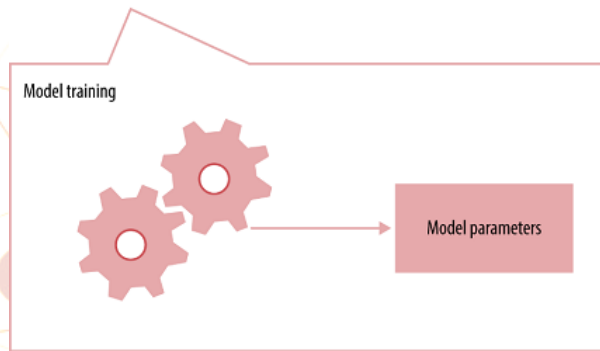
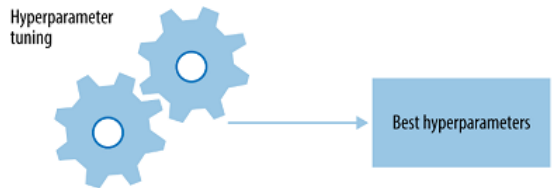
# Introduction

- **What is a hyperparameter?**

They are neural network “presets” like network architecture, learning rate, batch size, and more.

- **Why do we need to optimize the hyperparameters?**

A poor choice of hyperparameters can cause a network’s accuracy to converge slowly or not at all.



# Introduction

- **What are some obstacles to optimizing hyperparameters?**
  - The Curse of Dimensionality
  - Highly irregular (nonconvex, nondifferentiable) search spaces
- **What are some standard hyperparameter optimization techniques?**
  - Classic Approaches: Grid Search, Random Search
  - Modern Approaches: Early Stopping, Evolutionary Algorithms

