# Recurrent Neural Networks

Jim Glass / MIT 6.806-6.864 / Spring 2021

# Today's RNN Story

- RNN language models
- Training RNNs
- Gated RNNs
  - Gated recurrent units (GRUs)
  - Long short-term memories (LSTMs)
- Other architectures
  - Deep RNNs
  - Bidirectional RNNs

# Language Modeling (Again)

- $n$-gram LMs represent history with the previous $n - 1$ words
  - Number of $n$-gram parameters increases exponentially with $n$

$$p(w_t | w_{t-(n-1)}, \ldots, w_{t-1})$$

- An alternative model represents word history with a latent variable

$$p(w_t | w_1, \ldots, w_{t-1}) \approx p(w_t | h_t)$$

- Recurrent neural networks use hidden states to capture history
  - Latent variable $h_t$ is computed based on input $x_t$ and $h_{t-1}$

$$h_t = f(x_t, h_{t-1})$$

3

# A Feed-Forward LM

- Softmax output

$$\mathbf{y} = softmax(\mathbf{U}\mathbf{h} + \mathbf{c}) \quad \mathbf{U} \in \mathbb{R}^{V \times h} \quad \mathbf{c} \in \mathbb{R}^{V}$$

- Non-linear hidden layer

$$\mathbf{h} = f(\mathbf{W}\mathbf{v} + \mathbf{b}) \quad \mathbf{W} \in \mathbb{R}^{h \times 4d} \quad \mathbf{b} \in \mathbb{R}^{h}$$
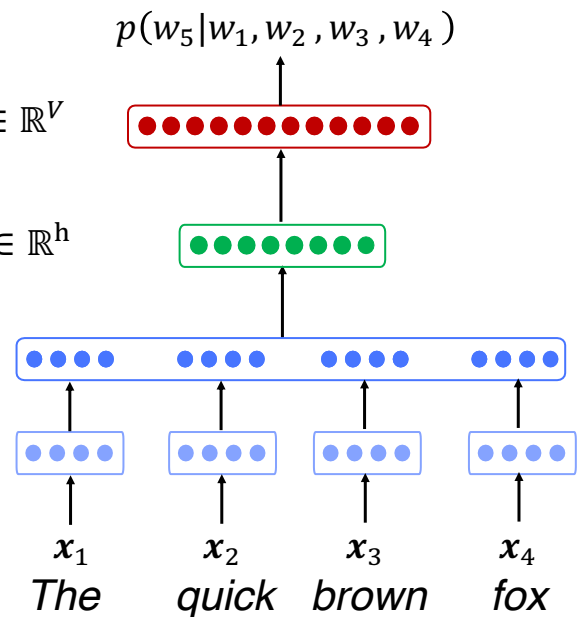
- Truncated context vector

$$\mathbf{v} = [\mathbf{v}_1; \mathbf{v}_2; \mathbf{v}_3; \mathbf{v}_4] \quad \mathbf{v} \in \mathbb{R}^{4d}$$

- Word embedding vectors

$$\mathbf{v}_i = \mathbf{E}\mathbf{x}_i \quad \mathbf{v}_i \in \mathbb{R}^{d} \quad \mathbf{E} \in \mathbb{R}^{d \times V}$$
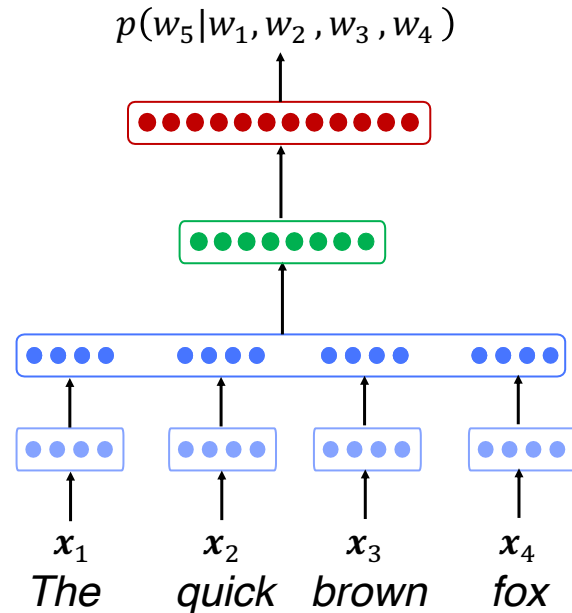
- One-hot input vectors

$$\mathbf{x}_i \in \{\mathbf{e}_i : 1 \leq i \leq V\}$$

$$p(w_5 | w_1, w_2, w_3, w_4)$$

$$\mathbf{x}_1 \quad \mathbf{x}_2 \quad \mathbf{x}_3 \quad \mathbf{x}_4$$

*The     quick     brown     fox*

4

# Feed-Forward LMs

- Advantages
  - No $n$-gram sparsity issues
  - Better memory usage
- Disadvantages
  - Fixed context window
  - Limits ability to capture history
  - No parameter sharing in $\mathbf{W}$

