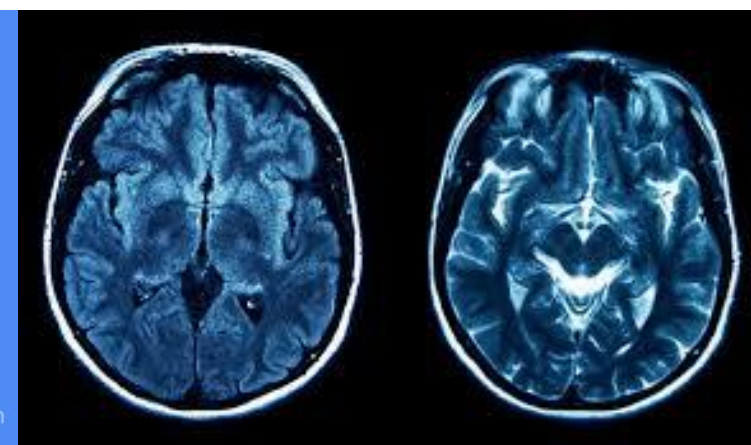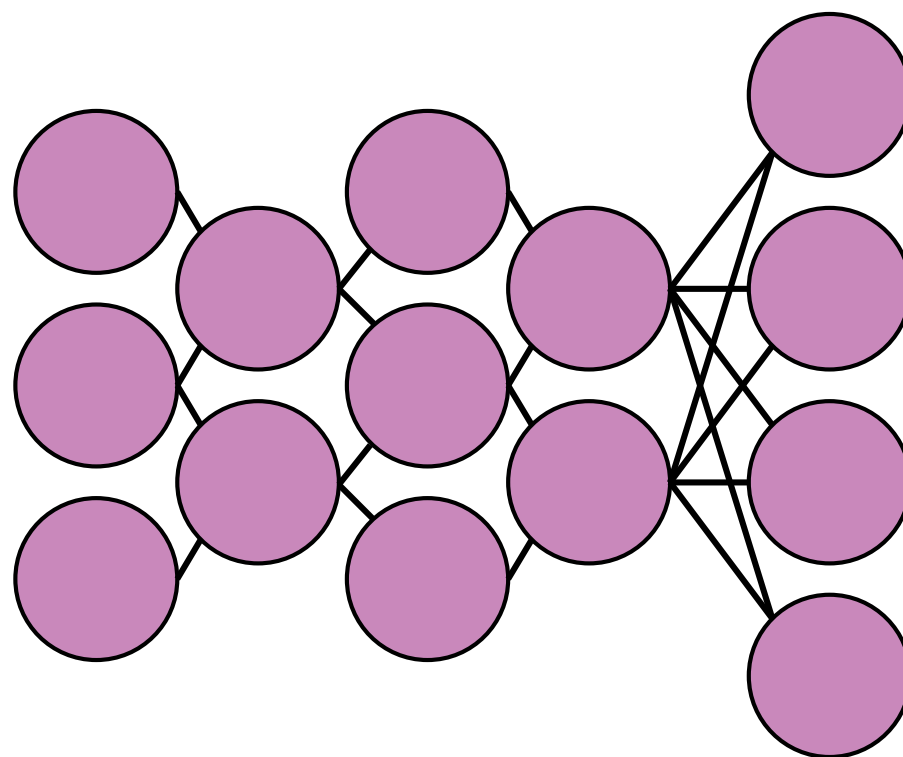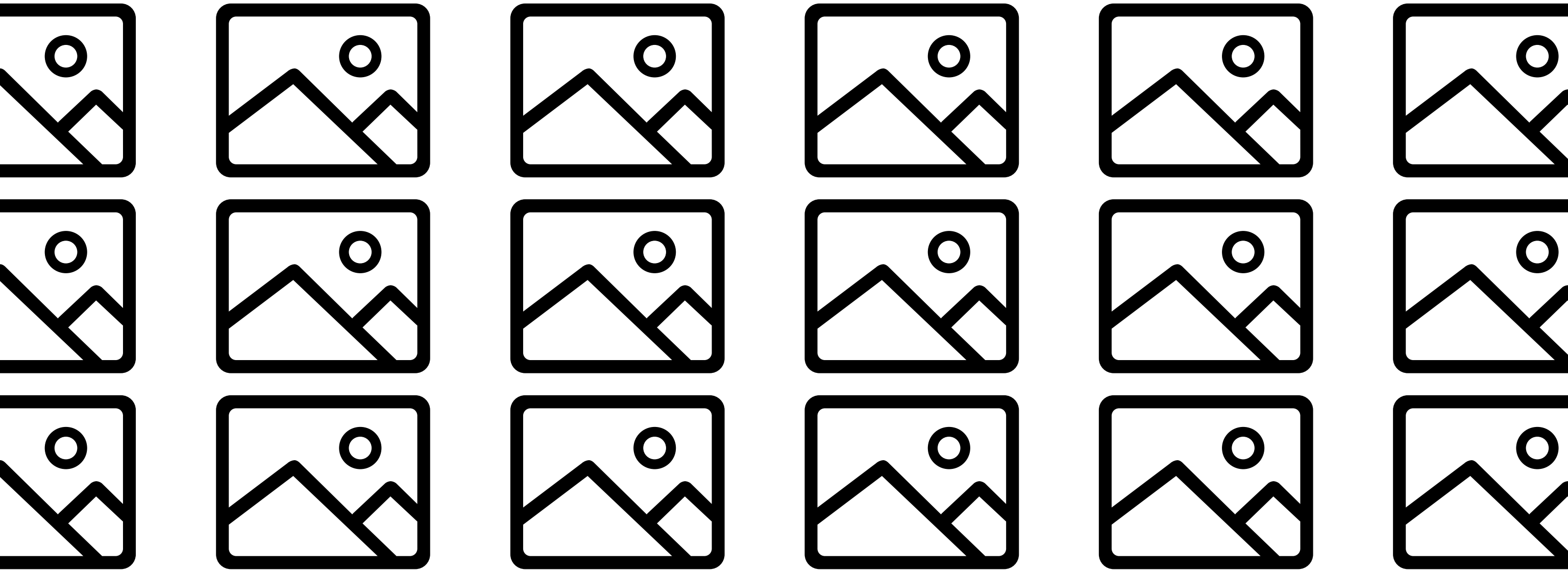# Accelerating Deep Learning with the Biggest Losers

