

# The Wonderful World of Embeddings

Franz Diebold, M.Sc. (TUM)



[github.com/franzdiebold](https://github.com/franzdiebold)



[linkedin.com/in/franz-diebold](https://linkedin.com/in/franz-diebold)

The screenshot shows the Google Translate interface. At the top, it says "Google Translate". Below that are buttons for "TURKISH" and "ENGLISH" with a double-headed arrow between them. A "Sign in" button is also present. The main input field contains the text "o bir doktor". Below the input field are microphone and speaker icons. The output section shows the translation "she is a doctor (feminine)" followed by "(masculine)". There is a blue button with a star icon and a note that translations are gender-specific with a link to "LEARN MORE".

More homes you may like



# Embeddings are all over the place!

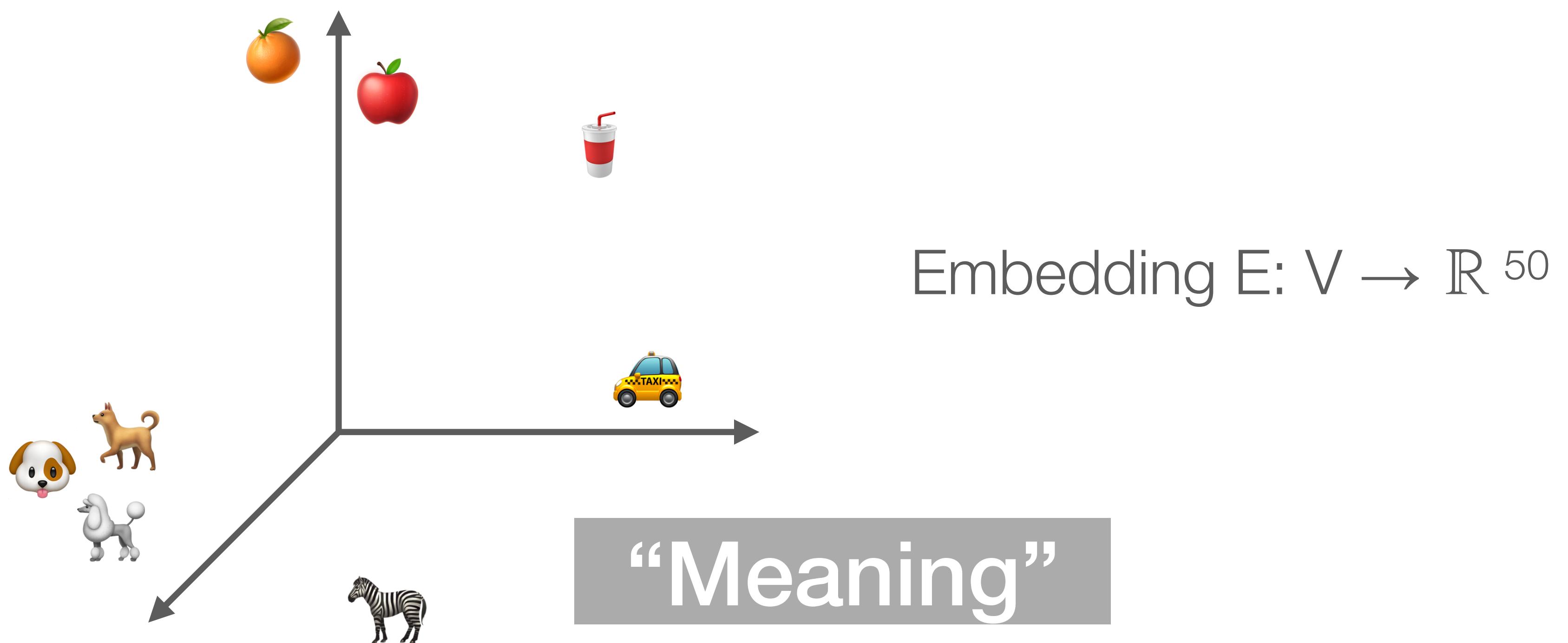


# Word Embedding

Demo

Vocabulary  $V = \{\text{apple}, \text{dog}, \text{chicken}, \text{poodle}, \text{orange}, \dots, \text{cup}, \text{taxi}, \text{zebra}\}$

$|V| = 100k$



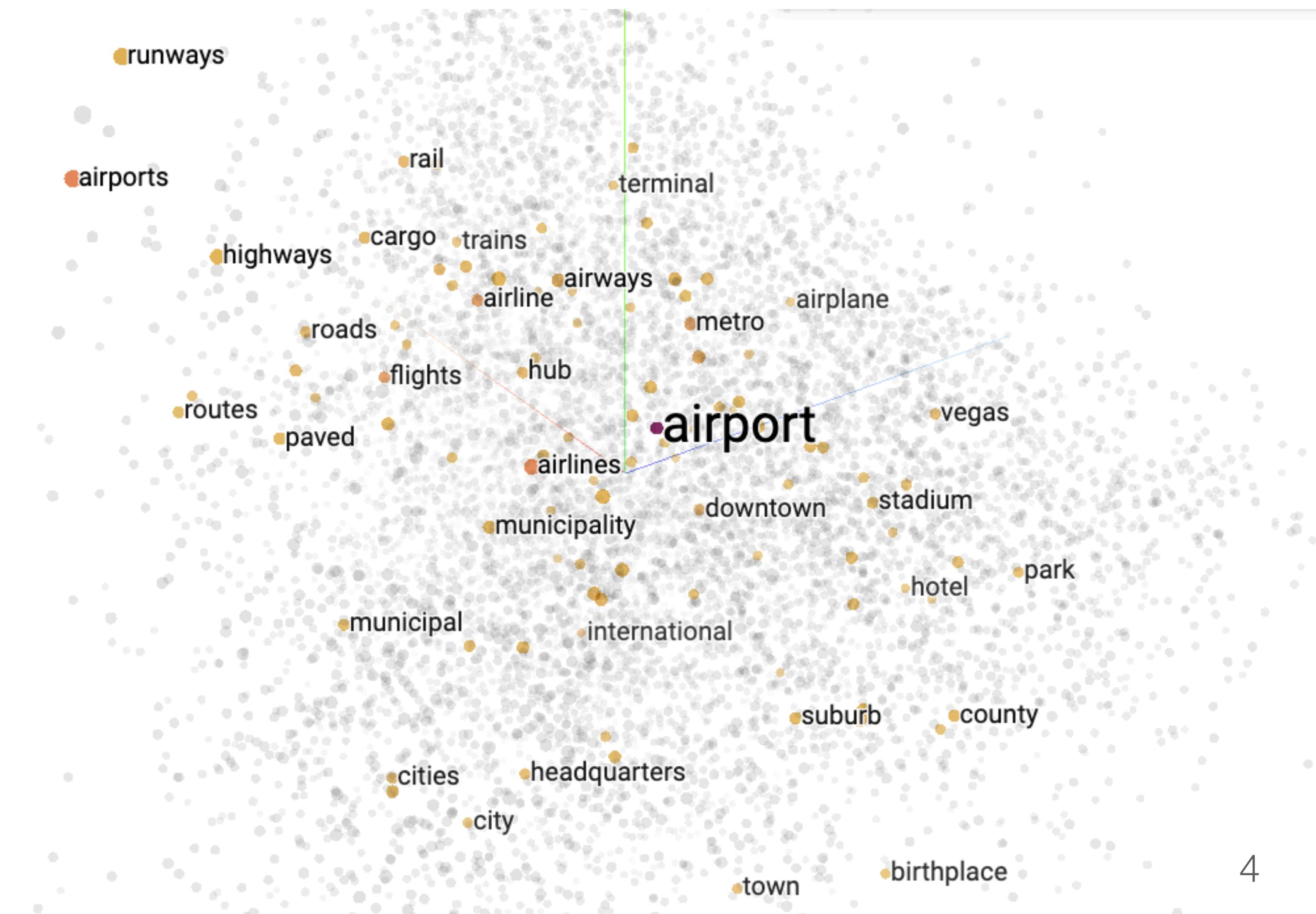
# Properties of Word Embeddings

## 1. Nearest neighbours

[Demo](#)

Using Euclidean distance or cosine similarity

Nearest points in the original space:	
airports	0.449
airlines	0.459
flights	0.493
airline	0.511
metro	0.525
international	0.546
ferry	0.550
downtown	0.556
airways	0.575
rail	0.578
passengers	0.579
hub	0.580

