

ENGR-E 533 “Deep Learning Systems”

Lecture 06: Adult Optimization

Minje Kim

Department of Intelligent Systems Engineering

Email: minje@indiana.edu

Website: <http://minjekim.com>

Research Group: <http://saige.sice.indiana.edu>

Meeting Request: <http://doodle.com/minje>



INDIANA UNIVERSITY

**SCHOOL OF INFORMATICS,
COMPUTING, AND ENGINEERING**

Initialization Methods, Activation Functions, and Batch Normalization

Yes, they are related

Initialization

-Why random?

- I remember a fellow student who didn't randomly initialize in his machine learning class a while ago
- What he did was to try out a bunch of different random numbers and pick up the best one
 - Best in the sense of test accuracy
- He was still using some MATLAB functions that generate random numbers
 - Why is not the random initialization?
 - What would be the correct random initialization?
- Random initialization means you generate random numbers to initialize your parameters
 - Therefore, there could be some good choices and bad choices depending on your luck
 - You never know until you see the results
- Eventually there's a technique called ensemble, but for now let's just bear with this uncertainty
- Remaining questions
 - What would be the right p.d.f. to sample from
 - How do we define the parameters of the p.d.f.?



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

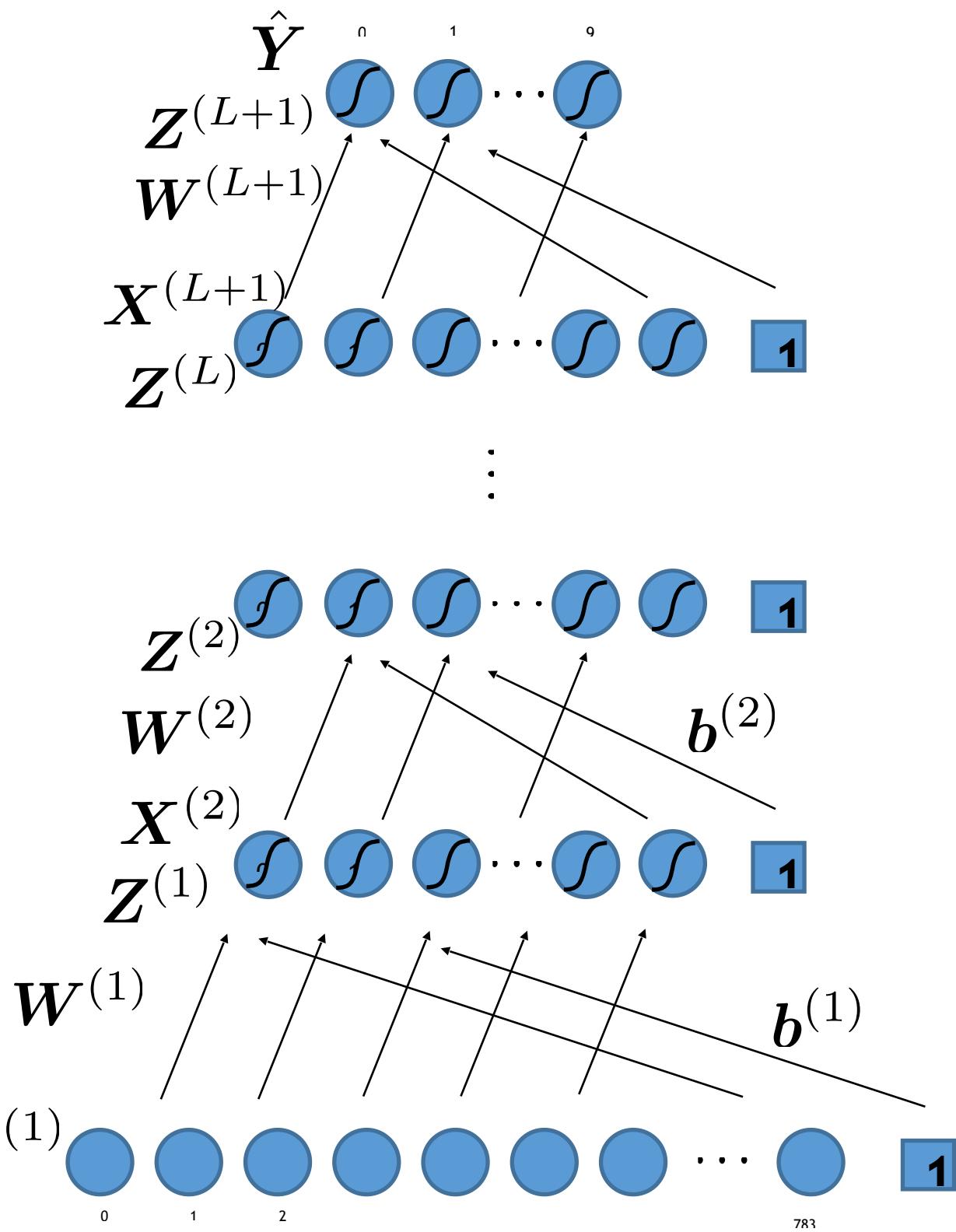
Initialization

-Choosing the best p.d.f.

- Optimal p.d.f. to sample your initial parameters from
 - Suppose the set of weights that give you the global optimum of the objective function
 - e.g. The weights that predict the test samples with the maximum possible accuracy
 - If we find a p.d.f. that's very similar to the sample distribution of these optimal weights
 - I have no theory behind this, but that would be the best one to start from
- Normally we should use some other distributions
 - Maybe some mismatching ones
 - What would be your choice?
 - Normal distribution $\mathcal{N}(\mu, \sigma)$ Many people also like to use uniform distributions, too
- What would be the mean and variance?
 - It might be safer to use zero mean
 - How about the variance?
 - Let's take a look at them
- Hold on.. You need to know which part of the network to look at:

$$\frac{\partial \mathcal{E}}{\partial \mathbf{W}^{(1)}} = \frac{\partial \mathcal{E}}{\partial \hat{\mathbf{Y}}} \frac{\partial \hat{\mathbf{Y}}}{\partial \mathbf{Z}^{(l+1)}} \prod_{l=1}^L \left[\frac{\partial \mathbf{Z}^{(l+1)}}{\partial \mathbf{X}^{(l+1)}} \frac{\partial \mathbf{X}^{(l+1)}}{\partial \mathbf{Z}^{(l)}} \right] \frac{\partial \mathbf{Z}^{(1)}}{\partial \mathbf{W}^{(1)}} = \left(\left(\mathbf{W}^{(2)\top} \left(\dots \left(\mathbf{W}^{(L+1)\top} (\hat{\mathbf{Y}} - \mathbf{Y}) \right) \odot \sigma'(\checkmark \mathbf{Z}^{(L)}) \dots \right) \right) \odot \sigma'(\checkmark \mathbf{Z}^{(1)}) \right) \mathbf{X}^{(1)\top}$$

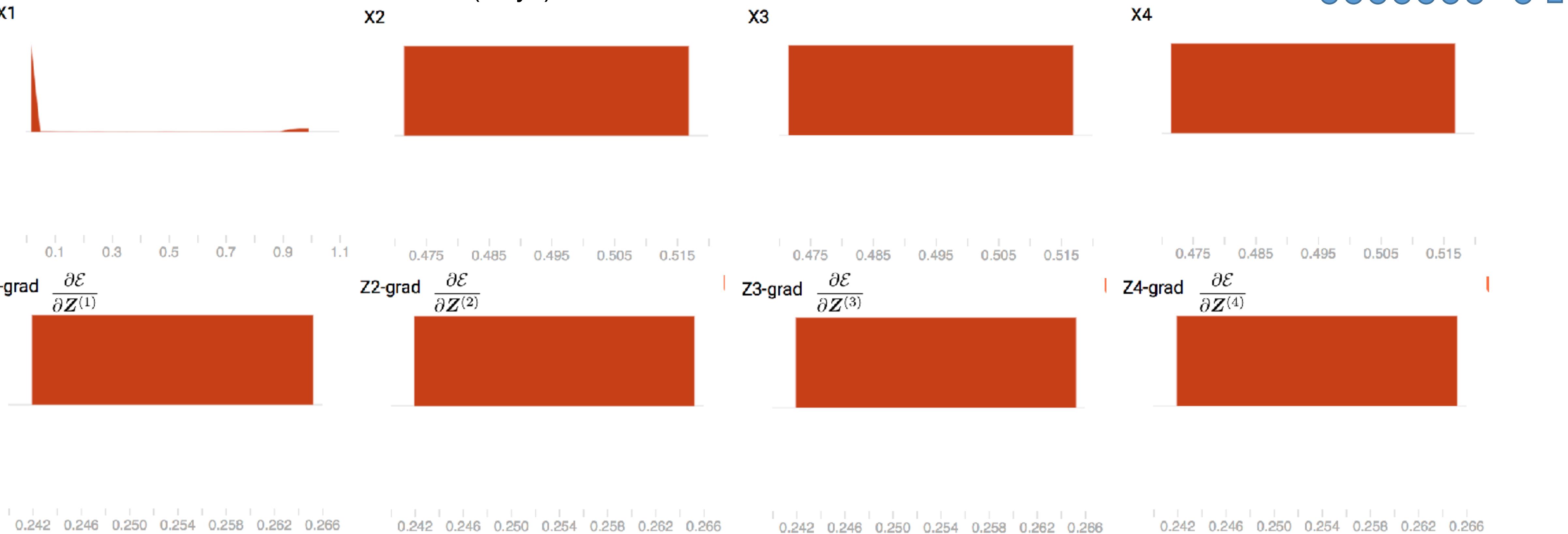
L speed bumps



Initialization

-Choosing the best p.d.f.

- Too small variance $\sigma = 0.0001$
 - All the hidden unit activations (features) are centered around a value
 - Not healthy (features should be distinctive)
 - Even worse with the tanh activations (why?)



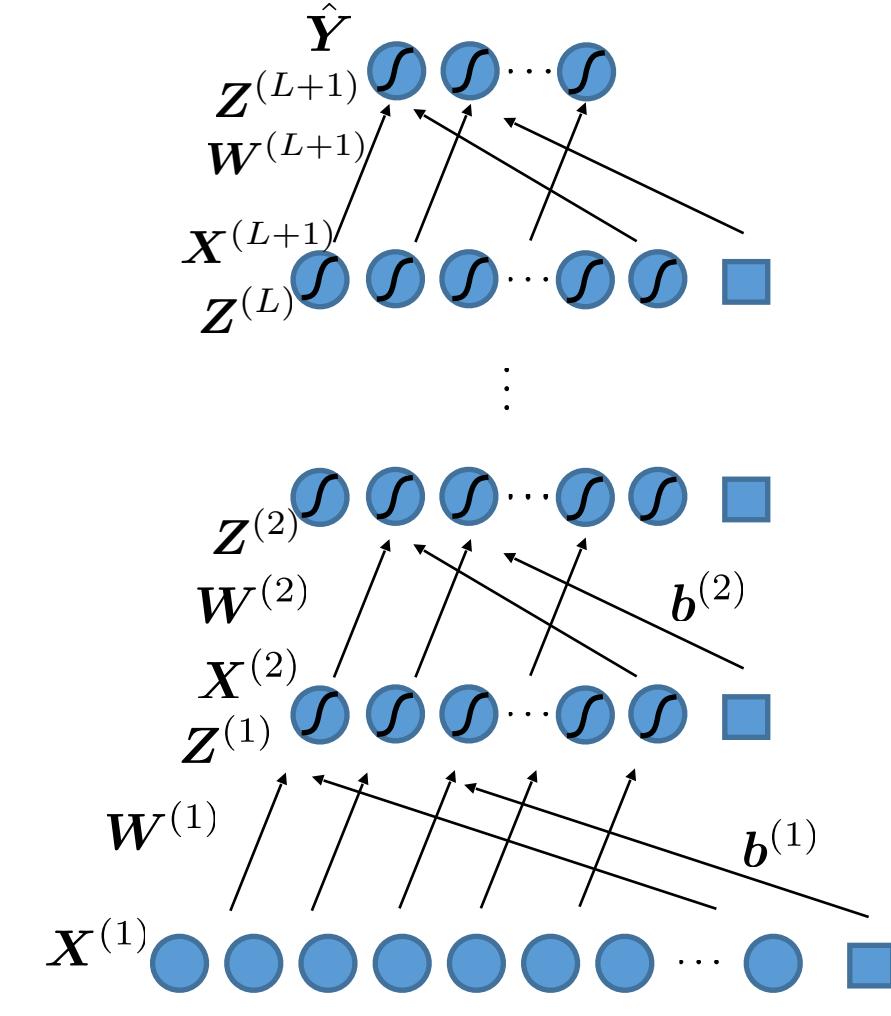
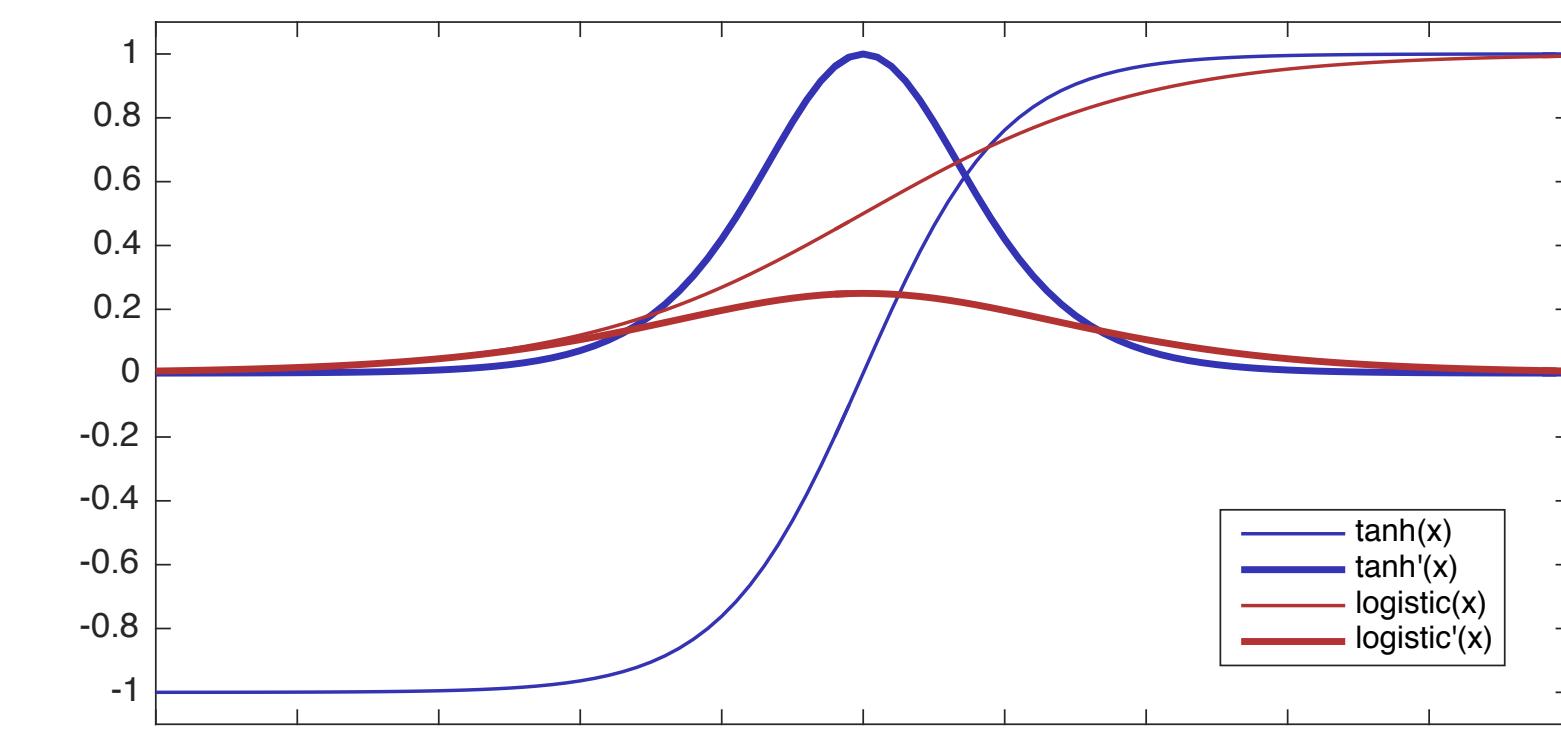
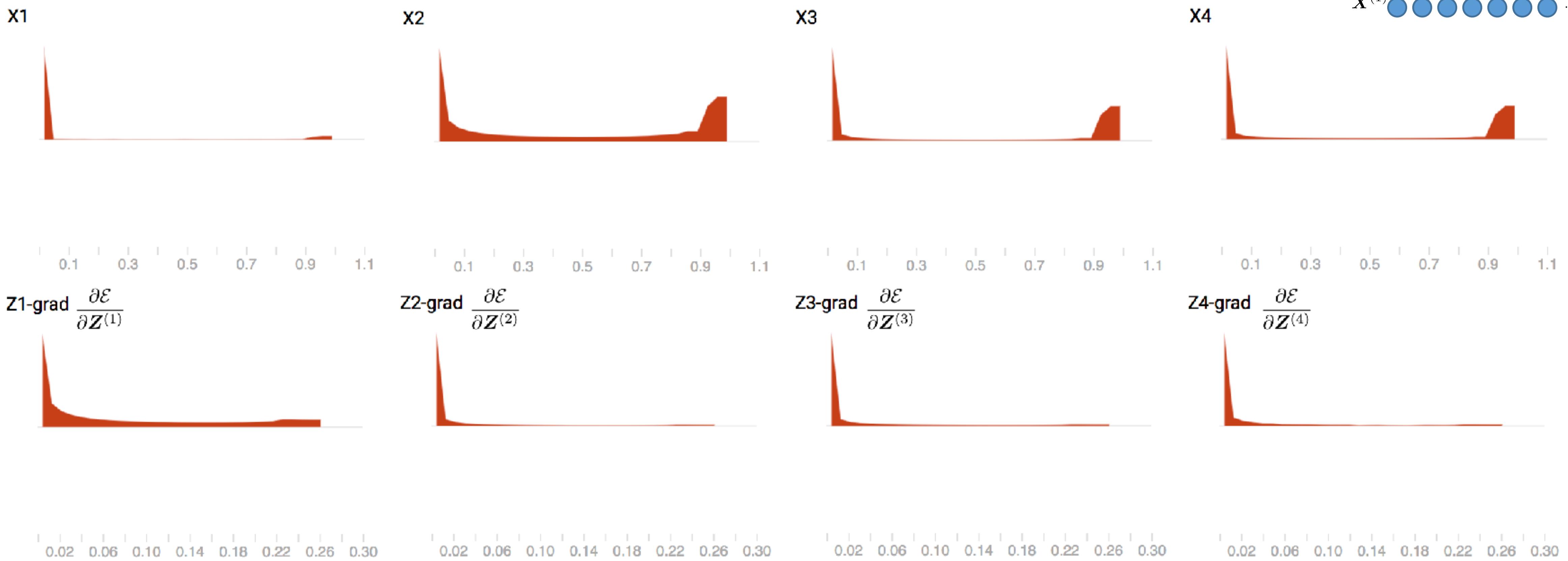
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Initialization

-Choosing the best p.d.f.

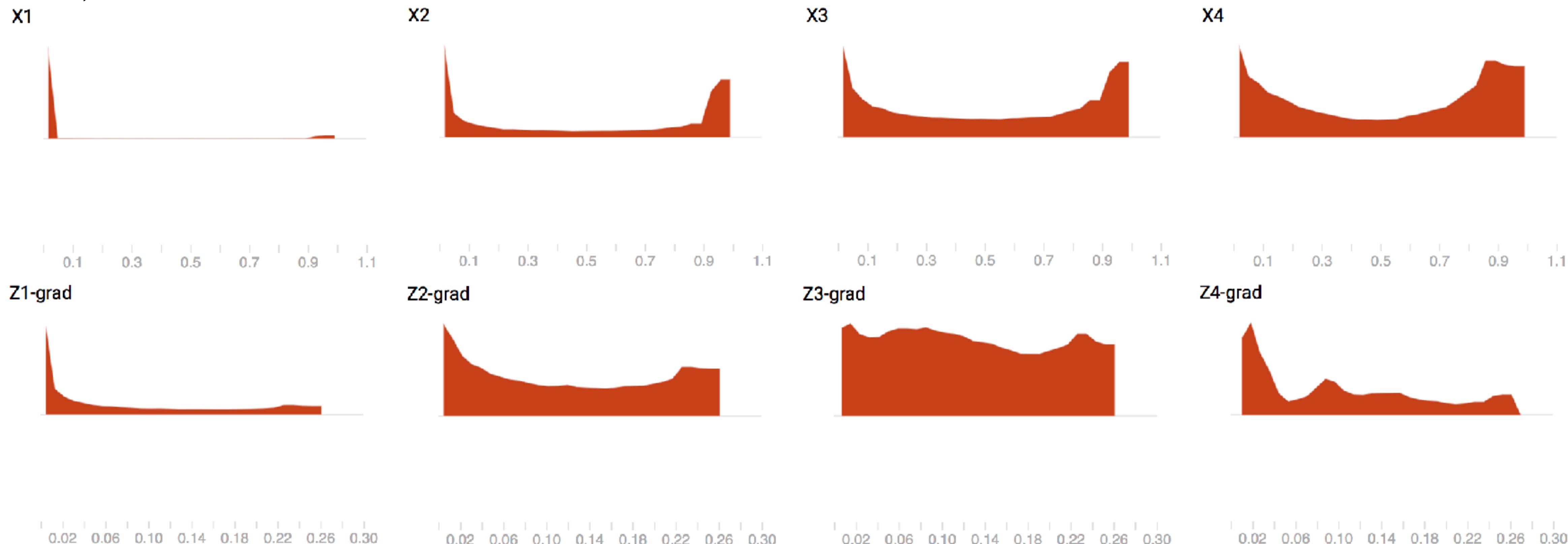
- Too large variance $\sigma = 0.5$
 - Hidden unit outputs are too extreme (why?)
 - Gradient vanishes



Initialization

-Dimension matters

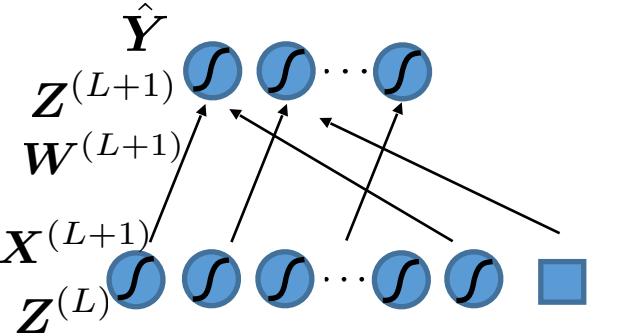
- Still too large variance $\sigma = 0.5$, but now are with only 50 hidden units per layer (previously it was 1,000)
 - Why does this happen?
 - $\text{var}(XY) = \text{var}(X)\text{var}(Y)$ (only if they are independent and zero mean)
 - So, what?



INDIANA UNIVERSITY

Frishman, Fred. "On the arithmetic means and variances of products and ratios of random variables." A Modern Course on Statistical Distributions in Scientific Work. Springer, Dordrecht, 1975. 401-406.

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING



Initialization

-Xavier Initialization

- Variance of the linear combination $\mathbf{Z}_{i,t} = \sum_{j=1}^{N_j} \mathbf{W}_{i,j} \mathbf{X}_{j,t}$

$$\text{var}(\mathbf{Z}_{i,t}) = N_j \text{var}(\mathbf{W}_{i,j}) \text{var}(\mathbf{X}_{j,t})$$

- To match the variance?

$$\text{var}(\mathbf{W}_{i,j}) = \frac{1}{N_j}$$

- We forgot about the other linear combination

$$\frac{\partial \mathcal{E}}{\partial \mathbf{W}^{(1)}} = \frac{\partial \mathcal{E}}{\partial \hat{\mathbf{Y}}} \frac{\partial \hat{\mathbf{Y}}}{\partial \mathbf{Z}^{(l+1)}} \prod_{l=1}^L \left[\frac{\partial \mathbf{Z}^{(l+1)}}{\partial \mathbf{X}^{(l+1)}} \frac{\partial \mathbf{X}^{(l+1)}}{\partial \mathbf{Z}^{(l)}} \right] \frac{\partial \mathbf{Z}^{(1)}}{\partial \mathbf{W}^{(1)}} = \left(\left(\mathbf{W}^{(2)\top} \left(\dots \left(\mathbf{W}^{(L+1)\top} (\hat{\mathbf{Y}} - \mathbf{Y}) \right) \odot \sigma'(\mathbf{Z}^{(L)}) \dots \right) \right) \odot \sigma'(\mathbf{Z}^{(1)}) \right) \mathbf{X}^{(1)\top}$$

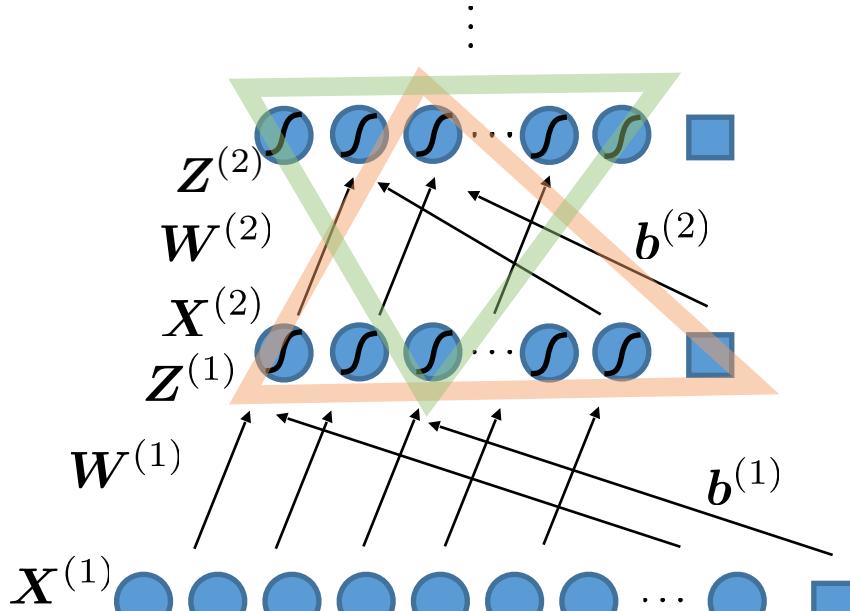
$$\frac{\partial \mathcal{E}}{\partial \mathbf{X}^{(l)}} = \mathbf{W}^{(l)\top} \Delta^{(l)}$$

$$\text{var} \left(\left[\frac{\partial \mathcal{E}}{\partial \mathbf{X}^{(l)}} \right]_{j,t} \right) = \text{var} \left(\sum_{i=1}^{N_i} \mathbf{W}_{i,j}^{(l)} \Delta_{i,t}^{(l)} \right)$$

- To match the variance? $\text{var}(\mathbf{W}_{i,j}) = \frac{1}{N_i}$

- Which one to choose?