**Angela H. Jiang**, Daniel L.-K. Wong, Giulio Zhou,
David G. Andersen, Jeffrey Dean,  Gregory R. Ganger,
Gauri Joshi, Michael Kaminsky, Michael A. Kozuch,
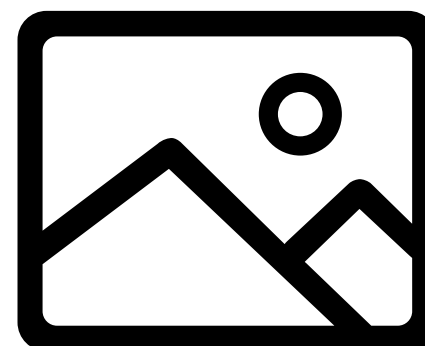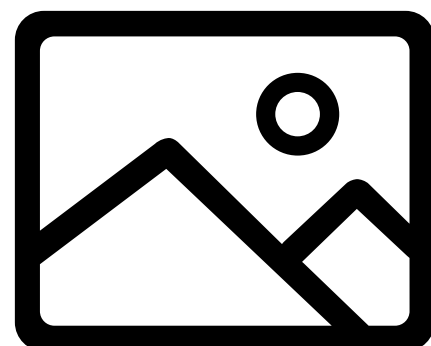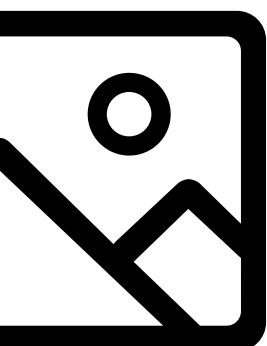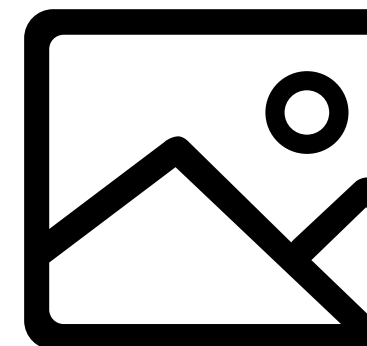Zachary C. Lipton, Padmanabhan Pillai

