

# Neural Network Language Model

Lecture # 4

# Neural Network Language Model

# Language Model

*You shall know a word by the company it keeps*

—Firth, J. R. 1957:11

# Language Model

Fill in the blank:

... a \_\_\_\_\_ ...

car

cars

water

cat

# Language Model

Fill in the blank:

... a \_\_\_\_\_ ...

*car* and *cat* both work

car

cars

water

cat

# Language Model

Fill in the blank:

... a \_\_\_\_\_ ...

*car* and *cat* both work

car

cars

water

cat

John is driving a \_\_\_\_\_ ...

# Language Model

Fill in the blank:

... a \_\_\_\_\_ ...

*car and cat both work*

car

cars

water

cat

John is driving a \_\_\_\_\_ ...

*only car works here*

# Language Model

Fill in the blank:

... a \_\_\_\_\_ ...

*car and cat both work*

car

cars

water

cat

John is driving a \_\_\_\_\_ ...

Similarly, machines use the context to predict the next words



# Language Model

You chose “driving a car” because you’ve seen that phrase more frequently

“driving a cat” is not a common phrase

# Language Model

Fill in the blank:

This \_\_\_\_\_ is going at 100 km/hours

car

bicycle

Car at 100km/hours is more probable than a bicycle

# Language Model

Language model defines  
“how probable a sentence is”

# Language Model

Let's look at the example again

How probable is:

John is driving a car vs. John is driving a cat

In other words, what is the probability to predict cat or car given the context "John is driving a"

# Language Model

**Q:** Can we see language modeling as a classification problem?

# Language Model

**Q:** Can we see language modeling as a classification problem?

**A:** Yes! We are just predicting which **word** (“class”) is coming next.

# Language Model

**Q:** Can we see language modeling as a classification problem?

**A:** Yes! We are just predicting which **word** (“class”) is coming next.

$$p(\langle s \rangle \text{ Dan likes ham } \langle /s \rangle) =$$

# Language Model

**Q:** Can we see language modeling as a classification problem?

**A:** Yes! We are just predicting which **word** (“class”) is coming next.

$$p(\langle s \rangle \text{ Dan likes ham } \langle /s \rangle) = p(\text{Dan} \mid \langle s \rangle)$$

Predict **Dan** given  $\langle s \rangle$



# Language Model

**Q:** Can we see language modeling as a classification problem?

**A:** Yes! We are just predicting which **word** (“class”) is coming next.

$$p(\langle s \rangle \text{ Dan likes ham } \langle /s \rangle) = p(\text{Dan} | \langle s \rangle)$$

Predict **Dan** given  **$\langle s \rangle$**

What is the probability of **Dan** given  **$\langle s \rangle$** ?

# Language Model

**Q:** Can we see language modeling as a classification problem?

**A:** Yes! We are just predicting which **word** (“class”) is coming next.

$$p(\langle s \rangle \text{ Dan likes ham } \langle /s \rangle) = p(\text{Dan} | \langle s \rangle) \cdot p(\text{likes} | \text{Dan})$$

Predict **likes** given **Dan**

# Language Model

**Q:** Can we see language modeling as a classification problem?

**A:** Yes! We are just predicting which **word** (“class”) is coming next.

$$p(\langle s \rangle \text{ Dan likes ham } \langle /s \rangle) = p(\text{Dan} | \langle s \rangle) \\ \cdot p(\text{likes} | \text{Dan}) \\ \cdot p(\text{ham} | \text{likes})$$

Predict **ham** given **likes**

# Language Model

**Q:** Can we see language modeling as a classification problem?

**A:** Yes! We are just predicting which **word** (“class”) is coming next.

$$p(\langle s \rangle \text{ Dan likes ham } \langle /s \rangle) = p(\text{Dan} | \langle s \rangle)$$

- $p(\text{likes} | \text{Dan})$
- $p(\text{ham} | \text{likes})$
- $p(\langle /s \rangle | \text{ham})$

# Language Model

Words represent classes that we want to predict!

