

Model Compression: Weights Pruning for RNNs

Luyu Wang

University of Waterloo & RBC Research Institute

July 17, 2017

Outlines

- Introduction to Model Compression
- RNN Model Compression via Weights Pruning
- Experimental Results
- Conclusion

The Need for Model Compression

- **Success of Deep Neural Network Models**
- **Vast datasets, GPU for training**
- Production environment
- Deep neural nets need lots of computational to make inferences

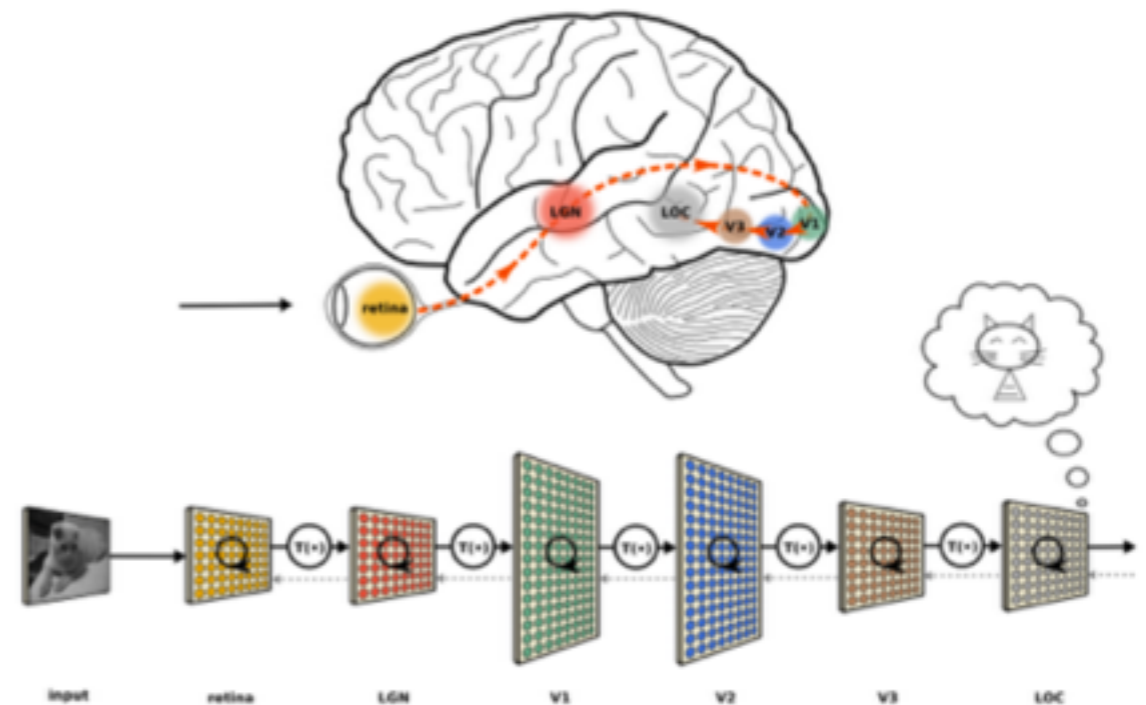


Image
Recognition



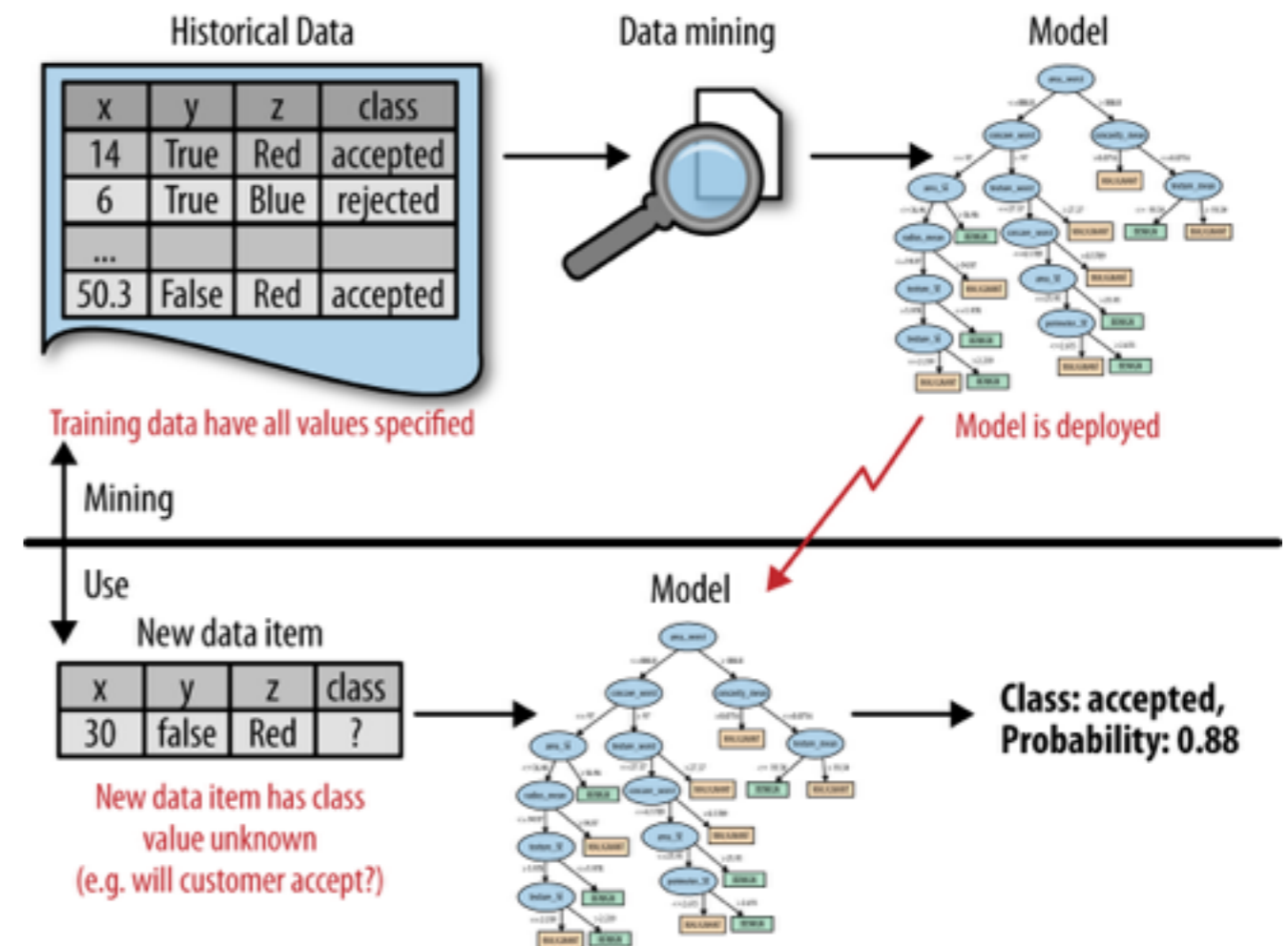
Speech
Recognition



Natural Language
Processing

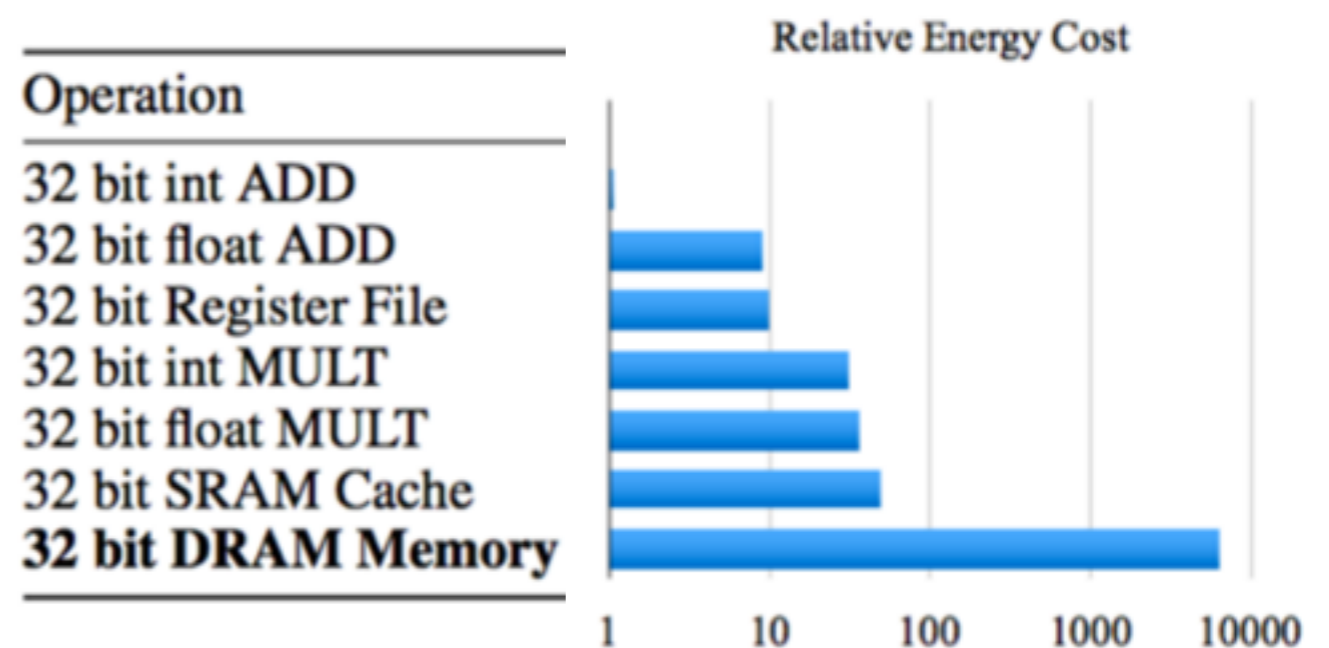
The Need for Model Compression

- Success of Deep Neural Network Models
- Vast datasets, GPU for training
- **Production environment**
- Deep neural nets need lots of computational to make inferences



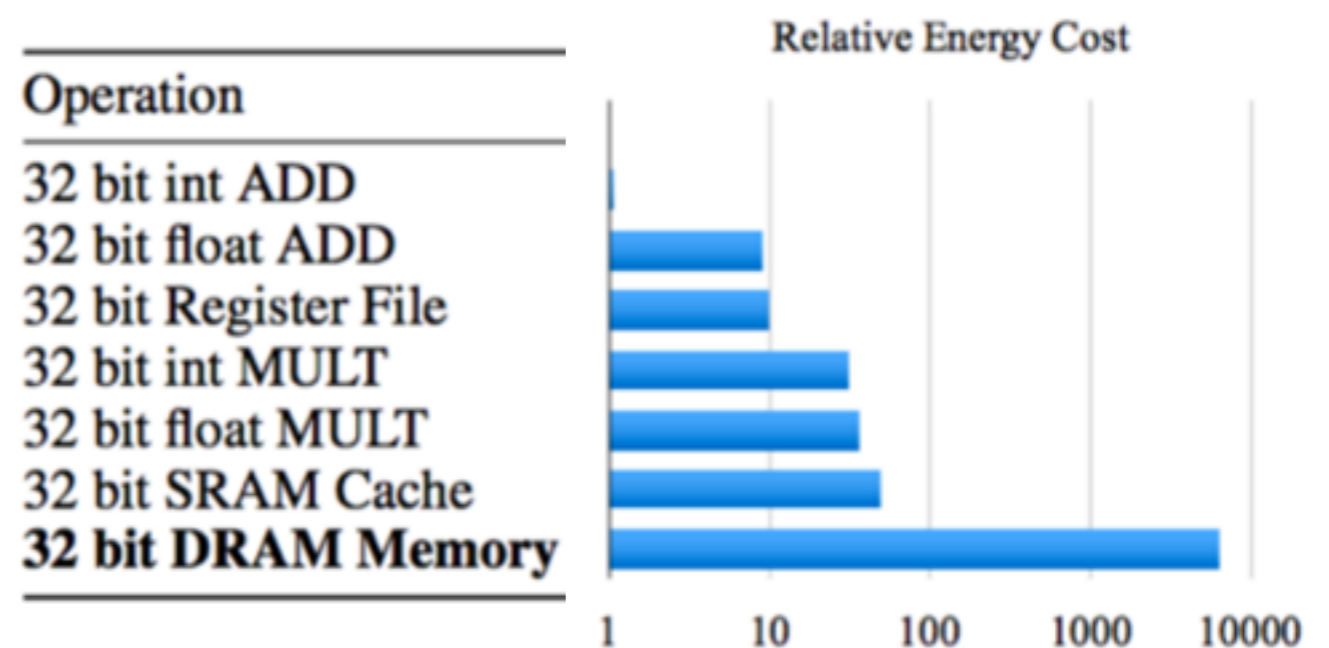
The Need for Model Compression

- Success of Deep Neural Network Models
- Vast datasets, GPU for training
- Production environment
- **Deep neural nets need lots of computational to make inferences**