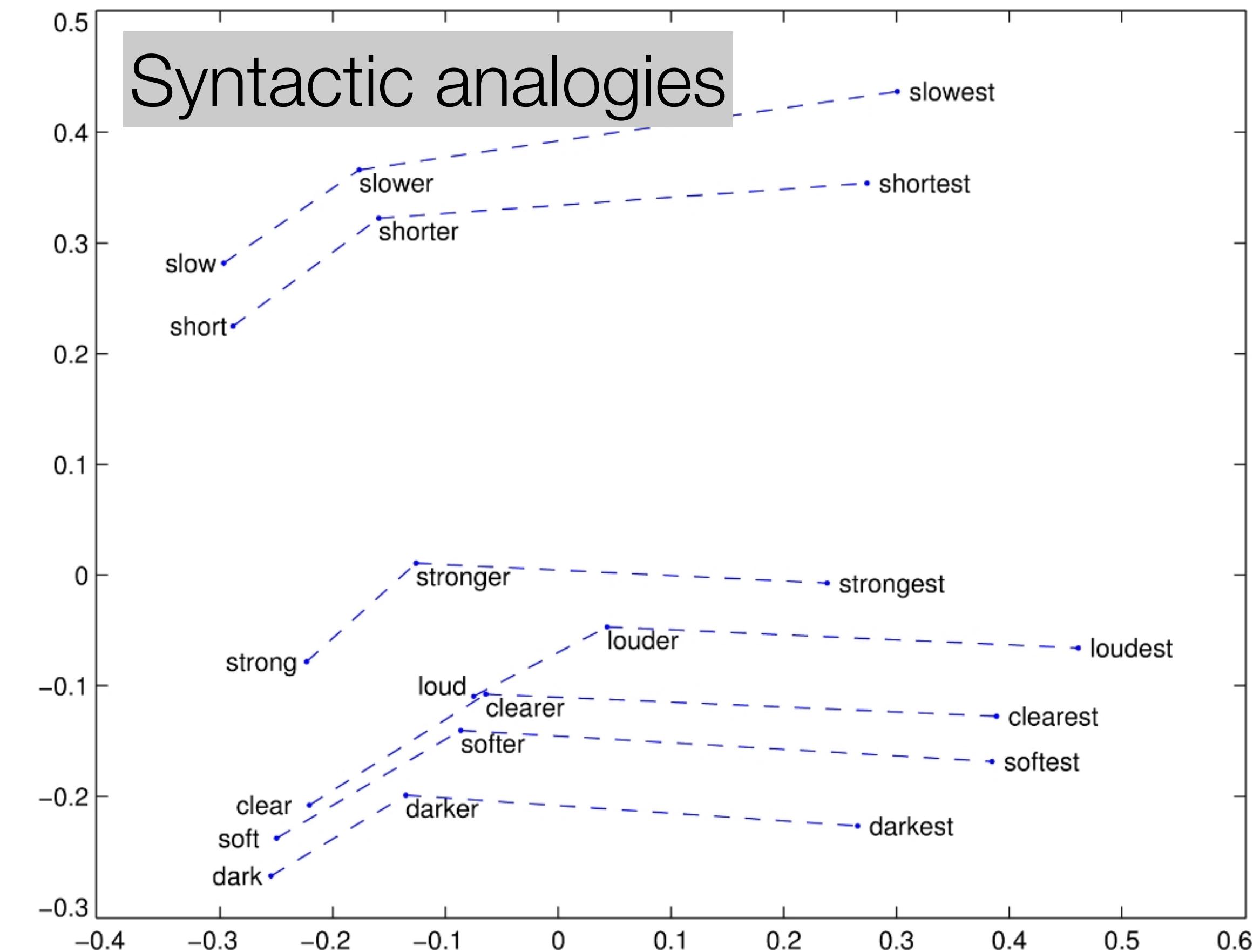
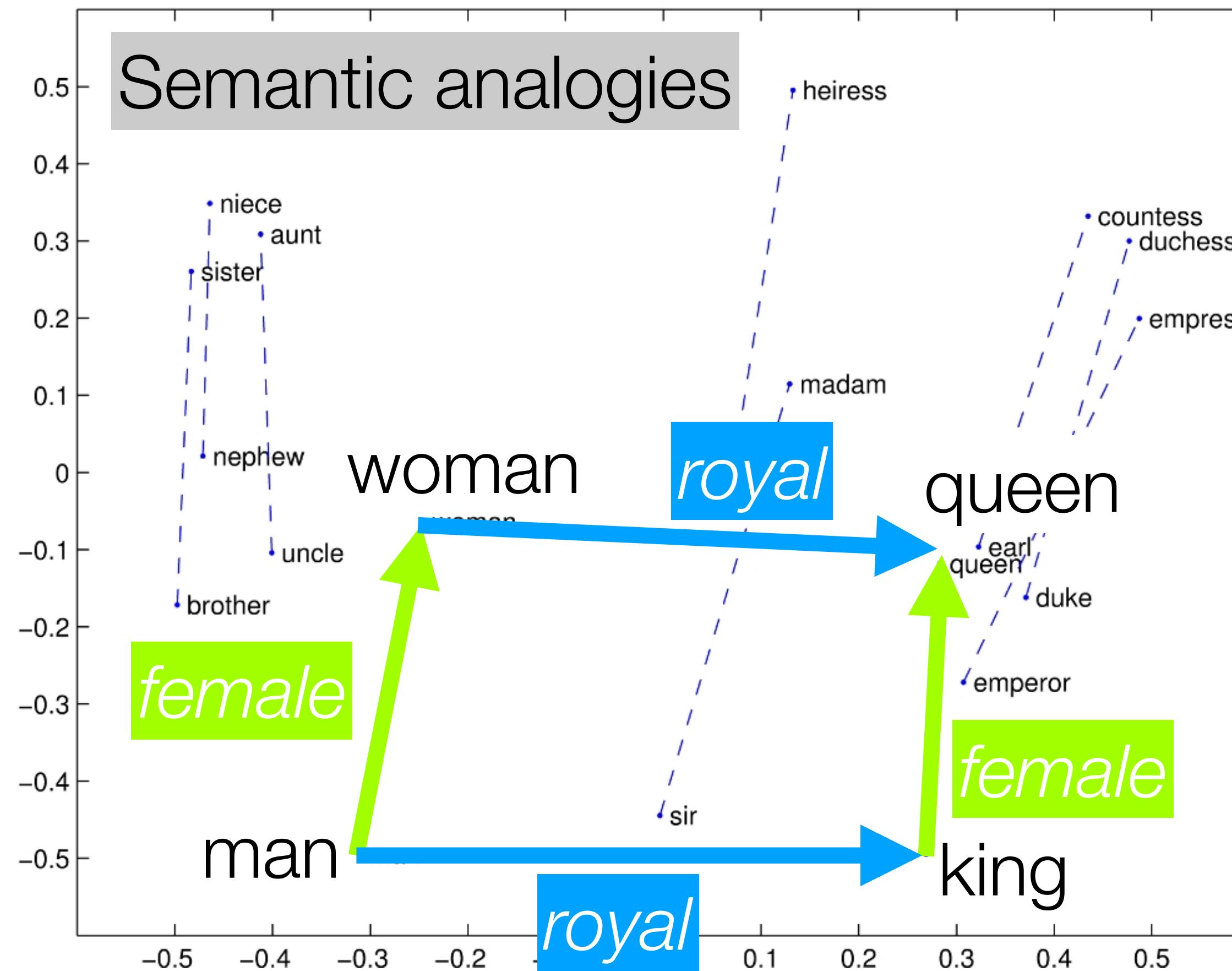


# Properties of Word Embeddings

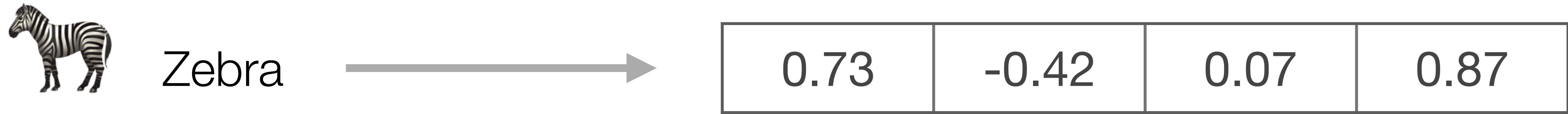
## 2. Linear substructures

Demo

For visualisation purposes: dimensions further reduced using or t-SNE



# Word2Vec

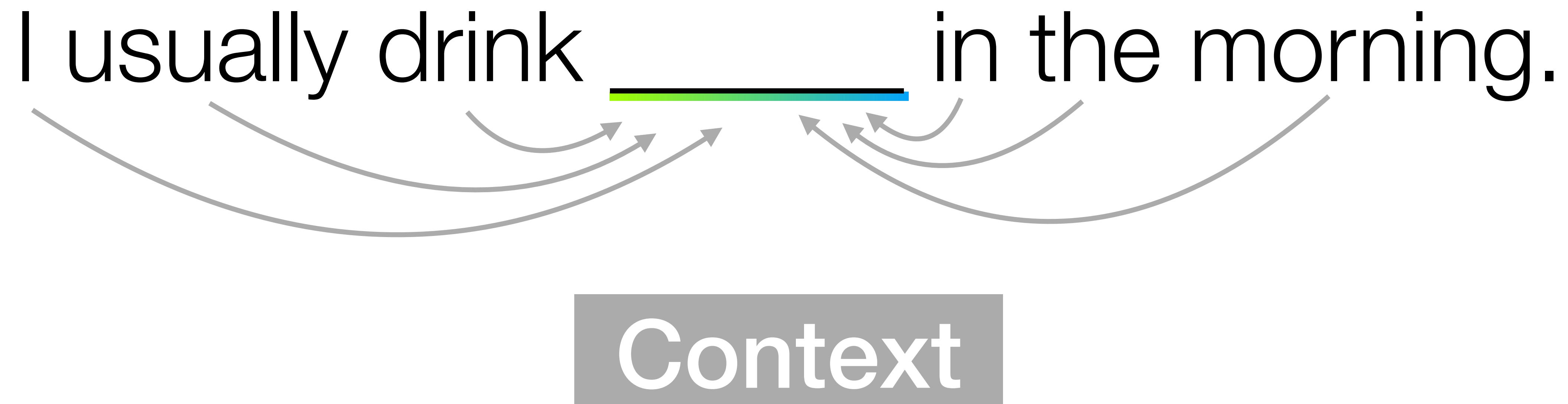


- Group of models to compute “*continuous vector representations of words*” (Word Embeddings)
- Two model architectures
  - Continuous Bag-of-Words (CBOW)
  - Continuous Skip-gram

