

# 자연어 처리 및 응용

The 13<sup>th</sup> KIAS CAC Summer School on  
the Parallel Computing and Artificial Intelligence

한양대학교 컴퓨터소프트웨어학부  
김은솔

# 대기업들이 뛰어드는 '초거대 AI'는 무엇

임영신 기자 | 입력 : 2021.07.08 10:23:10 수정 : 2021.07.08 10:23:26



## 글로벌 초거대 AI 성능 비교

초거대 AI	개발사	주요 기능	파라미터 개수	발표 시점
RoBERTa	페이스북	언어 생성 · 번역 · 검색 · 기사 작성 등	3억5500만	2019년 7월
GPT-2	오픈 AI	언어 생성 · 번역 · 검색 · 기사 작성 등	15억	2019년 8월
T5	구글	언어 생성 · 번역 · 검색 · 기사 작성 등	110억	2020년 2월
GPT-3	오픈 AI	기존 모든 기능의 고도화 · 프로그래밍	1750억	2020년 6월
하이퍼클로바	네이버	기존 모든 기능의 고도화 · 한국어 문장 생성 탁월	2040억	2021년 5월
우다오 2.0	베이징 지우안 인공지능연구원	기존 모든 기능의 고도화 · 중국어 문장 및 이미지 생성 탁월	1조7500억	2021년 6월
LG 초거대 AI	LG그룹	언어 · 이미지 이해 및 생성 · 데이터 추론	6000억	올 하반기(예정)
GPT-4	오픈 AI	GPT-3 초월 전망	100조	2023년(예정)

## 국내 기업의 주요 초거대 인공지능(AI) 기술

자료: 각 사

	초거대 AI	특징
카카오 브레인	코지피티(KoGPT)	<ul style="list-style-type: none"> <li>한국어 특화 AI 언어모델</li> <li>구글 텐서 처리장치 활용, 연산속도 고도화</li> </ul>
	민달리(minDALL-E)	<ul style="list-style-type: none"> <li>1400만 장의 텍스트·이미지 세트 사전 학습</li> <li>텍스트 명령어 입력하면 실시간 이미지 생성</li> </ul>
네이버	하이퍼클로바(HyperCLOVA)	<ul style="list-style-type: none"> <li>2040억 개에 이르는 매개변수(파라미터)</li> <li>학습 데이터의 한글 비중 97%, 한국어 집중 교육</li> </ul>
LG	엑사원(EXAONE)	<ul style="list-style-type: none"> <li>언어·이미지·영상 등을 다루는 멀티 모델리티 능력</li> <li>제조·연구·교육·금융 분야 상위 1% 전문가 목표</li> </ul>

# Transformer in Biological Science

## Highly accurate protein structure prediction with AlphaFold


<https://doi.org/10.1038/s41586-021-03819-2>

Received: 11 May 2021

Accepted: 12 July 2021

Published online: 15 July 2021

Open access

 Check for updates

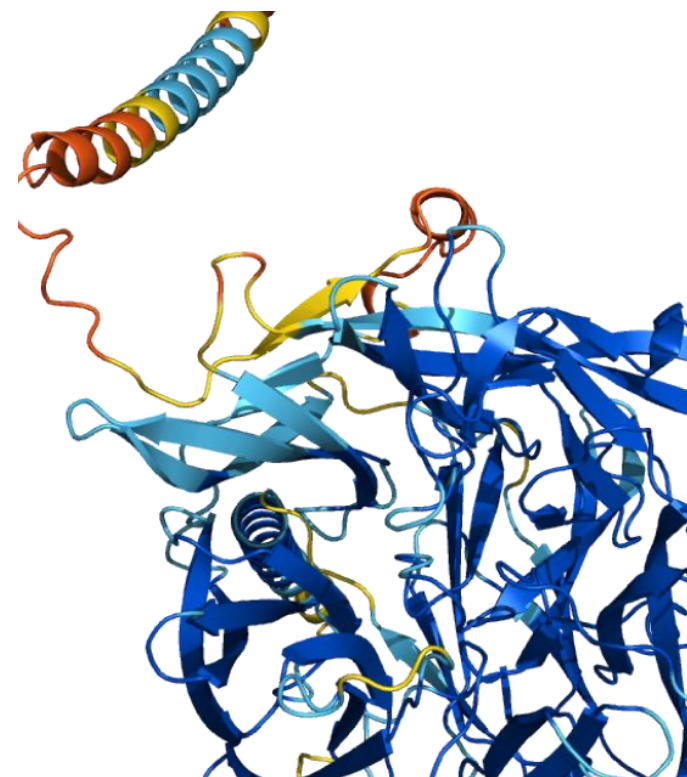
John Jumper<sup>1,4</sup>, Richard Evans<sup>1,4</sup>, Alexander Pritzel<sup>1,4</sup>, Tim Green<sup>1,4</sup>, Michael Figurnov<sup>1,4</sup>, Olaf Ronneberger<sup>1,4</sup>, Kathryn Tunyasuvunakool<sup>1,4</sup>, Russ Bates<sup>1,4</sup>, Augustin Židek<sup>1,4</sup>, Anna Potapenko<sup>1,4</sup>, Alex Bridgland<sup>1,4</sup>, Clemens Meyer<sup>1,4</sup>, Simon A. A. Kohl<sup>1,4</sup>, Andrew J. Ballard<sup>1,4</sup>, Andrew Cowie<sup>1,4</sup>, Bernardino Romera-Paredes<sup>1,4</sup>, Stanislav Nikolov<sup>1,4</sup>, Rishub Jain<sup>1,4</sup>, Jonas Adler<sup>1</sup>, Trevor Back<sup>1</sup>, Stig Petersen<sup>1</sup>, David Reiman<sup>1</sup>, Ellen Clancy<sup>1</sup>, Michal Zielinski<sup>1</sup>, Martin Steinegger<sup>2,3</sup>, Michalina Pacholska<sup>1</sup>, Tamas Berghammer<sup>1</sup>, Sebastian Bodenstein<sup>1</sup>, David Silver<sup>1</sup>, Oriol Vinyals<sup>1</sup>, Andrew W. Senior<sup>1</sup>, Koray Kavukcuoglu<sup>1</sup>, Pushmeet Kohli<sup>1</sup> & Demis Hassabis<sup>1,4</sup>

**nature**

## Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences

Alexander Rives<sup>a,b,1,2</sup>, Joshua Meier<sup>a,1</sup>, Tom Sercu<sup>a,1</sup>, Siddharth Goyal<sup>a,1</sup>, Zeming Lin<sup>b</sup>, Jason Liu<sup>a</sup>, Demi Guo<sup>c,3</sup>, Myle Ott<sup>a</sup>, C. Lawrence Zitnick<sup>a</sup>, Jerry Ma<sup>d,e,3</sup>, and Rob Fergus<sup>b</sup>

**PNAS**



# Contents

- Word Embedding
  - One-hot embedding
  - Word2Vec (Skip-gram, CBOW)
  - GloVe
- Language Model
  - n-gram
  - Recurrent Neural Network
  - Attention Methods
- Transformer
- Large-scale Language Models
- Applications
  - Machine Translation
  - Question Answering

# Word Embedding

# Data Representation

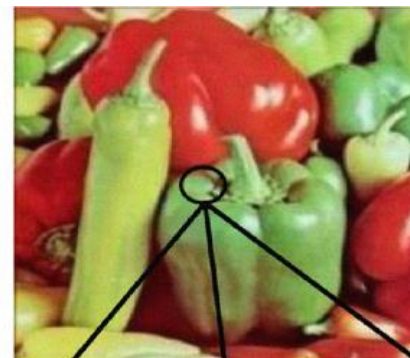
Table of baby-name data  
(baby-2010.csv)

name	rank	gender	year
Jacob	1	boy	2010
Isabella	1	girl	2010
Ethan	2	boy	2010
Sophia	2	girl	2010
Michael	3	boy	2010

Field  
names

One row  
(4 fields)

2000 rows  
all told



240 241 241  
240 237 238  
239 240 240  
238 237 240  
240 240 239  
239 240 240

207 199 196  
183 163 195  
183 166 184  
176 172 181  
184 167 176  
182 180 170

234 231 225  
223 213 225  
219 211 195  
176 205 189  
168 141 117  
160 142 117

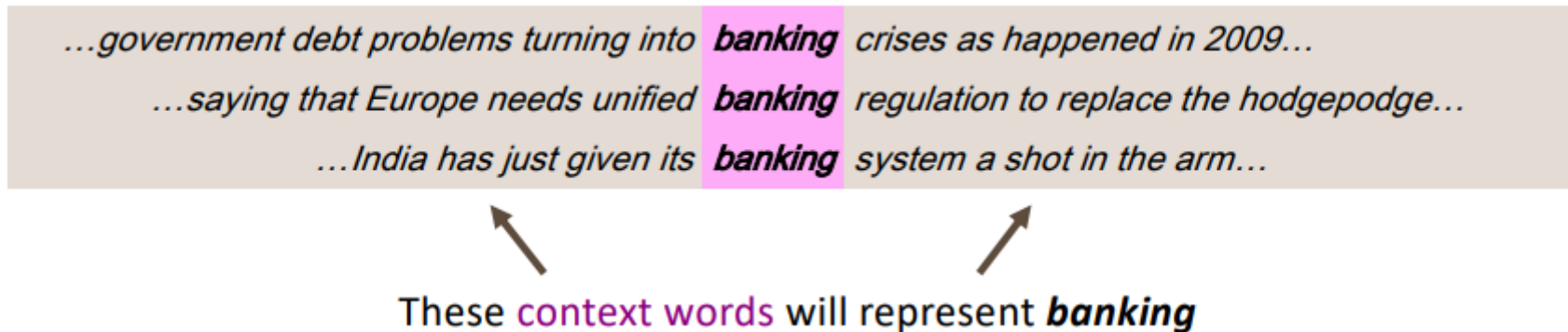
# Data Representation – Text

- Conventional Word Representations
  - 이미지, 음성과 달리 언어 데이터는 discrete
  - One-hot encoding
    - 데이터에 포함된 단어로 사전을 만들고, 이를 기반으로 one-hot encoding을 하여 단어를 표현
    - Discrete, Sparse
- All vectors are orthogonal
  - There are no natural notion of similarity for one-hot vectors

Word	One-hot encoding
economic	000010...
growth	001000...
has	100000...
slowed	000001...

