# Towards Better DL Frameworks

**Yangqing Jia**

Research Lead on AI Platforms, Facebook

Source: XKCD, [Girshick et al. CVPR 2014]

# The Needs
## Two sides of the same coin

- Researchers: "I will need to reproduce the ResNet paper."

- Companies: "I need to apply DL to drive cars."

# Democratizing Deep Learning w/ Caffe

Getting AlexNet running in 10 mins

- A grad student driven project
- Started by doing one job really well: image classification
- Adopted by industry participants
- Popular deep learning framework run by a non-profit.

# http://caffe.berkeleyvision.org/



| Maximally accurate | Maximally specific |
| --- | --- |
| cat | 1.79559 |
| feline | 1.74239 |
| domestic cat | 1.71551 |
| tabby | 0.95449 |
| domestic animal | 0.77145 |

# What makes a better DL library?

???

# "MAPS"

# "MAPS"
# -
# Scalability

# Scalability
## Run fast, run far

"How do I train on
multiple GPUs and machines?"

- Probably the most question we got from Caffe users

# Scalability

Run fast, run far

| L1 | L2 | L3 | L3b | L2b | L1b | U3 | U2 | U1 |

# Scalability

## Run fast, run far

# Scalability

Run fast, run far

# The Return of MPI
"I'm your father", said Allreduce.



Allreduce
Tree based - O(MlogN)
Ring based - O(M)
etc.

"Weak Scaling": Speedup in Throughput

Speedup

# GPUs

Resnet-50  Inception  Ideal

# Scalability
## Sitting on top of giants



... and many more

# "MAPS"
# -
# Portability

# Portable System
## Cloud, Mobile, IoT, Cars, Drones, Coffee makers

AI Math and Algorithms

Deployment Platforms

# Portable System
## Cloud, Mobile, IoT, Cars, Drones, Coffee makers



Model

auto predictor =
    caffe2::Predictor(model_file)

public class Predictor implements
    Caffe2ModelInterface;

# **Portable System Challenges**
## Still, a lot of thoughts needed

- Limited computation
- Battery life is a thing
- Our models may be luxurious
- Ecosystem less developed

# "MAPS"
-
# Augmented Comp Patterns

# **Augmented Comp Patterns**
Forget about float dense math, the world is bigger

- Quantized Computation
- Sparse Math Libraries
- Model Compression
- Rethinking Existing Operations

# Quantized Computation

## Forget about float, the world is bigger

# Quantized Computation

## Forget about float, the world is bigger

float add — 0.9
fp16 add — 0.4
fixed16 add — 0.05
fixed8 add — 0.03

float mul — 4.0
fp16 mul — 1.0
fixed8 mul — 0.2

# Why?



40x Efficiency vs CPU, 8x Efficiency vs FPGA

■ CPU   ■ FPGA   ■ 1x M4 (FP32)   ■ 1x P4 (INT8)

AlexNet

# Rethinking Existing Operations
## ResNEXT is coming to town

AlexNet Group
Conv

ResNext

# Augmented Math Challenges

Forget about float, the world is bigger

- Solutions
  - Eigen fp16
  - CuDNN
  - NNPack
  - gemmlowp

- Challenges
  - Seamless conversion?
  - Model training?
  - Performance tuning?
  - ...

# "**M**APS"
# -
# Modularity

# A Repeated Pattern

Many key components in deep learning
are
reusable
across frameworks.

# In 2013 it used to be...

Caffe

Torch

Theano

...

# Unix Philosophy?

## or, "UnFramework"

**Applications**

**Caffe, Torch, TF, MXNet, etc...**

| DataBases | Core Math | Comms | Low Level | Compilers |
|-----------|-----------|-------|-----------|-----------|
| LevelDB | Eigen | NCCL | CUDA | |
| RocksDB | CuDNN | MPI | OpenGL | |
| Hadoop | NNPack | ZeroMQ | OpenCL | |
| Amazon S3 | THNN | Redis | Vulkan | |
| your old disk | MKL | ... | ... | |

# **MAPS for a good framework**

| **Modular** Designs | **Augmented** Mathematics | **Portable** System | **Scalability** |
|---|---|---|---|
| Interface to Existing Toolkits | Optimized Math Libraries | Efficient Mobile Runtimes | Tuned Collective Primitives |

+

**Flexible Framework Design**

# No Silver Bullet?

# There is no silver bullet

D4J etc.

TensorFlow

Theano

Caffe

Torch

← ────────────────────────────── →

**Industry**:
Stability
Scale & speed
Data Integration
Relatively Fixed

**Research**:
Flexible
Fast Iteration
Debuggable
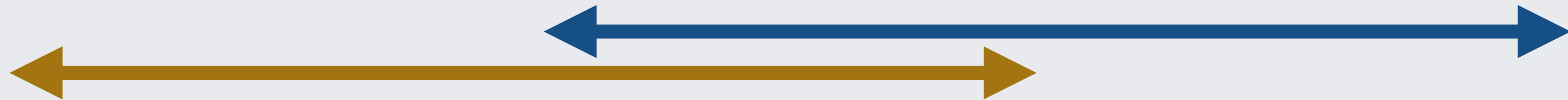Relatively bare-

# There is no silver bullet

Caffe                               Torch

**Industry**:                       **Research**:
Stability                               Flexible
Scale & speed                       Fast Iteration
Data Integration                      Debuggable
Relatively Fixed                   Relatively bare-

"In open source, we feel strongly that
to really do something well,
you have to get a lot of people involved."

— Linus Torvalds

# Thank you!

**Towards Better Deep Learning Frameworks**

Yangqing Jia, Research Lead on AI Platforms, Facebook