Language Technologies Institute

Carnegie Mellon University

# Multimodal Machine Learning

## Lecture 3.2: Language Representations and RNNs

**Louis-Philippe Morency**

* Original version co-developed with Tadas Baltrusaitis

# Administrative Stuff

# Piazza Live Q&A – Reminder

# Lecture Highlights - Reminder

https://forms.gle/W1VJDWaitEumn4XSA

New set of instructions…

　　…but same deadline: **Saturday 10:40am**

**IMPORTANT:** Be sure you received an email after your submission (or revisit the form and your answers should be there).

Language Technologies Institute

Carnegie Mellon University

# Reading Assignments – Reminder

Week 3 reading assignment was posted

1. **Friday 8pm:** Post your summary
2. **Monday 8pm:** End of the reading assignment

**Be sure to post your discussion comments, questions and answers before Monday 8pm!**

➡️ Start the discussion early ☺

➡️ Late submissions will be penalized

# Grades on Canvas



Grades are now on CMU Canvas for:

❑ Lecture highlights

❑ Reading assignments

Grades for the project assignments will be on gradescope

Language Technologies Institute

Carnegie Mellon University

**Multimodal Machine Learning**

**Lecture 3.2: Language Representations and RNNs**

**Louis-Philippe Morency**

**\* Original version co-developed with Tadas Baltrusaitis**

# Lecture Objectives

- Word representations
  - Distributional hypothesis
  - Learning neural representations

- Sentence representations and sequence modeling
  - Recurrent neural networks
  - Gated recurrent neural networks
  - Backpropagation through time

- Syntax and language structure
  - Phrase-structure and dependency grammars
  - Recursive neural network
    - Tree-based RNN, Stack LSTM

# Word Representations

# What is the meaning of "bardiwac"?

- He handed her 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

# 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

He was walking away because …
He was running away because …

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

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

# 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 $R^n$

|        | 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*: $\mathbf{x}_{dog}$ = (115, 10)

- similarity = spatial proximity (Euclidean distance)

- location depends on frequency of noun ($f_{dog} \approx 2.7 \cdot f_{cat}$)



Two dimensions of English V−Obj DSM

# Angle and similarity

- direction more important than location

- normalise "length" $\|\mathbf{x}_{dog}\|$ of vector

- or use angle $\alpha$ as distance measure



Two dimensions of English V–Obj DSM

Stefan Evert 2010

# How to learn (word) features/representations?

walking

x  (100 000d)

$W_1$

(300d)

$W_2$

(300d)

y  (100 000d)

He
Was

Away
because

[0; 0; 0; 0;…., 0; 0; 1; 0;…; 0; 0]

He was walking away because …
He was running away because …

[0; 1; 0; 0;…., 0; 0; 0; 0;…; 0; 0]

[0; 0; 0; 1;…., 0; 0; 0; 0;…; 0; 0]

[0; 0; 0; 0;…., 1; 0; 0; 0;…; 0; 0]

[0; 0; 0; 0;…., 0; 0; 0; 0;…; 0; 1]

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

Language Technologies Institute

Carnegie Mellon University

# How to use these word representations

If we would have a vocabulary of 100 000 words:

Classic NLP:     ⟵ 100 000 dimensional vector ⟶

Walking:     [0; 0; 0; 0;….; 0; 0; 1; 0;…; 0; 0]

Running:     [0; 0; 0; 0;….; 0; 0; 0; 0;…; 1; 0]

➡ Similarity = 0.0

⬇ Transform: x'=x*W

Goal:     ⟵ 300 dimensional vector ⟶

Walking:     [0,1; 0,0003; 0;….; 0,02; 0.08; 0,05]

Running:     [0,1; 0,0004; 0;….; 0,01; 0.09; 0,05]

➡ Similarity = 0.9

100 000d     x     $W_1$

300d

# 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:

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

# Vector space models of words: semantic relationships



Trained on the Google news corpus with over 300 billion words

# Word Representation Resources

Word-level representations:

    Word2Vec (Google, 2013)

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

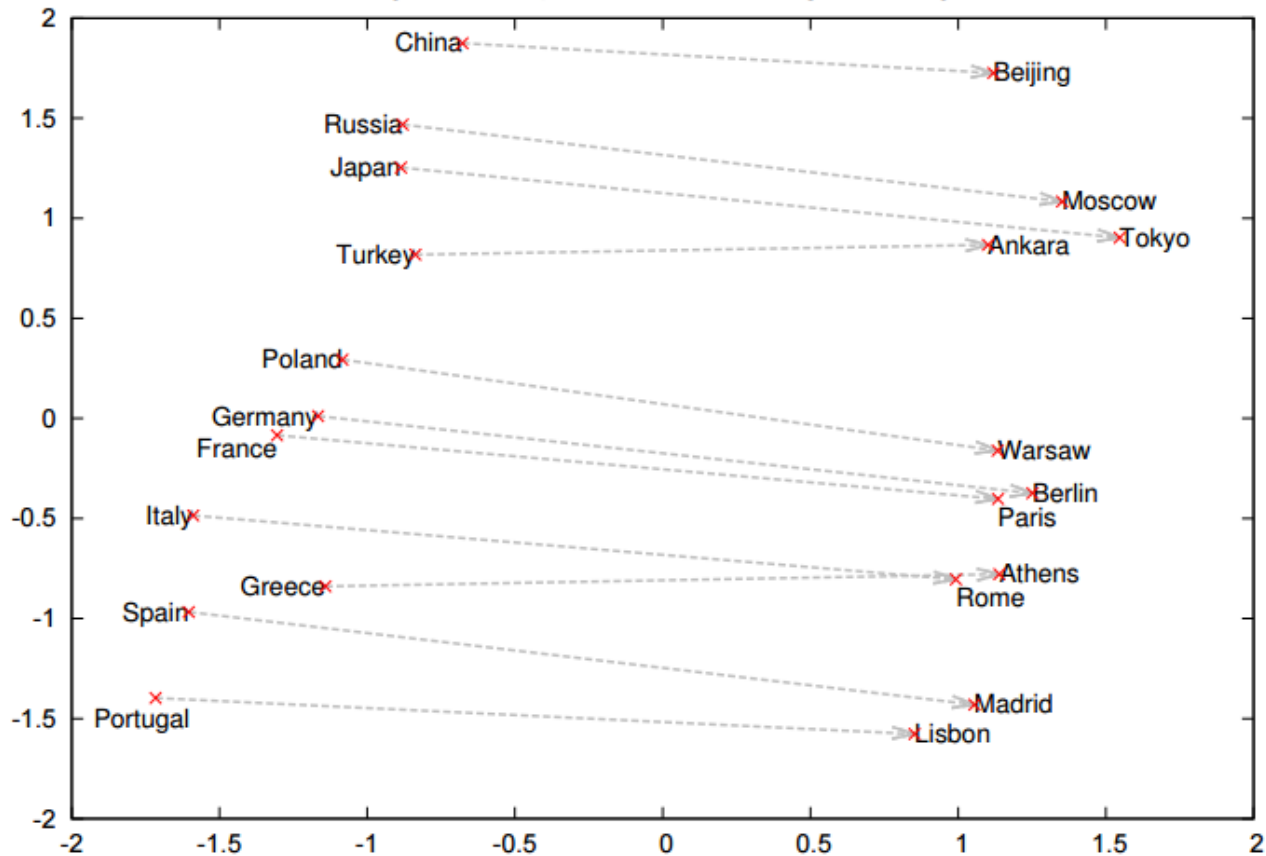    Glove (Stanford, 2014)

        https://nlp.stanford.edu/projects/glove/

    FastText (Facebook, 2017)

        https://fasttext.cc/

Sentence-level representations:

    ELMO (Allen Institute for AI, 2018)

        https://allennlp.org/elmo

    BERT (Google, 2018)

        https://github.com/google-research/bert

    RoBERTa (Facebook, 2019)

        https://github.com/pytorch/fairseq

Word representations are contextualized using all the words in the sentence.