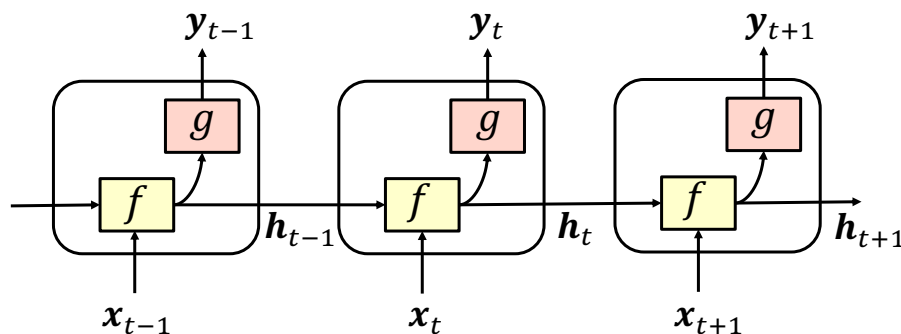- Need LM model to handle variable length sequential inputs

$$p(w_5|w_1, w_2, w_3, w_4)$$

$$\mathbf{x}_1 \quad \mathbf{x}_2 \quad \mathbf{x}_3 \quad \mathbf{x}_4$$

*The    quick   brown   fox*

# Recurrent Neural Networks (RNNs)

- RNNs accept variable length input sequences $\boldsymbol{x}_t$
- Use a hidden layer that incorporates current input and prior state

$$\boldsymbol{h}_t = f(\boldsymbol{W}_{hh}\boldsymbol{h}_{t-1} + \mathbf{W}_{xh}\boldsymbol{x}_t + \mathbf{b}_h)$$

- Optional outputs can be produced at every step

$$\boldsymbol{y}_t = g(\mathbf{W}_{hy}\boldsymbol{h}_t + \mathbf{b}_y)$$

$$\boldsymbol{y}_{t-1} \qquad \boldsymbol{y}_t \qquad \boldsymbol{y}_{t+1}$$

$$\boldsymbol{h}_{t-1} \qquad \boldsymbol{h}_t \qquad \boldsymbol{h}_{t+1}$$

$$\boldsymbol{x}_{t-1} \qquad \boldsymbol{x}_t \qquad \boldsymbol{x}_{t+1}$$

# Recurrent Neural Networks (RNNs)

- Uses a hidden layer that incorporates current input and prior hidden state
- Accepts variable length sequences
- Optional outputs can be produced at every step
- Inherent parameter sharing

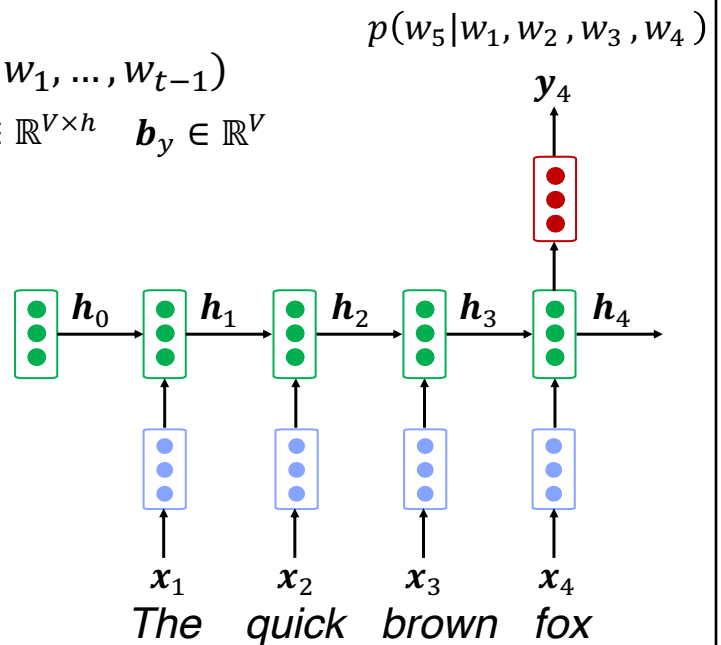$$y_1 \quad y_2 \quad y_3 \quad y_4$$

$$h_0 \quad h_1 \quad h_2 \quad h_3 \quad h_4$$

$$x_1 \quad x_2 \quad x_3 \quad x_4$$

# An RNN LM

$$p(w_5|w_1, w_2, w_3, w_4)$$

- Softmax output computes $p(w_t|w_1, \ldots, w_{t-1})$

$$\boldsymbol{y}_t = softmax(\boldsymbol{W}_{hy}\boldsymbol{h}_t + \boldsymbol{b}_y) \quad \boldsymbol{W}_{hy} \in \mathbb{R}^{V \times h} \quad \boldsymbol{b}_y \in \mathbb{R}^V$$

- Non-linear hidden layer

$$\boldsymbol{h}_t = f(\boldsymbol{W}_{hh}\boldsymbol{h}_{t-1} + \boldsymbol{W}_{xh}\boldsymbol{v}_t + \boldsymbol{b}_h)$$

$$\boldsymbol{W}_{hh} \in \mathbb{R}^{h \times h} \quad \boldsymbol{W}_{xh} \in \mathbb{R}^{h \times d} \quad \boldsymbol{b}_h \in \mathbb{R}^h$$

- Word embedding vectors

$$\boldsymbol{v}_t = \mathbf{E}\boldsymbol{x}_t \quad \boldsymbol{v}_t \in \mathbb{R}^d \quad \boldsymbol{E} \in \mathbb{R}^{d \times V}$$

