



# XNOR-Nets: Binarized Neural Networks

---

*Sai Krishna Aditya Parvatha*

*Piyush Bhatt*

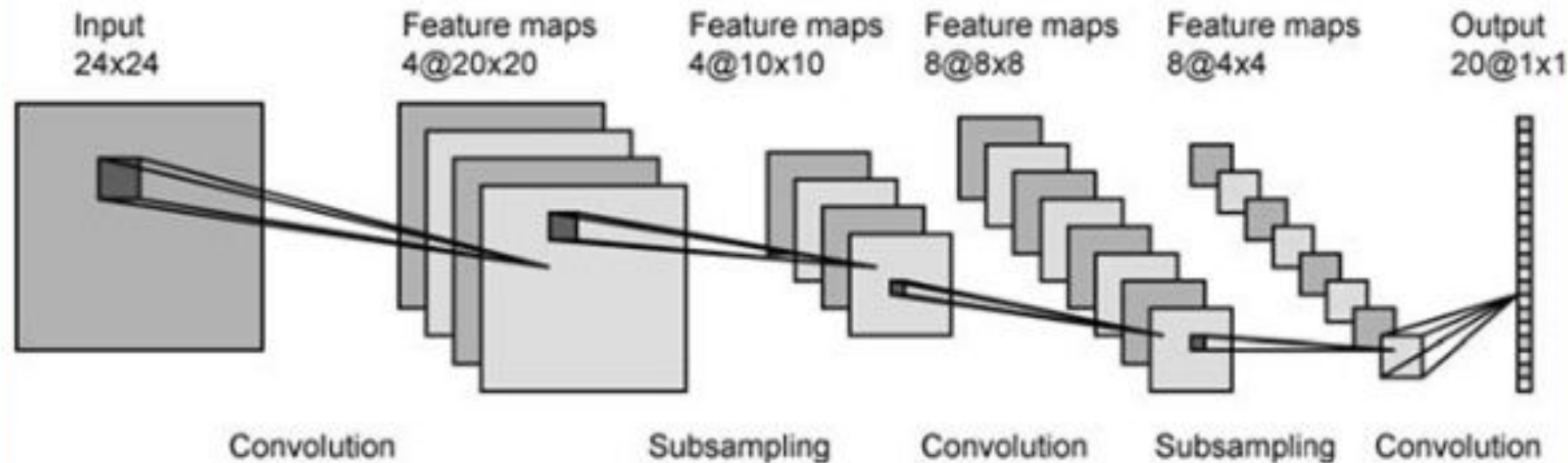
*Sameer Kumar Behera*

*Tanmay Verma*

# Convolutional Neural Network



TEXAS A&M  
UNIVERSITY.



# Convolution Operations



TEXAS A&M  
UNIVERSITY

Model	LeNet-5	AlexNet	VGG-16	GoogleNet v1	ResNet-5
Weights	60K	61M	138M	7M	25.5M
MAC	341K	724M	15.5G	1.43G	3.9G

**Source:** V. Sze, Y. Chen, T. Yang, and J. S. Emer. Efficient processing of deep neural networks: A tutorial and survey. Proceedings of the IEEE. 105(12):2295–2329, Dec 2017

---



# **INTRODUCTION TO BINARIZED NEURAL NETWORKS**

# Binarized Convolutional Neural Network



TEXAS A&M  
UNIVERSITY

- Weights are constrained to -1 (represented by bit value 0) or +1 (represented by bit value 1)
  - Pre-activations are also constrained to -1 and 1
  - Different popular attempts: Bitwise Neural Network, BinConnect, BinaryNet and XNOR-Net
-

# Matrix Multiplication vs XNOR



TEXAS A&M  
UNIVERSITY

$\begin{matrix} -1 & +1 & +1 \\ -1 & -1 & -1 \\ +1 & +1 & -1 \end{matrix}$	$\begin{matrix} -1 \\ +1 \\ +1 \end{matrix}$	$=$	$\begin{matrix} (-1 \times 1) + (1 \times 1) + (1 \times 1) \\ (-1 \times -1) + (1 \times -1) + (1 \times -1) \\ (-1 \times 1) + (1 \times 1) + (1 \times -1) \end{matrix}$	$=$	$\begin{matrix} 3 \\ -1 \\ -1 \end{matrix}$
$\begin{matrix} 0 & +1 & +1 \\ 0 & 0 & 0 \\ +1 & +1 & 0 \end{matrix}$	$\begin{matrix} 0 \\ +1 \\ +1 \end{matrix}$	$=$	$\begin{matrix} \text{popcnt}(\text{xnor}(011, 011)) \\ \text{popcnt}(\text{xnor}(011, 000)) \\ \text{popcnt}(\text{xnor}(011, 110)) \end{matrix}$	$=$	$\begin{matrix} 3 \\ -1 \\ -1 \end{matrix}$

$$\text{popcnt}(\text{xnor}(011, 110)) = \text{popcnt}(\text{xnor}(0,1), \text{xnor}(1,1), \text{xnor}(1,0)) = \text{popcnt}(010) = -1 + 1 - 1 = -1$$

- Small-sized model (32-bit weights replaced with a single-bit)
  - Fast inference
  - Power-efficient operations
-



# XNOR-Net



- First Binary Neural Network that scales up to **ImageNet**
- A simple and efficient way of training and running inference on Binary Convolution Layers.

$$\alpha^*, \mathbf{B}^*, \beta^*, \mathbf{H}^* = \underset{\alpha, \mathbf{B}, \beta, \mathbf{H}}{\operatorname{argmin}} \|\mathbf{X}^T \mathbf{W} - \beta \alpha \mathbf{H}^T \mathbf{B}\|$$

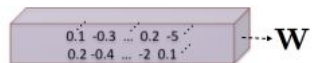
- **BOTTLENECK:** However, real-valued weights are still stored to calculate gradients and updating the weights.

# Convolution with XNOR-Net

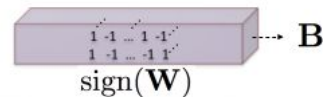


TEXAS A&M  
UNIVERSITY

## (1) Binarizing Weight

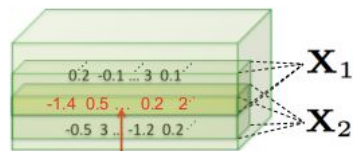


$$\frac{1}{n} \|W\|_{\ell_1} = \alpha$$



## (2) Binarizing Input

*Inefficient*

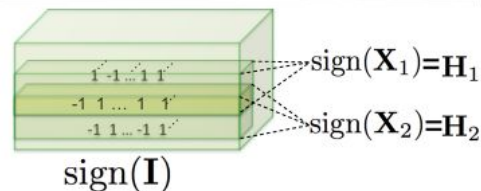


Redundant computations in overlapping areas

$$\frac{1}{n} \|X_1\|_{\ell_1} = \beta_1$$

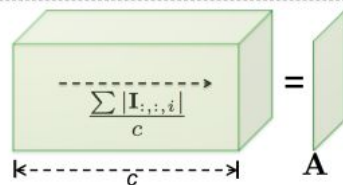
$$\frac{1}{n} \|X_2\|_{\ell_1} = \beta_2$$

$K$



## (3) Binarizing Input

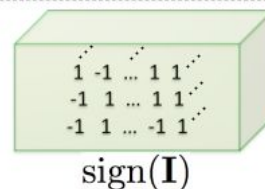
*Efficient*



$$A * K = \beta_1$$

$$A * K = \beta_2$$

$K$



## (4) Convolution with XNOR-Bitcount

$$I * W \approx \left[ \text{sign}(I) * \text{sign}(W) \right] \odot K \odot \alpha$$

$I$   $W$   $\text{sign}(I)$   $\text{sign}(W)$   $K$   $\alpha$

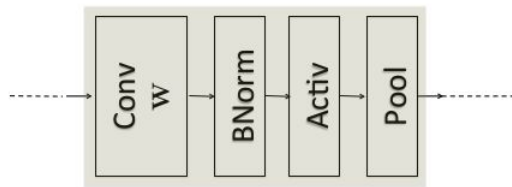
**Source:** Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet classification using binary convolutional neural networks. CoRR, abs/1603.05279, 2016.