More details later in this lecture and during Week 5

# Lexicon-based Word Representation

**LIWC: Language Inquiry & Word Count**

Manually created dictionaries for different topics and categories:

- Function words: *pronouns, preposition, negation…*
- Affect words: *positive, negative emotions*
- Social words: *family, friends, referents*
- Cognitive processes: *Insight, cause, …*
- Perceptual processes: *Seeing, hearing, feeling*
- Biological processes: *Body, health/illness,…*
- Drives and needs*: Affiliation, achievement, …*
- Time orientation: *past, present, future*
- Relativity: *motion, space, time*
- Personal concerns: *work, leisure, money, religion …*
- Informal speech: *swear words, fillers, assent,…*

LIWC can encode individual words or full sentences.

https://liwc.wpengine.com/

Commercial software. Contact TAs in advance if you would like to use it.

# Other Lexicon Resources

Lexicons

- General Inquirer (Stone et al., 1966)
- OpinionFinder lexicon (Wiebe & Riloff, 2005)
- SentiWordNet (Esuli & Sebastiani, 2006)
- LIWC (Pennebaker)

Other Tools

- LightSIDE
- Stanford NLP toolbox
- IBM Watson Tone Analyzer
- Google Cloud Natural Language
- Microsoft Azure Text Analytics

# Sentence Modeling

# Sentence 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 humourous manner.

0 of 4 people found this review helpful

Prediction → Sentiment ?
(positive or negative)

**Sentiment label?**

Ideal   for   anyone   with   an   interest   in   disguises

# Sentence 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 humourous 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 |

# Sentence Modeling: Sequence Representation

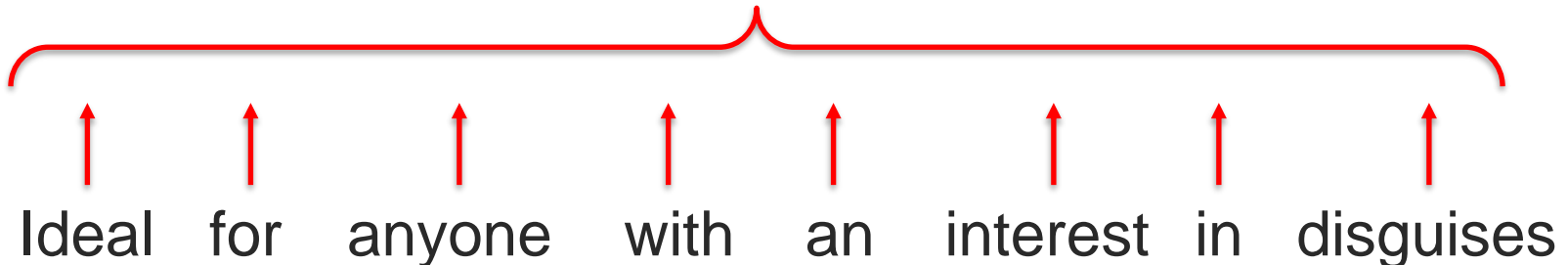⭐⭐⭐⭐⭐ **Masterful!**

By Antony Witheyman - January 12, 2006

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

0 of 4 people found this review helpful

Learning ➡ Sequence representation

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

Ideal   for   anyone   with   an   interest   in   disguises

Language Technologies Institute

Carnegie Mellon University

# Sentence 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 humourous manner.
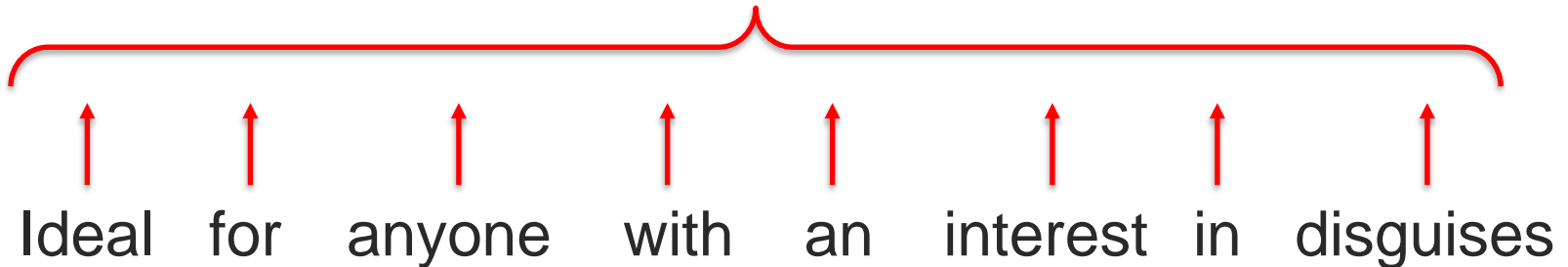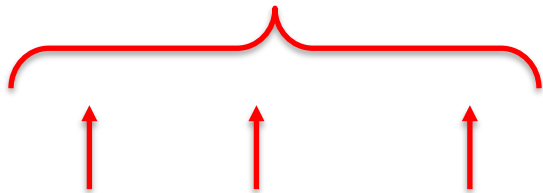
0 of 4 people found this review helpful

Prediction →

Language Model

**Next word?**

Ideal    for    anyone    with    an    interest    in    disguises

# Language Model Application: Language Generation

**Embedding**
[0,1;
0,0004;
….;
0.09;
0,05]

Generation →

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

## Example: Image captioning



→

[0,1;
0,0004;
….;
0.09;
0,05]

→

The man at bat readies to swing at the pitch while the umpire looks on.

Language Technologies Institute

Carnegie Mellon University

# Language Model Application: Speech Recognition

$$\underset{wordsequence}{\arg\max}\, P(wordsequence \mid acoustics) =$$

$$\underset{wordsequence}{\arg\max}\, \frac{P(acoustics \mid wordsequence) \times P(wordsequence)}{P(acoustics)}$$

$$\underset{wordsequence}{\arg\max}\, P(acoustics \mid wordsequence) \times P(wordsequence)$$

**Language model**

# Challenges in 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 humourous manner.

0 of 4 people found this review helpful

→ Model →

- 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

Carnegie Mellon University
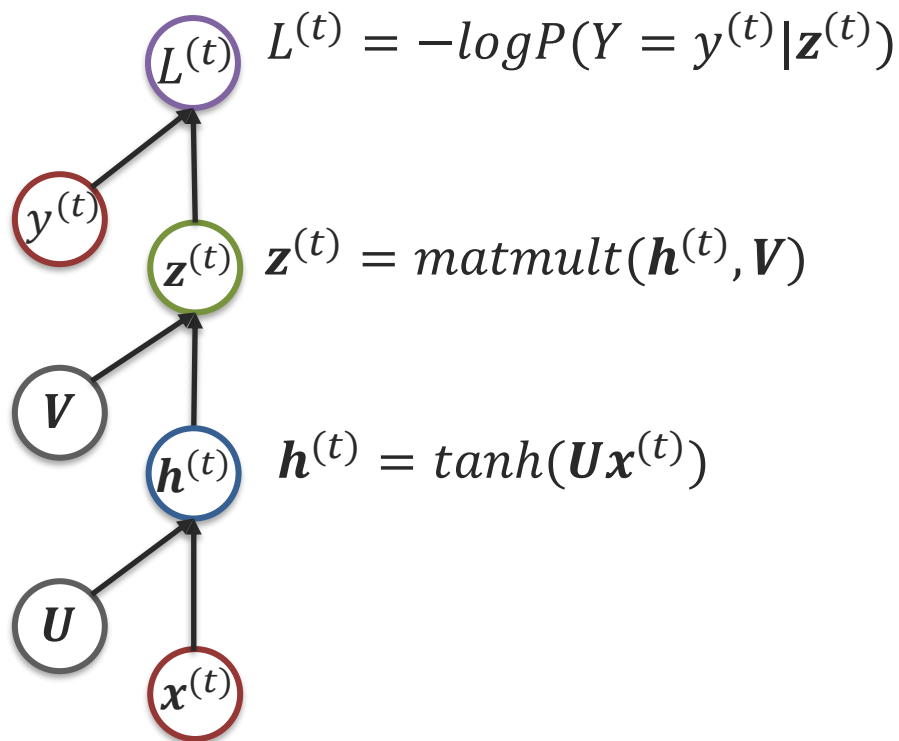
# Recurrent Neural Networks

