

Towards Scalable Parallel Training of Deep Neural Networks

Sam Adé Jacobs, Nikoli Dryden, Roger Pearce & Brian Van Essen

November 13, 2017

MLHPC 2017

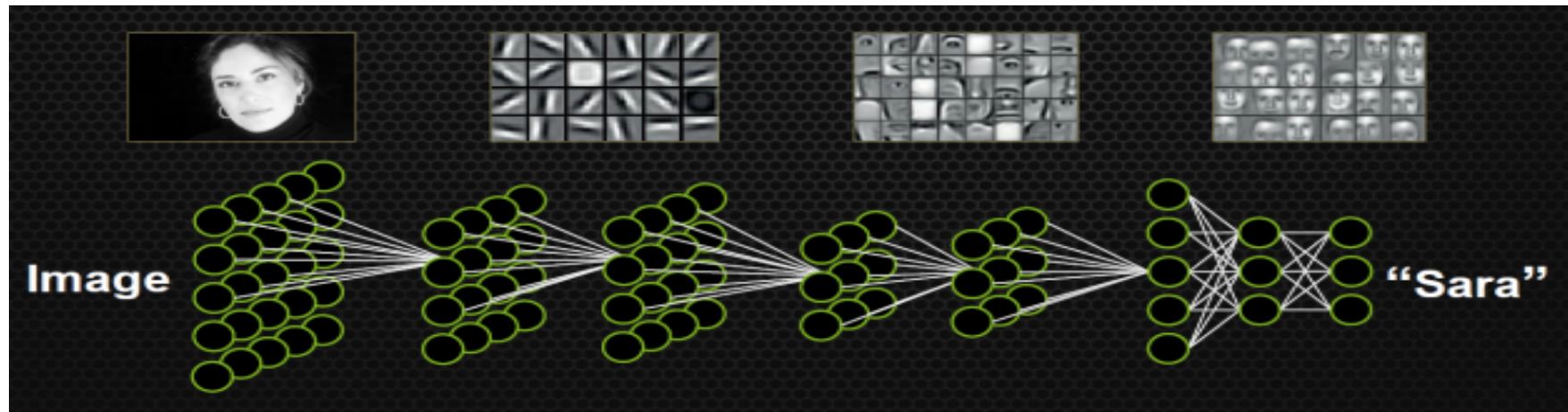


LLNL-PRES-741592

This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under contract DE-AC52-07NA27344. Lawrence Livermore National Security, LLC

Motivation

Scale Matters!

