

# Machine Learning to Language Model

## Topic 02 - Word Embedding

**Jaihua Yen**

<https://jaihuayen.github.io/>

# Contents

- Deep Neural Network
- Word Embedding
- Transfer Learning
- Wrap Up

# Review Bi-Gram Model

# Bi-Gram Model

## A Very Simple Language Model

- **Intuition:** We only predict the next word based on the previous word.
- Model the predicted probability of a certain word based on a given word.

$$P(c_i | c_{i-1})$$

$c_i$  is the word in the  $i$  position.

- Here we give an example of a sentence:

<s> AI could finally be introduced into practice in general tasks <e>

Start Token

End Token

# Bi-Gram Model

## A Very Simple Language Model

<s> AI could finally be introduced into practice in general tasks <e>

$c_{i-1}$

<s>

$c_i$

AI

# Bi-Gram Model

## A Very Simple Language Model

<s> Al could finally be introduced into practice in general tasks <e>

$c_{i-1}$

Al

$c_i$

could

# Bi-Gram Model

## A Very Simple Language Model

<s> AI could finally be introduced into practice in general tasks <e>

$c_{i-1}$

tasks

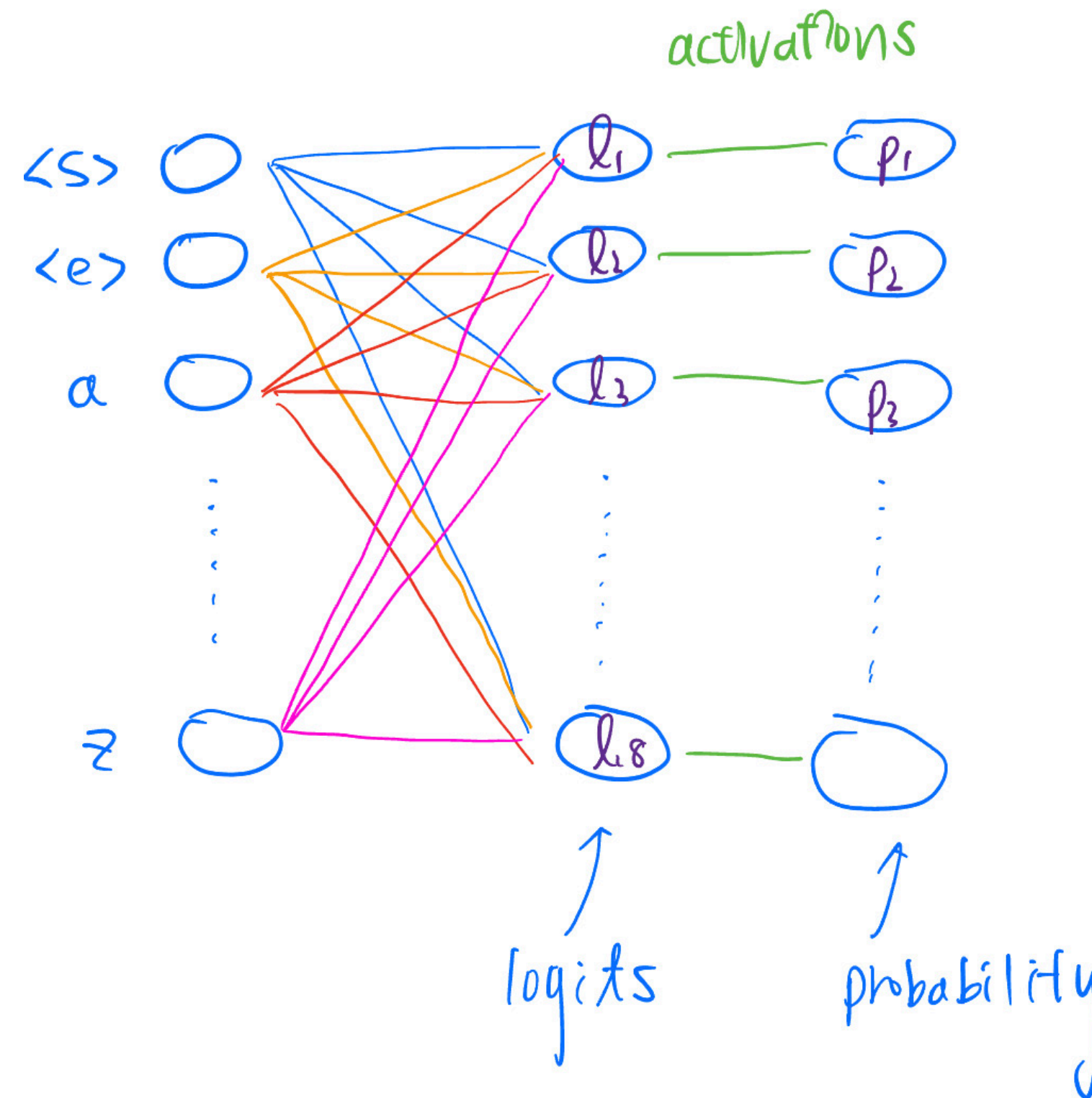
$c_i$

<e>

# Bi-Gram Model

## A Very Simple Language Model

- What if the features cannot be extracted only by one layer of neural network?