The Need for Model Compression

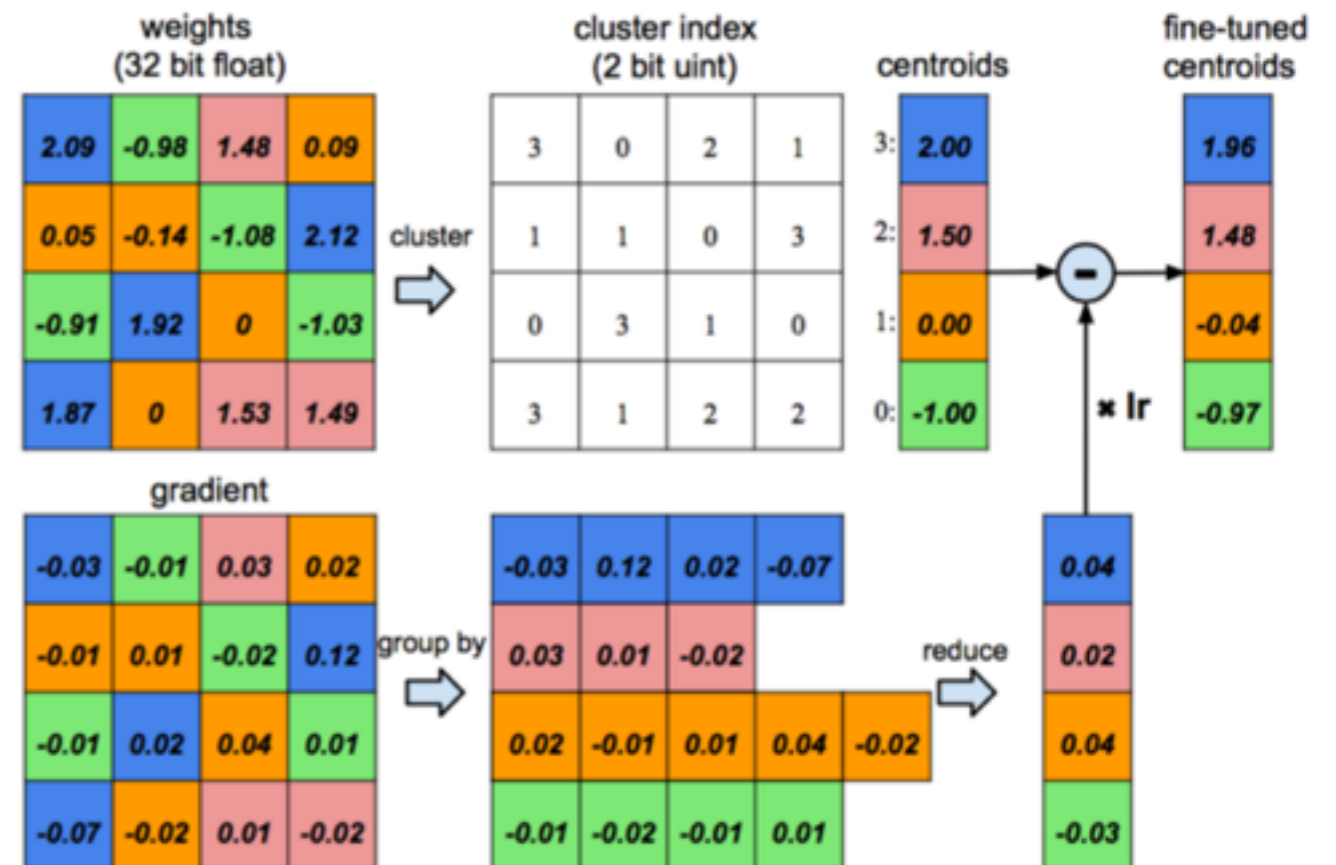
- Success of Deep Neural Network Models
- Vast datasets, GPU for training
- Production environment
- **Deep neural nets need lots of computational to make inferences**



Question: less cost for deployed models?

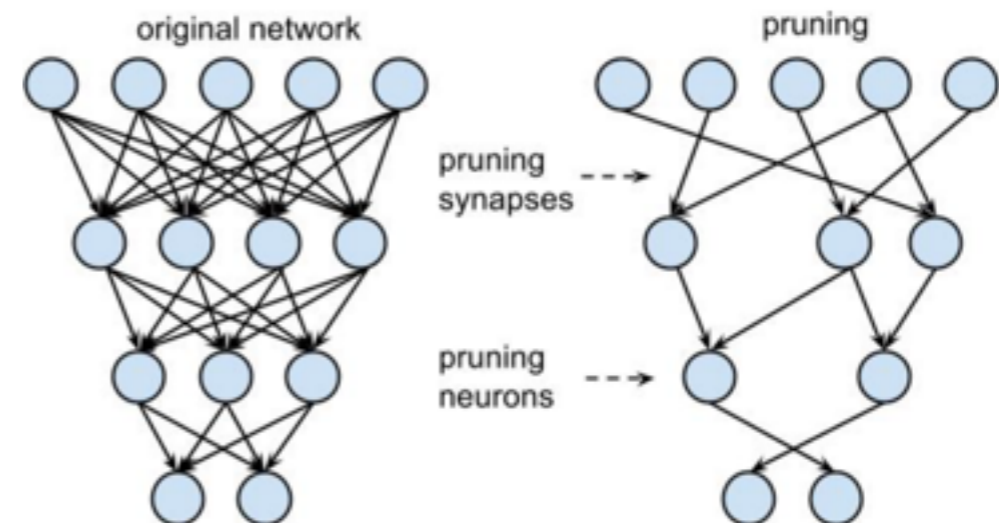
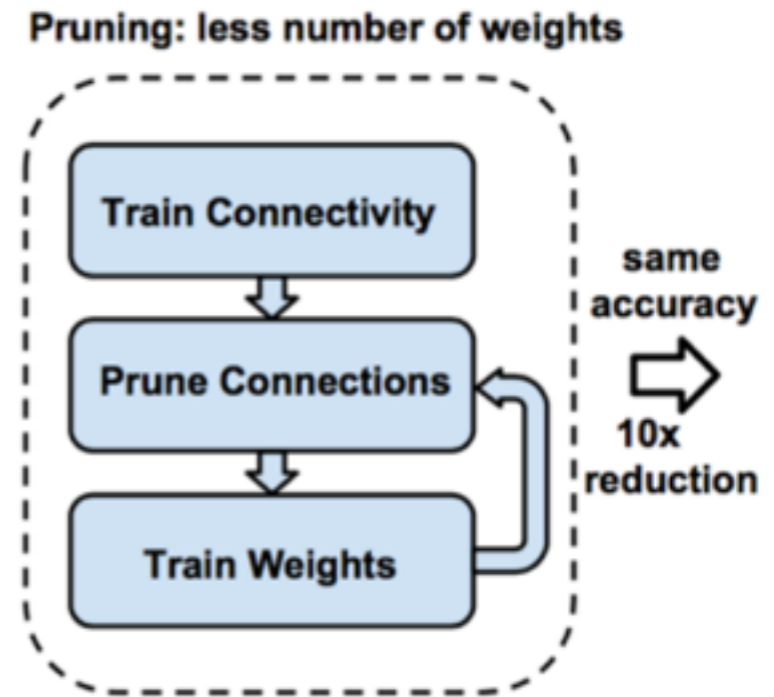
Methods for Model Compression

- Weights quantization (or even binarization)
- Model distillation
- Weights pruning



Weights Pruning

- Learn the connectivity via normal network training
- Prune the low-weight connections
- Retrain the sparse network



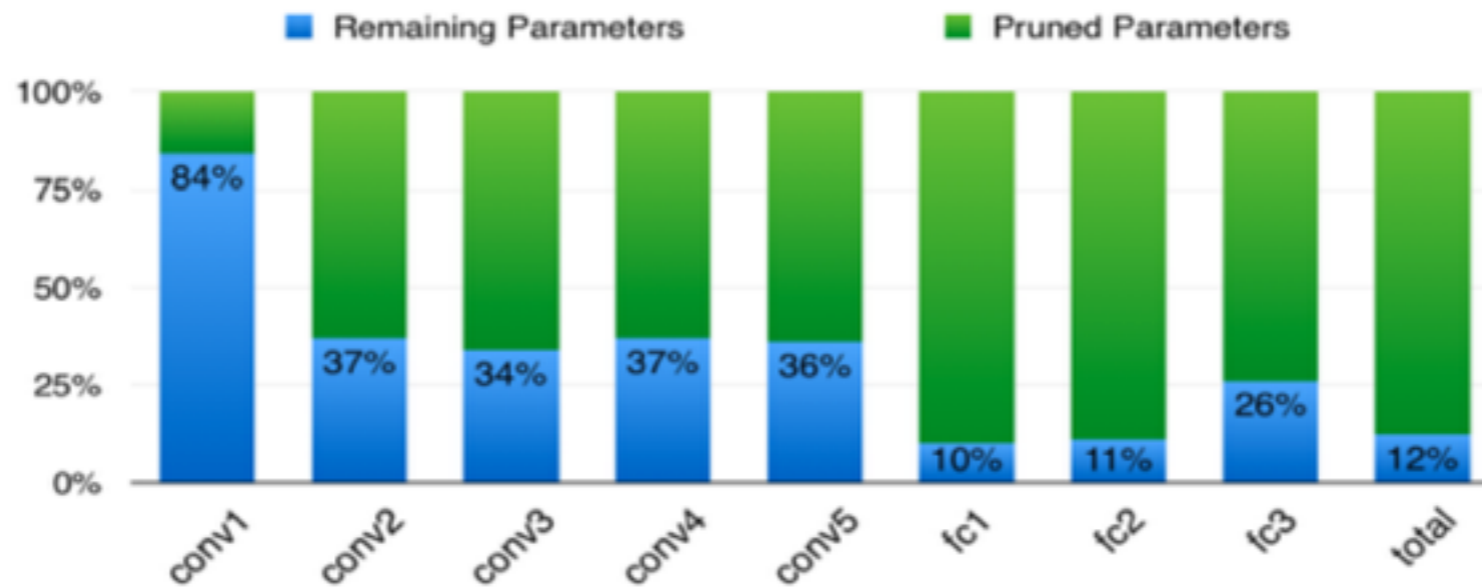
Han, Song, et al. "Learning both weights and connections for efficient neural network." Advances in Neural Information Processing Systems. 2015.

Weights Pruning

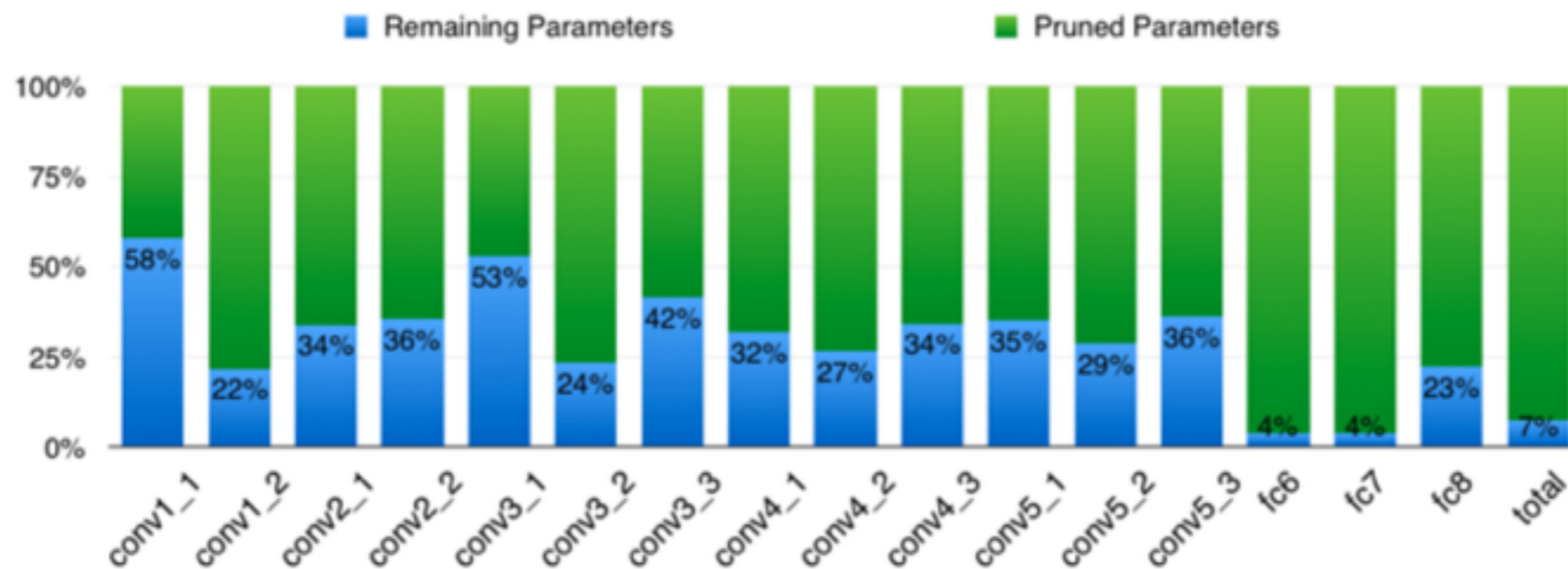
Network	Top-1 Error	Top-5 Error	Parameters	Compression Rate
LeNet-300-100 Ref	1.64%	-	267K	
LeNet-300-100 Pruned	1.59%	-	22K	12×
LeNet-5 Ref	0.80%	-	431K	
LeNet-5 Pruned	0.77%	-	36K	12×
AlexNet Ref	42.78%	19.73%	61M	
AlexNet Pruned	42.77%	19.67%	6.7M	9×
VGG16 Ref	31.50%	11.32%	138M	
VGG16 Pruned	31.34%	10.88%	10.3M	13×

