

ELEC/COMP 576: Gated Recurrent Neural Network Architectures & Applications

Ankit B. Patel

Baylor College of Medicine (Neuroscience Dept.)

Rice University (ECE Dept.)

RNN Training

Why Training is Unstable

$$x^{(l)} = W^{(l-1)}y^{(l-1)} + b^{(l-1)}$$

$$y^{(l)} = f(x^{(l)})$$

Let the activation function $f(x) = \alpha x + \beta$,

$$\text{Var} \left(y^{(l)} \right) = \alpha^2 n_{l-1} \sigma_{l-1}^2 \left(\text{Var} \left(y^{(l-1)} \right) + \beta^2 I_{n_l} \right) .$$

$$\text{Var} \left(\frac{\partial \text{cost}}{\partial y^{(l-1)}} \right) = \alpha^2 n_l \sigma_{l-1}^2 \text{Var} \left(\frac{\partial \text{cost}}{\partial y^{(l)}} \right) .$$

Variance of activations/gradients grows multiplicatively

Interesting Question

- Are there modifications to an RNN such that it can combat these gradient problems?

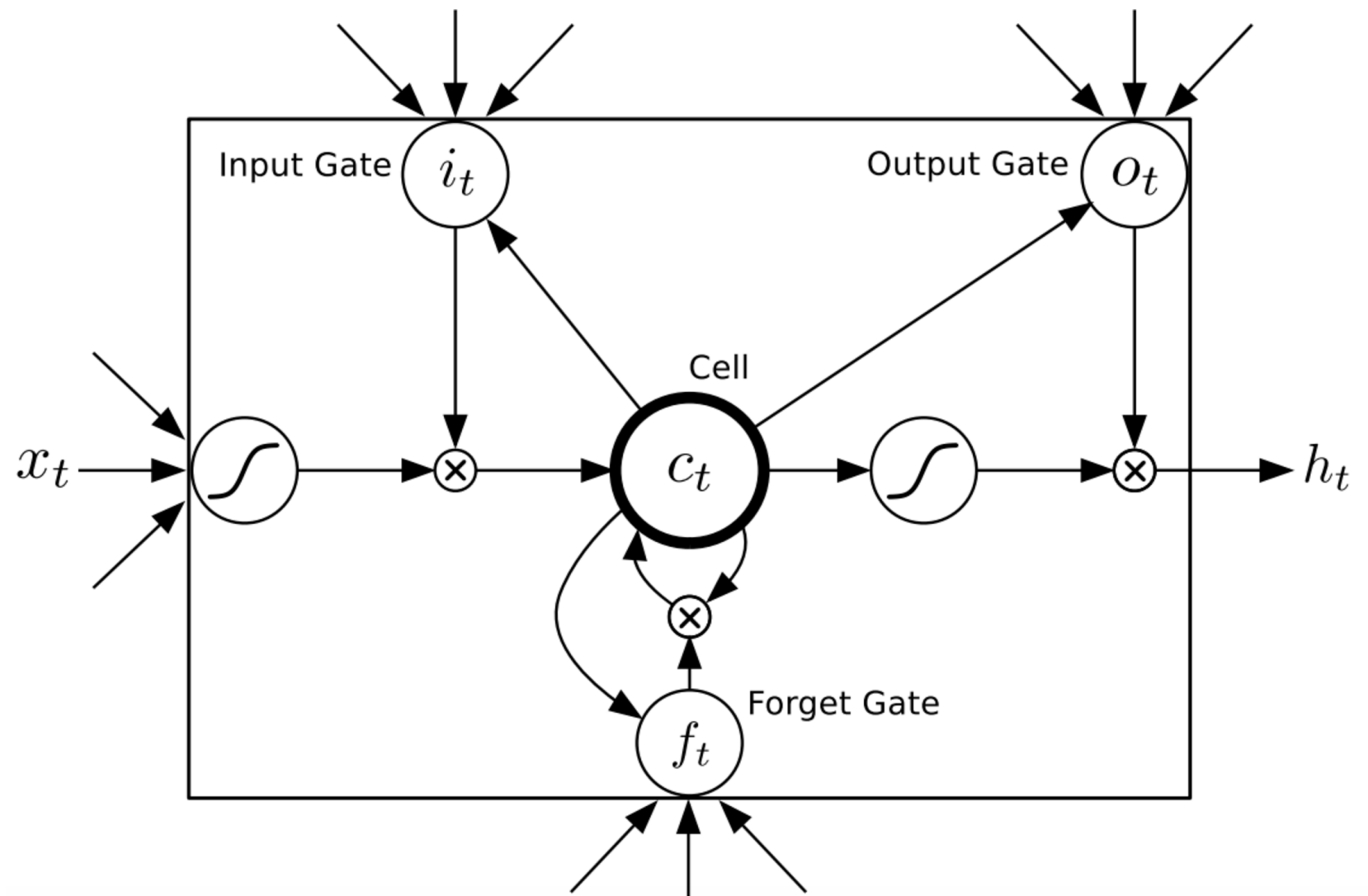
RNNs with Longer Term Memory

Motivation

- The need to remember certain events for arbitrarily long periods of time (Non-Markovian)
- The need to forget certain events

Long Short Term Memory

- 3 gates
 - Input
 - Forget
 - Output



[Zygmunt Z.]

LSTM Formulation

$$i_t = \sigma (W_{xi}x_t + W_{hi}h_{t-1} + W_{ci}c_{t-1} + b_i)$$

$$f_t = \sigma (W_{xf}x_t + W_{hf}h_{t-1} + W_{cf}c_{t-1} + b_f)$$

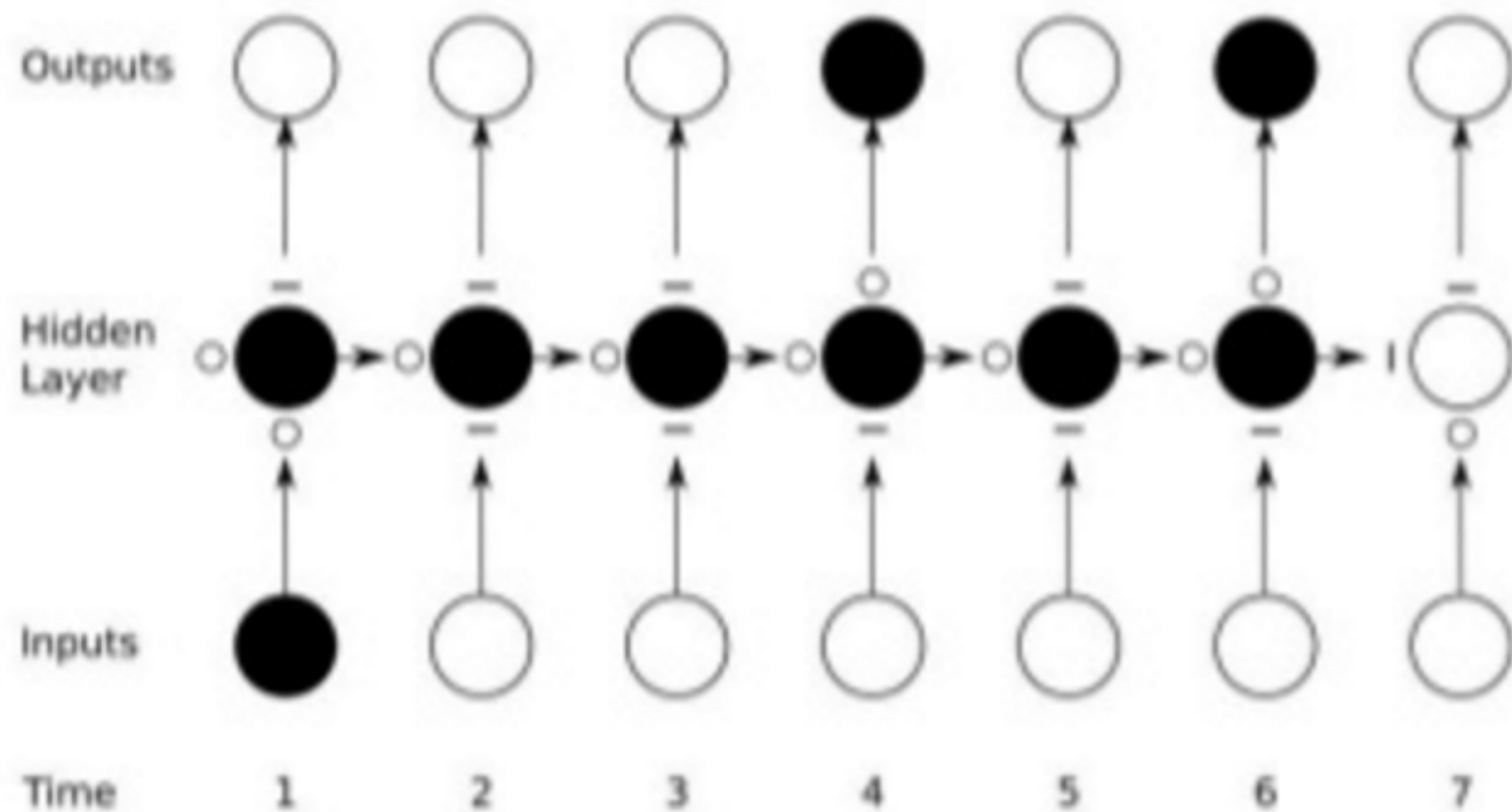
$$c_t = f_t c_{t-1} + i_t \tanh (W_{xc}x_t + W_{hc}h_{t-1} + b_c)$$

