



Language
Technologies
Institute

Carnegie
Mellon
University

Advanced Multimodal Machine Learning

Lecture 4.1: Recurrent Networks

Louis-Philippe Morency

* Original version co-developed with Tadas Baltrusaitis

Administrative Stuff



Upcoming Schedule

- First project assignment:
 - Proposal presentation (10/3 and 10/5)
 - First project report (Sunday 10/8)
- Second project assignment
 - Midterm presentations (11/7 and 11/9)
 - Midterm report (Sunday 11/12)
- Final project assignment
 - Final presentation (12/4 & 12/5)
 - Final report (Sunday 12/10)



Proposal Presentation (10/3/2017 and 10/5/2017)

- 5 minutes (about 5-10 slides)
- All team members should be involved in the presentation
- Will receive feedback from instructors and other students
 - 1-2 minutes between presentations reserved for written feedback
- Main presentation points (similar to pre-proposal)
 - General research problem and motivation
 - Dataset and input modalities
 - Multimodal challenges and prior work
- You need to submit a copy of your slides (PDF or PPT)
 - Deadline: Friday 10/6 (on Gradescope)

Project Proposal Report – Due on 10/8/17

- Part 1 (updated version of your pre-proposal)
 - **Research problem:**
 - Describe and motivate the research problem
 - Define in generic terms the main computational challenges
 - **Dataset and Input Modalities:**
 - Describe the dataset(s) you are planning to use for this project.
 - Describe the input modalities and annotations available in this dataset.

Project Proposal Report – Due on 10/8/17

- Part 2
 - **Related Work:**
 - Include 12-15 paper citations which give an overview of the prior work
 - Present in more details the 3-4 research papers most related to your work
 - **Research Challenges and Hypotheses:**
 - Describe your specific challenges and/or research hypotheses
 - Highlight the novel aspect of your proposed research

Project Proposal Report – Due on 10/8/17

- Part 3
 - **Language Modality Exploration:**
 - Explore neural language models on your dataset (using Keras)
 - Train at least two different language models (e.g., using SimpleRNN, GRU or LSTM) on your dataset and compare their perplexity.
 - Include qualitative examples of successes and failure cases.
 - **Visual Modality Exploration:**
 - Explore pre-trained Convolutional Neural Networks (CNNs) on your dataset
 - Load a pre-existing CNN model trained for object recognition (e.g., AlexNet or VGG-Net) and process your test images.
 - Extract features at different network layers in the network and visualize them (using t-sne visualization) with overlaid class labels with different colors.



Lecture Objectives

- Word representations & distributional hypothesis
 - Learning neural representations (e.g., Word2vec)
- Language models
- Sequence modeling tasks
- Recurrent neural networks
- Backpropagation through time
- Gates recurrent neural networks
 - Long Short-Term Memory (LSTM) model

Distributed Semantics



Possible ways of representing words

Given a text corpus containing 100,000 unique words

- ➡ Classic binary word representation: $[0; 0; 0; 0; \dots; 0; 0; 1; 0; \dots; 0; 0]$
100,000d vector
 - ➡ Only non-zero at the index of the word
- ➡ Classic word feature representation: $[5; 1; 0; 0; \dots; 0; 20; 1; 0; \dots; 3; 0]$
300d vector
 - ➡ Manually define 300 “good” features (e.g., ends on –ing)
- ➡ Learned word representation: $[0,1; 0,0003; 0; \dots; 0,02; 0.08; 0,05]$
300d vector
 - ➡ This 300-dimension vector should approximate the “meaning” of the word



The Distributional Hypothesis

- Distribution Hypothesis (DH) [Lenci 2008]
 - At least certain aspects of the meaning of lexical expressions depend on their distributional properties in the linguistic contexts
 - The degree of semantic similarity between two linguistic expressions α and β is a function of the similarity of the linguistic contexts in which α and β can appear
- Weak and strong DH
 - Weak view as a quantitative method for semantic analysis and lexical resource induction
 - Strong view as a cognitive hypothesis about the form and origin of semantic representations; assuming that word distributions in context play a specific *causal role* in forming meaning representations.




What is the meaning of “bardiwac”?

- He handed her glass of **bardiwac**.
 - Beef dishes are made to complement the **bardiwacs**.
 - Nigel staggered to his feet, face flushed from too much **bardiwac**.
 - Malbec, one of the lesser-known **bardiwac** grapes, responds well to Australia’s sunshine.
 - I dined off bread and cheese and this excellent **bardiwac**.
 - The drinks were delicious: blood-red **bardiwac** as well as light, sweet Rhenish.
- ⇒ **bardiwac** is a heavy red alcoholic beverage made from grapes

Geometric interpretation

- row vector \mathbf{x}_{dog} describes usage of word *dog* in the corpus
- can be seen as coordinates of point in n -dimensional Euclidean space \mathbb{R}^n

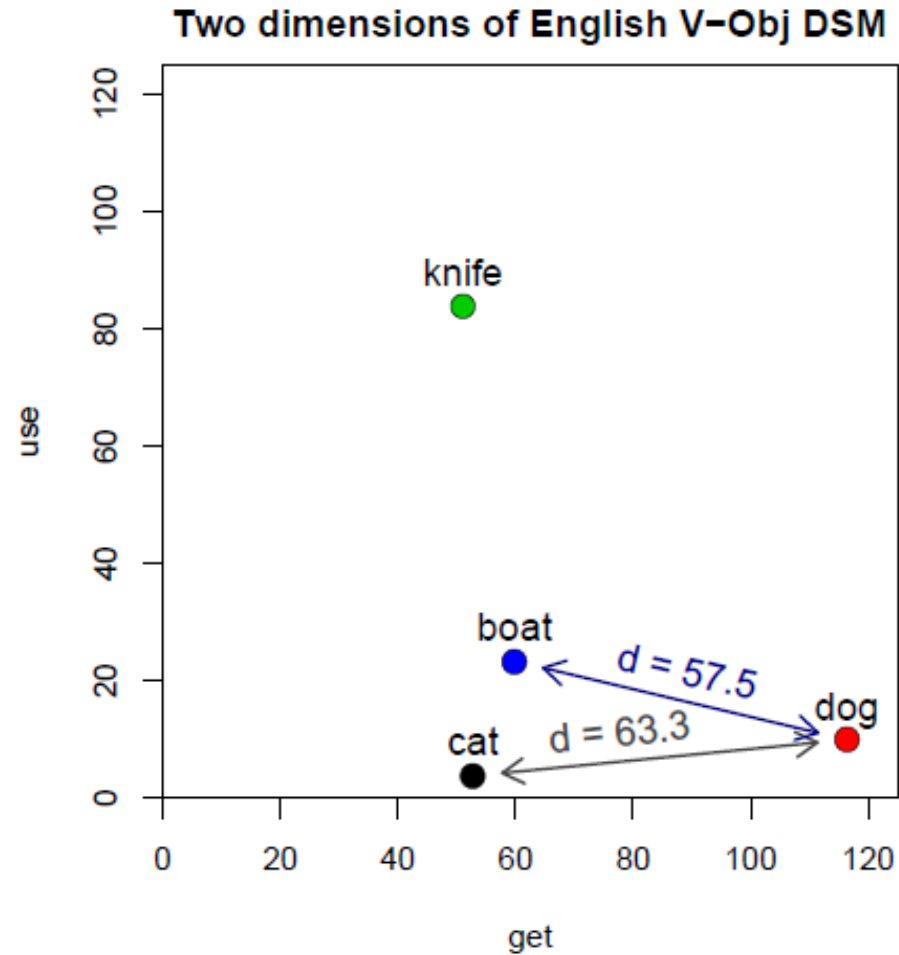


	get	see	use	hear	eat	kill
knife	51	20	84	0	3	0
cat	52	58	4	4	6	26
dog	115	83	10	42	33	17
boat	59	39	23	4	0	0
cup	98	14	6	2	1	0
pig	12	17	3	2	9	27
banana	11	2	2	0	18	0

co-occurrence matrix M

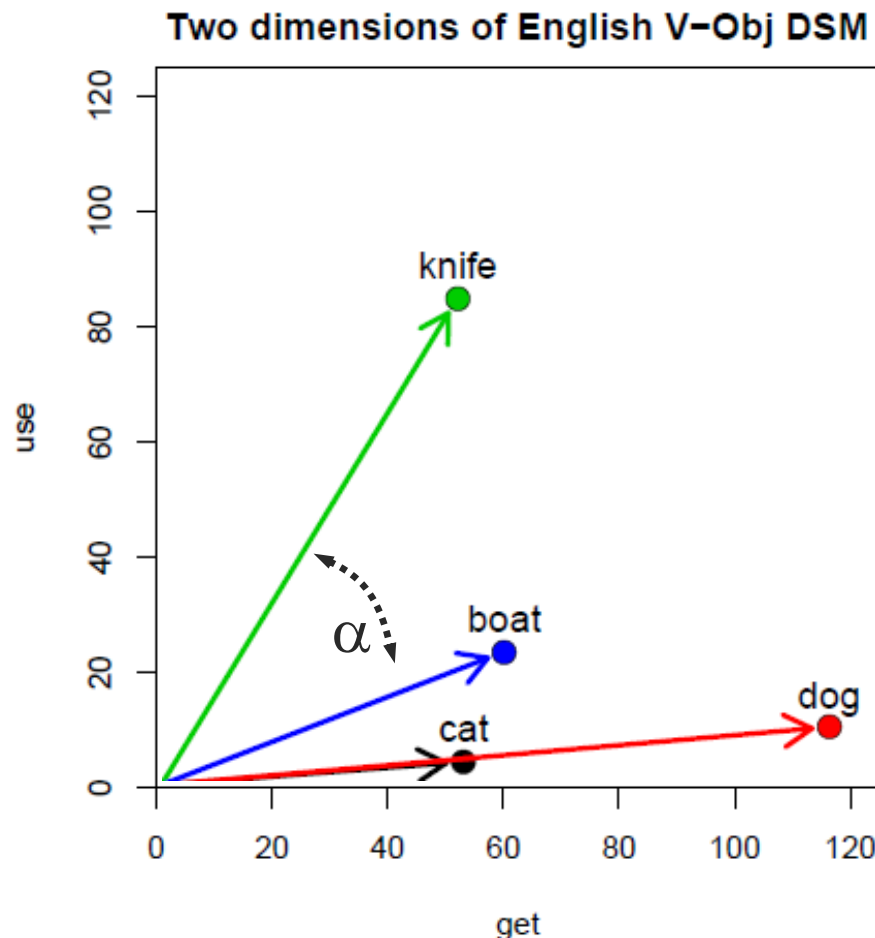
Distance and similarity

- illustrated for two dimensions: *get* and *use*: $\mathbf{x}_{\text{dog}} = (115, 10)$
- similarity = spatial proximity (Euclidean distance)
- location depends on frequency of noun ($f_{\text{dog}} \approx 2.7 \cdot f_{\text{cat}}$)

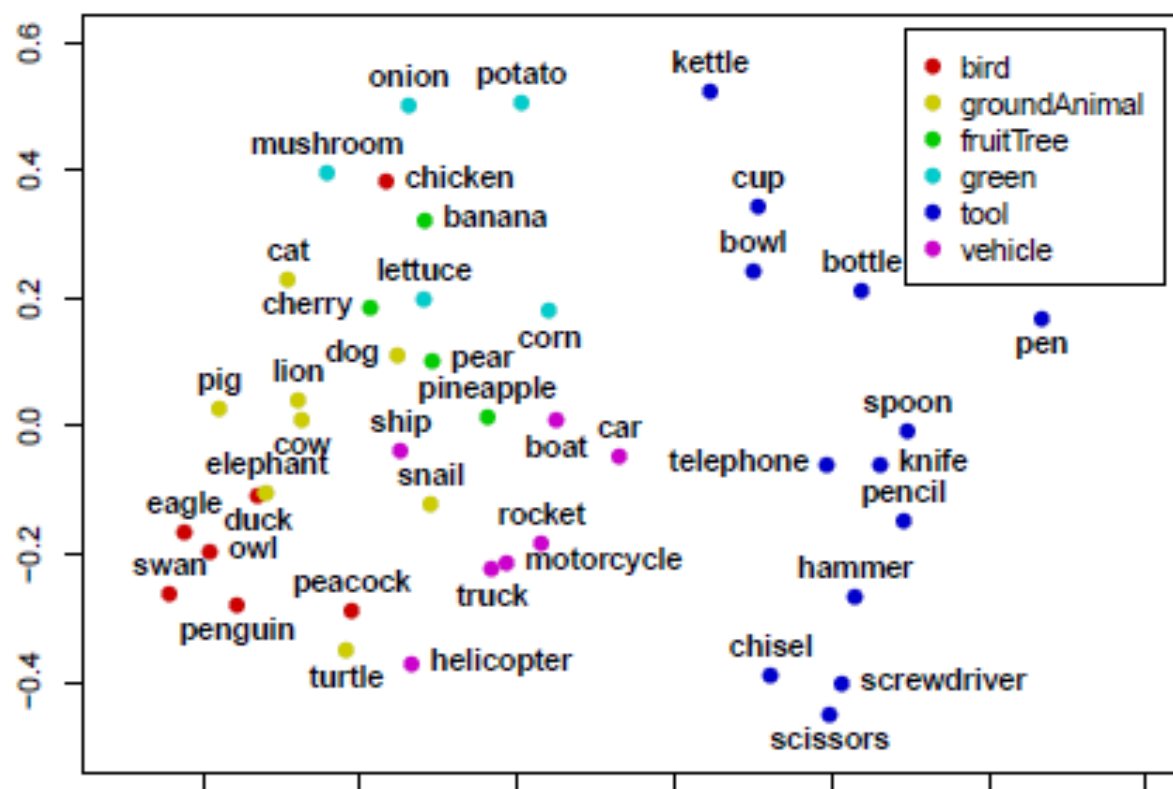


Angle and similarity

- direction more important than location
- normalise “length”
 $\|\mathbf{x}_{\text{dog}}\|$ of vector
- or use angle α as distance measure



Semantic maps



Learning Neural Representations of Words



How to learn (word) features/representations?

