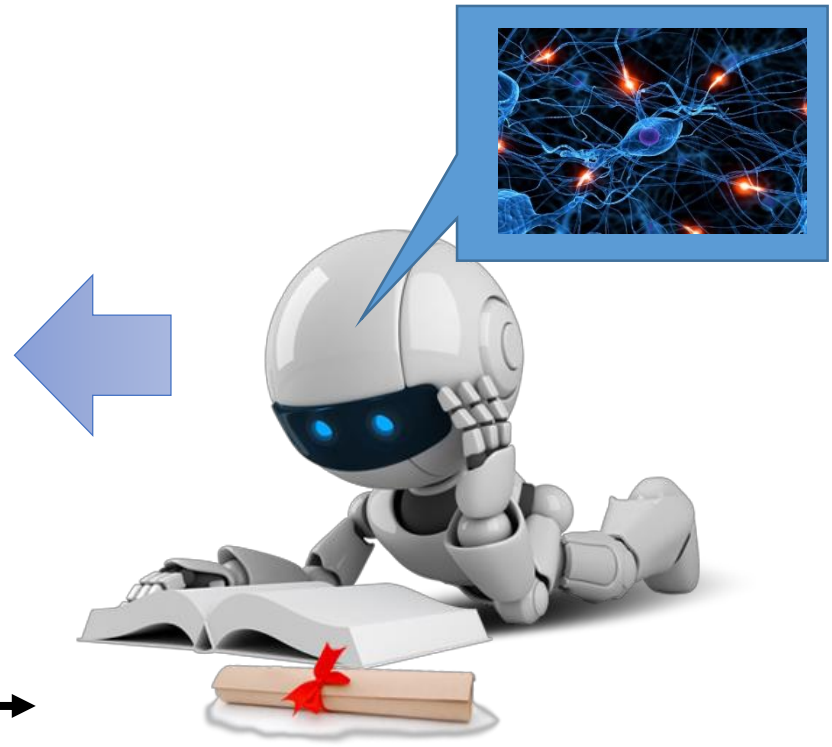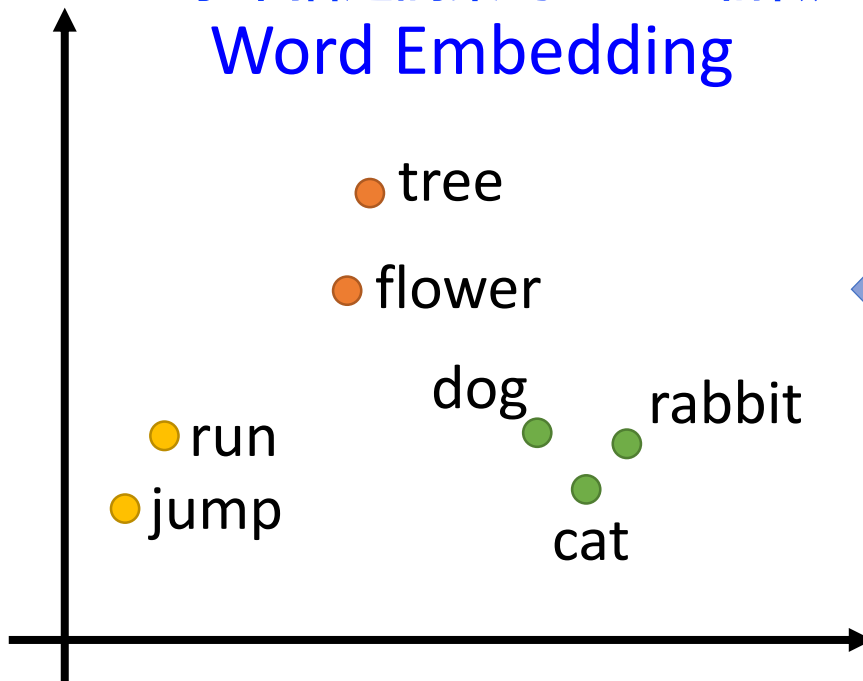# Unsupervised Learning: Word Embedding

**特別用在文字上的dimension reduction**

# Word Embedding

• Machine learns the meaning of words from reading a lot of documents without supervision

拿來描述詞彙的vector稱做
Word Embedding

tree

flower

dog    rabbit

run

jump

cat

用lexicon size的vector來描述一個詞彙

## *1-of-N Encoding*

apple = [ 1  0  0  0  0]
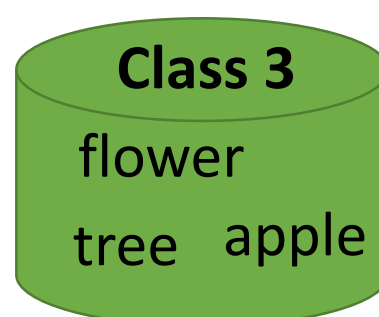
bag   = [ 0  1  0  0  0]

cat   = [ 0  0  1  0  0]

dog   = [ 0  0  0  1  0]

elephant  = [ 0  0  0  0  1]

壞處是詞彙間的關係無
法藉由vector傳遞出來

## *Word Embedding*



dog
rabbit
run
jump
cat
tree
flower

所以選擇continuous的
vector (word embedding)

## *Word Class*



**class 1**
dog
cat
bird

**Class 2**
ran
jumped
walk

**Class 3**
flower
tree  apple

因此我們可以先做clustering，但這樣的分類方式太粗糙

# Word Embedding

**想法：每一個word皆可以用上下文來判斷其語意**

- Machine learns the meaning of words from reading a lot of documents without supervision