# Deep learning enables emerging applications

# DNN training analyzes many examples

# Selective-Backprop prioritizes informative examples

# DNN basics

**How to use and train a DNN**

# Example task: Image classification

# DNN inference: From image to "Dog"

# Training DNNs relies on a labeled dataset



**Class:
Dog**

**Class:
Bird**

**Class:
Lizard**

**Class:
Dog**

# DNN training: Determining the weights



**Class 1**

# DNN training: Determining the weights via backpropagation



Class
3

# DNN training: Determining the weights via backpropagation



**Class 2**

# DNN training analyzes an example many times

# DNN training analyzes an example many times

**Epoch 250**

Forward pass

Predicted Class

Target Class

Loss

Backwards pass

# SelectiveBackprop targets slowest part of training

**Legend:**
- Backwards (training)
- Forwards (selecting)
- Forwards (training)
- Other (setup, overhead)

# Not all examples are equally useful



data_batch_1.mat index 1681
data_batch_1.mat index 2362
data_batch_1.mat index 8559
data_batch_1.mat index 8825

data_batch_1.mat index 3786
data_batch_1.mat index 2163

data_batch_2.mat index 3239
data_batch_2.mat index 9450
data_batch_3.mat index 1137
data_batch_3.mat index 6315

data_batch_3.mat index 9429
data_batch_4.mat index 7995
data_batch_4.mat index 8571
data_batch_5.mat index 1091

data_batch_1.mat index 1851
data_batch_1.mat index 9049
data_batch_4.mat index 5925

data_batch_5.mat index 4273
data_batch_5.mat index 7588
data_batch_5.mat index 9096
data_batch_5.mat index 9875

# Prioritize examples with high loss



**Examples with low loss**

**Examples with high loss**

# *Selective Backprop algorithm*

# DNN training analyzes an example many times

# Bad idea #1:
## Deciding with a hard threshold

*if loss > threshold: backprop()*

# Bad idea #2:
## Deciding probabilistically with absolute loss

*P(backprop) = normalize(loss, 0, 1)*

# Good idea:
Use relative probabilistic calculation

$$P(backprop) =$$
$$Percentile(loss, recent\ losses)^B$$

# Example of probability calculation



CDF(loss) = 0.9

Loss = 2.3

Probability = .96

CDF(loss) =0.9

# Selective-Backprop approach

**Forward propagate example through the network**



**Calculate usefulness of backpropping example based on its accuracy**

$$P(\text{Backprop}) = \text{L2 Dist}^2\left(\text{Output}, \text{Target}\right)$$

**Decide probabilistically if we should backprop**

Output → Loss →

# StaleSB reduces forward passes

**Forward propagate example through the network** *every n epochs*



**Calculate usefulness of backpropping example based on its accuracy**

$$P(\text{Backprop}) = \text{L2 Dist}^2(\text{Output}, \text{Target})$$

**Decide probabilistically if we should backprop**

Output → Loss →

# *Evaluation of Selective Backprop*

# Datasets



**CIFAR10**
**60,000 Training Images**

**CIFAR100**
**60,000 Training Images**

**SVHN**
**604,388 Training Images**

# Train CIFAR10 to 4.14% (1.4x Traditional's final error)

*SB: 1.5X faster*

*StaleSB: 2X faster*

Legend:
- Backwards (training)
- Forwards (selecting)
- Forwards (training)
- Other (setup, overhead)

X-axis: Traditional Epoch=120, Selective-Backprop Epoch=120, Stale-SB (n=3) Epoch=120

Y-axis: Wall-clock Time (hours)

Train CIFAR100 to 25.5% (1.4x Traditional's final error)

SB: 1.2X faster

StaleSB: 1.6X faster

Train SVHN to 1.72% (1.4x Traditional's final error)