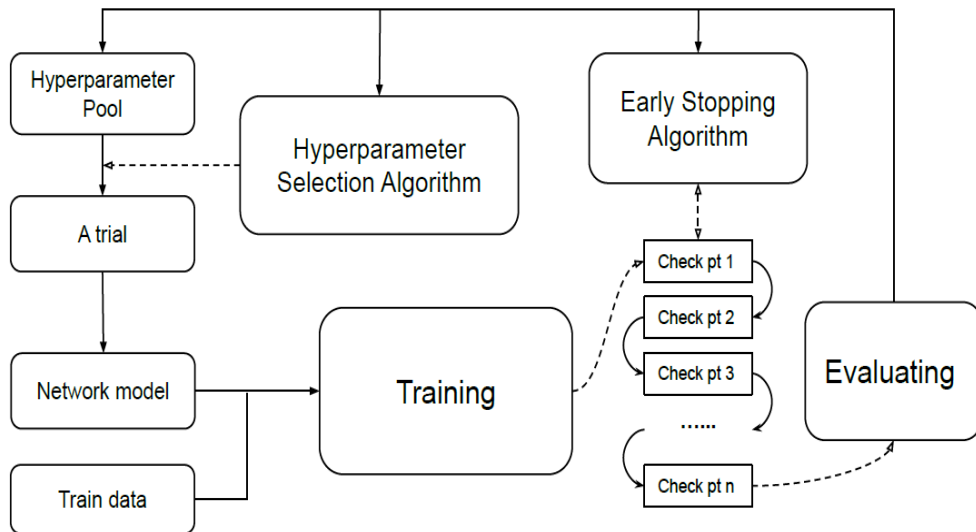
# Part I

## An Early Stopping Algorithm Based on Learning Curve Matching

Chris Ouyang

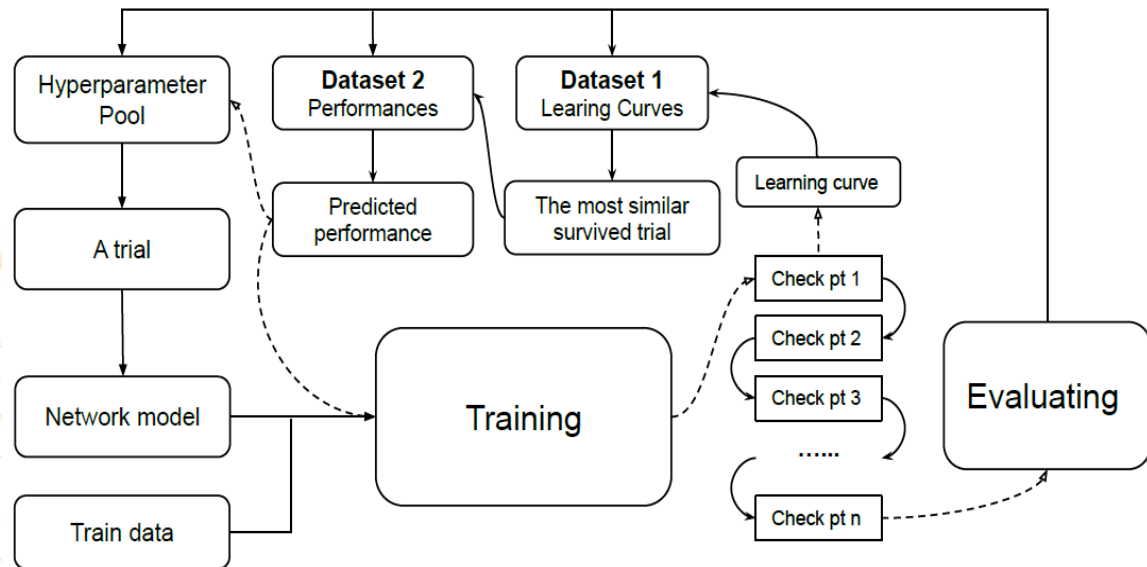
# Hyperparameter Algorithms

- **Hyperparameter Selection:** Random search, grid search and Bayesian optimization
- **Early stopping:** Successive Halving Algorithm (SHA) and Hyperband
- **Advanced Algorithm:** Evolutionary Algorithm, such as population based training (PBT) and swarm optimization.



# LCM Algorithm: Flow Chart and Terms

- ***Trials***: Sets contain a single sample for every hyperparameter.
- ***Learning Curves***: arrays of the numerical values of loss function in some certain stages during a single training.
- ***Check Points***: points where apply LCM to decide whether abort the training



# LCM Algorithm: Cumulation Stage

Learning Curve with performance

[Loss\_1, Loss\_2, Loss\_3, Loss\_4, Loss\_5, ....., Loss\_n, Performance]

