

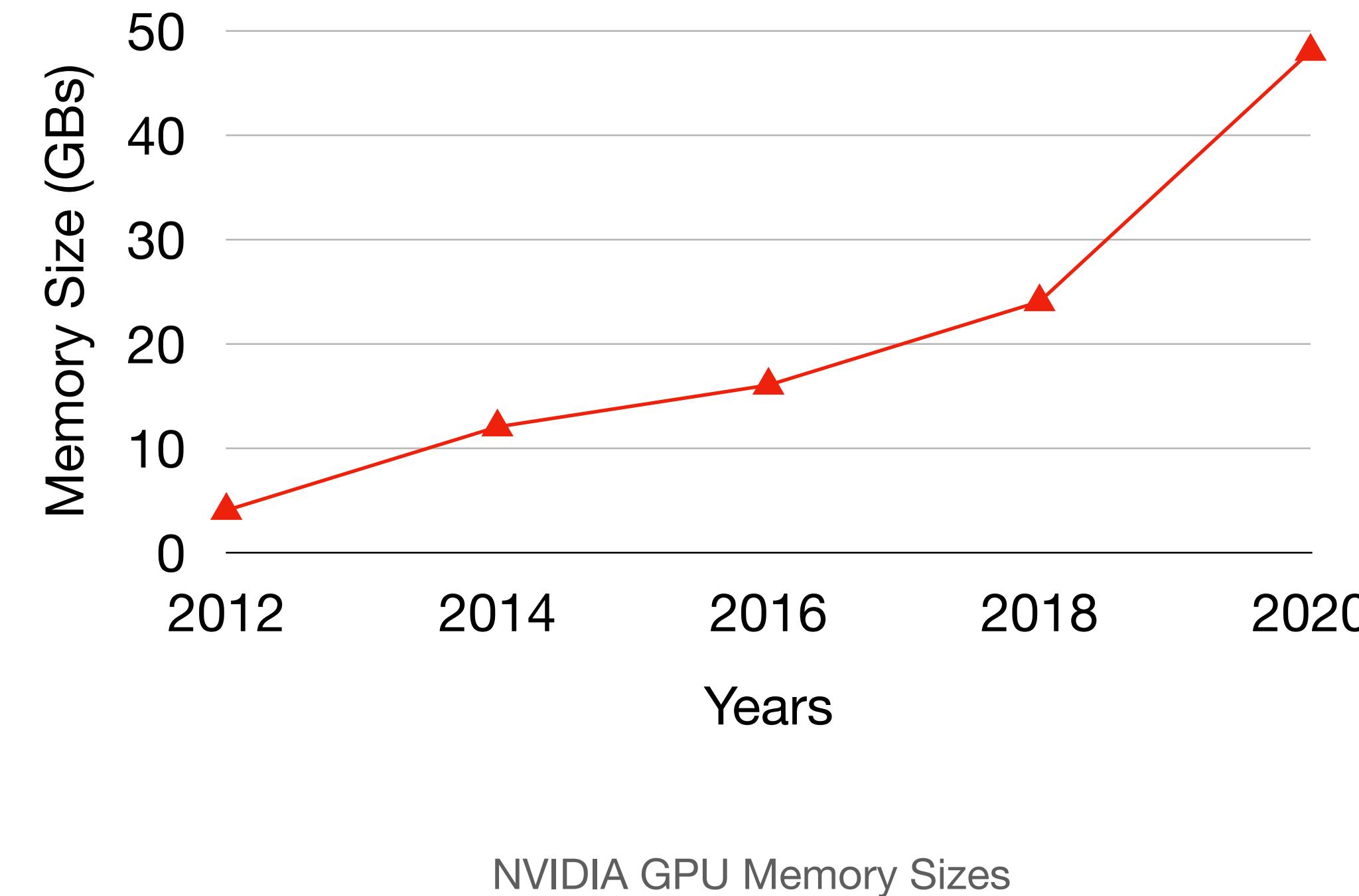
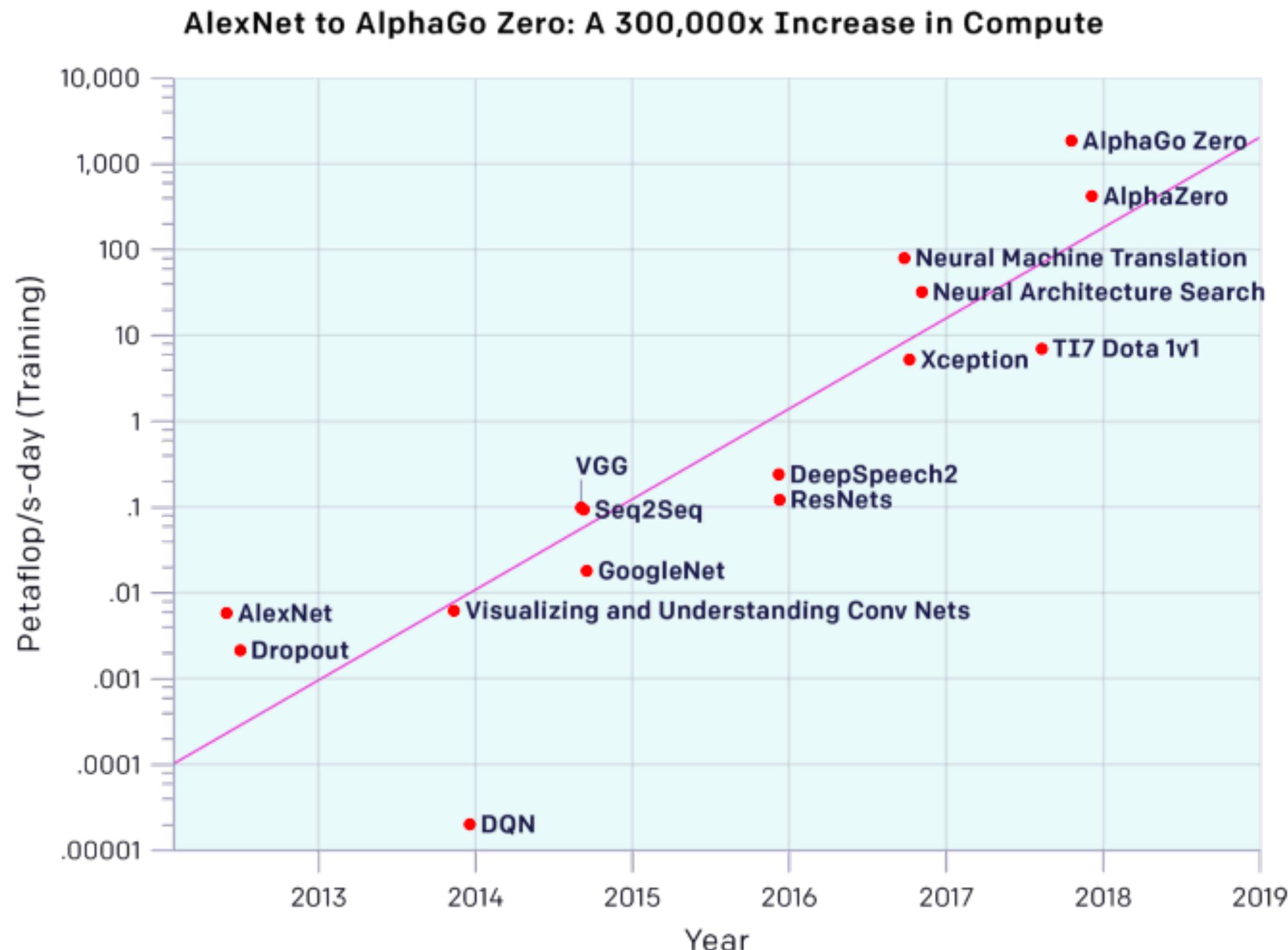
ML Systems Project