- A word can be understood by its context

蔡英文、馬英九 are something very similar

You shall know a word by the company it keeps

馬英九 520宣誓就職

有類似的content，可能vector某個dimension會一樣

蔡英文 520宣誓就職

# How to exploit the context?

- **Count based** 原則：若wi,wj常常一起出現，則他們的 vector繼算相似度應該要很接近
  - If two words $w_i$ and $w_j$ frequently co-occur, $V(w_i)$ and $V(w_j)$ would be close to each other
  - E.g. Glove Vector: http://nlp.stanford.edu/projects/glove/

越接近越好

$$V(w_i) . V(w_j) \longleftrightarrow N_{i,j}$$

N(i,j): 一起出現的次數

Inner product

Number of times $w_i$ and $w_j$ in the same document

- **Perdition based**

# Prediction-based – Training

首先讓machine predict接下來出現的詞彙

Collect data: 首先要先做斷詞

潮水 退了 就 知道 誰 …
不爽 不要 買 …
公道價 八萬 一 …
………

**Minimizing cross entropy**

潮水
退了 → Neural Network → m i n i m i z e cross entropy → 就

退了
就 → Neural Network → m i n i m i z e cross entropy → 知道

就
知道 → Neural Network → m i n i m i z e cross entropy → 誰

# Prediction-based - 推文接話

推 louisee :話說十幾年前我念公立國中時,老師也曾做過這種事,但

https://www.ptt.cc/bbs/Teacher/M.1317226791.A.558.html

推 AO56789: 我同學才扯好不好，他有一次要交家政料理報告
→ AO56789:其中一個是要寫一樣水煮料理的食譜，他居然給我寫

著名簽名檔 (出處不詳)

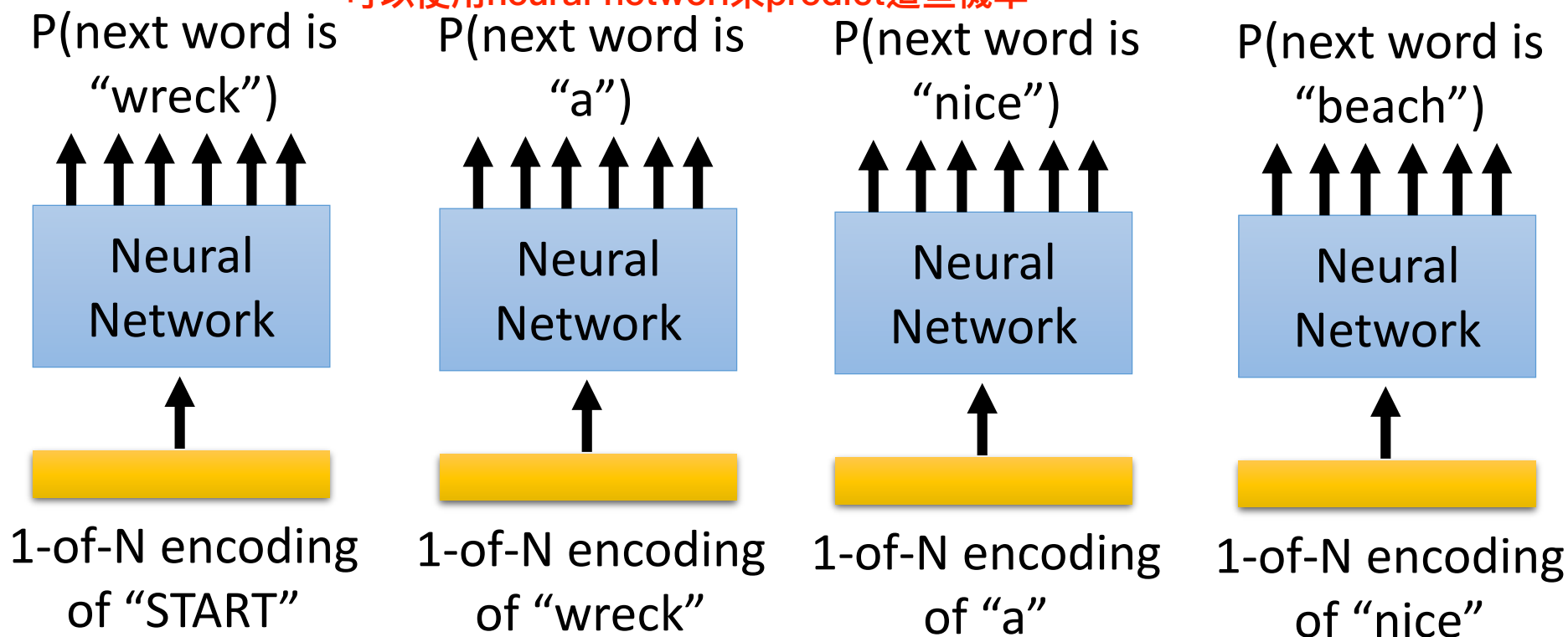# Prediction-based – Language Modeling

P("wreck a nice beach")

即使搜集了大量文章再去算這個句子出現的機率，仍然有機會等於零

=P(wreck|START)P(a|wreck)P(nice|a)P(beach|nice)

P(b|a): the probability of NN predicting the next word.