Language Technologies Institute

**Carnegie Mellon University**

# Recurrent Neural Network

## Feedforward Neural Network



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

$$\mathbf{z}^{(t)} = matmult(\mathbf{h}^{(t)}, \mathbf{V})$$

$$\mathbf{h}^{(t)} = tanh(\mathbf{U}\mathbf{x}^{(t)})$$

Language Technologies Institute

Carnegie Mellon University

# Recurrent Neural Networks

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

$L^{(t)}$

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

$y^{(t)}$

$\mathbf{z}^{(t)}$

$$\mathbf{z}^{(t)} = matmult(\mathbf{h}^{(t)}, \mathbf{V})$$

$W$

$V$

$\mathbf{h}^{(t)}$

$U$

$$\mathbf{h}^{(t)} = tanh(\mathbf{U}\mathbf{x}^{(t)} + \mathbf{W}\mathbf{h}^{(t-1)})$$

$\mathbf{x}^{(t)}$

Language Technologies Institute

Carnegie Mellon University

# Recurrent Neural Networks - Unrolling

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

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

$$\mathbf{z}^{(t)} = matmult(\boldsymbol{h}^{(t)}, \boldsymbol{V})$$

$$\boldsymbol{h}^{(t)} = tanh(\boldsymbol{U}\boldsymbol{x}^{(t)} + \boldsymbol{W}\boldsymbol{h}^{(t-1)})$$

**Same model parameters are used for all time parts.**

Language Technologies Institute

Carnegie Mellon University

# Backpropagation Through Time

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

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

Gradient = "backprop" gradient
x "local" Jacobian

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

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

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

$L^{(\tau)}$

$y^{(\tau)}$  $\mathbf{z}^{(\tau)}$

$\mathbf{h}^{(\tau)}$

$\mathbf{x}^{(\tau)}$

Language Technologies Institute
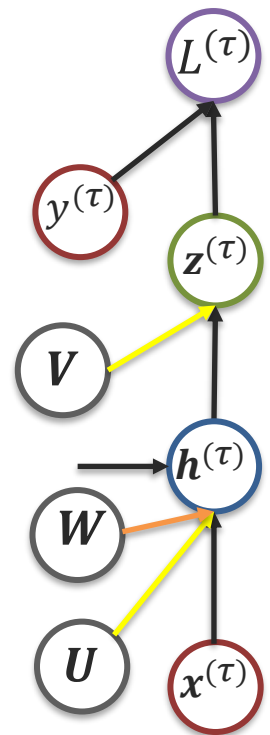
Carnegie Mellon University

# Backpropagation Through Time

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

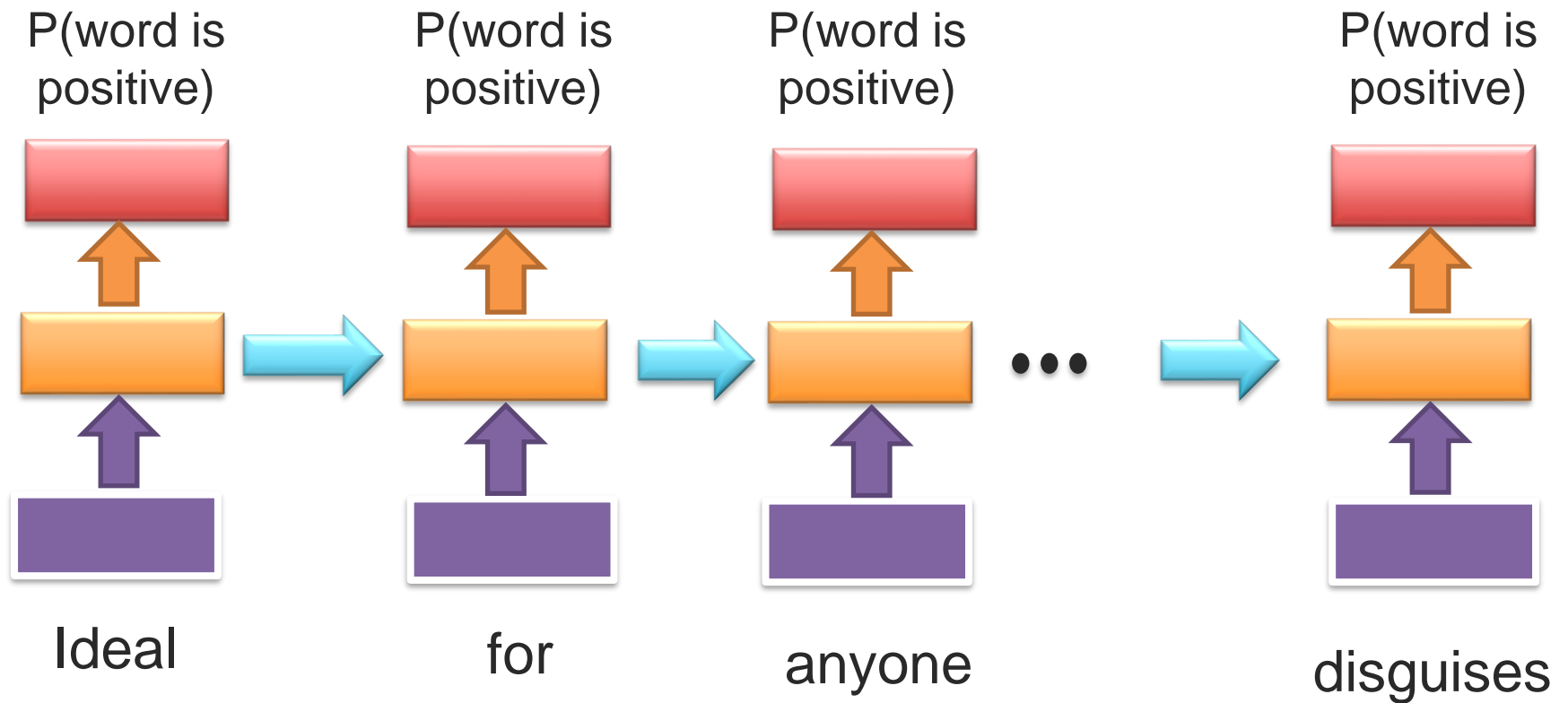> Gradient = "backprop" gradient
>
> x "local" Jacobian

$V$    $\nabla_{\mathbf{V}} L = \sum_t (\nabla_{\mathbf{z}^{(t)}} L) \dfrac{\partial \mathbf{z}^{(t)}}{\partial \mathbf{V}}$

$W$    $\nabla_{\mathbf{W}} L = \sum_t (\nabla_{\mathbf{h}^{(t)}} L) \dfrac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{W}}$

$U$    $\nabla_{\mathbf{U}} L = \sum_t (\nabla_{\mathbf{h}^{(t)}} L) \dfrac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{U}}$

Language Technologies Institute

Carnegie Mellon University

# RNN for Sequence Prediction

P(word is
positive)

P(word is
positive)

P(word is
positive)

P(word is
positive)

Ideal

for

anyone

disguises

What is the loss?

