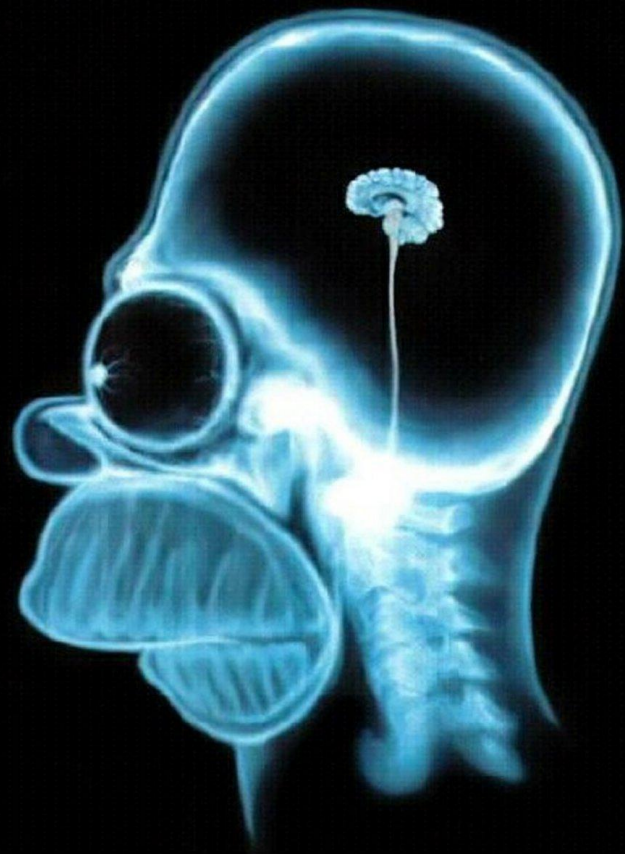




# Transparent parallelization of neural network training

Cyprien Noel

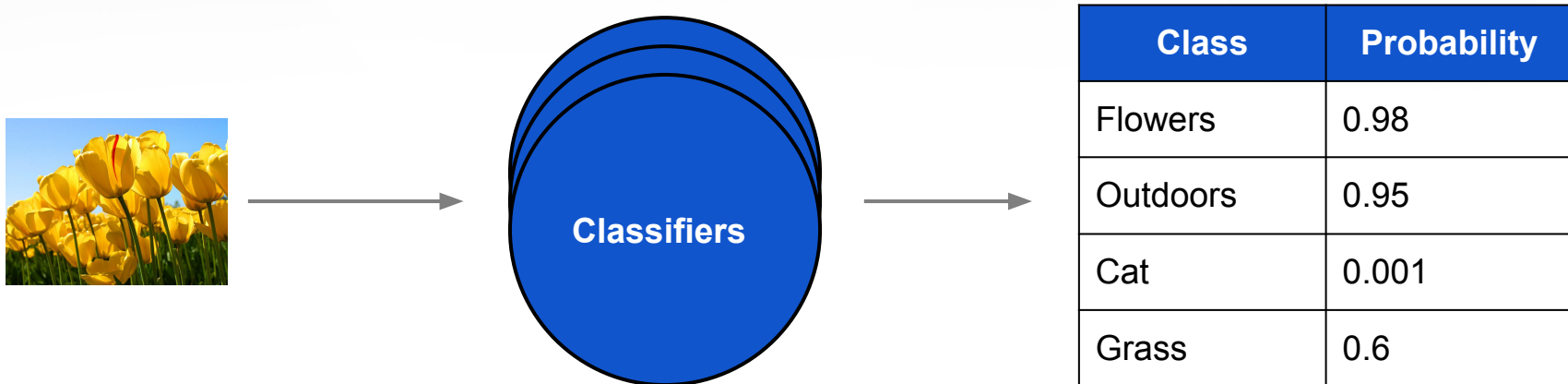
Flickr / Yahoo - GTC 2015



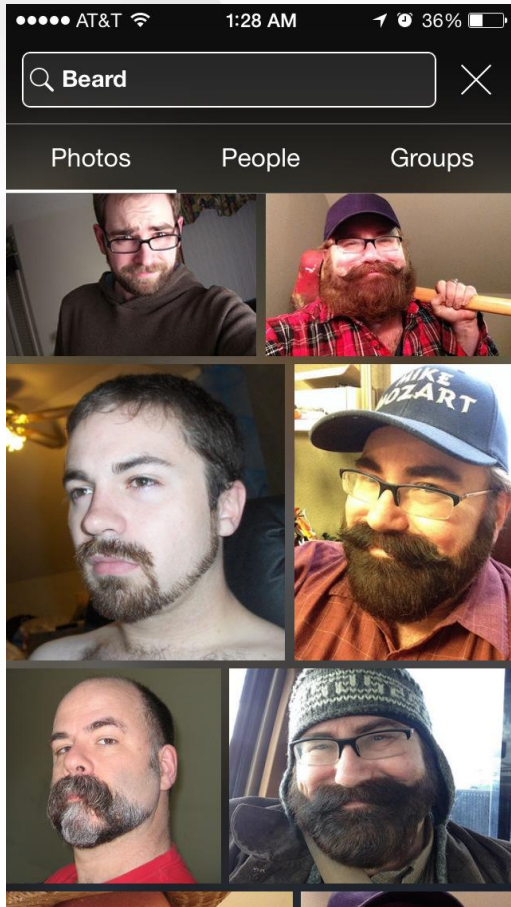
# Outline

- **Neural Nets at Flickr**
- Training Fast
  - Parallel
  - Distributed
- Q&A

# Tagging Photos



Any photo on Flickr is classified using computer vision



# Auto Tags Feeding Search

- Flowers
- Buildings
- Outdoors
- Beach
- Beards
- and more!