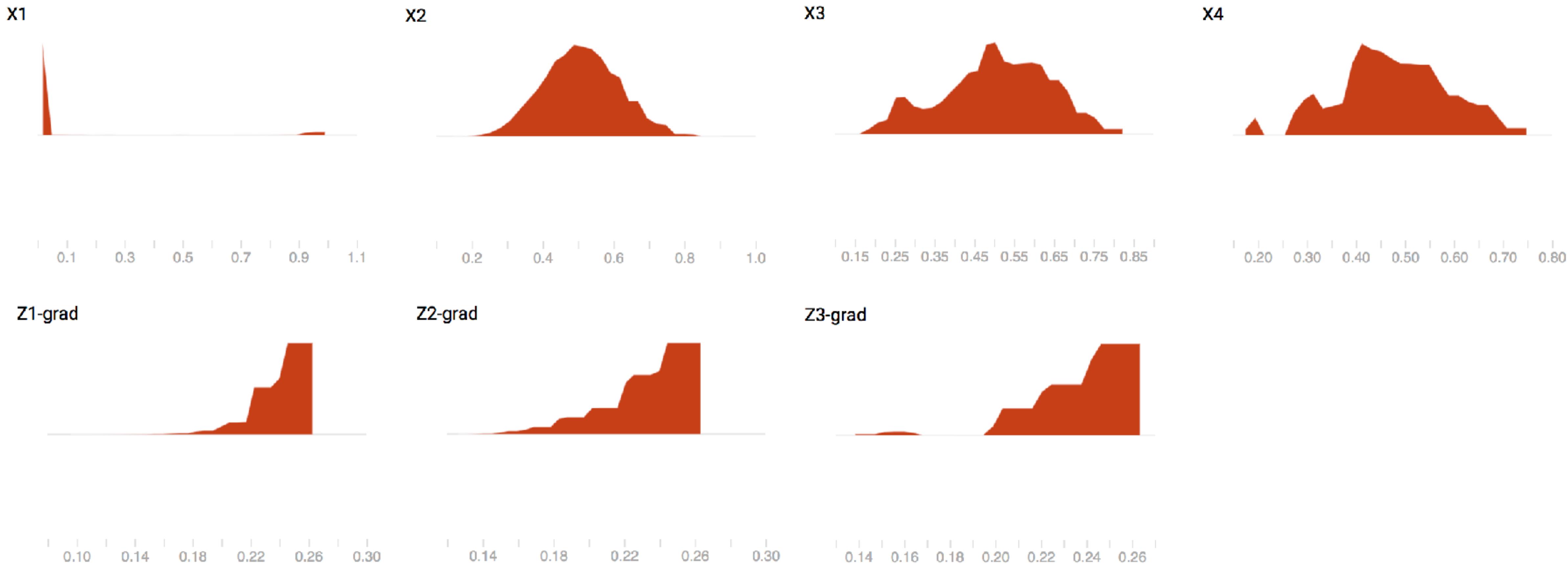
- Average number of units: $(N_j + N_i)/2$
- Hence, $\text{var}(\mathbf{W}_{i,j}) = \frac{2}{N_i + N_j}$



Initialization

-Xavier Initialization

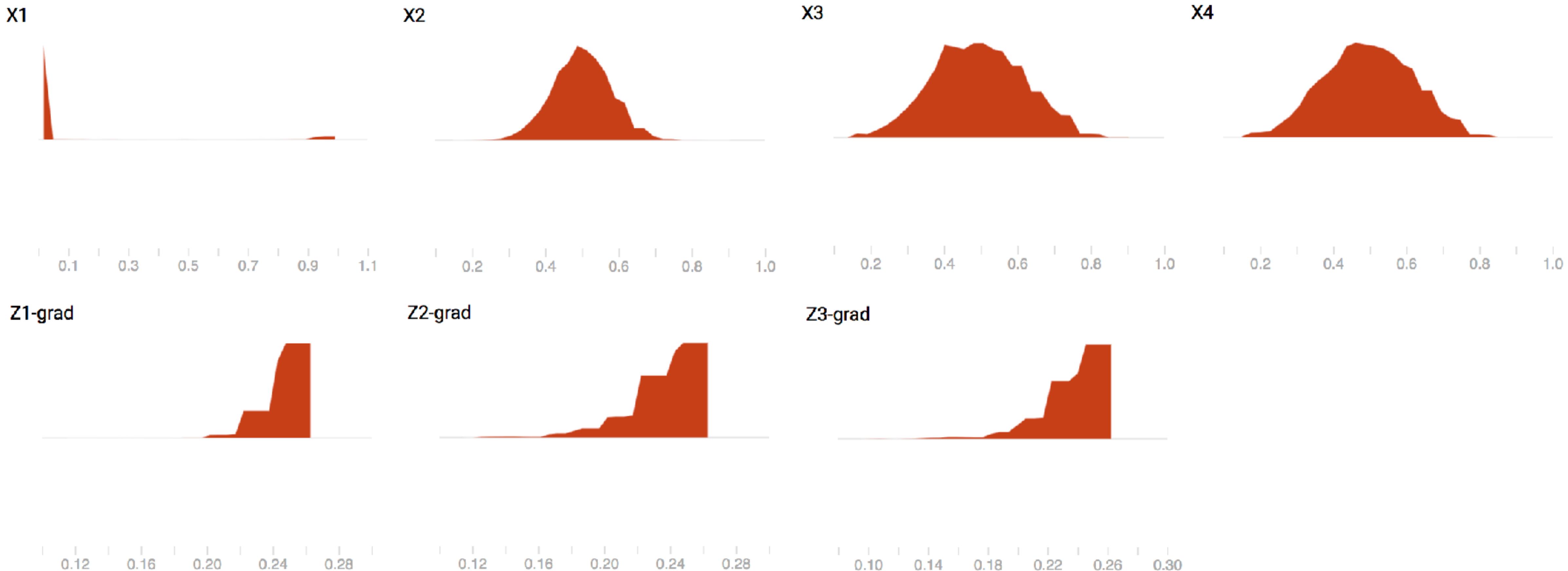
- Xavier initialization with 100 hidden units



Initialization

-Xavier Initialization

- Xavier initialization with 1000 hidden units
 - Properly adapts to the change of the network topology



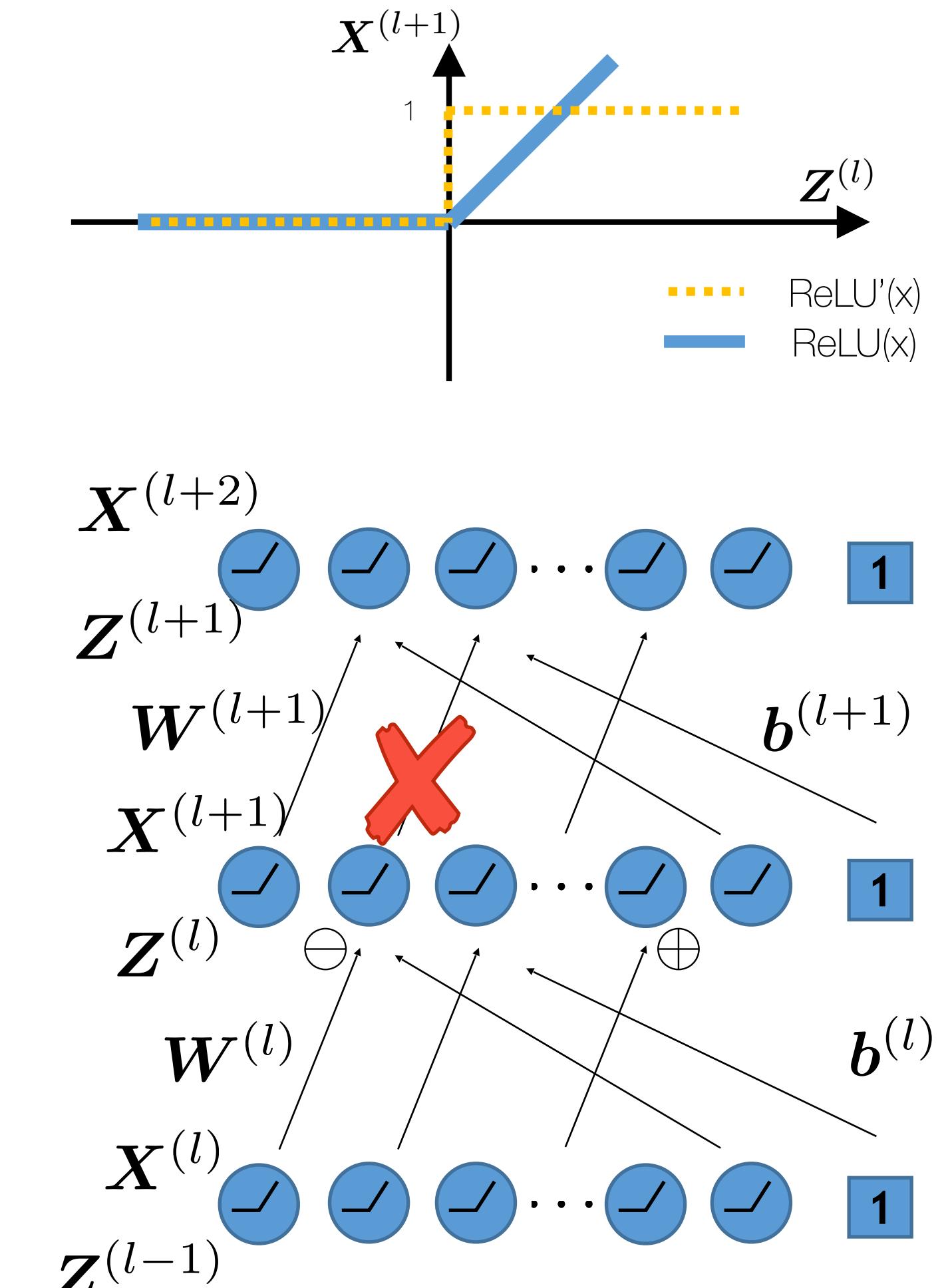
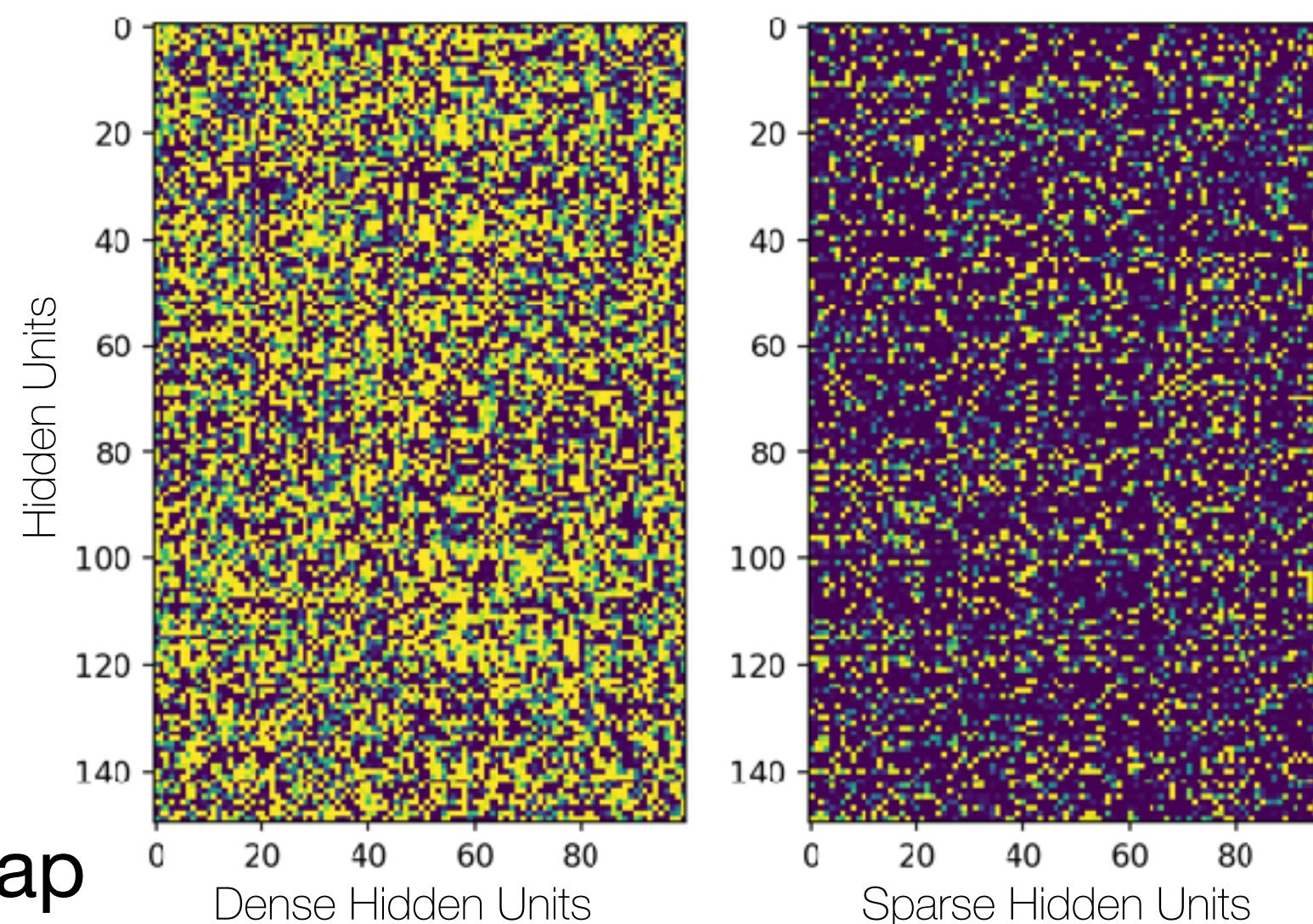
Activation Functions

-Rectified Linear Units (ReLU)

- We all know that sigmoid functions are not a good choice for deep learning
 - Their derivatives block the flow of the BP error
- ReLU doesn't (entirely) block the BP error
 - It actually does block about half of them, but the other half survive
 - It actually blocks about half of the FP (forwardpropagation), too, but the other half survive
- Blocking FP path is actually good
 - A natural way to enforce sparsity
- Blocking BP path is NOT good
 - But we do have half the pathways survived

$$\frac{\partial \mathbf{X}_{i,t}^{(l+1)}}{\partial \mathbf{Z}_{i,t}^{(l)}} = \begin{cases} 1 & \mathbf{Z}_{i,t}^{(l)} \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

- Another thing to note is that ReLU is cheap to calculate both for FP and BP
- Let's see what happens



INDIANA UNIVERSITY

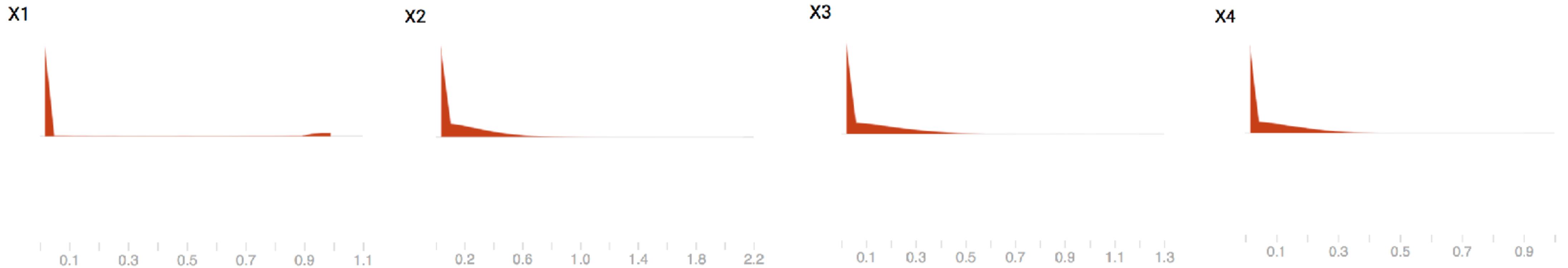
SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Xavier Glorot, Antoine Bordes, and Yoshua Bengio. "Deep sparse rectifier neural networks." Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics. 2011.

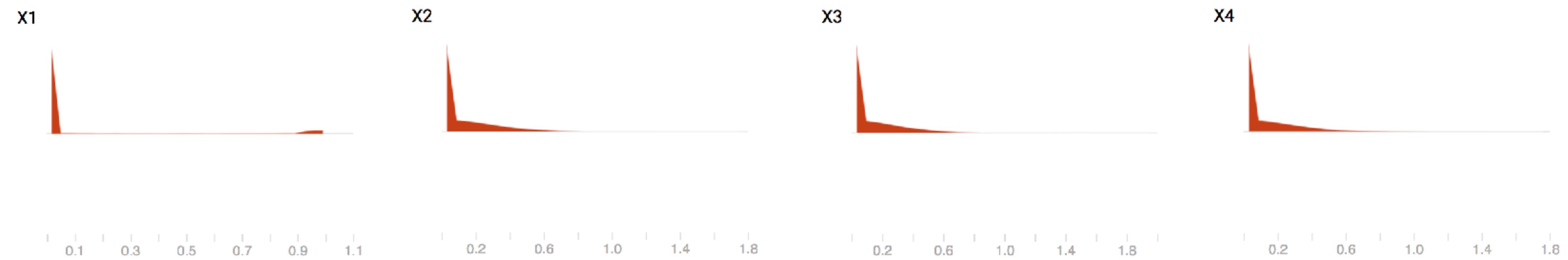
Activation Functions

-Rectified Linear Units (ReLU)

- ReLU with Xavier initialization (notice the different xtick ranges)



- ReLU with He initialization



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

He, Kaiming, et al. "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification." Proceedings of the IEEE international conference on computer vision. 2015.

Activation Functions

-He initialization for ReLU

- Variance of the linear combination $Z_{i,t}^{(l+1)} = \sum_{j=1}^{N_j} W_{i,j}^{(l+1)} X_{j,t}^{(l+1)}$

$$\text{var}(Z_{i,t}^{(l+1)}) = N_j \text{var}(W_{i,j}^{(l+1)}) \text{var}(X_{j,t}^{(l+1)}) + N_j \text{var}(W_{i,j}^{(l+1)}) [E(X_{j,t}^{(l+1)})]^2 + N_j \text{var}(X_{j,t}^{(l+1)}) [E(W_{i,j}^{(l+1)})]^2$$

This is correct only if your r.v.'s on the right hand side is centered

- We have no control over the units, but can still center the weights $\therefore [E(X_{j,t}^{(l+1)})]^2 = 0$

$$\text{var}(Z_{i,t}^{(l+1)}) = N_j \text{var}(W_{i,j}^{(l+1)}) \left(\text{var}(X_{j,t}^{(l+1)}) - [E(X_{j,t}^{(l+1)})]^2 \right) = N_j \text{var}(W_{i,j}^{(l+1)}) [E(X_{j,t}^{(l+1)})^2] \stackrel{?}{=} \text{var}(Z_{j,t}^{(l)})$$

- We look for $\text{var}(W_{i,j}^{(l+1)})$ that meets the final equation

- Let's see what happens in the lower layer

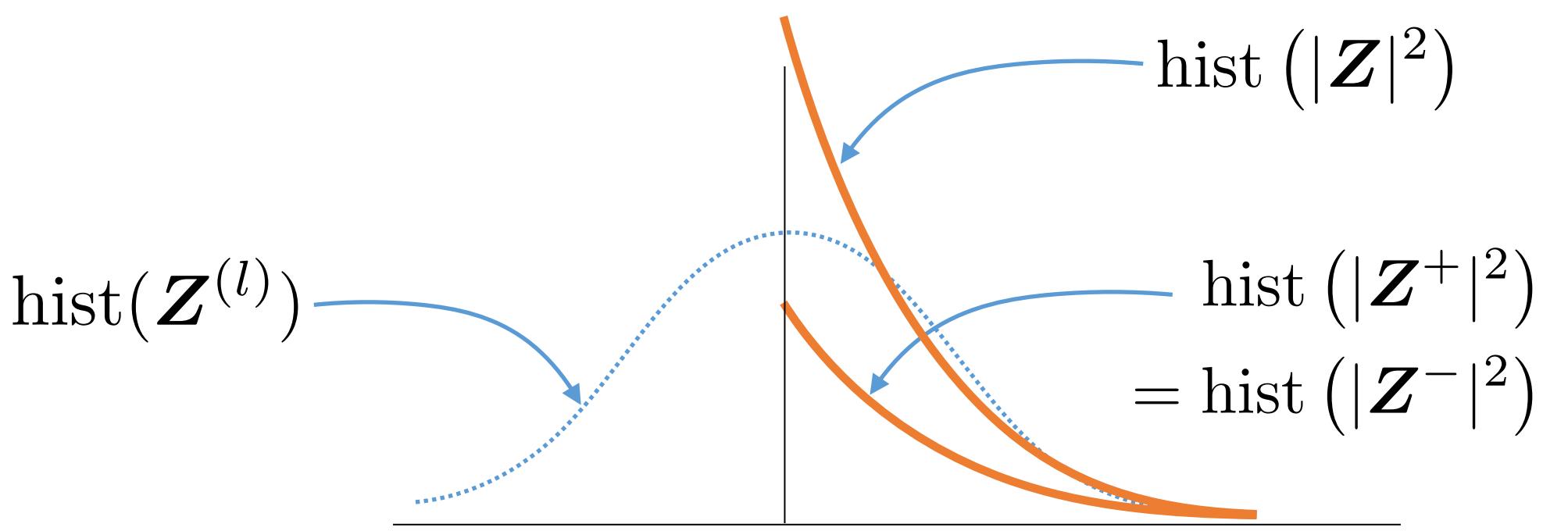
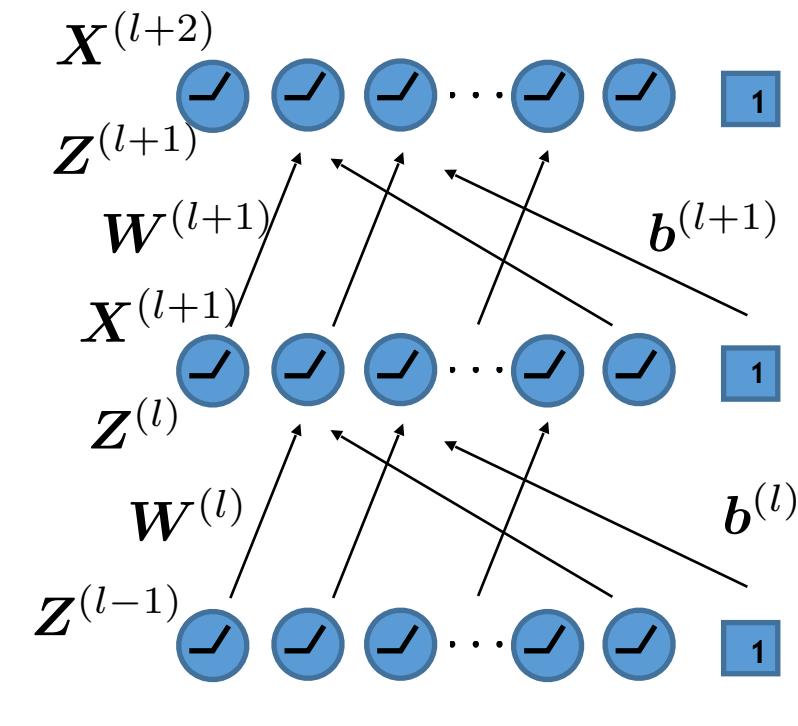
- Since we assume a centered symmetric distribution for $W_{i,j}^{(l)}$, should $P(Z_{i,j}^{(l)})$ be, too

$$\begin{aligned} E(X_{j,t}^{(l+1)2}) &= E(\text{ReLU}(Z_{j,t}^{-(l)})^2)/2 + E(\text{ReLU}(Z_{j,t}^{+(l)})^2)/2 \\ &= E((Z_{j,t}^{+(l)})^2)/2 = E(Z_{j,t}^{(l)2})/2 = \text{var}(Z_{j,t}^{(l)})/2 \end{aligned}$$