**Data Set:**

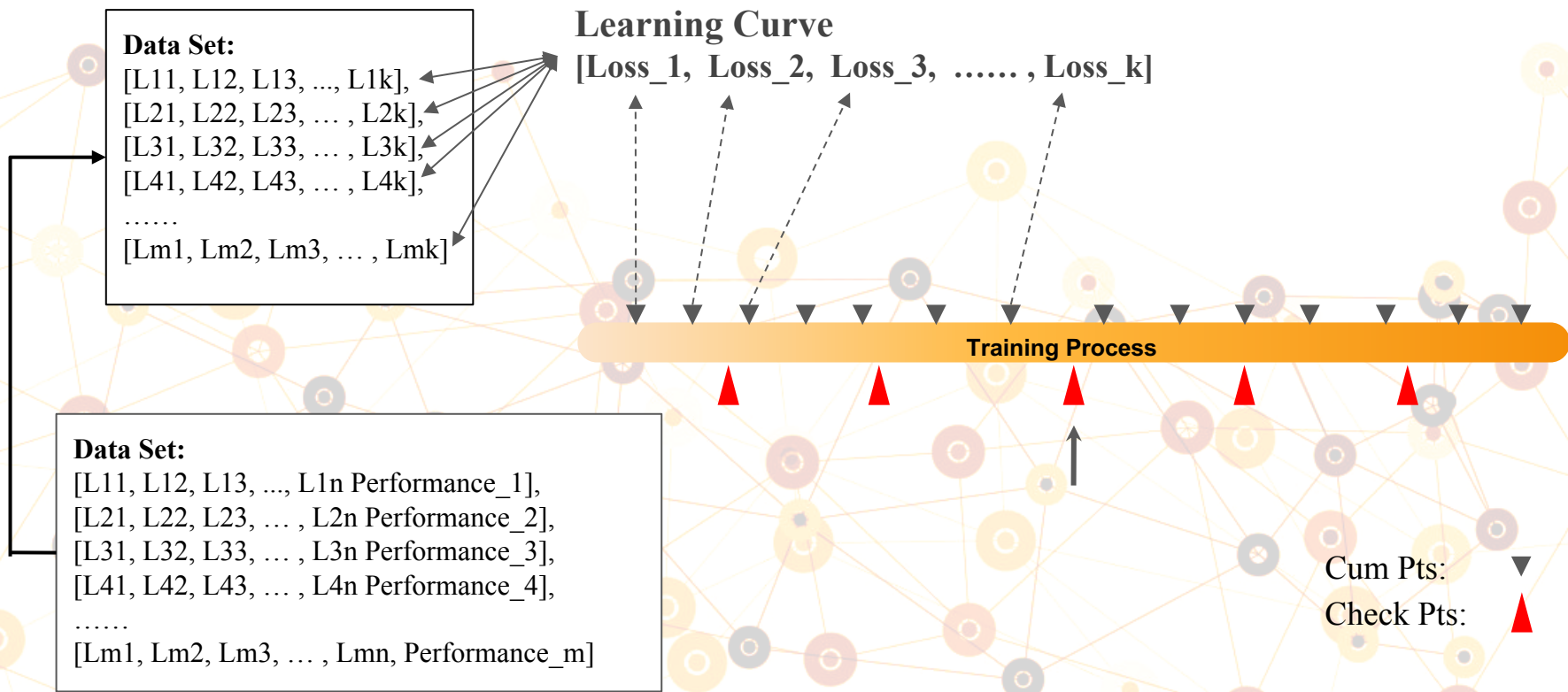
[LC\_1, Performance\_1],  
[LC\_2, Performance\_2],  
[LC\_3, Performance\_3],  
[LC\_4, Performance\_4],  
.....  
[LC\_m, Performance\_m]