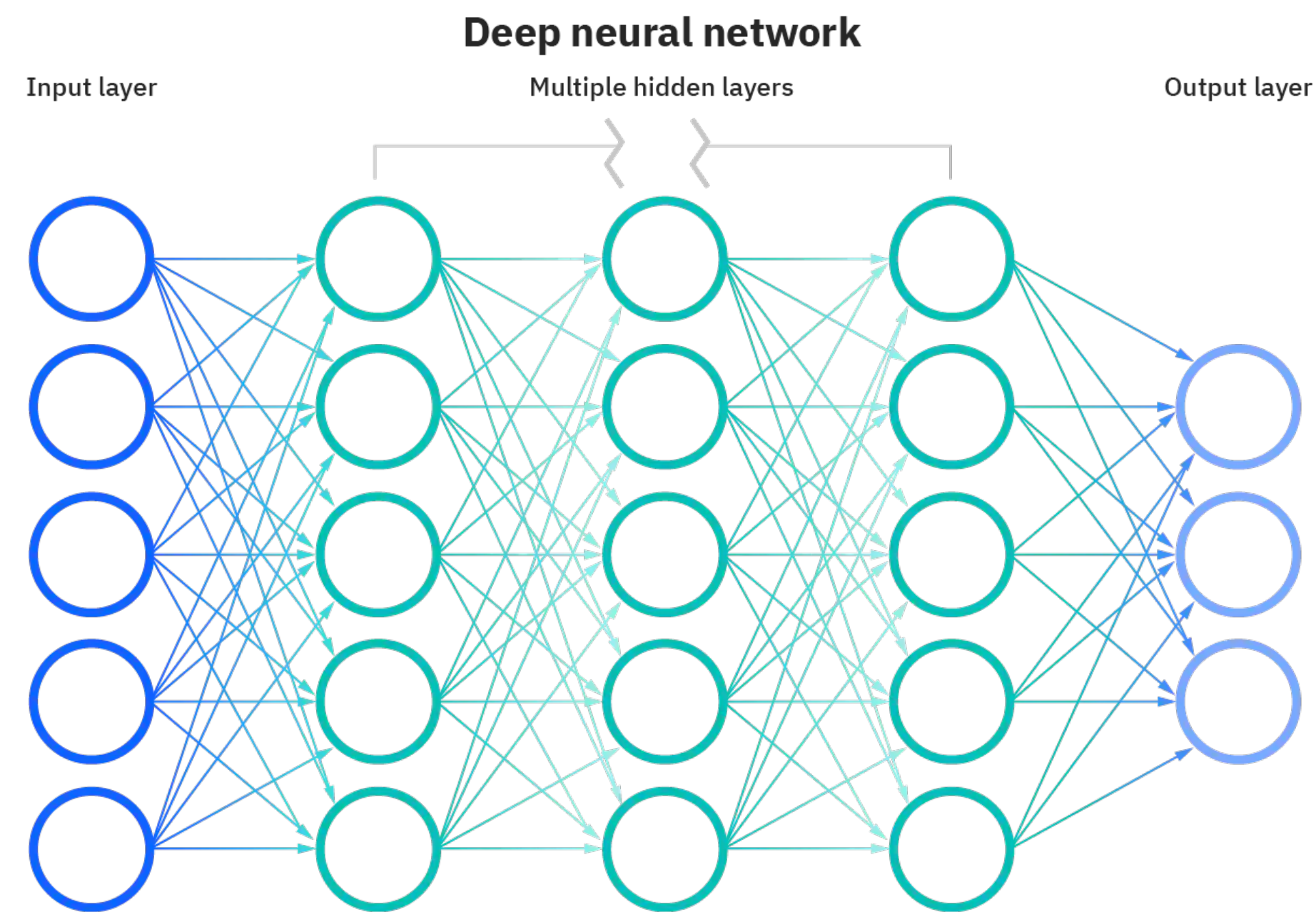
$$p_i = \frac{e^{l_i}}{\sum_{i=1}^{28} e^{l_i}}$$



# Deep Neural Network

# What Deeper?

- Deep Neural Networks (DNN) extract text semantics meanings on a deeper level.



# Word Embedding

# Word Embedding

## Word Representation

- From all the experiments above, we all use one-hot encodings.

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

# Word Embedding

## Word Representation

- However, we cannot extract the meaning between those two words.

