# Deep Learning
## Do machines know the meaning of a word?

Hung-yi Lee

# Language Technology

**_spam detection_**



(http://spam-filter-review.toptenreviews.com/)

**_Part-of-speech Tagging_**

John   saw   the   saw.

↓         ↓         ↓        ↓

PN      V        D        N

**_Name Entity Recognition_**

這　位　是　李　宏　毅

Name of People

**_Sentiment Analysis_**

這部電影太糟了

Negative (負雷)

**_Translation_**

"Machine learning ……"

⬍

"機器學習 ……"

**_Summarization_**



document        summary

**_Retrieval_**



**_Speech Recognition_**



大家好……

**_Syntactic Analysis_**

# Do machine really understand human language?



http://cse3521.artifice.cc/chinese-room.html

# Meaning Representation

Do machine know the meaning of a word or word sequence?

# Meaning of Word

# Predicting the next word
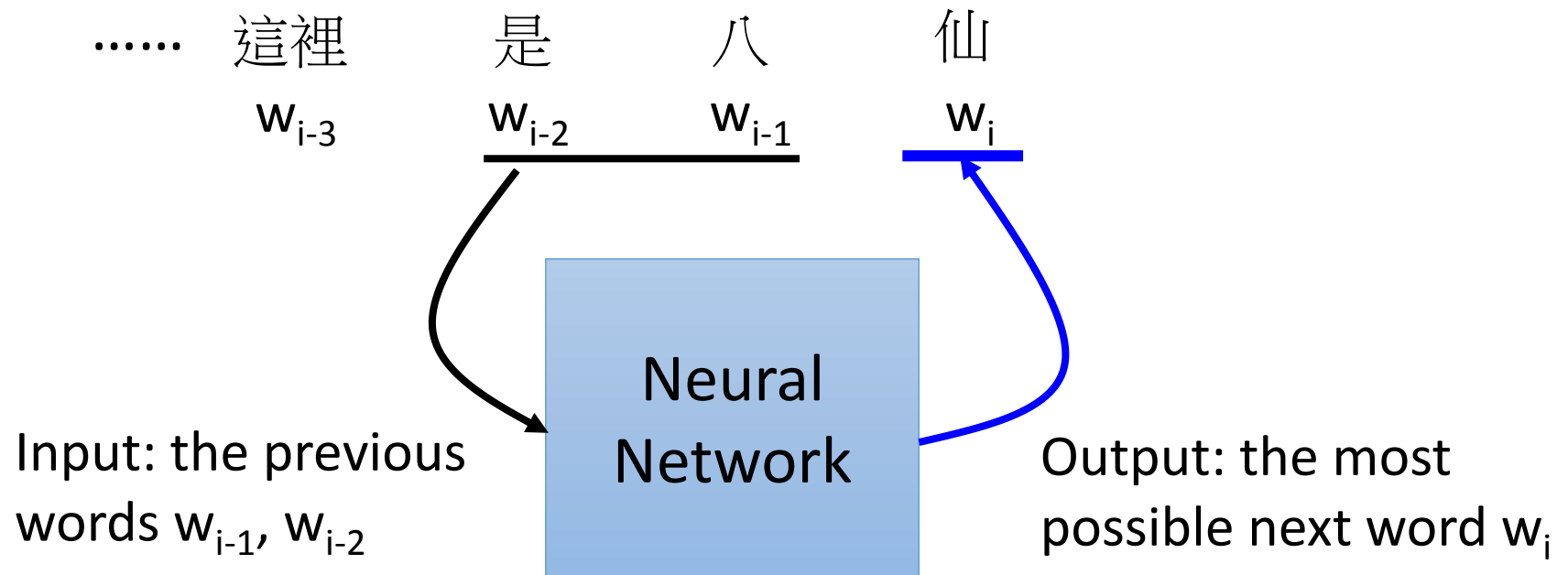
- Given a sequence of words, predict the next word

麻煩這系列的請到政黑或其他地方討論好嗎?這裡是八          04/27 00:40

# Predicting the next word

- Given a sequence of words, predict the next word

...... 這裡　　是　　八　　仙

$w_{i-3}$　　$w_{i-2}$　　$w_{i-1}$　　$w_i$

Neural Network

Input: the previous words $w_{i-1}$, $w_{i-2}$

Output: the most possible next word $w_i$

Each word should be represented as a feature vector.

# Predicting the next word

**_1-of-N Encoding_**

lexicon = {apple, bag, cat, dog, elephant}

apple = [ 1  0  0  0  0]     The vector is lexicon size.

bag    = [ 0  1  0  0  0]
                                    Each dimension corresponds
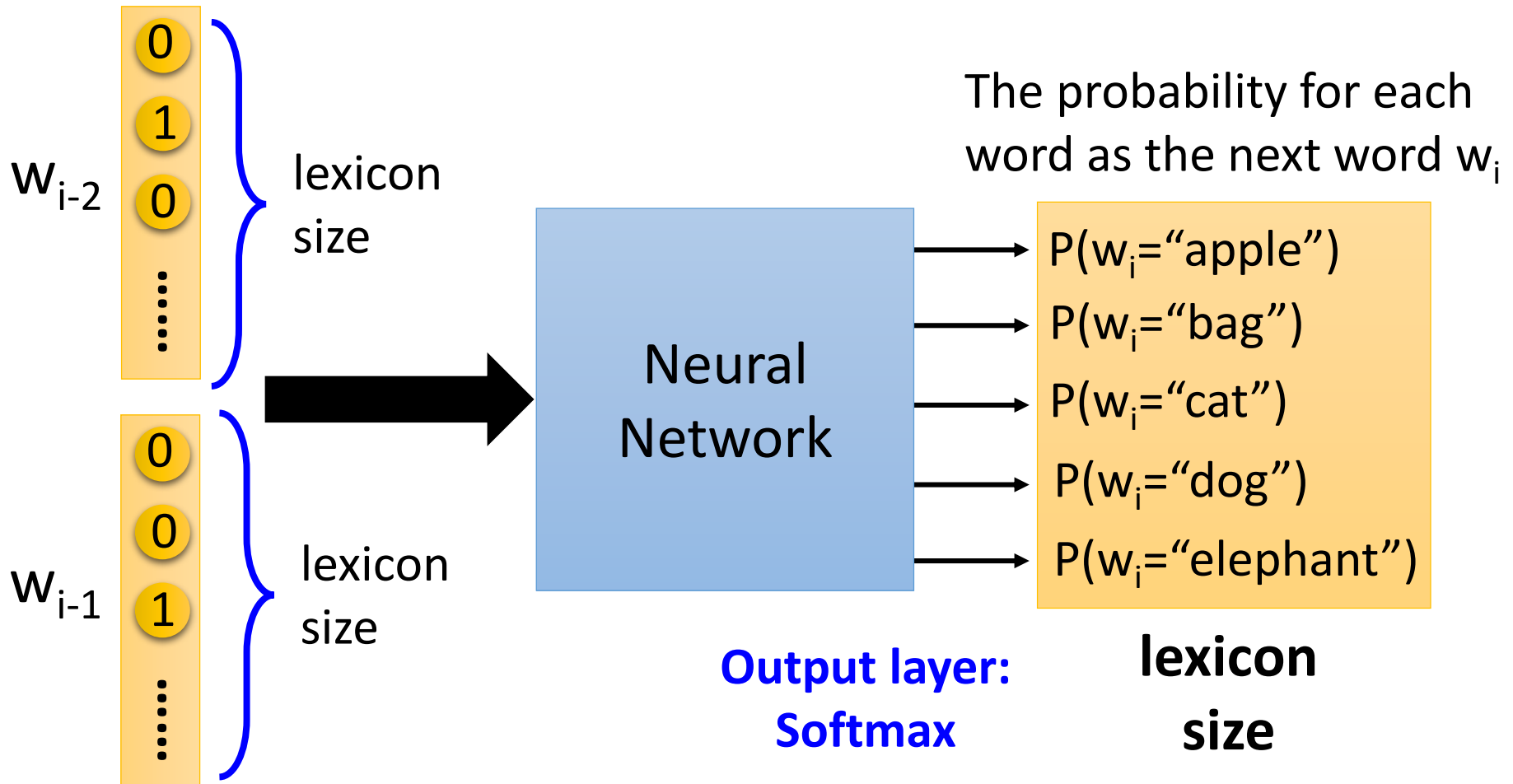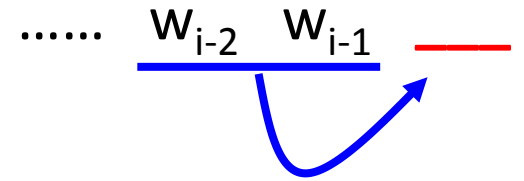cat    = [ 0  0  1  0  0]     to a word in the lexicon

dog    = [ 0  0  0  1  0]
                                    The dimension for the word
elephant  = [ 0  0  0  0  1]     is 1, and others are 0