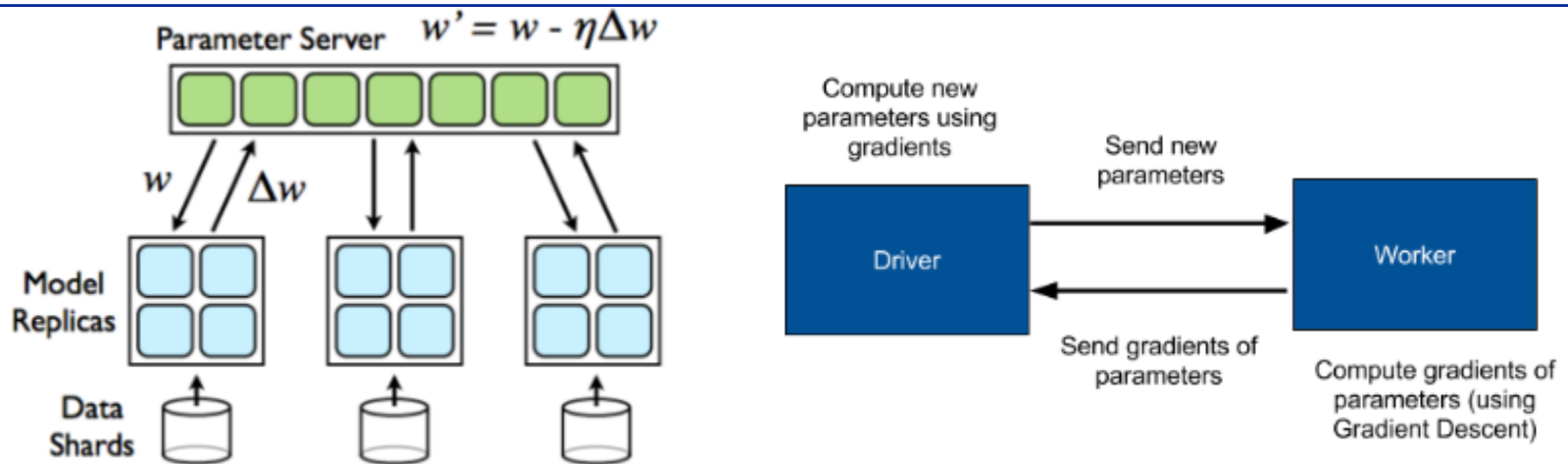


Accelerate deep learning training as model and data scale grow

The image above and subsequent images are exclusive copyright of their respective sources

Prior Arts

Parameter server (manager-worker architecture) and its variants



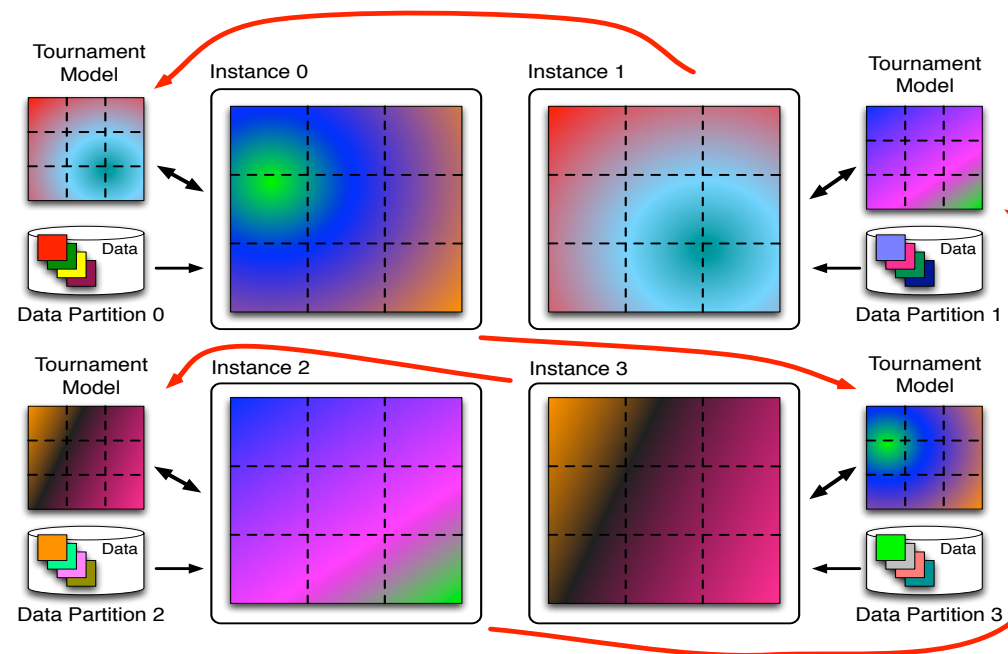
Scalability is limited by (global) communication overhead

Livermore Tournament Fast Batch Learning (LTFB)

Multi-Level Tournament Voting with Random Pairing

Framework

- Multiple independent trainers
- Periodically exchange model with random peer
 - Multiple mini batches or epochs of training
- Run local tournament to select current or exchanged model
- Continue training using winning model



Livermore Tournament Fast Batch Learning (LTFB) Example

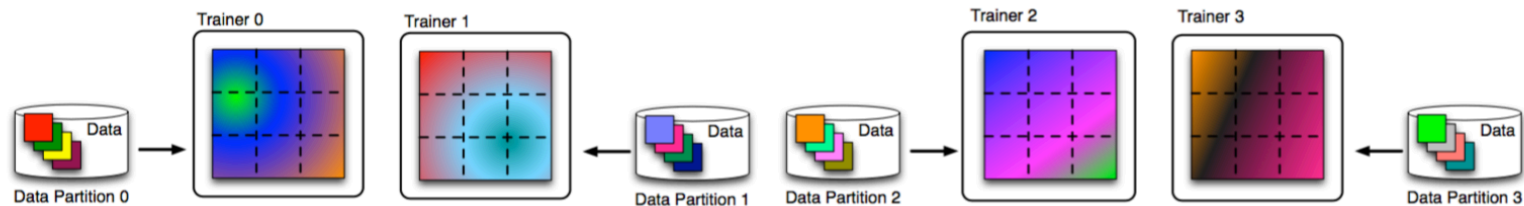


Figure 1: LTFB Architecture: training data is partitioned among p (4) trainers. Each trainer has an instance of the model and a partition of the training data set. For this example, each trainer shown is internally parallelized over 9 HPC nodes. Multi-color patterns represent the trainer's computed gradient for its weight matrices.