**Input to the classifier:** previous words i.e. context

**Output:** probability distribution over all possible words, i.e. our vocabulary

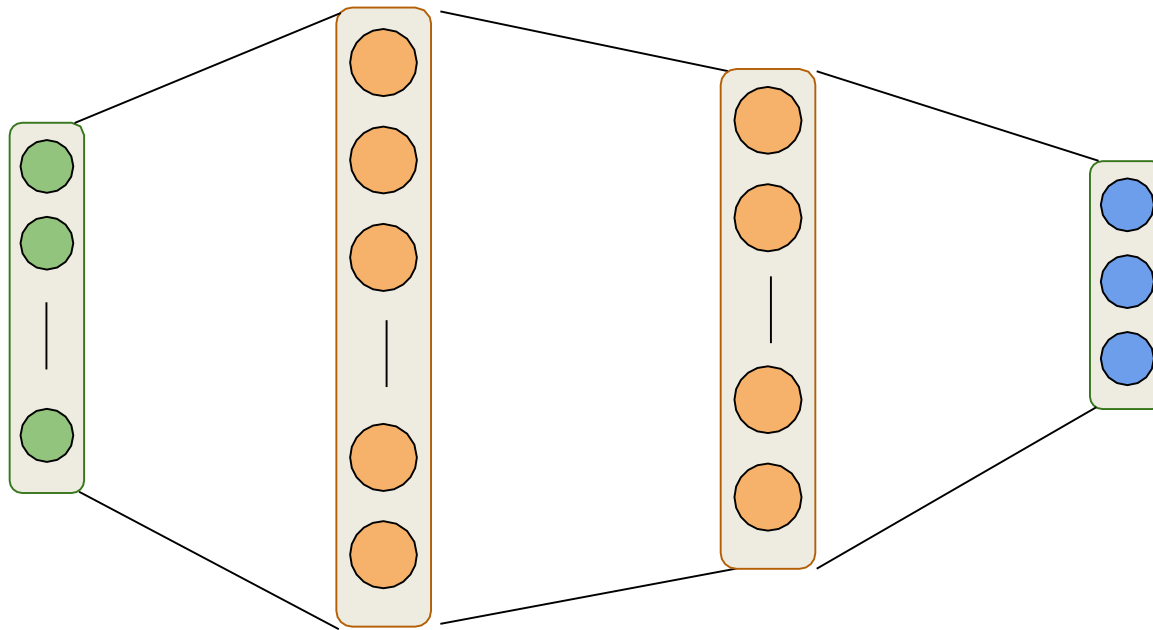
# Neural Network Language Model

Input

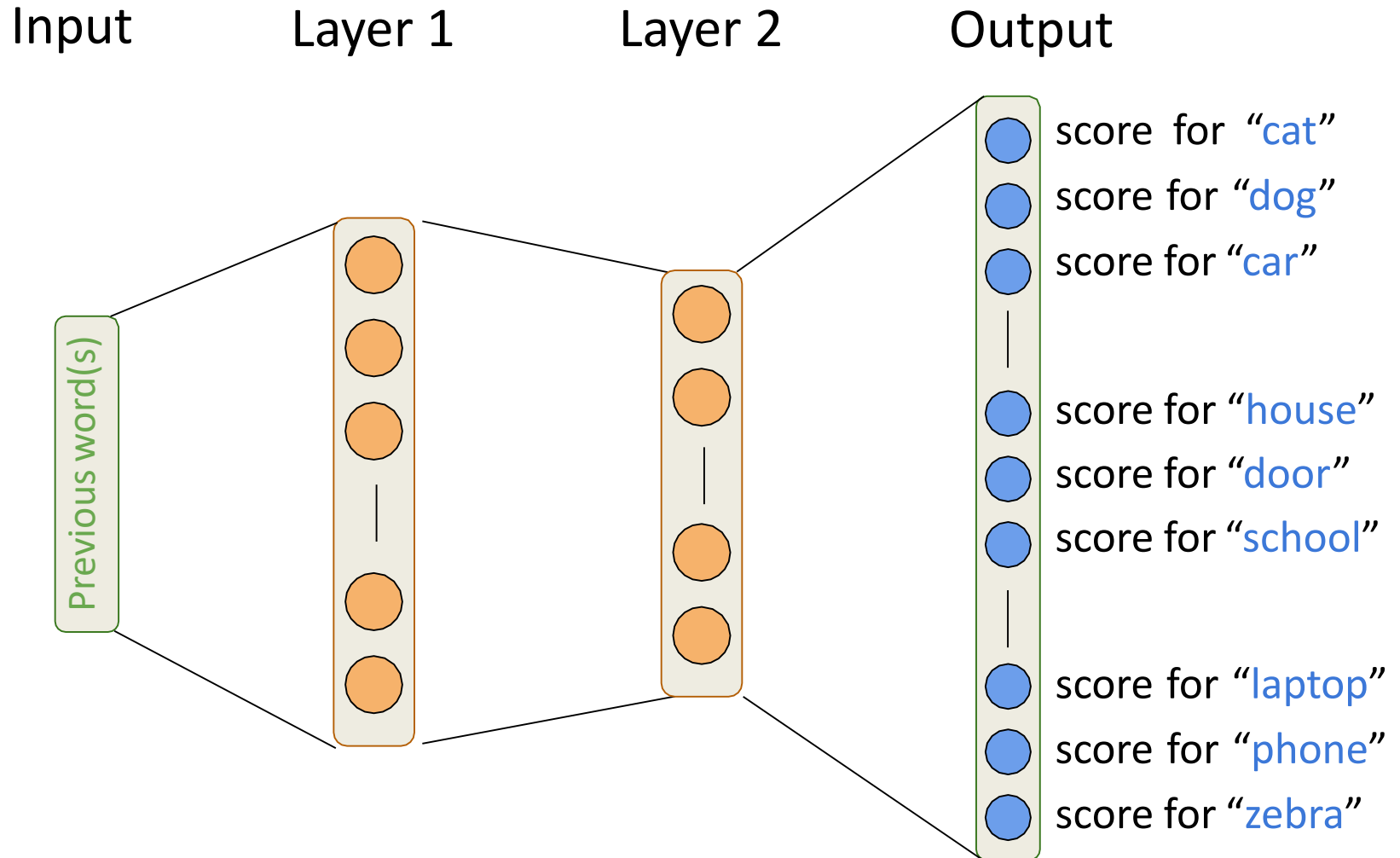
Layer 1

Layer 2

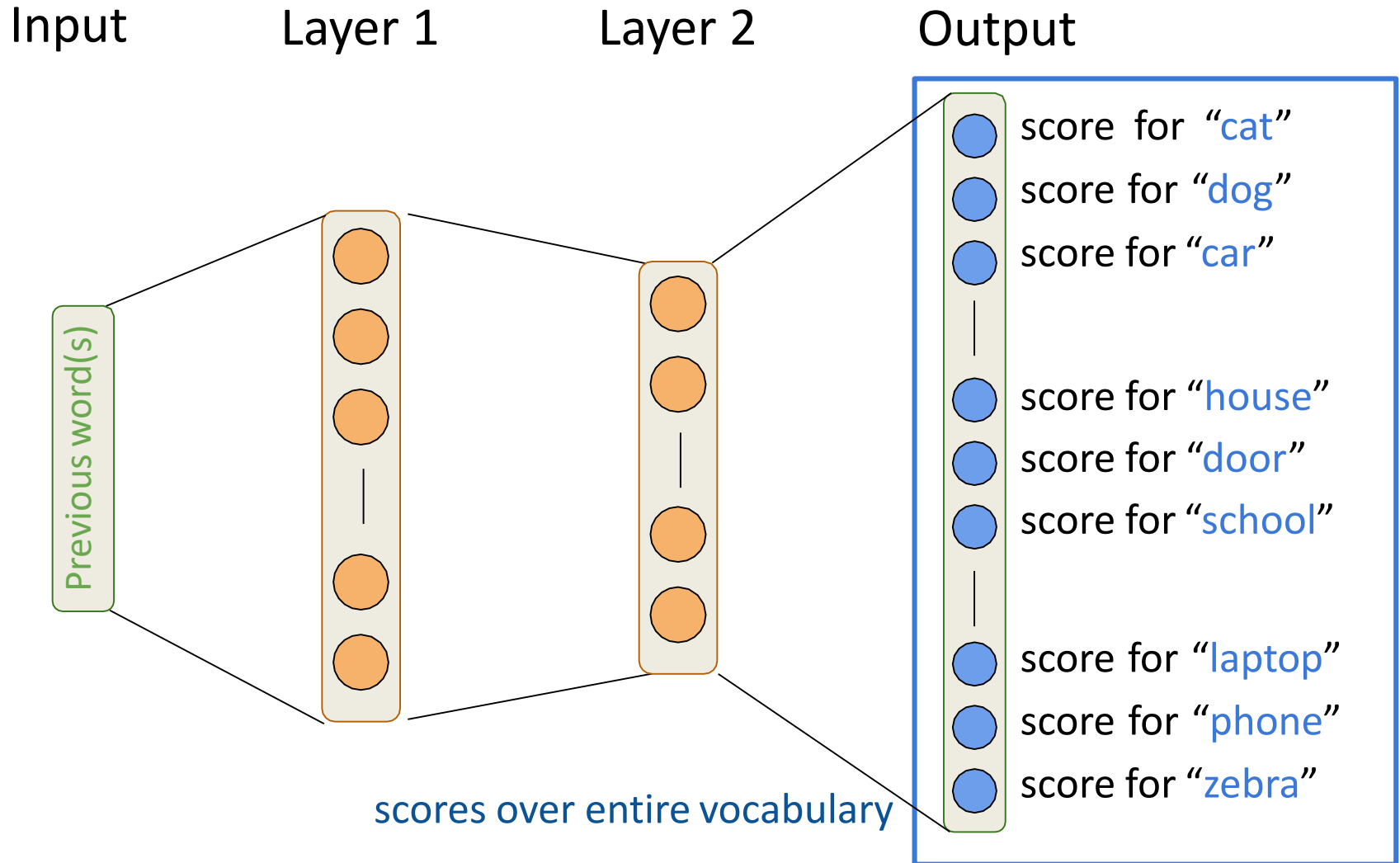
Output



# Neural Network Language Model



# Neural Network Language Model





# Input Representation

Input

Previously we've used a vector as input, where each element of the vector represented some "feature" of the input

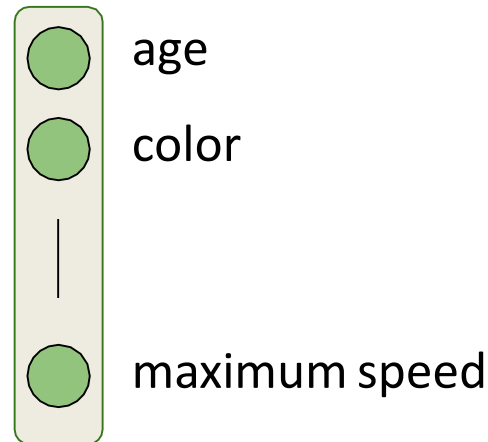
Previous word(s)

# Input Representation

Input

Previously we've used a vector as input, where each element of the vector represented some "feature" of the input

Previous word(s)



Car

# Input Representation

Input

Can we represent a word as a feature vector?

Previous word(s)

“Universität”



?



# One Hot Vector Representation

- Every word can be represented as a ***one hot vector***

# One Hot Vector Representation

- Every word can be represented as a ***one hot vector***
- Suppose the total number of unique words in the corpus is 10,000

# One Hot Vector Representation

- Every word can be represented as a ***one hot vector***
- Suppose the total number of unique words in the corpus is 10,000
- Assign each word a unique index:

```
Universität: 1  
cat: 2  
house: 3  
car: 4  
⋮  
apple: 10,000
```

# One Hot Vector Representation

- Every word can be represented as a ***one hot vector***
- Suppose the total number of unique words in the corpus is 10,000
- Assign each word a unique index:

```
Universität: 1  
cat: 2  
house: 3  
car: 4  
⋮  
apple: 10,000
```

Dictionary

```
cat =  $\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$ 
```

One-hot representation

# One Hot Vector Representation

- Every word can be represented as a ***one hot vector***
- Suppose the total number of unique words in the corpus is 10,000
- Assign each word a unique index:

Only index that represents the input word will be one

$$cat = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

One-hot representation



# One Hot Vector Representation

- Every word can be represented as a ***one hot vector***
- Suppose the total number of unique words in the corpus is 10,000
- Assign each word a unique index:

Only index that represents the input word will be one

$$\begin{array}{l} \text{cat} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \text{car} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \end{array}$$

One-hot representation

# One Hot Vector Representation

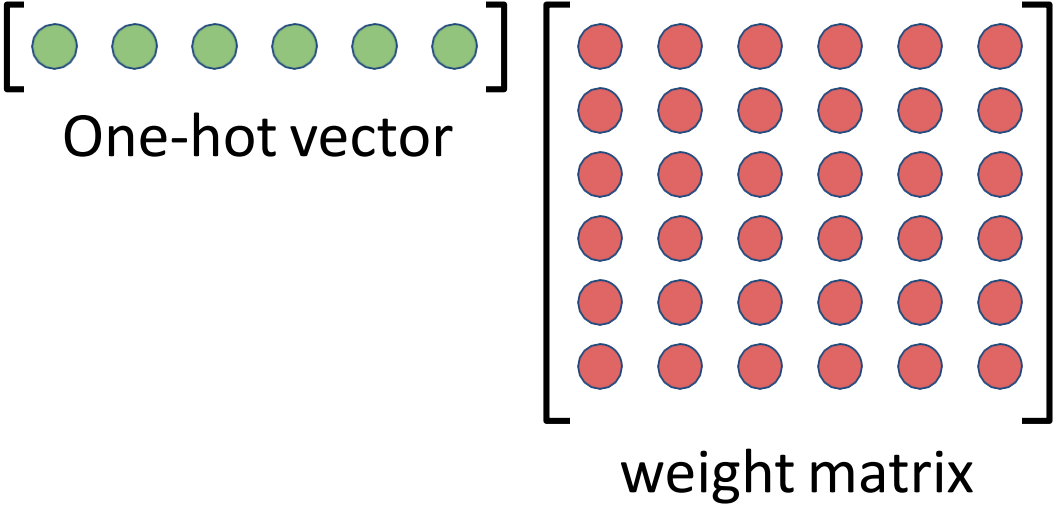
- Every word can be represented as a ***one hot vector***
- Suppose the total number of unique words in the corpus is 10,000
- Assign each word a unique index:

Vector size will be the size of the vocabulary, i.e. 10,000 in this case

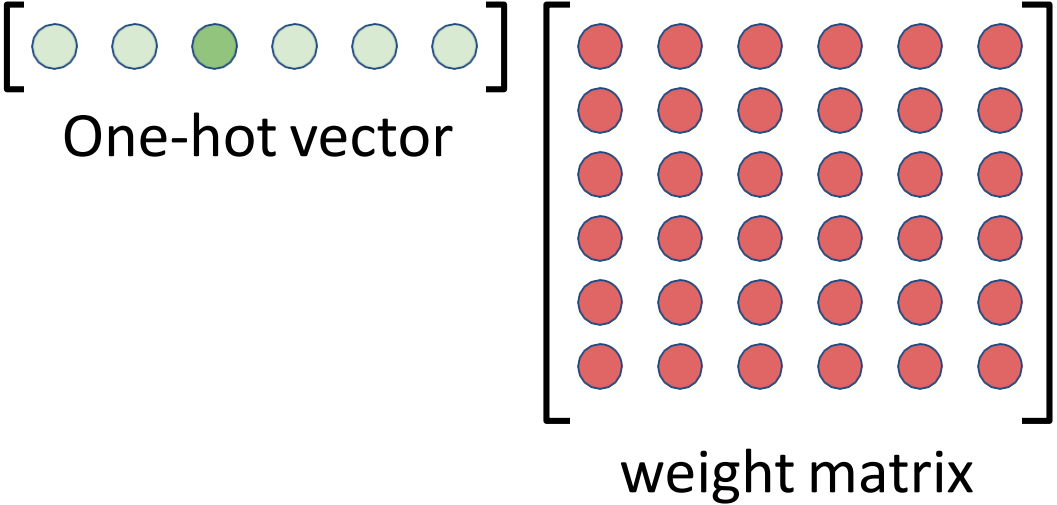
$$\begin{array}{l} \text{cat} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \text{car} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \end{array}$$

One-hot representation

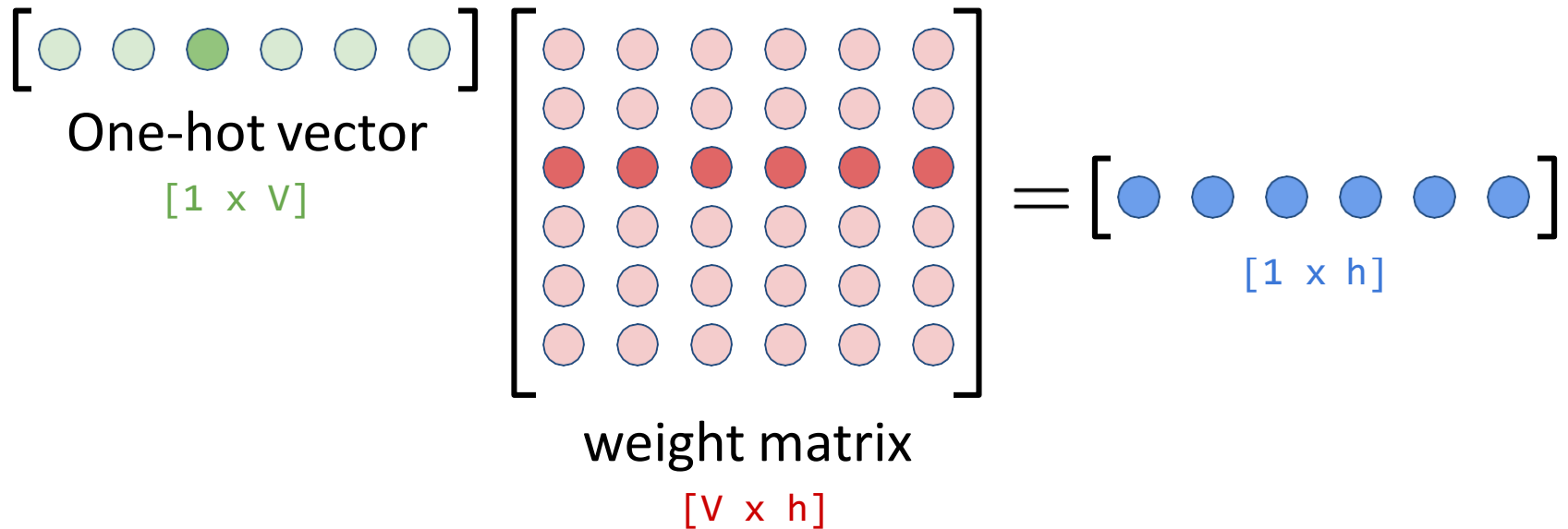
# One Hot Vector Representation



# One Hot Vector Representation



# One Hot Vector Representation



One-hot vector will “turn on” one row of weights

# Higher ngram Vector Representation

- What about representing multiple words?