$$o_t = \sigma (W_{xo}x_t + W_{ho}h_{t-1} + W_{co}c_t + b_o)$$

$$h_t = o_t \tanh(c_t)$$

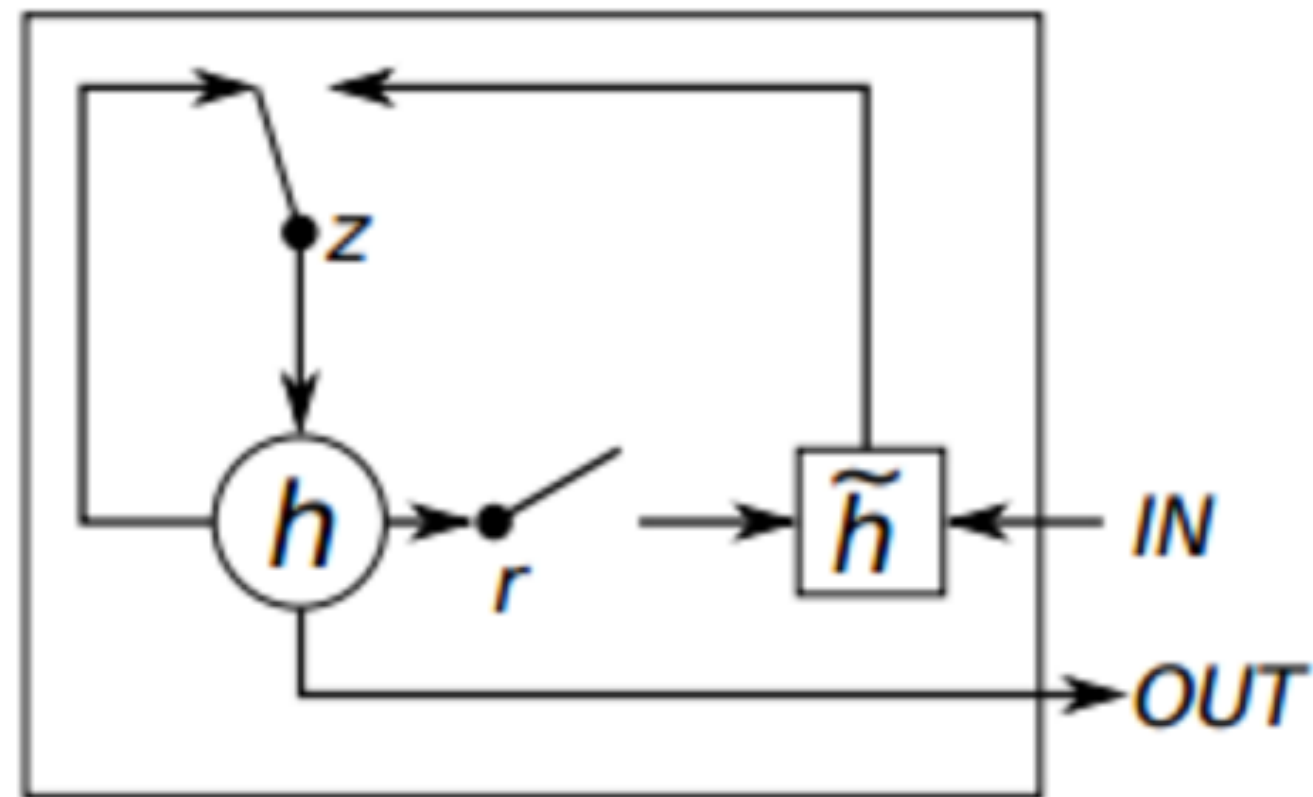
$$y_t = W_{ho}h_t + b_o$$

Preserving Gradients



Gated Recurrent Unit

- 2 gates
 - Reset
 - Combine new input with previous memory
 - Update
 - How long the previous memory should stay



[Zygmunt Z.]

GRU Formulation

$$z = \sigma(x_t U^z + s_{t-1} W^z)$$

$$r = \sigma(x_t U^r + s_{t-1} W^r)$$

$$h = \tanh(x_t U^h + (s_{t-1} \circ r) W^h)$$

$$s_t = (1 - z) \circ h + z \circ s_{t-1}$$

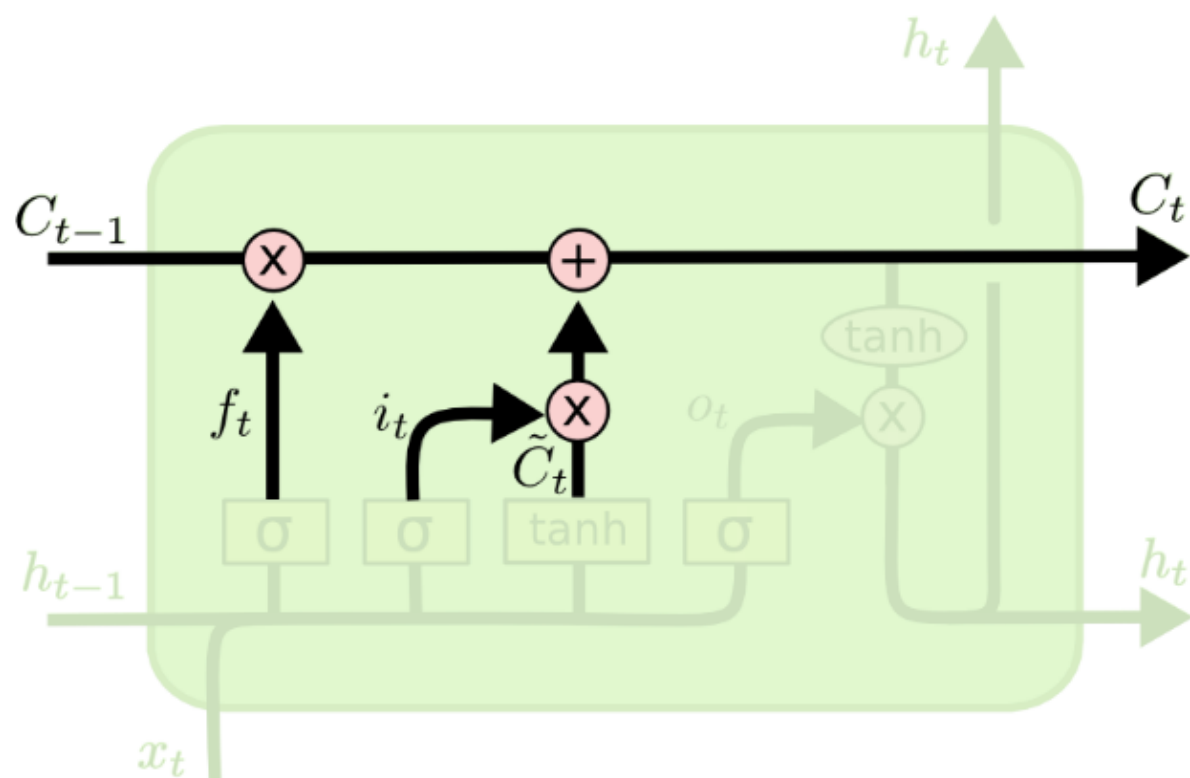
LSTM & GRU Benefits

- Remember for longer temporal durations
 - RNN has issues for remembering longer durations
- Able to have feedback flow at different strengths depending on inputs

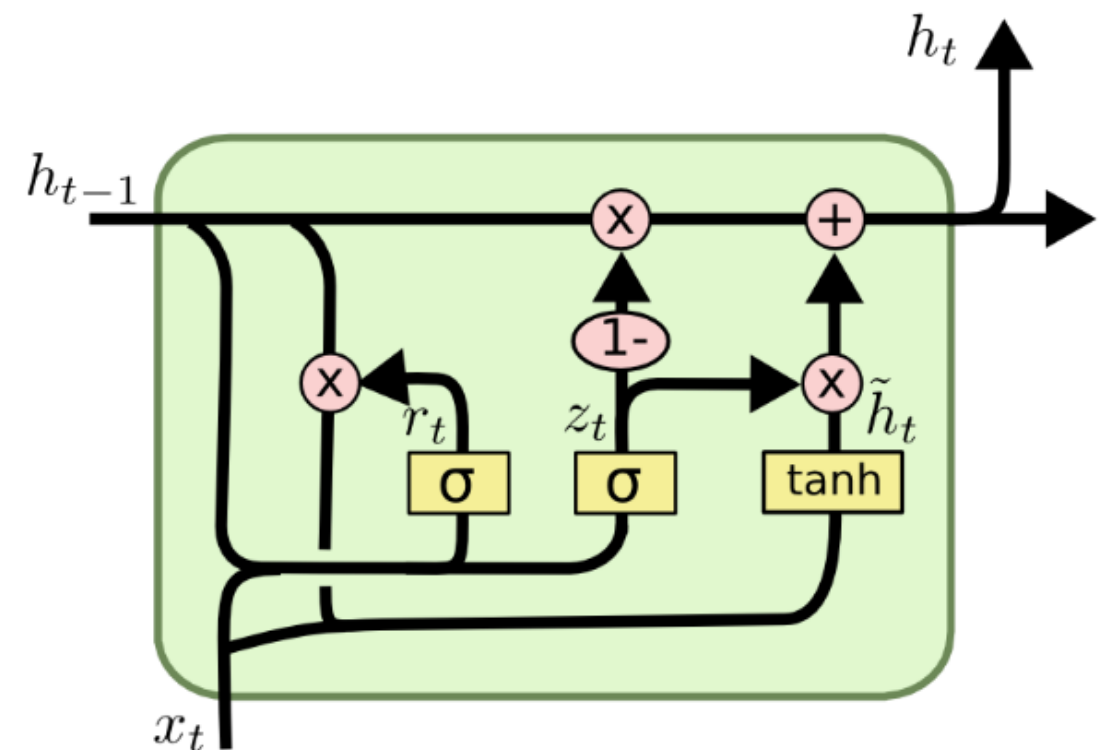
Differences between LSTM & GRU

- GRU has two gates, while LSTM has three gates
- GRU does not have internal memory
- GRU does not use a second nonlinearity for computing the output

Visual Difference of LSTM & GRU



LSTM



GRU

LSTM vs GRU Results

			tanh	GRU	LSTM
Music Datasets	Nottingham	train	3.22	2.79	3.08
		test	3.13	3.23	3.20
	JSB Chorales	train	8.82	6.94	8.15
		test	9.10	8.54	8.67
	MuseData	train	5.64	5.06	5.18
		test	6.23	5.99	6.23
	Piano-midi	train	5.64	4.93	6.49
		test	9.03	8.82	9.03
Ubisoft Datasets	Ubisoft dataset A	train	6.29	2.31	1.44
		test	6.44	3.59	2.70
	Ubisoft dataset B	train	7.61	0.38	0.80
		test	7.62	0.88	1.26

Other Methods for Stabilizing RNN Training

Why Training is Unstable

$$x^{(l)} = W^{(l-1)}y^{(l-1)} + b^{(l-1)}$$

$$y^{(l)} = f(x^{(l)})$$

Let the activation function $f(x) = \alpha x + \beta$,

$$\text{Var} \left(y^{(l)} \right) = \alpha^2 n_{l-1} \sigma_{l-1}^2 \left(\text{Var} \left(y^{(l-1)} \right) + \beta^2 I_{n_l} \right) .$$

$$\text{Var} \left(\frac{\partial \text{cost}}{\partial y^{(l-1)}} \right) = \alpha^2 n_l \sigma_{l-1}^2 \text{Var} \left(\frac{\partial \text{cost}}{\partial y^{(l)}} \right) .$$

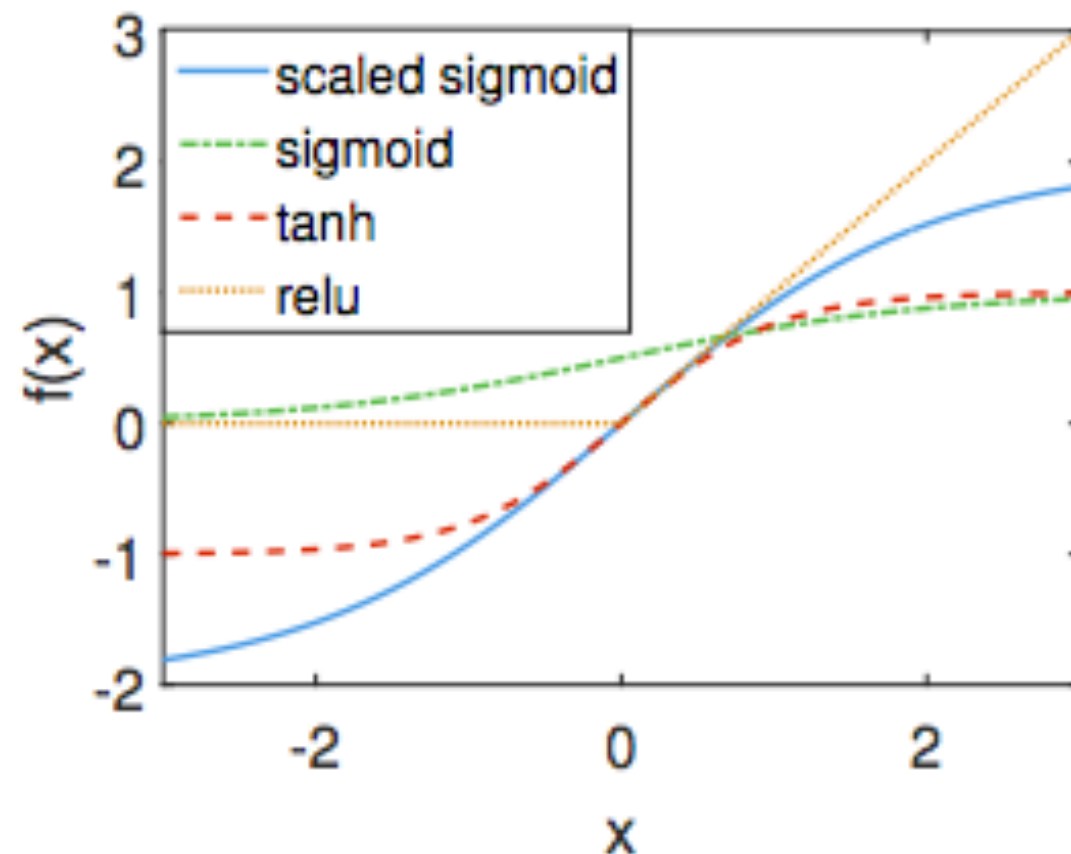
Variance of activations/gradients grows multiplicatively