# Important Guidelines to Preserve Accuracy

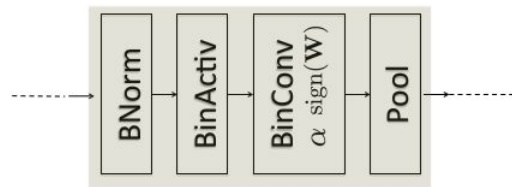


TEXAS A&M  
UNIVERSITY

- The first convolution layer and last layer in the network are not binarized.
- The sequence of operations is slightly modified from traditional CNN block



A typical block in CNN



A block in XNOR-Net

# Our Project Includes



TEXAS A&M  
UNIVERSITY

- Implementation of XNOR-Nets to compare accuracy loss against the full-precision Neural Networks
  - Filter and Activation Visualization to see whether Binarization actually works and gain better insights and validate our implementation
  - Evaluate the discussed Binarization Technique on Textual Dataset (Sentiment Analysis)
-



# **CLASSIFICATION PROBLEM IN A SMALL IMAGE DATASET (MNIST)**

# PROBLEM STATEMENT



TEXAS A&M  
UNIVERSITY

- Hand-written Digit Recognition on MNIST Dataset
  - MNIST consists of:
    - Training Set of 60000 examples
    - Test Set of 10000 examples
  - Considered as 'Hello World' of Machine Learning
  - Will visualize the activations and filters to validate our implementation
-

# Our LeNet-5 Model



TEXAS A&M  
UNIVERSITY

```
LeNet_5(  
  (conv1): Conv2d(1, 20, kernel_size=(5, 5), stride=(1, 1))  
  (bn_conv1): BatchNorm2d(20, eps=0.0001, momentum=0.1, affine=False, track_running_stats=True)  
  (relu_conv1): ReLU(inplace)  
  (pool1): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)  
  (bin_conv2): BinConv2d(  
    (bn): BatchNorm2d(20, eps=0.0001, momentum=0.1, affine=True, track_running_stats=True)  
    (conv): Conv2d(20, 50, kernel_size=(5, 5), stride=(1, 1))  
    (k_conv): Conv2d(1, 1, kernel_size=(5, 5), stride=(1, 1)) <= Fixed-filter convolution to obtain the Beta values  
    (relu): ReLU(inplace)  
  )  
  (pool2): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)  
  (bin_ip1): BinConv2d(  
    (bn): BatchNorm2d(50, eps=0.0001, momentum=0.1, affine=True, track_running_stats=True)  
    (linear): Linear(in_features=800, out_features=500, bias=True)  
    (relu): ReLU(inplace)  
  )  
  (ip2): Linear(in_features=500, out_features=10, bias=True)
```

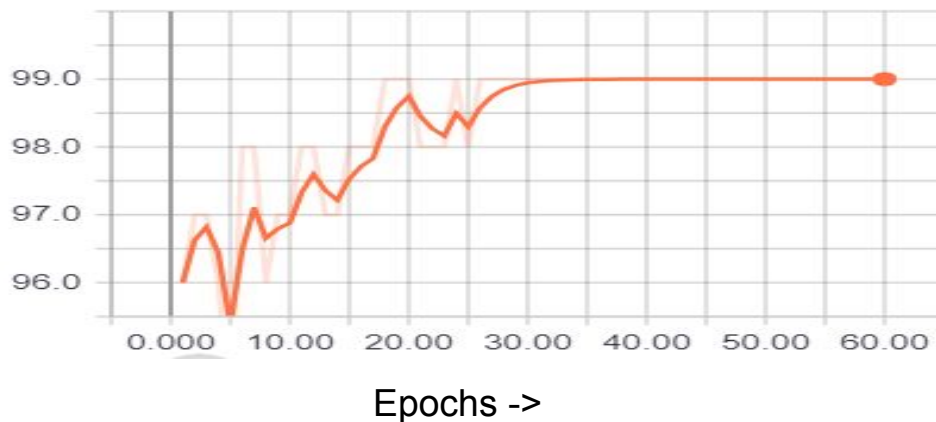
**Source:** <http://yann.lecun.com/exdb/lenet/>

# Validation Accuracy and Loss



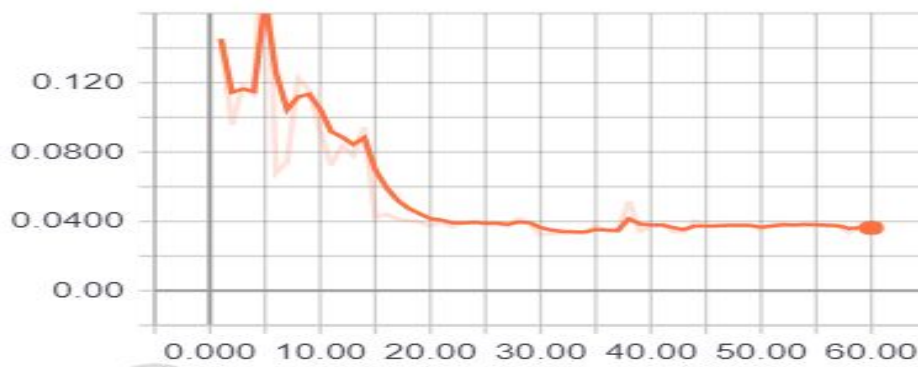
TEXAS A&M  
UNIVERSITY

Accuracy



Max: 99.24%

Loss



Min: 0.036



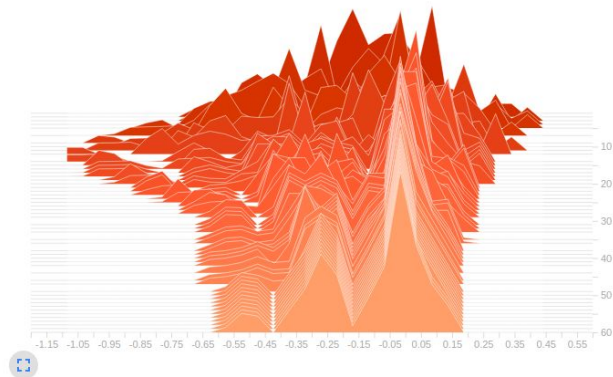
# Bias and Weights Histograms



TEXAS A&M  
UNIVERSITY

bin\_conv2D/bias

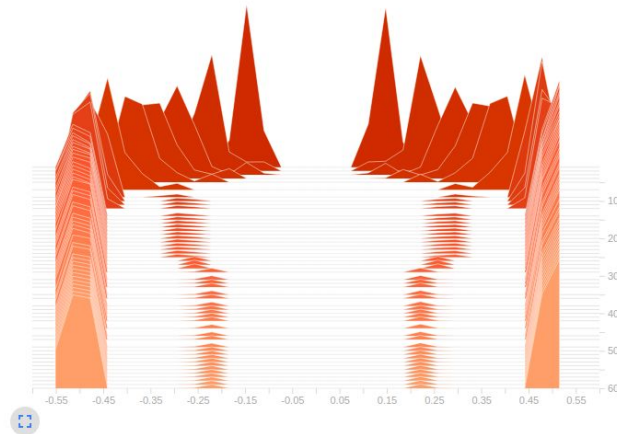
exp1/LeNet\_5



Bias

bin\_conv2D/weights

exp1/LeNet\_5



Weights

Epochs

# Learned Kernels in bin\_conv2



TEXAS A&M  
UNIVERSITY.

bin\_conv2 is a 2D convolutional layer that takes a 28x28x20 input and produces a 10x10x50 output. The input is a 28x28x20 volume, where the first two dimensions are spatial (height and width) and the third is the number of input channels. The output is a 10x10x50 volume, where the first two dimensions are spatial and the third is the number of output channels. The convolution is performed using a 5x5x20 kernel, which is learned during training. The kernel is applied to the input volume to produce the output volume. The output volume is then passed to the next layer in the network.

