

# Natural Language Processing

**CSE 447 / 547 M**

## Pre-training

Lecturer: Kabir Ahuja

*Slides adapted from Liwei Jiang, John Hewitt, Anna Goldie*

# Major Paradigms in NLP

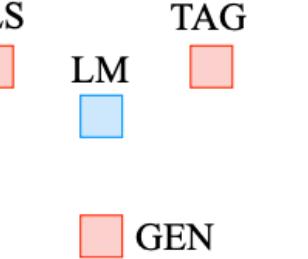
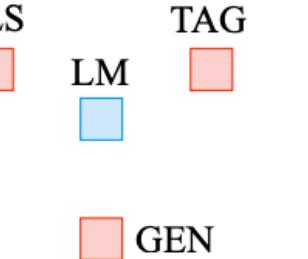
Liu et al. 2021

# Major Paradigms in NLP

Paradigm	Engineering	Task Relation
a. Fully Supervised Learning (Non-Neural Network)	Features (e.g. word identity, part-of-speech, sentence length)	 CLS (Red square) LM (Blue square) TAG (Red square) GEN (Red square)

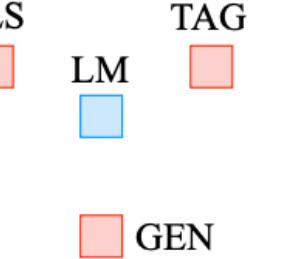
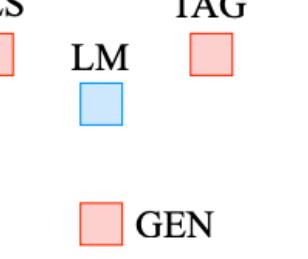
Liu et al. 2021

# Major Paradigms in NLP

Paradigm	Engineering	Task Relation
a. Fully Supervised Learning (Non-Neural Network)	Features (e.g. word identity, part-of-speech, sentence length)	
b. Fully Supervised Learning (Neural Network)	Architecture (e.g. convolutional, recurrent, self-attentional)	

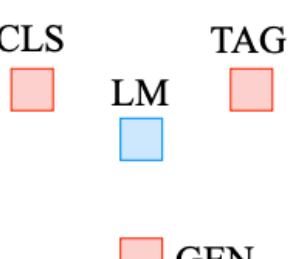
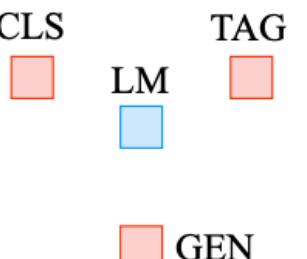
Liu et al. 2021

# Major Paradigms in NLP

Paradigm	Engineering	Task Relation
a. Fully Supervised Learning (Non-Neural Network)	Features (e.g. word identity, part-of-speech, sentence length)	
b. Fully Supervised Learning (Neural Network)	Architecture (e.g. convolutional, recurrent, self-attentional)	

Liu et al. 2021

# Major Paradigms in NLP

Paradigm	Engineering	Task Relation
a. Fully Supervised Learning (Non-Neural Network)	Features (e.g. word identity, part-of-speech, sentence length)	 <p>Task Relation diagram showing three nodes: CLS (red square), LM (blue square), and TAG (red square). A red arrow points from CLS to TAG. A blue arrow points from LM to TAG.</p>
b. Fully Supervised Learning (Neural Network)	Architecture (e.g. convolutional, recurrent, self-attentional)	 <p>Task Relation diagram showing three nodes: CLS (red square), LM (blue square), and TAG (red square). A red arrow points from CLS to TAG. A blue arrow points from LM to TAG.</p>

Pre 2017

Liu et al. 2021

# Major Paradigms in NLP

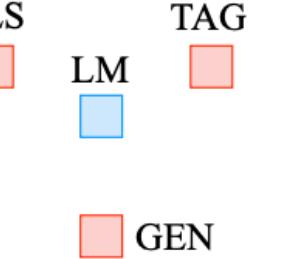
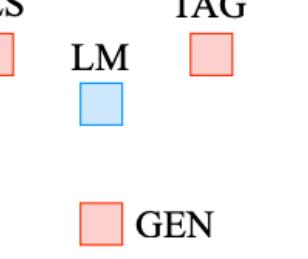
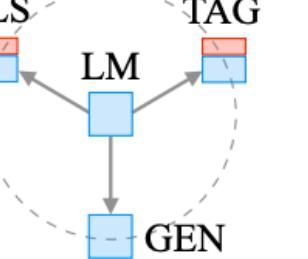
Paradigm	Engineering	Task Relation
a. Fully Supervised Learning (Non-Neural Network)	Features (e.g. word identity, part-of-speech, sentence length)	<p>Diagram illustrating task relations:</p> <ul style="list-style-type: none"><li>CLS (red square)</li><li>LM (blue square)</li><li>TAG (red square)</li><li>GEN (red square)</li></ul>
b. Fully Supervised Learning (Neural Network)	Architecture (e.g. convolutional, recurrent, self-attentional)	<p>Diagram illustrating task relations:</p> <ul style="list-style-type: none"><li>CLS (red square)</li><li>LM (blue square)</li><li>TAG (red square)</li><li>GEN (red square)</li></ul>

Pre 2017

What we have seen so far

Liu et al. 2021

# Major Paradigms in NLP

Paradigm	Engineering	Task Relation
a. Fully Supervised Learning (Non-Neural Network)	Features (e.g. word identity, part-of-speech, sentence length)	
b. Fully Supervised Learning (Neural Network)	Architecture (e.g. convolutional, recurrent, self-attentional)	
c. Pre-train, Fine-tune	Objective (e.g. masked language modeling, next sentence prediction)	

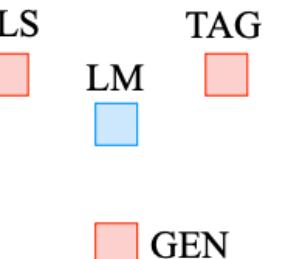
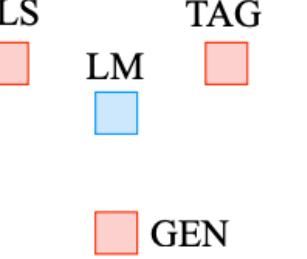
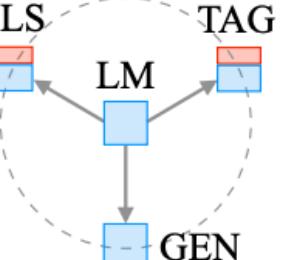
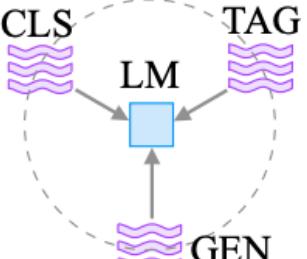
What we have seen so far

Pre 2017

2017-2019

Liu et al. 2021

# Major Paradigms in NLP

Paradigm	Engineering	Task Relation
a. Fully Supervised Learning (Non-Neural Network)	Features (e.g. word identity, part-of-speech, sentence length)	
b. Fully Supervised Learning (Neural Network)	Architecture (e.g. convolutional, recurrent, self-attentional)	
c. Pre-train, Fine-tune	Objective (e.g. masked language modeling, next sentence prediction)	
d. Pre-train, Prompt, Predict	Prompt (e.g. cloze, prefix)	

What we have seen so far

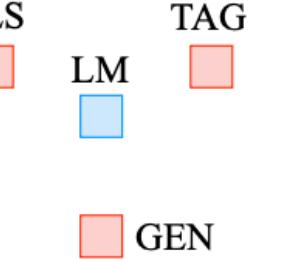
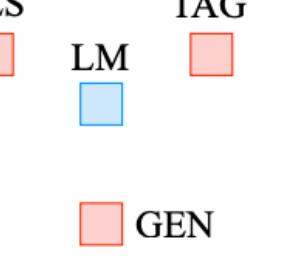
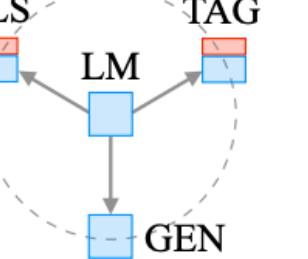
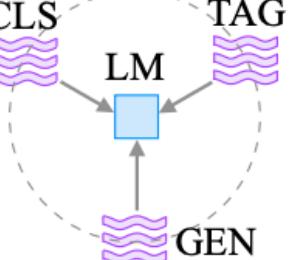
Pre 2017

2017-2019

2021- Present?

Liu et al. 2021

# Major Paradigms in NLP

Paradigm	Engineering	Task Relation
a. Fully Supervised Learning (Non-Neural Network)	Features (e.g. word identity, part-of-speech, sentence length)	
b. Fully Supervised Learning (Neural Network)	Architecture (e.g. convolutional, recurrent, self-attentional)	
c. Pre-train, Fine-tune	Objective (e.g. masked language modeling, next sentence prediction)	
d. Pre-train, Prompt, Predict	Prompt (e.g. cloze, prefix)	
e. Pre-train, Alignment, (Fine-tune), Predict		

What we have seen so far

Pre 2017

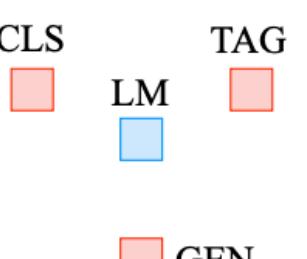
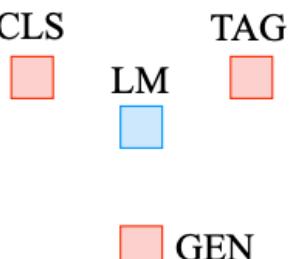
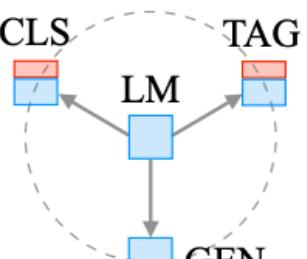
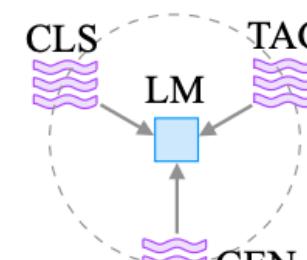
2017-2019

2021- Present?

2022- Present

Liu et al. 2021

# Major Paradigms in NLP

Paradigm	Engineering	Task Relation
a. Fully Supervised Learning (Non-Neural Network)	Features (e.g. word identity, part-of-speech, sentence length)	
b. Fully Supervised Learning (Neural Network)	Architecture (e.g. convolutional, recurrent, self-attentional)	
c. Pre-train, Fine-tune	Objective (e.g. masked language modeling, next sentence prediction)	
d. Pre-train, Prompt, Predict	Prompt (e.g. cloze, prefix)	
e. Pre-train, Alignment, (Fine-tune), Predict		

Liu et al. 2021

Pre 2017

**What we have seen so far**

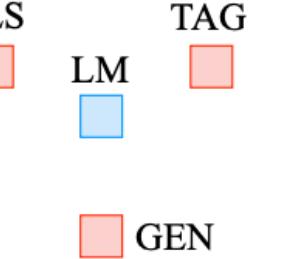
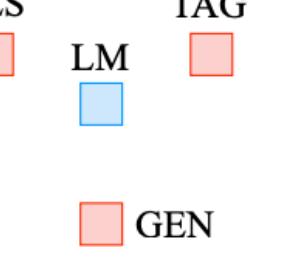
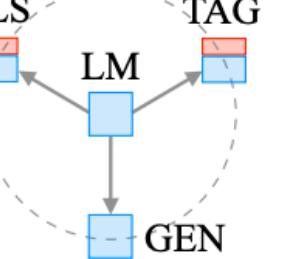
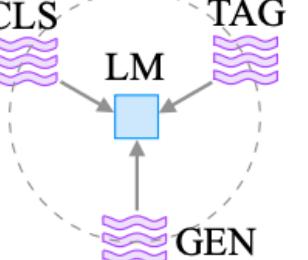
2017-2019

**What we will see in the coming lectures**

2021- Present?

2022- Present

# Major Paradigms in NLP

Paradigm	Engineering	Task Relation
a. Fully Supervised Learning (Non-Neural Network)	Features (e.g. word identity, part-of-speech, sentence length)	 A diagram showing three boxes labeled 'CLS' (red), 'LM' (blue), and 'TAG' (red). 'CLS' has a solid arrow pointing to 'LM'. 'LM' has a solid arrow pointing to 'TAG'. There are no arrows between 'CLS' and 'TAG'.
b. Fully Supervised Learning (Neural Network)	Architecture (e.g. convolutional, recurrent, self-attentional)	 A diagram showing three boxes labeled 'CLS' (red), 'LM' (blue), and 'TAG' (red). 'CLS' has a dashed arrow pointing to 'LM'. 'LM' has a solid arrow pointing to 'TAG'. There are no arrows between 'CLS' and 'TAG'.
c. Pre-train, Fine-tune	Objective (e.g. masked language modeling, next sentence prediction)	 A diagram showing three boxes labeled 'CLS' (red), 'LM' (blue), and 'TAG' (red). 'CLS' has a dashed arrow pointing to 'LM'. 'LM' has a dashed arrow pointing to 'TAG'. There are no arrows between 'CLS' and 'TAG'.
d. Pre-train, Prompt, Predict	Prompt (e.g. cloze, prefix)	 A diagram showing three boxes labeled 'CLS' (red), 'LM' (blue), and 'TAG' (red). 'CLS' has a dashed arrow pointing to 'LM'. 'LM' has a dashed arrow pointing to 'TAG'. There are no arrows between 'CLS' and 'TAG'.
e. Pre-train, Alignment, (Fine-tune), Predict		

Liu et al. 2021

Pre 2017

What we have seen so far

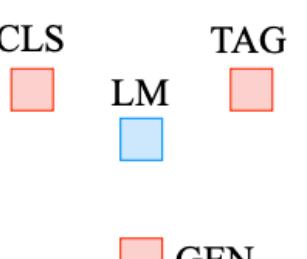
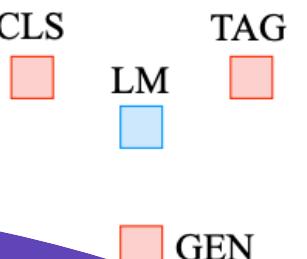
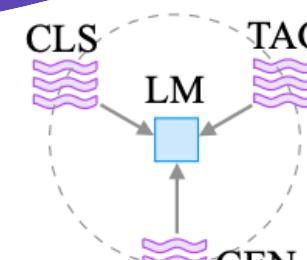
2017-2019

What we will see in the coming lectures

2021- Present?

2022- Present

# Major Paradigms in NLP

Paradigm	Engineering	Task Relation
a. Fully Supervised Learning (Non-Neural Network)	Features (e.g. word identity, part-of-speech, sentence length)	 A diagram showing three boxes labeled 'CLS' (red), 'LM' (blue), and 'TAG' (red) arranged horizontally. Below them is a red box labeled 'GEN'. Arrows point from 'CLS' to 'LM', 'LM' to 'TAG', and 'TAG' to 'GEN'.
b. Fully Supervised Learning (Neural Network)	Architecture (e.g. convolutional, recurrent, self-attentional)	 A diagram showing three boxes labeled 'CLS' (red), 'LM' (blue), and 'TAG' (red) arranged horizontally. Below them is a red box labeled 'GEN'. Arrows point from 'CLS' to 'LM', 'LM' to 'TAG', and 'TAG' to 'GEN'.
c. Pre-train, Fine-tune	Pre-training common across all major paradigms post 2017	
d. Pre-train, Prompt, Predict	Prompt (e.g. cloze, prefix)	 A diagram showing three boxes labeled 'CLS' (red), 'LM' (blue), and 'TAG' (red) arranged horizontally. Below them is a red box labeled 'GEN'. Arrows point from 'CLS' to 'LM', 'LM' to 'TAG', and 'TAG' to 'GEN'. There are also dashed arrows pointing from 'CLS' to 'GEN' and from 'TAG' to 'GEN'.
e. Pre-train, Alignment, (Fine-tune), Predict		

Liu et al. 2021

Pre 2017

What we have seen so far

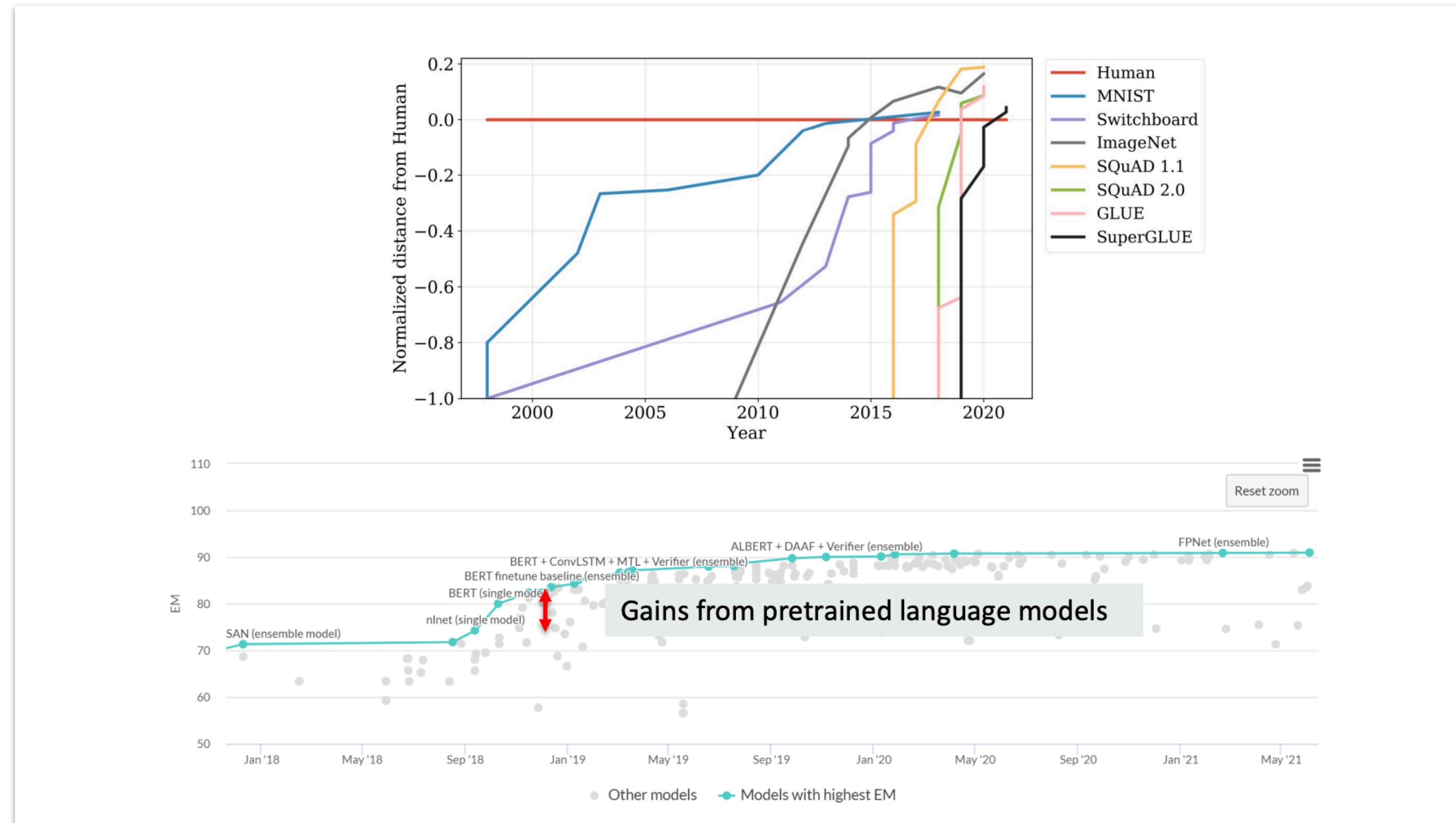
2017-2019

What we will see in the coming lectures

2021- Present?

2022- Present

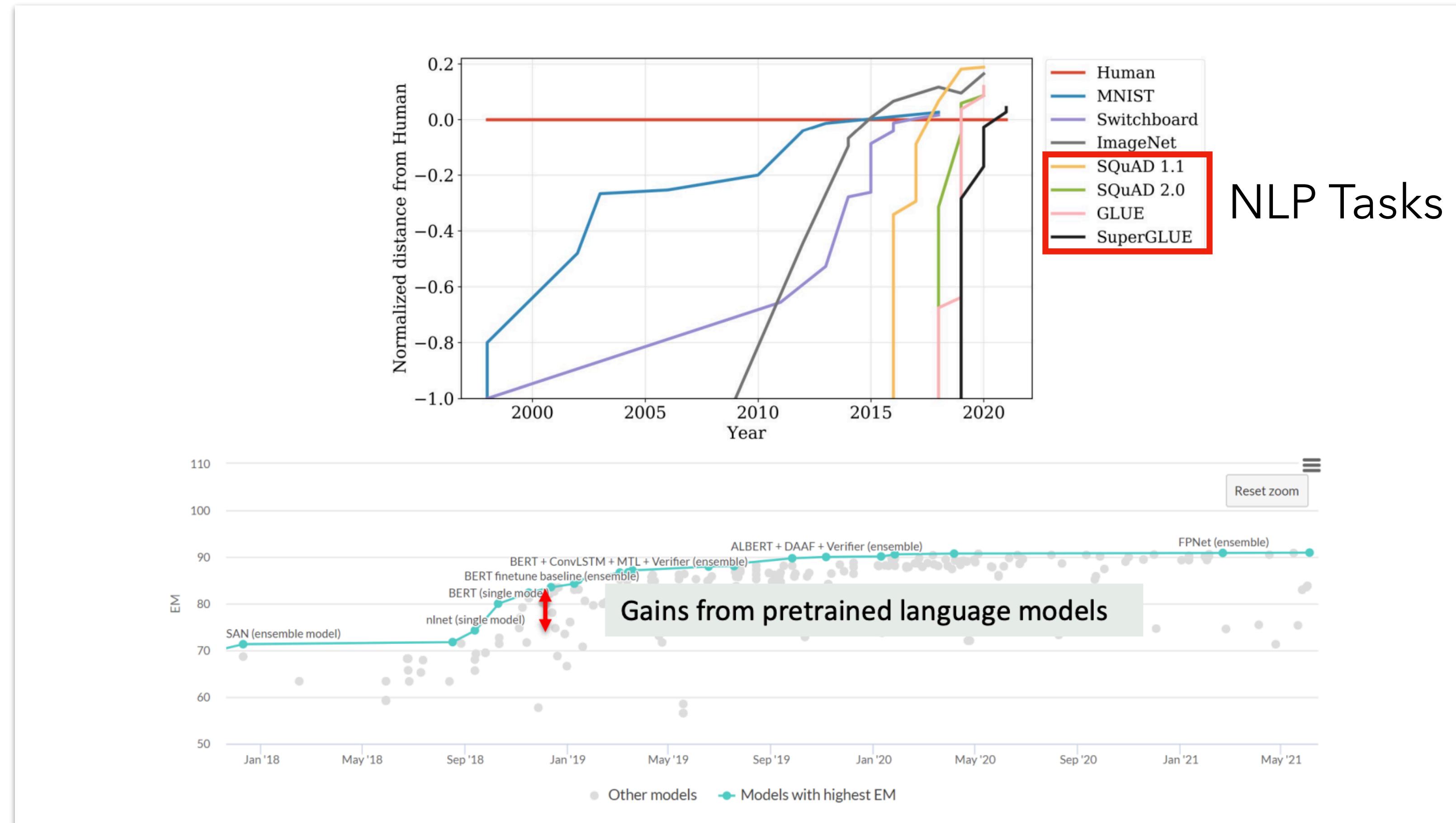
# The Pre-training Revolution



Pre-training has had a major, tangible impact on how well NLP systems work

Slide from Chris Manning, Lecture 9: Pre-training, CS224n Spring 2024

# The Pre-training Revolution



Pre-training has had a major, tangible impact on how well NLP systems work

Slide from Chris Manning, Lecture 9: Pre-training, CS224n Spring 2024

# Lecture Outline

1. Motivating Pre-training, aka Self-supervised Learning
2. Pre-training Architectures and Training Objectives
  1. Encoders
  2. Encoder-Decoders
  3. Decoder

# Lecture Outline

1. Motivating Pre-training, aka Self-supervised Learning
2. Pre-training Architectures and Training Objectives
  1. Encoders
  2. Encoder-Decoders
  3. Decoder

# Issues with Fully Supervised Learning Approaches



**Food Review:** "I recently had the pleasure of dining at Fusion Bites, and the experience was nothing short of spectacular. The menu boasts an exciting blend of global flavors, and each dish is a masterpiece in its own right."