Stabilizing Activations & Gradients

$$\text{Var} \left(y^{(l)} \right) = \text{Var} \left(y^{(l-1)} \right) \quad \text{and} \quad \text{Var} \left(\frac{\partial \text{cost}}{\partial y^{(l)}} \right) = \text{Var} \left(\frac{\partial \text{cost}}{\partial y^{(l-1)}} \right) ;$$
$$n_l \sigma_{l-1}^2 \cong 1 \quad \text{and} \quad n_{l-1} \sigma_{l-1}^2 \cong 1;$$

We want $\alpha = 1$ and $\beta = 0$.

Taylor Expansions of Different Activation Functions



$$\text{sigmoid}(x) = \frac{1}{2} + \frac{x}{4} - \frac{x^3}{48} + O(x^5)$$

$$\tanh(x) = 0 + x - \frac{x^3}{3} + O(x^5)$$

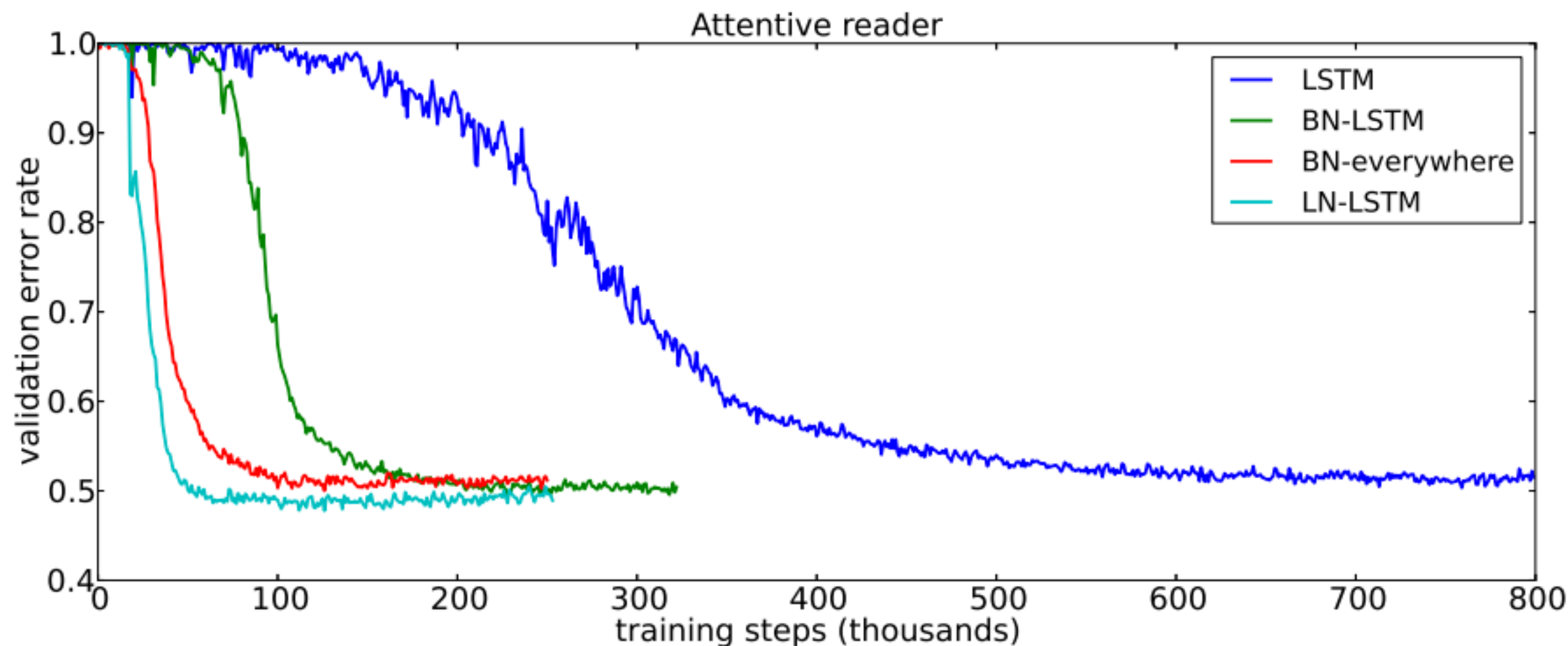
$$\text{relu}(x) = 0 + x \quad \text{for } x \geq 0.$$

Layer Normalization

- Similar to batch normalization
- Apply it to RNNs to stabilize the hidden state dynamics

$$\mathbf{h}^t = f \left[\frac{\mathbf{g}}{\sigma^t} \odot (\mathbf{a}^t - \mu^t) + \mathbf{b} \right] \quad \mu^t = \frac{1}{H} \sum_{i=1}^H a_i^t \quad \sigma^t = \sqrt{\frac{1}{H} \sum_{i=1}^H (a_i^t - \mu^t)^2}$$

Layer Normalization Results



[Ba, Kiros, Hinton]

Streamlining Normalization

Algorithm 5 Streaming Normalization Layer: Forward

Require: layer input \mathbf{x} (a mini-batch), NormOP $N(.,.)$, function $S(.)$ to compute NormStats for every element of \mathbf{x} , running estimates of NormStats and/or related information packed in a structure/table H_1 , function F to update H_1 and generate current estimate of NormStats \hat{s} .

Ensure: layer output \mathbf{y} (a mini-batch) and H_1 (it is stored in this layer, instead of feeding to other layers), always maintain the latest \hat{s} in case of testing

if *training* **then**

$s = S(\mathbf{x})$

$\{H_1, \hat{s}\} = F(H_1, s)$

$\mathbf{y} = N(\mathbf{x}, \hat{s})$

else *testing*

$\mathbf{y} = N(\mathbf{x}, \hat{s})$

end if

Streamlining Normalization

Algorithm 6 Streaming Normalization Layer: Backpropagation

Require: $\frac{\partial E}{\partial y}$ (a mini-batch) where E is objective, layer input \mathbf{x} (a mini-batch), NormOP $N(.,.)$, function $S(.)$ to compute NormStats for every element of \mathbf{x} , running estimates of NormStats \hat{s} , running estimates of gradients and/or related information packed in a structure/table H_2 , function G to update H_2 and generate the current estimates of gradients of NormStats $\widehat{\frac{\partial E}{\partial \hat{s}}}$.

Ensure: $\frac{\partial E}{\partial x}$ (a mini-batch) and H_2 (it is stored in this layer, instead of feeding to other layers)

$\frac{\partial E}{\partial \hat{s}}$ is calculated using chain rule.

$\{H_2, \widehat{\frac{\partial E}{\partial \hat{s}}}\} = G(H_2, \frac{\partial E}{\partial \hat{s}})$

Use $\widehat{\frac{\partial E}{\partial \hat{s}}}$ for further backpropagation, instead of $\frac{\partial E}{\partial \hat{s}}$

$\frac{\partial E}{\partial x}$ is calculated using chain rule.