# Tagging the Flickr corpus

- Classify millions of new photos per day
- Apply new models to billions of photos
- Train new models using Caffe



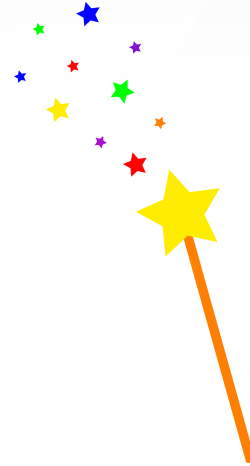
# Training new models

- Manual experimentation
- Hyperparameter search
- Limitation is training time

→ **Parallelize Caffe**

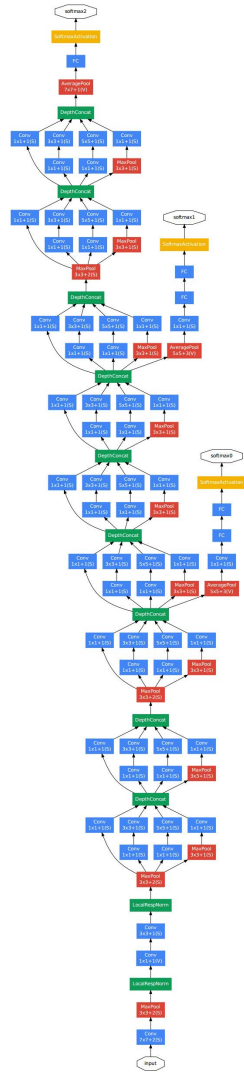
# Goals

- “Transparent”
  - Code Isolation
  - Existing Models
  - Globally connected layers
- Existing Infrastructure