Weights Pruning

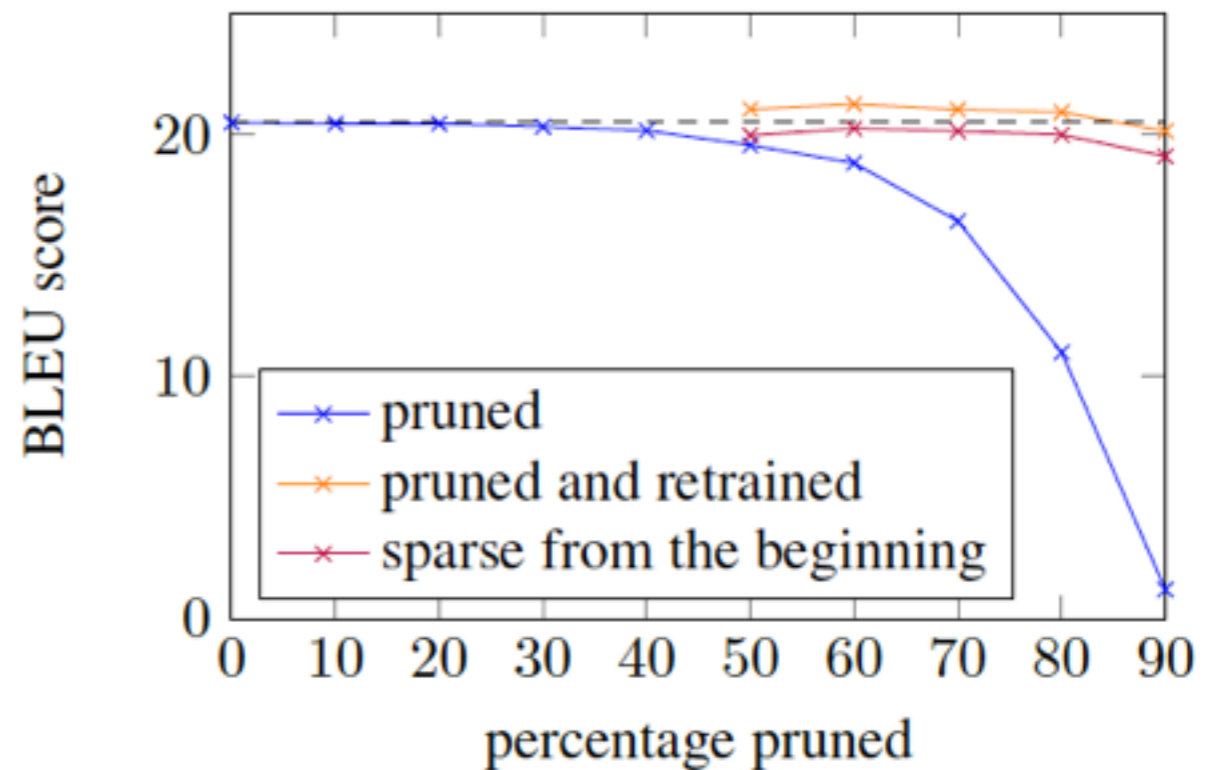
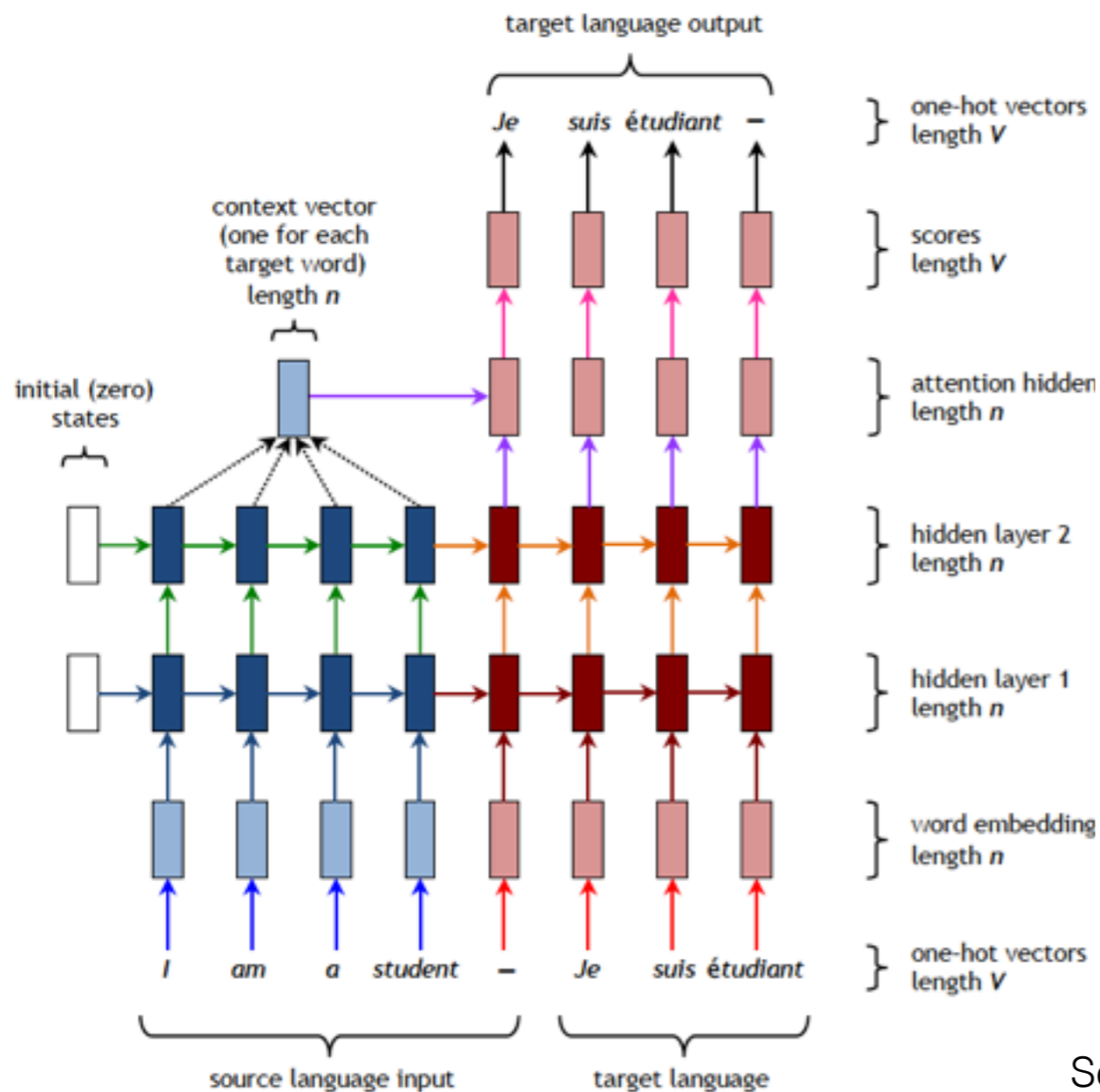
AlexNet



VGG



Neural Machine Translation Models via Pruning



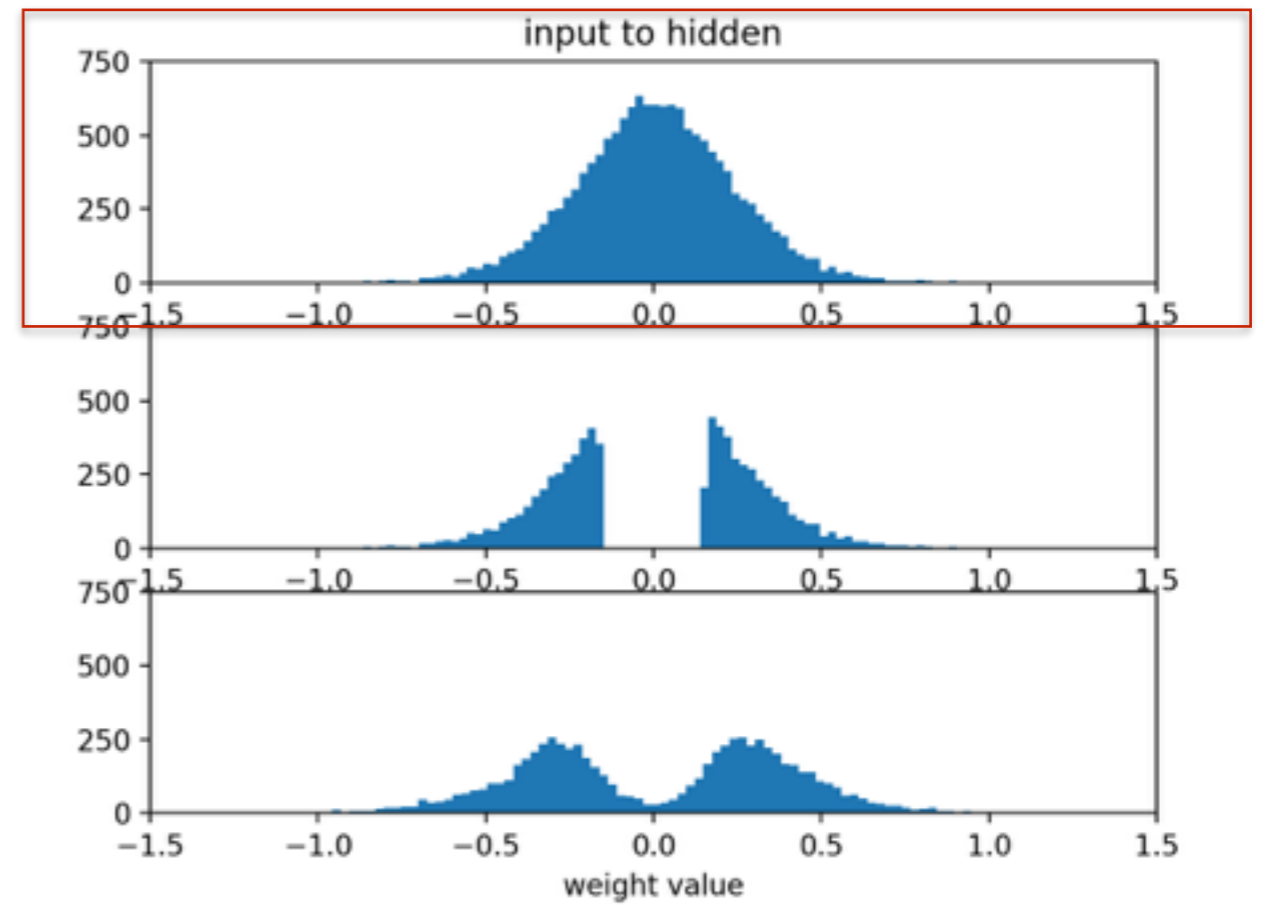
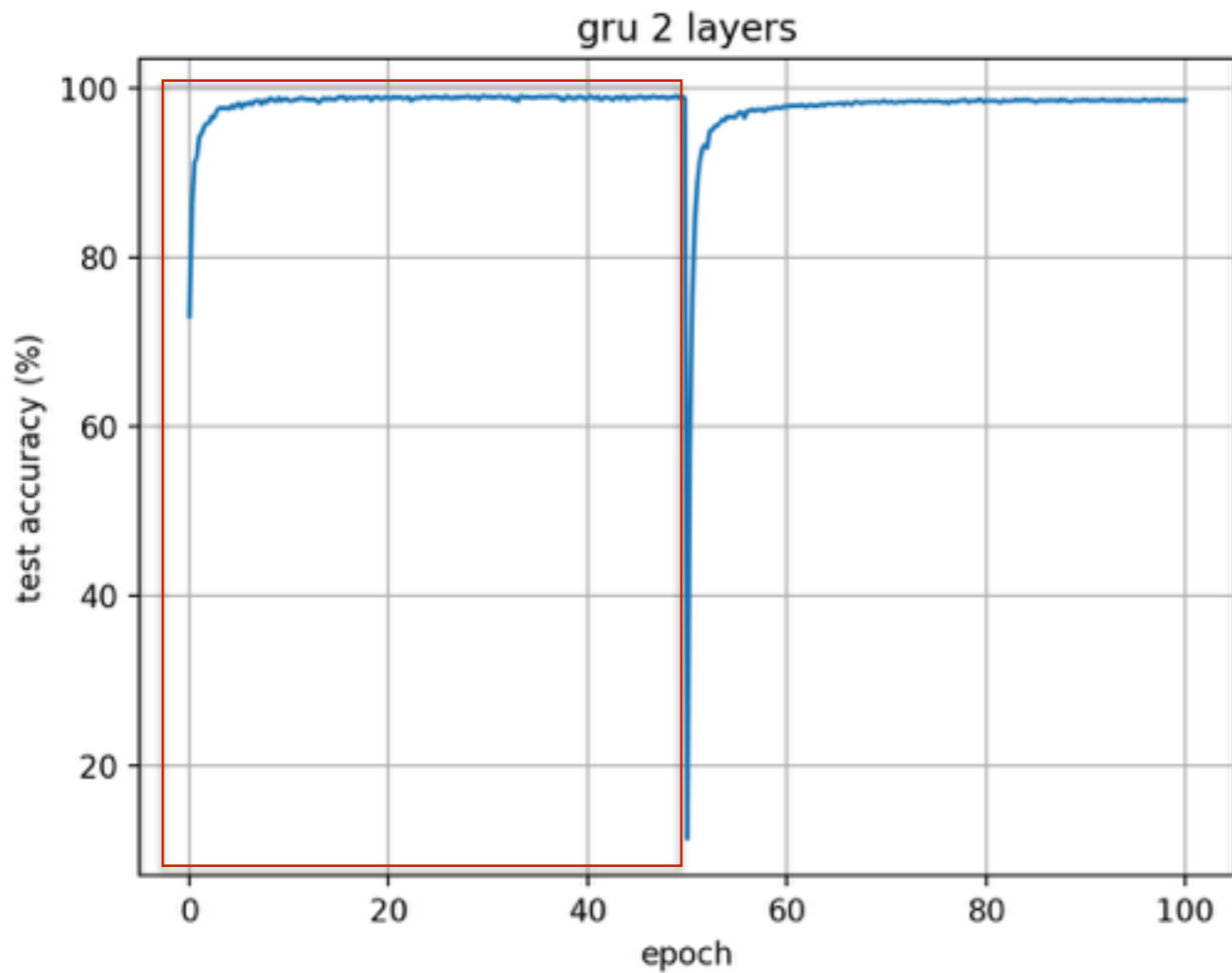
See, Abigail, Minh-Thang Luong, and Christopher D. Manning. "Compression of Neural Machine Translation Models via Pruning." CoNLL 2016 (2016): 291.

