

CS310 Natural Language Processing

自然语言处理

Lecture 14 - Course Review

Instructor: Yang Xu

主讲人：徐炆

xuyang@sustech.edu.cn

A Historic View of NLP

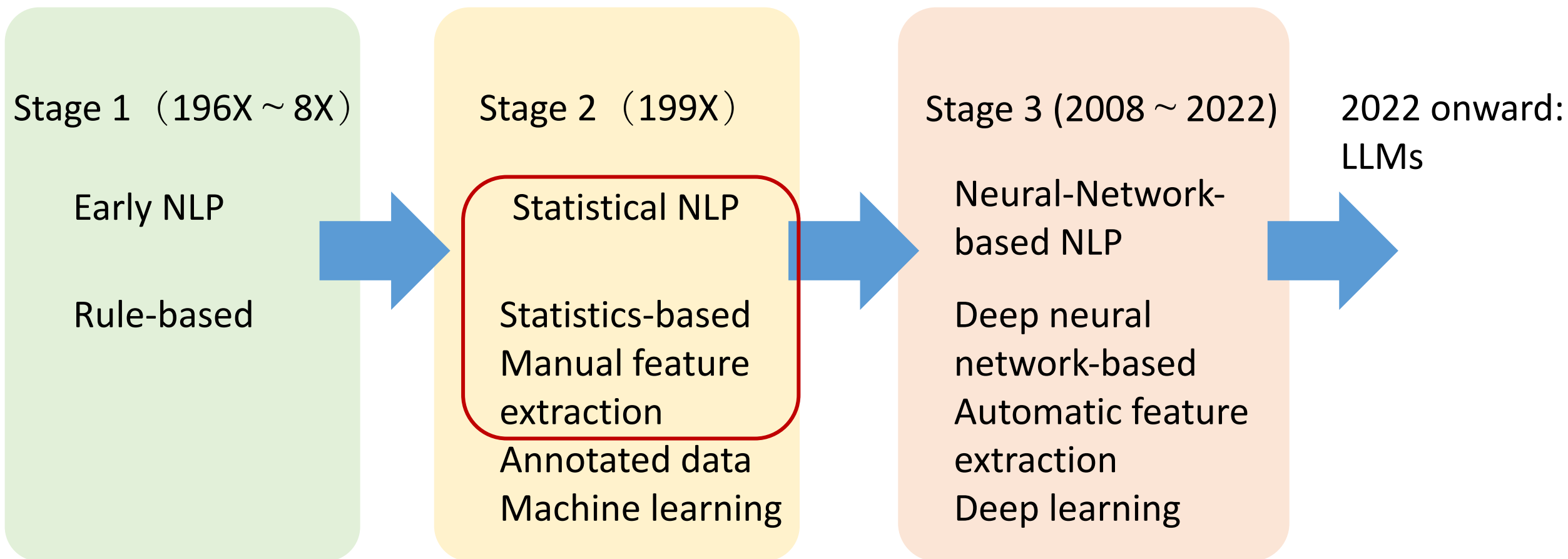


Table of Content Lecture 02 - Word Vectors

- Motivation
- Documents and Counts-based Method
- Neural Network-based Method -- word2vec
- Evaluation and Applications

Build Word-Document Matrix (term-document matrix)^[1]

- Build matrix $\mathbf{A} \in \mathbb{R}^{V \times C}$, which contains the count of each word in each document

- Example:**

x_1 : 学 而 时 习 之

x_2 : 学 而 不 思 则 罔

x_3 : 思 而 不 学 则 殆

Entry $\mathbf{A}_{v,c} = \text{count}_{x_c}(v)$, count of word v in the c th document

		C		
		x_1	x_2	x_3
V	学	1	1	1
	而	1	1	1
	不	0	1	1
	思	0	1	1
	则	0	1	1
	时	1	0	0
	习	1	0	0
	之	1	0	0
	罔	0	1	0
	殆	0	0	1

[1] https://en.wikipedia.org/wiki/Term-document_matrix

Pointwise Mutual Information

$$PMI = [\mathbf{A}]_{v,c} = \left[\log \frac{\text{count}_{x_c}(v)}{\frac{\text{count}_x(v)}{N} \cdot \ell_c} \right]_+$$

- If a word v has nearly same frequency in every document, then its row $[\mathbf{A}]_{v,*}$ will be nearly all zeros
- If a word v only occurs in one document c , then its PMI will be large and positive
- Thus, PMI is sensitive to rare words; usually need to smooth the frequencies by filtering rare words

	x_1	x_2	x_3
学	1	1	1
而	1	1	1
不	0	1	1
思	0	1	1
则	0	1	1
时	1	0	0
习	1	0	0
之	1	0	0
罔	0	1	0
殆	0	0	1

Improvement: Latent Semantic Analysis

(Deerwester et al., 1990)

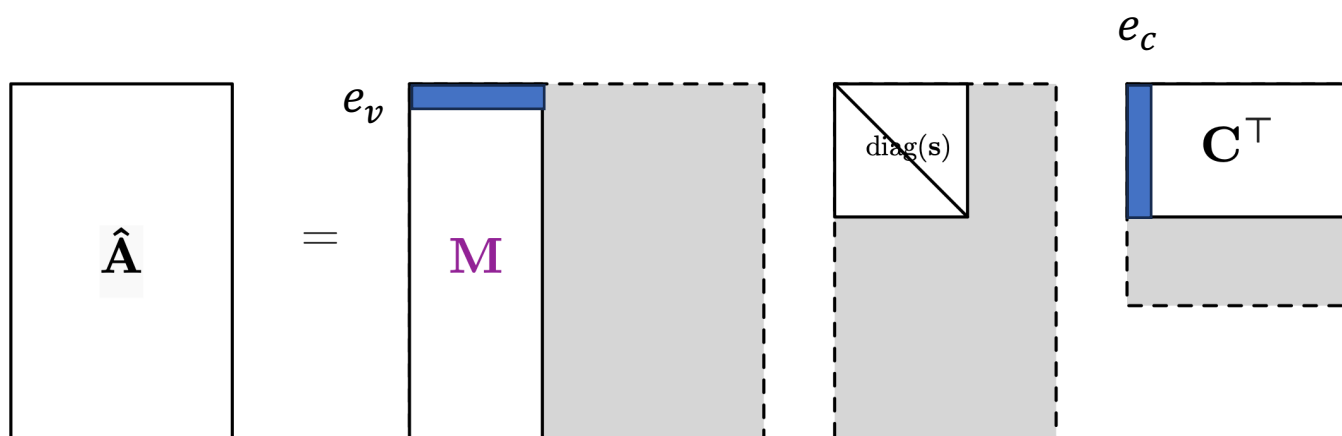
- LSA seeks to find a more compact (low rank) representation of document-word matrix \mathbf{A}

$$\underset{V \times C}{\mathbf{A}} \approx \underset{V \times d}{\hat{\mathbf{A}}} = \underset{V \times d}{\mathbf{M}} \times \underset{d \times d}{\text{diag}(\mathbf{s})} \times \underset{d \times C}{\mathbf{C}^T}$$

- Can be solved by applying singular value decomposition to \mathbf{A} , and then truncating to d dimensions ($\hat{\mathbf{A}}$)
- \mathbf{M} contains left singular vectors of \mathbf{A}
- \mathbf{C} contains right singular vectors of \mathbf{A}
- \mathbf{s} are singular values of \mathbf{A}

Truncated SVD => word vectors

$$\mathbf{A} \approx \hat{\mathbf{A}} = \mathbf{M} \times \text{diag}(\mathbf{s}) \times \mathbf{C}^\top$$

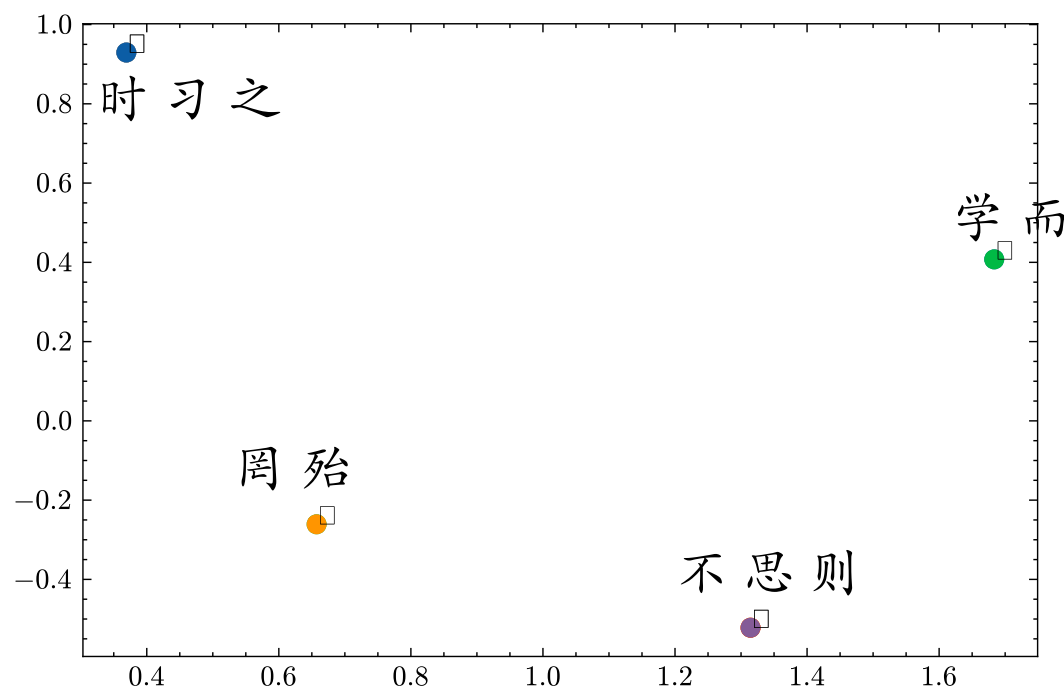


- v th column in \mathbf{M} is the embedding vector for word v
- c th column in \mathbf{C} is the embedding vector for document c

