

Language Representation and Recurrent Nets

Kapil Thadani
kapil@cs.columbia.edu



Outline

- o Language representation
 - Bag of words
 - Distributional hypothesis
 - Word embeddings: word2vec, GloVe
 - Beyond words: paragraph vector
- o Recurrent neural networks
 - Backpropagation through time
 - Long short-term memory (LSTM)
 - Gated recurrent units (GRU)
- o NLP scenarios
 - Classification
 - Tagging
 - Generation
 - Text-to-text

Formal language

(i) Set of sequences over symbols from an alphabet

sentences

words

vocabulary

utterances

morphemes

documents

MWEs

{

{

(ii) Rules for valid sequences

spelling orthography, morphology

grammar syntax

meaning semantics, discourse, pragmatics, ...

Natural language

- (i) Set of sequences over symbols from an alphabet

sentences	words	vocabulary
utterances	morphemes	
documents	MWEs	
⋮	⋮	

- (ii) Rules for valid sequences

spelling	orthography, morphology
grammar	syntax
meaning	semantics, discourse, pragmatics, ⋯

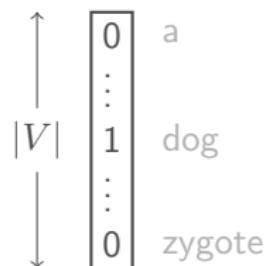
Sparse representations

Word: one-hot ($1\text{-of-}V$) vectors

Document: "bag of words"

Emphasize rare words with inverse document frequency (IDF)

Compare documents with cosine similarity



- + Simple and interpretable
- No notion of word order
- No implicit semantics
- Curse of dimensionality with large $|V|$

Sparse representations

Consider a neural network to read documents with:

- at most T words
- drawn from vocabulary V
- into a hidden layer with H units

How many parameters in the input layer for:

- Tweets?
- News stories?

Sparse representations

Consider a neural network to read documents with:

- at most T words
- drawn from vocabulary V
- into a hidden layer with H units

How many parameters in the input layer for:

- Tweets?

$$T = 50, |V| = 100K, H = 100 \rightarrow 0.5B$$

- News stories?

Sparse representations

Consider a neural network to read documents with:

- at most T words
- drawn from vocabulary V
- into a hidden layer with H units

How many parameters in the input layer for:

- Tweets?

$$T = 50, |V| = 100K, H = 100 \rightarrow 0.5B$$

- News stories?

Sparse representations

Consider a neural network to read documents with:

- at most T words
- drawn from vocabulary V
- into a hidden layer with H units

How many parameters in the input layer for:

- Tweets?

$$T = 50, |V| = 100K, H = 100 \rightarrow \text{0.5B}$$

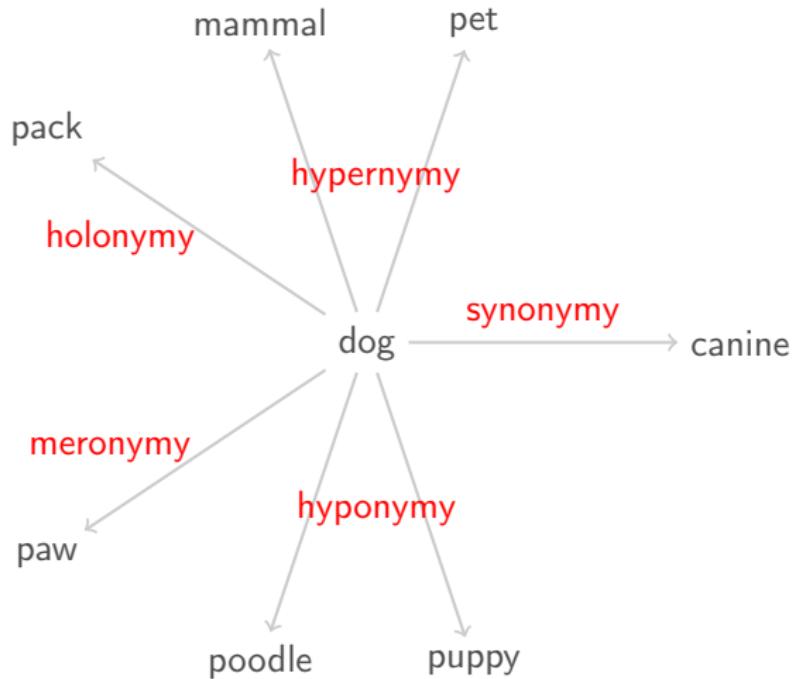
- News stories?

$$T = 2000, |V| = 200K, H = 100 \rightarrow \text{40B}$$

Lexical semantics

dog

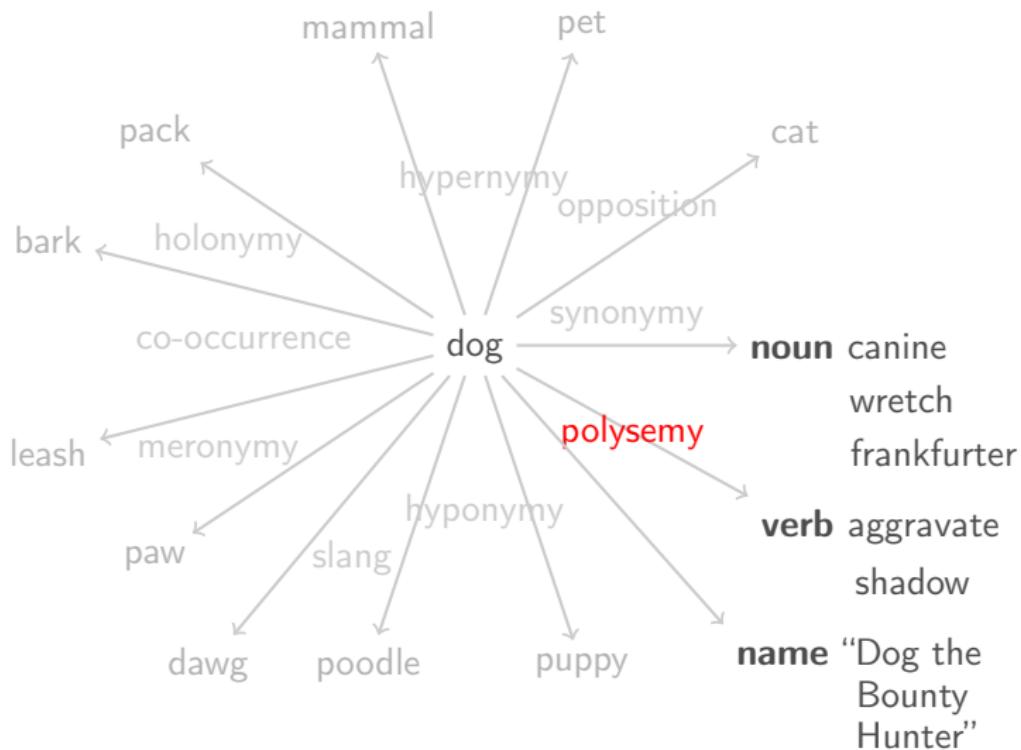
Lexical semantics



Lexical semantics



Lexical semantics



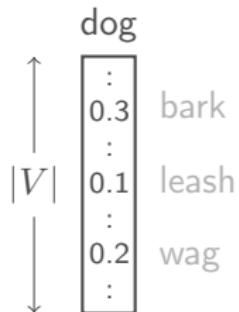
Distributional approaches

Words that occur in similar contexts have similar meanings

e.g., record word co-occurrence within a context window over a large corpus

Weight association with pointwise mutual information (PMI), etc

$$PMI(w_1, w_2) = \log_2 \frac{p(w_1, w_2)}{p(w_1)p(w_2)}$$



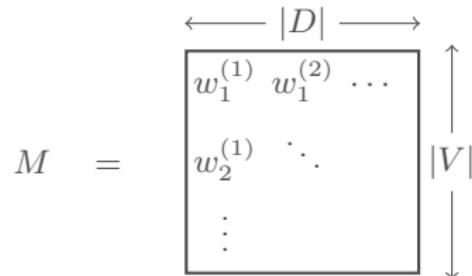
- + Implicit semantics, i.e., related words have similar representations
- Domain dependence on training corpus
- Curse of dimensionality with large $|V|$

Latent Semantic Analysis

Deerwester et al. (1990)

Indexing by Latent Semantic Analysis

Construct term-document matrix



Singular value decomposition

$$M \approx \begin{matrix} & \boxed{u_1} & \boxed{u_2} & \boxed{u_3} & \dots & & \\ M \approx & & & & & & \\ & k & & & & & \\ & & & & \boxed{\lambda_1} & \boxed{\lambda_2} & \lambda_3 & \ddots & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \end{matrix}$$

Select top k singular vectors for k -dim embeddings of words/docs

word2vec

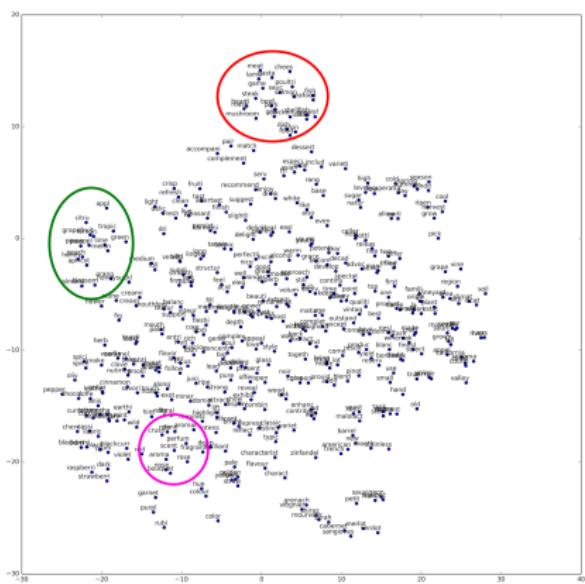
Mikolov et al. (2013)

Efficient Estimation of Word Representations in Vector Space

word2vec

Mikolov et al. (2013)

Efficient Estimation of Word Representations in Vector Space

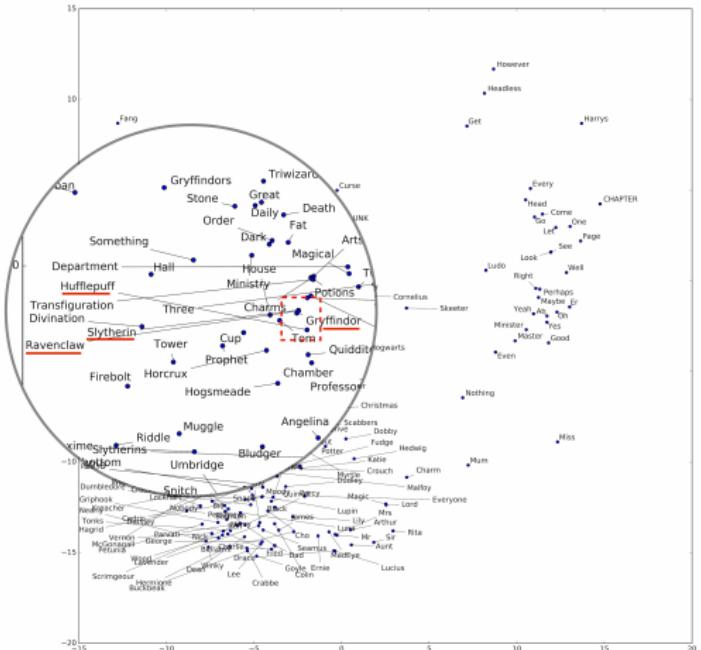


t-SNE projection of word embeddings from 58k Winemaker's Notes

<http://methodmatters.blogspot.com/2017/11/using-word2vec-to-analyze-word.html>

word2vec

Efficient Estimation of Word Representations in Vector Space



t-SNE projection of name embeddings from all 7 Harry Potter books

<https://github.com/nchah/word2vec4everything>

word2vec

Mikolov et al. (2013)