# Word2Vec

## Continuous Bag-of-Words (CBOW)

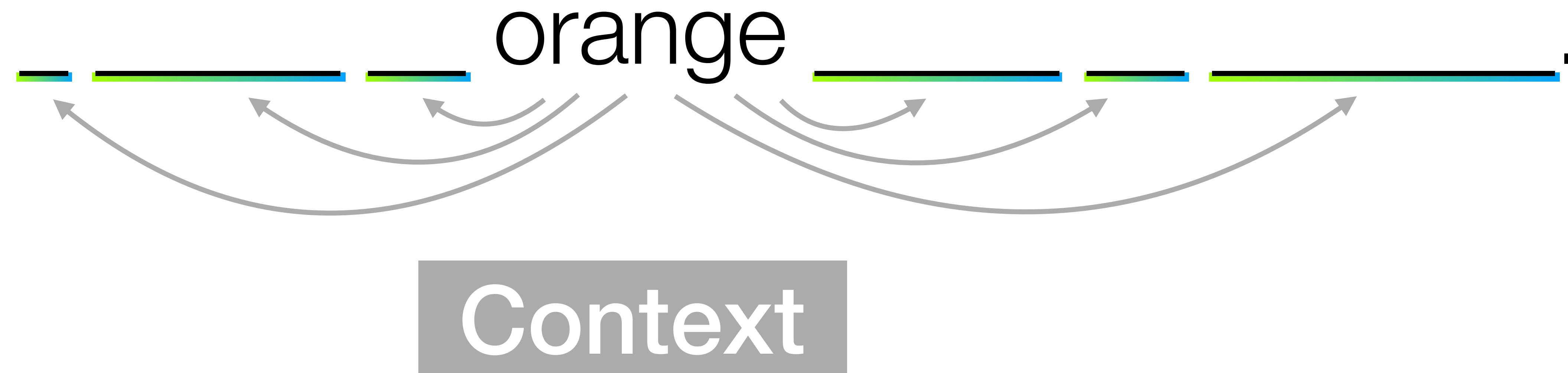
- Given the surrounding words, predict the middle word.



# Word2Vec

## Continuous Skip-gram

- Given a word, predict the surrounding words.



# Word2Vec

## Continuous Skip-gram

- Context window size: 2

Jupiter	is	the	largest	planet	in	the	Solar	System
---------	----	-----	---------	--------	----	-----	-------	--------

# Word2Vec

## Continuous Skip-gram

- Context window size: 2



Input	Label
Jupiter	is
Jupiter	the

# Word2Vec

## Continuous Skip-gram

- Context window size: 2



Input	Label
Jupiter	is
Jupiter	the
is	Jupiter
is	the
is	largest

# Word2Vec

## Continuous Skip-gram

- Context window size: 2



Input	Label
Jupiter	is
Jupiter	the
is	Jupiter
is	the
is	largest
the	Jupiter
the	is
the	largest
the	planet

# Word2Vec

## Continuous Skip-gram

- Context window size: 2



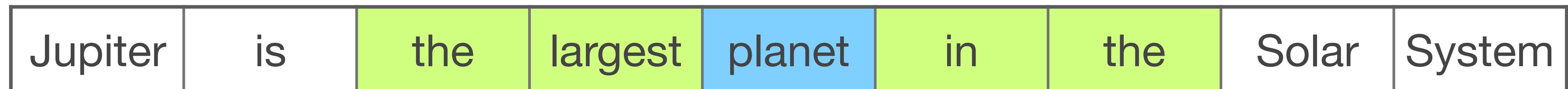
Input	Label
Jupiter	is
Jupiter	the
is	Jupiter
is	the
is	largest
the	Jupiter
the	is
the	largest
the	planet

Input	Label
largest	is
largest	the
largest	planet
largest	in

# Word2Vec

## Continuous Skip-gram

- Context window size: 2



Input	Label
Jupiter	is
Jupiter	the
is	Jupiter
is	the
is	largest
the	Jupiter
the	is
the	largest
the	planet

Input	Label
largest	is
largest	the
largest	planet
largest	in
planet	the
planet	largest
planet	in
planet	the

# Word2Vec

## Continuous Skip-gram

- Context window size: 2

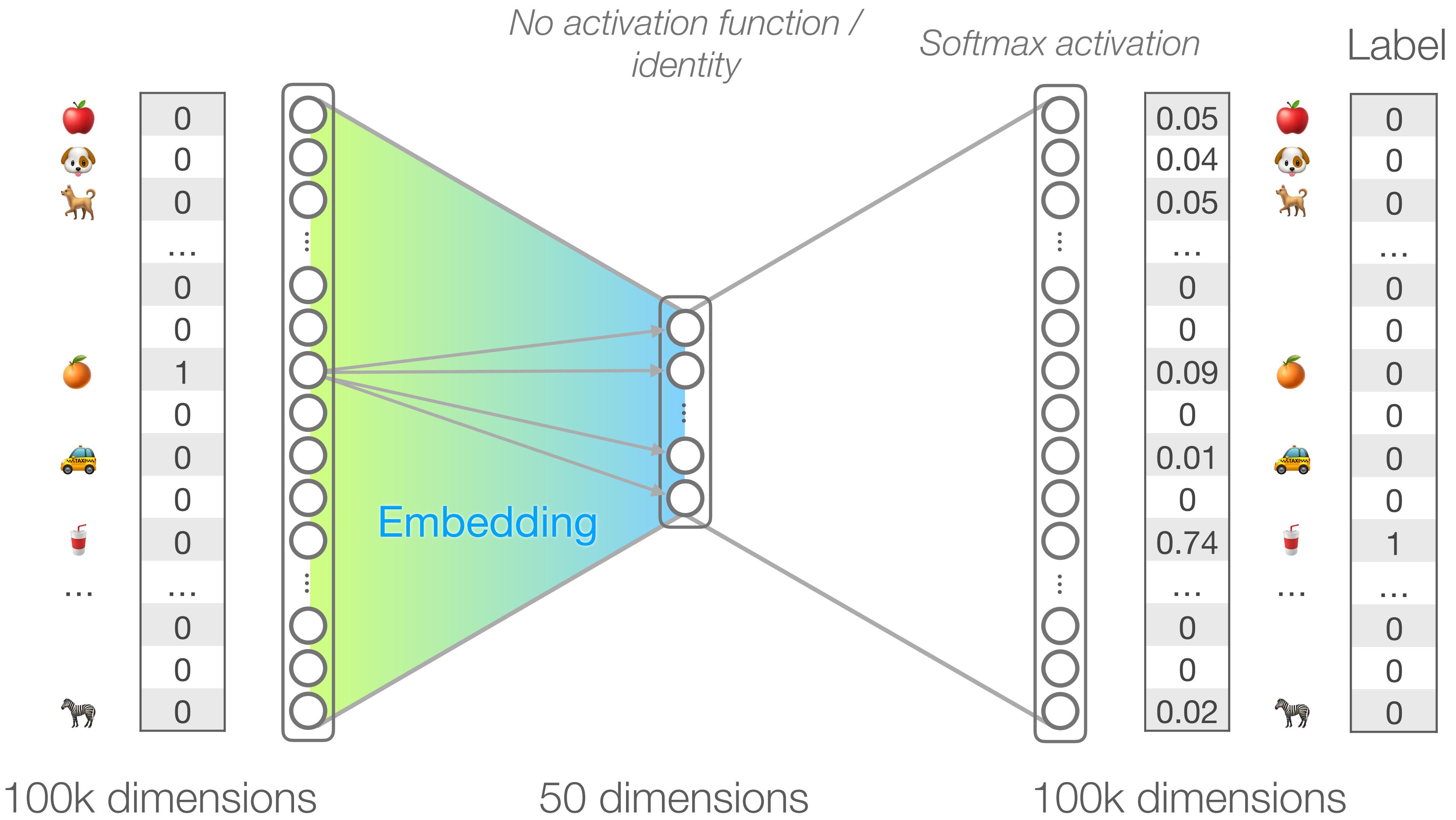


Input	Label
Jupiter	is
Jupiter	the
is	Jupiter
is	the
is	largest
the	Jupiter
the	is
the	largest
the	planet

Input	Label
largest	is
largest	the
largest	planet
largest	in
planet	the
planet	largest
planet	in
planet	the
...	...