**Bag of words approach**

# Higher ngram Vector Representation

- What about representing multiple words?

## Bag of words approach

**Bigram:** indices of the *two previous words* are 1 in the vector

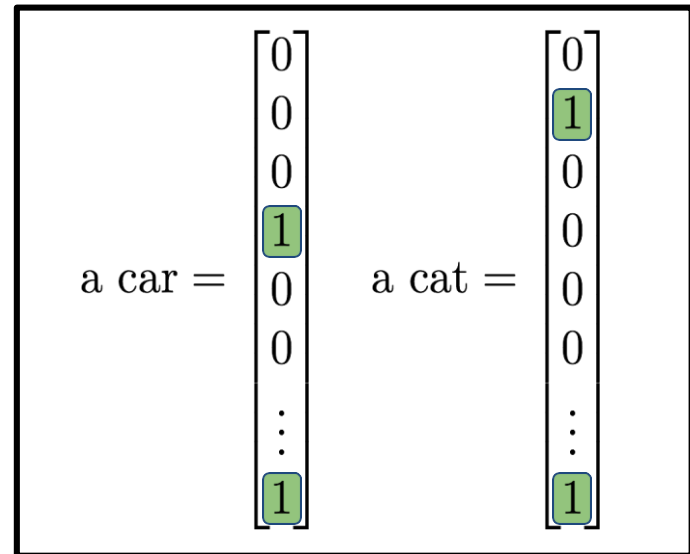
$$\begin{array}{l} \text{a car} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix} \quad \text{a cat} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix} \end{array}$$

# Higher ngram Vector Representation

- What about representing multiple words?

## Bag of words approach

**Bigram:** indices of the *two previous words* are 1 in the vector



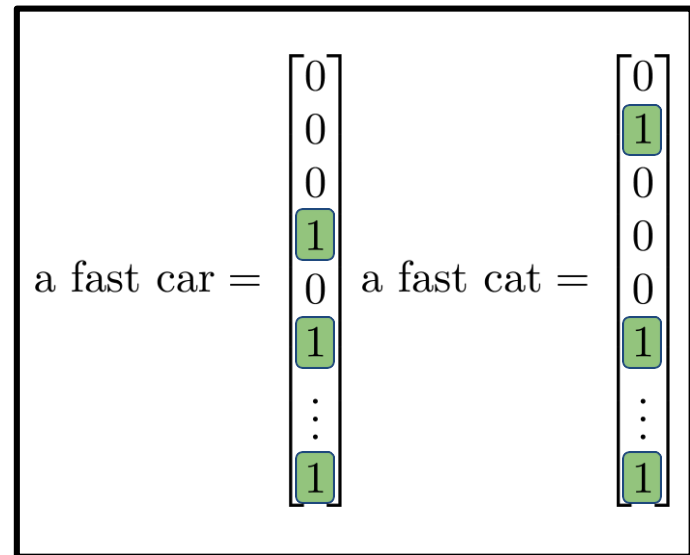


# Higher ngram Vector Representation

- What about representing multiple words?

## Bag of words approach

**Trigram:** indices of the *three previous words* are 1 in the vector



# Higher ngram Vector Representation

- What about representing multiple words?

## Context-aware approach

In the **bag of words** approach, order information is lost!

# Higher ngram Vector Representation

- What about representing multiple words?

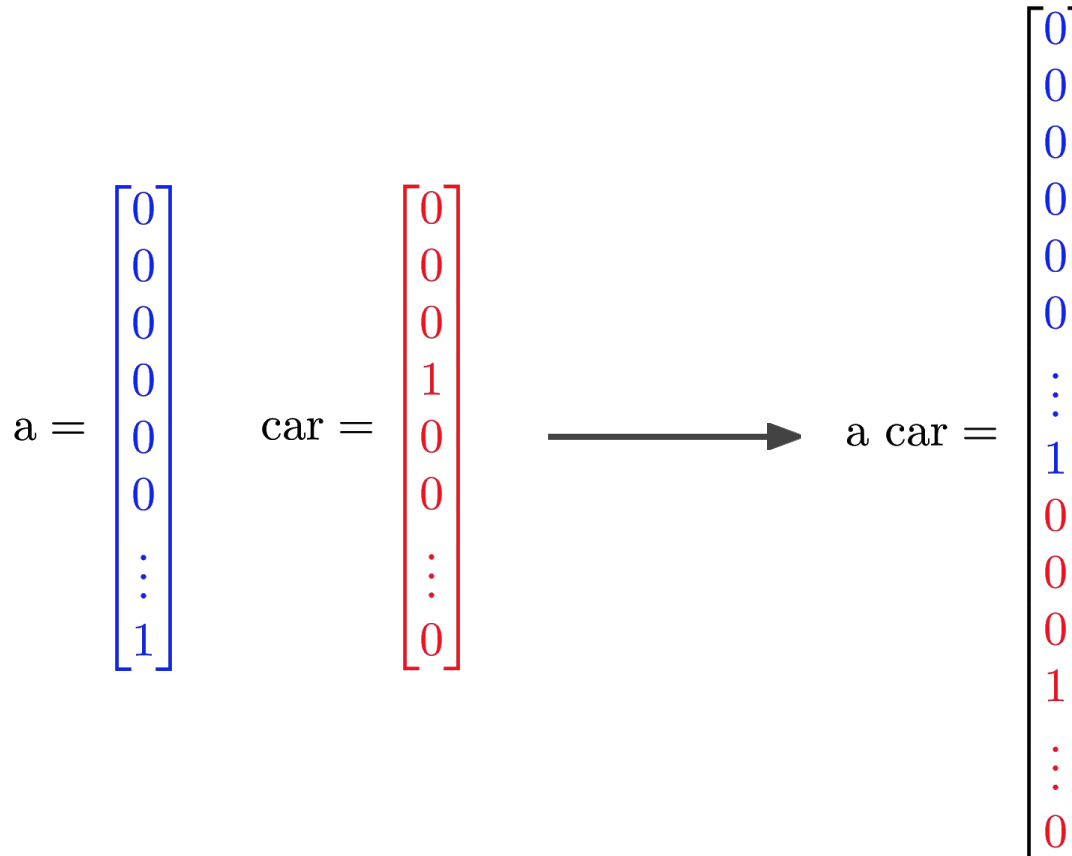
## Context-aware approach

In the **bag of words** approach, order information is lost!

**Solution:** for  $N$  words, concatenate one-hot vectors for each of the words in the correct order

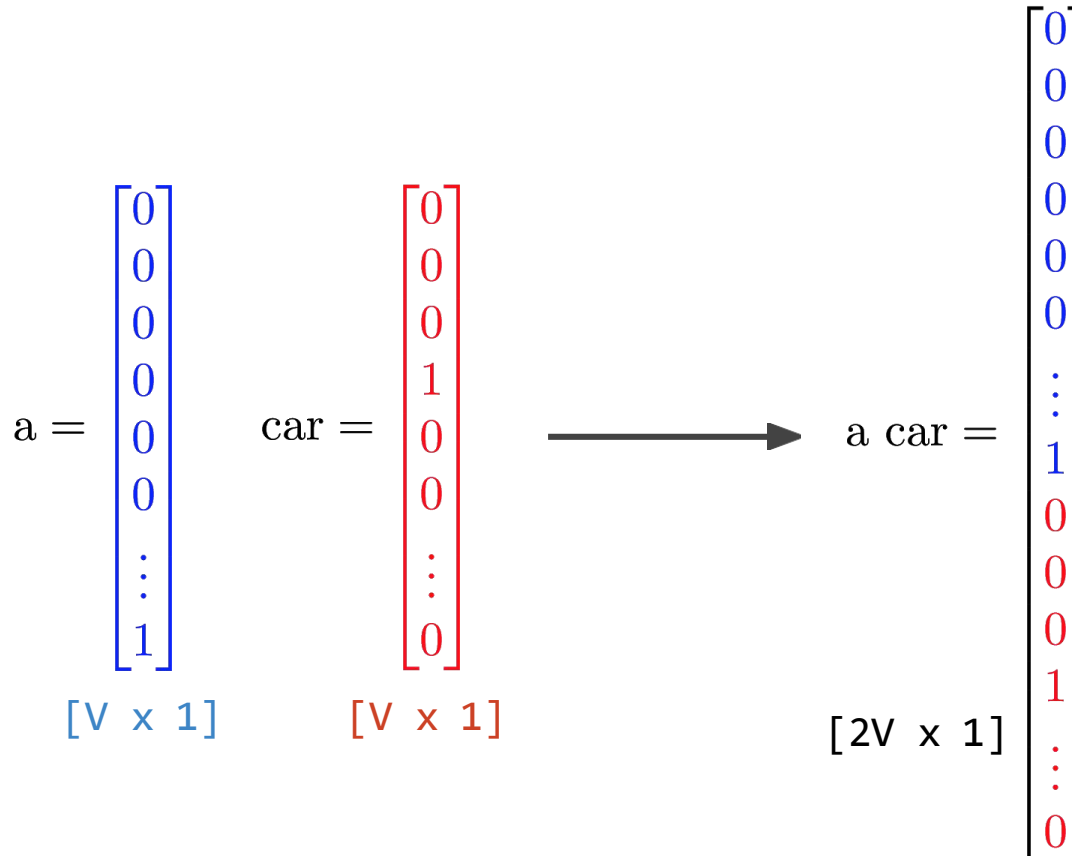
# Higher ngram Vector Representation

## Context-aware approach



# Higher ngram Vector Representation

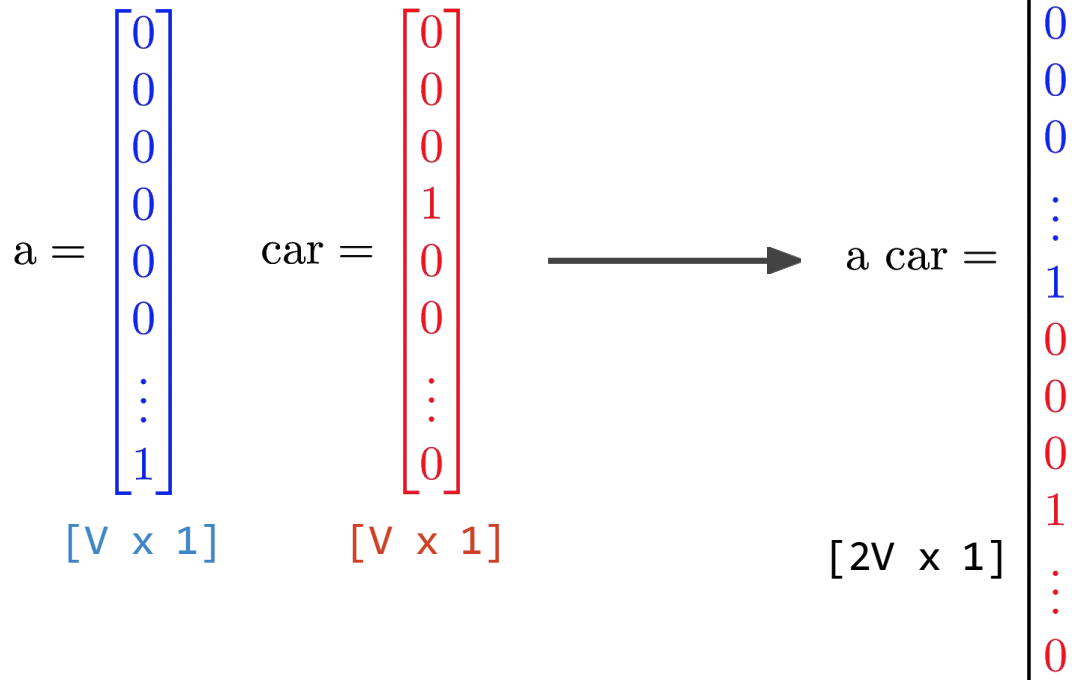
## Context-aware approach



# Higher ngram Vector Representation

## Context-aware approach

input vector length has increased



# Higher ngram Vector Representation

## Context-aware approach

input vector length has increased

- order information is available for the training

Advantages

- long vectors in case of large context size
- number of parameters increases with context size

Disadvantages

# Higher ngram Vector Representation

- Bag of words vs. context-aware approach?
  - Given the disadvantages of the context-aware approach, Bag of words is more commonly used
  - Works well in practice



# Input Representation

Generally, the size of the vocabulary is very large

- Results in very large one-hot vectors!