Efficient Estimation of Word Representations in Vector Space

◦ sir

◦ king

◦ madam

◦ queen

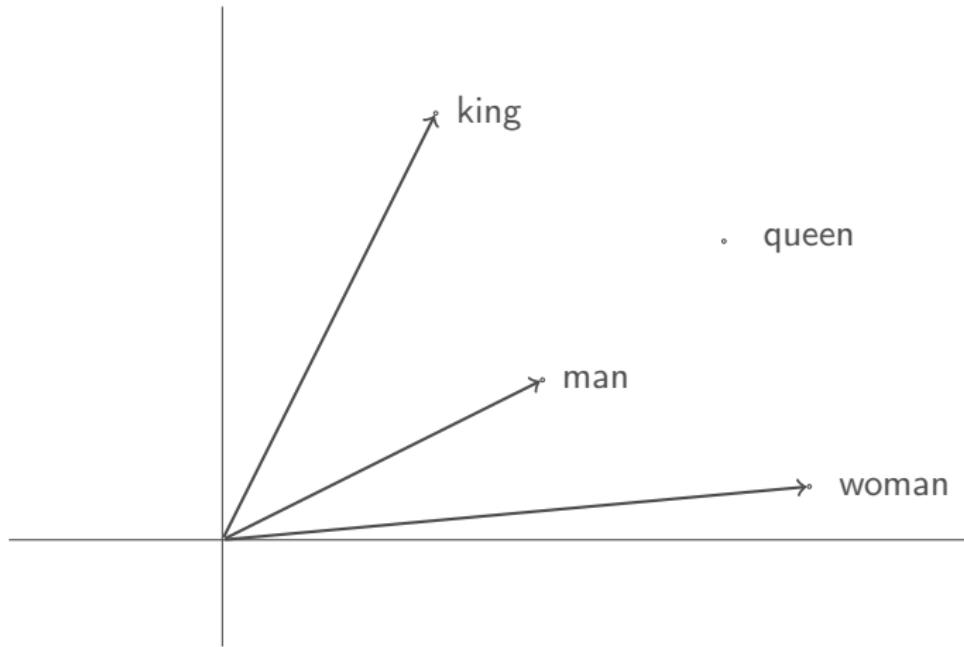
◦ man

◦ uncle

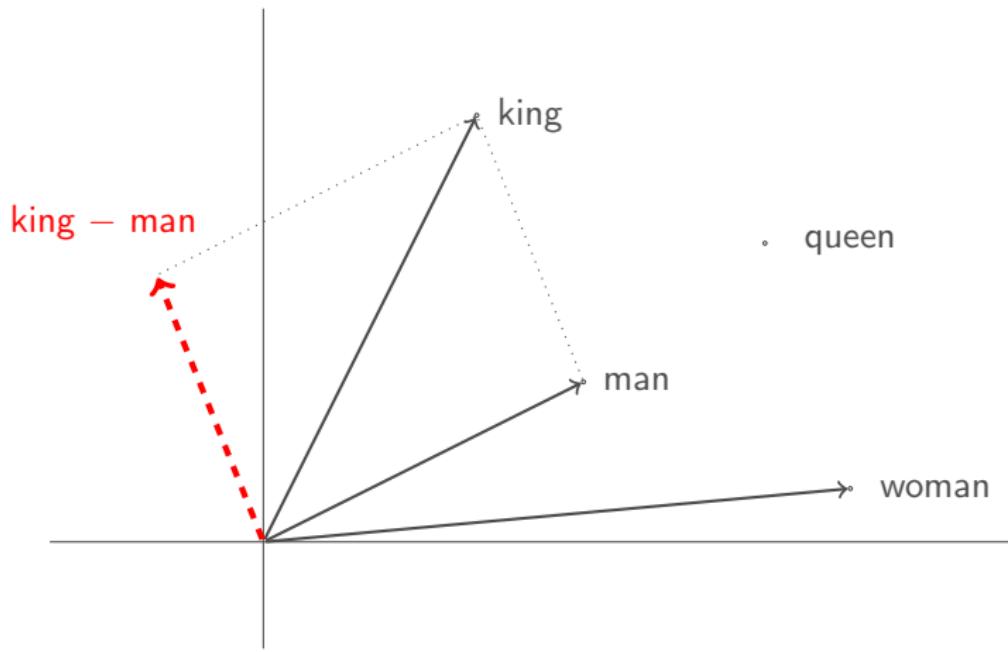
◦ woman

◦ aunt

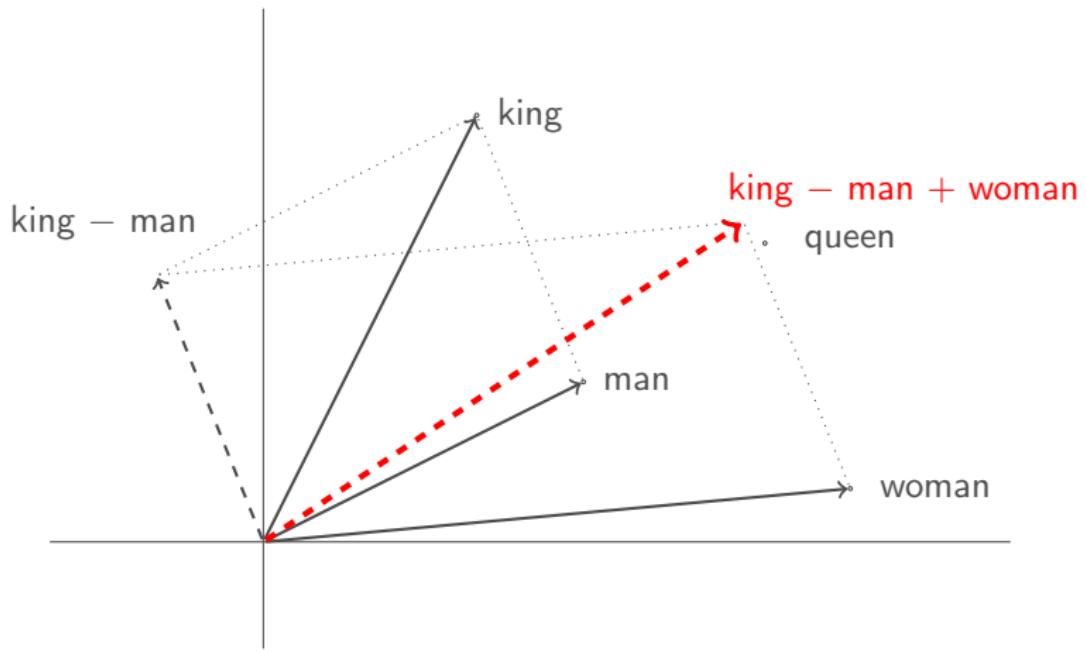
Linear substructure for related words



Linear substructure for related words



Linear substructure for related words



Linear substructure for related words

Efficient Estimation of Word Representations in Vector Space

Relationship	Example 1	Example 2	Example 3
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee
big - bigger	small: larger	cold: colder	quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy - France	Berlusconi: Italy	Merkel: Germany	Koizumi: Japan
copper - Cu	zinc: Zn	gold: Au	uranium: plutonium
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone
Microsoft - Ballmer	Google: Yahoo	IBM: McNealy	Apple: Jobs
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza

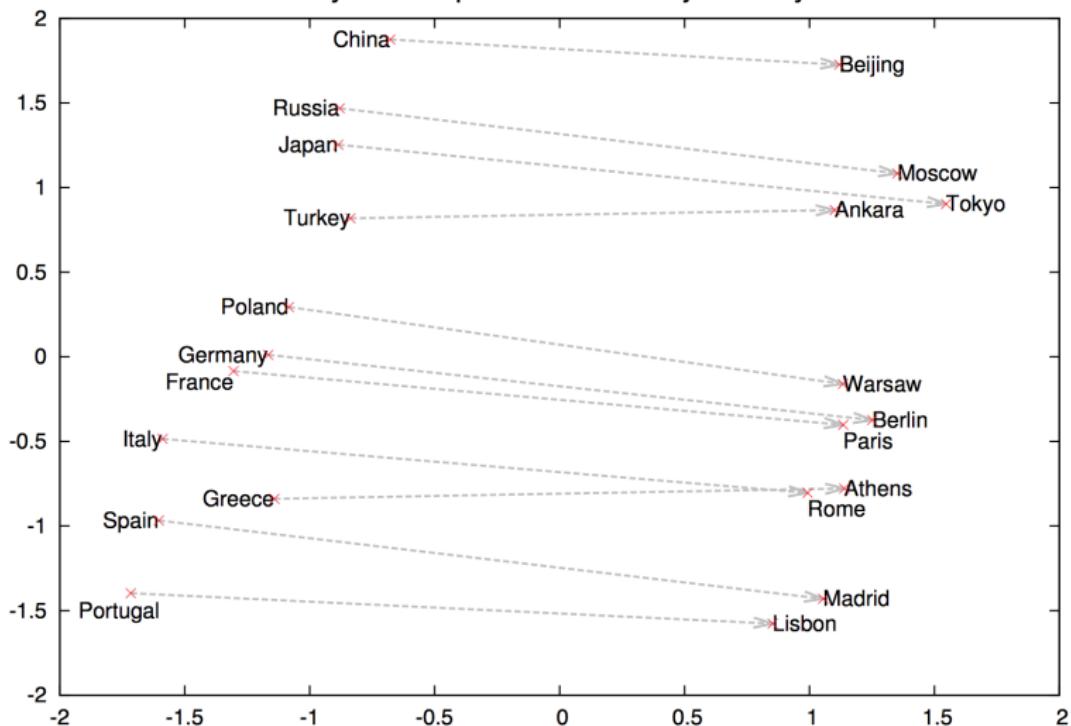
Analogical reasoning

word2vec

Mikolov et al. (2013)

Distributed Representations of Words and Phrases and their Compositionality

Country and Capital Vectors Projected by PCA



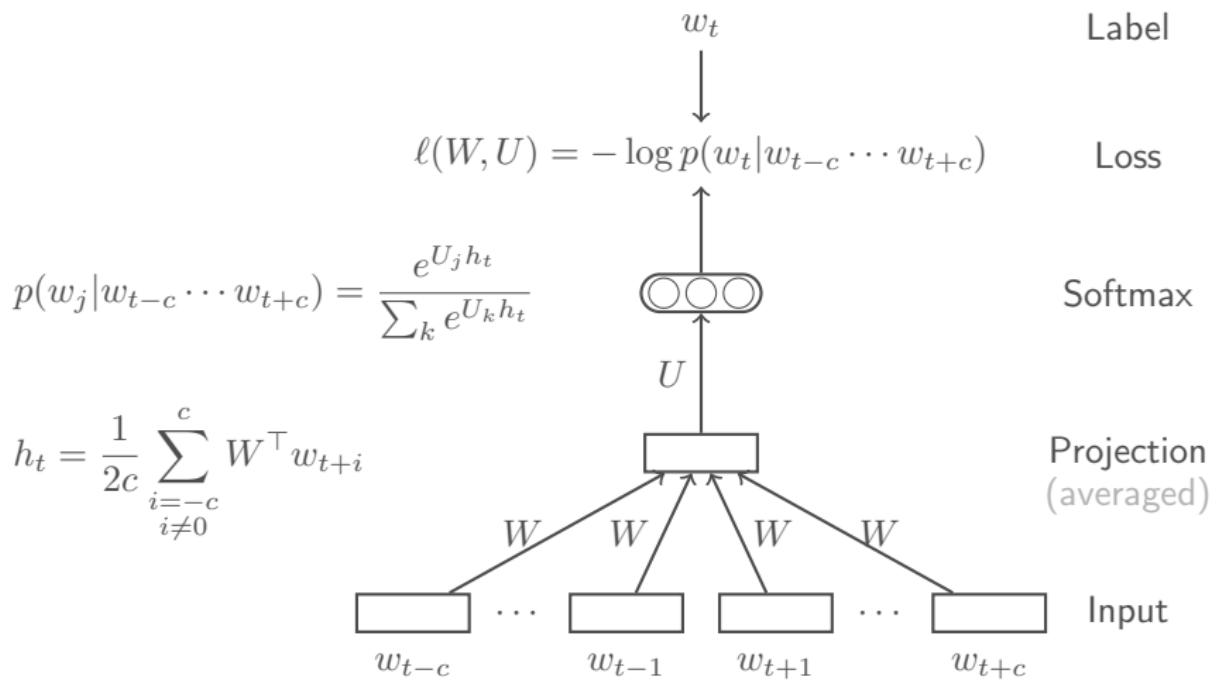
Visualizing lexical relationships

word2vec

Mikolov et al. (2013)

Continuous Bag-of-Words (CBOW)

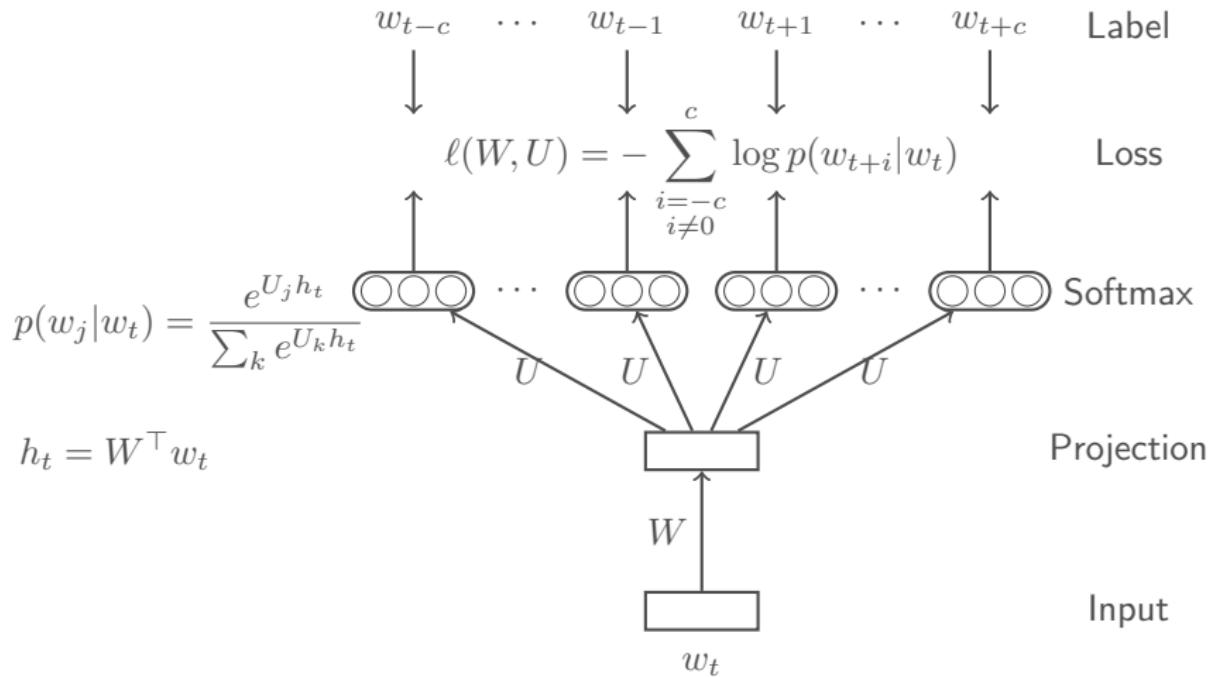
- Predict target w_t given context $w_{t-c}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+c}$



word2vec

Skip-gram

- Predict context $w_{t-c}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+c}$ given target w_t

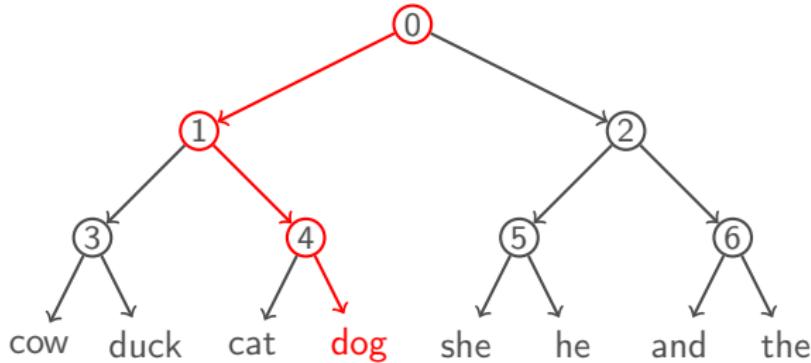


Cost of computing $\nabla p(w_j | \dots)$ is proportional to $V!$

Alternative 1: Hierarchical softmax

- Predict path in binary tree representation of output layer
- Reduces to $\log_2(V)$ binary decisions

$$p(w_t = \text{"dog"} | \dots) = (1 - \sigma(U_0 h_t)) \times \sigma(U_1 h_t) \times \sigma(U_4 h_t)$$



Cost of computing $\nabla p(w_j | \dots)$ is proportional to $V!$

Alternative 2: Negative sampling

- Change objective to differentiate target vector from noisy samples with logistic regression

$$\max \log \sigma(u_j^\top h_t) + \sum_{k=1}^K \mathbb{E}_{w_m \sim \Psi} \log \sigma(-u_m^\top h_t)$$

where $u_j = U_j = j$ 'th column of U

and $w_j \in \text{context}(w_t)$

- Noise distribution Ψ typically unigram, uniform or in between
- Number of samples K typically 5–20

Skip-gram with negative sampling increases $u_j^\top h_t$ for real word-context pairs $\langle w_t, w_j \rangle$ and decreases it for noise pairs

Given:

- a matrix of d -dim word vectors W ($|V_w| \times d$)
- a matrix of d -dim context vectors U ($|V_u| \times d$)

Skip-gram is implicitly factorizing the matrix $M = WU^\top$

What is M ?

