# Thesis Progress report

## Ching-Yuan, Tsai

2018,1,17
Presented by Ching-Yuan, Tsai

Introduction
Prune
Gradients request scheduling
Timetable

Motivation
Objective

# Outline

Introduction
Prune
Gradients request scheduling
Timetable

Motivation
Objective

## Motivation

- In large module, distributed training may be slower than single-machine training.
- Parameter server uses a lot of time on synchronous stage.
- Network is bottleneck in our system.

Introduction
Prune
Gradients request scheduling
Timetable

Motivation
Objective

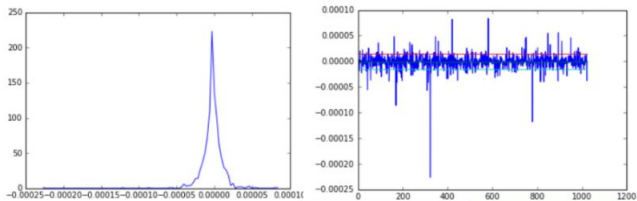## Objective

- Use software method to reduce network load.
    - Prune gradients by threshold.
        - static threshold
        - standard deviation threshold
        - dynamic threshold
    - Gradients request scheduling.

Introduction
Prune
Gradients request scheduling
Timetable

Observation
Prune
Standard deviation threshold
Static threshold
Dynamic threshold

# Outline

1. **Introduction**

2. **Prune**

3. **Gradients request scheduling**

4. **Timetable**

Introduction
Prune
Gradients request scheduling
Timetable

Observation
Prune
Standard deviation threshold
Static threshold
Dynamic threshold

## Observation

Introduction
Prune
Gradients request scheduling
Timetable

Observation
Prune
Standard deviation threshold
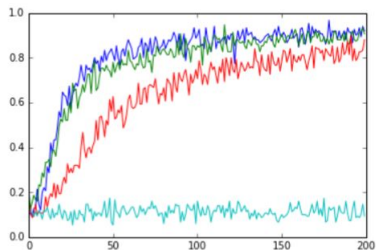Static threshold
Dynamic threshold

## Prune

- Prune gradients: Send gradients which are absolutely greater than threshold.
    - Static threshold: 10%, 1%, 0.1%.
    - Standard deviation threshold: 1, 2, 3.
    - Dynamic threshold: mean of the gradients which are greater than last threshold.

Introduction
Prune
Gradients request scheduling
Timetable

Observation
Prune
Standard deviation threshold
Static threshold
Dynamic threshold

## Standard deviation threshold

- Compute standard deviation on gpu.
- Faster 30 times than cpu.
- Use 25% computation time.

Introduction
Prune
Gradients request scheduling
Timetable

Observation
Prune
Standard deviation threshold
Static threshold
Dynamic threshold

# experiment 1

Introduction
Prune
Gradients request scheduling
Timetable

Observation
Prune
Standard deviation threshold
Static threshold
Dynamic threshold

## experiment 2

- To do: accuracy for each time.

Introduction
Prune
Gradients request scheduling
Timetable

Observation
Prune
Standard deviation threshold
Static threshold
Dynamic threshold

# Static threshold

- Compute on gpu.
- Selection algorithm.

Introduction
Prune
Gradients request scheduling
Timetable

Observation
Prune
Standard deviation threshold
Static threshold
Dynamic threshold

## experiment 1

- To do: accuracy for each interation.

Introduction
Prune
Gradients request scheduling
Timetable

Observation
Prune
Standard deviation threshold
Static threshold
Dynamic threshold

experiment 2

- To do: accuracy for each time.

Introduction
**Prune**
Gradients request scheduling
Timetable

Observation
Prune
Standard deviation threshold
Static threshold
**Dynamic threshold**

# Dynamic threshold

- Compute on gpu.
- Next threshold $=$ mean of the gradients which are greater than current threshold

Introduction
**Prune**
Gradients request scheduling
Timetable

Observation
Prune
Standard deviation threshold
Static threshold
**Dynamic threshold**

## experiment 1

- To do: accuracy for each interation.

Introduction
Prune
Gradients request scheduling
Timetable

Observation
Prune
Standard deviation threshold
Static threshold
Dynamic threshold

## experiment 2

- To do: accuracy for each time.

Introduction
Prune
Gradients request scheduling
Timetable

Introduction
Some examples
Problem description
Scheduling

# Outline

1. **Introduction**

2. **Prune**

3. **Gradients request scheduling**

4. **Timetable**

Introduction
Prune
Gradients request scheduling
Timetable

Introduction
Some examples
Problem description
Scheduling

# Introduction

- Client changes gradient request order to get some benefits.

Introduction
Prune
Gradients request scheduling
Timetable

Introduction
Some examples
Problem description
Scheduling

# Reduce suspend

Introduction
Prune
Gradients request scheduling
Timetable

Introduction
Some examples
Problem description
Scheduling

# Reduce staleness

Introduction
Prune
Gradients request scheduling
Timetable

Introduction
Some examples
Problem description
Scheduling

## Input

- $B$ = staleness bound.
- $L$ = number of layers.
- $I$ = number of iterations.
- $W$ = number of workers.
- $Cp$ = computation time : f1, b1, f2, b2, ...
- $Cm$ = communication time : t1, t2, ...

Introduction
Prune
Gradients request scheduling
Timetable

Introduction
Some examples
Problem description
Scheduling

## Intermediate data

- Server maintains a set of variables indicating the minimum iteration of each layer.
- $c_1$, $c_2$, ... , $c_L$

Introduction
Prune
Gradients request scheduling
Timetable

Introduction
Some examples
Problem description
Scheduling

## Objective

- Minimize training time and staleness.

Introduction
Prune
Gradients request scheduling
Timetable

Introduction
Some examples
Problem description
Scheduling

# Request pool

- Pull request.
- Push request.

Introduction
Prune
Gradients request scheduling
Timetable

Introduction
Some examples
Problem description
Scheduling

# Scheduling

- First in first out.
- First in last out.
    - Starvation.

Introduction
Prune
Gradients request scheduling
Timetable

Introduction
Some examples
Problem description
Scheduling

## Optimal solution

- Pull
  - Pull the stalest layer.
- Push
  - Push the layer which can update intermediate data.

# Outline

1. **Introduction**

2. **Prune**

3. **Gradients request scheduling**

4. **Timetable**

# Timetable

- -3/31 : Finish all experiment.
- -4/30 : Write paper.
- 6/15- : Oral.