Say that we are given a dataset of 100K food reviews with sentiment labels, **how do we train a model to perform sentiment analysis over unseen food reviews?**

# Issues with Fully Supervised Learning Approaches



**Food Review:** "I recently had the pleasure of dining at Fusion Bites, and the experience was nothing short of spectacular. The menu boasts an exciting blend of global flavors, and each dish is a masterpiece in its own right."

Say that we are given a dataset of 100K food reviews with sentiment labels, **how do we train a model to perform sentiment analysis over unseen food reviews?**

**We can directly train a randomly initialized model to take in food review texts and output “positive” or “negative” sentiment labels.**

# Issues with Fully Supervised Learning Approaches



**Food Review:** "I recently had the pleasure of dining at Fusion Bites, and the experience was nothing short of spectacular. The menu boasts an exciting blend of global flavors, and each dish is a masterpiece in its own right."



**Movie Review:** "The narrative unfolds with a steady pace, showcasing a blend of various elements. While the performances are competent, and the cinematography captures the essence of the story, the overall impact falls somewhere in the middle."

If we are instead given **movie reviews** to classify, can we use the same system trained from food reviews to predict the sentiment?

# Issues with Fully Supervised Learning Approaches



**Food Review:** "I recently had the pleasure of dining at Fusion Bites, and the experience was nothing short of spectacular. The menu boasts an exciting blend of global flavors, and each dish is a masterpiece in its own right."



**Movie Review:** "The narrative unfolds with a steady pace, showcasing a blend of various elements. While the performances are competent, and the cinematography captures the essence of the story, the overall impact falls somewhere in the middle."

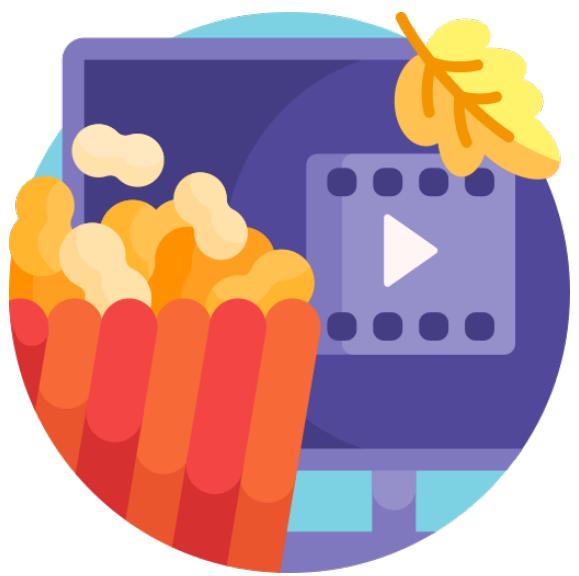
If we are instead given **movie reviews** to classify, can we use the same system trained from food reviews to predict the sentiment?

**May NOT generalize well due to distributional shift!**

# Issues with Fully Supervised Learning Approaches



**Food Review:** "I recently had the pleasure of dining at Fusion Bites, and the experience was nothing short of spectacular. The menu boasts an exciting blend of global flavors, and each dish is a masterpiece in its own right."



**Movie Review:** "The narrative unfolds with a steady pace, showcasing a blend of various elements. While the performances are competent, and the cinematography captures the essence of the story well, the plot feels somewhat predictable and lacks depth somewhere in the middle."

If we are instead given **movie reviews** to train a model, what would happen if we used a system trained from food reviews to

Fully Supervised  
Learning

Collect a labeled dataset for movie reviews and train a model from scratch on this new dataset

**May NOT generalize well due to distributional shift!**

# Transfer Learning: A History Lesson from Computer Vision

- Instead of training a randomly initialized neural network every time we encounter a new task or domain,
  - can we re-use the learned representations from one task/domain for another?

Image from Lecture 7 CS231n slides by Fei-Fei Li, Ehsan Adeli, Zane Durante

# Transfer Learning: A History Lesson from Computer Vision

- Instead of training a randomly initialized neural network every time we encounter a new task or domain,
  - can we re-use the learned representations from one task/domain for another?

**Idea:** Train a **(very) deep neural network** on a **large-scale dataset** and re-use the learned representations from this network to adapt to new tasks

Image from Lecture 7 CS231n slides by Fei-Fei Li, Ehsan Adeli, Zane Durante

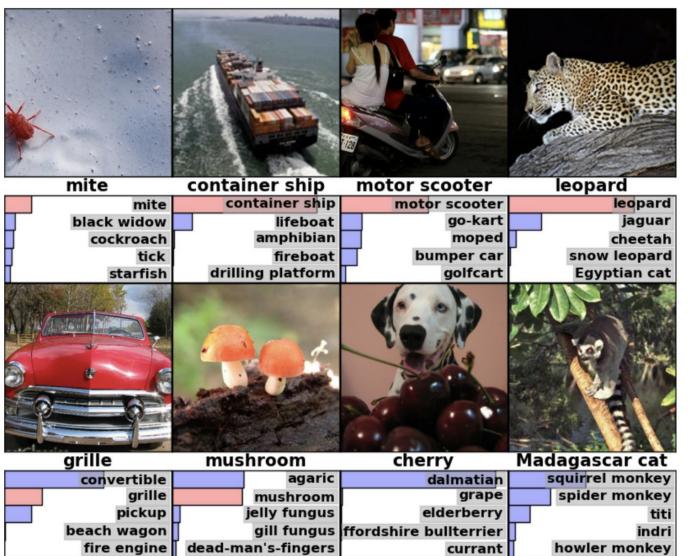
# Transfer Learning: A History Lesson from Computer Vision

- Instead of training a randomly initialized neural network every time we encounter a new task or domain,
  - can we re-use the learned representations from one task/domain for another?

**Idea:** Train a **(very) deep neural network** on a **large-scale dataset** and re-use the learned representations from this network to adapt to new tasks

## ImageNet Challenge

IMAGENET



- 1,000 object classes (categories).
- Images:
  - 1.2 M train
  - 100k test.

Image from Lecture 7 CS231n slides by Fei-Fei Li, Ehsan Adeli, Zane Durante

# Transfer Learning: A History Lesson from Computer Vision

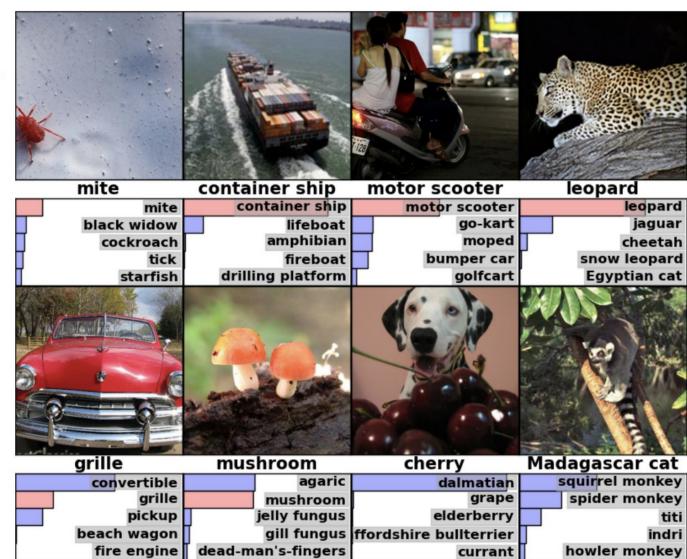
- Instead of training a randomly initialized neural network every time we encounter a new task or domain,
- can we re-use the learned representations from one task/domain for another?

**Idea:** Train a **(very) deep neural network** on a **large-scale dataset** and re-use the learned representations from this network to adapt to new tasks

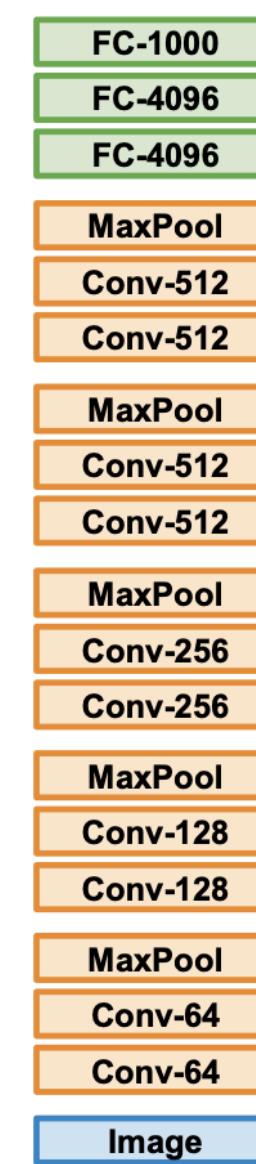
## ImageNet Challenge

IMAGENET

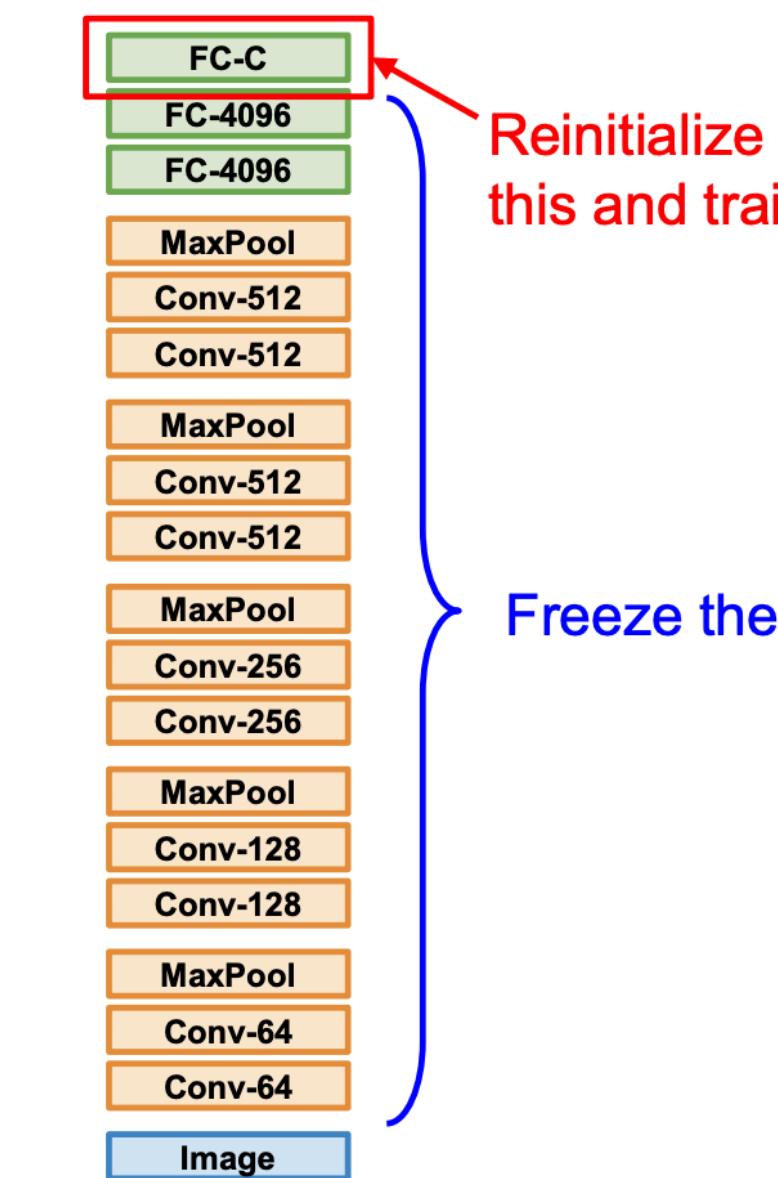
- 1,000 object classes (categories).
- Images:
  - 1.2 M train
  - 100k test.



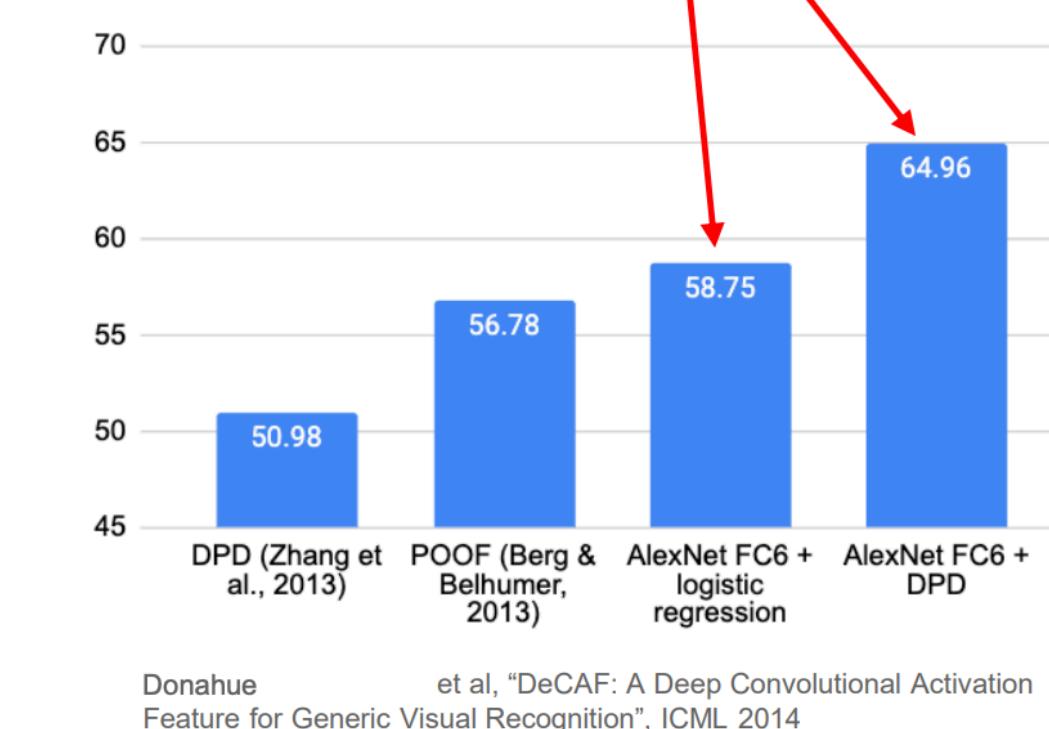
## 1. Train on Imagenet



## 2. Small Dataset ( $C$ classes)



Finetuned from AlexNet



Donahue et al., "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014

This is called Fine-tuning!

Image from Lecture 7 CS231n slides by Fei-Fei Li, Ehsan Adeli, Zane Durante

# Transfer Learning: A History Lesson from Computer Vision

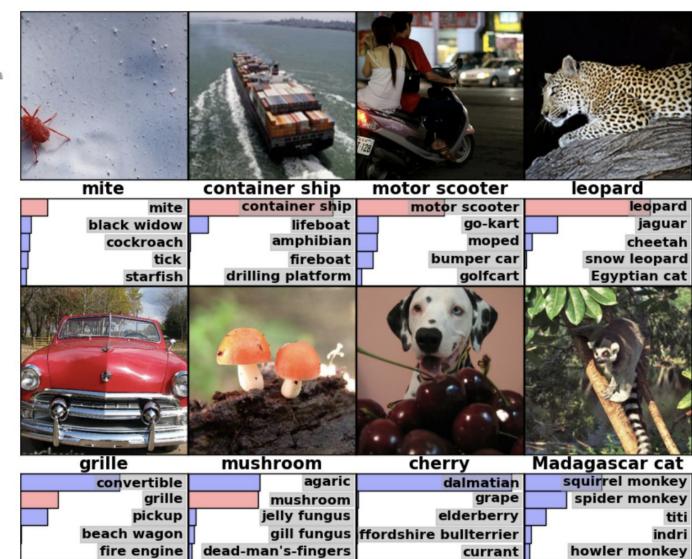
- Instead of training a randomly initialized neural network every time we encounter a new task or domain,
- can we re-use the learned representations from one task/domain for another?

**Idea:** Train a **(very) deep neural network** on a **large-scale dataset** and re-use the learned representations from this network to adapt to new tasks

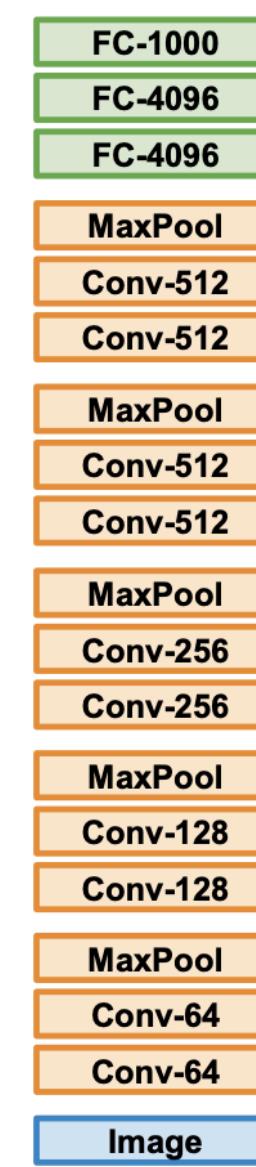
## ImageNet Challenge

IMAGENET

- 1,000 object classes (categories).
- Images:
  - 1.2 M train
  - 100k test.



## 1. Train on Imagenet



## 2. Small Dataset ( $C$ classes)



Finetuned from AlexNet

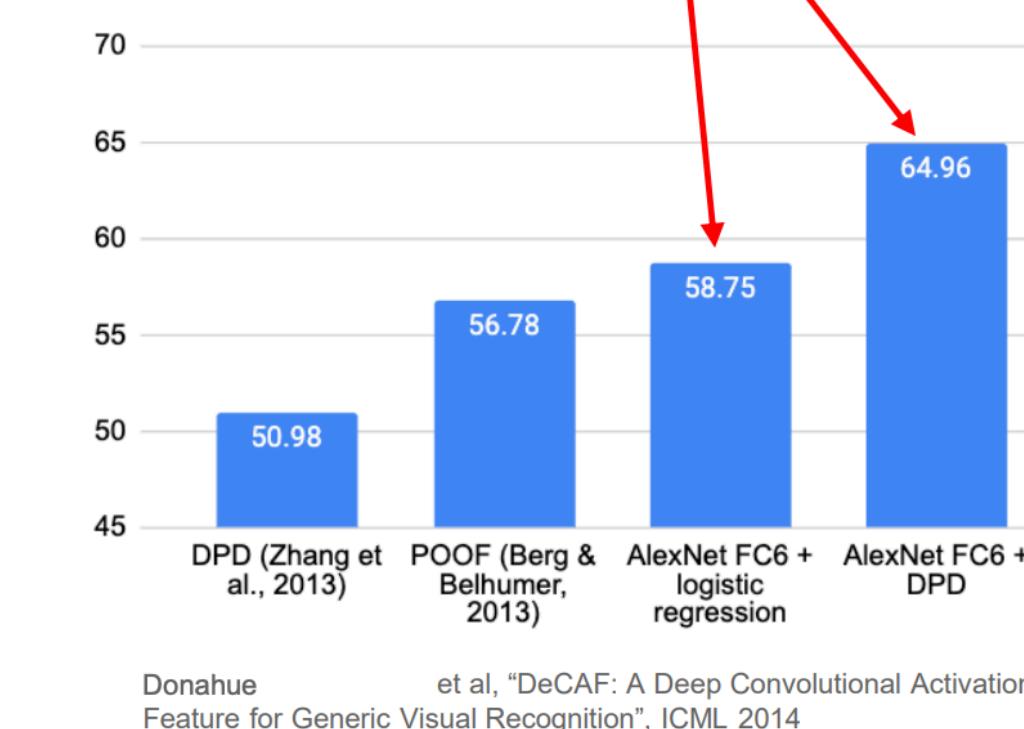


Image from Lecture 7 CS231n slides by Fei-Fei Li, Ehsan Adeli, Zane Durante

- A very successful recipe for adapting to different vision tasks like object detection, semantic segmentation, pose estimation, etc.

- Also, reduced the reliance on large training datasets to achieve good performance

This is called Fine-tuning!

# Why it took so long for NLP?

- Since 2014, it had become common practice in the Computer Vision community to download a pre-trained (on Image Net) deep neural network model and “fine-tune” it on the problem at hand instead of starting from scratch.
- This wasn’t the case in NLP till late 2017s.
- It was common to use pre-trained word vectors like word2vec, GloVe for NLP tasks, and while those would help boost performance, most often it was a marginal improvement.

# Why it took so long for NLP?

- Since 2014, it had become common practice in the Computer Vision community to download a pre-trained (on Image Net) deep neural network model and “fine-tune” it on the problem at hand instead of starting from scratch.
- This wasn’t the case in NLP till late 2017s.
- It was common to use pre-trained word vectors like word2vec, GloVe for NLP tasks, and while those would help boost performance, most often it was a marginal improvement.



You might have  
seen this already  
while attempting  
HW2

# Why it took so long for NLP?

# Why it took so long for NLP?

**We can mostly boil down this delay to two factors**

# Why it took so long for NLP?

**We can mostly boil down this delay to two factors**

1. Lack of a large-scale general dataset
  1. It wasn't clear what would be a suitable NLP task most representative of the space of NLP tasks (classification, QA, NLI, Parsing, Language Modeling?). Getting high-quality label at such a large scale was also a challenge.

# Why it took so long for NLP?

**We can mostly boil down this delay to two factors**

1. Lack of a large-scale general dataset
  1. It wasn't clear what would be a suitable NLP task most representative of the space of NLP tasks (classification, QA, NLI, Parsing, Language Modeling?). Getting high-quality label at such a large scale was also a challenge.
2. Neural Network Models for NLP were usually very shallow
  1. Pre-2017, dominant models used in NLP were recurrent neural networks e.g. LSTMs
  2. These models were usually 1-2 hidden layers, and scaling them to a large number of layers was non-trivial as these models were notoriously hard to train

# Why it took so long for NLP?

What changed starting from 2017?

We can mostly boil down this delay to two factors

1. Lack of a large-scale general dataset
  1. It wasn't clear what would be a suitable NLP task most representative of the space of NLP tasks (classification, QA, NLI, Parsing, Language Modeling?). Getting high-quality label at such a large scale was also a challenge.
2. Neural Network Models for NLP were usually very shallow
  1. Pre-2017, dominant models used in NLP were recurrent neural networks e.g. LSTMs
  2. These models were usually 1-2 hidden layers, and scaling them to a large number of layers was non-trivial as these models were notoriously hard to train

# Why it took so long for NLP?

What changed starting from 2017?

We can mostly boil down this delay to two factors

1. Lack of a large-scale general dataset

## Self-supervised Learning

1. It wasn't clear what would be a suitable NLP task most representative of the space of NLP tasks (classification, QA, NLI, Parsing, Language Modeling?). Getting high-quality label at such a large scale was also a challenge.

2. Neural Network Models for NLP were usually very shallow

1. Pre-2017, dominant models used in NLP were recurrent neural networks e.g. LSTMs

2. These models were usually 1-2 hidden layers, and scaling them to a large number of layers was non-trivial as these models were notoriously hard to train

# Why it took so long for NLP?

What changed starting from 2017?

We can mostly boil down this delay to two factors

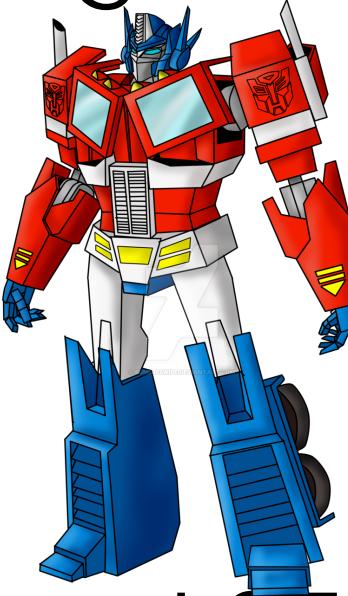
1. Lack of a large-scale general dataset

**Self-supervised Learning**

1. It wasn't clear what would be a suitable NLP task most representative of the space of NLP tasks (classification, QA, NLI, Parsing, Language Modeling?). Getting high-quality label at such a large scale was also a challenge.