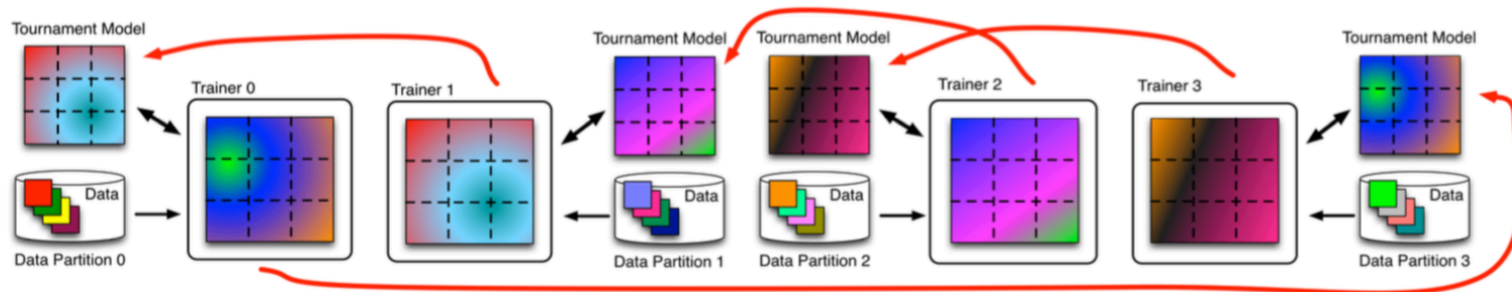


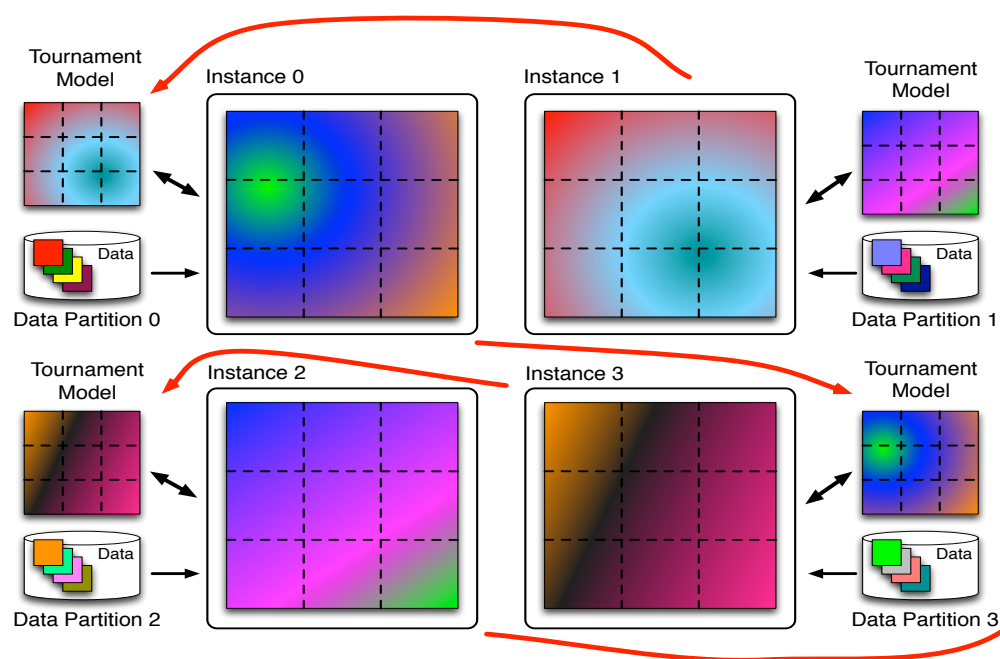
Figure 2: LTFB Architecture: models are exchanged between trainers. Each trainer evaluates the state of its current model and model received from another trainer against their private holdout tournament dataset. The winning model is kept for next round of training and loser is discarded.