➡ **Distribution hypothesis:** Approximate the word meaning by its surrounding words

➡ Words used in a similar context will lie close together

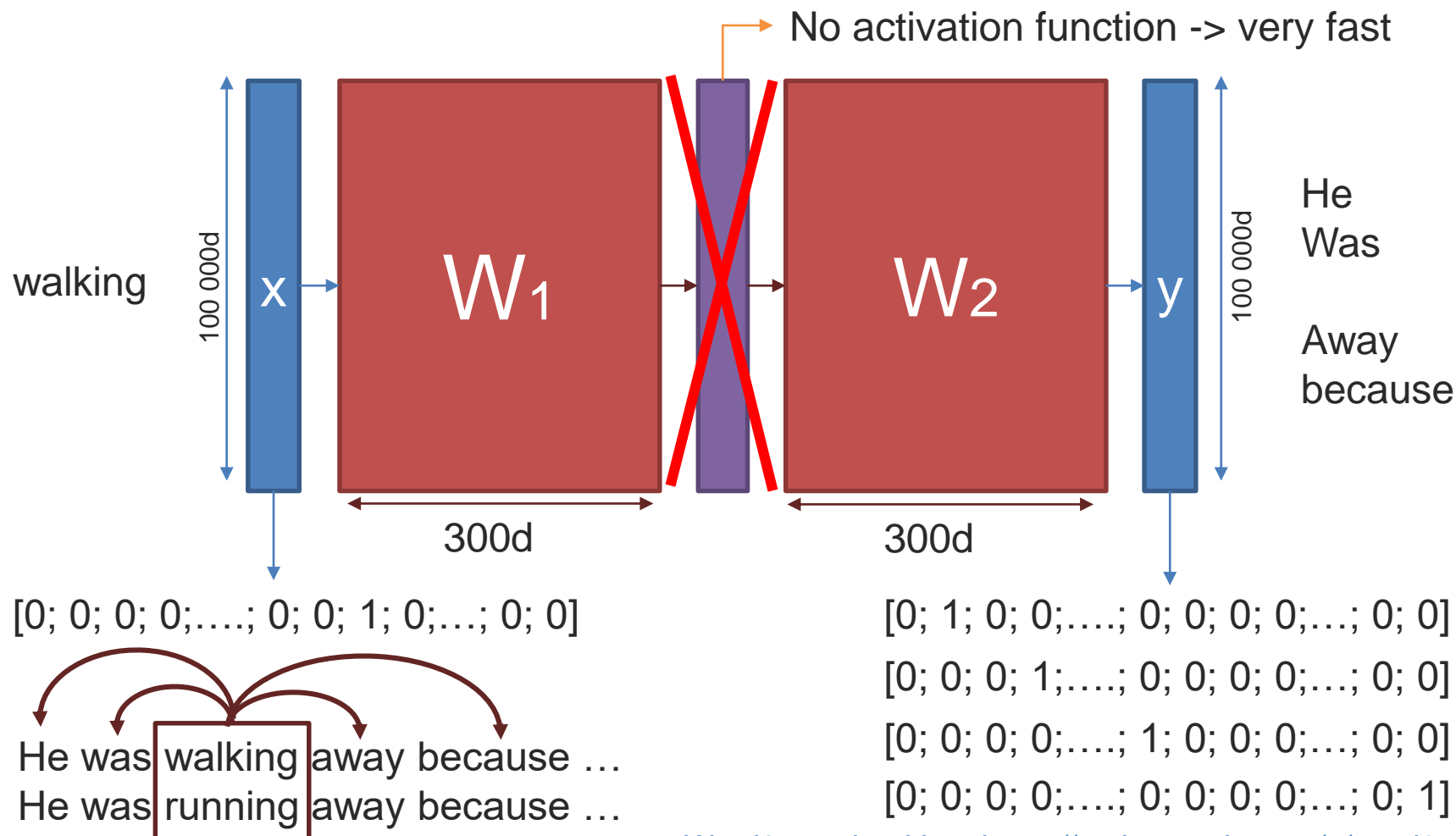


➡ **Instead of capturing co-occurrence counts directly, predict surrounding words of every word**

$$\frac{1}{T} \sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log p(w_{t+j} | w_t)$$



How to learn (word) features/representations?



Word2vec algorithm: <https://code.google.com/p/word2vec/>

How to use these word representations

If we would have a vocabulary of 100 000 words:

Classic NLP: $\xleftarrow{100\,000 \text{ dimensional vector}}$

Walking: $[0; 0; 0; 0; \dots; 0; 0; 1; 0; \dots; 0; 0]$

Running: $[0; 0; 0; 0; \dots; 0; 0; 0; 0; \dots; 1; 0]$

➡ Similarity = 0.0



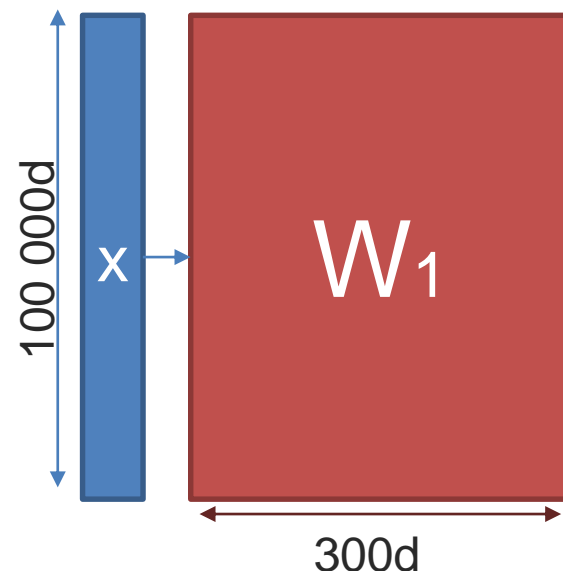
Transform: $x' = x * W$

Goal: $\xleftarrow{300 \text{ dimensional vector}}$

Walking: $[0,1; 0,0003; 0; \dots; 0,02; 0,08; 0,05]$

Running: $[0,1; 0,0004; 0; \dots; 0,01; 0,09; 0,05]$

➡ Similarity = 0.9



Word representations: examples

Word similarity test:

➡ Trained on 400 million tweets having 5 billion words

Input: running	Cosine similarity	Input: :)	Cosine similarity
runnin	0.758099	:))	0.885355
runing	0.702119	=)	0.836011
Running	0.69014	:D	0.818340
runnnng	0.669039	;))	0.814380
sprinting	0.587385	(:	0.809806
runnung	0.578426	:)))	0.808298
run	0.576671	:-)	0.798115
walking/running	0.563114	:))))	0.777765
runin	0.556682	:)	0.772422
walking	0.542137	:-))	0.758584



Vector space models of words

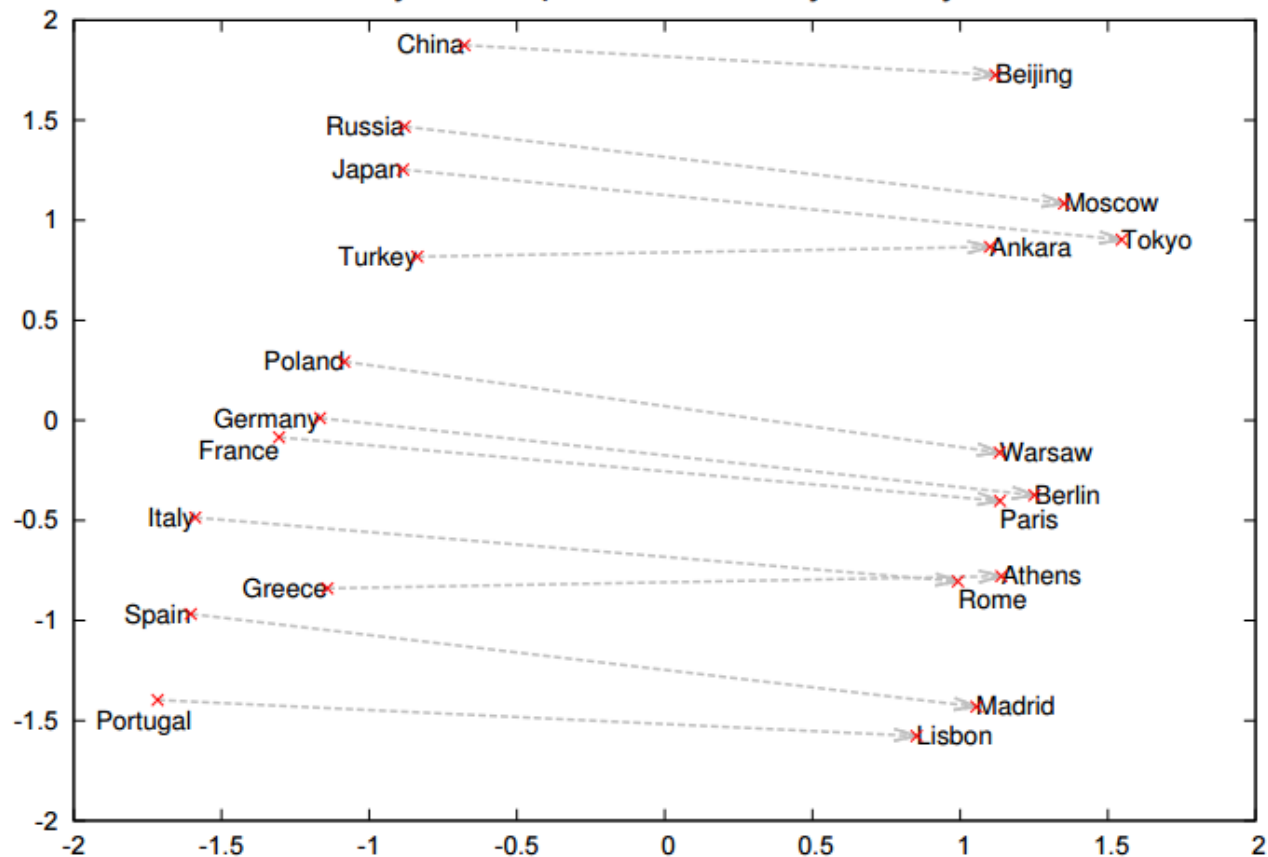
- ➔ While learning these word representations, we are actually building a vector space in which all words reside with certain relationships between them
- ➔ Encodes both syntactic and semantic relationships
- ➔ This vector space allows for algebraic operations:

$$\text{Vec}(\text{king}) - \text{vec}(\text{man}) + \text{vec}(\text{woman}) \approx \text{vec}(\text{queen})$$

Why linear algebra is working?