# Word Embedding

- Assumption: Distributional semantics (hypothesis)
  - *Linguistic items with similar distributions have similar meanings*
  - Representing words by their context





# Word Embedding

- Build a dense vector for each word
- similar to vectors of words that appear in similar contexts

$$\textit{banking} = \begin{pmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 0.109 \\ -0.542 \\ 0.349 \\ 0.271 \end{pmatrix}$$



# Word2Vec

- Efficient Estimation of Word Representations in Vector Space
  - T. Mikolov et al., ICLR Workshop, 2013
- Distributed Representations of Words and Phrases and their Compositionality
  - T. Mikolov et al., 2013, NeurIPS

---

## **Distributed Representations of Words and Phrases and their Compositionality**

---

**Tomas Mikolov**  
Google Inc.  
Mountain View  
mikolov@google.com

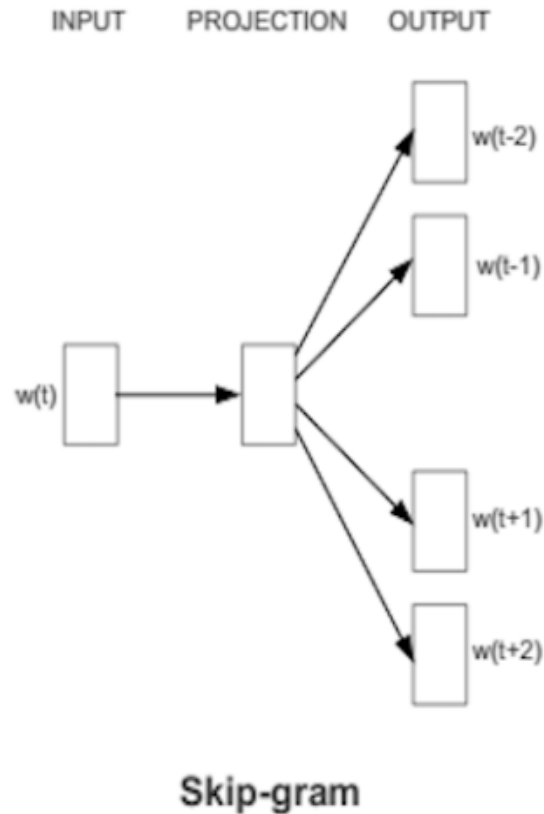
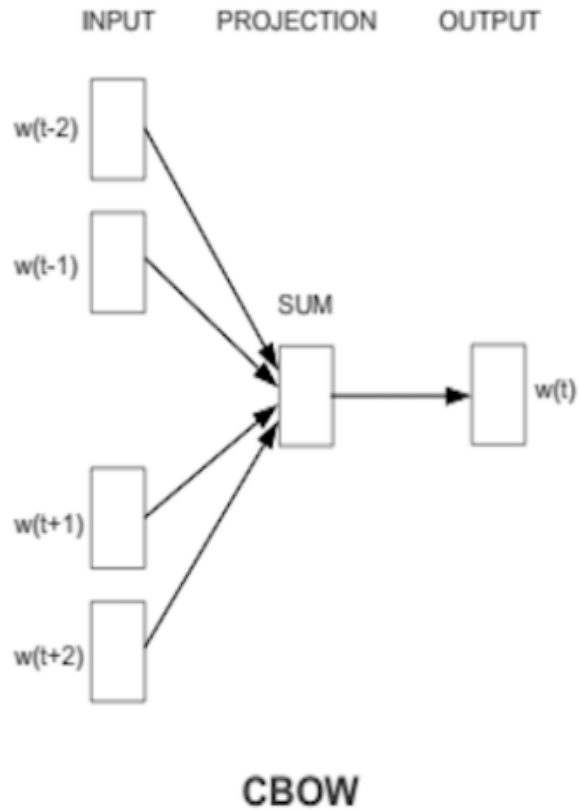
**Ilya Sutskever**  
Google Inc.  
Mountain View  
ilyasu@google.com

**Kai Chen**  
Google Inc.  
Mountain View  
kai@google.com

**Greg Corrado**  
Google Inc.  
Mountain View  
gcorrado@google.com

**Jeffrey Dean**  
Google Inc.  
Mountain View  
jeff@google.com

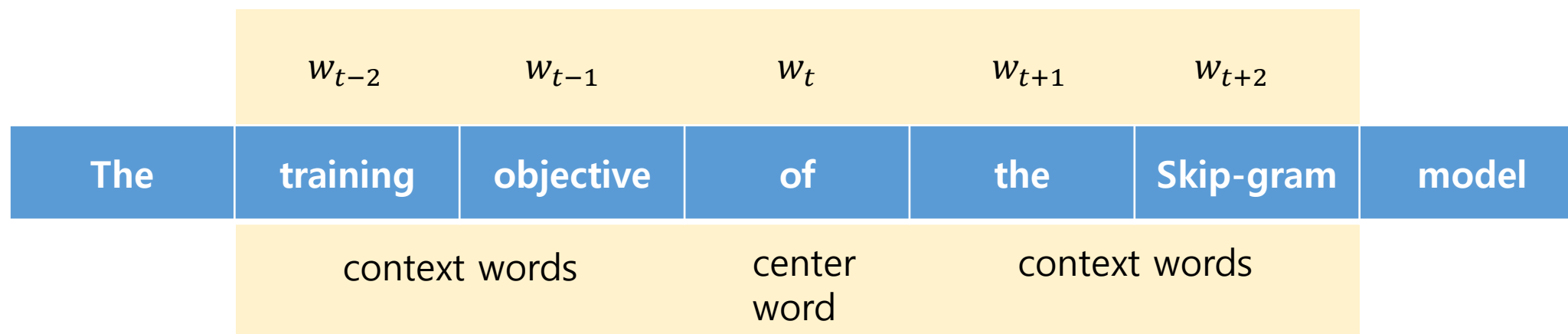
# Word2Vec



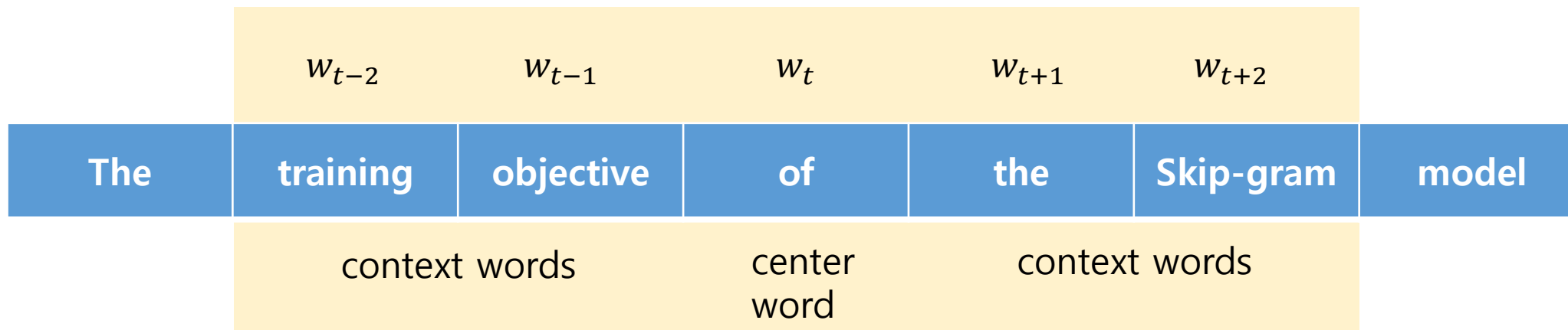
- Key Idea

- Find word representations that are useful for predicting the surrounding words
- Use the similarity of the word vectors to calculate the probability