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

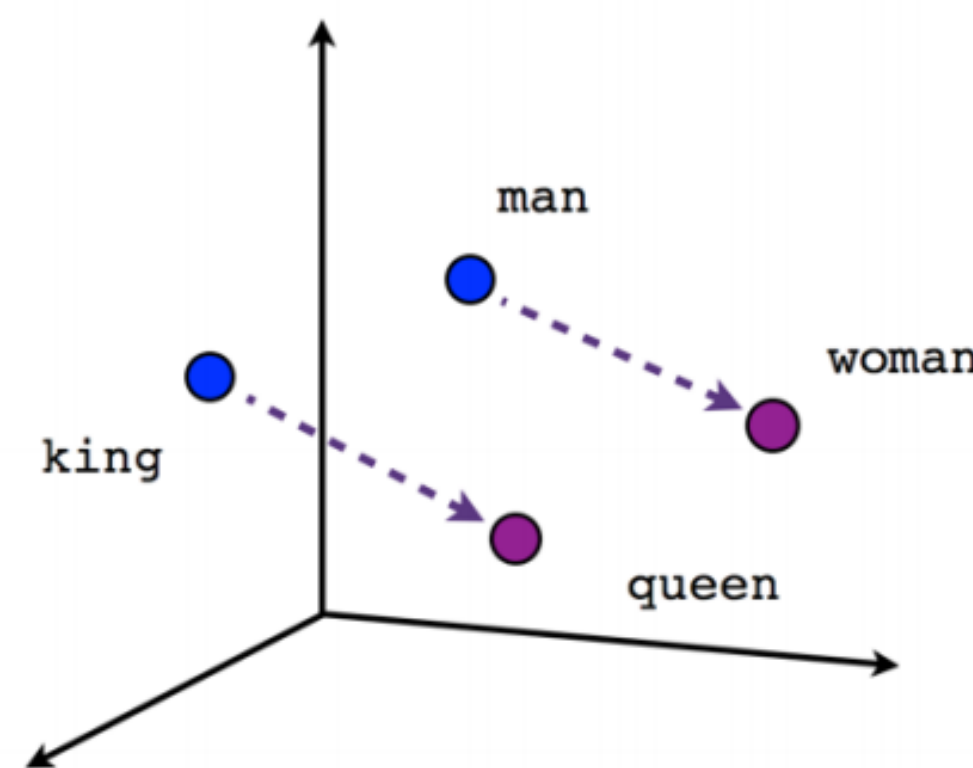
$$\cos(\theta) = \frac{v_{dog}^T v_{cat}}{\|v_{dog}\| \|v_{cat}\|} = 0$$