Vector space models of words: semantic relationships

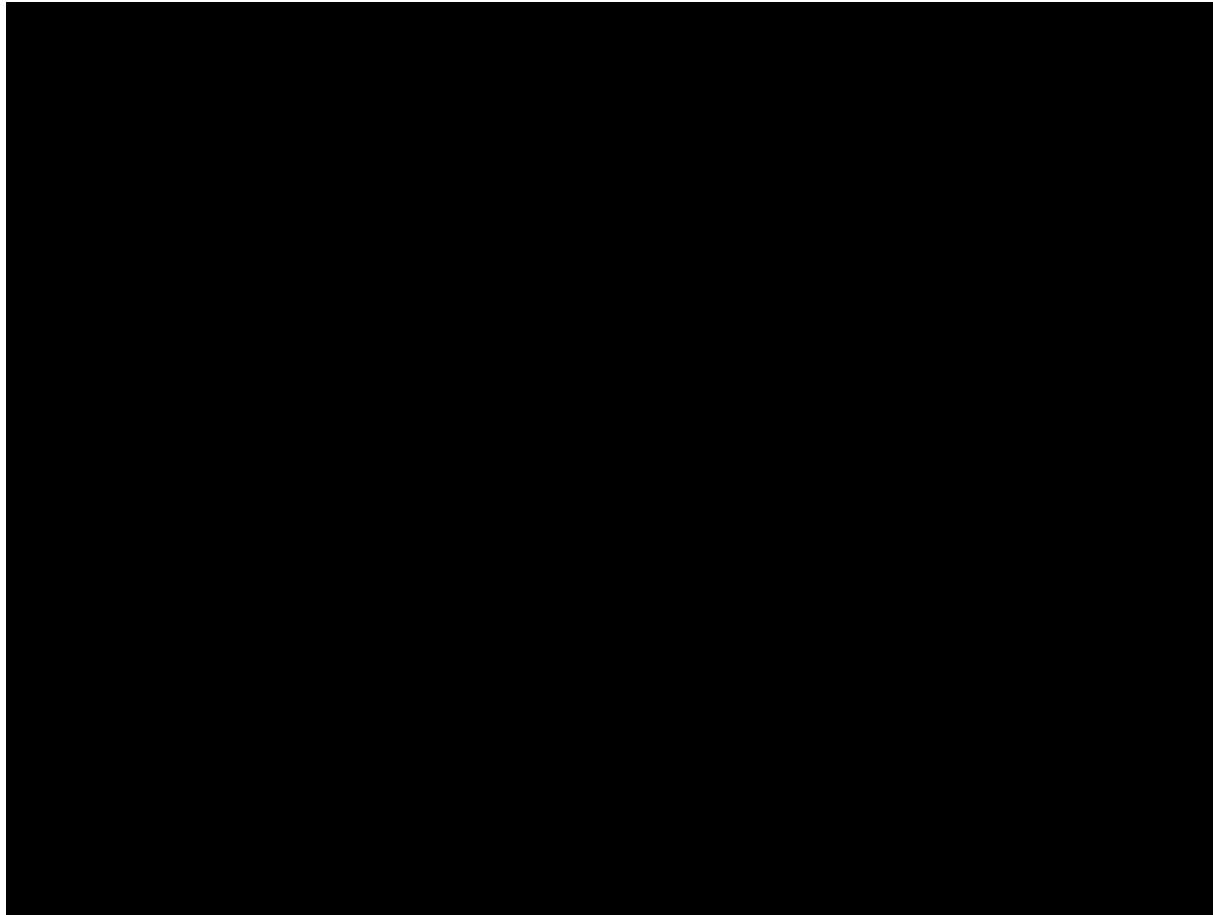


Trained on the Google news corpus with over 300 billion words

Language Models



Example of Language Model



Transcription

- **Couric:** You've cited Alaska's proximity to Russia as part of your foreign policy experience. What did you mean by that?
- **Sarah Palin:** That Alaska has a very narrow maritime border between a foreign country, Russia, and, on our other side, the land-boundary that we have with Canada. It's funny that a comment like that was kinda made to ... I don't know, you know ... reporters.
- **Couric:** Mocked?
- **Palin:** Yeah, mocked, I guess that's the word, yeah.
- **Couric:** Well, explain to me why that enhances your foreign-policy credentials.



Language Models: N-Grams

- Estimate probability of each word given prior context.
 - $P(\text{mock} \mid \text{that was kinda made to})$
- An N-gram model uses only $N-1$ words of prior context.
 - Unigram: $P(\text{mock})$
 - Bigram: $P(\text{mock} \mid \text{to})$
 - Trigram: $P(\text{mock} \mid \text{made to})$
- The *Markov assumption* is the presumption that the future behavior of a dynamical system only depends on its recent history.



N-Gram Model Formulas

- Word sequences

$$w_1^n = w_1 \dots w_n$$

- Chain rule of probability

$$P(w_1^n) = P(w_1)P(w_2 | w_1)P(w_3 | w_1^2) \dots P(w_n | w_1^{n-1}) = \prod_{k=1}^n P(w_k | w_1^{k-1})$$

- Bigram approximation

$$P(w_1^n) = \prod_{k=1}^n P(w_k | w_{k-1})$$

- N-gram approximation

$$P(w_1^n) = \prod_{k=1}^n P(w_k | w_{k-N+1}^{k-1})$$



Estimating Probabilities

- N-gram conditional probabilities can be estimated from raw text based on the *relative frequency* of word sequences.

Bigram:
$$P(w_n | w_{n-1}) = \frac{C(w_{n-1}w_n)}{C(w_{n-1})}$$

N-gram:
$$P(w_n | w_{n-N+1}^{n-1}) = \frac{C(w_{n-N+1}^{n-1}w_n)}{C(w_{n-N+1}^{n-1})}$$



Application: Speech Recognition

$$\arg \max_{wordsequence} P(wordsequence | acoustics) =$$

$$\arg \max_{wordsequence} \frac{P(acoustics | wordsequence) \times P(wordsequence)}{P(acoustics)}$$

$$\arg \max_{wordsequence} P(acoustics | wordsequence) \times P(wordsequence)$$



Language model

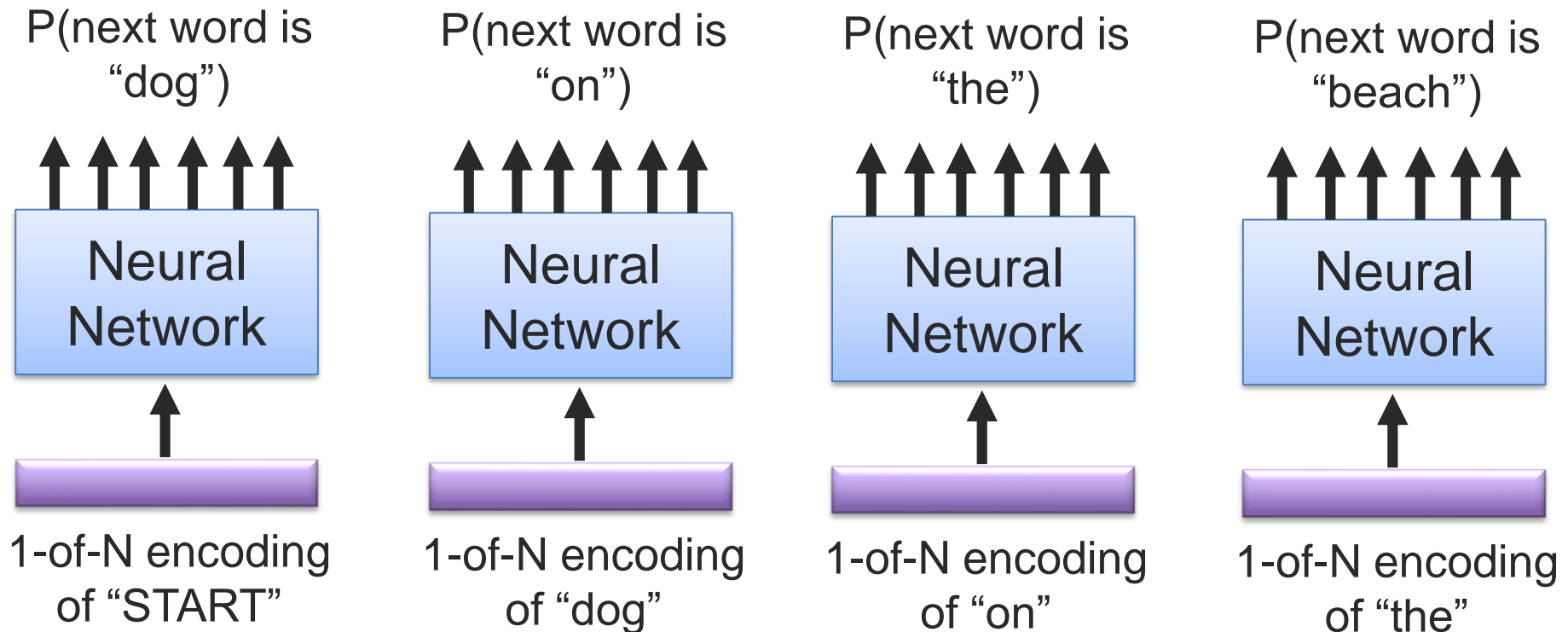


Neural-based Unigram Language Model (LM)

$P(\text{"dog on the beach"})$

$= P(\text{dog}|\text{START})P(\text{on}|\text{dog})P(\text{the}|\text{on})P(\text{beach}|\text{the})$

$P(b|a)$: not from count, but the NN that can predict the next word.

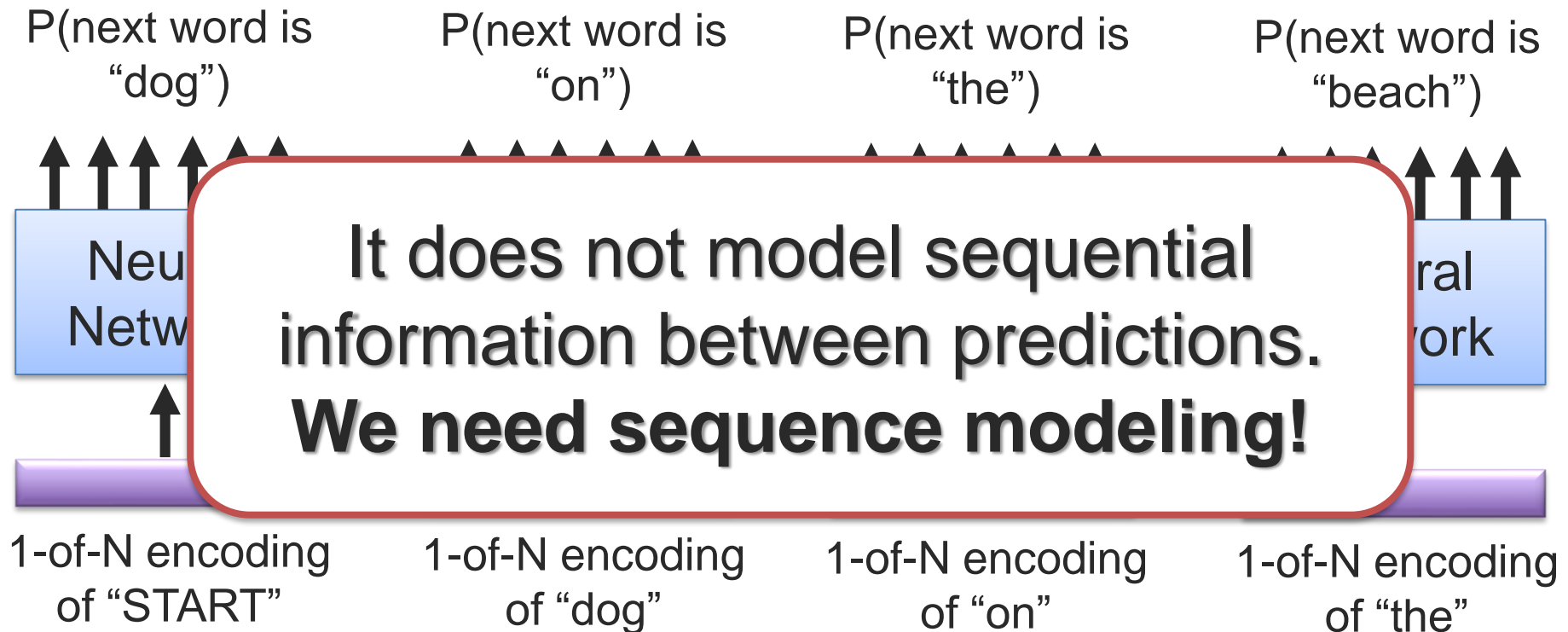


Neural-based Unigram Language Model (LM)

$P(\text{"dog on the beach"})$

$=P(\text{dog}|\text{START})P(\text{on}|\text{dog})P(\text{the}|\text{on})P(\text{beach}|\text{the})$

$P(b|a)$: not from count, but the NN that can predict the next word.



Sequence Modeling Tasks



Sequence Modeling: Language Model

★★★★★ **Masterful!**

By Antony Witheyman - January 12, 2006

Ideal for anyone with an interest in disguises who likes to see the subject tackled in a humorous manner.

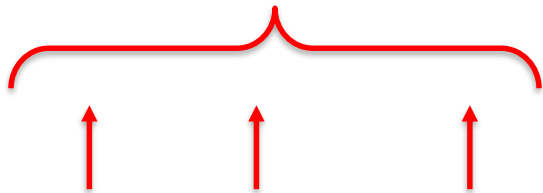
0 of 4 people found this review helpful

Prediction



Language Model

Next word?



Ideal for anyone with an interest in disguises



Sequence Modeling: Sequence Prediction

★★★★★ **Masterful!**

By Antony Witheyman - January 12, 2006

Ideal for anyone with an interest in
disguises who likes to see the subject
tackled in a humorous manner.