CS265 *Spring 2025*



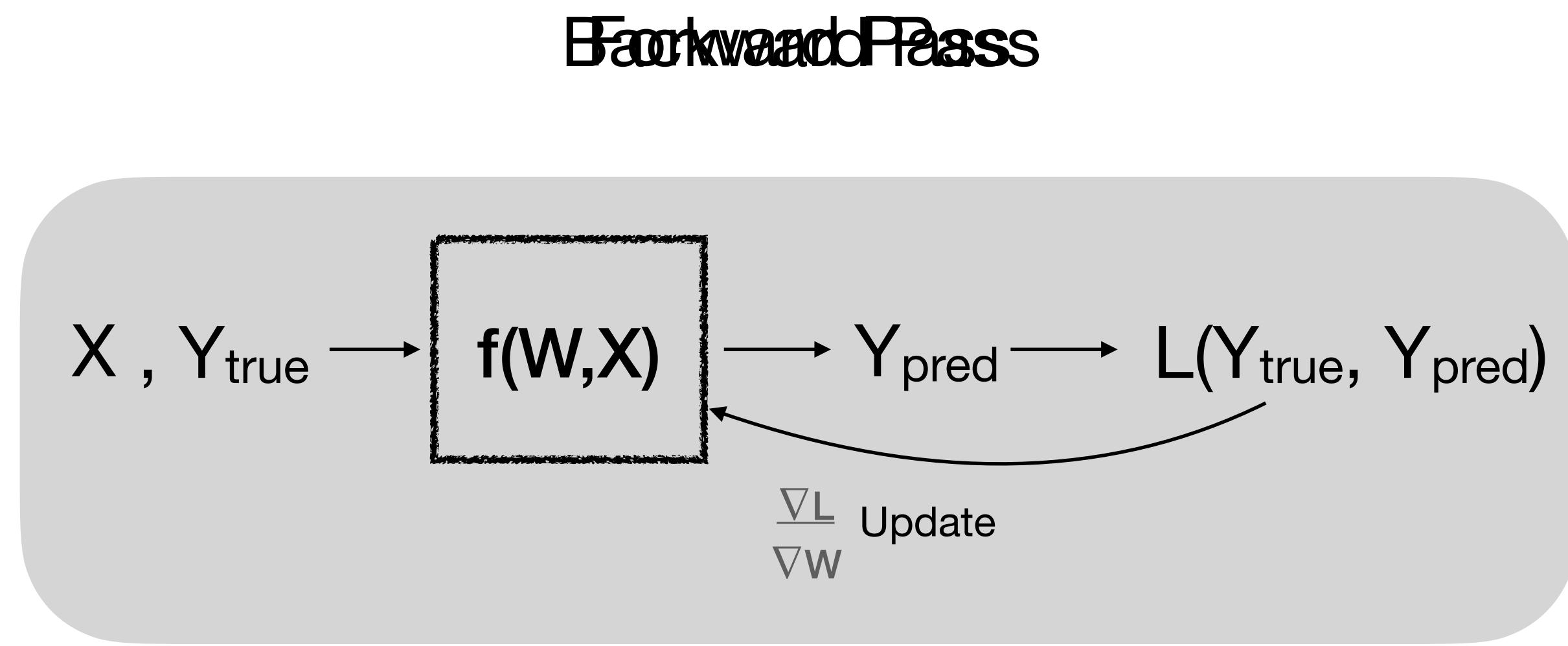
DASlab
@ Harvard SEAS

Memory is the Bottleneck in NN Training



Neural Network Training

Black-box function: f
Input: X
Model Parameters: W



Loss Function: L
Prediction: Y_{pred}
True Label: Y_{true}

Loss (L): Distance from True Value

$\frac{\nabla L}{\nabla W}$: Gradient of Loss (L) wrt weights (W)

Forward Pass

Input X

$$Z_1 = W_1 X$$

$$Z_2 = \sigma(Z_1)$$

$$Z_3 = W_3 Z_2$$

$$Z_4 = \sigma(Z_3)$$

$$L = \frac{1}{2}(Z_4 - Y)^2$$

True Label Y

$$W_i = W_i - \alpha \nabla W_i$$

Update Rule

Function Composition

$$\begin{aligned} y &= f(x) & \frac{\delta z}{\delta x} &= \frac{\delta g(y)}{\delta y} \frac{\delta y}{\delta x} \\ z &= g(y) & & \\ z &= g(f(x)) & \frac{\delta z}{\delta x} &= \frac{\delta g(f(x))}{\delta f(x)} \frac{\delta f(x)}{\delta x} \end{aligned}$$

Chain Rule

Backward Pass

$$\nabla W_1 = \nabla Z_1 X$$

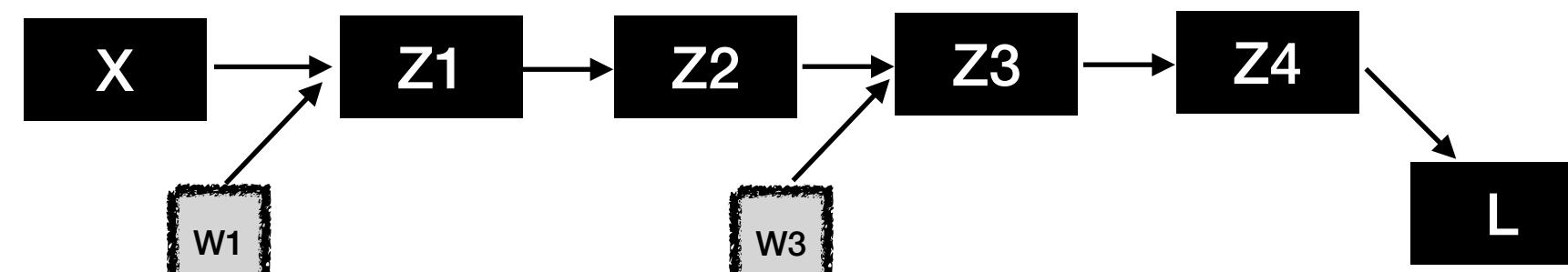
$$\nabla Z_1 = \nabla Z_2 Z_2 (1 - Z_2)$$