SB: 3.5X faster

StaleSB: 5X faster

# SB on CIFAR10 targets hard examples



**3% correct w/ Traditional**

**29% correct w/ SB**

**Output** = [0.1, 0.3, 0.6]

**Target confidence** = 0.3

✓ **Selective-Backprop accelerates training**

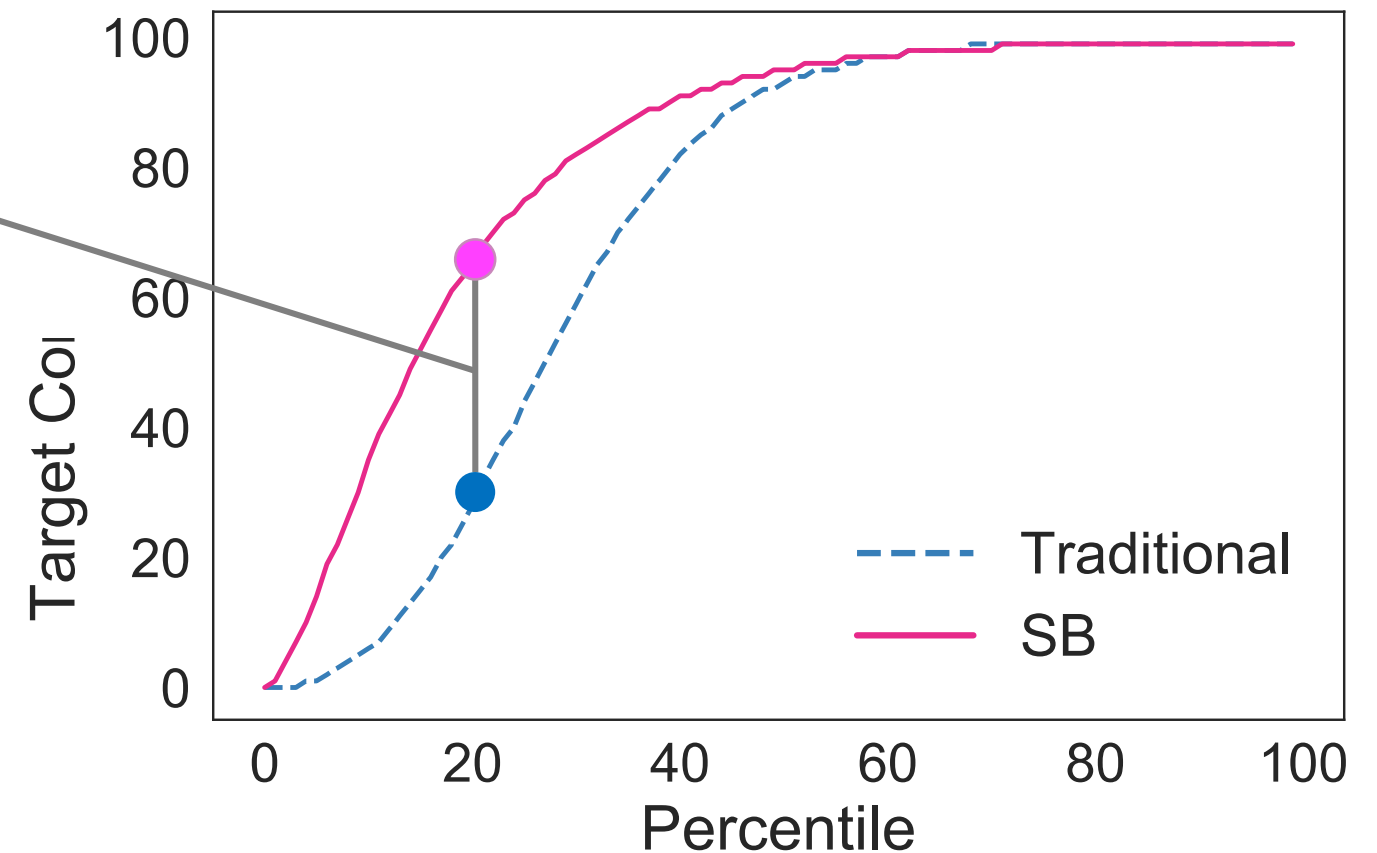Reduces time spent in the backwards pass by prioritizing high-loss examples