Livermore Tournament Fast Batch Learning (LTFB)

Multi-Level Tournament Voting with Random Pairing

Key Benefits

- Scalable peer to peer communication
- Use parallel resources to :
 - reduce total time to train
 - to achieve higher quality solution via more extensive training
- Streamlines hyper parameter exploration



Experimental Setup

Dataset & DNN Architectures

- CIFAR10 and ImageNet1K dataset
- Network architectures taken from original (Berkeley) Caffe
 - GoogleNet for ImageNet1K
- Same hyper parameter (learning rate, convolutional filter sizes, optimizer etc) as in original Caffe

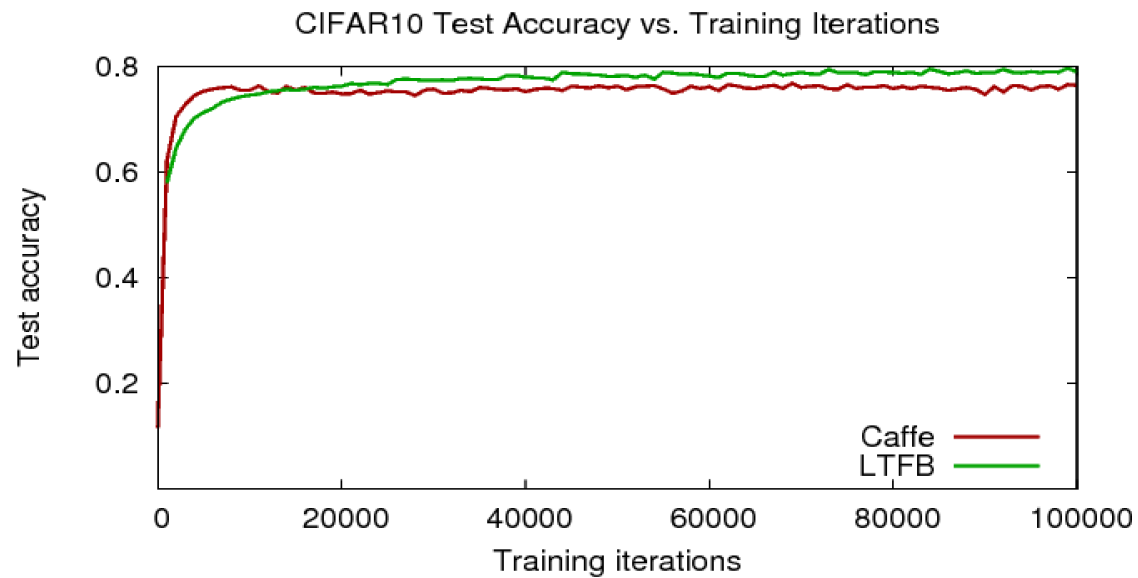
HPC Machines

- Surface
 - 156 Intel Xeon (Sandybridge) compute nodes
 - 16 CPU cores, 256GB memory, and 2 Tesla K40 GPUs per node
- Ray
 - Sierra (CORAL) early access system
 - 54 IBM Power8+ compute nodes
 - 20 CPU cores, 256GB memory, and 4 Tesla P100 (Pascal) GPUs per node