2. Neural Network Models for NLP were usually very shallow

**Transformers**



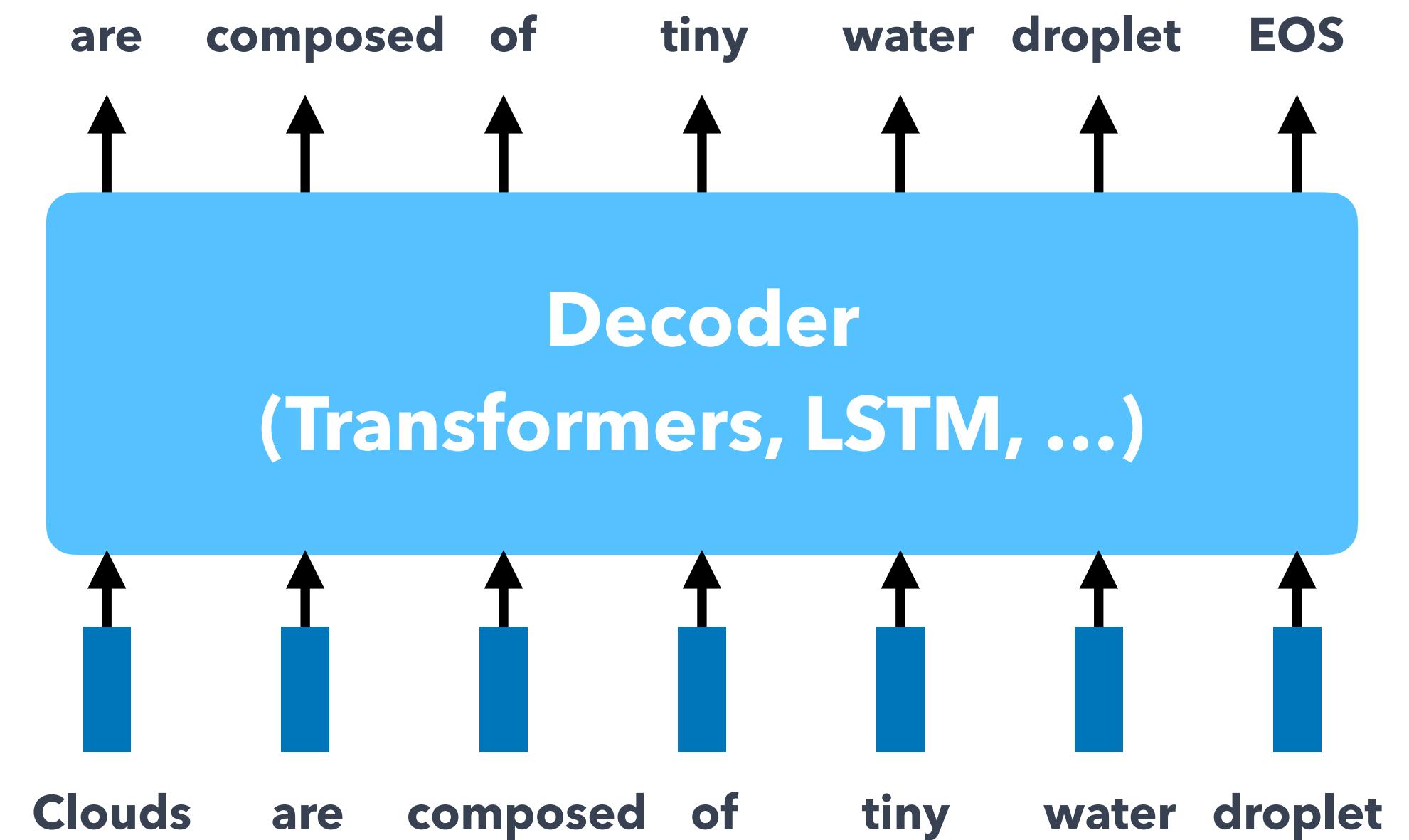
1. Pre-2017, dominant models used in NLP were recurrent neural networks e.g. LSTMs

2. These models were usually 1-2 hidden layers, and scaling them to a large number of layers was non-trivial as these models were notoriously hard to train

# **Self-supervised Pre-training for Learning Underlying Patterns, Structures, and Semantic Knowledge**

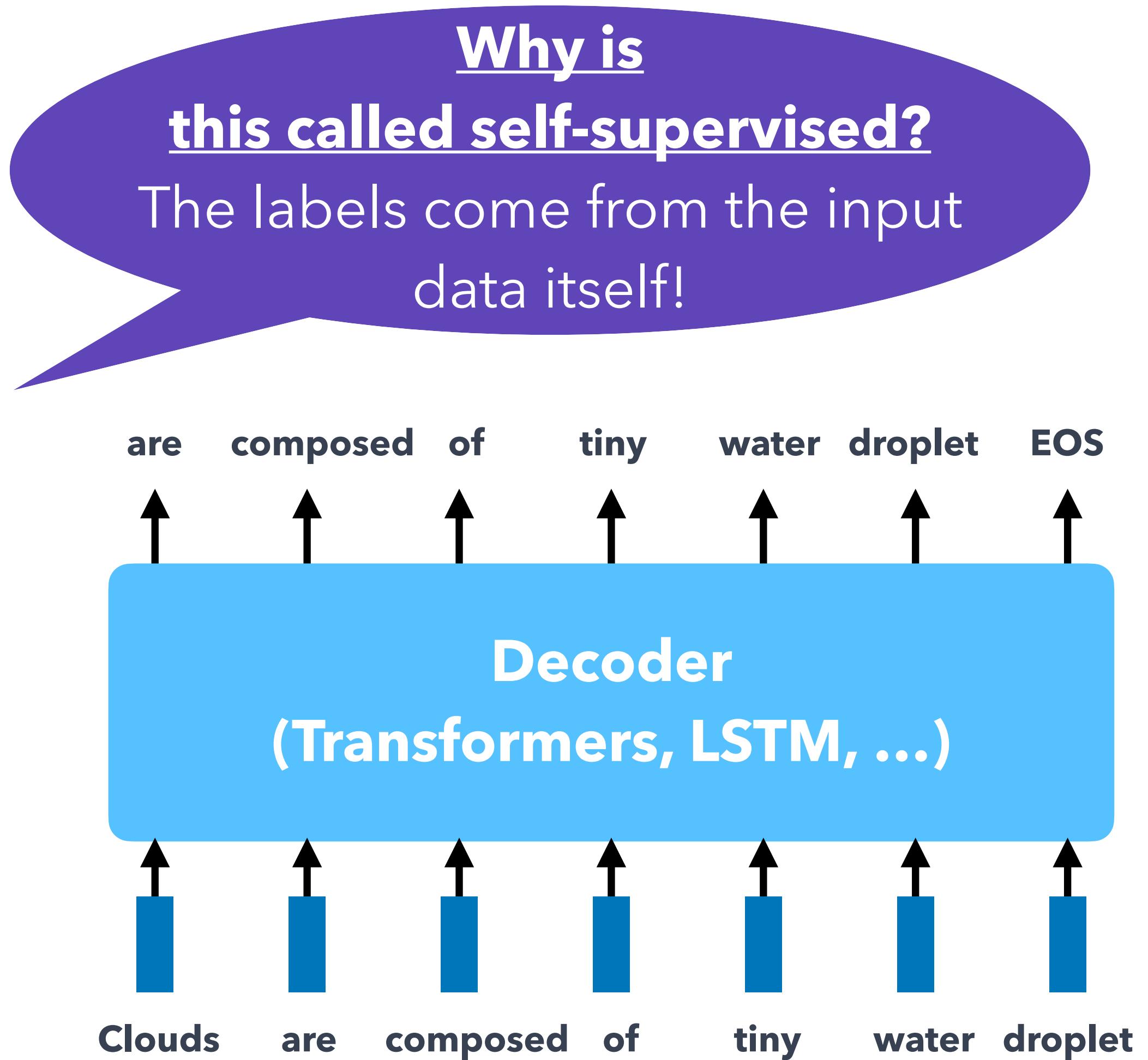
# Self-supervised Pre-training for Learning Underlying Patterns, Structures, and Semantic Knowledge

- Pre-training through **language modeling** [Dai and Le, 2015]
  - Model  $P_{\theta}(w_t | w_{1:t-1})$ , the probability distribution of the next word given previous contexts.
  - There's lots of (English) data for this!** E.g., books, websites.
  - Self-supervised** training of a neural network to perform the language modeling task with massive raw text data.
  - Save the network parameters to reuse later.



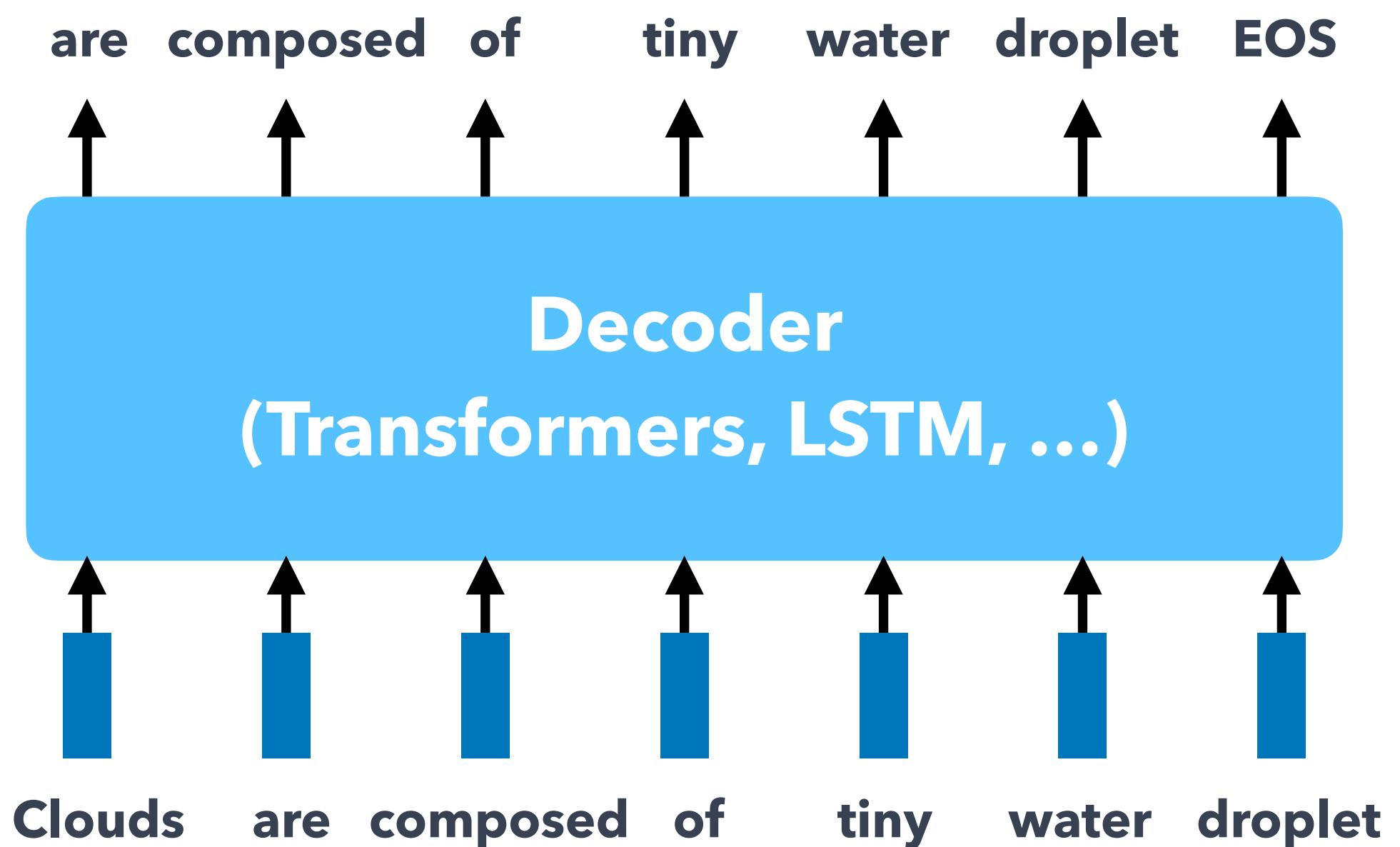
# Self-supervised Pre-training for Learning Underlying Patterns, Structures, and Semantic Knowledge

- Pre-training through **language modeling** [Dai and Le, 2015]
  - Model  $P_{\theta}(w_t | w_{1:t-1})$ , the probability distribution of the next word given previous contexts.
  - **There's lots of (English) data for this!** E.g., books, websites.
  - **Self-supervised** training of a neural network to perform the language modeling task with massive raw text data.
  - Save the network parameters to reuse later.



# Supervised Fine-tuning for Specific Tasks

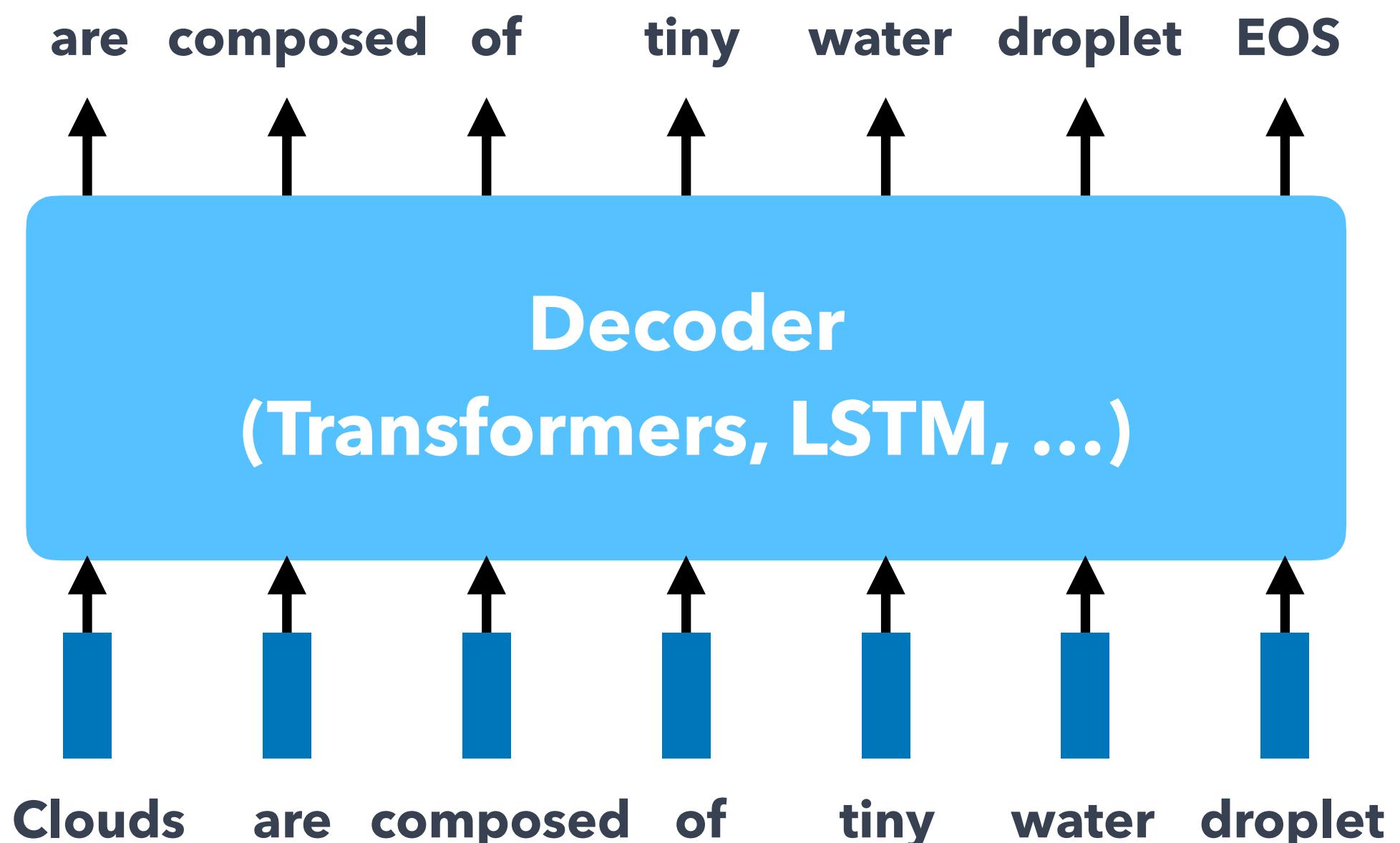
## Step 1: Pre-training



Abundant data; learn general language

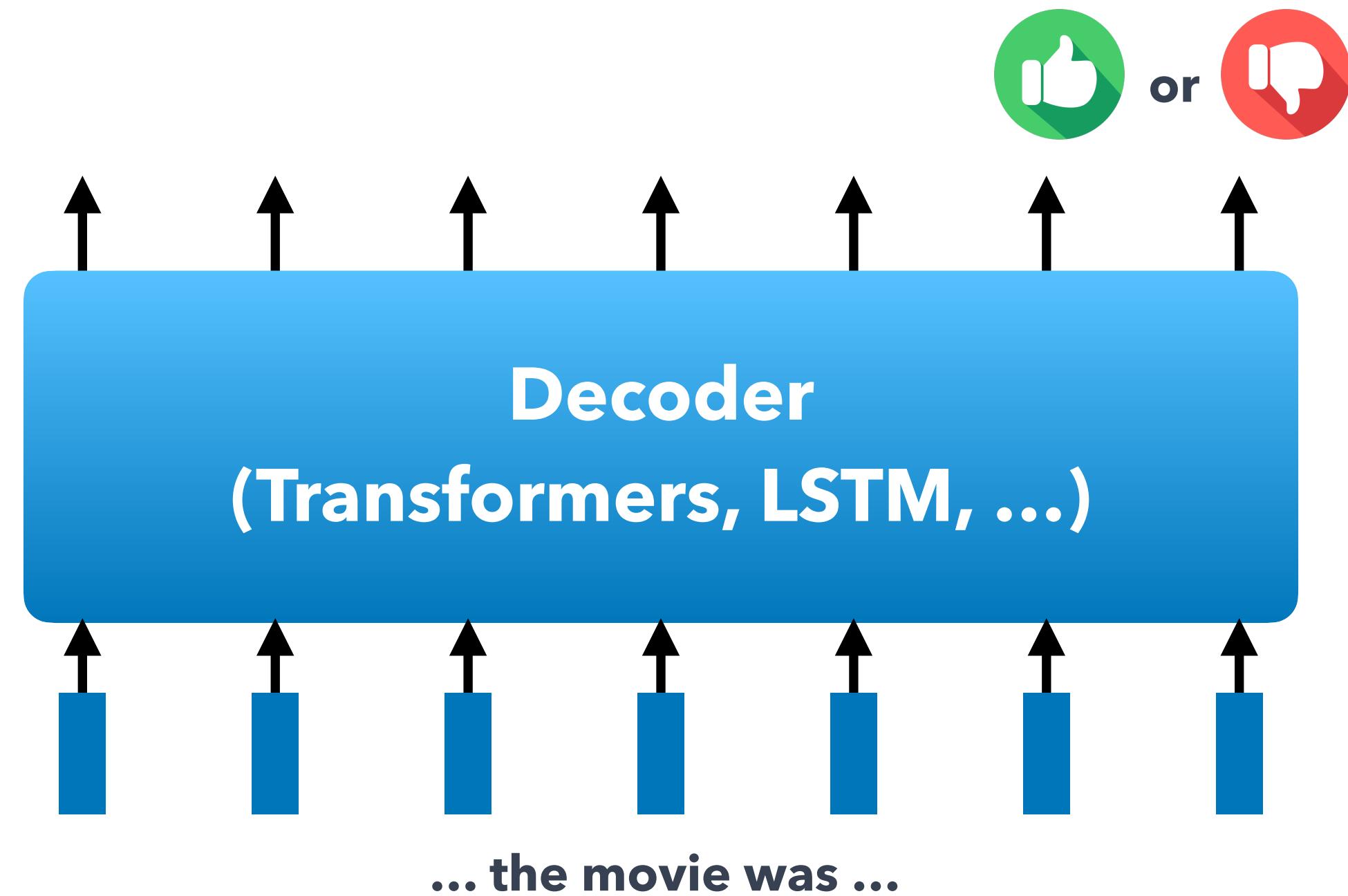
# Supervised Fine-tuning for Specific Tasks