Streamlining Normalization for RNNs

Normalized RNN

$$h_t = \text{NonLinear}(\text{Norm}(W_x * x_t) + \text{Norm}(W_h * h_{t-1}))$$

Normalized GRU

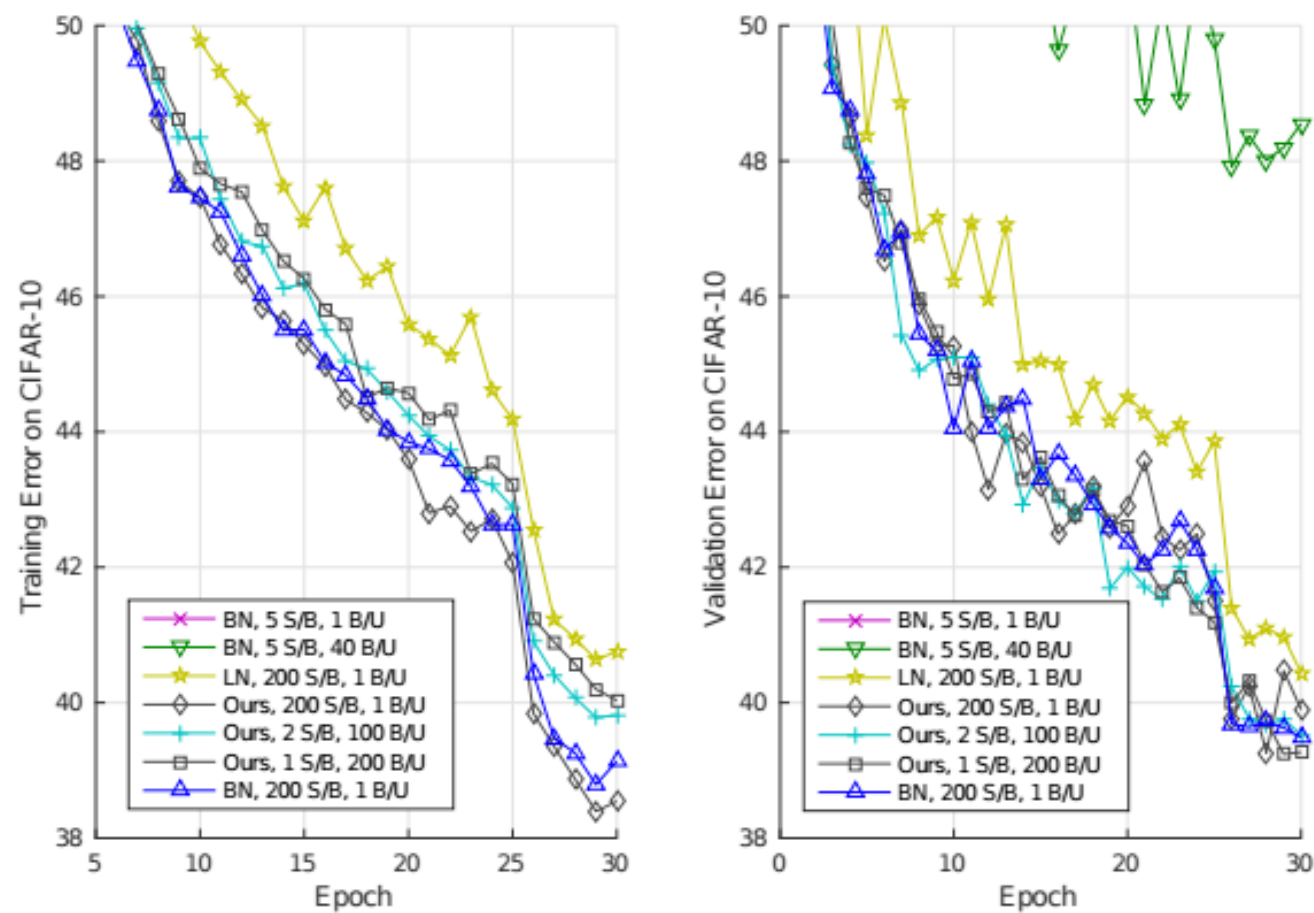
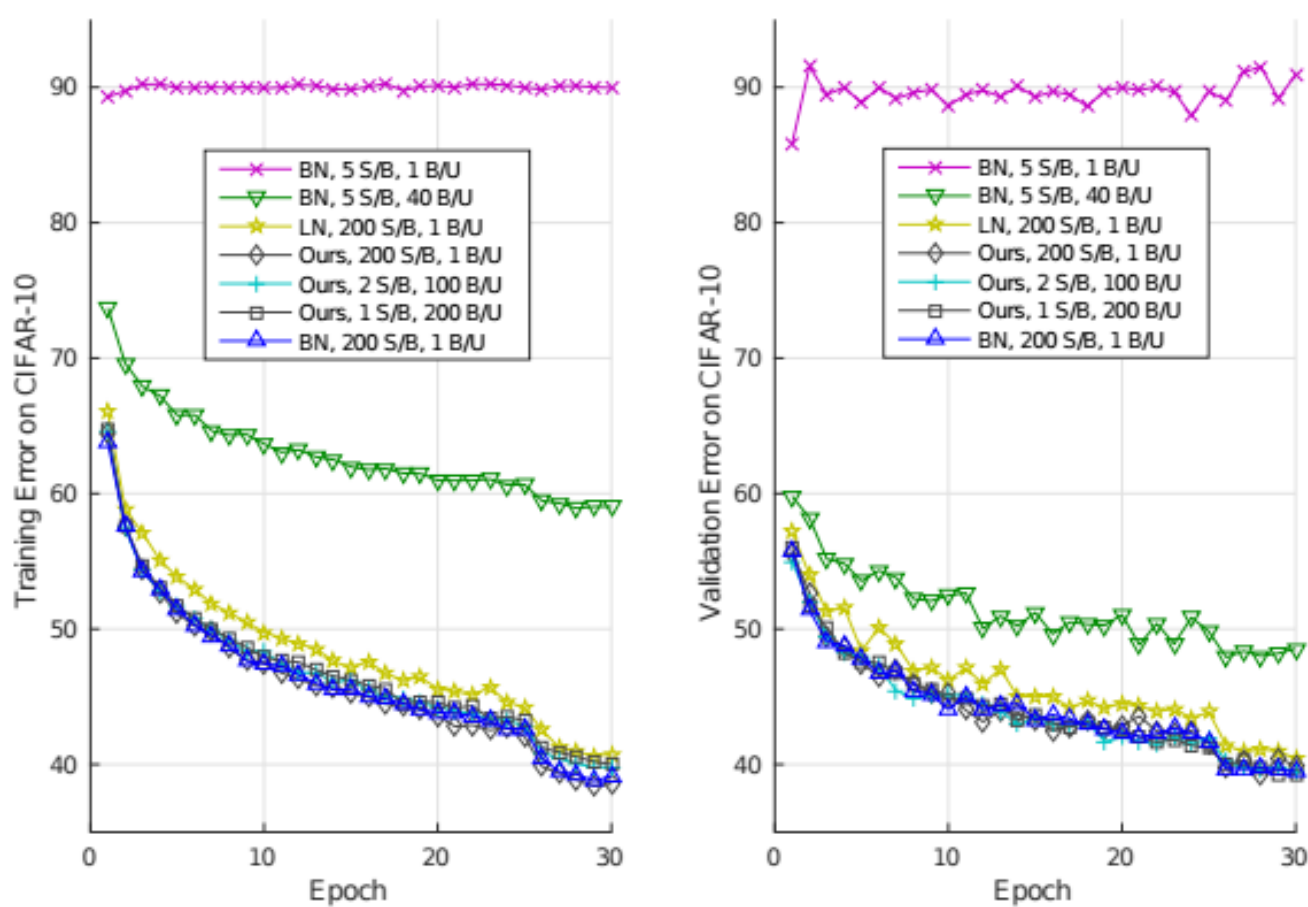
$$g_r = \text{Sigmoid}(\text{Norm}(W_{xr} * x_t) + \text{Norm}(W_{hr} * h_{t-1}))$$

$$g_z = \text{Sigmoid}(\text{Norm}(W_{xz} * x_t) + \text{Norm}(W_{hz} * h_{t-1}))$$

$$h_{\text{new}} = \text{NonLinear}(\text{Norm}(W_{xh} * x_t) + \text{Norm}(W_{hh} * (h_{t-1} \odot g_r)))$$

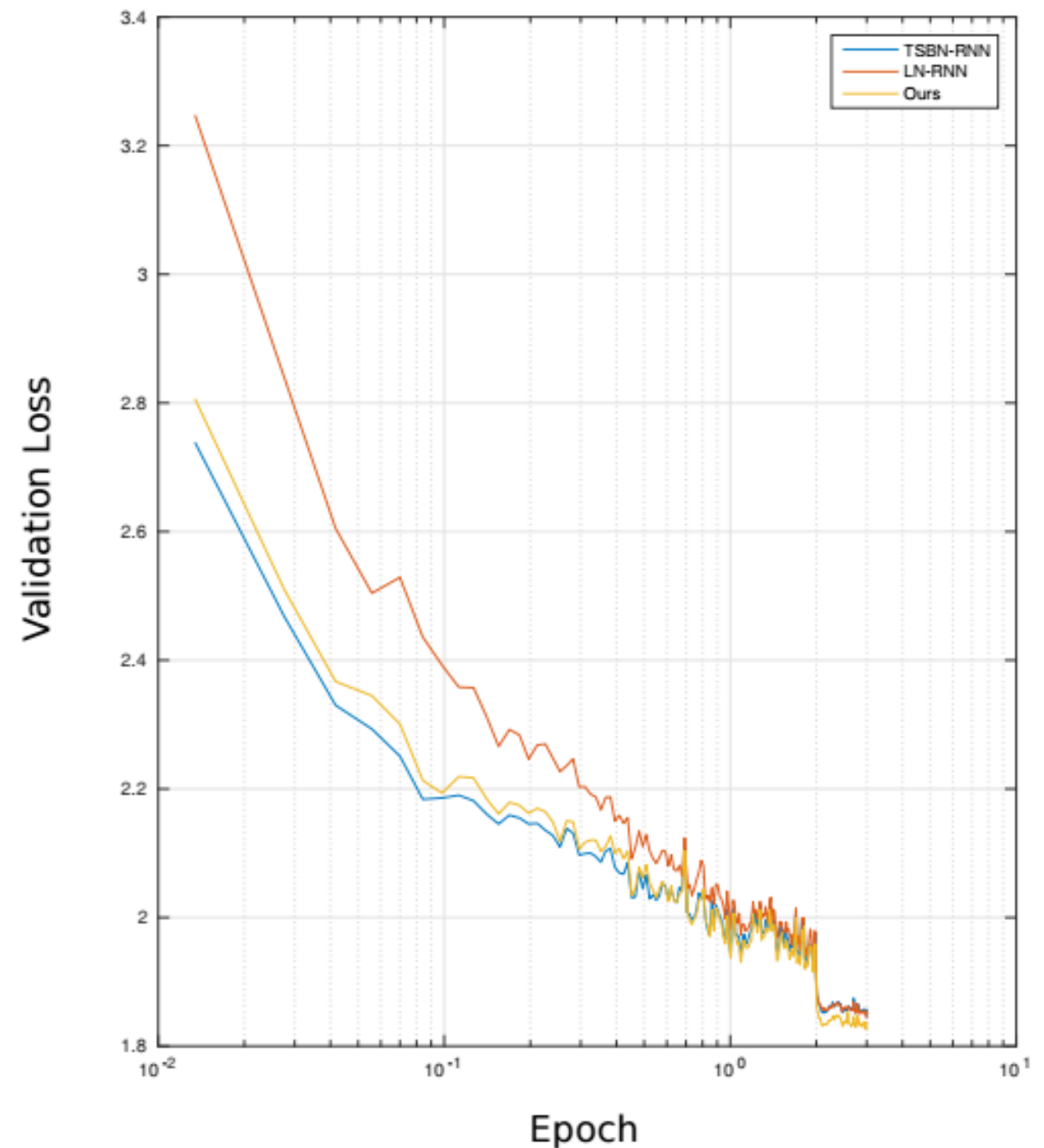
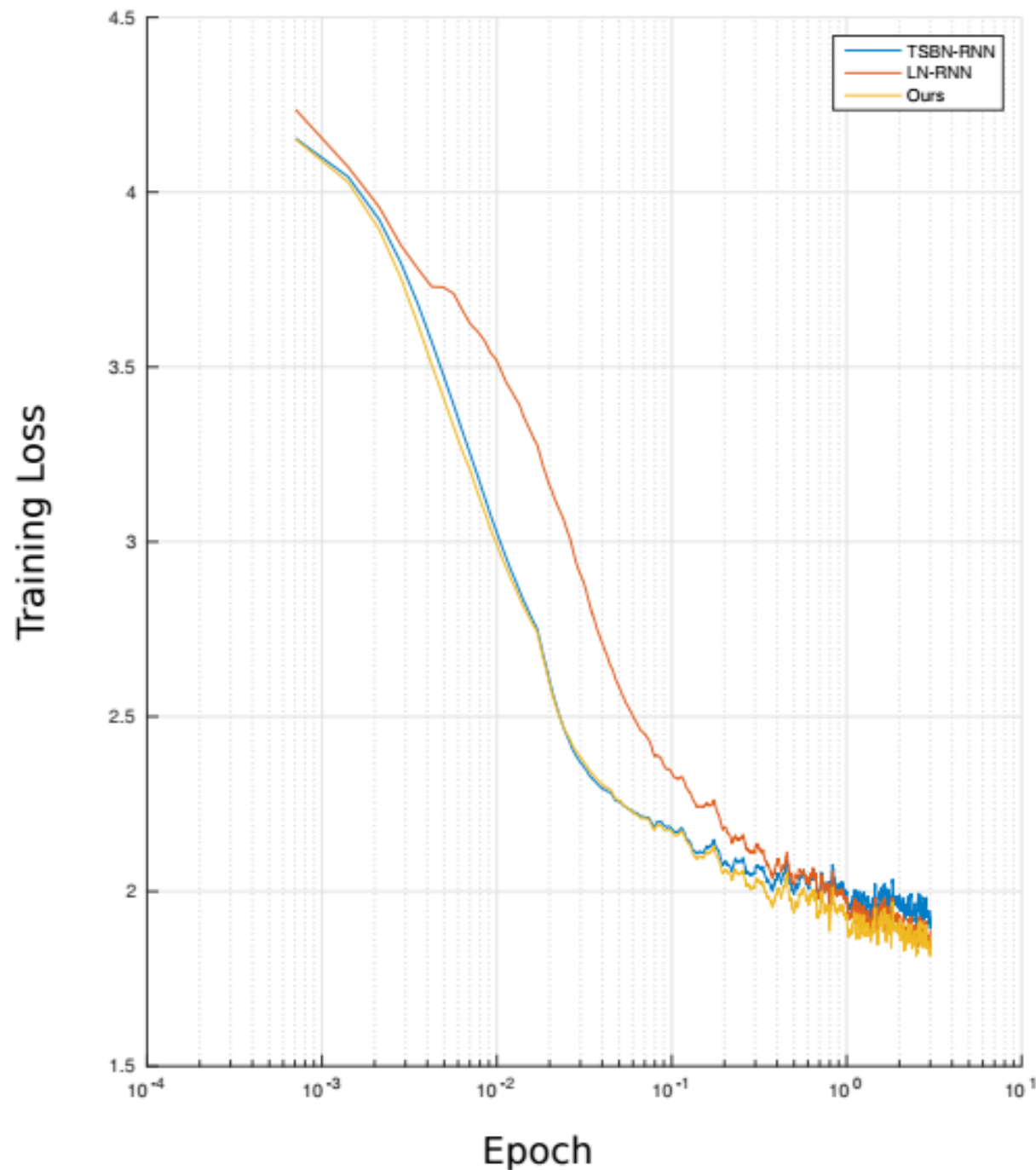
$$h_t = g_z \odot h_{\text{new}} + (1 - g_z) \odot h_{t-1}$$