Pruning RNNs

- Vanilla RNN
- LSTM
- GRU
- MNIST Dataset 28x28
- Hidden states = 128
- RMSprop with fixed learning rate 0.001
- PyTorch on AWS GPU server

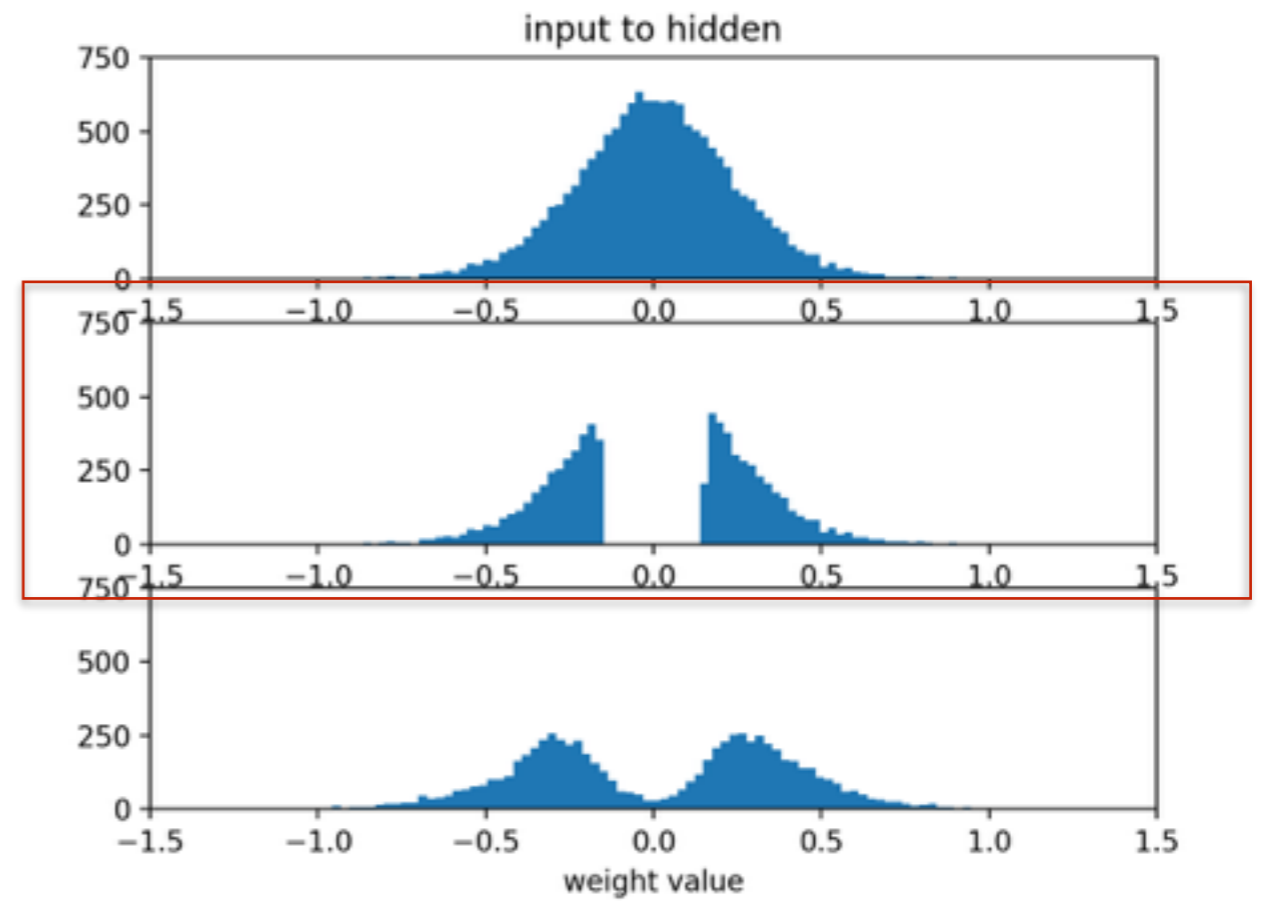
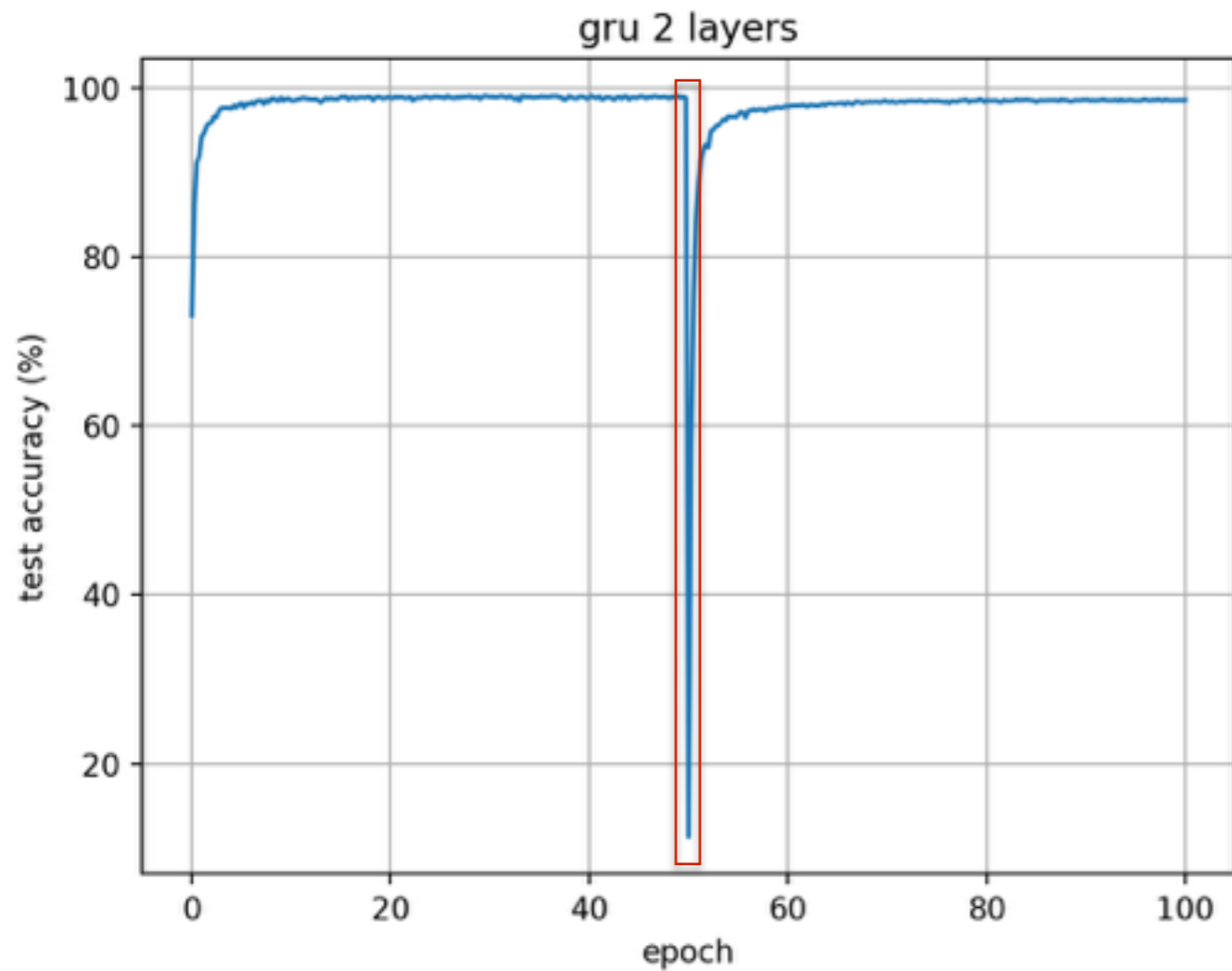
Pruning RNNs