可以使用neural networl來predict這些機率

P(next word is "wreck")

P(next word is "a")

P(next word is "nice")

P(next word is "beach")

| Neural Network | Neural Network | Neural Network | Neural Network |

1-of-N encoding of "START"

1-of-N encoding of "wreck"

1-of-N encoding of "a"

1-of-N encoding of "nice"

# 利用neural network做出一個language model



$i$-th output $= P(w_t = i \mid context)$

softmax

most computation here

tanh

$C(w_{t-n+1})$   $C(w_{t-2})$   $C(w_{t-1})$   降維

Table look−up in $C$

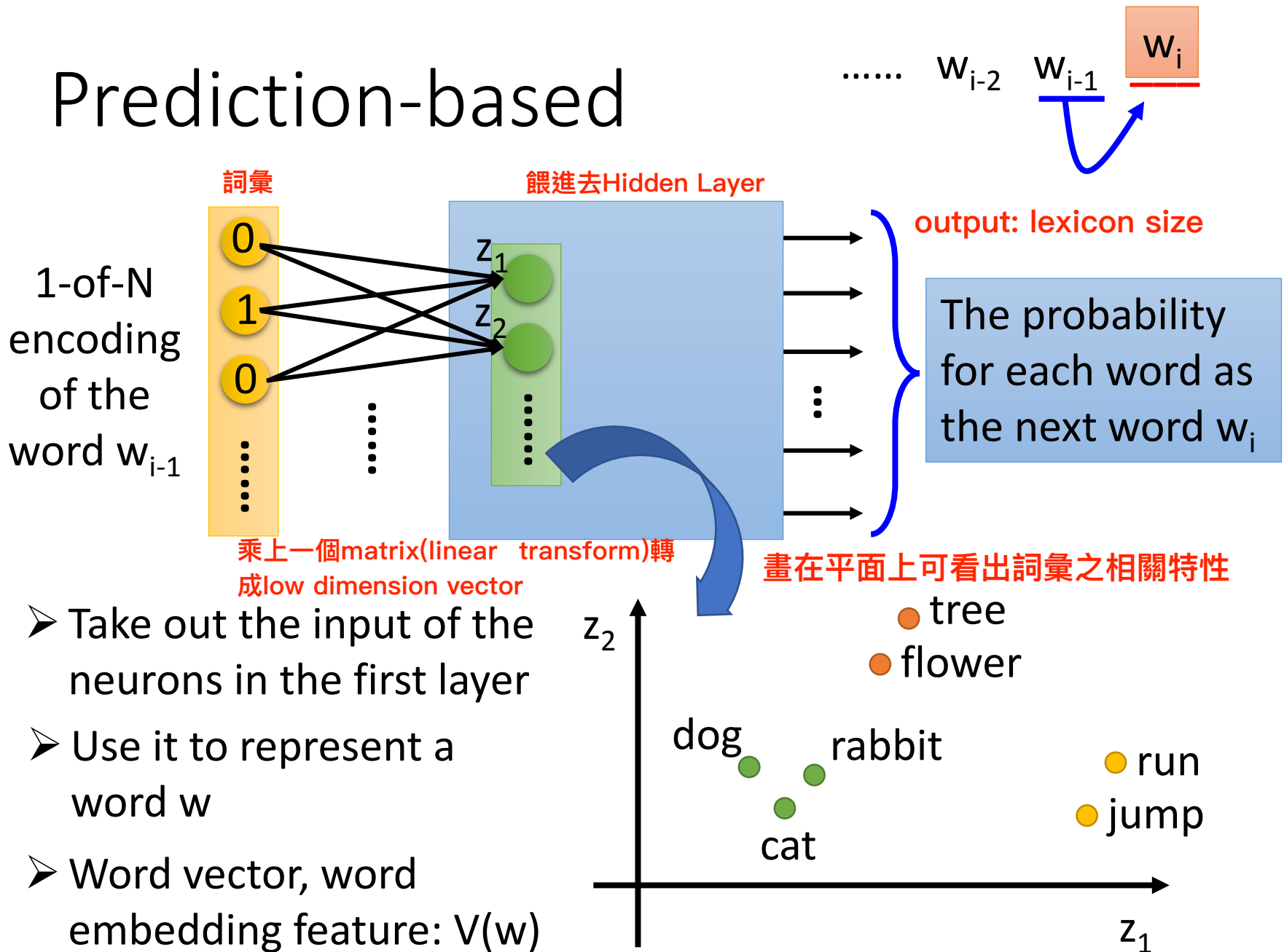Matrix $C$ shared parameters across words