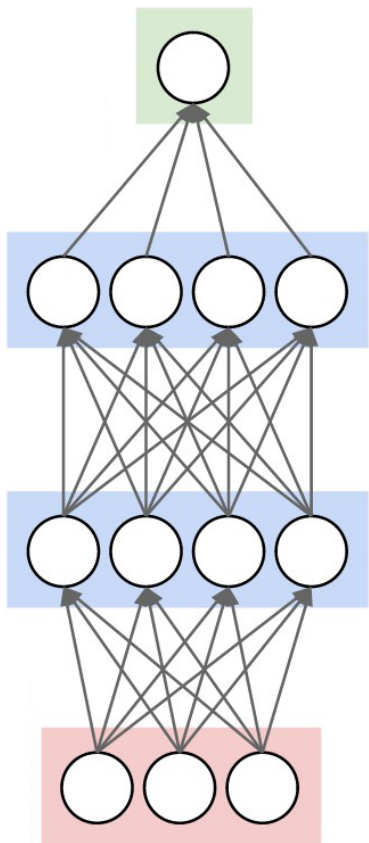
# Outline

- Neural Nets at Flickr
- **Training Fast**
  - Parallel
  - Distributed
- Q&A



GoogLeNet, 2014

# Ways to Parallelize



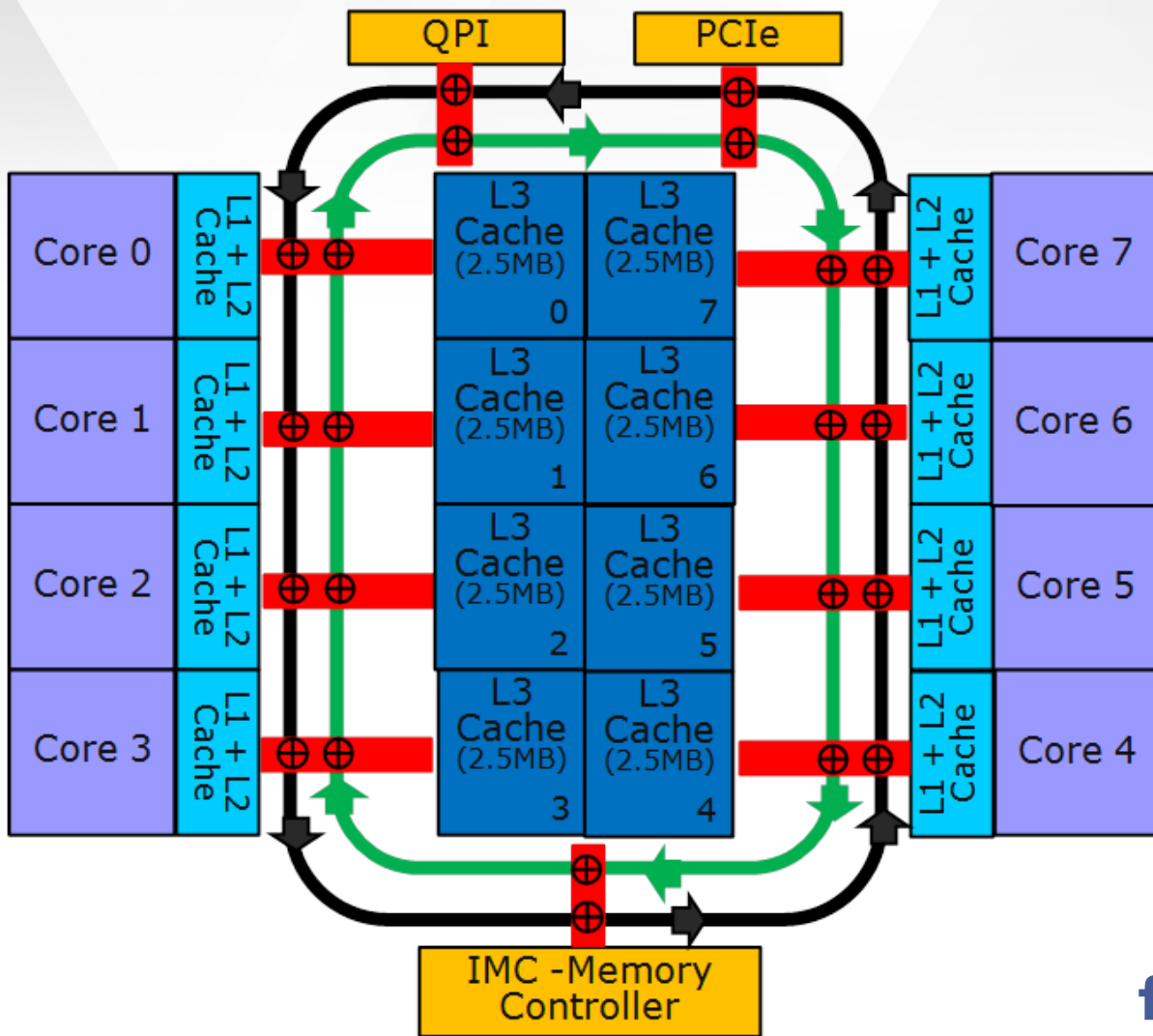
- Model
  - Caffe team enabling this now
- Data
  - Synchronous
  - Asynchronous