Streamlining Normalization Results



Streamlining Normalization

RNN Results



Normalization Techniques

Approach	FF & FC	FF & Conv	Rec & FC	Rec & Conv	Online Learning	Small Batch	All Combined
Original Batch Normalization(BN)	✓	✓	✗	✗	✗	Suboptimal	✗
Time-specific BN	✓	✓	Limited	Limited	✗	Suboptimal	✗
Layer Normalization	✓	✗*	✓	✗*	✓	✓	✗*
Streaming Normalization	✓	✓	✓	✓	✓	✓	✓

Table 1: An overview of normalization techniques for different tasks. ✓: works well. ✗: does not work well. FF: Feedforward. Rec: Recurrent. FC: Fully-connected. Conv: convolutional. Limited: time-specific BN requires recording normalization statistics for each timestep and thus may not generalize to novel sequence length. *Layer normalization does not fail on these tasks but perform significantly worse than the best approaches.

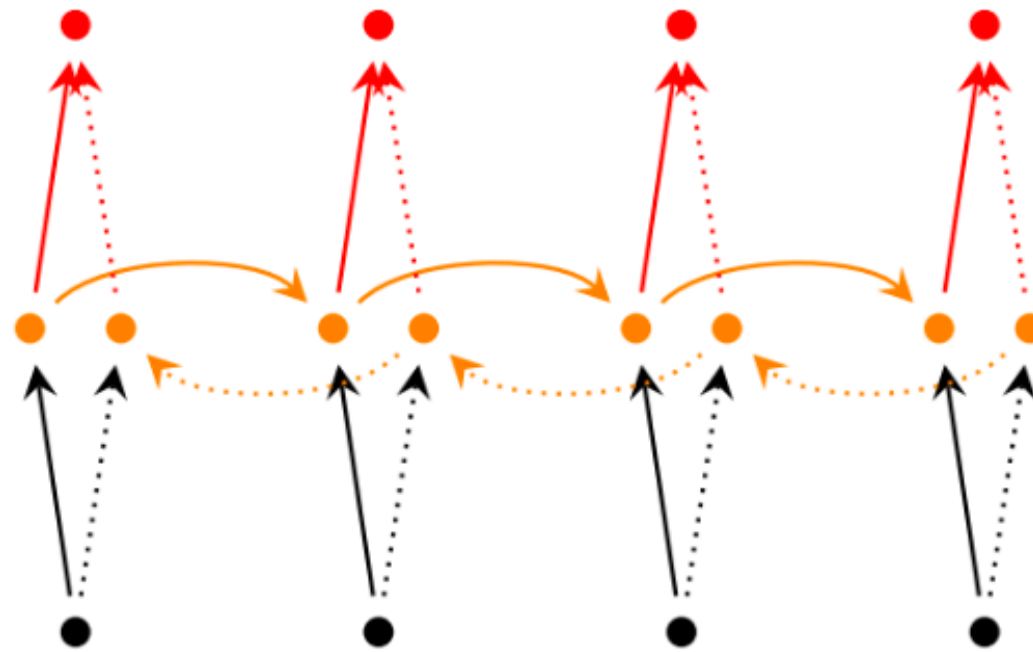
Applications

RNN Applications

- Speech Recognition
- Natural Language Processing
- Action Recognition
- Machine Translation
- Many more to come

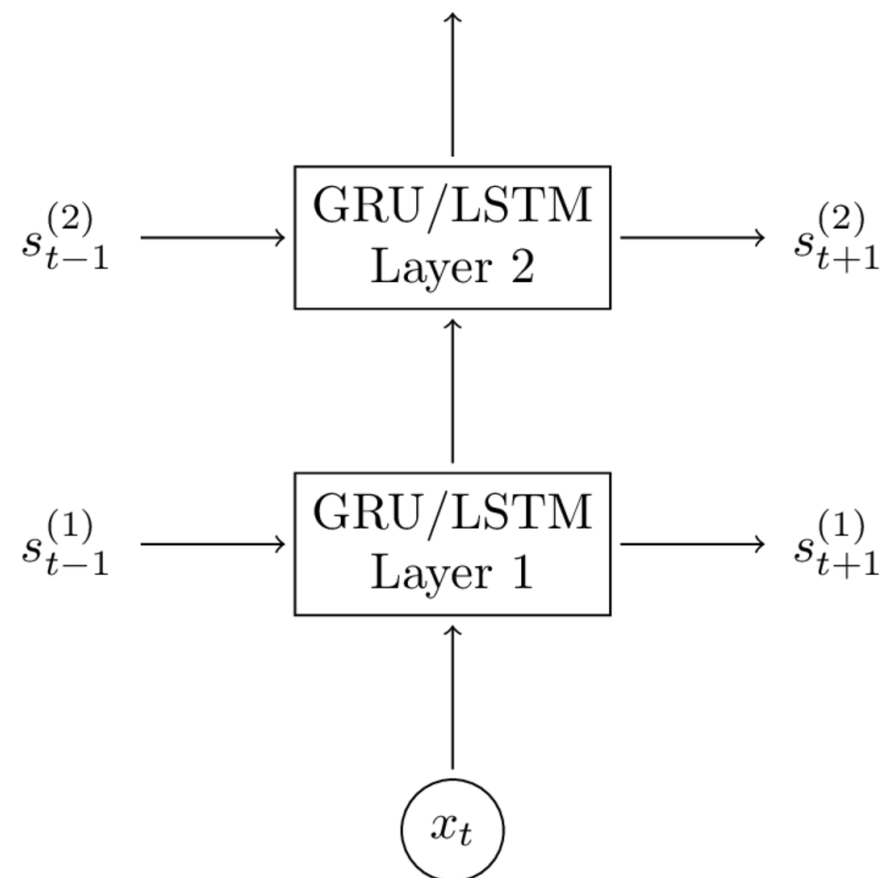
Bidirectional RNNs

- The output at time t does not depend on previous time steps but also the future
- Two RNNs stacked on top of each other



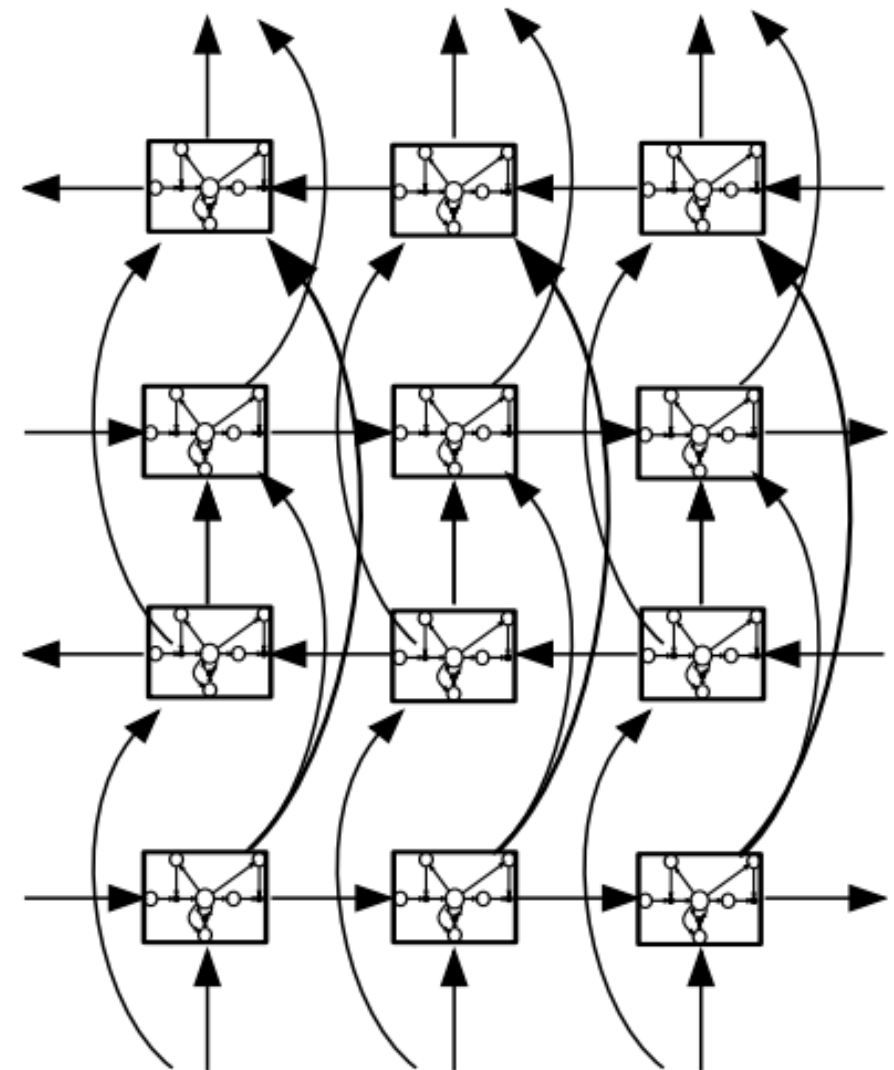
Deep RNNs

- Stack them on top of each other
- The output of the previous RNN is the input to the next one