- Word-context matrix where each cell (i, j) contains $PMI(w_i, w_j)$
- If number of negative samples $K > 1$, this is shifted by a constant $-\log K$
- (Assuming large enough d and iterations)

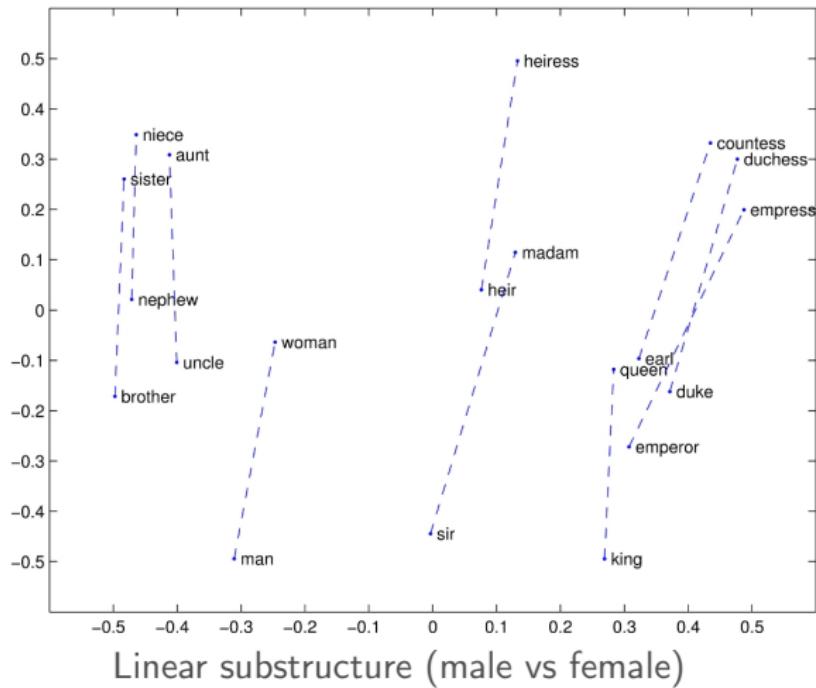
Similar words have similar **ratios of co-occurrence probabilities** for context words

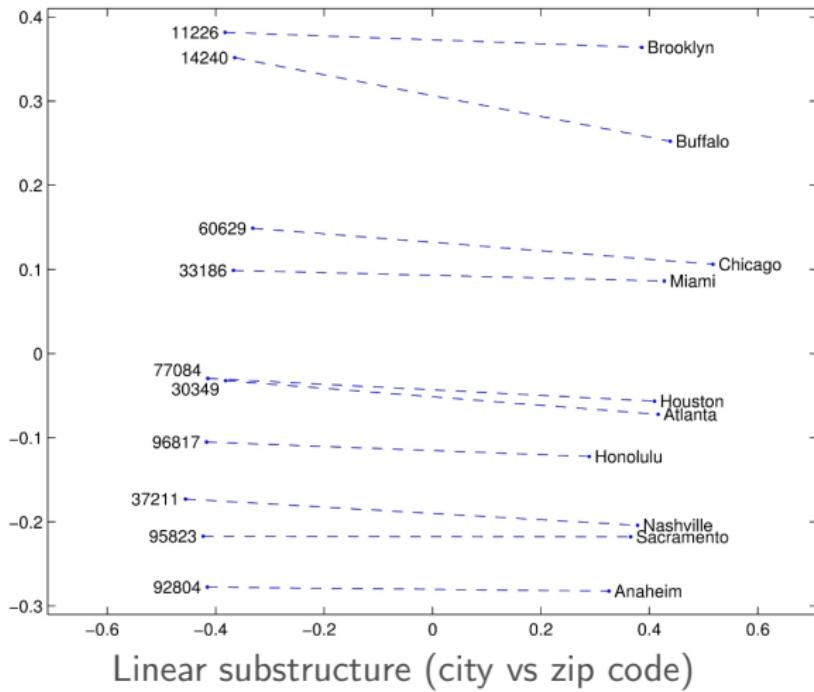
$$\begin{aligned} f((u_i - u_j)^\top \tilde{u}_k) &= \frac{p(w_k|w_i)}{p(w_k|w_j)} \\ \Rightarrow u_i^\top \tilde{u}_k + b_i + \tilde{b}_k &= \log \frac{\text{count}(w_i, w_k)}{\text{count}(w_i)} \end{aligned}$$

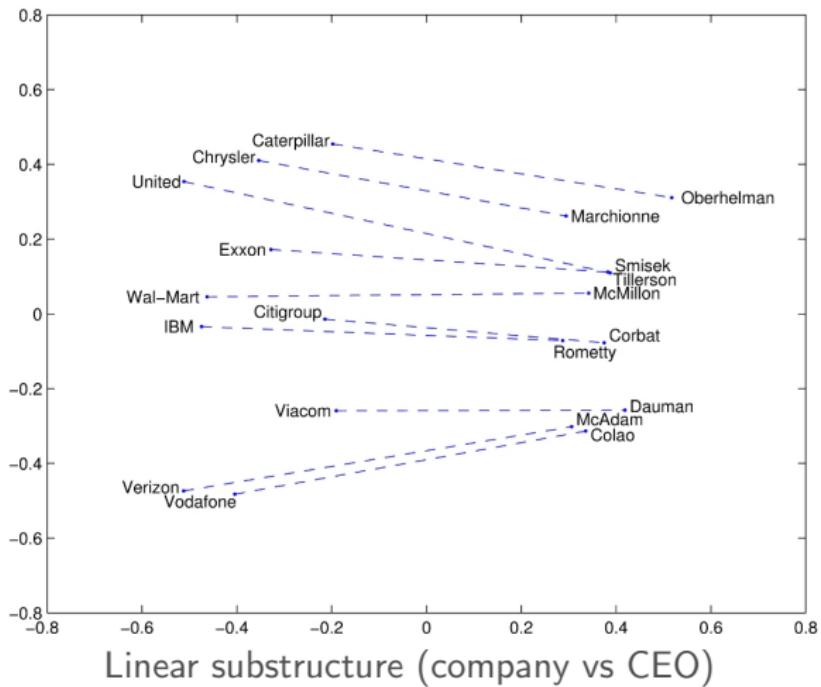
- + Explicitly encodes linear substructure between similar words
- + Scales to huge corpora

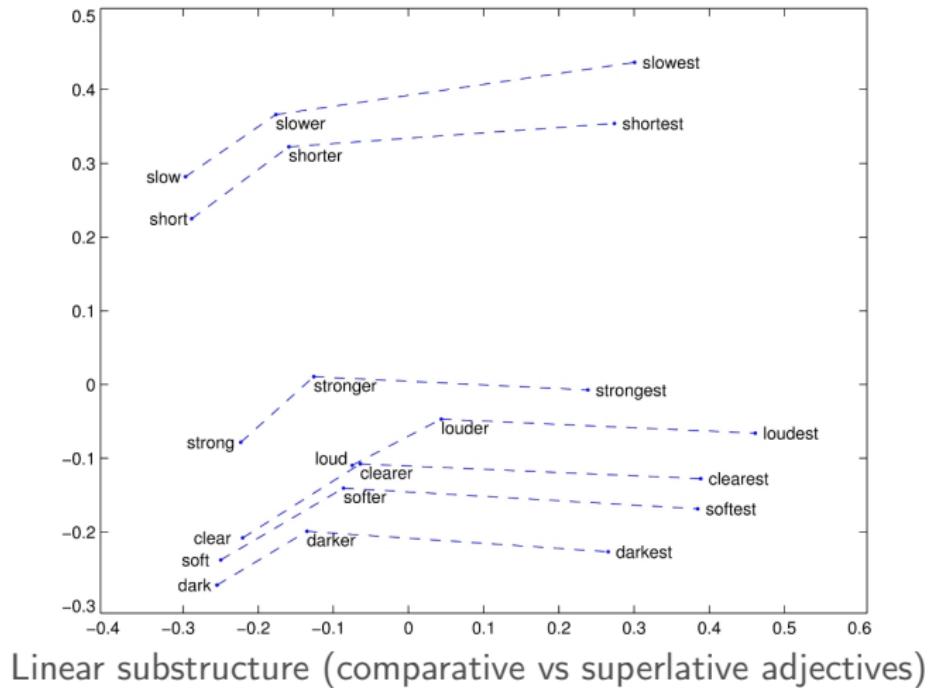
Pre-trained embeddings at <https://nlp.stanford.edu/projects/glove/>

- Common Crawl: 840B tokens, $|V| = 2.2M$, $d = 300$
- Twitter: 27B tokens, $|V| = 1.2M$, $d = 25 - 200$









word2vec with phrases

Mikolov et al. (2013)

Distributed Representations of Words and Phrases and their Compositionality

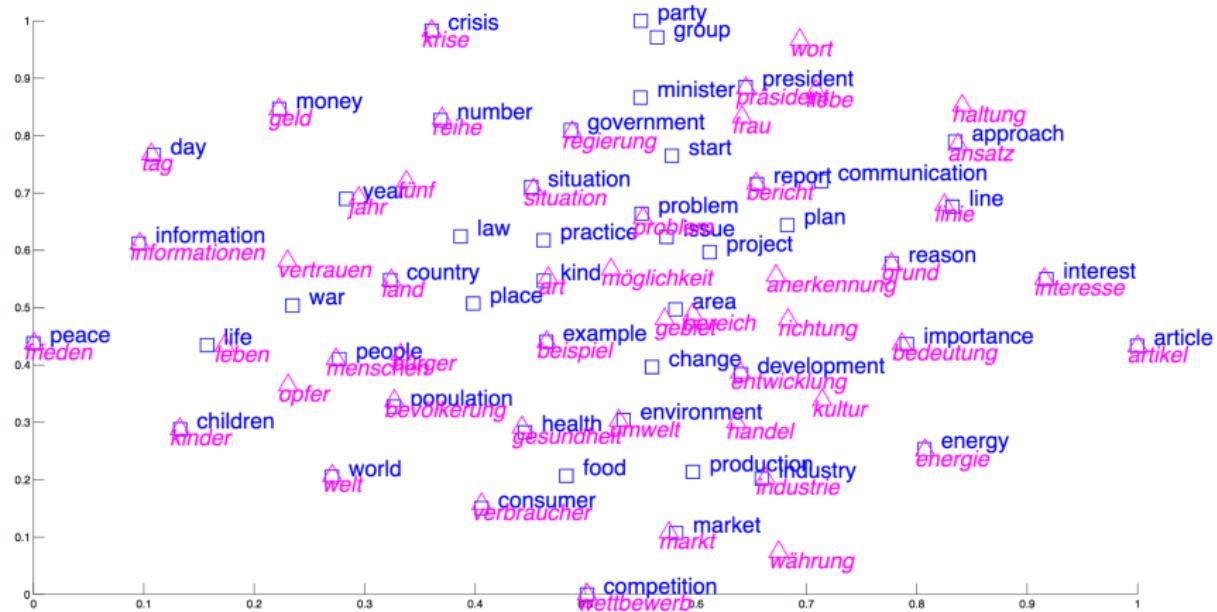
Newspapers			
New York	New York Times	Baltimore	Baltimore Sun
San Jose	San Jose Mercury News	Cincinnati	Cincinnati Enquirer
NHL Teams			
Boston	Boston Bruins	Montreal	Montreal Canadiens
Phoenix	Phoenix Coyotes	Nashville	Nashville Predators
NBA Teams			
Detroit	Detroit Pistons	Toronto	Toronto Raptors
Oakland	Golden State Warriors	Memphis	Memphis Grizzlies
Airlines			
Austria	Austrian Airlines	Spain	Spainair
Belgium	Brussels Airlines	Greece	Aegean Airlines
Company executives			
Steve Ballmer	Microsoft	Larry Page	Google
Samuel J. Palmisano	IBM	Werner Vogels	Amazon

Phrase analogies

Czech + currency	Vietnam + capital	German + airlines	Russian + river	French + actress
koruna	Hanoi	airline Lufthansa	Moscow	Juliette Binoche
Check crown	Ho Chi Minh City	carrier Lufthansa	Volga River	Vanessa Paradis
Polish zolty	Viet Nam	flag carrier Lufthansa	upriver	Charlotte Gainsbourg
CTK	Vietnamese	Lufthansa	Russia	Cecile De

Additive compositionality

Bilingual word embeddings



Aligned embeddings for English and German (Luong et al., 2015)

Resources for word embeddings

Original code and pre-trained embeddings:

<https://code.google.com/archive/p/word2vec/>

Python library for word and document embeddings:

<https://radimrehurek.com/gensim/models/word2vec.html>

Tensorflow tutorial and implementation:

<https://www.tensorflow.org/tutorials/representation/word2vec>

FastText library and pre-trained embeddings for 157 languages:

<https://github.com/facebookresearch/fastText>

Beyond words

Can we add word vectors to make sentence/paragraph/doc vectors?

$$\text{doc } A = a_1 + a_2 + a_3$$

$$\text{doc } B = b_1 + b_2 + b_3$$

$$\begin{aligned} \cos(A, B) &= \frac{A \cdot B}{\|A\| \cdot \|B\|} \\ &= \frac{1}{\|A\| \cdot \|B\|} (a_1 \cdot b_1 + a_1 \cdot b_2 + a_1 \cdot b_3 + \\ &\quad a_2 \cdot b_1 + a_2 \cdot b_2 + a_2 \cdot b_3 + \\ &\quad a_3 \cdot b_1 + a_3 \cdot b_2 + a_3 \cdot b_3) \\ &= \text{weighted all-pairs similarity over } A \text{ and } B \end{aligned}$$

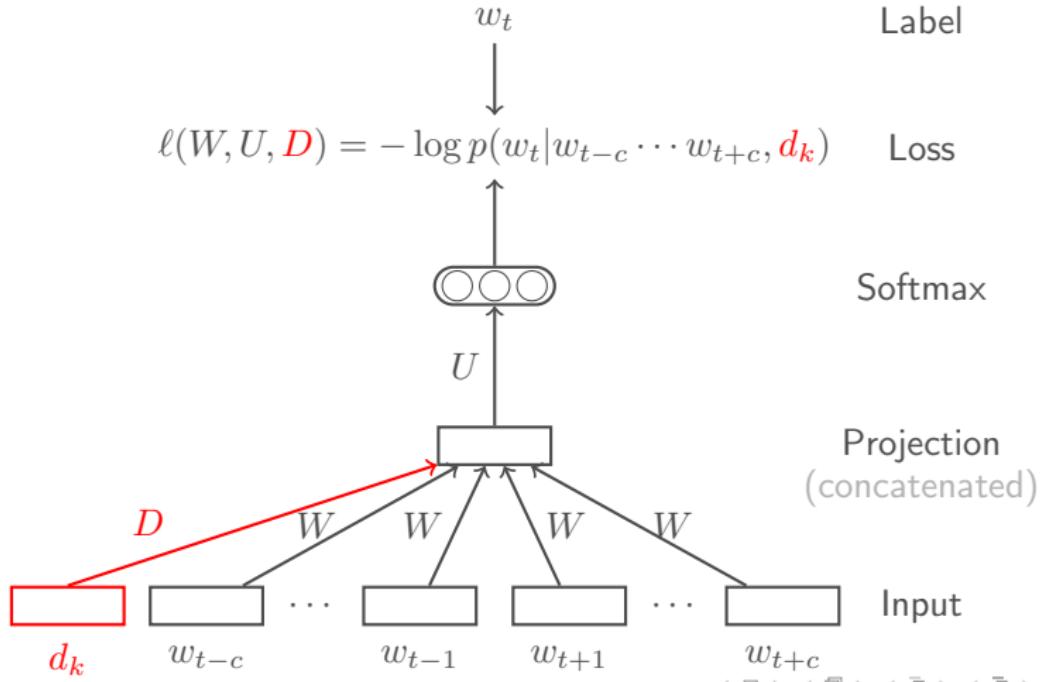
Paragraph vector (a.k.a doc2vec)

Le & Mikolov (2014)

Distributed memory

Distributed Representations of Sentences and Documents

- Predict target w_t given context $w_{t-c}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+c}$ and doc label d_k
 - At test time, hold U, W fixed and back-prop into expanded D



Paragraph vector (a.k.a doc2vec)

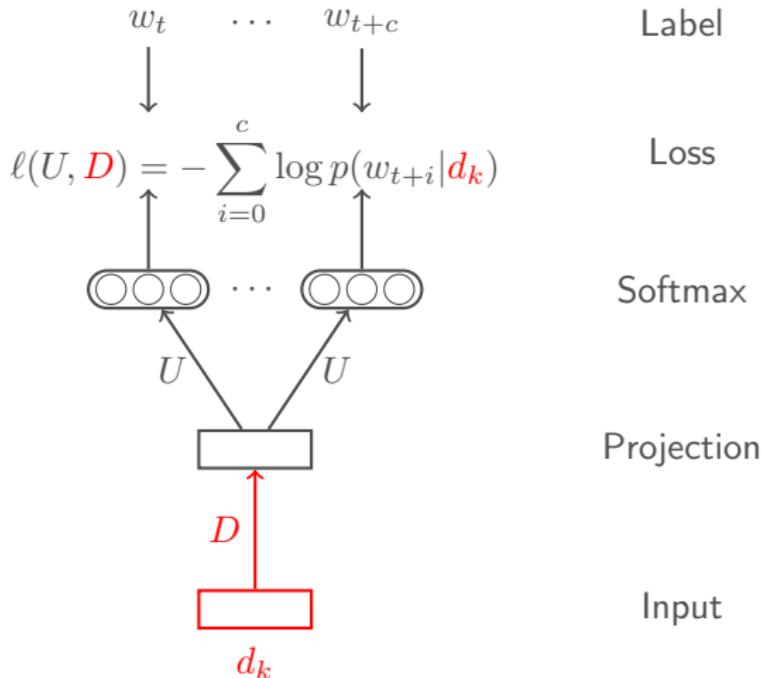
Le & Mikolov (2014)

22

Distributed Representations of Sentences and Documents

Distributed Bag-of-Words (DBOW)

- Predict target n-grams w_t, \dots, w_{t+c} given doc label d_k
 - At test time, hold U fixed and back-prop into expanded D



Semantics are elusive

Visitors saw her duck with binoculars .