# Outline

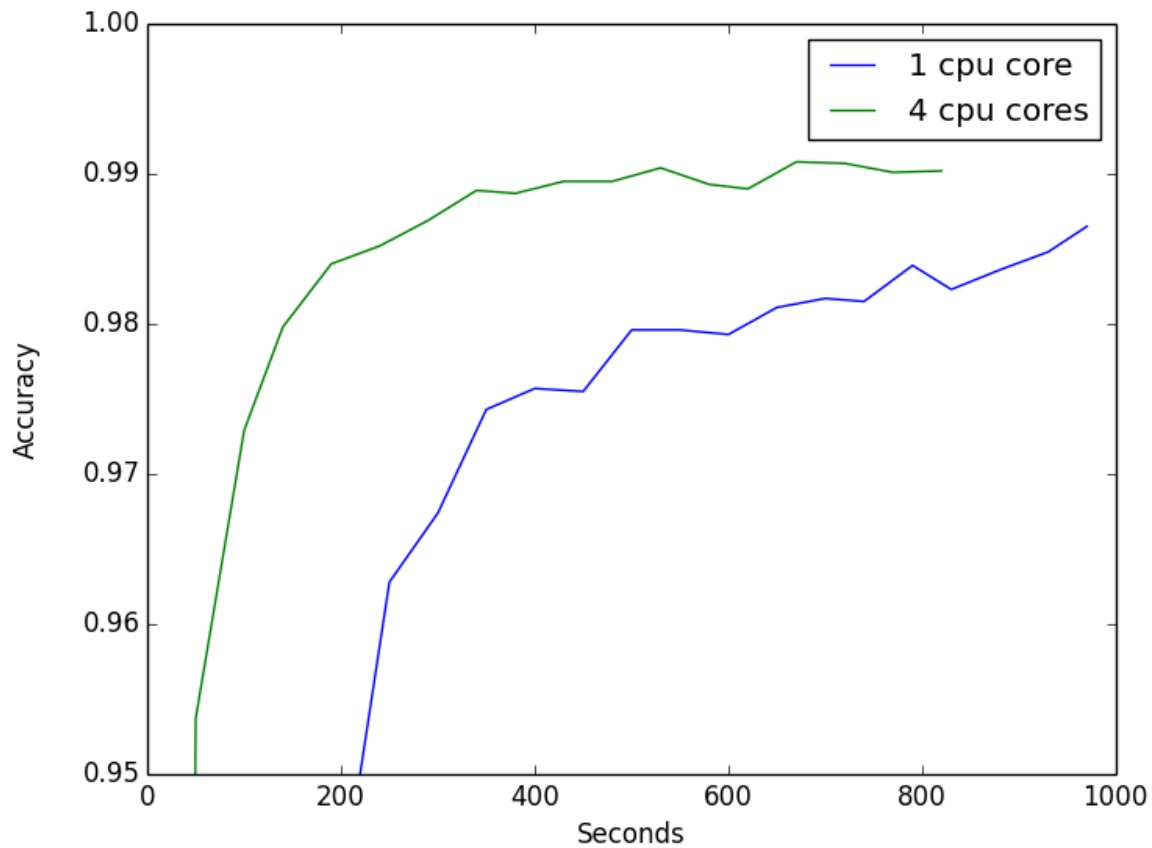
- Neural Nets at Flickr
- Training Faster
  - **Parallel**
  - Distributed
- Q&A