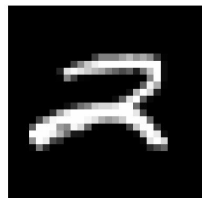
bin\_conv2 have 20  
input channels and  
50 output channels.  
Hence, there are  
1000 kernels.

**NOTE:** All the  
weights are binary,  
either white or black.

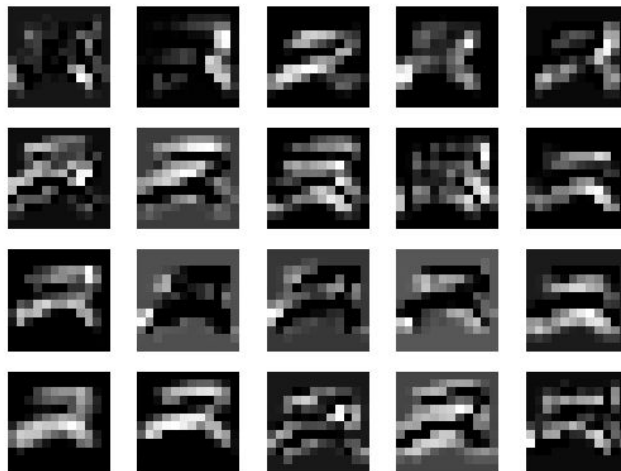
# Output of conv1 (Full-precision)



TEXAS A&M  
UNIVERSITY



Input

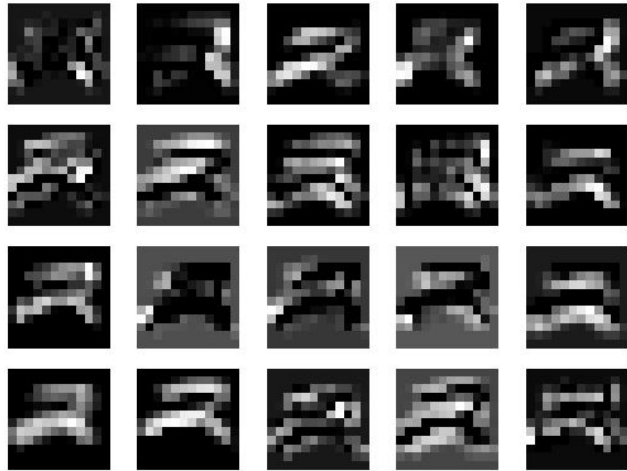


Output of conv1 (full-precision convolution)

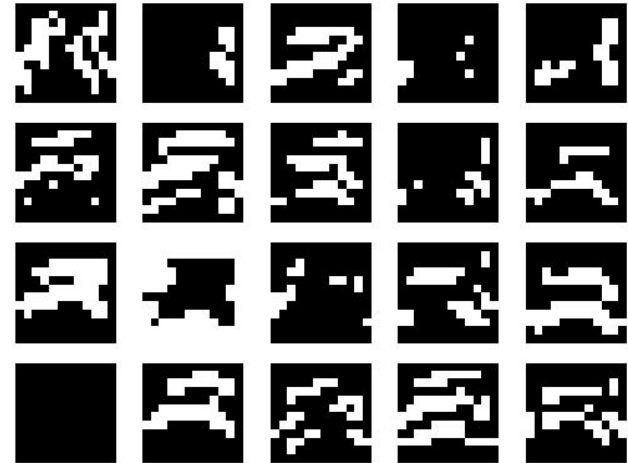
# Binary Activation



TEXAS A&M  
UNIVERSITY



Output of first convolution block



Input to bin\_conv2 which is either 0 or 1

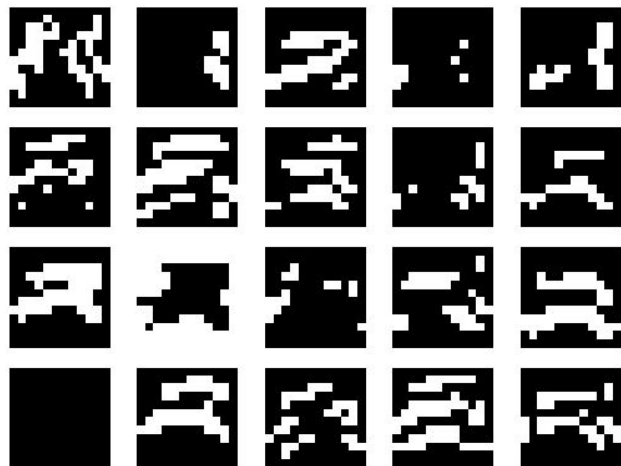
**NOTE:** At this point we have verified that both the inputs and the kernel weights to bin\_conv2 are either 0 or 1

---

# Output of bin\_conv2



TEXAS A&M  
UNIVERSITY.



Input to bin\_conv2 which is  
either 0 or 1

The output of bin\_conv2



# **CLASSIFICATION PROBLEM IN A LARGE IMAGE DATASET (CIFAR-10)**



# PROBLEM STATEMENT



TEXAS A&M  
UNIVERSITY.

- Image Classification on CIFAR-10
- The CIFAR-10 dataset consists of 60000 32x32 colour images in 10 classes, with 6000 images per class.
- There are 50000 training images and 10000 test images.
- The dataset is divided into 5 training batches and 1 test batch.
- The classes are completely mutually exclusive.

airplane



automobile



bird



cat



deer



dog



frog



horse



ship



truck



# Network-In-Network (NIN) Model



TEXAS A&M  
UNIVERSITY

```
(module): Net(
  (xnor): Sequential(
    (0): Conv2d(3, 192, kernel_size=(5, 5), stride=(1, 1), padding=(2, 2))
    (1): BatchNorm2d(192, eps=0.0001, momentum=0.1, affine=False, track_running_stats=True)
    (2): ReLU(inplace)
    (3): BinConv2d(
      (bn): BatchNorm2d(192, eps=0.0001, momentum=0.1, affine=True, track_running_stats=True)
      (conv): Conv2d(192, 160, kernel_size=(1, 1), stride=(1, 1))
      (relu): ReLU(inplace)
    )
    (4): BinConv2d(
      (bn): BatchNorm2d(160, eps=0.0001, momentum=0.1, affine=True, track_running_stats=True)
      (conv): Conv2d(160, 96, kernel_size=(1, 1), stride=(1, 1))
      (relu): ReLU(inplace)
    )
    (5): MaxPool2d(kernel_size=3, stride=2, padding=1, dilation=1, ceil_mode=False)
```