index for $w_{t-n+1}$   index for $w_{t-2}$   index for $w_{t-1}$

Bengio, Y., Ducharme, R., Vincent, P., & Jauvin, C. (2003). A neural probabilistic language model. *Journal of machine learning research*, *3*(Feb), 1137-1155.
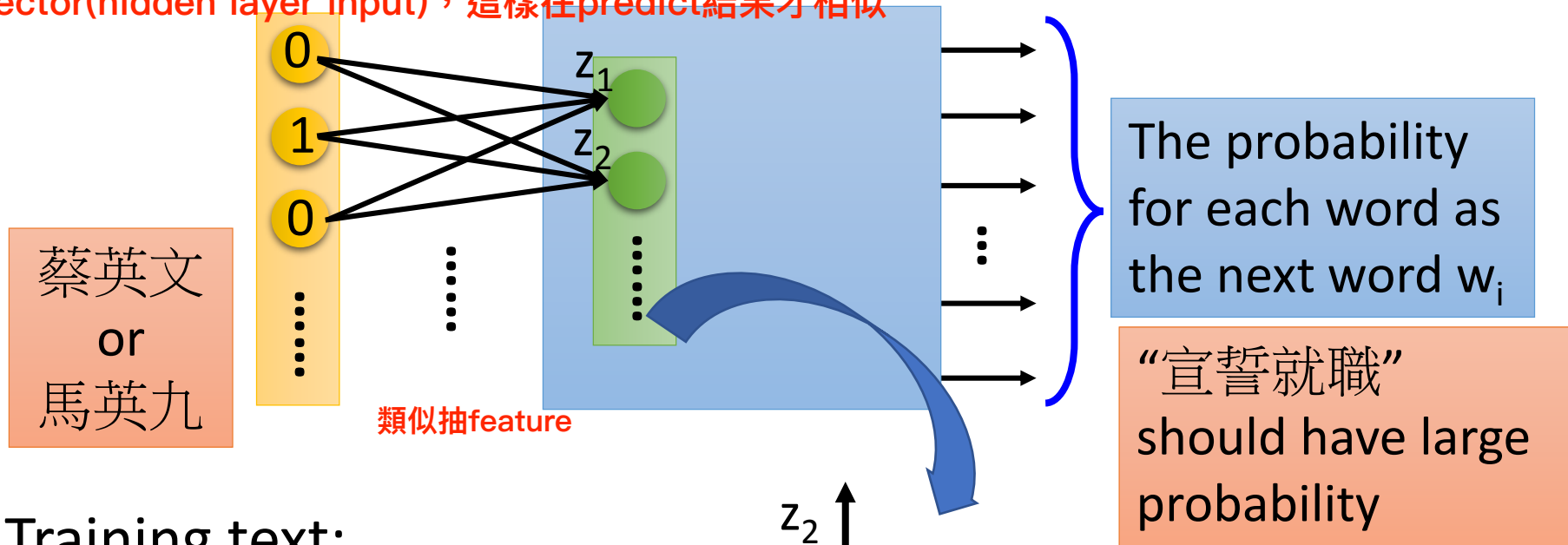
第一篇用NN解language model

# Prediction-based

$$w_{i-2} \quad w_{i-1} \quad \boxed{w_i}$$

……

**詞彙**

**餵進去Hidden Layer**

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

$$\begin{array}{c} 0 \\ 1 \\ 0 \\ \vdots \end{array}$$

$z_1$

$z_2$

$\vdots$

**output: lexicon size**

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

**乘上一個matrix(linear transform)轉成low dimension vector**

**畫在平面上可看出詞彙之相關特性**

➢ 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

$z_1$

# Prediction-based

在做word embedding與其做一個deep的model，不如用一個shallow的linear model，類似抽feature（PCA）

將不同的input做transform後能夠得到相似的vector(hidden layer input)，這樣在predict結果才相似

類似抽feature

蔡英文
or
馬英九

$z_1$
$z_2$

You shall know a word by the company it keeps

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

"宣誓就職" should have large probability

Training text:

...... 蔡英文  宣誓就職 ......
$w_{i-1}$     $w_i$

...... 馬英九  宣誓就職 ......
$w_{i-1}$     $w_i$

比較靠近的word embedding
蔡英文
馬英九

$z_2$
$z_1$

不同的詞彙如果有同樣的attribute，在做word embedding，他們在某幾個dimension會有接近的值

# Prediction-based – Sharing Parameters

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

$W_1$

$x_{i-2}$

$z_1$
$z_2$
$z$

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

|V| = lexicon size

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|X|V| matrices.  將w1,w2 share 參數

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

考慮長一點的word才能有好的結果(至少10~20)     共用參數，使得vector dimension不會增加

12

# Prediction-based
## – 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 $w_i$

每個dimension所對應的weight是相同的

The weights with the same color should be the same.

Or, one word would have two word vectors.

由於將參數tight在一起，因此同個詞彙放在不同位置得到的結果是相同的

13

# Prediction-based – Various Architectures

- Continuous bag of word (CBOW) model

看context predict中間的詞彙

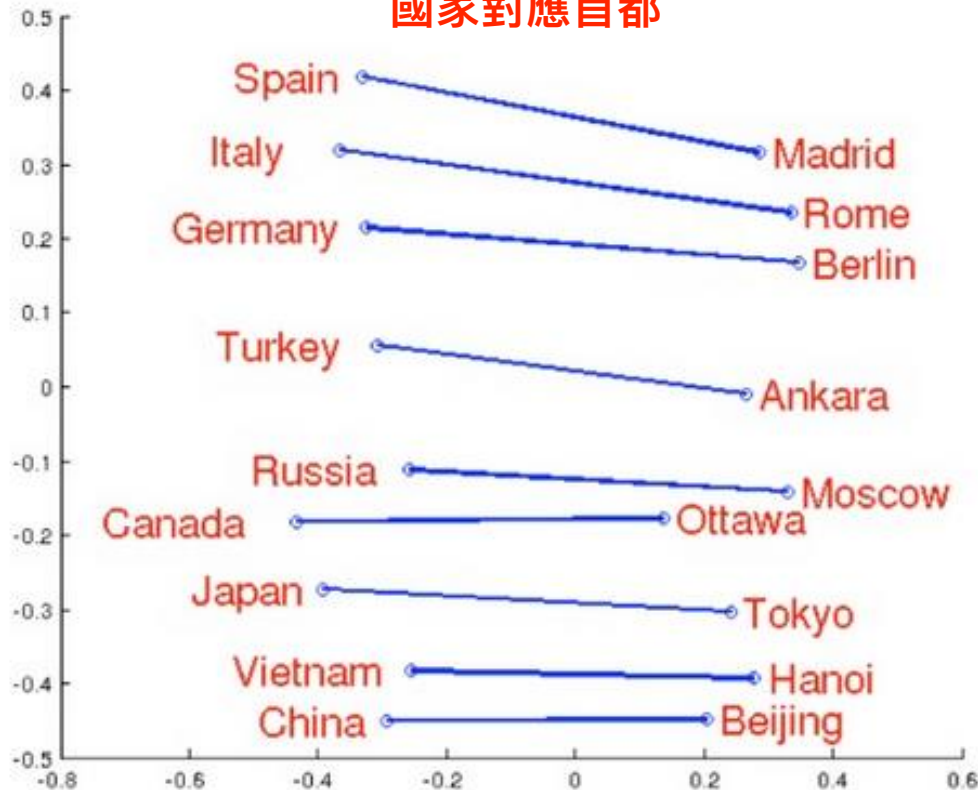$\ldots\ldots$ $w_{i-1}$ ____ $w_{i+1}$ $\ldots\ldots$