# Predicting the next word

$\cdots\cdots \quad \dfrac{w_{i-2} \quad w_{i-1}}{\qquad} \quad \text{\textcolor{red}{\rule{1cm}{0.4pt}}}$

$w_{i-2}$

| 0 |
| 1 |
| 0 |
| ⋮ |

lexicon size

$w_{i-1}$

| 0 |
| 0 |
| 1 |
| ⋮ |

lexicon size

Neural Network

**Output layer: Softmax**

The probability for each word as the next word $w_i$

$P(w_i=\text{"apple"})$

$P(w_i=\text{"bag"})$

$P(w_i=\text{"cat"})$

$P(w_i=\text{"dog"})$

$P(w_i=\text{"elephant"})$

**lexicon size**

# Predicting the next word

- Training:

Collect data:

這裡　是　八　仙　樂園

………

………

………

**Minimizing cross entropy**



這裡 是 → Neural Network → 八

是 八 → Neural Network → 仙

八 仙 → Neural Network → 樂園

# Word Vector

1-of-N encoding of the word $w_{i-1}$

$z_1$
$z_2$

The probability for each word as the next word $w_i$
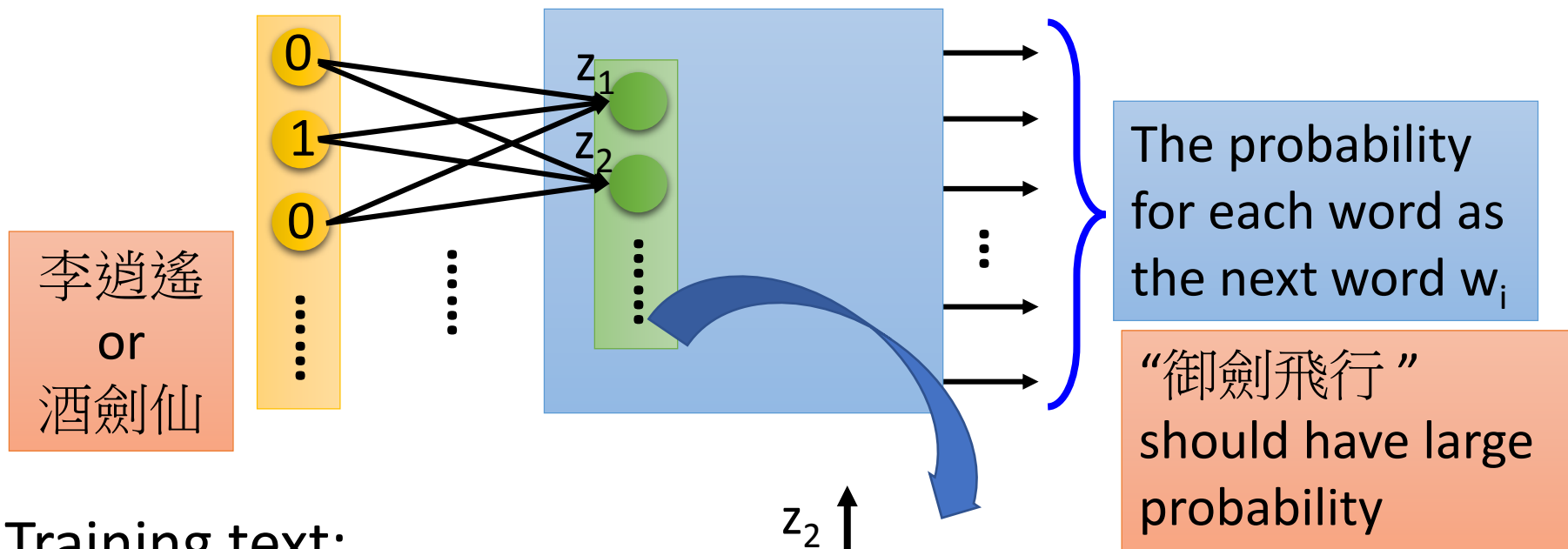
➤ Take out the input of the neurons in the first layer

➤ Use it to represent a word w

➤ Word vector, word embedding feature: V(w)

$z_2$

tree
flower

dog        rabbit

cat

run
jump

$z_1$

# Word Vector

You shall know a word by the company it keeps

$z_1$
$z_2$

The probability for each word as the next word $w_i$

"御劍飛行" should have large probability
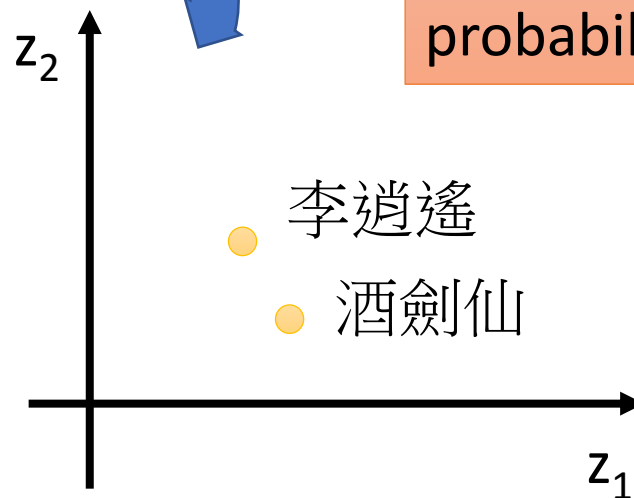
李逍遙 or 酒劍仙

Training text:

...... 李逍遙 御劍飛行 ......
$w_{i-1}$     $w_i$

...... 酒劍仙 御劍飛行 ......
$w_{i-1}$     $w_i$

$z_2$

李逍遙
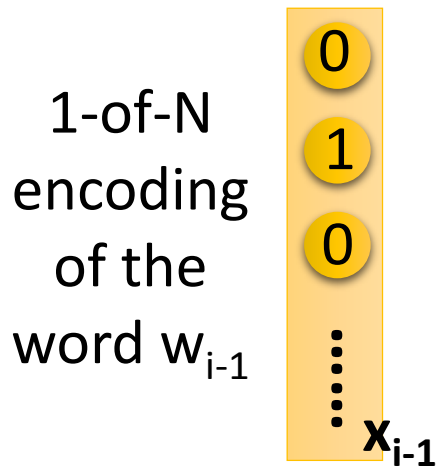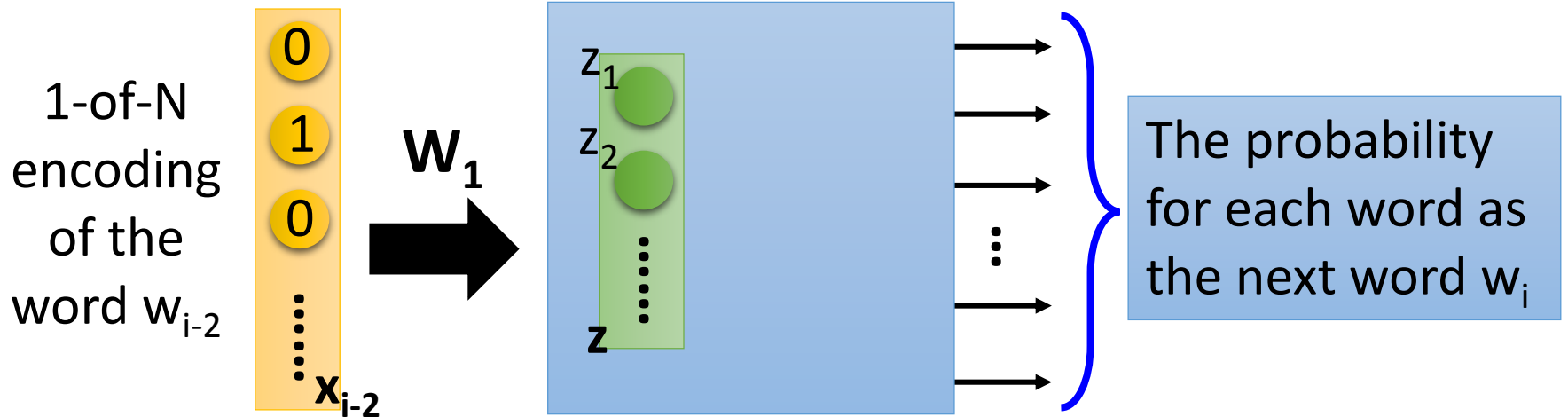酒劍仙

$z_1$

# Word Vector – Sharing Parameters



1-of-N encoding of the word $w_{i-2}$

1-of-N encoding of the word $w_{i-1}$

$z_1$

$z_2$

The probability for each word as the next word $w_i$

The weights with the same color should be the same.