Did she duck or does she *have* a duck?

Who has the binoculars?

How many pairs of binoculars are there?

Semantics are elusive

Visitors saw her duck with binoculars .

Did she duck or does she *have* a duck?

Who has the binoculars?

How many pairs of binoculars are there?

Semantics are elusive

Visitors saw her duck with binoculars .
NNS VBD PRP\$ NN IN NNS .

Did she duck or does she *have* a duck?

Who has the binoculars?

How many pairs of binoculars are there?

Semantics are elusive

Visitors saw her duck with binoculars .
NNS VBD PRP VB IN NNS .

Did she duck or does she *have* a duck?

Who has the binoculars?

How many pairs of binoculars are there?

Semantics are elusive

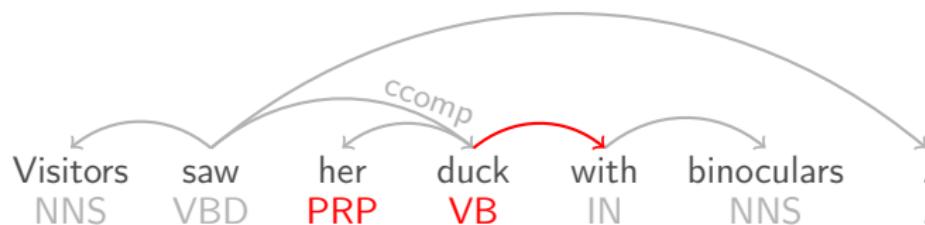
Visitors saw her duck with binoculars .
NNS VBD PRP VB IN NNS .

Did she duck or does she *have* a duck?

Who has the binoculars?

How many pairs of binoculars are there?

Semantics are elusive

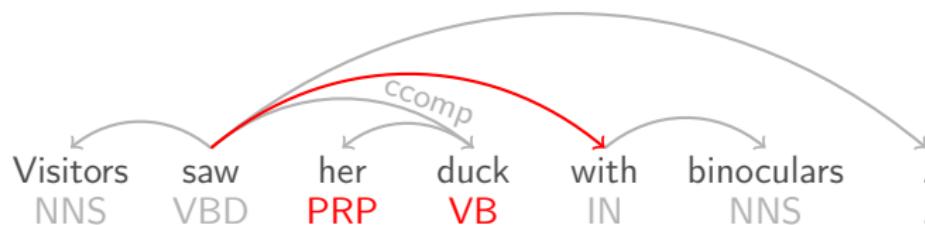


Did she duck or does she *have* a duck?

Who has the binoculars?

How many pairs of binoculars are there?

Semantics are elusive

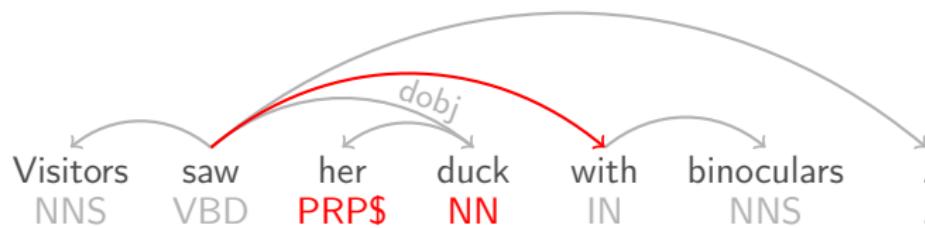


Did she duck or does she *have* a duck?

Who has the binoculars?

How many pairs of binoculars are there?

Semantics are elusive

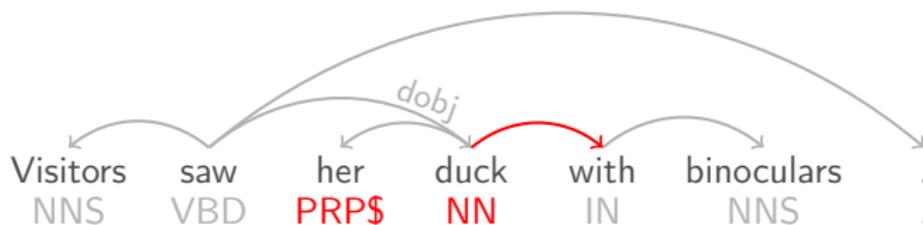


Did she duck or does she *have* a duck?

Who has the binoculars?

How many pairs of binoculars are there?

Semantics are elusive

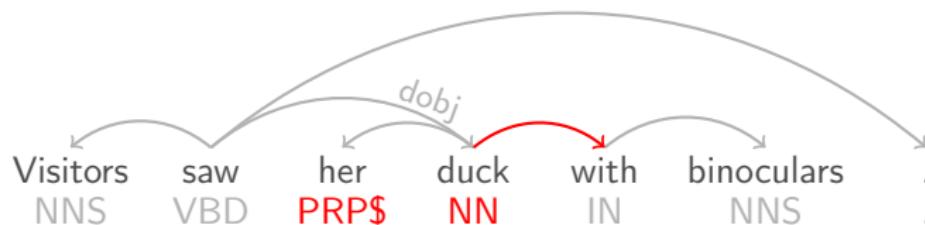


Did she duck or does she *have* a duck?

Who has the binoculars?

How many pairs of binoculars are there?

Semantics are elusive



Did she duck or does she *have* a duck?

Who has the binoculars?

How many pairs of binoculars are there?

Recurrent connections

Output vector



Hidden state

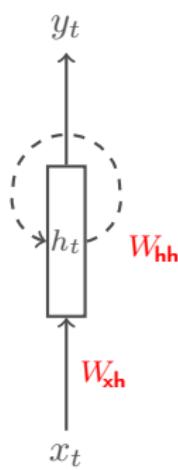
Input vector

Recurrent connections

Output vector

Hidden state

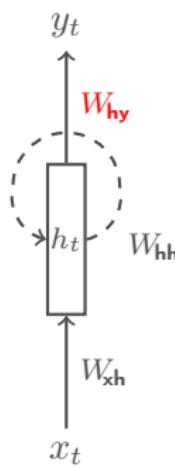
Input vector



$$h_t = \phi_{\text{h}}(W_{\text{xh}} x_t + W_{\text{hh}} h_{t-1})$$

Recurrent connections

Output vector



$$y_t = \phi_y(W_{hy} h_t)$$

Hidden state

$$h_t = \phi_h(W_{xh} x_t + W_{hh} h_{t-1})$$

Input vector

Recurrent connections: Unfolding

x_1

x_2

x_3

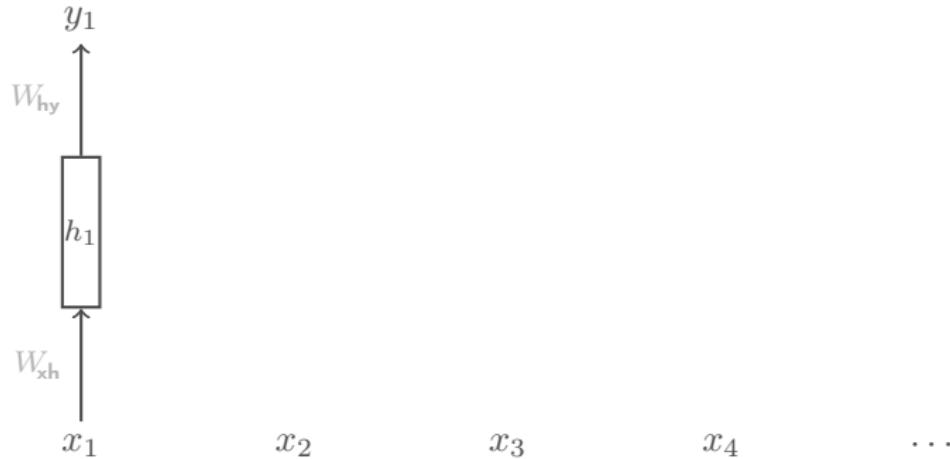
x_4

...