# Network-In-Network (NIN) Model



TEXAS A&M  
UNIVERSITY

**(6): BinConv2d(**

**(bn): BatchNorm2d(96, eps=0.0001, momentum=0.1, affine=True, track\_running\_stats=True)**

**(dropout): Dropout(p=0.5)**

**(conv): Conv2d(96, 192, kernel\_size=(5, 5), stride=(1, 1), padding=(2, 2))**

**(relu): ReLU(inplace)**

**)**

**(7): BinConv2d(**

**(bn): BatchNorm2d(192, eps=0.0001, momentum=0.1, affine=True, track\_running\_stats=True)**

**(conv): Conv2d(192, 192, kernel\_size=(1, 1), stride=(1, 1))**

**(relu): ReLU(inplace)**

**)**

**(8): BinConv2d(**

**(bn): BatchNorm2d(192, eps=0.0001, momentum=0.1, affine=True, track\_running\_stats=True)**

**(conv): Conv2d(192, 192, kernel\_size=(1, 1), stride=(1, 1))**

**(relu): ReLU(inplace)**

**)**

**(9): AvgPool2d(kernel\_size=3, stride=2, padding=1)**

# Network-In-Network (NIN) Model



TEXAS A&M  
UNIVERSITY

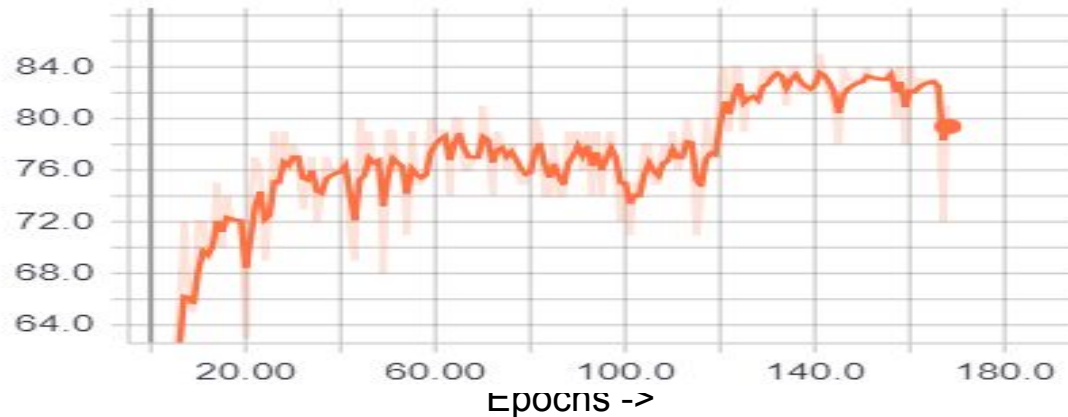
```
(10): BinConv2d(  
(bn): BatchNorm2d(192, eps=0.0001, momentum=0.1, affine=True, track_running_stats=True)  
(dropout): Dropout(p=0.5)  
(conv): Conv2d(192, 192, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
(relu): ReLU(inplace)  
)  
(11): BinConv2d(  
(bn): BatchNorm2d(192, eps=0.0001, momentum=0.1, affine=True, track_running_stats=True)  
(conv): Conv2d(192, 192, kernel_size=(1, 1), stride=(1, 1))  
(relu): ReLU(inplace)  
)  
(12): BatchNorm2d(192, eps=0.0001, momentum=0.1, affine=False, track_running_stats=True)  
(13): Conv2d(192, 10, kernel_size=(1, 1), stride=(1, 1))  
(14): ReLU(inplace)  
(15): AvgPool2d(kernel_size=8, stride=1, padding=0)  
)  
)
```