- \mathbf{M} contains useful word vectors (“embeddings”) of d dimensions
- \mathbf{C} contains document vectors

LSA Example $d = 2$

- Word vectors \mathbf{M} plotted
- Note that some words are in the same spot. Why?



	x_1	x_2	x_3
学	1	1	1
而	1	1	1
不	0	1	1
思	0	1	1
则	0	1	1
时	1	0	0
习	1	0	0
之	1	0	0
罔	0	1	0
殆	0	0	1

$\mathbf{M} =$

$\mathbf{A} =$

	x_1	x_2	x_3
学	1	1	1
而	1	1	1
不	0	1	1
思	0	1	1
则	0	1	1
时	1	0	0
习	1	0	0
之	1	0	0
罔	0	1	0
殆	0	0	1

Overview Lecture 12 - Question Answering

- **Question Answering (QA)**
 - What is QA?
 - Information Retrieval; Tf-idf
 - Retriever-based QA; Datasets
 - Answer Span Extraction
 - Retrieval-Augmented Generation

Brief Overview of Information Retrieval (IR)

- **Information retrieval, IR:** Retrieval of all kinds of media based on user information needs. IR system \approx **search engine**
- We focus on **ad hoc (临时) retrieval**: a user poses a query to an IR system, which then returns an ordered set of documents from some collection

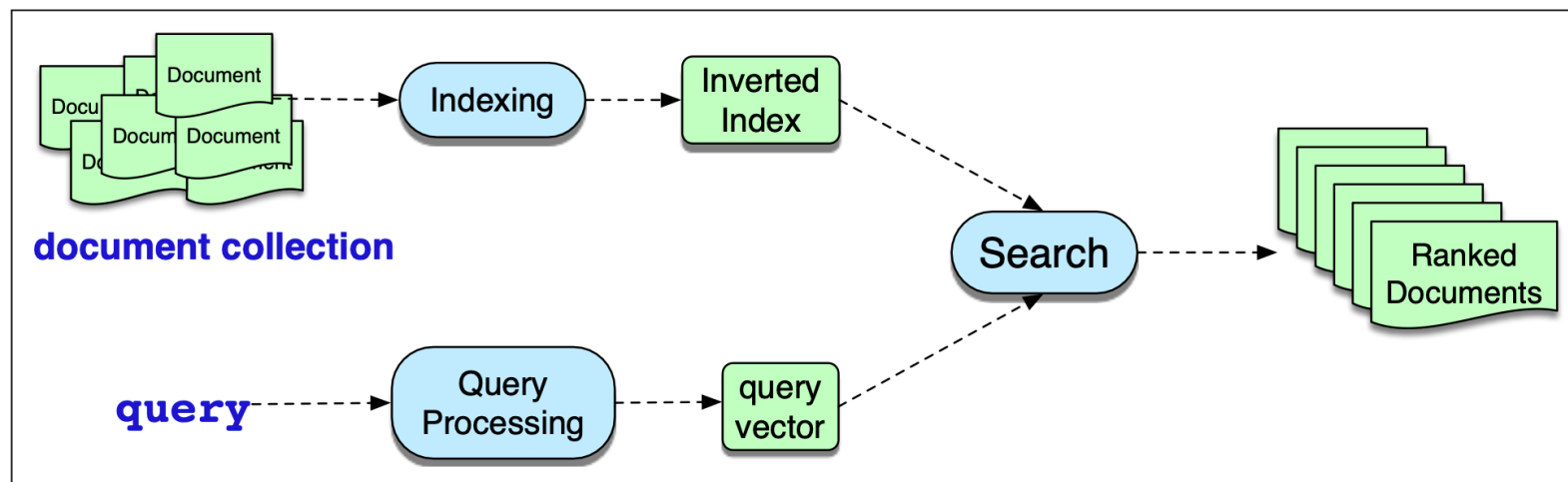


Figure 14.1 The architecture of an ad hoc IR system.

Query: a user's information need expressed as a set of **terms**

Term refers to a word/phrase in a collection of documents

Figure from SLP3, Ch 14

How to match a document a query?

- Compute a term weight for each document term
- Common method: **tf-idf** and BM25
 - **tf**: term frequency
 - **idf**: inverse document frequency
- $\text{tf-idf} \triangleq \text{tf} \times \text{idf}$ (product of the two)

term t ; document d

$$\text{tf}_{t,d} = \begin{cases} 1 + \log_{10} \text{count}(t, d) & \text{if } \text{count}(t, d) > 0 \\ 0 & \text{otherwise} \end{cases}$$

- **tf**: words that occur more often in a document are likely to be informative about the document's content
- Use the \log_{10} of word frequency count rather than raw count
- Why? A word appearing 100 times doesn't make it 100 times more likely

Tf-idf

$$\text{tf}_{t,d} = \begin{cases} 1 + \log_{10} \text{count}(t,d) & \text{if } \text{count}(t,d) > 0 \\ 0 & \text{otherwise} \end{cases}$$

term t ; document d

term occurs 0 times in document: $\text{tf} = 0$
term occurs 1 times in document: $\text{tf} = 1$
term occurs 10 times in document: $\text{tf} = 2, \dots$

- **document frequency** df_t of a term t is the number of documents it occurs in
- Terms that occur in only **a few** documents are useful for discriminating those documents from the rest of the collection;
- terms that occur across the entire collection aren't as helpful (*the, a, an, ...*)
- **inverse document frequency** or **idf** is defined as:

$$\text{idf}_t = \log_{10} \frac{N}{\text{df}_t}$$

N : total number of documents
The fewer documents in which t occurs, the higher idf_t

Inverse document frequency example

- Some idf values for some words in the corpus of Shakespeare plays

Word	df	idf
Romeo	1	1.57
salad	2	1.27
Falstaff	4	0.967
forest	12	0.489
battle	21	0.246
wit	34	0.037
fool	36	0.012
good	37	0
sweet	37	0

Extremely informative words that occur in only one play like *Romeo*

good or *sweet* are completely non-discriminative since they occur in all 37 plays

Scoring with tf-idf

- We can score document d by the cosine of its vector \vec{d} with the query vector \vec{q} :

$$\text{score}(q, d) = \cos(\vec{q}, \vec{d}) = \frac{\vec{q} \cdot \vec{d}}{|\vec{q}| \cdot |\vec{d}|}$$

- in which \vec{q} and \vec{d} are vectors of query length n , whose values are the **tf-idf** values (normalized):

$$\begin{aligned} \vec{q} &= \frac{[\text{tf-idf}(t_1, q), \dots, \text{tf-idf}(t_n, q)]}{\sqrt{\sum_{t \in q} \text{tf-idf}^2(t, q)}} \\ \vec{d} &= \frac{[\text{tf-idf}(t_1, d), \dots, \text{tf-idf}(t_n, d)]}{\sqrt{\sum_{t \in d} \text{tf-idf}^2(t, d)}} \end{aligned} \quad \Rightarrow \quad \text{score}(q, d) = \sum_{t_i \in q} \frac{\text{tf-idf}(t_i, q)}{\sqrt{\sum_{t \in q} \text{tf-idf}^2(t, q)}} \cdot \frac{\text{tf-idf}(t_i, d)}{\sqrt{\sum_{t \in d} \text{tf-idf}^2(t, d)}}$$

Table of Content

- **Language Modeling**
- Neural Language Models
- Recurrent Neural Networks for LM
- Evaluate LMs

Lecture 03 - Recurrent Neural Networks and Language Modeling

How to Learn an LM?

- Pre- Neural network solution: ***n*-gram Language Model**
- **Def.** An *n*-gram is a chunk of *n* consecutive words

the Sun rises every _____

- **Unigrams** ($n=1$): “the”, “Sun”, “rises”, “every”
 - **Bigrams** ($n=2$): “the Sun”, “Sun rises”, “rises every”
 - **Trigrams** ($n=3$): “the Sun rises”, “Sun rises every”
 - **Four-grams** ($n=4$): “the Sun rises every”
-
- **Idea:** Count the frequencies of different *n*-grams and use these to predict the next word



Andrey Andreyevich Markov
(14 June 1856 – 20 July 1922)

n -grams LM: Markov assumption

- **Markov assumption**: a word at only depends on its preceding $n - 1$ words

$$P(x^{(t+1)} | x^{(1)}, \dots, x^{(t)}) = P(x^{(t+1)} | \underbrace{x^{(t-n+2)}, \dots, x^{(t)}}_{n-1 \text{ words}})$$