0 of 4 people found this review helpful

Prediction


Part-of-speech ?
(noun, verb,...)

POS?	POS?	POS?	POS?	POS?	POS?	POS?	POS?
↑	↑	↑	↑	↑	↑	↑	↑
Ideal	for	anyone	with	an	interest	in	disguises

Sequence Modeling: Sequence Label Prediction

★★★★★ **Masterful!**

By Antony Witheyman - January 12, 2006

Ideal for anyone with an interest in
disguises who likes to see the subject
tackled in a humorous manner.

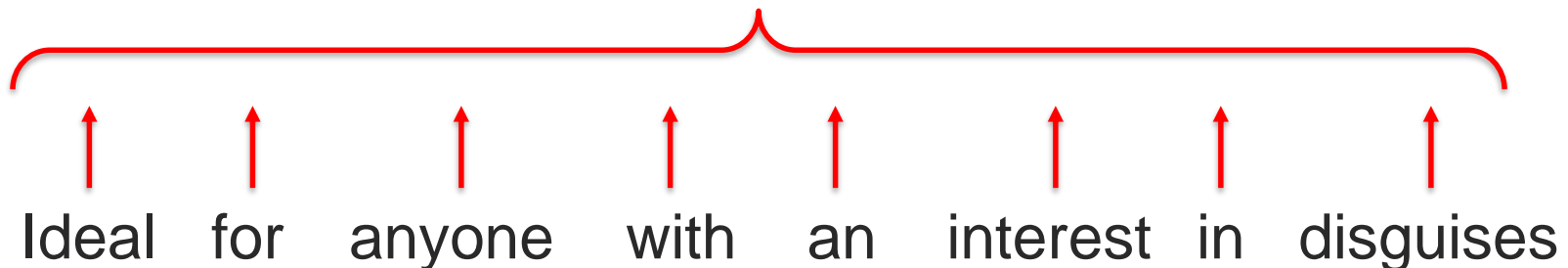
0 of 4 people found this review helpful

Prediction



Sentiment ?
(positive or negative)

Sentiment label?



Sequence Modeling: Sequence Representation

★★★★★ **Masterful!**

By Antony Witheyman - January 12, 2006

Ideal for anyone with an interest in disguises who likes to see the subject tackled in a humorous manner.

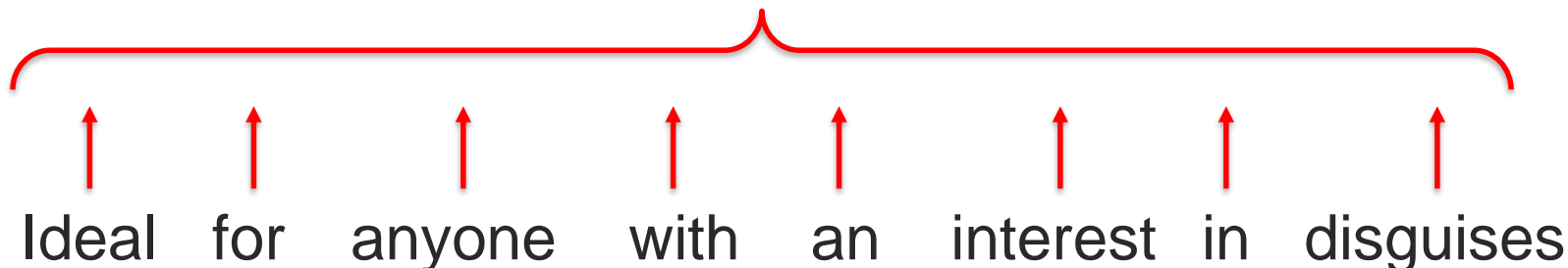
0 of 4 people found this review helpful

Learning



Sequence representation

[0,1; 0,0004; 0;....; 0,01; 0.09; 0,05]



Sequence Modeling

★★★★★ **Masterful!**

By Antony Witheyman - January 12, 2006

Ideal for anyone with an interest in disguises who likes to see the subject tackled in a humorous manner.

0 of 4 people found this review helpful



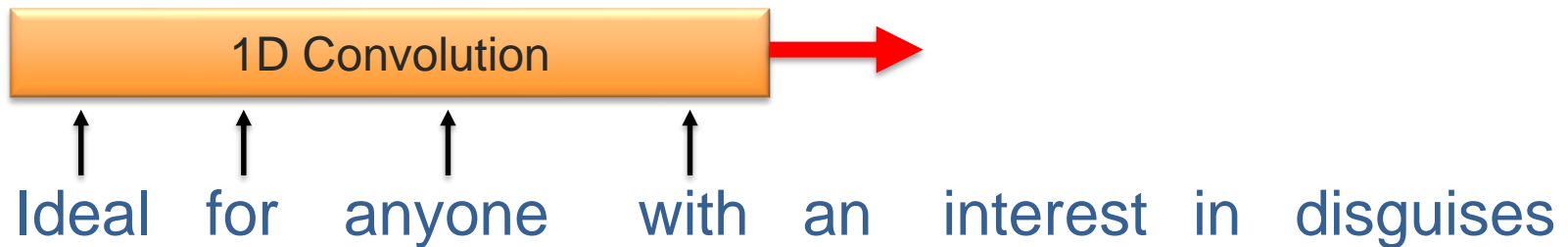
- Part-of-speech ?
(noun, verb,...)
- Sentiment ?
(positive or negative)
- Language Model
- Sequence representation

Main Challenges:

- Sequences of variable lengths (e.g., sentences)
- Keep the number of parameters at a minimum
- Take advantage of possible redundancy



Time-Delay Neural Network



Main Challenges:

- Sequences of variable lengths (e.g., sentences)
- Keep the number of parameters at a minimum
- Take advantage of possible redundancy



Recurrent Neural Networks

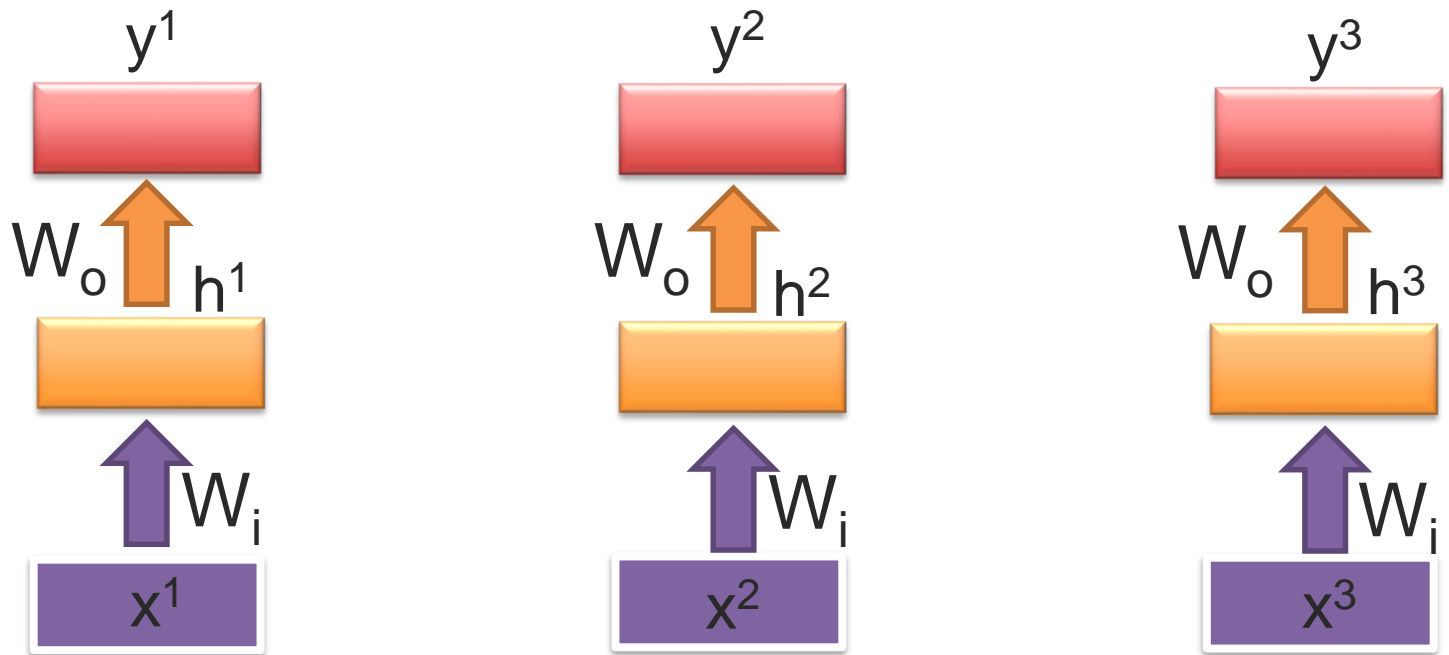


Sequence Prediction

(or Unigram Language Model)

Input data: x^1 x^2 x^3 (x^i are vectors)

Output data: y^1 y^2 y^3 (y^i are vectors)



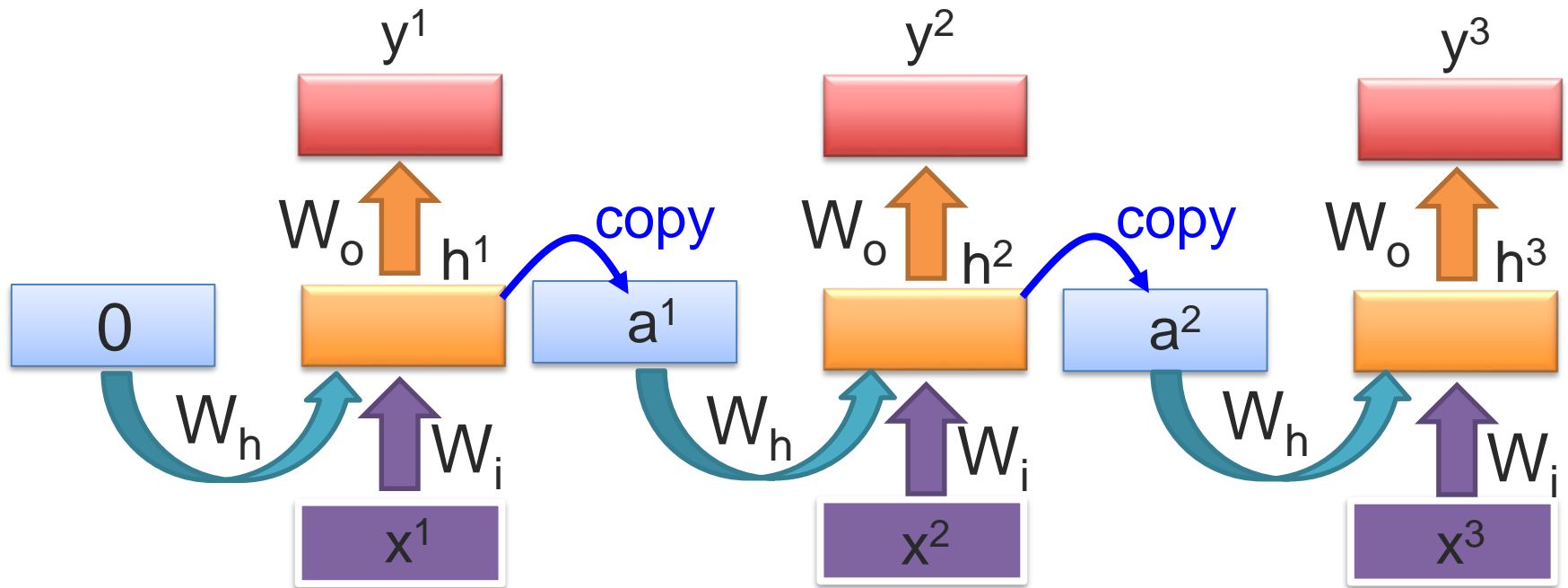
How can we include temporal dynamics?