## Step 1: Pre-training



Abundant data; learn general language

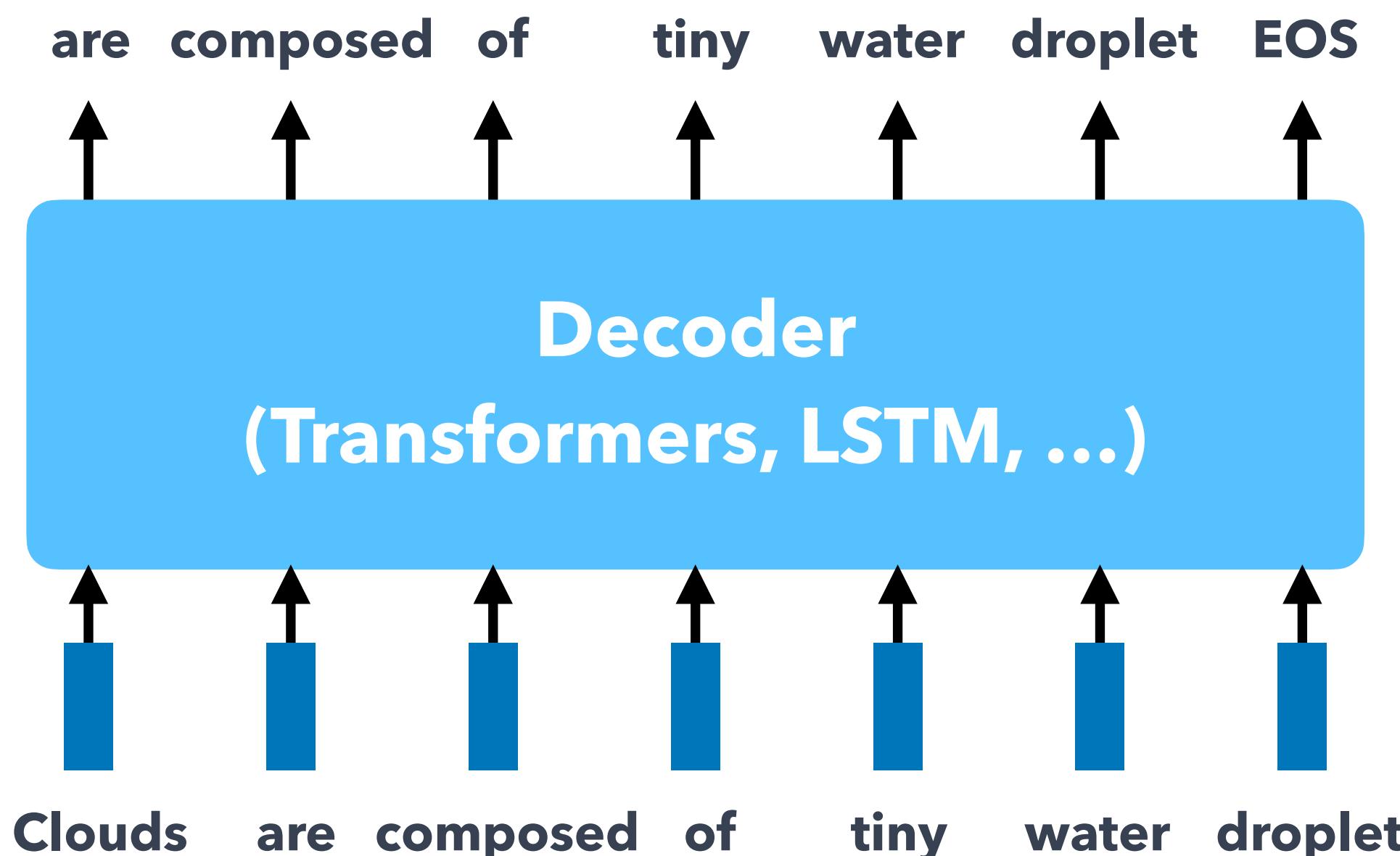
## Step 2: Fine-tuning



Limited data; adapt to the task

# Supervised Fine-tuning for Specific Tasks

## Step 1: Pre-training

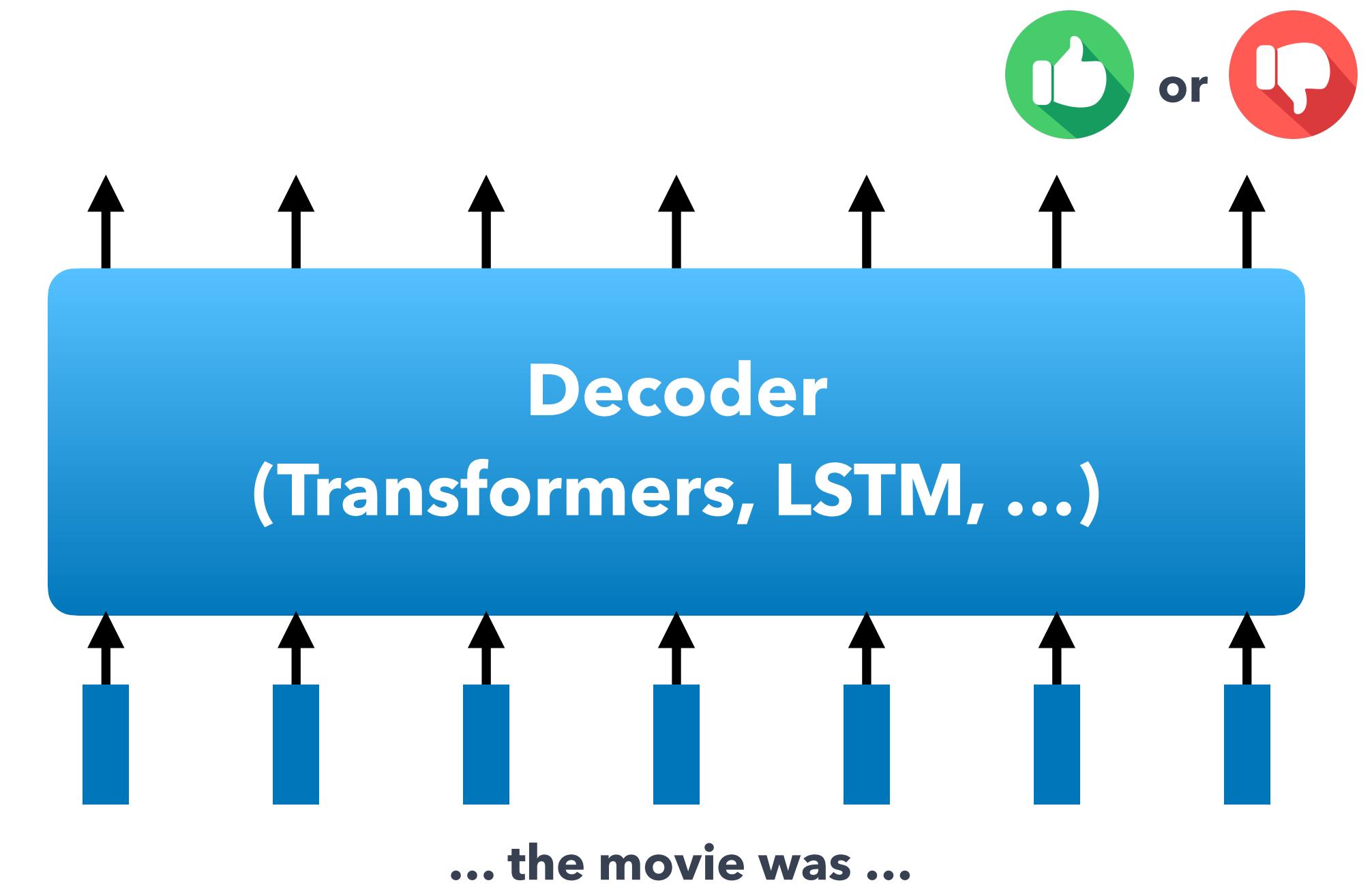


Abundant data; learn general language

**Remember this is paradigm 3 from before**

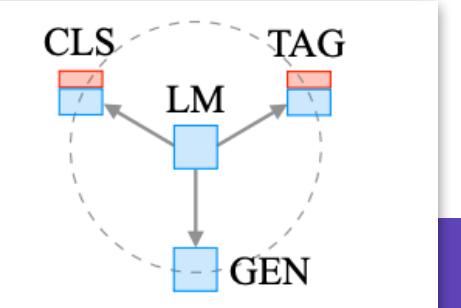
c. Pre-train, Fine-tune

## Step 2: Fine-tuning



Limited data; adapt to the task

Objective  
(e.g. masked language modeling,  
next sentence prediction)



Pre-training



# **Why this works?**

# Lots of Information in Raw Texts

I went to Hawaii for snorkeling, hiking, and whale \_\_\_\_\_.

I walked across the street, checking for traffic \_\_\_\_\_ my shoulders.

I use \_\_\_\_\_ and fork to eat steak.

Ruth Bader Ginsburg was born in \_\_\_\_\_.

University of Washington is located at \_\_\_\_\_, Washington.

I was thinking about the sequence that goes 1, 1, 2, 3, 5, 8, 13, 21, \_\_\_\_\_.

Sugar is composed of carbon, hydrogen, and \_\_\_\_\_.

# Lots of Information in Raw Texts

## Verb

I went to Hawaii for snorkeling, hiking, and whale watching.

I walked across the street, checking for traffic \_\_\_\_\_ my shoulders.

I use \_\_\_\_\_ and fork to eat steak.

Ruth Bader Ginsburg was born in \_\_\_\_\_.

University of Washington is located at \_\_\_\_\_, Washington.

I was thinking about the sequence that goes 1, 1, 2, 3, 5, 8, 13, 21, \_\_\_\_\_.

Sugar is composed of carbon, hydrogen, and \_\_\_\_\_.

# Lots of Information in Raw Texts

**Verb**

I went to Hawaii for snorkeling, hiking, and whale watching.

**Preposition**

I walked across the street, checking for traffic over my shoulders.

I use \_\_\_\_\_ and fork to eat steak.

Ruth Bader Ginsburg was born in \_\_\_\_\_.

University of Washington is located at \_\_\_\_\_, Washington.

I was thinking about the sequence that goes 1, 1, 2, 3, 5, 8, 13, 21, \_\_\_\_\_.

Sugar is composed of carbon, hydrogen, and \_\_\_\_\_.

# Lots of Information in Raw Texts

**Verb**

I went to Hawaii for snorkeling, hiking, and whale watching.

**Preposition**

I walked across the street, checking for traffic over my shoulders.

**Commonsense**

I use knife and fork to eat steak.

Ruth Bader Ginsburg was born in \_\_\_\_\_.

University of Washington is located at \_\_\_\_\_, Washington.

I was thinking about the sequence that goes 1, 1, 2, 3, 5, 8, 13, 21, \_\_\_\_\_.

Sugar is composed of carbon, hydrogen, and \_\_\_\_\_.

# Lots of Information in Raw Texts

**Verb**

I went to Hawaii for snorkeling, hiking, and whale watching.

**Preposition**

I walked across the street, checking for traffic over my shoulders.

**Commonsense**

I use knife and fork to eat steak.

**Time**

Ruth Bader Ginsburg was born in 1933.

University of Washington is located at \_\_\_\_\_, Washington.

I was thinking about the sequence that goes 1, 1, 2, 3, 5, 8, 13, 21, \_\_\_\_\_.

Sugar is composed of carbon, hydrogen, and \_\_\_\_\_.

# Lots of Information in Raw Texts

**Verb**

I went to Hawaii for snorkeling, hiking, and whale watching.

**Preposition**

I walked across the street, checking for traffic over my shoulders.

**Commonsense**

I use knife and fork to eat steak.

**Time**

Ruth Bader Ginsburg was born in 1933.

**Location**

University of Washington is located at Seattle, Washington.

I was thinking about the sequence that goes 1, 1, 2, 3, 5, 8, 13, 21, \_\_\_\_\_.

Sugar is composed of carbon, hydrogen, and \_\_\_\_\_.

# Lots of Information in Raw Texts

**Verb**

I went to Hawaii for snorkeling, hiking, and whale \_watching\_.

**Preposition**

I walked across the street, checking for traffic \_over\_ my shoulders.

**Commonsense**

I use \_\_knife\_\_ and fork to eat steak.

**Time**

Ruth Bader Ginsburg was born in \_\_1933\_\_.

**Location**

University of Washington is located at \_\_Seattle\_\_, Washington.

**Math**

I was thinking about the sequence that goes 1, 1, 2, 3, 5, 8, 13, 21, \_\_34\_\_.

Sugar is composed of carbon, hydrogen, and \_\_\_\_\_.

# Lots of Information in Raw Texts

**Verb**

I went to Hawaii for snorkeling, hiking, and whale \_watching\_.

**Preposition**

I walked across the street, checking for traffic \_over\_ my shoulders.

**Commonsense**

I use \_\_knife\_\_ and fork to eat steak.

**Time**

Ruth Bader Ginsburg was born in \_\_1933\_\_.

**Location**

University of Washington is located at \_\_Seattle\_\_, Washington.

**Math**

I was thinking about the sequence that goes 1, 1, 2, 3, 5, 8, 13, 21, \_\_34\_\_.

**Chemistry**

Sugar is composed of carbon, hydrogen, and \_\_oxygen\_\_.

# Lots of Information in Raw Texts

**Verb**

I went to Hawaii for snorkeling, hiking, and whale watching.

**Preposition**

I walked across the street, checking for traffic over my shoulders.

**Commonsense**

I use knife and fork to eat steak.

**Time**

Ruth Bader Ginsburg was born in 1933.

**Location**

University of Washington is located at Seattle, Washington.

**Math**

I was thinking about the sequence that goes 1, 1, 2, 3, 5, 8, 13, 21, 34.

**Chemistry**

Sugar is composed of carbon, hydrogen, and oxygen.

...

# The Stochastic Gradient Descent Angle

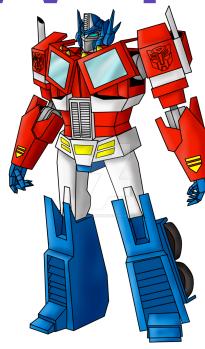
## Why should pre-training and then fine-tuning help?

- Providing parameters  $\hat{\theta}$  by approximating the pre-training loss,  
$$\min_{\theta} \mathcal{L}_{\text{pretrain}}(\theta).$$
- Then, starting with parameters  $\hat{\theta}$ , approximating fine-tuning loss,  
$$\min_{\theta} \mathcal{L}_{\text{finetune}}(\theta).$$
- **Stochastic gradient descent sticks (relatively) close to  $\hat{\theta}$  during fine-tuning.**
  - So, maybe the fine-tuning local minima near  $\hat{\theta}$  tend to generalize well!
  - And/or, maybe the gradients of fine-tuning loss near  $\hat{\theta}$  propagate nicely!

# Advantages of Pre-training & Fine-tuning

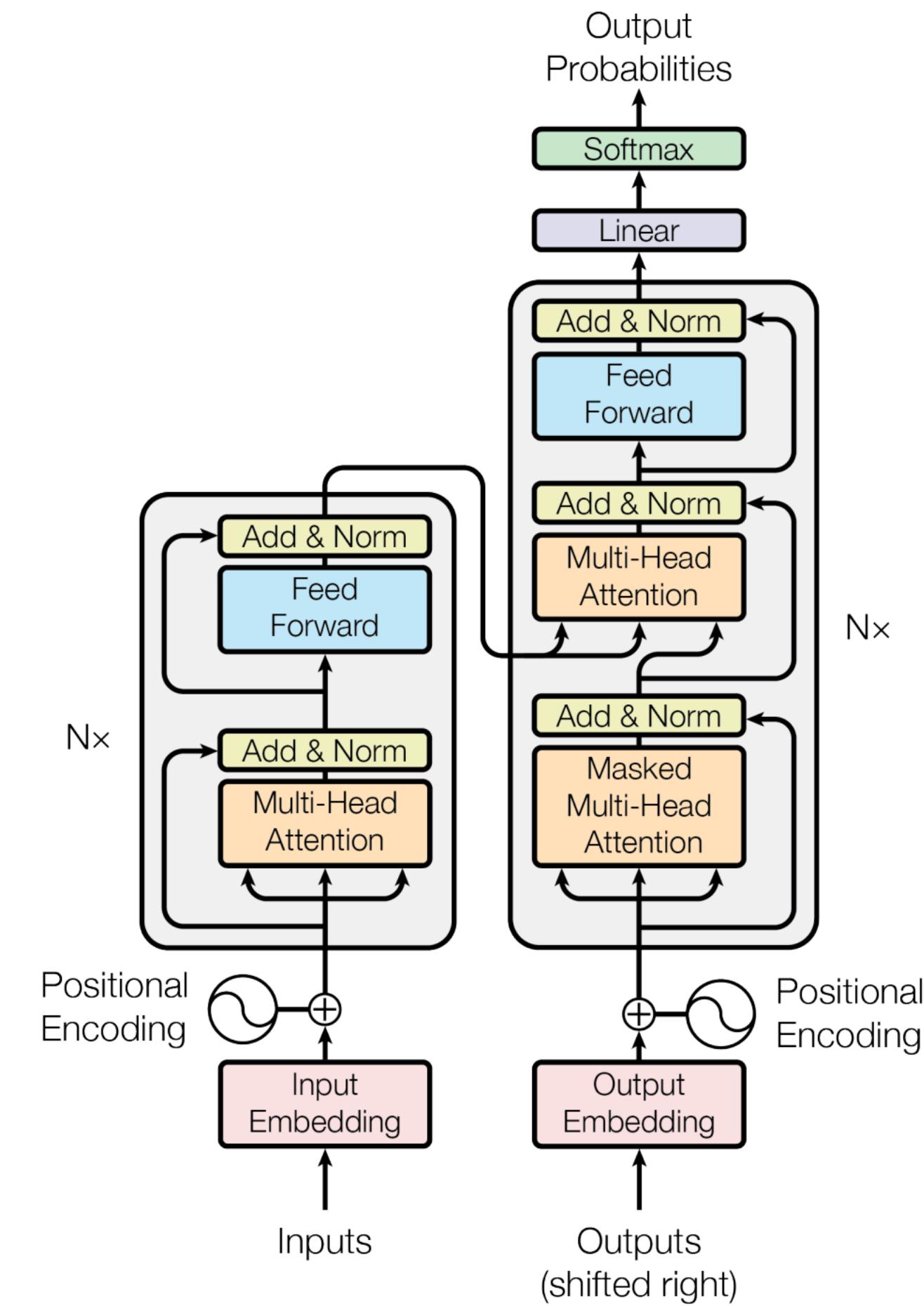
- **Leveraging rich underlying information** from abundant raw texts.
- **Reducing the reliance of task-specific labeled data** that is difficult or costly to obtain.
- **Initializing model parameters** for more **generalizable** NLP applications.
- **Saving training cost** by providing a reusable model checkpoints.
- **Providing robust representation** of language contexts.

# Solving Shallow Networks Problem in NLP: Enter Transformers



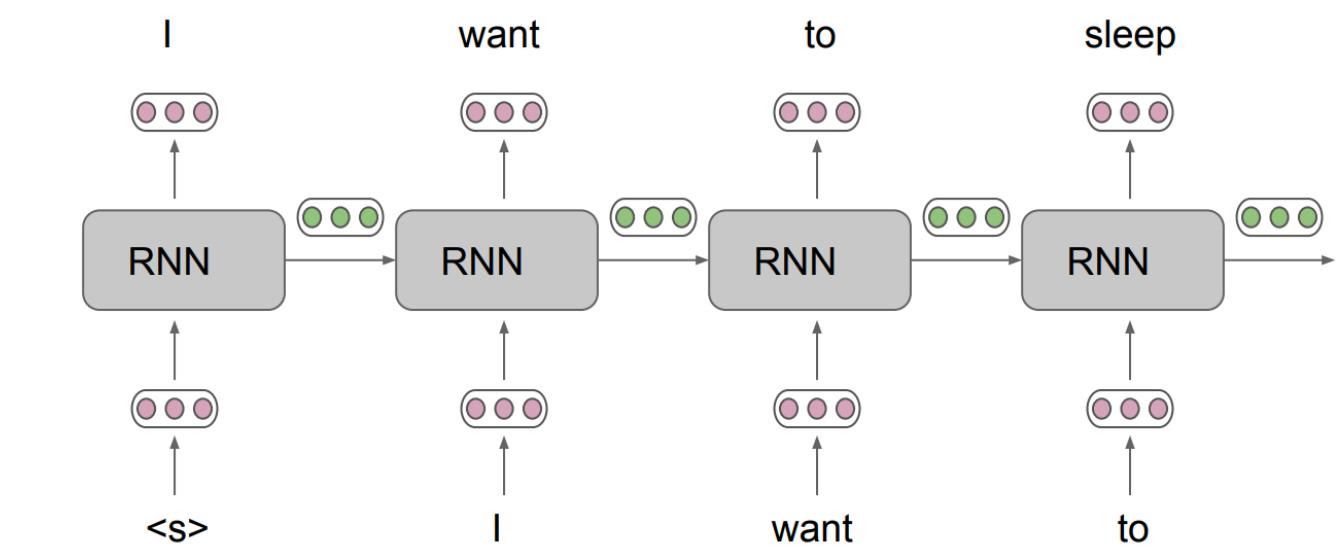
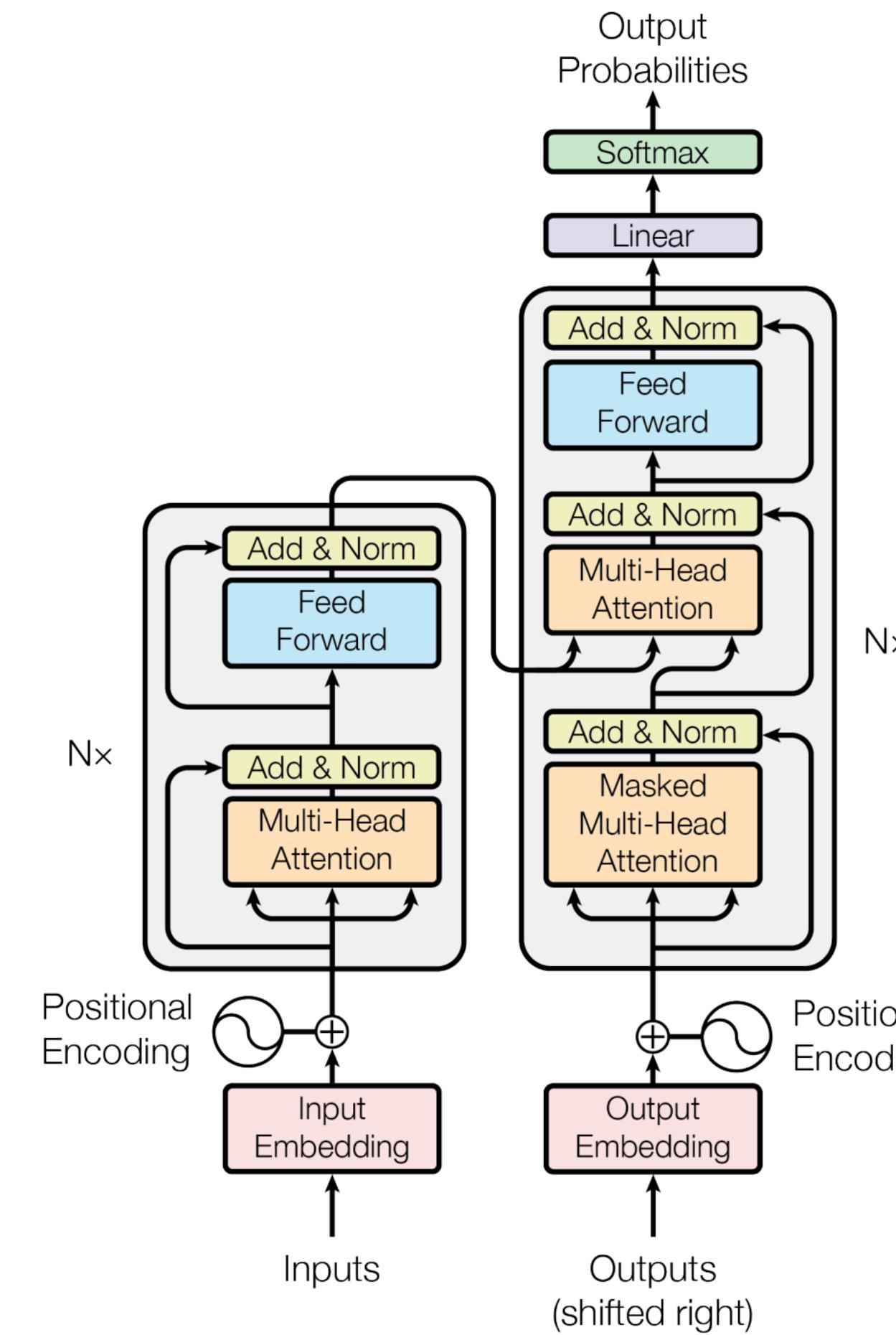
Attention is all You Need. 2017.

# Solving Shallow Networks Problem in NLP: Enter Transformers



Attention is all You Need. 2017.

# Solving Shallow Networks Problem in NLP: Enter Transformers

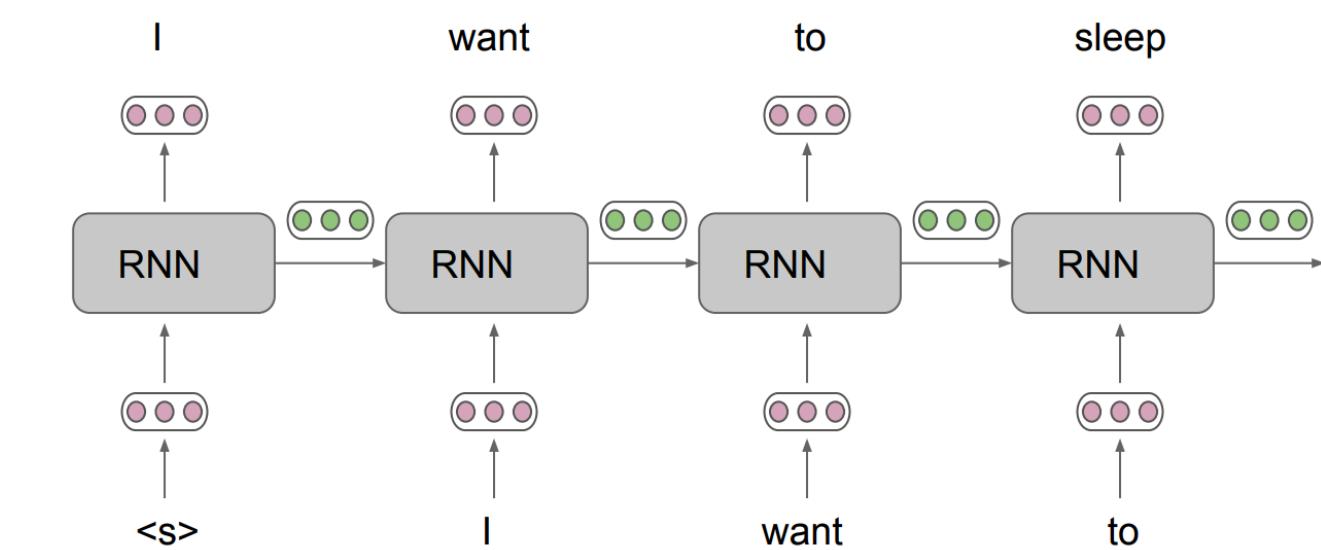
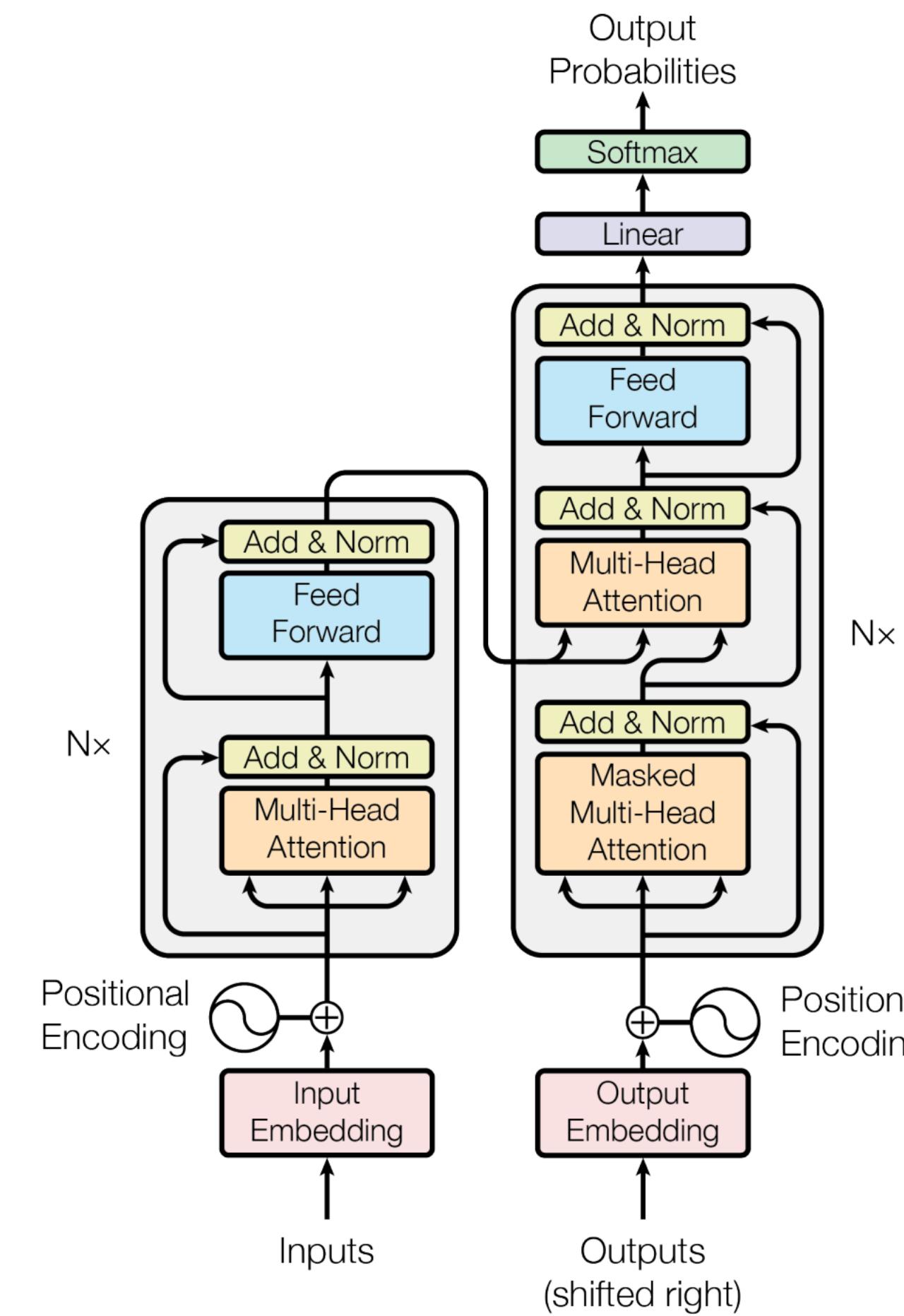


Attention is all You Need. 2017.