Or, one word would have two word vectors.

# Word Vector – Sharing Parameters

1-of-N encoding of the word $w_{i-2}$

$W_1$

$z_1$
$z_2$
...
$z$

$x_{i-2}$

The probability for each word as the next word $w_i$

1-of-N encoding of the word $w_{i-1}$

$W_2$

$x_{i-1}$

The length of $x_{i-1}$ and $x_{i-2}$ are both $|V|$.

The length of $z$ is $|Z|$.

$z = W_1 x_{i-2} + W_2 x_{i-1}$

The weight matrix $W_1$ and $W_2$ are both $|Z| \times |V|$ matrices.

$W_1 = W_2 = W$ ➡ $z = W ( x_{i-2} + x_{i-1} )$

# Word Vector – Sharing Parameters

1-of-N encoding of the word $w_{i-2}$

$w_i$

1-of-N encoding of the word $w_{i-1}$

$w_j$

$z_1$
$z_2$

The probability for each word as the next word $w_i$

How to make $w_i$ equal to $w_j$

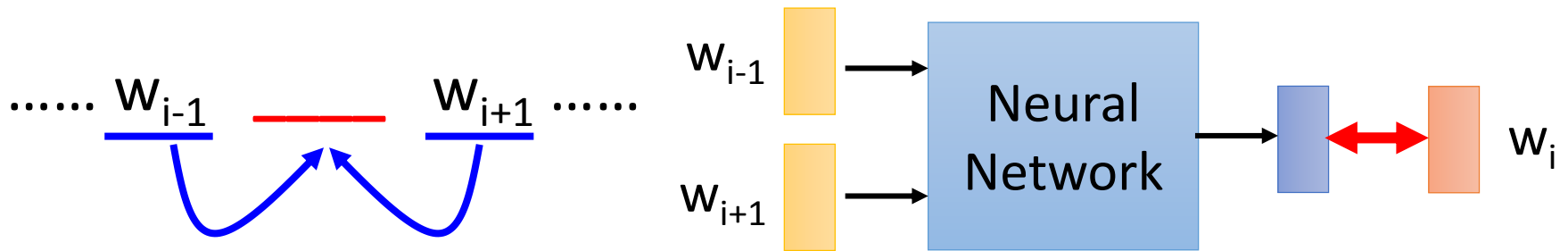Given $w_i$ and $w_j$ the same initialization

$$w_i \leftarrow w_i - \eta \frac{\partial C}{\partial w_i} - \eta \frac{\partial C}{\partial w_j}$$

$$w_j \leftarrow w_j - \eta \frac{\partial C}{\partial w_j} - \eta \frac{\partial C}{\partial w_i}$$
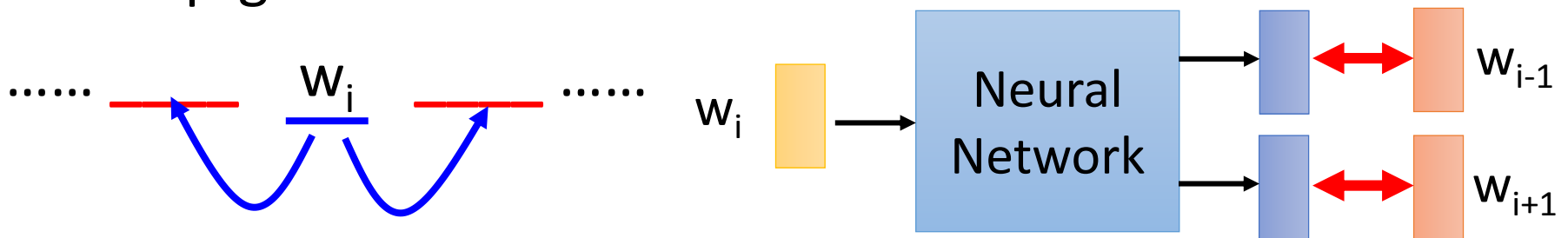
# Word Vector – Various Architectures