$$\cos(\theta) = \frac{v_{dog}^T v_{table}}{\|v_{dog}\| \|v_{table}\|} = 0$$

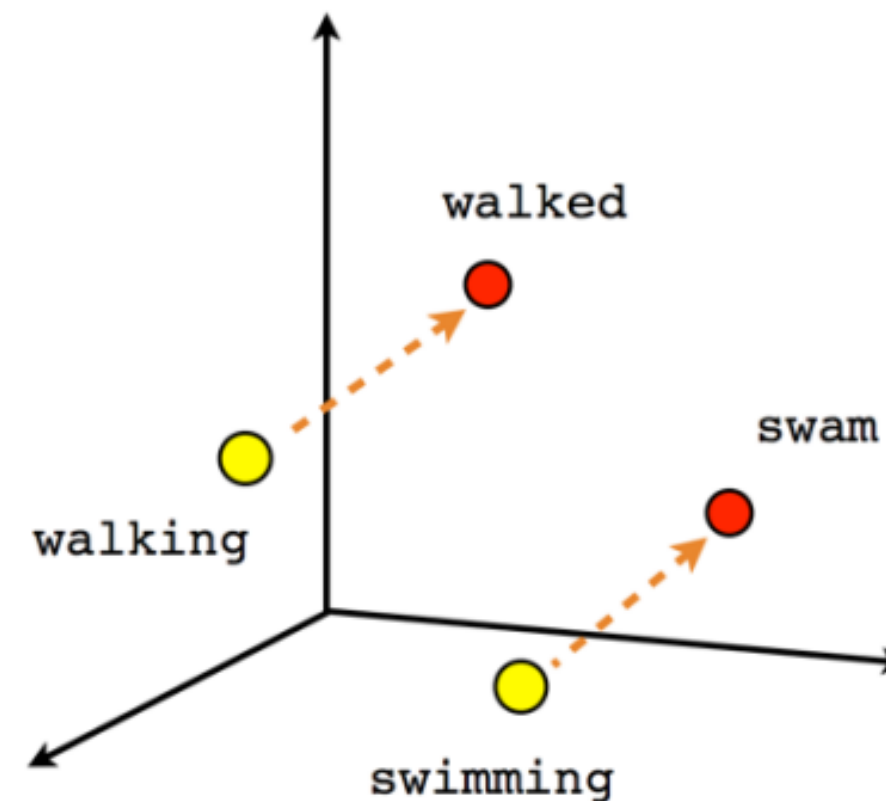
# Word Embedding

## Word Representation

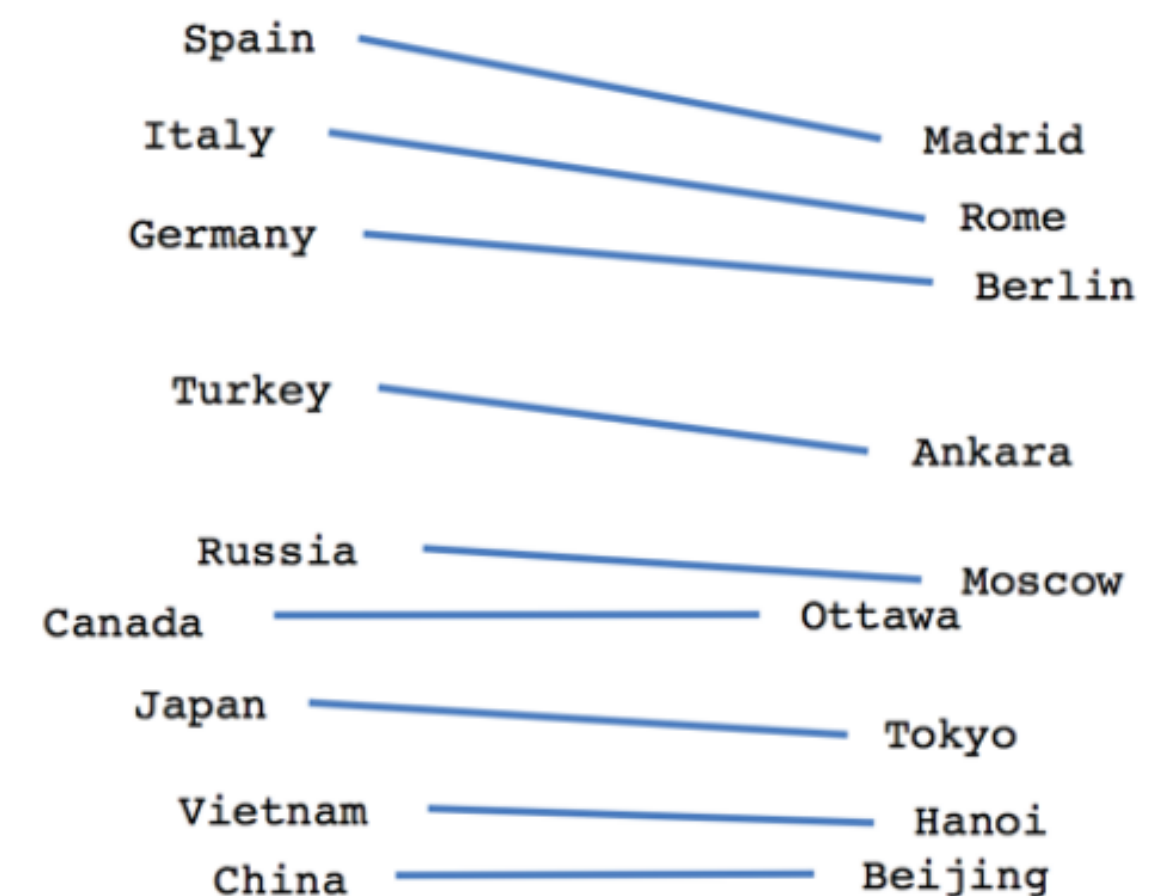
- Word embedding uses a vector representation which could indicate the semantic relationship between words.



Male-Female



Verb tense



Country-Capital

<https://leemeng.tw/find-word-semantic-by-using-word2vec-in-tensorflow.html>

# Word Embedding

## Advantages using word embeddings

- Find the semantic relationship between words.
- Map a high-dimensional one-hot encoding vector to a lower-dimensional word embedding vector

One-hot encoding

	cat	mat	on	sat	the
<b>the</b> =>	0	0	0	0	1
<b>cat</b> =>	1	0	0	0	0
<b>sat</b> =>	0	0	0	1	0
...					

A 4-dimensional embedding

<b>cat</b> =>	1.2	-0.1	4.3	3.2
<b>mat</b> =>	0.4	2.5	-0.9	0.5
<b>on</b> =>	2.1	0.3	0.1	0.4
...				