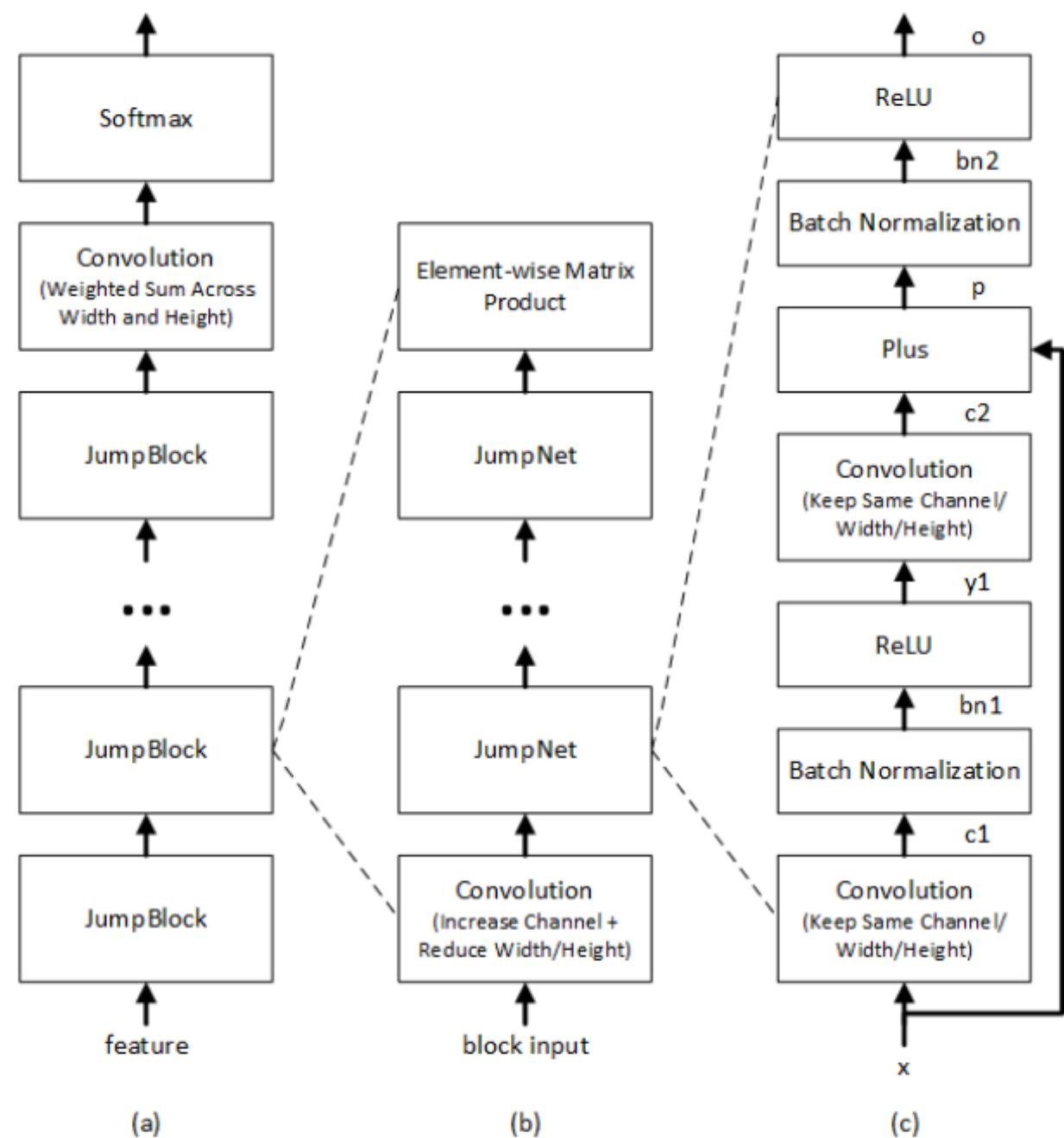
Speech Recognition

- Infer phonemes from the past and the future
- Bidirectional Stacked LSTM



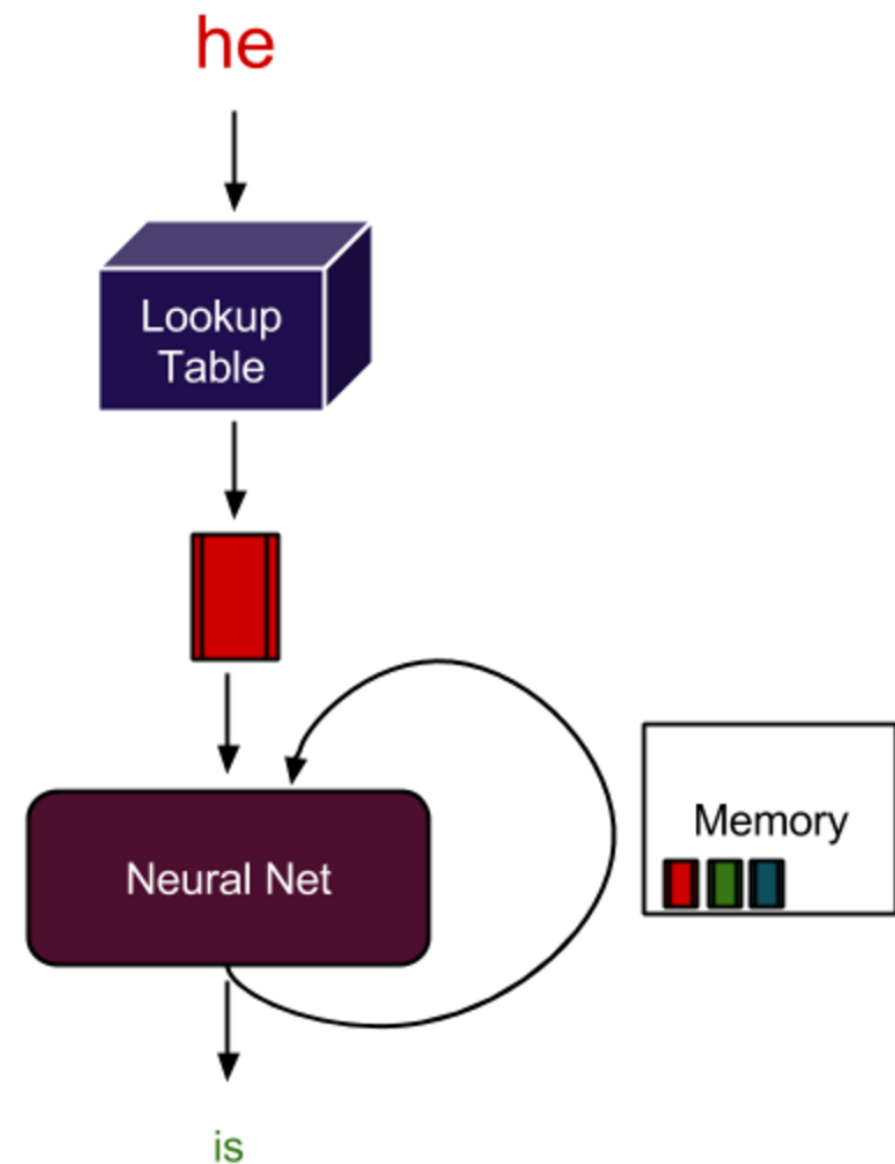
Conversational Speech Recognition

- Achieving human parity
- Combining RNN and CNN



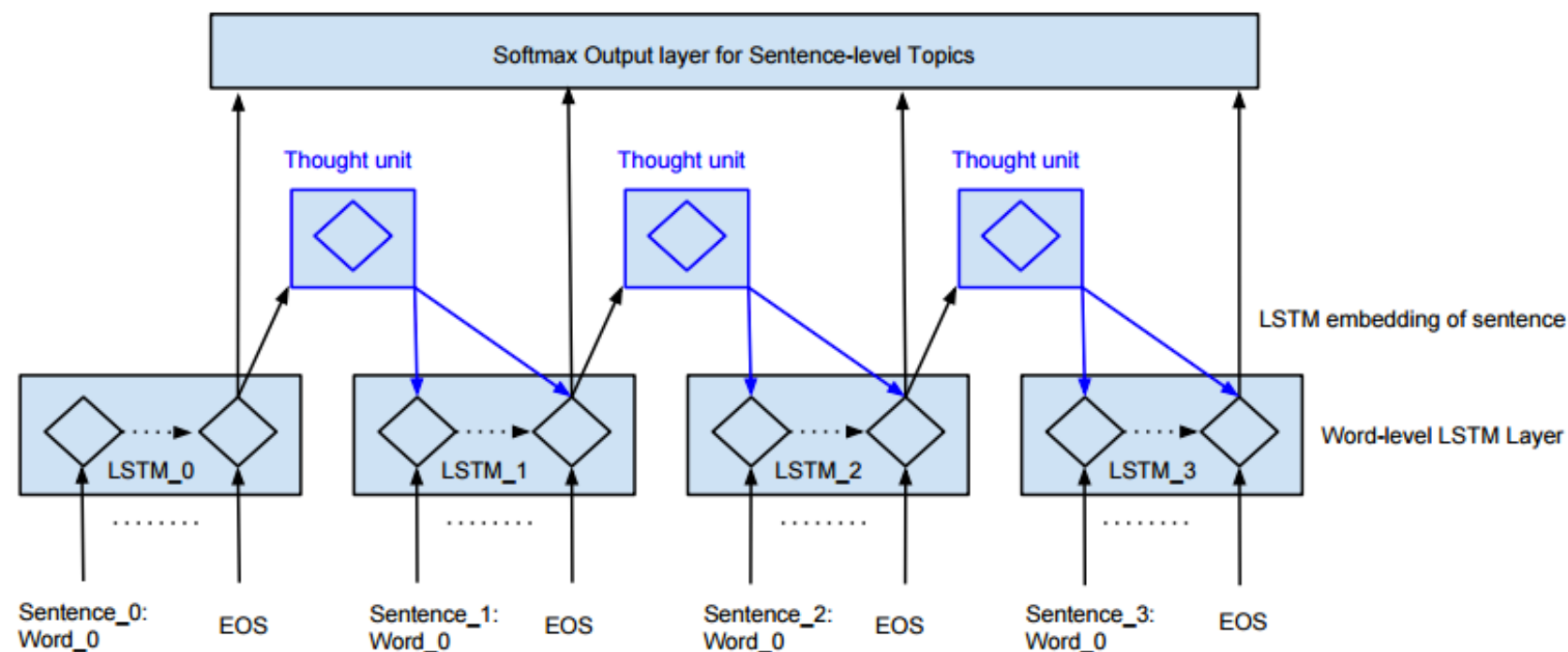
Natural Language Processing

- Infer character distribution from past characters in the sequence
- Remember for longer durations