Experimental Results

CIFAR10

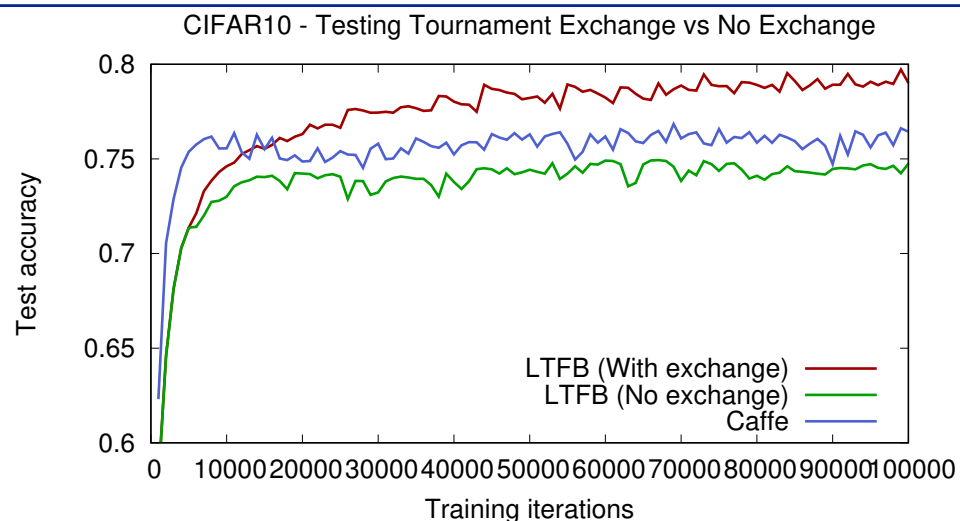
4 LTFB Trainers vs Caffe



Experimental Results

CIFAR10

Digging Deeper : LTFB with(out) Tournament Exchange

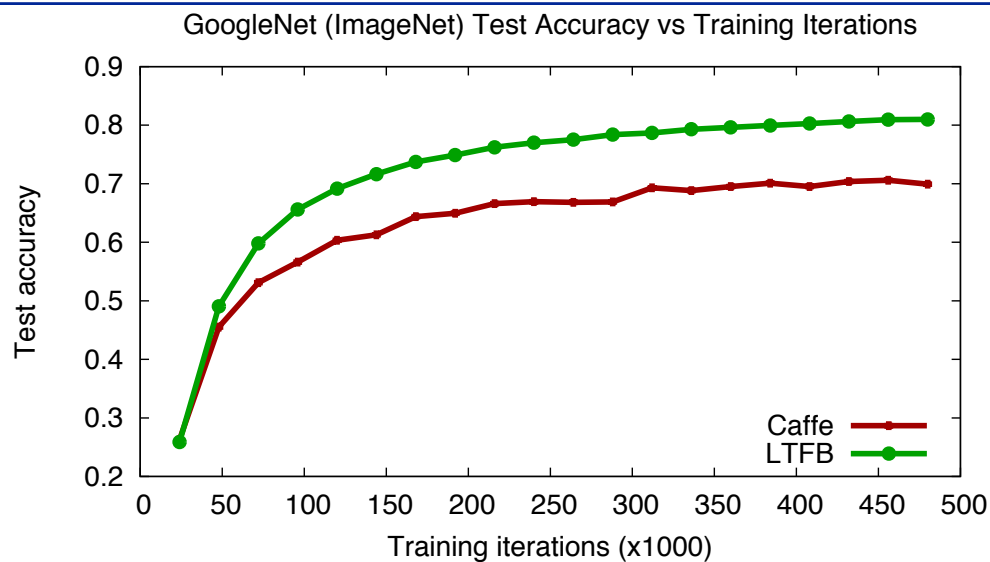


Without tournament exchange (knowledge sharing), Caffe beats LTFB

Experimental Results

ImageNet (GoogleNet)

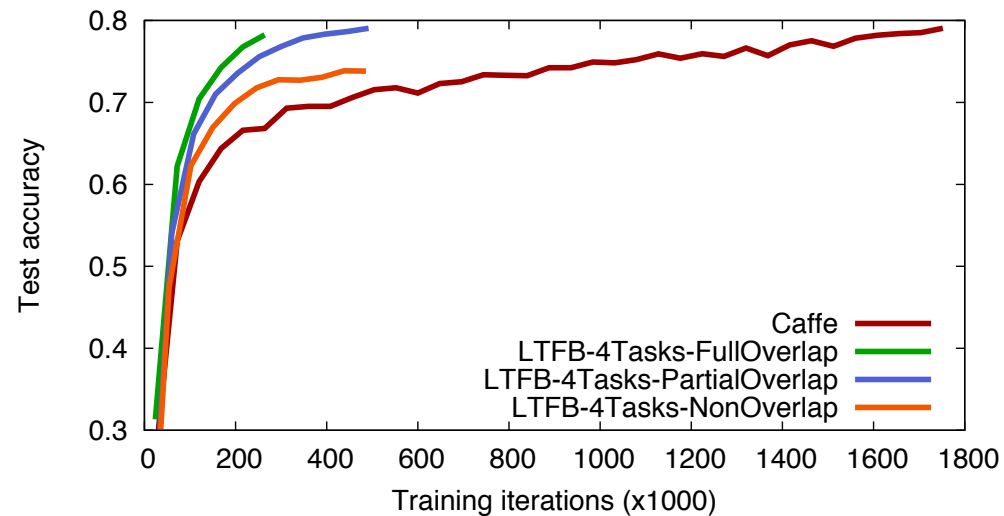
4 LTFB Trainers vs Caffe



Experimental Results

ImageNet (GoogLeNet)