- Therefore,

$$\text{var}(Z_{i,t}^{(l+1)}) = \frac{N_j}{2} \text{var}(W_{i,j}^{(l+1)}) \text{var}(Z_{j,t}^{(l)}) \stackrel{?}{=} \text{var}(Z_{j,t}^{(l)})$$

$$\text{var}(W_{i,j}^{(l+1)}) = \frac{2}{N_j}$$



INDIANA UNIVERSITY

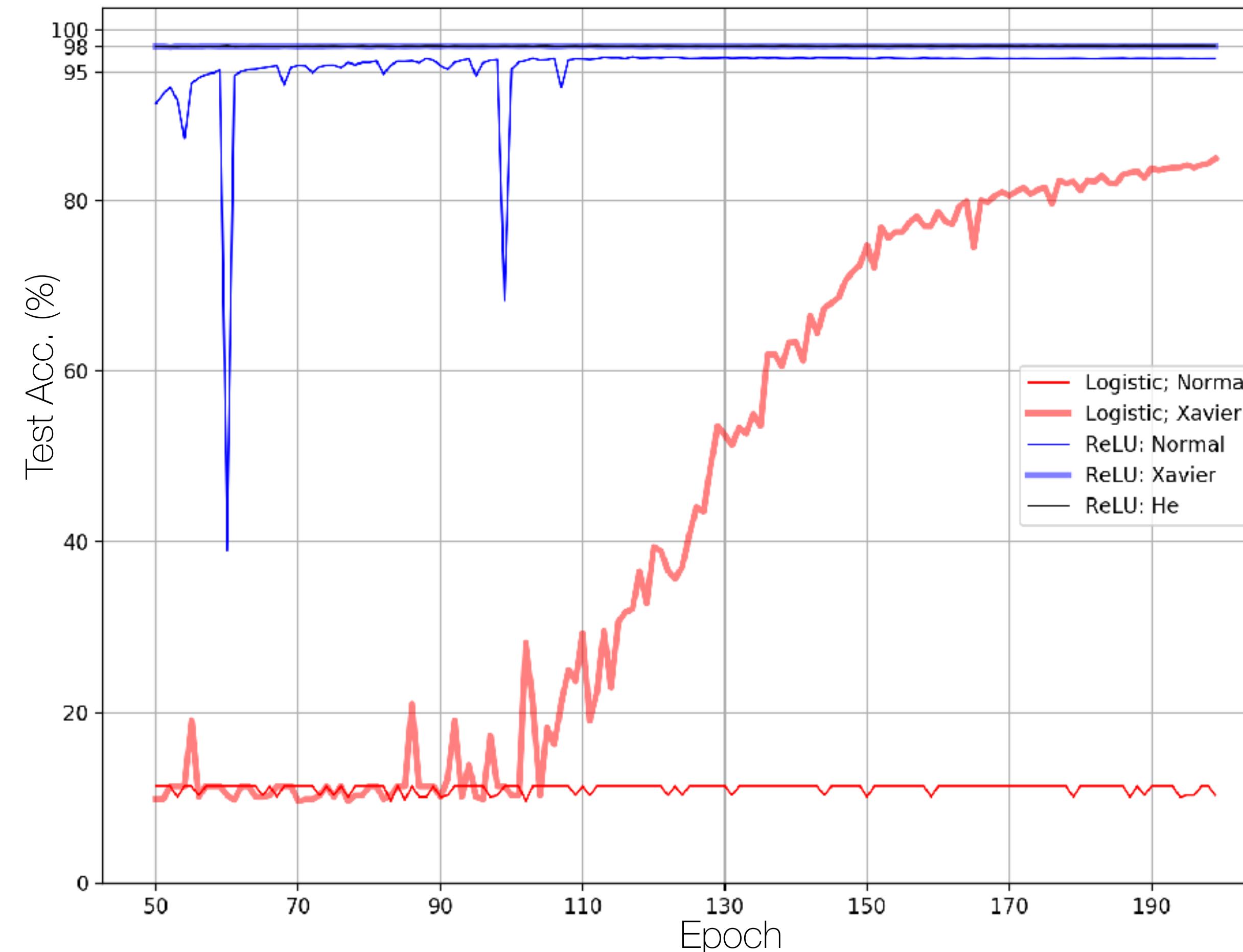
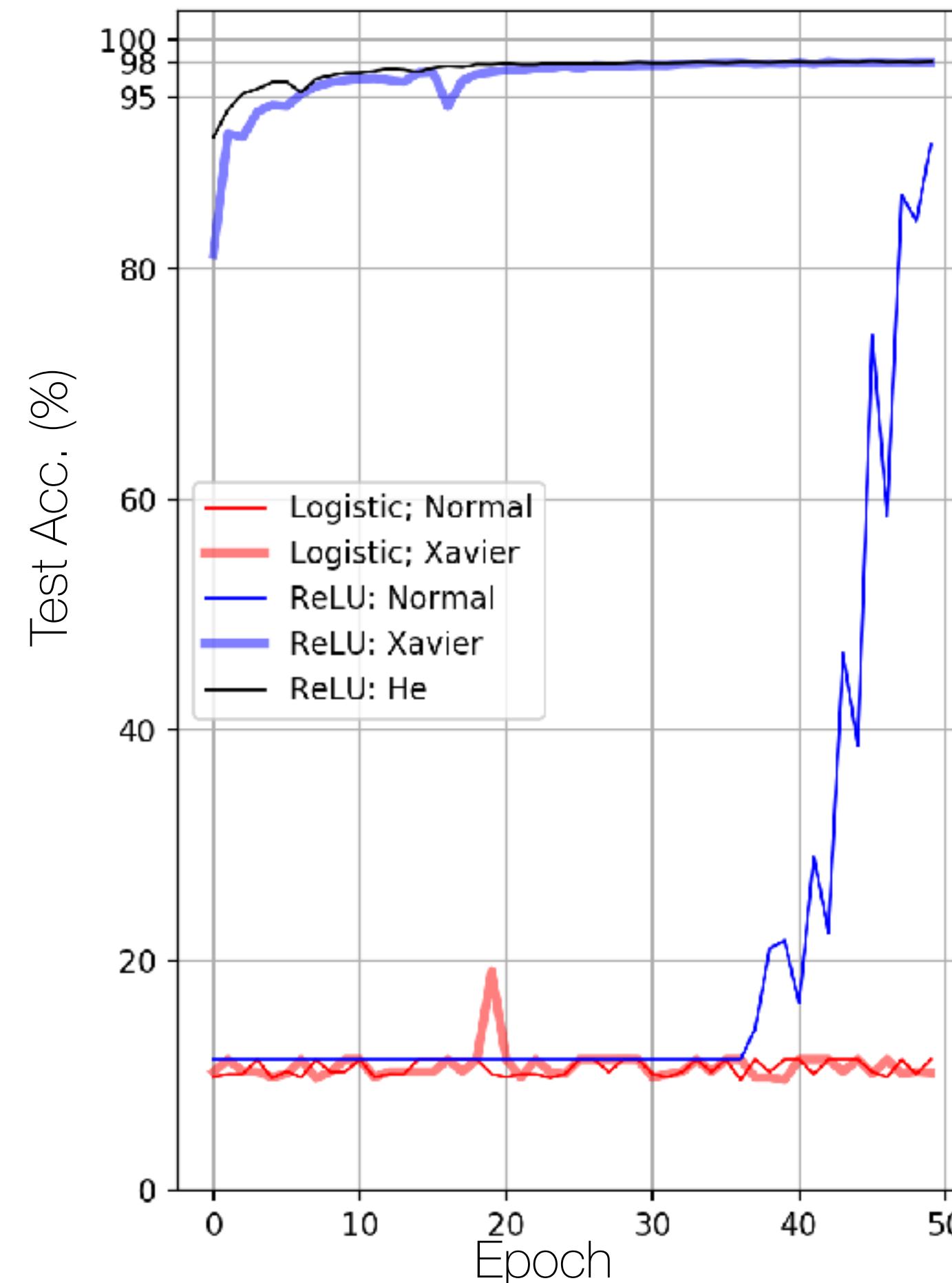
SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

He, Kaiming, et al. "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification." Proceedings of the IEEE international conference on computer vision. 2015.

Activation Functions

-He initialization for ReLU

- For a 512X5 network for MNIST



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Activation Functions

-The other activation functions

Leaky ReLU

$$f(x) = \max(x, 0) + a \min(0, x)$$

Parametric ReLU

$$f(x) = \max(x, 0) + a \min(0, x)$$

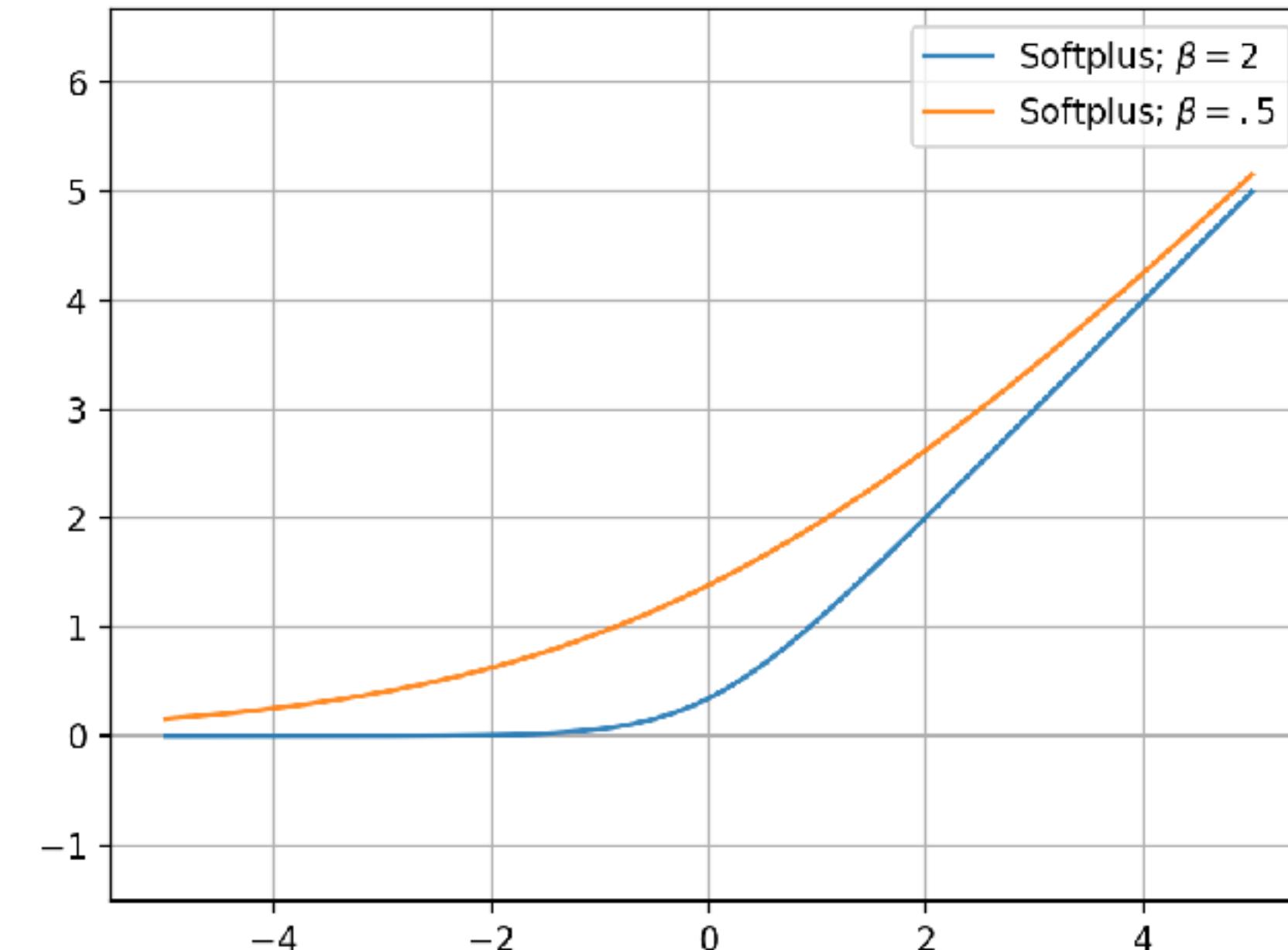
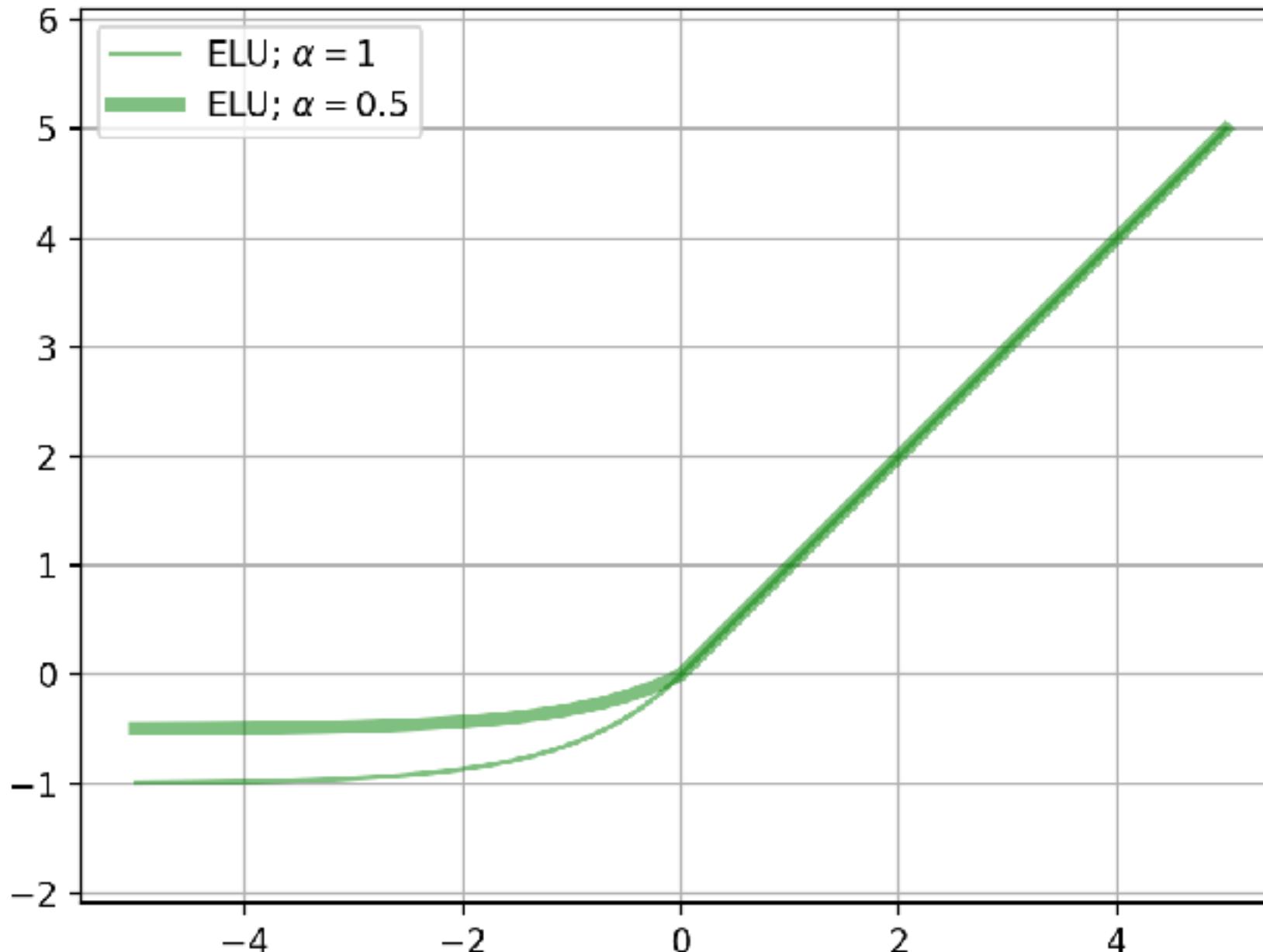
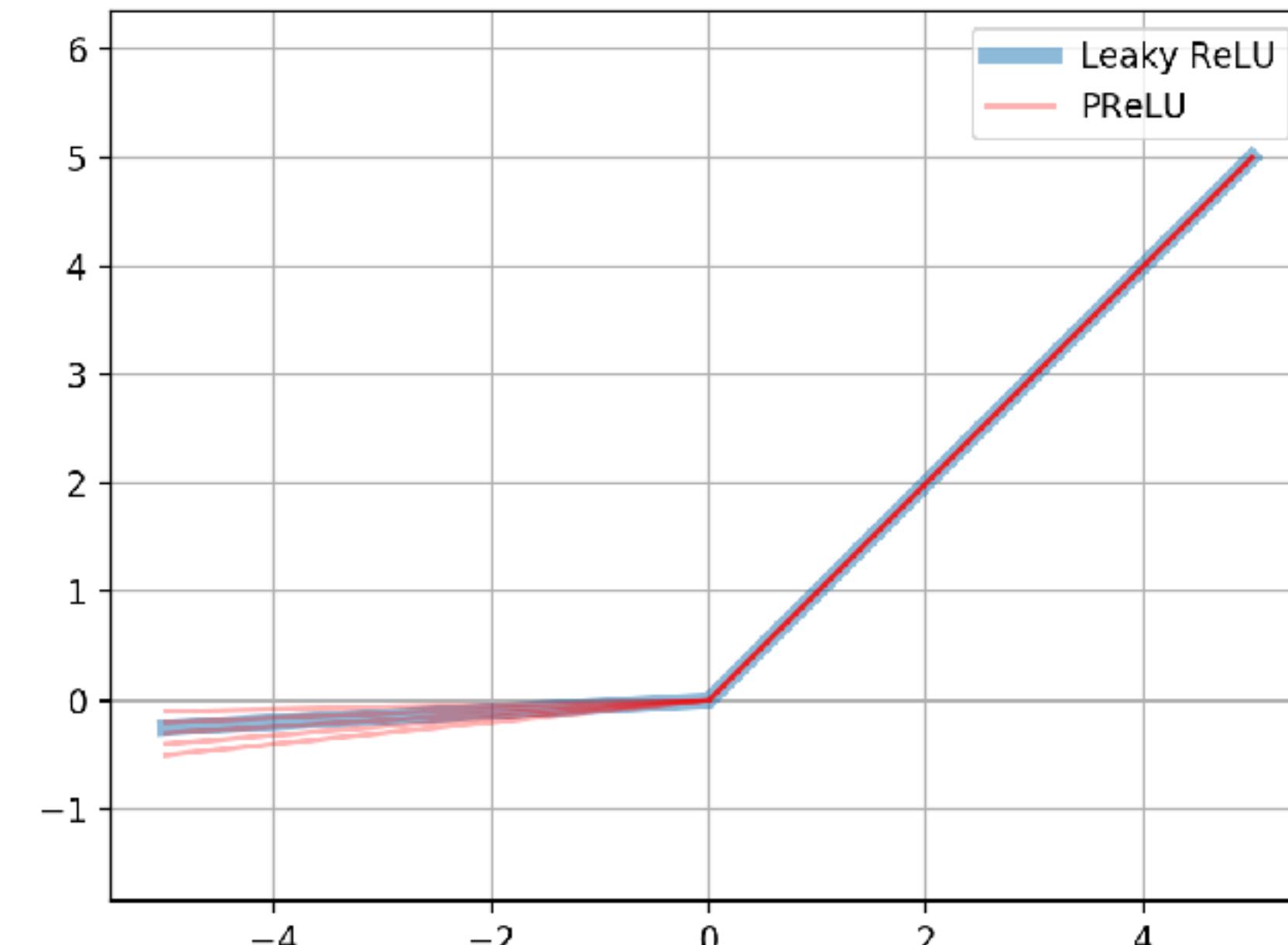
Softplus

$$f(x) = \frac{1}{\beta} \log(1 + \exp(\beta x))$$

Exponential LU

$$f(x) = \max(x, 0) + \min(0, \alpha(\exp(x) - 1))$$

You name it!



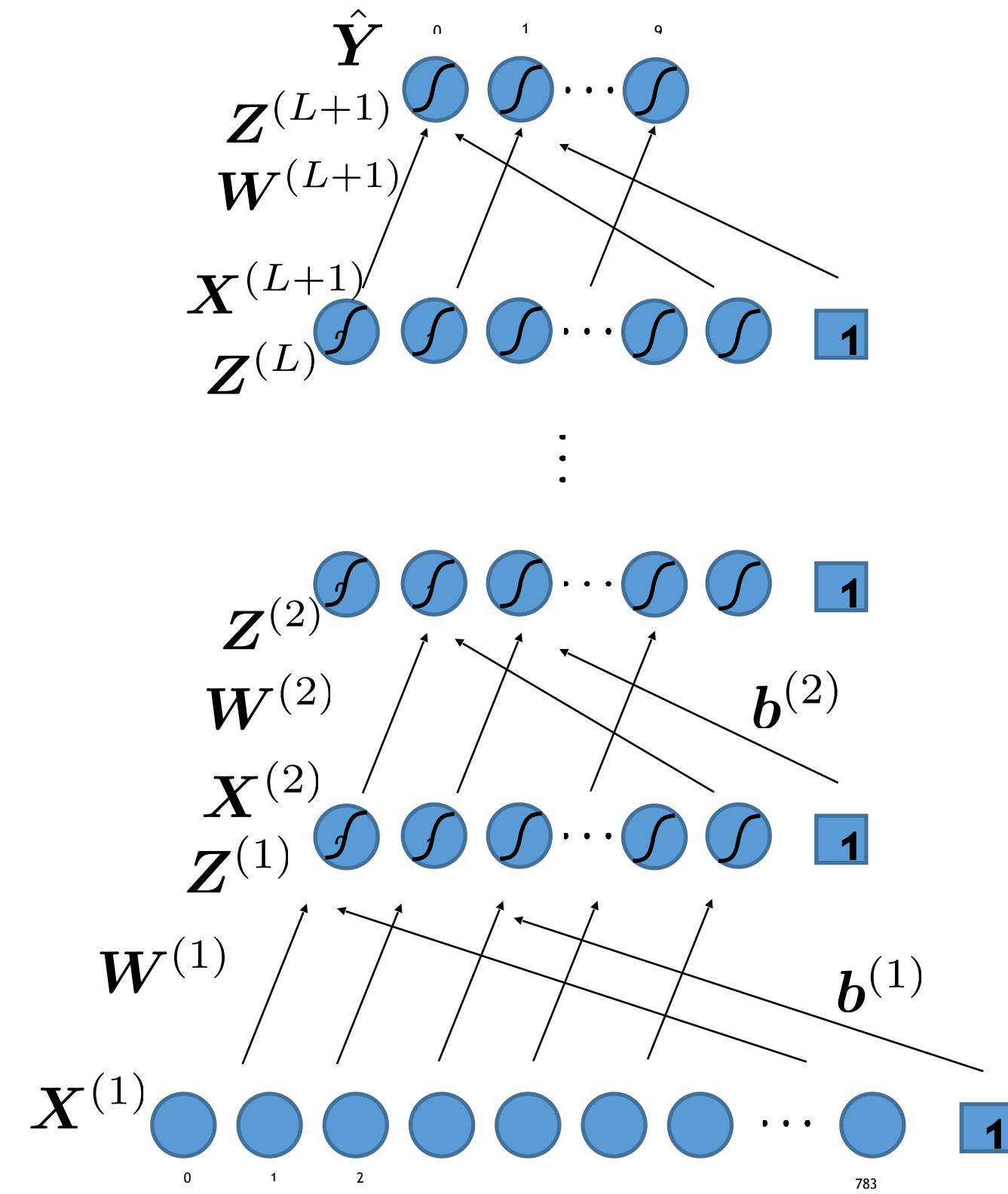
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Batch Normalization

-A straightforward way to keep the distribution organized

- What is a batch?
- Let me do some calculation first
 - I've got a 1024X3 DNN for MNIST classification
 - What's the size of the input training matrix $X^{(1)}$?
 - $784 \times 55,000 = 43,120,000$
 - In byte?
 - + $784 \times 55,000 \times 32 / 8 = 172.48\text{MB}$
 - What's the size of $X^{(2)}$?
 - $1024 \times 55,000 \times 32 / 8 = 225.28\text{MB}$
 - What's the size of $W^{(2)}$?
 - $1024 \times 1024 \times 32 / 8 = 4.2\text{MB}$
 - These days this amount of calculation can be done in one shot in some good computers
 - But was too big back in the days
 - Also, there are many datasets larger than MNIST out there
 - Mini-batches is a compromised version of batch processing and SGD
 - When the mini-batch size is 128: $X^{(1)} \in \mathbb{R}^{784 \times 128}, X^{(2)} \in \mathbb{R}^{1024 \times 128}, \dots$