$w_{i-1}$ → Neural Network → ↔ $w_i$
$w_{i+1}$ →

***predicting the word given its context***

- Skip-gram

$\ldots\ldots$ ____ $w_i$ ____ $\ldots\ldots$

看中間的詞彙predict前後

$w_i$ → Neural Network → ↔ $w_{i-1}$
→ ↔ $w_{i+1}$

***predicting the context given a word***

# Word Embedding



國家對應首都

動詞三態變化

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

# Word Embedding



找出詞彙間的從屬關係

狗-警犬
dog - police dog

鸡-公鸡
chicken - cock

兔-长毛兔
rabbit - wool rabbit

驴-野驴
donkey - wild ass

羊-小尾寒羊
sheep - small-tail Han sheep

马-斑马
equus - zebra

羊-公羊
sheep - ram

蟹-海蟹
crab - sea crab

虾-对虾
shrimp - prawn

海豚-白鳍豚
dolphin - white-flag dolphin

鱼-鲨鱼
fish - shark

鱼-金鱼
fish - gold fish

鱼-热带鱼
fish - tropical fish

运动员-足球球员
sportsman - footballer

职员-售货员
staff - salesclerk

职员-售票员
staff - conductor

职员-空姐
airline hostess

职员-公务员
staff - civil servant

工人-木匠
laborer - carpenter

工人-园丁
laborer - gardener

工人-临时工
laborer - temporary worker

海员-领航员
seaman - navigator

职位-校长
position - headmaster

职位-总领事
position – consul general

演员-歌手
actor - singer

演员-小丑
actor - clown

演员-主角
actor - protagonist

演员-斗牛士
actor - matador
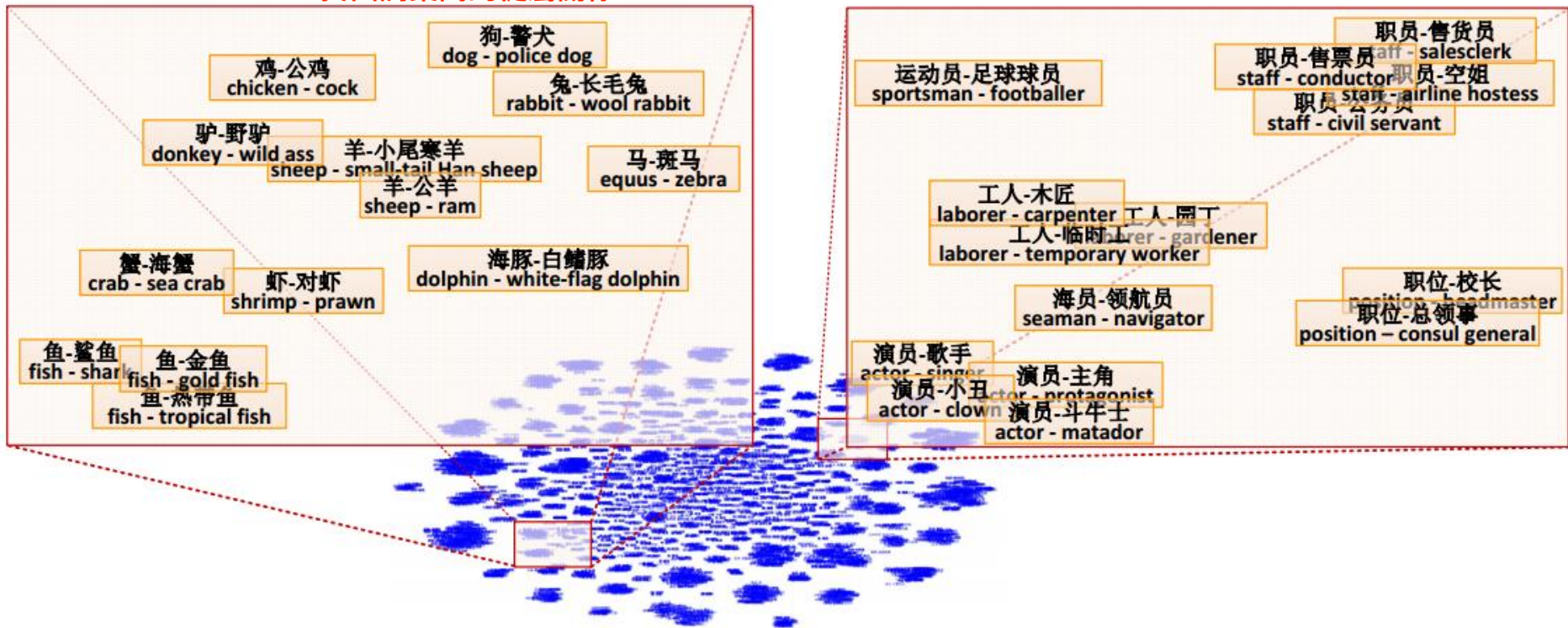
Fu, Ruiji, et al. "Learning semantic hierarchies via word embeddings."*Proceedings of the 52th Annual Meeting of the Association for Computational Linguistics: Long Papers*. Vol. 1. 2014.