Elman Network for Sequence Prediction

(or Unigram Language Model)

Input data: x^1 x^2 x^3 (x^i are vectors)

Output data: y^1 y^2 y^3 (y^i are vectors)

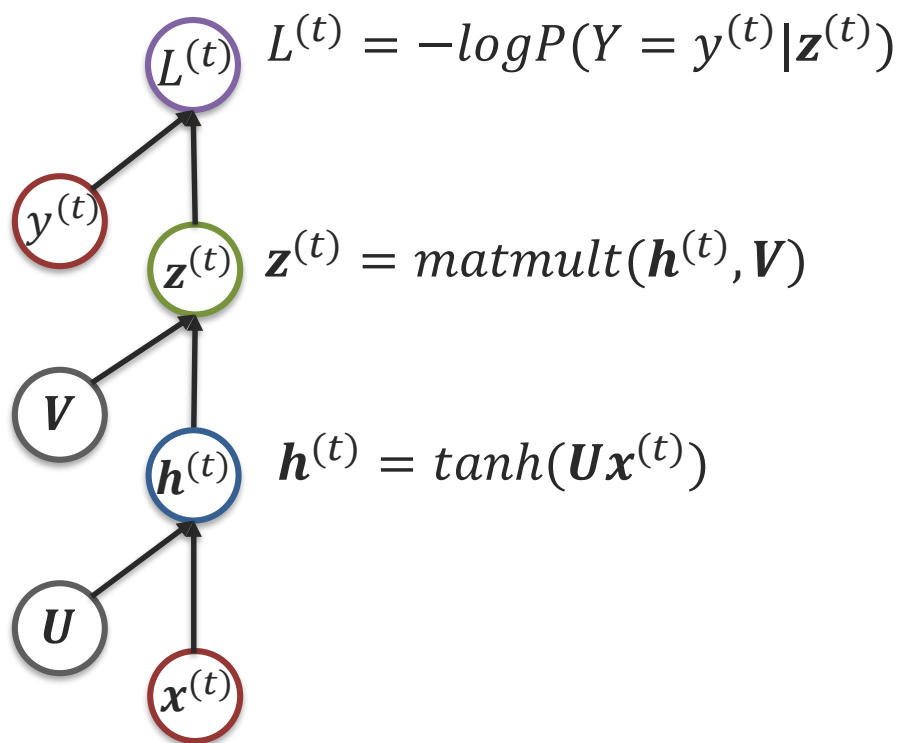


The same model parameters are used again and again.

Can be trained using backpropagation

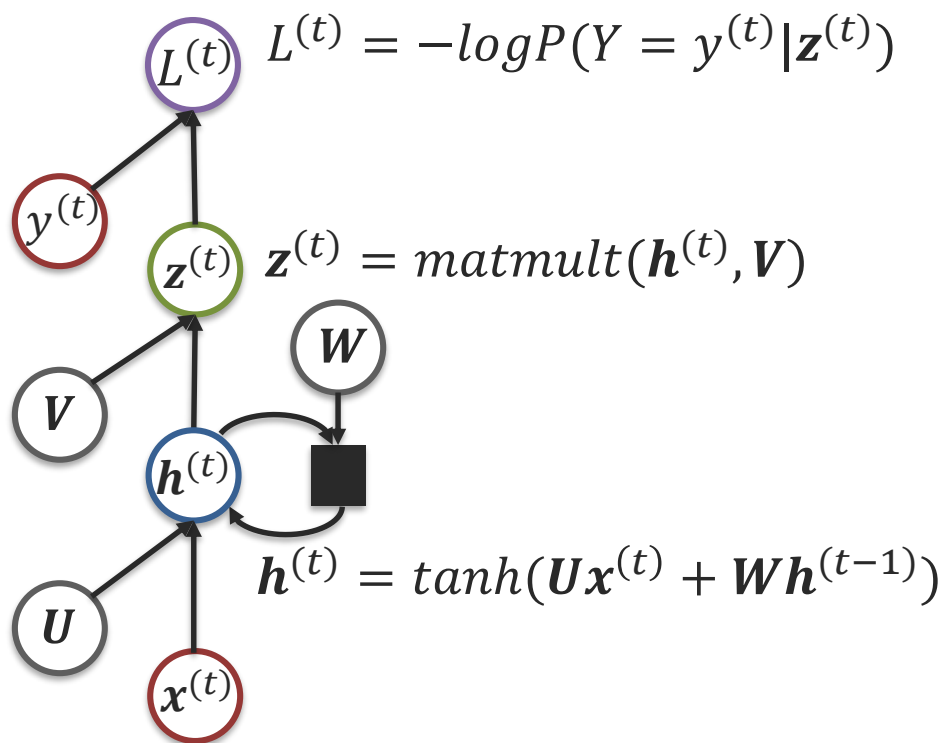
Recurrent Neural Network

Feedforward Neural Network



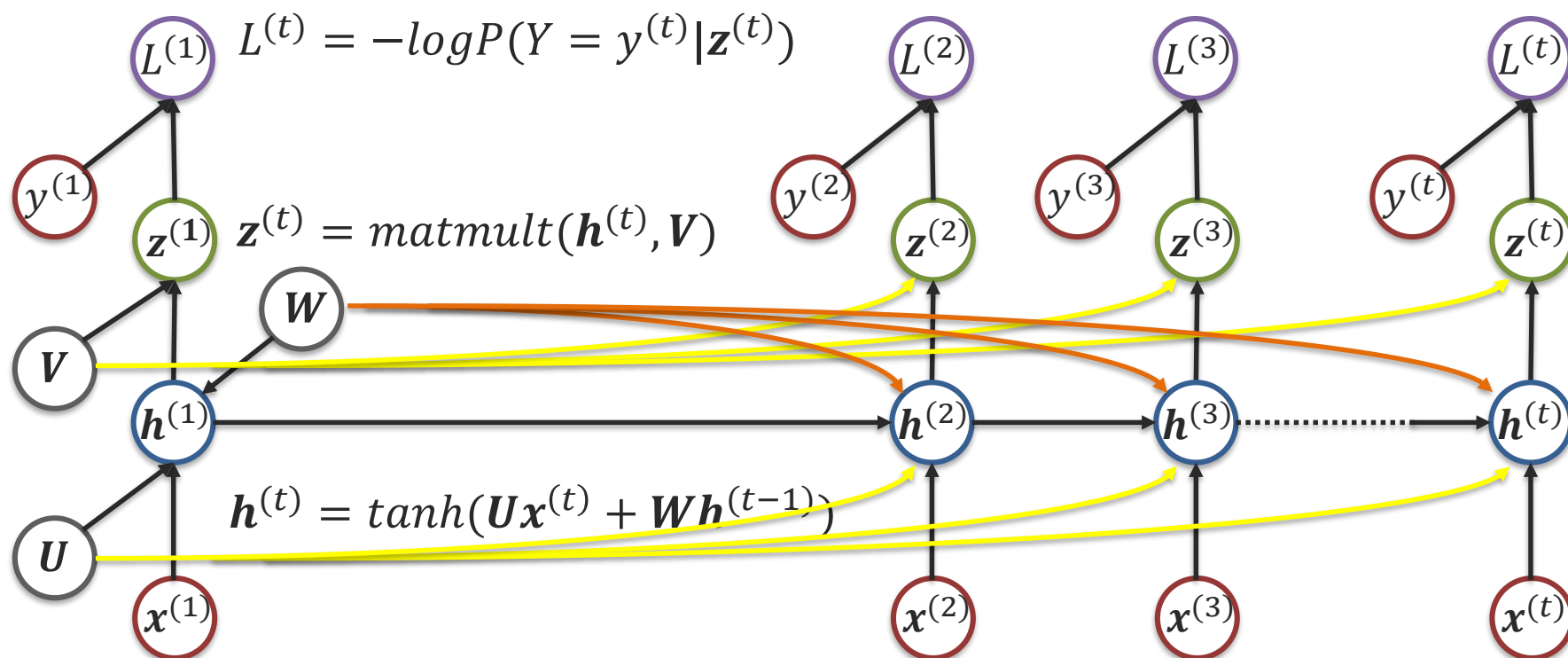
Recurrent Neural Networks

$$L = \sum_t L^{(t)}$$



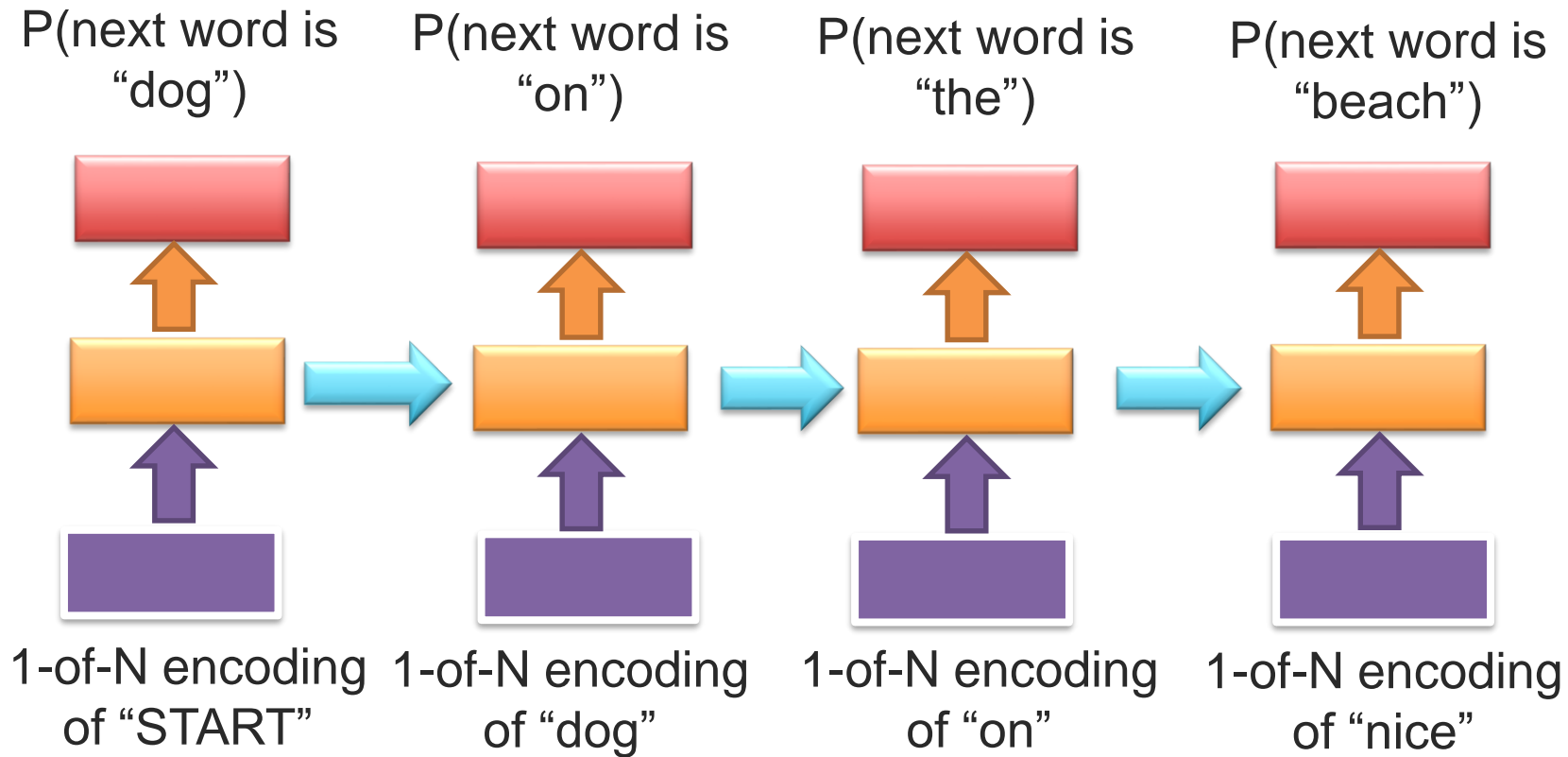
Recurrent Neural Networks - Unrolling

$$L = \sum_t L^{(t)}$$



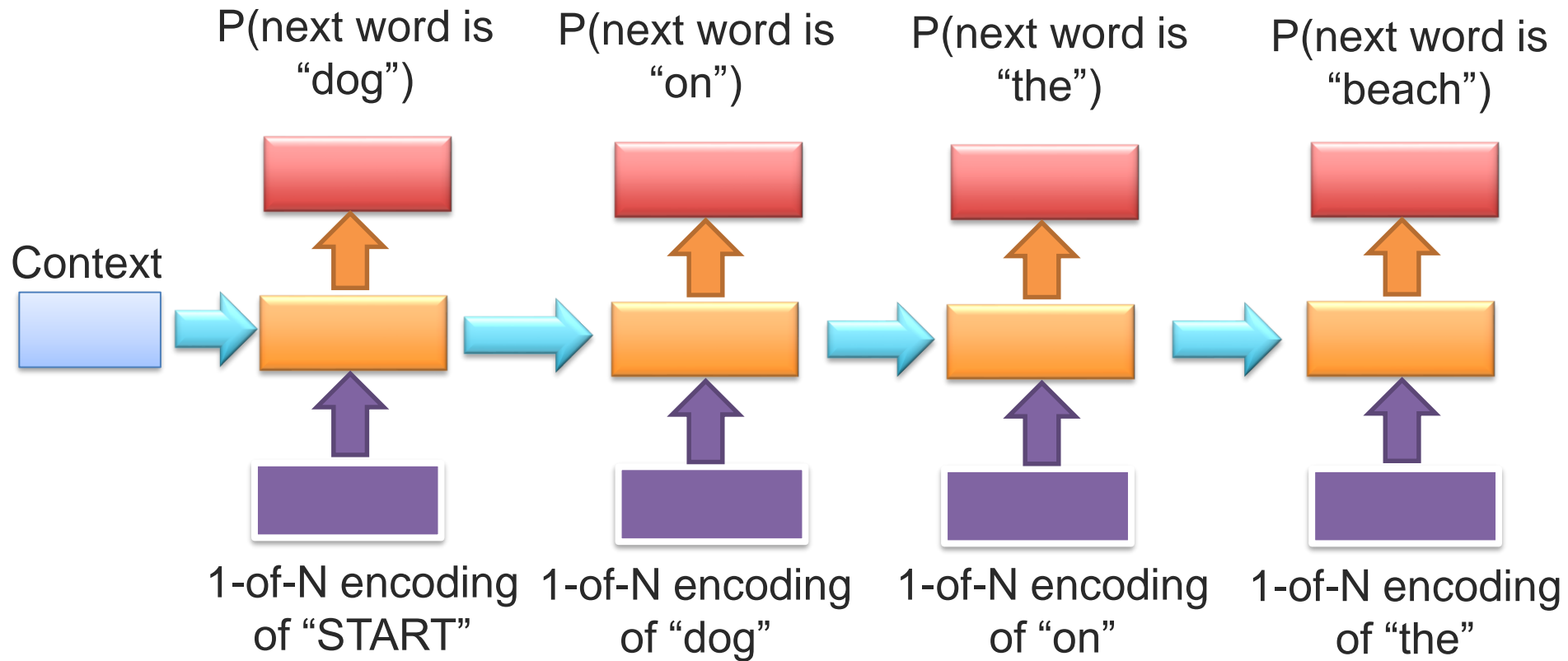
Same model parameters are used for all time parts.