# Word2Vec



# Word2Vec



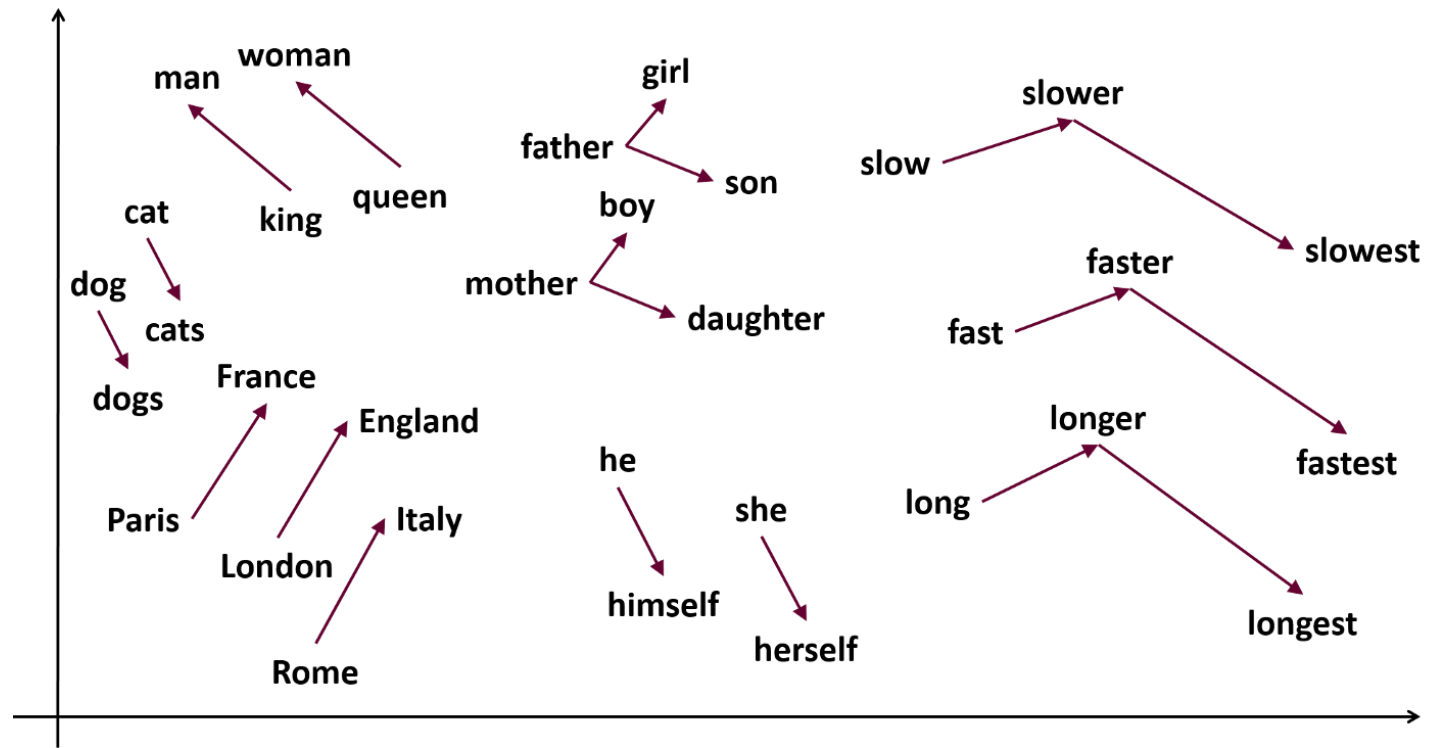
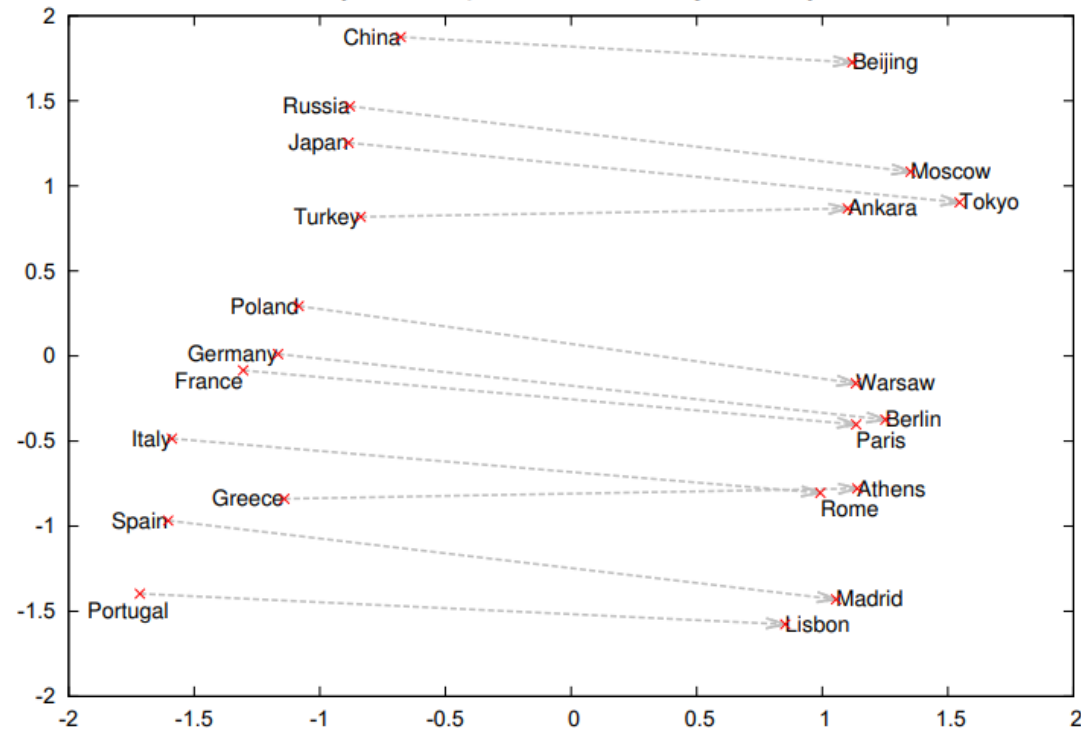
training objective function

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

$$p(w_O | w_I) = \frac{\exp(v'_{w_O}{}^\top v_{w_I})}{\sum_{w=1}^W \exp(v'_w{}^\top v_{w_I})}$$

# Interesting Results of the Word2Vec

Country and Capital Vectors Projected by PCA



# Another Word (Neural) Embedding

- Using co-occurrence information

	<i>a</i>	<i>as</i>	<i>chuck</i>	<i>could</i>	<i>how</i>	<i>if</i>	<i>much</i>	<i>wood</i>	<i>woodch.</i>	<i>would</i>	<i>,</i>	<i>.</i>	<i>?</i>	<i>a</i>	<i>as</i>	<i>chuck</i>	<i>could</i>	<i>how</i>	<i>if</i>	<i>much</i>	<i>wood</i>	<i>woodch.</i>	<i>would</i>	<i>,</i>	<i>.</i>	<i>?</i>
<i>a</i>	13	24	12	3	9	20	22	31	16	23	18	0	7	13	7	31	26	0	14	4	21	50	9	16	7	7
<i>as</i>	7	8	15	11	0	5	9	25	10	0	3	0	17	24	8	2	3	0	9	10	10	20	13	11	0	0
<i>chuck</i>	31	2	5	20	5	14	6	9	36	15	12	0	0	12	15	5	6	0	9	8	30	10	2	11	9	12
<i>could</i>	26	3	6	0	0	16	2	4	30	9	14	0	0	3	11	20	0	0	0	6	23	2	1	0	8	8
<i>how</i>	0	0	0	0	0	0	0	0	0	0	0	0	0	9	0	5	0	0	3	10	9	7	8	4	0	0
<i>if</i>	14	9	9	0	3	0	8	11	16	15	20	0	2	20	5	14	16	0	0	3	14	18	0	0	5	5
<i>much</i>	4	10	8	6	10	3	0	8	5	0	2	0	9	22	9	6	2	0	8	0	20	18	15	10	0	0
<i>wood</i>	21	10	30	23	9	14	20	7	26	5	11	0	8	31	25	9	4	0	11	8	7	26	20	14	10	10
<i>woodch.</i>	50	20	10	2	7	18	18	26	13	20	16	0	5	16	10	36	30	0	16	5	26	13	10	18	9	9
<i>would</i>	9	13	2	1	8	0	15	20	10	0	0	0	4	23	0	15	9	0	15	0	5	20	0	17	3	0
<i>,</i>	16	11	11	0	4	0	10	14	18	17	0	0	3	18	3	12	14	0	20	2	11	16	0	0	4	4
<i>.</i>	7	0	9	8	0	5	0	10	9	3	4	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
<i>?</i>	7	0	12	8	0	5	0	10	9	0	4	0	0	7	17	0	0	0	2	9	8	5	4	3	0	0

# Co-occurrence based word vectors

- Singular Value Decomposition of co-occurrence matrix  $X$
- Factorize  $X$  into  $U\Sigma V^T$ 
  - $U, V$  are orthogonal

$$\begin{array}{c}
 \begin{array}{ccc}
 \begin{array}{c} m \\ \boxed{\phantom{X}} \\ n \end{array} & = & \begin{array}{c} r \\ \boxed{\begin{array}{c} | \\ U_1 \\ | \\ U_2 \\ | \\ U_3 \\ | \\ \vdots \end{array}} \\ n \end{array} \begin{array}{c} r \\ \boxed{\begin{array}{c} S_1 \quad S_2 \quad S_3 \quad \dots \quad 0 \\ 0 \quad \quad \quad \ddots \quad S_i \end{array}} \\ r \end{array} \begin{array}{c} m \\ \boxed{\begin{array}{c} \text{---} V_1 \text{---} \\ \text{---} V_2 \text{---} \\ \text{---} V_3 \text{---} \\ \vdots \end{array}} \\ r \end{array} \\
 X & & U \quad S \quad V^T
 \end{array} \\
 \\
 \begin{array}{ccc}
 \begin{array}{c} m \\ \boxed{\phantom{\hat{X}}} \\ n \end{array} & = & \begin{array}{c} k \\ \boxed{\begin{array}{c} | \\ U_1 \\ | \\ U_2 \\ | \\ U_3 \\ | \\ \vdots \end{array}} \\ n \end{array} \begin{array}{c} k \\ \boxed{\begin{array}{c} S_1 \quad S_2 \quad S_3 \quad \dots \quad 0 \\ 0 \quad \quad \quad \ddots \quad S_i \end{array}} \\ k \end{array} \begin{array}{c} m \\ \boxed{\begin{array}{c} \text{---} V_1 \text{---} \\ \text{---} V_2 \text{---} \\ \text{---} V_3 \text{---} \\ \vdots \end{array}} \\ k \end{array} \\
 \hat{X} & & \hat{U} \quad \hat{S} \quad \hat{V}^T
 \end{array}
 \end{array}$$

An Improved Model of Semantic Similarity Based on Lexical Co-Occurrence

**Douglas L. T. Rohde**

Massachusetts Institute of Technology, Department of Brain and Cognitive Sciences

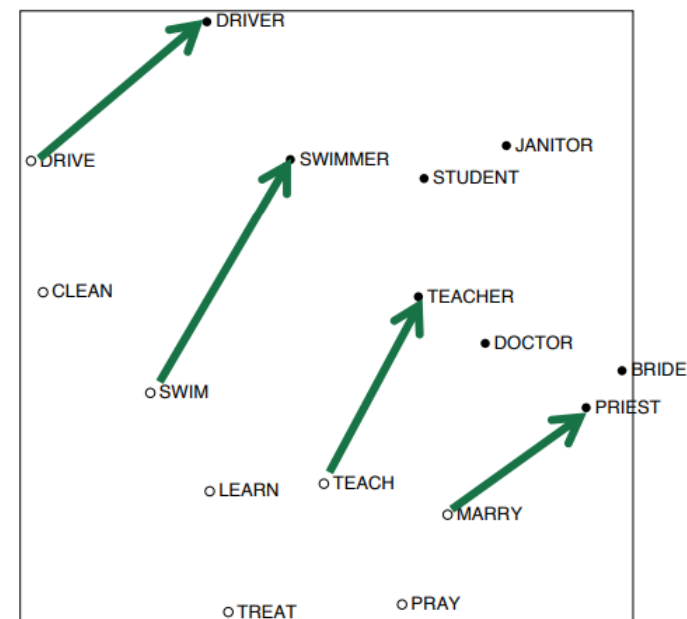
**Laura M. Gonnerman**

Lehigh University, Department of Psychology

**David C. Plaut**

Carnegie Mellon University, Department of Psychology,  
and the Center for the Neural Basis of Cognition