# Validation Accuracy and Loss



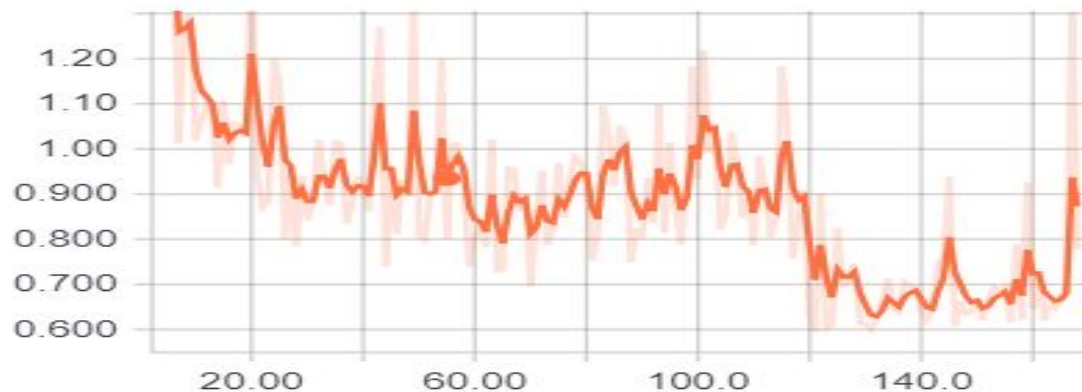
TEXAS A&M  
UNIVERSITY

Accuracy



**Max: 85.02%**

Loss



**Min: 0.595**

# Weights and Bias Histograms

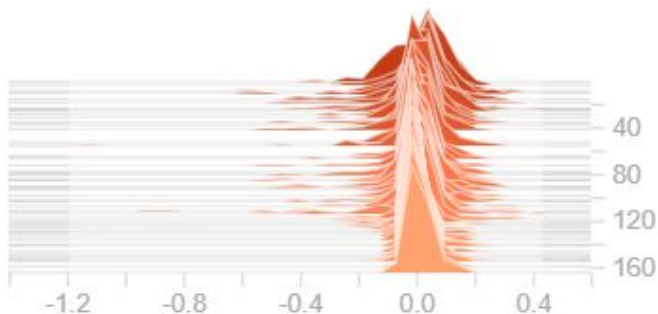


TEXAS A&M  
UNIVERSITY

## Layer 3 - 1st Binary Convolution

Layer\_3\_bias

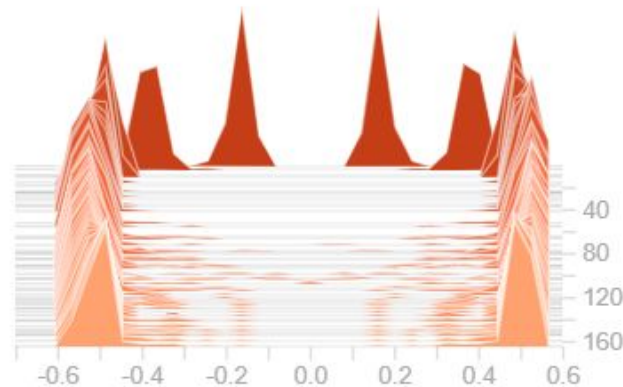
exp1/nin



Bias

Layer\_3\_weights

exp1/nin



Weights

Epochs

# Weights and Bias Histograms

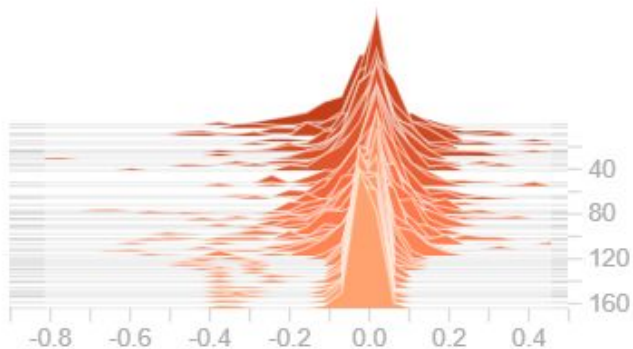


TEXAS A&M  
UNIVERSITY

## Layer 4 - 2nd Binary Convolution

Layer\_4\_bias

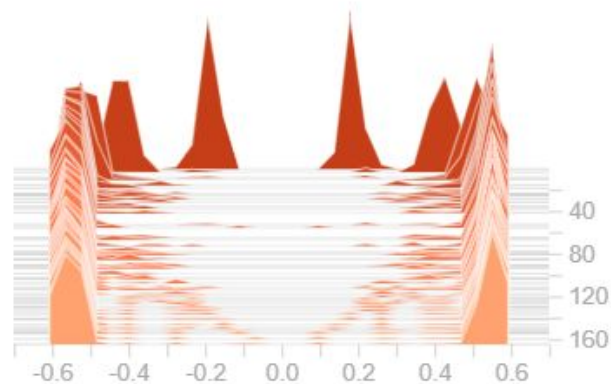
exp1/nin



Bias

Layer\_4\_weights

exp1/nin



Weights

Epochs

# Weights and Bias Histograms

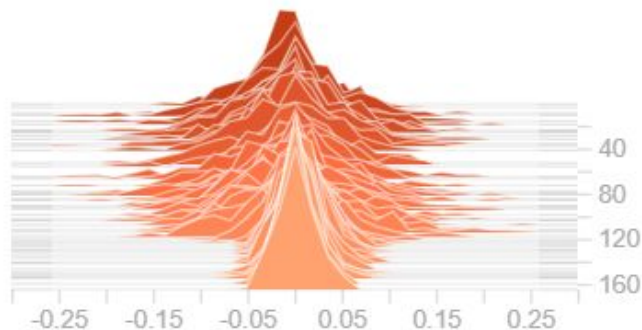


TEXAS A&M  
UNIVERSITY

## Layer 6- 3rd Binary Convolution

Layer\_6\_bias

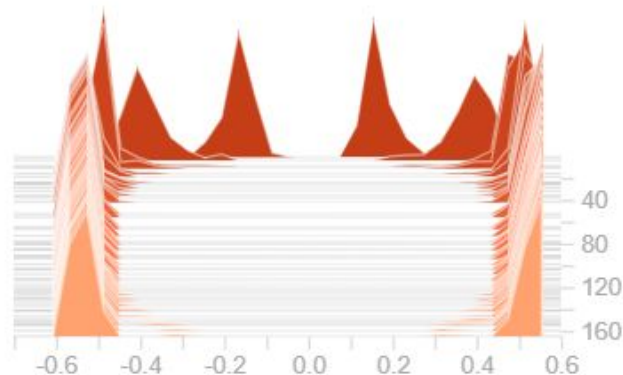
exp1/nin



Bias

Layer\_6\_weights

exp1/nin



Weights

Epochs

# Weights and Bias Histograms

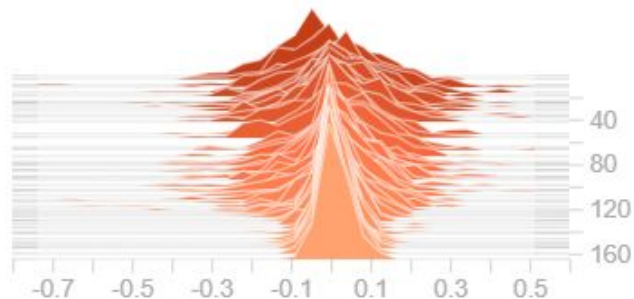


TEXAS A&M  
UNIVERSITY

## Layer 7 - 4th Binary Convolution

Layer\_7\_bias

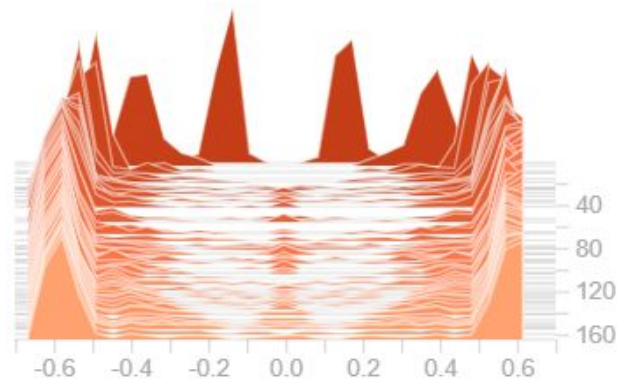
exp1/nin



Bias

Layer\_7\_weights

exp1/nin



Weights

Epochs

# Weights and Bias Histograms

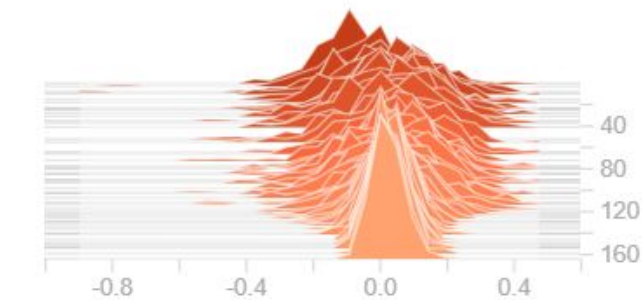