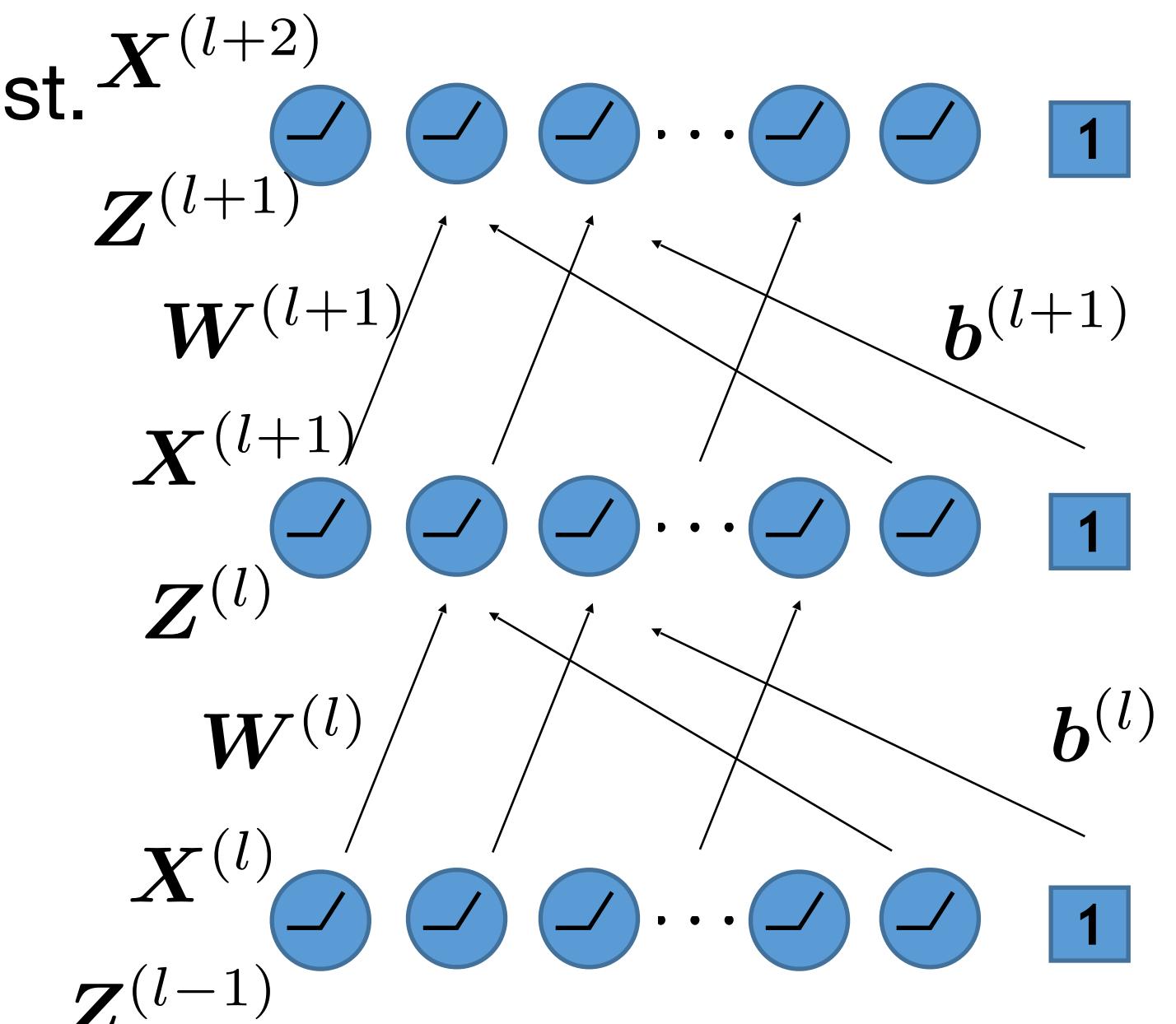
Batch Normalization

-A straightforward way to keep the distribution organized

- So far we've seen that some wrong initialization can damage the convergence
 - The intermediate variables can have drastically different distributions across layers
 - Due to the wrong choice of initialization, and also because of a particular choice of activation functions
 - But it can happen anytime during training in deep learning
- Batch normalization: makes sure a variable (unit) follow the standard normal dist.
 $X_{i,t}^{(l)} \sim [\text{Some distribution we don't like}] \xrightarrow{\text{Batch Normalization}} \overline{X}_{i,t}^{(l)} \sim \mathcal{N}(0, 1)$

$t = [1, \dots, N_b]$ is for the sample index for the batch
 $b = [1, \dots, B]$ denotes batches

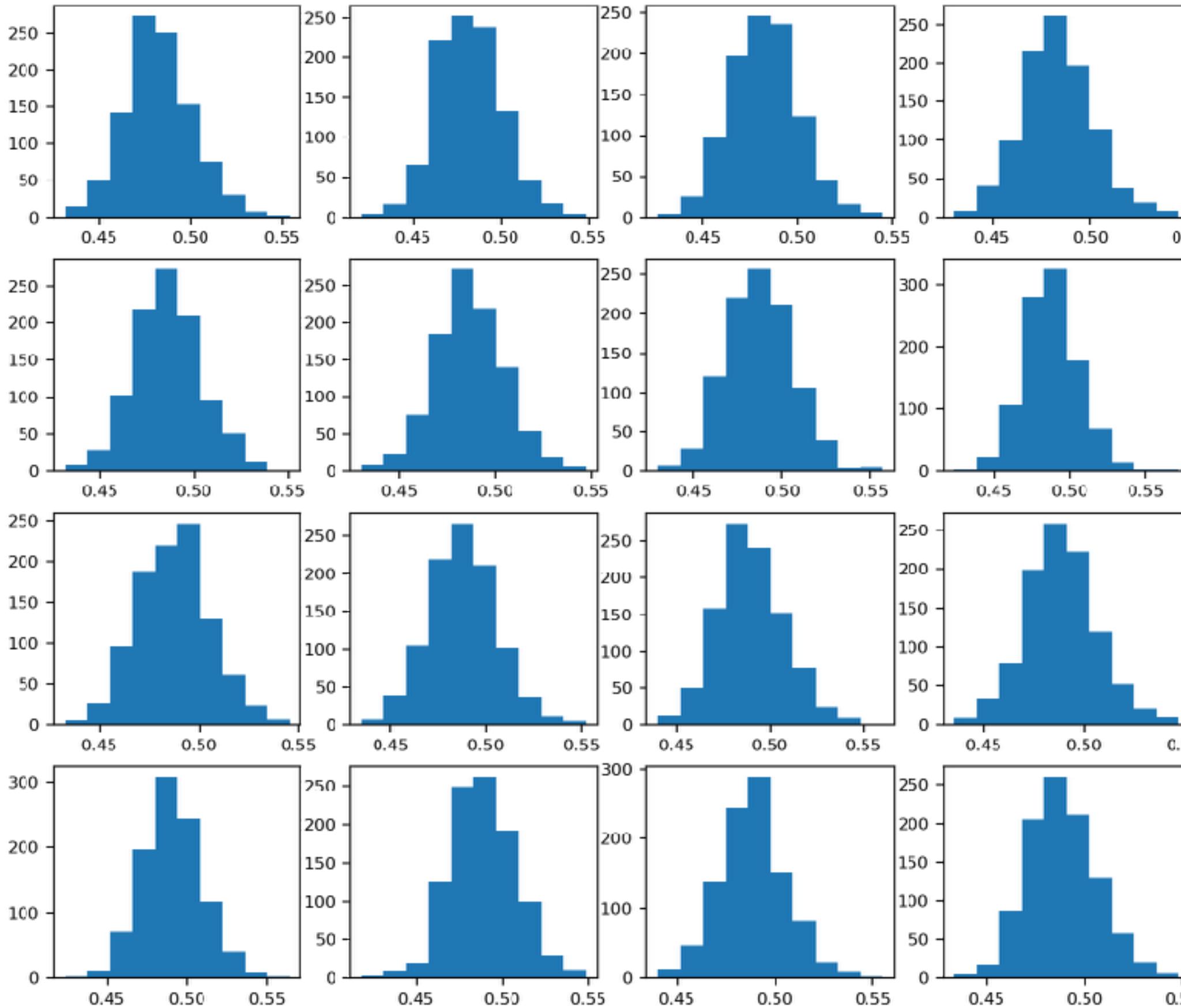
$$\mu_{i,b}^{(l)} = \frac{1}{N_b} \sum_{t=1}^{N_b} X_{i,t}^{(l)}$$
$$\sigma_{i,b}^{(l)} = \frac{1}{N_b} \sum_{t=1}^{N_b} (X_{i,t}^{(l)} - \mu_{i,b}^{(l)})$$
$$\overline{X}_{i,t}^{(l)} = \frac{X_{i,t}^{(l)} - \mu_{i,b}^{(l)}}{\sigma_{i,b}^{(l)}}$$
 - Note that N_b is the number of samples in a (mini)batch
- This is a bit problematic
 - All input units of all layers are forced to produce the same sample distribution
 - Let's give them some flexibility $\widehat{X}_{i,t}^{(l)} = \gamma_i^{(l)} \overline{X}_{i,t}^{(l)} + \beta_i^{(l)}$
- Eventually,
batch normalization standardizes mini-batch statistics, but keeps the empirical statistics of the entire dataset



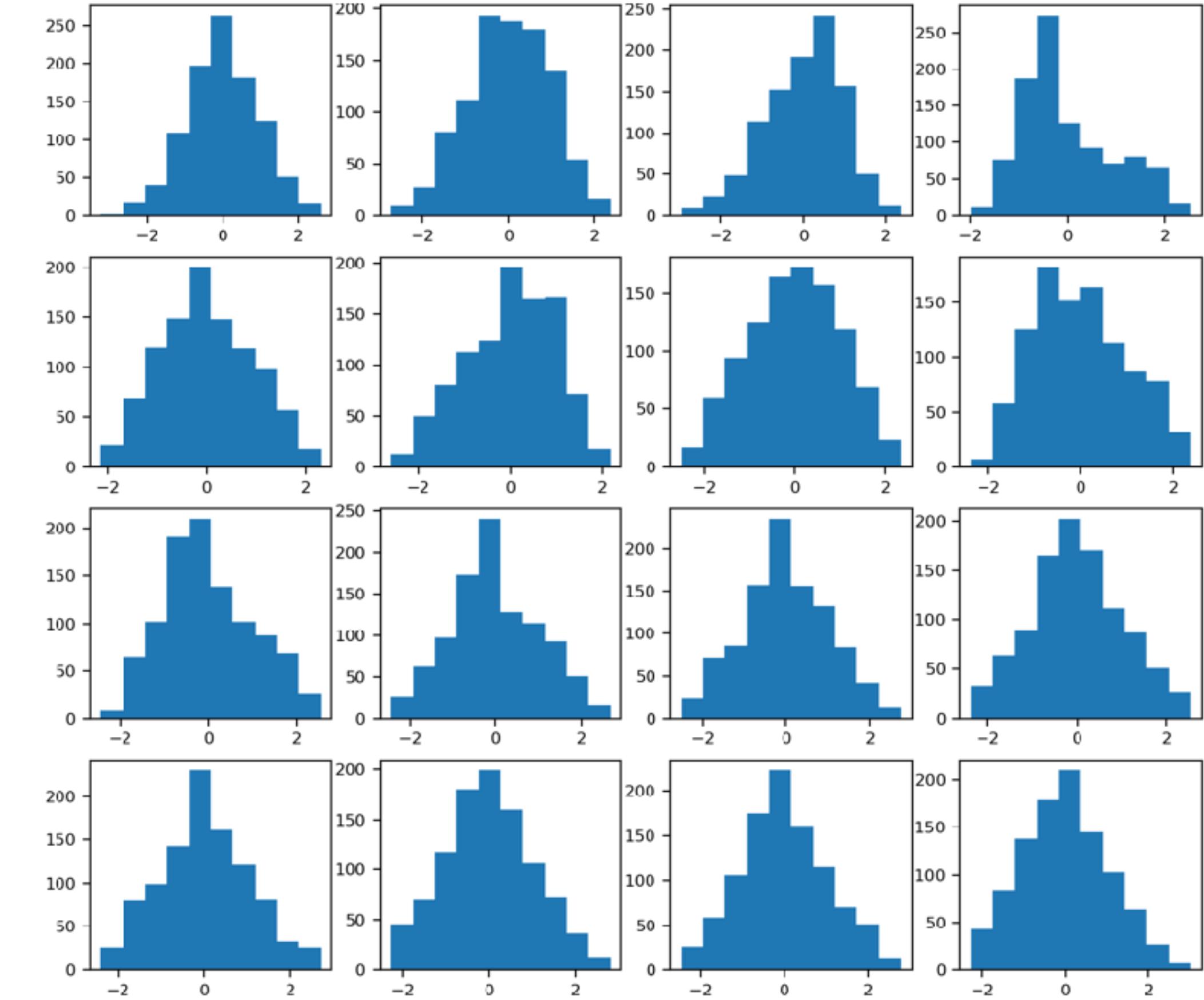
Batch Normalization

-A straightforward way to keep the distribution organized

- 1st layer, 1st epoch; a histogram is from within a mini-batch for a fixed hidden unit output



No BN



BN



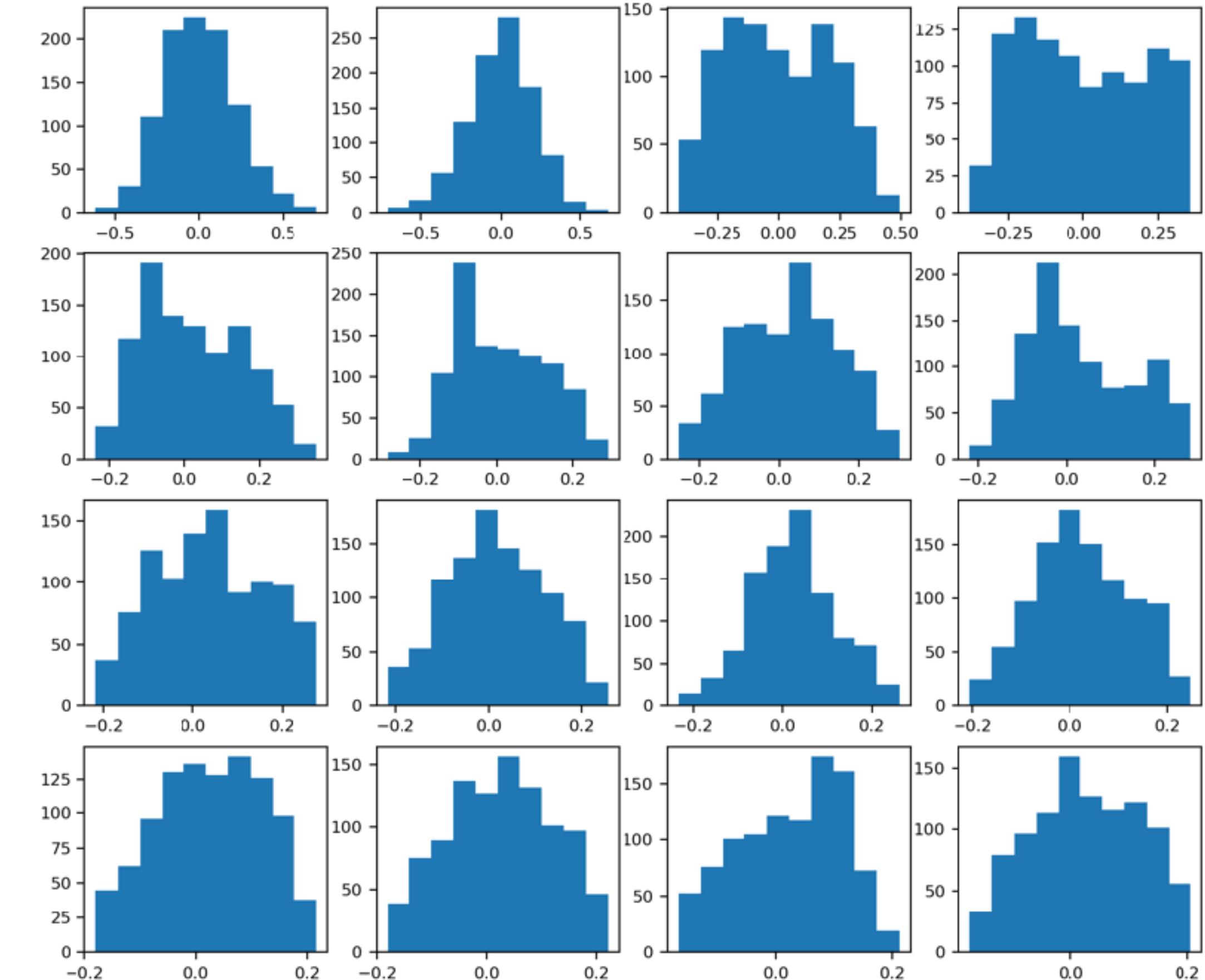
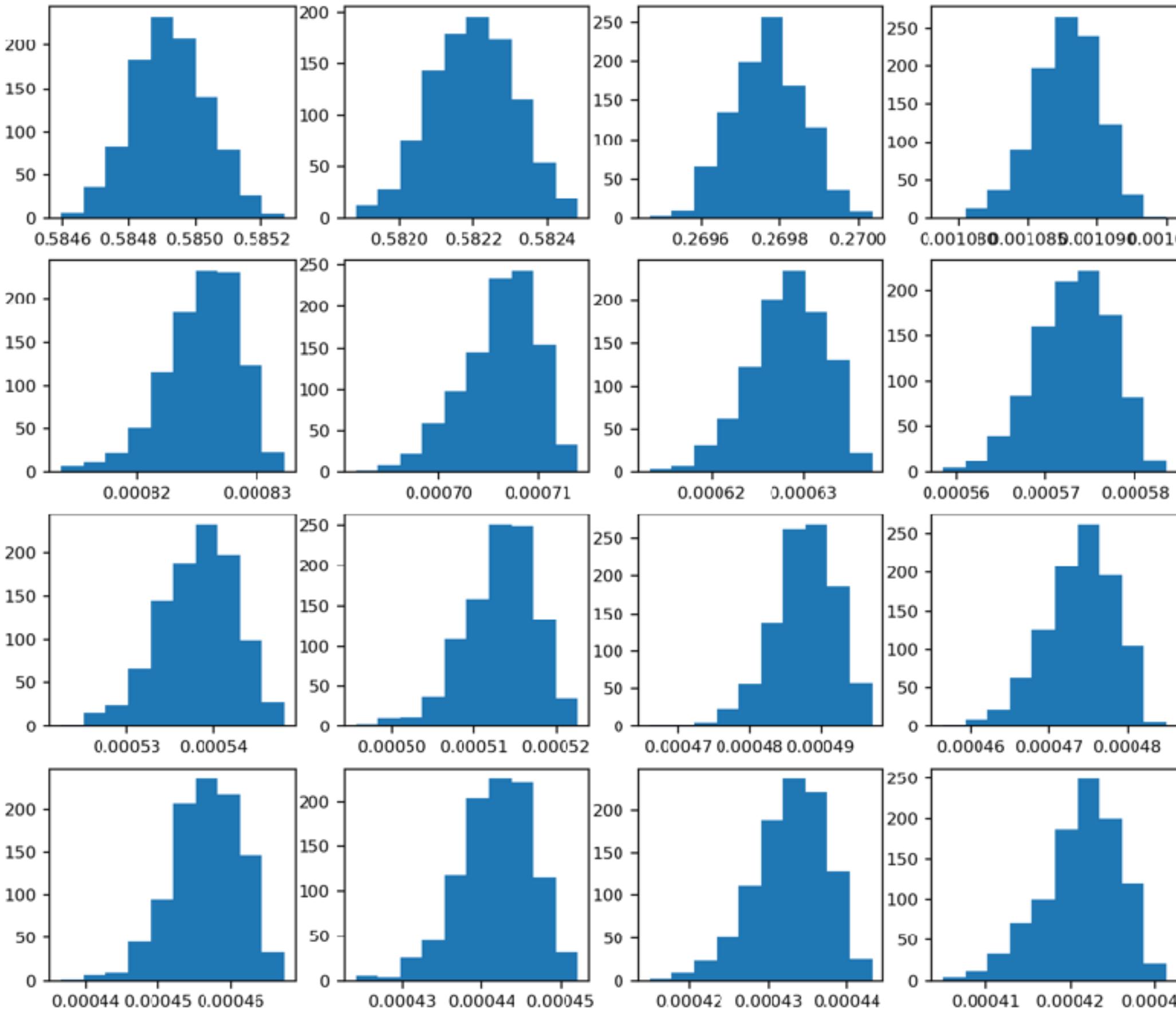
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Batch Normalization

-A straightforward way to keep the distribution organized

- 3rd layer, 1st epoch; a histogram is from within a mini-batch for a fixed hidden unit output



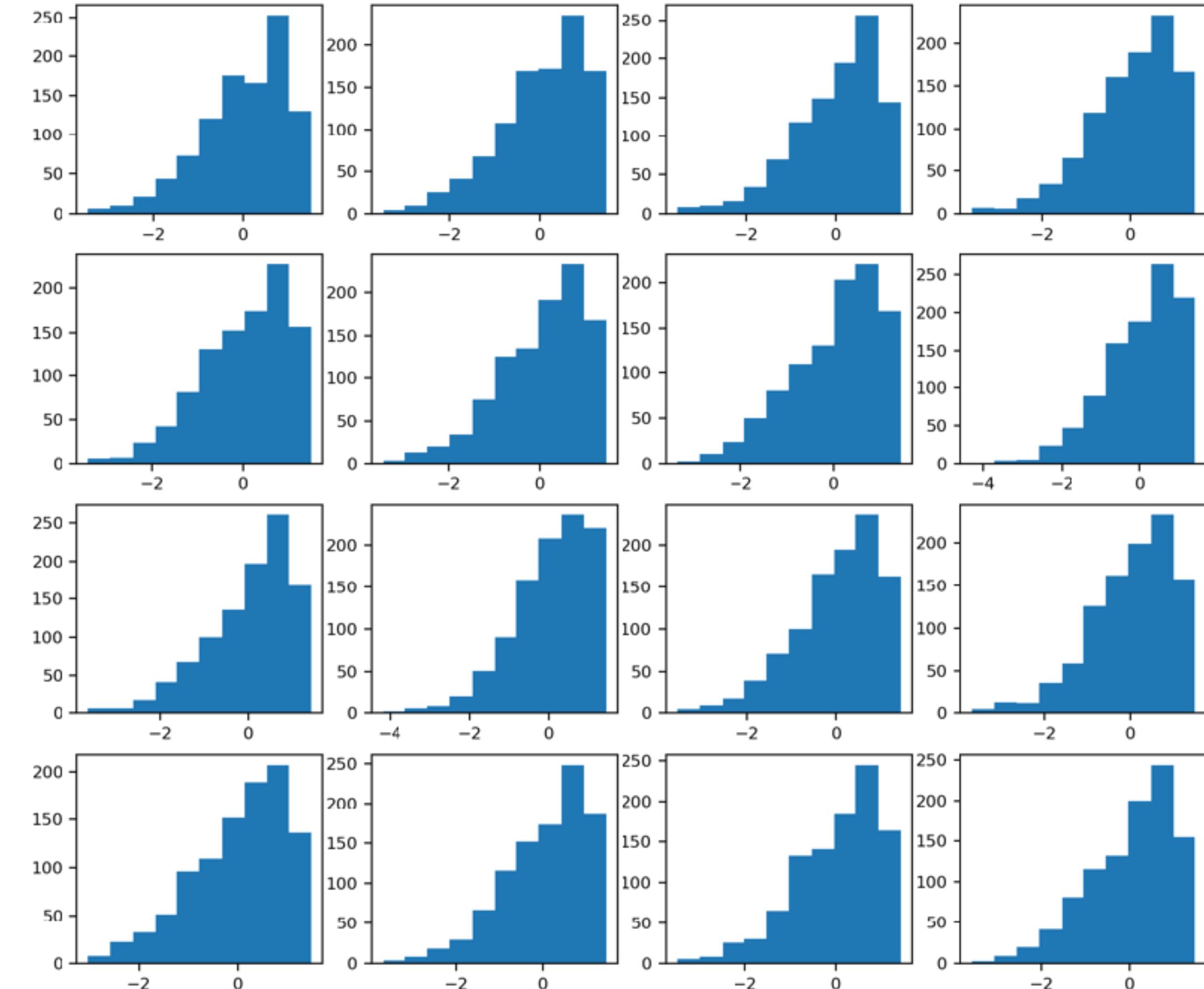
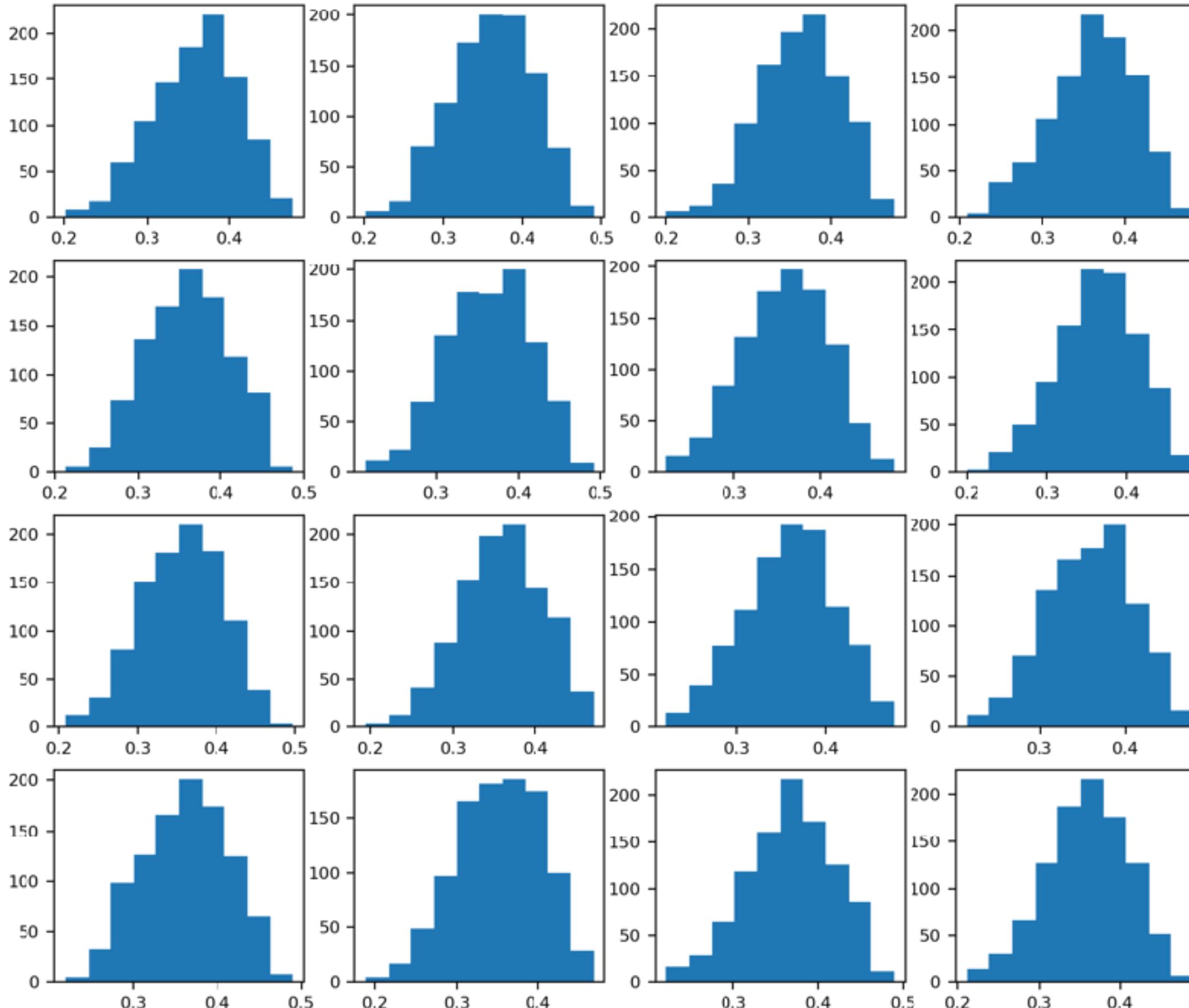
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Batch Normalization