Training Process

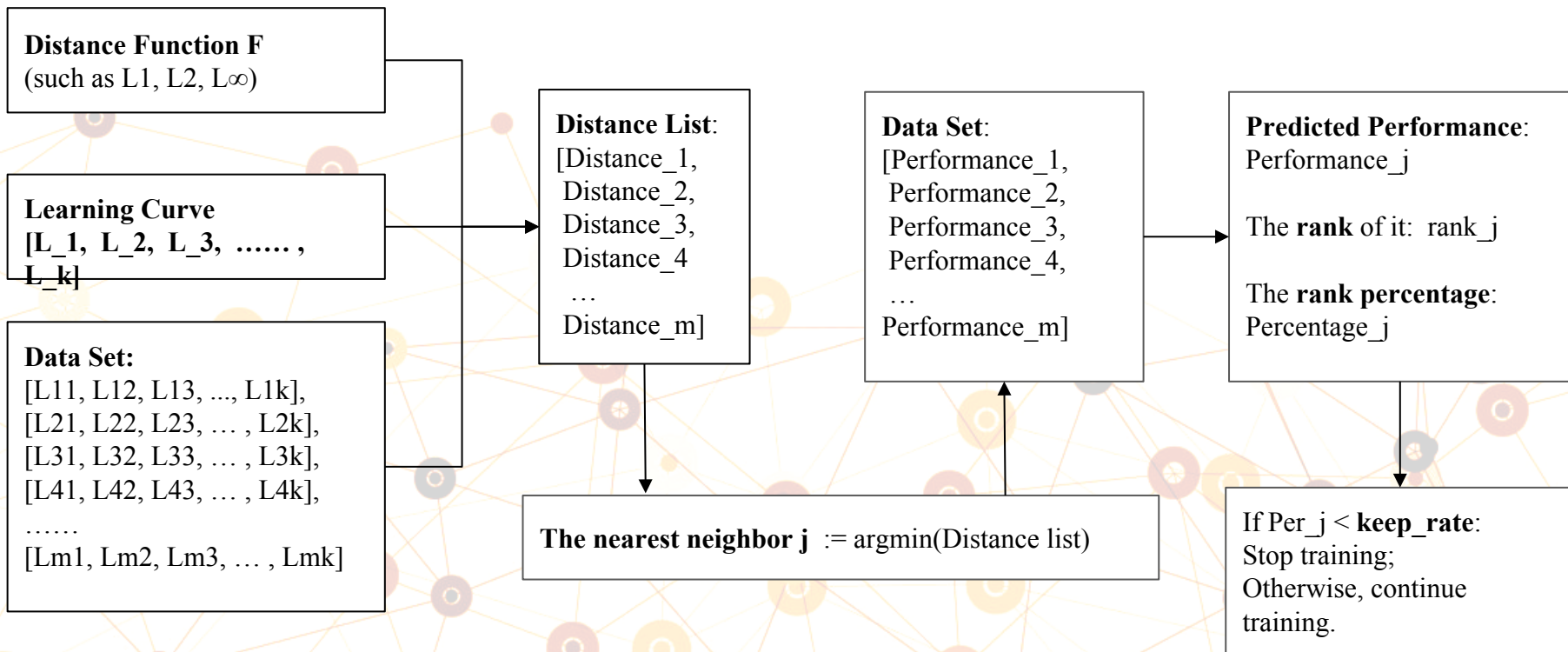
Cum Pts: ▼  
Check Pts: ▲



# LCM Algorithm: Checking Stage



# LCM Algorithm: Checking Stage



# LCM Algorithm: Comparisons

- **Network:** Only one dense layer
- **Dataset:** MNIST
- **Optimizer:** stochastic gradient descent
- **Hyperparameter:** Epochs, batch sizes, learning rate, momentum and decay
- **Benchmark:** Random search
- **Times:** 9

|               | Trials | Computer Time (S) | Best Performance (%) |
|---------------|--------|-------------------|----------------------|
| <b>LCM</b>    | 100    | 778.50            | 97.10                |
| <b>Random</b> | 100    | 3657.75           | 97.41                |

**Remark:** In 5 of 9 experiments, two algorithms got the same optimal hyperparameters.

# LCM Algorithm: Comparisons

- **Network:** Only one dense layer
- **Dataset:** MNIST
- **Optimizer:** stochastic gradient descent
- **Hyperparameter:** Epochs, batch sizes, learning rate, momentum and decay
- **Benchmark:** Random search
- **Times:** 6

|               | Trials | Computer Time (S) | Best Performance (%) |
|---------------|--------|-------------------|----------------------|
| <b>LCM</b>    | 37.67  | 4800              | 97.82                |
| <b>Random</b> | 67.33  | 4800              | 97.69                |

**Remark:** In 4 of 6 experiments, two algorithms got the same optimal hyperparameters.

# LCM Algorithm: Comparisons

- **Network:** Four CNN layers and several dense layers
- **Dataset:** CIFAR10
- **Optimizer:** Adam
- **Hyperparameter:** More than 10 hyperparameters
- **Benchmark:** Random search
- **Times:** 12

|               | Trials | Computer Time (S) | Best Performance (%) |
|---------------|--------|-------------------|----------------------|
| <b>LCM</b>    | 100    | 8069.08           | 67.05                |
| <b>Random</b> | 100    | 26498.00          | 67.26                |

