

Lattice Recurrent Unit

Chaitanya Ahuja

Language Technologies Institute, CMU

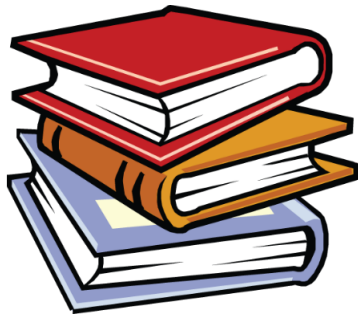
May 30, 2017

A plethora of data is **sequential** in nature

A plethora of data is **sequential** in nature
⇒ The **order** is important

Motivation

A plethora of data is **sequential** in nature
⇒ The **order** is important



A plethora of data is **sequential** in nature
⇒ The **order** is important

and because the permutations increase **factorially** ($n!$)
with the number of data points (n)

A plethora of data is **sequential** in nature
⇒ The **order** is important

and because the permutations increase **factorially** ($n!$)
with the number of data points (n)

We might lose information if the order is not preserved

Language Model

Language Model

CH A R

Language Model

CH A R

$$P(R | CHA)$$

Language Model

C H A R

$$P\left(R \mid C H A\right)$$

$$P\left(x_t \mid x_{1:t-1}\right)$$

Rise of **Neural Networks**

- **Function Approximator:** A model that can approximate such distributions.

Rise of **Neural Networks**

- **Function Approximator:** A model that can approximate such distributions.

⇒ **Neural Networks**

Rise of **Neural Networks**

- **Function Approximator:** A model that can approximate such distributions.

⇒ **Neural Networks**

- **Adaptable:** Generalizable over large datasets by increasing number of parameters.

Rise of **Neural Networks**

- **Function Approximator:** A model that can approximate such distributions.

⇒ **Neural Networks**

- **Adaptable:** Generalizable over large datasets by increasing number of parameters.

⇒ Increase **depth** of Neural Networks.

Rise of **Neural Networks**

- **Function Approximator:** A model that can approximate such distributions.

⇒ **Neural Networks**

- **Adaptable:** Generalizable over large datasets by increasing number of parameters.

⇒ Increase **depth** of Neural Networks.

- **Order Preserving:** Use **order-information** as well to predict.

Rise of **Neural Networks**

- **Function Approximator:** A model that can approximate such distributions.

⇒ **Neural Networks**

- **Adaptable:** Generalizable over large datasets by increasing number of parameters.

⇒ Increase **depth** of Neural Networks.

- **Order Preserving:** Use **order-information** as well to predict.

⇒ Deep Recurrent Neural Networks (**Deep RNNs**)

Rise of **Neural Networks**

- **Function Approximator:** A model that can approximate such distributions.

⇒ **Neural Networks**

- **Adaptable:** Generalizable over large datasets by increasing number of parameters.

⇒ Increase **depth** of Neural Networks.

- **Order Preserving:** Use **order-information** as well to predict.

⇒ Deep Recurrent Neural Networks (**Deep RNNs**)

- **Trainable:** Estimate *accurate* functions in a **practical** time-frame.

Rise of **Neural Networks**

- **Function Approximator:** A model that can approximate such distributions.

⇒ **Neural Networks**

- **Adaptable:** Generalizable over large datasets by increasing number of parameters.

⇒ Increase **depth** of Neural Networks.

- **Order Preserving:** Use **order-information** as well to predict.

⇒ Deep Recurrent Neural Networks (**Deep RNNs**)

- **Trainable:** Estimate *accurate* functions in a **practical** time-frame.
Use **gradient** based approaches.

Rise of **Neural Networks**

- **Function Approximator:** A model that can approximate such distributions.

⇒ **Neural Networks**

- **Adaptable:** Generalizable over large datasets by increasing number of parameters.

⇒ Increase **depth** of Neural Networks.

- **Order Preserving:** Use **order-information** as well to predict.

⇒ Deep Recurrent Neural Networks (**Deep RNNs**)

- **Trainable:** Estimate *accurate* functions in a **practical** time-frame.

Use **gradient** based approaches.

Are they **easy to train**?

Rise of **Recurrent** Neural Networks

$$x_t = \text{RNN}(x_{t-1}, c_{t-1})$$

Rise of **Recurrent** Neural Networks

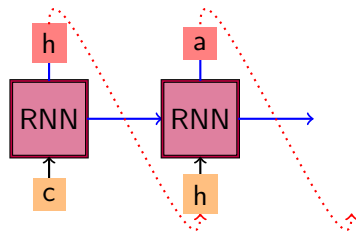
$$\begin{aligned}x_t &= \text{RNN}(x_{t-1}, c_{t-1}) \\&= \text{RNN}(x_{t-1}, x_{t-2} \dots, x_1) \\&= \text{RNN}(x_{1:t-1})\end{aligned}$$

Rise of **Recurrent** Neural Networks

$$\begin{aligned}x_t &= \text{RNN}(x_{t-1}, c_{t-1}) \\&= \text{RNN}(x_{t-1}, x_{t-2} \dots, x_1) \\&= \text{RNN}(x_{1:t-1}) \\ \Rightarrow \text{RNN} &\sim \text{P}\left(x_t \mid x_{1:t-1}\right)\end{aligned}$$

Rise of **Recurrent** Neural Networks

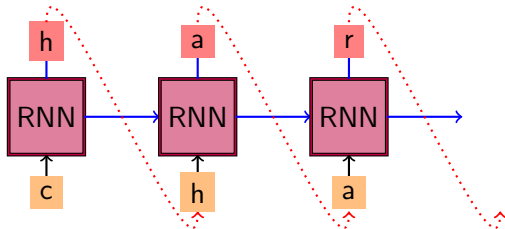
Language Modeling



- Note: Blue Arrows correspond to **hidden states** (c_{t-1}).

Rise of **Recurrent** Neural Networks

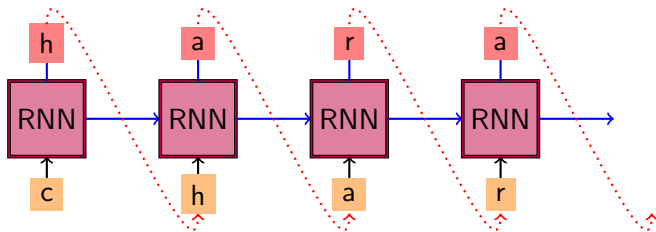
Language Modeling



- Note: Blue Arrows correspond to **hidden states** (c_{t-1}).

Rise of **Recurrent** Neural Networks

Language Modeling



- Note: Blue Arrows correspond to **hidden states** (c_{t-1}).