-A straightforward way to keep the distribution organized

- 1st layer, 20th epoch; a histogram is from within a mini-batch for a fixed hidden unit output



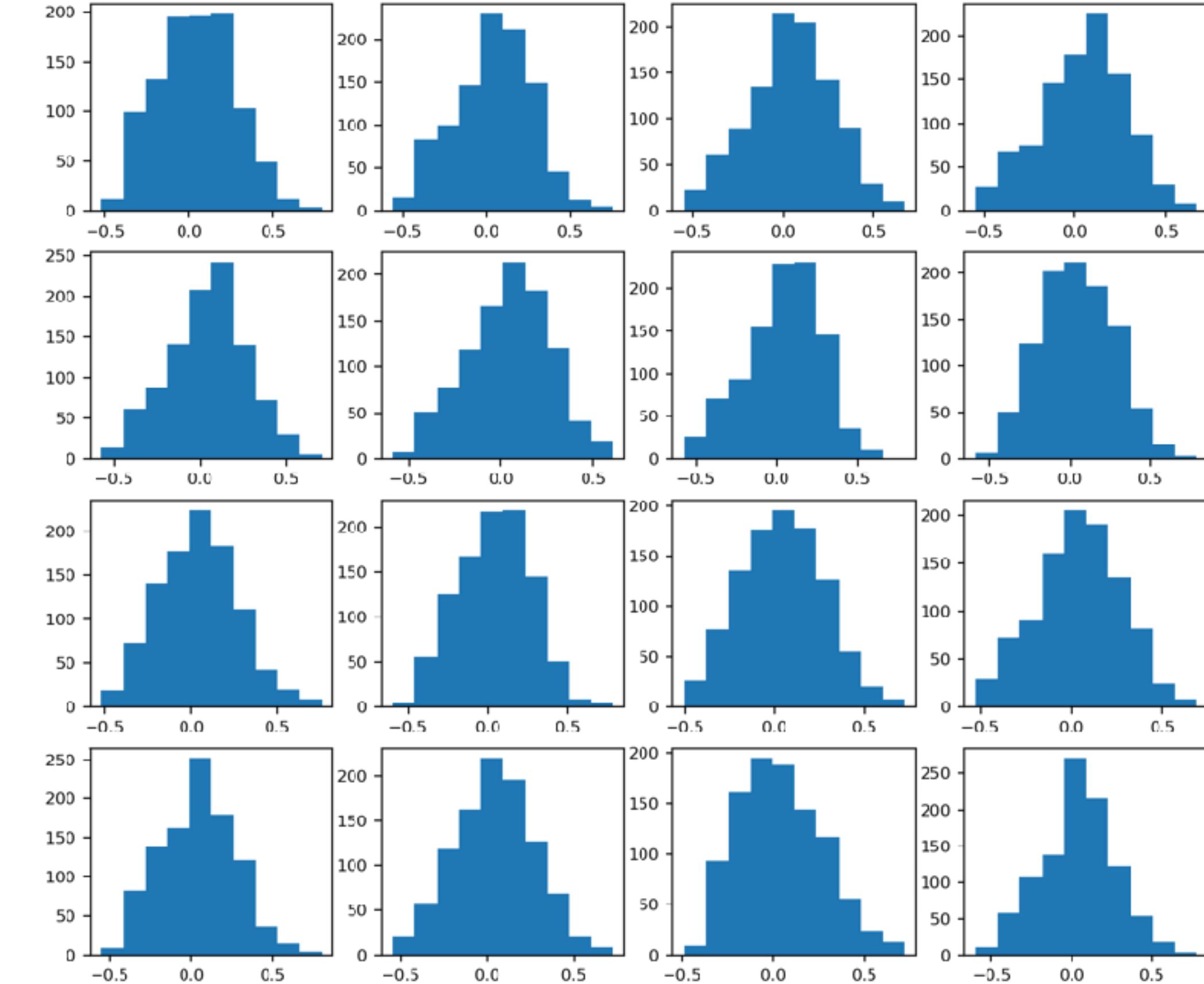
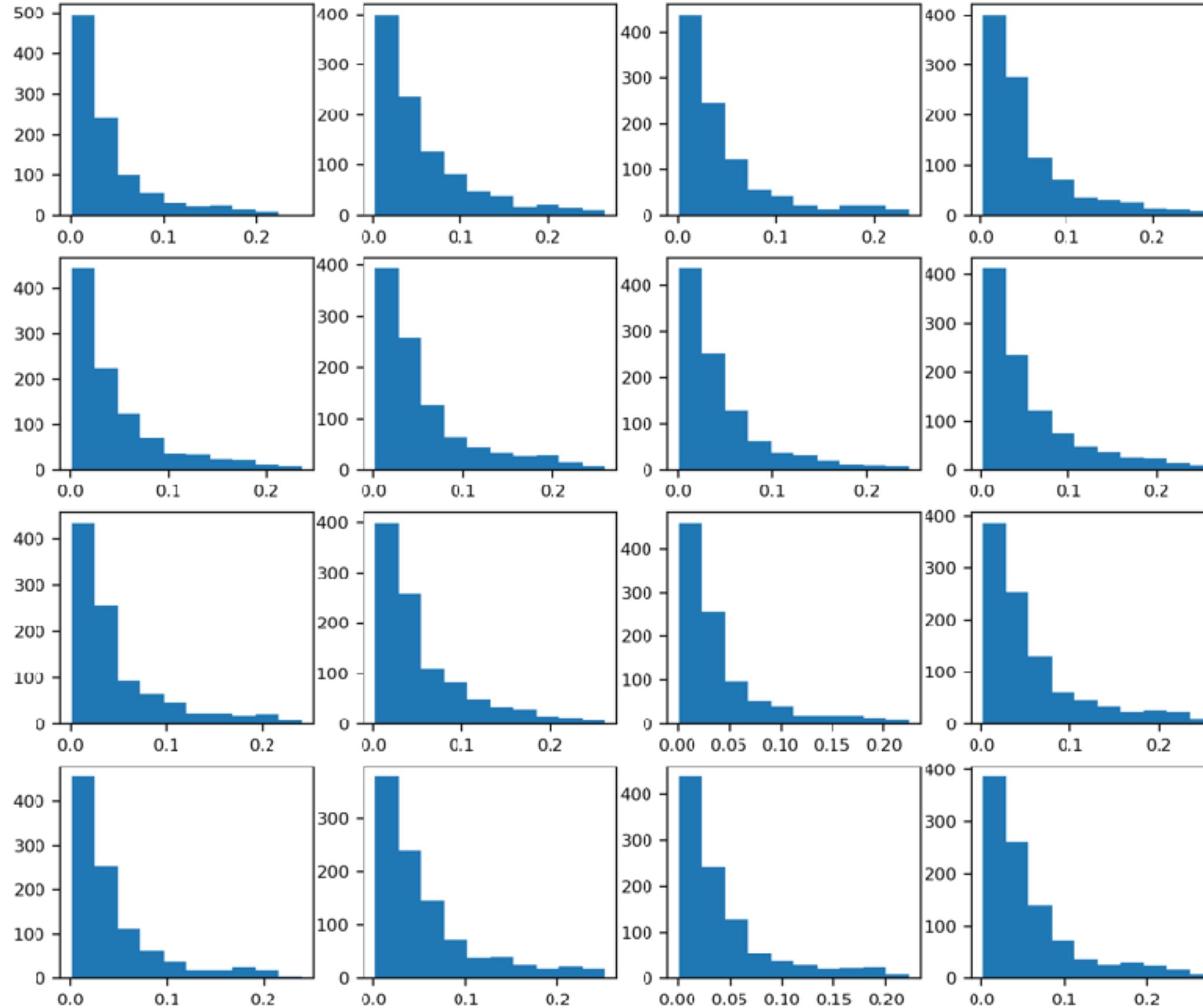
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Batch Normalization

-A straightforward way to keep the distribution organized

- 3rd layer, 20th epoch; a histogram is from within a mini-batch for a fixed hidden unit output



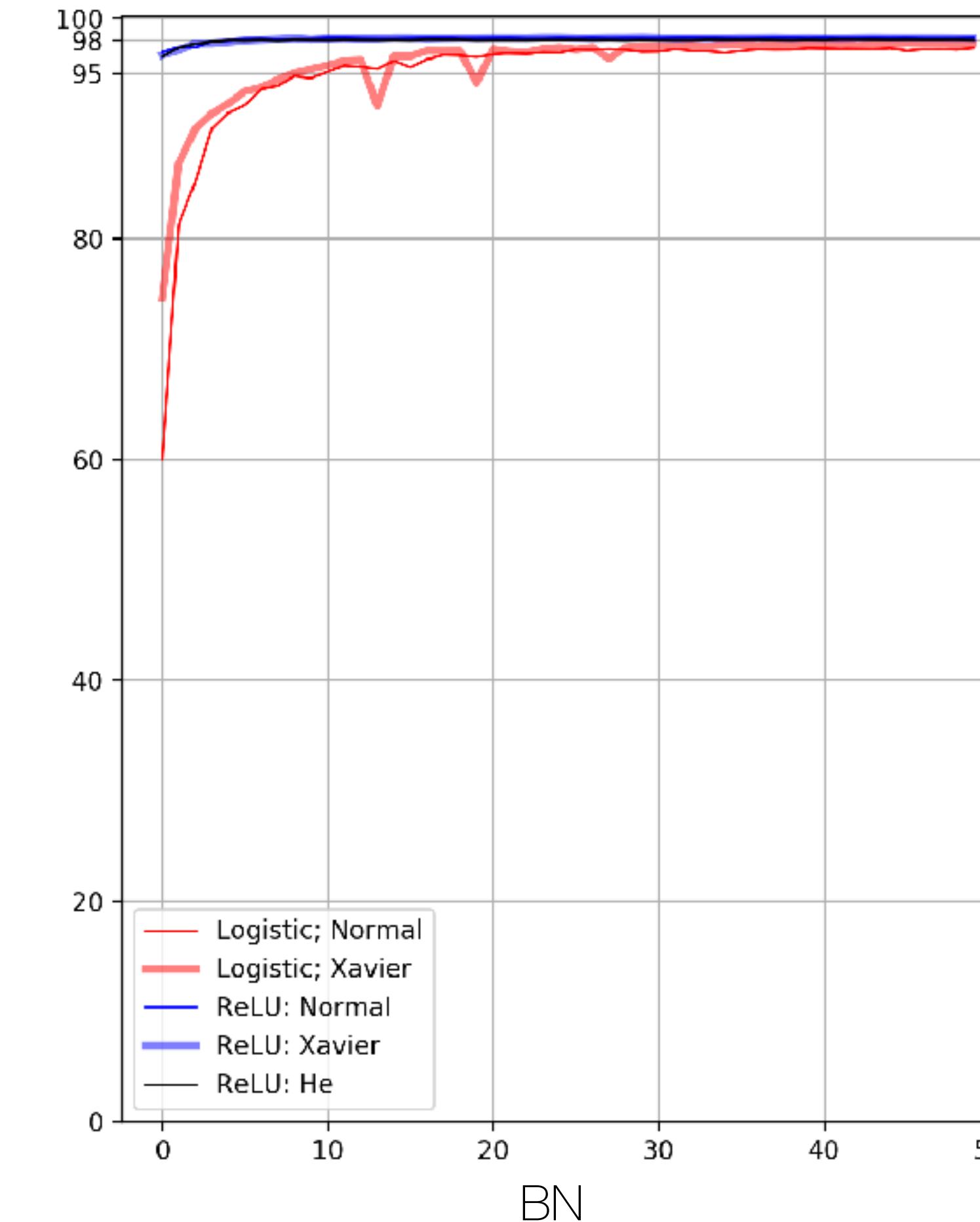
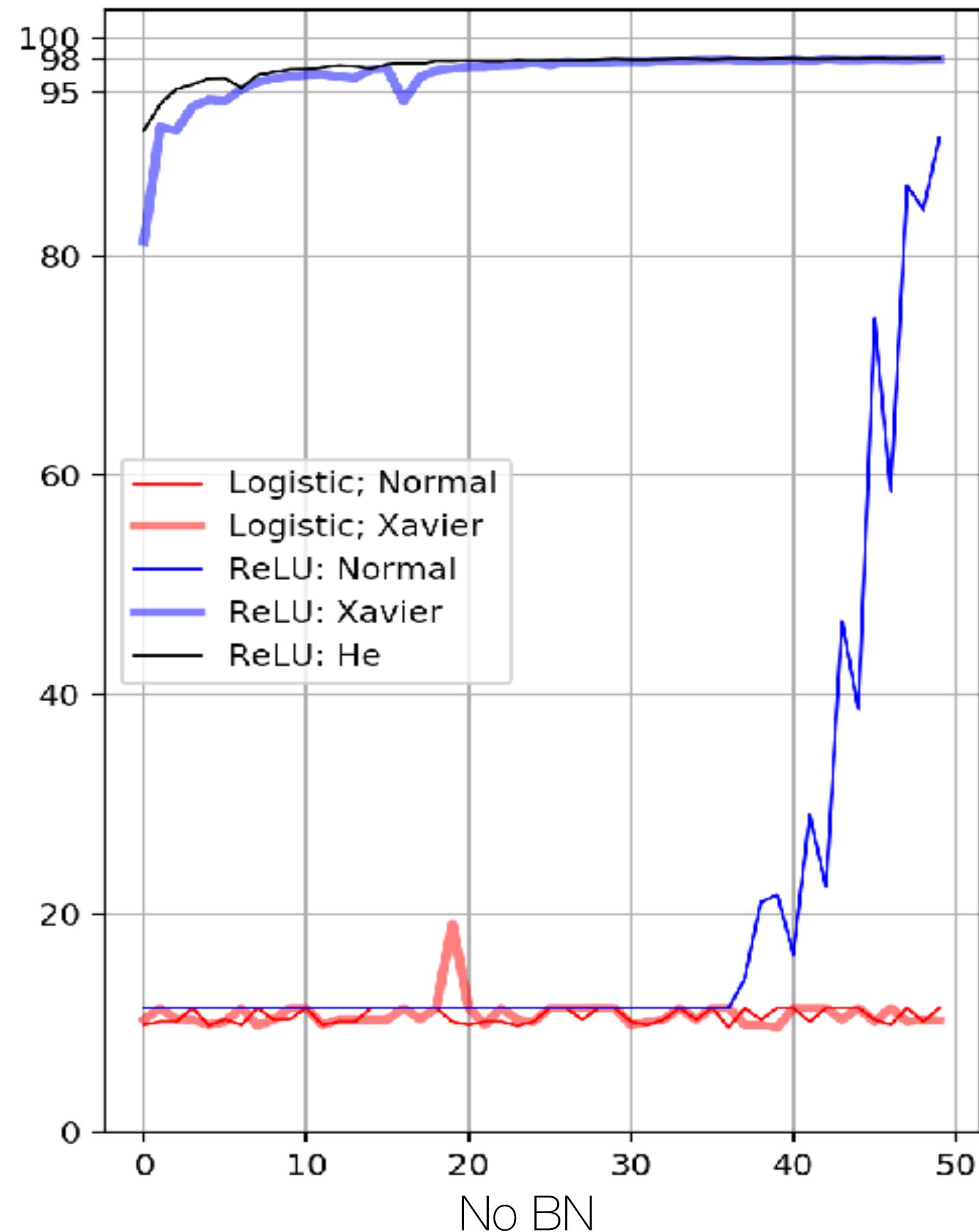
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Batch Normalization

-A straightforward way to keep the distribution organized

- Batch normalization reduces the influence of the initialization methods and activation functions



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Momentum and Adaptive Learning Rates

Yes, they are related, too.

Gradient Descent

-Aren't you already tired of it?

- The (stochastic) gradient descent method has some problems
 - It assumes the objective function is locally linear
 - Newton's method can fix this, but not suitable for deep learning. Why?
 - Inverse of Hessian matrix is too expensive because of the too many parameters to estimate
 - Quasi-Newton's method is doable, but not so popular choice
 - A fixed learning rate can be either too large or too small
 - Too large LR can overlook an important local minimum
 - Too small LR can stuck in an unimportant local minimum
 - A decaying learning rate can be a reasonable solution, but scheduling is not easy
 - A fixed scalar learning rate might not be best for all parameters
- So far unless I confessed in advance I've used SGD
 - So my previous arguments about initialization and batch normalization can change if I choose another optimization technique



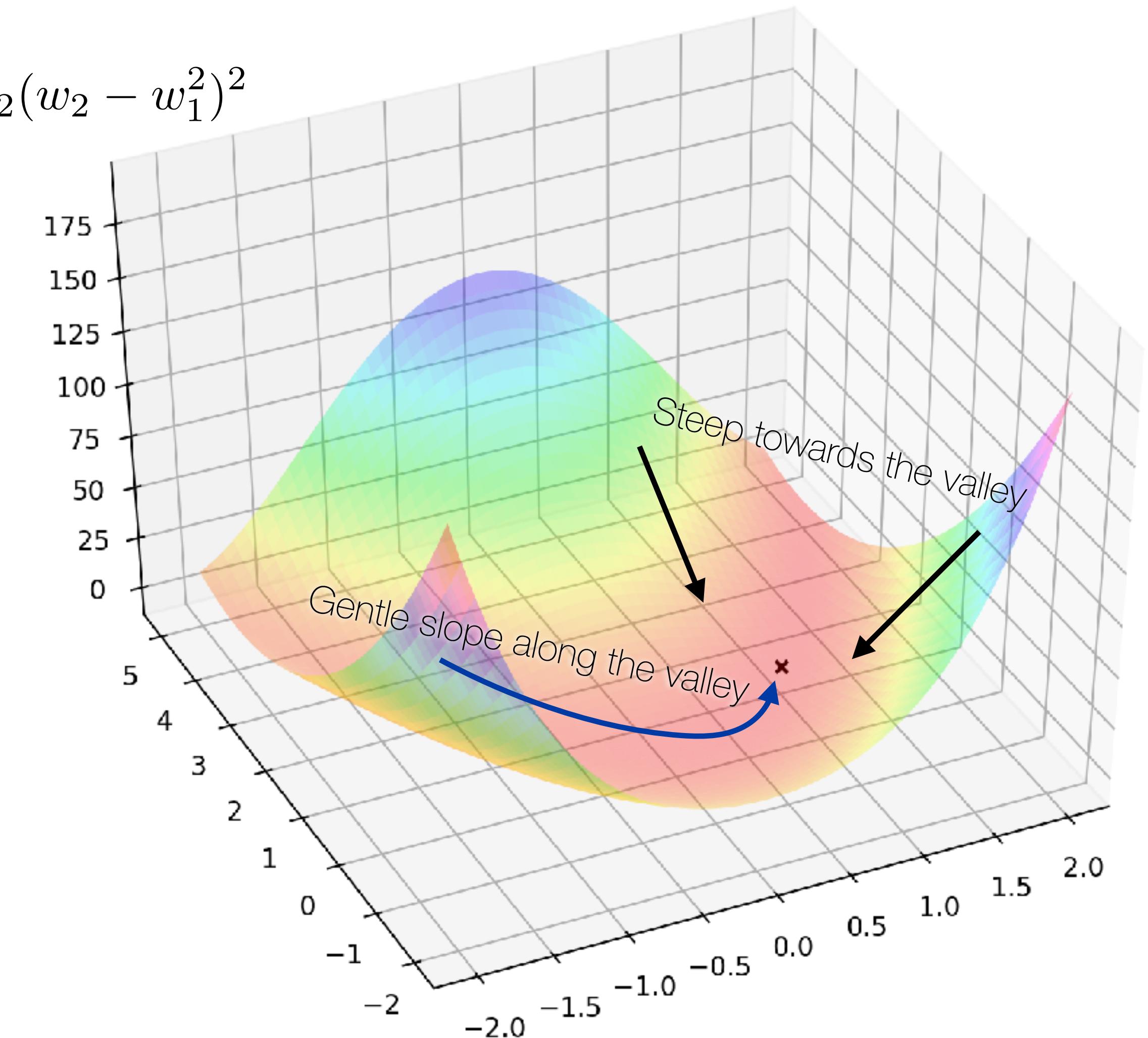
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Gradient Descent