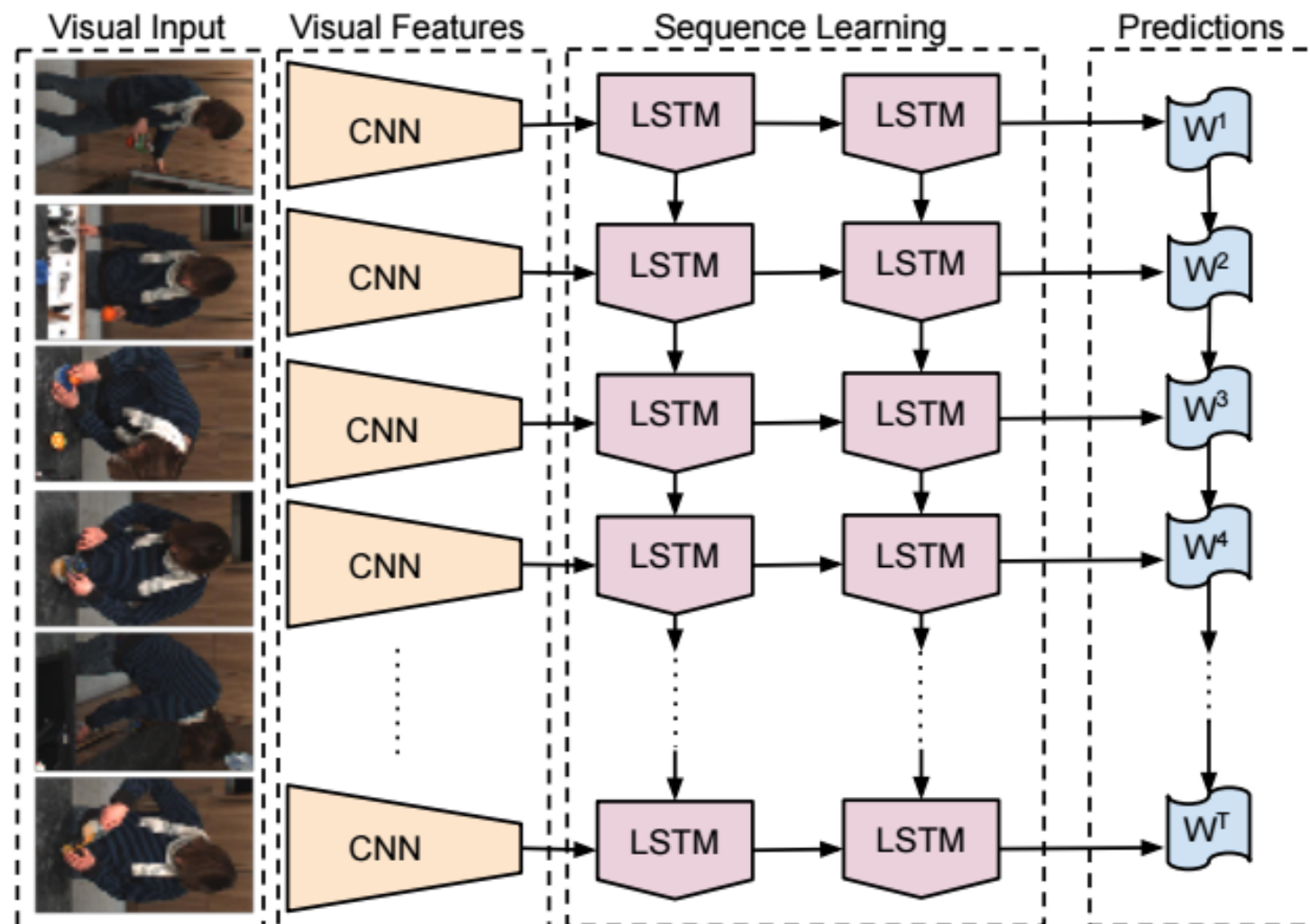
Contextual LSTM for NLP Tasks

- Using a word embedding with LSTM layers to learn sentiments
- Incorporating a thought unit



Action Recognition

- Using LSTMs and CNN for videos
- CNN creates a feature vector that is feed into LSTM



Google's Neural Machine Translation System

- Encoder and Decoder LSTMs
- Attention model

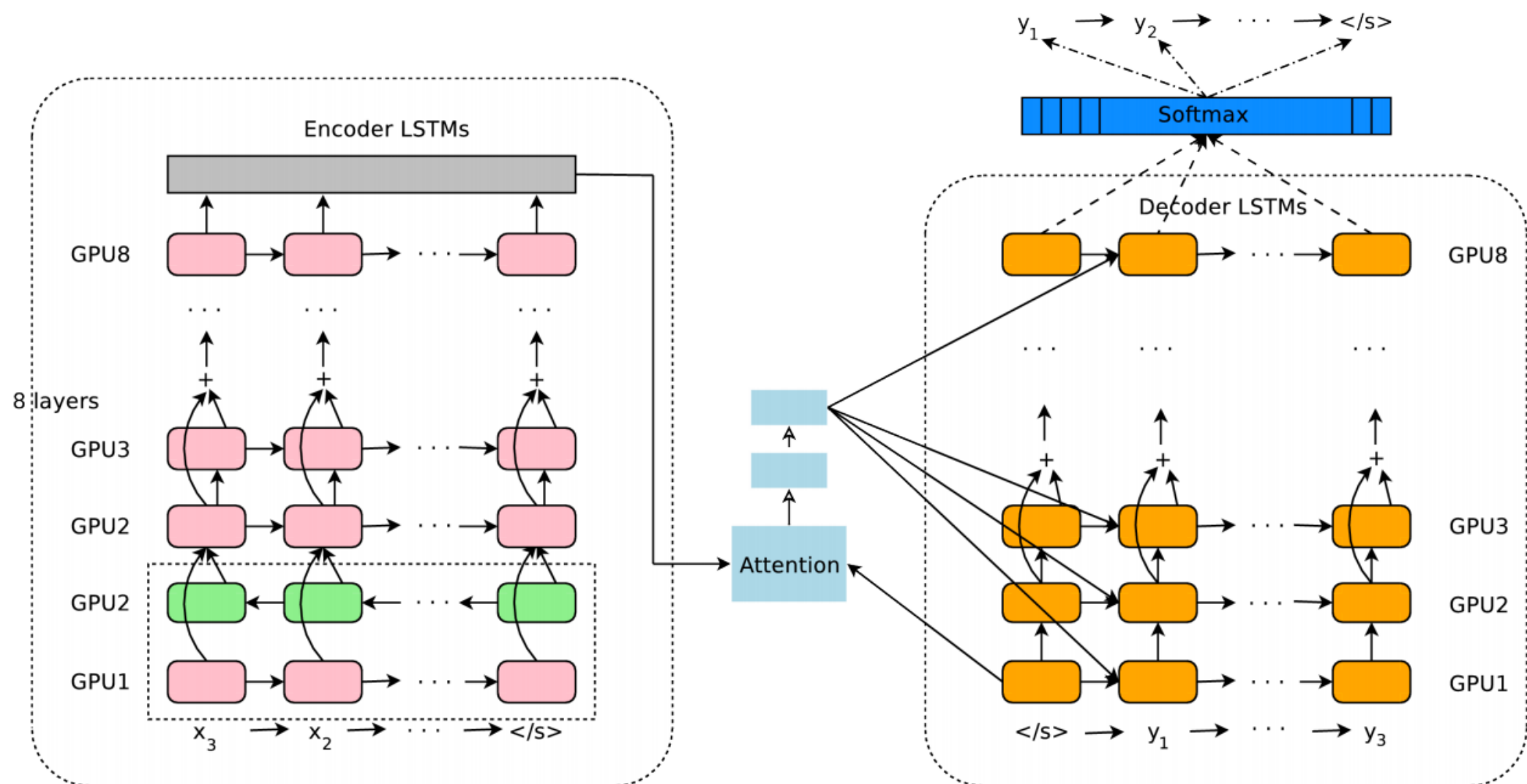


Image Captioning Pt 1

- Combination of CNN and LSTM to caption images
- Using a pretrained CNN for visual features

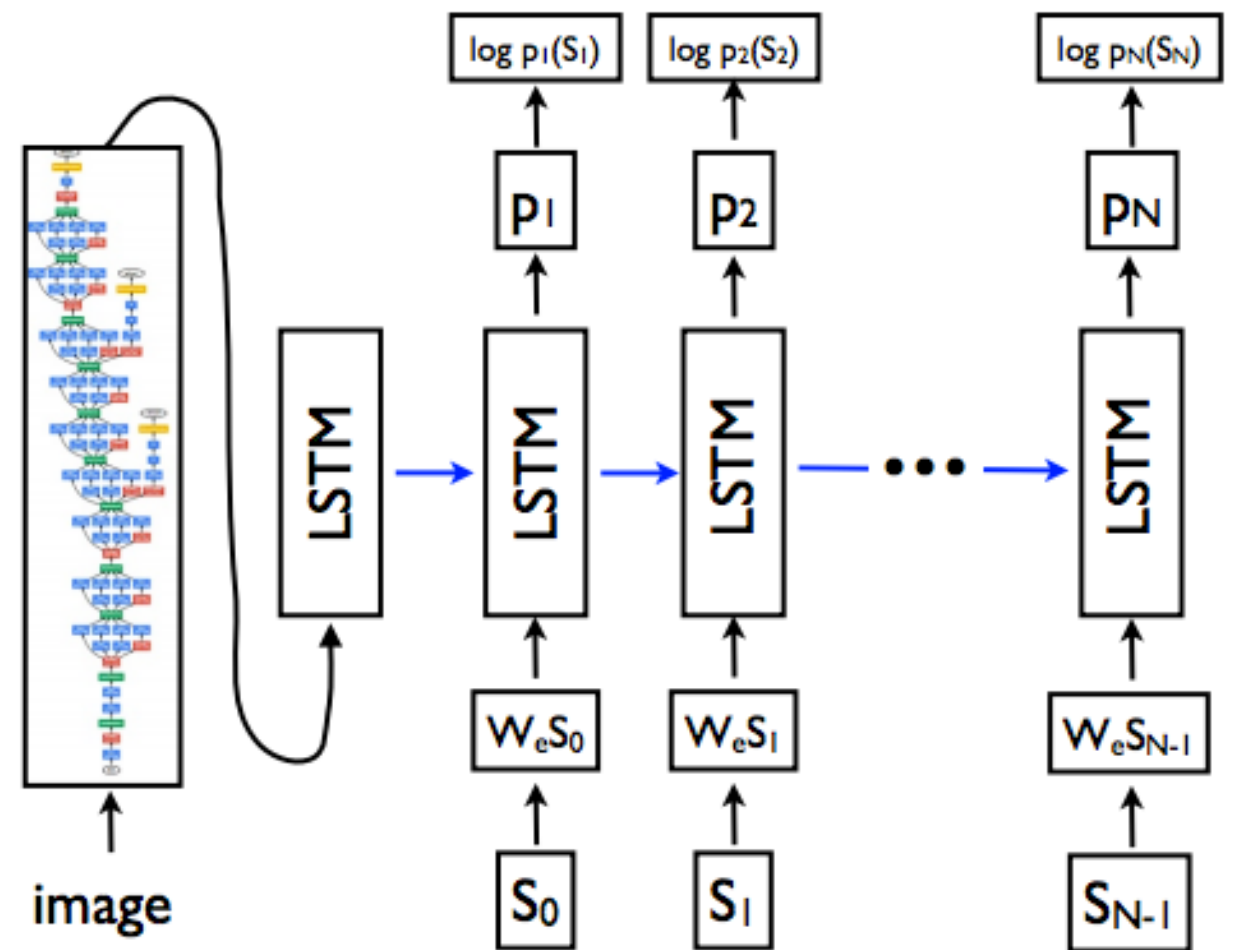


Image Captioning Pt 2

A person riding a motorcycle on a dirt road.



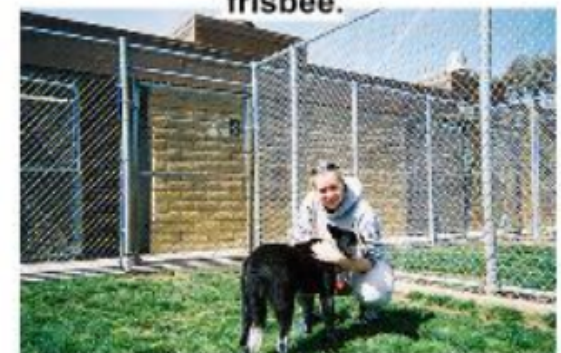
Two dogs play in the grass.



A skateboarder does a trick on a ramp.



A dog is jumping to catch a frisbee.



A group of young people playing a game of frisbee.