$$\nabla Z_2 = \nabla Z_3 W_3$$

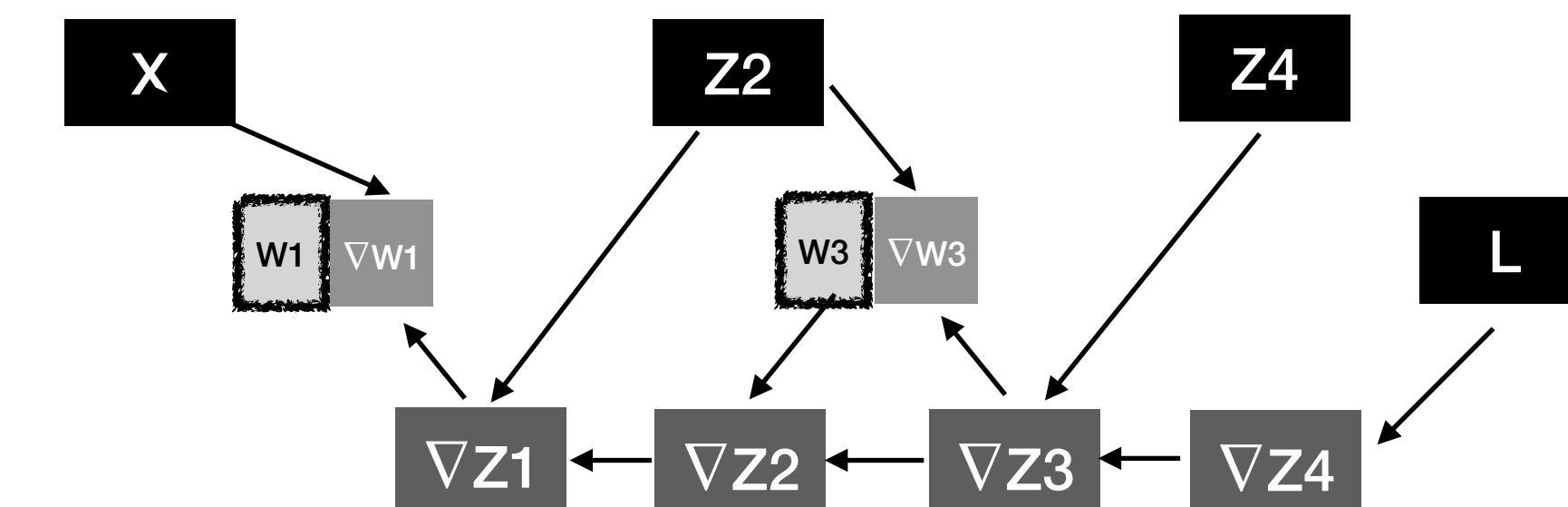
$$\nabla W_3 = \nabla Z_3 Z_2$$

$$\nabla Z_3 = \nabla Z_4 Z_4 (1 - Z_4)$$

$$\nabla Z_4 = (Z_4 - Y)$$



Forward Computational Graph



Backward Computational Graph

Forward Pass

Input X

$$Z_1 = W_1 X$$

$$Z_2 = \sigma(Z_1)$$

$$Z_3 = W_3 Z_2$$

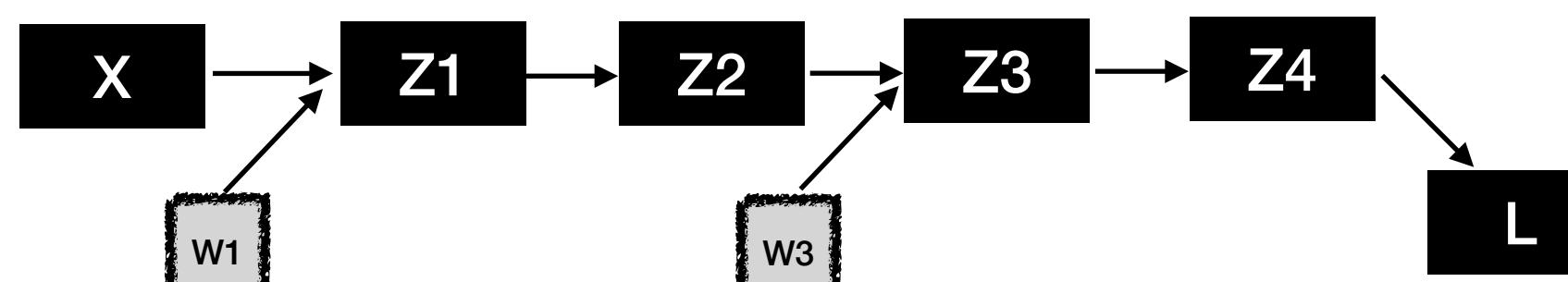
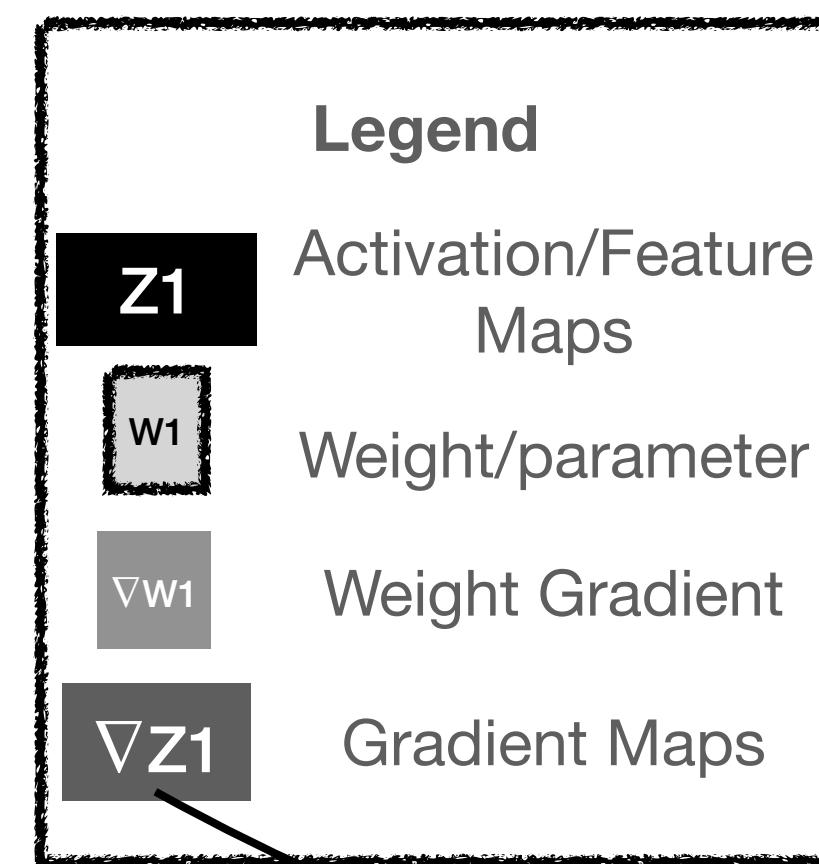
$$Z_4 = \sigma(Z_3)$$