- One-hot input vectors

$$\boldsymbol{x}_t \in \{\boldsymbol{e}_i : 1 \le i \le V\}$$

$$\boldsymbol{y}_4$$

$$h_0 \quad h_1 \quad h_2 \quad h_3 \quad h_4$$

$$x_1 \quad x_2 \quad x_3 \quad x_4$$

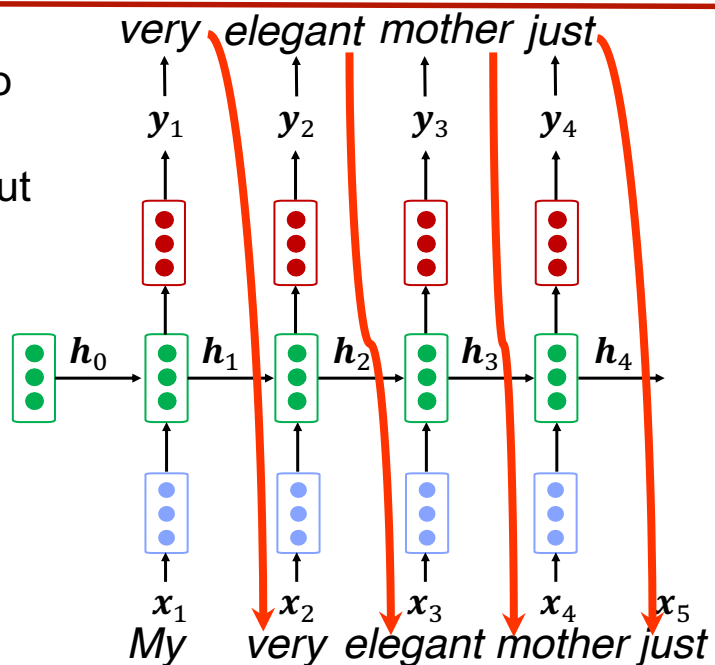*The   quick   brown   fox*

# RNN LMs

- Advantages
  - Can process variable length input
  - No truncated history
  - Model size context-independent
  - Sharing among weights
- Disadvantages
  - Recurrent computation is slow
  - Limitations on how far back it can incorporate context

$y_4$

$h_0$ $h_1$ $h_2$ $h_3$ $h_4$
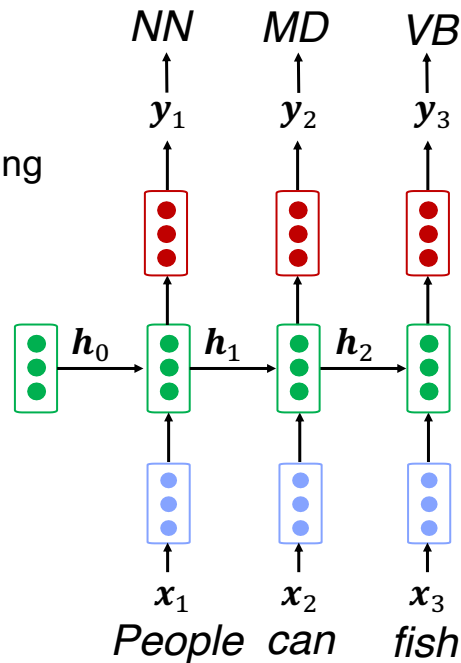
$x_1$ $x_2$ $x_3$ $x_4$

*The quick brown fox*

# RNN-based Language Generation

*very elegant mother just*

- An RNN LM can be used to *generate* text by sampling
- Sampled output is next input

$y_1$ $y_2$ $y_3$ $y_4$

$h_0$ $h_1$ $h_2$ $h_3$ $h_4$

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$

*My very elegant mother just*

# RNN-based Sequence Labeling

- An RNN can be trained to output tags for each word
  - Part-of-speech (POS) tagging
  - Named Entity Recognition

$$NN \qquad MD \qquad VB$$

$$y_1 \qquad y_2 \qquad y_3$$

$h_0 \quad h_1 \quad h_2$

$x_1 \qquad x_2 \qquad x_3$

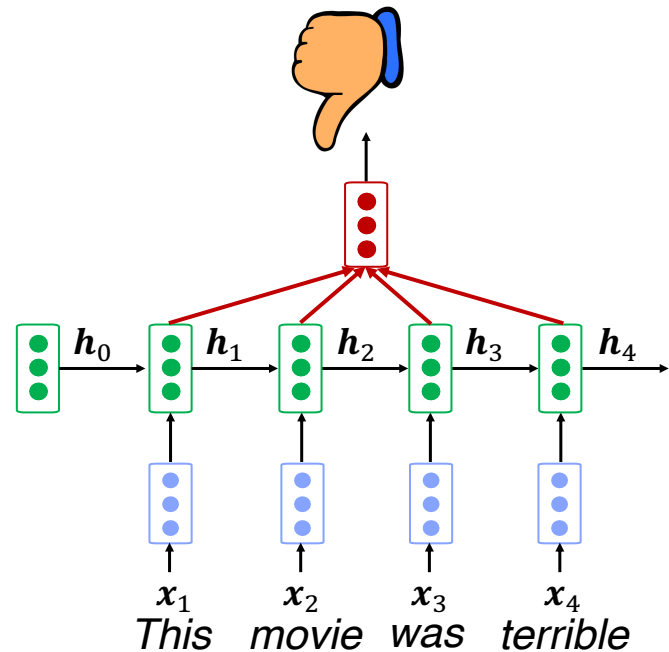*People   can      fish*

# RNN-based Sentence Classification

- Represent sentence as vector
- Use final hidden state as a representation of sentence
- Feed into penultimate layer for sentence classification