Rise of **Recurrent** Neural Networks

- **Trainability:** Vanilla RNNs suffer from the **Vanishing Gradient** problem,

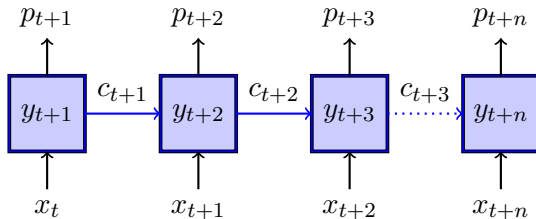


Figure 1: Forward Pass

Rise of **Recurrent** Neural Networks

- **Trainability:** Vanilla RNNs suffer from the **Vanishing Gradient** problem,

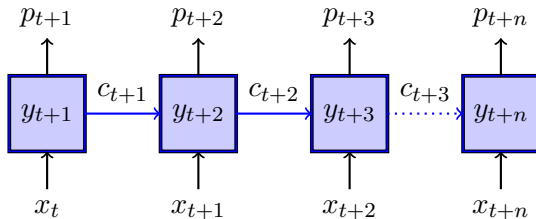


Figure 1: Forward Pass

$$y_{t+1} = \mathbf{X}x_t + \mathbf{C}c_t$$

Rise of **Recurrent** Neural Networks

- **Trainability:** Vanilla RNNs suffer from the **Vanishing Gradient** problem,

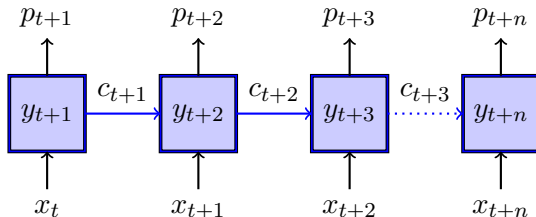


Figure 1: Forward Pass

$$\begin{aligned}y_{t+1} &= \mathbf{X}x_t + \mathbf{C}c_t \\c_{t+1} &= \sigma(y_{t+1})\end{aligned}$$

Rise of **Recurrent** Neural Networks

- **Trainability:** Vanilla RNNs suffer from the **Vanishing Gradient** problem,

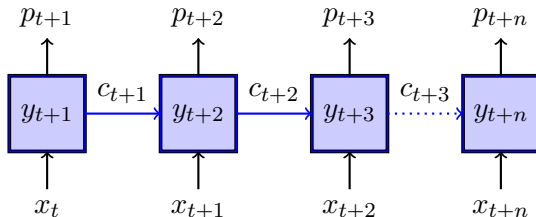


Figure 1: Forward Pass

$$\begin{aligned}y_{t+1} &= \mathbf{X}x_t + \mathbf{C}c_t \\c_{t+1} &= \sigma(y_{t+1}) \\p_{t+1} &= \text{softmax}(\mathbf{W}c_{t+1})\end{aligned}$$

Rise of **Recurrent** Neural Networks

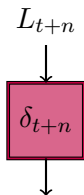


Figure 2: Gradients in Vanilla RNN

Rise of **Recurrent** Neural Networks

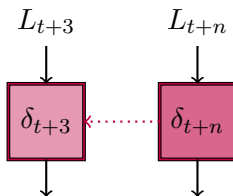


Figure 2: Gradients in Vanilla RNN

Rise of **Recurrent** Neural Networks

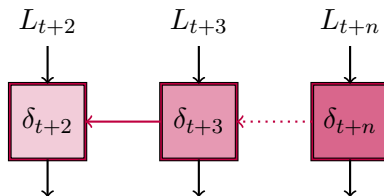


Figure 2: Gradients in Vanilla RNN

Rise of **Recurrent** Neural Networks

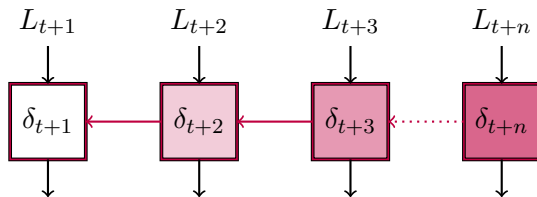


Figure 2: Gradients in Vanilla RNN

Rise of **Recurrent** Neural Networks

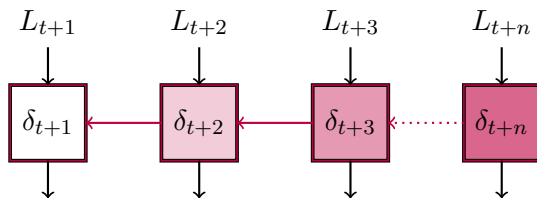


Figure 2: Gradients in Vanilla RNN

$$L_j = (p_j^{gt} - p_j)^2$$

Rise of **Recurrent** Neural Networks

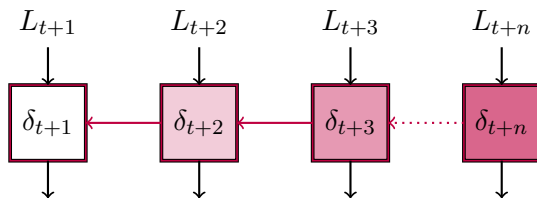


Figure 2: Gradients in Vanilla RNN

$$L_j = (p_j^{gt} - p_j)^2$$
$$\frac{\partial L_{t+n}}{\partial y_{t+1}} = \frac{\partial L_{t+n}}{\partial y_{t+n}} \cdot \frac{\partial y_{t+n}}{\partial y_{t+n-1}} \cdot \dots \cdot \frac{\partial y_{t+2}}{\partial y_{t+1}}$$

Rise of **Recurrent** Neural Networks

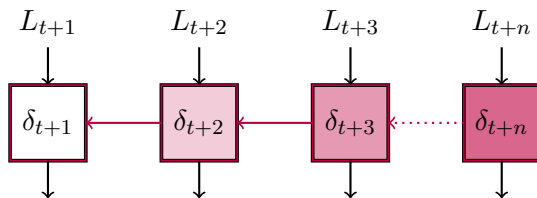


Figure 2: Gradients in Vanilla RNN

$$\begin{aligned} L_j &= (p_j^{gt} - p_j)^2 \\ \frac{\partial L_{t+n}}{\partial y_{t+1}} &= \frac{\partial L_{t+n}}{\partial y_{t+n}} \cdot \frac{\partial y_{t+n}}{\partial y_{t+n-1}} \cdots \frac{\partial y_{t+2}}{\partial y_{t+1}} \\ &= \frac{\partial L_{t+n}}{\partial y_{t+n}} \cdot \prod_{\tau=t+2}^{t+n} \frac{\partial y_{\tau}}{\partial y_{\tau-1}} \end{aligned}$$