November 7, 2005





## GloVe: Global Vectors for Word Representation

# GloVe

Jeffrey Pennington, Richard Socher, Christopher D. Manning  
Computer Science Department, Stanford University, Stanford, CA 94305  
jpennin@stanford.edu, richard@socher.org, manning@stanford.edu

- Key idea
  - Ratios of co-occurrence probabilities can encode meaning components

	$x = \text{solid}$	$x = \text{gas}$	$x = \text{water}$	$x = \text{random}$
$P(x \text{ice})$	large	small	large	small
$P(x \text{steam})$	small	large	large	small
$\frac{P(x \text{ice})}{P(x \text{steam})}$	large	small	$\sim 1$	$\sim 1$

## GloVe: Global Vectors for Word Representation

# GloVe

Jeffrey Pennington, Richard Socher, Christopher D. Manning  
Computer Science Department, Stanford University, Stanford, CA 94305  
jpennin@stanford.edu, richard@socher.org, manning@stanford.edu

- Key idea
  - Ratios of co-occurrence probabilities can encode meaning components

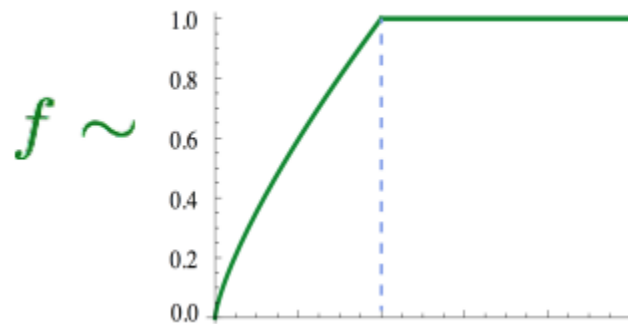
	$x = \text{solid}$	$x = \text{gas}$	$x = \text{water}$	$x = \text{fashion}$
$P(x \text{ice})$	$1.9 \times 10^{-4}$	$6.6 \times 10^{-5}$	$3.0 \times 10^{-3}$	$1.7 \times 10^{-5}$
$P(x \text{steam})$	$2.2 \times 10^{-5}$	$7.8 \times 10^{-4}$	$2.2 \times 10^{-3}$	$1.8 \times 10^{-5}$
$\frac{P(x \text{ice})}{P(x \text{steam})}$	8.9	$8.5 \times 10^{-2}$	1.36	0.96

# GloVe

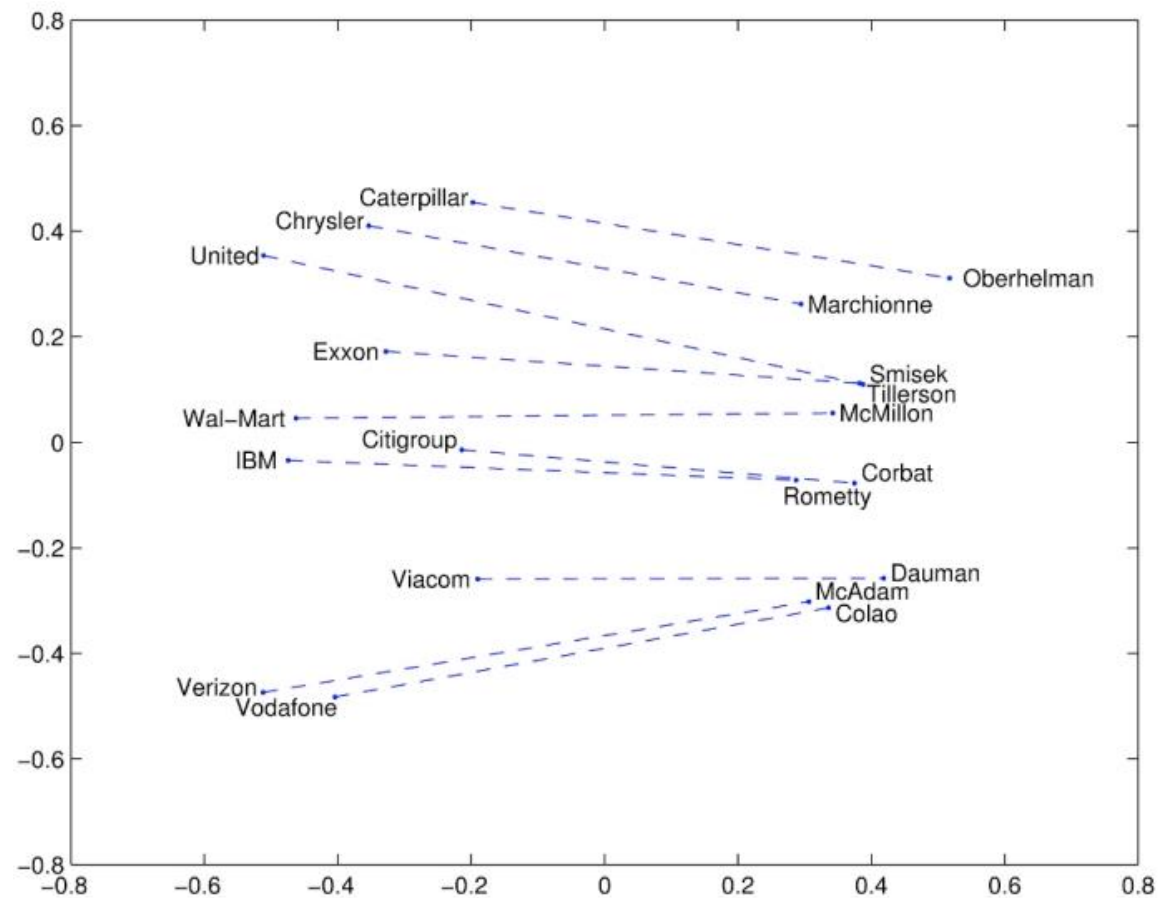
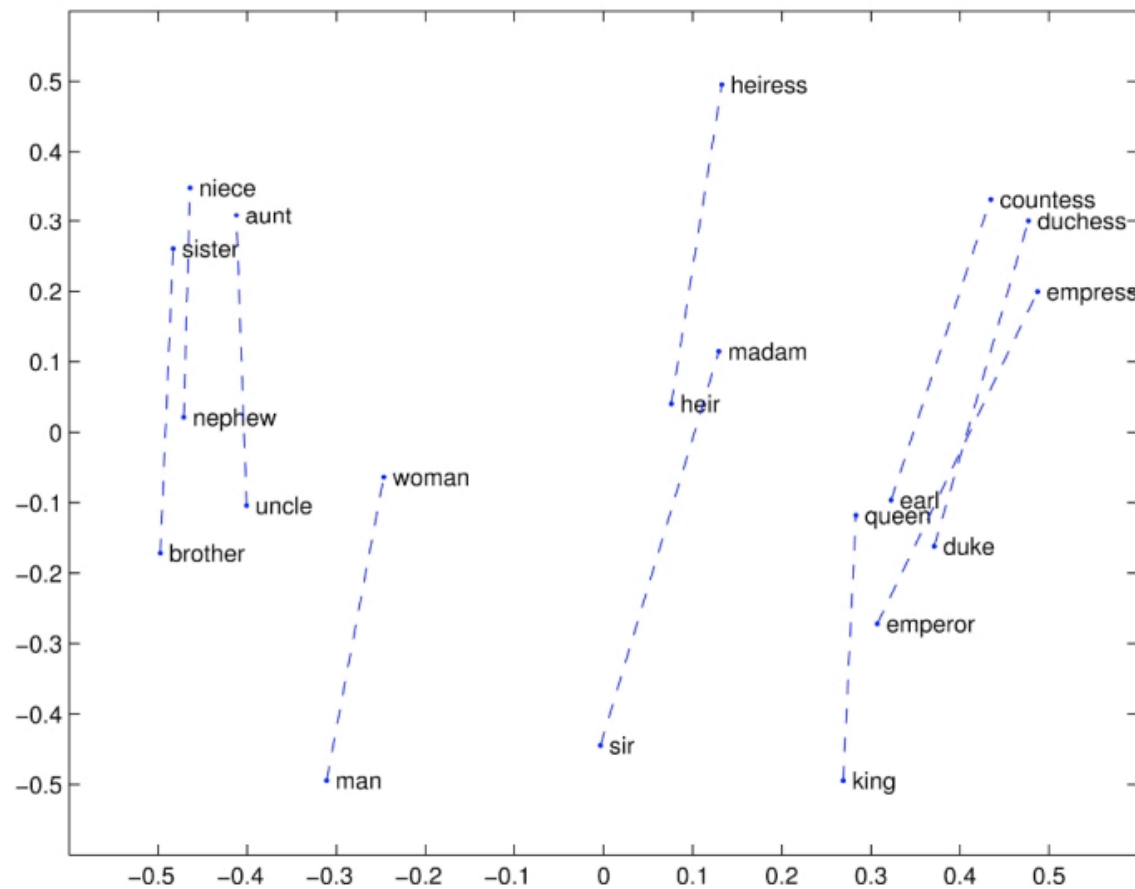
- Training Objective

$$w_i \cdot w_j = \log P(i|j)$$

$$J = \sum_{i,j=1}^V f(X_{ij}) \left( w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij} \right)^2$$