$$L = \frac{1}{N}\sum_t L^{(t)} = \frac{1}{N}\sum_t -logP(Y = y^{(t)}|z^{(t)})$$

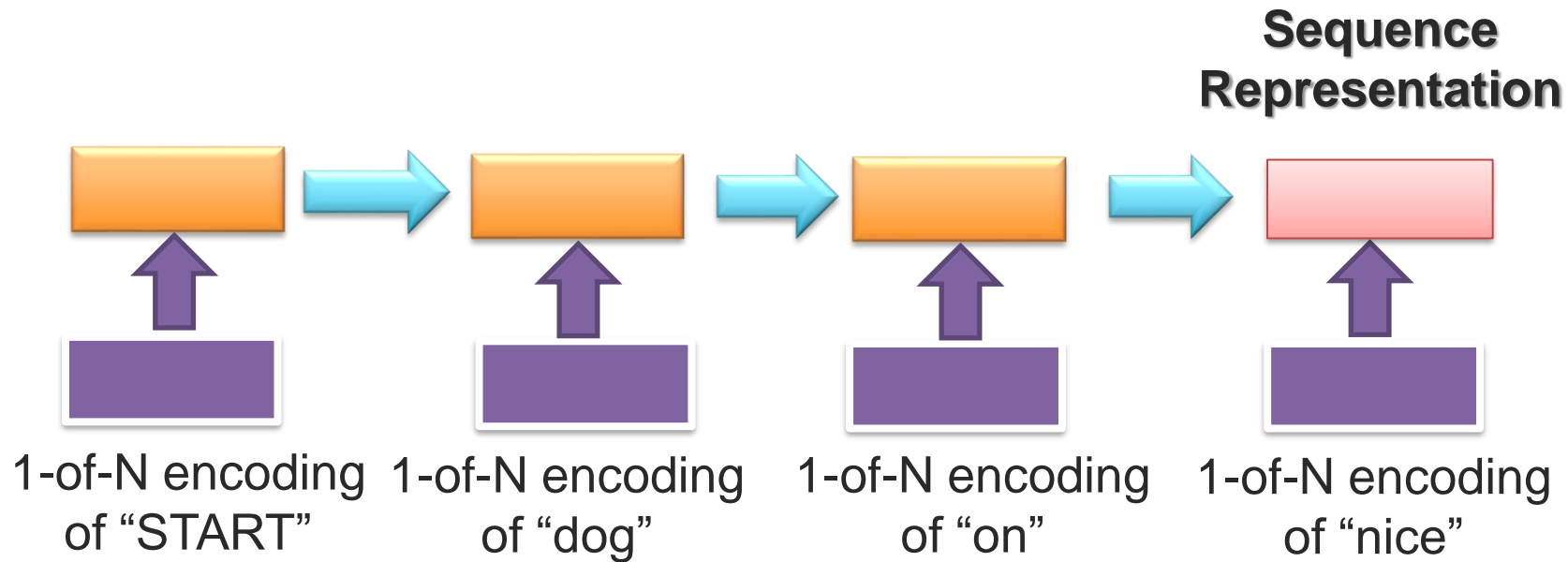# RNN for Sequence Prediction

P(sequence is positive)



Ideal       for       anyone       disguises

What is the loss?    $L = L^{(N)} = -log P(Y = y^{(N)} | z^{(N)})$
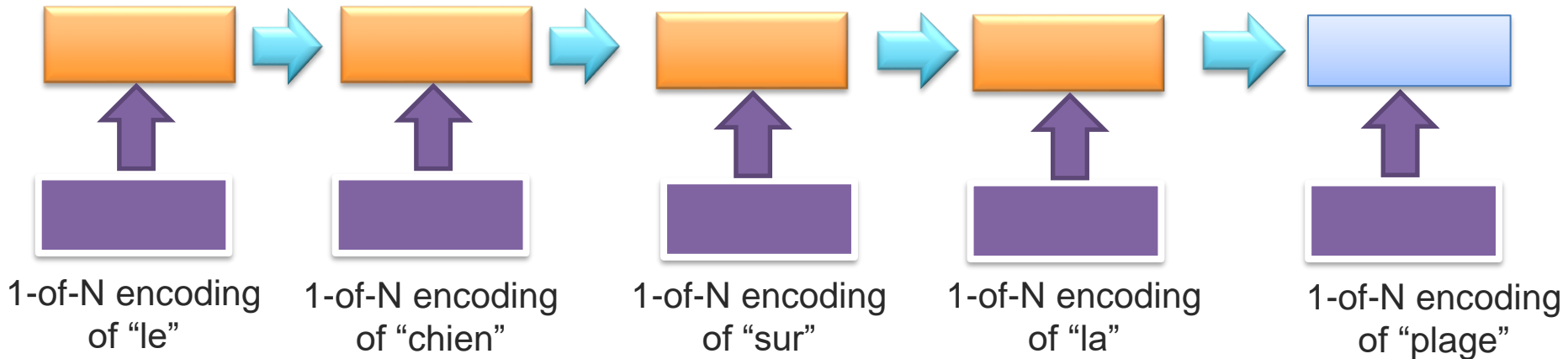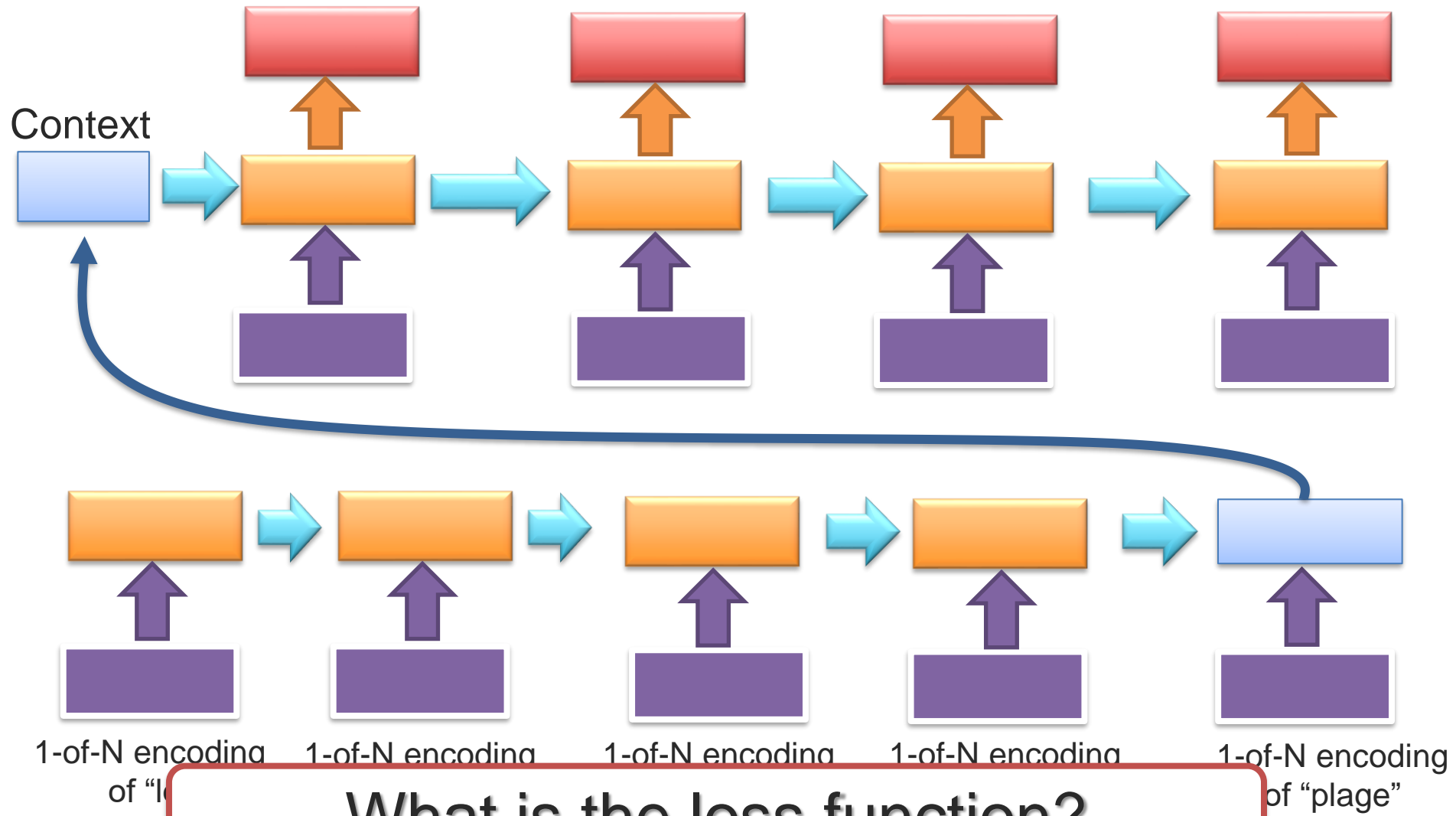
# RNN for Sequence Representation (Encoder)

**Sequence Representation**



1-of-N encoding of "START"   1-of-N encoding of "dog"   1-of-N encoding of "on"   1-of-N encoding of "nice"

# RNN-based for Machine Translation

Le chien sur la plage  ➡  The dog on the beach



1-of-N encoding of "le"  1-of-N encoding of "chien"  1-of-N encoding of "sur"  1-of-N encoding of "la"  1-of-N encoding of "plage"

# Encoder-Decoder Architecture

Context

1-of-N encoding
of "lo"

1-of-N encoding

1-of-N encoding

1-of-N encoding

1-of-N encoding
of "plage"

What is the loss function?