Probability of a n -gram

Probability of a $(n-1)$ -gram

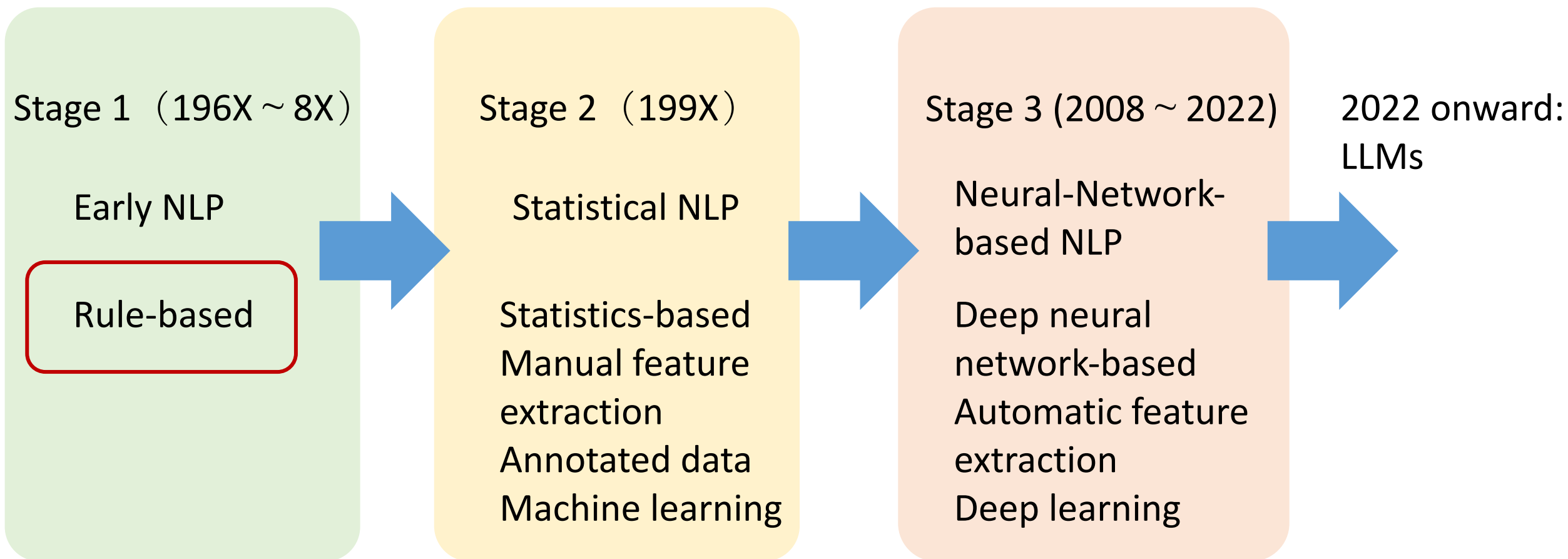
$$= \frac{P(x^{(t-n+2)}, \dots, x^{(t)}, x^{(t+1)})}{P(x^{(t-n+2)}, \dots, x^{(t)})}$$

By the
definition of
conditional
probability

- **Question**: How to obtain the probabilities?
- **Answer**: By counting them from some large enough corpora (statistical approximation)

$$\approx \frac{\text{count}(x^{(t-n+2)}, \dots, x^{(t)}, x^{(t+1)})}{\text{count}(x^{(t-n+2)}, \dots, x^{(t)})}$$

A Historic View of NLP



Overview

- Context Free Grammars (CFG)
- Constituency Parsing
- Neural Constituency Parsing

Lecture 05 - Context-Free Grammars and Constituency Parsing

Context-Free Grammars (CFG)

- Also called Phrase-Structure Grammar
- Wilhelm Wundt (1900), Chomsky (1956), Backus (1959), Naur (1963)
- A set of **rules** or **productions** that express the ways symbols can be grouped and ordered together.

Wundt, German psychologist
Chomsky, American linguist
Backus & Naur, American
computer scientists (FORTRAN,
ALGO)

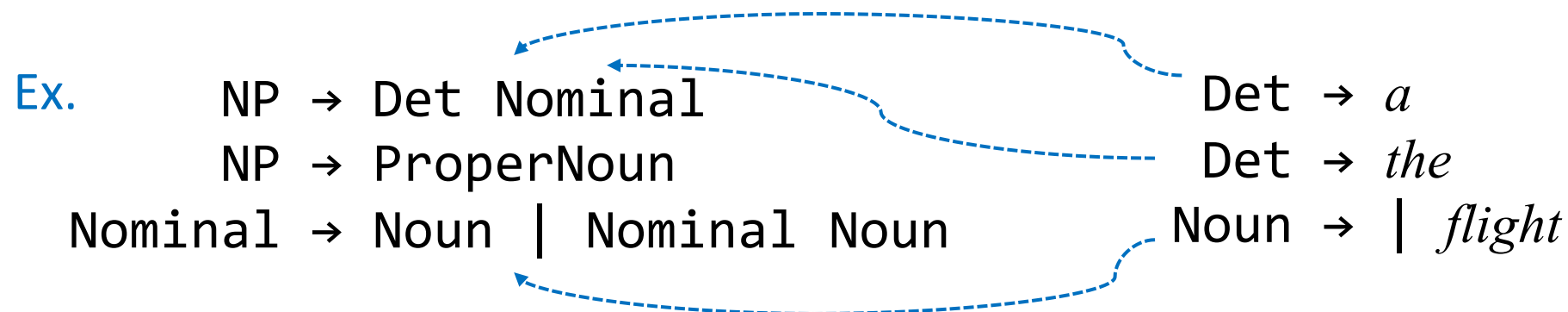
Ex. NP \rightarrow Det Nominal
 NP \rightarrow ProperNoun
 Nominal \rightarrow Noun | Nominal Noun

A **NP (noun phrase)** can be composed of

- either a *ProperNoun* (专有名词)
- or a determiner (*Det* 冠词)
 followed by a *Nominal* (名词性)

A *Nominal* can be one or more *Nouns* (名词)

Rules can be hierarchically embedded



Two classes of symbols: **terminal** and **non-terminal**

- **terminal**: symbols that are actual words in language (“a”, “the”, ...)
- **non-terminal**: symbols that express abstractions over terminals

CKY Recognition

- For a sentence of length n , work with an $(n + 1) \times (n + 1)$ **matrix** m , using the upper triangular portion \Rightarrow **parse table**
- Cell $m[i, j]$ contains the set of non-terminals representing all constituents that span position i through j .

Index: 0 *Book* 1 *the* 2 *flight* 3 *through* 4 *Houston* 5

Ex. Cell $m[0,3]$ contains the set of non-terminals for all constituents spanning 0 to 3

Task: Fill in the parse table correctly in a bottom-up way

		Column index					
		<i>Book</i>	<i>the</i>	<i>flight</i>	<i>through</i>	<i>Houston</i>	
		0	1	2	3	4	5
Row index	0		[0,1]				
	1						
	2						
	3						
	4						
	5						

CKY Recognition: Task breakdown

- Because grammar is in **CNF** \Rightarrow each non-terminal entry has exactly **two children**.
- For each constituent represented by entry $[i, j]$, there must be a position in the input, k , where it can be split into two parts ($i < k < j$)
- The first component $[i, k]$ is to the left of $[i, j]$, at the same row i
- The second component $[k, j]$ is beneath $[i, j]$, at the same column j

- Then, $[i, j]$ is the **combination** of $[i, k]$ and $[k, j]$ for all possible k

The smallest problems are $[j - 1, j]$, i.e., spans of single word

In order to solve the **problem**, we need to solve **smaller problems** first

- Very typical dynamic programming approach - *optimal substructure*

CKY Recognition: Example

Index: 0 *Book* 1 *the* 2 *flight* 3

Solve span [0, 1]

	0	1	2	3
0		S, VP, Verb, Nominal, Noun		
1				
2				
3				

5 production rules match with "**Book**"

Store them in **cell [0,1]**

$m[0,1]$
 $= \{S, \text{Nominal}, VP, \text{Noun}, \text{Verb}\}$

\mathcal{L}_1 in CNF

$S \rightarrow NP VP$

$S \rightarrow X1 VP$

$X1 \rightarrow Aux NP$

$S \rightarrow book \mid include \mid prefer$

$S \rightarrow Verb NP$

$S \rightarrow X2 PP$

$S \rightarrow Verb PP$

$S \rightarrow VP PP$

$NP \rightarrow I \mid she \mid me$

$NP \rightarrow TWA \mid Houston$

$NP \rightarrow Det Nominal$

$Nominal \rightarrow book \mid flight \mid meal \mid money$

$Nominal \rightarrow Nominal Noun$

$Nominal \rightarrow Nominal PP$

$VP \rightarrow book \mid include \mid prefer$

$VP \rightarrow Verb NP$

$VP \rightarrow X2 PP$

$X2 \rightarrow Verb NP$

$VP \rightarrow Verb PP$

$VP \rightarrow VP PP$

$PP \rightarrow Preposition NP$

$Det \rightarrow that \mid this \mid the \mid a$

$Noun \rightarrow book \mid flight \mid meal \mid money$

$Verb \rightarrow book \mid include \mid prefer$

Overview

Lecture 06 - Dependency Parsing

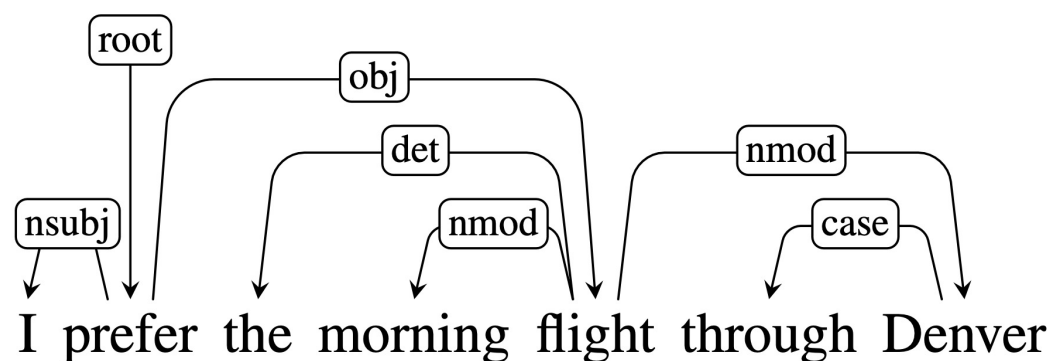
- Dependency Grammars
- Transition-Based Dependency Parsing
- Graph-Based Dependency Parsing
- Evaluation