TEXAS A&M  
UNIVERSITY

## Layer 8 - 5th Binary Convolution

Layer\_8\_bias

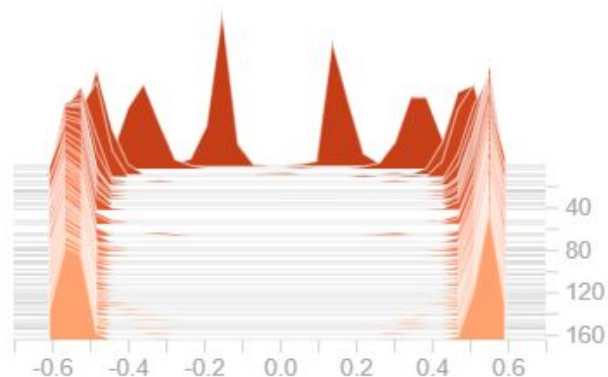
exp1/nin



Bias

Layer\_8\_weights

exp1/nin



Weights

Epochs



# Weights and Bias Histograms

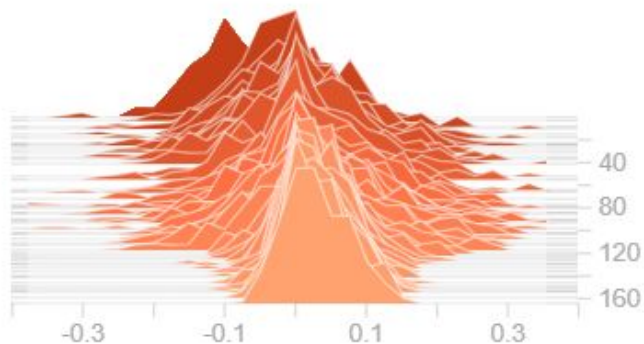


TEXAS A&M  
UNIVERSITY

## Layer 10- 6th Binary Convolution

Layer\_10\_bias

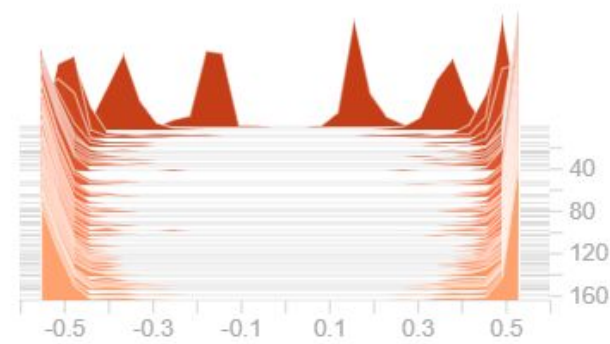
exp1/nin



Bias

Layer\_10\_weights

exp1/nin



Weights

Epochs

# Weights and Bias Histograms

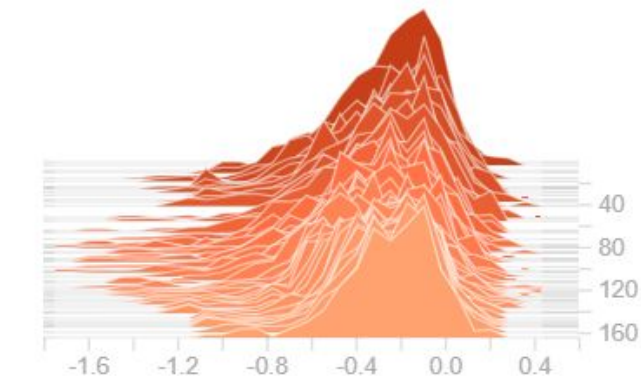


TEXAS A&M  
UNIVERSITY

## Layer 11 - 7th Binary Convolution

Layer\_11\_bias

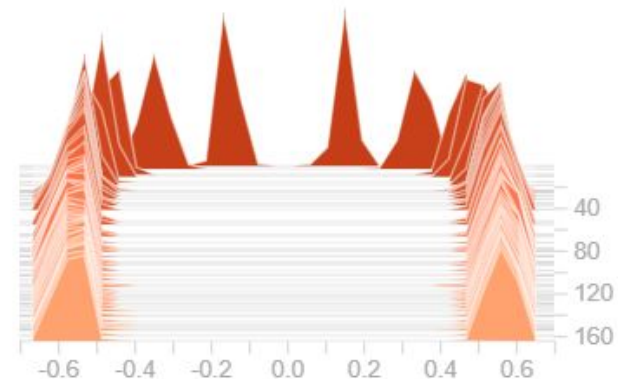
exp1/nin



Bias

Layer\_11\_weights

exp1/nin



Weights

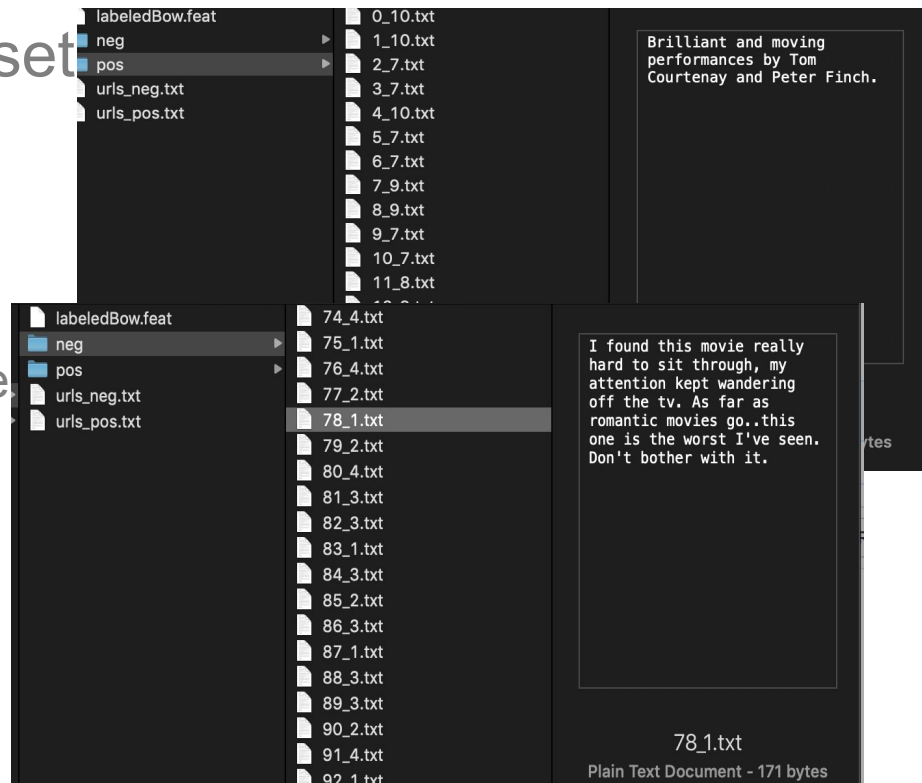
Epochs



# **CLASSIFICATION PROBLEM IN A TEXT DATASET (IMDB)**

## IMDB Movie Review Text Dataset

- Movie Reviews
- Labels : Binary Sentiment
- 50k labelled reviews
- 50k unlabelled reviews
- At most 30 reviews for a single movie
  - Correlation of reviews
- Stratified train and test set
  - Disjoint on the basis of movies
- Sentiment Analysis



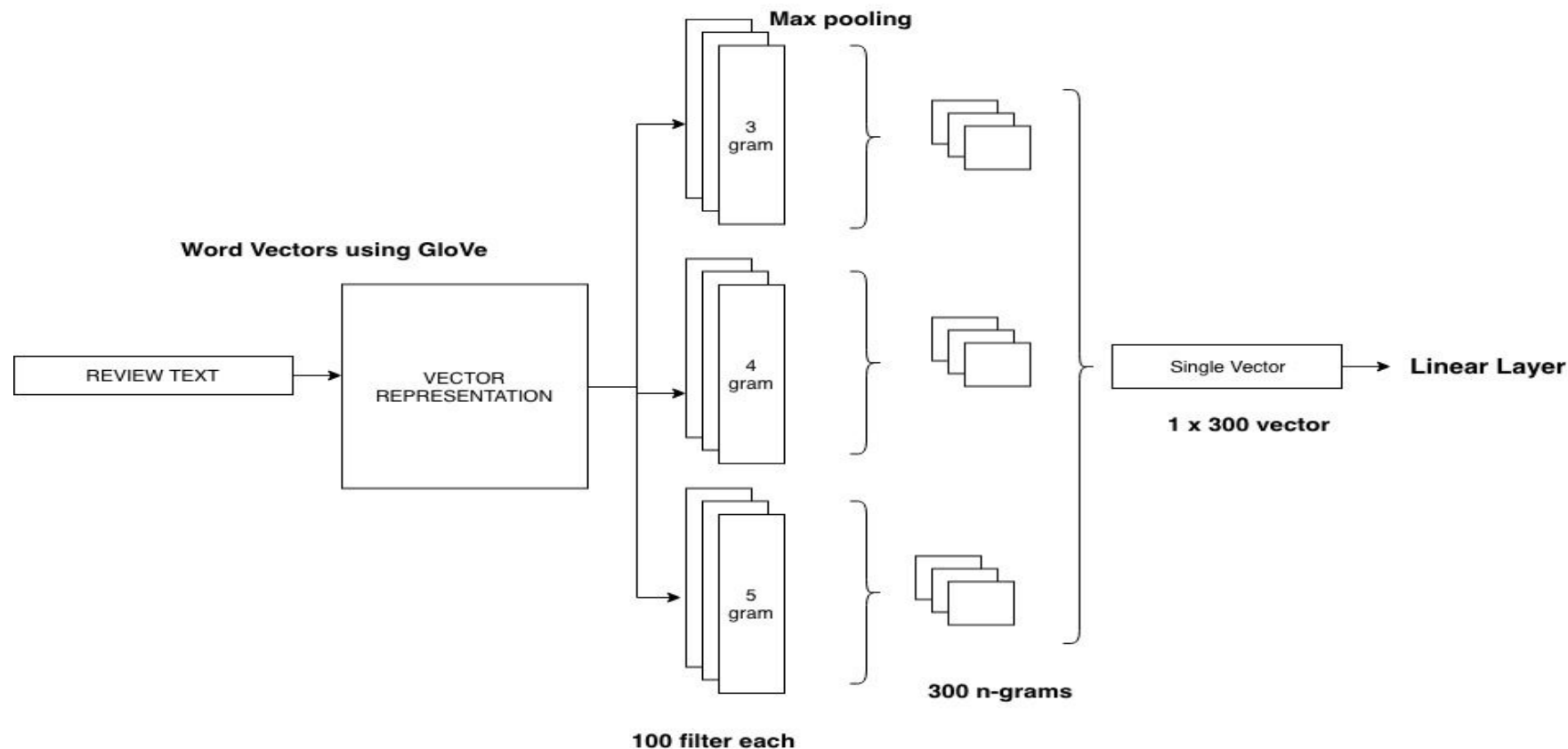
## GloVe - **G**lobal **V**ectors for Word Representation