Rise of **Recurrent** Neural Networks

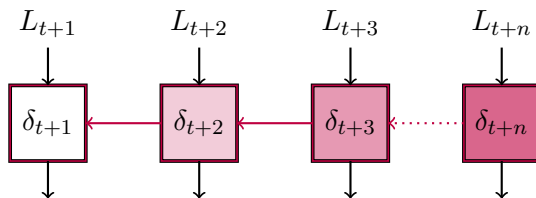


Figure 2: Gradients in Vanilla RNN

$$\prod_{\tau=t+2}^{t+n} \frac{\partial y_{\tau}}{\partial y_{\tau-1}}$$

Rise of **Recurrent** Neural Networks

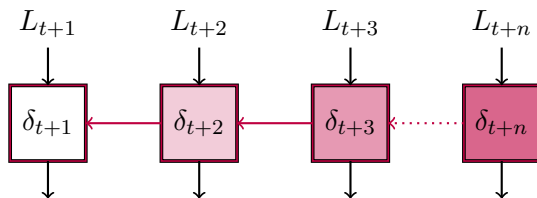


Figure 2: Gradients in Vanilla RNN

$$\prod_{\tau=t+2}^{t+n} \frac{\partial y_{\tau}}{\partial y_{\tau-1}}$$

$$\frac{\partial y_{\tau}}{\partial y_{\tau-1}} = \mathbf{C}^T \sigma'(y_{\tau-1})$$

Rise of **Recurrent** Neural Networks

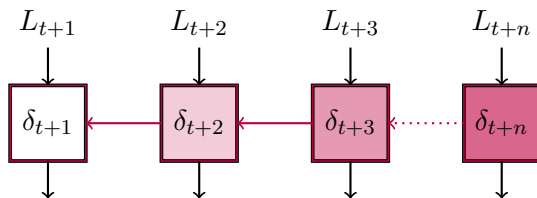


Figure 2: Gradients in Vanilla RNN

$$\prod_{\tau=t+2}^{t+n} \frac{\partial y_{\tau}}{\partial y_{\tau-1}}$$

$$\begin{aligned} \frac{\partial y_{\tau}}{\partial y_{\tau-1}} &= \mathbf{C}^T \sigma'(y_{\tau-1}) \\ \left\| \frac{\partial y_{\tau}}{\partial y_{\tau-1}} \right\| &\leq \|\mathbf{C}\| \frac{1}{4} \end{aligned}$$

Rise of **Recurrent** Neural Networks

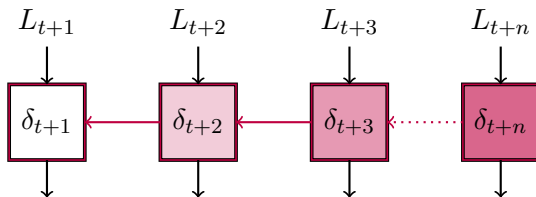


Figure 2: Gradients in Vanilla RNN

$$\prod_{\tau=t+2}^{t+n} \frac{\partial y_{\tau}}{\partial y_{\tau-1}}$$

$$\frac{\partial y_{\tau}}{\partial y_{\tau-1}} = \mathbf{C}^T \sigma'(y_{\tau-1})$$

$$\left\| \frac{\partial y_{\tau}}{\partial y_{\tau-1}} \right\| \leq \|\mathbf{C}\| \frac{1}{4}$$

$$\prod_{\tau=t+2}^{t+n} \left\| \frac{\partial y_{\tau}}{\partial y_{\tau-1}} \right\| \leq \left(\|\mathbf{C}\| \frac{1}{4} \right)^{n-1}$$

Rise of **Recurrent** Neural Networks

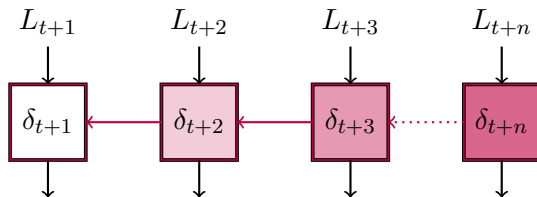


Figure 2: Gradients in Vanilla RNN

$$\prod_{\tau=t+2}^{t+n} \frac{\partial y_{\tau}}{\partial y_{\tau-1}}$$

Exponentially Decaying
Gradients

$$\begin{aligned} \frac{\partial y_{\tau}}{\partial y_{\tau-1}} &= \mathbf{C}^T \sigma'(y_{\tau-1}) \\ \left\| \frac{\partial y_{\tau}}{\partial y_{\tau-1}} \right\| &\leq \|\mathbf{C}\| \frac{1}{4} \\ \prod_{\tau=t+2}^{t+n} \left\| \frac{\partial y_{\tau}}{\partial y_{\tau-1}} \right\| &\leq \left(\|\mathbf{C}\| \frac{1}{4} \right)^{n-1} \end{aligned}$$

Rise of **Recurrent** Neural Networks

- GRUs/LSTMs alleviate this problem by using gates on inputs, outputs and hidden states [Hochreiter and Schmidhuber, 1997].

$$y_{\tau} = \mathbf{X}x_{\tau} + \mathbf{C}\sigma(y_{\tau-1}) \quad (1)$$

Let $\mathbf{C} = \mathbf{I}$ and $\sigma(x) = x$. Then,

Rise of **Recurrent** Neural Networks

- GRUs/LSTMs alleviate this problem by using gates on inputs, outputs and hidden states [Hochreiter and Schmidhuber, 1997].

$$y_{\tau} = \mathbf{X}x_{\tau} + \mathbf{C}\sigma(y_{\tau-1}) \quad (1)$$

Let $\mathbf{C} = \mathbf{I}$ and $\sigma(x) = x$. Then,

$$\frac{\partial y_{\tau}}{\partial y_{\tau-1}} = \mathbf{C}^T \sigma'(y_{\tau-1})$$

Rise of **Recurrent** Neural Networks

- GRUs/LSTMs alleviate this problem by using gates on inputs, outputs and hidden states [Hochreiter and Schmidhuber, 1997].

$$y_\tau = \mathbf{X}x_\tau + \mathbf{C}\sigma(y_{\tau-1}) \quad (1)$$

Let $\mathbf{C} = \mathbf{I}$ and $\sigma(x) = x$. Then,

$$\begin{aligned} \frac{\partial y_\tau}{\partial y_{\tau-1}} &= \mathbf{C}^T \sigma'(y_{\tau-1}) \\ &\approx \mathbf{I} \end{aligned}$$