16

# Word Embedding

- Characteristics

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

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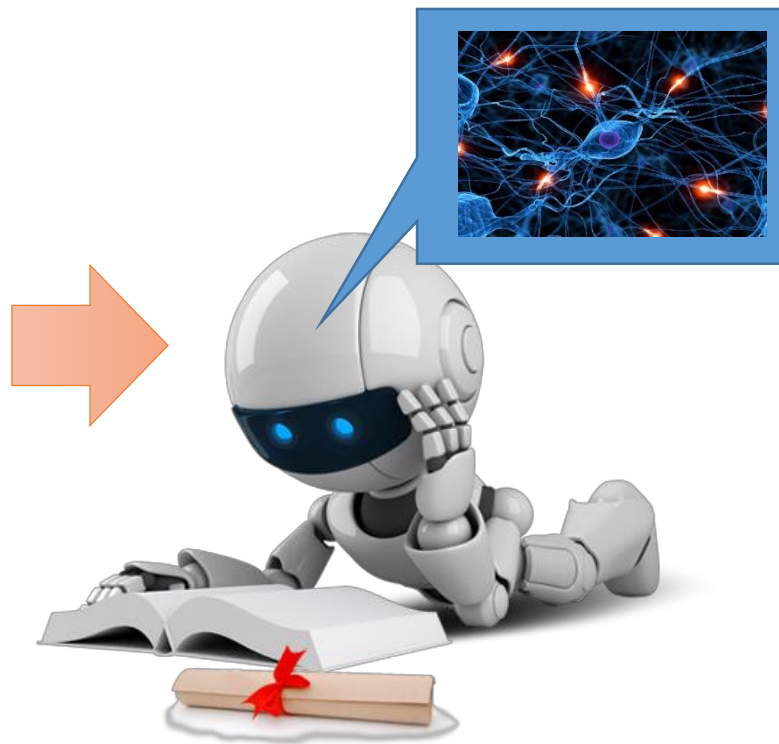
- Solving analogies

羅馬之於意大利=柏林之於？

Rome : Italy = Berlin : ?

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

Find the word w with the closest V(w)

算出vector後找最相近的vector

# Demo

- Machine learns the meaning of words from reading a lot of documents without supervision
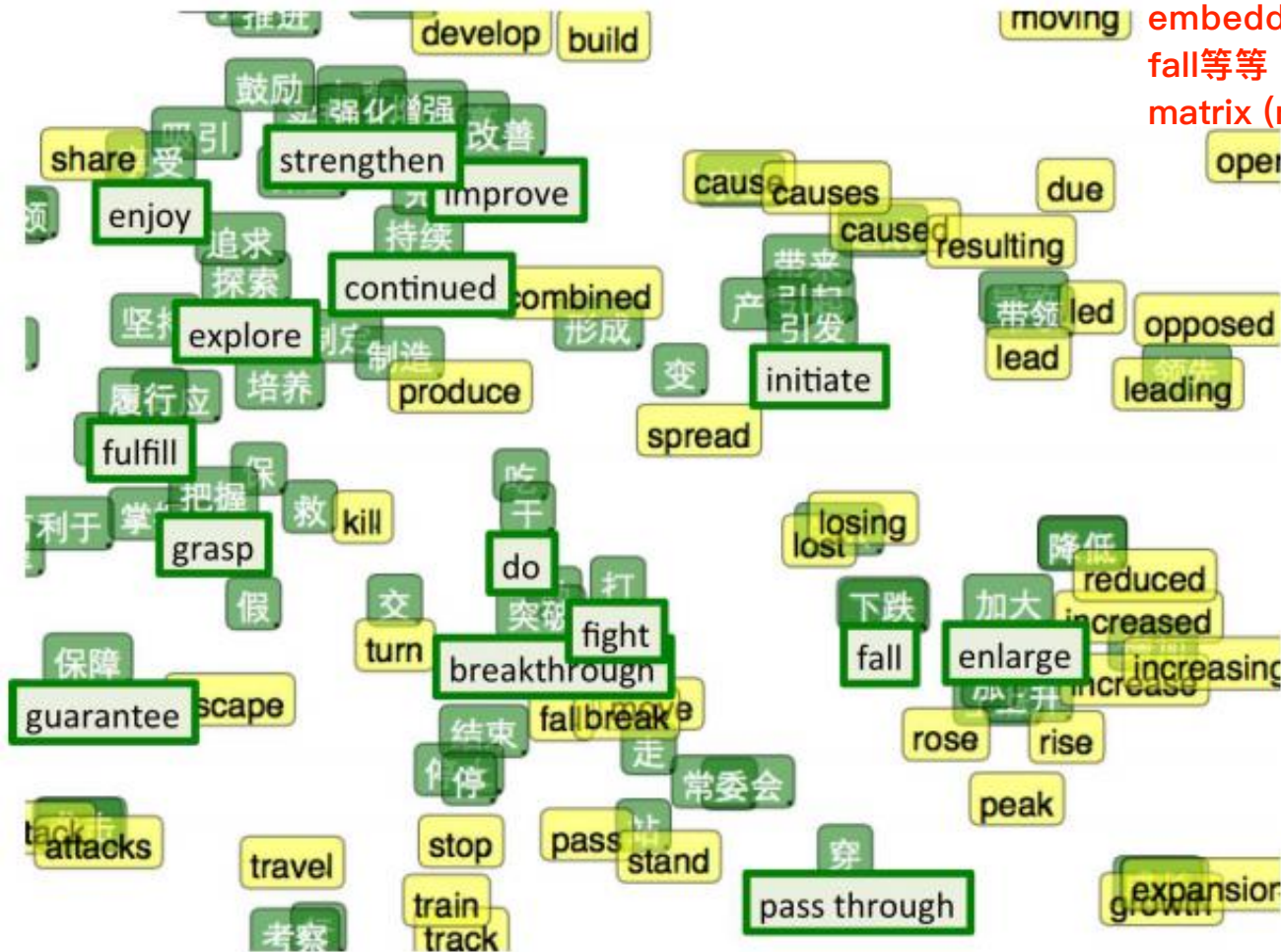
# Demo