$h_0 \quad h_1 \quad h_2 \quad h_3 \quad h_4$

$x_1 \qquad x_2 \qquad x_3 \qquad x_4$

*This   movie   was   terrible*

# An Alternate RNN-based Sentence Classifier

- Represent sentence as vector
- Take element-wise max or mean of all hidden states
  - A simple form of *attention*
- Feed into penultimate layer for sentence classification



$$h_0 \quad h_1 \quad h_2 \quad h_3 \quad h_4$$

$$x_1 \quad x_2 \quad x_3 \quad x_4$$

*This   movie   was   terrible*

# RNN-based Encoder

- Represent sentence as vector
- Take element-wise max or mean of all hidden states
  - A simple form of *attention*
- Feed into subsequent layers for downstream processing
  - Question-answering
  - Machine translation
  - Etc.



$$e$$

$$h_0 \quad h_1 \quad h_2 \quad h_3 \quad h_4$$

$$x_1 \quad x_2 \quad x_3 \quad x_4$$

*Are   there   any   weather*

**Sequence to Sequence Learning with Neural Networks**

Ilya Sutskever
Google
ilyasu@google.com

Oriol Vinyals
Google
vinyals@google.com

Quoc V. Le
Google
qvl@google.com

NeurIPS 2014

**15**

---

# Today's RNN Story

- RNN language models
- **Training RNNs**
- Gated RNNs
  - Gated recurrent units (GRUs)
  - Long short-term memories (LSTMs)
- Other architectures
  - Deep RNNs
  - Bidirectional RNNs

**16**

# Training RNNs

- RNN parameters are learned on a training corpus
- Overall loss is average step-by-step cross-entropy loss

$$L(\boldsymbol{\theta}) = \frac{1}{T}\sum_{t=1}^{T} L_t$$

- For LMs, $L_t$ equivalent to negative log likelihood of next true word
  - Computed by taking dot-product with next one-hot vector

$$L_t = -\log p(w_t|w_1, \cdots, w_{t-1}) = -\log(\boldsymbol{y}_i \cdot \boldsymbol{x}_{t+1})$$

- Parameters optimized via back-propagation and SGD

# RNN LM Training

- The total loss is the sum of all word-by-word losses

$$L(\boldsymbol{\theta}) = \frac{1}{T}\sum_{t=1}^{T} L_t$$

- For SGD, losses are typically accumulated in batches of sentences

# Learning Long-Term Dependencies with Gradient Descent is Difficult

Yoshua Bengio, Patrice Simard, and Paolo Frasconi, *Student Member*, *IEEE*

---

## On the difficulty of training recurrent neural networks

---

**Razvan Pascanu**         PASCANUR@IRO.UMONTREAL.CA
Université de Montréal, 2920, chemin de la Tour, Montréal, Québec, Canada, H3T 1J8

**Tomas Mikolov**         T.MIKOLOV@GMAIL.COM
Speech@FIT, Brno University of Technology, Brno, Czech Republic

**Yoshua Bengio**         YOSHUA.BENGIO@UMONTREAL.CA
Université de Montréal, 2920, chemin de la Tour, Montréal, Québec, Canada, H3T 1J8

ICML 2013

**19**

---

# Training Neural Networks via Backpropagation

- Gradients propagated backwards through network to minimize loss
  - Gradients are accumulated for each parameter in a training batch
- Numerical stability issues for many layered networks
  - Exploding gradients, vanishing gradients
  - Initialization, non-linearity choices affect results
  - Techniques developed for clipping exploding gradients
  - Residual and highway connections etc. help with vanishing gradients

**20**

# Computational Graphs

- Neural networks can be represented as a computational graph, e.g.,
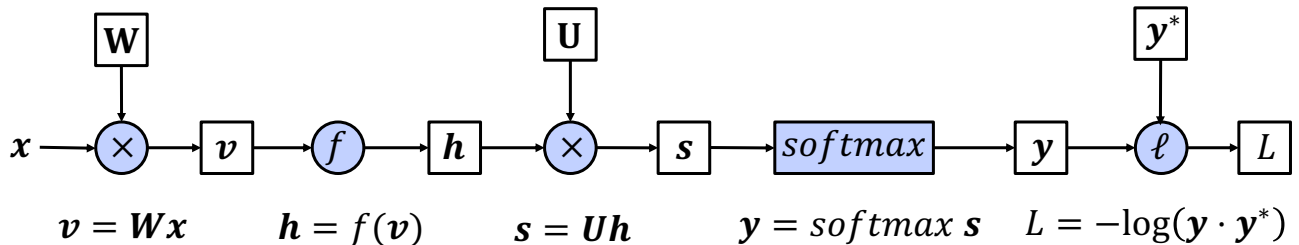
$$v = Wx \qquad h = f(v) \qquad s = Uh$$

$$y = softmax\ s \qquad L = -\log(y \cdot y^*)$$

- For training, each data point takes a forward and backward pass through graph

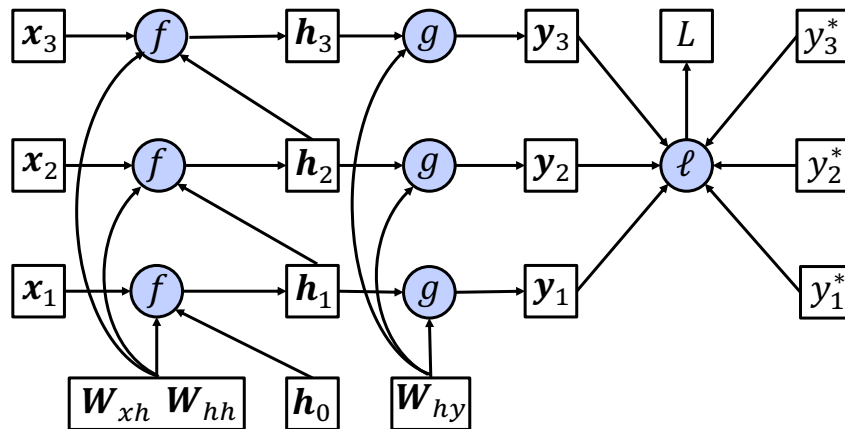  - Gradients are accumulated for each parameter