$$L = \frac{1}{2}(Z_4 - Y)^2$$

True Label Y

$$W_i = W_i - \alpha \nabla W_i$$

Update Rule



Forward Computational Graph

Backward Pass

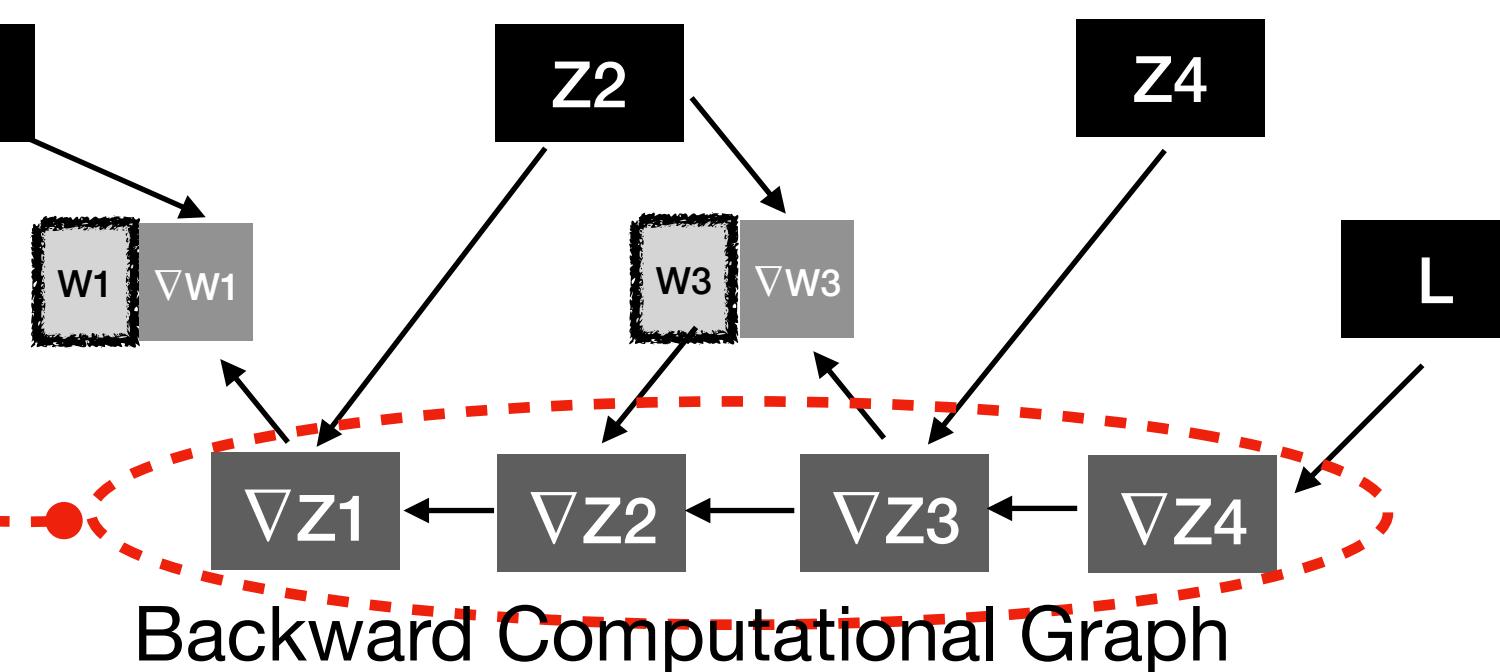
$$\nabla W_1 = \nabla Z_1 X$$

$$\nabla Z_1 = \nabla Z_2 Z_2 (1 - Z_2)$$

$$\nabla Z_2 = \nabla Z_3 W_3$$

$$\nabla Z_3 = \nabla Z_4 Z_4 (1 - Z_4)$$

$$\nabla Z_4 = (Z_4 - Y)$$



Backward Computational Graph

Forward Pass

Input X

$$Z_1 = W_1 X$$

$$Z_2 = \sigma(Z_1)$$

$$Z_3 = W_3 Z_2$$

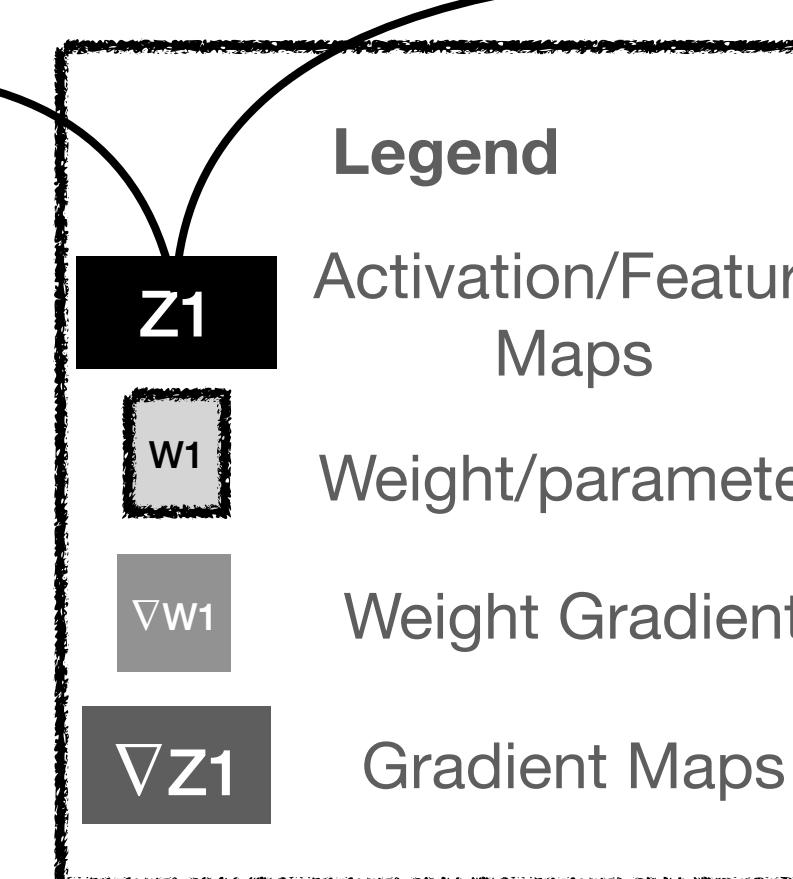
$$Z_4 = \sigma(Z_3)$$

$$L = \frac{1}{2} (Z_4 - Y)^2$$

True Label Y

$$W_i = W_i - \alpha \nabla W_i$$

Update Rule



Backward Pass

$$\nabla W_1 = \nabla Z_1 X$$

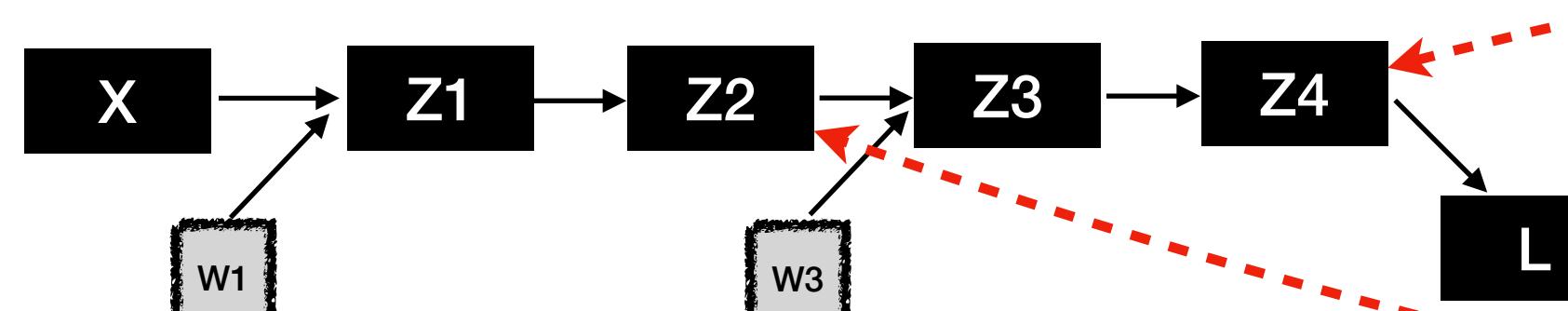
$$\nabla Z_1 = \nabla Z_2 Z_2 (1 - Z_2)$$