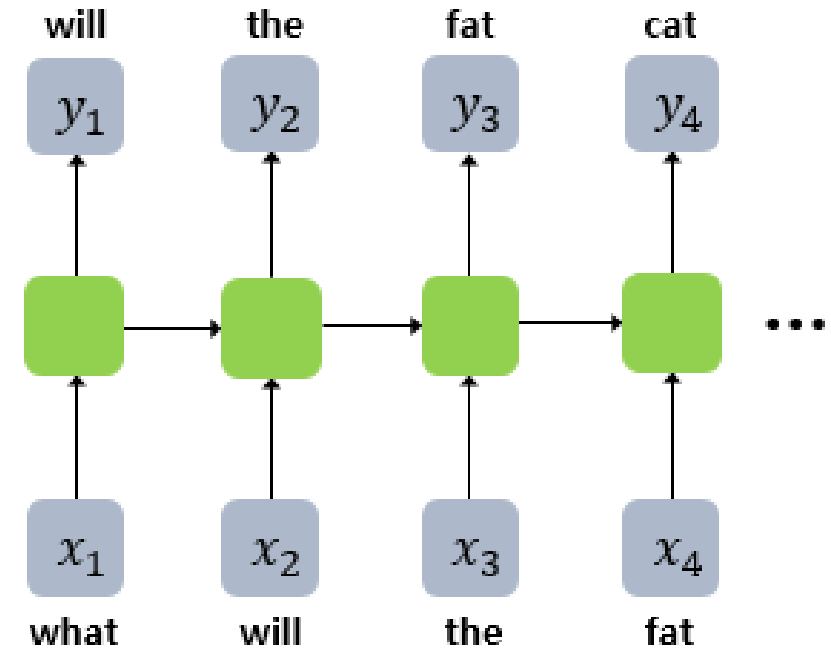


# GloVe



# Language Models

# Language Models



The

training

objective

of

the

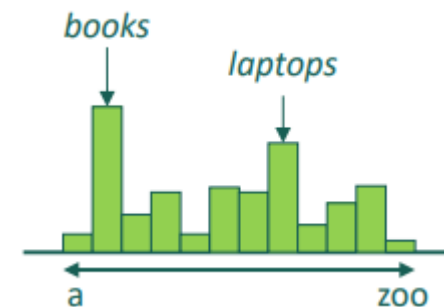
Skip-gram

model

# Language Models

- Language modeling is the task of **predicting what word comes next.**

$$P(x^{t+1} | x^t, \dots, x^1)$$



$$V = \{w_1, w_2, \dots, w_{|V|}\}$$

$$P(x^1, \dots, x^t) = P(x^1) \times P(x^2 | x^1) \times \dots \times P(x^t | x^{t-1}, \dots, x^1)$$

$$P(\text{This is a sentence}) = P(\text{This}) \times P(\text{is} | \text{This}) \times P(\text{a} | \text{This is}) \times P(\text{sentence} | \text{This is a})$$

# n-gram Language Models

- Modeling with Markov assumption
  - $x^t$  depends only on the preceding (n-1) words

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

- For example, if n=3

$P(\text{This is a sentence from AAA}) = P(\text{This})$

x  $P(\text{is} \mid \text{This})$   
x  $P(\text{a} \mid \text{This is})$   
x  $P(\text{sentence} \mid \text{This is a})$   
x  $P(\text{from} \mid \text{This is a sentence})$   
x  $P(\text{AAA} \mid \text{This is a sentence from})$

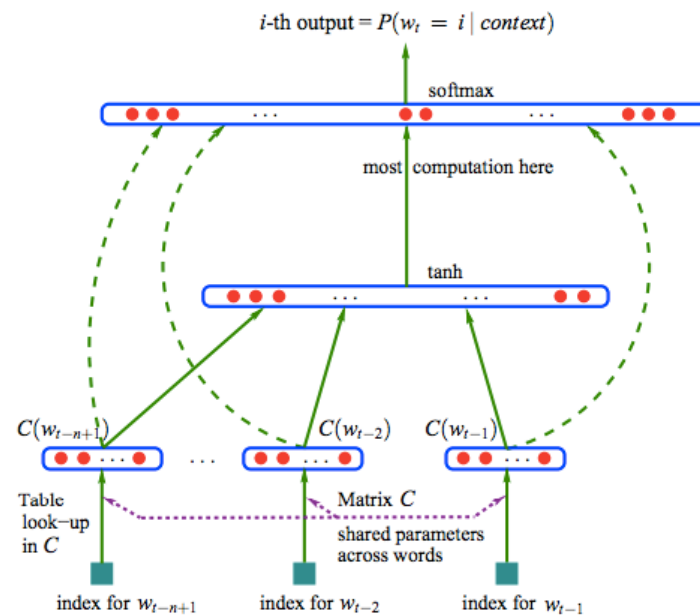
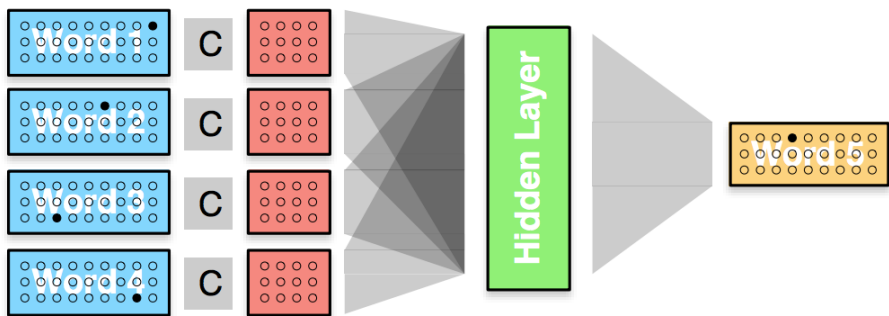
calculate by counting  
the phrases in large  
corpus of text



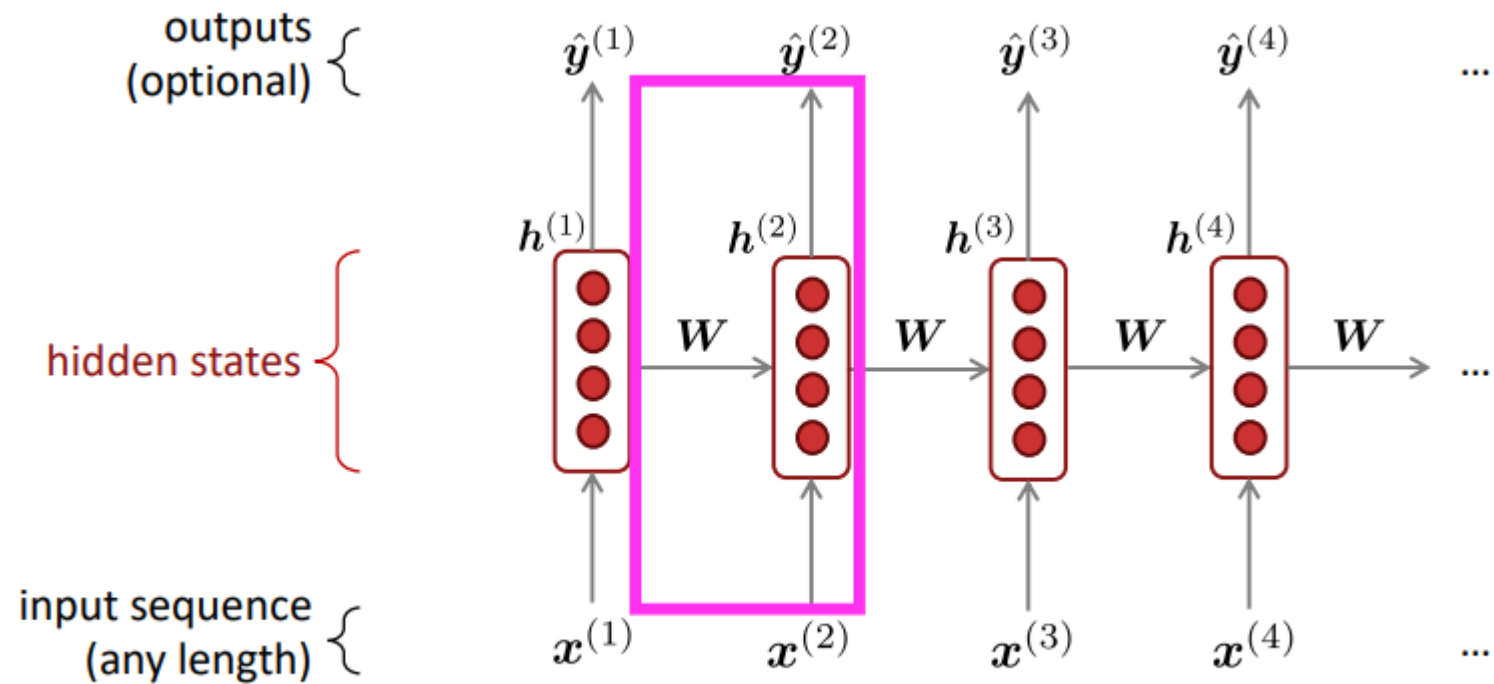
# Neural Word Embedding

- Neural network의 hidden vector 값을 이용하여 단어의 의미를 표현하는 기법
- Neural network language model (NNLM, 2003)