# Solving Shallow Networks Problem in NLP: Enter Transformers

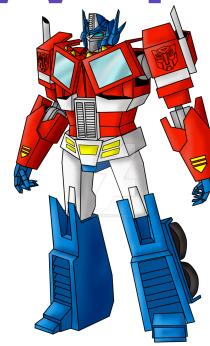


Transformers managed to avoid the two major problems that made Recurrent Neural Networks hard to scale on larger depths:



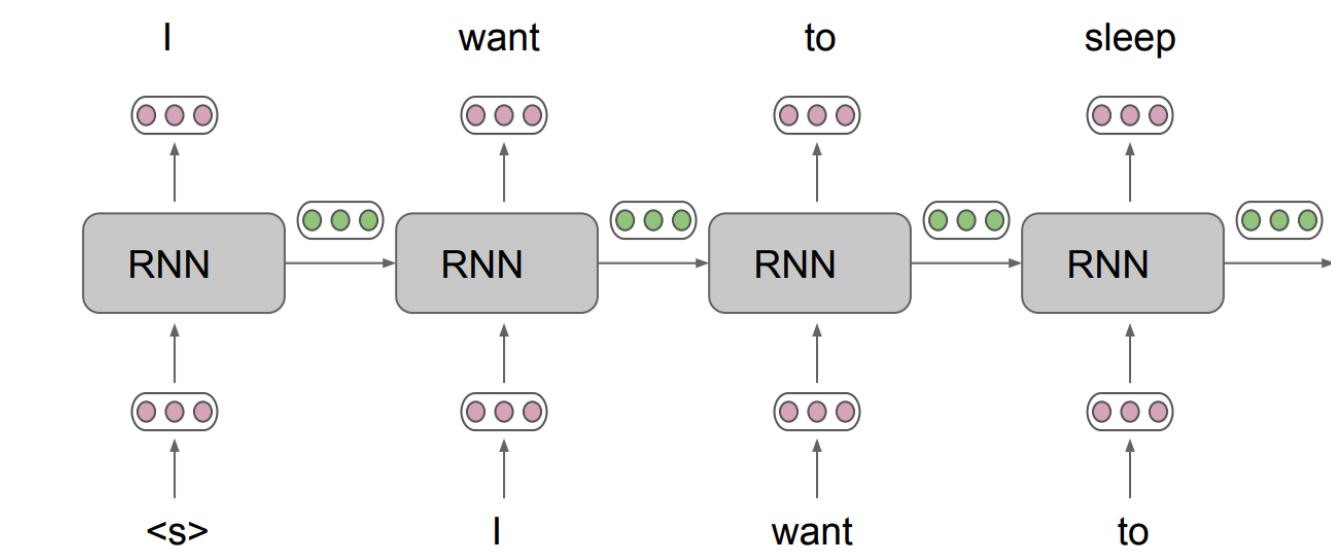
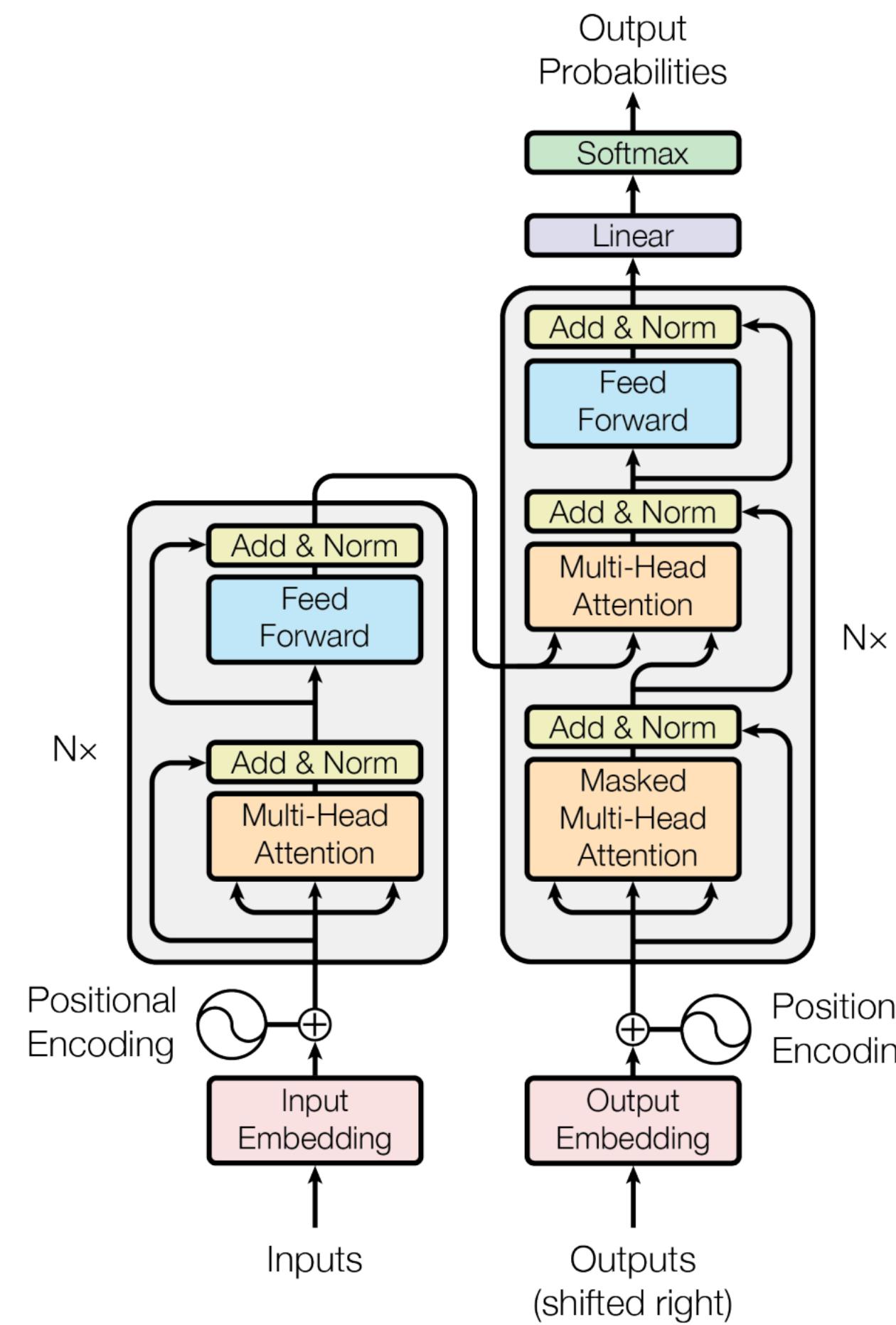
Attention is all You Need. 2017.

# Solving Shallow Networks Problem in NLP: Enter Transformers



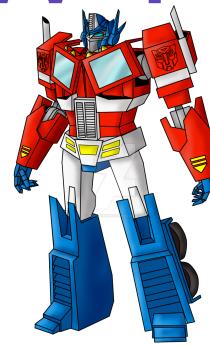
Transformers managed to avoid the two major problems that made Recurrent Neural Networks hard to scale on larger depths:

- **Highly Parallelizable During Training:**  
Need not wait for the computation at the previous time step to complete to execute the next step



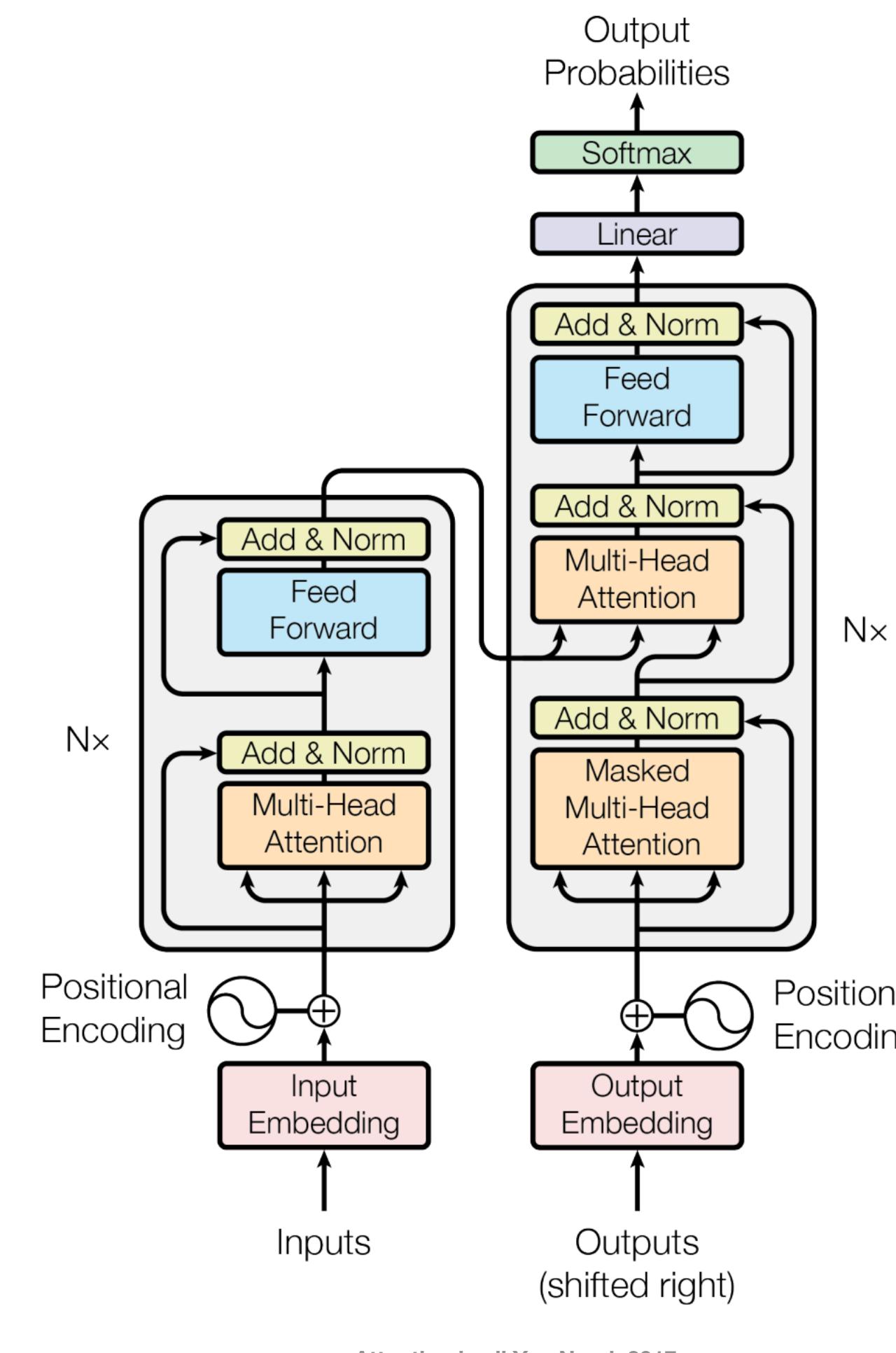
Attention is all You Need. 2017.

# Solving Shallow Networks Problem in NLP: Enter Transformers

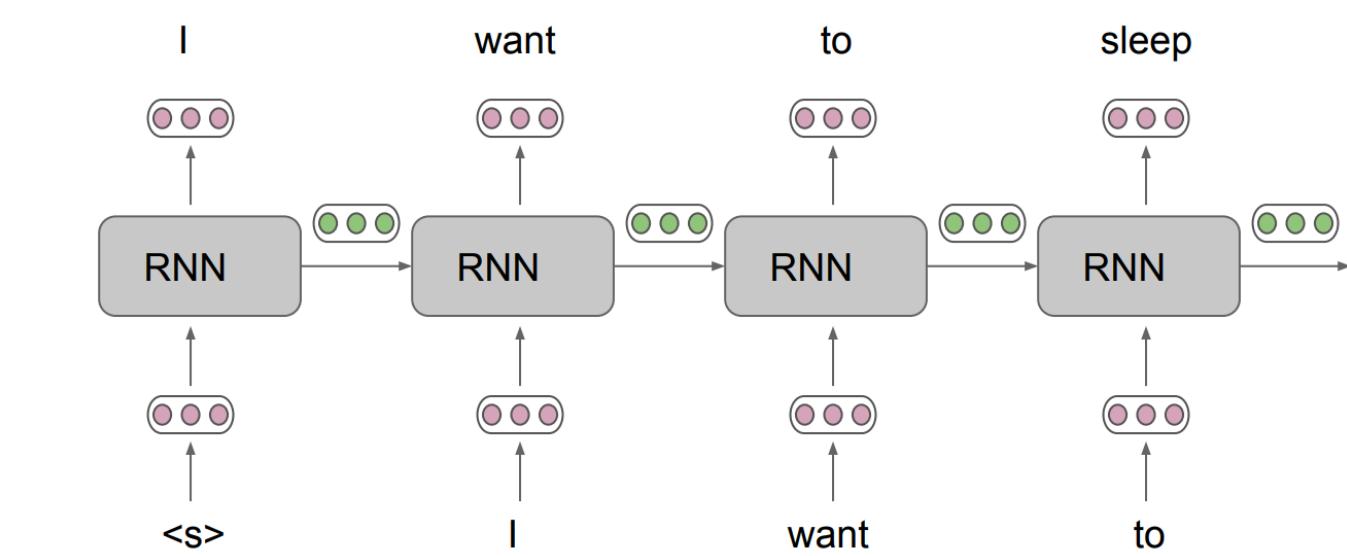


Transformers managed to avoid the two major problems that made Recurrent Neural Networks hard to scale on larger depths:

- **Highly Parallelizable During Training:**  
Need not wait for the computation at the previous time step to complete to execute the next step
- **Avoids Training Complications like Vanishing Gradients:** Unlike RNNs, which have a fixed state that gets updated repeatedly, transformers have dynamic memory, which also avoids issues such as vanishing gradients



Attention is all You Need. 2017.



# Lecture Outline

1. Motivating Pre-training, aka Self-supervised Learning
2. Pre-training Architectures and Training Objectives
  1. Encoders
  2. Encoder-Decoders
  3. Decoder

# 3 Pre-training Paradigms/Architectures

**Encoder**

- E.g., BERT, RoBERTa, DeBERTa, ...
- **Autoencoder** model
- **Masked** language modeling

**Encoder-Decoder**

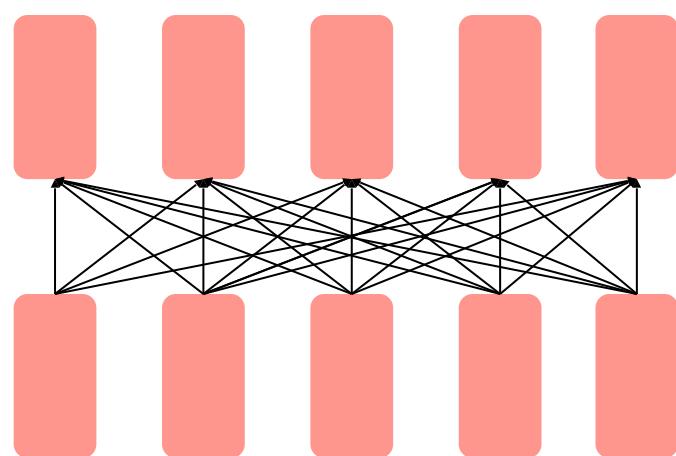
- E.g., T5, BART, ...
- **seq2seq** model

**Decoder**

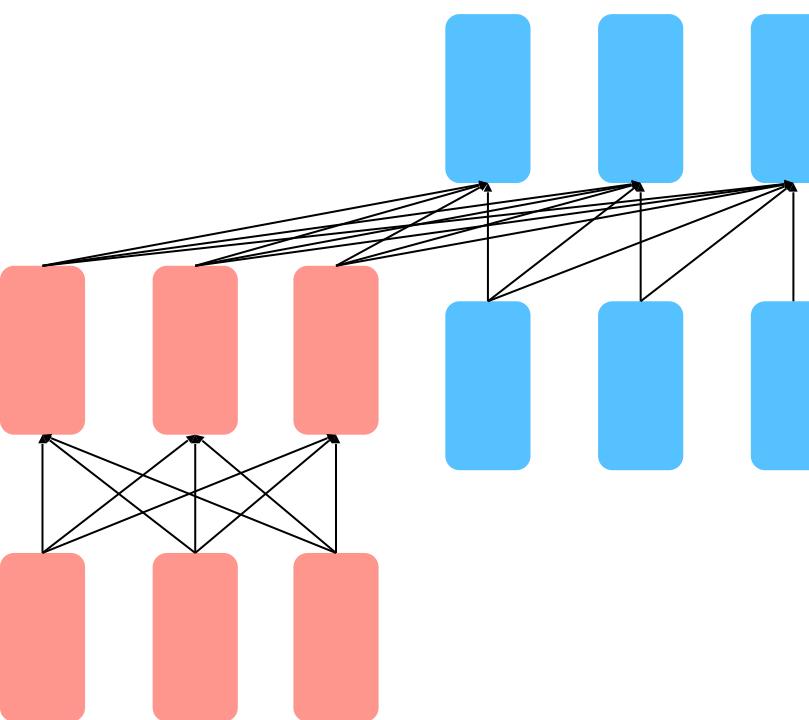
- E.g., GPT, GPT2, GPT3, ...
- **Autoregressive** model
- **Left-to-right** language modeling

# 3 Pre-training Paradigms/Architectures

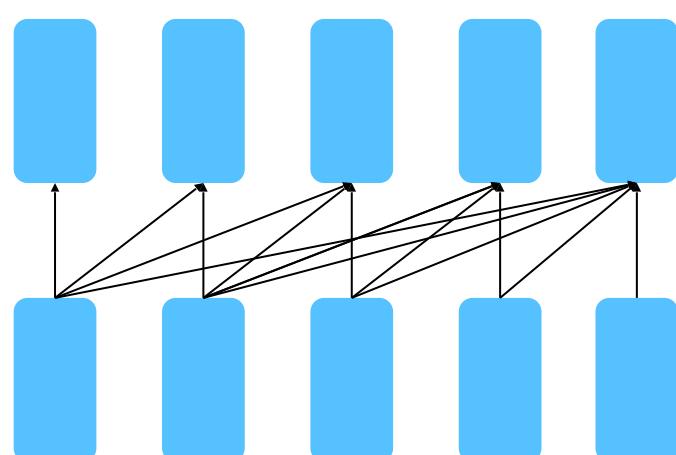
**Encoder**



**Encoder-Decoder**



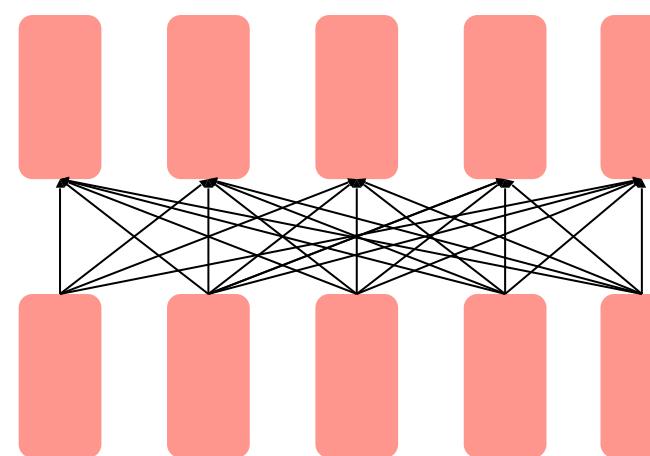
**Decoder**



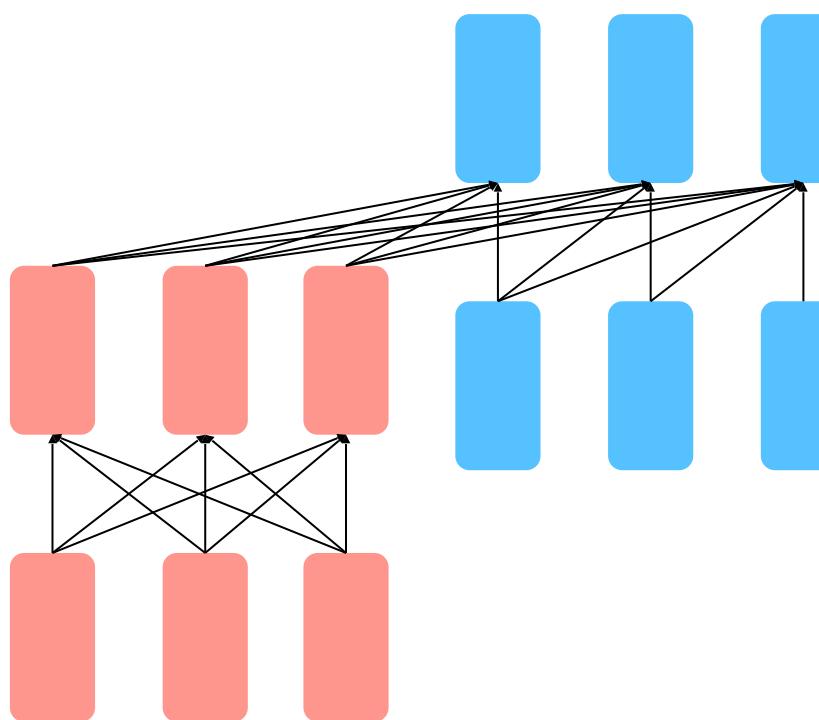
- Bidirectional; can condition on the future context
- Map two sequences of different length together
- Language modeling; can only condition on the past context