# Input Representation

Generally, the size of the vocabulary is very large

- Results in very large one-hot vectors!

Some tricks to reduce vocabulary size:

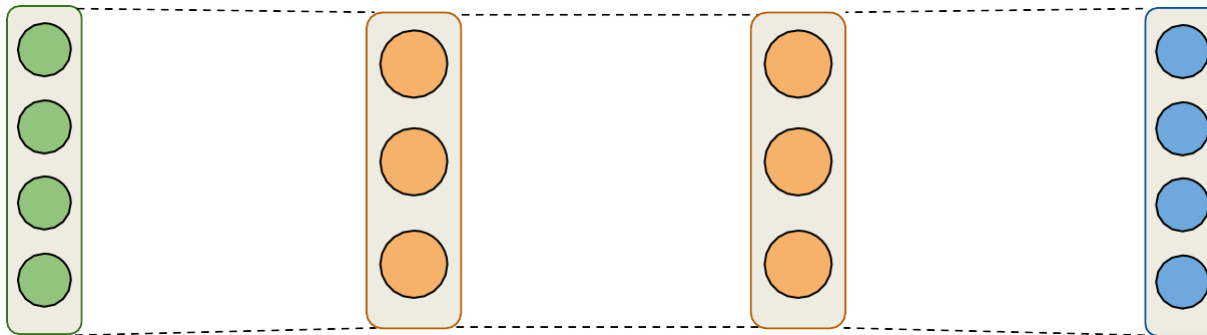
- 1) Take most frequent top words. For example, consider only 10,000 most frequent words and map the rest to a unique token `<UNK>`
- 2) Cluster words
  - a) based on context
  - b) based on linguistic properties

# Neural Network Language Model

Let us look at a complete example:

**Vocabulary:** {"how", "you", "hello", "are"}

**Network Architecture:** 2 hidden layers of size 3 each

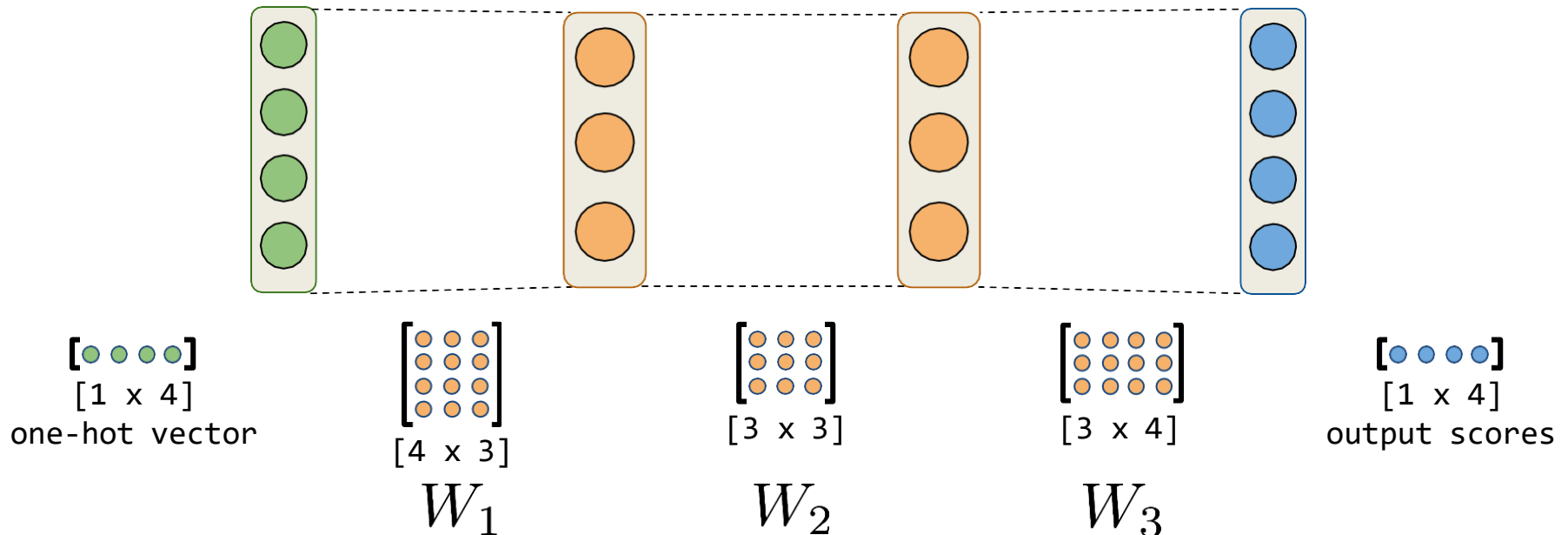


# Neural Network Language Model

Let us look at a complete example:

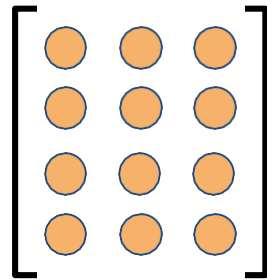
**Vocabulary:** {"how", "you", "hello", "are"}

**Network Architecture:** 2 hidden layers of size 3 each

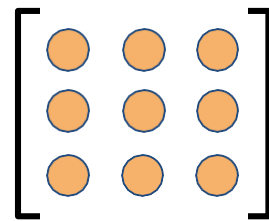


# Neural Network Language Model

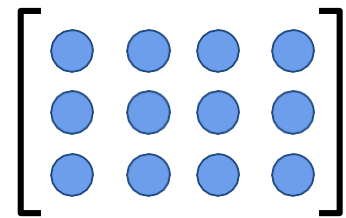
**Vocabulary:** {"how", "you", "hello", "are"}



$W_1$



$W_2$



$W_3$

# Neural Network Language Model

**Vocabulary:** {"how", "you", "hello", "are"}

[○ ○ ● ○]

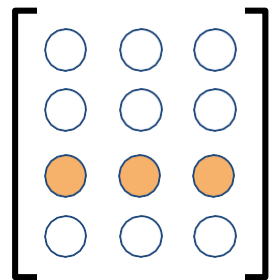
"hello"

# Neural Network Language Model

**Vocabulary:** {"how", "you", "hello", "are"}



"hello"



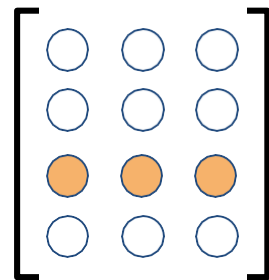
$W_1$

# Neural Network Language Model

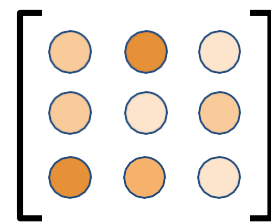
**Vocabulary:** {"how", "you", "hello", "are"}



"hello"



$W_1$

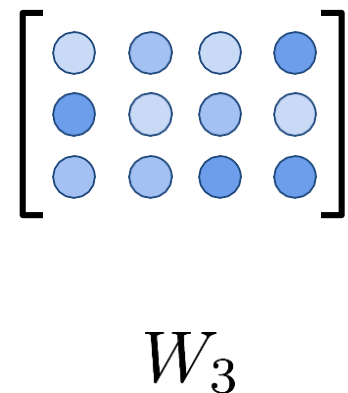
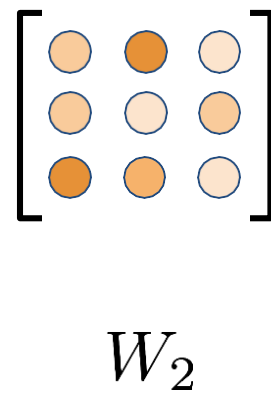
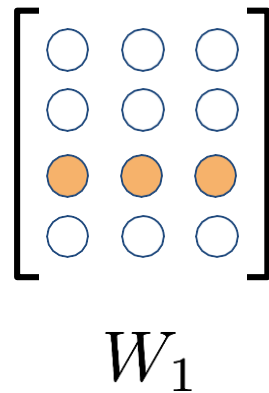
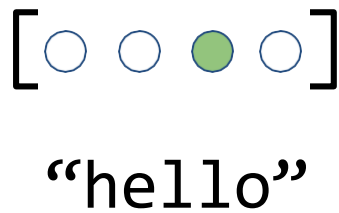


$W_2$



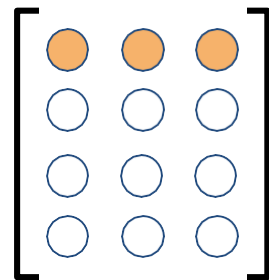
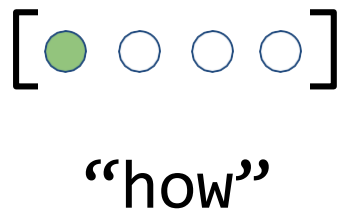
# Neural Network Language Model

**Vocabulary:** {"how", "you", "hello", "are"}

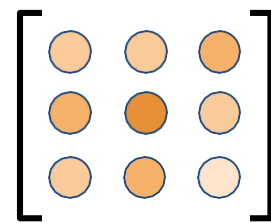


# Neural Network Language Model

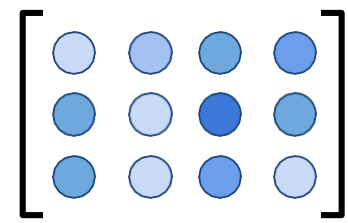
**Vocabulary:** {"how", "you", "hello", "are"}