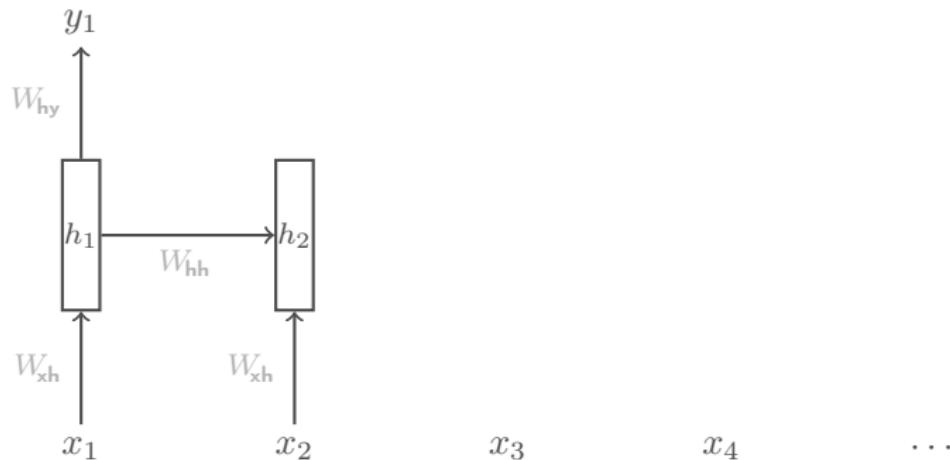
Recurrent connections: Unfolding



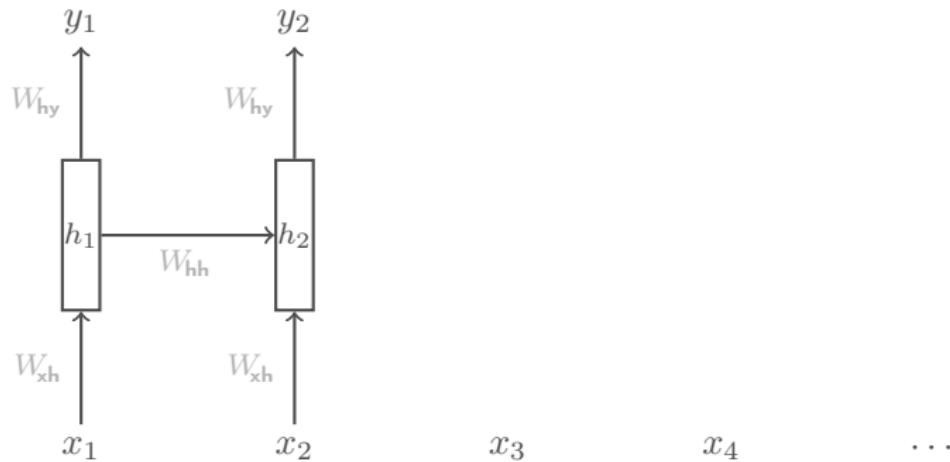
Recurrent connections: Unfolding



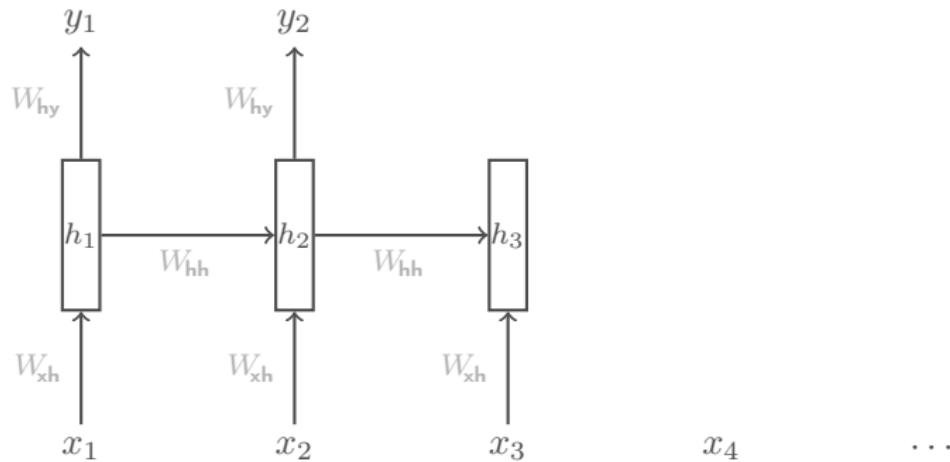
Recurrent connections: Unfolding



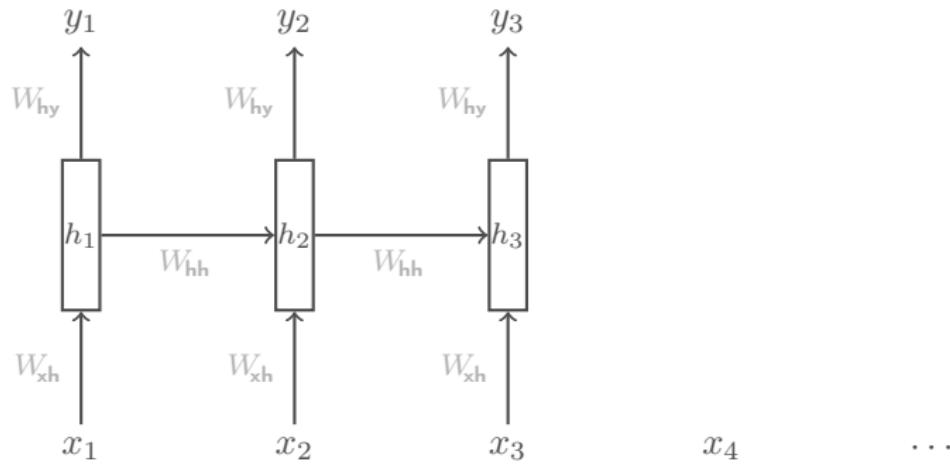
Recurrent connections: Unfolding



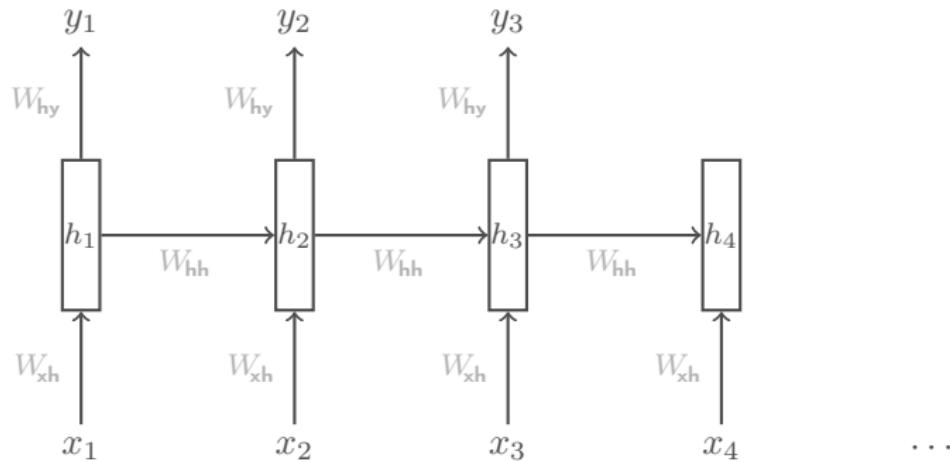
Recurrent connections: Unfolding



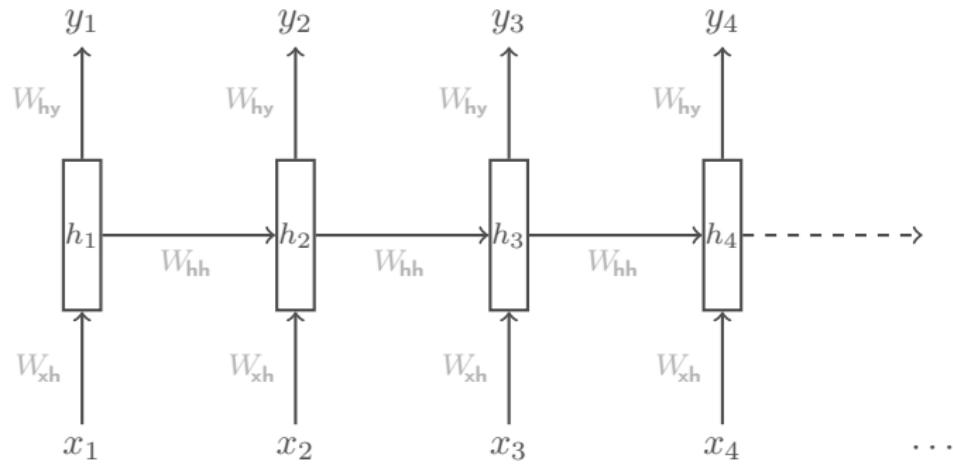
Recurrent connections: Unfolding



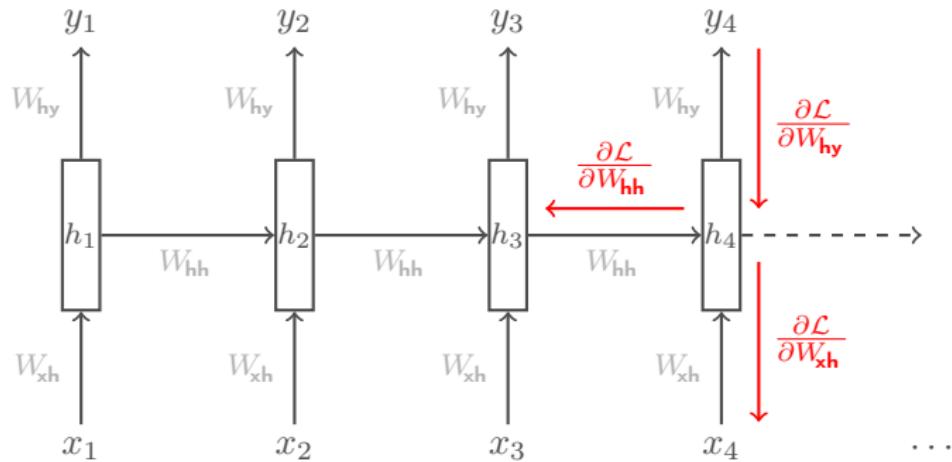
Recurrent connections: Unfolding



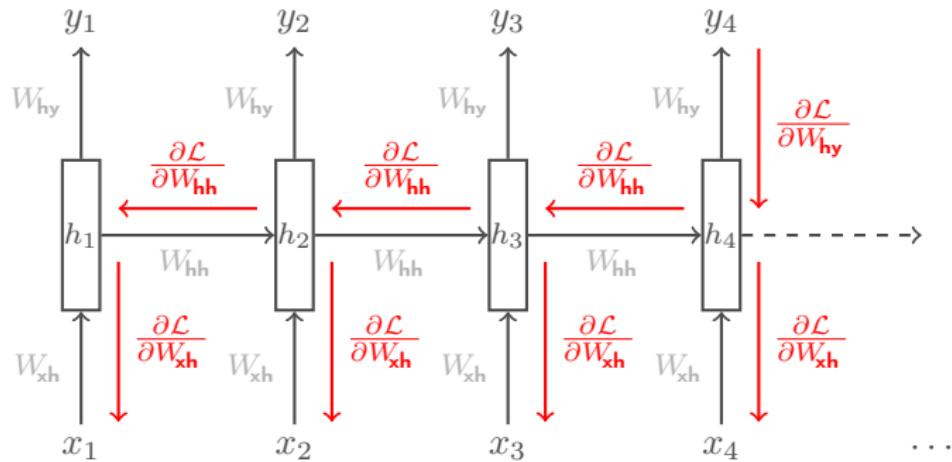
Recurrent connections: Unfolding



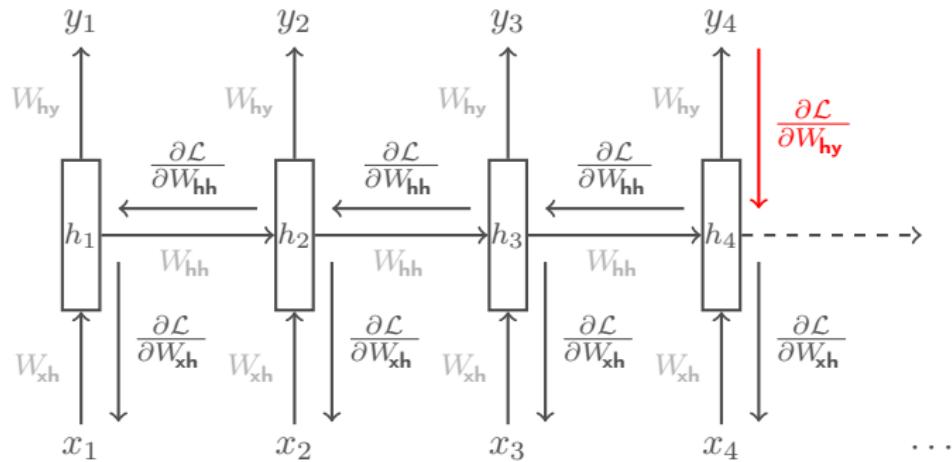
Recurrent connections: Backprop through time



Recurrent connections: Backprop through time

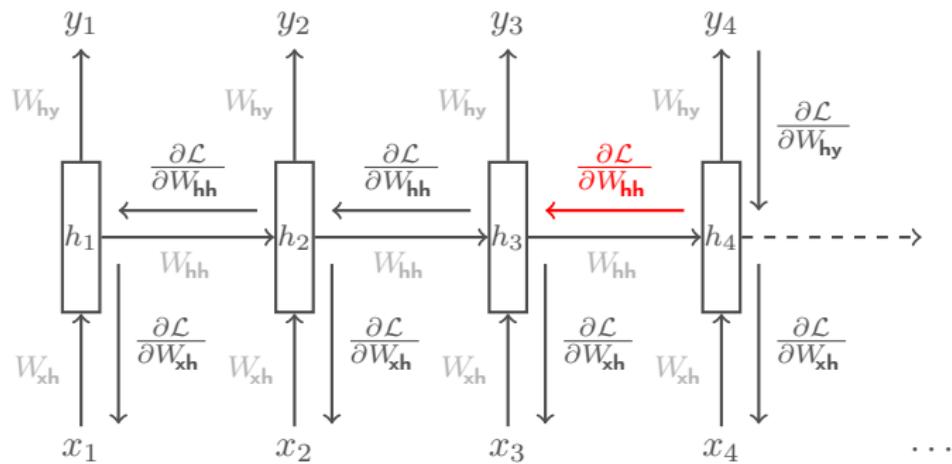


Recurrent connections: Backprop through time



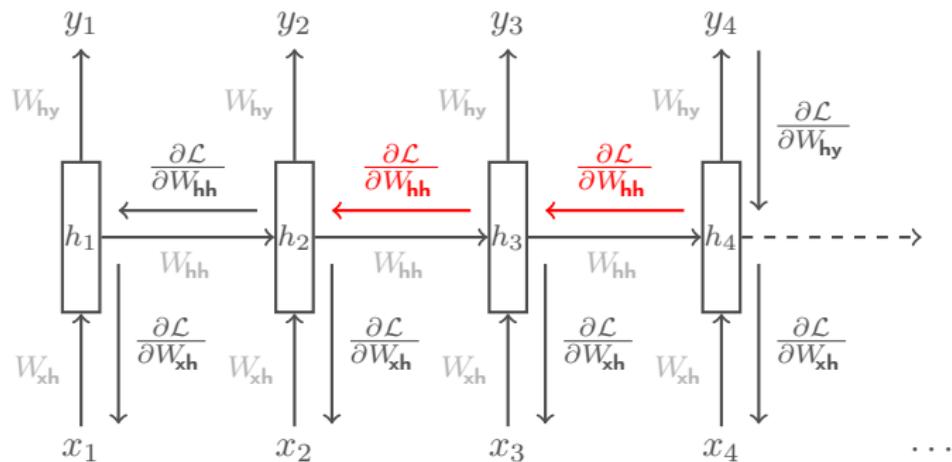
$$\frac{\partial \mathcal{L}}{\partial W_{hy}} = \frac{\partial \mathcal{L}}{\partial y_4} \frac{\partial y_4}{\partial W_{hy}}$$

Recurrent connections: Backprop through time



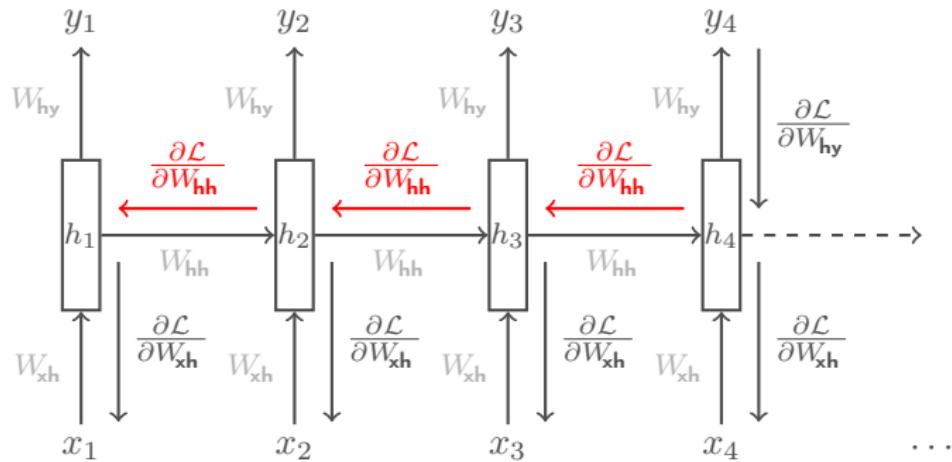
$$\frac{\partial \mathcal{L}}{\partial W_{hh}} = \frac{\partial \mathcal{L}}{\partial y_4} \frac{\partial y_4}{\partial h_4} \frac{\partial h_4}{\partial W_{hh}}$$

Recurrent connections: Backprop through time



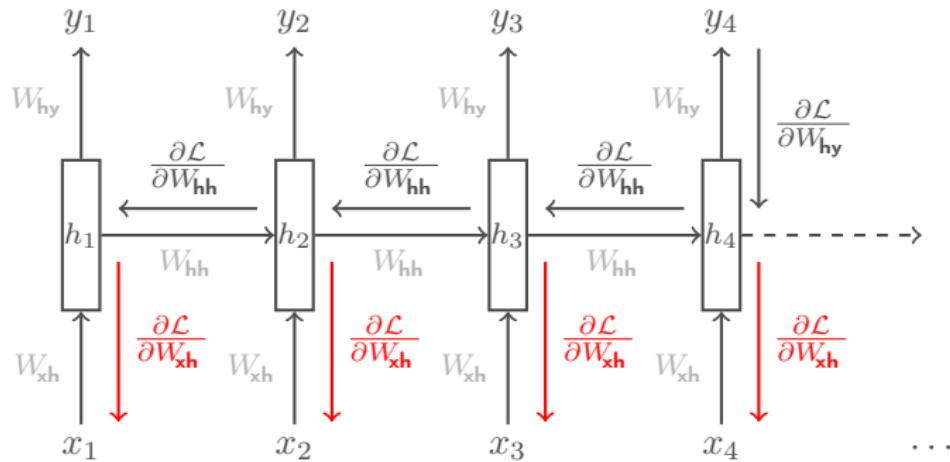
$$\frac{\partial \mathcal{L}}{\partial W_{\text{hh}}} = \frac{\partial \mathcal{L}}{\partial y_4} \frac{\partial y_4}{\partial h_4} \frac{\partial h_4}{\partial W_{\text{hh}}} + \frac{\partial \mathcal{L}}{\partial y_4} \frac{\partial y_4}{\partial h_4} \frac{\partial h_4}{\partial h_3} \frac{\partial h_3}{\partial W_{\text{hh}}}$$

Recurrent connections: Backprop through time



$$\frac{\partial \mathcal{L}}{\partial W_{hh}} = \frac{\partial \mathcal{L}}{\partial y_4} \frac{\partial y_4}{\partial h_4} \frac{\partial h_4}{\partial W_{hh}} + \frac{\partial \mathcal{L}}{\partial y_4} \frac{\partial y_4}{\partial h_4} \frac{\partial h_4}{\partial h_3} \frac{\partial h_3}{\partial W_{hh}} + \frac{\partial \mathcal{L}}{\partial y_4} \frac{\partial y_4}{\partial h_4} \frac{\partial h_4}{\partial h_3} \frac{\partial h_3}{\partial h_2} \frac{\partial h_2}{\partial W_{hh}}$$

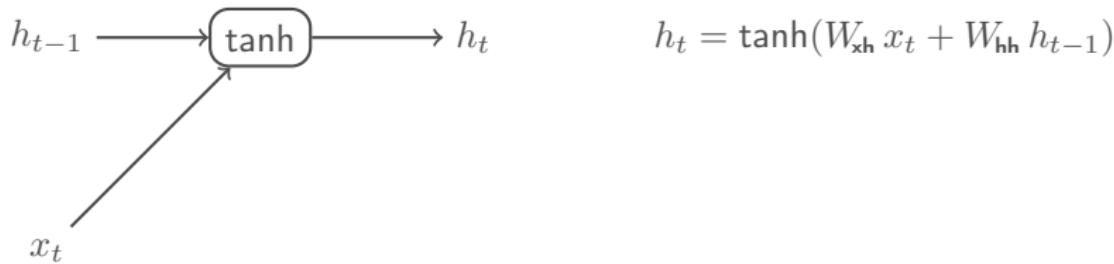
Recurrent connections: Backprop through time