Dependency Grammars

- Different from context-free grammars and constituency-based representations
- **Dependency Grammars** describe syntactic structure of a sentence solely in terms of directed grammatical *relations between words*

arc: *n.* 弧线

Labeled arcs from **heads** to **dependents**



$prefer \xrightarrow{\text{nsubj}} I$

prefer is the head of *I*

$prefer \xrightarrow{\text{obj}} flight$

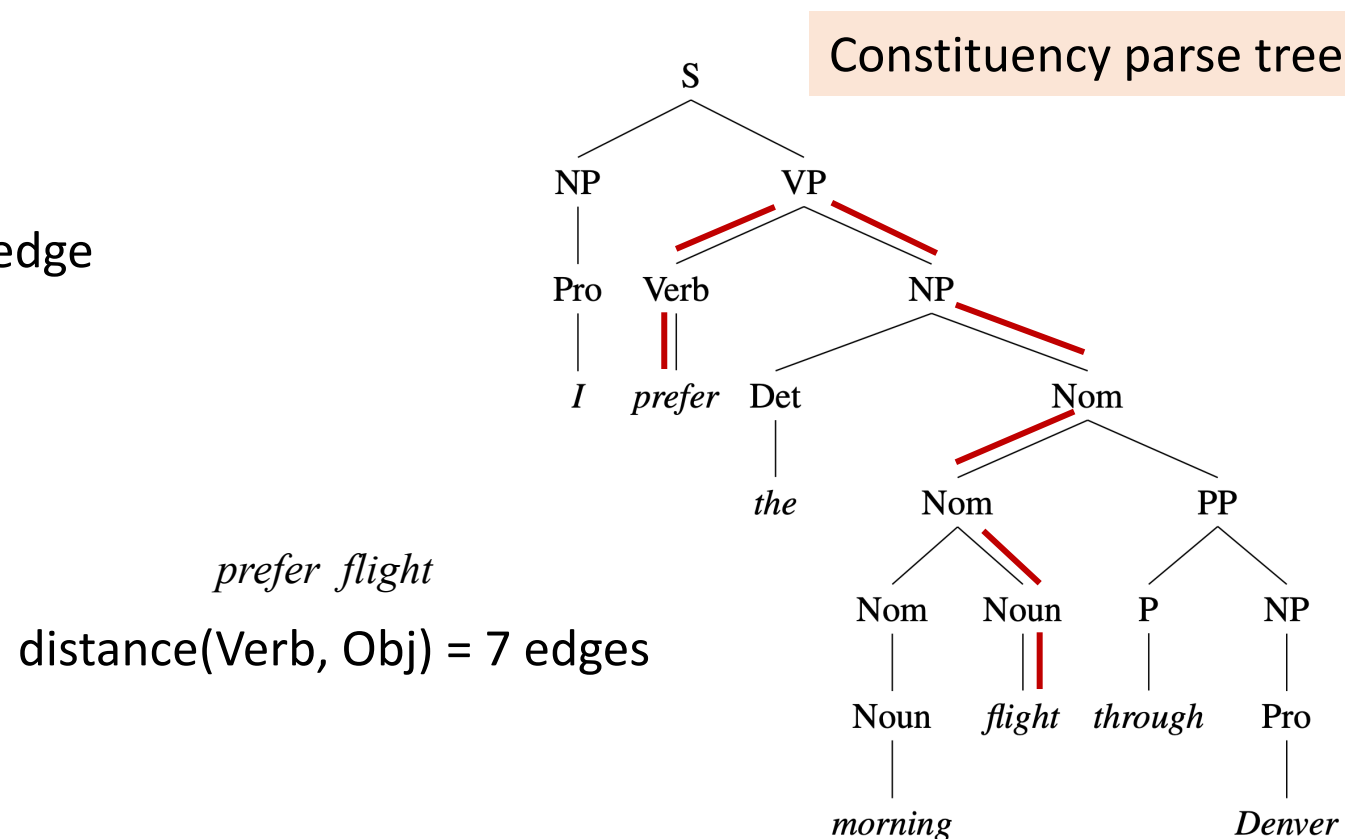
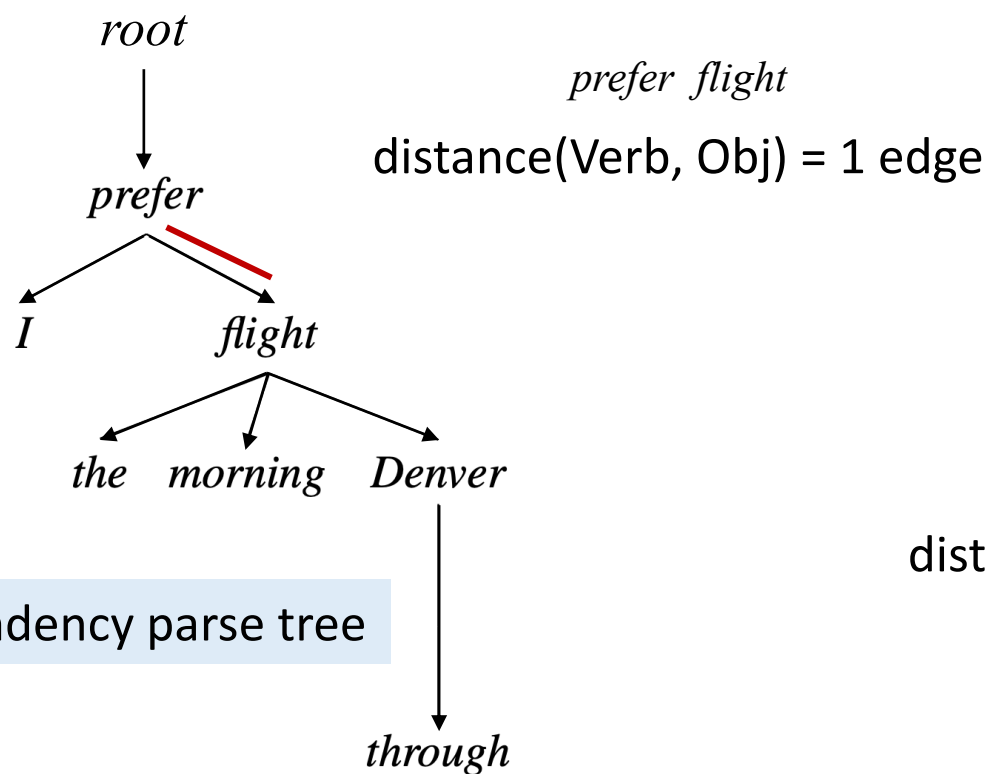
flight is the dependent of *prefer*

$root \rightarrow prefer$

root is the head of *prefer* and also the head of the entire structure (root of the tree)

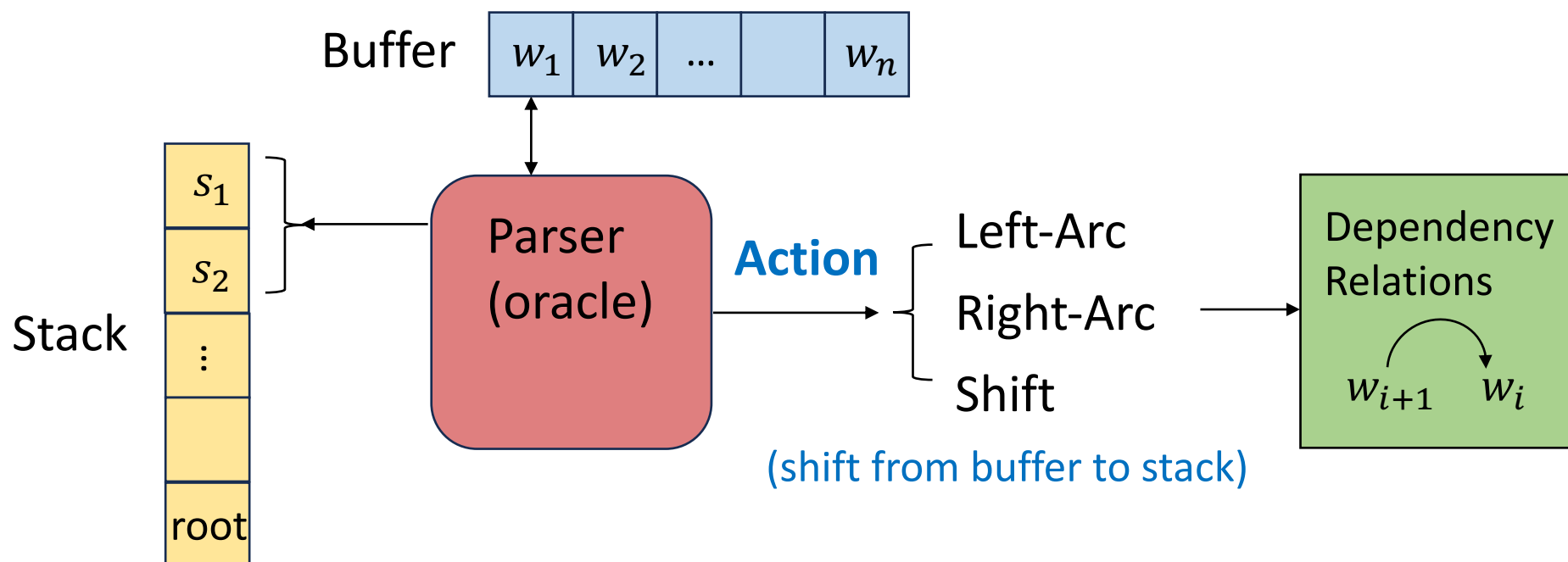
Compare to Context-Free Grammars

- Head-dependent relations **directly** encode important information that is often **buried** in the more complex constituency parses (by CFG)



Transition-Based Dependency Parsing

- An architecture that draws on shift-reduce parsing (a paradigm for analyzing programming languages)
- Key components: **stack**, **buffer**, and **oracle**.