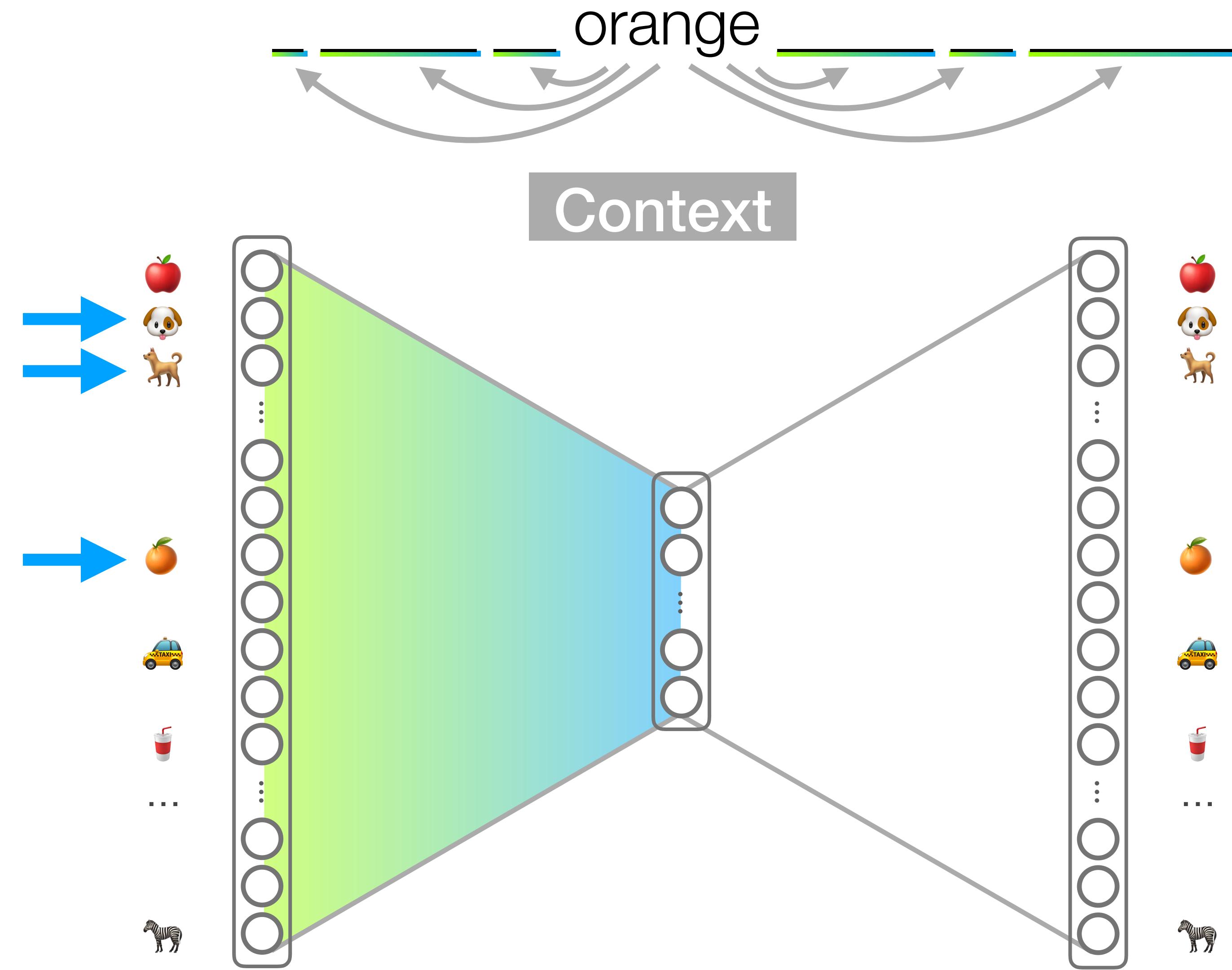
# Word2Vec

## Continuous Skip-gram - The Model



# (Word) Embeddings

## Intuition



# Word2Vec

## Hyperparameters

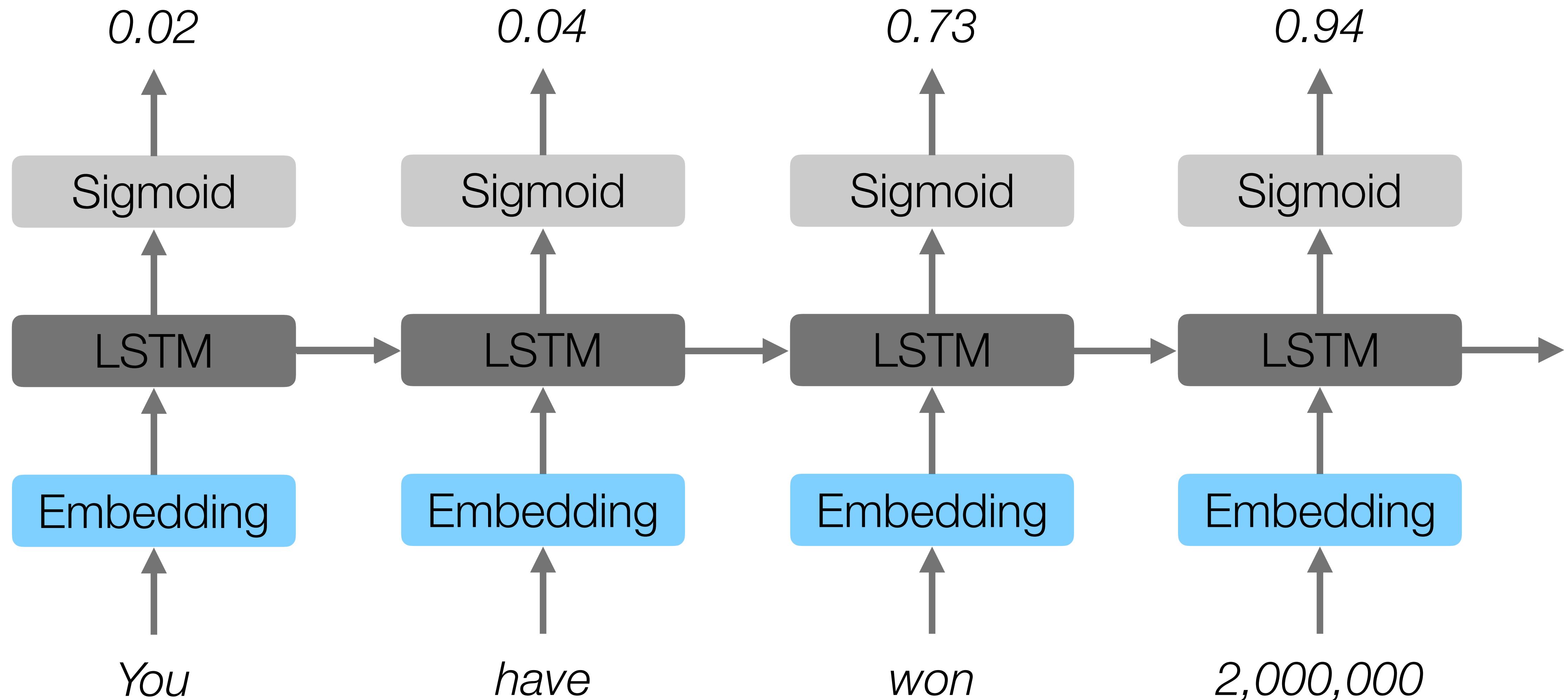
- Model architecture
  - Continuous Bag-of-Words (CBOW)
  - Continuous Skip-gram
- Vocabulary
- Embedding dimension
  - Typically 25 - 300
- Context window size
  - Typically 5 - 10
- Sampling techniques
  - Subsampling of frequent words
  - Negative sampling

# Embeddings in Machine Learning Frameworks

```
from keras.layers import Embedding  
  
# ...  
model.add(Embedding(input_dim=len(vocabulary), output_dim=50))  
# ...
```

```
from torch import nn  
  
# ...  
self.embedding = nn.Embedding(num_embeddings=len(vocabulary),  
                           embedding_dim=50)  
# ...
```

# NLP Use Case: Spam Detection



# Airbnb



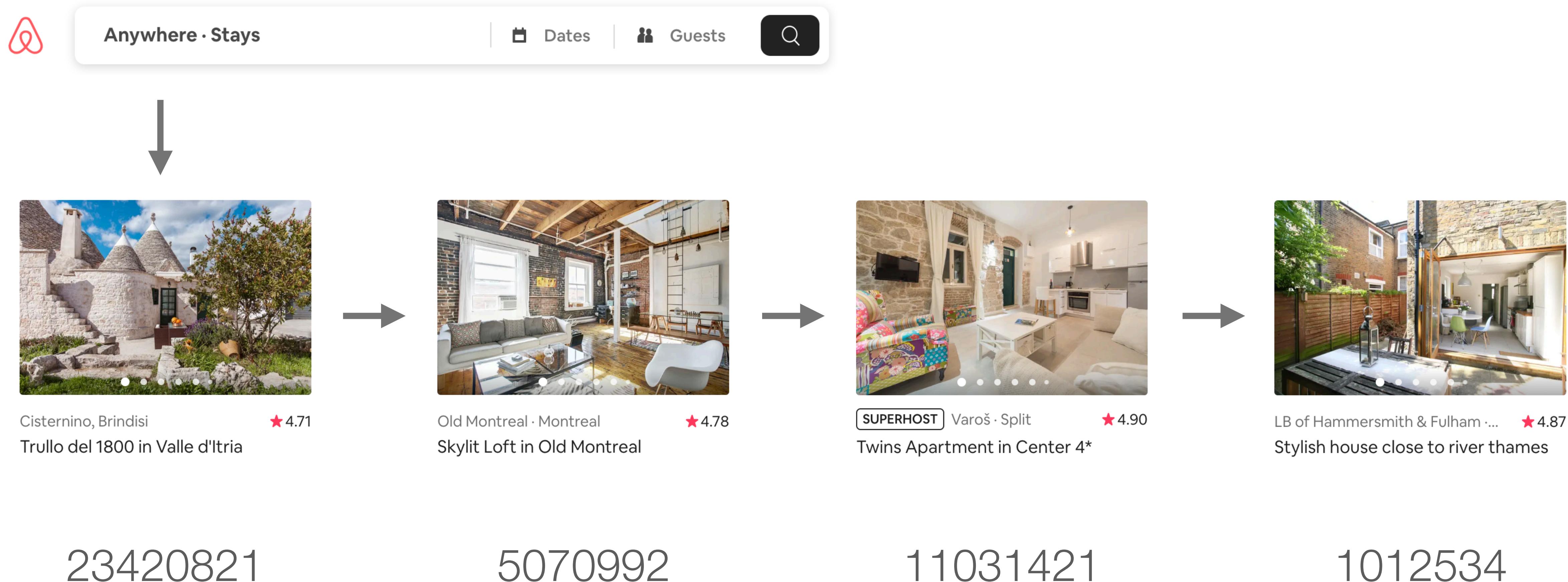
I SETTE CONI - TRULLO EDERA  
Entire home/apt · 2 beds · 4 guests · Business Travel Ready