$$y = softmax(\mathbf{U}h)$$

$$h = f(\mathbf{W}x)$$

# Computational Graph MLP Example



$$v = Wx \qquad h = f(v) \qquad s = Uh \qquad y = softmax\ s \qquad L = -\log(y \cdot y^*)$$

$$\frac{\partial v}{\partial W} \qquad \frac{\partial h}{\partial v} \qquad \frac{\partial s}{\partial h} \quad \frac{\partial s}{\partial U} \qquad \frac{\partial y}{\partial s} \qquad \frac{\partial L}{\partial y}$$

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial s} \times \frac{\partial s}{\partial h} \times \frac{\partial h}{\partial v} \times \frac{\partial v}{\partial W} \qquad\qquad \frac{\partial L}{\partial U} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial s} \times \frac{\partial s}{\partial U}$$

- Gradients are accumulated for each parameter over batch
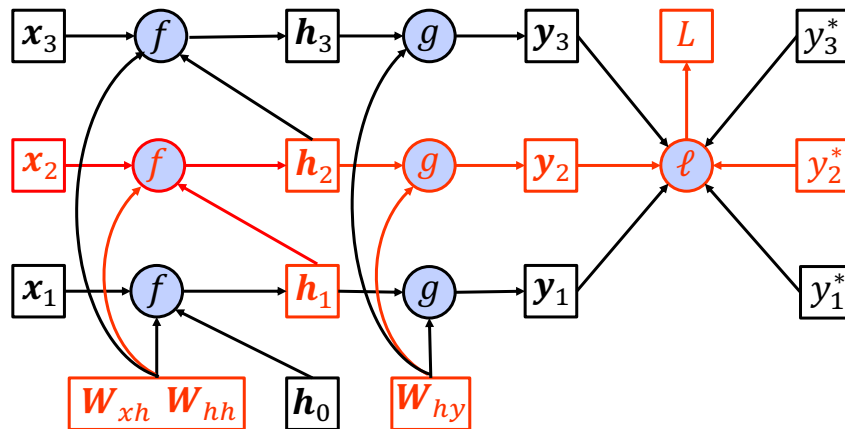
# An RNN Computational Graph

# RNN Forward Pass Step 1
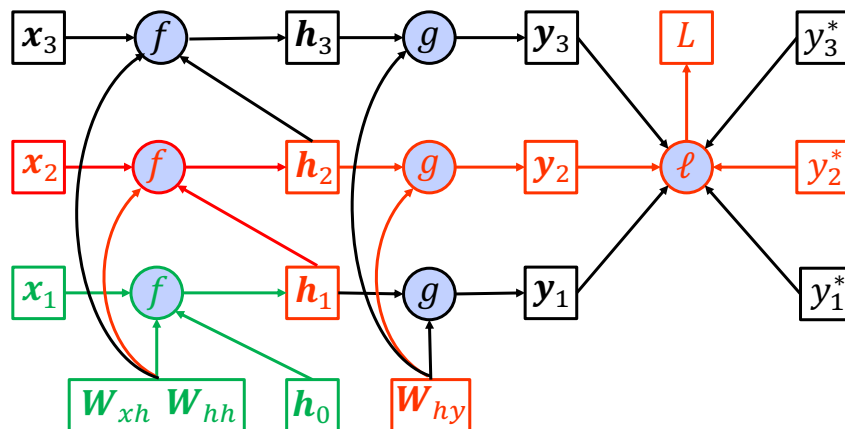


- Backpropagation for step 1 touches the same parameters
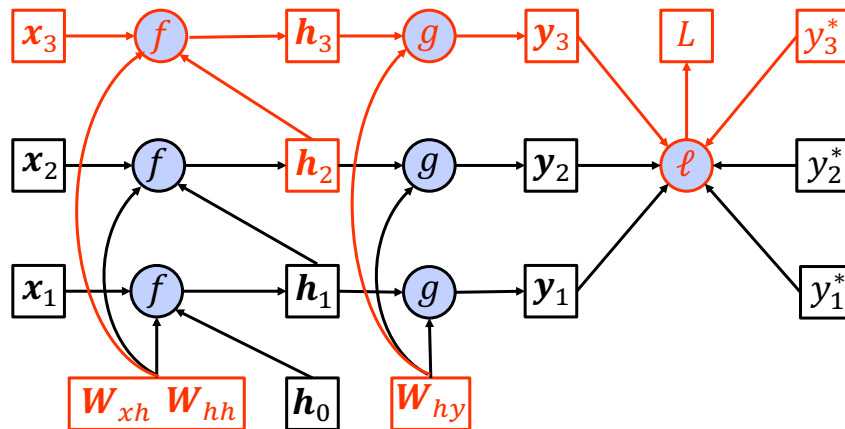
# RNN Forward Pass Step 2