# Gated Recurrent Neural Networks

Carnegie Mellon University

# Long-term Dependencies

Vanishing gradient problem for RNNs:

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



➢ 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.

# Recurrent Neural Networks

$L^{(t)}$

$y^{(t)}$

$z^{(t)}$

$h^{(t)}$

$x^{(t)}$

$h^{(t)}$

tanh
+1
-1

$x^{(t)}$

$h^{(t+1)}$

Language Technologies Institute

Carnegie Mellon University

# LSTM ideas: (1) "Memory" Cell and Self Loop

[Hochreiter and Schmidhuber, 1997]

Long Short-Term Memory (LSTM)

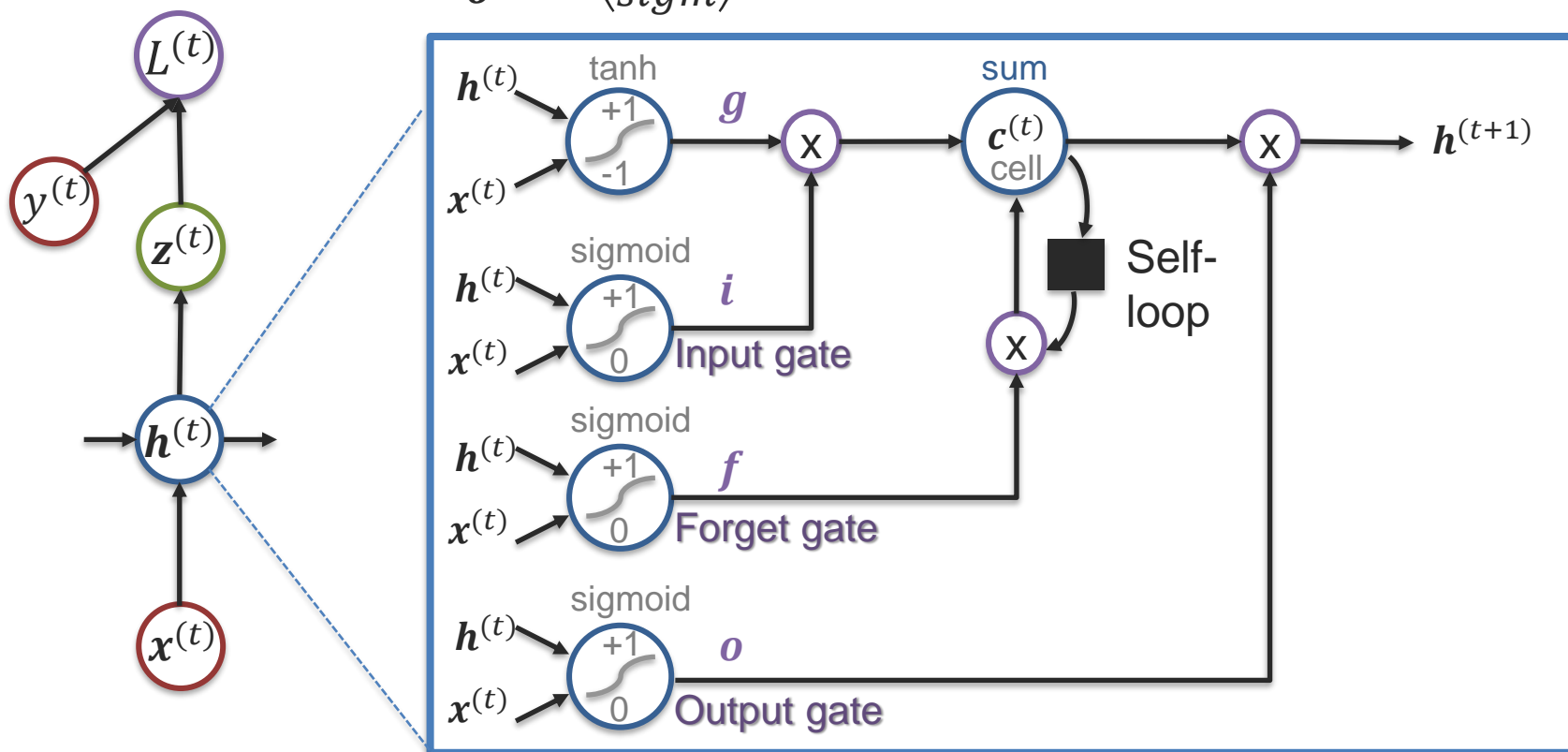# LSTM Ideas: (2) Input and Output Gates

[Hochreiter and Schmidhuber, 1997]

Language Technologies Institute

Carnegie Mellon University

# LSTM Ideas: (3) Forget Gate [Gers et al., 2000]

$$\begin{pmatrix} \boldsymbol{g} \\ \boldsymbol{i} \\ \boldsymbol{f} \\ \boldsymbol{o} \end{pmatrix} = \begin{pmatrix} tanh \\ sigm \\ sigm \\ sigm \end{pmatrix} W \begin{pmatrix} \boldsymbol{h}^{(t)} \\ \boldsymbol{x}^{(t)} \end{pmatrix} \qquad \boldsymbol{c}^{(t)} = \boldsymbol{f} \odot \boldsymbol{c}^{(t-1)} + \boldsymbol{i} \odot \boldsymbol{g}$$

$$\boldsymbol{h}^{(t)} = \boldsymbol{o} \odot \tanh(\boldsymbol{c}^{(t)})$$

Language Technologies Institute

Carnegie Mellon University

# Recurrent Neural Network using LSTM Units



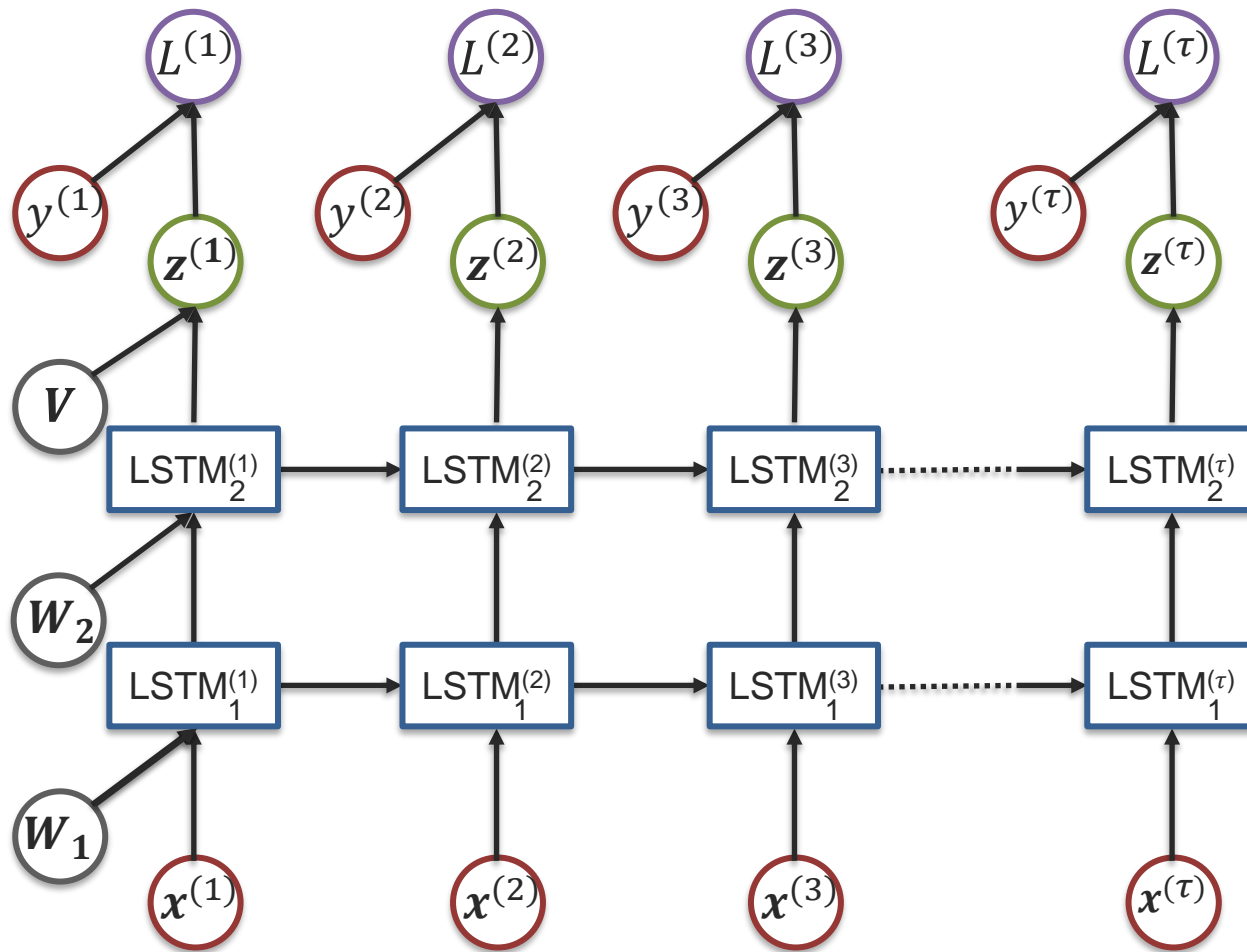Gradient can still be computer using backpropagation!