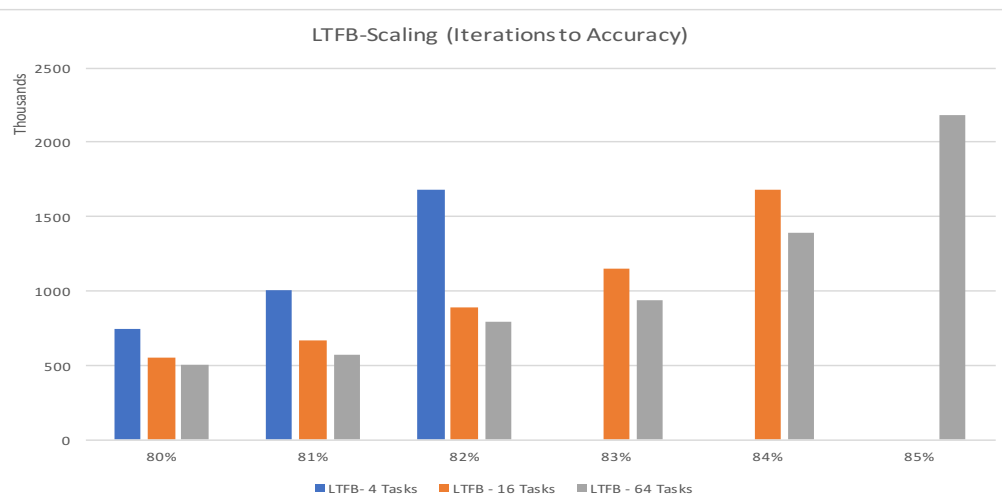
Varying Data Partitioning



Experimental Results

ImageNet (GoogleNet)

LTFB Scaling Study: Accelerating Time to Solution



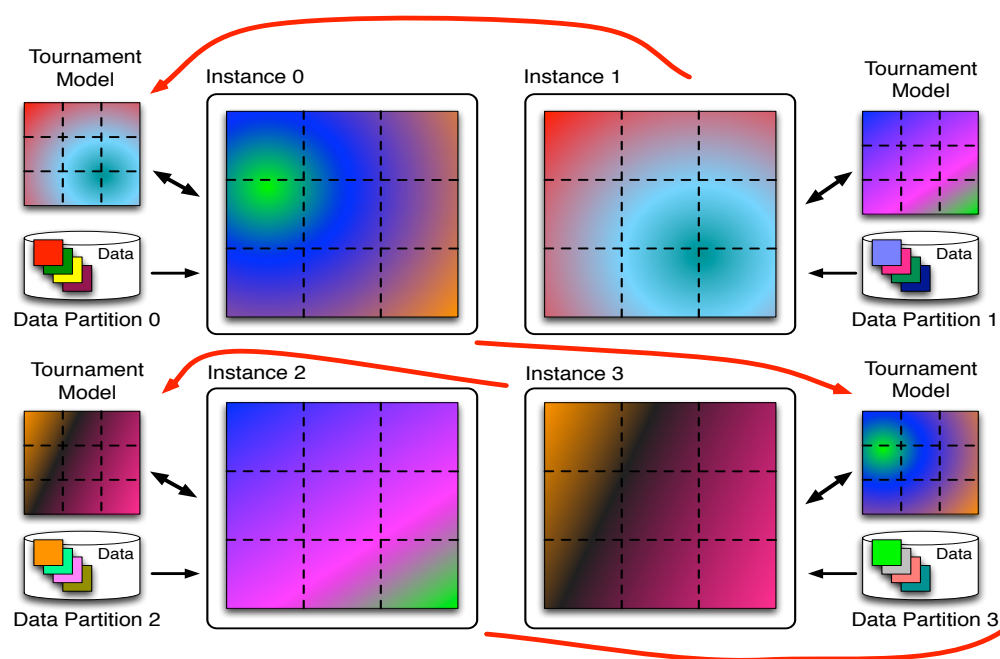
Increased concurrency in LTFB pays off at higher accuracy levels, every round of voting explores multiple paths from the strongest models

Livermore Tournament Fast Batch Learning (LTFB)