- Continuous bad of word (CBOW) model
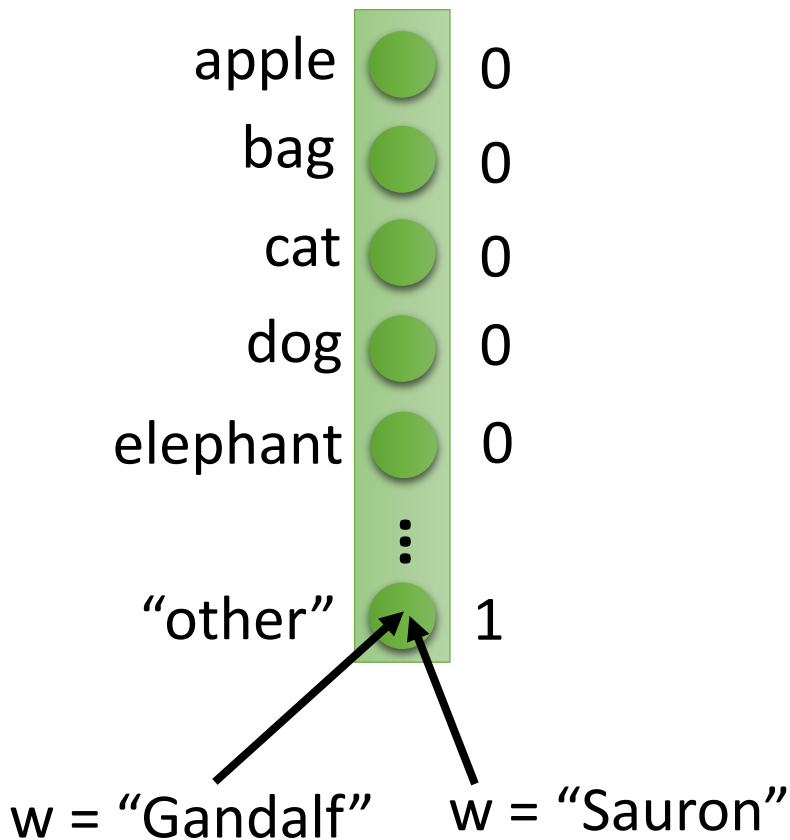


*predicting the word given its context*

- Skip-gram



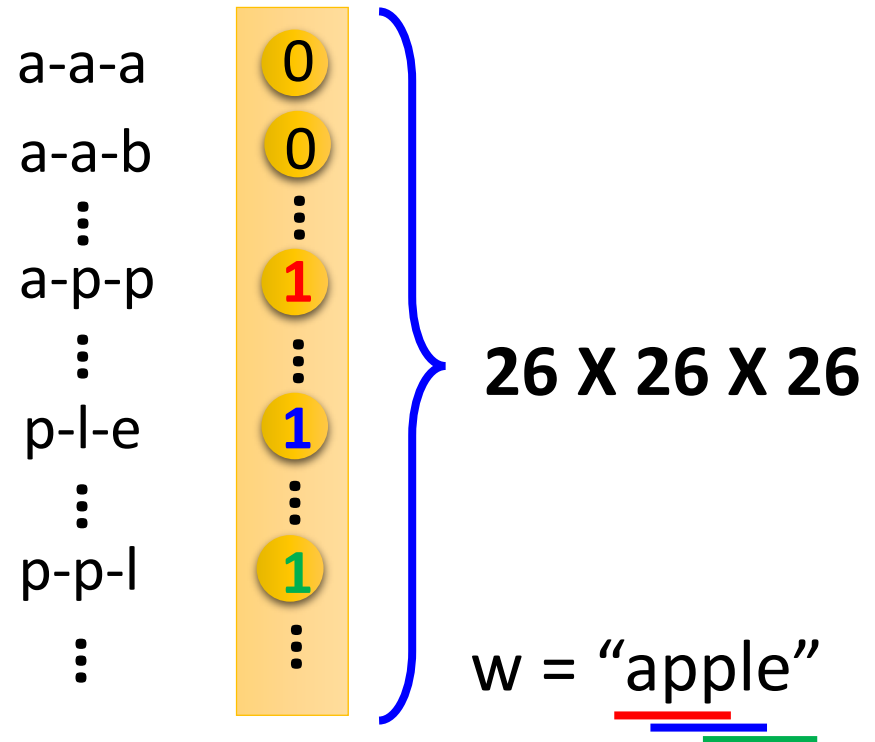*predicting the context given a word*
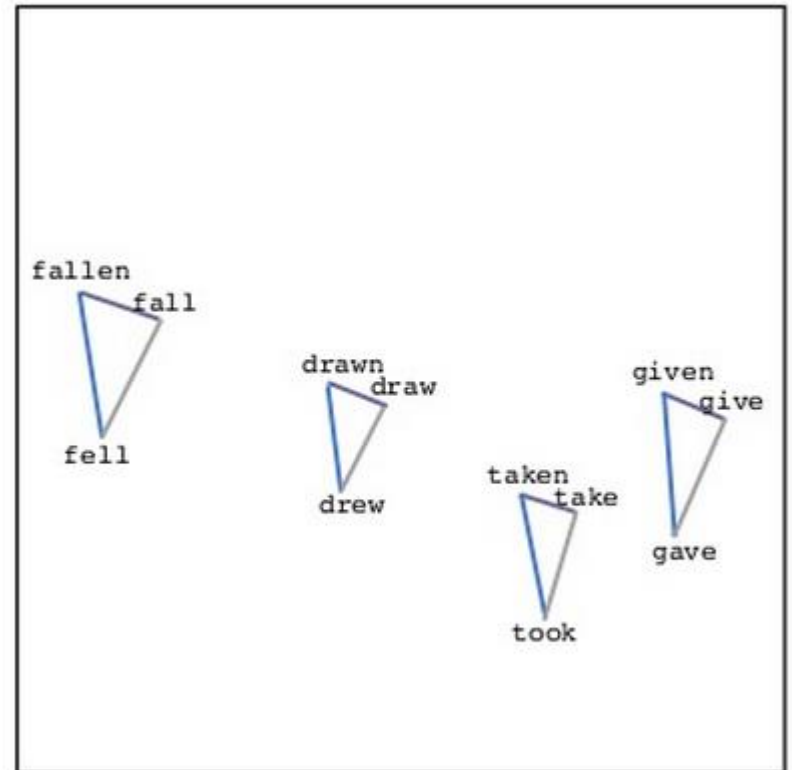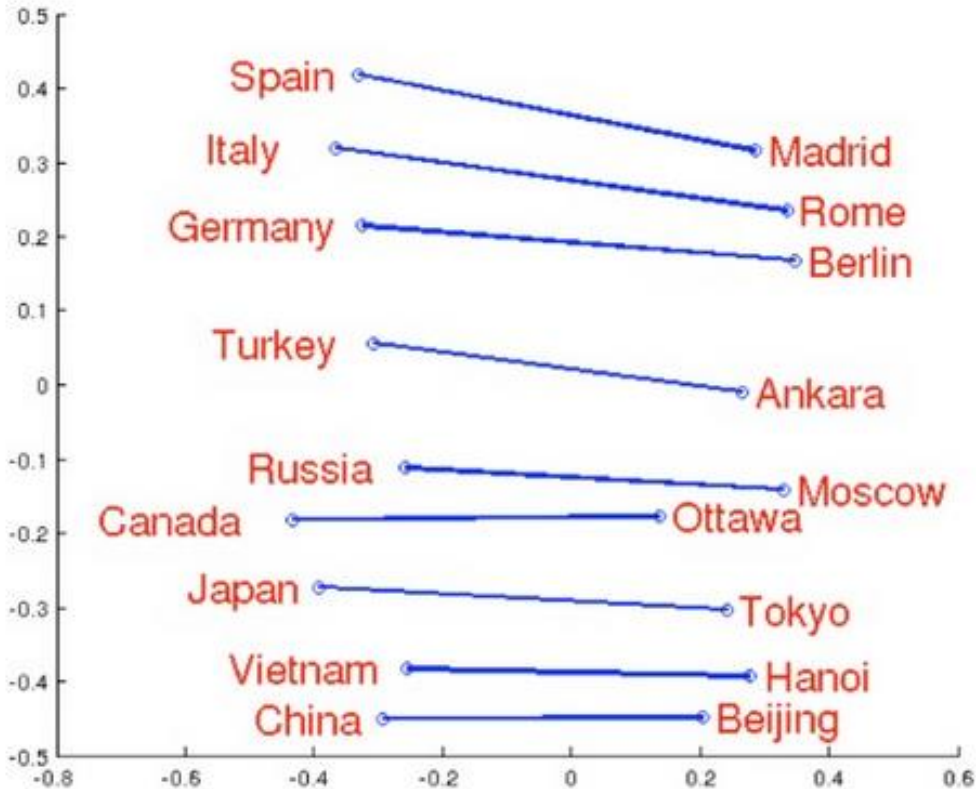
# Beyond 1-of-N encoding

**_Dimension for "Other"_**

| | | |
|---|---|---|
| apple | ● | 0 |
| bag | ● | 0 |
| cat | ● | 0 |
| dog | ● | 0 |
| elephant | ● | 0 |
| ⋮ | ⋮ | |
| "other" | ● | 1 |

w = "Gandalf"   w = "Sauron"

**_Word hashing_**

| | |
|---|---|
| a-a-a | 0 |
| a-a-b | 0 |
| ⋮ | ⋮ |
| a-p-p | 1 |
| ⋮ | ⋮ |
| p-l-e | 1 |
| ⋮ | ⋮ |
| p-p-l | 1 |
| ⋮ | |

26 X 26 X 26

w = "apple"

# Word Vector



Source: http://www.slideshare.net/hustwj/cikm-keynotenov2014

# Word Vector

$$V(Germany)$$
$$\approx V(Berlin) - V(Rome) + V(Italy)$$

- Characteristics

$$V(hotter) - V(hot) \approx V(bigger) - V(big)$$
$$V(Rome) - V(Italy) \approx V(Berlin) - V(Germany)$$
$$V(king) - V(queen) \approx V(uncle) - V(aunt)$$

- Solving analogies

Rome : Italy = Berlin : ?

Compute $V(Berlin) - V(Rome) + V(Italy)$

Find the word w with the closest V(w)

# Demo

- Model used in demo is provided by 陳仰德
  - Part of the project done by 陳仰德、林資偉
  - TA: 劉元銘
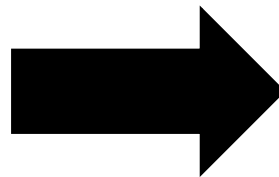  - Training data is from PTT (collected by 葉青峰)

# Meaning of Word Sequence

# Meaning of Word Sequence

- word sequences with different lengths → the vector with the same length
  - The vector representing the meaning of the word sequence
  - A word sequence can be a document or a paragraph

**word sequence**
(a document or paragraph)

# Meaning of Word Sequence
- Outline

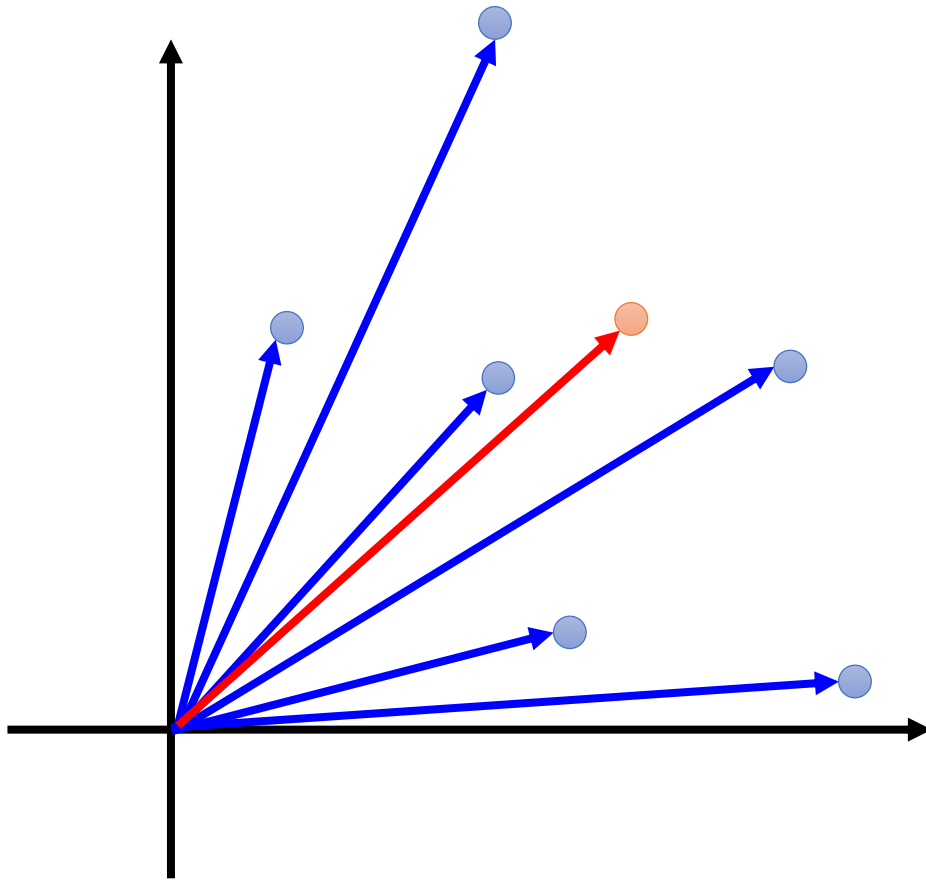| Deep Structured Semantic Model (DSSM) | • Application: Information Retrieval (IR) |
| Recursive Deep Model | • Application: Sentiment Analysis |
| Paragraph Vector | • Unsupervised |

Reference: http://www.msr-waypoint.net/pubs/198202/cikm2013_DSSM_fullversion.pdf

# Information Retrieval (IR)



## Vector Space Model

The documents are vectors in the space.