- 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)

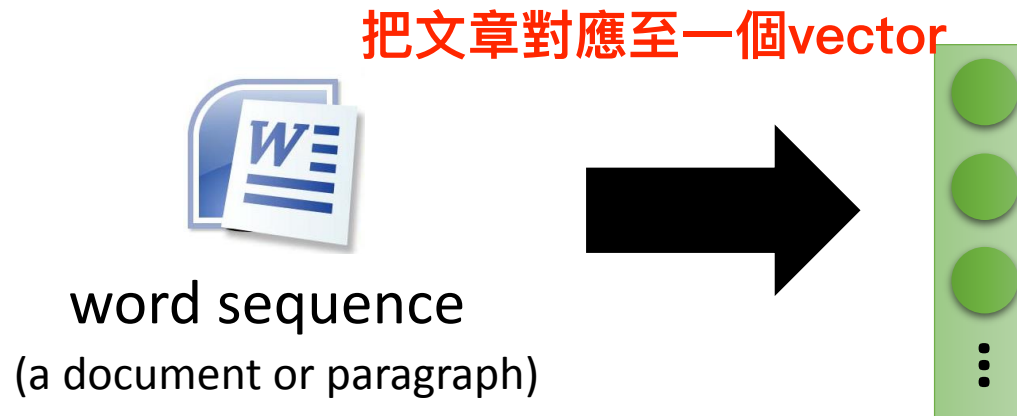word embedding是unsupervised的，因此我們並不清楚每個dimension代表什麼意思

# Multi-lingual Embedding

中文跟英文無法一起train，舉例來說中文的第三維可能是動物英文的第七維才是，因此會亂掉
因此要先有一些已經可以相互對應的 word embedding，如 下跌v.s. fall等等，找出transform的 matrix (mapping)
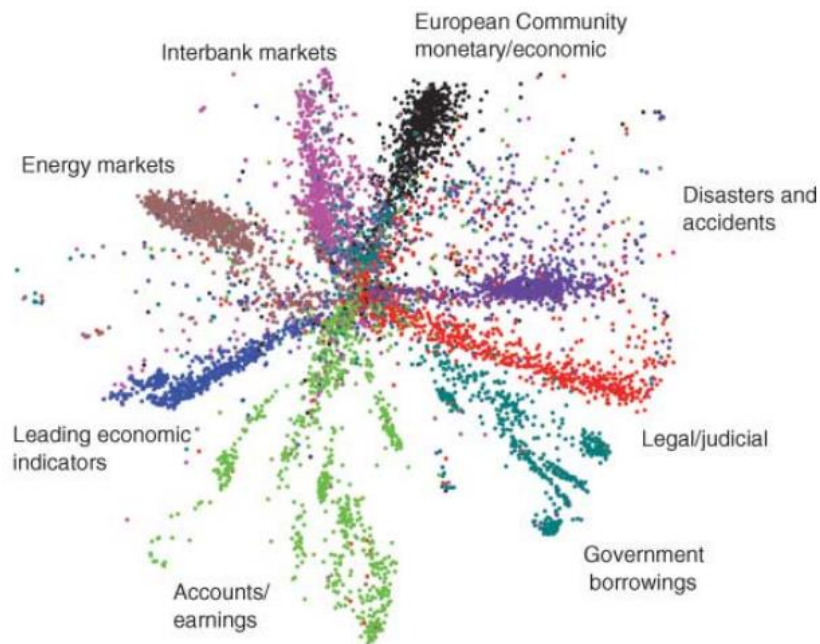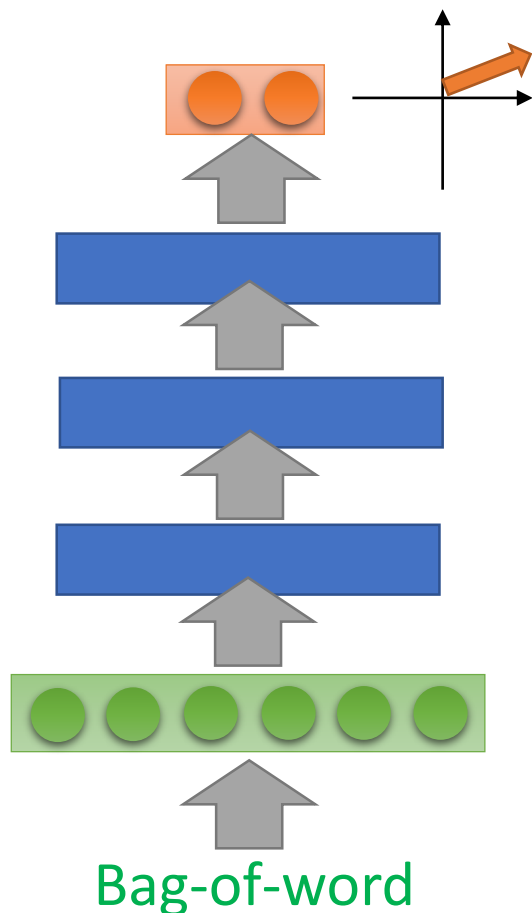