RNN-based Language Model



➤ Models long-term information

RNN-based Sentence Generation (Decoder)



➤ Models long-term information

Sequence Modeling: Sequence Prediction

★★★★★ **Masterful!**

By Antony Witheyman - January 12, 2006

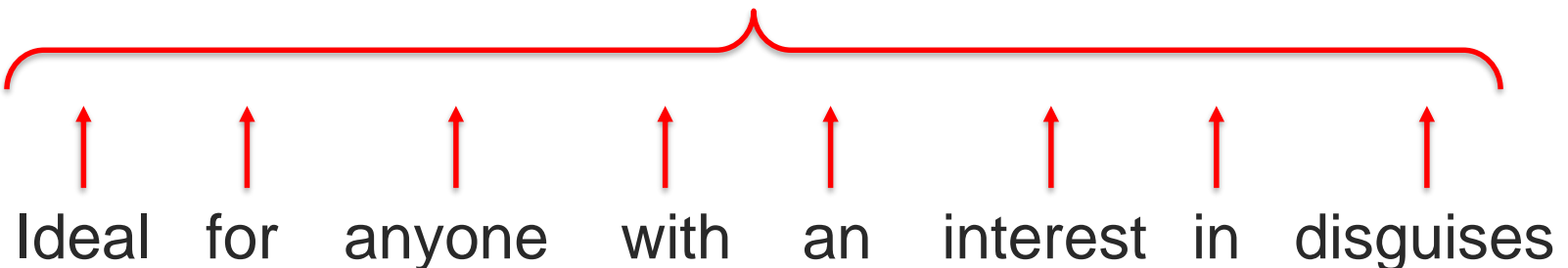
Ideal for anyone with an interest in
disguises who likes to see the subject
tackled in a humorous manner.

0 of 4 people found this review helpful

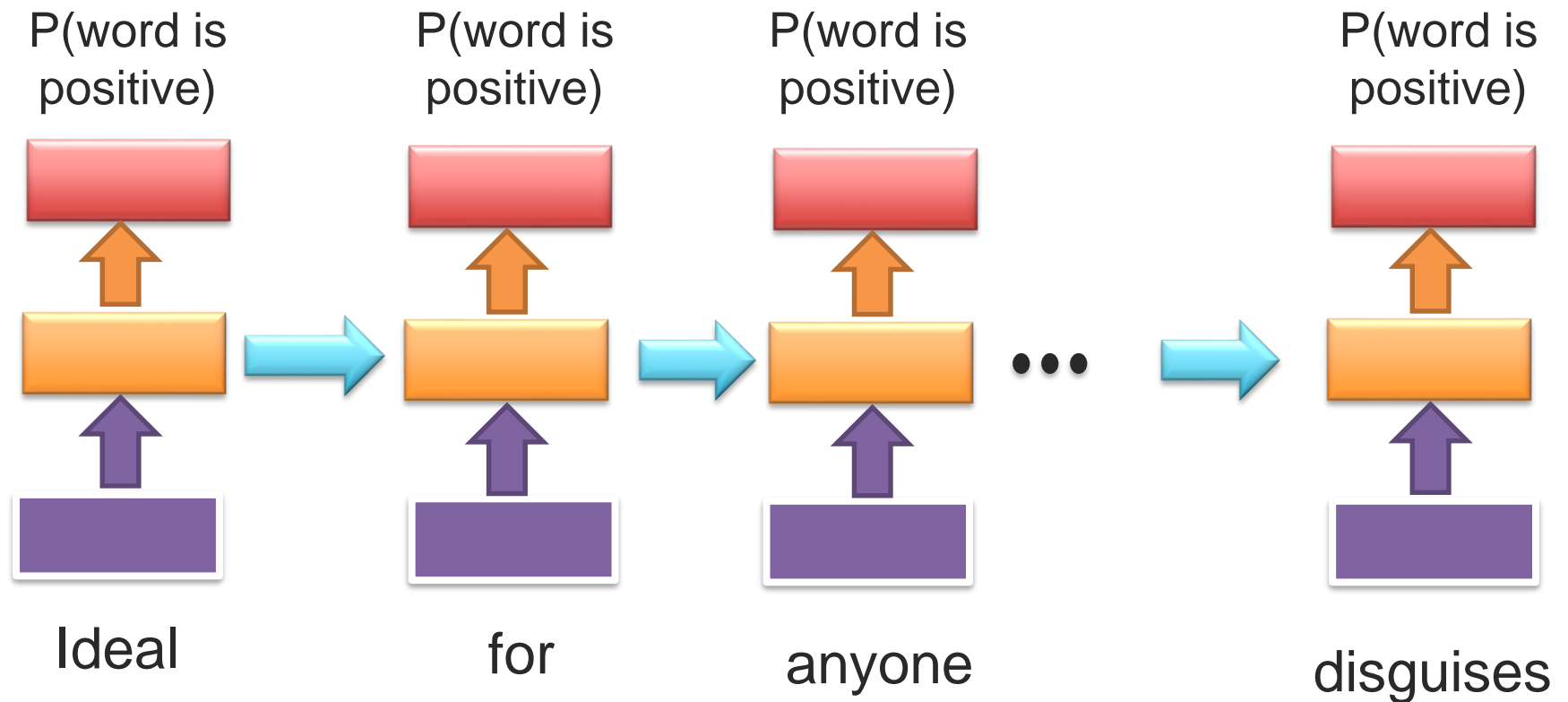
Prediction


Sentiment ?
(positive or negative)

Sentiment label?

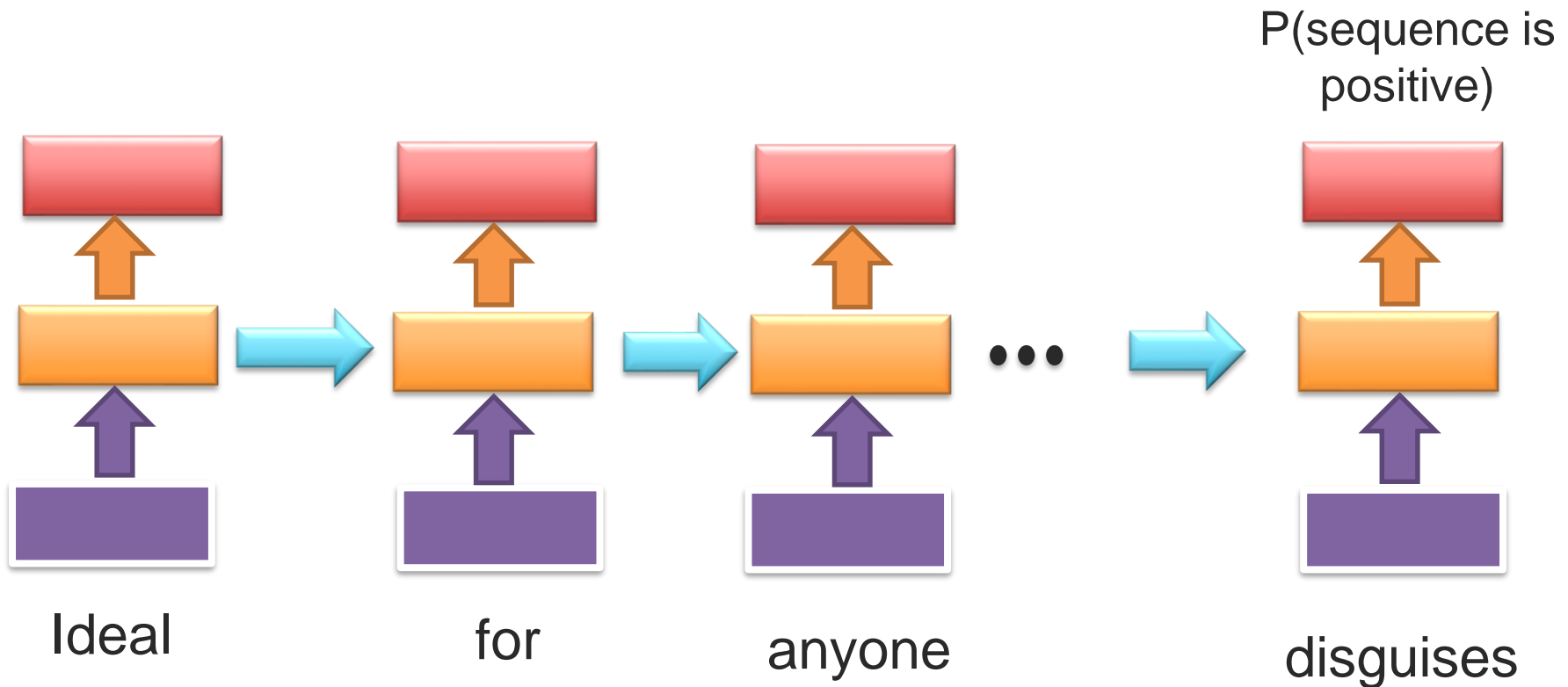


RNN for Sequence Prediction



$$L = \sum_t L^{(t)} = \sum_t -\log P(Y = y^{(t)} | \mathbf{z}^{(t)})$$

RNN for Sequence Prediction



$$L = L^{(N)} = -\log P(Y = y^{(N)} | \mathbf{z}^{(N)})$$

Sequence Modeling: Sequence Representation

★★★★★ **Masterful!**

By Antony Witheyman - January 12, 2006

Ideal for anyone with an interest in disguises who likes to see the subject tackled in a humorous manner.