A Historic View of NLP

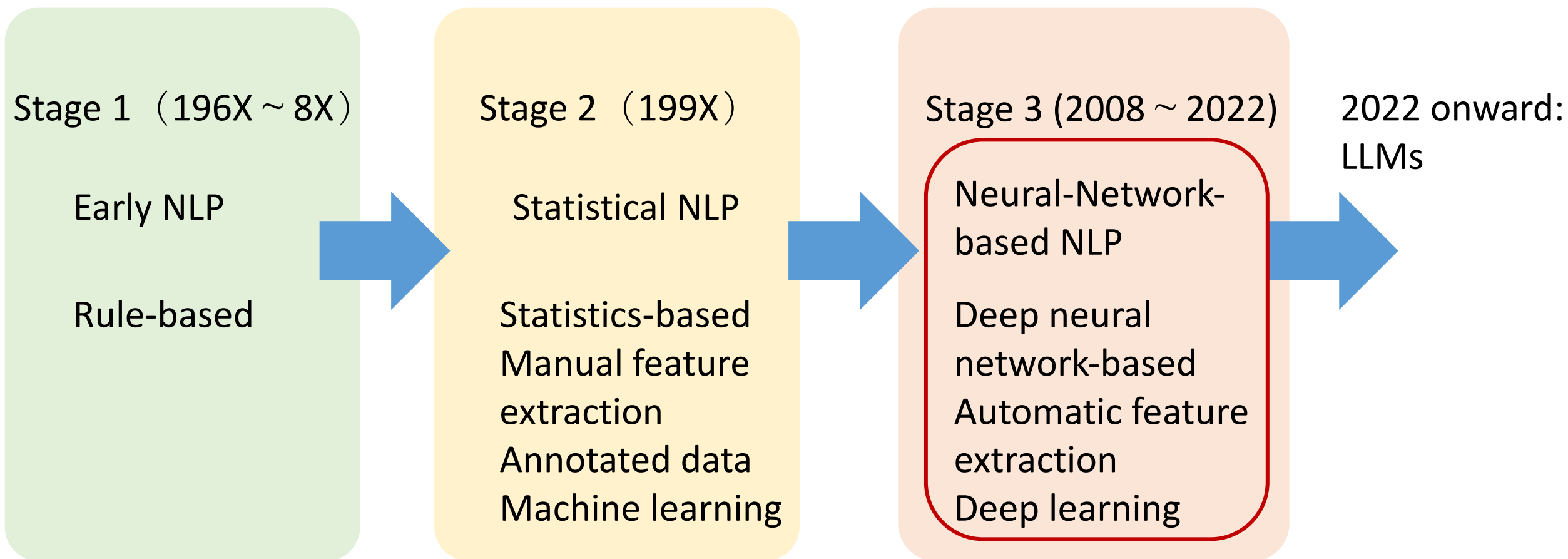


Table of Content

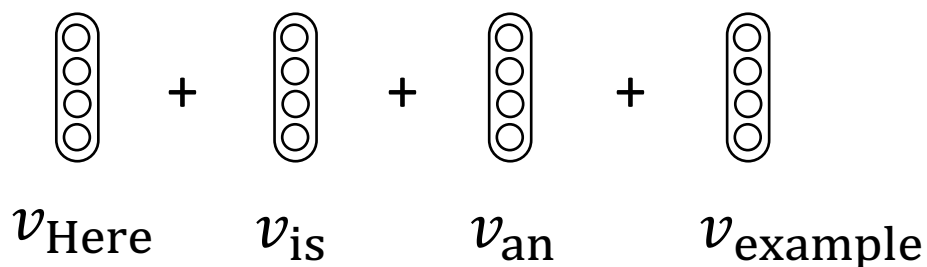
- **Motivation**
- Word Vectors
- Neural Networks
- Neural Text Classification

Lecture 01 - Word Vectors and Neural Networks

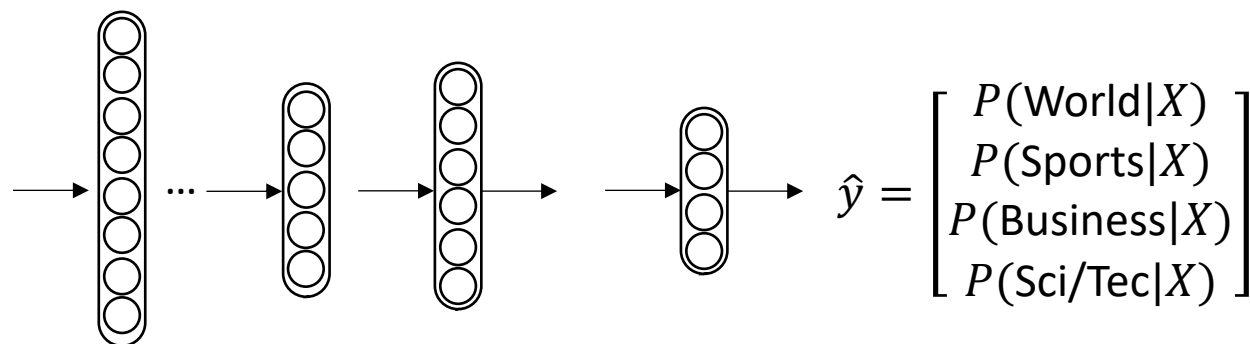
Bag-of-Words Neural Net

Task: News text classification

X : ["Here", "is", "an", "example"]



sentence vector v_X



A typical PyTorch implementation

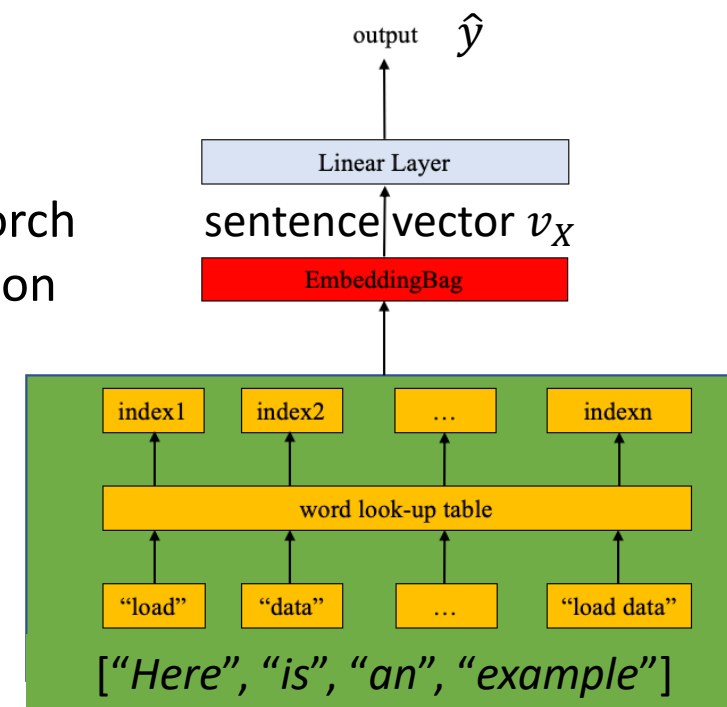


Table of Content Lecture 02 - Word Vectors

- Motivation
- Documents and Counts-based Method
- Neural Network-based Method -- word2vec
- Evaluation and Applications

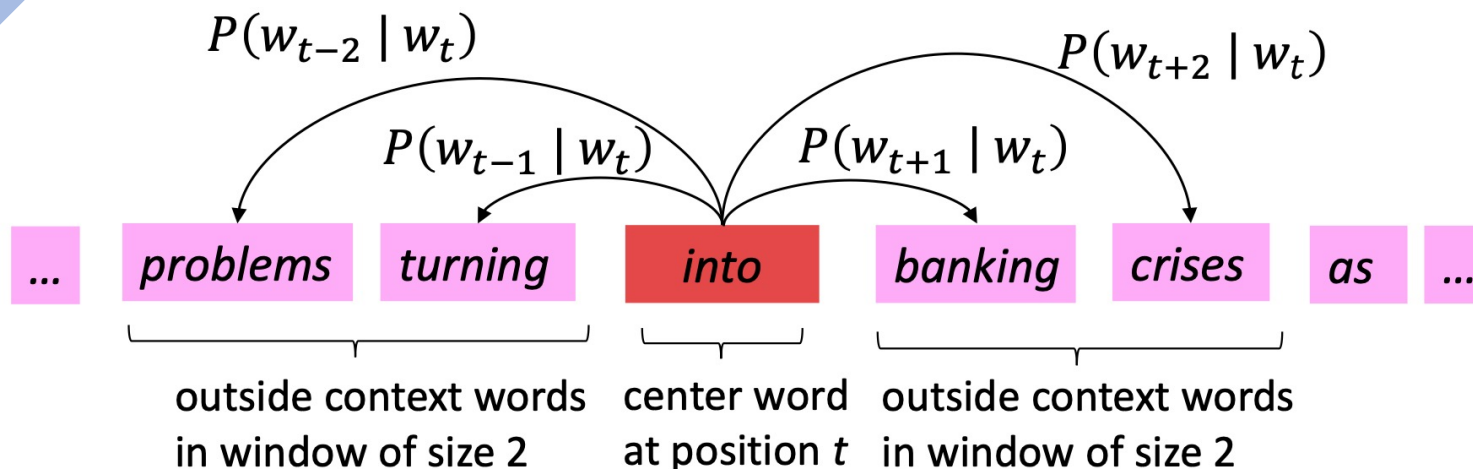
Two architectures of Word2vec

Word2vec { **Skip-gram**: Maximize $P(o|c)$

- “o” for outside (context) words
- “c” for center word

Continuous Bag-of-words (**CBOW**): Maximize $P(c|o)$

Compute probability
 $P(w_{t+j}|w_t)$,
for $j \in \{-2, -1, 1, 2\}$
when window size is 2