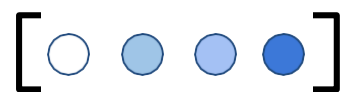
$W_1$



$W_2$



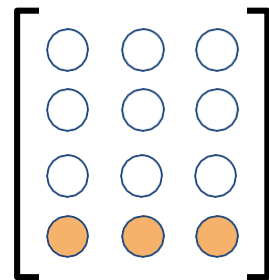
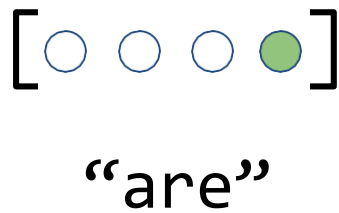
$W_3$



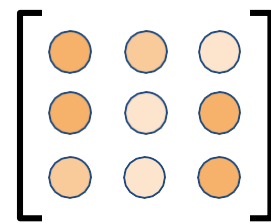
output scores  
max: "are"

# Neural Network Language Model

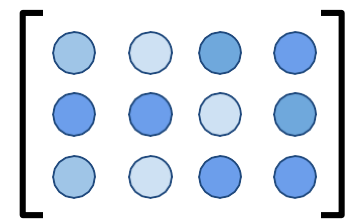
**Vocabulary:** {"how", "you", "hello", "are"}



$W_1$



$W_2$

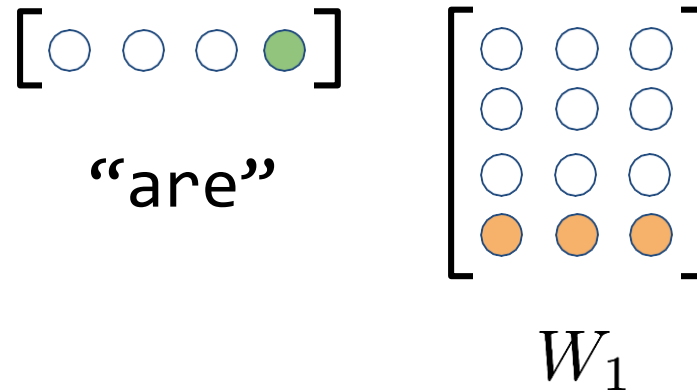


$W_3$



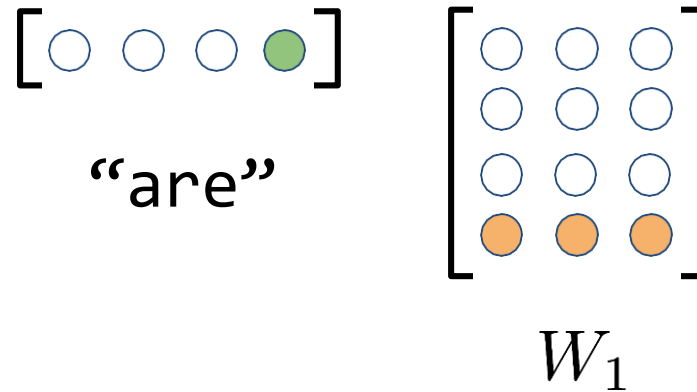
output scores  
max: "you"

# Neural Network Language Model



Each one-hot vector turns on one row in the weight matrix and results in  $[1 \times 3]$  vector

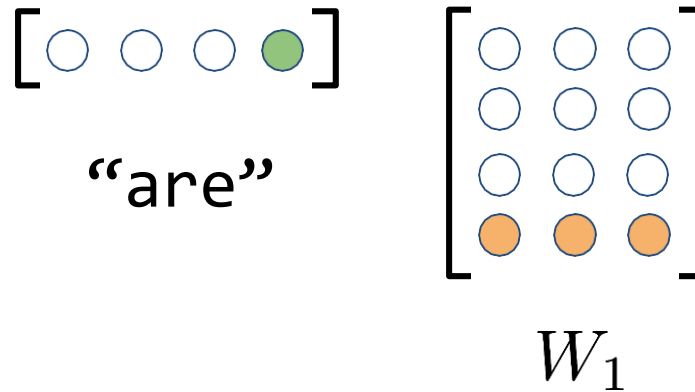
# Neural Network Language Model



Each one-hot vector turns on one row in the weight matrix and results in  $[1 \times 3]$  vector

*Can we say that the  $[1 \times 3]$  vector represents the input word?*

# Neural Network Language Model



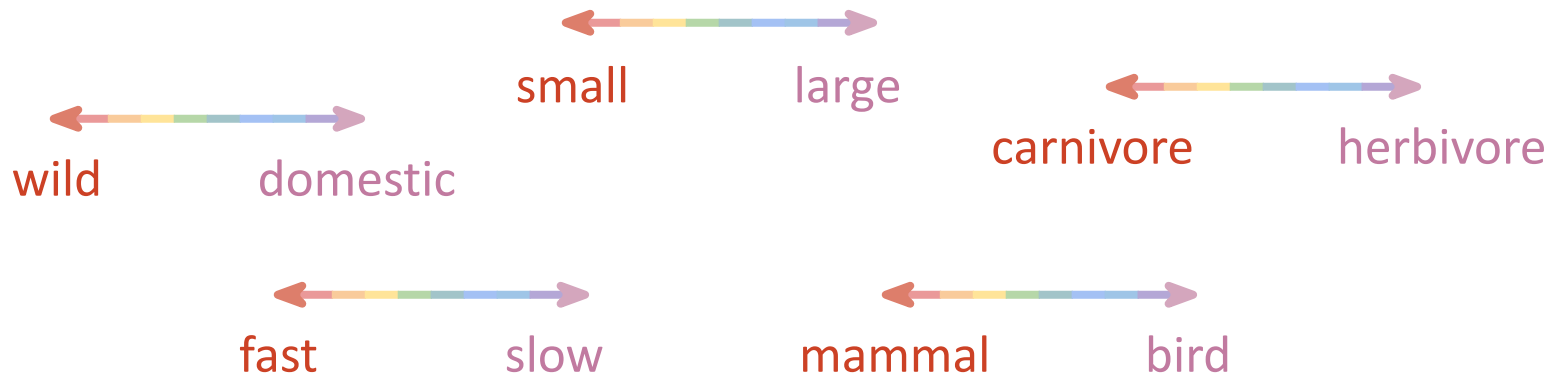
Each one-hot vector turns on one row in the weight matrix and results in  $[1 \times 3]$  vector

*Can we say that the  $[1 \times 3]$  vector represents the input word? **Yes***

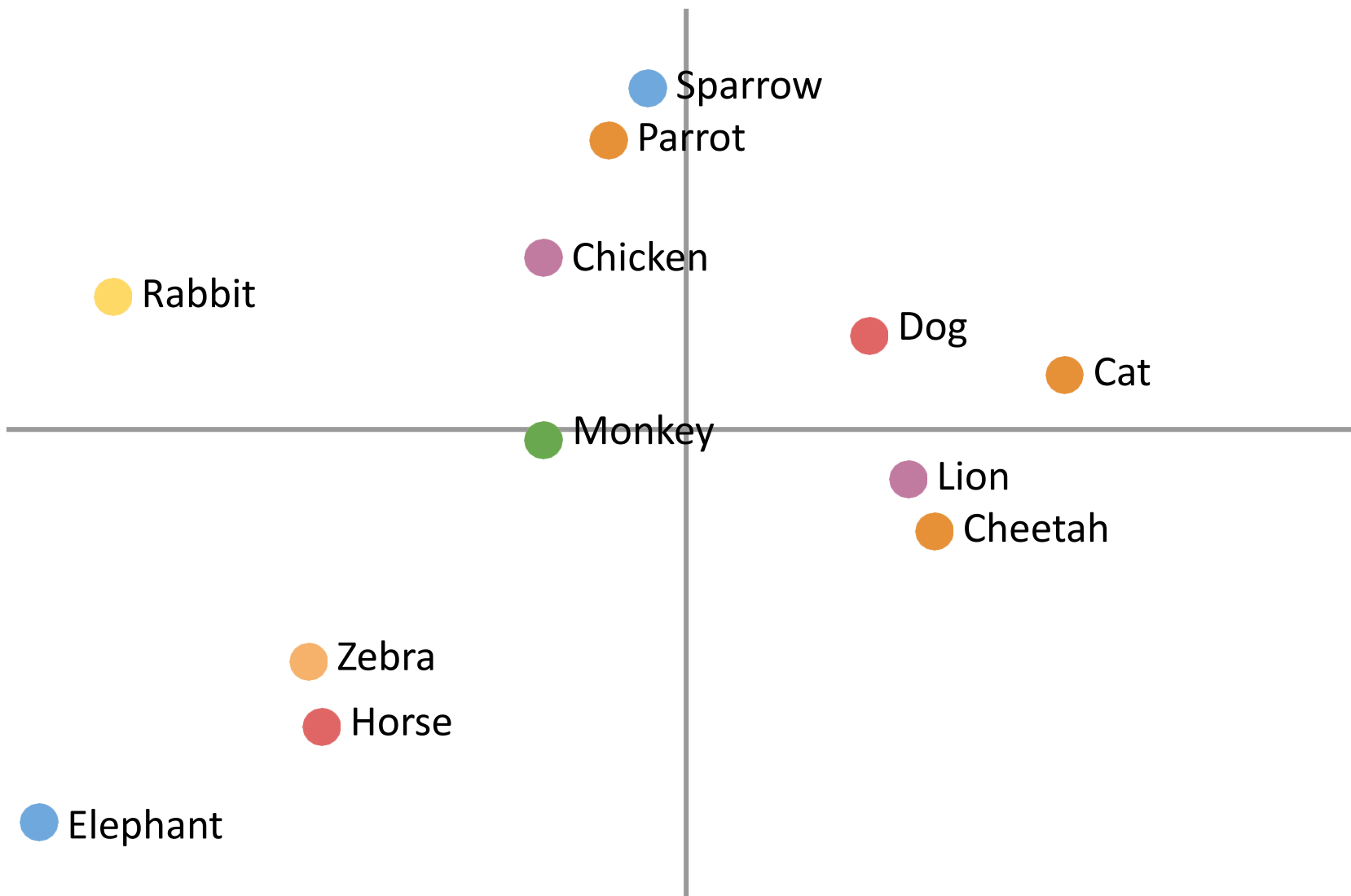
# Exercise

Create a 2D vector space representation of the following words:

dog, lion, cat, rabbit, horse, zebra, cheetah, parrot, sparrow, elephant, chicken, monkey