The query is also a vector.

How to use a vector to represent word sequences

# Information Retrieval (IR)

### *Bag-of-word*

| | | |
|---|---|---|
| this | | 1 |
| is | | 1 |
| word string s1:    a | | 0 |
| "This is an apple"    an | | 1 |
| apple | | 1 |
| pen | | 0 |
| ⋮ | | |

| | | |
|---|---|---|
| this | | 1 |
| is | | 1 |
| word string s2:    a | | 1 |
| "This is a pen"    an | | 0 |
| apple | | 0 |
| pen | | 1 |
| ⋮ | | |

Weighted by IDF

# Information Retrieval (IR)

## *Vector Space Model + Bag-of-word*

Retrieved

Bag-of-word

Query q

Document $d_1$

Document $d_2$

All documents in the database

➢ All the words are treated as discrete tokens

➢ Never consider different words can have the same meaning, and the same word can have different meanings

# IR - Semantic Embedding



Source:
http://www.cs.toronto.edu/~hinton/science.pdf

How to achieve that? (No target ......)

Auto-encoder is one solution, but not today

Bag-of-word

word string
(document or query)

# DSSM

Click-through data:  $q_1 \longrightarrow d_1 : +$   $d_2 : -$

$q_2 \longrightarrow d_3 : -$   $d_4 : +$

......

Training:



query $q_1$      document $d_1$ +      document $d_2$ -

query $q_2$      document $d_3$ -      document $d_4$ +

# DSSM v.s. Typical DNN

**_Typical DNN_**

**_DSSM_**



reference

input

query q

document d **+**

query q

document d **+**

Click-through data:   $q_1$ ⟶ $d_1$ : **+**   $d_2$ : **-**

$q_2$ ⟶ $d_3$ : **-**   $d_4$ : **+**

……

- How to do retrieval?

Retrieved

New Query q'          Document $d_1$          Document $d_2$

# More …

- **Convolutional DSSM**: http://www.iro.umontreal.ca/~lisa/pointeurs/ir0895-he-2.pdf



Take max at each dimension across all word-trigram features

# Meaning of Word Sequence - Outline

| | |
|---|---|
| **Deep Structured Semantic Model (DSSM)** | • Application: Information Retrieval (IR) |
| **Recursive Deep Model** | • Application: Sentiment Analysis |
| **Paragraph Vector** | • Unsupervised |

Reference: http://nlp.stanford.edu/~socherr/EMNLP2013_RNTN.pdf
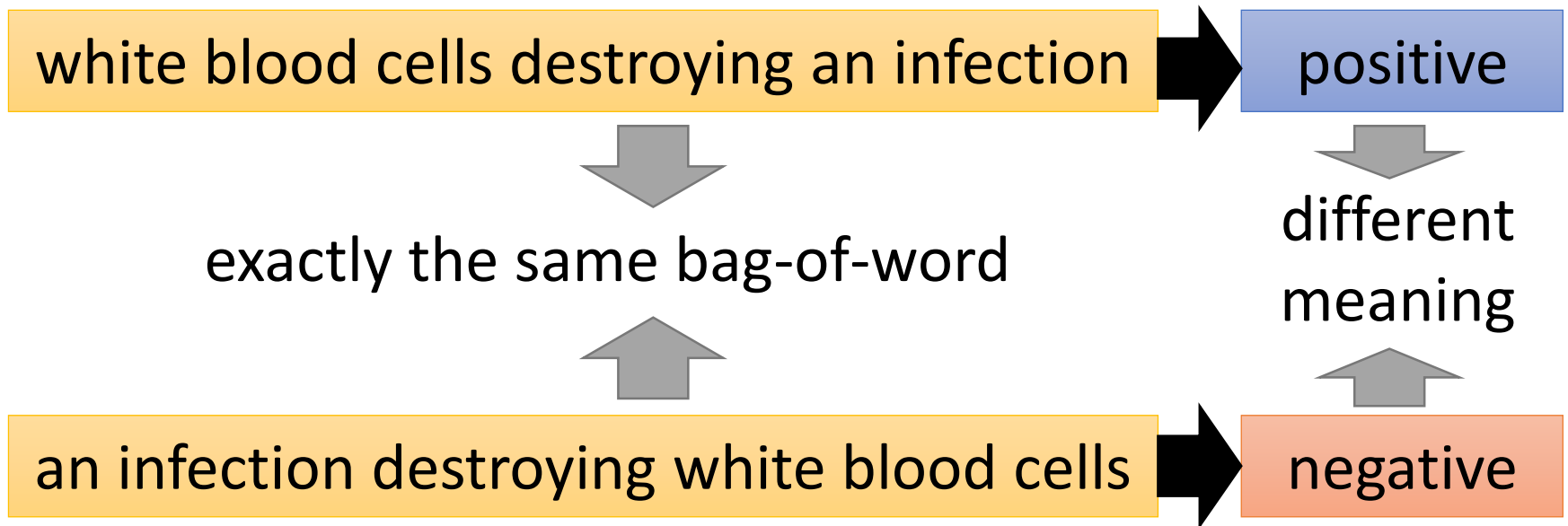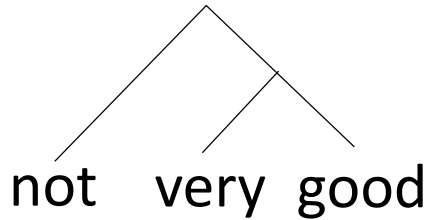
# Recursive Deep Model

- When understanding the meaning of a word sequence, the order of the words can not be ignore.

| white blood cells destroying an infection | ➡ | positive |

exactly the same bag-of-word

different meaning

| an infection destroying white blood cells | ➡ | negative |

# Recursive Deep Model

syntactic structure

How to do it is out
of the scope

not   very  good

word sequence:

not                              very                              good

# Recursive Deep Model

syntactic structure

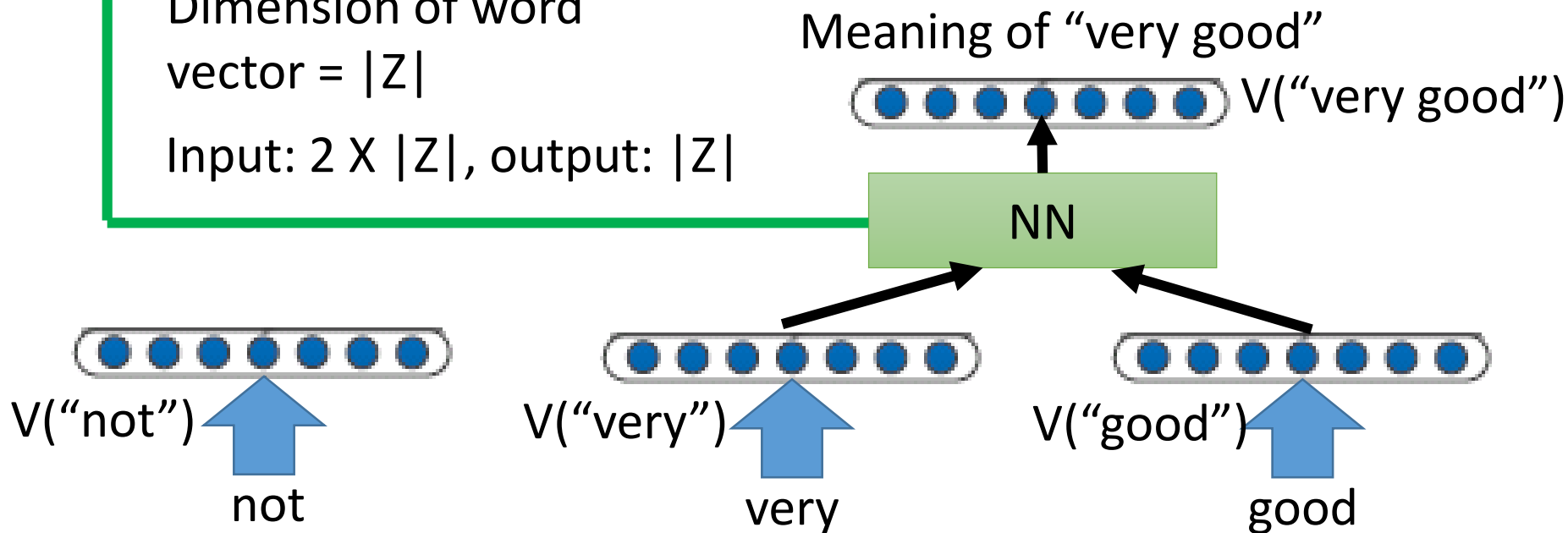not  very  good

By composing the two meaning, what should the meaning be.