## Similar Listings

Similar Listings	
A thumbnail image of a trullo with a white heart icon overlaid. Navigation arrows are present on either side.	\$99 ⚡ Trullo of 1800 in the Itria Valley Entire home/apt · 5 beds · 5 guests
A thumbnail image of a trullo with a white heart icon overlaid. Navigation arrows are present on either side.	\$76 ⚡ 🌟 Trulli Tramonti d'Itria - Trullo Luna Entire home/apt · 3 beds · 4 guests
A thumbnail image of a trullo with a white heart icon overlaid. Navigation arrows are present on either side.	\$67 ⚡ 🌟 I TRULLINI OSTUNI - MARTINA FRANCA Entire home/apt · 1 bed · 2 guests

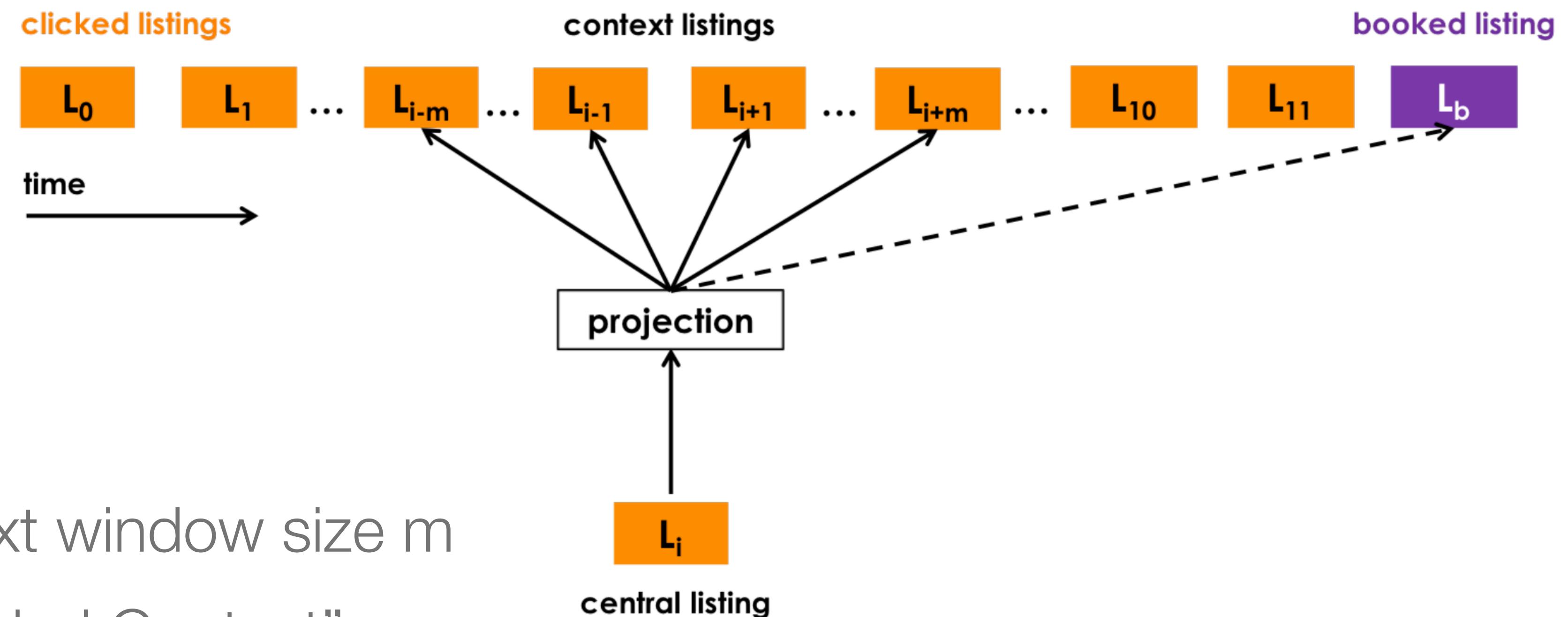
# Airbnb



# Listing Embeddings @ Airbnb



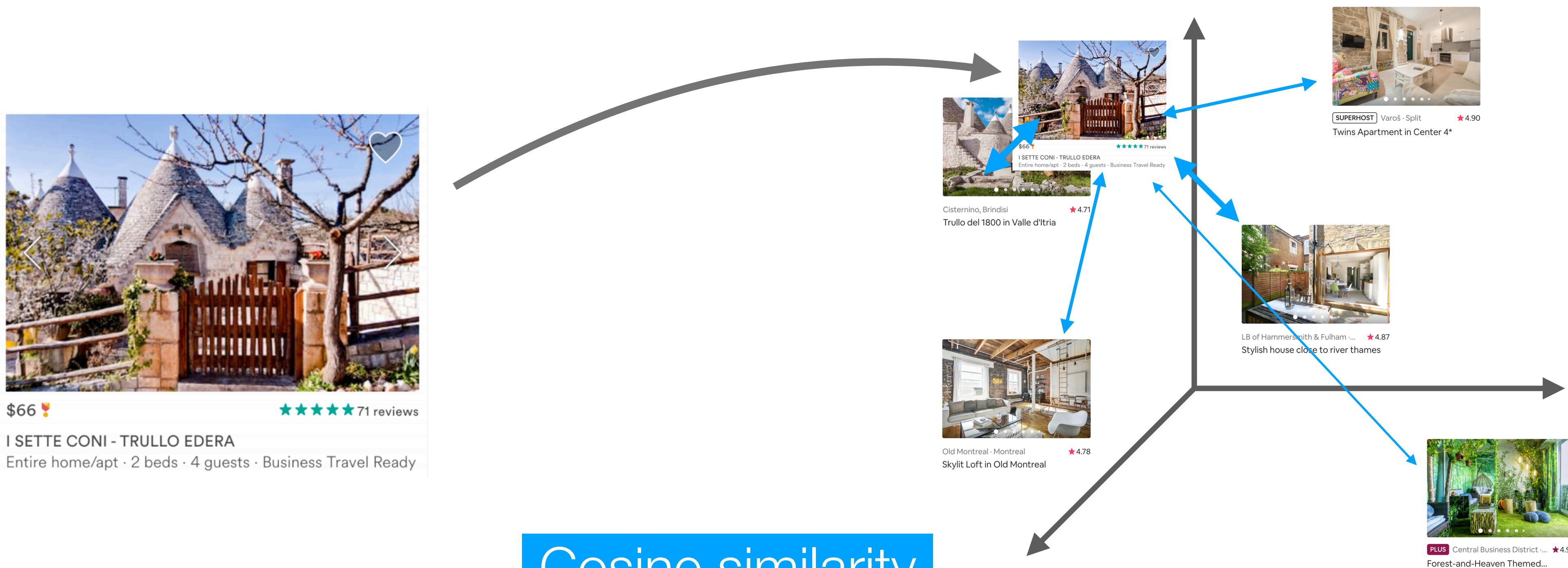
“Click Session”