# RNN Backpropagation Pass Step 2



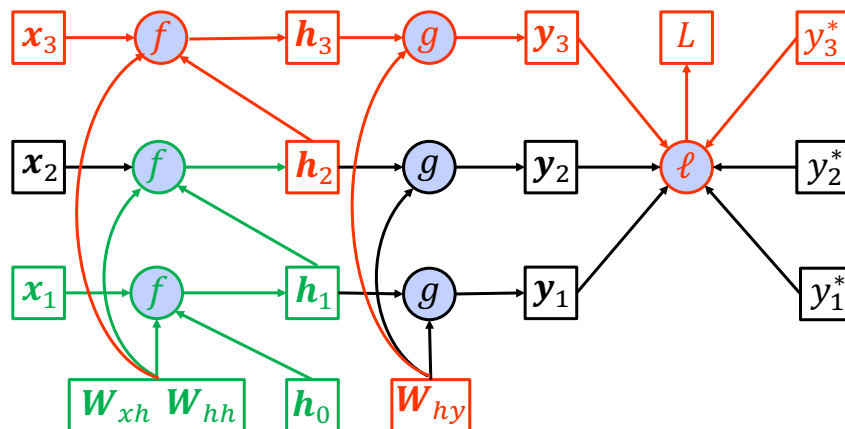- Backpropagation for step 2 must also consider gradients for $h_1$

# RNN Forward Pass Step 3

# RNN Backpropagation Pass Step 3



- Backpropagation for step 3 must also consider partials for $\boldsymbol{h}_2$ and $\boldsymbol{h}_1$

# Backpropagation Through Time

- SGD for RNNs must consider the impact of past inputs and states
  - This process is known as *Backpropagation Through Time* (BPTT)
- The gradients for longer time spans are exponential, e.g.,

$$\frac{\partial \boldsymbol{h}_t}{\partial \boldsymbol{W}_{hh}} = \sum_{i=1}^{t} \left(\boldsymbol{W}_{hh}^T\right)^{t-i} \boldsymbol{h}_i \qquad \frac{\partial \boldsymbol{h}_t}{\partial \boldsymbol{W}_{xh}} = \sum_{i=1}^{t} \left(\boldsymbol{W}_{hh}^T\right)^{t-i} \boldsymbol{v}_i$$

  - Potential for exploding gradients or vanishing gradients
- Since BPTT is computationally intensive for long sequences, sometimes truncated BPTT is used to save computation

29

# Impact of Vanishing Gradients

- Long distance gradients are weaker and have less impact than local gradients
- Model parameters primarily learn local dependencies
- This motivated the search for RNNs that could better model long distance dependencies by some internal memory state



30

# Today's RNN Story

- RNN language models
- Training RNNs
- **Gated RNNs**
  - Gated recurrent units (GRUs)
  - Long short-term memories (LSTMs)
- Other architectures
  - Deep RNNs
  - Bidirectional RNNs

# Gated RNNs

- Conventional RNNs can be challenging to train
  - Long products of matrices lead to vanishing or divergent gradients
  - Effect of BPTT focuses attention on recent history
- Desiderata:
  - A *memory* mechanism to store important information over long distances
  - A *forgetting* mechanism to erase unimportant information from the model
  - A mechanism to *reset* the internal state representation
- A number of alternative RNNs attempt to address these issues
  - Gated RNNs are far more commonly used for sequence labeling
  - One of the earliest is *Long Short-Term Memory* (LSTM) RNNs
  - *Gated Recurrent Unit* (GRU) RNNs are more streamlined and faster

**Learning Phrase Representations using RNN Encoder–Decoder
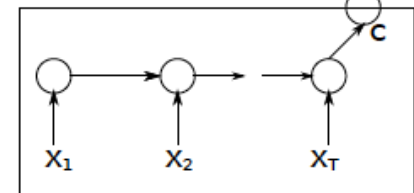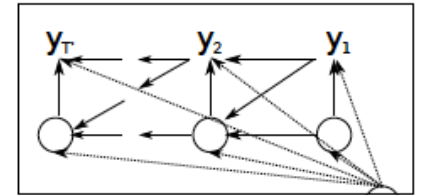for Statistical Machine Translation**

Kyunghyun Cho
Bart van Merriënboer    Caglar Gulcehre
Université de Montréal
firstname.lastname@umontreal.ca

Dzmitry Bahdanau
Jacobs University, Germany
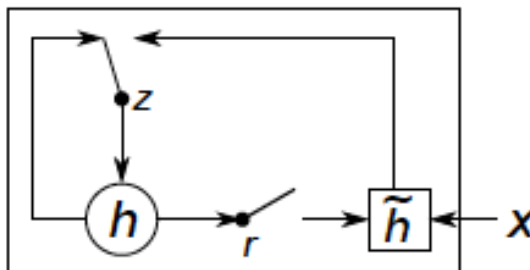d.bahdanau@jacobs-university.de

Fethi Bougares    Holger Schwenk
Université du Maine, France
firstname.lastname@lium.univ-lemans.fr

Yoshua Bengio
Université de Montréal, CIFAR Senior Fellow
find.me@on.the.web

33

# Gated Recurrent Units (GRUs)

- Key distinction between regular RNNs and GRUs is gating
- Dedicated mechanisms for updating and resetting hidden state
  - *Reset* gate controls how much prior state information to remember
  - *Update* gate controls how much new state retains of old state
- Gating mechanisms are a function of current input and prior state