Dimension of word vector = |Z|

Input: 2 X |Z|, output: |Z|

Meaning of "very good"

V("very good")

NN

V("not")

not

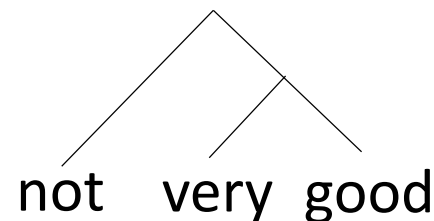V("very")

very

V("good")

good

# Recursive Deep Model
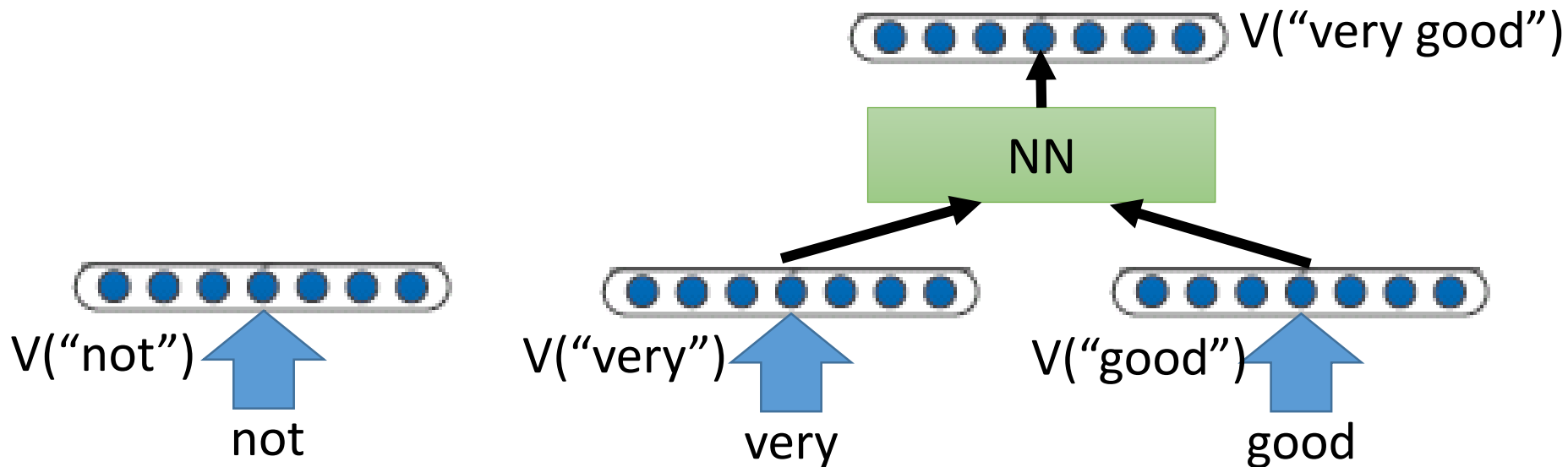
$V(w_A\ w_B) \neq V(w_A) + V(w_B)$

"not": neutral

"good": positive

"not good": negative

syntactic structure

Meaning of "very good"
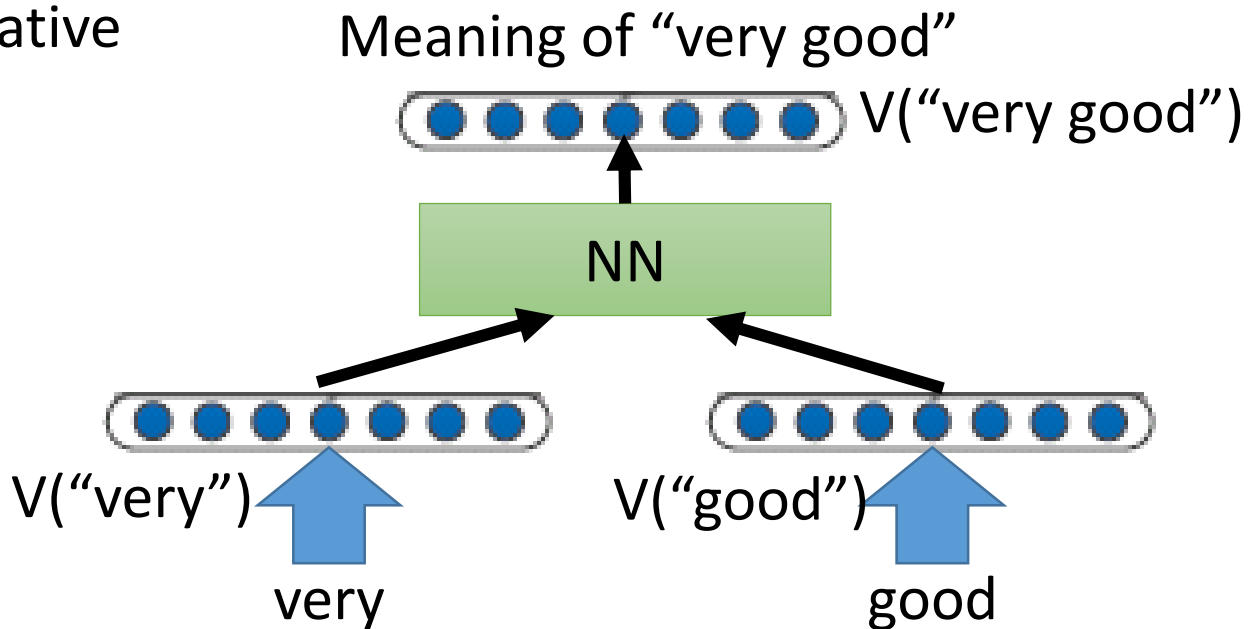
# Recursive Deep Model

$V(w_A \, w_B) \neq V(w_A) + V(w_B)$

"棒": positive

"好棒": positive

"好棒棒": negative

syntactic structure

not  very  good



V("not")

not

Meaning of "very good"

V("very good")

NN

V("very")

very

V("good")

good

# Recursive Deep Model



"not good"

"not bad"

syntactic structure

not     very     good

NN

NN

"not"     "good"

"not"     "bad"

Meaning of "very good"

V("very good")

: "reverse" another input

"not"

NN

V("not")

V("very")

V("good")

not

very

good

# Recursive Deep Model

"very good"



"very bad"



NN

"very"  "good"

NN

"very"  "bad"

syntactic structure

not   very   good

Meaning of "very good"

V("very good")

NN

: "emphasize" another input

"very"

V("not")

not

V("very")

very

V("good")