$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial W_{xh}} = & \frac{\partial \mathcal{L}}{\partial y_4} \frac{\partial y_4}{\partial h_4} \frac{\partial h_4}{\partial W_{xh}} + \frac{\partial \mathcal{L}}{\partial y_4} \frac{\partial y_4}{\partial h_4} \frac{\partial h_4}{\partial h_3} \frac{\partial h_3}{\partial W_{xh}} + \frac{\partial \mathcal{L}}{\partial y_4} \frac{\partial y_4}{\partial h_4} \frac{\partial h_4}{\partial h_3} \frac{\partial h_3}{\partial h_2} \frac{\partial h_2}{\partial W_{xh}} + \\ & \frac{\partial \mathcal{L}}{\partial y_4} \frac{\partial y_4}{\partial h_4} \frac{\partial h_4}{\partial h_3} \frac{\partial h_3}{\partial h_2} \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial W_{xh}} \end{aligned}$$

Activation functions

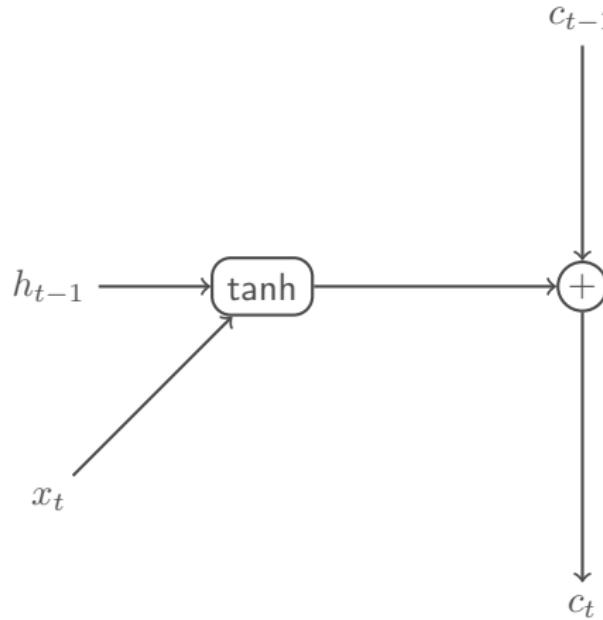
ϕ_h is typically a smooth, bounded function, e.g., σ , \tanh



- Susceptible to vanishing gradients
- Can fail to capture long-term dependencies

Long short-term memory (LSTM)

Hochreiter & Schmidhuber (1997)

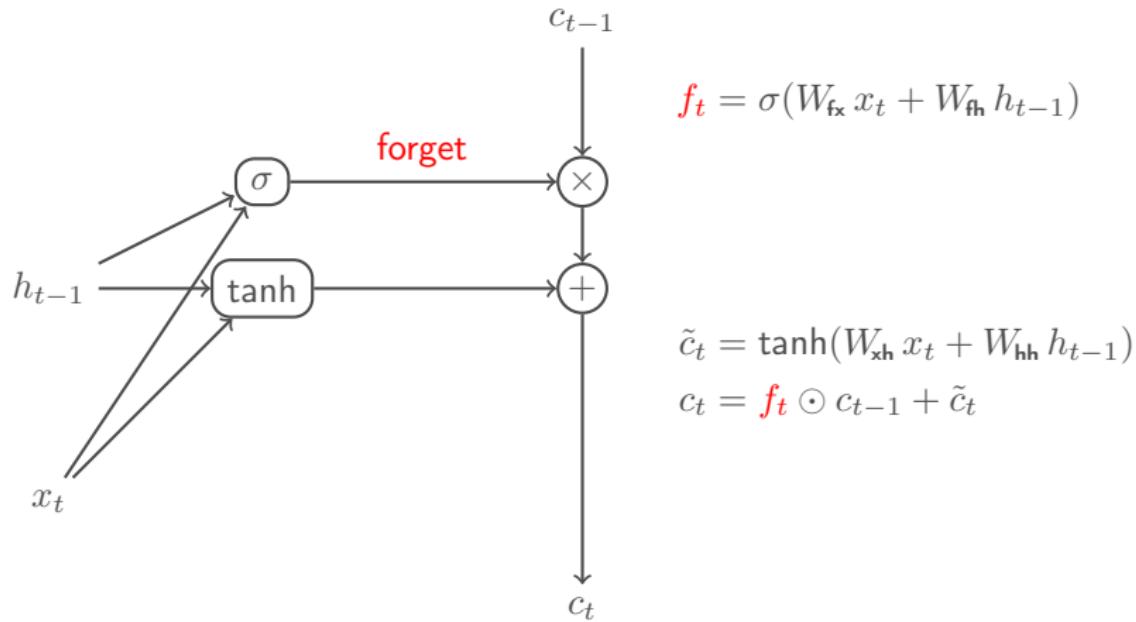


$$\begin{aligned}\tilde{c}_t &= \tanh(W_{\text{xt}} x_t + W_{\text{hh}} h_{t-1}) \\ c_t &= c_{t-1} + \tilde{c}_t\end{aligned}$$

Long short-term memory (LSTM)

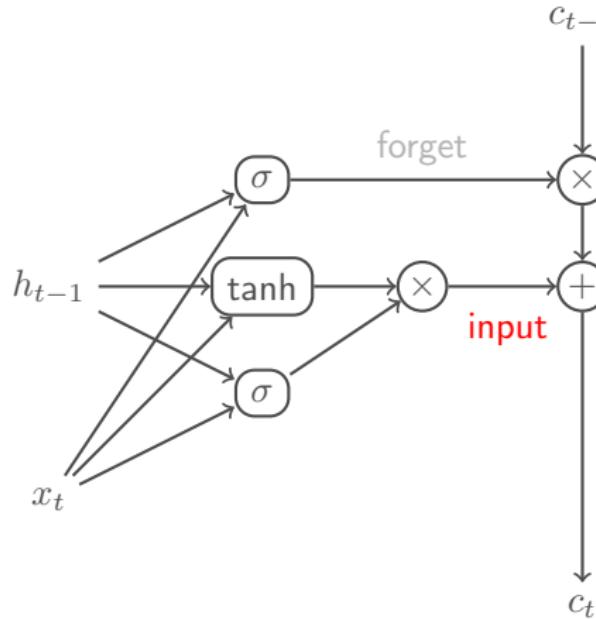
27

Hochreiter & Schmidhuber (1997)



Long short-term memory (LSTM)

Hochreiter & Schmidhuber (1997)



$$f_t = \sigma(W_{fx} x_t + W_{fh} h_{t-1})$$

$$i_t = \sigma(W_{ix} x_t + W_{ih} h_{t-1})$$

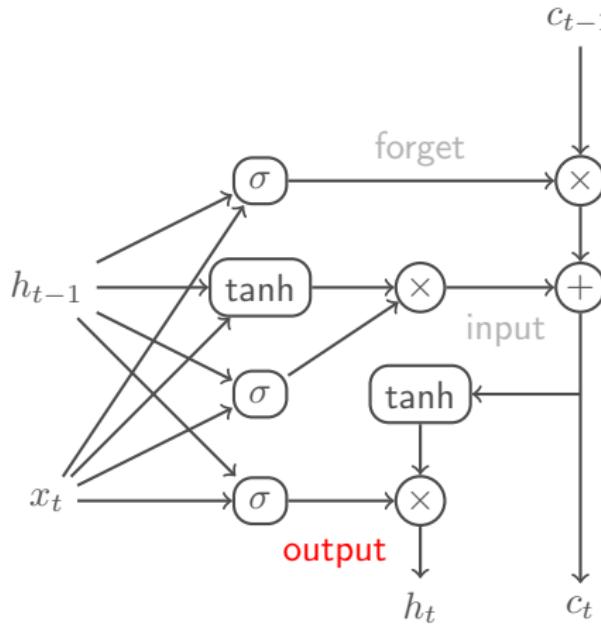
$$\tilde{c}_t = \tanh(W_{xh} x_t + W_{hh} h_{t-1})$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$$

Long short-term memory (LSTM)

27

Hochreiter & Schmidhuber (1997)



$$f_t = \sigma(W_{\text{fx}} x_t + W_{\text{fh}} h_{t-1})$$

$$i_t = \sigma(W_{\text{ix}} x_t + W_{\text{ih}} h_{t-1})$$

$$o_t = \sigma(W_{\text{ex}} x_t + W_{\text{eh}} h_{t-1})$$

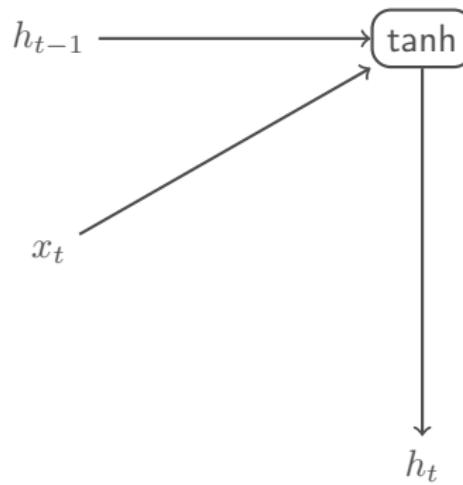
$$\tilde{c}_t = \tanh(W_{\mathbf{x}h} x_t + W_{hh} h_{t-1})$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$$

$$h_t = \textcolor{red}{o}_t \odot \tanh(c_t)$$

Gated Recurrent Unit (GRU)

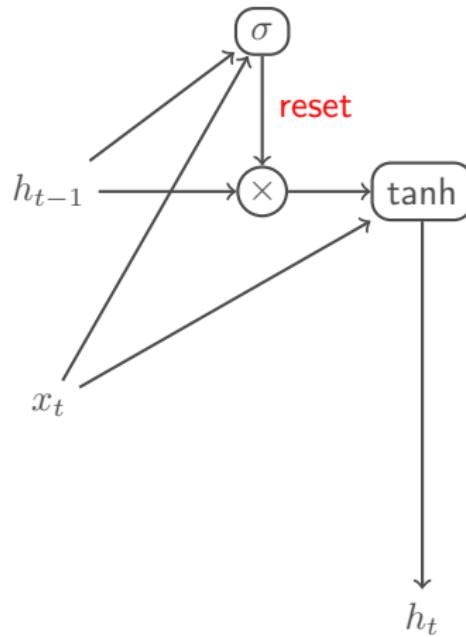
Cho et al. (2014)



$$\begin{aligned}\tilde{h}_t &= \tanh(W_{\text{xh}} x_t + W_{\text{hh}} h_{t-1}) \\ h_t &= \tilde{h}_t\end{aligned}$$

Gated Recurrent Unit (GRU)

Cho et al. (2014)



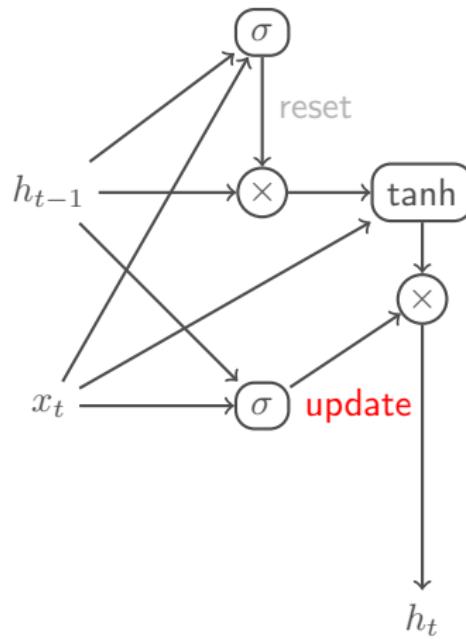
$$r_t = \sigma(W_{rx} x_t + W_{rh} h_{t-1})$$

$$\tilde{h}_t = \tanh(W_{xh} x_t + W_{hh} (r_t \odot h_{t-1}))$$

$$h_t = \tilde{h}_t$$

Gated Recurrent Unit (GRU)

Cho et al. (2014)



$$r_t = \sigma(W_{rx} x_t + W_{rh} h_{t-1})$$

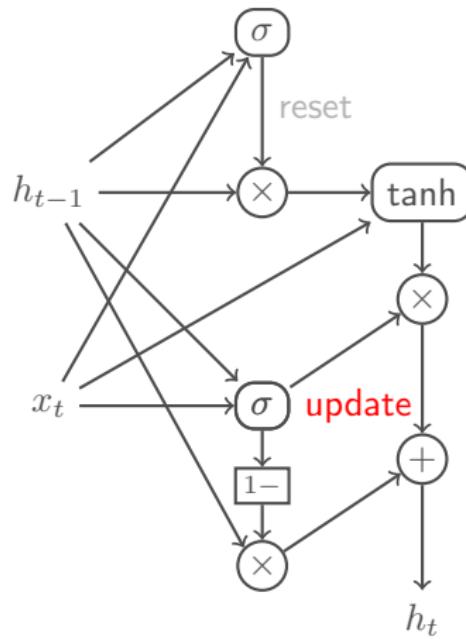
$$z_t = \sigma(W_{zx} x_t + W_{zh} h_{t-1})$$

$$\tilde{h}_t = \tanh(W_{xh} x_t + W_{hh} (r_t \odot h_{t-1}))$$

$$h_t = z_t \odot \tilde{h}_t$$

Gated Recurrent Unit (GRU)

Cho et al. (2014)



$$r_t = \sigma(W_{rx} x_t + W_{rh} h_{t-1})$$

$$z_t = \sigma(W_{zx} x_t + W_{zh} h_{t-1})$$

$$\tilde{h}_t = \tanh(W_{xh} x_t + W_{hh} (r_t \odot h_{t-1}))$$

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t$$

NLP scenarios: Classification

Given variable-length text $w_1 \dots w_n$ (sentence, document, etc), find label y

Normal discriminative approach:

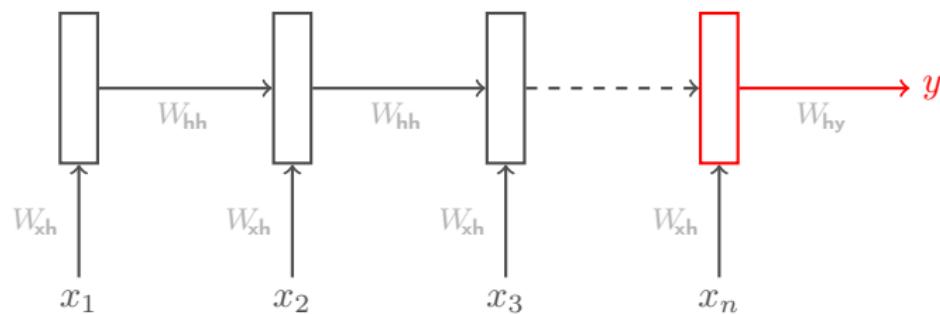
- Extract features over the input text
- Train a linear classifier

Examples:

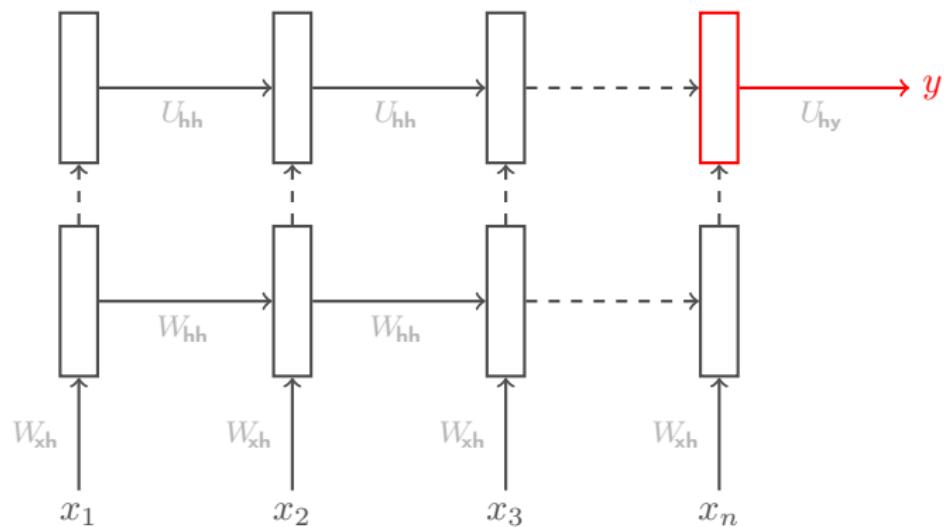
- Topic classification
- Sentiment analysis
- Entailment recognition

⋮

Classification with RNNs

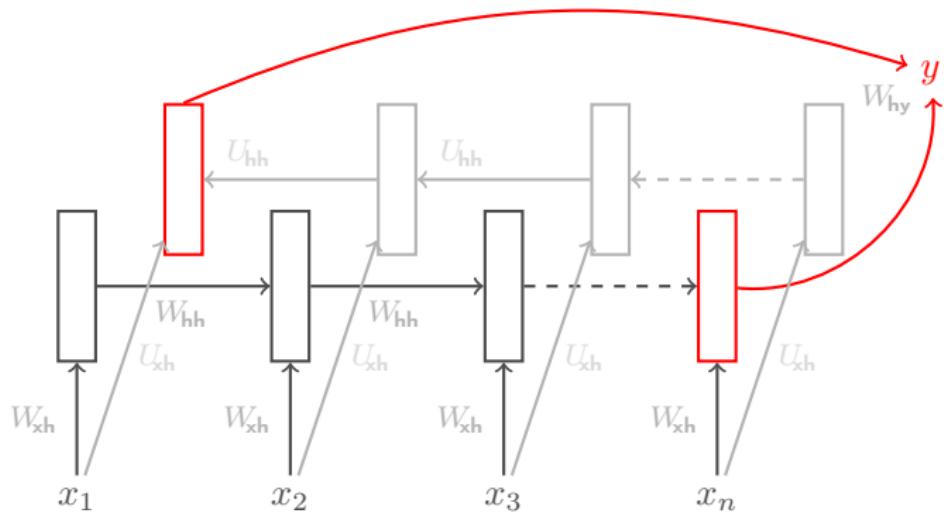


Classification with **deep** RNNs



- + Can learn more abstract representations
 - Slow computation because of recurrent dependencies

Classification with bidirectional RNNs

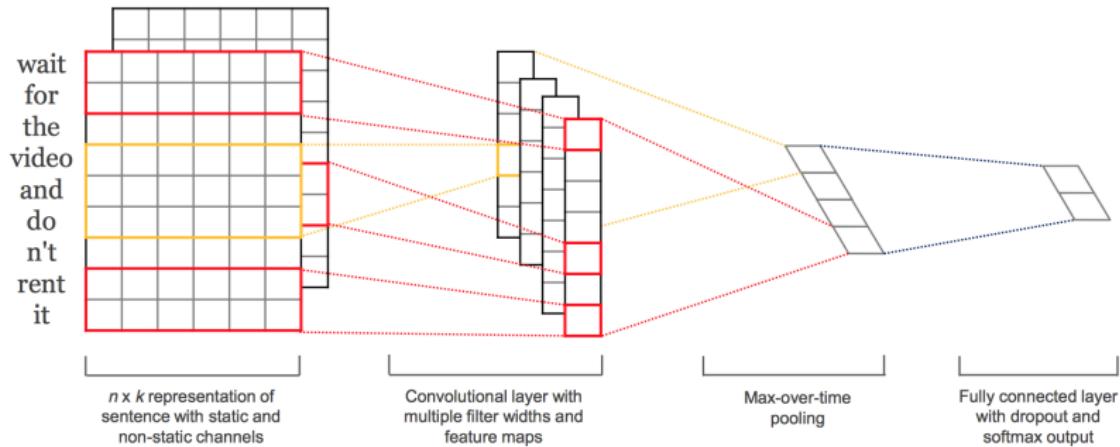


- + Less sensitive to vanishing gradients for long sequences

Classification with CNNs

Kim (2014)

Convolutional Neural Networks for Sentence Classification

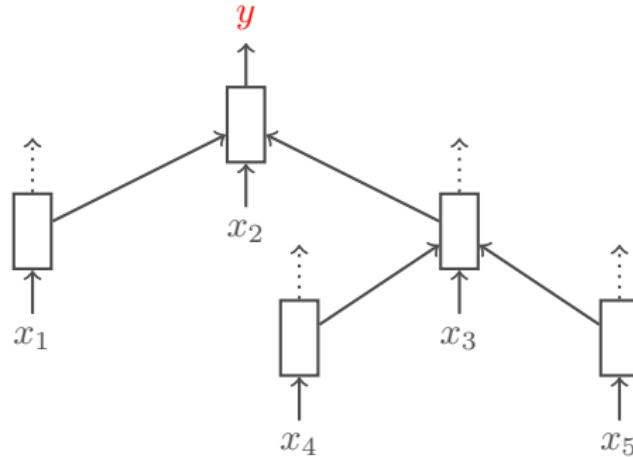


- Max-over-time pooling
 - Two input embedding “channels” — one updated during training

Classification with Tree-RNNs

Tai et al (2015)

Improved Semantic Representations From Tree-Structured Long Short-Term Memory Networks

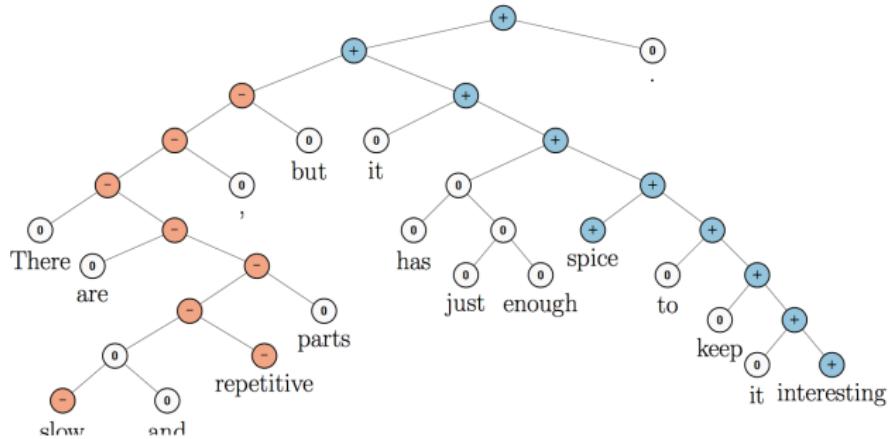


- Computation graph follows dependency or constituent parse
 - Child-sum:
 - Good for arbitrary fan-out or unordered children
 - Suited to dependency trees (input x_i is head word)
 - N -ary:
 - Fixed number of children, each parameterized separately
 - Suited to binarized constituency parses (leaves take word inputs x_i)

Classification with Tree-RNNs

Socher et al (2013)

Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank



- Computation graph follows dependency or constituent parse
 - Child-sum:
 - Good for arbitrary fan-out or unordered children
 - Suited to dependency trees (input x_i is head word)
 - N -ary:
 - Fixed number of children, each parameterized separately
 - Suited to binarized constituency parses (leaves take word inputs x_i)

NLP scenarios: Tagging

Given variable-length text $w_1 \dots w_n$, label spans z_1, \dots, z_m

Normal discriminative approach:

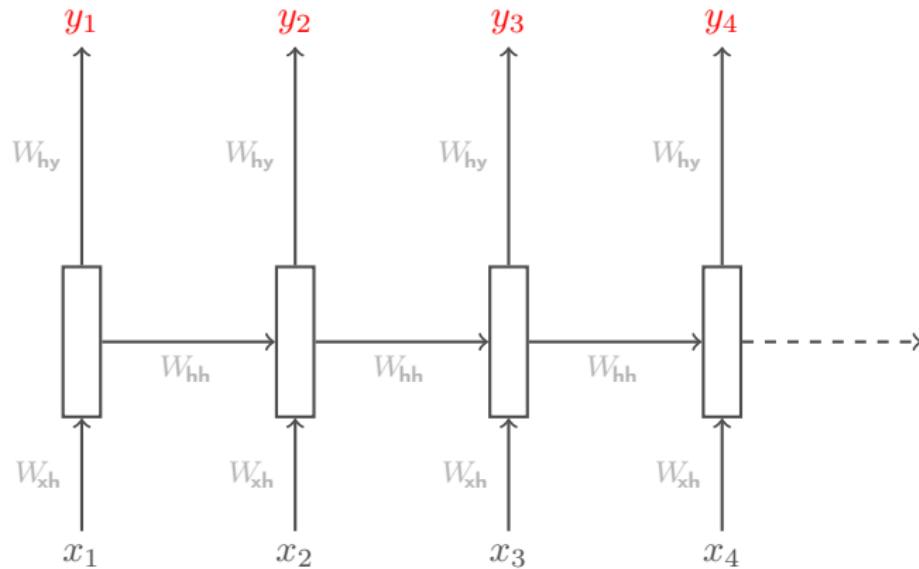
- Distribute labels over input words (e.g., BIO, BILOU encodings) to produce per-word labels y_1, \dots, y_n
- Extract features over input words
- Train a linear-chain conditional random field

Examples:

- Part-of-speech tagging
- Chunking
- Named entity recognition
- Semantic role labeling

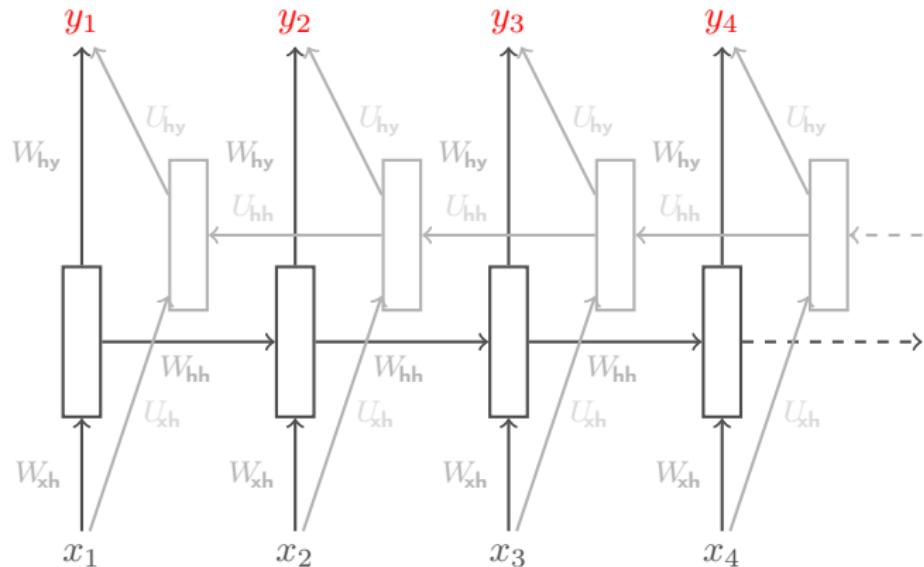
⋮

Tagging with RNNs



- + Effective with BIO/BILOU label encodings

Tagging with bidirectional RNNs



- + Effective with BIO/BILOU label encodings
- + Less sensitive to vanishing gradients for long sequences

NLP scenarios: Generation

Probabilistic language modeling

- Distribution over sequences of words $p(w_1, \dots, w_T)$ in a language
- Typically made tractable via conditional independence assumptions

$$p(w_1, \dots, w_n) = \prod_{t=1}^T p(w_t | w_{t-1}, \dots, w_{t-n})$$

- n -gram counts estimated from large corpora
- Distributions *smoothed* to tolerate data sparsity, e.g., Laplace (add-one) smoothing, Kneser-Ney smoothing
- Evaluate on *perplexity* over held-out data

$$2^{\frac{1}{N} \sum_{i=1}^N \log_2 p(w_1^{(i)} \dots w_{T_i}^{(i)})}$$

NLP scenarios: Generation

Bengio et al (2003)

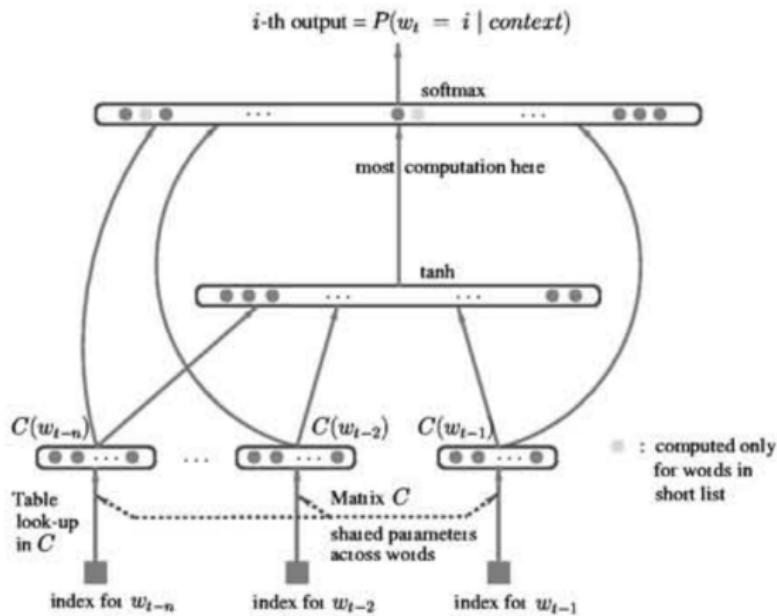
A Neural Probabilistic Language Model

Discriminative language modeling

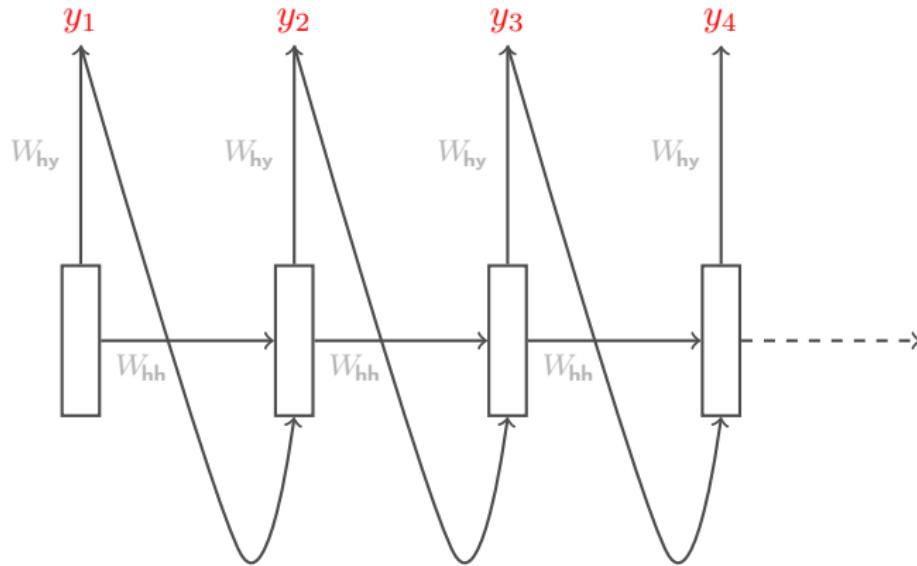
- Estimate n-gram probabilities with a discriminative model