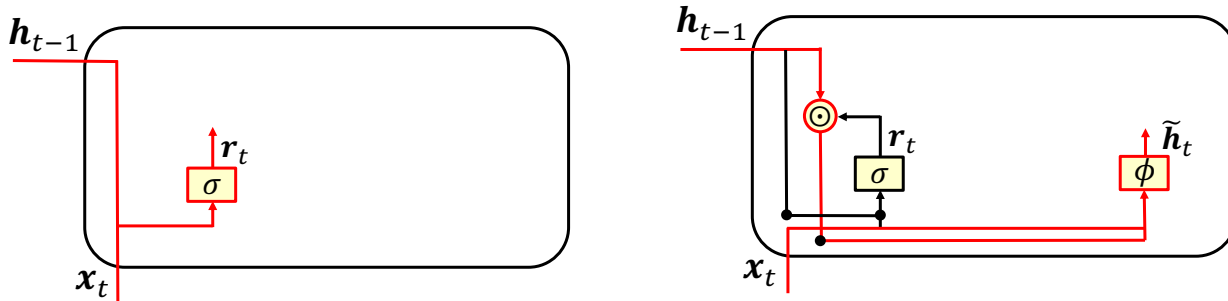⊙ Element-wise multiplication

34

# GRU Reset Gates

- Reset gate is used to reduce influence of $h_{t-1}$ (i.e., reset past)
$$r_t = \sigma(W_{hr}h_{t-1} + W_{xr}x_t + b_r) \qquad r_t \in \mathbb{R}^h$$

- Produces a candidate hidden state by de-weighting prior state
$$\widetilde{h}_t = \tanh(W_{hh}(r_t \odot h_{t-1}) + W_{xh}x_t + b_h)$$

- For $r_t \sim 1$ GRU behaves as RNN; for $r_t \sim 0$ GRU behaves as MLP

# GRU Update Gates

- Determines extent new state is old state vs new candidate state
$$u_t = \sigma(W_{xu}x_t + W_{hu}h_{t-1} + b_u) \qquad u_t \in \mathbb{R}^h$$

- Update gate $u_t$ applied in convex combination with $h_{t-1}$ and $\widetilde{h}_t$
$$h_t = u_t \odot h_{t-1} + (1 - u_t) \odot \widetilde{h}_t$$

- When $u_t \sim 1$ we essentially skip time step t and remember prior state

Sepp Hochreiter
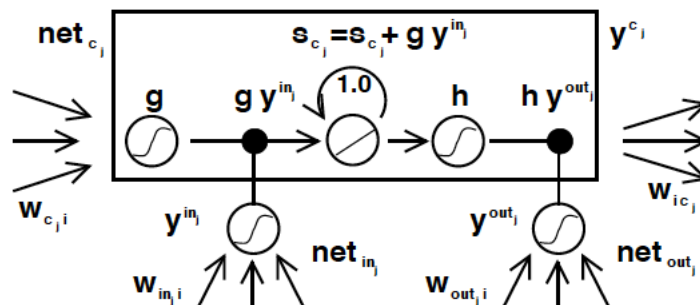Fakultät für Informatik
Technische Universität München
80290 München, Germany
hochreit@informatik.tu-muenchen.de
http://www7.informatik.tu-muenchen.de/˜hochreit

Jürgen Schmidhuber
IDSIA
Corso Elvezia 36
6900 Lugano, Switzerland
juergen@idsia.ch
http://www.idsia.ch/˜juergen
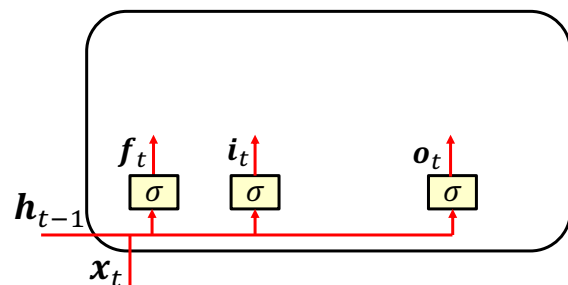
---

# Long Short-Term Memories (LSTMs)

- LSTMs pre-date GRUs, but are slightly more complex (3 gates)
- Inspired by logic gates to control a *memory cell*
  - An *output* gate reads out entries from the cell
  - An input gate is used to read data into a cell
  - A forget gate is used to reset cell contents
- All three gates are a function of current input and prior state

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

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

$$o_t = \sigma(W_{ho}h_{t-1} + W_{xo}x_t + +b_o)$$
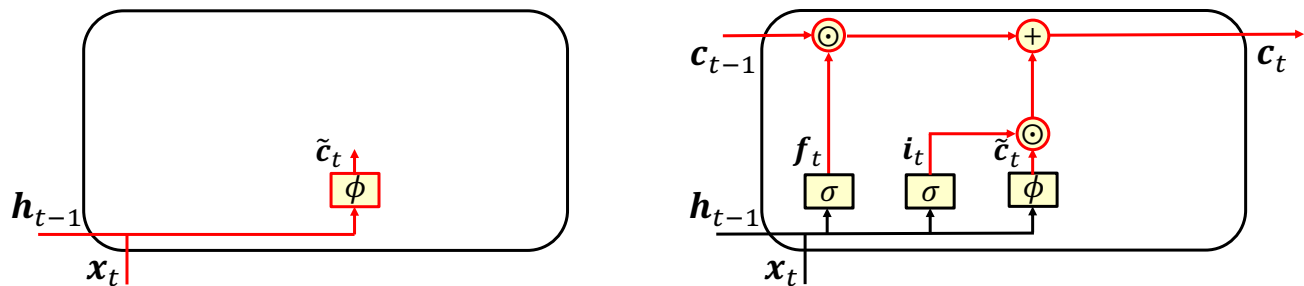
$$i_t, f_t, o_t \in \mathbb{R}^h$$