Pre-train



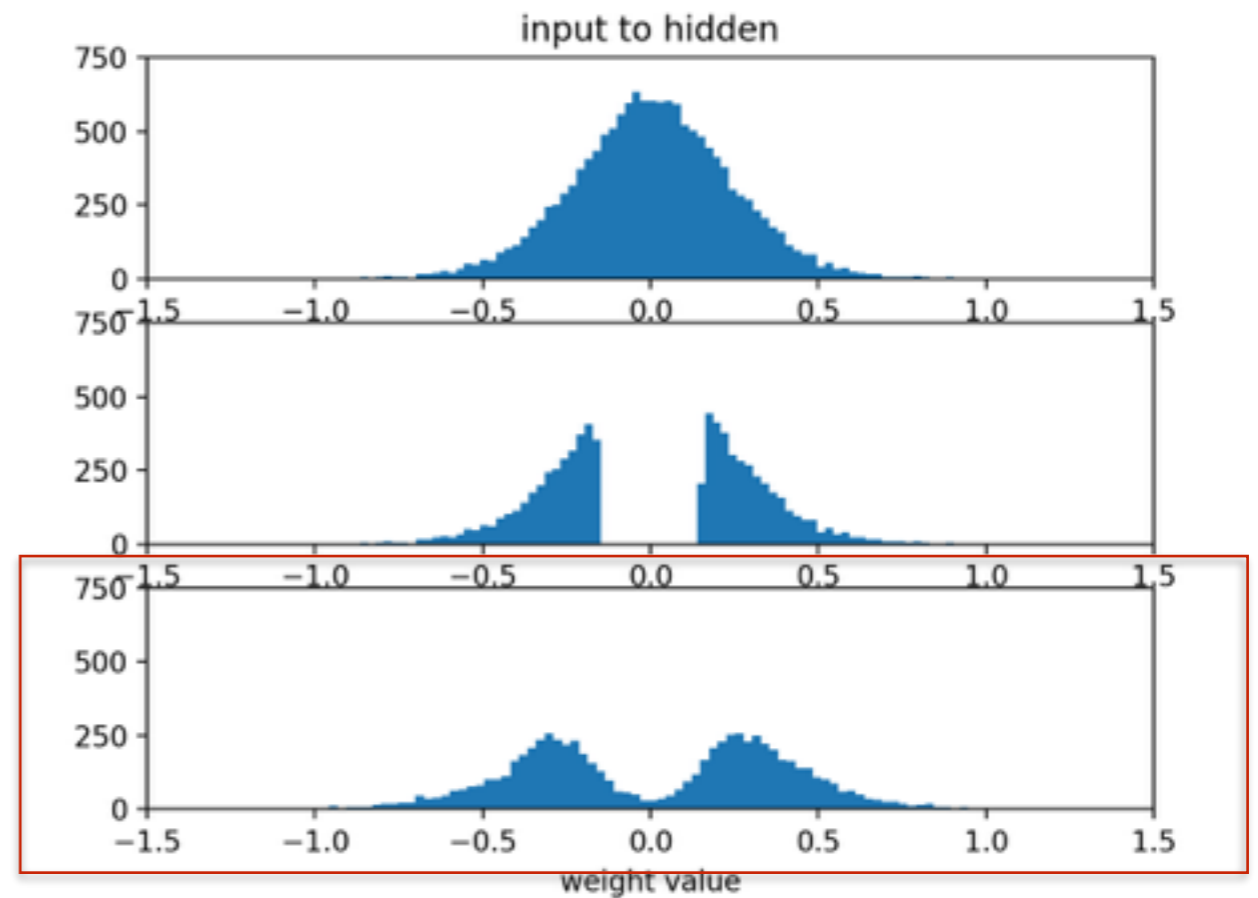
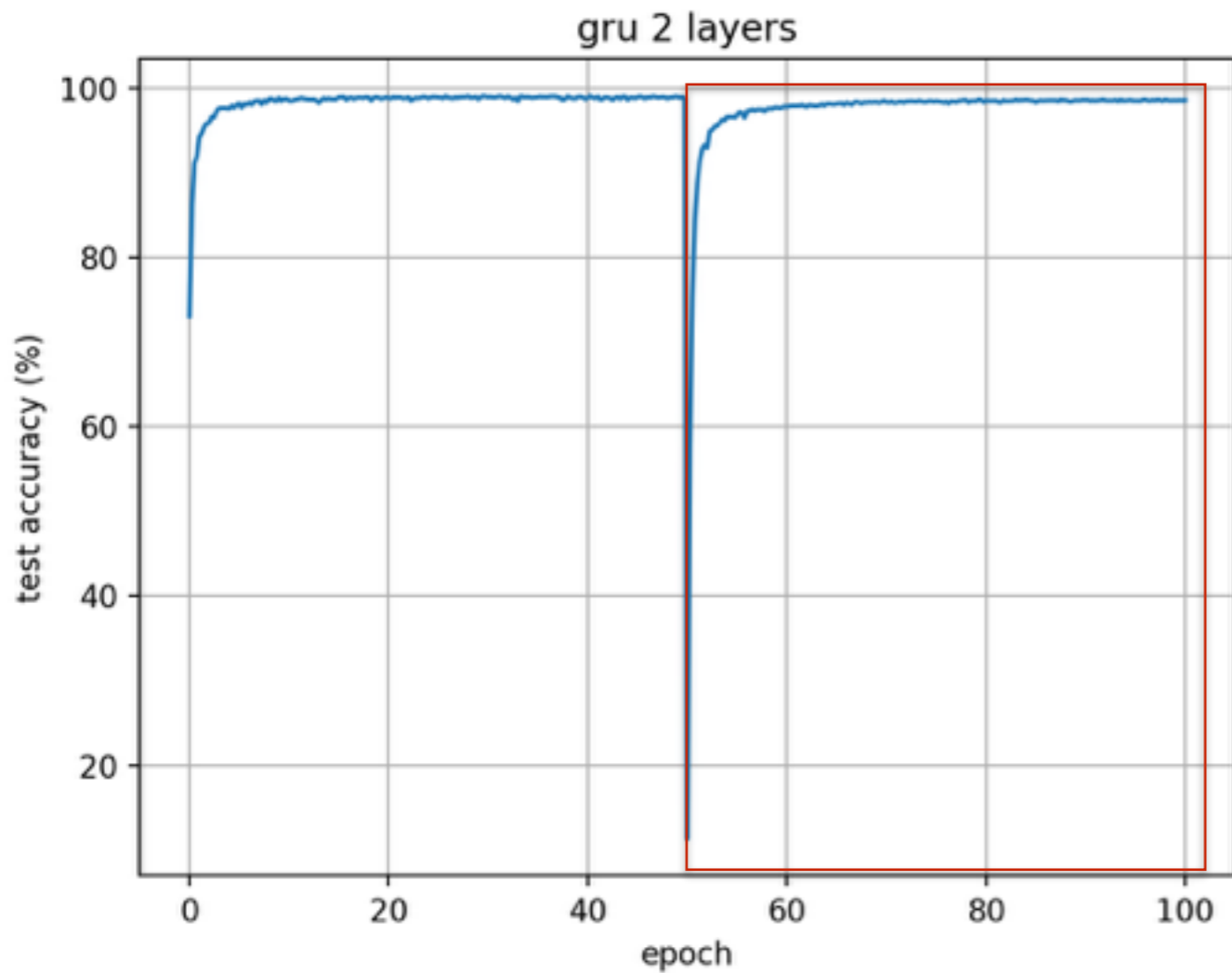
Pruning RNNs

Pruning



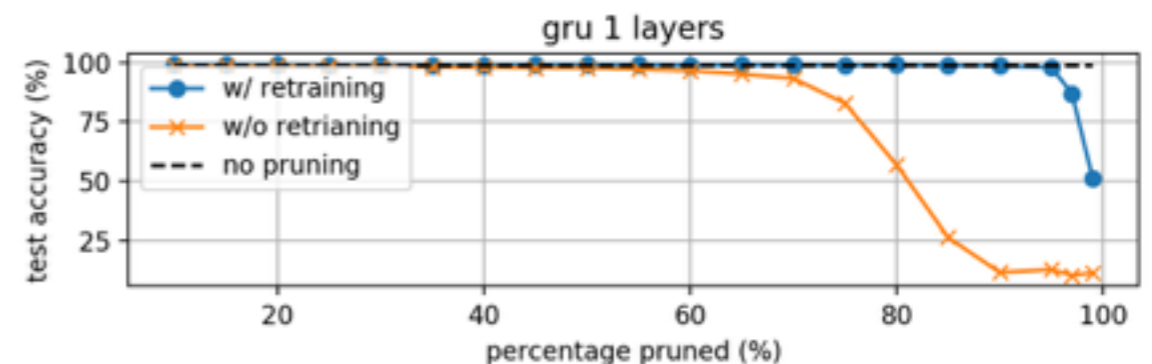
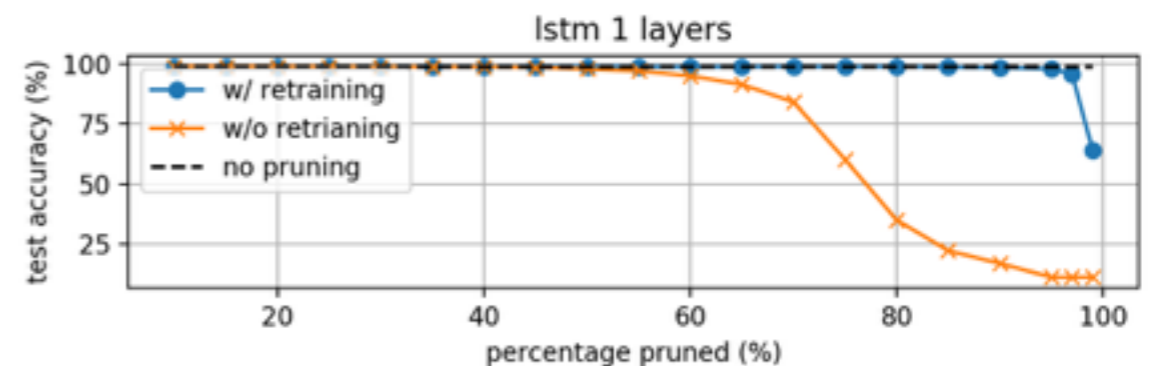
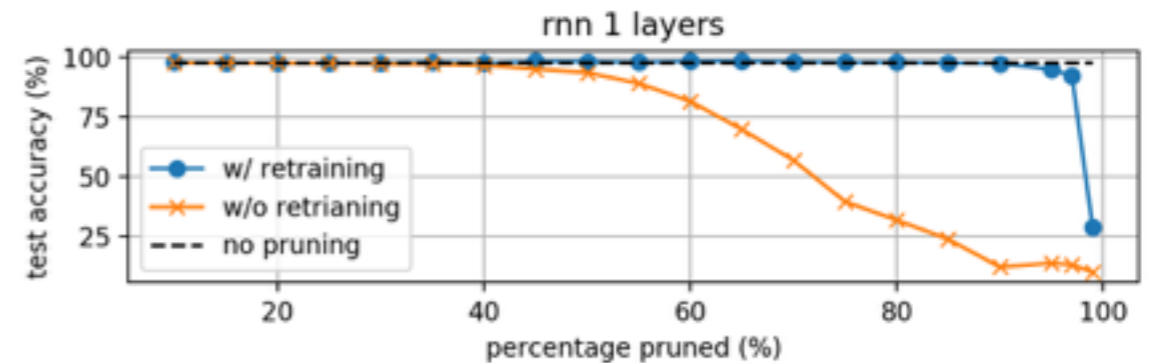
Pruning RNNs

Re-train



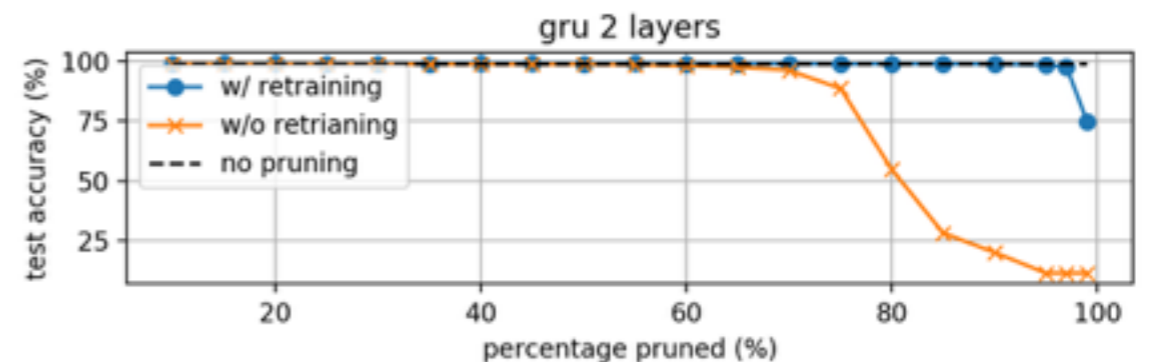
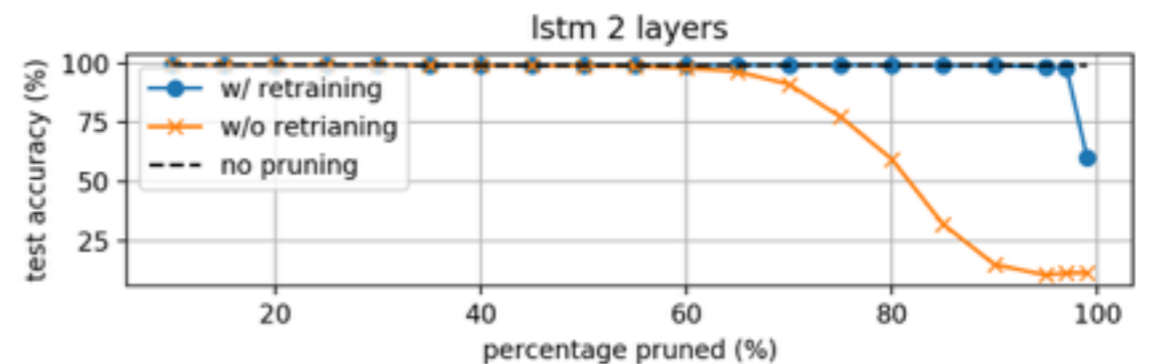
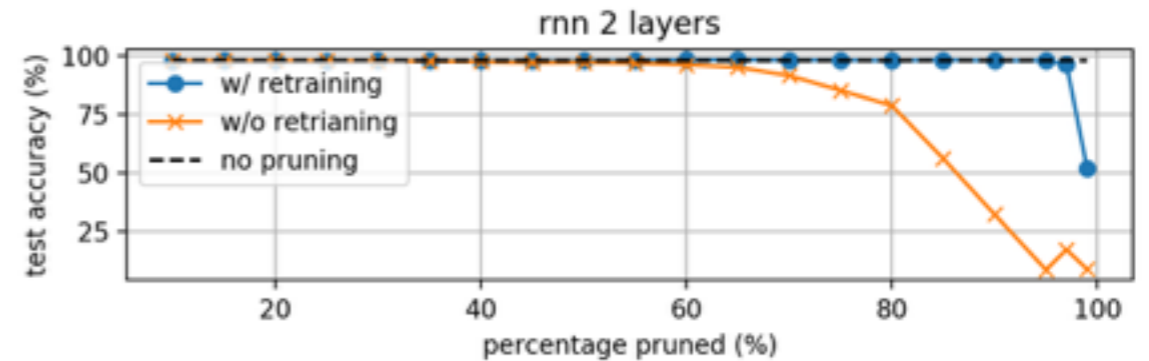
Pruning RNNs

- Retraining closes the gap
- 95% pruned - accuracy loss smaller than 1%
- LSTM and GRU more resilient to pruning - additional redundancy in gates