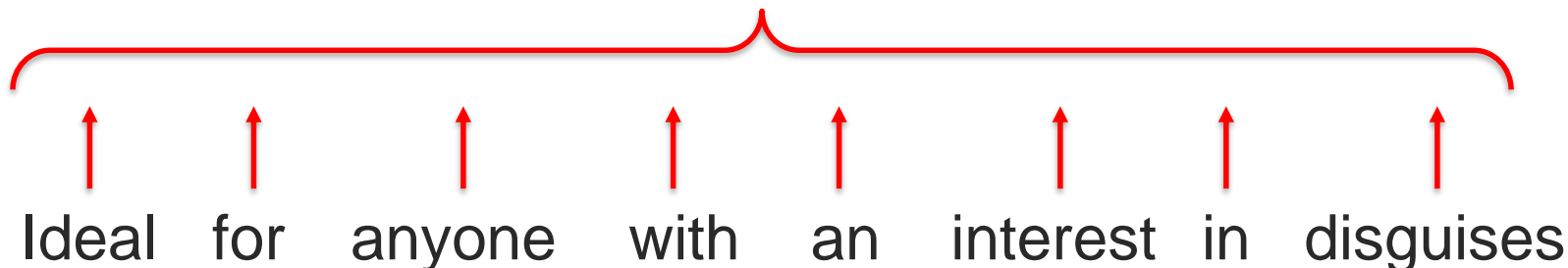
0 of 4 people found this review helpful

Learning

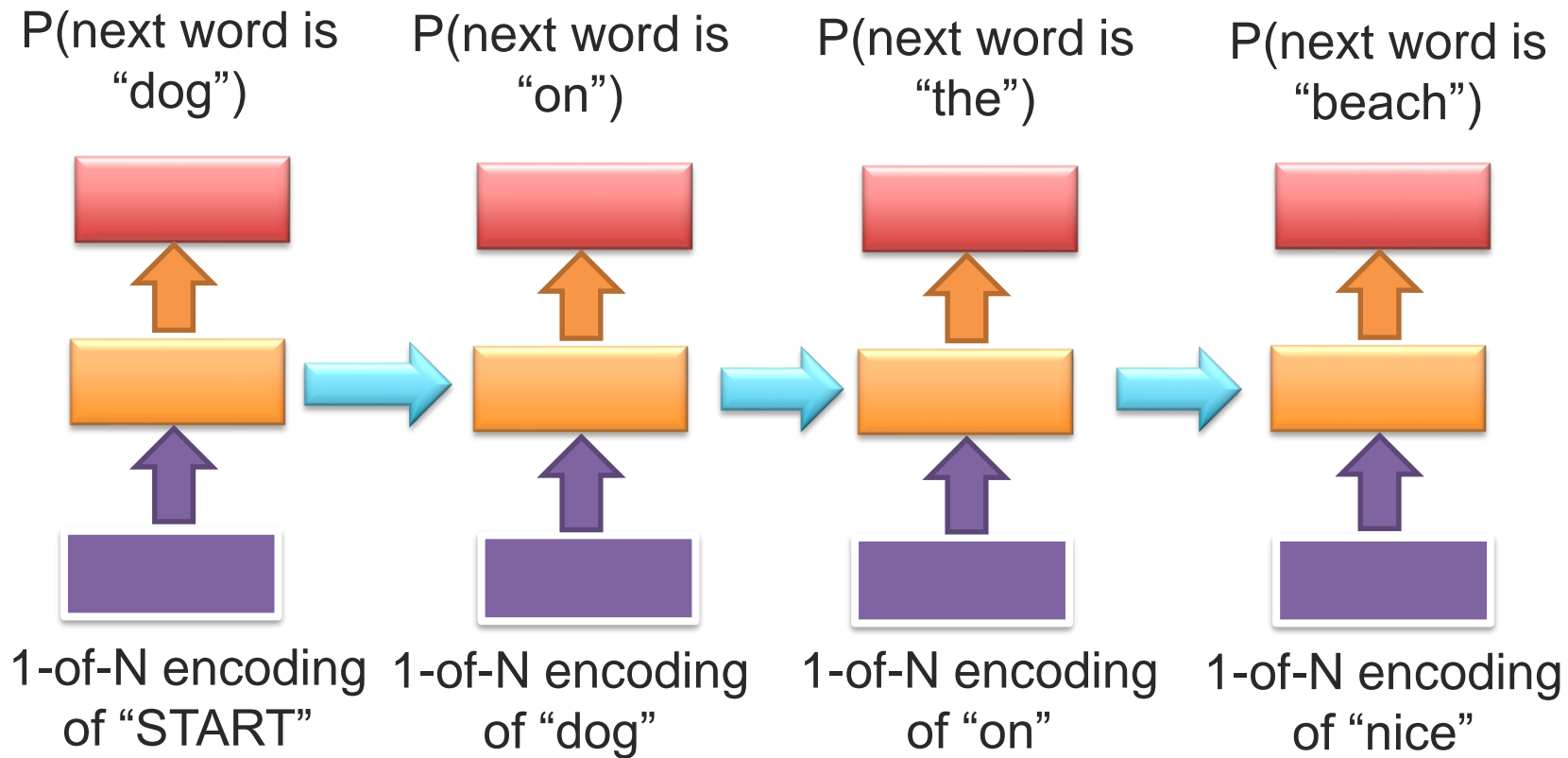


Sequence representation

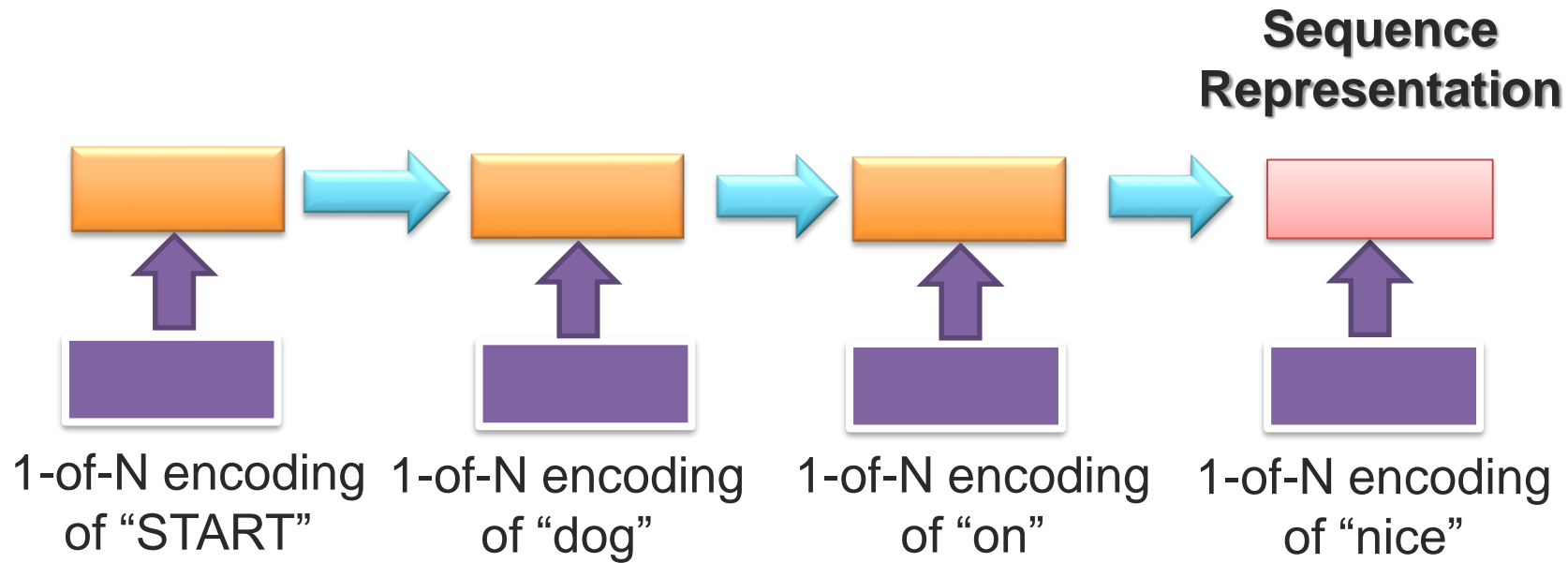
[0,1; 0,0004; 0;....; 0,01; 0.09; 0,05]



RNN for Sequence Representation

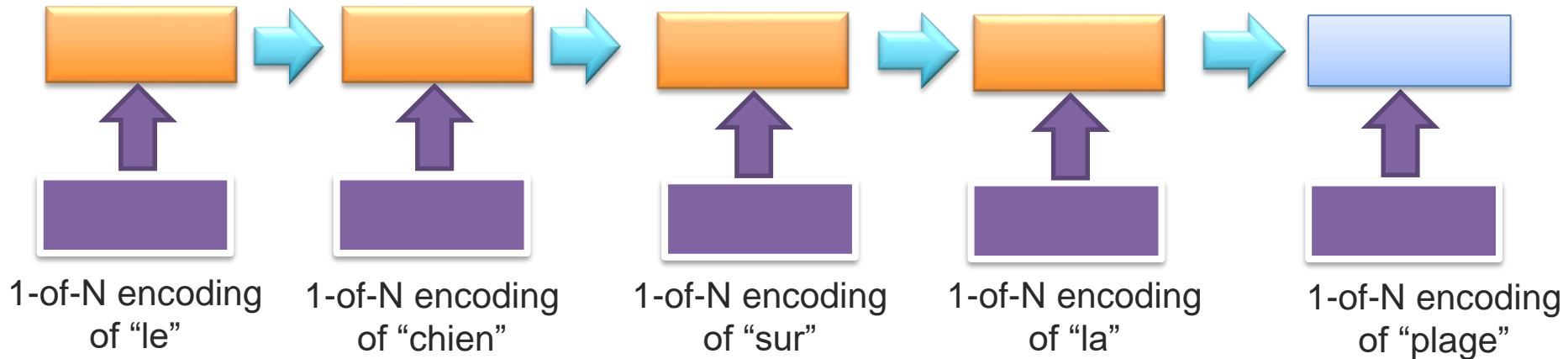


RNN for Sequence Representation (Encoder)

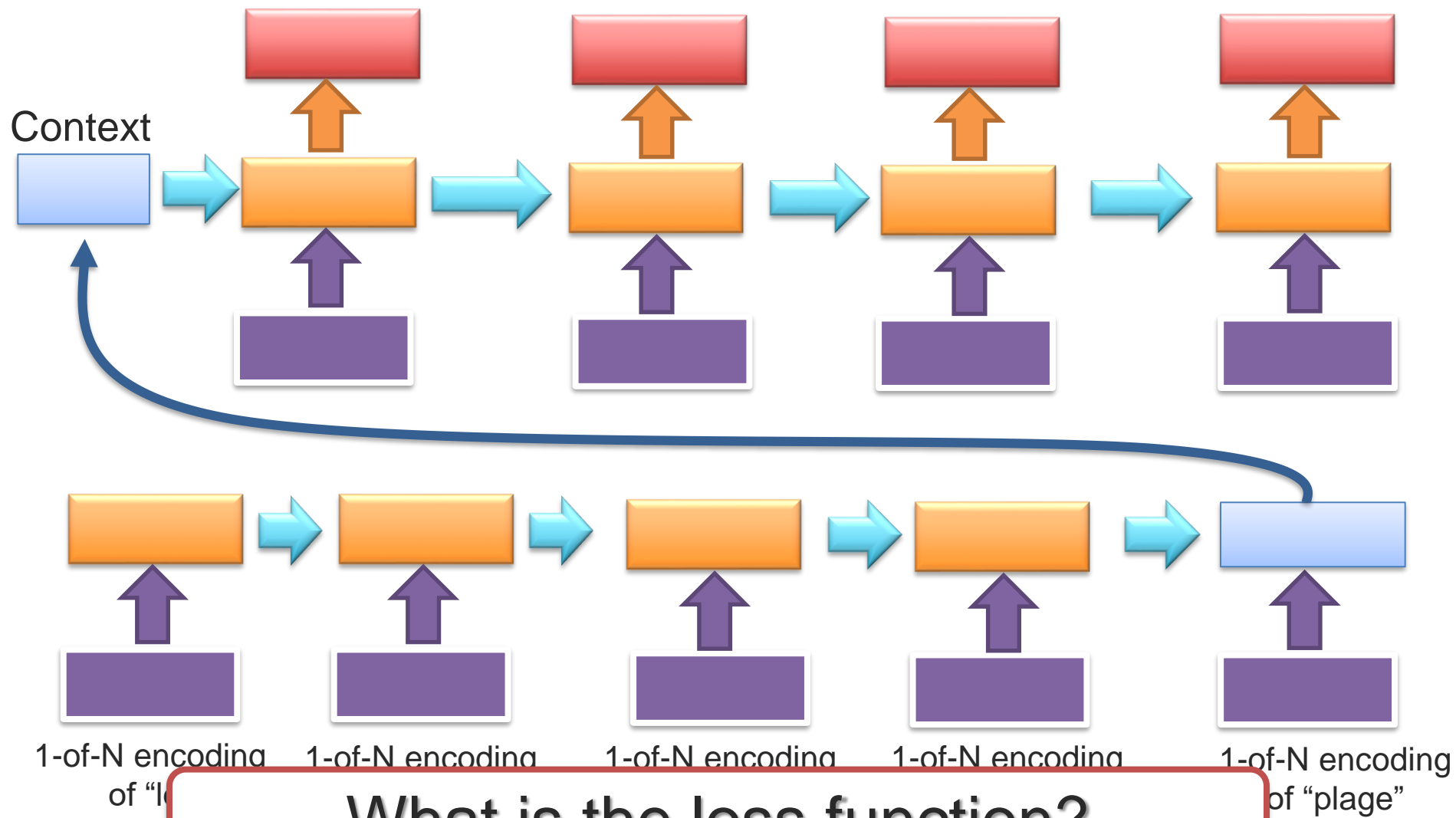


RNN-based for Machine Translation

Le chien sur la plage → The dog on the beach



Encoder-Decoder Architecture



What is the loss function?