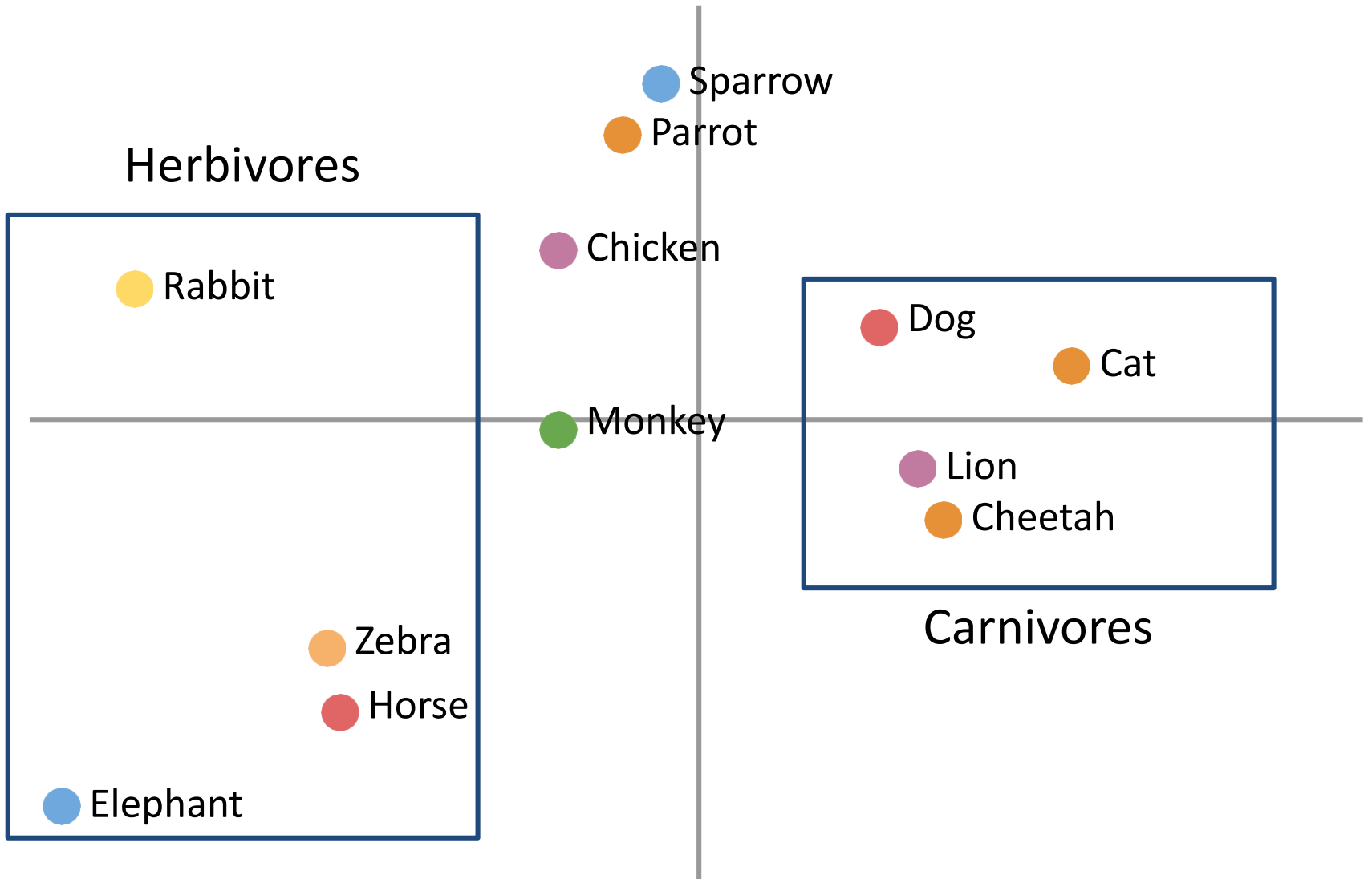


# Exercise



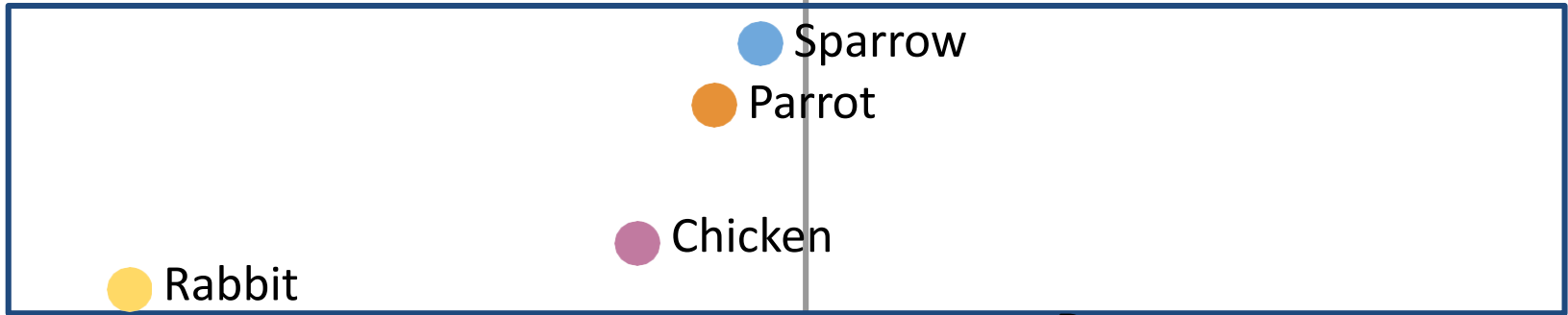


# Exercise



# Exercise

Small Animals



Monkey

Lion

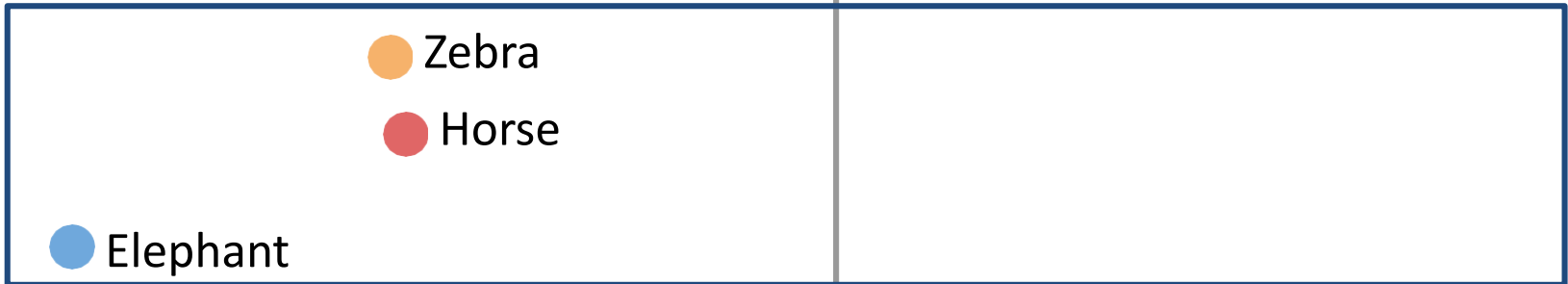
Cheetah

Large Animals

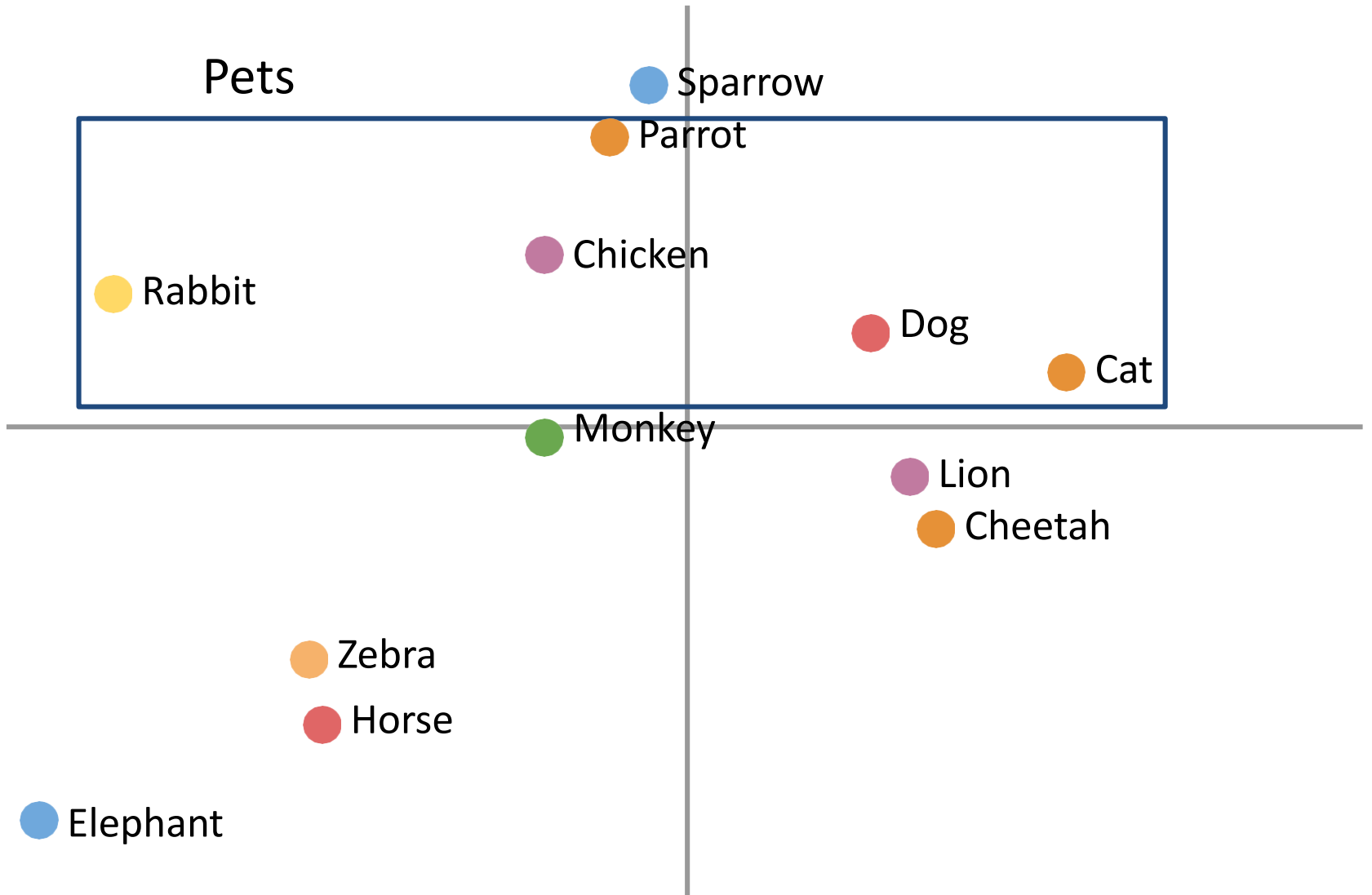
Zebra

Horse

Elephant



# Exercise



# Word Embeddings

How did you decide which animals need to be closer?

How did you handle conflicts between animals that belong to multiple groups?

How does having this kind of vector space representation help us?

# Word Embeddings

- In one-hot vector representation, a word is represented as one large *sparse* vector

only one element is 1 in the entire vector

vectors of different words do not give us any information about the potential relations between the words!