To retain modeling power, introduce a non-linear term

Rise of **Recurrent** Neural Networks

- GRUs/LSTMs alleviate this problem by using gates on inputs, outputs and hidden states [Hochreiter and Schmidhuber, 1997].

$$y_\tau = \mathbf{X}x_\tau + \mathbf{C}\sigma(y_{\tau-1}) \quad (1)$$

Let $\mathbf{C} = \mathbf{I}$ and $\sigma(x) = x$. Then,

$$\begin{aligned} \frac{\partial y_\tau}{\partial y_{\tau-1}} &= \mathbf{C}^T \sigma'(y_{\tau-1}) \\ &\approx \mathbf{I} \quad \text{(Constant Error Carrousel)} \end{aligned}$$

To retain modeling power, introduce a non-linear term

Rise of Recurrent Neural Networks

- GRUs/LSTMs alleviate this problem by using gates on inputs, outputs and hidden states [Hochreiter and Schmidhuber, 1997].

$$y_{\tau} = \mathbf{X}x_{\tau} + \mathbf{C}\sigma(y_{\tau-1}) \quad (1)$$

Let $\mathbf{C} = \mathbf{I}$ and $\sigma(x) = x$. Then,

$$\begin{aligned} \frac{\partial y_{\tau}}{\partial y_{\tau-1}} &= \mathbf{C}^T \sigma'(y_{\tau-1}) \\ &\approx \mathbf{I} \quad \text{(Constant Error Carrousel)} \end{aligned}$$

To retain modeling power, introduce a non-linear term

$$y_{\tau} = \mathbf{X}x_{\tau} + y_{\tau-1} + \mathbf{z} \cdot \phi(y_{\tau-1})$$

Rise of Recurrent Neural Networks

- GRUs/LSTMs alleviate this problem by using gates on inputs, outputs and hidden states [Hochreiter and Schmidhuber, 1997].

$$y_\tau = \mathbf{X}x_\tau + \mathbf{C}\sigma(y_{\tau-1}) \quad (1)$$

Let $\mathbf{C} = \mathbf{I}$ and $\sigma(x) = x$. Then,

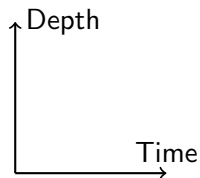
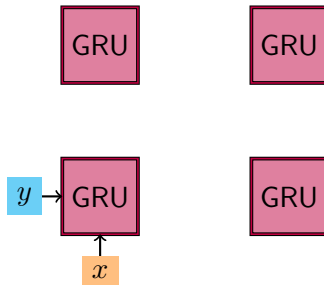
$$\begin{aligned} \frac{\partial y_\tau}{\partial y_{\tau-1}} &= \mathbf{C}^T \sigma'(y_{\tau-1}) \\ &\approx \mathbf{I} \quad \text{(Constant Error Carrousel)} \end{aligned}$$

To retain modeling power, introduce a non-linear term

$$\begin{aligned} y_\tau &= \mathbf{X}x_\tau + y_{\tau-1} + \mathbf{z} \cdot \phi(y_{\tau-1}) \\ \frac{\partial y_\tau}{\partial y_{\tau-1}} &\approx \mathbf{I} + \mathbf{z} \cdot \phi'(y_{\tau-1}) \end{aligned}$$

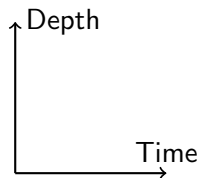
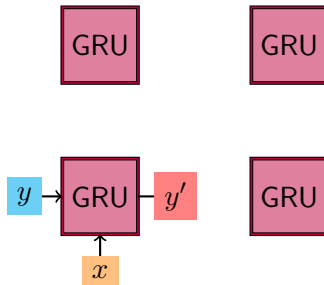
Rise of **Recurrent** Neural Networks

GRU Equations



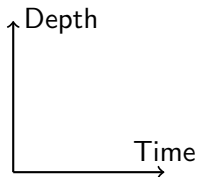
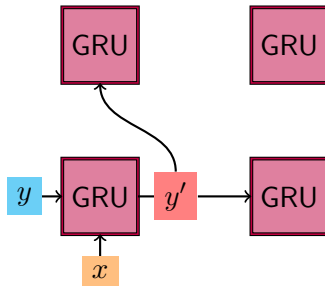
Rise of **Recurrent** Neural Networks

GRU Equations



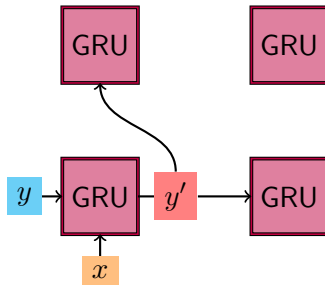
Rise of **Recurrent** Neural Networks

GRU Equations

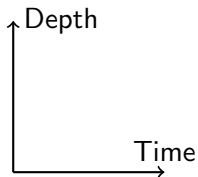


Rise of **Recurrent** Neural Networks

GRU Equations

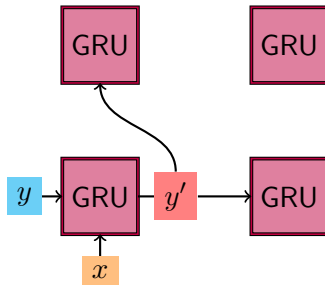


$$z = \sigma(\mathbf{X}_z x + \mathbf{Y}_z y)$$



Rise of **Recurrent** Neural Networks

GRU Equations



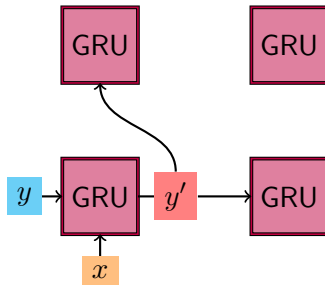
$$z = \sigma(\mathbf{X}_z x + \mathbf{Y}_z y)$$

$$r = \sigma(\mathbf{X}_r x + \mathbf{Y}_r y)$$

Depth
Time

Rise of **Recurrent** Neural Networks

GRU Equations



$$z = \sigma(\mathbf{X}_z \text{orange } x + \mathbf{Y}_z \text{blue } y)$$

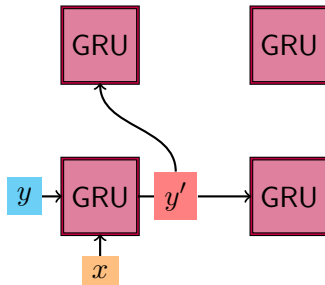
$$r = \sigma(\mathbf{X}_r \text{orange } x + \mathbf{Y}_r \text{blue } y)$$

$$\tilde{y} = \tanh(\mathbf{X} \text{orange } x + \mathbf{Y}(r \cdot \text{blue } y))$$