**Remark:** in 7 of 12 experiments, two algorithms got the same optimal hyperparameters.

# Part II

## Population Based Training with MagmaDNN and OpenDIEL

Daniel McBride



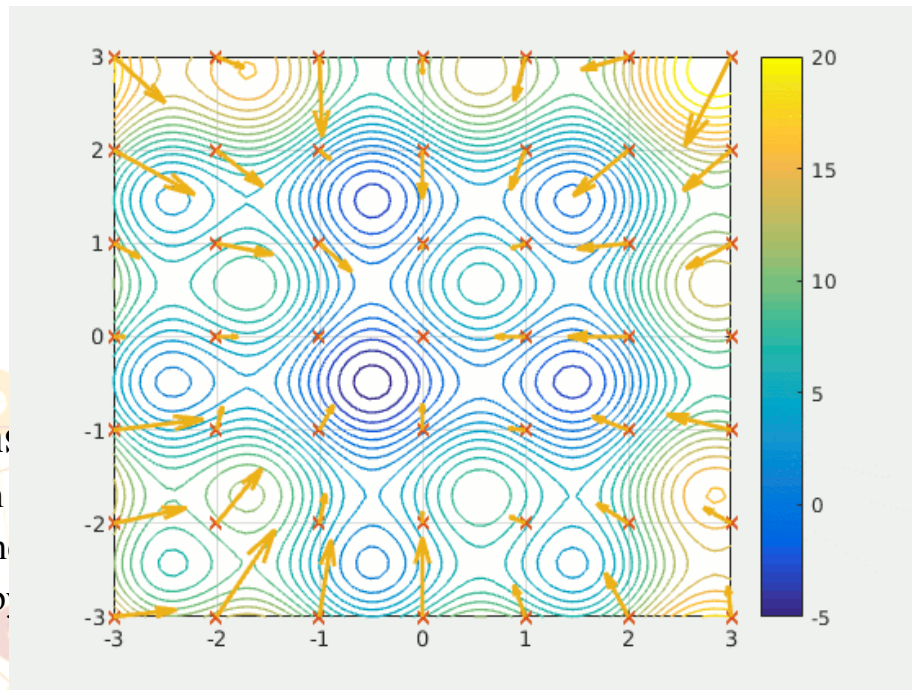
# PBT: Background

- **What is Population Based Training (PBT)?**

PBT is an evolutionary hyperparameter optimization algorithm.

- **Evolutionary optimization algorithms** use natural models to inspire a particular approach to traversing a search space. One classic case is the Particle Swarm Optimization algorithm, inspired by the swarming behavior of bees.

Particle Swarm Optimization



# PBT: Background

- **What are the benefits of PBT?**

PBT outperforms the standard hyperparameter tuning benchmarks. These benchmark algorithms, **Grid Search and Random Search**, each have their own limitations, which PBT overcomes.

- **Why should we implement it on MagmaDNN and OpenDIEL?**

- MagmaDNN and OpenDIEL are engineered for supercomputers.
- The current standard implementation (Ray-Tune: shared memory model) has a scalability bottleneck.