Multi-Level Tournament Voting with Random Pairing

Key Benefits

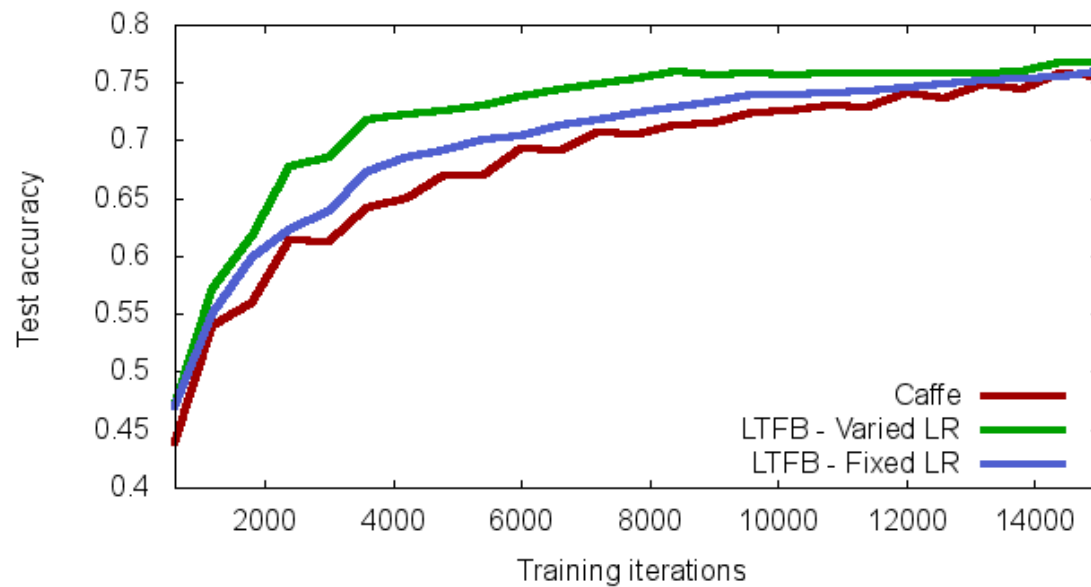
- Scalable peer to peer communication
- Use parallel resources to reduce total time to train
- Streamlines hyper parameter exploration



Experimental Results

Hyper parameter Exploration

CIFAR10: 4 LTFB trainers with different LR vs fixed LR using default mini batch



Experimental Results

Hyper parameter Exploration

Setup

- Task: explore mini batch / learning rate pairing mismatch
- ImageNet dataset, GoogleNet
- A set of 16 learning rates, each given to an LTFB trainer
- Caffe has one of the learning rate from the set
- Tournaments help to dynamically select the best hyper parameter for the current state of training

LTFB achieves ~66% accuracy in 4 epochs, Caffe did not learn

	Round 1	Round 2	Round 3	Round 4
Trainer 0	23.8633	42.9233	53.8433	64.4033
Trainer 1	25.81	47.2767	58.2667	66.1967
Trainer 2	14.4367	47.2767	57.3967	65.2867
Trainer 3	1.26333	43.3733	51.8733	64.4033
Trainer 4	24.1833	48.6967	60.0233	65.3467
Trainer 5	0.473333	0.483333	56.2866	66.1967
Trainer 6	0.49	31.1167	54.8067	65.1033
Trainer 7	0.48	35.1933	54.2467	63.6833
Trainer 8	1.26333	43.3733	54.2467	63.6801
Trainer 9	0.46	35.1933	53.7067	65.0367
Trainer 10	0.493333	0.516666	0.52333	60.6334
Trainer 11	23.8633	43.3733	53.8433	61.1934
Trainer 12	0.493333	0.48	0.48	58.24
Trainer 13	0.49	48.6967	54.9966	60.6334
Trainer 14	0.49	0.473333	54.9966	63.6833
Trainer 15	0.496666	0.48	53.7067	58.24

Summary

Our Contributions and Future Work

- Present a new multi-level tournament voting parallel training algorithm that uses scalable peer to peer communication
- Our framework streamlines hyper parameter exploration by allowing diversity in each independently trained model, and minimizing the time spent training with suboptimal parameters
- Demonstration of feasibility of the new approach on image classification tasks using HPC clusters
- **Future Work:** Big Science Big Data HPC Deep Learning
 - Extremely large (infinite streams of scientific) data sets