# LSTM Memory Cell

- A *candidate* memory cell is based on a regular RNN hidden state
$$\tilde{c}_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + b_h)$$
- The input and forget gates are used to create the new memory cell
$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$$
  - $i_t$ controls how much new information to take into account
  - $f_t$ controls how much old information to retain
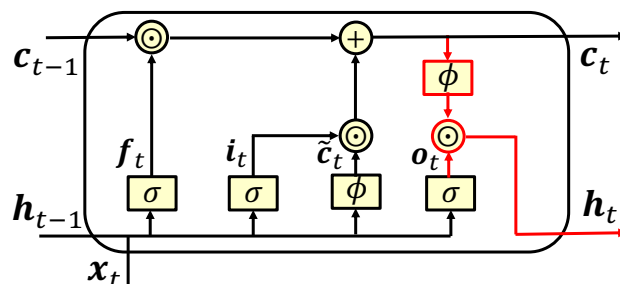  - If $f_t \sim 1$ and $i_t \sim 0$ then old information will be retained

# LSTM Hidden State

- The new hidden state is partially read from the new memory cell
$$h_t = o_t \odot \tanh(c_t)$$
  - Amount retained in $h_t$ is controlled by output gate
  - tanh ensures $h_t$ spans interval (-1,1)
- If $o_t \sim 1$ pass all information through to prediction for next time step

# GRUs vs LSTMs

- Both gated RNNs are much better able to maintain information over many timesteps compared to a vanilla RNN
- Both models have been very effective on many NLP tasks
  - LSTMs attained state-of-the-art results in the 2013-2015 time frame
  - GRUs are newer models, but have also achieved good results
- The LSTM memory cell is not bounded like the hidden state, and has demonstrated an excellent ability to count etc.
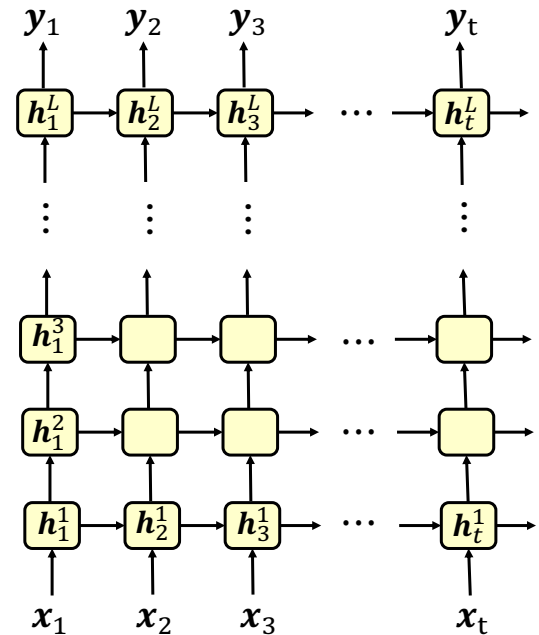- The GRU has fewer parameters and is faster than the LSTM

41

# Today's RNN Story

- RNN language models
- Training RNNs
- Gated RNNs
  - Gated recurrent units (GRUs)
  - Long short-term memories (LSTMs)
- Other architectures
  - Deep RNNs
  - Bidirectional RNNs

42

# Deep RNNs

- The conventional DNN is a feed-forward network
- RNNs can be stacked in multiple layers on top of each other
  - The hidden state $\boldsymbol{h}_t$ from lower RNN becomes input to the upper RNN
  - The topmost layer is responsible for generating any outputs $\boldsymbol{y}_t$
- Deep RNN layers can potentially focus on different information
  - Much more flexibility than HMMs

---

# Bidirectional RNNs

- Bidirectional RNNs consist of two RNNs
  - $\overrightarrow{RNN}$ runs in a forward direction starting from the beginning
  $$\overrightarrow{\boldsymbol{h}}_t = \phi(\overrightarrow{\mathbf{W}}_{hh}\overrightarrow{\boldsymbol{h}}_{t-1} + \overrightarrow{\mathbf{W}}_{xh}\boldsymbol{x}_t + \overrightarrow{\mathbf{b}}_h)$$
  - $\overleftarrow{RNN}$ runs in a backward direction starting from the end
  $$\overleftarrow{\boldsymbol{h}}_t = \phi\left(\overleftarrow{\mathbf{W}}_{hh}\overleftarrow{\boldsymbol{h}}_{t+1} + \overleftarrow{\mathbf{W}}_{xh}\boldsymbol{x}_t + \overleftarrow{\mathbf{b}}_h\right)$$
- At time step $t$, the hidden state is the concatenation of $\overrightarrow{\boldsymbol{h}}_t$ and $\overleftarrow{\boldsymbol{h}}_t$
  $$\boldsymbol{h}_t = \left[\overrightarrow{\boldsymbol{h}}_t; \overleftarrow{\boldsymbol{h}}_t\right] \qquad \boldsymbol{h}_t \in \mathbb{R}^{2h}$$
- The bidirectional RNN output is computed like vanilla RNNs
  $$\boldsymbol{y}_t = \phi(\mathbf{W}_{hq}\boldsymbol{h}_t + \mathbf{b}_q) \qquad \mathbf{W}_{hq} \in \mathbb{R}^{q \times 2h}$$

## Bidirectional RNNs

- Provides fuller context for each input token
- Bidirectional RNNs have achieved very good performance
- Requires access to entire input label sequence

## References

- Readings:
  - Jurafsky & Martin, "Speech and Language Processing," 2020 (RNNs 9.2-9.3)