# 3 Pre-training Paradigms/Architectures

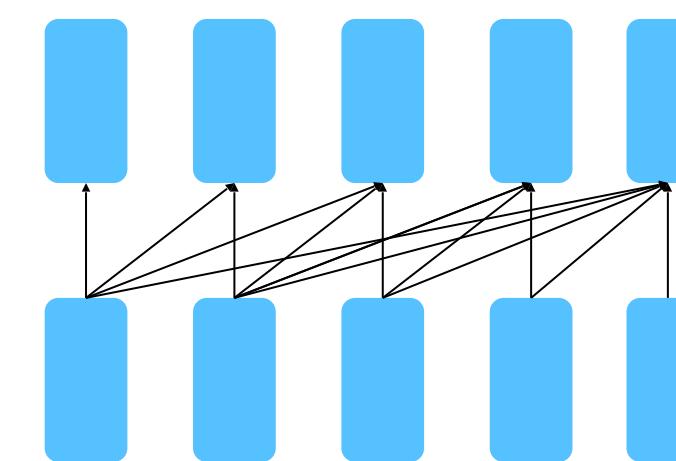
**Encoder**



**Encoder-Decoder**



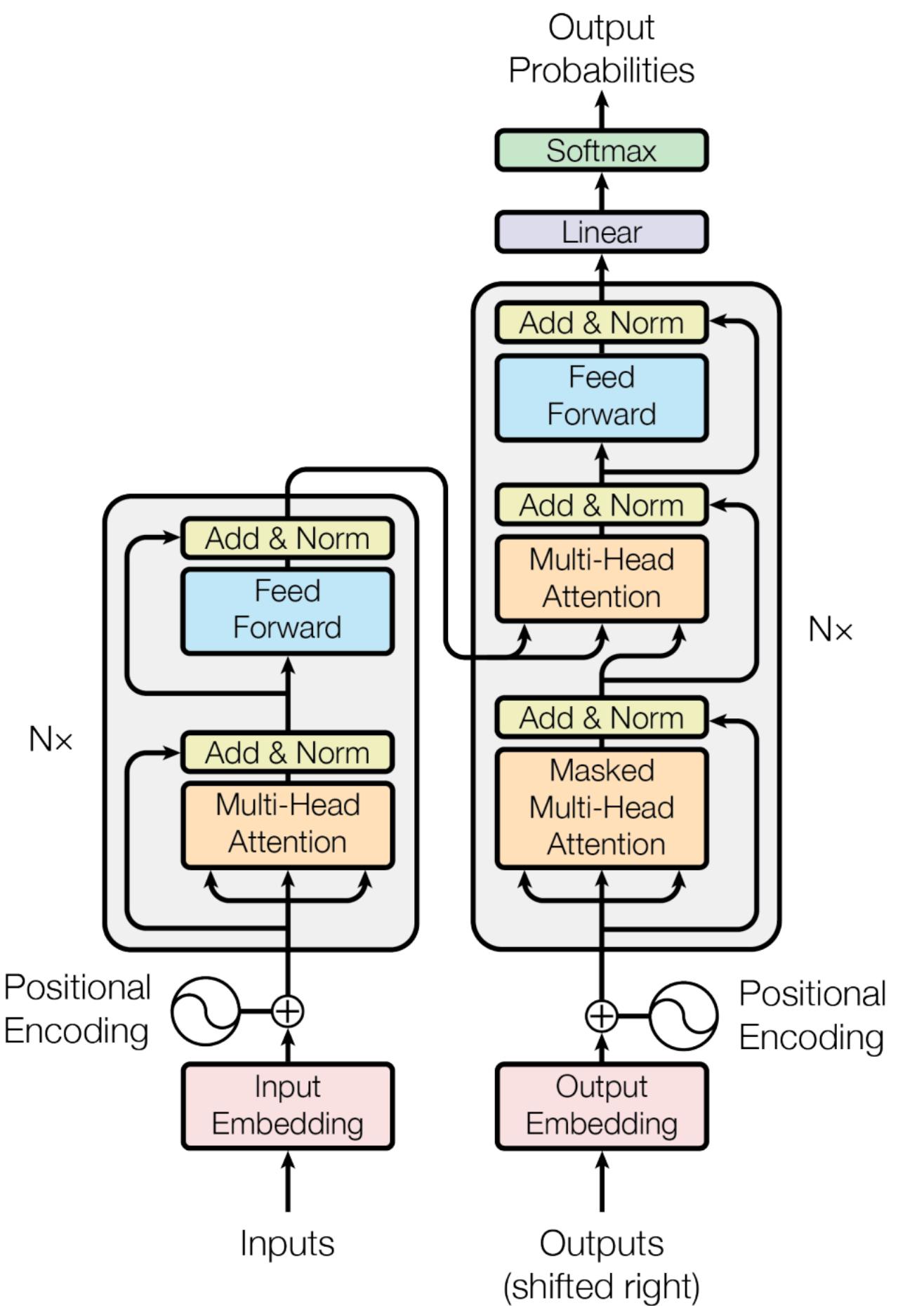
**Decoder**



- Bidirectional; can condition on the future context
- Map two sequences of different length together
- Language modeling; can only condition on the past context