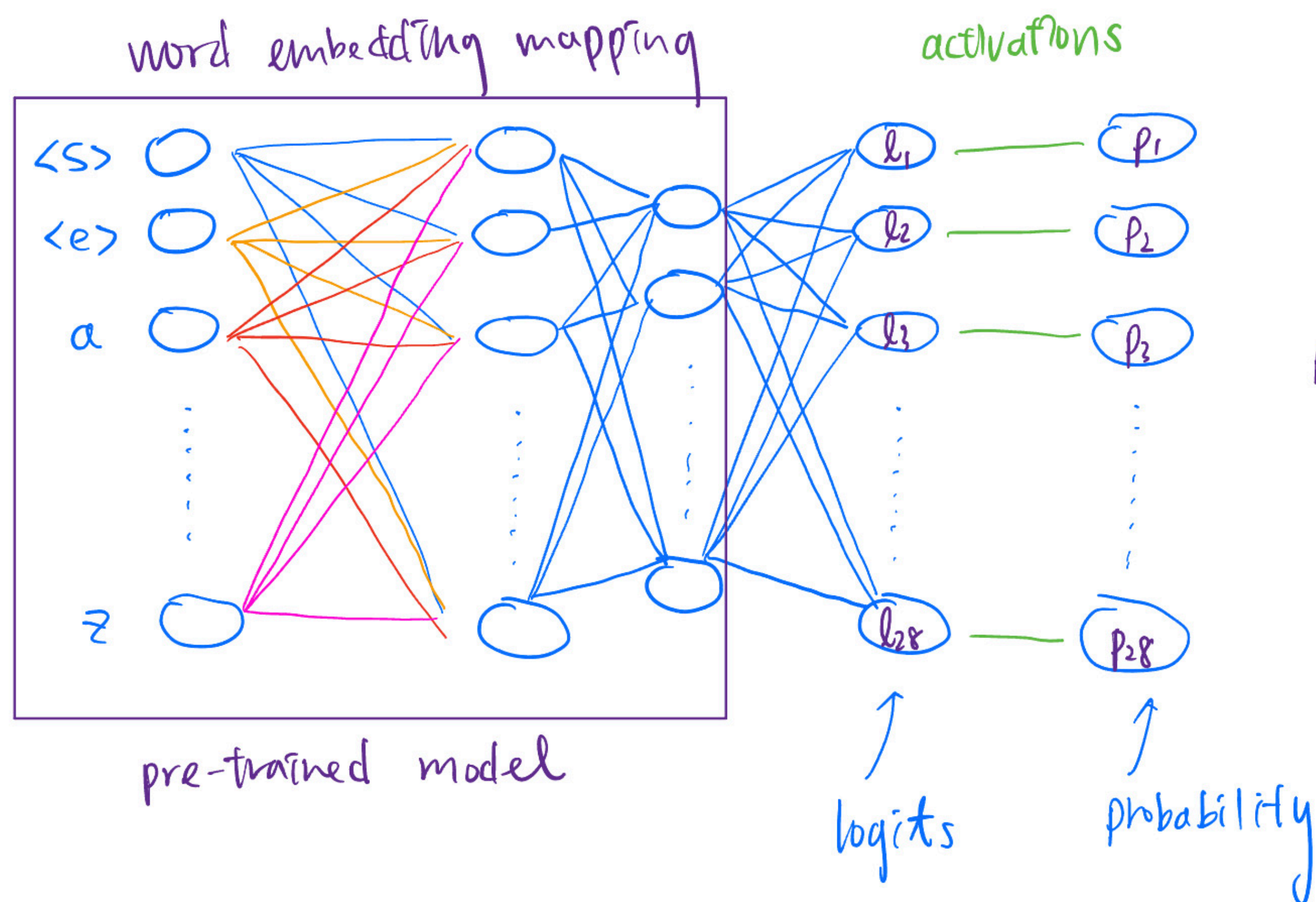
[https://www.tensorflow.org/text/guide/word\\_embeddings](https://www.tensorflow.org/text/guide/word_embeddings)

# How to Train Word Embedding?



# Deep Neural Network

Word embedding layer is in the hidden layer of DNN



Let's do this in Colab!

# Questions

- What if we have a massive dataset that cannot fit in memory?
- How can we compute gradient with more than one hidden layer?

# Further Reading

## A Neural Probabilistic Language Model

---

Journal of Machine Learning Research 3 (2003) 1137–1155

Submitted 4/02; Published 2/03

### A Neural Probabilistic Language Model

**Yoshua Bengio**  
**Réjean Ducharme**  
**Pascal Vincent**  
**Christian Jauvin**

*Département d'Informatique et Recherche Opérationnelle*  
*Centre de Recherche Mathématiques*  
*Université de Montréal, Montréal, Québec, Canada*

BENGIOY@IRO.UMONTREAL.CA  
DUCHARME@IRO.UMONTREAL.CA  
VINCENTP@IRO.UMONTREAL.CA  
JAUVINC@IRO.UMONTREAL.CA

**Editors:** Jaz Kandola, Thomas Hofmann, Tomaso Poggio and John Shawe-Taylor

### Abstract

A goal of statistical language modeling is to learn the joint probability function of sequences of words in a language. This is intrinsically difficult because of the **curse of dimensionality**: a word sequence on which the model will be tested is likely to be different from all the word sequences seen during training. Traditional but very successful approaches based on n-grams obtain generalization