# Listing Embeddings @ Airbnb



Finding similar Listings

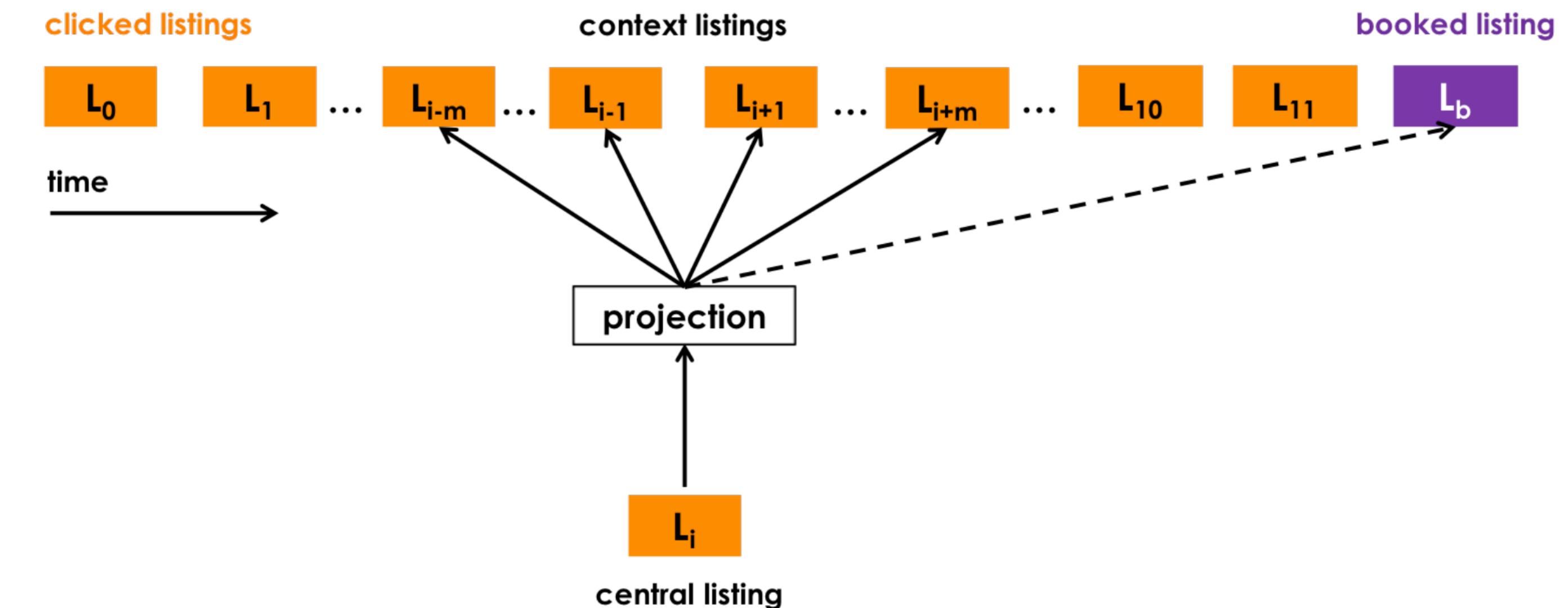


$$\text{similarity}(L_i, L_j) = \cos(\theta_{L_i, L_j}) = \frac{\mathbf{L}_i \cdot \mathbf{L}_j}{\|\mathbf{L}_i\| \|\mathbf{L}_j\|} \in [-1, 1]$$

# Listing Embeddings @ Airbnb



- Context window size  $m = 5$
- Embedding dimension: 32
- Vocabulary: 4.5M listings
- Training data: 800M click sessions
- Daily re-training from scratch with sliding window training data



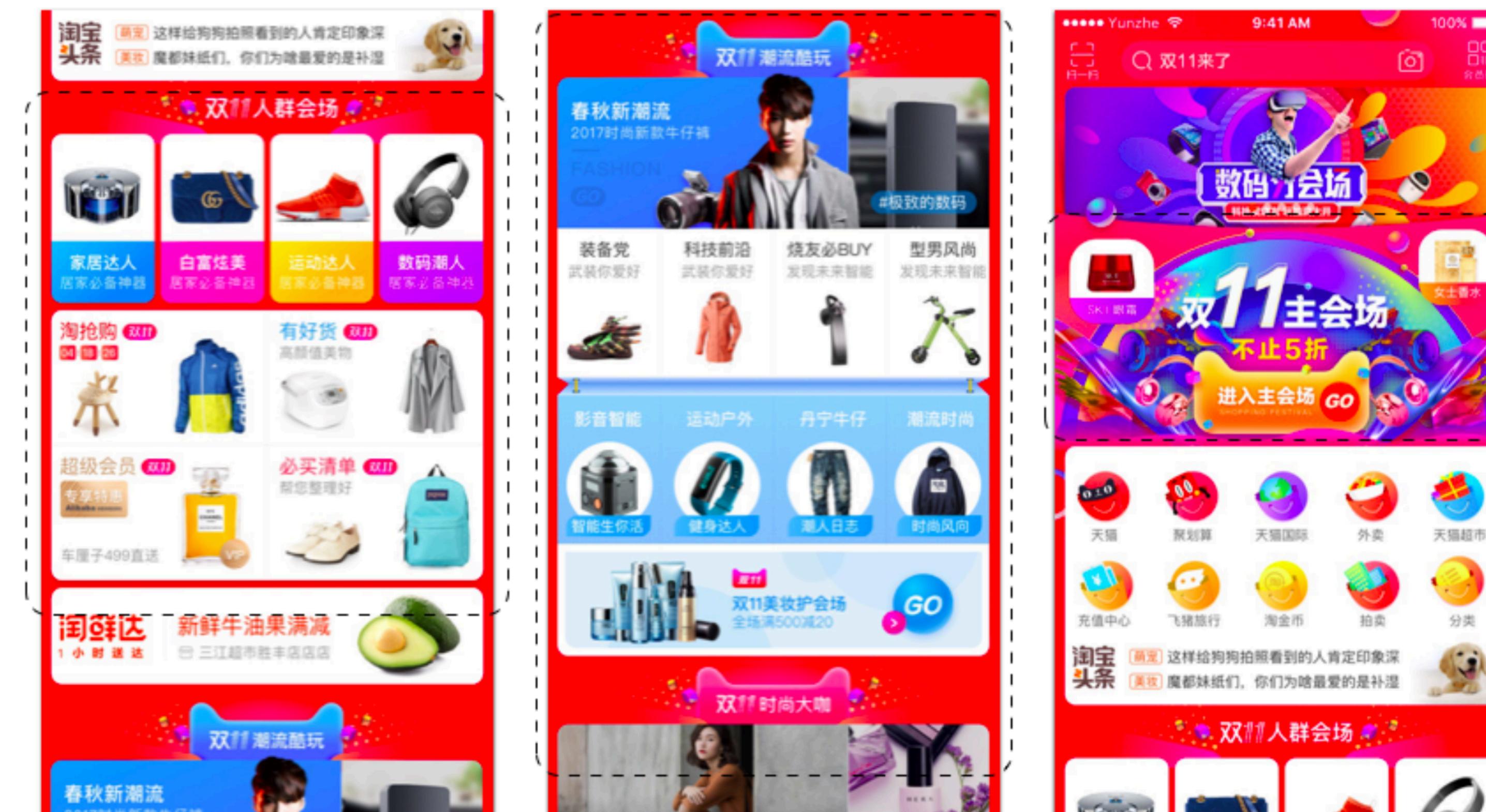
⇒ “21% increase in Similar Listing carousel CTR”