# First Approach: CPU

- Hogwild! (2011)
- Cores read and write from shared buffer
- No synchronization
- Data races surprisingly low



# MNIST CPU




# Hogwild

- Plateaus with core counts
- Some Potential
  - On a grid
  - With model parallelism

But we are at GTC



# GPU Cluster

- A lot of time spent preparing experiments
- Code Deployment
  -  docker
- Data Handling
  - On the fly datasets for “big data”

# Outline

- Neural Nets at Flickr
- Training Fast
  - Parallel
  - **Distributed**
- Q&A

## Second Approach: Lots of Boxes



## Second Approach: Lots of Boxes

- Exchange gradients between nodes
- Parameter server setup
- Easy: move data fast

# GPU memory - PCI - Ethernet



## Second Approach: Lots of Boxes

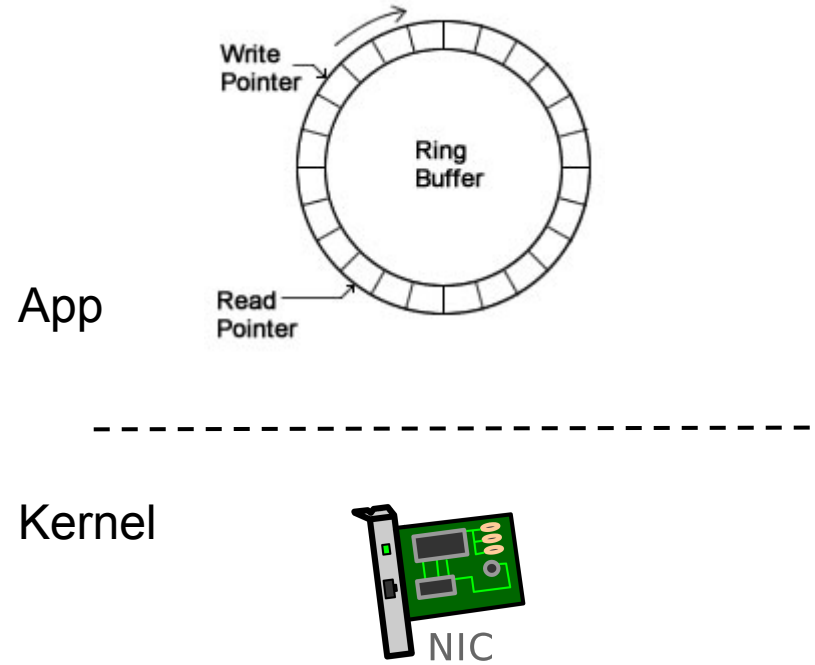
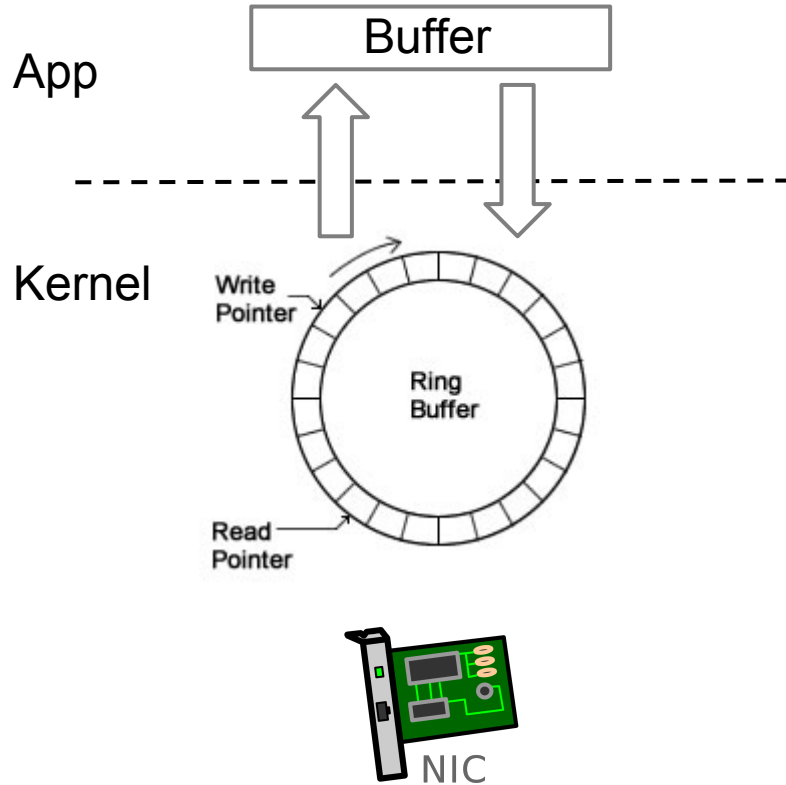
- $230\text{MB} * 2 * N$  per batch
- TCP/UDP chokes
  - Machines unreachable
- No InfiniBand or RoCE