# Further Reading

## Efficient Estimation of Word Representations in Vector Space (Word2Vec)

---

### Efficient Estimation of Word Representations in Vector Space

---

**Tomas Mikolov**

Google Inc., Mountain View, CA  
tmikolov@google.com

**Kai Chen**

Google Inc., Mountain View, CA  
kaichen@google.com

**Greg Corrado**

Google Inc., Mountain View, CA  
gcorrado@google.com

**Jeffrey Dean**

Google Inc., Mountain View, CA  
jeff@google.com

### Abstract

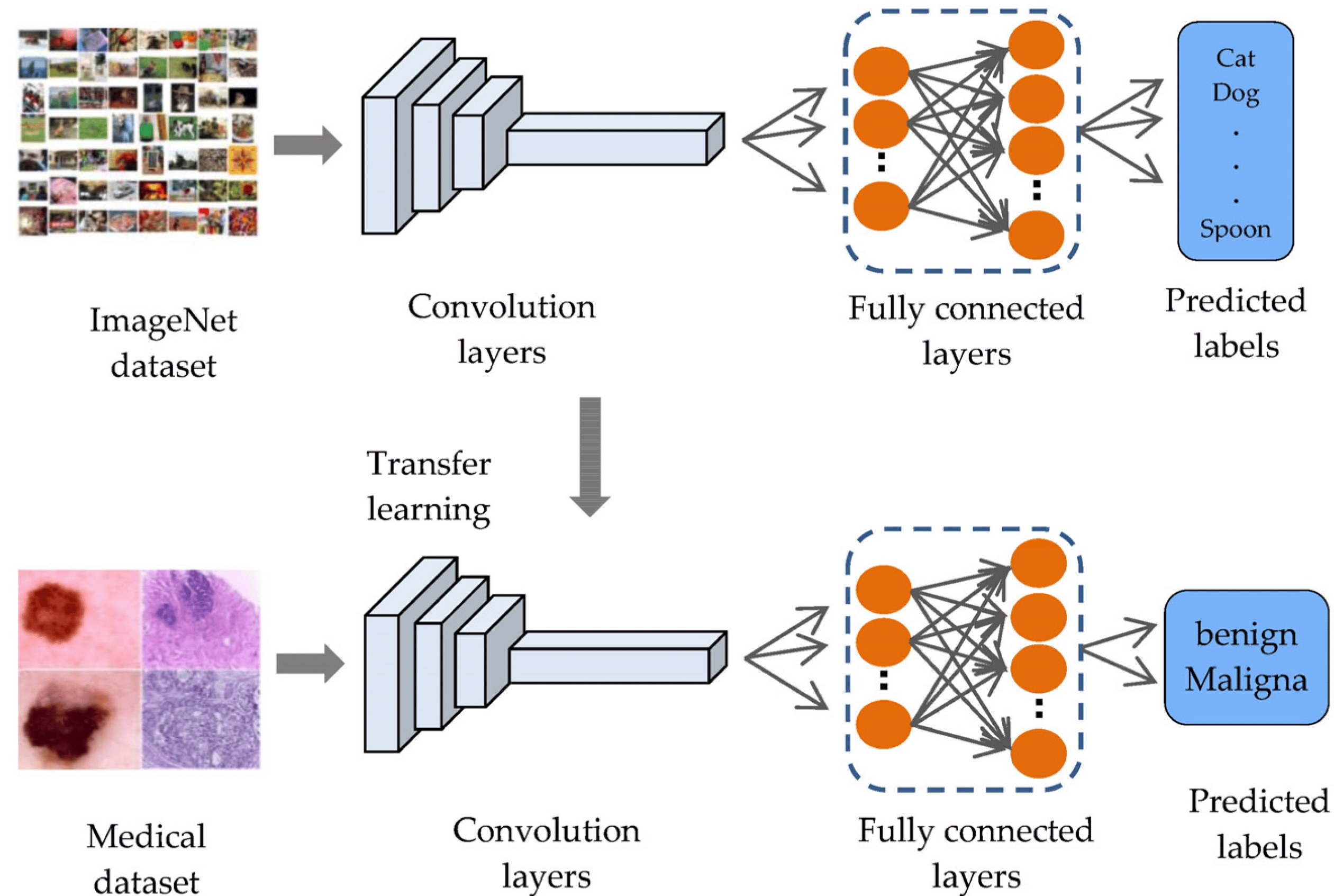
We propose two novel model architectures for computing continuous vector representations of words from very large data sets. The quality of these representations