Two hockey players are fighting over the puck.



A little girl in a pink hat is blowing bubbles.



A refrigerator filled with lots of food and drinks.



A herd of elephants walking across a dry grass field.



A close up of a cat laying on a couch.



A red motorcycle parked on the side of the road.



A yellow school bus parked in a parking lot.



Describes without errors

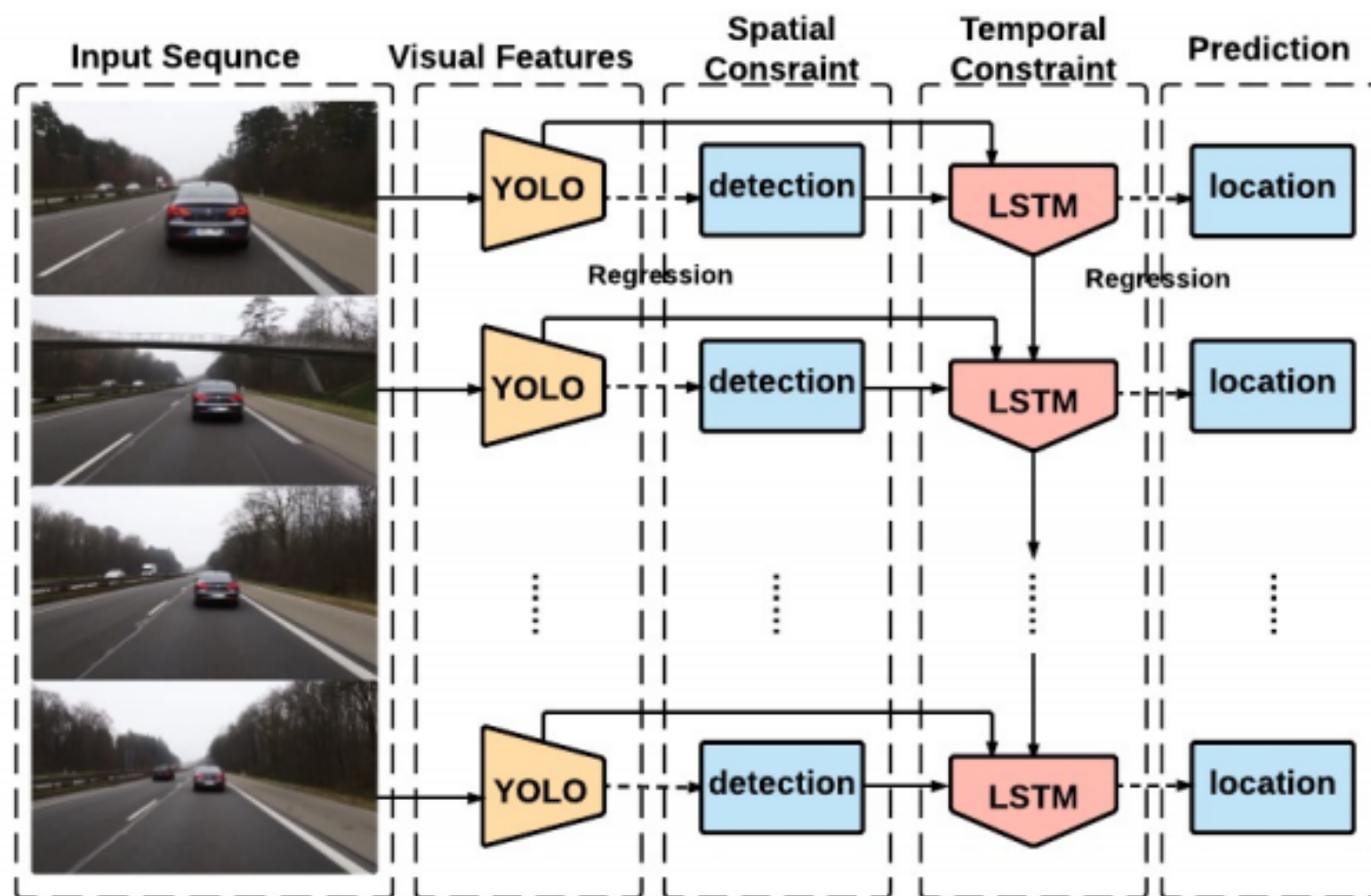
Describes with minor errors

Somewhat related to the image

Unrelated to the image

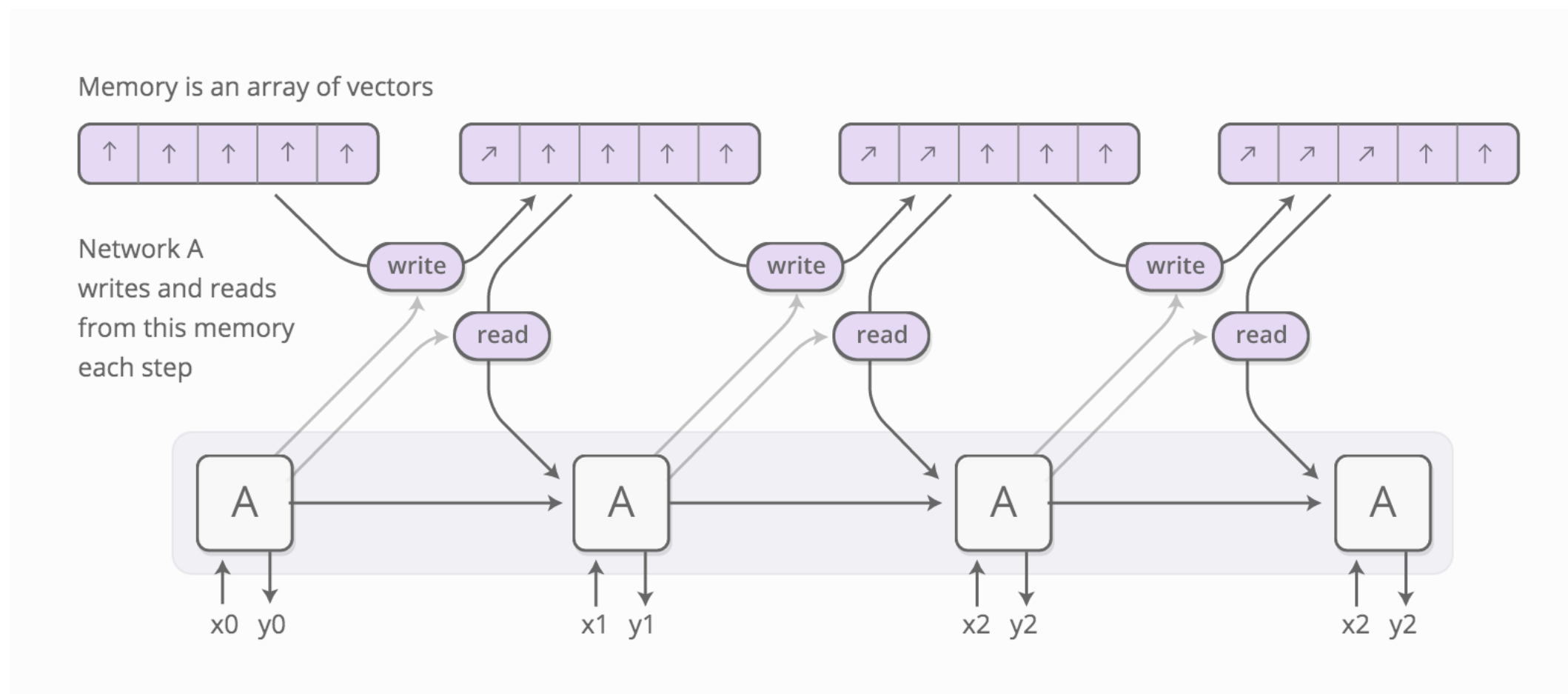
Object Tracking

- Using an object detector with an LSTM to track objects
- Model the dynamics of video



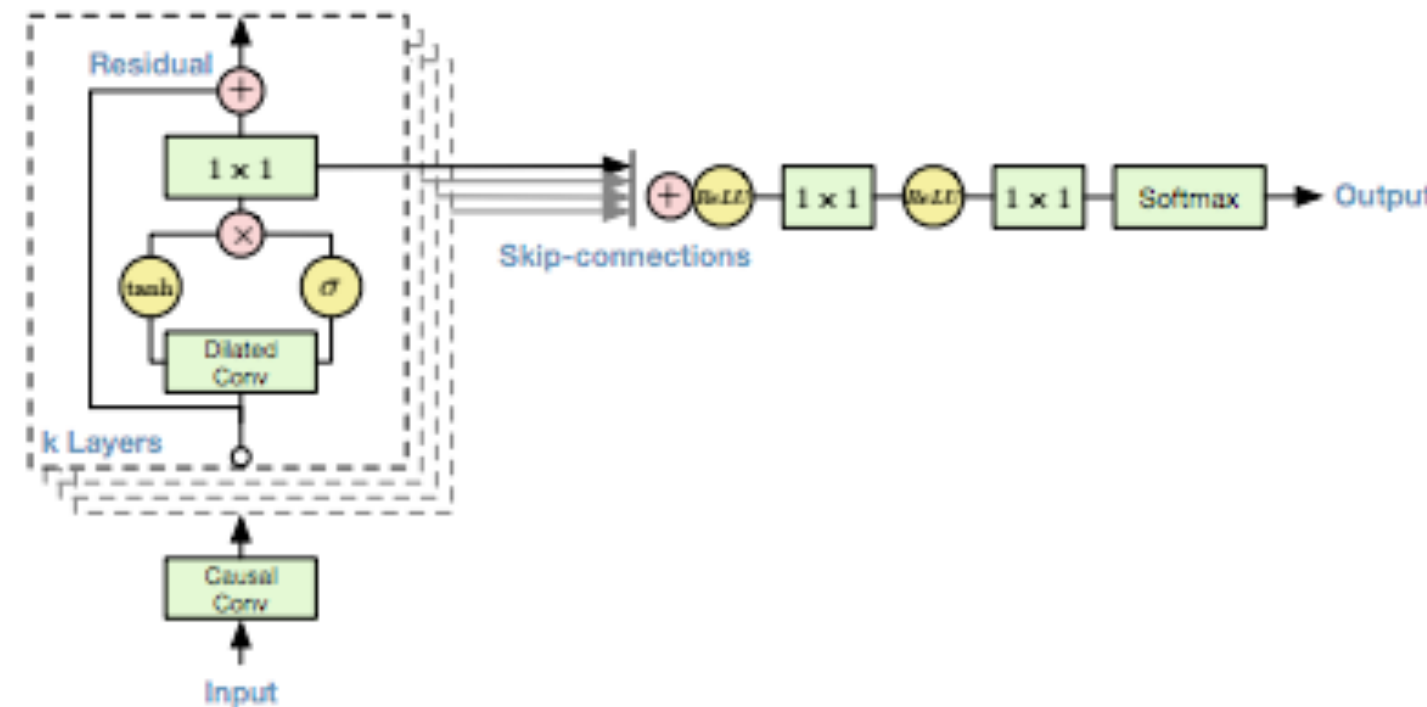
Neural Turing Machines

- LSTM with external memory
- Analogous to a Turing Machine



WaveNet

- Using stack of diluted layers
- To generate next sample, it models conditional probability given previous samples



DoomBot

- Doom Competition
- Facebook won 1st place (F1)
- <https://www.youtube.com/watch?v=94EPSjQH38Y>