✓ **SelectiveBackprop outperforms static approaches**

Trains up to 3.5x faster compared to standard SGD

Trains 1.02-1.8X faster than state-of-the-art importance sampling approach

✓ **Stale-SB further accelerates training**

Trains on average 26% faster compared to SB

*www.github.com/angelajiang/SelectiveBackprop*

# Compared approaches

**Traditional**

**Classic SGD with no filtering**

**Katharopoulos18**

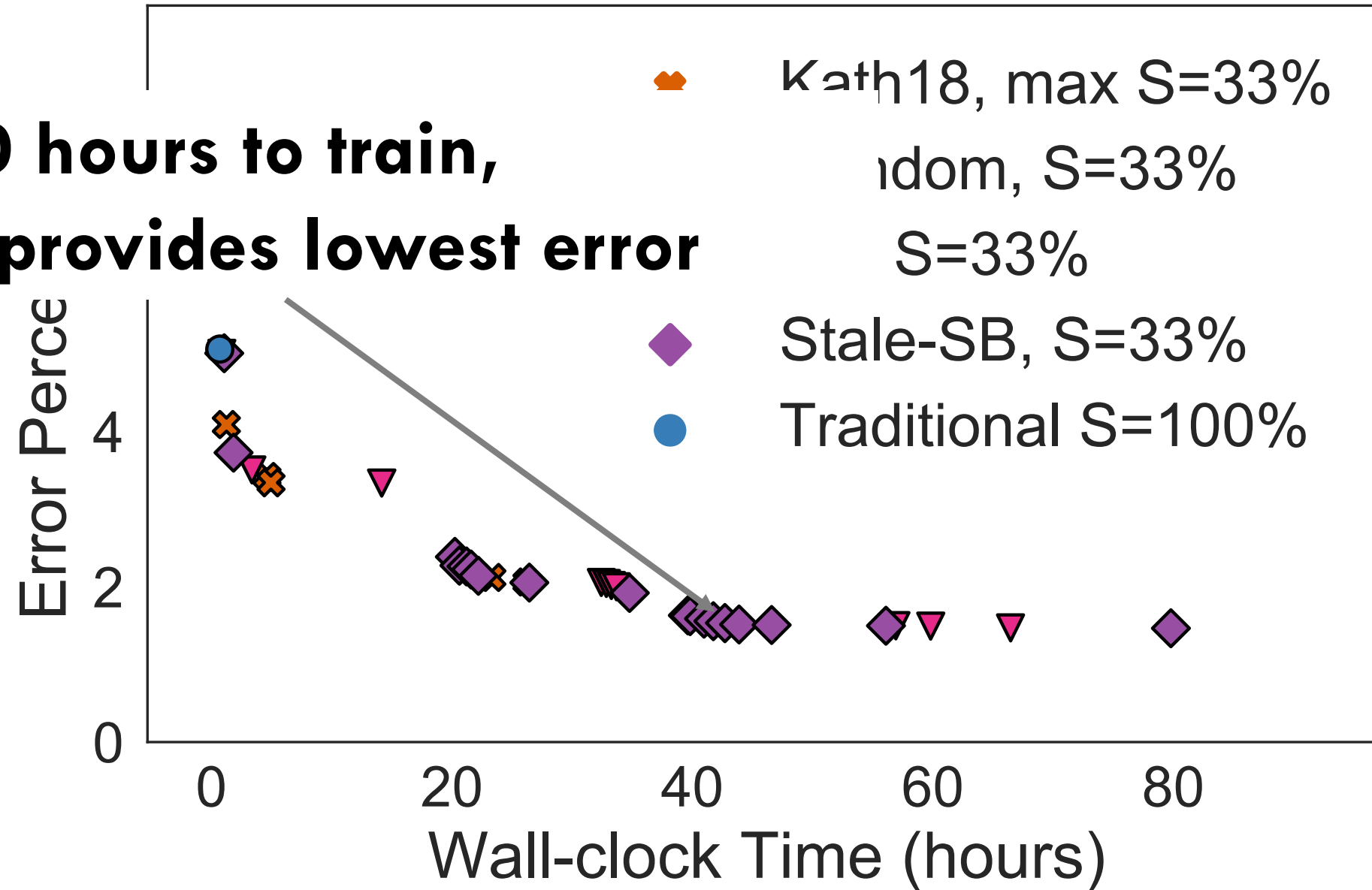**State of the art importance sampling approach**