# Transfer Learning



# Transfer Learning

Use Pre-trained Model in other tasks

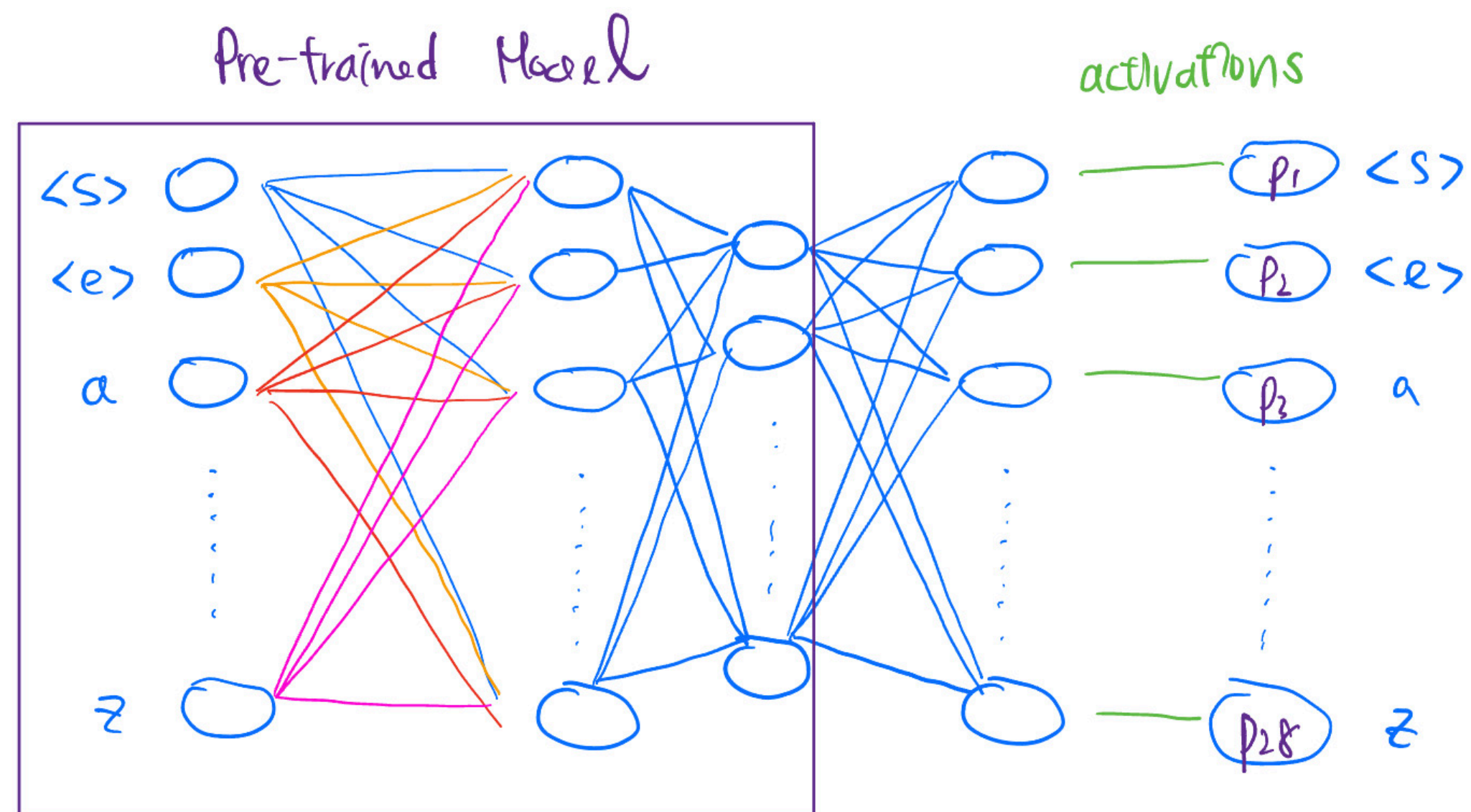


<https://www.mdpi.com/1424-8220/23/2/570>

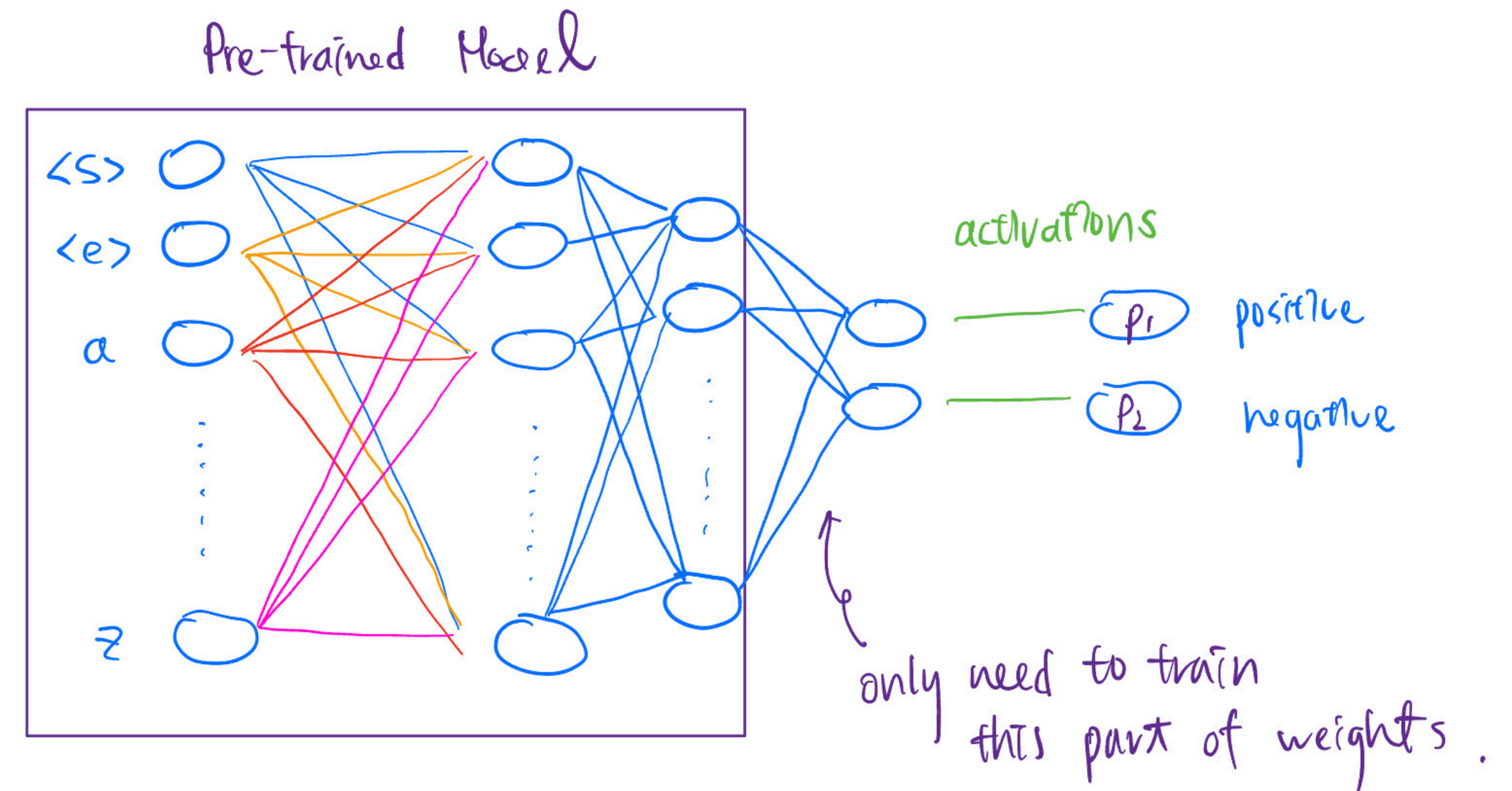
# Transfer Learning

## Use Pre-trained Model in other tasks

Original Task