# Bi-directional LSTM Network



ELMO: Two bi-directional LSTMs are used to contextualize the word embeddings
https://allennlp.org/elmo
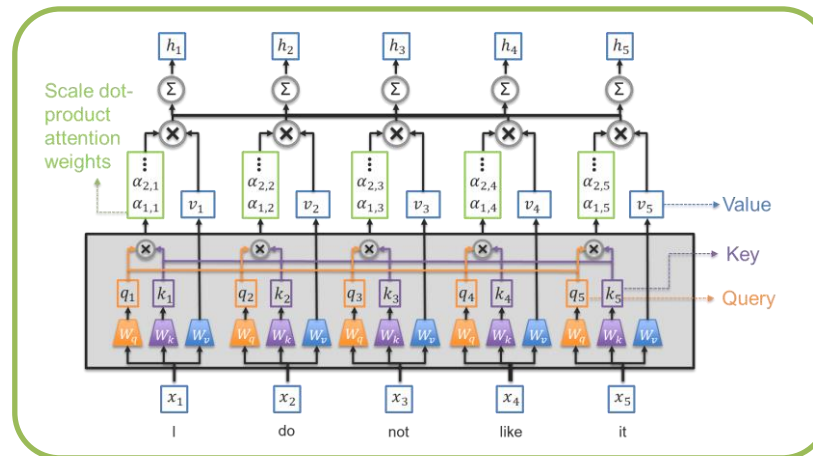
# Deep LSTM Network

Language Technologies Institute

Carnegie Mellon University

# And There Are More Ways To Model Sequences…



… in Week 5!

## Self-attention Models
(e.g., BERT, RoBERTa)

Language Technologies Institute

Carnegie Mellon University

# Syntax and Language Structure

# Syntax and Language Structure

What can you tell about this sentence?



Phrase-structure Grammar

② Syntactic parse tree

① Part-of-speech tags

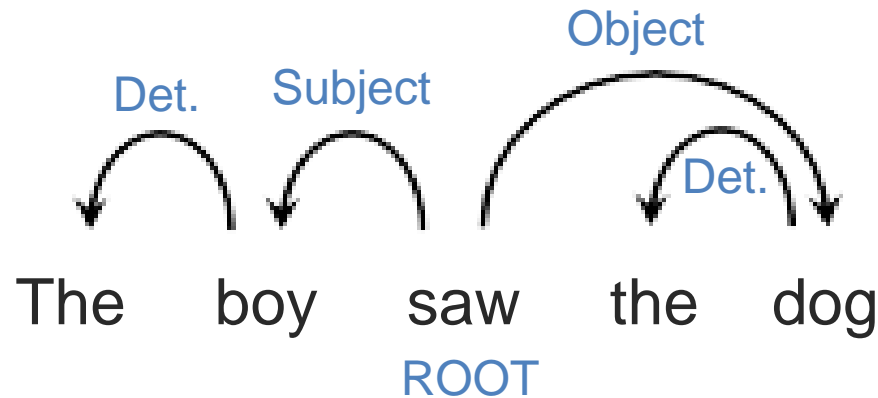# Syntax and Language Structure

What can you tell about this sentence?

# Dependency Grammar

**Main idea:** Syntactic structure consists of *lexical items*, linked by binary asymmetric relations called *dependencies*

➢ Easier to convert to predicate-argument structure

➢ You can try to convert one representation into another

❑ But, in general, these formalisms are not equivalent



The boy saw the dog

# Ambiguity in Syntactic Parsing

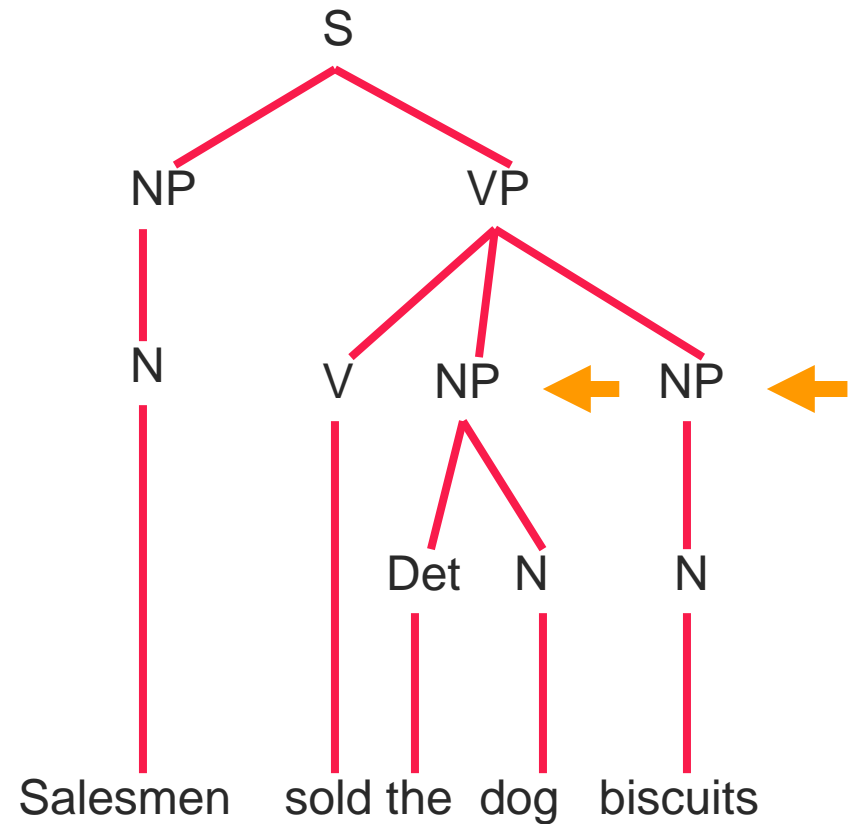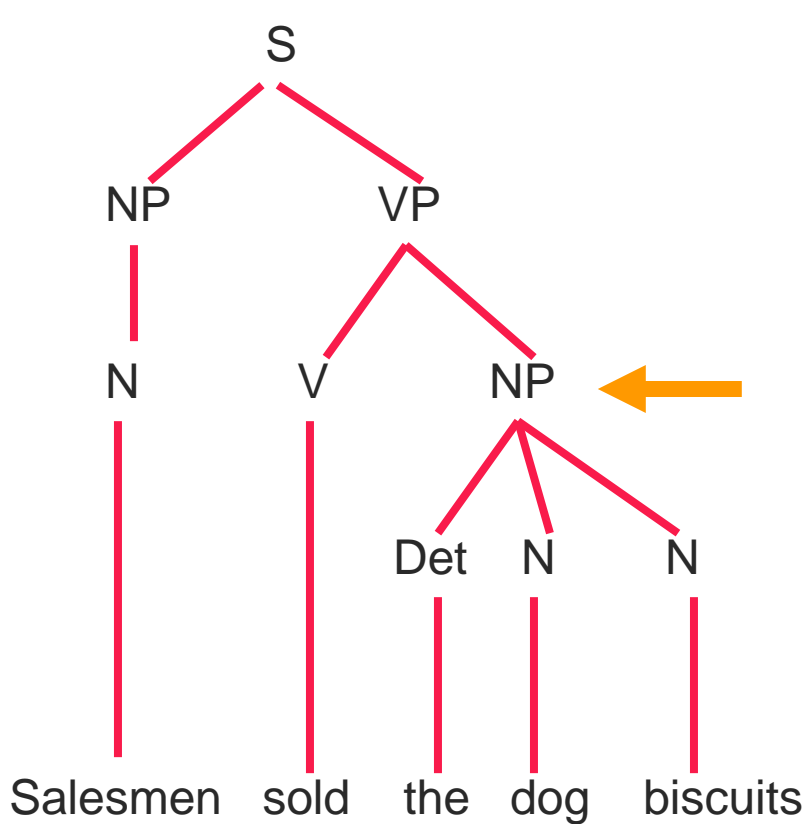**"Like" can be a verb or a preposition**

- I like/VBP candy.
- Time flies like/IN an arrow.