-Let's see how bad it can be

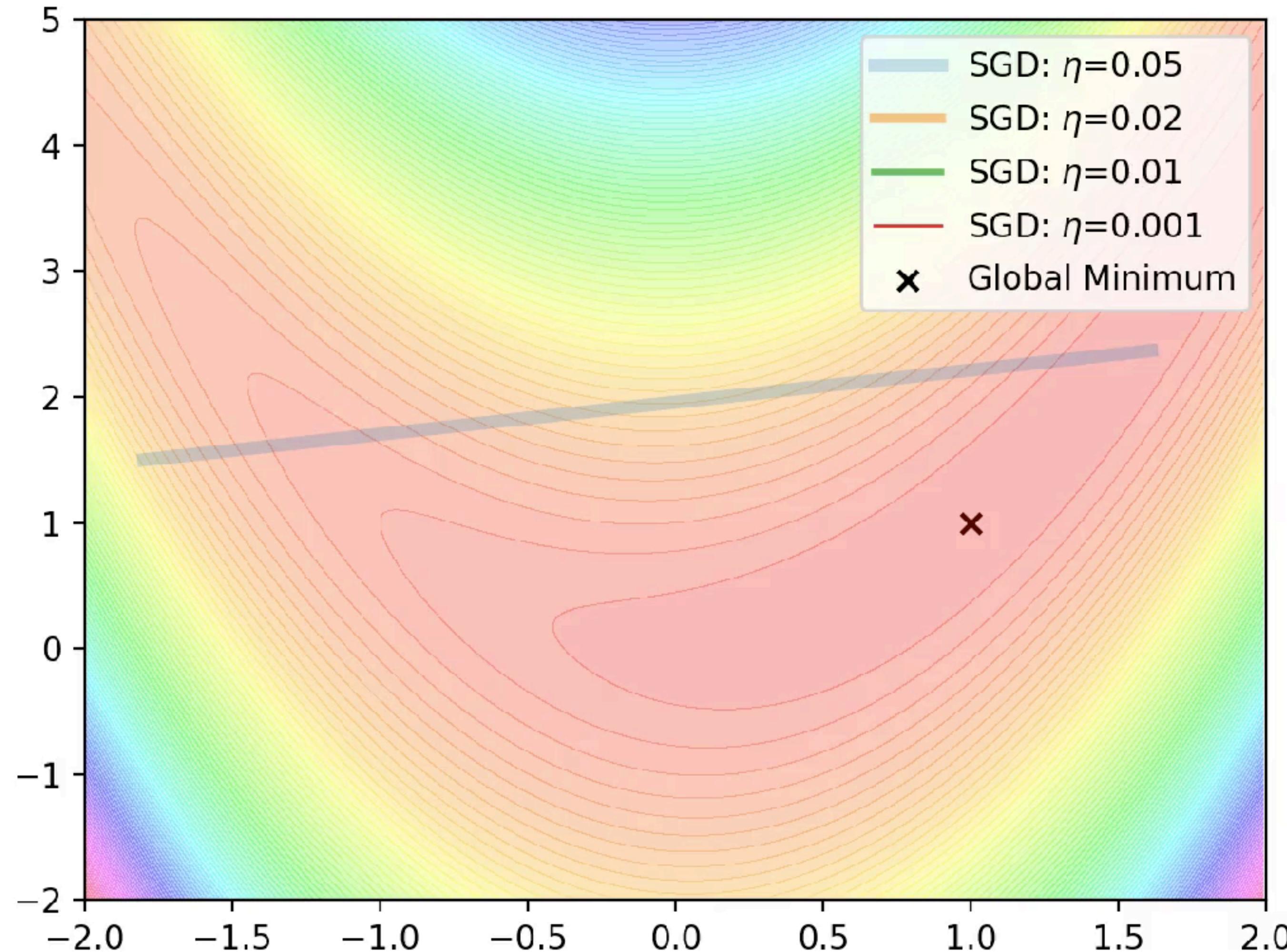
- The Rosenbrock's banana function $f(\mathbf{w}) = (x_1 - w_1)^2 + x_2(w_2 - w_1^2)^2$
 - x_1, x_2 : constants (input data)
 - w_1, w_2 : parameters to update
 - $f(\mathbf{w})$: the objective function
- The Rosenbrock function is a non-convex function, but it has a global minimum at
 $\mathbf{w}^* = [1, 1]^\top \quad f(\mathbf{w}^*) = 0$
 - It's just difficult to find it
- Gradient descent upadates
$$\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla_{\mathbf{w}} f(\mathbf{x}; \mathbf{w})$$
- This function is weird
 - It's something like a steep valley
 - where water runs very slow



Gradient Descent

-Let's see how bad it can be

- Choosing the best learning rate is difficult
- The best learning rate is not best everywhere



Momentum

-Let's give it some speed

- Velocity

$$\mathbf{v} \leftarrow \gamma \mathbf{v} - \eta \nabla_{\mathbf{w}} f(\mathbf{x}; \mathbf{w})$$

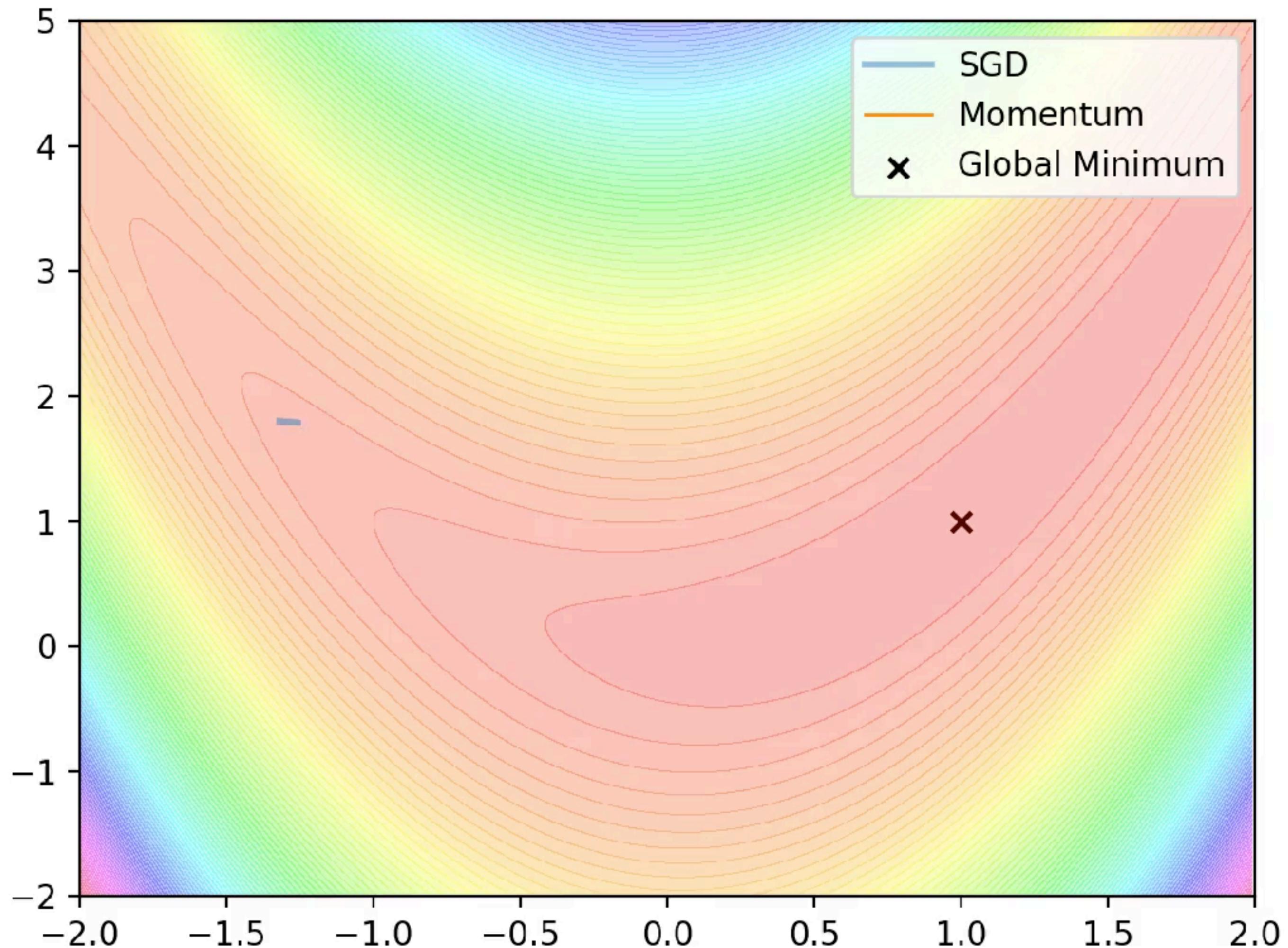
- Initially a zero vector,
but eventually accumulates all the SGD directions

- Update

$$\mathbf{w} \leftarrow \mathbf{w} + \mathbf{v}$$

- Meaning?

- If the update has been done
with a large velocity,
sudden changes in the SGD directions
don't affect the update too much
- Because usually γ is a large number, e.g. 0.9
- See the video for a comparison with SGD



INDIANA UNIVERSITY

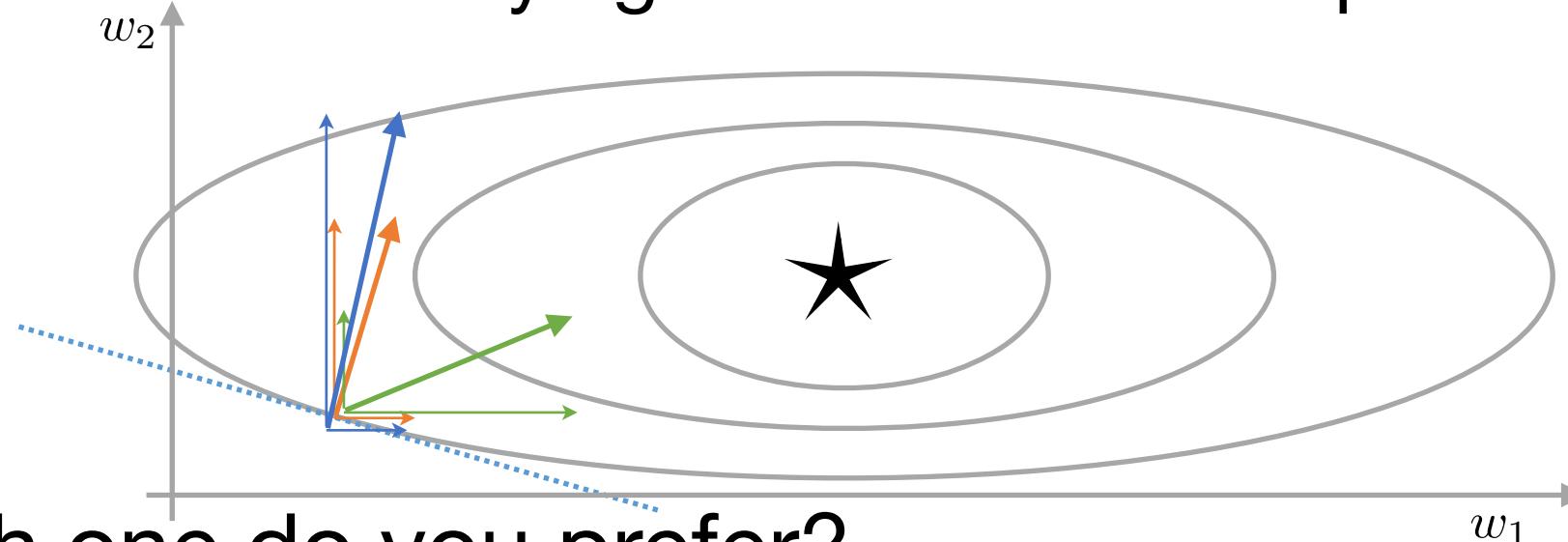
SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Sutskever, Ilya, et al. "On the importance of initialization and momentum in deep learning." International conference on machine learning. 2013.

AdaGrad

-Per-parameter adaptive learning rate

- It's to gradually decay the learning rate
 - But different decaying rates for different parameters



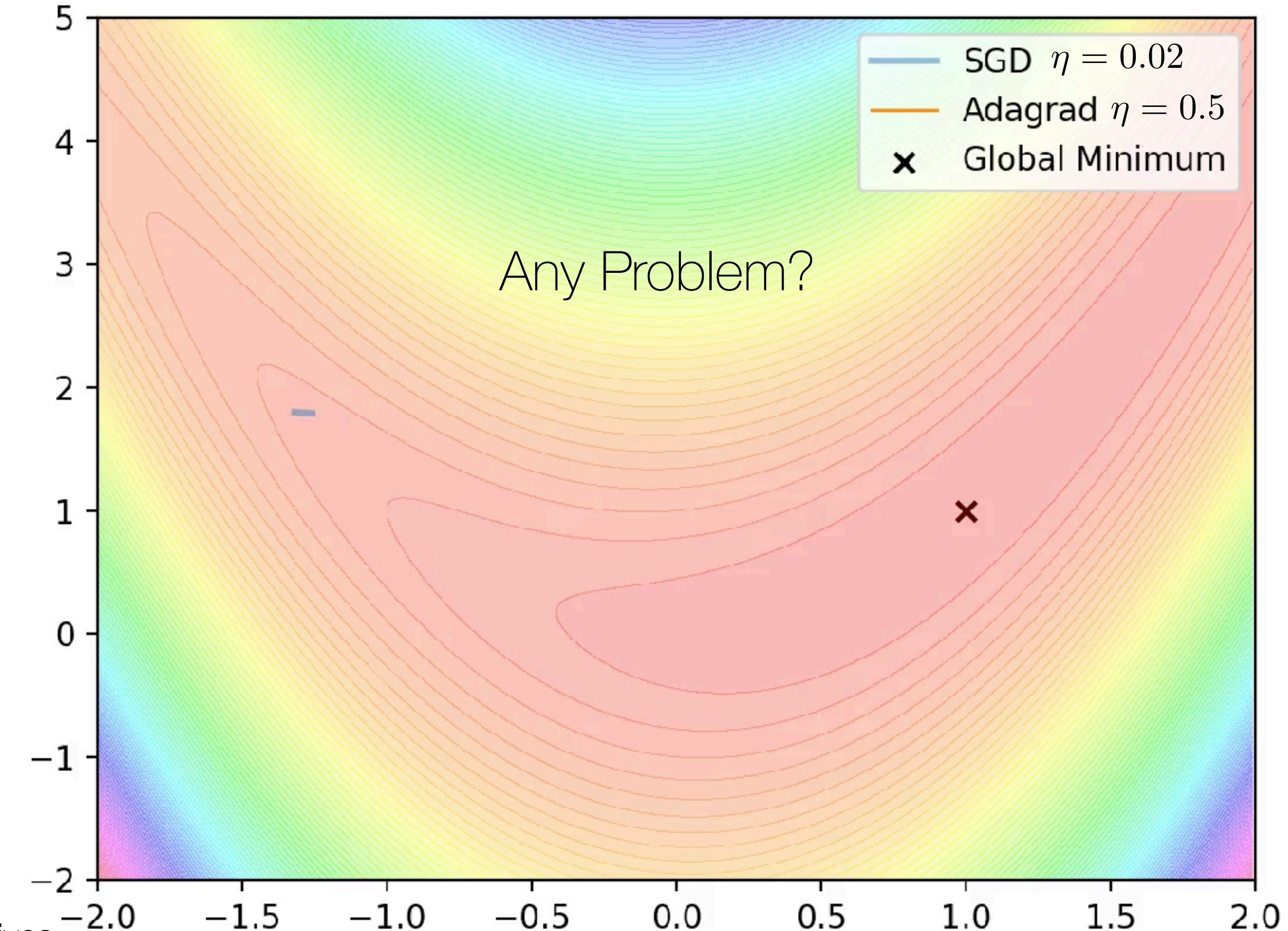
- Which one do you prefer?

- Orange arrow: Original GD direction
 $-\nabla_{\mathbf{w}} = -\nabla_{w_1} - \nabla_{w_2}$
- Green arrow: Larger LR to the gentle slope
 $-\nabla_{\mathbf{w}} = -\eta_1 \nabla_{w_1} - \eta_2 \nabla_{w_2}, \quad \eta_1 > \eta_2$
- Blue arrow: Larger LR to the steeper slope
 $-\nabla_{\mathbf{w}} = -\eta_1 \nabla_{w_1} - \eta_2 \nabla_{w_2}, \quad \eta_1 < \eta_2$

- Update

- Gradient $\mathbf{g} \leftarrow \eta \nabla_{\mathbf{w}} f(\mathbf{x}; \mathbf{w})$
- Accumulated squared gradient $\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g}$
- Update $\mathbf{w} \leftarrow \mathbf{w} - \frac{\eta}{\epsilon + \sqrt{\mathbf{r}}} \odot \mathbf{g}$

Large partial derivatives
denote LR



Any Problem?



INDIANA UNIVERSITY

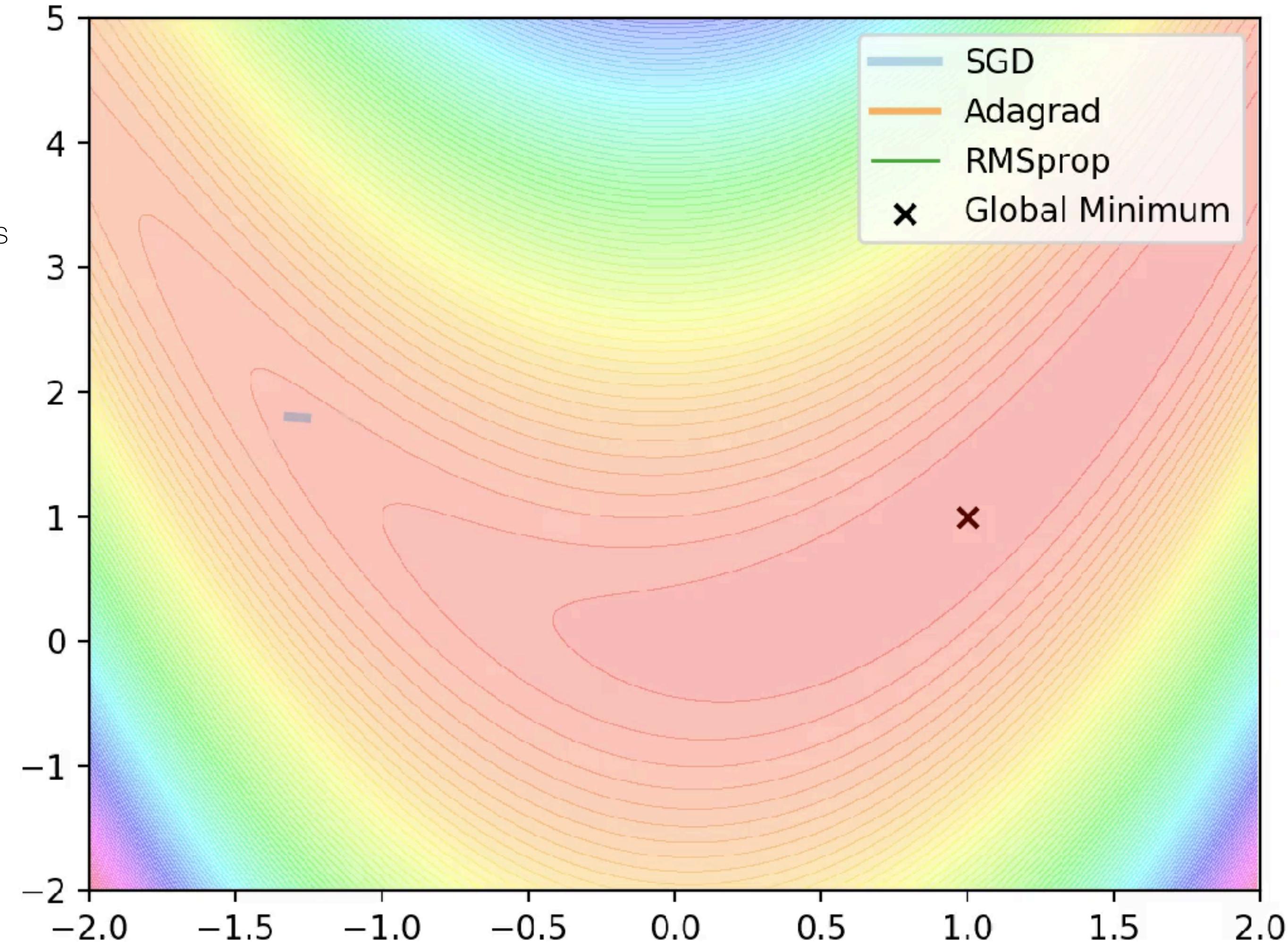
SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Duchi, John, Elad Hazan, and Yoram Singer. "Adaptive subgradient methods for online learning and stochastic optimization." Journal of Machine Learning Research 12.Jul (2011): 2121-2159.

RMSprop

-To improve AdaGrad

- A potential problem with AdaGrad
 - Gradient $\mathbf{g} \leftarrow \eta \nabla_{\mathbf{w}} f(\mathbf{x}; \mathbf{w})$
 - Accumulated squared gradient $\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g}$
 - Update $\mathbf{w} \leftarrow \mathbf{w} - \frac{\eta}{\epsilon + \sqrt{\mathbf{r}}} \odot \mathbf{g}$
 - Large partial derivatives denote LR
 - \mathbf{r} keeps getting large, and the LR can vanish
- RMSprop
 - $\mathbf{r} \leftarrow \gamma \mathbf{r} + (1 - \gamma) \mathbf{g} \odot \mathbf{g}$
 - Meaning?
 - The contribution of the first \mathbf{r} exponentially decays over iterations:
 - At i -th iteration $\gamma^i \mathbf{r}^{(0)}$
 - $\mathbf{g} \odot \mathbf{g}$ in the first iteration
 - Prevents \mathbf{r} from becoming too large



INDIANA UNIVERSITY

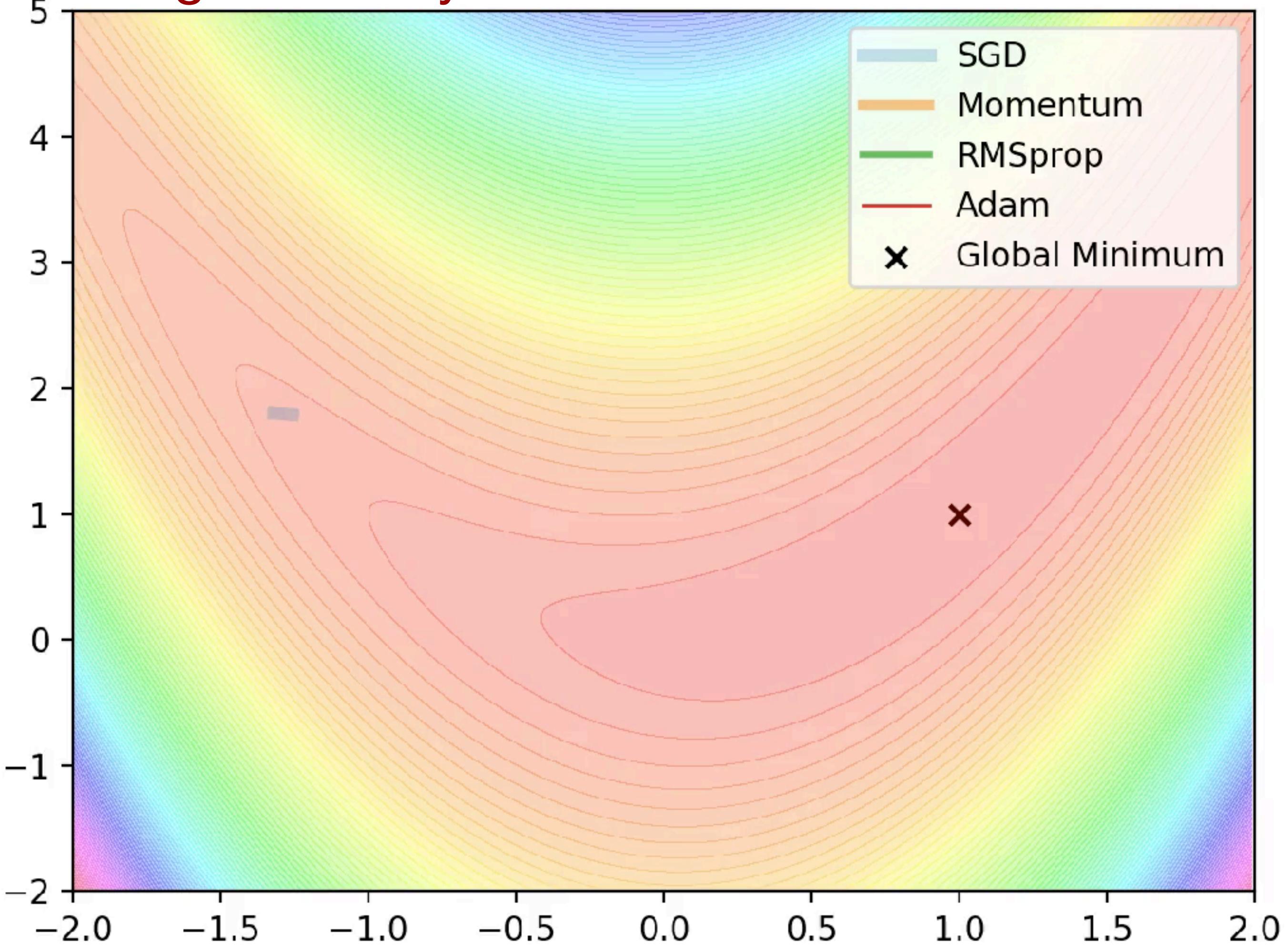
SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

G. Hinton, "Neural Networks for machine learning," Coursera Lecture (2012)

Adam

-Perhaps the most popular choice in deep learning these days

- Combination of RMSprop and momentum
- Momentum (review)
 - Velocity $\mathbf{v} \leftarrow \gamma\mathbf{v} - \eta \nabla_{\mathbf{w}} f(\mathbf{x}; \mathbf{w})$
 - Update $\mathbf{w} \leftarrow \mathbf{w} + \mathbf{v}$
- RMSprop (review)
 - Gradient $\mathbf{g} \leftarrow \eta \nabla_{\mathbf{w}} f(\mathbf{x}; \mathbf{w})$
 - Exponentially accumulated squared gradient
$$\mathbf{r} \leftarrow \gamma\mathbf{r} + (1 - \gamma)\mathbf{g} \odot \mathbf{g}$$
 - Update $\mathbf{w} \leftarrow \mathbf{w} - \frac{\eta}{\epsilon + \sqrt{\mathbf{r}}} \odot \mathbf{g}$
- Adam
 - Gradient $\mathbf{g} \leftarrow \eta \nabla_{\mathbf{w}} f(\mathbf{x}; \mathbf{w})$
 - Velocity $\mathbf{v} \leftarrow \beta_1 \mathbf{v} + (1 - \beta_1) \mathbf{g}$
 - Exponentially accumulated squared gradient $\hat{\mathbf{v}} \leftarrow \frac{\mathbf{v}}{1 - \beta_1^i}$
$$\mathbf{r} \leftarrow \beta_2 \mathbf{r} + (1 - \beta_2) \mathbf{g} \odot \mathbf{g}$$
$$\hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}}{1 - \beta_2^i}$$
 - Update $\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{\hat{\mathbf{v}}}{\sqrt{\hat{\mathbf{r}}} + \epsilon}$