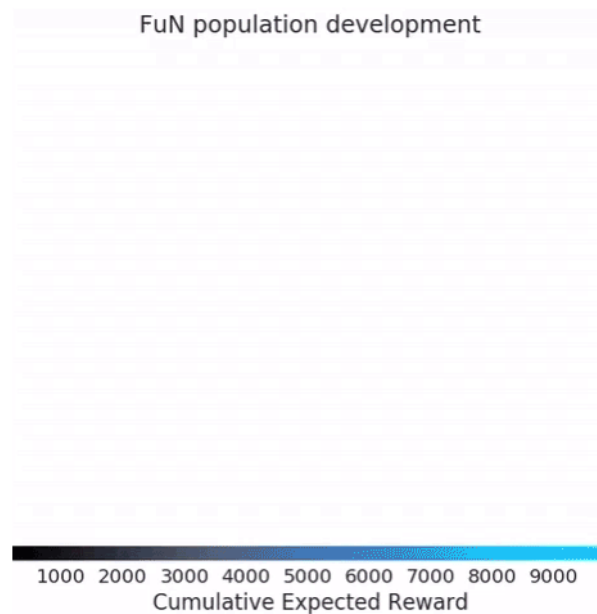
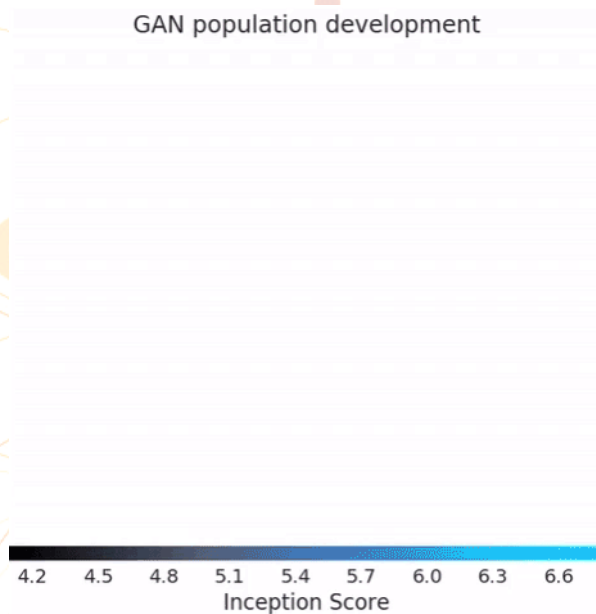
# PBT: Algorithm

## How does the PBT Algorithm work?

- Population Model
- Stochasticity
- Explore / Exploit
- Early Stopping
- Evolution
- Adaptive Hyperparameter Scheduling

# PBT: Algorithm

## How does the PBT Algorithm work?



# PBT: Algorithm

Does PBT's functionality improve on the benchmark algorithms?

|                          | Grid Search | Random Search | PBT |
|--------------------------|-------------|---------------|-----|
| Parallelizability        | ✓           | ✓             | ✓   |
| Stochasticity            | ✗           | ✓             | ✓   |
| Early Stopping           | ✗           | ✗             | ✓   |
| Adaptive Hyperparameters | ✗           | ✗             | ✓   |

# PBT: Analysis - Dynamic Learning Rate

- **Data: MNIST**

- 60k images of handwritten digits 0-9
- 256 greyscale pixels per image
- 10 categories (0-9)