$$\nabla Z_2 = \nabla Z_3 W_3$$

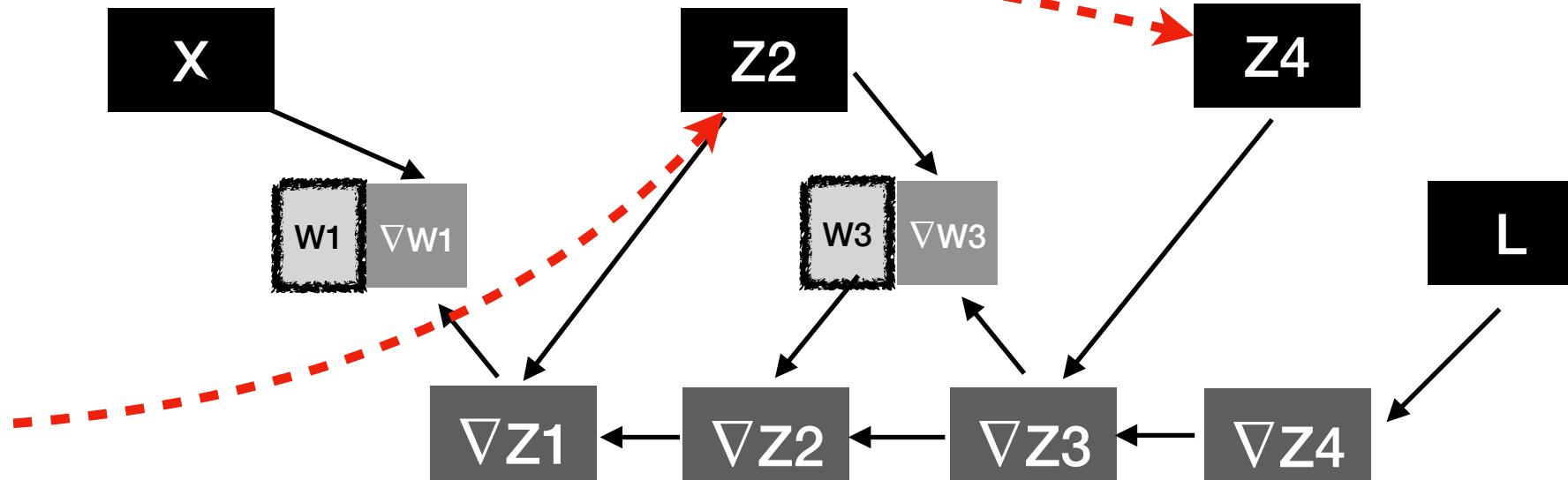
$$\nabla W_3 = \nabla Z_3 Z_2$$

$$\nabla Z_3 = \nabla Z_4 Z_4 (1 - Z_4)$$

$$\nabla Z_4 = (Z_4 - Y)$$



Forward Computational Graph



Backward Computational Graph



Forward Pass

Input X

$$Z_1 = W_1 X$$

$$Z_2 = \sigma(Z_1)$$

$$Z_3 = W_3 Z_2$$

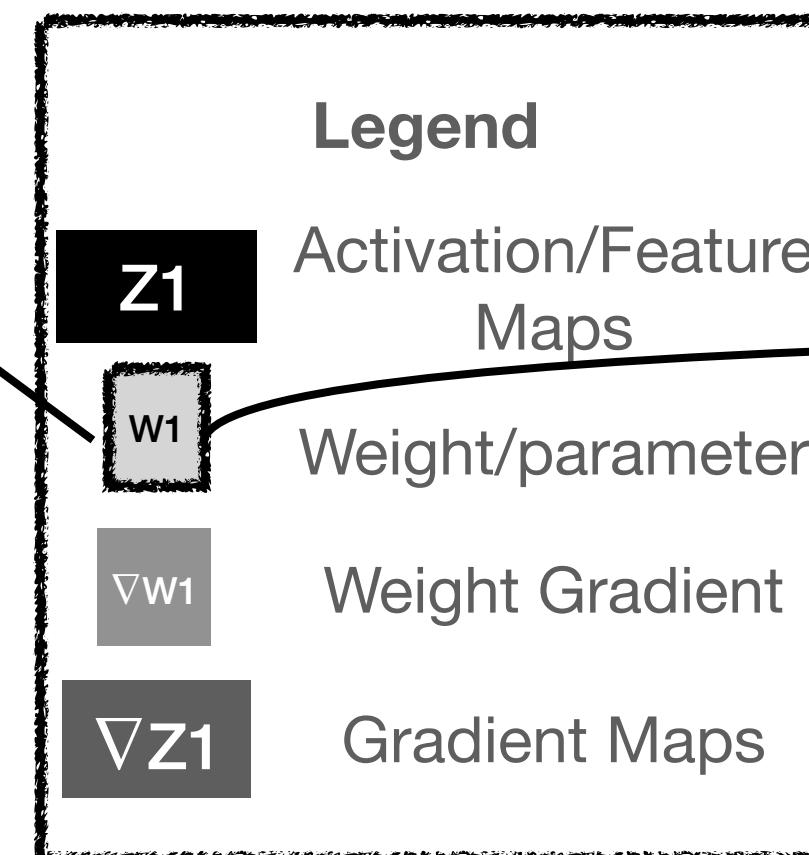
$$Z_4 = \sigma(Z_3)$$

$$L = \frac{1}{2}(Z_4 - Y)^2$$

True Label Y

$$W_i = W_i - \alpha \nabla W_i$$

Update Rule



Backward Pass

$$\nabla W_1 = \nabla Z_1 X$$

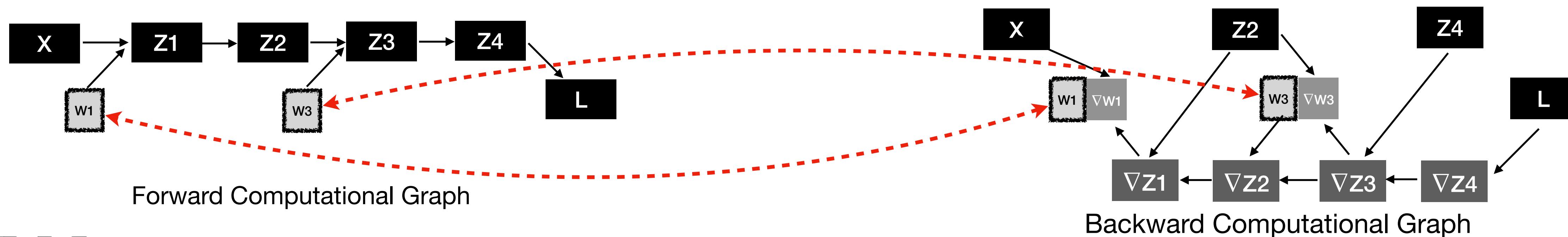
$$\nabla Z_1 = \nabla Z_2 Z_2 (1 - Z_2)$$