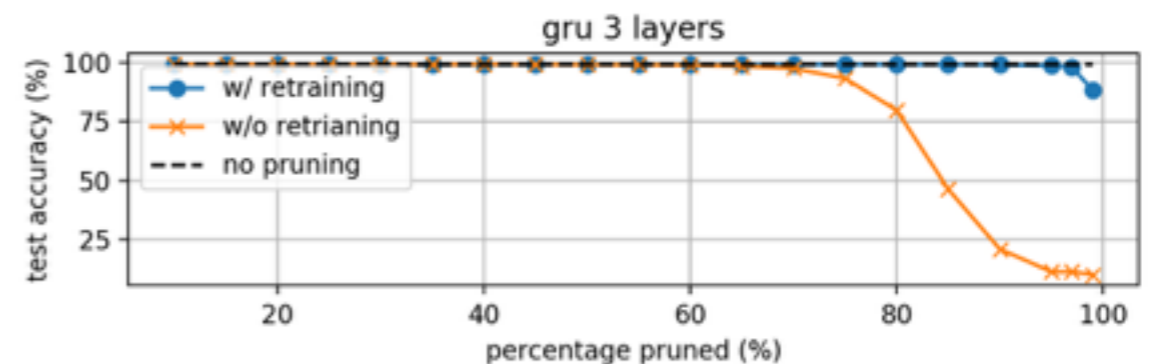
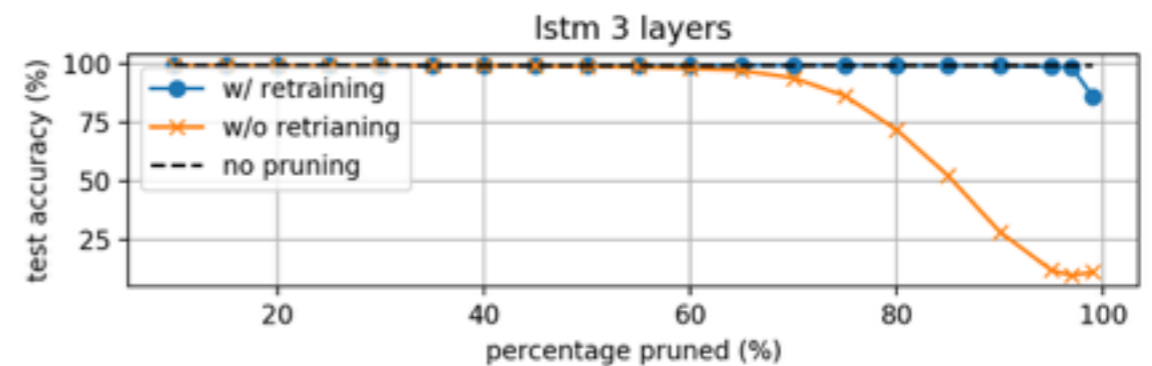
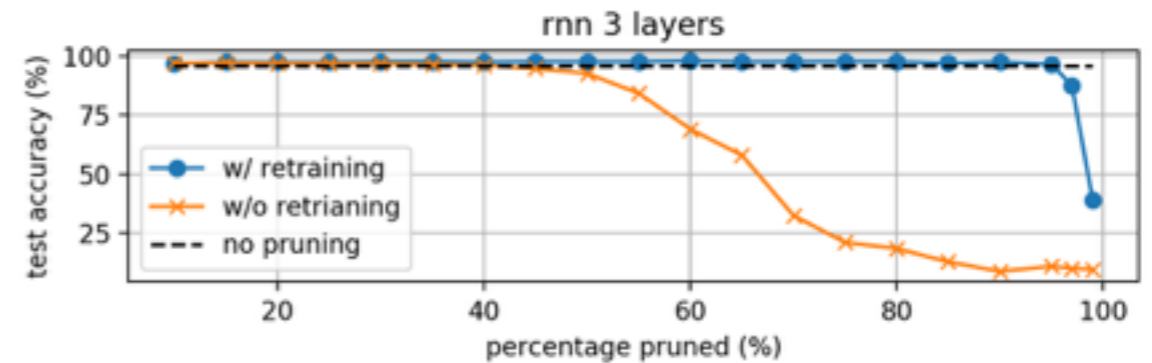
Pruning RNNs

- Retraining closes the gap
- 95% pruned - accuracy loss smaller than 1%
- LSTM and GRU more resilient to pruning - additional redundancy in gates

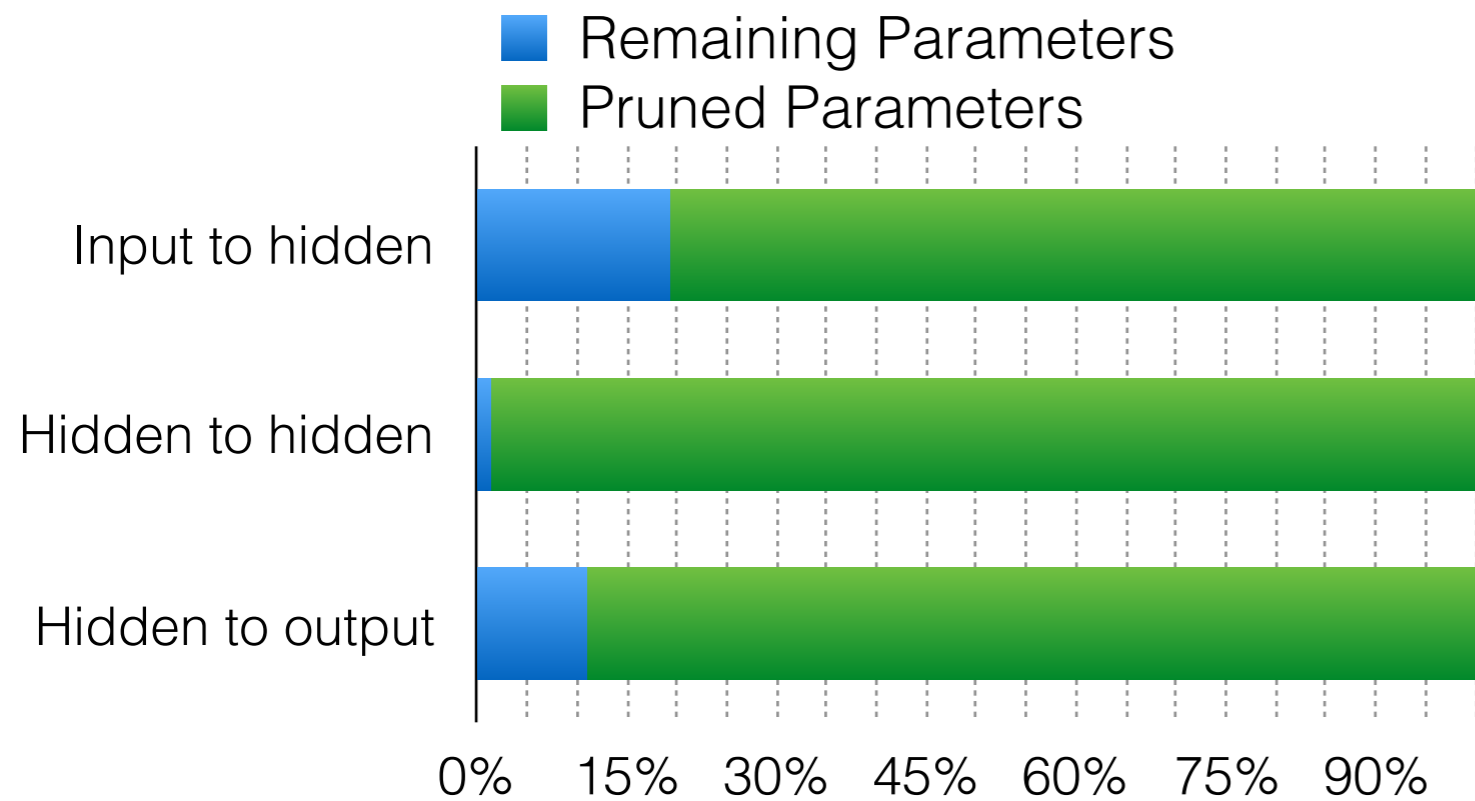


Pruning RNNs

- Retraining closes the gap
- 95% pruned - accuracy loss smaller than 1%
- LSTM and GRU more resilient to pruning - additional redundancy in gates
- More layers - more can be pruned for LSTM & GRU (RNN?)



1 Layer RNN - 95% Pruned



Hidden to hidden more redundancy

$$\mathbf{a}^{(t)} = \mathbf{b} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)}$$

$$\mathbf{h}^{(t)} = \psi(\mathbf{a}^{(t)})$$

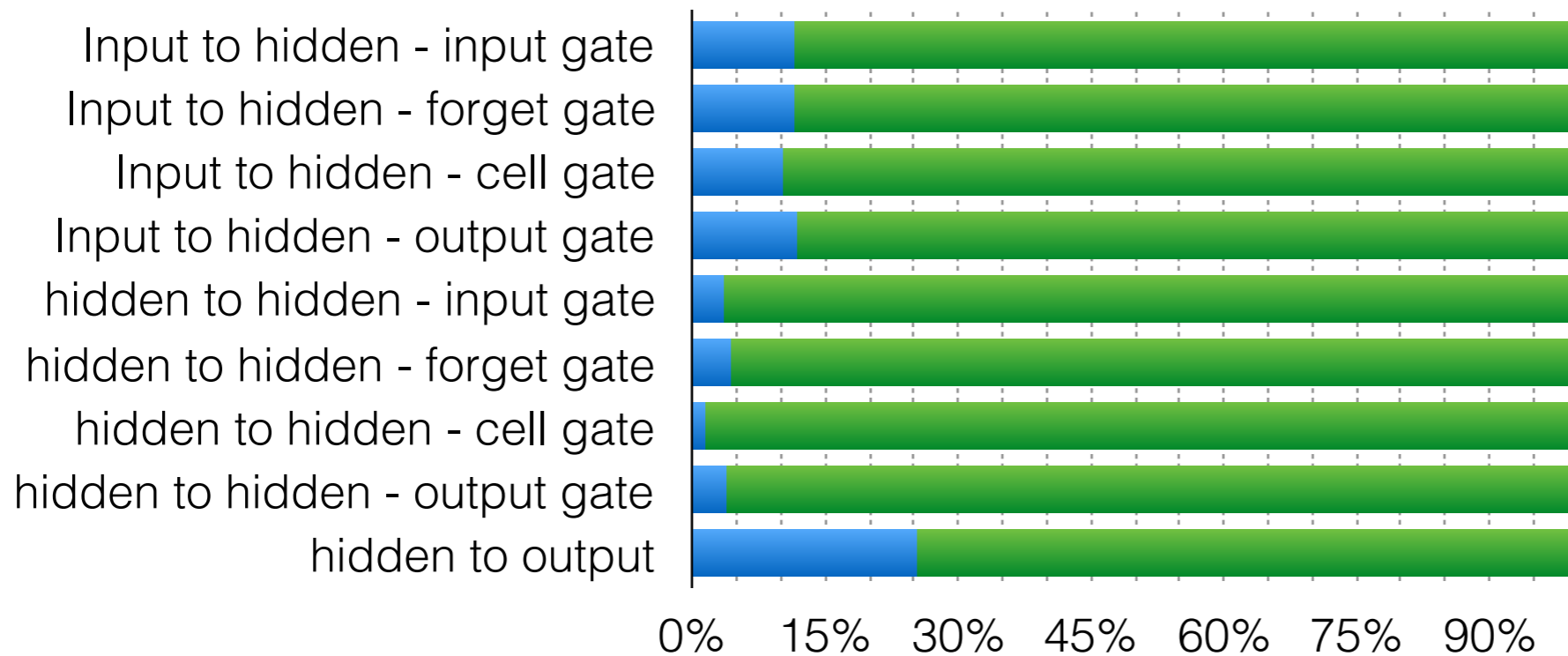
$$\mathbf{o}^{(t)} = \mathbf{c} + \mathbf{V}\mathbf{h}^{(t)}$$