Depth
Time

Rise of **Recurrent** Neural Networks

GRU Equations



Depth
↑
Time
→

$$z = \sigma(\mathbf{X}_z \text{orange} + \mathbf{Y}_z \text{blue})$$

$$r = \sigma(\mathbf{X}_r \text{orange} + \mathbf{Y}_r \text{blue})$$

$$\tilde{y} = \tanh(\mathbf{X} \text{orange} + \mathbf{Y}(r \cdot \text{blue}))$$

$$y_{t+1} = \text{red } y' = \text{green } (1 - z) \cdot y + z \cdot \tilde{y}$$

Rise of **Deep** Recurrent Neural Networks

- Even though the gating mechanisms (i.e. CEC) resolve vanishing gradient issues along time,

Rise of **Deep** Recurrent Neural Networks

- Even though the gating mechanisms (i.e. CEC) resolve vanishing gradient issues along time,
- Gradients have been shown to vanish along depth.

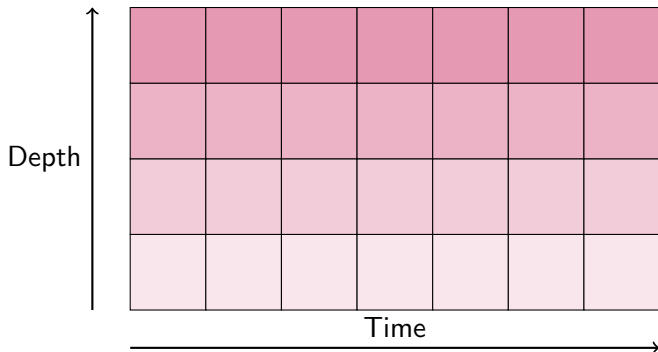


Figure 3: Gradients along Depth in a GRU

Rise of **Deep** Recurrent Neural Networks

- We observe the gradients while training a Deep GRU network on Character-Level Language Model.

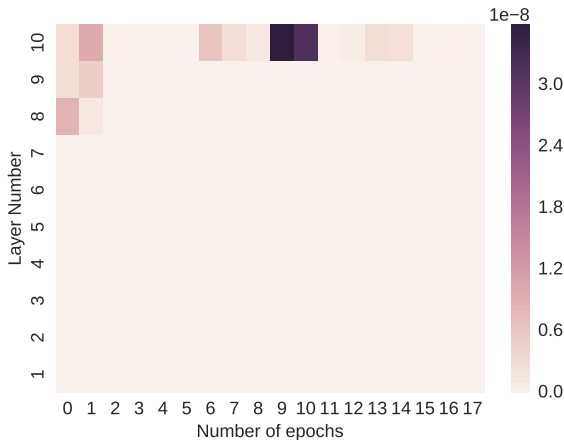
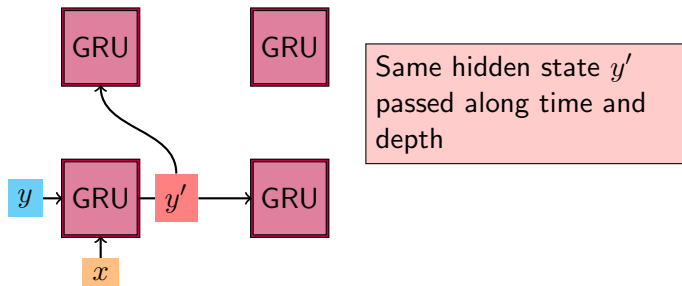


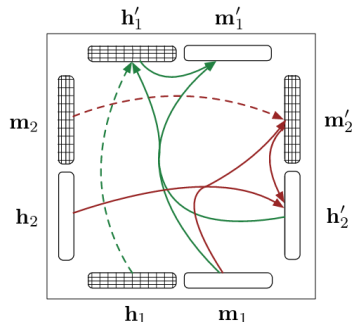
Figure 4: Gradient-Norms across depth in a 10-layered GRU

Rise of **Deep** Recurrent Neural Networks

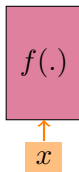


Related Work

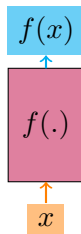
- Grid-LSTMs [Kalchbrenner et al., 2015] is to have **different hidden states along time and depth**.



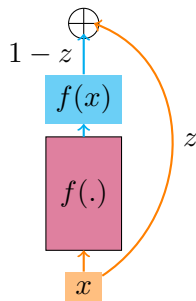
- Recurrent Highway Networks [Zilly et al., 2016] use gating in the depth to enforce **Constant Error Carousel**.



- Recurrent Highway Networks [Zilly et al., 2016] use gating in the depth to enforce **Constant Error Carrousel**.

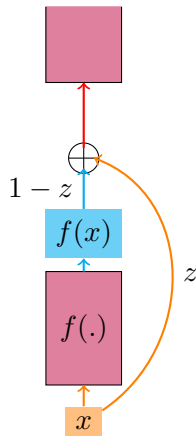


- Recurrent Highway Networks [Zilly et al., 2016] use gating in the depth to enforce **Constant Error Carrousel**.



Related Work

- Recurrent Highway Networks [Zilly et al., 2016] use gating in the depth to enforce **Constant Error Carrousel**.



Is it possible to design a **Recurrent Unit**

Is it possible to design a **Recurrent Unit**
that passes
different hidden states along depth and time

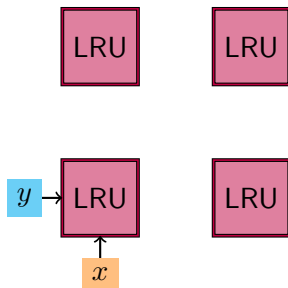
Is it possible to design a **Recurrent Unit**
that passes
different hidden states along depth and time
while
enforcing **Constant Error Carrousel** ?

Is it possible to design a **Recurrent Unit**
that passes
different hidden states along depth and time
while
enforcing **Constant Error Carrousel** ?

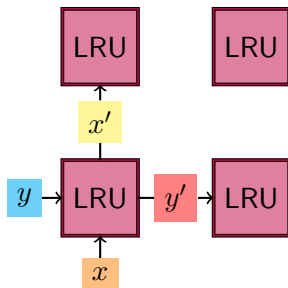
YES!!!