# Word Embeddings

- In one-hot vector representation, a word is represented as one large *sparse* vector
- Instead, **word embeddings** are *dense* vectors in some vector space

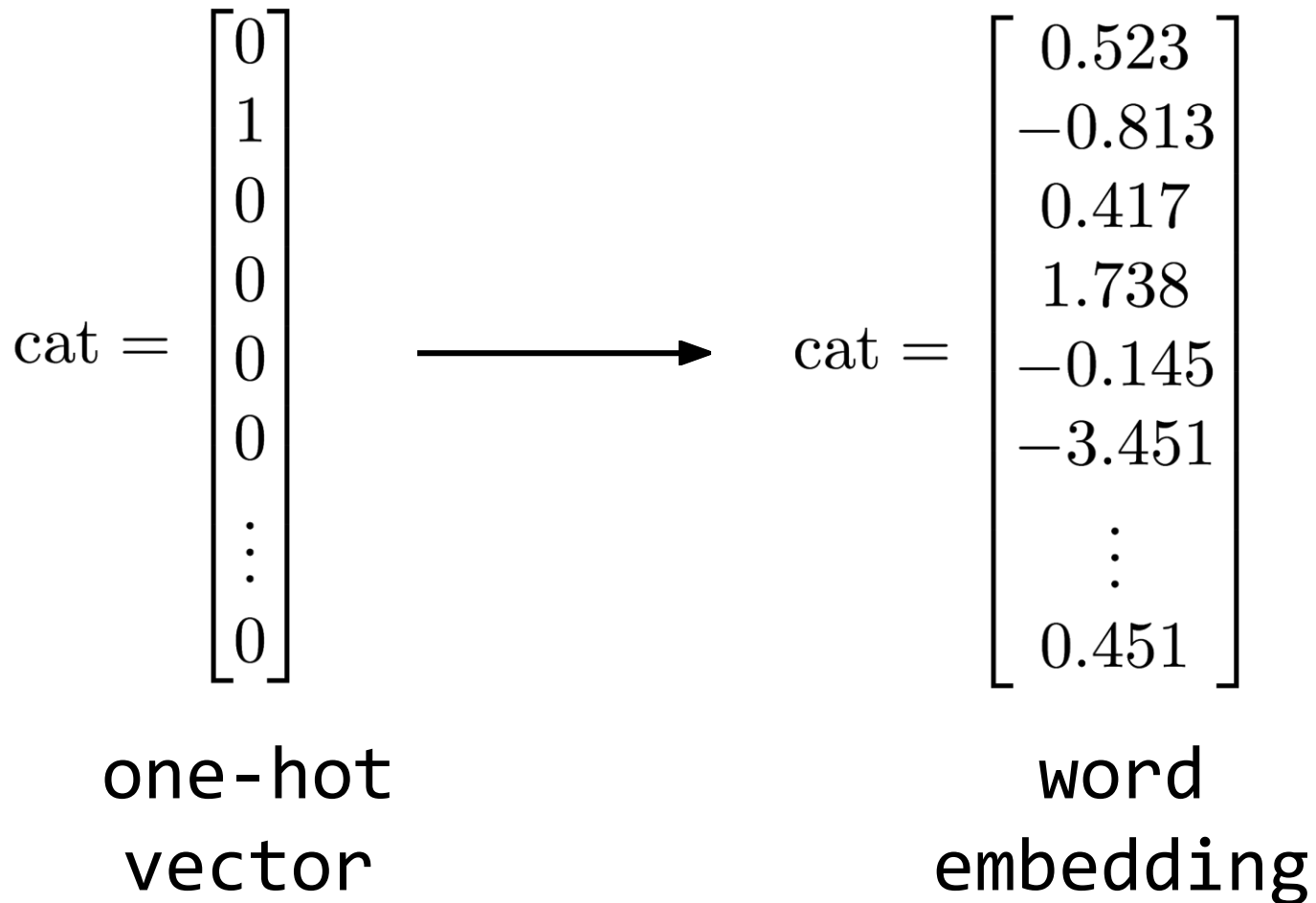
# Word Embeddings

- In one-hot vector representation, a word is represented as one large *sparse* vector
- Instead, **word embeddings** are *dense* vectors in some vector space

word vectors are *continuous* representations of words

vectors of different words give us information about the potential relations between the words - words closer together in meaning have vectors closer to each other

# Word Embeddings





# Word Embeddings

*“Representation of words in continuous space”*

## Inherit benefits

- Reduce dimensionality
- Semantic relatedness
- Increase expressiveness
  - one word is represented in the form of several features (numbers)

# Word Embeddings

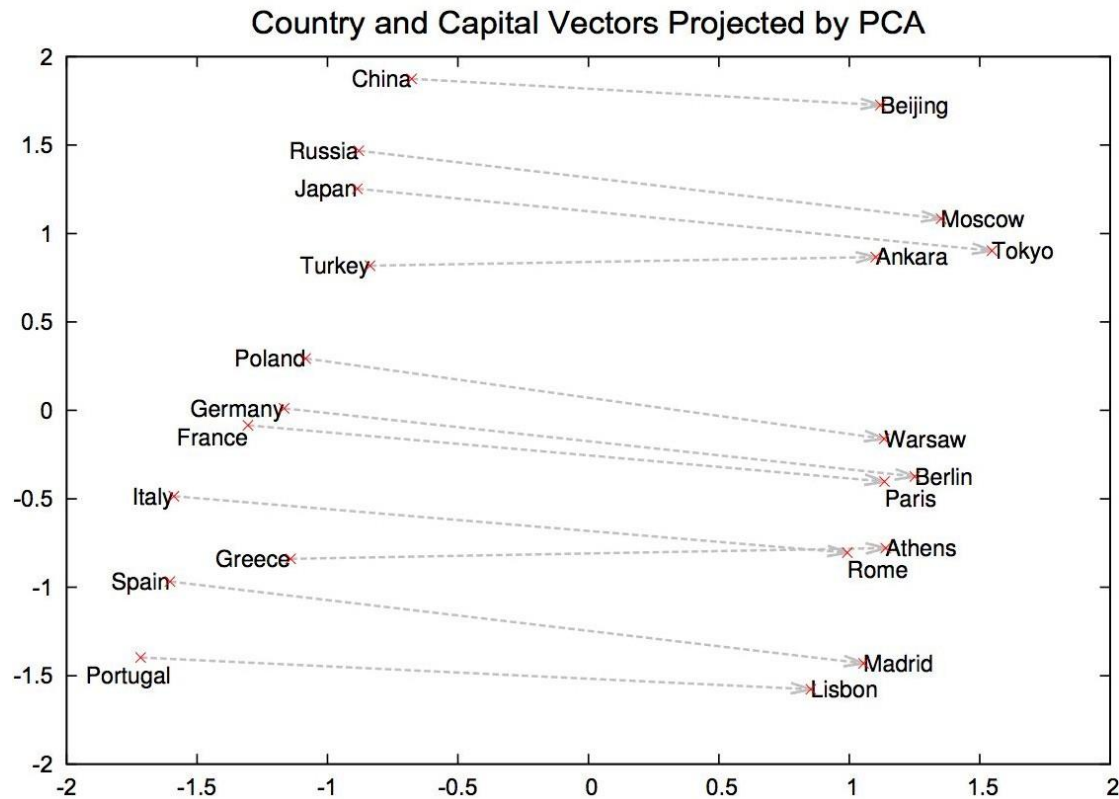
Play with some embeddings!

[https://rare-technologies.com/word2vec-tutorial/#bonus\\_app](https://rare-technologies.com/word2vec-tutorial/#bonus_app)

Try various relationships...