**"Around" can be a preposition, particle, or adverb**

- I bought it at the shop around/IN the corner.
- I never got around/RP to getting a car.
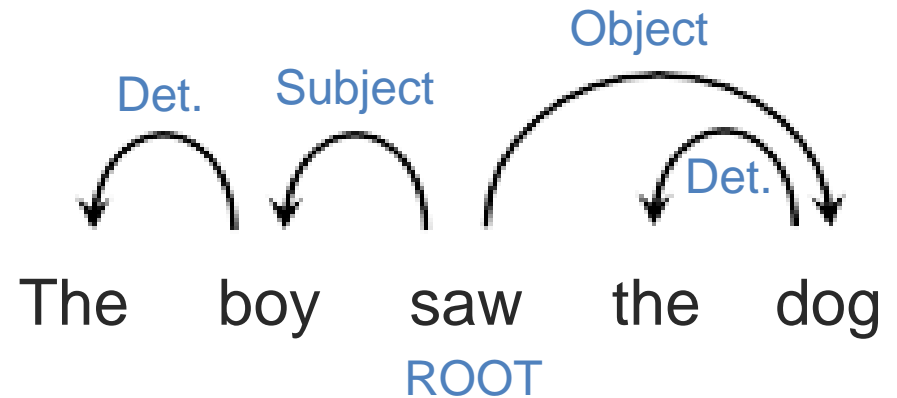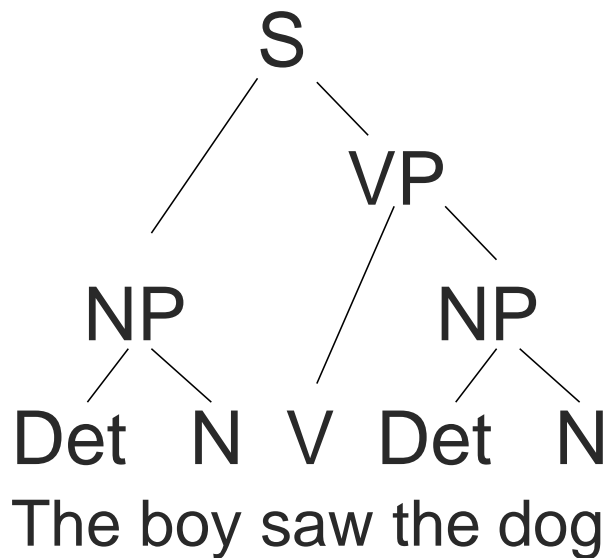- A new Prius costs around/RB $25K.

# Language Ambiguity

# Language Syntax – Examples

| Det | Noun | Verb | Det | Noun | Prep | Det | Noun |
|-----|------|------|-----|------|------|-----|------|
| The | boy  | saw  | the | dog  | in   | the | park |

**Part of Speech tagging**

```
              S
             / \
              VP
             /  \
   NP          NP
  /  \        /  \
Det  N  V  Det  N
The boy saw the dog
```

**Constituency Parsing**

Det.   Subject   Object

Det.

The    boy    saw    the    dog

ROOT

**Dependency Parsing**

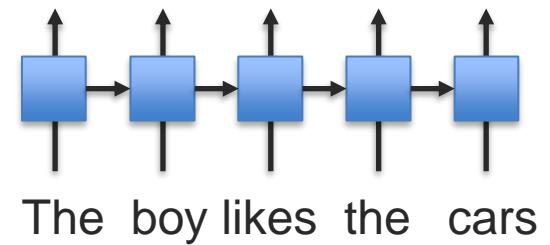How to take advantage of syntax when modeling language with neural networks?

Language

University
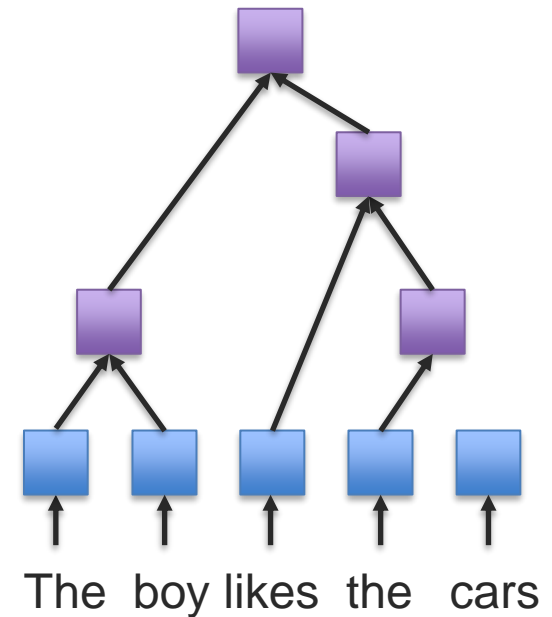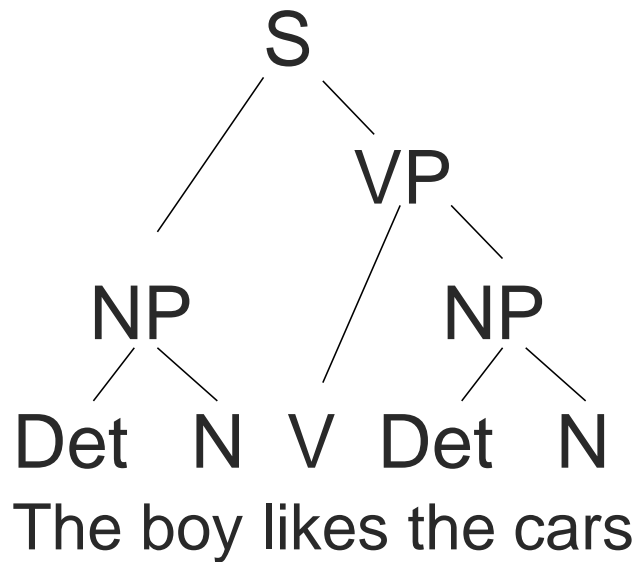
# Recursive Neural Network

# How to Model Syntax with RNNs?

S

VP

NP NP

Det N V Det N

The boy likes the cars

**?**


The boy likes the cars

We could use Part-of-Speech tags.
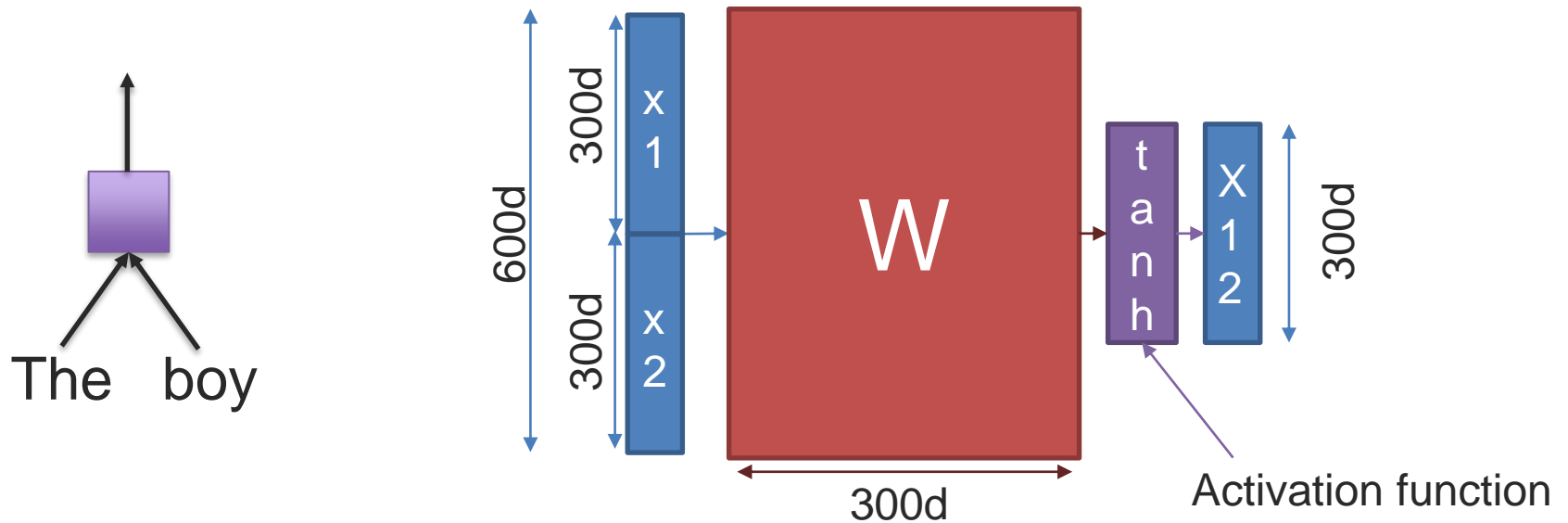
Language Technologies Institute