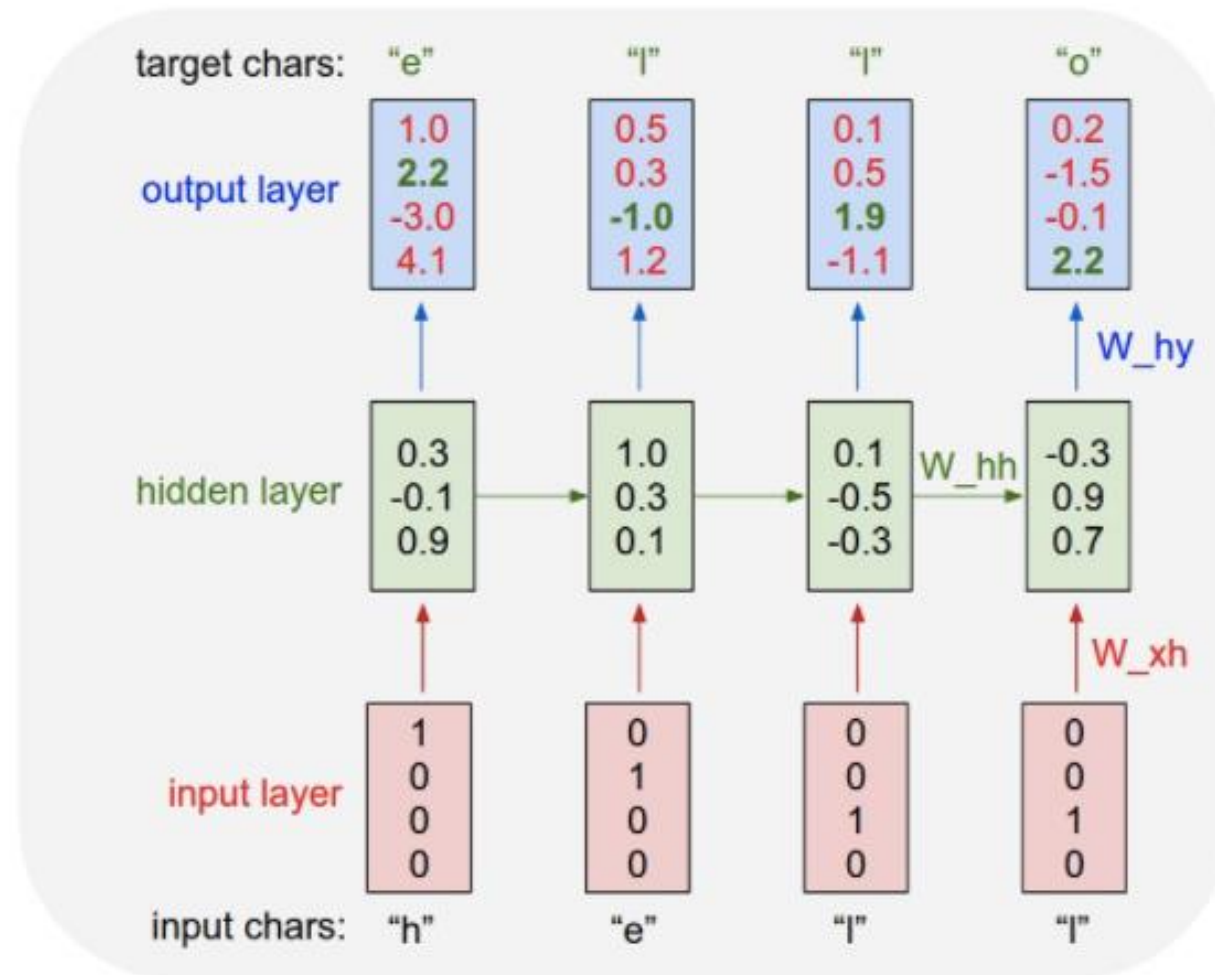
$$p(W) = \sum_i p(w_i | w_1, \dots, w_{i-1})$$
$$p(w_i | w_1, \dots, w_{i-1}) \simeq p(w_i | w_{i-4}, w_{i-3}, w_{i-2}, w_{i-1})$$



# Recurrent Neural Network (RNN)



# An illustrative example



# Recurrent Neural Network (RNN)

output distribution

$$\hat{y}^{(t)} = \text{softmax} \left( U h^{(t)} + b_2 \right) \in \mathbb{R}^{|V|}$$