$$\nabla Z_2 = \nabla Z_3 W_3$$

$$\nabla W_3 = \nabla Z_3 Z_2$$

$$\nabla Z_3 = \nabla Z_4 Z_4 (1 - Z_4)$$

$$\nabla Z_4 = (Z_4 - Y)$$



Forward Pass

Input X

$$Z_1 = W_1 X$$

$$Z_2 = \sigma(Z_1)$$

$$Z_3 = W_3 Z_2$$

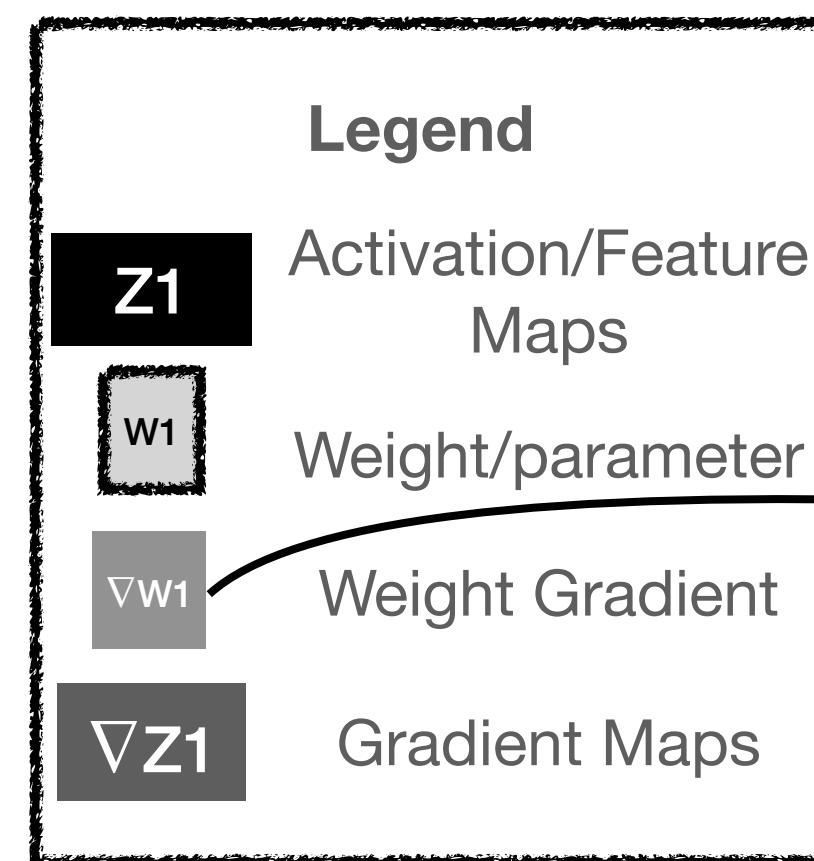
$$Z_4 = \sigma(Z_3)$$

$$L = \frac{1}{2}(Z_4 - Y)^2$$

True Label Y

$$W_i = W_i - \alpha \nabla W_i$$

Update Rule



Backward Pass

$$\nabla W_1 = \nabla Z_1 X$$

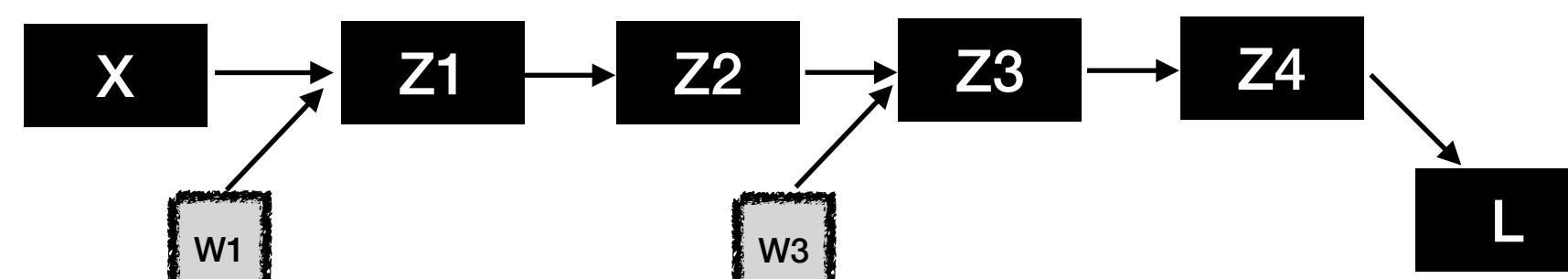
$$\nabla Z_1 = \nabla Z_2 Z_2 (1 - Z_2)$$