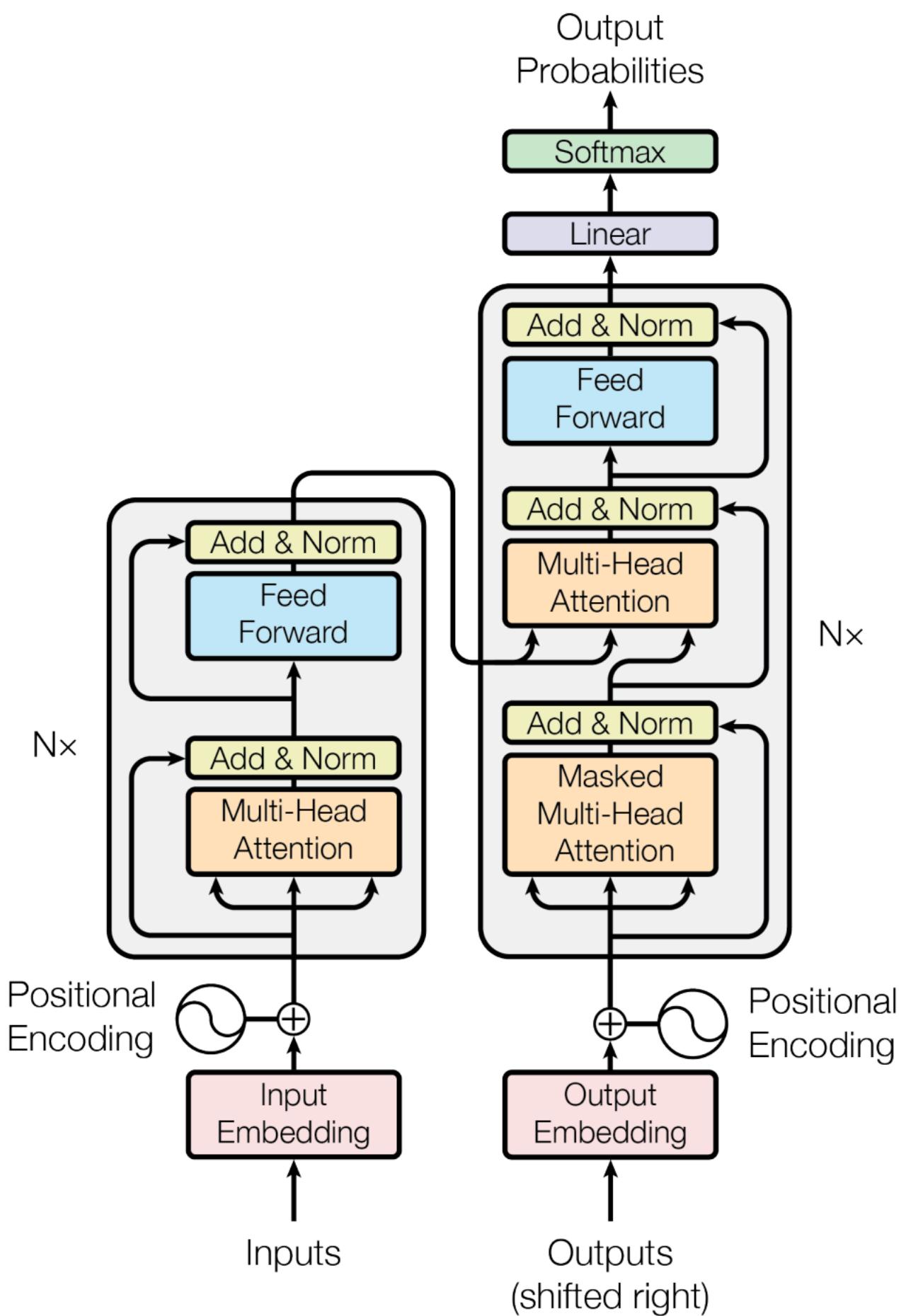
# Encoder: Architecture

## Full-Transformer Architecture (Encoder-Decoder)

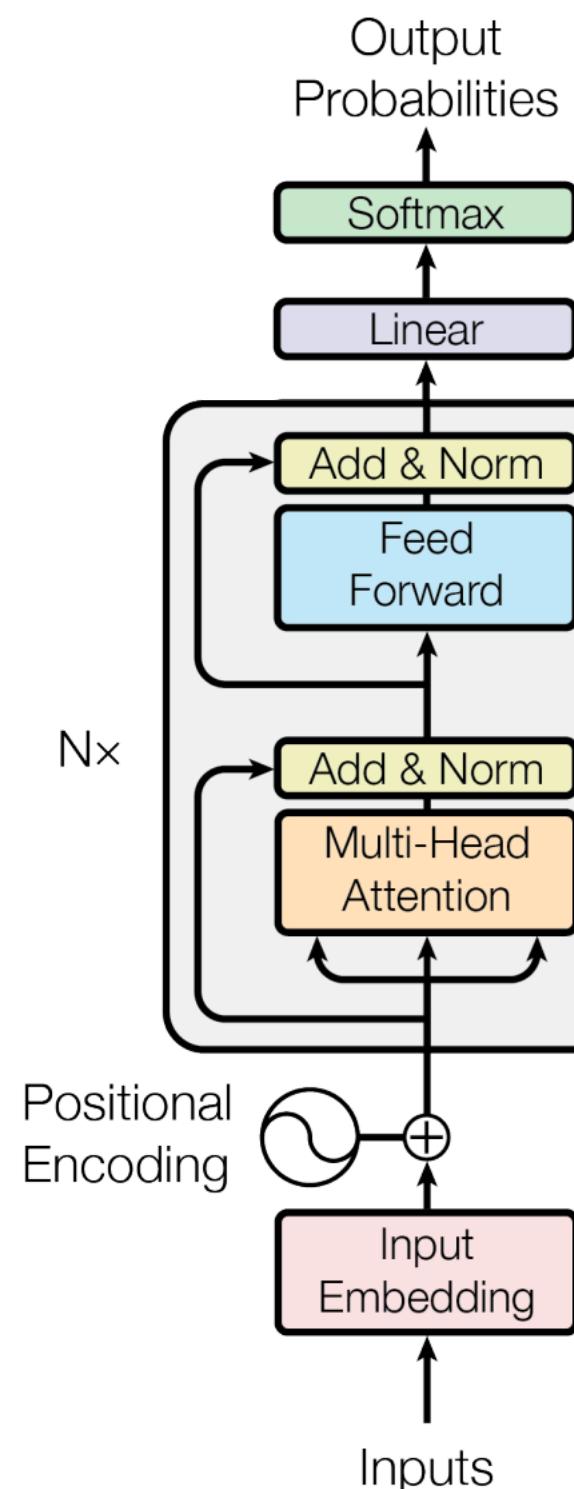


# Encoder: Architecture

**Full-Transformer Architecture  
(Encoder-Decoder)**

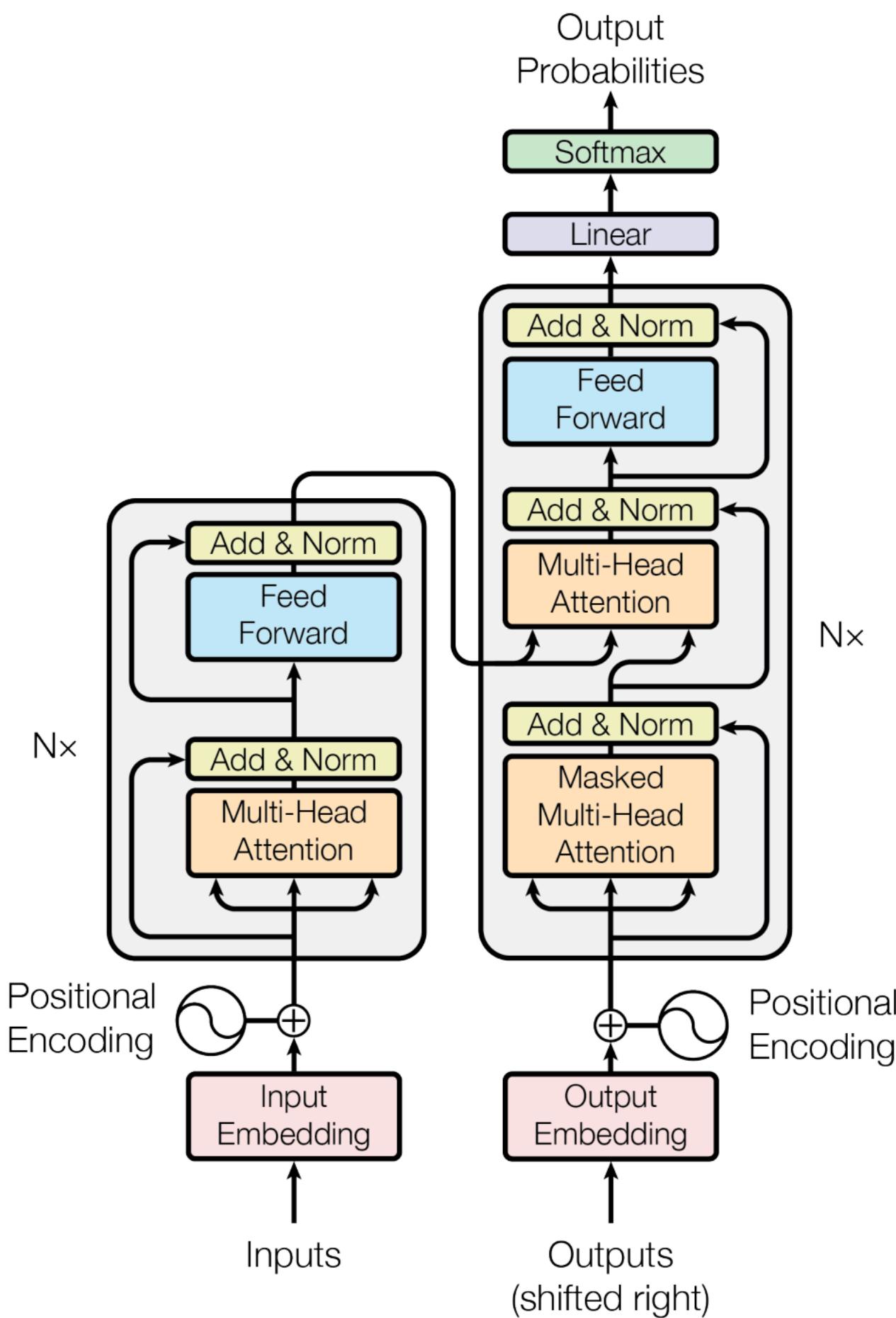


**Encoder-Only Transformer  
Architecture**

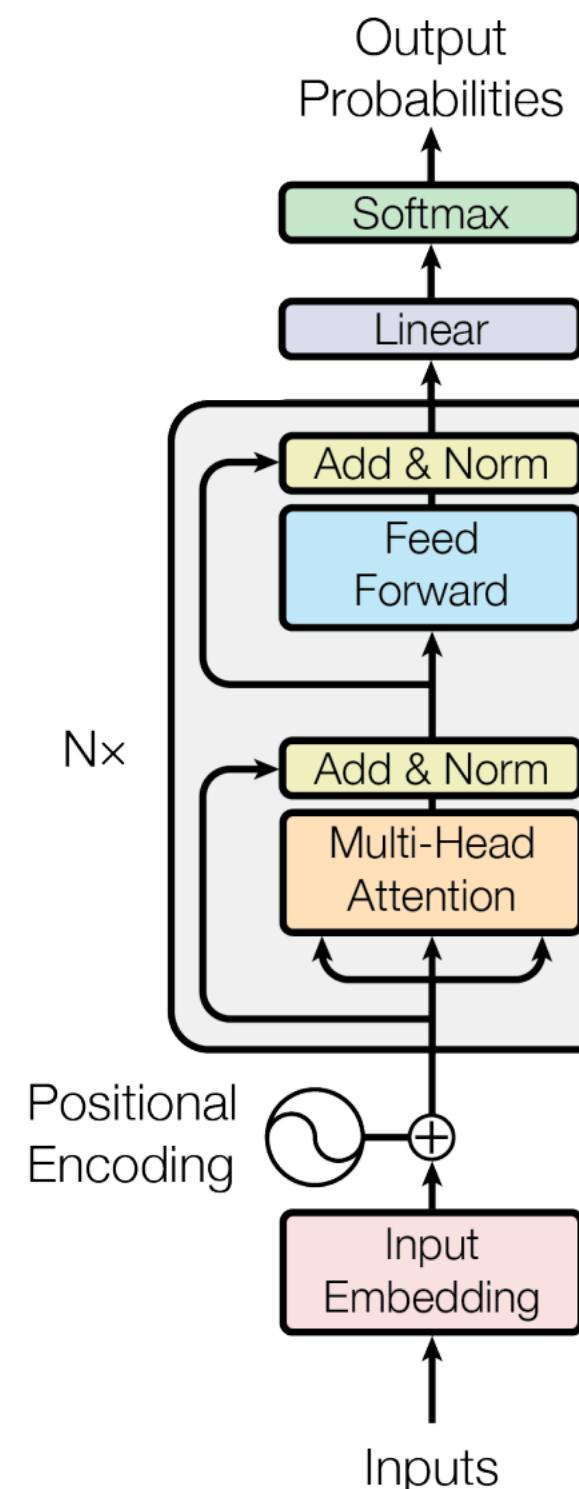


# Encoder: Architecture

**Full-Transformer Architecture  
(Encoder-Decoder)**

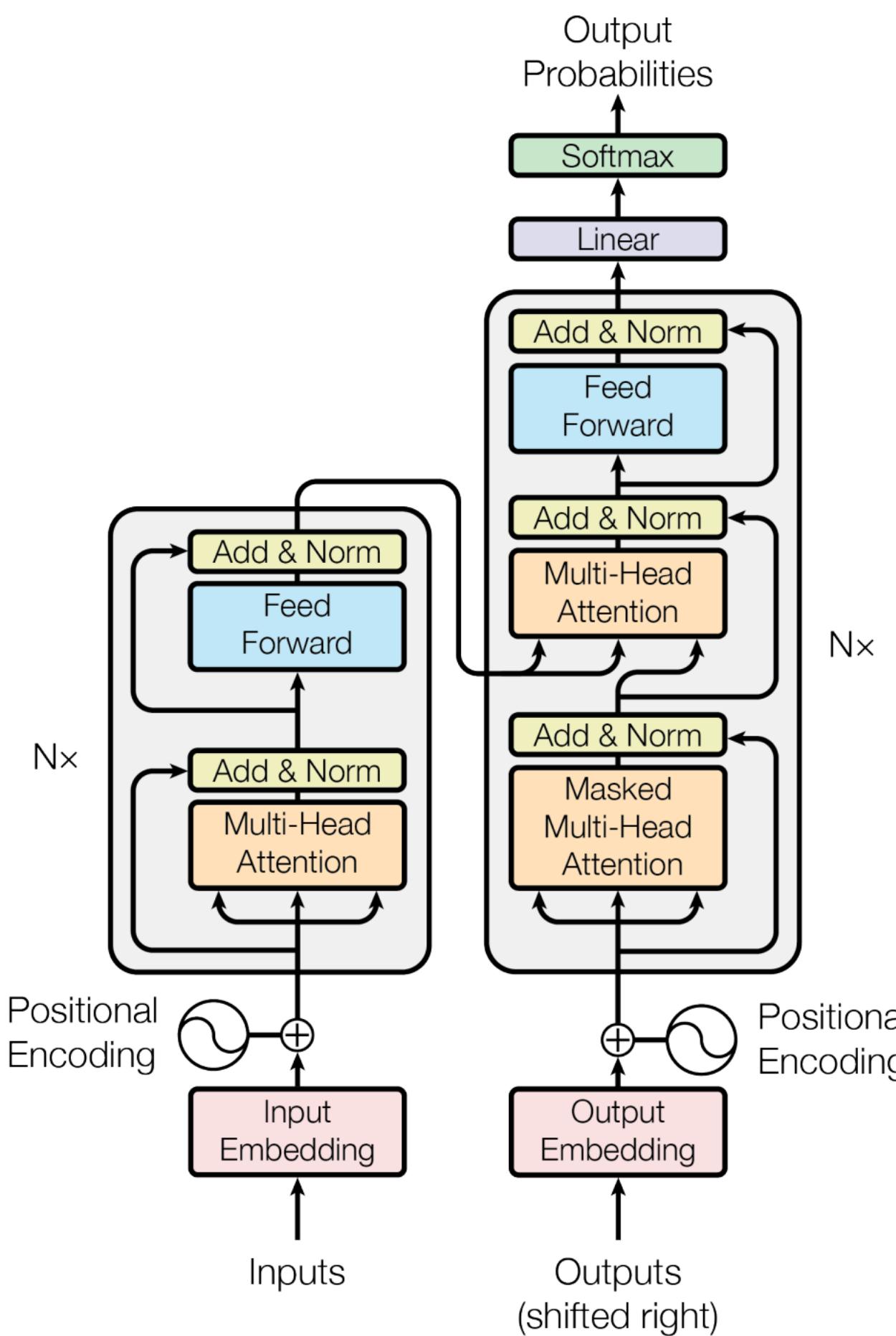


**Encoder-Only Transformer  
Architecture**

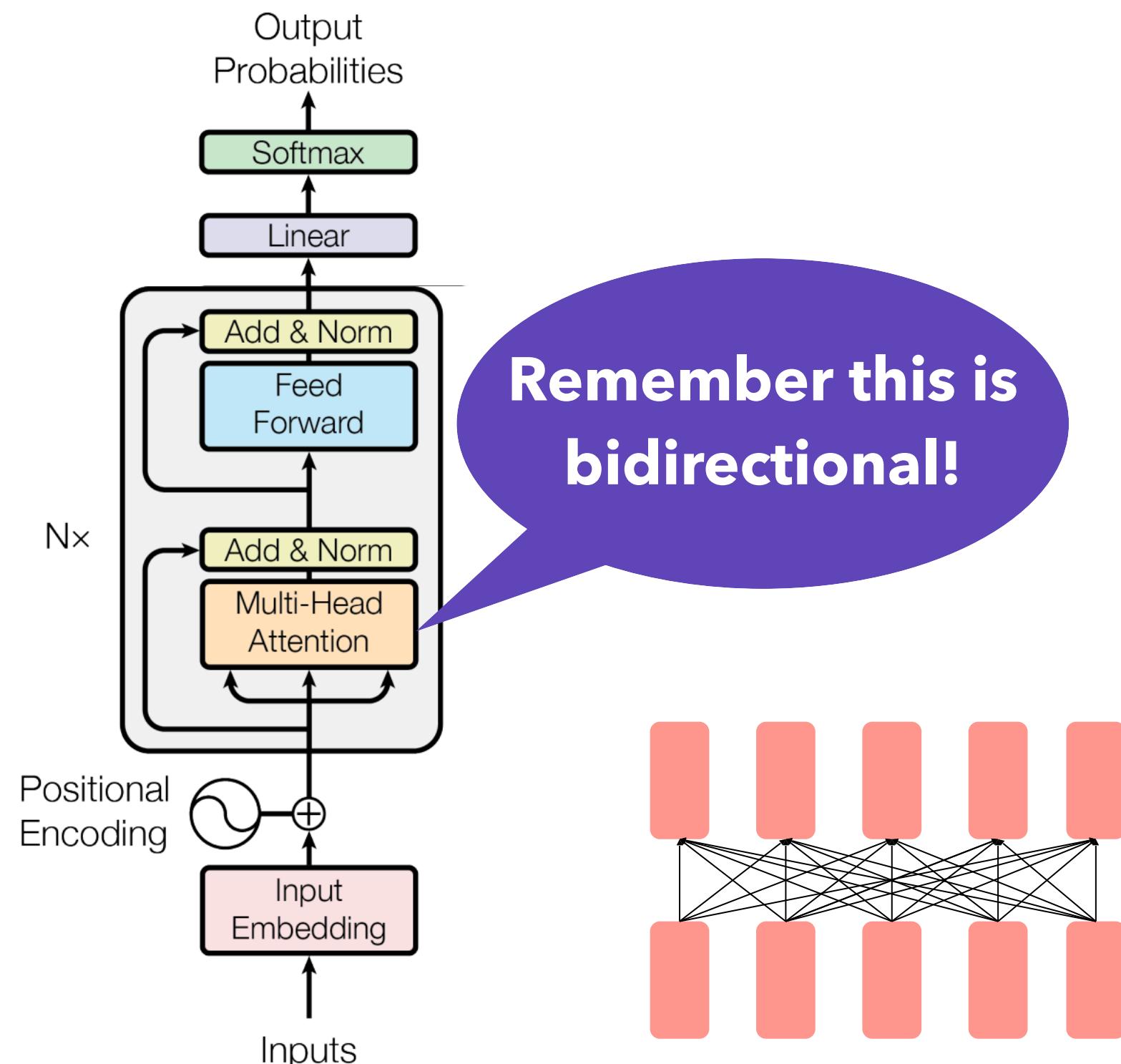


# Encoder: Architecture

**Full-Transformer Architecture  
(Encoder-Decoder)**



**Encoder-Only Transformer  
Architecture**

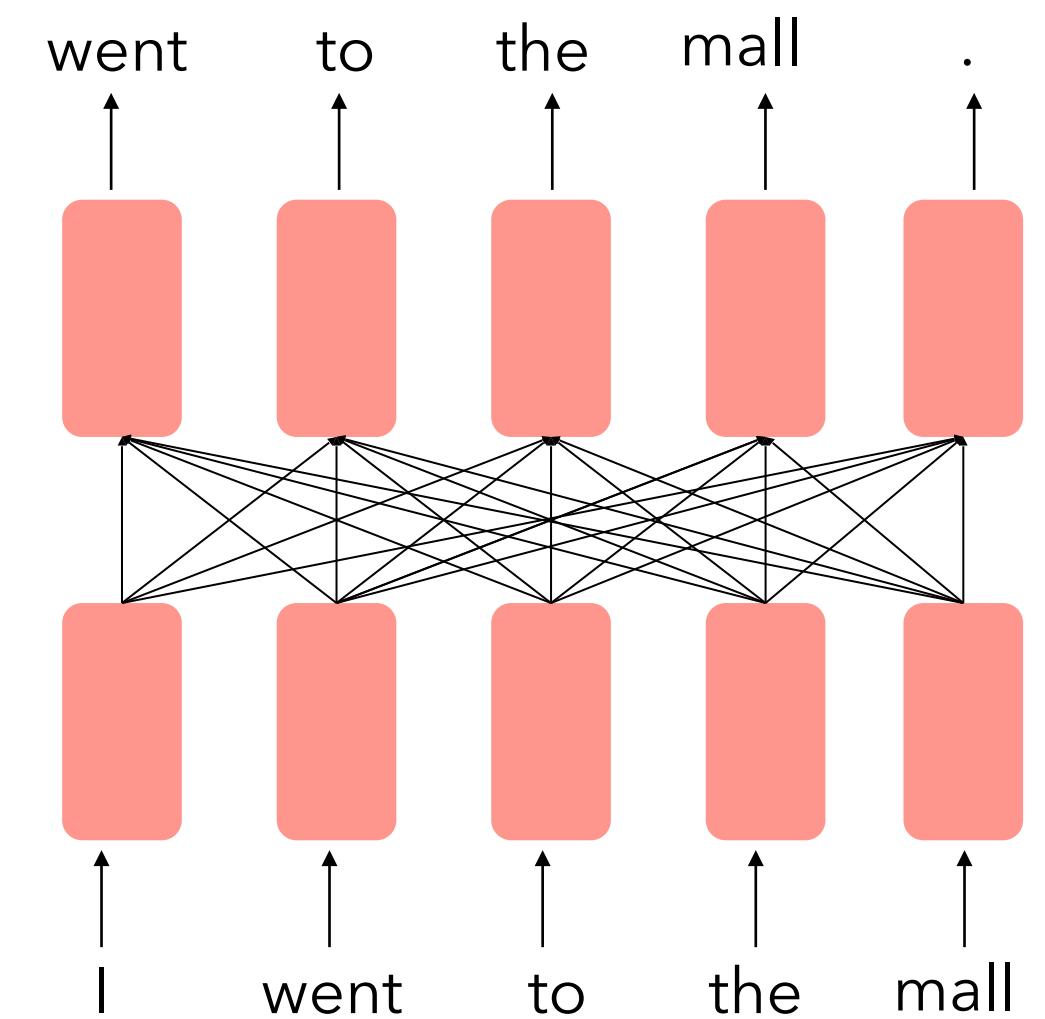


# Encoder: Training Objective

- So far, we've looked at language modeling for pre-training.
- Language Model Pretraining is problematic for encoders
- Why?
  - **Encoders get bidirectional contexts**
  - The model can cheat by just looking at the next token when predicting it without actually learning anything about language!

# Encoder: Training Objective

- So far, we've looked at language modeling for pre-training.
- Language Model Pretraining is problematic for encoders
- Why?
  - **Encoders get bidirectional contexts**
  - The model can cheat by just looking at the next token when predicting it without actually learning anything about language!



# Encoder: Training Objective

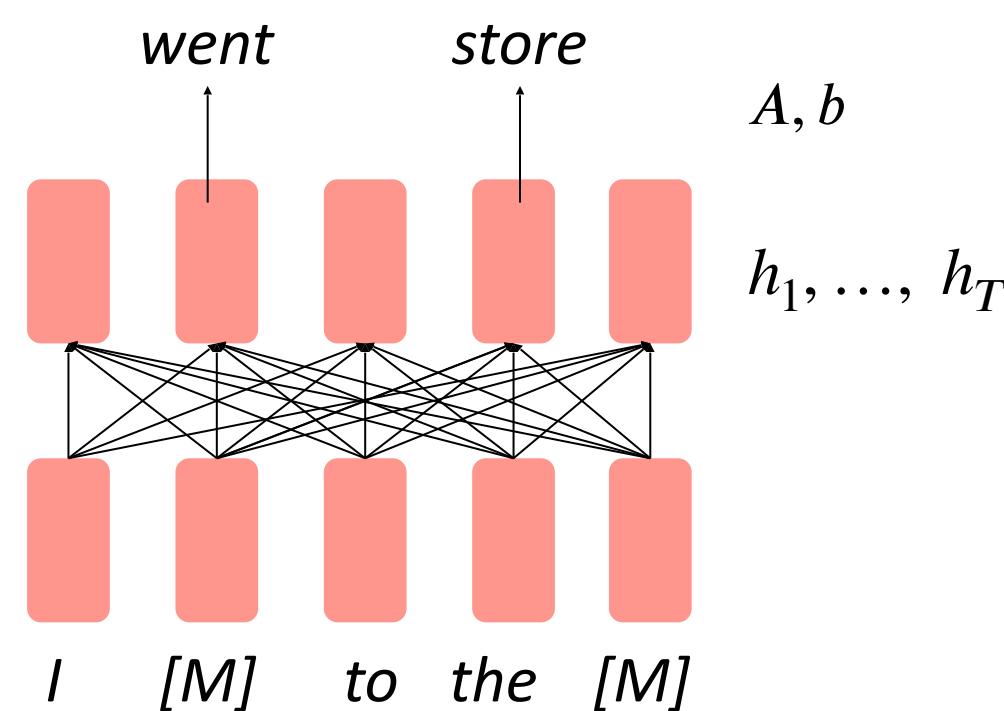
[Devlin et al., 2018]

- How to encode information from both **bidirectional** contexts?
- General Idea: **text reconstruction!**
  - Your time is [MASK], so don't [MASK] it living someone else's life.  
Don't be trapped by [MASK], which is [MASK] with the results of  
other [MASK]'s thinking. – [MASK] Jobs

# Encoder: Training Objective

[Devlin et al., 2018]

- How to encode information from both **bidirectional** contexts?
- General Idea: **text reconstruction!**
  - Your time is [MASK], so don't [MASK] it living someone else's life.  
Don't be trapped by [MASK], which is [MASK] with the results of other [MASK]'s thinking. – [MASK] Jobs



$$h_1, \dots, h_T = \text{Encoder}(w_1, \dots, w_T)$$

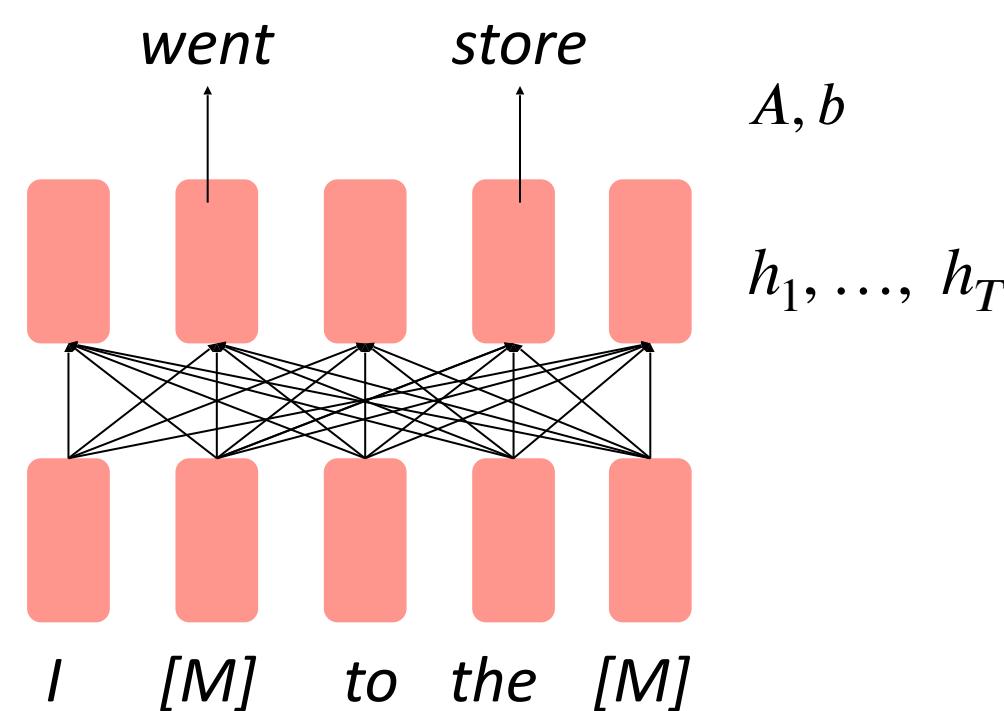
$$y_i \sim Aw_i + b$$

Only add loss terms from the masked tokens. If  $\tilde{x}$  is the masked version of  $x$ , we're learning  $p_\theta(x | \tilde{x})$ . Called **Masked Language model (MLM)**.

# Encoder: Training Objective

[Devlin et al., 2018]

- How to encode information from both **bidirectional** contexts?
- General Idea: **text reconstruction!**
  - Your time is limited so don't [MASK] it living someone else's life.  
Don't be trapped by [MASK], which is [MASK] with the results of other [MASK]'s thinking. – [MASK] Jobs



$$h_1, \dots, h_T = \text{Encoder}(w_1, \dots, w_T)$$

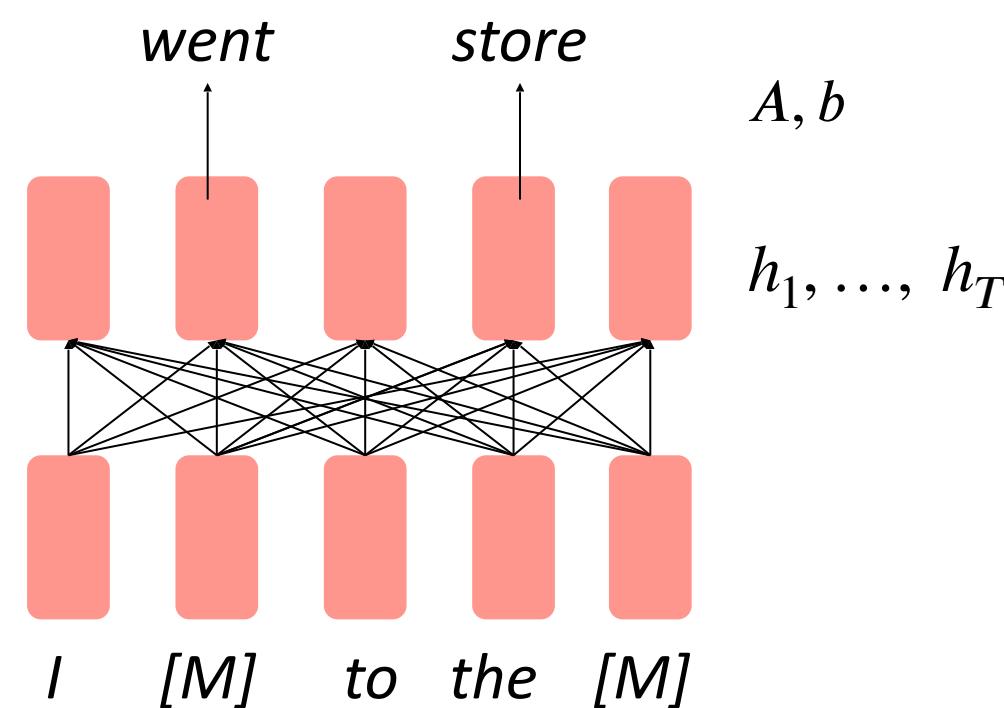
$$y_i \sim Aw_i + b$$

Only add loss terms from the masked tokens. If  $\tilde{x}$  is the masked version of  $x$ , we're learning  $p_\theta(x | \tilde{x})$ . Called **Masked Language model (MLM)**.

# Encoder: Training Objective

[Devlin et al., 2018]

- How to encode information from both **bidirectional** contexts?
- General Idea: **text reconstruction!**
  - Your time is limited so don't waste it living someone else's life.  
Don't be trapped by [MASK], which is [MASK] with the results of other [MASK]'s thinking. – [MASK] Jobs

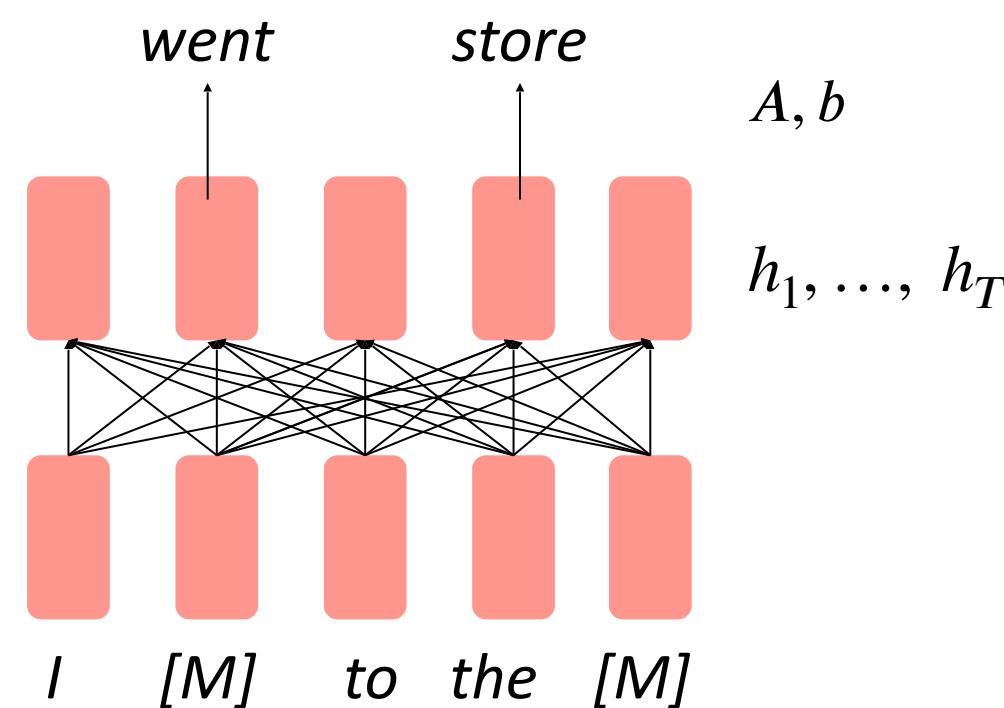


Only add loss terms from the masked tokens. If  $\tilde{x}$  is the masked version of  $x$ , we're learning  $p_\theta(x | \tilde{x})$ . Called **Masked Language model (MLM)**.

# Encoder: Training Objective

[Devlin et al., 2018]

- How to encode information from both **bidirectional** contexts?
- General Idea: **text reconstruction!**
  - Your time is limited so don't waste it living someone else's life.  
Don't be trapped by dogma which is [MASK] with the results of other [MASK]'s thinking. – [MASK] Jobs



$$h_1, \dots, h_T = \text{Encoder}(w_1, \dots, w_T)$$

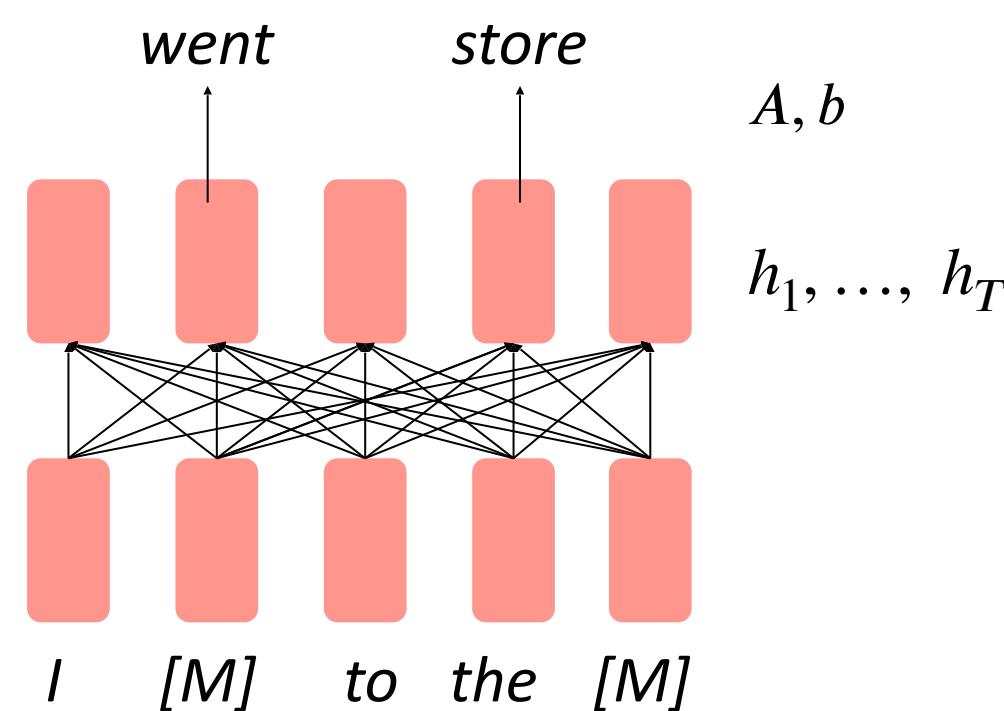
$$y_i \sim Aw_i + b$$

Only add loss terms from the masked tokens. If  $\tilde{x}$  is the masked version of  $x$ , we're learning  $p_\theta(x | \tilde{x})$ . Called **Masked Language model (MLM)**.

# Encoder: Training Objective

[Devlin et al., 2018]

- How to encode information from both **bidirectional** contexts?
- General Idea: **text reconstruction!**
  - Your time is limited so don't waste it living someone else's life.  
Don't be trapped by dogma which is living with the results of other [MASK]'s thinking. – [MASK] Jobs



$$h_1, \dots, h_T = \text{Encoder}(w_1, \dots, w_T)$$

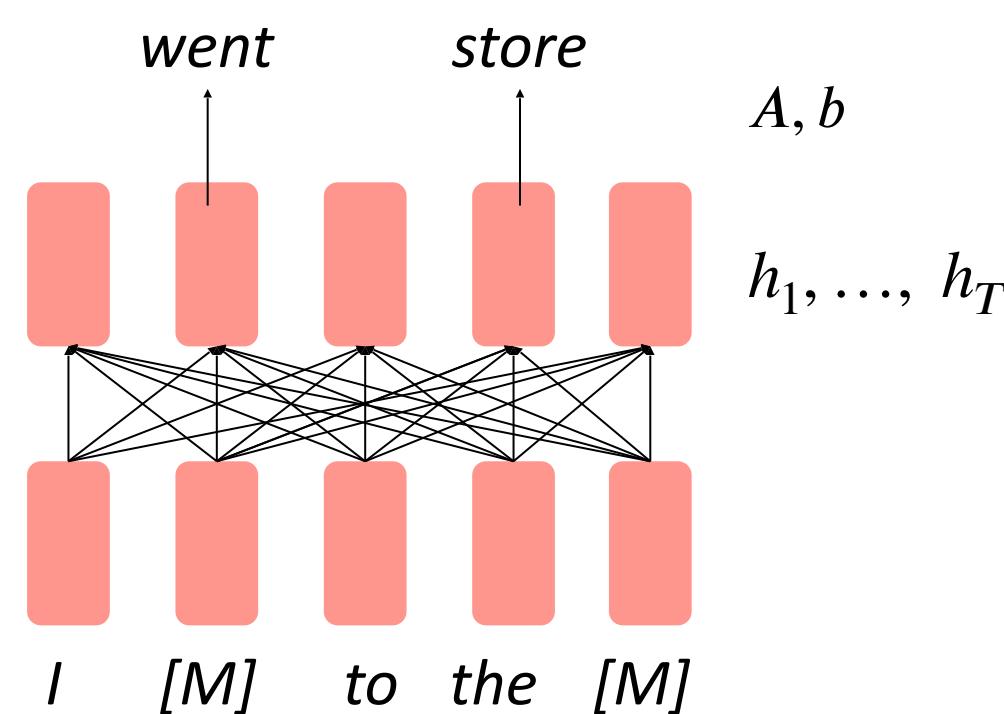
$$y_i \sim Aw_i + b$$

Only add loss terms from the masked tokens. If  $\tilde{x}$  is the masked version of  $x$ , we're learning  $p_\theta(x | \tilde{x})$ . Called **Masked Language model (MLM)**.

# Encoder: Training Objective

[Devlin et al., 2018]

- How to encode information from both **bidirectional** contexts?
- General Idea: **text reconstruction!**
  - Your time is limited so don't waste it living someone else's life.  
Don't be trapped by dogma which is living with the results of other people's thinking. – [MASK] Jobs



$$h_1, \dots, h_T = \text{Encoder}(w_1, \dots, w_T)$$

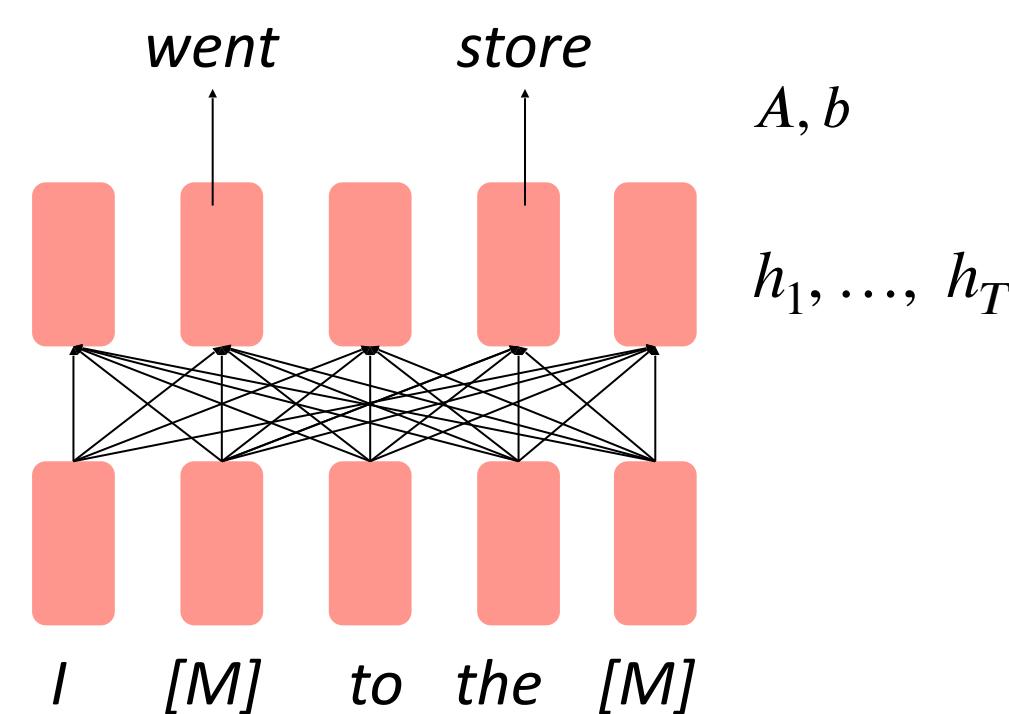
$$y_i \sim Aw_i + b$$

Only add loss terms from the masked tokens. If  $\tilde{x}$  is the masked version of  $x$ , we're learning  $p_\theta(x | \tilde{x})$ . Called **Masked Language model (MLM)**.