hidden states

$$h^{(t)} = \sigma \left( W_h h^{(t-1)} + W_e e^{(t)} + b_1 \right)$$

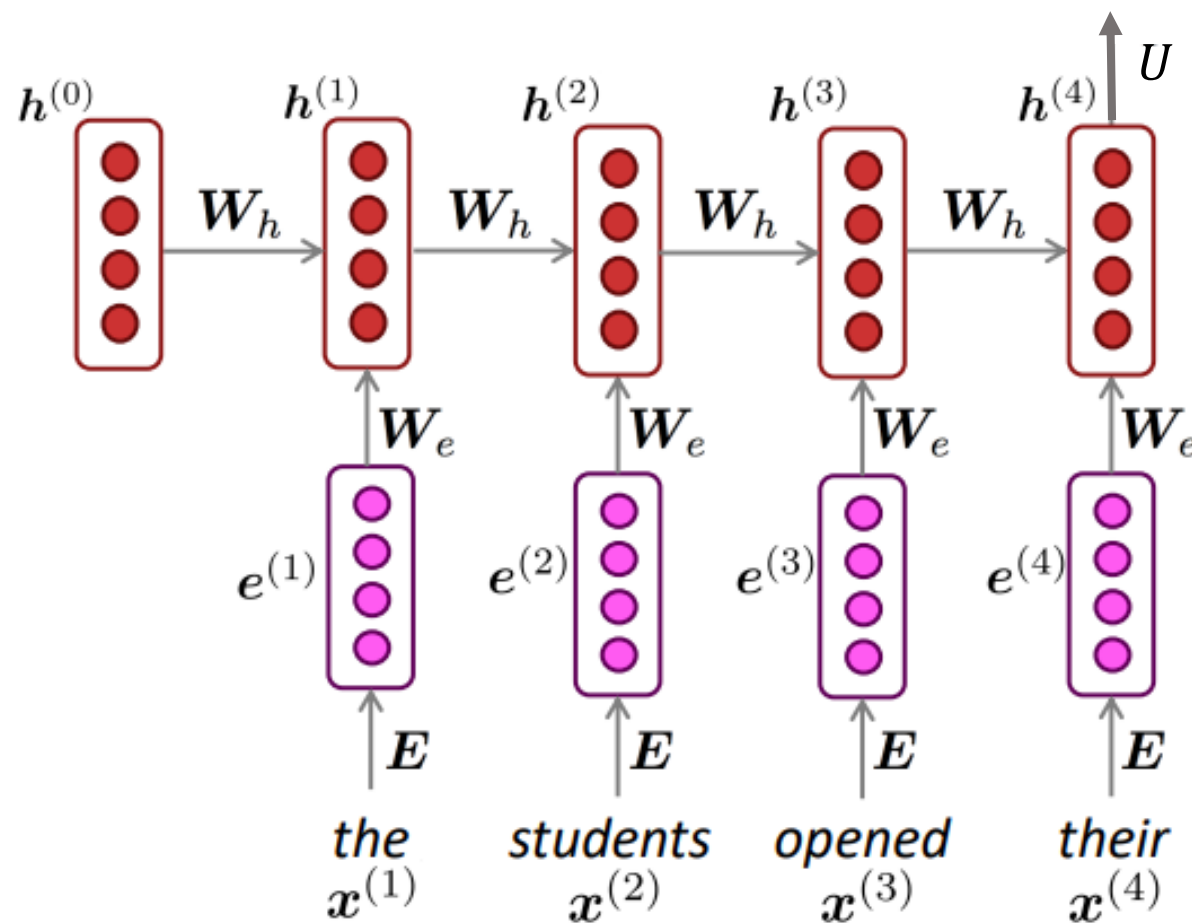
$h^{(0)}$  is the initial hidden state

word embeddings

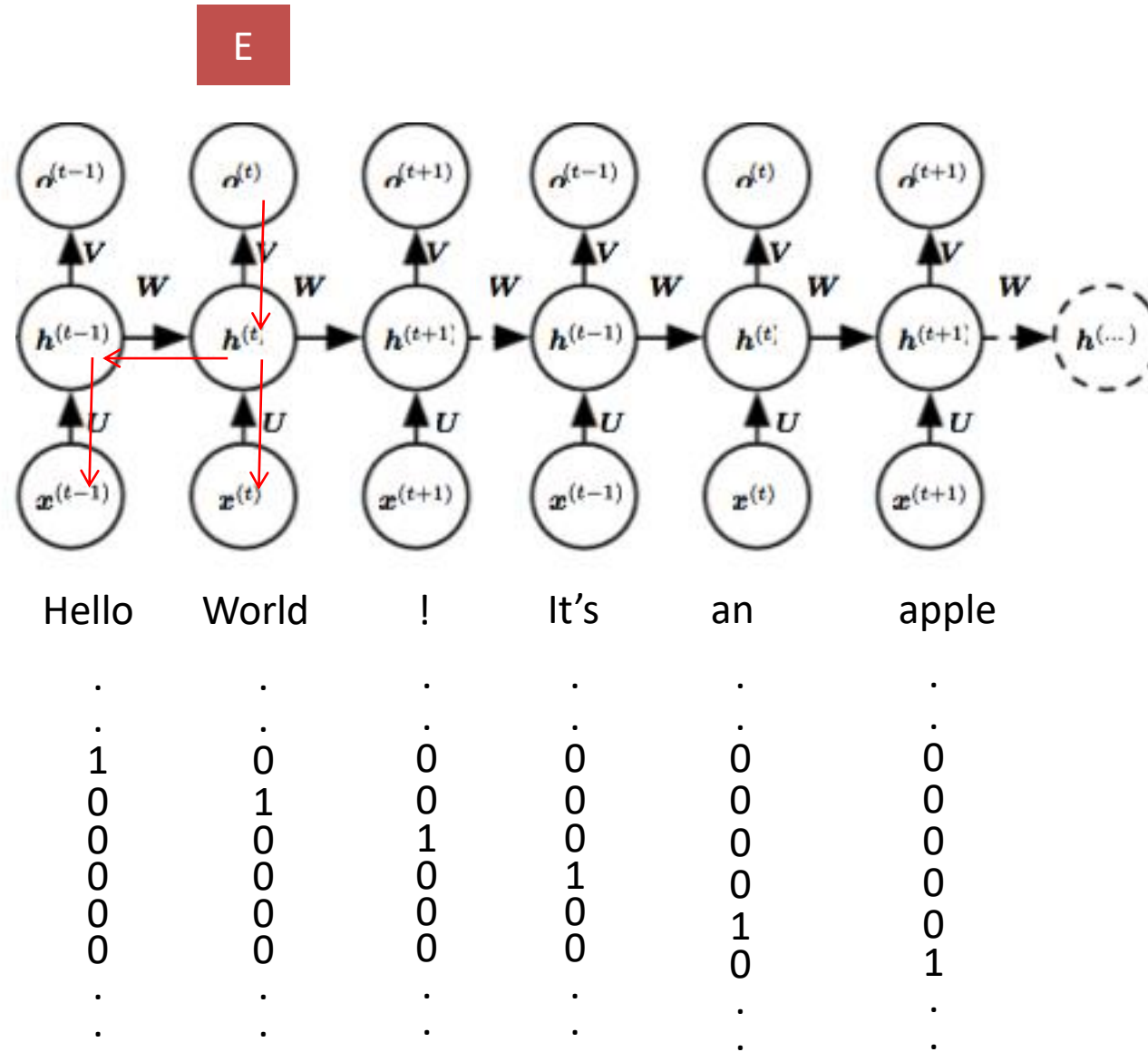
$$e^{(t)} = E x^{(t)}$$

words / one-hot vectors

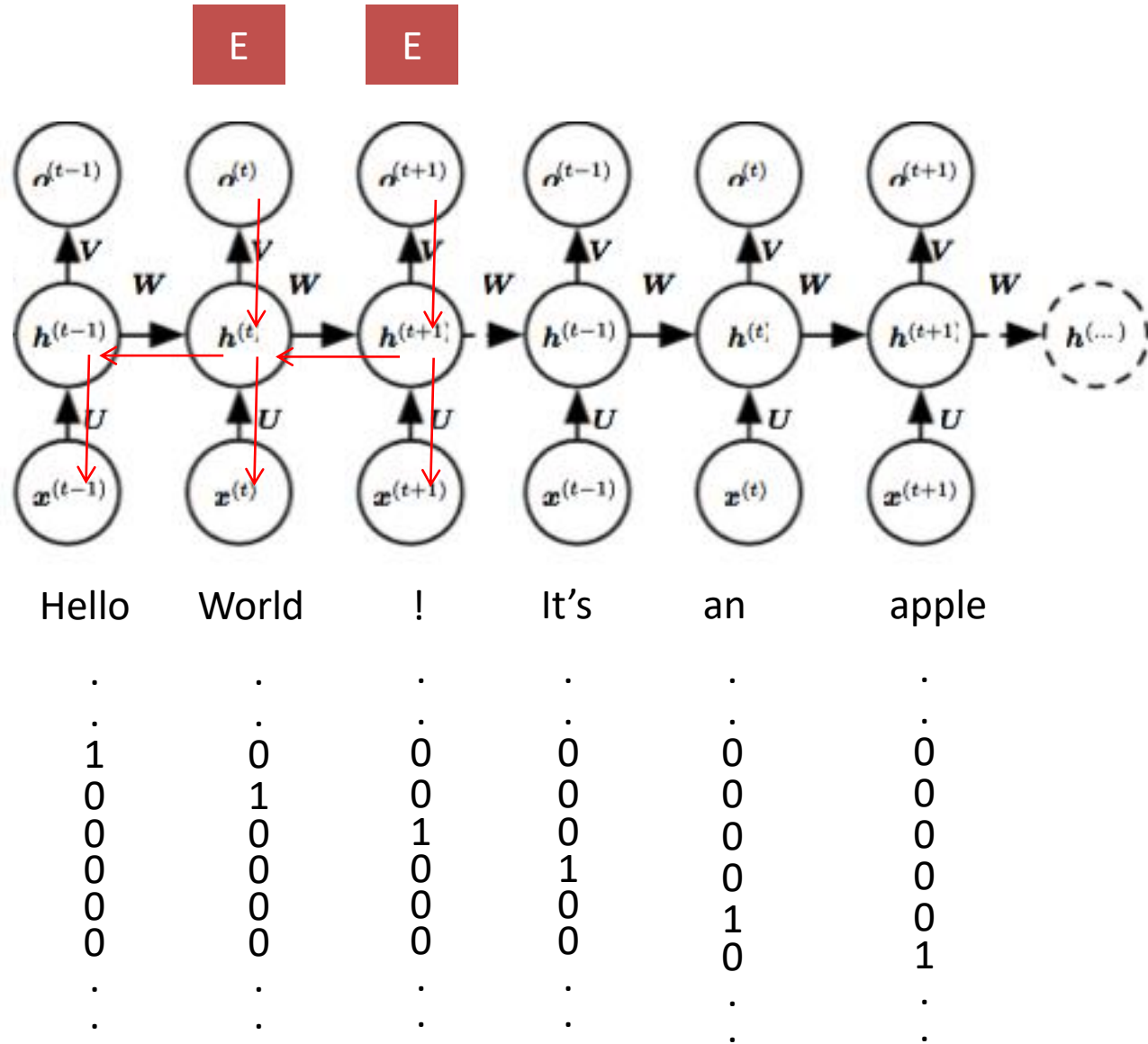
$$x^{(t)} \in \mathbb{R}^{|V|}$$



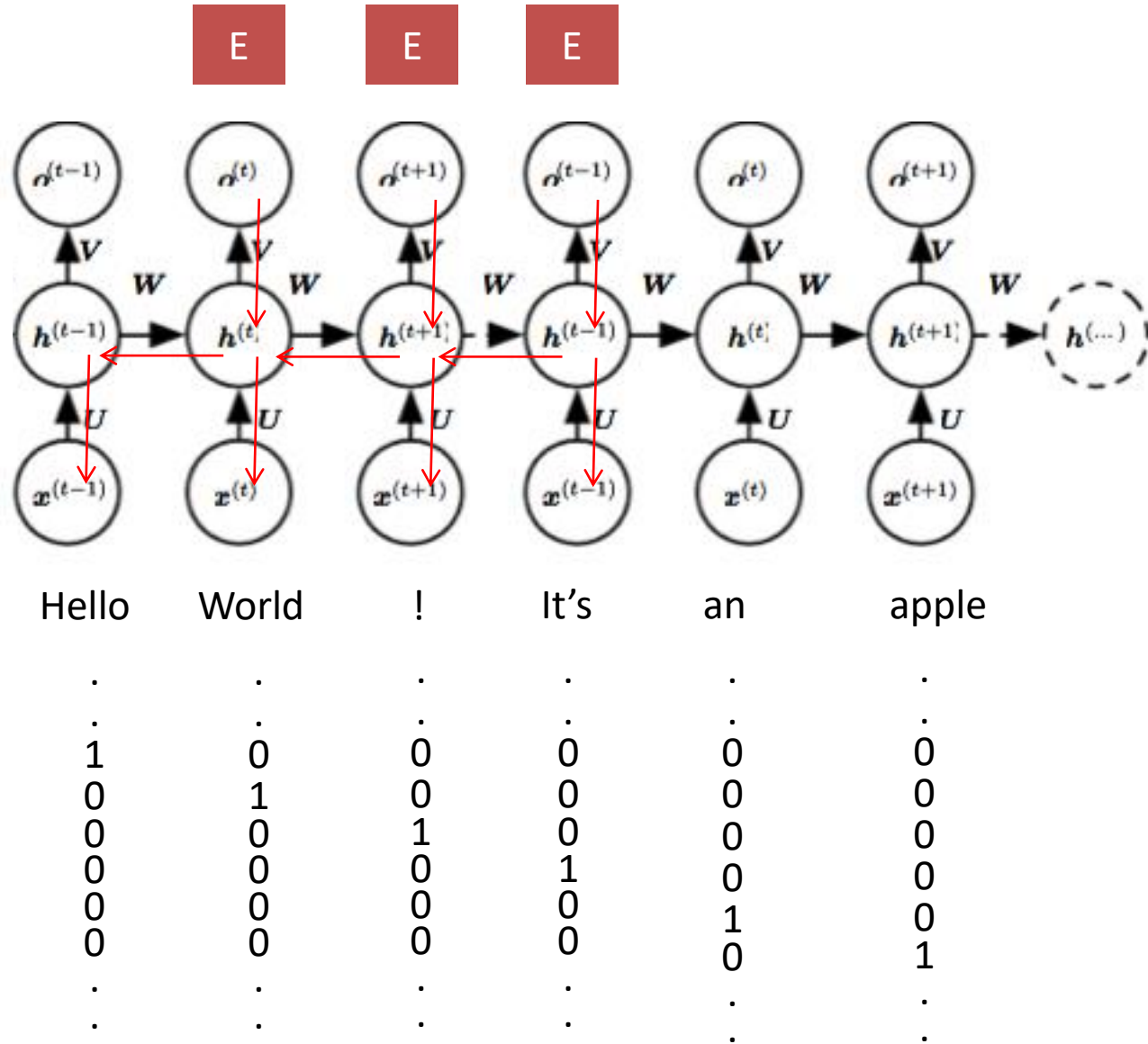
# Process – Next word prediction



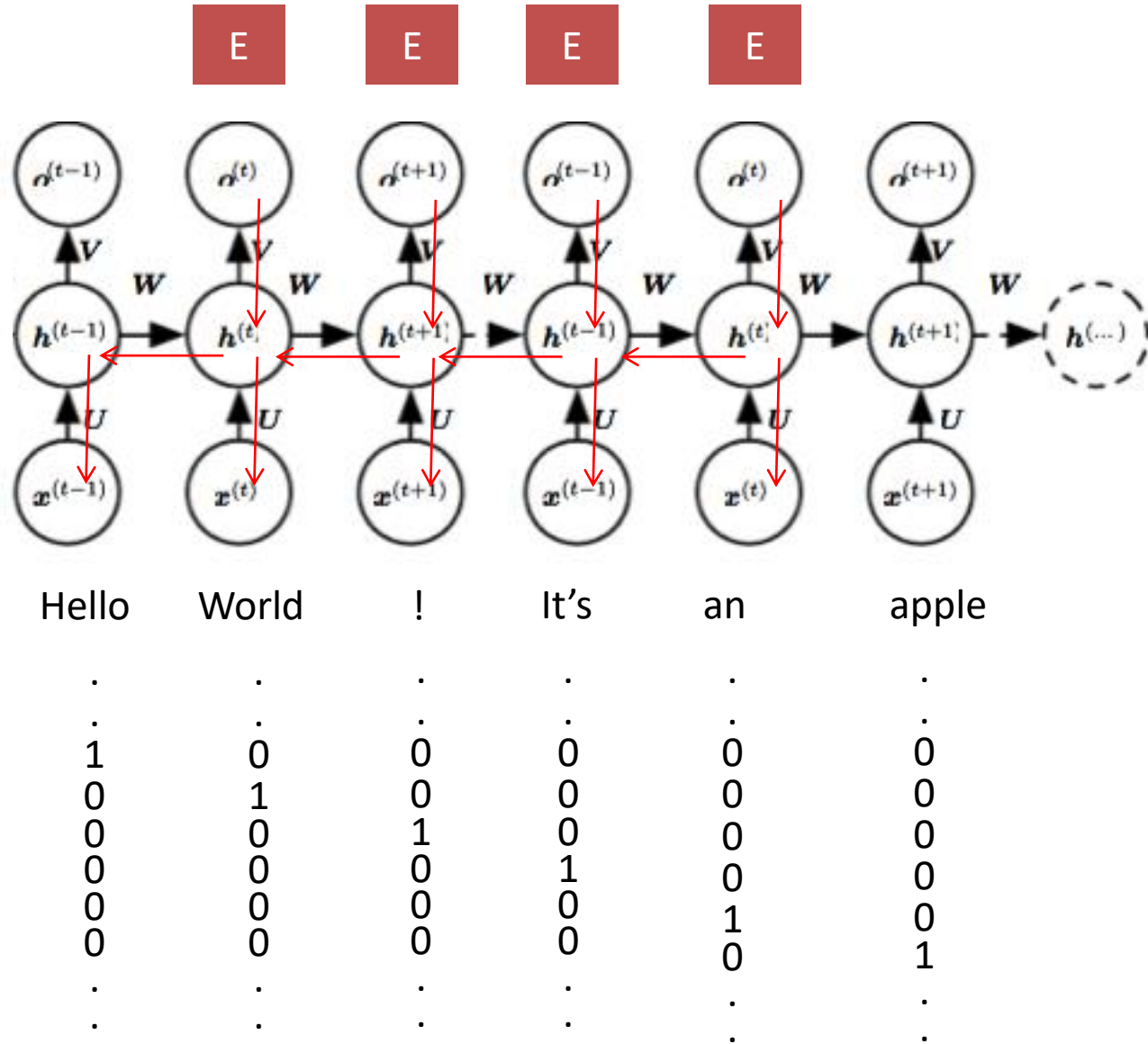
# Process – Next word prediction



# Process – Next word prediction

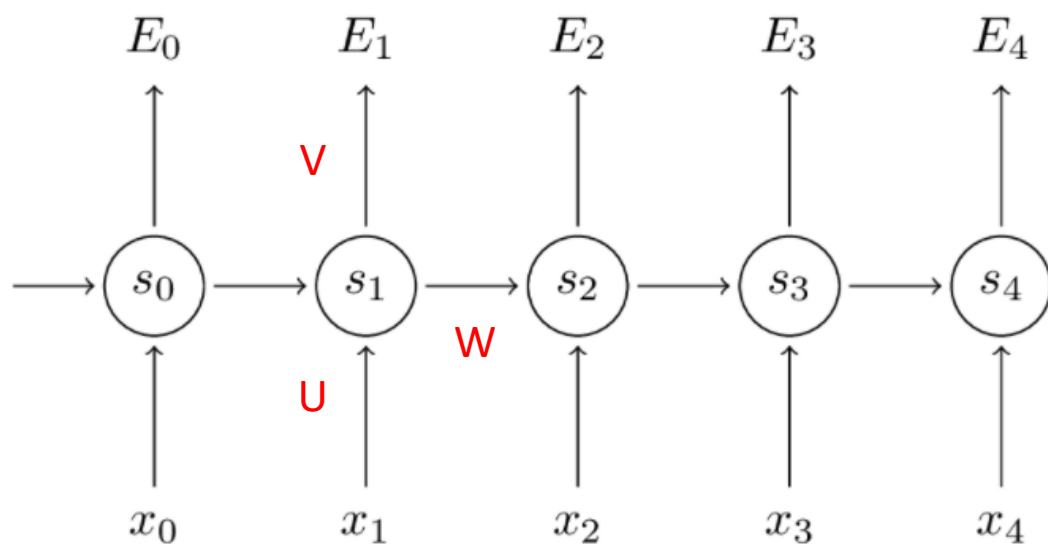


# Process – Next word prediction





# Model



$$s_t = \tanh(Ux_t + Ws_{t-1})$$

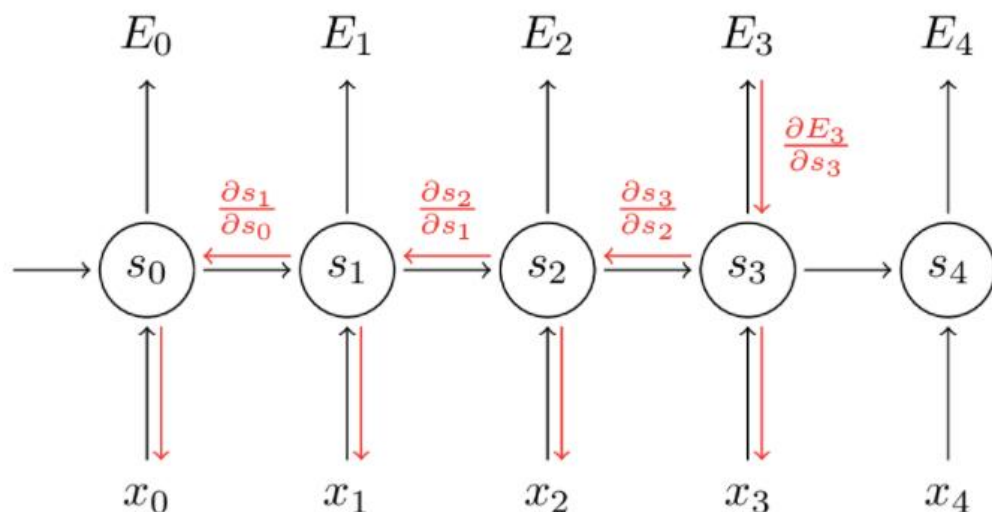
$$\hat{y}_t = \text{softmax}(Vs_t)$$

$$E(y_t, \hat{y}_t) = -y_t \log \hat{y}_t$$

$$E(y, \hat{y}) = - \sum_t E_t(y_t, \hat{y}_t)$$

$$= - \sum_t -y_t \log \hat{y}_t$$

# Learning



$$\begin{aligned}\frac{\partial E_3}{\partial V} &= \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial V} \\ &= \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial z_3} \frac{\partial z_3}{\partial V} \\ &= (\hat{y}_3 - y_3) \otimes s_3\end{aligned}$$

$$\frac{\partial E_3}{\partial W} = \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \frac{\partial s_3}{\partial W}$$