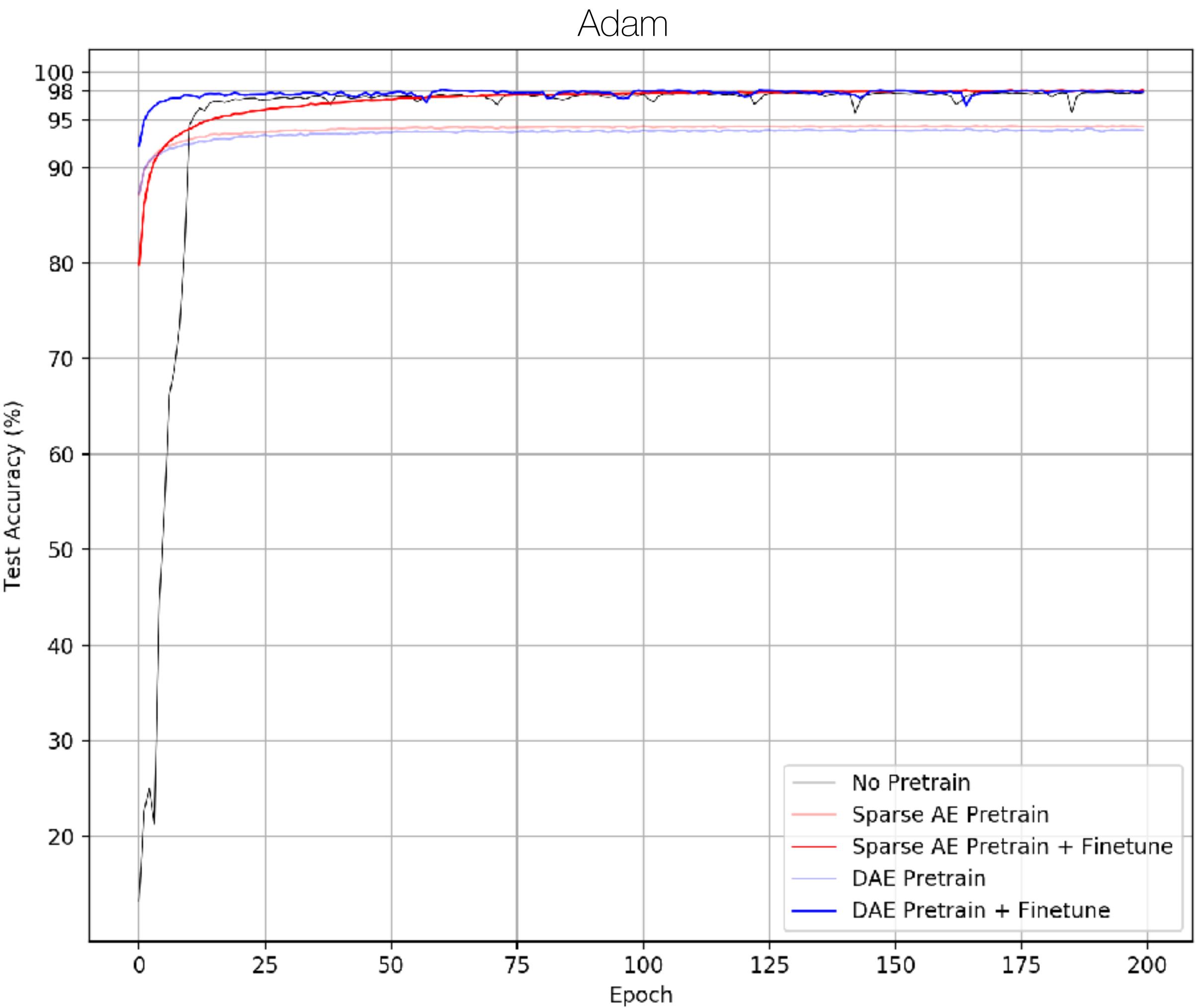
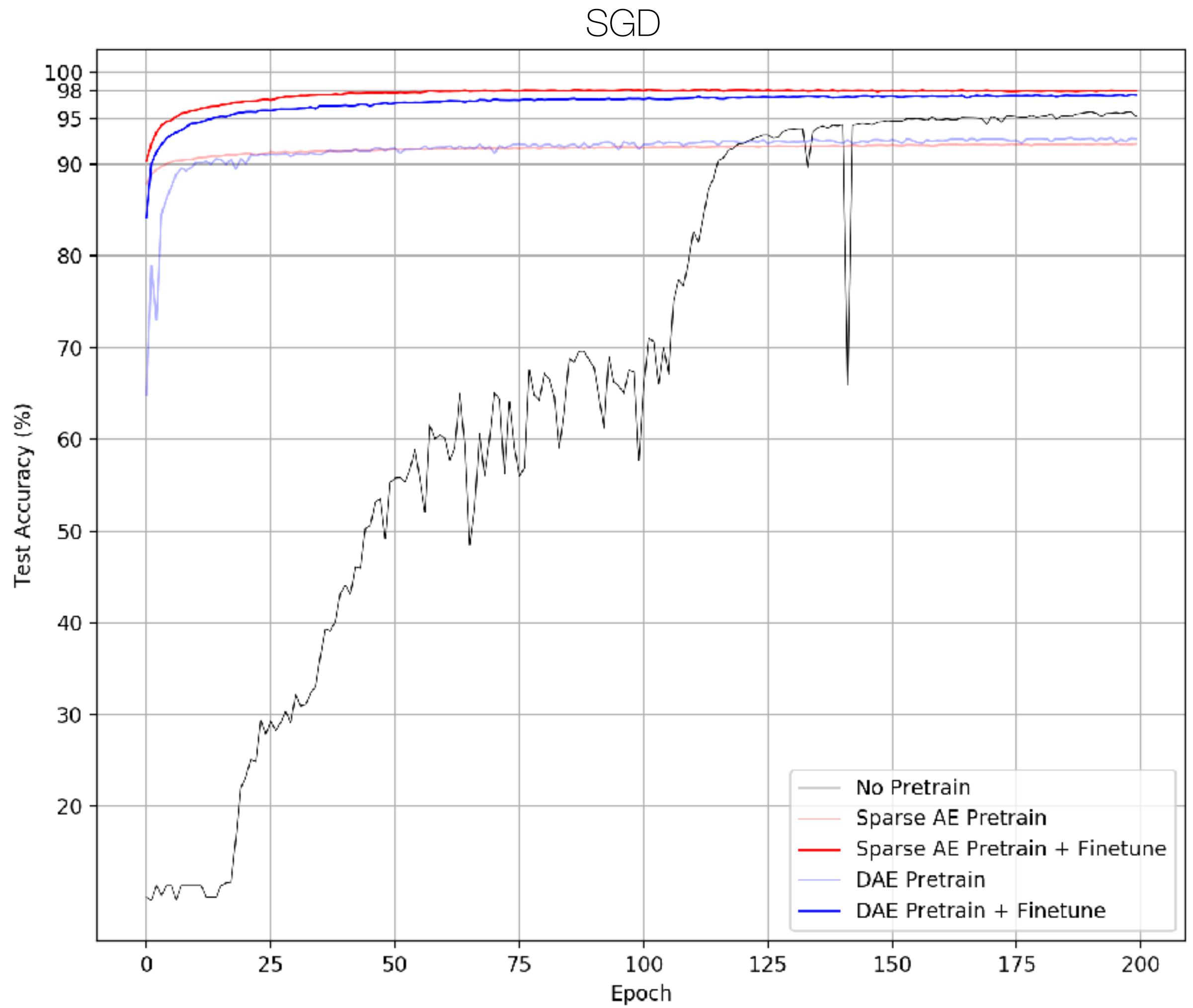
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Kingma, Diederik P., and Jimmy Ba. "Adam: A method for stochastic optimization." arXiv preprint arXiv:1412.6980 (2014).

Adam

-On pretrained networks



INDIANA UNIVERSITY

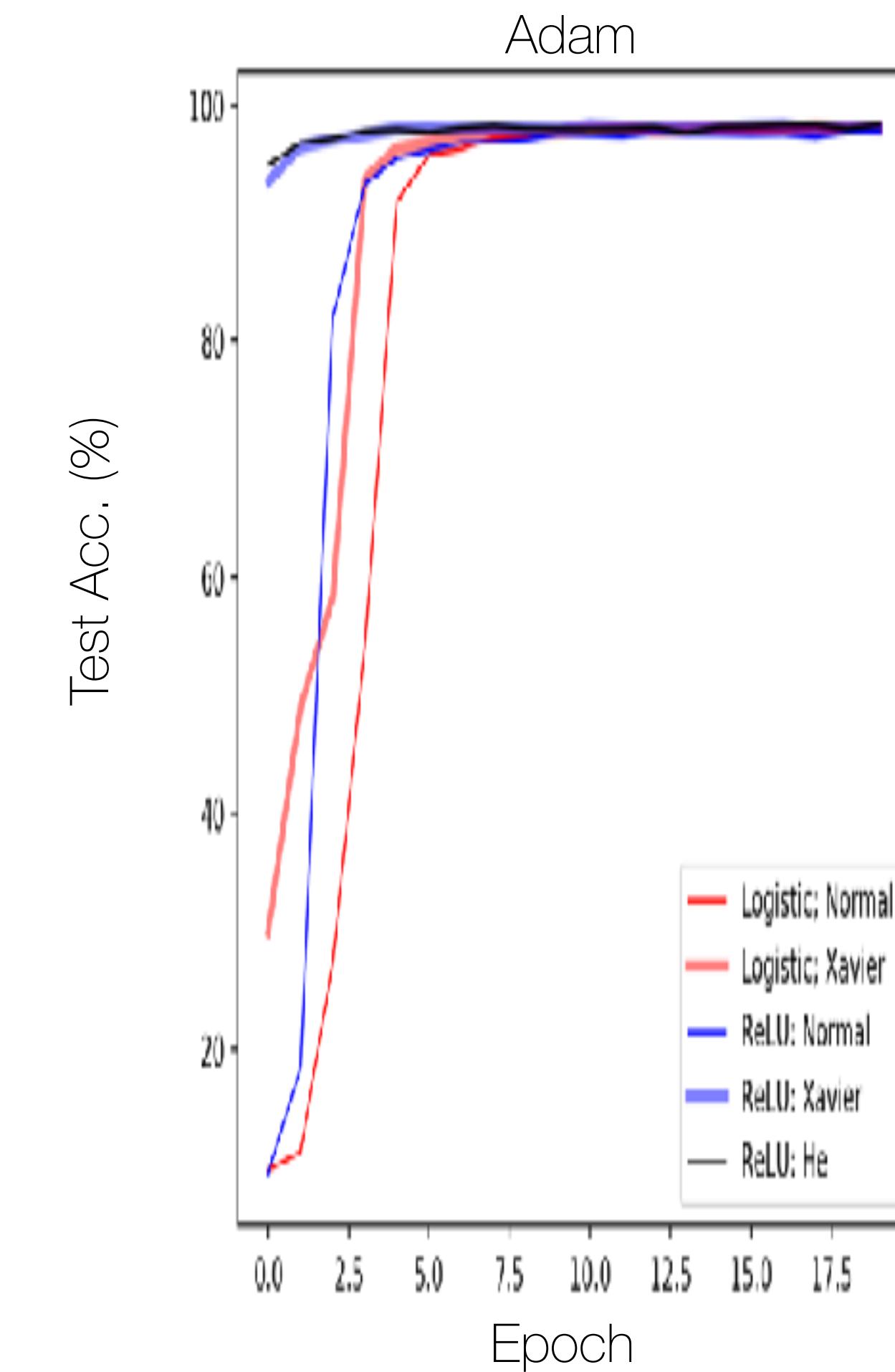
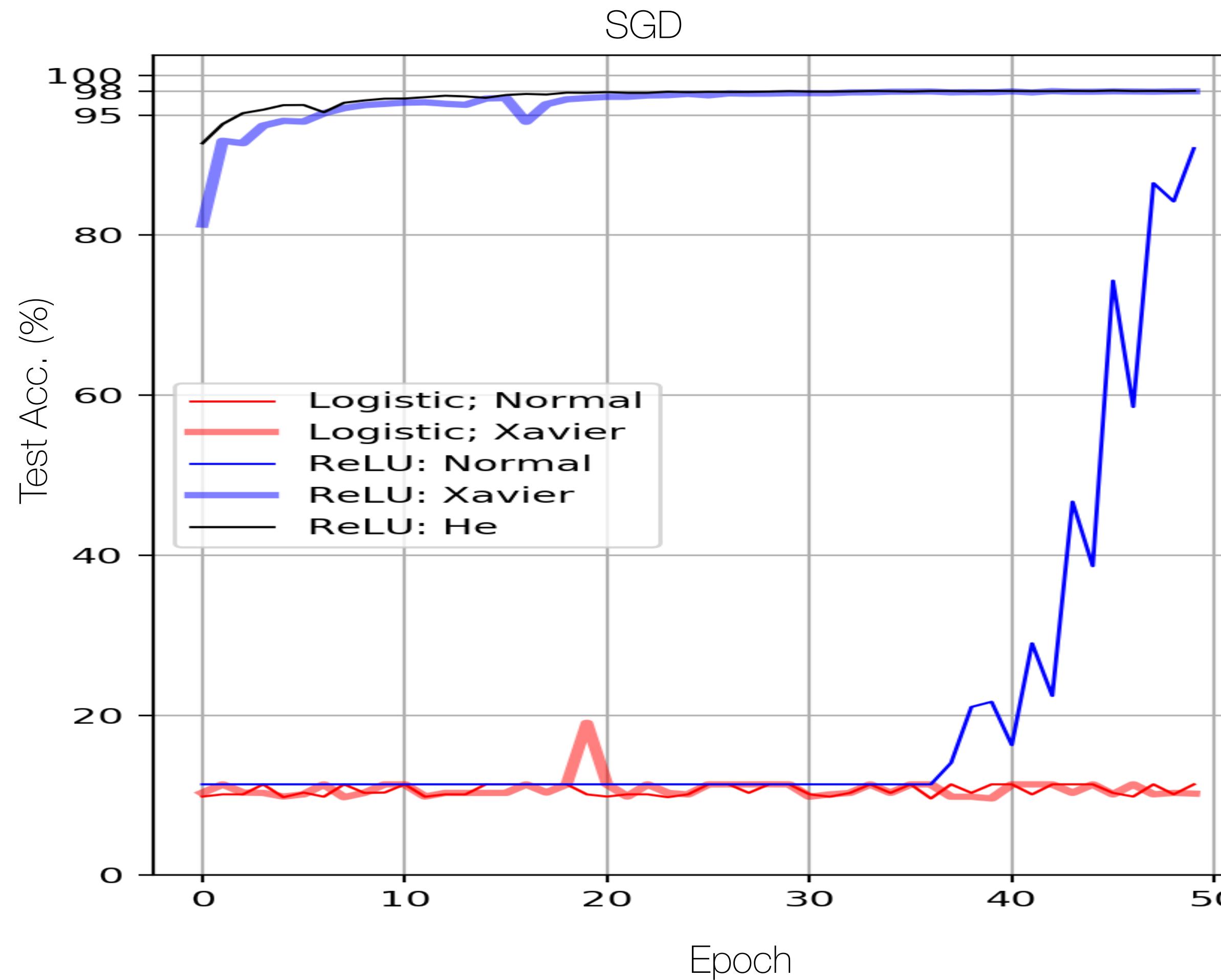
SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Kingma, Diederik P., and Jimmy Ba. "Adam: A method for stochastic optimization." arXiv preprint arXiv:1412.6980 (2014).

Adam

-On deeper networks with no pretraining

- For a 512X5 network for MNIST



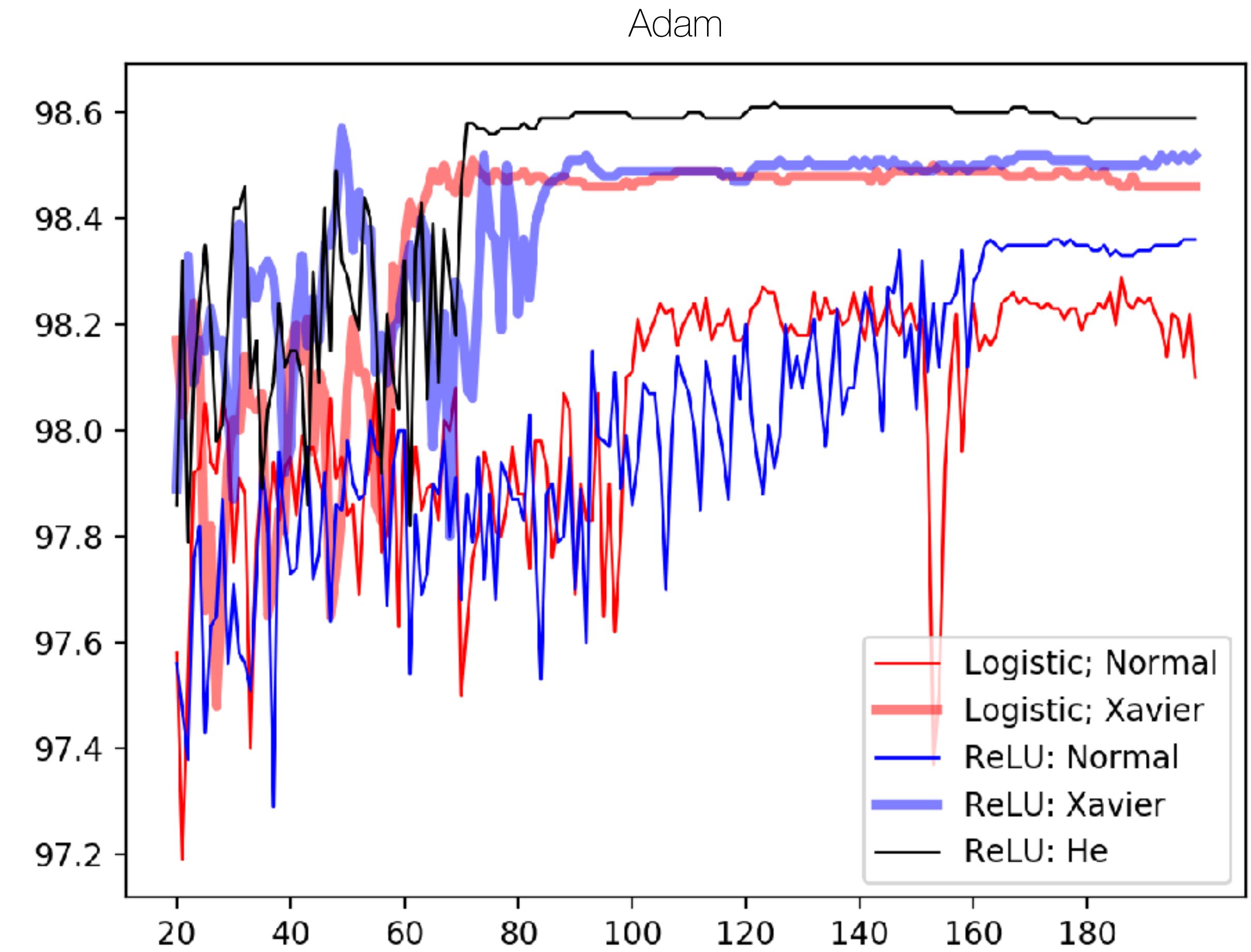
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Kingma, Diederik P., and Jimmy Ba. "Adam: A method for stochastic optimization." arXiv preprint arXiv:1412.6980 (2014).

Adam

-On deeper networks with no pretraining



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Kingma, Diederik P., and Jimmy Ba. "Adam: A method for stochastic optimization." arXiv preprint arXiv:1412.6980 (2014).

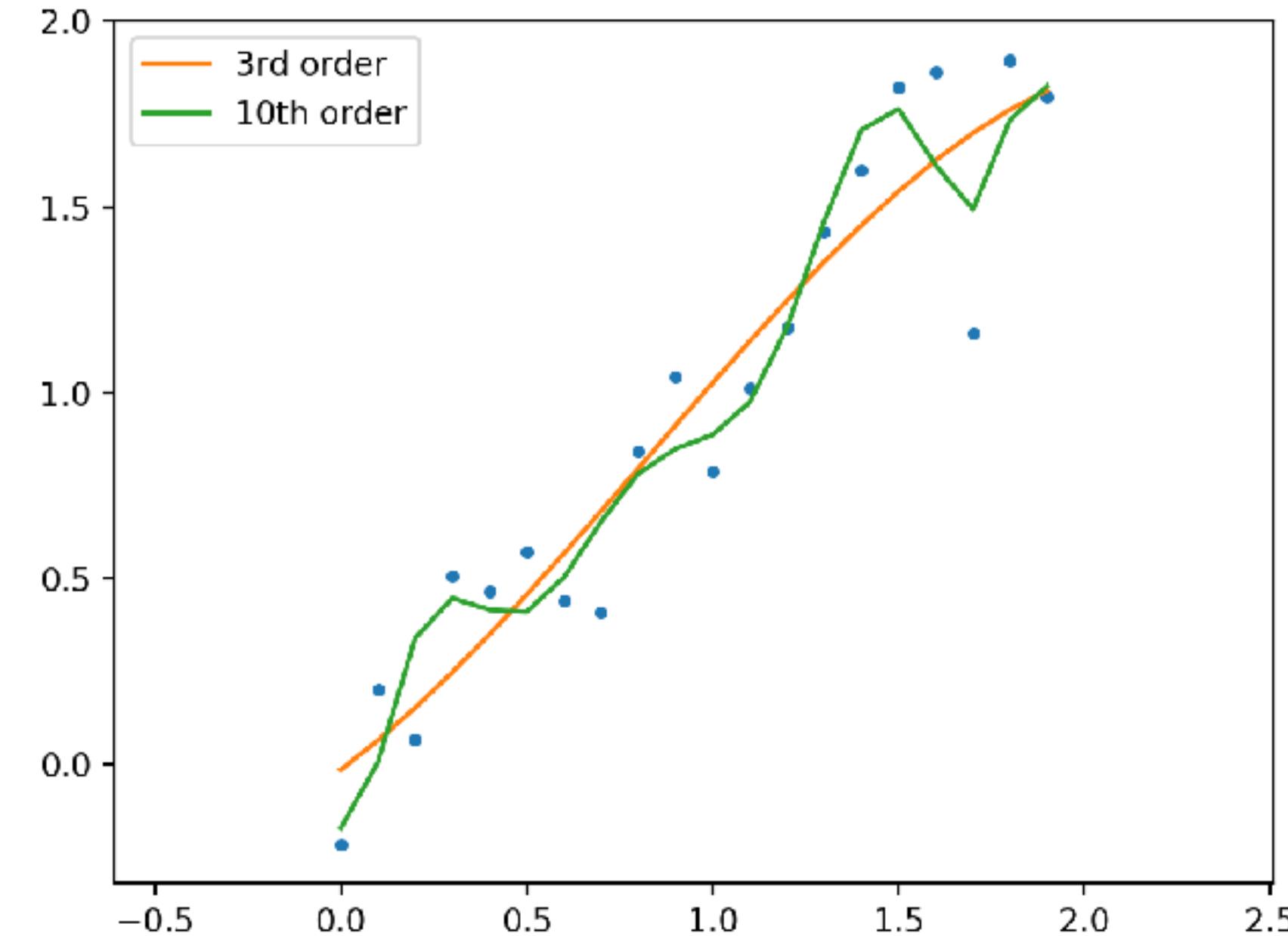
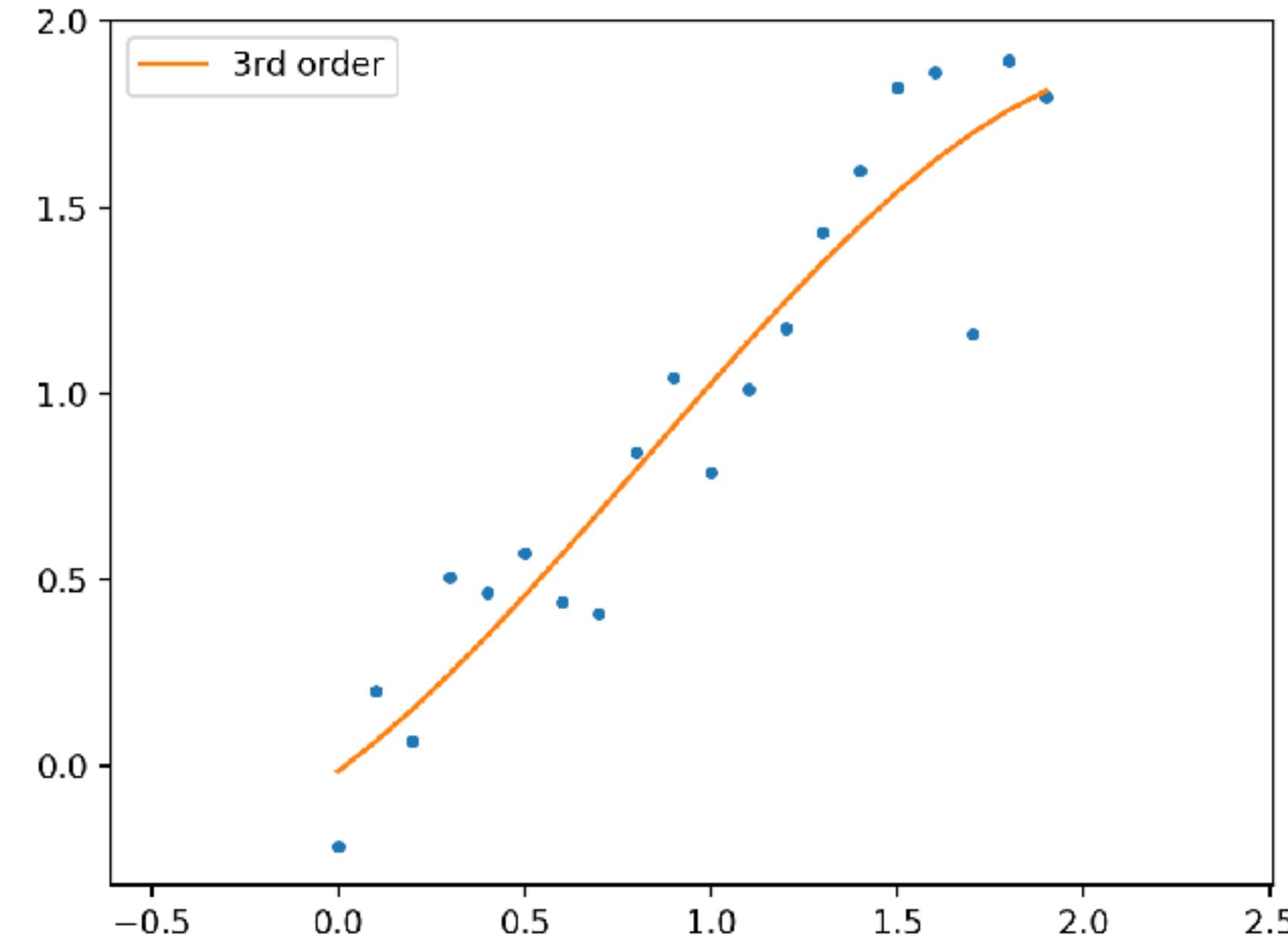
Regularization, Dropout, and Validation

Once again, they are related.