- **Network: MagmaDNN**

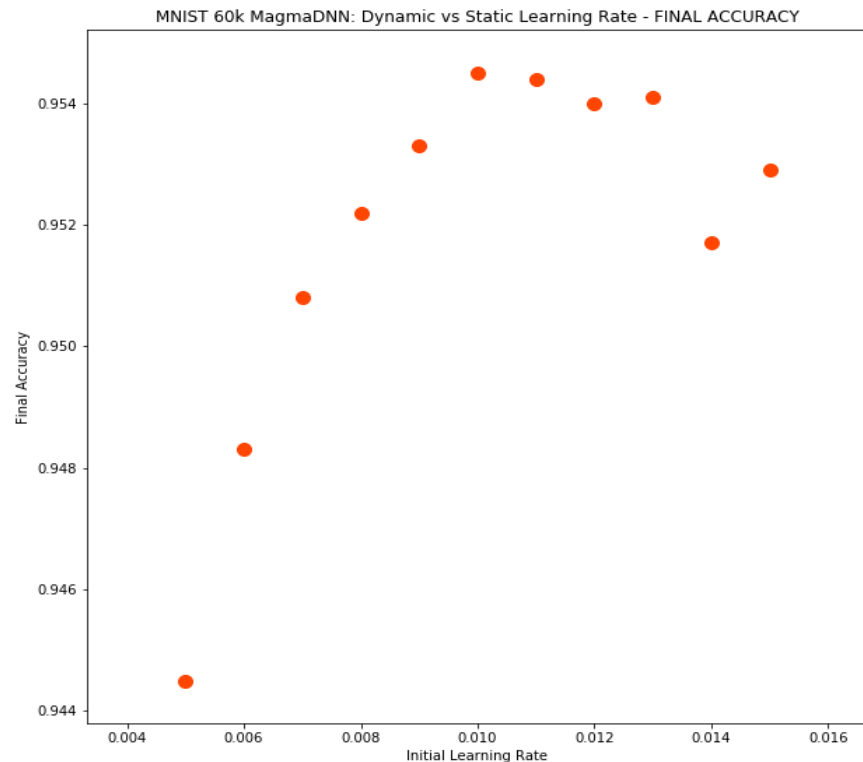
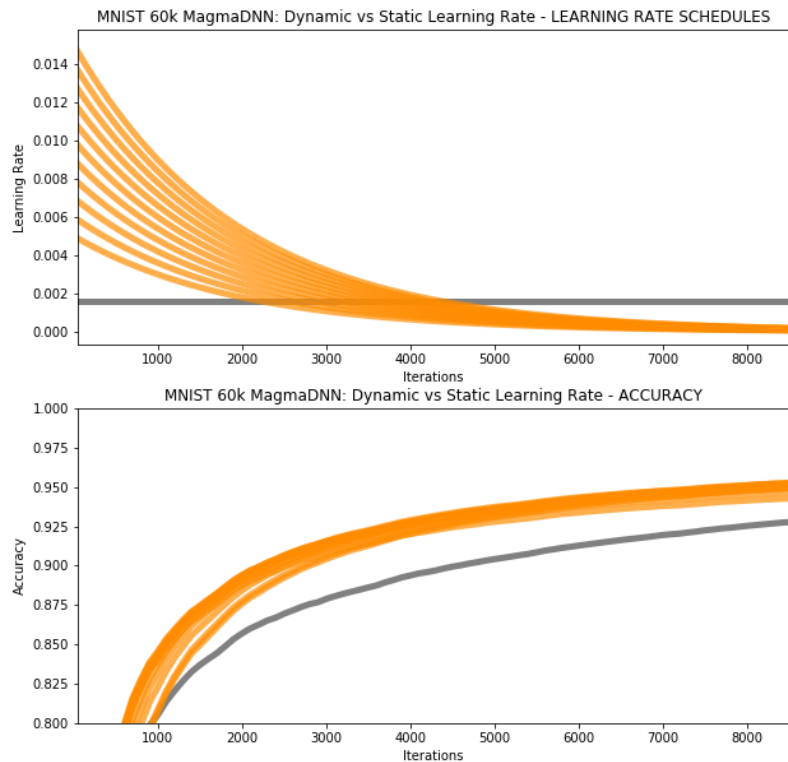
- Network Structure: In -> FCB -> Sig -> FCB -> Sig -> FCB -> Out
- Weight Optimizer: Stochastic Gradient Descent
- Number of Epochs = 5
- Batch Size = 32

- **Benchmark:** constant learning rate = .0016

- **Experiments:** dynamic learning rate schedules with variable initial values

\*FCB := Fully Connected Layer with Bias  
\*Sig := Sigmoid

# PBT: Analysis - Dynamic Learning Rate



# PBT: Goals

- Extend the OpenDIEL Grid Search Application to have PBT functionality, i.e. stochasticity and evolution.
- Program more custom MagmaDNN classes to explore the effect of tuning Convolutional Neural Network hyperparameters.
- Implement PBT on MagmaDNN and OpenDIEL with a distributed Worker, and overcome the Ray-Tune bottleneck.

**Thanks for listening!**  
*-The Hyperparameter Team*