good

The word order is considered.
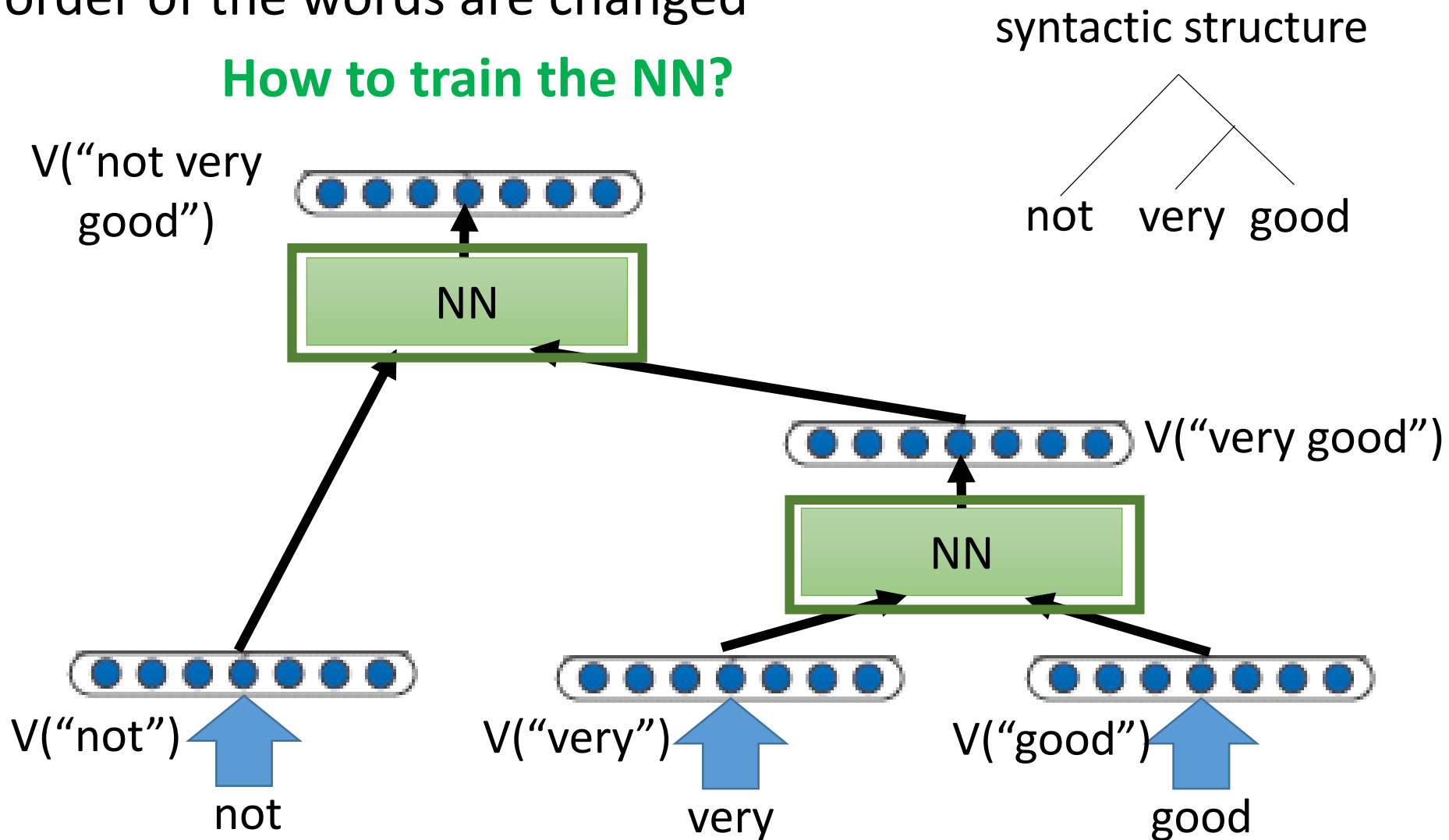
The representation of the sequence will change if the order of the words are changed

syntactic structure

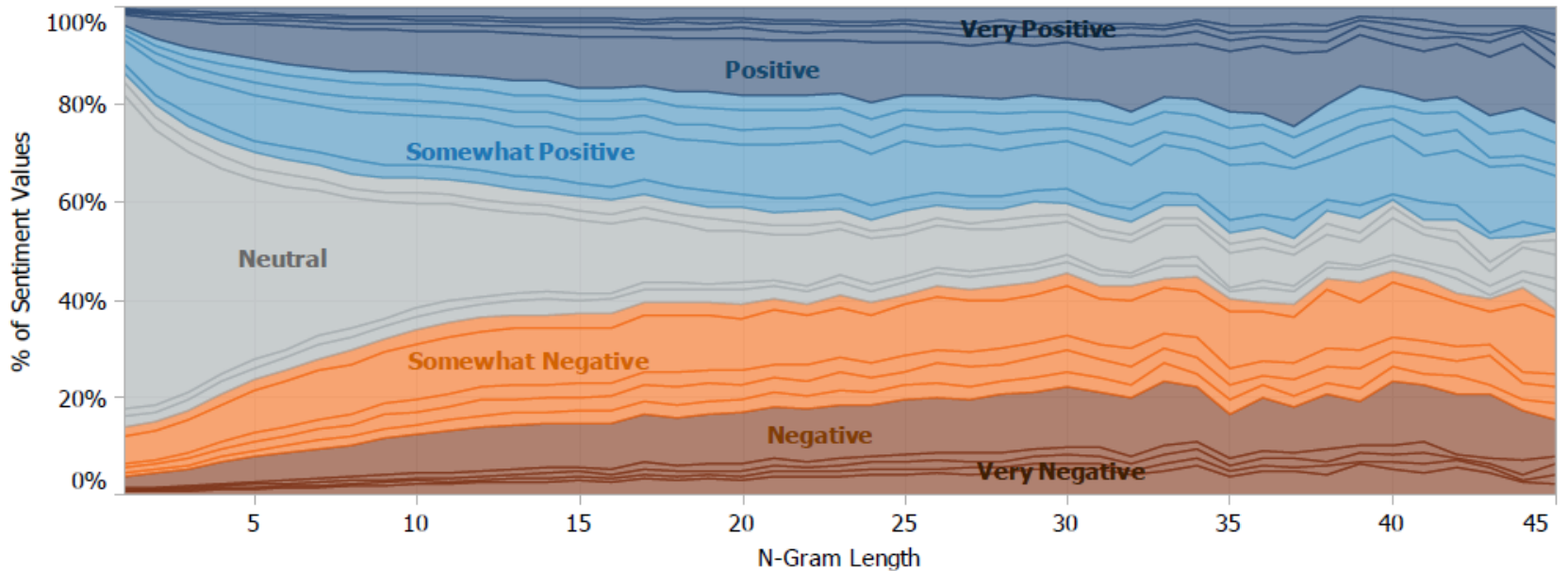**How to train the NN?**

```
     not   very  good
```

V("not very good")

NN

V("very good")

NN

V("not")

not

V("very")

very

V("good")

good

# Training Data



5-class sentiment classification
( -- , - , 0 , + , ++ )

- ref

output

5 classes
( -- , - , 0 , + , ++ )

Train both ...

NN

NN

NN

++ ref

output

NN

0 ref

output

NN

V("not")

not

0 ref

output

NN

V("very")

very

NN

NN

+ ref

output

NN

V("good")

good

# More …

- Demo
  - http://nlp.stanford.edu:8080/sentiment/rntnDemo.html

This course is a little difficult, but it is fun still.

# Meaning of Word Sequence
- Outline

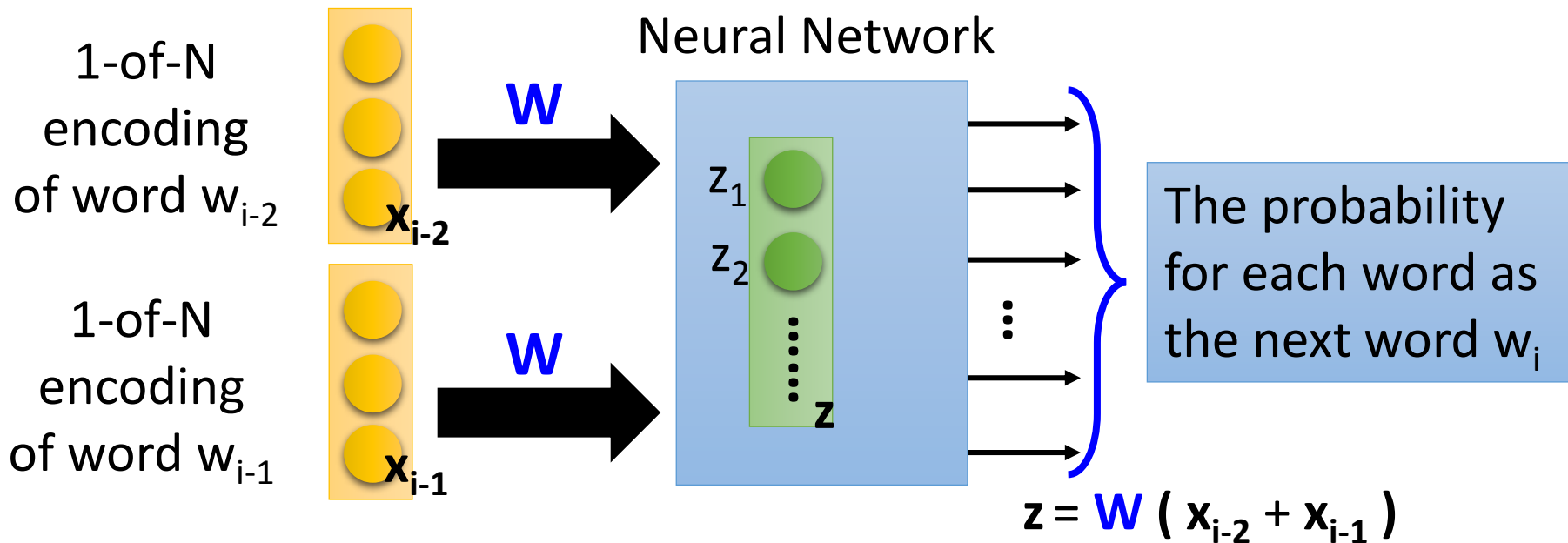| Deep Structured Semantic Model (DSSM) | • Application: Information Retrieval (IR) |
|---|---|
| Recursive Deep Model | • Application: Sentiment Analysis |
| Paragraph Vector | • Unsupervised |