$$\hat{\mathbf{y}}^{(t)} = \phi(\mathbf{o}^{(t)})$$

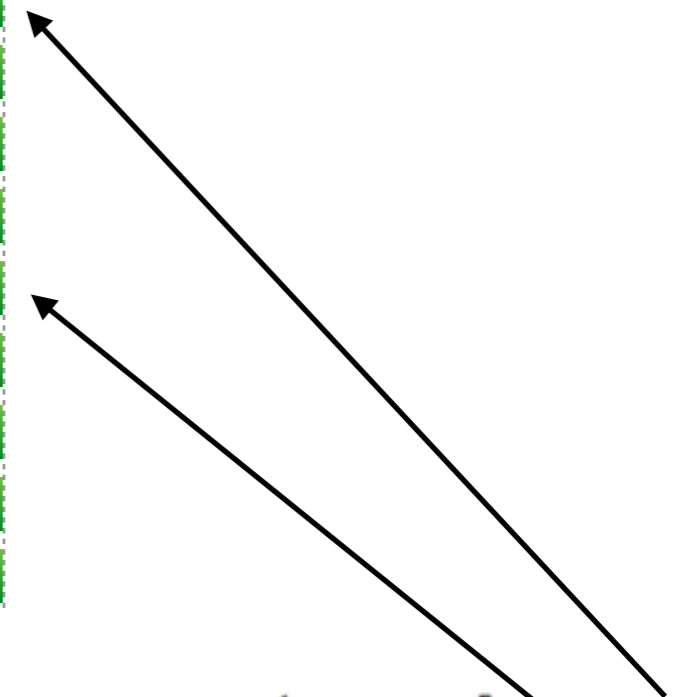
1 Layer LSTM - 95% Pruned Input Gate

■ Remaining Parameters

■ Pruned Parameters



Hidden to hidden more redundancy

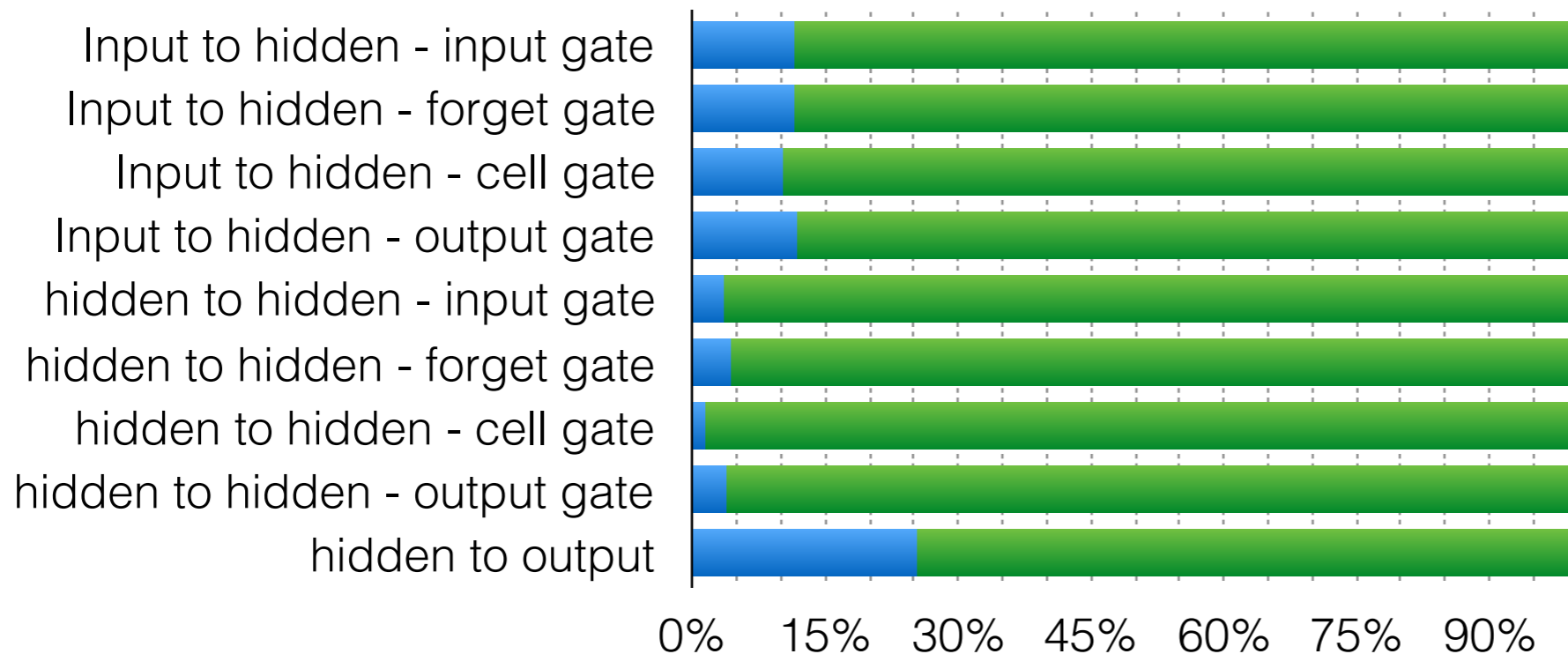


$$\begin{aligned}
 \mathbf{i}^{(t)} &= \sigma \left(\mathbf{b}^i + \mathbf{W}^i \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 \mathbf{g}^{(t)} &= \tanh \left(\mathbf{b}^g + \mathbf{W}^g \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 \mathbf{f}^{(t)} &= \sigma \left(\mathbf{b}^f + \mathbf{W}^f \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 \mathbf{c}^{(t)} &= \mathbf{f}^{(t)} \odot \mathbf{c}^{(t-1)} + \mathbf{i}^{(t)} \odot \mathbf{g}^{(t)} \\
 \mathbf{o}^{(t)} &= \sigma \left(\mathbf{b}^o + \mathbf{W}^o \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 \mathbf{h}^{(t)} &= \mathbf{o}^{(t)} \odot \tanh(\mathbf{c}^{(t)})
 \end{aligned}$$

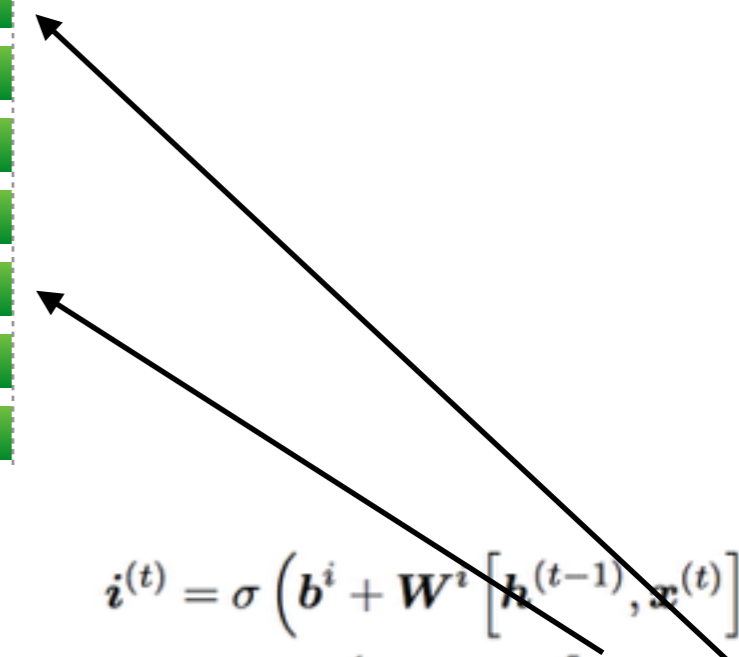
1 Layer LSTM - 95% Pruned Cell Gate

■ Remaining Parameters

■ Pruned Parameters



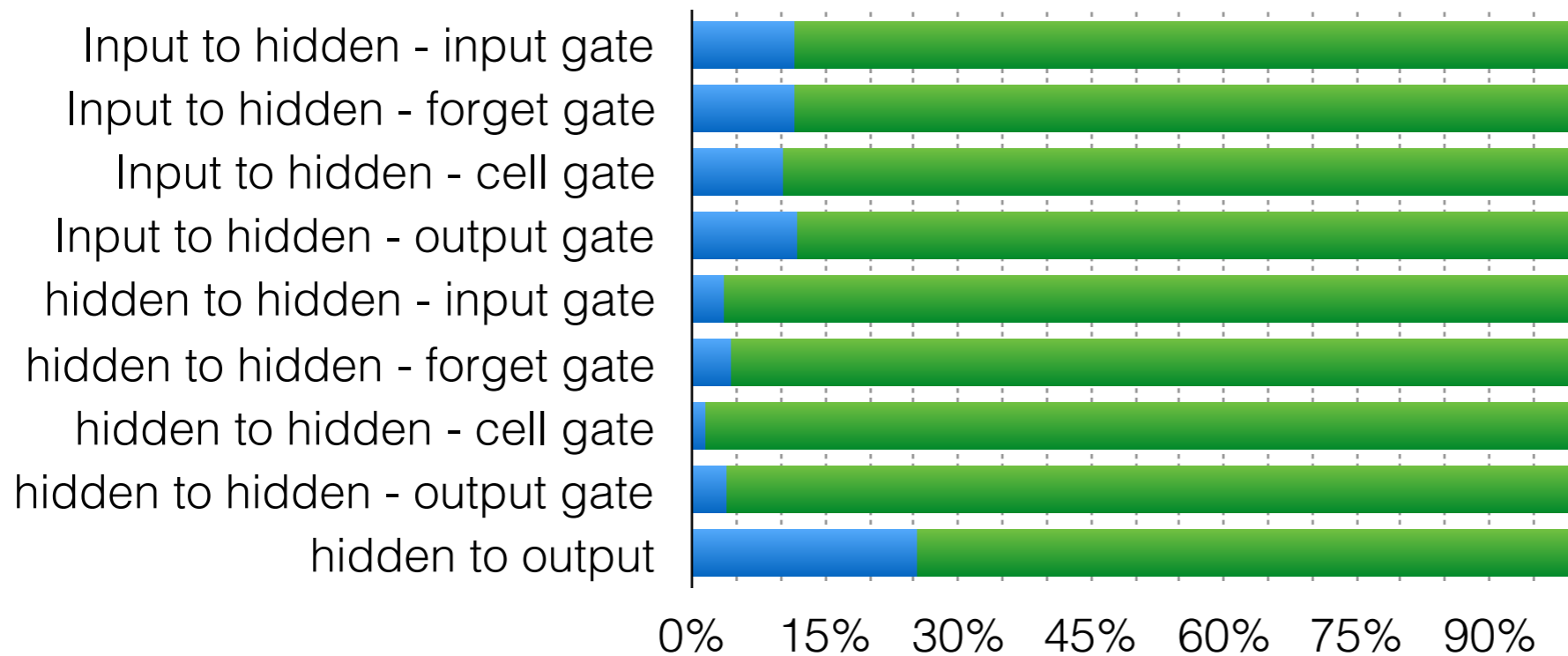
Hidden to hidden more redundancy



$$\begin{aligned}
 i^{(t)} &= \sigma \left(\mathbf{b}^i + \mathbf{W}^i \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 g^{(t)} &= \tanh \left(\mathbf{b}^g + \mathbf{W}^g \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 f^{(t)} &= \sigma \left(\mathbf{b}^f + \mathbf{W}^f \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 \mathbf{c}^{(t)} &= f^{(t)} \odot \mathbf{c}^{(t-1)} + i^{(t)} \odot g^{(t)} \\
 \mathbf{o}^{(t)} &= \sigma \left(\mathbf{b}^o + \mathbf{W}^o \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 \mathbf{h}^{(t)} &= \mathbf{o}^{(t)} \odot \tanh(\mathbf{c}^{(t)})
 \end{aligned}$$

1 Layer LSTM - 95% Pruned Forget Gate

■ Remaining Parameters ■ Pruned Parameters



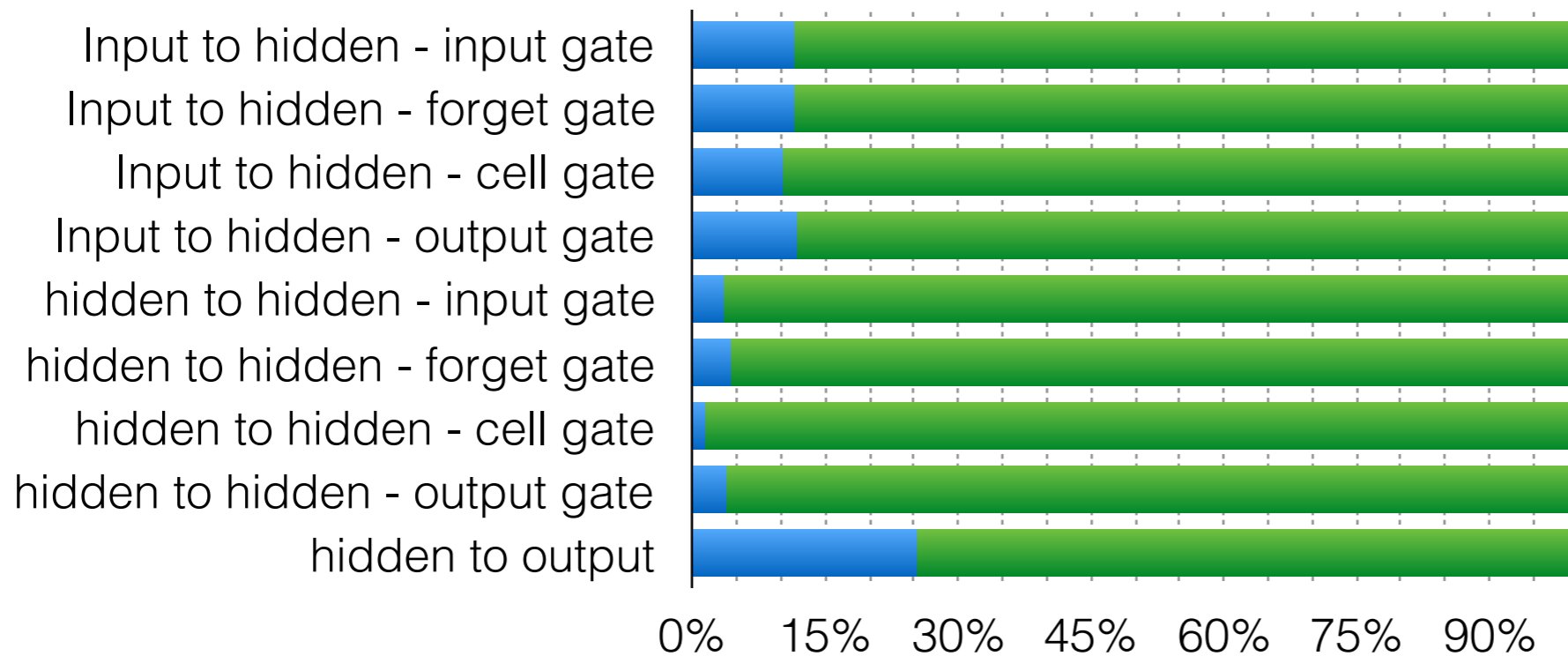
Hidden to hidden more redundancy

$$\begin{aligned}
 i^{(t)} &= \sigma \left(\mathbf{b}^i + \mathbf{W}^i \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 g^{(t)} &= \tanh \left(\mathbf{b}^g + \mathbf{W}^g \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 f^{(t)} &= \sigma \left(\mathbf{b}^f + \mathbf{W}^f \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 \mathbf{c}^{(t)} &= f^{(t)} \odot \mathbf{c}^{(t-1)} + i^{(t)} \odot g^{(t)} \\
 \mathbf{o}^{(t)} &= \sigma \left(\mathbf{b}^o + \mathbf{W}^o \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 \mathbf{h}^{(t)} &= \mathbf{o}^{(t)} \odot \tanh(\mathbf{c}^{(t)})
 \end{aligned}$$