**Random**

**Random sampling approach**
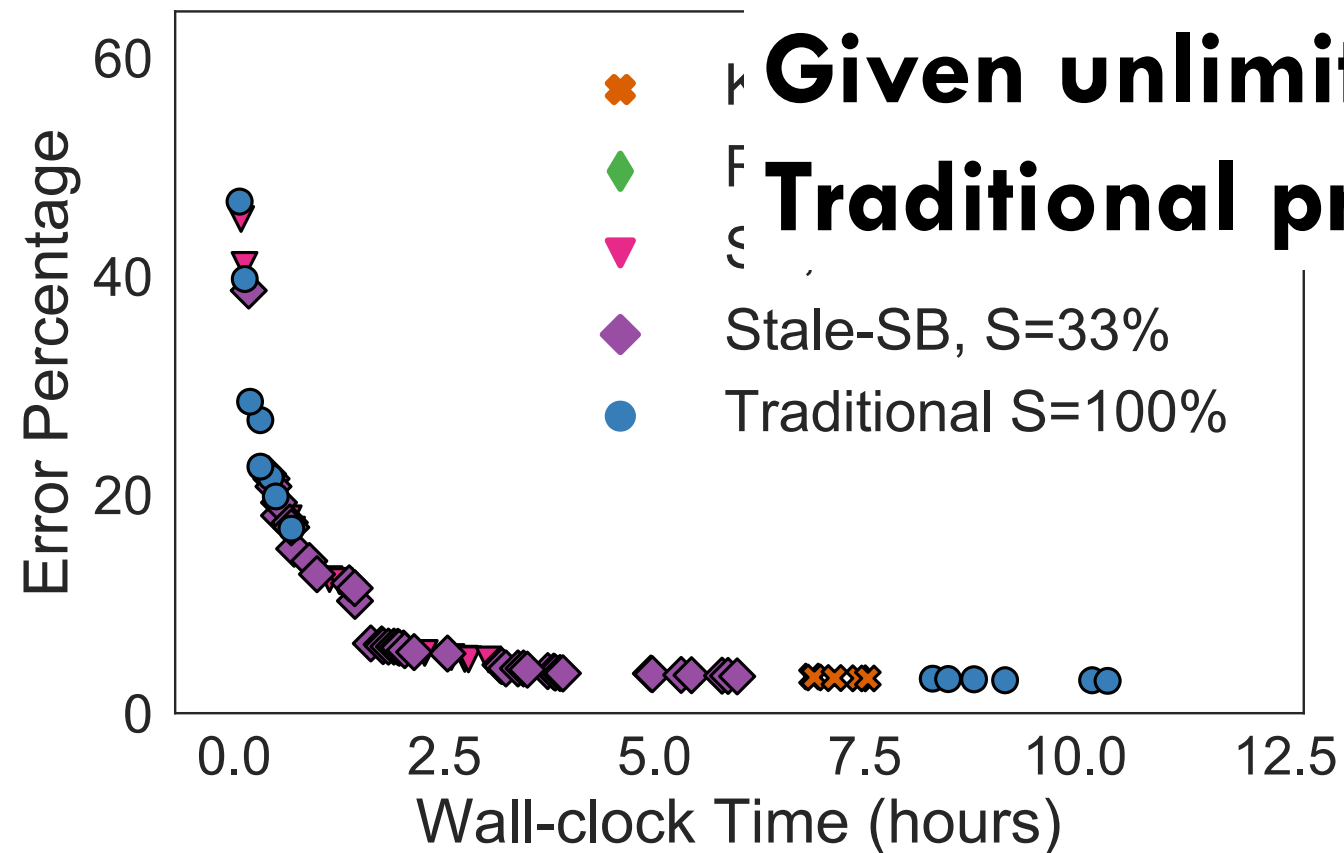
**Selective-Backprop (Us)**

# Most Pareto optimal points are SB or StaleSB



**CIFAR10**

**CIFAR100**

**Given unlimited time to train, Traditional provides lowest error**

# SB is robust to modest amounts of error

## 0.1% Randomized

SB, S=33%, Err=3.73%
Traditional S=100%, Err=3.41%

Test Error

Num Images Backpropped (millions)

## 10% Randomized

SB, S=33%, Err=6.72%
Traditional S=100%, Err=6.92%

Num Images Backpropped (millions)

## 20% Randomized

SB, S=33%, Err=9.33%
Traditional S=100%, Err=8.12%

Num Images Backpropped (millions)