Reference : http://cs.stanford.edu/~quocle/paragraph_vector.pdf

1-of-N encoding of word $w_{i-2}$

**W**

$\mathbf{x_{i-2}}$

Neural Network

$z_1$

$z_2$

$\mathbf{z}$

The probability for each word as the next word $w_i$

1-of-N encoding of word $w_{i-1}$

**W**

$\mathbf{x_{i-1}}$

$z = $ **W** ( $\mathbf{x_{i-2}}$ + $\mathbf{x_{i-1}}$ )

Paragraph $d_1$: (The paragraph is related to "The lord of the ring")

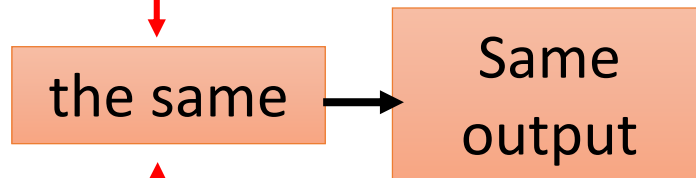...... 魔君　名叫　索倫 (Sauron) ......

$w_{i-2}$　　$w_{i-1}$　　$w_i$

$z = $ **W** ( $\mathbf{x_{i-2}}$ + $\mathbf{x_{i-1}}$ )

the same → Same output

Paragraph $d_2$: (The paragraph is related to "仙五")

...... 魔君　名叫　姜世離　......

$w_{i-2}$　　$w_{i-1}$　　$w_i$

$z = $ **W** ( $\mathbf{x_{i-2}}$ + $\mathbf{x_{i-1}}$ )

# Paragraph Vector

$d_1$ **1** **0** **0** ⋮   $d_2$ **0** **1** **0** ⋮   $d_3$ **0** **0** **1** ⋮

1-of-N encoding of paragraph d

**W'**

Original word vector: $z = W ( x_{i-2} + x_{i-1} )$

Paragraph vector:

$z = W ( x_{i-2} + x_{i-1} ) + W' d$

**d**

1-of-N encoding of word $w_{i-2}$

$x_{i-2}$

**W**

Neural Network

$z_1$
$z_2$
⋮
**z**

1-of-N encoding of word $w_{i-1}$

$x_{i-1}$

**W**

The probability for each word as the next word $w_i$

# Paragraph Vector

Original word vector:
$$z = W ( x_{i-2} + x_{i-1} )$$
Paragraph vector:
$$z = W ( x_{i-2} + x_{i-1} ) + W' d$$

Then error of the prediction can be explained by the meaning of the paragraphs.

Paragraph $d_1$: (The paragraph is related to "The lord of the ring")

...... 魔君　名叫　索倫 (Sauron) ......
$w_{i-2}$　　$w_{i-1}$　　$w_i$

$$z = W ( x_{i-2} + x_{i-1} ) + W' d_1$$
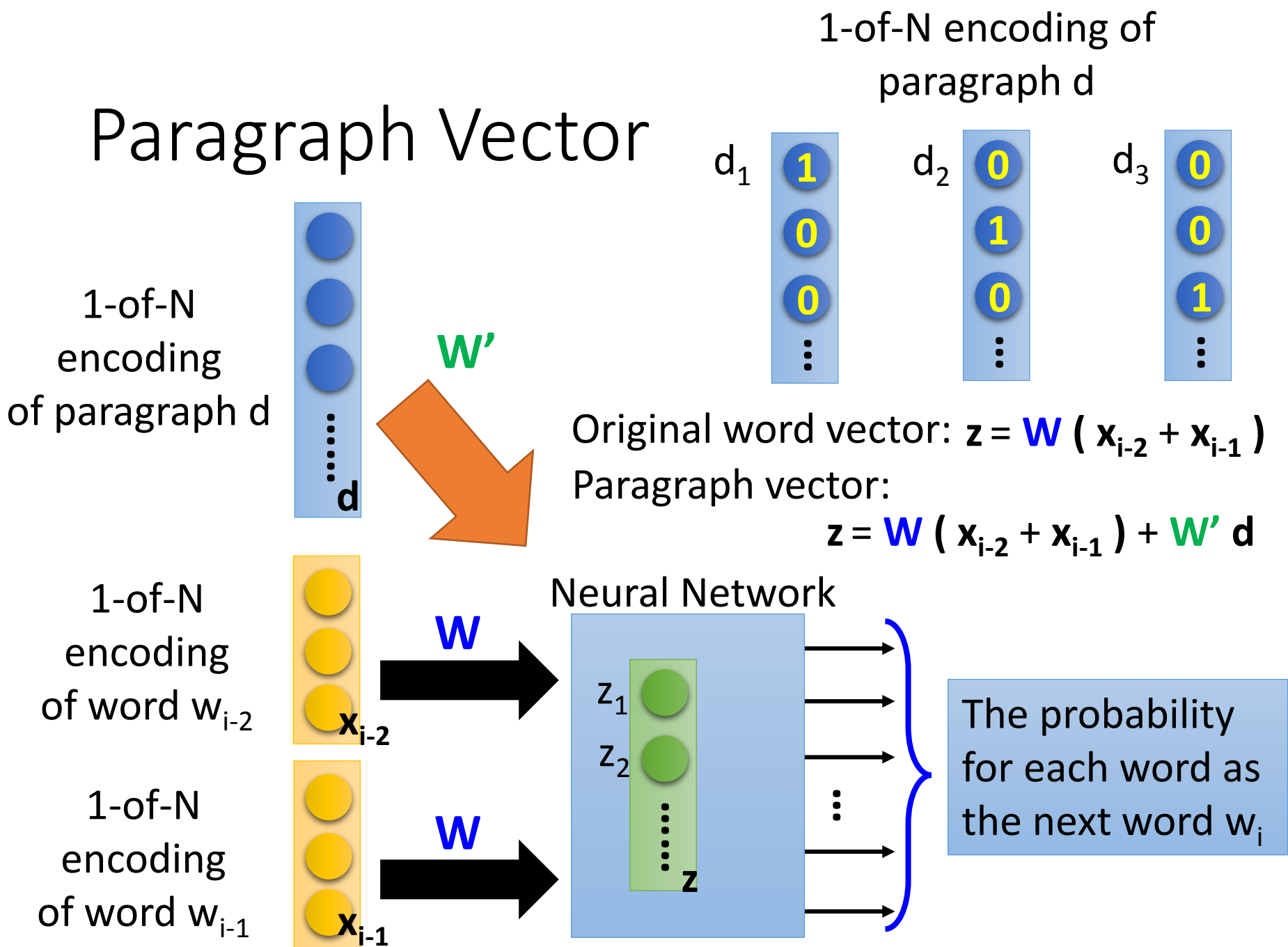
Paragraph $d_2$: (The document is related to "仙五")

...... 魔君　名叫　姜世離 ......
$w_{i-2}$　　$w_{i-1}$　　$w_i$

$$z = W ( x_{i-2} + x_{i-1} ) + W' d_2$$

different

***Paragraph vector*** of d: $V(d) = W' d$ ➡ Meaning of the paragraph

# Meaning of Word Sequence - Summary

| Deep Structured Semantic Model (DSSM) | • Application: Information Retrieval (IR) |
|---|---|
| Recursive Deep Model | • Application: Sentiment Analysis |
| Paragraph Vector | • Unsupervised |

Thank You

# Appendix

# Chinese Room



https://www.youtube.com/watch?feature=player_embedded&v=0F3-j-GQcts

# Demo in the paper

# More ……

- The paragraph vector can also be used in retrieval
  - Demo: http://www.logos.t.u-tokyo.ac.jp/~hassy/implementations/paragraph_vector/
- Toolkit: https://github.com/klb3713/sentence2vec

# Word classes

- One of the most successful NLP concepts in practice
- Similar words should share parameter estimation, which leads to generalization
- Example:
$$Class_1 = (yellow, green, blue, red)$$
$$Class_2 = (Italy, Germany, France, Spain)$$

- Usually, each vocabulary word is mapped to a single class (similar words share the same class)

# Word classes

- There are many ways how to compute the classes – usually, it is assumed that similar words appear in similar contexts

- Instead of using just counts of words, we can use also counts of classes, which leads to generalization (better performance on novel data)

*Class-based n-gram models of natural language* (Brown, 1992)