## Second Approach: Lots of Boxes

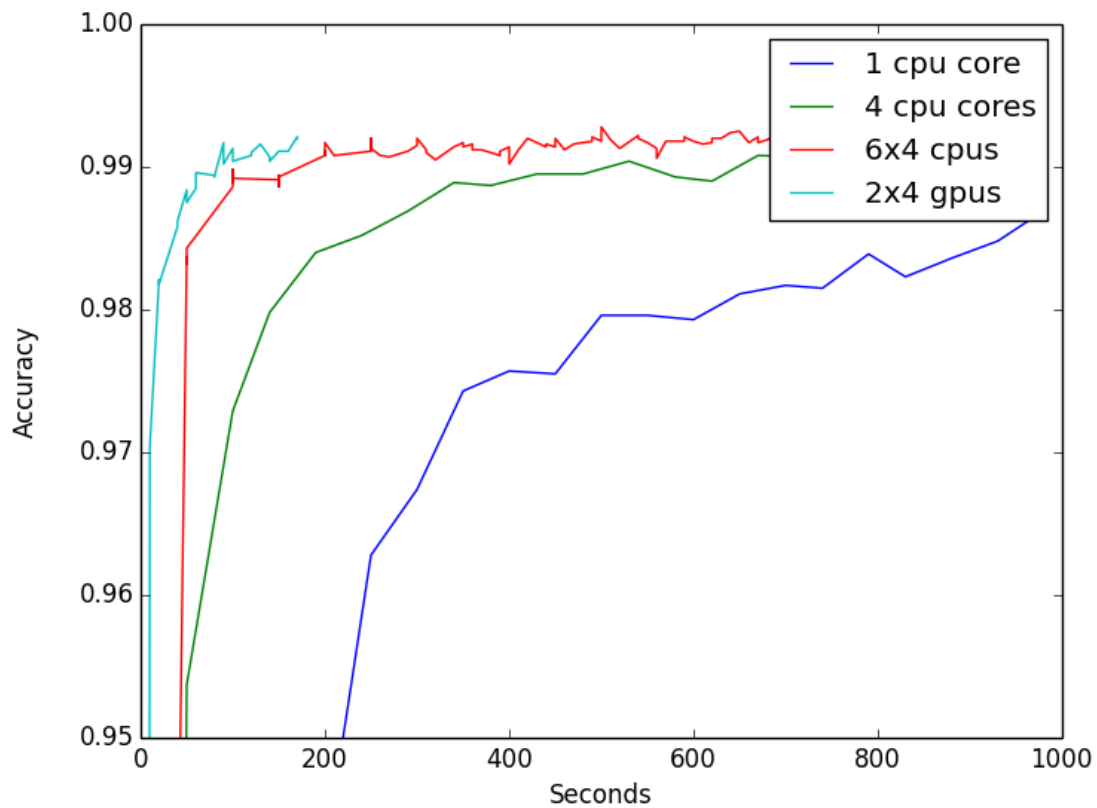
- Modify Caffe: chunk parameters



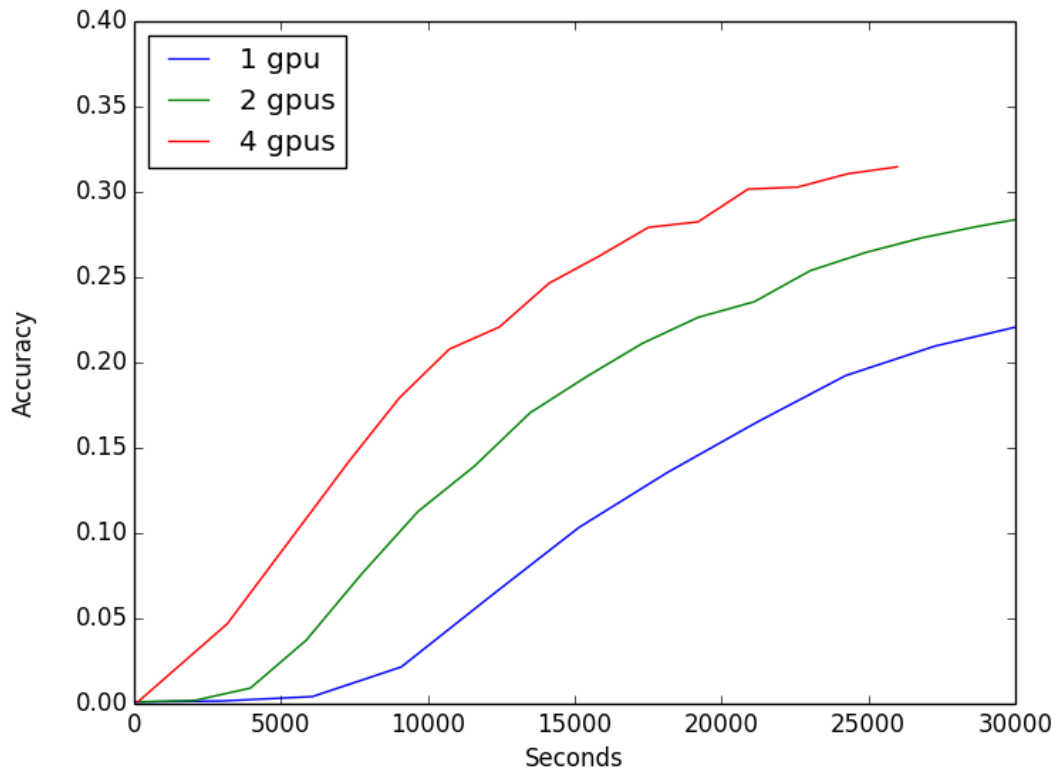
# Packet\_mmap



# MNIST



# ImageNet



# NVIDIA

- Large Machines
  - 4 or 8 GPUs
  - Root PCI switches
- InfiniBand

## Third Approach: CUDA P2P

- GPUs on single machine
- Data Feeding
  - Caffe Pipeline
  - Async Streams

# State of Things

- Async  $\sim 8x$  but no momentum
- Sync  $\sim 2x$
- Combining both, and model parallelism
- Working on auto tuning of params (batch, rate)
- Different ratios of compute vs. IO

# Takeaway

- Check Caffe, including Flickr's contributions
- CUDA + Docker = Love
- Small SOC servers might be interesting for ML

# Thanks!

Flickr vision team  
Flickr backend team  
Yahoo labs

[cypof@yahoo-inc.com](mailto:cypof@yahoo-inc.com)