Negative Sampling: Objective Function

- For token at position t , maximize the log-likelihood:

Word o is the positive sample

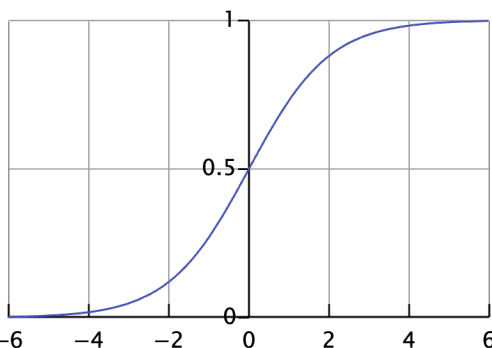
$$J_t(\theta) = \log \sigma(u_o^\top v_c) + \sum_{i=1}^k \mathbb{E}_{w_i \sim P(w)} [\log \sigma(-u_{w_i}^\top v_c)]$$

The k words w_i ($i = 1 \dots k$) are the negative samples

- Sigmoid function $\sigma(u_o^\top v_c)$ outputs the probability of o in the context window of c

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

a monotone increasing function



Maximizing this term will push the dot product $u_o^\top v_c$ to **larger** values, i.e., making o and c closer in semantic space

Maximizing this term will push the dot product $u_{w_i}^\top v_c$ to **smaller** values, i.e., making w_i and c farther apart in semantic space

Table of Content

- **Language Modeling**
- Neural Language Models
- Recurrent Neural Networks for LM
- Evaluate LMs

Lecture 03 - Recurrent Neural Networks and Language Modeling

Fixed-Window Neural Language Model

- Idea:** Represent words with embedding vectors; predict the next word using the concatenated embeddings from a fixed context window

concatenated word embeddings

$$e = \underbrace{[e^{(1)}; e^{(2)}; e^{(3)}; e^{(4)}]}_{4 \times d} \Bigg\}^d$$

Input tokens: $x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$

(window size = 4)

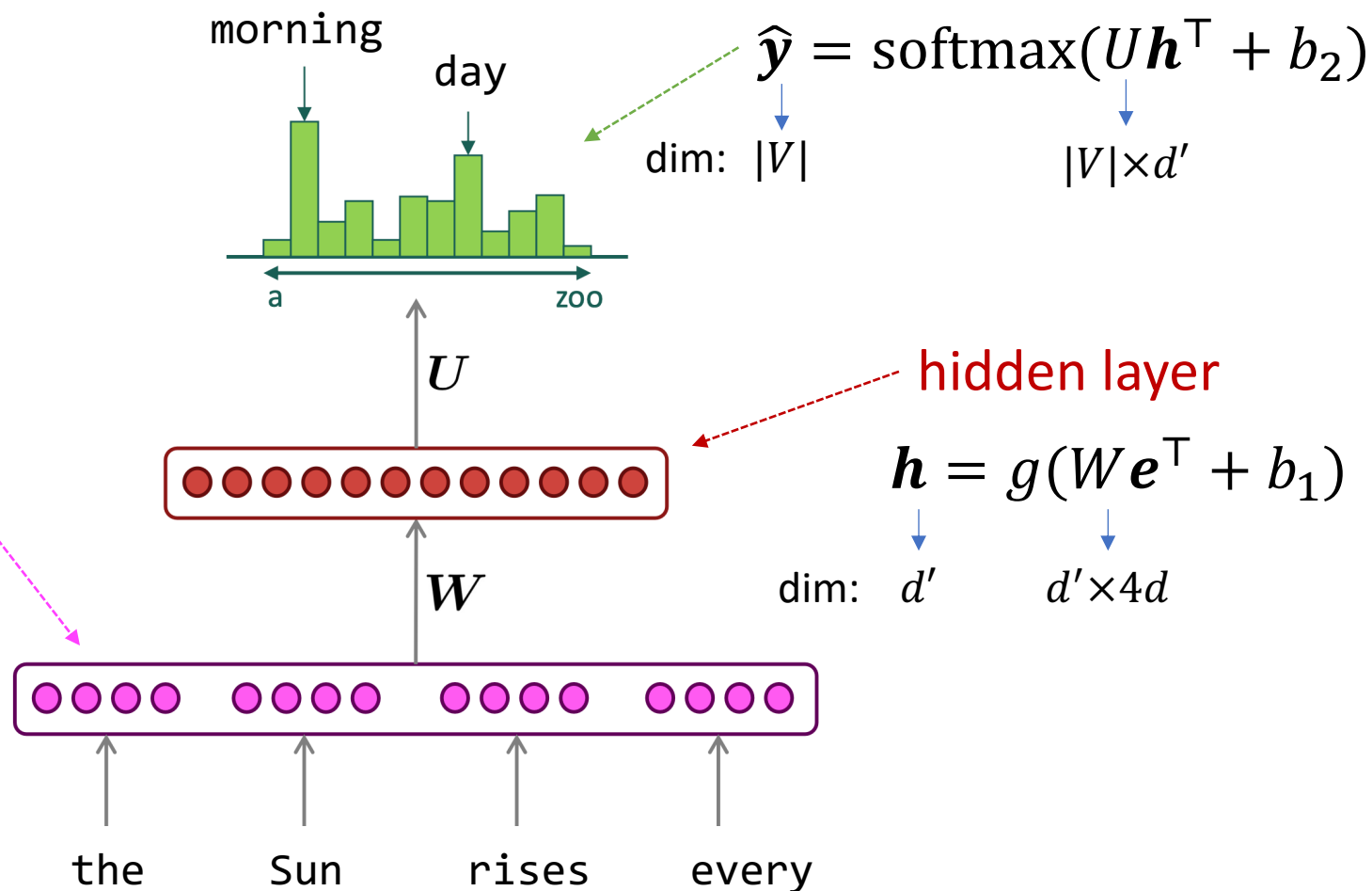
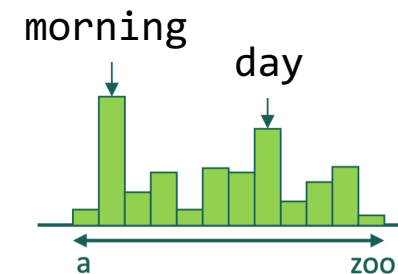


Figure from: <https://web.stanford.edu/class/archive/cs/cs224n/cs224n.1224/>

A Simple RNN Language Model



$$\hat{y}^{(4)} = P(x^{(5)} | \text{the Sun rises every})$$

output (optional)

$$\hat{y}^{(t)} = \text{softmax}(\mathbf{U}\mathbf{h}^{(t)} + b_2)$$

hidden state

$$\mathbf{h}^{(t)} = g(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_e \mathbf{e}^{(t)} + b_1)$$

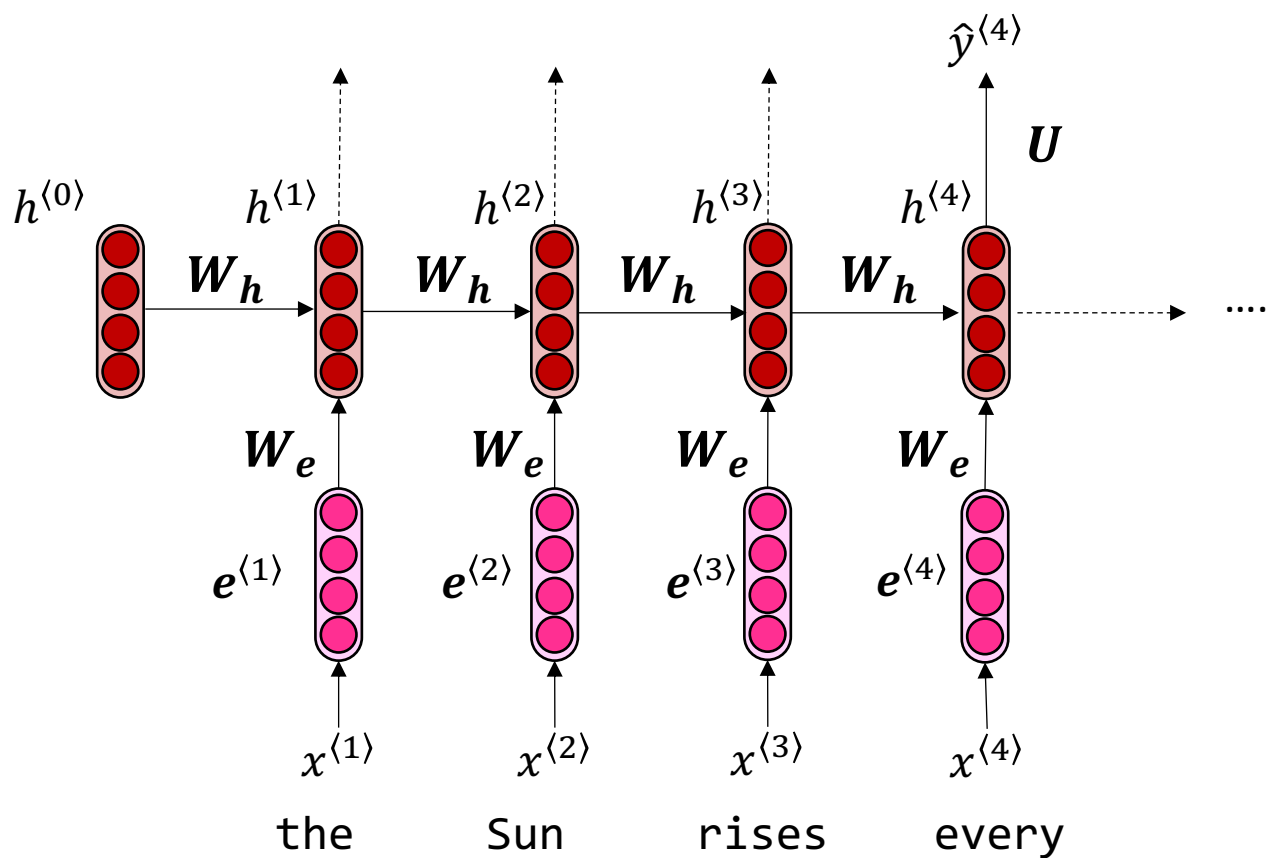
($\mathbf{h}^{(0)}$ is the initial hidden state)