$$\nabla Z_2 = \nabla Z_3 W_3$$

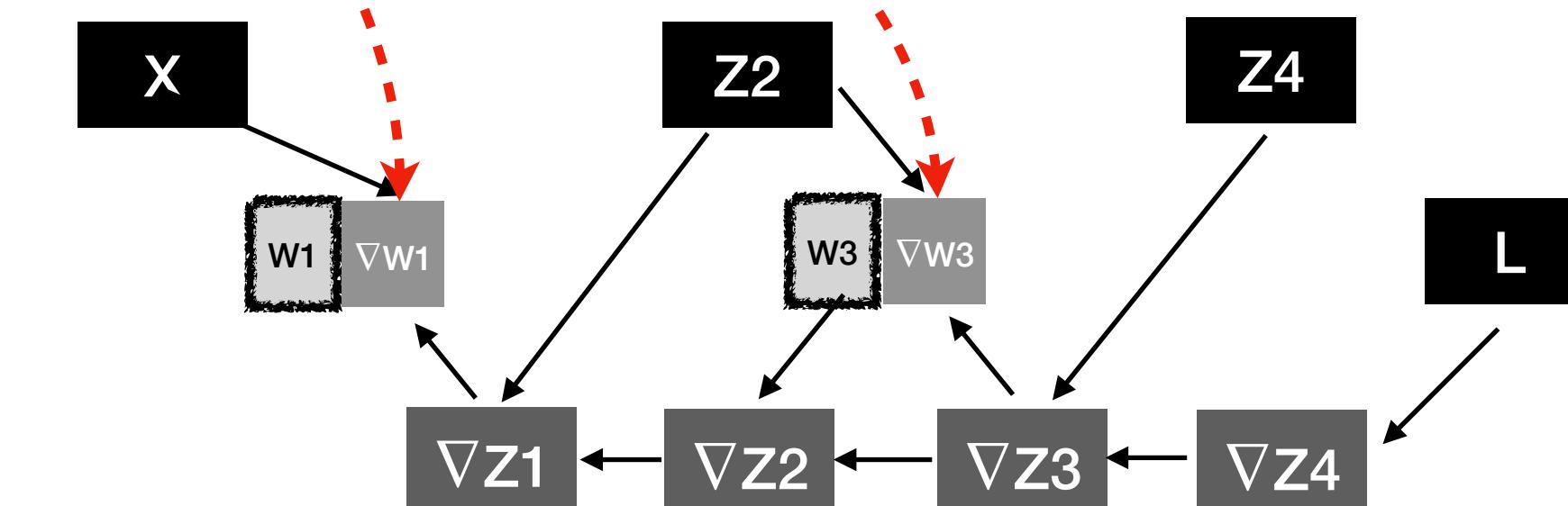
$$\nabla W_3 = \nabla Z_3 Z_2$$

$$\nabla Z_3 = \nabla Z_4 Z_4 (1 - Z_4)$$

$$\nabla Z_4 = (Z_4 - Y)$$



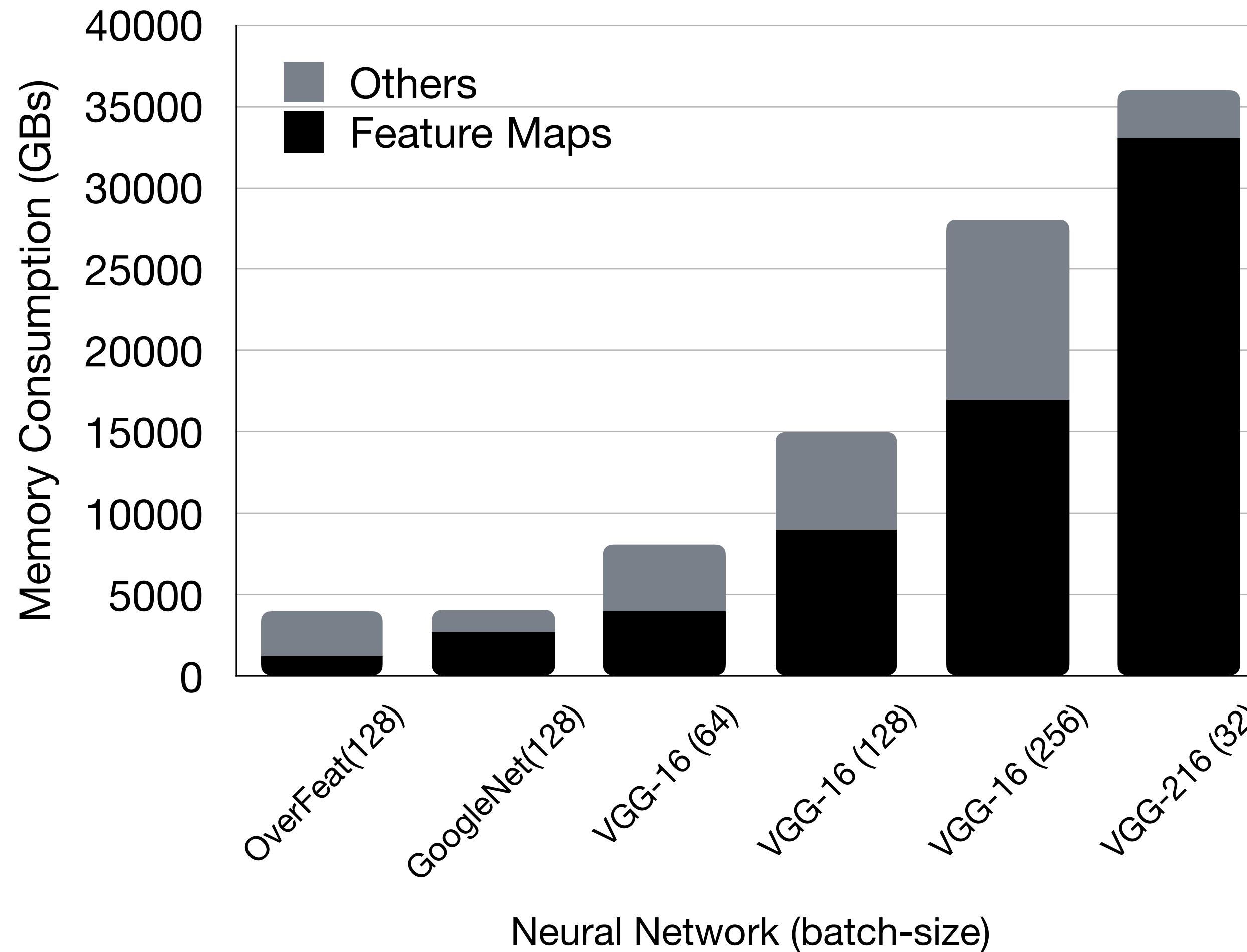
Forward Computational Graph



Backward Computational Graph

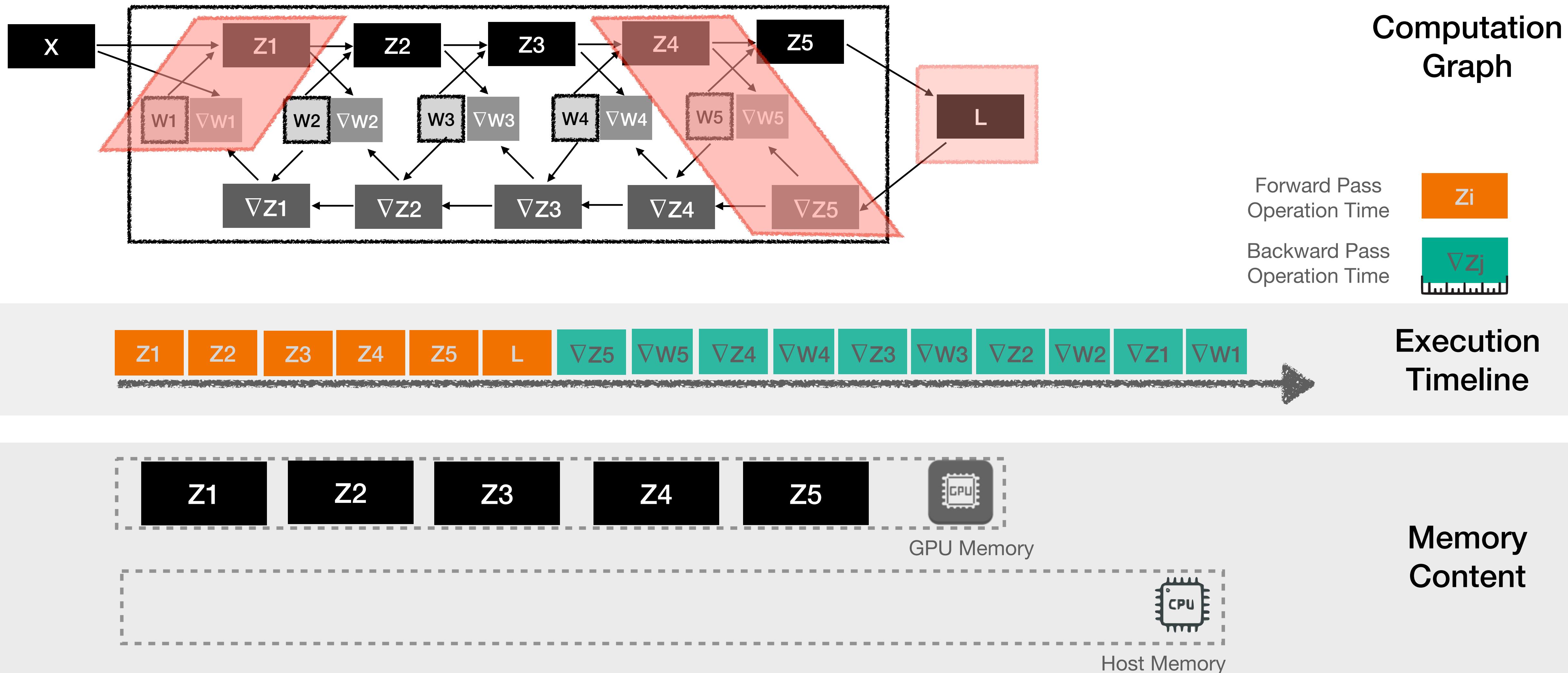


Reduce Memory Footprint

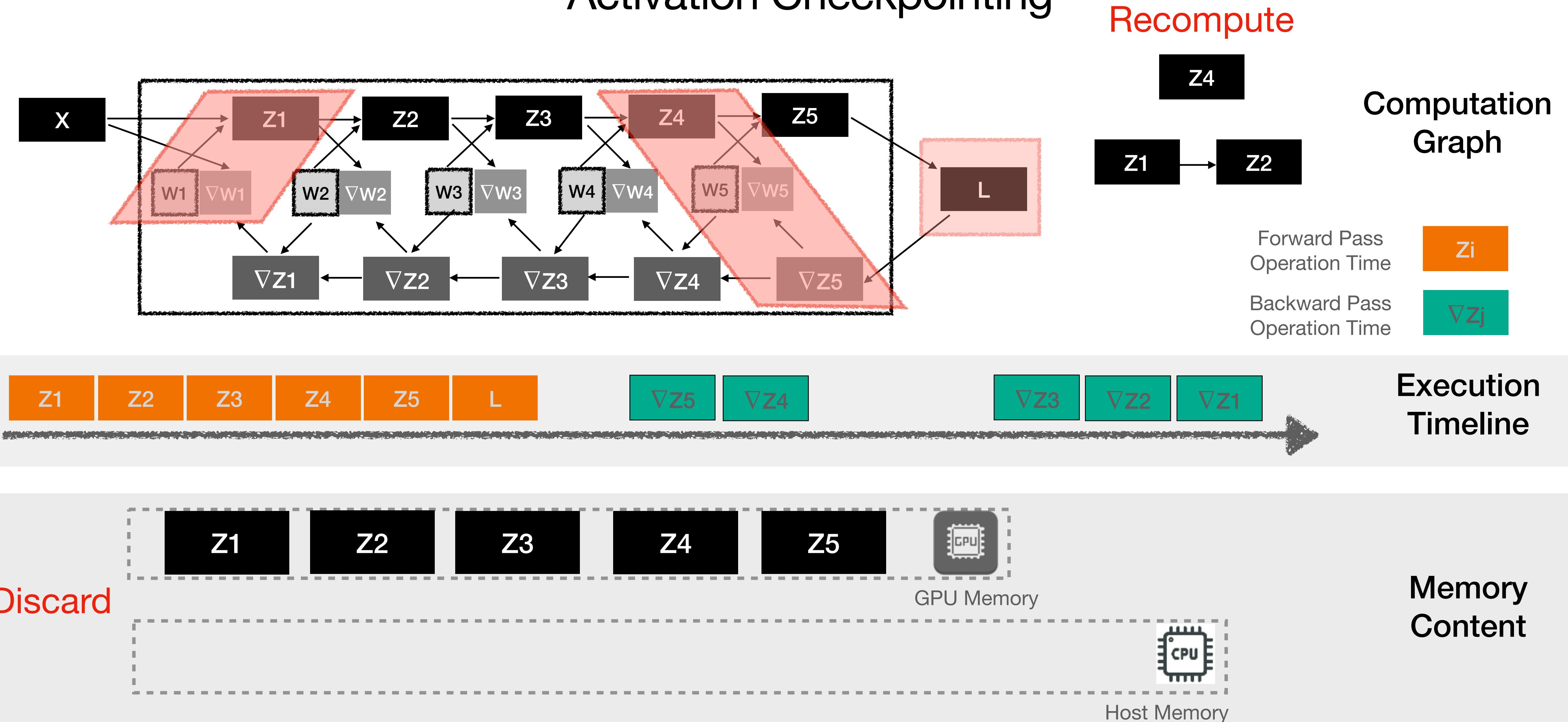


Adapted From: Minsoo Rhu et al. vDNN, MICRO 2016

Execution Framework

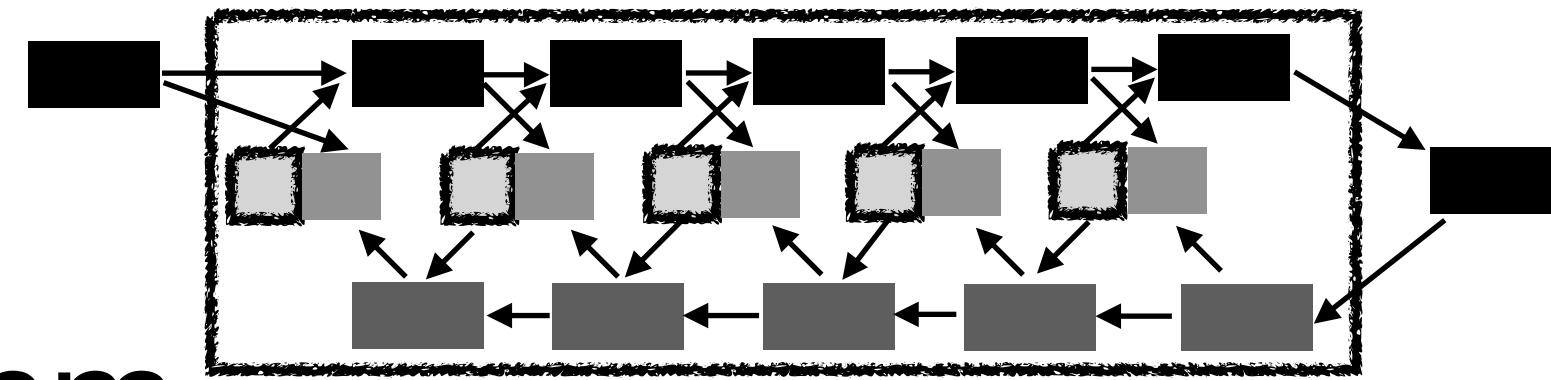


Activation Checkpointing



The ML Systems Project

1. construct and profile the computational graph on PyTorch



2. design the activation checkpointing (AC) algorithm

3. rewrite the PyTorch computational graph to deploy activation
checkpointing

[http://daslab.seas.harvard.edu/classes/cs265/files/CS 265 Systems Project Description.pdf](http://daslab.seas.harvard.edu/classes/cs265/files/CS%20265%20Systems%20Project%20Description.pdf)

How to Start

1. play with the starter code (instructions in README.md)

<https://github.com/qtwang/CS265-mlsys-project>

2. check the reference paper

Sanket Purandare, Abdul Wasay, Animesh Jain, Stratos Idreos: μ -TWO: 3 \times Faster Multi-Model Training with Orchestration and Memory Optimization. MLSys 2023.

3. come to the labs with questions!

Midway Check-in

1. complete Phase 1: construct and profile the computational graph
2. a document with experimental analysis consisting of the following two deliverables:
 - computation and memory profiling statistics and static analysis
 - peak memory consumption vs. mini-batch size bar graph without AC

Final Deliverables

1. code deliverable + code review + demo = 50%
2. a document with experimental analysis = 50%
 - computation and memory profiling statistics and static analysis
 - peak memory consumption vs mini-batch size bar graph (w and w/o AC)
 - iteration latency vs mini-batch size performance graph (w and w/o AC)

Enjoy!