- Used to obtain vector representation for words
  - Pre-trained word vector
    - **glove.6B**
      - 50d
      - **100d (400K vocab, 6B tokens)**
      - 200d
      - 300d
    - glove.42B.300d
    - glove.840B.300d
    - glove.twitter.27B
-

# Model



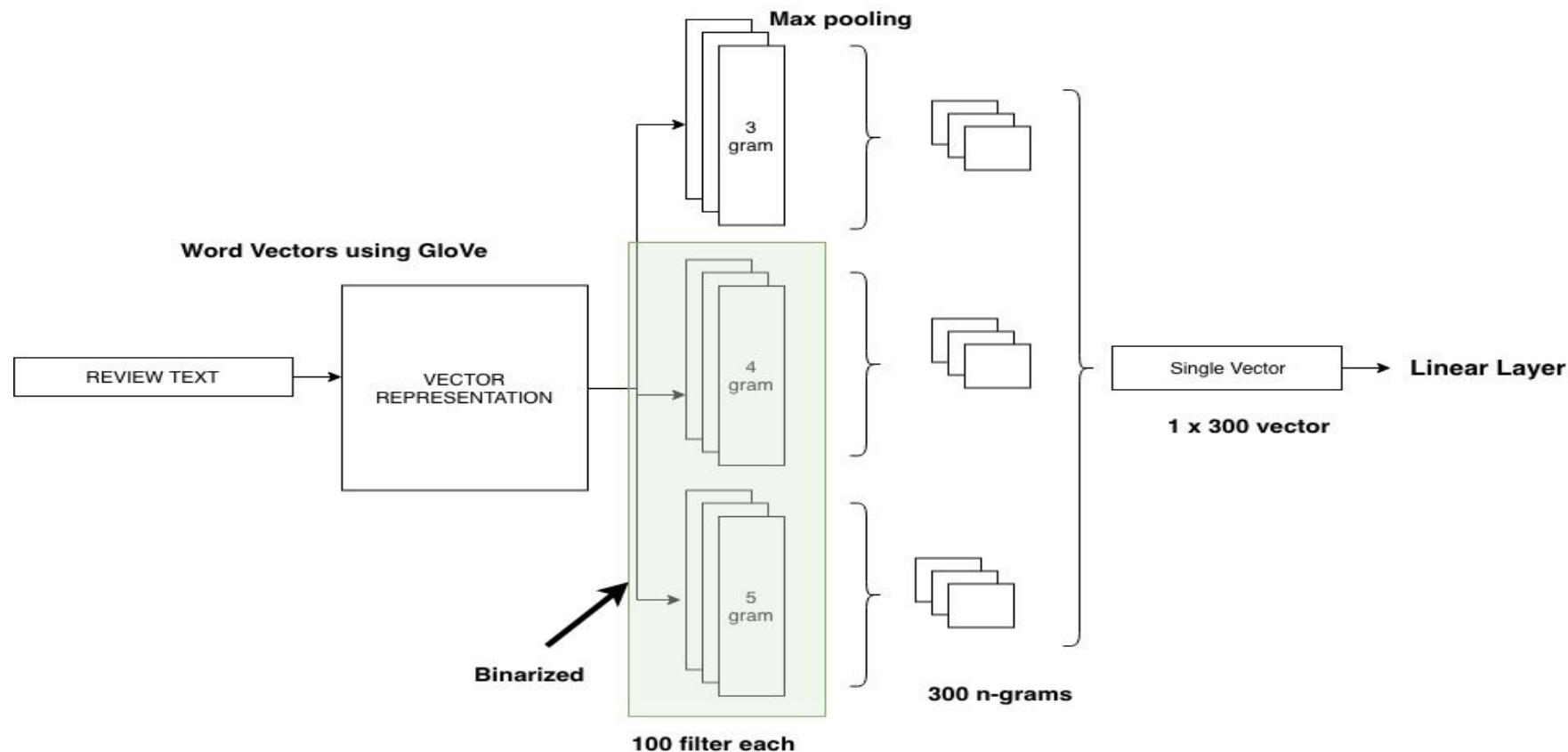
TEXAS A&M  
UNIVERSITY



# Binarization



TEXAS A&M  
UNIVERSITY

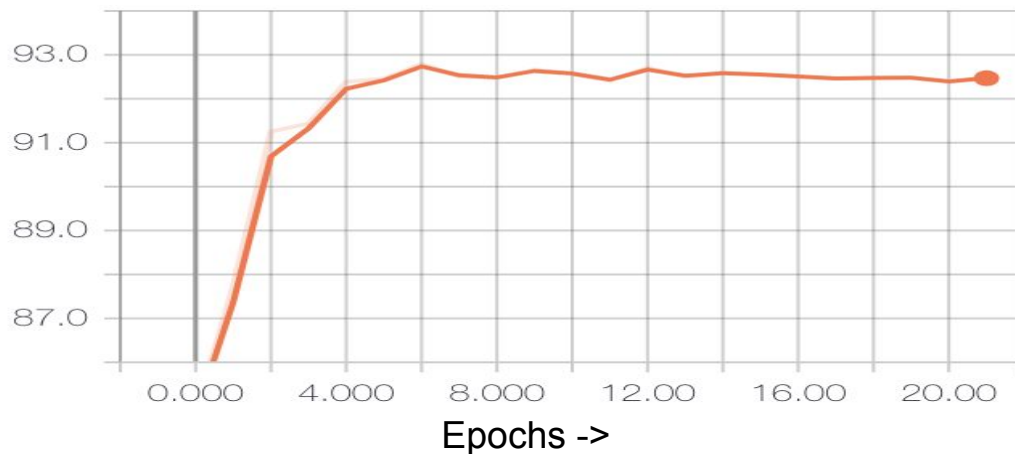


# Validation Accuracy and Loss

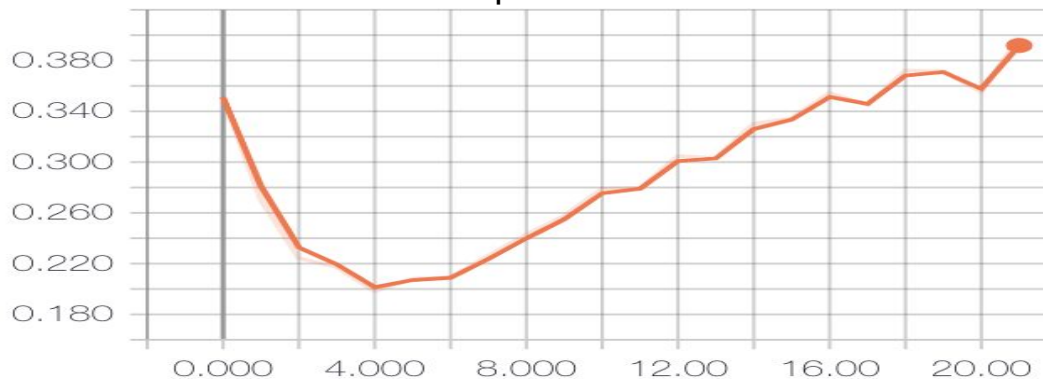


TEXAS A&M  
UNIVERSITY

Accuracy



Loss



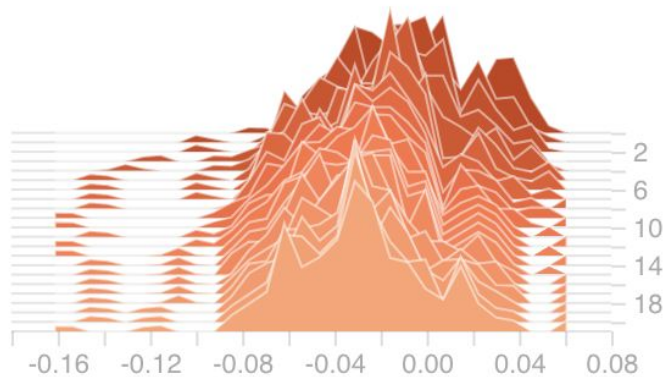