input embedding

$$\mathbf{e}^{(t)} \in \mathbb{R}^d$$

input sequence

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

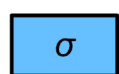
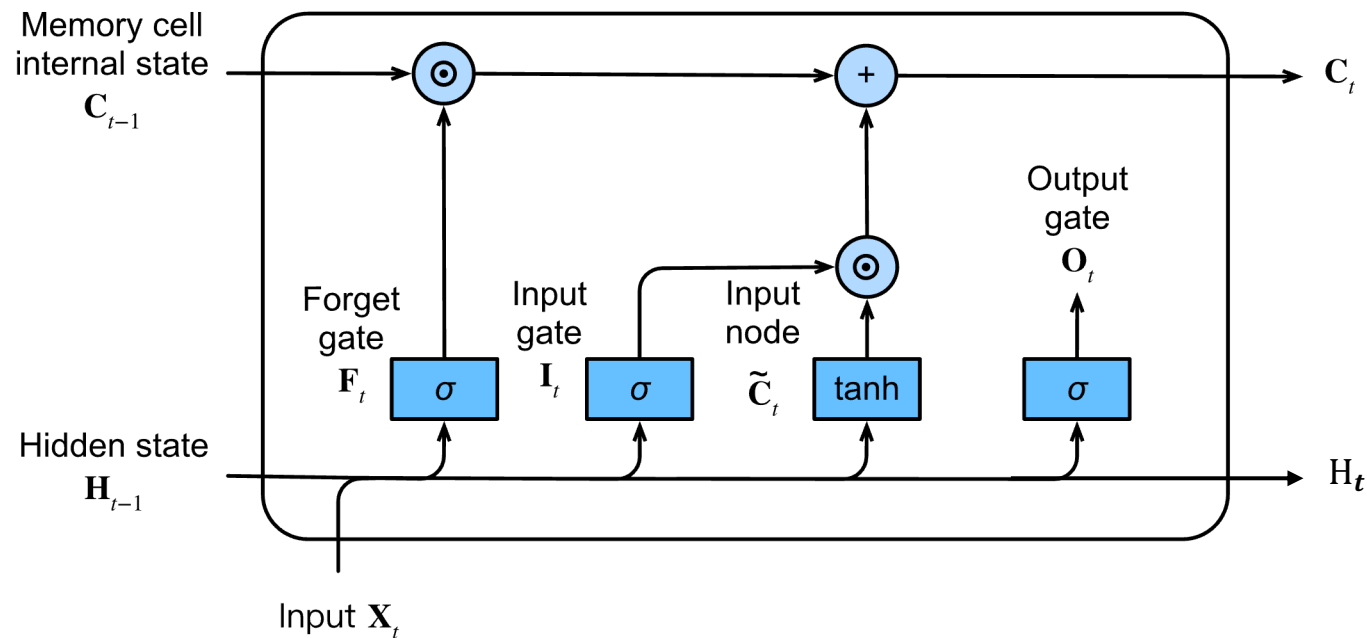


Overview

- Long Short-Term Memory RNNs (LSTMs)
- Bidirectional and multi-layer RNNs
- Sequence Labeling Task

Lecture 04 - Recurrent Neural Networks and Sequence Labeling

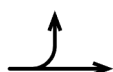
LSTM Computational Graph



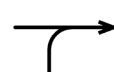
FC layer with
activation function



Elementwise
operator



Copy



Concatenate

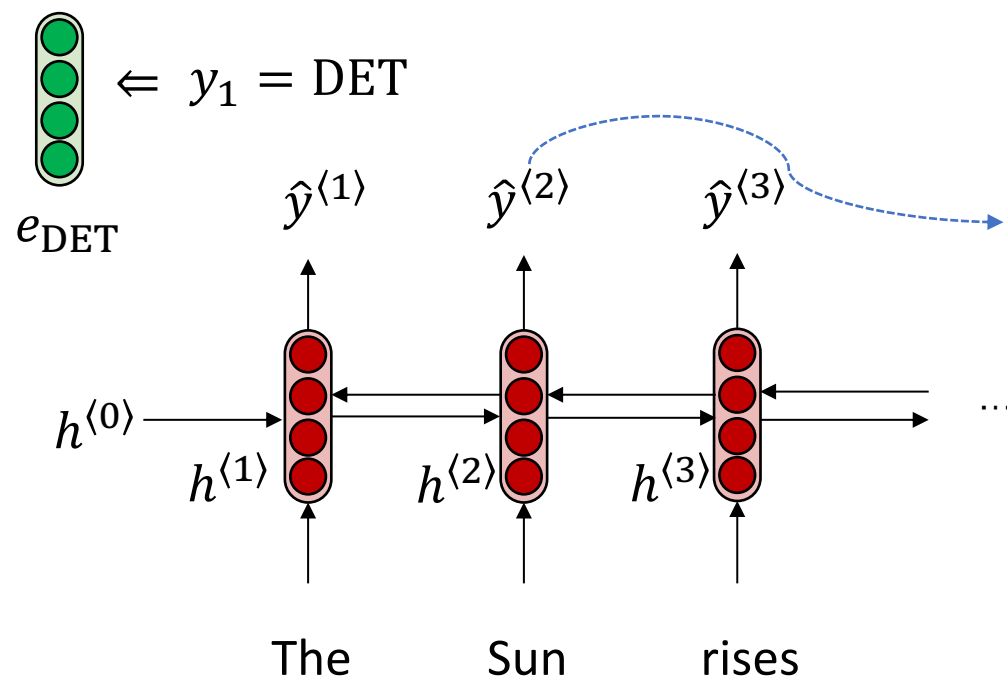
Figure from: https://d2l.ai/chapter_recurrent-modern/lstm.html

Level 3: MEMM by RNN

- Choice for $\phi(\mathbf{x}, y_i, y_{i+1}) \Rightarrow$ it should measures the likelihood $p(y_{i+1}|\mathbf{x}, y_i)$ and be in the form of a valid probability:

$$\phi = \frac{\exp(\text{feat}(\mathbf{x}, y_i, y_{i+1}))}{\sum_{y_{i+1} \in \mathcal{L}} \exp(\text{feat}(\mathbf{x}, y_i, y_{i+1}))}$$

$\text{feat}()$ is a parameterized feature function



$$= \frac{\exp(\text{fc}([\mathbf{h}^{(2)}; e_{\text{DET}}]) \odot e_{\text{NOUN}})}{\sum_{y \in \mathcal{L}} \exp(\text{fc}([\mathbf{h}^{(2)}; e_{\text{DET}}]) \odot e_y)}$$