# Alibaba (Taobao)

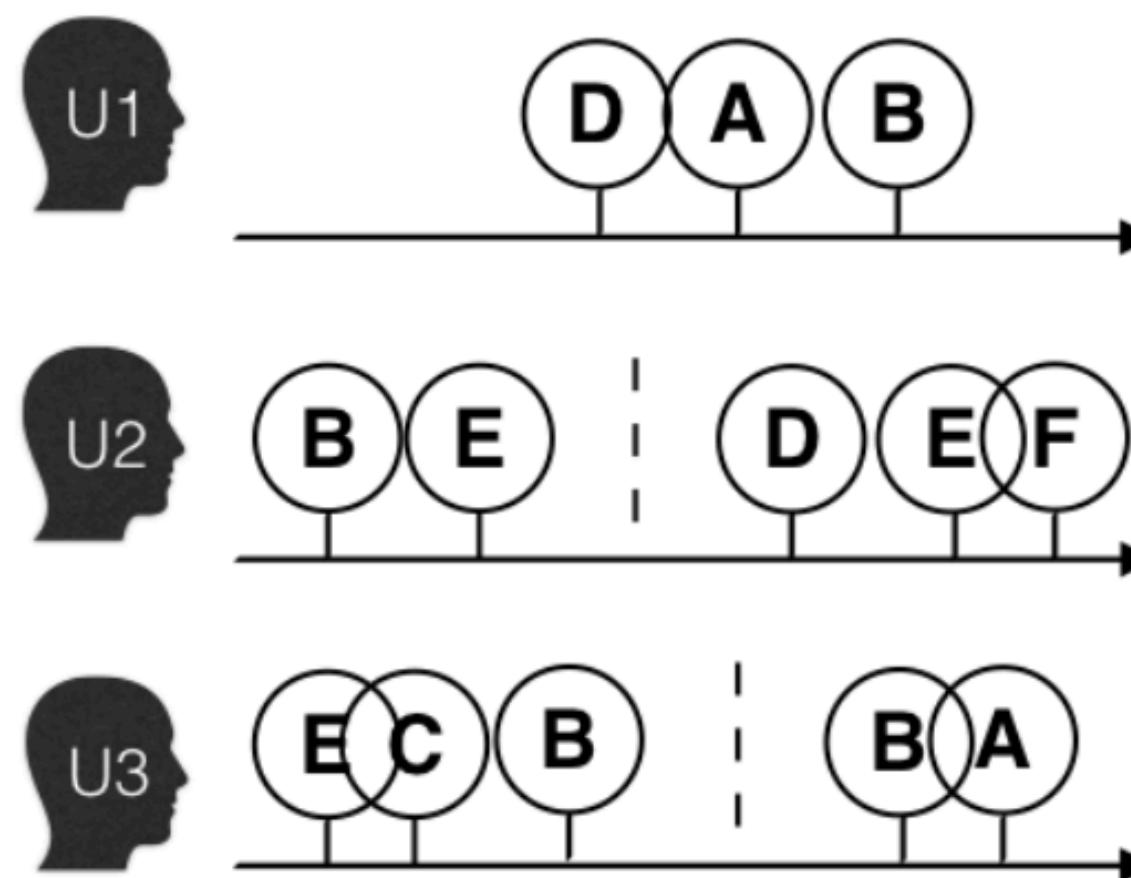
- “eBay” of China



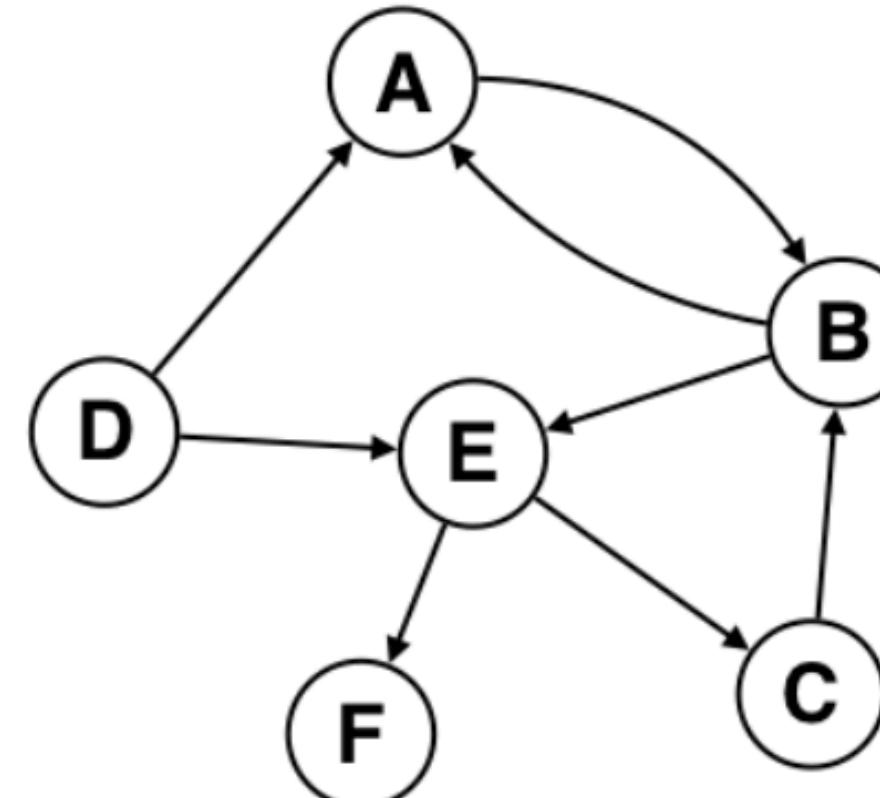
淘宝网  
Taobao.com



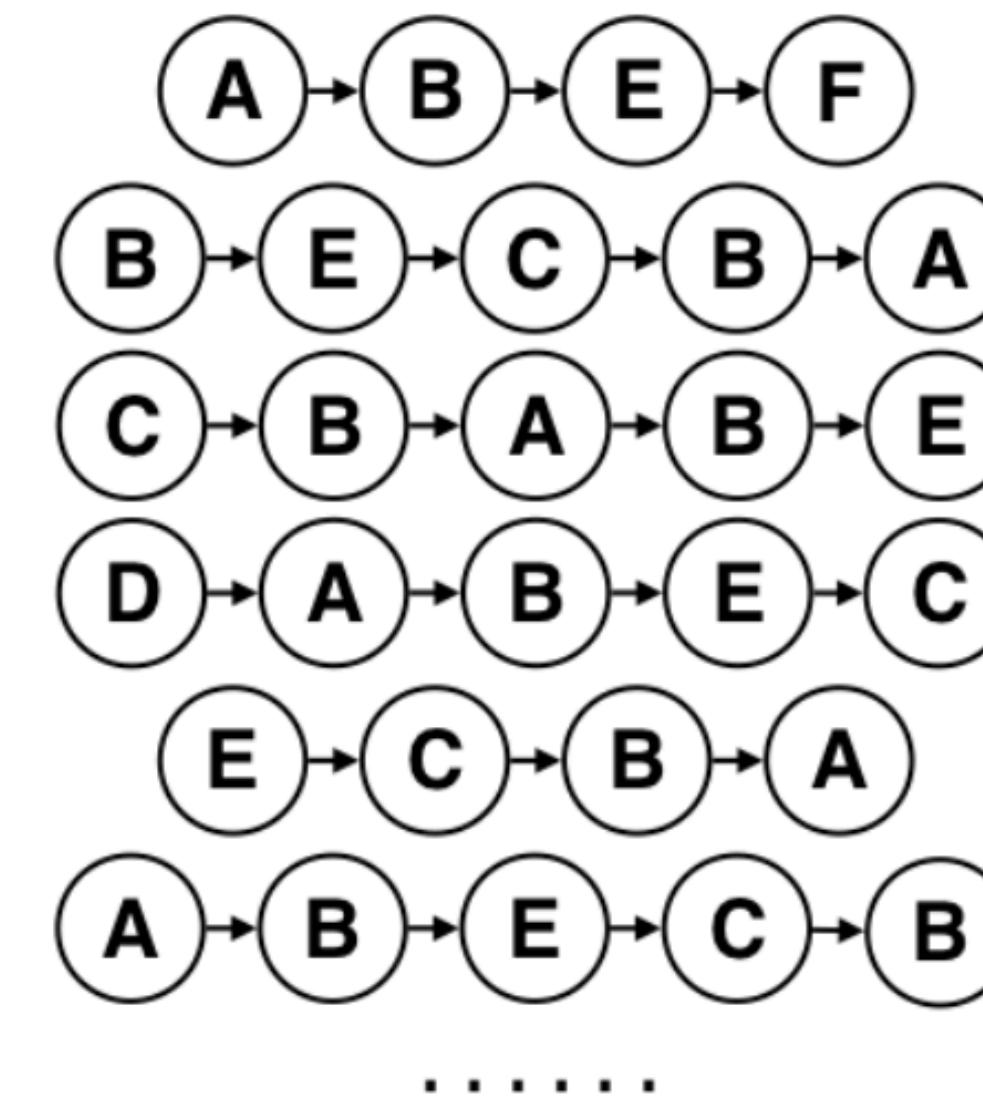
# Graph Embedding (Deep Walk)