$$\frac{\partial E_3}{\partial W} = \sum_{k=0}^3 \frac{\partial E_3}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial s_3} \frac{\partial s_3}{\partial s_k} \frac{\partial s_k}{\partial W}$$

# RNN Applications

## Automatic Text Generation

PANDARUS:

Alas, I think he shall be come approached and the day  
When little strain would be attain'd into being never fed,  
And who is but a chain and subjects of his death,  
I should not sleep.

Second Senator:

They are away this miseries, produced upon my soul,  
Breaking and strongly should be buried, when I perish  
The earth and thoughts of many states.

DUKE VINCENTIO:

Well, your wit is in the care of side and that.

Second Lord:

They would be ruled after this chamber, and  
my fair nudes begun out of the fact, to be conveyed,  
Whose noble souls I'll have the heart of the wars.

Clown:

Come, sir, I will make did behold your worship.

VIOLA:

I'll drink it.

# RNN Applications

## Automatic Image Caption Generation



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"girl in pink dress is jumping in air."

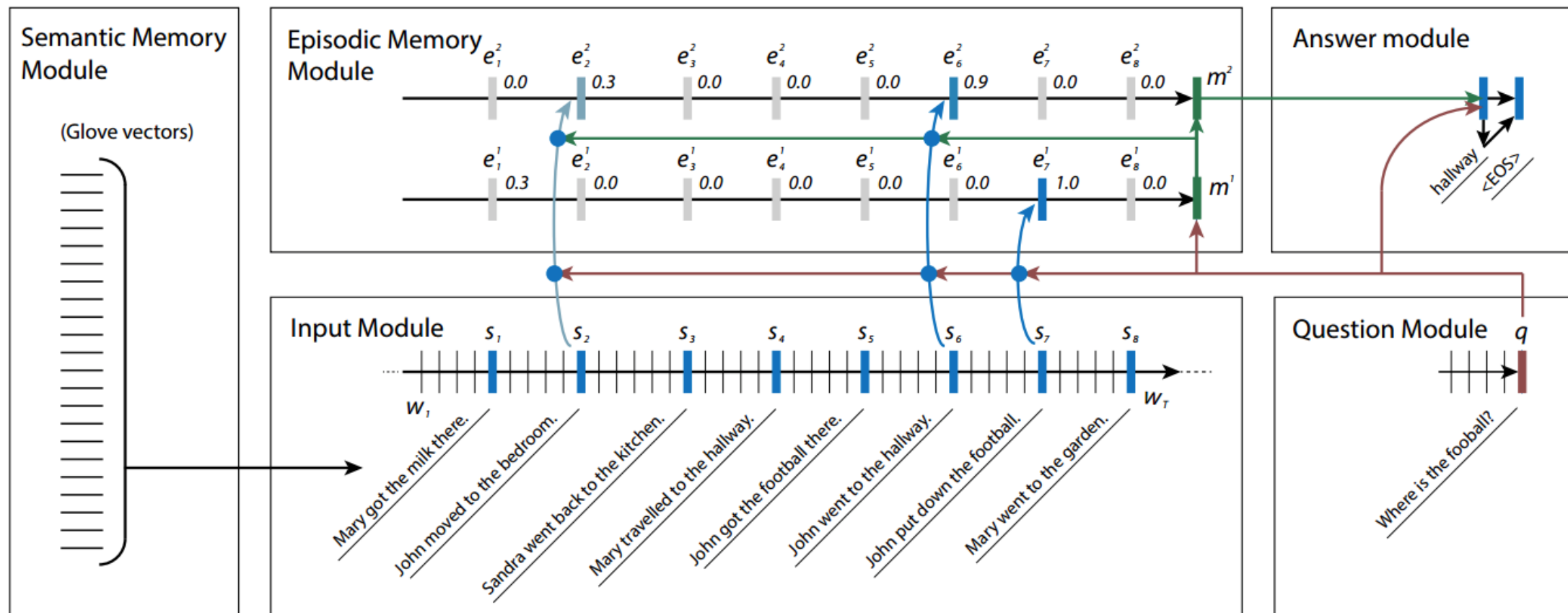


"black and white dog jumps over bar."



"young girl in pink shirt is swinging on swing."

# RNN Applications– Question Answering



# RNN Applications – Machine Translation