Formulation of Lattice Recurrent Unit

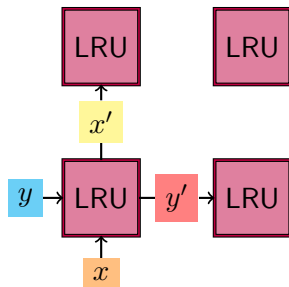
Formulation of Lattice Recurrent Unit



Formulation of Lattice Recurrent Unit

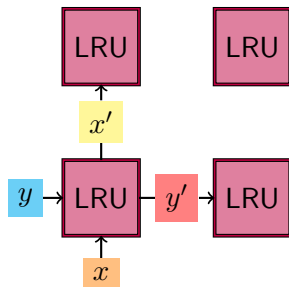


Formulation of Lattice Recurrent Unit



$$z = \sigma(\mathbf{X}_z \textcolor{brown}{x} + \mathbf{Y}_z \textcolor{blue}{y})$$

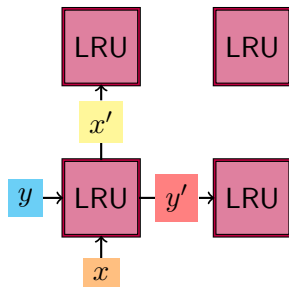
Formulation of Lattice Recurrent Unit



$$z = \sigma(\mathbf{X}_z \text{orange } x + \mathbf{Y}_z \text{blue } y)$$

$$r = \sigma(\mathbf{X}_r \text{orange } x + \mathbf{Y}_r \text{blue } y)$$

Formulation of Lattice Recurrent Unit

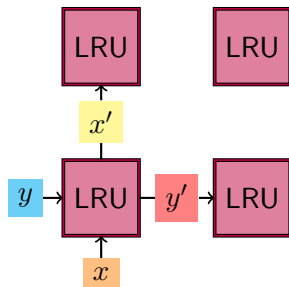


$$z = \sigma(\mathbf{X}_z \text{orange } x + \mathbf{Y}_z \text{blue } y)$$

$$r = \sigma(\mathbf{X}_r \text{orange } x + \mathbf{Y}_r \text{blue } y)$$

$$q = \sigma(\mathbf{X}_q \text{orange } x + \mathbf{Y}_q \text{blue } y)$$

Formulation of Lattice Recurrent Unit



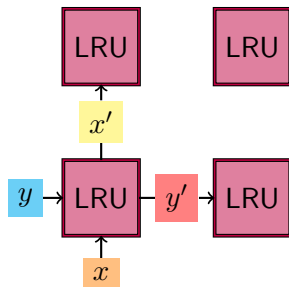
$$z = \sigma(\mathbf{X}_z \text{orange } x + \mathbf{Y}_z \text{blue } y)$$

$$r = \sigma(\mathbf{X}_r \text{orange } x + \mathbf{Y}_r \text{blue } y)$$

$$q = \sigma(\mathbf{X}_q \text{orange } x + \mathbf{Y}_q \text{blue } y)$$

$$\tilde{x} = \tanh\left(\mathbf{X}_x \text{orange } x + \mathbf{Y}_x (r \cdot \text{blue } y)\right)$$

Formulation of Lattice Recurrent Unit



$$z = \sigma(\mathbf{X}_z \text{orange } x + \mathbf{Y}_z \text{blue } y)$$

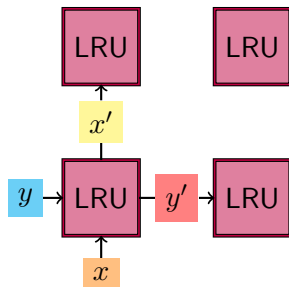
$$r = \sigma(\mathbf{X}_r \text{orange } x + \mathbf{Y}_r \text{blue } y)$$

$$q = \sigma(\mathbf{X}_q \text{orange } x + \mathbf{Y}_q \text{blue } y)$$

$$\tilde{x} = \tanh\left(\mathbf{X}_x \text{orange } x + \mathbf{Y}_x (r \cdot \text{blue } y)\right)$$

$$\tilde{y} = \tanh\left(\mathbf{Y}_y \text{blue } y + \mathbf{X}_y (q \cdot \text{orange } x)\right)$$

Formulation of Lattice Recurrent Unit



$$z = \sigma(\mathbf{X}_z \text{orange } x + \mathbf{Y}_z \text{blue } y)$$

$$r = \sigma(\mathbf{X}_r \text{orange } x + \mathbf{Y}_r \text{blue } y)$$

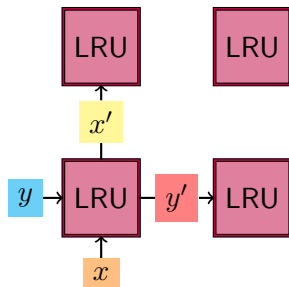
$$q = \sigma(\mathbf{X}_q \text{orange } x + \mathbf{Y}_q \text{blue } y)$$

$$\tilde{x} = \tanh\left(\mathbf{X}_x \text{orange } x + \mathbf{Y}_x (r \cdot \text{blue } y)\right)$$

$$\tilde{y} = \tanh\left(\mathbf{Y}_y \text{blue } y + \mathbf{X}_y (q \cdot \text{orange } x)\right)$$

$$x_{t+1} = \text{yellow } x' = z \cdot \tilde{y} + \text{green } (1 - z) \cdot x$$

Formulation of Lattice Recurrent Unit



$$z = \sigma(\mathbf{X}_z \text{orange } x + \mathbf{Y}_z \text{blue } y)$$

$$r = \sigma(\mathbf{X}_r \text{orange } x + \mathbf{Y}_r \text{blue } y)$$

$$q = \sigma(\mathbf{X}_q \text{orange } x + \mathbf{Y}_q \text{blue } y)$$

$$\tilde{x} = \tanh\left(\mathbf{X}_x \text{orange } x + \mathbf{Y}_x (r \cdot \text{blue } y)\right)$$

$$\tilde{y} = \tanh\left(\mathbf{Y}_y \text{blue } y + \mathbf{X}_y (q \cdot \text{orange } x)\right)$$

$$x_{t+1} = \text{yellow } x' = z \cdot \tilde{y} + \text{green } (1 - z) \cdot x$$

$$y_{t+1} = \text{red } y' = z \cdot \tilde{x} + \text{green } (1 - z) \cdot y$$

Task: **Character Level** Language Modeling

- Out of Vocabulary Words can potentially be modeled.

Task: **Character Level** Language Modeling

- Out of Vocabulary Words can potentially be modeled.
- Character Aware Neural Language Models [Kim et al., 2016]

Task: **Character Level** Language Modeling

- Out of Vocabulary Words can potentially be modeled.
- Character Aware Neural Language Models [Kim et al., 2016]
- Speech Synthesis [Wang et al., 2017]

- We take 2 datasets **Penn Tree Bank** and **War and Peace**.