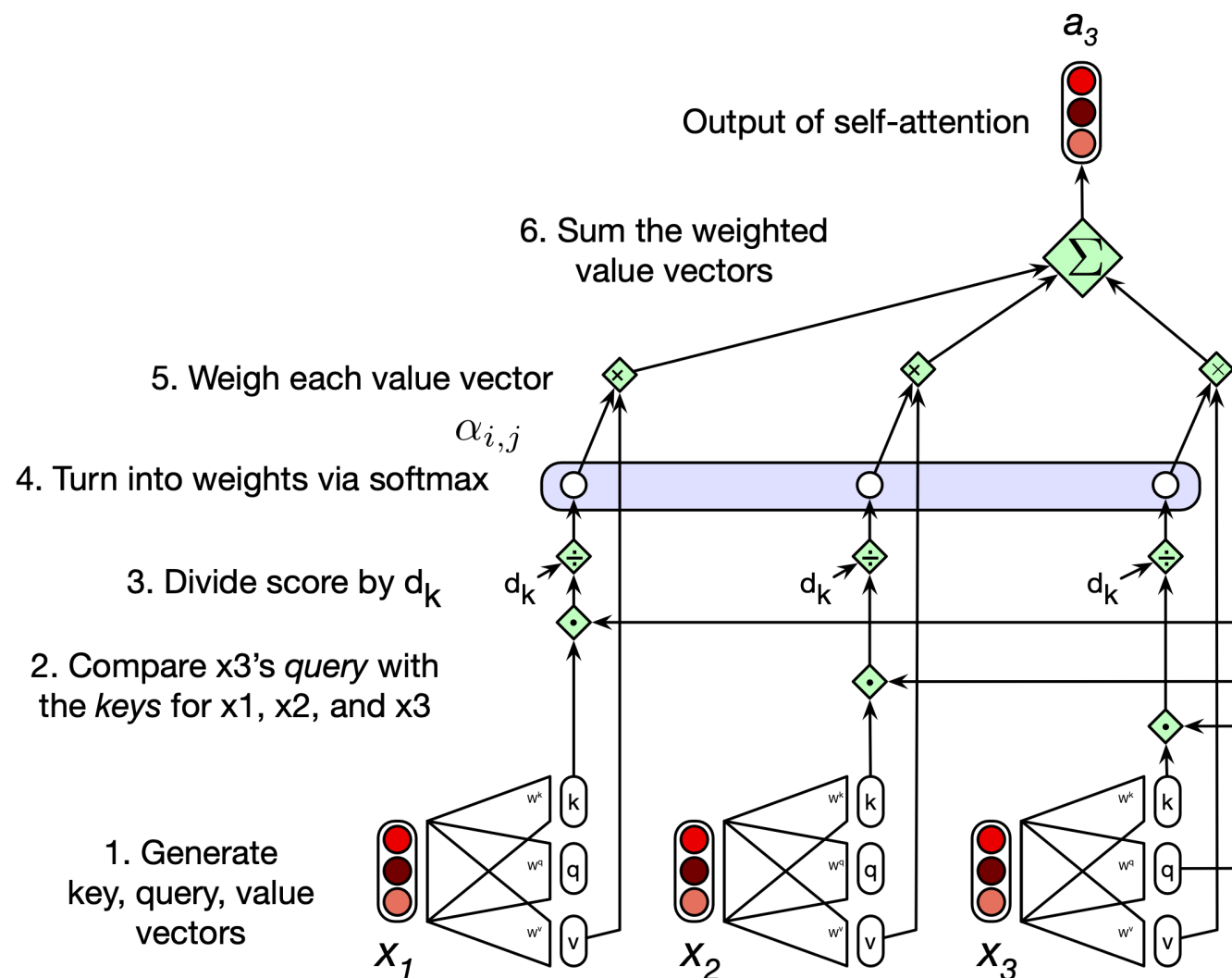
- $\text{fc}()$ is a fully-connected layer
- $[\mathbf{h}^{(2)}; e_{\text{DET}}]$ is the concatenation of $\mathbf{h}^{(2)}$ and e_{DET}
- \odot is dot product

Overview

Lecture 07 - Transformer

- Motivation
- Transformer Model
- Achievements and Drawbacks

Self-Attention Final Version



$$\begin{aligned} q_i &= x_i W^Q \\ k_i &= x_i W^K \\ v_i &= x_i W^V \end{aligned}$$

$$\text{score}(x_i, x_j) = \frac{q_i \cdot k_j}{\sqrt{d_k}}$$

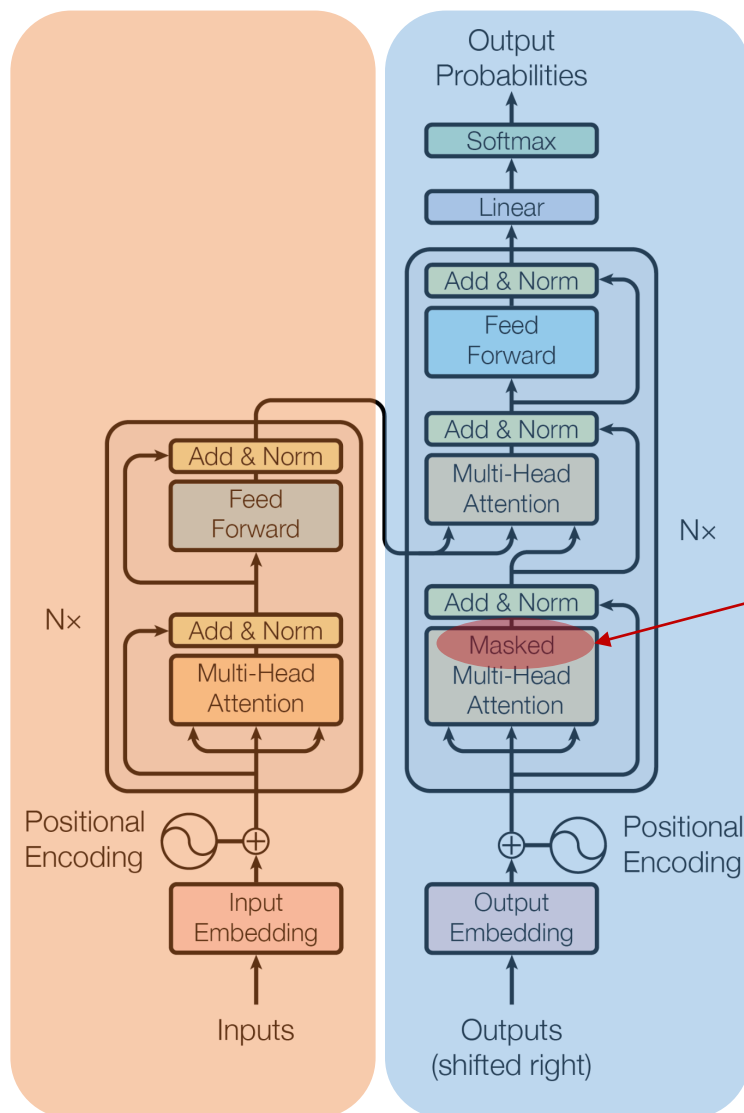
$$\alpha_{ij} = \text{softmax}(\text{score}(x_i, x_j))$$

$$\forall j \leq i$$

$$a_i = \sum_{j \leq i} \alpha_{ij} v_j$$

Transformer Model (encoder-decoder)

Encoder



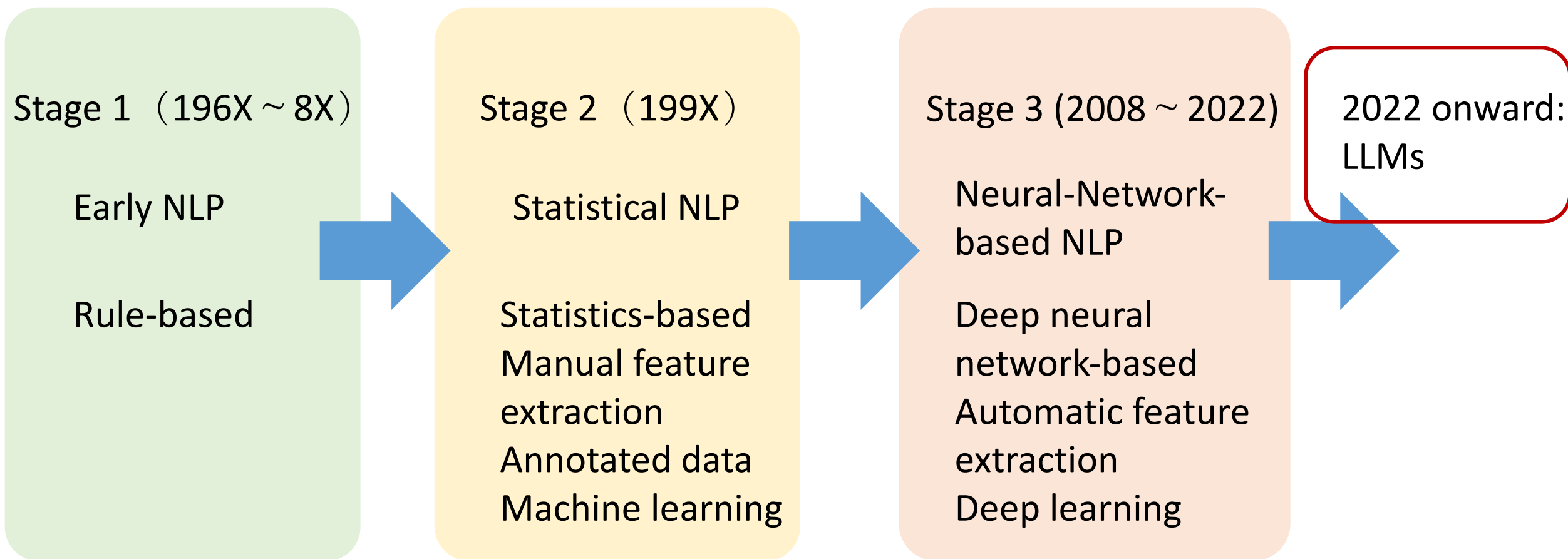
Decoder

Vaswani et al. (2017)'s original work is for machine translation task

- In their encoder, the self-attention is bidirectional;
- In the decoder, the self-attention is causal, i.e., future words are masked out

It was only later that the paradigm for causal language model was defined using only the decoder part

A Historic View of NLP



Overview

- Motivation: What is NLG?
- Review of Language Models
- Decoding from NLG models
- Evaluating NLG Systems

Lecture 09 - Natural Language Generation

Overview

Lecture 10 - Instruction Tuning and Human Feedbacks

- Motivation
- Instruction Tuning
- Reinforcement Learning from Human Feedback (RLHF)
- Evaluation of RLHF
- What's Next?

Overview

- Prompting
- Parameter Efficient Fine-Tuning

Lecture 11 - Prompting and Parameter Efficient Fine-Tuning