Carnegie Mellon University

# Tree-based RNNs (or Recursive Neural Network)

# Recursive Neural Unit

➡️ Pair-wise combination of two input features



Activation function
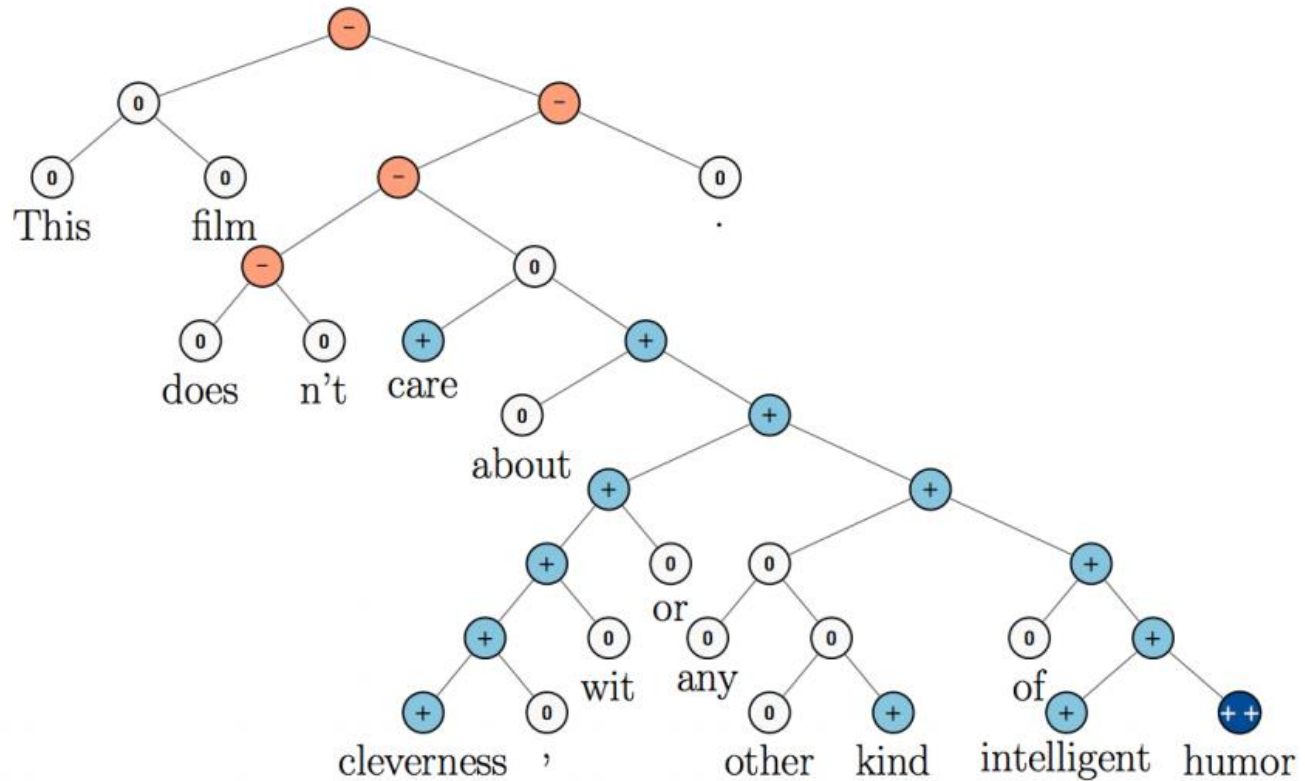
Language Technologies Institute

Carnegie Mellon University

# Recursive Neural Network for Sentiment Analysis



$p_2 = g(a, p_1)$

$p_1 = g(b, c)$

... not     very     good ...

a          b         c

Socher et al., Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank, EMNLP 2013

Language Technologies Institute

Carnegie Mellon University

# Recursive Neural Network for Sentiment Analysis

Classification of a sentence using tree-based compositionality of words



Demo: http://nlp.stanford.edu/sentiment/

Socher et al., Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank, EMNLP 2013

Language Technologies Institute

Carnegie Mellon University

# Stack LSTM



stack of partially constructed dependency subtrees

buffer of words remaining to be processed

stack representing the history of actions taken by the parser
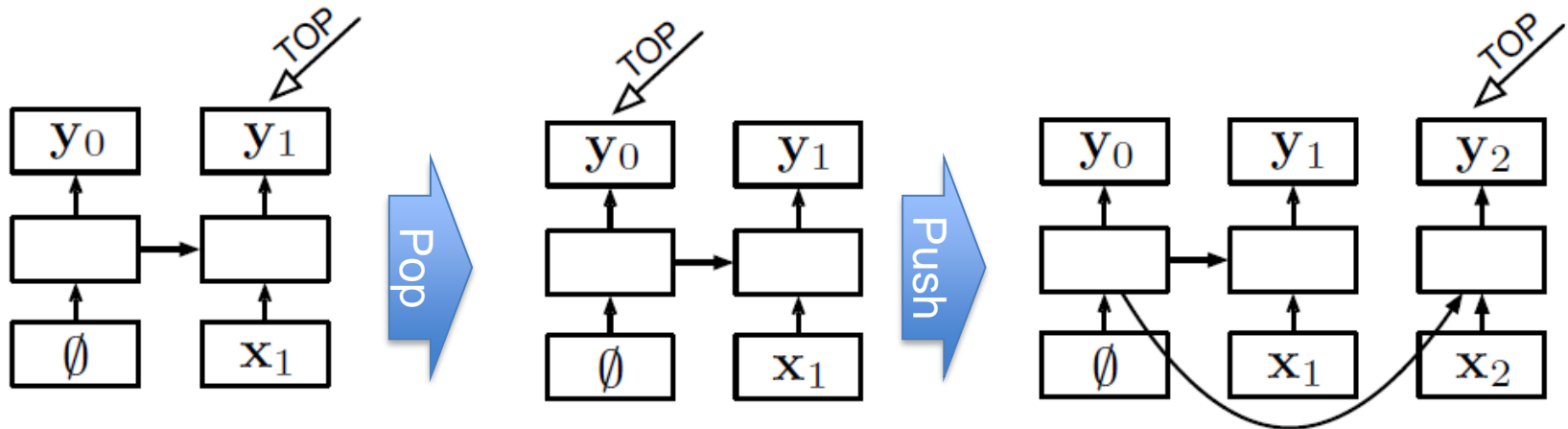
Dyer et al., Transition-Based Dependency Parsing with Stack Long Short-Term Memory, 2015

# Stack LSTM



Dyer et al., Transition-Based Dependency Parsing with Stack Long Short-Term Memory, 2015

# Resources

- Stanford NLP software

https://nlp.stanford.edu/software/

  - Stanford Parser
  - Stanford POS Tagger

- UC Berkeley Parser

https://github.com/slavpetrov/berkeleyparser

- Parsers by Kenji Sagae (syntactic parsers)
  http://www.sagae.org/software.html