1 Layer LSTM - 95% Pruned Output Gate

■ Remaining Parameters

■ Pruned Parameters

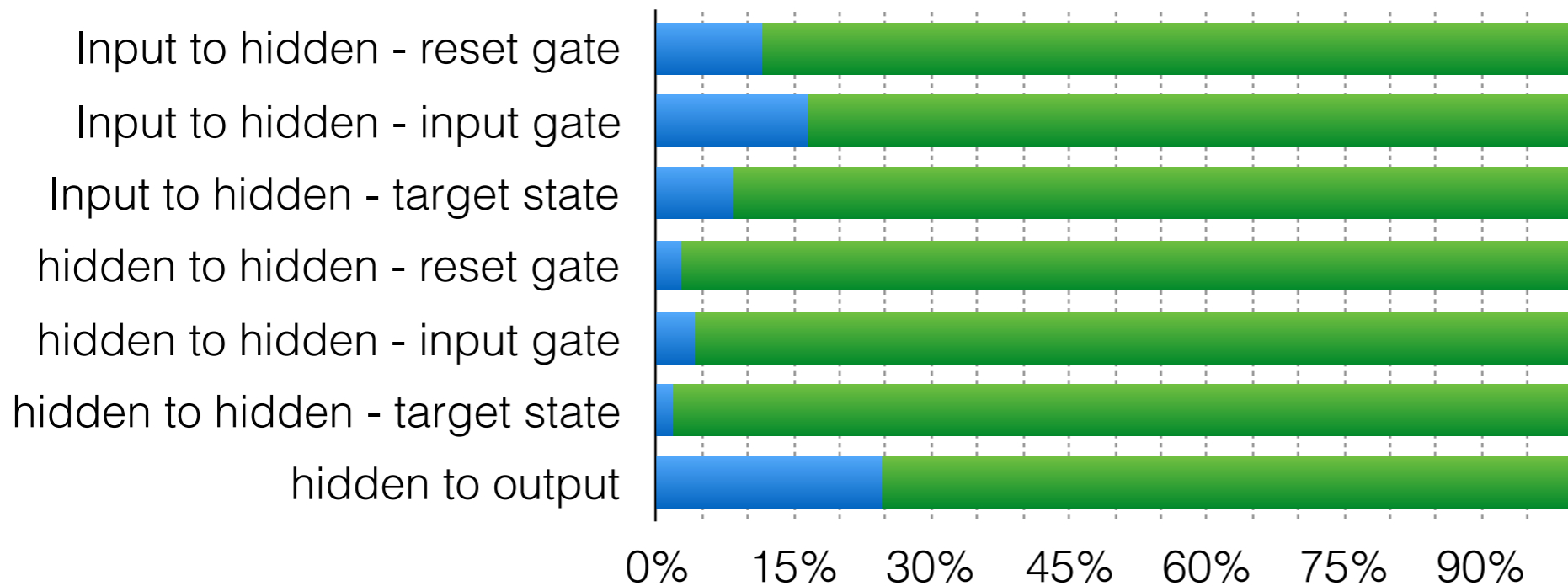


Hidden to hidden more redundancy

$$\begin{aligned}
 i^{(t)} &= \sigma \left(\mathbf{b}^i + \mathbf{W}^i \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 g^{(t)} &= \tanh \left(\mathbf{b}^g + \mathbf{W}^g \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 f^{(t)} &= \sigma \left(\mathbf{b}^f + \mathbf{W}^f \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 \mathbf{c}^{(t)} &= f^{(t)} \odot \mathbf{c}^{(t-1)} + i^{(t)} \odot g^{(t)} \\
 \mathbf{o}^{(t)} &= \sigma \left(\mathbf{b}^o + \mathbf{W}^o \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 \mathbf{h}^{(t)} &= \mathbf{o}^{(t)} \odot \tanh(\mathbf{c}^{(t)})
 \end{aligned}$$

1 Layer GRU - 95% Pruned Input Gate

■ Remaining Parameters ■ Pruned Parameters

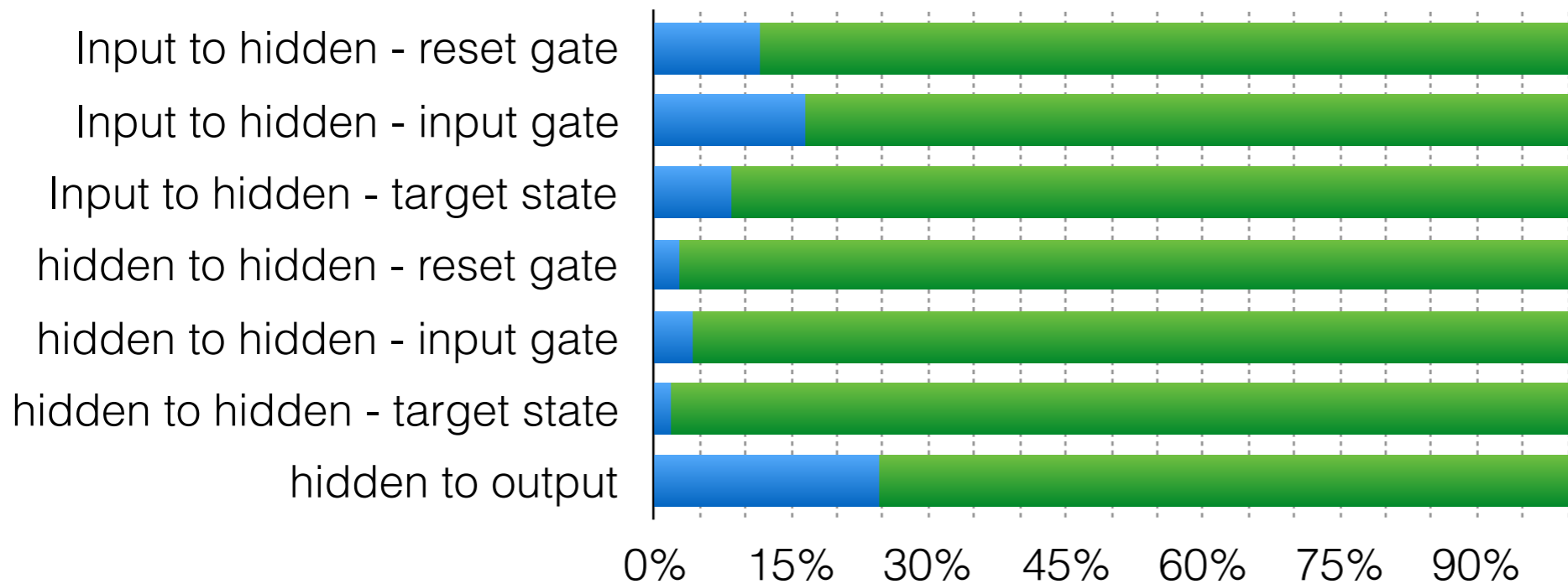


Hidden to hidden more redundancy
Input gate less redundancy

$$\begin{aligned}
 \mathbf{z}^{(t)} &= \sigma \left(\mathbf{b}^z + \mathbf{W}^z \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 \mathbf{r}^{(t)} &= \sigma \left(\mathbf{b}^r + \mathbf{W}^r \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 \mathbf{g}^{(t)} &= \tanh \left(\mathbf{b}^g + \mathbf{W}^g \left[\mathbf{r}^{(t)} \odot \mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 \mathbf{h}^{(t)} &= (1 - \mathbf{z}^{(t)}) \odot \mathbf{h}^{(t-1)} + \mathbf{z}^{(t)} \odot \mathbf{g}^{(t)}
 \end{aligned}$$

1 Layer GRU - 95% Pruned Reset Gate

■ Remaining Parameters ■ Pruned Parameters

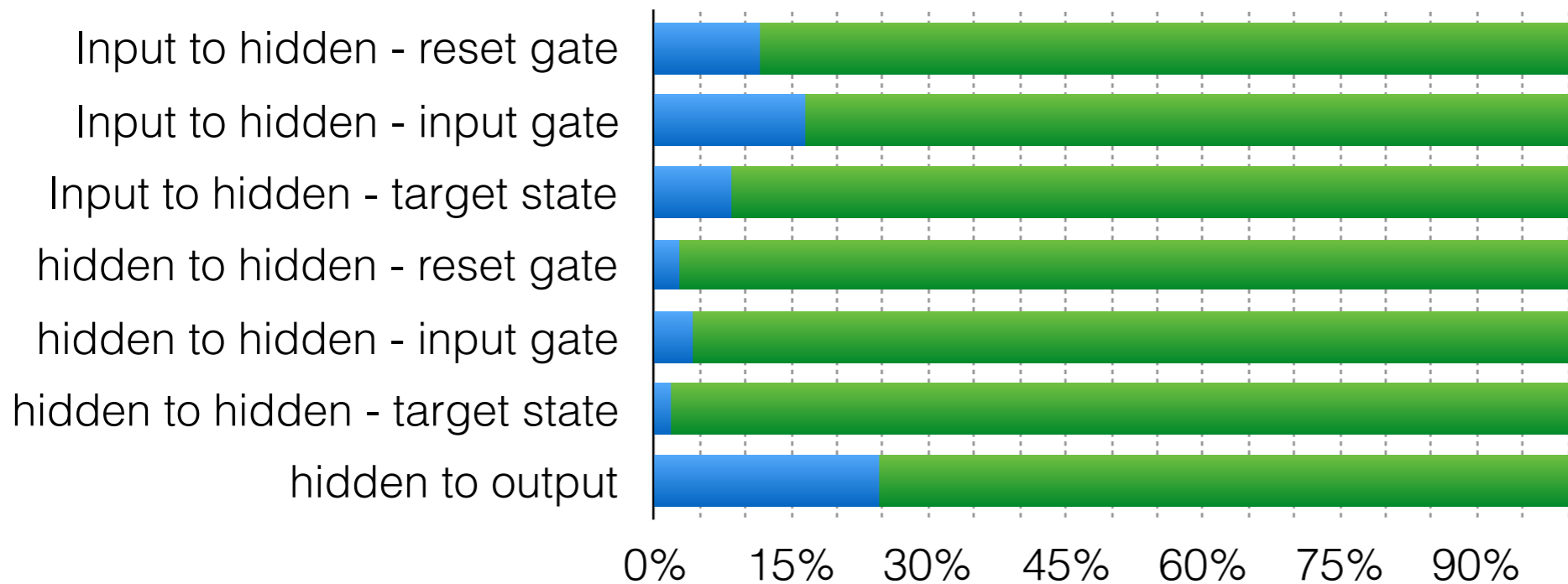


Hidden to hidden more redundancy
Input gate less redundancy

$$\begin{aligned}
 \mathbf{z}^{(t)} &= \sigma \left(\mathbf{b}^z + \mathbf{W}^z \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 \mathbf{r}^{(t)} &= \sigma \left(\mathbf{b}^r + \mathbf{W}^r \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 \mathbf{g}^{(t)} &= \tanh \left(\mathbf{b}^g + \mathbf{W}^g \left[\mathbf{r}^{(t)} \odot \mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 \mathbf{h}^{(t)} &= (1 - \mathbf{z}^{(t)}) \odot \mathbf{h}^{(t-1)} + \mathbf{z}^{(t)} \odot \mathbf{g}^{(t)}
 \end{aligned}$$

1 Layer GRU - 95% Pruned Target State

■ Remaining Parameters ■ Pruned Parameters



Hidden to hidden more redundancy
Input gate less redundancy

$$\begin{aligned}
 \mathbf{z}^{(t)} &= \sigma \left(\mathbf{b}^z + \mathbf{W}^z \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 \mathbf{r}^{(t)} &= \sigma \left(\mathbf{b}^r + \mathbf{W}^r \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 \mathbf{g}^{(t)} &= \tanh \left(\mathbf{b}^g + \mathbf{W}^g \left[\mathbf{r}^{(t)} \odot \mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right) \\
 \mathbf{h}^{(t)} &= (1 - \mathbf{z}^{(t)}) \odot \mathbf{h}^{(t-1)} + \mathbf{z}^{(t)} \odot \mathbf{g}^{(t)}
 \end{aligned}$$

Summary

- On MNIST, RNN, LSTM, GRU can be pruned by 95% without significant accuracy loss
- LSTM and GRU have more redundancy than RNN
- Hidden to hidden layers have more redundancy