Bilingual Word Embeddings for Phrase-Based Machine Translation, Will Zou, Richard Socher, Daniel Cer and Christopher Manning, EMNLP, 2013

# Document Embedding

- word sequences with different lengths → the vector with the same length
  - The vector representing the meaning of the word sequence  不同文章的長度不同
  - A word sequence can be a document or a paragraph

把文章對應至一個vector

word sequence
(a document or paragraph)

# Semantic Embedding



**Bag-of-word**

Reference: Hinton, Geoffrey E., and Ruslan R. Salakhutdinov. "Reducing the dimensionality of data with neural networks." *Science* 313.5786 (2006): 504-507
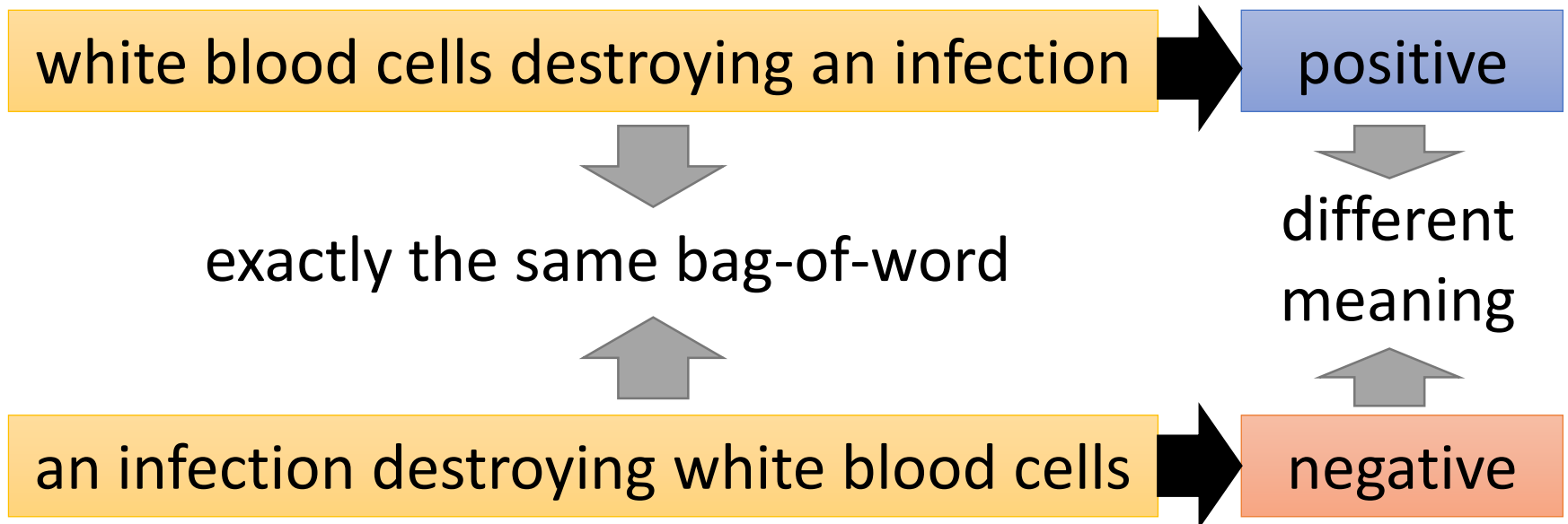
先將詞彙統計成Bag of Word，然後利用auto-encoder降維至二維平面

# Beyond Bag of Word

- To understand the meaning of a word sequence, the order of the words can not be ignored.

  即使Bag of Word相同，但是詞彙順序不同會影響其語意

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

exactly the same bag-of-word

different meaning

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

23

# Beyond Bag of Word

- ***Paragraph Vector***: Le, Quoc, and Tomas Mikolov. "Distributed Representations of Sentences and Documents." ICML, 2014

- ***Seq2seq Auto-encoder***: Li, Jiwei, Minh-Thang Luong, and Dan Jurafsky. "A hierarchical neural autoencoder for paragraphs and documents." arXiv preprint, 2015

- ***Skip Thought***: Ryan Kiros, Yukun Zhu, Ruslan Salakhutdinov, Richard S. Zemel, Antonio Torralba, Raquel Urtasun, Sanja Fidler, "Skip-Thought Vectors" arXiv preprint, 2015.