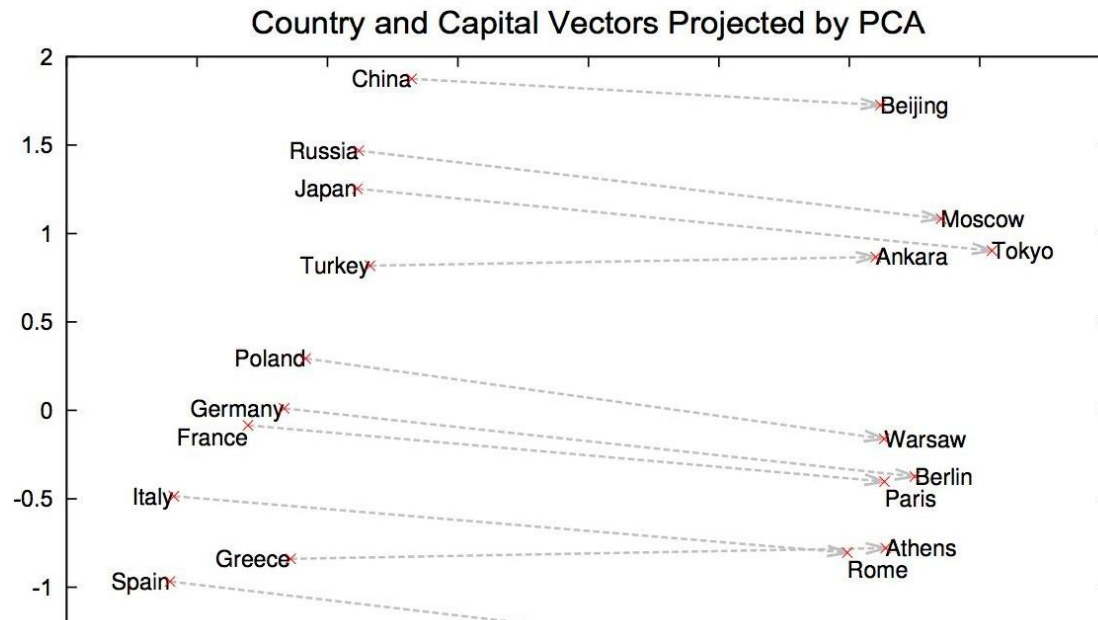
# Word Embeddings

- Plot the embedding vectors



# Word Embeddings

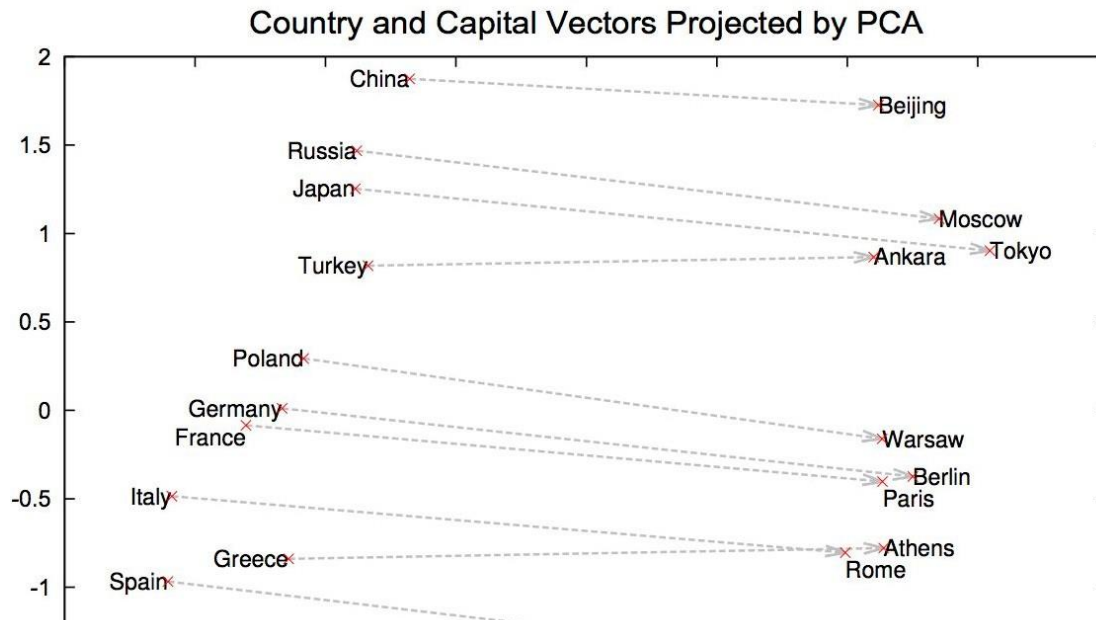
- Plot the embedding vectors



Plot shows the relationship between vectors representing related concepts

# Word Embeddings

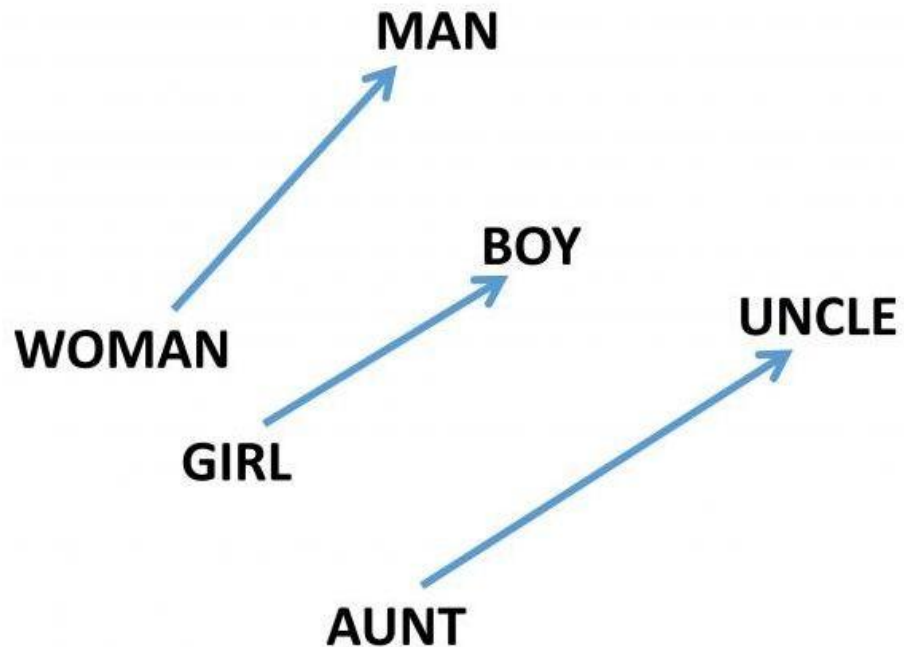
- Plot the embedding vectors



The vectors from countries to capitals point roughly in the same direction

# Word Embeddings

- Similarly, learning the gender relationship



# Word Embeddings

**Q:** How can we learn these embeddings automatically?

# Word Embeddings

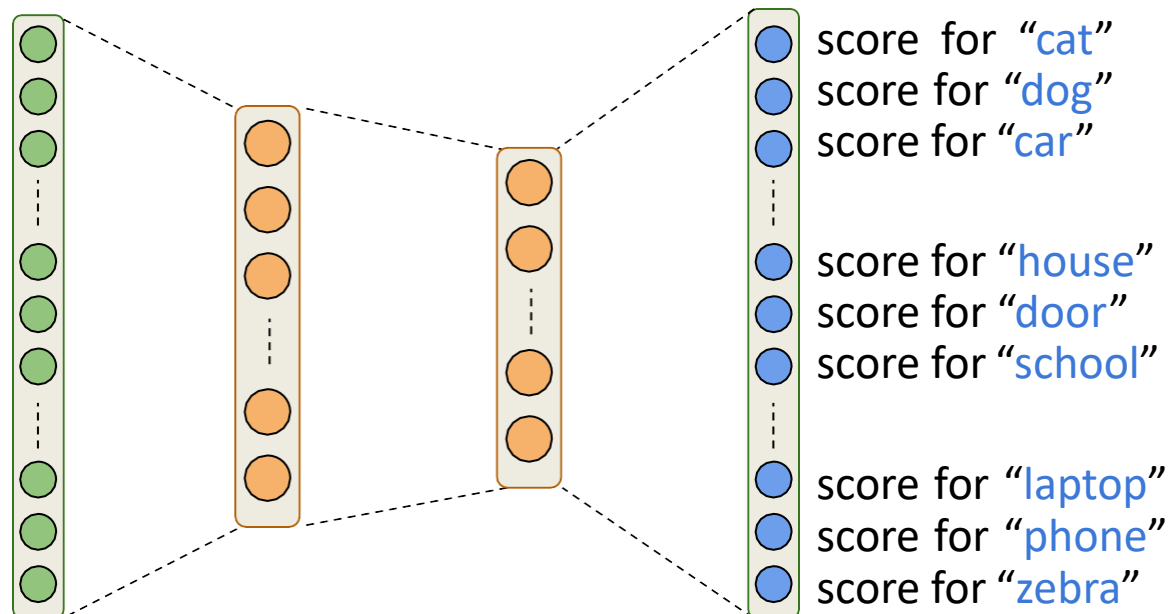
**Q:** How can we learn these embeddings automatically?

**A:** Neural Networks are a step ahead - embeddings are already learned as “richer” features



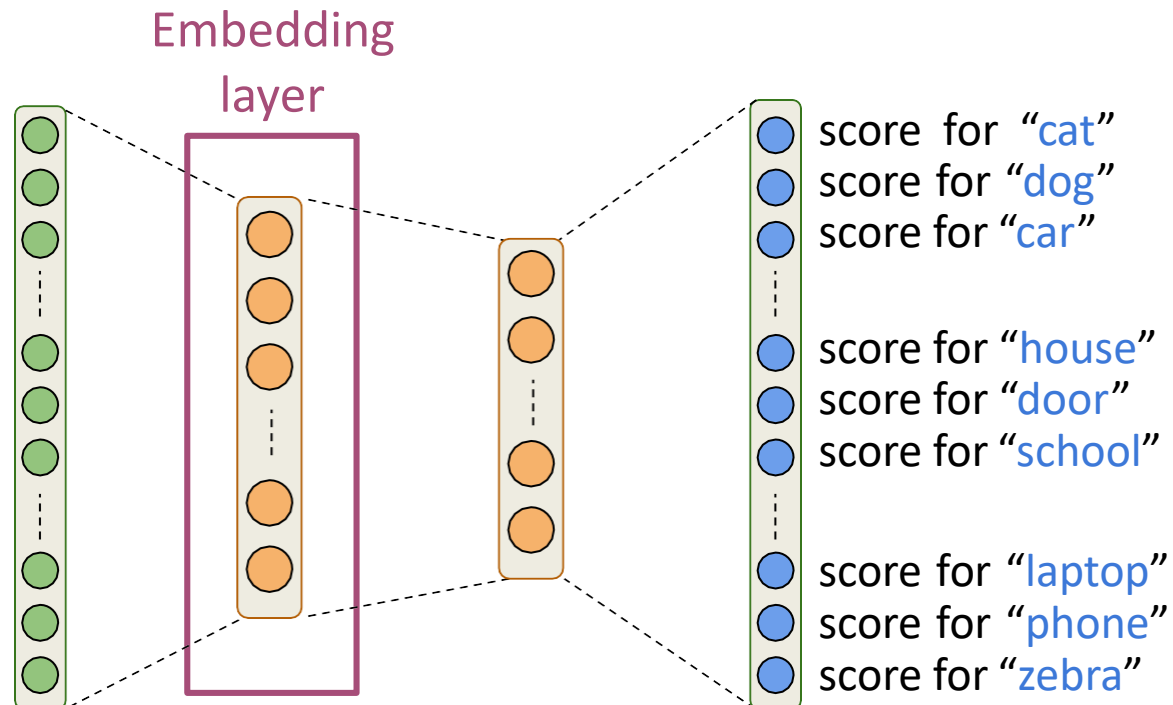
# Word Embeddings

Neural Networks are a step ahead -  
embeddings are already learned as “richer”  
features



# Word Embeddings

Neural Networks are a step ahead - embeddings are already learned as “richer” features



# Word Embeddings

The overall **training task** defines the relationships which will be learned by the model

For example:

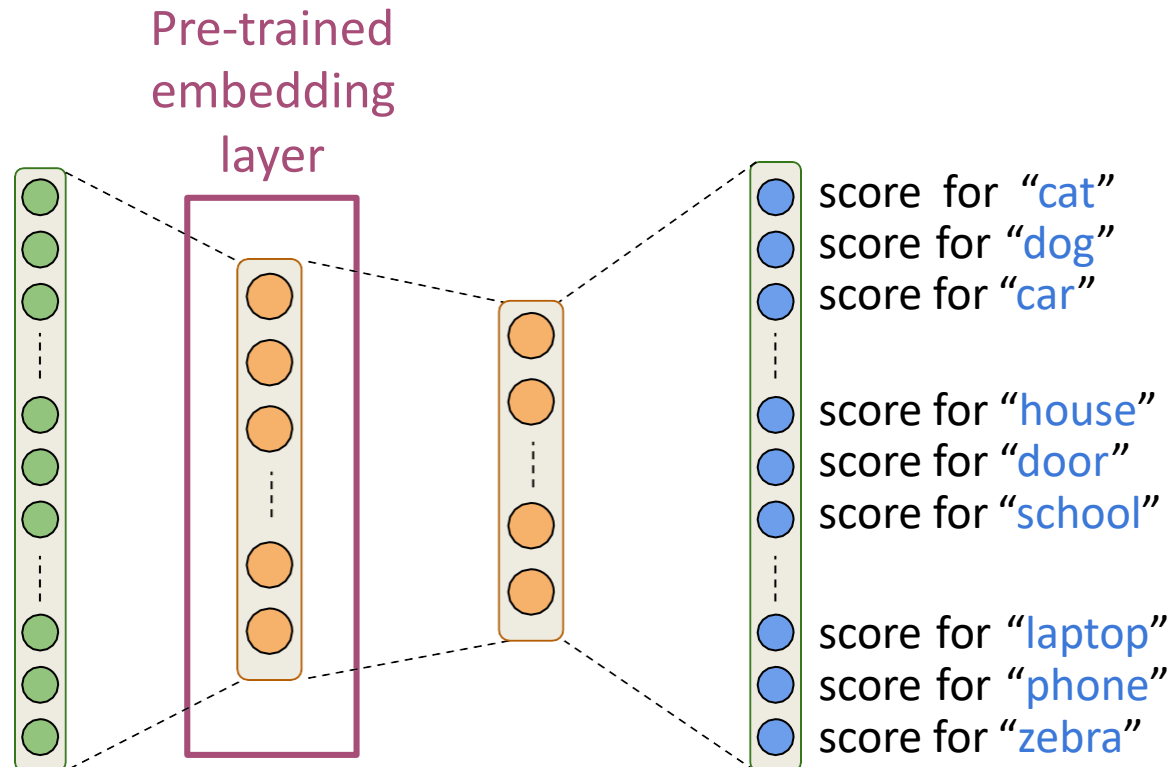
- In language modeling, the model uses neighboring context thus bringing words with similar context closer
- In doing POS tagging task, words with similar POS tags will come close to each other
- If our network is doing machine translation, the embeddings will be tuned for translation

# Word Embeddings

- Generally, task specific embeddings are better than generic embeddings
- In case of small amount of training data, generic embeddings learned on large amount of data works better
- Generic embeddings can also be used as a starting point

# Word Embeddings

We can use pre-trained embeddings as well - just initialize the weights in the first layer with some learned embeddings



# Word Embedding Tools

Some tools to learn word embeddings:

- Word2Vec (from Google)
- FastText (from Facebook)
- GloVe (from Stanford)

# Word Embeddings

## A few pre-trained word embeddings

- GloVe: Wikipedia plus Gigaword <https://goo.gl/1XYZhc>
- FastText: Wikipedia of 294 languages <https://goo.gl/1v423g>
- Dependency-based <https://goo.gl/tpgw4R>

## Using pre-trained embeddings in keras:

<https://blog.keras.io/using-pre-trained-word-embeddings-in-a-keras-model.html>