Datasets

- We take 2 datasets **Penn Tree Bank** and **War and Peace**.
- Each of them has around **5 million** characters.

War and Peace (WP)

Well, Prince, so Genoa and Lucca are now just family estates of the Bonapartes. But I warn you, if you don't tell me that this means war, if you still try to defend the infamies and horrors perpetrated by that Antichrist—I really believe he is Antichrist—I will have nothing more to do with you and you are no longer my friend, no longer my 'faithful slave,' as you call yourself!

Penn Tree Bank (PTB)

the asbestos fiber <unk> is unusually <unk> once it enters the <unk> with even brief exposures to it causing symptoms that show up decades later researchers said, <unk> inc. the unit of new york-based <unk> corp. that makes kent cigarettes stopped using <unk> in its <unk> cigarette filters in N.

- Multi-Class Classification Problem

Experimental Setup

- Multi-Class Classification Problem
- Categorical Cross-Entropy Loss
 - $\text{loss} = \sum_{\forall i \in \mathcal{C}} p_i \log(\hat{p}_i)$
where \mathcal{C} is the set of all classes.

Experimental Setup

- Multi-Class Classification Problem
- Categorical Cross-Entropy Loss
 - $\text{loss} = \sum_{\forall i \in \mathcal{C}} p_i \log(\hat{p}_i)$
where \mathcal{C} is the set of all classes.
- RNN unrolled 50 time-steps in time.

Experimental Questions

- **Accuracy:** For equal number of parameters does the **loss improve**?

Experimental Questions

- **Accuracy:** For equal number of parameters does the **loss improve**?
- **Convergence Rate:** How many epochs does it take for the model to converge?

Experimental Questions

- **Accuracy:** For equal number of parameters does the **loss improve**?
- **Convergence Rate:** How many epochs does it take for the model to converge?
- **Trainability:** How are the **gradient norms** distributed across layers?

Results - Accuracy

Table 1: Penn Treebank Dataset and losses are in bits per character (BPC). Lower is better. l is the number of layers.

Model	Tess Loss (BPC)				Hidden Units
	$l = 2$	$l = 4$	$l = 6$	$l = 8$	
GRU	1.107	1.092	1.125	2.99	387

Results - Accuracy

Table 1: Penn Treebank Dataset and losses are in bits per character (BPC). Lower is better. l is the number of layers.

Model	Tess Loss (BPC)				Hidden Units
	$l = 2$	$l = 4$	$l = 6$	$l = 8$	
GRU	1.107	1.092	1.125	2.99	387
LSTM	1.110	1.13	1.254	1.311	335

Results - Accuracy

Table 1: Penn Treebank Dataset and losses are in bits per character (BPC). Lower is better. l is the number of layers.

Model	Tess Loss (BPC)				Hidden
	$l = 2$	$l = 4$	$l = 6$	$l = 8$	Units
GRU	1.107	1.092	1.125	2.99	387
LSTM	1.110	1.13	1.254	1.311	335
GLSTM	1.119	1.106	1.111	1.105	237
Highway	1.110	1.102	1.107	-	-

Results - Accuracy

Table 1: Penn Treebank Dataset and losses are in bits per character (BPC). Lower is better. l is the number of layers.

Model	Tess Loss (BPC)				Hidden
	$l = 2$	$l = 4$	$l = 6$	$l = 8$	Units
GRU	1.107	1.092	1.125	2.99	387
LSTM	1.110	1.13	1.254	1.311	335
GLSTM	1.119	1.106	1.111	1.105	237
Highway	1.110	1.102	1.107	-	-
LRU	1.097	1.102	1.101	1.103	300

Results - Accuracy

Table 2: War and Peace Dataset and losses are in bits per character (BPC). Lower is better. l is the number of layers.

Model	Tess Loss (BPC)				Hidden Units
	$l = 2$	$l = 4$	$l = 6$	$l = 8$	
GRU	1.285	1.293	1.344	3.104	387

Results - Accuracy

Table 2: War and Peace Dataset and losses are in bits per character (BPC). Lower is better. l is the number of layers.

Model	Tess Loss (BPC)				Hidden Units
	$l = 2$	$l = 4$	$l = 6$	$l = 8$	
GRU	1.285	1.293	1.344	3.104	387
LSTM	1.297	1.353	1.530	3.100	335

Results - Accuracy

Table 2: War and Peace Dataset and losses are in bits per character (BPC). Lower is better. l is the number of layers.

Model	Tess Loss (BPC)				Hidden Units
	$l = 2$	$l = 4$	$l = 6$	$l = 8$	
GRU	1.285	1.293	1.344	3.104	387
LSTM	1.297	1.353	1.530	3.100	335
GLSTM	1.319	1.317	1.312	1.31	237
Highway	1.294	1.298	1.306	1.305	-

Results - Accuracy

Table 2: War and Peace Dataset and losses are in bits per character (BPC). Lower is better. l is the number of layers.

Model	Tess Loss (BPC)				Hidden
	$l = 2$	$l = 4$	$l = 6$	$l = 8$	Units
GRU	1.285	1.293	1.344	3.104	387
LSTM	1.297	1.353	1.530	3.100	335
GLSTM	1.319	1.317	1.312	1.31	237
Highway	1.294	1.298	1.306	1.305	-
LRU	1.279	1.280	1.281	1.287	300

Results - Convergence Rates

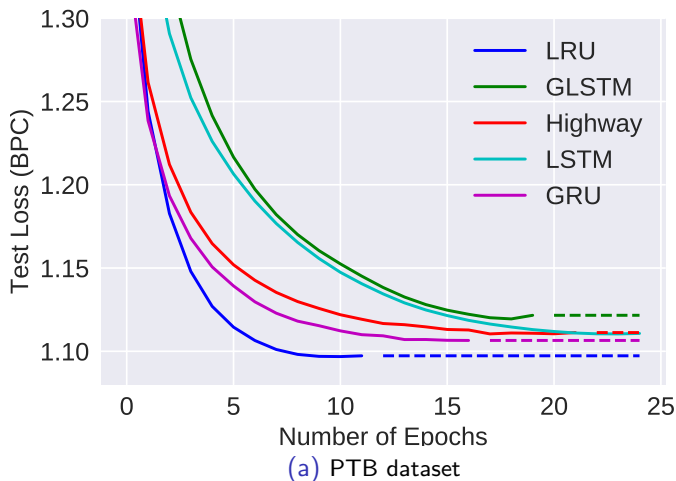


Figure 5: Convergence rates of various models on PTB and WP datasets.