Weight Decaying

- You don't want your weights to be too large

- If your model is reasonably complex (e.g. 3rd order polynomial for a linear dataset)
 - The model tries to fit the data samples by coming up with small weights $[-0.20, 0.50, 0.74, -0.02]$
- If the model is too complex with too many free parameters, then overfitting can happen
 - As a result, the parameters tend to become large $[-25.01, 190.83, -560.18, 727.14, -185.85, -598.32, 741.70, -359.62, 72.98, -2.61, -0.17]$



- L2 regularization $\arg \min_{\mathbb{W}} \sum_t \mathcal{D}(\mathbf{Y}_{:,t} || \mathcal{G}(\mathbf{X}_{:,t}; \mathbb{W})) + \lambda \sum_{w \in \mathbb{W}} w^2$



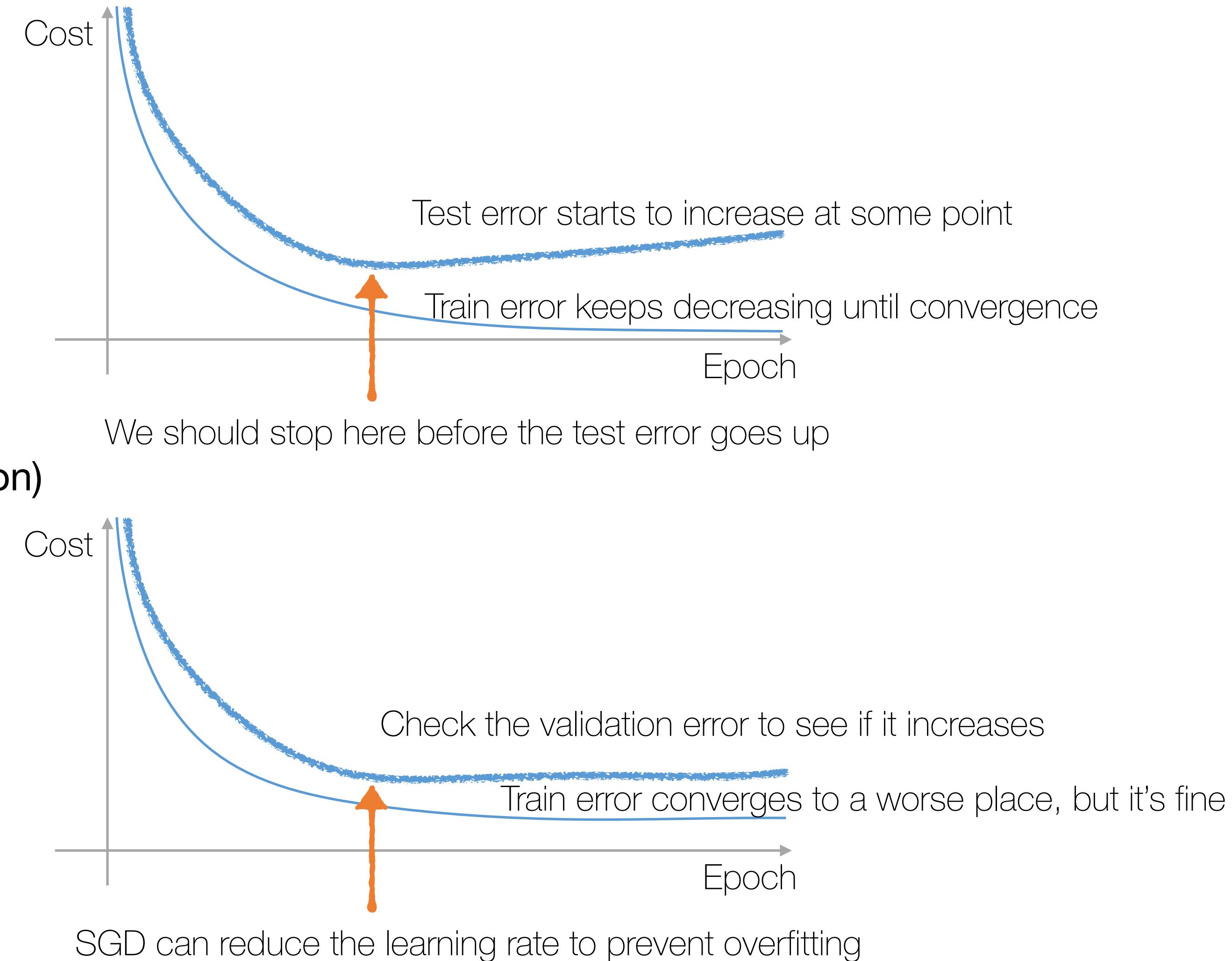
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Early Stopping

-Use validation error to prevent overfitting

- So, what happens if overfitting happens?
- What can we do without using the test set?
 - We can set aside some parts of the training set as the validation set
 - Check to see if the validation error increases
- What can we do if the validation error starts to increase?
 - In SGD, we can reduce the learning rate to stay at the local minimum (of the validation error function)
- An indirect way to predict the performance during the test time
- Disadvantage: you can't fully make use of the entire training set
- Cross validation can be also used to find out the optimal hyperparameters
 - e.g. initial LRs, network structures, etc.



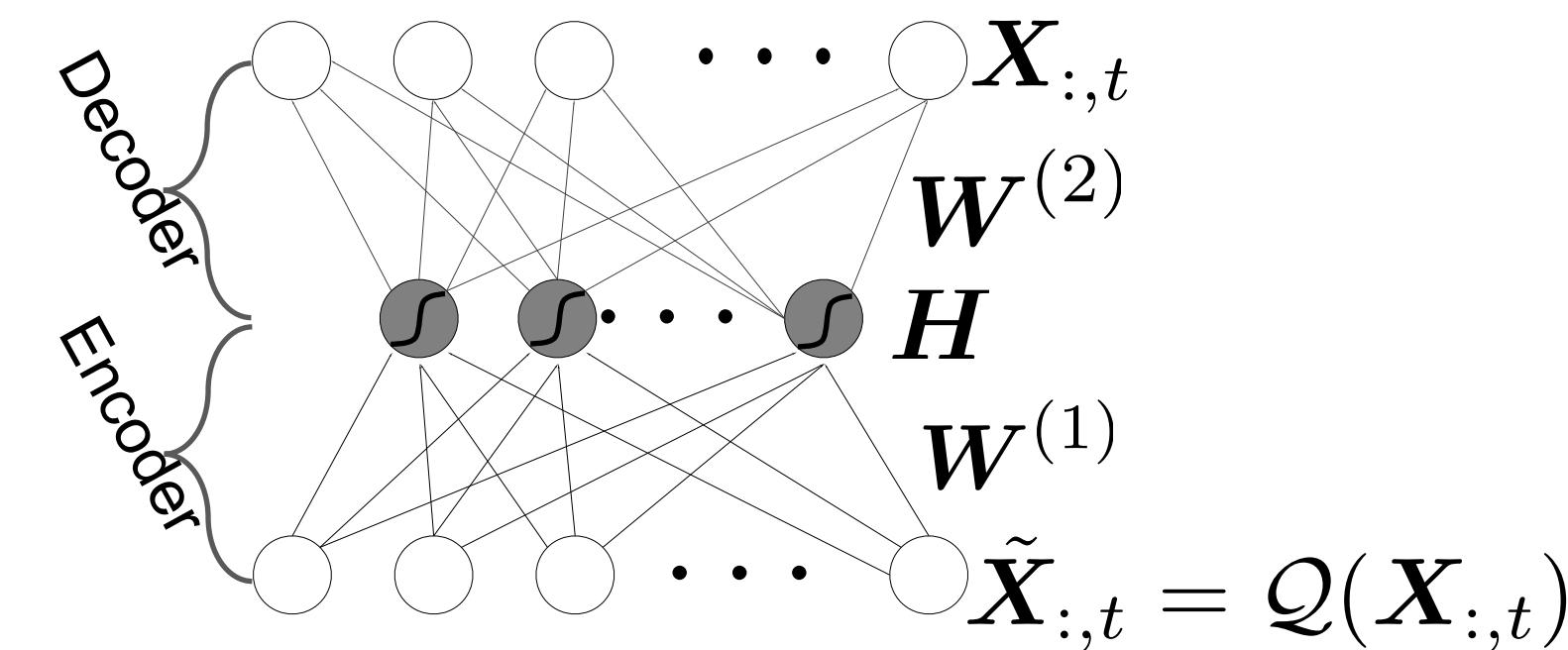
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Dropout

-Structural noise injection

- Denoising AutoEncoders
 - Contaminate the input with noise $\tilde{\mathbf{X}}_{:,t} = \mathcal{Q}(\mathbf{X}_{:,t})$
 - Train the autoencoder to predict the clean input
- e.g. Masking noise $\tilde{\mathbf{X}}_{f,t} = b\mathbf{X}_{f,t}$, $b \sim \text{Bern}(p)$
 - At each iteration you start from a corrupted version of the dataset
 - Still has to do your best to minimize the error
 - This can be seen as noise injection to prevent overfitting, too



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Vincent, Pascal, et al. "Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion." Journal of Machine Learning Research 11.Dec (2010): 3371-3408.

Dropout

-Structural noise injection

- For each data sample, and its feedforward

- Dropout randomly turns off the input units to the layer

$$Z_{i,t}^{(l)} \leftarrow \sum_j W_{i,j}^{(l)} M_{j,t} X_{j,t}^{(l)} + b_i^{(l)} \quad M_{j,t} \sim \text{Bern}(p)$$

- So, the masking noise is injected to every layer

- Kinda similar to stacked DAE
 - But, this is done during the full network training, not pretraining
 - Network is robust to the masking noise

- Network is robust to the structural changes

- Every feedforward done in a thinned-out network

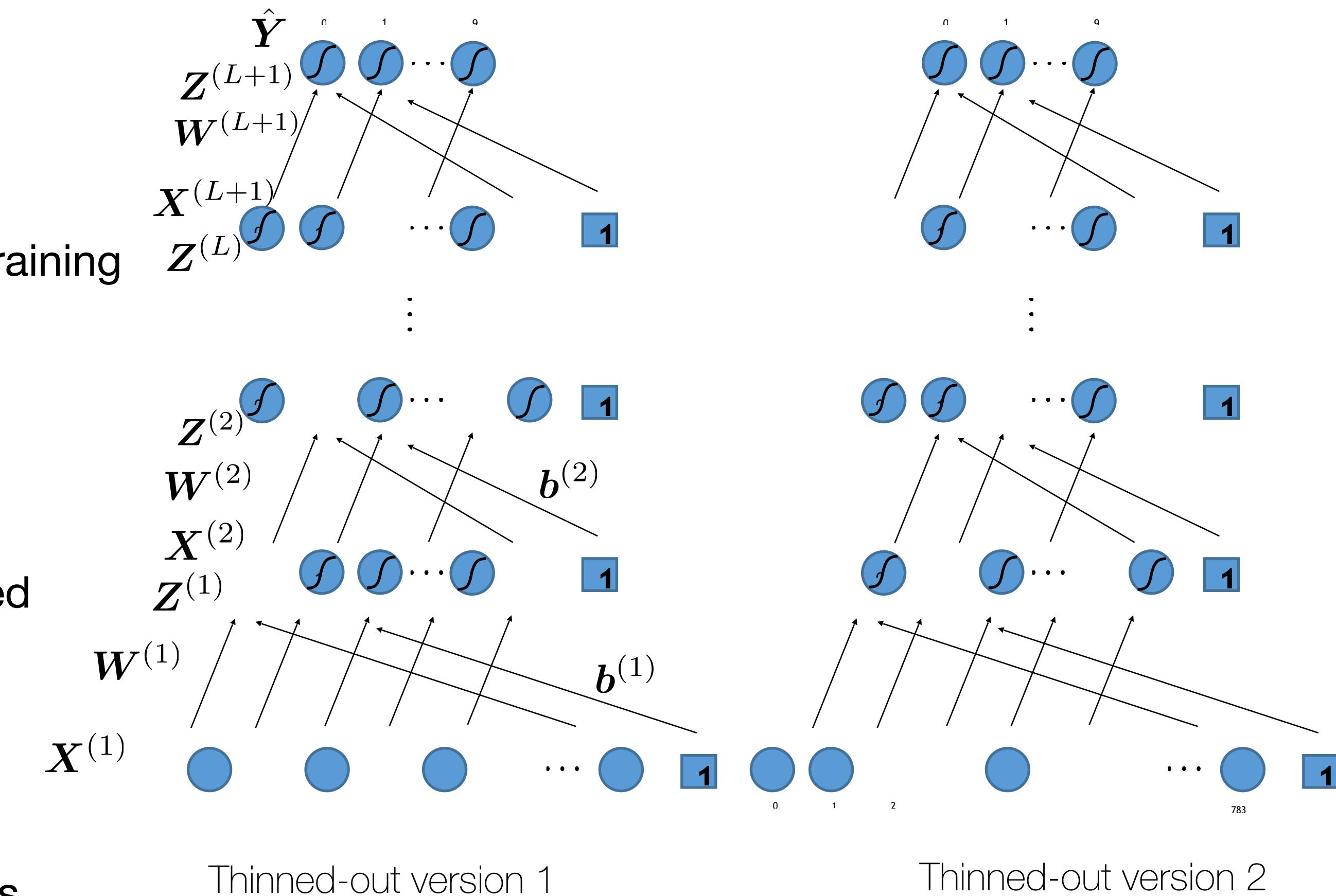
- Ensemble

- Each data sample predicts the output using a subsampled version of the network
 - Can be seen as an ensemble technique

- Testing time

$$Z_{i,t}^{(l)} \leftarrow \sum_j p W_{i,j}^{(l)} X_{j,t}^{(l)} + b_i^{(l)}$$

- Because the weights are trained to handle less input units



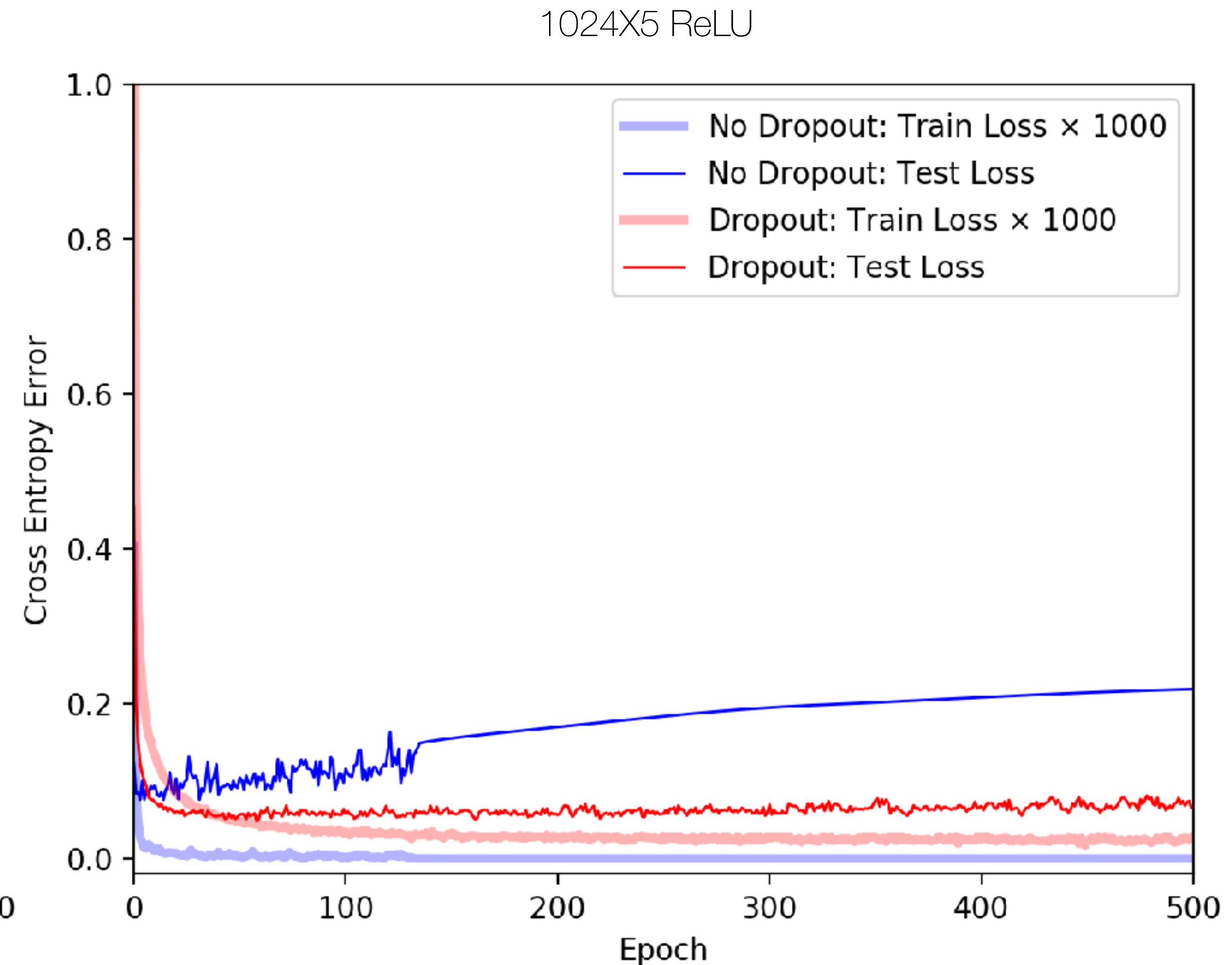
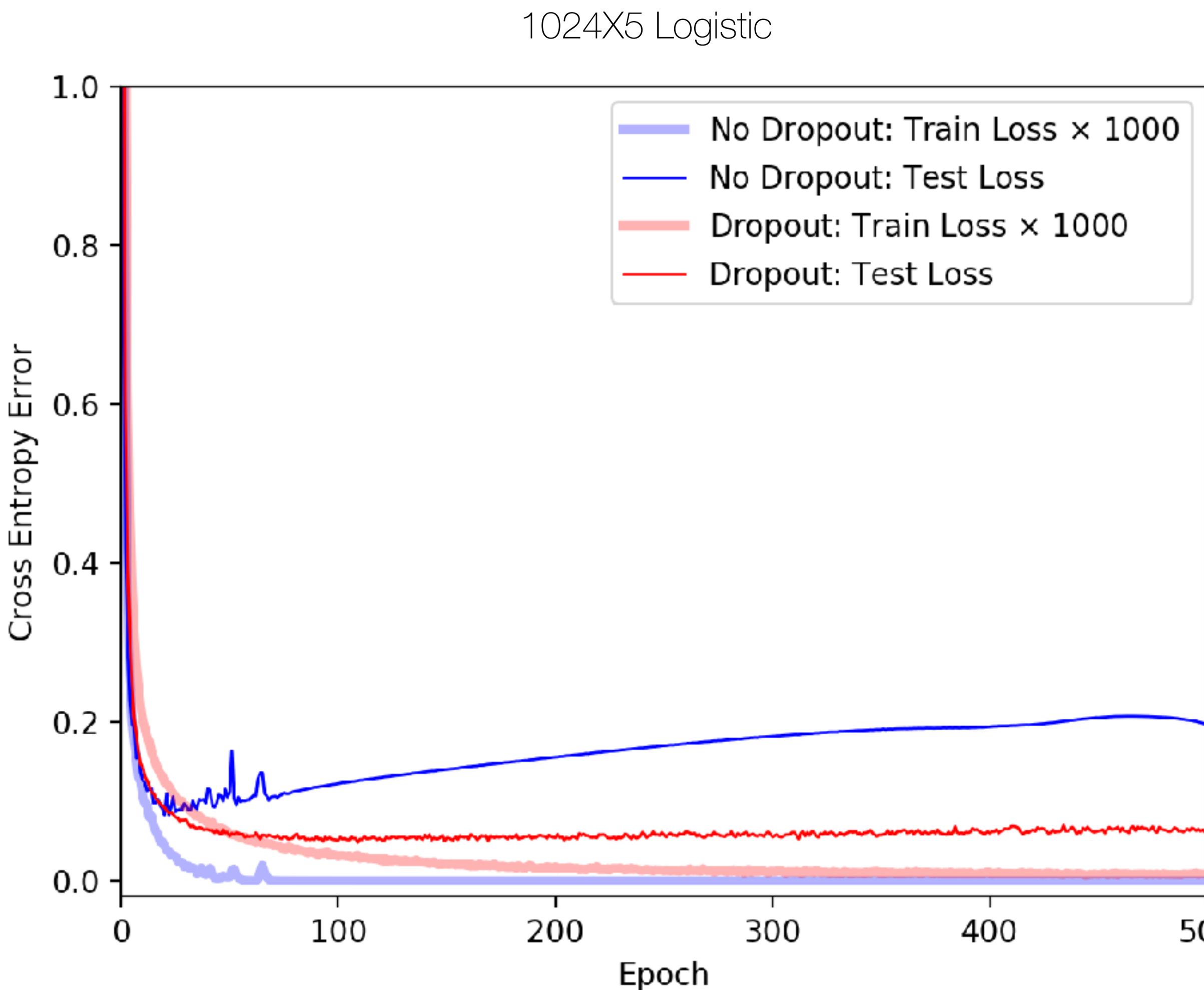
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Srivastava, Nitish, et al. "Dropout: A simple way to prevent neural networks from overfitting." The Journal of Machine Learning Research 15.1 (2014): 1929-1958.

Dropout

-Structural noise injection



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Srivastava, Nitish, et al. "Dropout: A simple way to prevent neural networks from overfitting." The Journal of Machine Learning Research 15.1 (2014): 1929-1958.

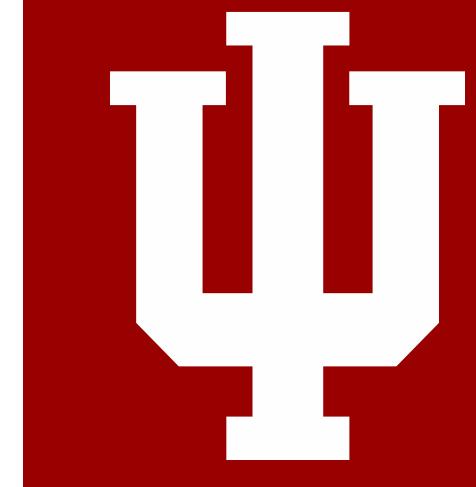
Reading

- All the paper cited
- Textbook Chapter 7, 8



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING



Thank You!