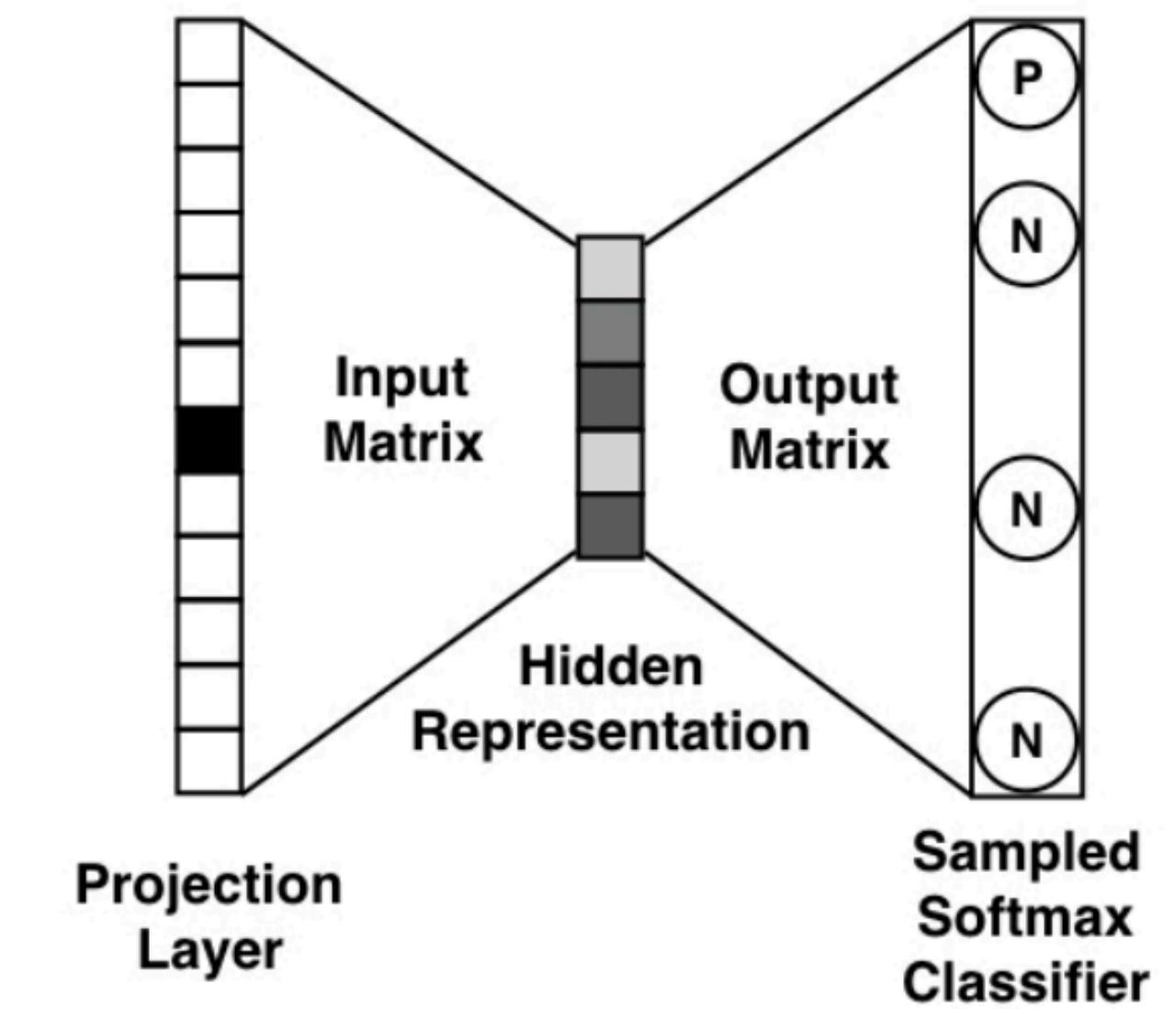
(a) Users' behaviours sequences



(b) Item graph construction



(c) Random walk generation

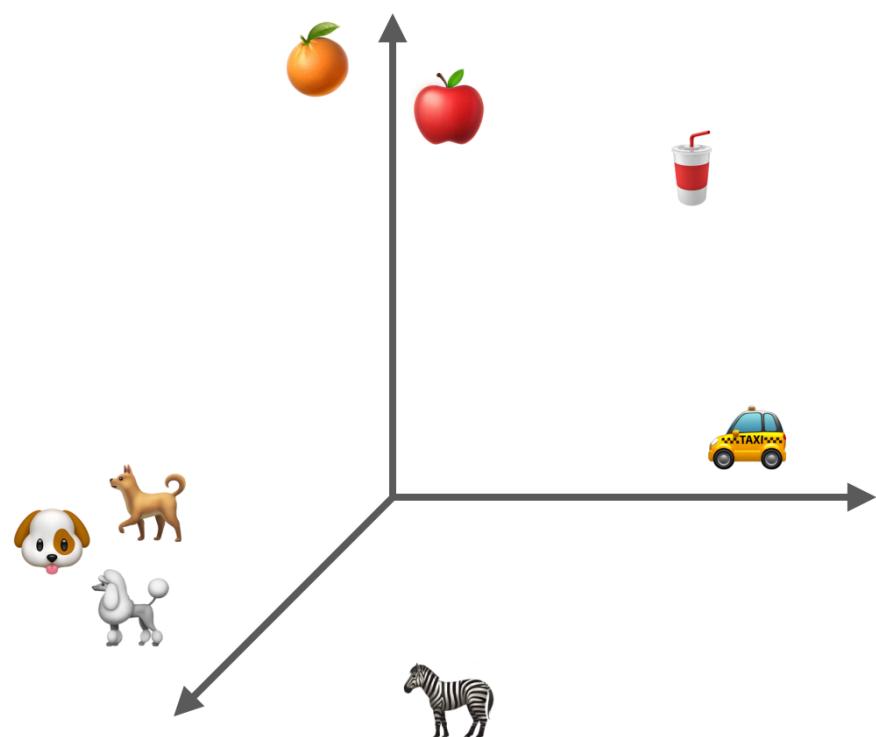


(d) Embedding with Skip-Gram

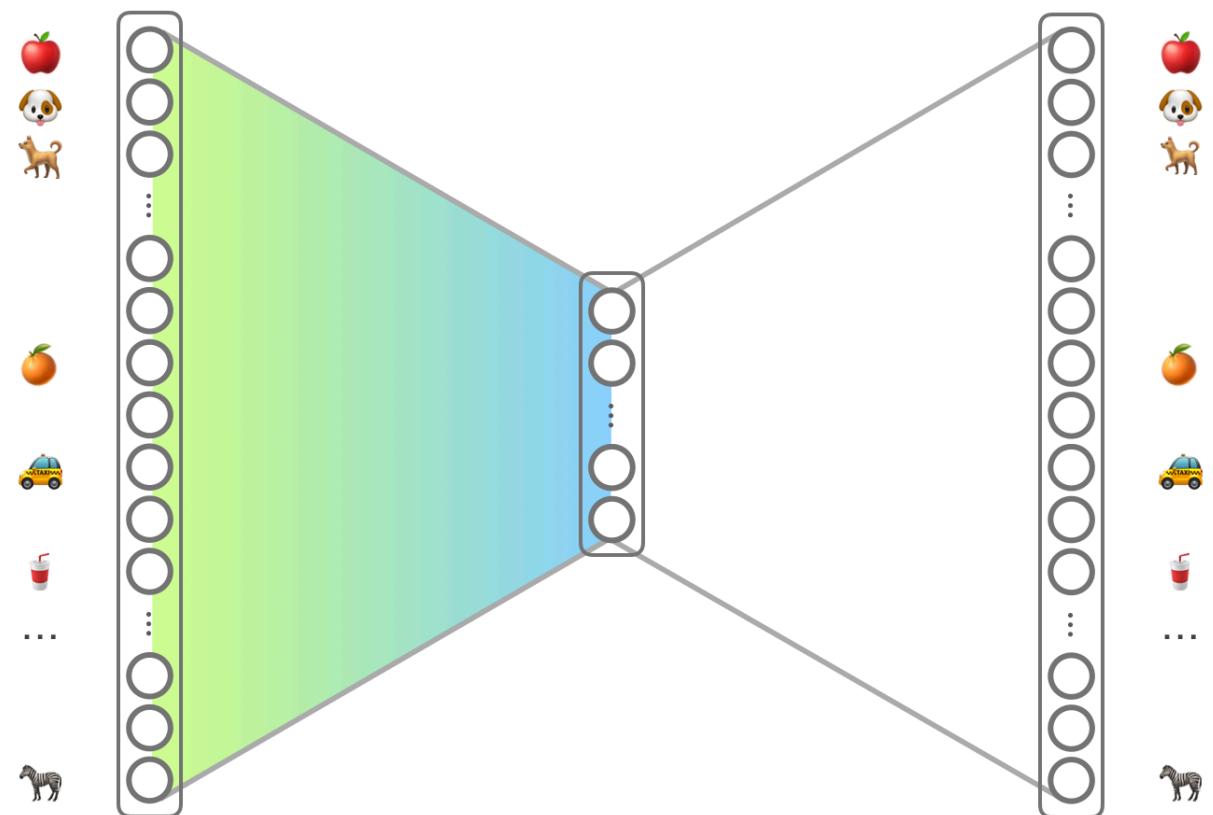
- Context window size: 5
- Embedding dimension: 160
- Vocabulary: ~2B items

# The Wonderful World of Embeddings

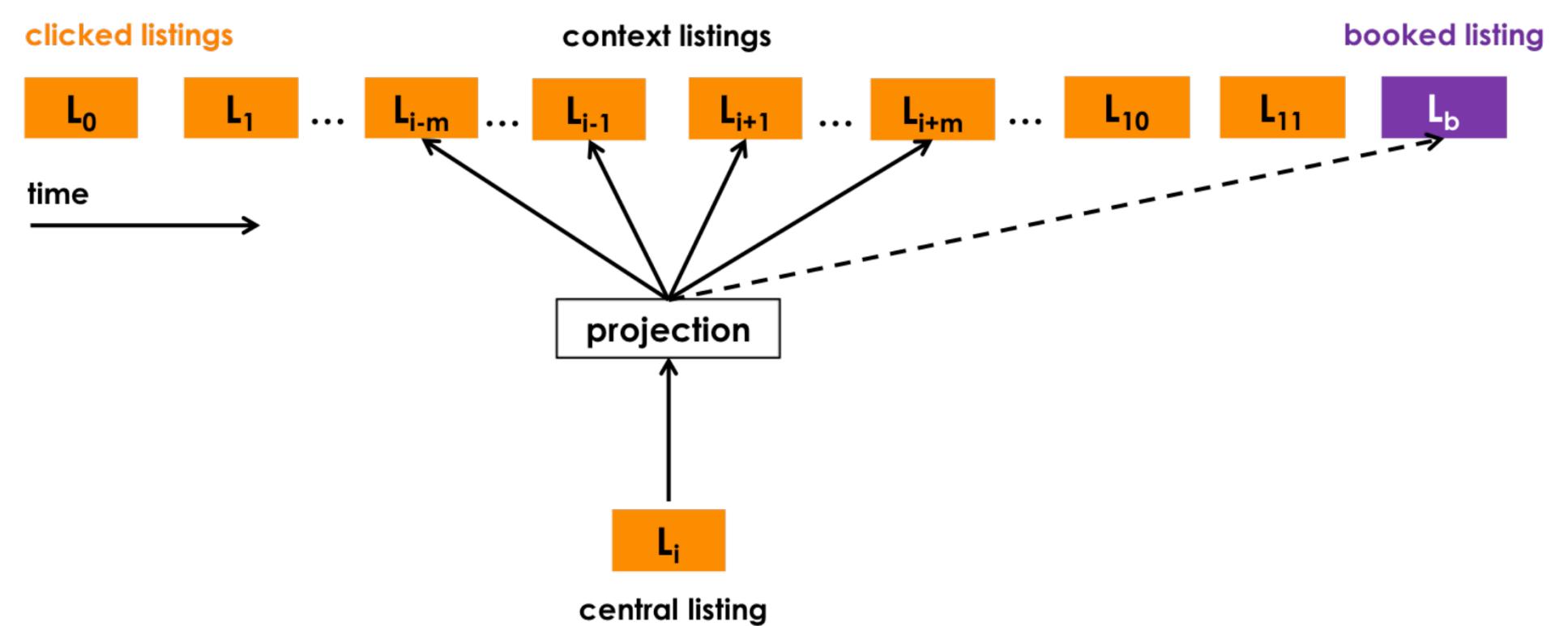
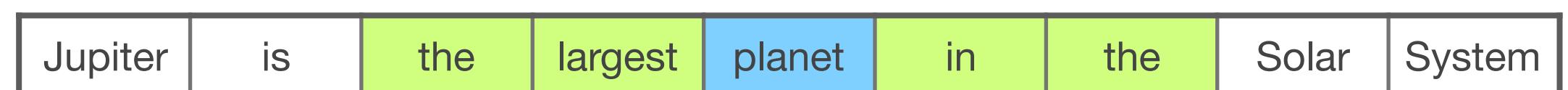
Embedding



Skip-gram Model Architecture



Skip-gram



Franz Diebold, M.Sc. (TUM)

[github.com/franzdiebold](https://github.com/franzdiebold)

[linkedin.com/in/franz-diebold](https://linkedin.com/in/franz-diebold)