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

e.g., model $p(w_t | w_{t-1}, \dots, w_{t-n})$ with a feed-forward neural net



Language modeling with auto-regressive RNNs



- Supply pre-softmax activations $\phi_y(W_{hy} h_t)$ as input to timestep $t + 1$
Optionally include y_t , e.g., if using beam search
 - *Curriculum learning* to overcome model initialization and speed convergence

Model $p(w_t | w_{t-1}, \dots, w_{t-n})$ with an RNN

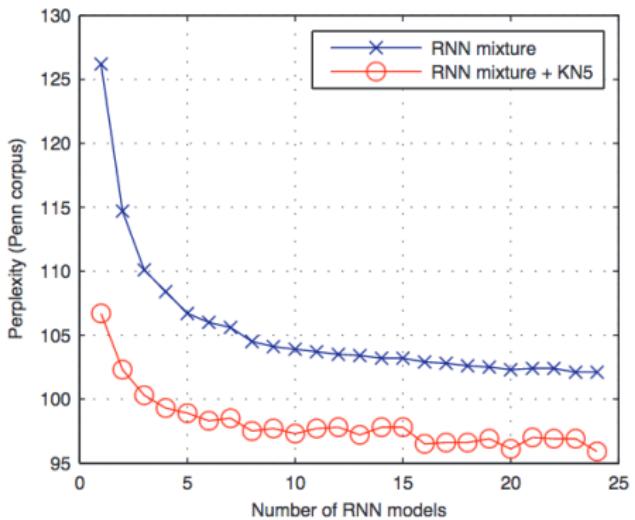
Or an ensemble of multiple RNNs, randomly initialized

Model	Perplexity
Kneser-Ney 5-gram	141
Random forest [Xu 2005]	132
Structured LM [Filimonov 2009]	125
Feedforward NN LM	116
Syntactic NN LM [Emami 2004]	110
RNN trained by BP	113
RNN trained by BPTT	106
4x RNN trained by BPTT	98

Results on Penn Treebank corpus

Model $p(w_t | w_{t-1}, \dots, w_{t-n})$ with an RNN

Or an ensemble of multiple RNNs, randomly initialized



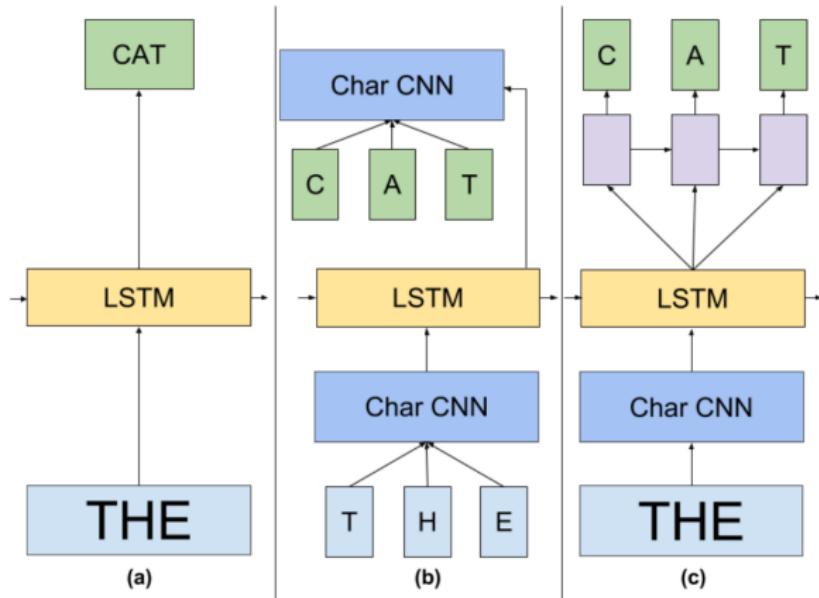
Comparison of single RNN vs RNN ensembles

RNNLMs with character CNNs

Jozefowicz et al (2015)

39

Exploring the Limits of Language Modeling



Recent models with character-CNN inputs and softmax alternatives

RNNLMs with character CNNs

Jozefowicz et al (2015)

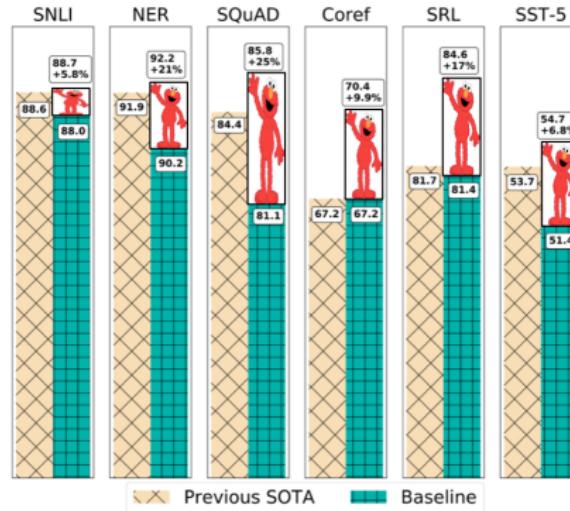
Exploring the Limits of Language Modeling

WORD	TOP-1	TOP-2	TOP-3
INCREDIBLE	INCREDIBLE	NONEDEBLE	EXTENDIBLE
WWW.A.COM	WWW.AA.COM	WWW.AAA.COM	WWW.CA.COM
7546	7646	7534	8566
TOWNHAL1	TOWNHALL	DJC2	MOODSWING360
KOMARSKI	KOHARSKI	KONARSKI	KOMANSKI

Nearest neighbors in character-CNN embedding space

Transfer learning with NNLMs

Pre-trained neural language models are useful in classification tasks



Performance gains with ELMo representations added to various models

Deep Contextualized Word Representations (Peters et al., 2018)

Universal Language Model Fine-tuning for Text Classification (Howard and Ruder, 2018)

Improving Language Understanding by Generative Pre-Training (Radford et al., 2018)

NLP scenarios: Text-to-text

Given variable-length text $x_1 \dots x_n$ (sentence, document, etc), produce output text y_1, \dots, y_m under some transformation

Traditional approaches:

- Pipelined models
- Constrained optimization

⋮

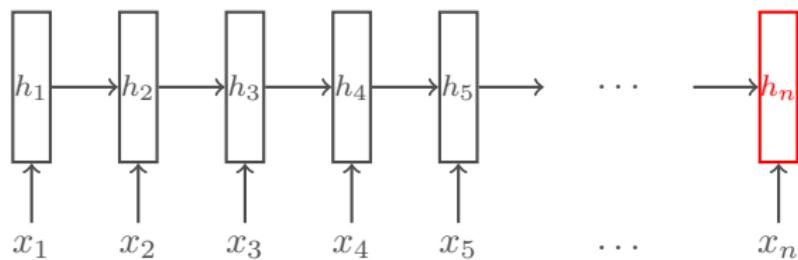
Examples:

- Machine translation
- Document summarization
- Sentence simplification
- Paraphrase generation

⋮

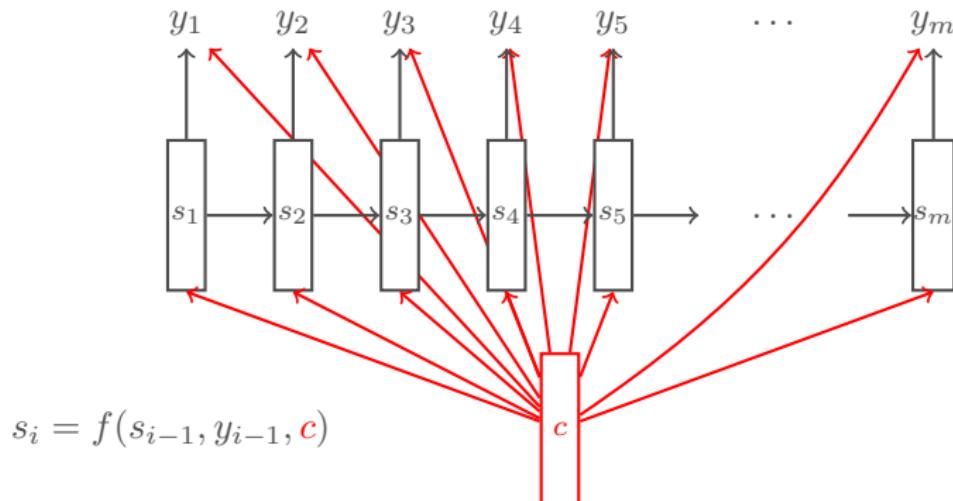
Encoding with RNNs

- Input: source words x_1, \dots, x_n
- Output: representation h_n



Decoding with RNNs

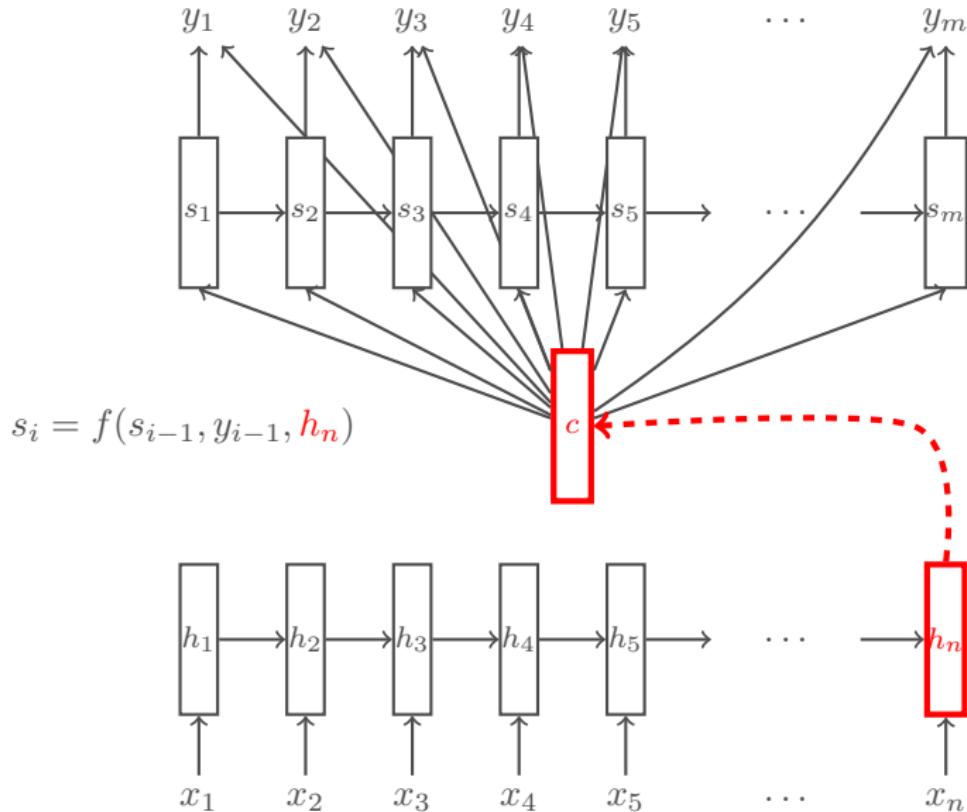
- Input: context vector c
- Output: words y_1, \dots, y_m



Sequence-to-sequence learning

Sutskever, Vinyals & Le (2014)

Sequence to Sequence Learning with Neural Networks



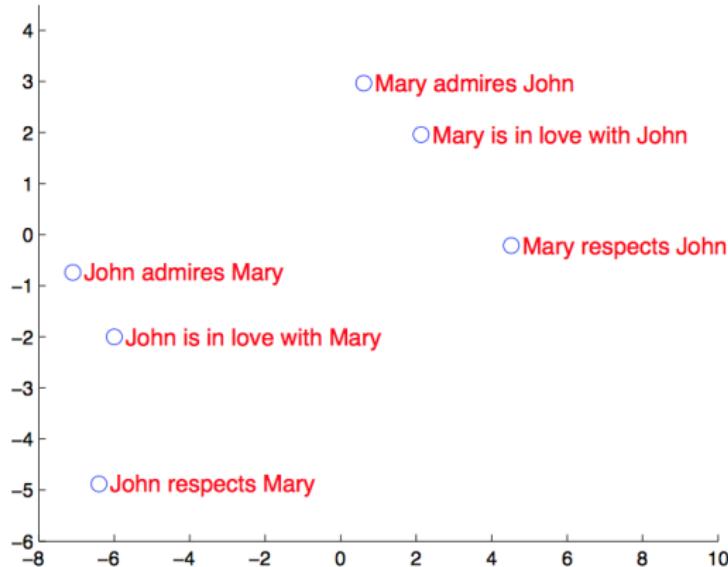
Sequence-to-sequence learning

Sutskever, Vinyals & Le (2014)

Sequence to Sequence Learning with Neural Networks

Produces a fixed length representation of input

- “sentence embedding” or “thought vector” ($\vec{\Theta}_\text{S}$)



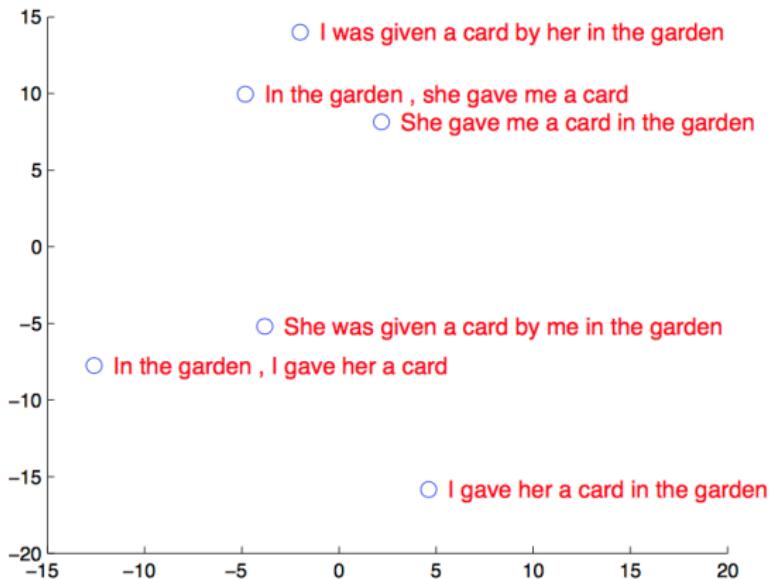
Sequence-to-sequence learning

Sutskever, Vinyals & Le (2014)

Sequence to Sequence Learning with Neural Networks

Produces a fixed length representation of input

- “sentence embedding” or “thought vector” ($\vec{\Theta}_\text{S}$)



Overview: processing text with RNNs

Input

- Word/sentence embeddings
- One-hot words/characters
- CNNs over characters/words/sentences, e.g., document modeling
- Absent, e.g., RNN-LMs

⋮

Recurrent layer

- Gated units: LSTMs, GRUs
- Forward, backward, bidirectional
- ReLUs initialized with identity matrix

⋮

Output

- Softmax over words/characters/labels, e.g., text generation
- Deeper RNN layers
- Absent, e.g., text encoders

⋮