# Weights and Bias Histograms



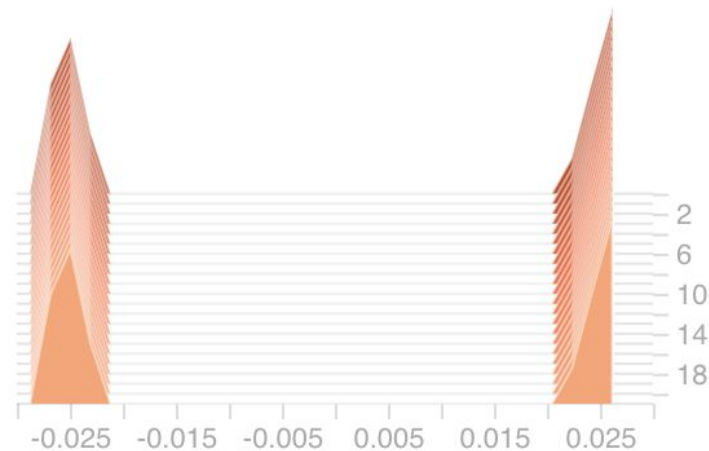
TEXAS A&M  
UNIVERSITY

4-gram



Bias

Epochs



Weights

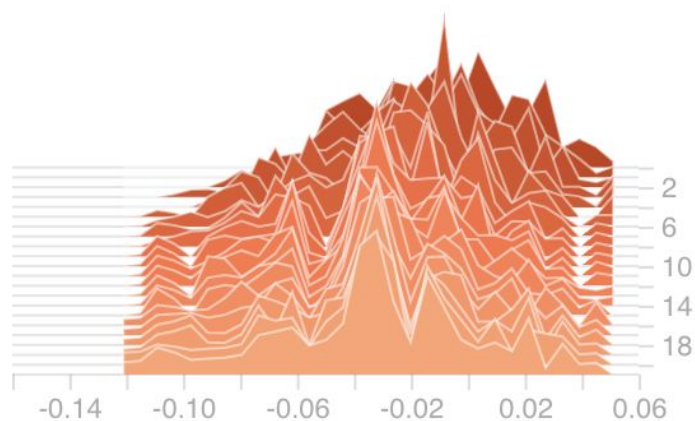
# Weights and Bias Histograms



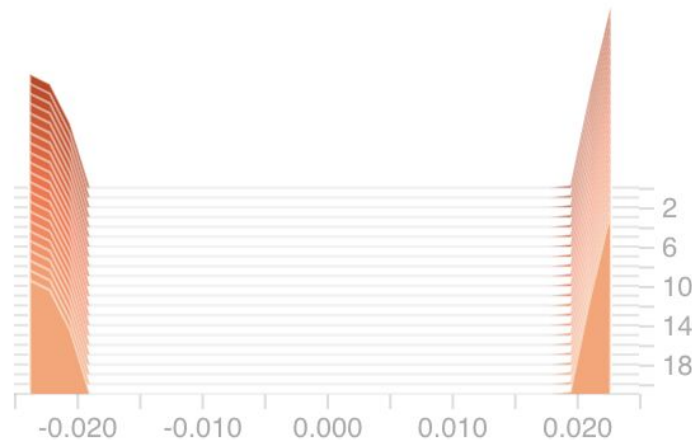
TEXAS A&M  
UNIVERSITY

5-gram

Epochs



Bias



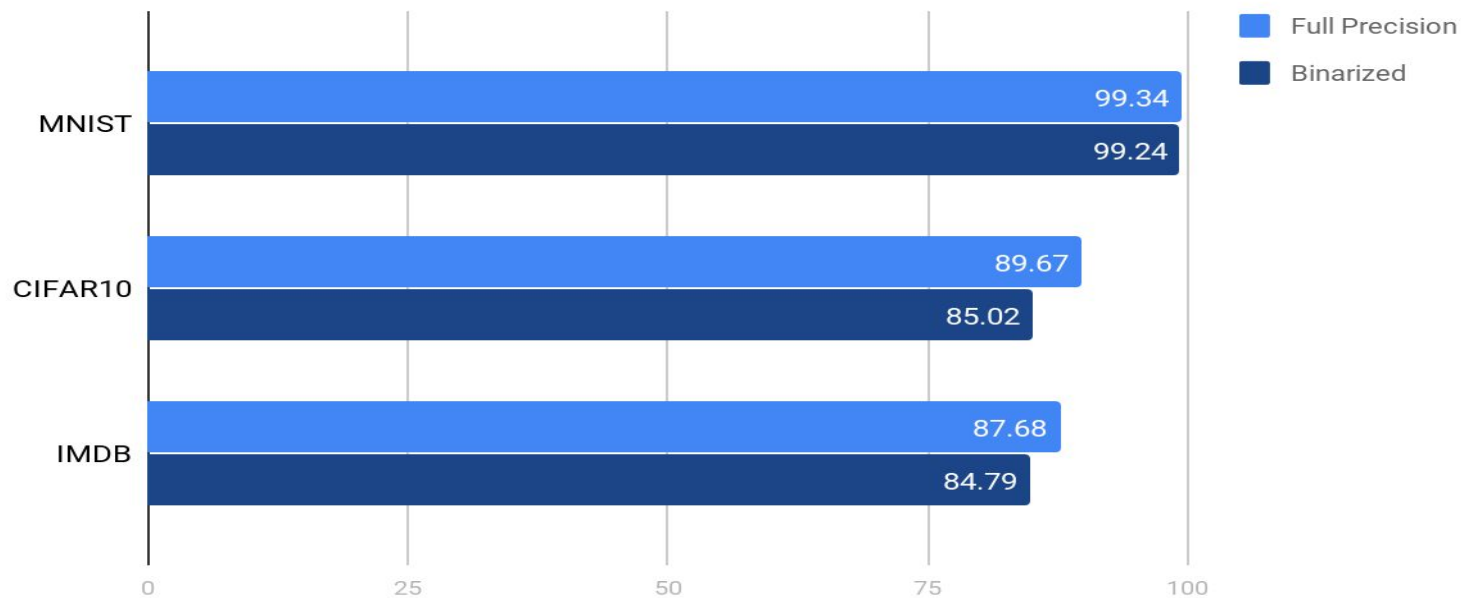
Weights

# Final Results



TEXAS A&M  
UNIVERSITY

## Accuracy





# Future Work

- ImageNet Dataset
  - I2B2 Dataset
  - Optimizing with actual Bitwise Operations on GPU and how the binary values are saved.
-



**THANK YOU**

**Questions?**