“Advanced” Topics

- Character-level “language models”
 - Xiang Zhang, Junbo Zhao and Yann LeCun, Character-level Convolutional Networks for Text Classification, NIPS 2015
<http://arxiv.org/pdf/1509.01626v2.pdf>
- Skip-thought: embedding at the sentence level
 - Ryan Kiros, Yukun Zhu, Ruslan Salakhutdinov, Richard S. Zemel, Antonio Torralba, Raquel Urtasun, Sanja Fidler. Skip-Thought Vectors, NIPS 2015
<http://arxiv.org/pdf/1506.06726v1.pdf>



Backpropagation Through Time



Optimization: Gradient Computation

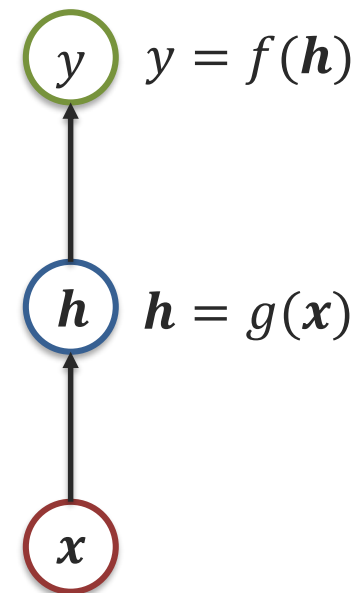
Vector representation:

Gradient $\nabla_x y = \left[\frac{\partial y}{\partial x_1}, \frac{\partial y}{\partial x_2}, \frac{\partial y}{\partial x_3} \right]$

$\nabla_x y = \left(\frac{\partial \mathbf{h}}{\partial \mathbf{x}} \right)^T \nabla_h y$

“local” Jacobian (matrix of size $|h| \times |x|$ computed using partial derivatives)

“backprop” Gradient



Backpropagation Algorithm

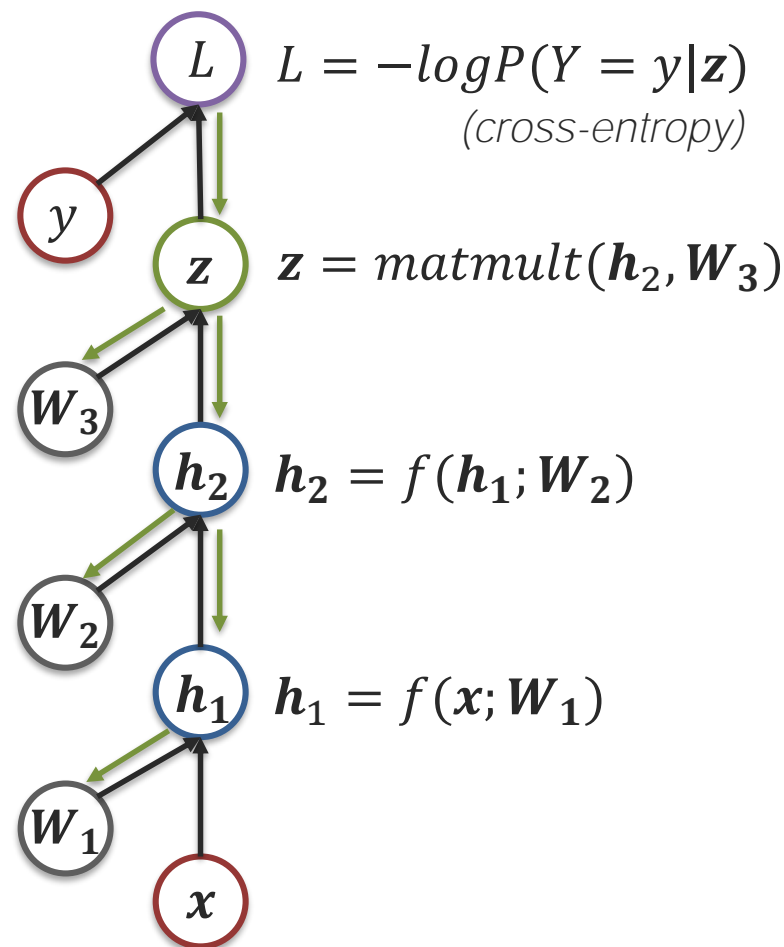
Forward pass

- Following the graph topology, compute value of each unit

Backpropagation pass

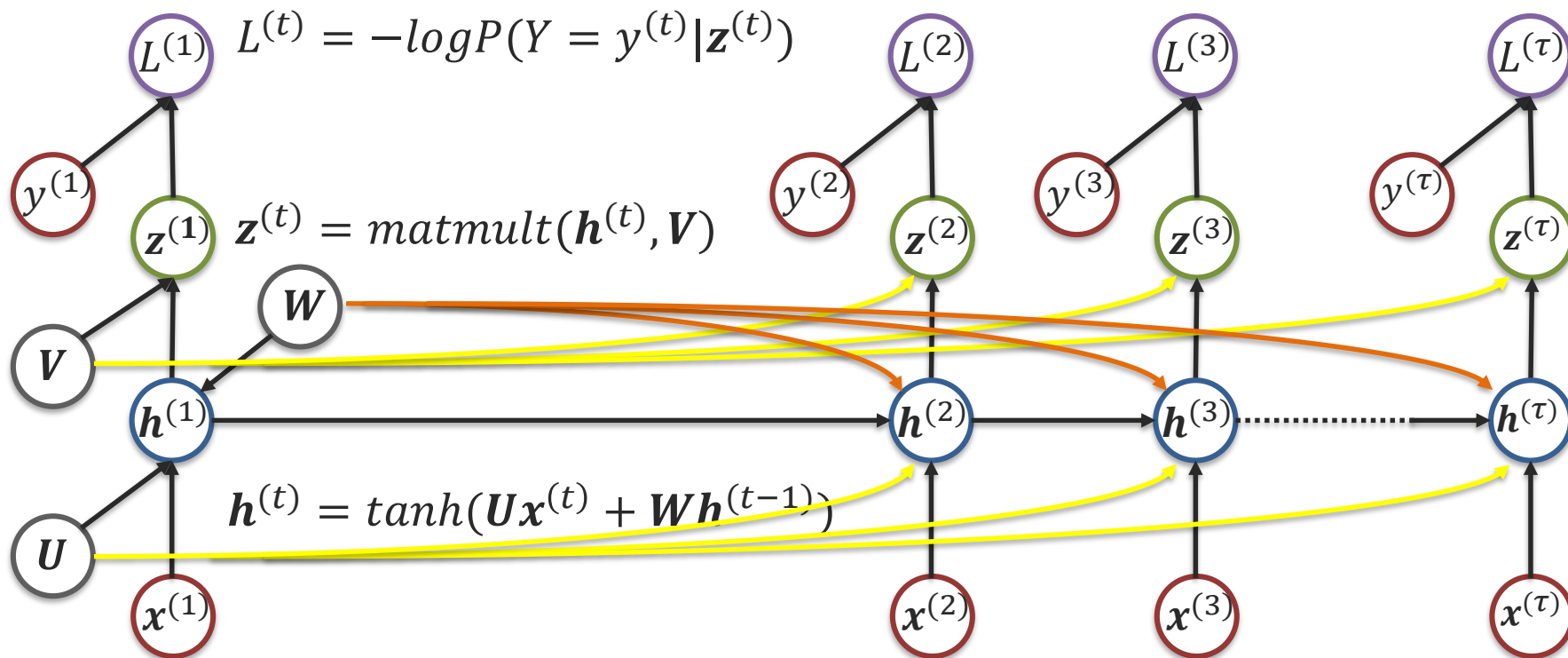
- Initialize output gradient = 1
- Compute “local” Jacobian matrix using values from forward pass
- Use the chain rule:

Gradient = “local” Jacobian \times
“backprop” gradient



Recurrent Neural Networks

$$L = \sum_t L^{(t)}$$



Backpropagation Through Time

$$L = \sum_t L^{(t)} = - \sum_t \log P(Y = y^{(t)} | \mathbf{z}^{(t)})$$

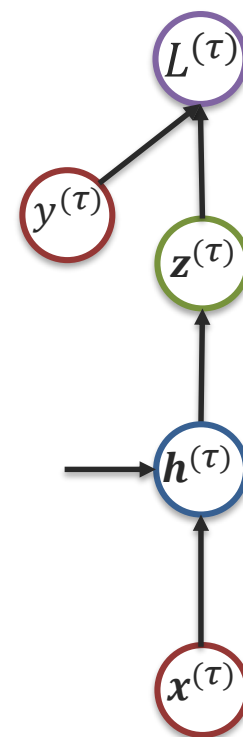
Gradient = “backprop” gradient
x “local” Jacobian

$$L^{(\tau)} \text{ or } L^{(t)} \quad \frac{\partial L}{\partial L^{(t)}} = 1$$

$$\mathbf{z}^{(\tau)} \text{ or } \mathbf{z}^{(t)} \quad (\nabla_{\mathbf{z}^{(t)}} L)_i = \frac{\partial L}{\partial z_i^{(t)}} = \frac{\partial L}{\partial L^{(t)}} \frac{\partial L^{(t)}}{\partial z_i^{(t)}} = \text{sigmoid}(z_i^t) - \mathbf{1}_{i,y^{(t)}}$$

$$\mathbf{h}^{(\tau)} \quad \nabla_{\mathbf{h}^{(\tau)}} L = \nabla_{\mathbf{z}^{(\tau)}} L \frac{\partial \mathbf{z}^{(\tau)}}{\partial \mathbf{h}^{(\tau)}} = \nabla_{\mathbf{z}^{(\tau)}} L \mathbf{V}$$

$$\mathbf{h}^{(t)} \rightarrow \mathbf{h}^{(t+1)} \quad \nabla_{\mathbf{h}^{(t)}} L = \nabla_{\mathbf{z}^{(t)}} L \frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{h}^{(t)}} + \nabla_{\mathbf{z}^{(t+1)}} L \frac{\partial \mathbf{h}^{(t+1)}}{\partial \mathbf{h}^{(t)}}$$



Backpropagation Through Time

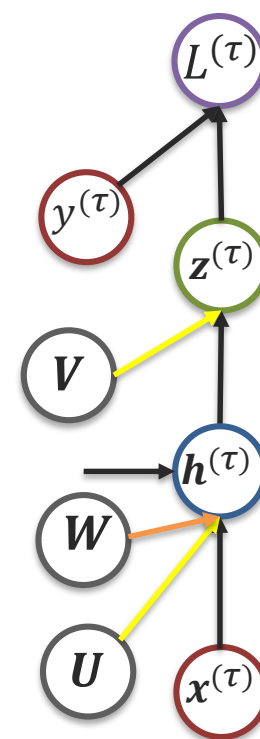
$$L = \sum_t L^{(t)} = - \sum_t \log P(Y = y^{(t)} | \mathbf{z}^{(t)})$$

Gradient = “backprop” gradient
x “local” Jacobian

$$\textcircled{V} \quad \nabla_V L = \sum_t (\nabla_{\mathbf{z}^{(t)}} L) \frac{\partial \mathbf{z}^{(t)}}{\partial V}$$

$$\textcircled{W} \quad \nabla_W L = \sum_t (\nabla_{\mathbf{h}^{(t)}} L) \frac{\partial \mathbf{h}^{(t)}}{\partial W}$$

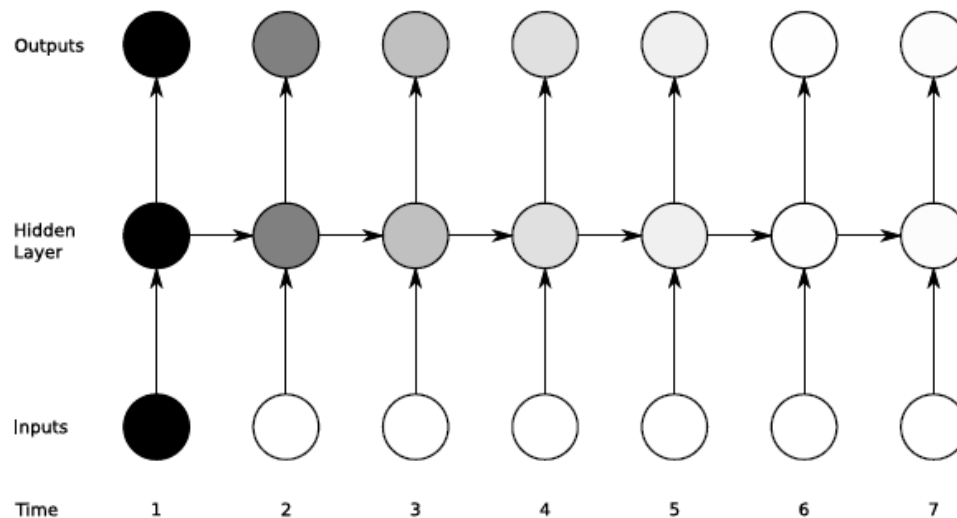
$$\textcircled{U} \quad \nabla_U L = \sum_t (\nabla_{\mathbf{h}^{(t)}} L) \frac{\partial \mathbf{h}^{(t)}}{\partial U}$$



Long-term Dependencies

Vanishing gradient problem for RNNs:

$$\mathbf{h}^{(t)} \sim \tanh(\mathbf{W}\mathbf{h}^{(t-1)})$$



- The influence of a given input on the hidden layer, and therefore on the network output, either decays or blows up exponentially as it cycles around the network's recurrent connections.