New Classification Task





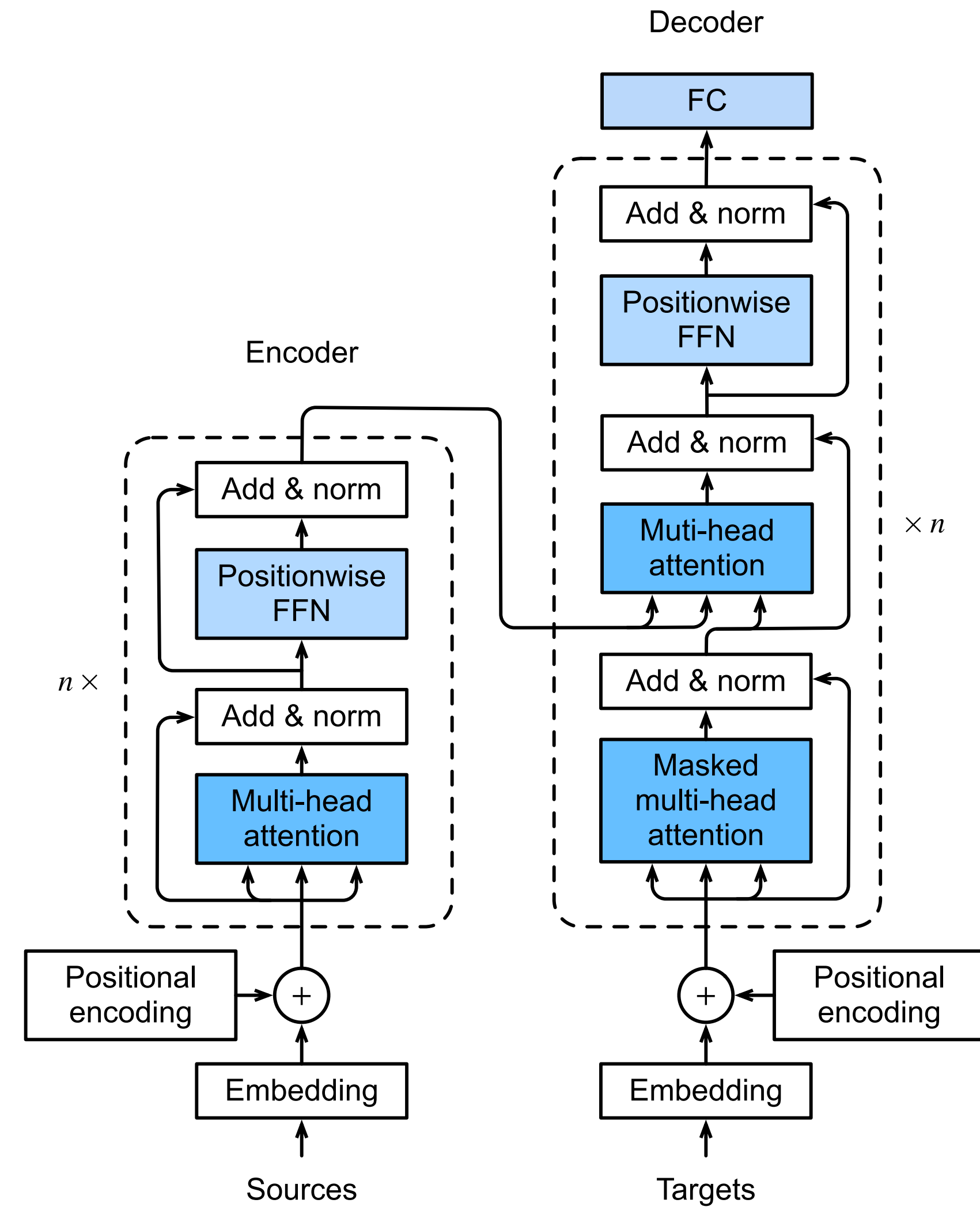
# Wrap Up

# What We Have Gone Through

- Deep Neural Network
- Word Embedding
- Transfer Learning

# What's Next

- Transformer - Self Attention



**Q & A**