Results - Convergence Rates

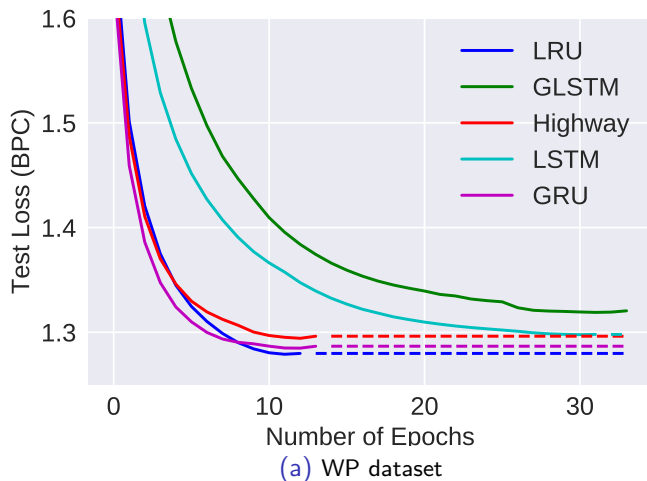
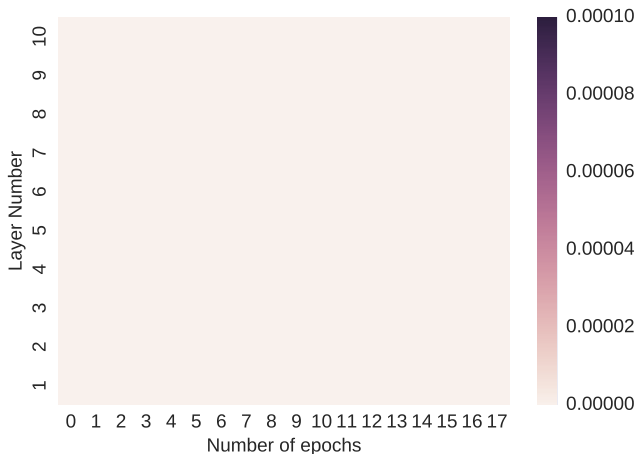


Figure 5: Convergence rates of various models on PTB and WP datasets.

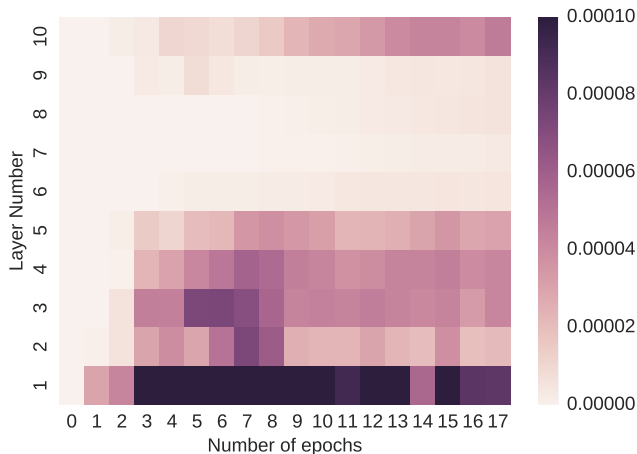
Results - Trainability



(a) GRU

Figure 6: Gradient Norms across layers and number of epochs as a heatmap. Darker values are bigger.

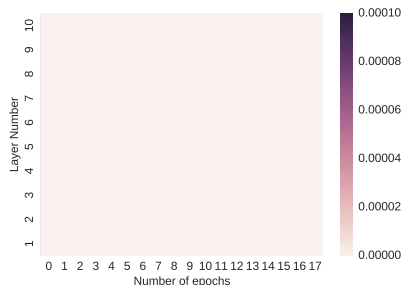
Results - Trainability



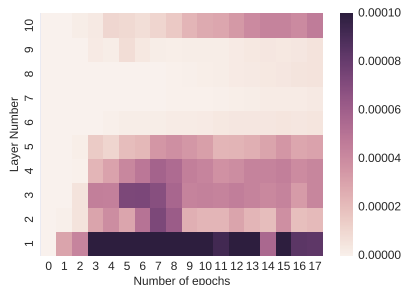
(a) LRU

Figure 6: Gradient Norms across layers and number of epochs as a heatmap. Darker values are bigger.

Results - Trainability



(a) GRU



(b) LRU

Figure 6: Gradient Norms across layers and number of epochs as a heatmap. Darker values are bigger.

- **Faster** convergence rates for **LRU**

- **Faster** convergence rates for **LRU**
- Gradients are spread out across the layers. Hence **alleviation** of vanishing gradients.

- **Faster** convergence rates for **LRU**
- Gradients are spread out across the layers. Hence **alleviation** of vanishing gradients.
- While **not aggravating** the performance.

References I



Hochreiter, S. and Schmidhuber, J. (1997).

Long short-term memory.

Neural computation, 9(8):1735–1780.



Kalchbrenner, N., Danihelka, I., and Graves, A. (2015).

Grid long short-term memory.

arXiv preprint arXiv:1507.01526.



Kim, Y., Jernite, Y., Sontag, D., and Rush, A. M. (2016).

Character-aware neural language models.

In *Thirtieth AAAI Conference on Artificial Intelligence*.



Wang, Y., Skerry-Ryan, R., Stanton, D., Wu, Y., Weiss, R. J., Jaitly, N., Yang, Z., Xiao, Y., Chen, Z., Bengio, S., et al. (2017).

Tacotron: Towards end-to-end speech syn.

arXiv preprint arXiv:1703.10135.



Zilly, J. G., Srivastava, R. K., Koutník, J., and Schmidhuber, J. (2016).

Recurrent highway networks.

arXiv preprint arXiv:1607.03474.

References II



Hochreiter, S. and Schmidhuber, J. (1997).

Long short-term memory.

Neural computation, 9(8):1735–1780.



Kalchbrenner, N., Danihelka, I., and Graves, A. (2015).

Grid long short-term memory.

arXiv preprint arXiv:1507.01526.



Kim, Y., Jernite, Y., Sontag, D., and Rush, A. M. (2016).

Character-aware neural language models.

In *Thirtieth AAAI Conference on Artificial Intelligence*.



Wang, Y., Skerry-Ryan, R., Stanton, D., Wu, Y., Weiss, R. J., Jaitly, N., Yang, Z., Xiao, Y., Chen, Z., Bengio, S., et al. (2017).

Tacotron: Towards end-to-end speech syn.

arXiv preprint arXiv:1703.10135.



Zilly, J. G., Srivastava, R. K., Koutník, J., and Schmidhuber, J. (2016).

Recurrent highway networks.

arXiv preprint arXiv:1607.03474.

Thank You