# Encoder: Training Objective

[Devlin et al., 2018]

- How to encode information from both **bidirectional** contexts?
- General Idea: **text reconstruction!**
  - Your time is limited so don't waste it living someone else's life.  
Don't be trapped by dogma which is living with the results of other people's thinking. – Steve Jobs



$$h_1, \dots, h_T = \text{Encoder}(w_1, \dots, w_T)$$

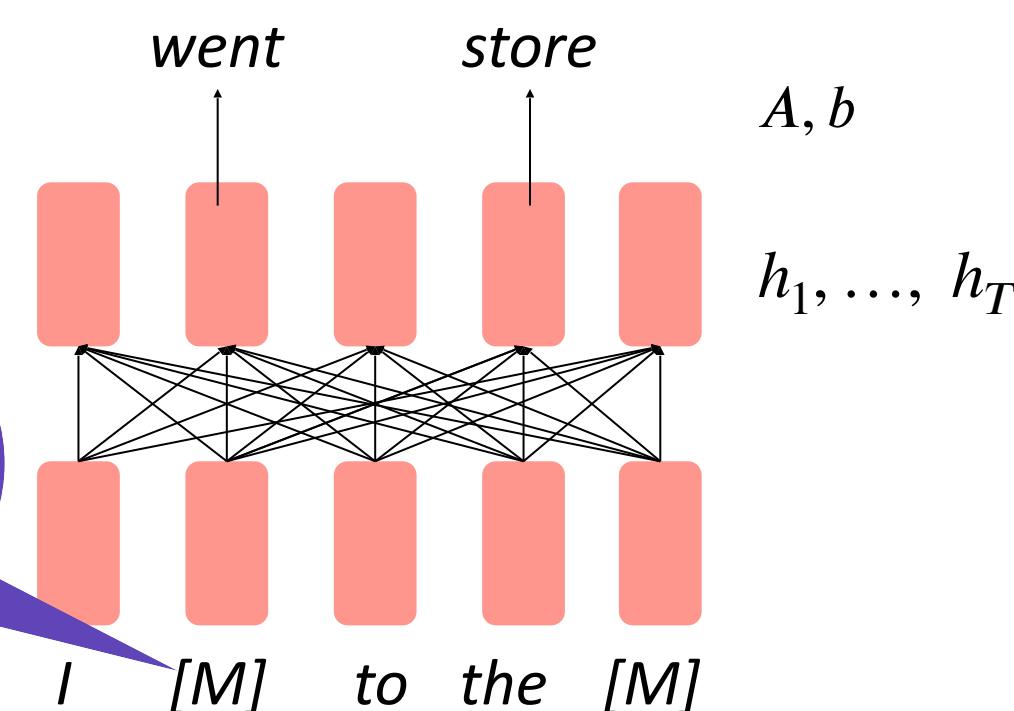
$$y_i \sim Aw_i + b$$

Only add loss terms from the masked tokens. If  $\tilde{x}$  is the masked version of  $x$ , we're learning  $p_\theta(x | \tilde{x})$ . Called **Masked Language model (MLM)**.

# Encoder: Training Objective

[Devlin et al., 2018]

- How to encode information from both **bidirectional** contexts?
- General Idea: **text reconstruction!**
  - Your time is limited so don't waste it living someone else's life.  
Don't be trapped by dogma which is living with the results of other people's thinking. – Steve Jobs



$$h_1, \dots, h_T = \text{Encoder}(w_1, \dots, w_T)$$

$$y_i \sim Aw_i + b$$

Only add loss terms from the masked tokens. If  $\tilde{x}$  is the masked version of  $x$ , we're learning  $p_\theta(x | \tilde{x})$ . Called **Masked Language model (MLM)**.

# Encoder: BERT

Bidirectional **E**ncoder  
Representations from **T**ransformers

[Devlin et al., 2018]

- **2 Pre-training Objectives:**

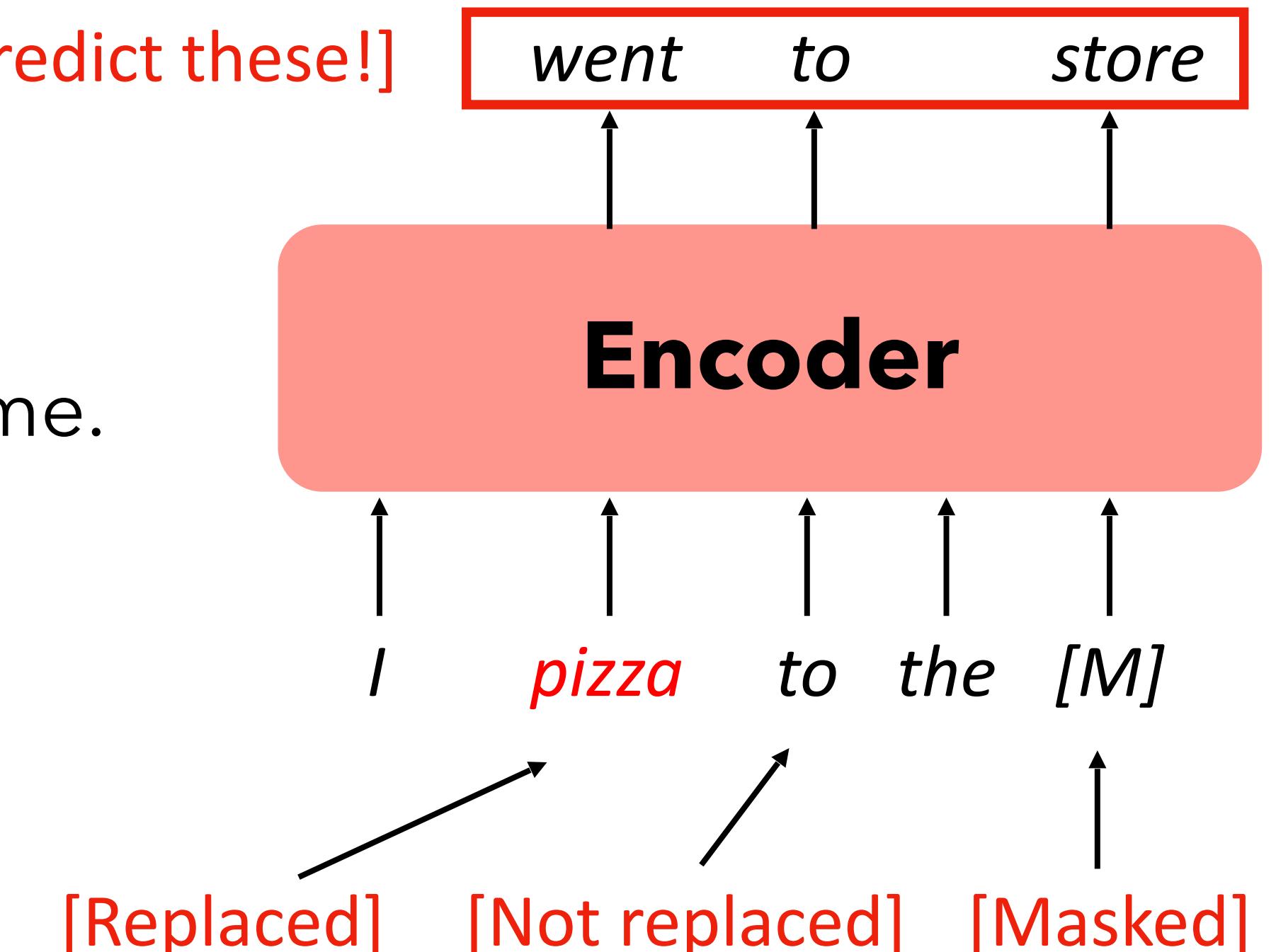
- **Masked LM: Choose a random 15% of tokens to predict.**

- For each chosen token:
  - Replace it with **[MASK]** 80% of the time.
  - Replace it with a **random token** 10% of the time.
  - Leave it **unchanged** 10% of the time (but still predict it!).

- **Next Sentence Prediction (NSP)**

- 50% of the time two adjacent sentences are in the correct order.

- **This actually hurts model learning based on later work!**



# Encoder: BERT

Bidirectional Encoder  
Representations from Transformers

[Devlin et al., 2018]

- **2 Pre-training Objectives:**

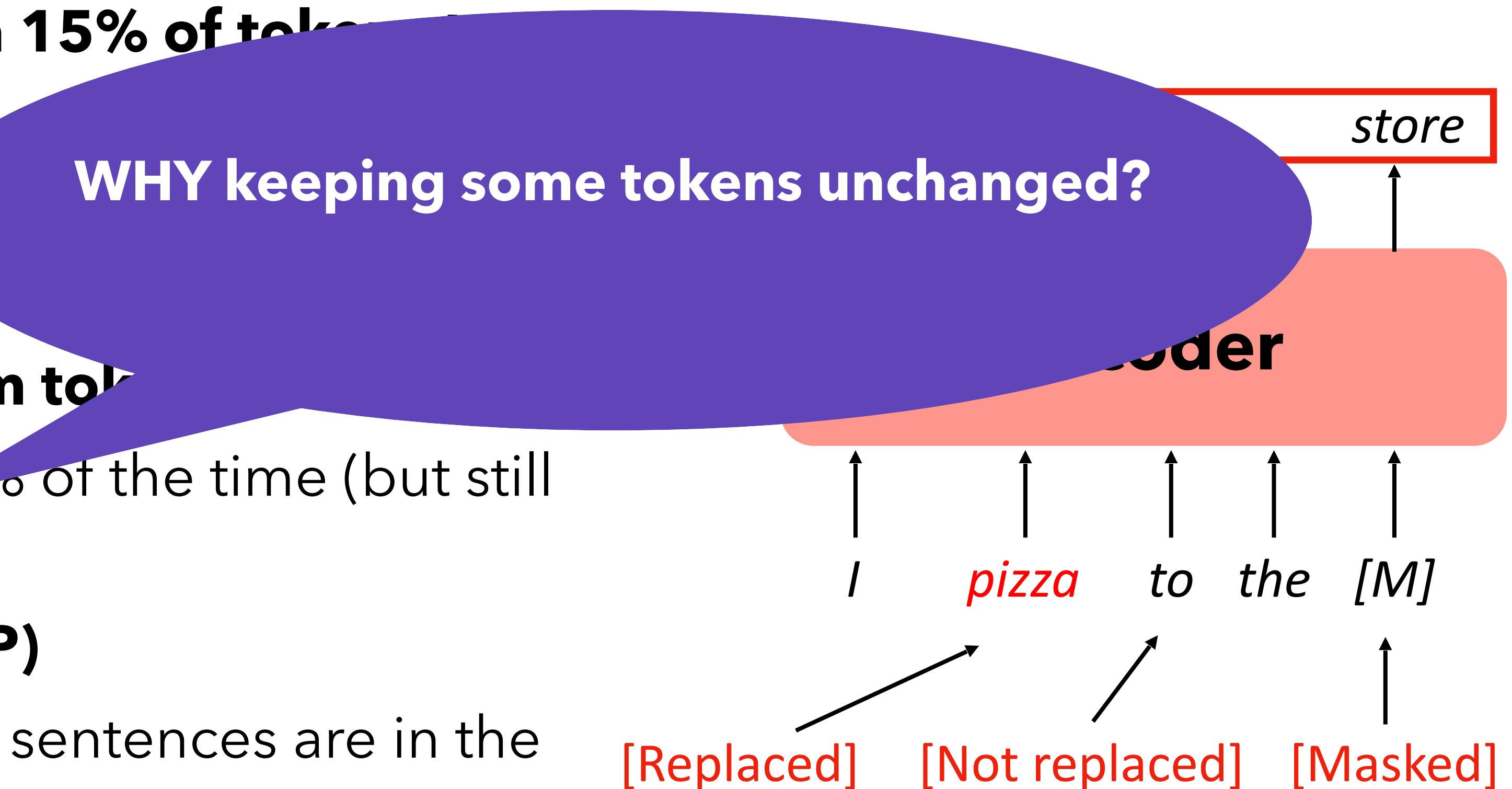
- **Masked LM: Choose a random 15% of tokens to predict.**

- For each chosen token:
  - Replace it with **[MASK]**
  - Replace it with a **random token**
  - Leave it **unchanged** 10% of the time (but still predict it!).

- **Next Sentence Prediction (NSP)**

- 50% of the time two adjacent sentences are in the correct order.

- **This actually hurts model learning based on later work!**



# Encoder: BERT

Bidirectional Encoder  
Representations from Transformers

[Devlin et al., 2018]

- **2 Pre-training Objectives:**

- **Masked LM: Choose a random 15% of tokens to predict.**

- For each chosen token:
  - Replace it with **[MASK]**
  - Replace it with a **random token**
  - Leave it **unchanged** 10% of the time (but still predict it!).

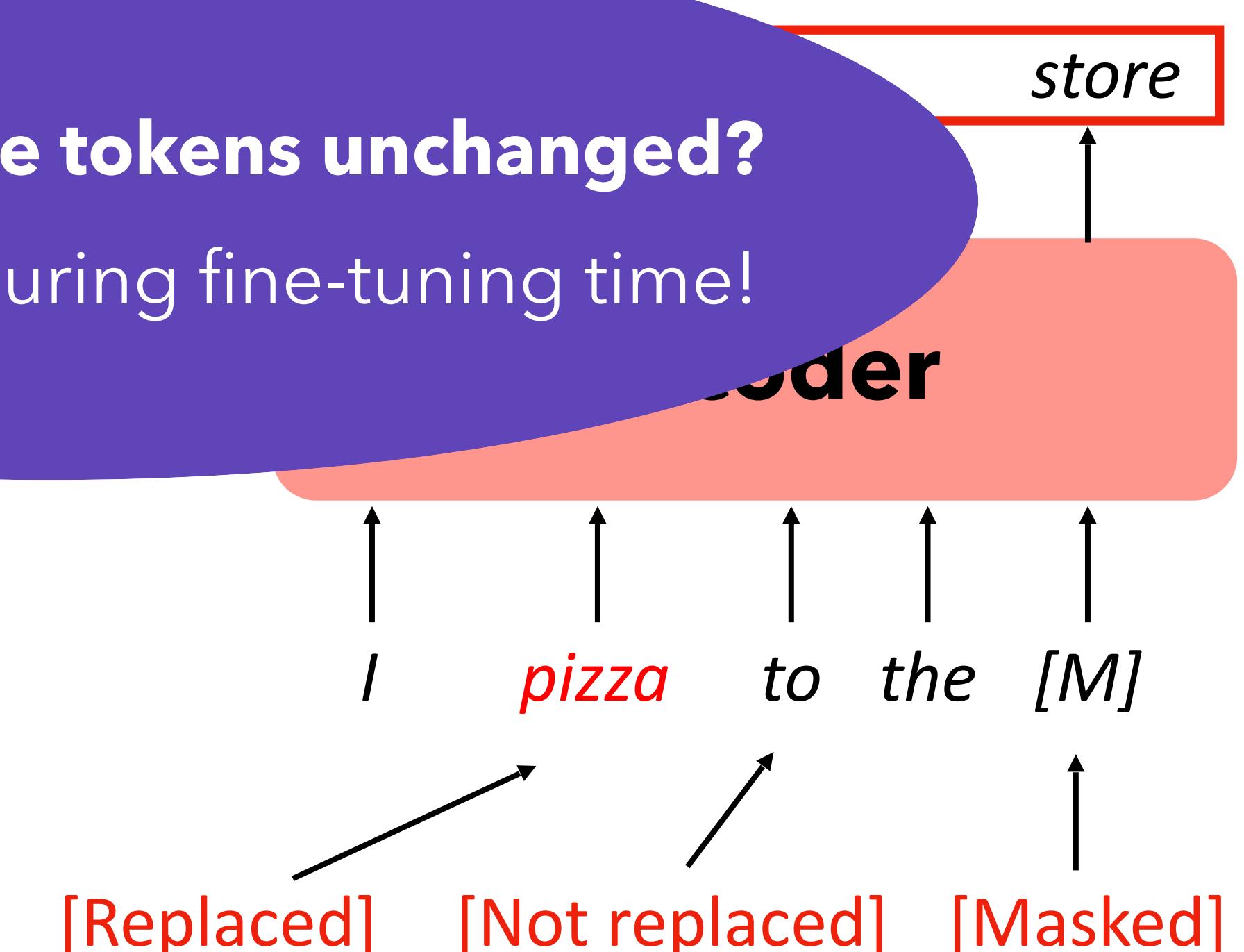
- **Next Sentence Prediction (NSP)**

- 50% of the time two adjacent sentences are in the correct order.

- **This actually hurts model learning based on later work!**

**WHY keeping some tokens unchanged?**

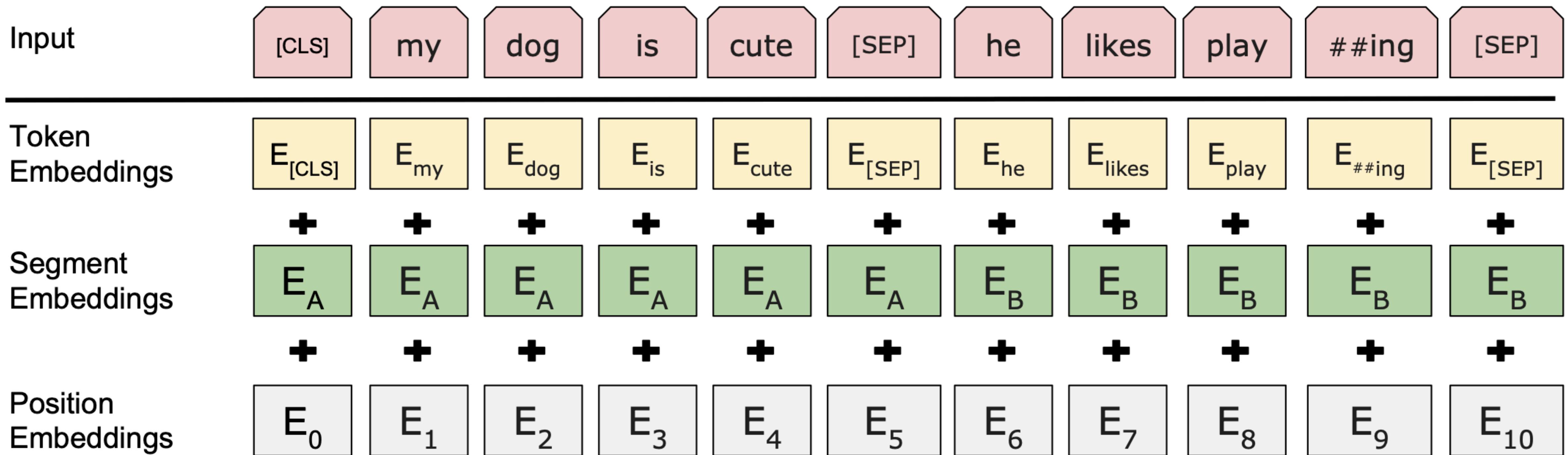
There's no [MASK] during fine-tuning time!



# Encoder: BERT

**Bidirectional Encoder Representations from Transformers**

[Devlin et al., 2018]

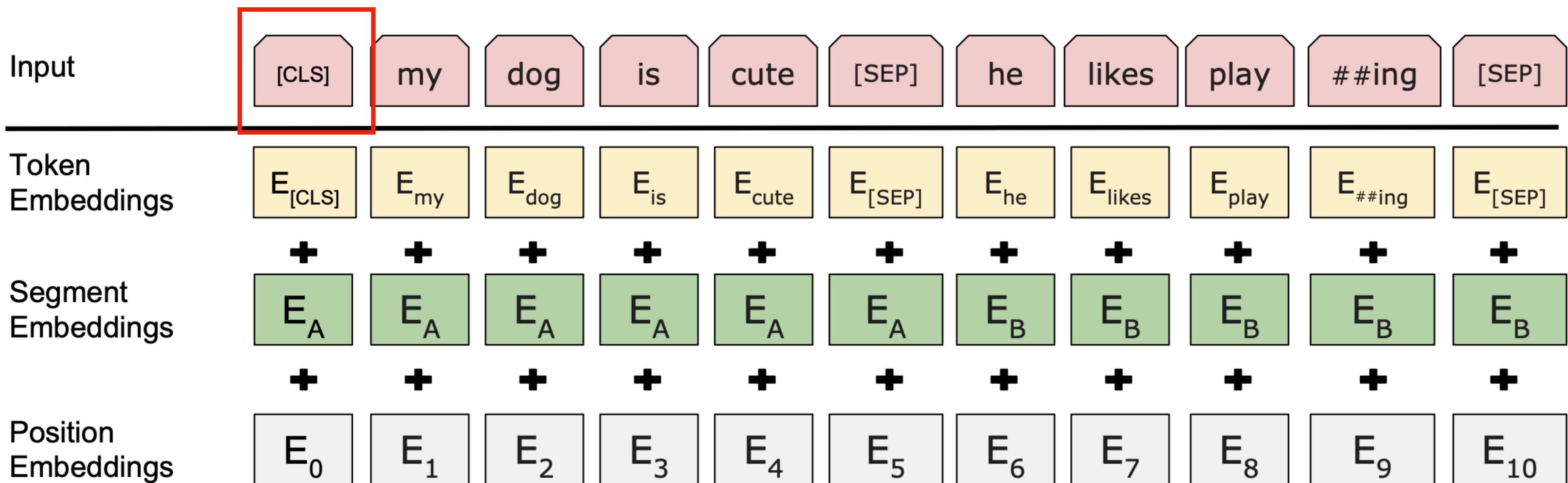


# Encoder: BERT

**Bidirectional Encoder Representations from Transformers**

[Devlin et al., 2018]

Special token added to the beginning of each input sequence



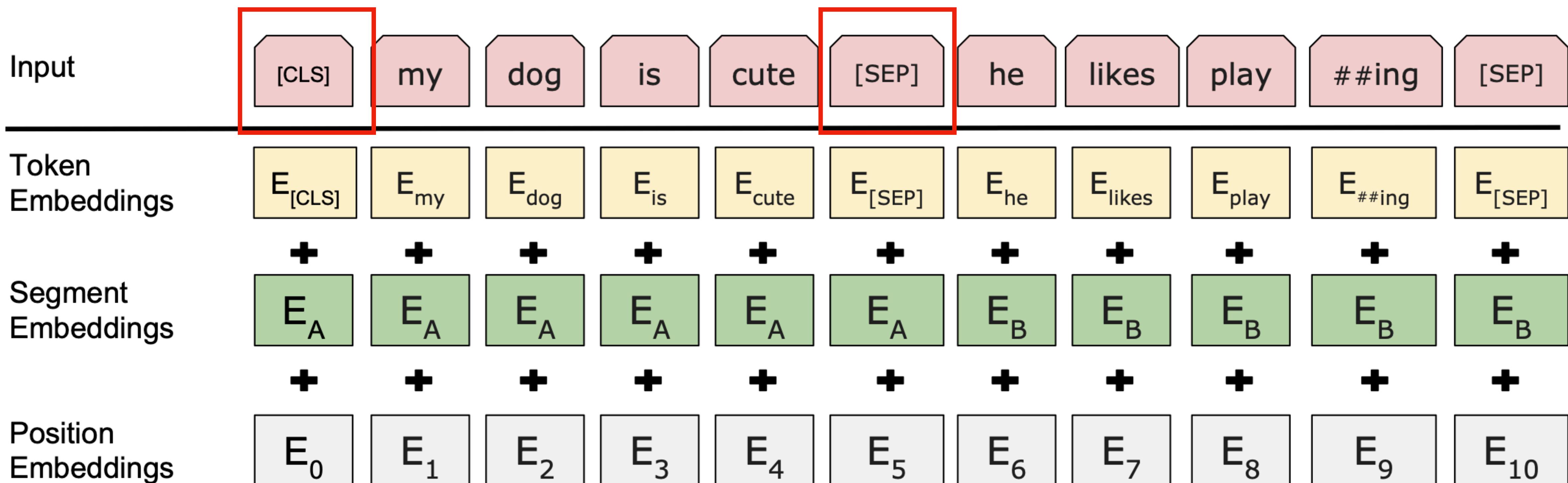
# Encoder: BERT

**Bidirectional Encoder Representations from Transformers**

[Devlin et al., 2018]

Special token added to the beginning of each input sequence

Special token to separate sentence A/B



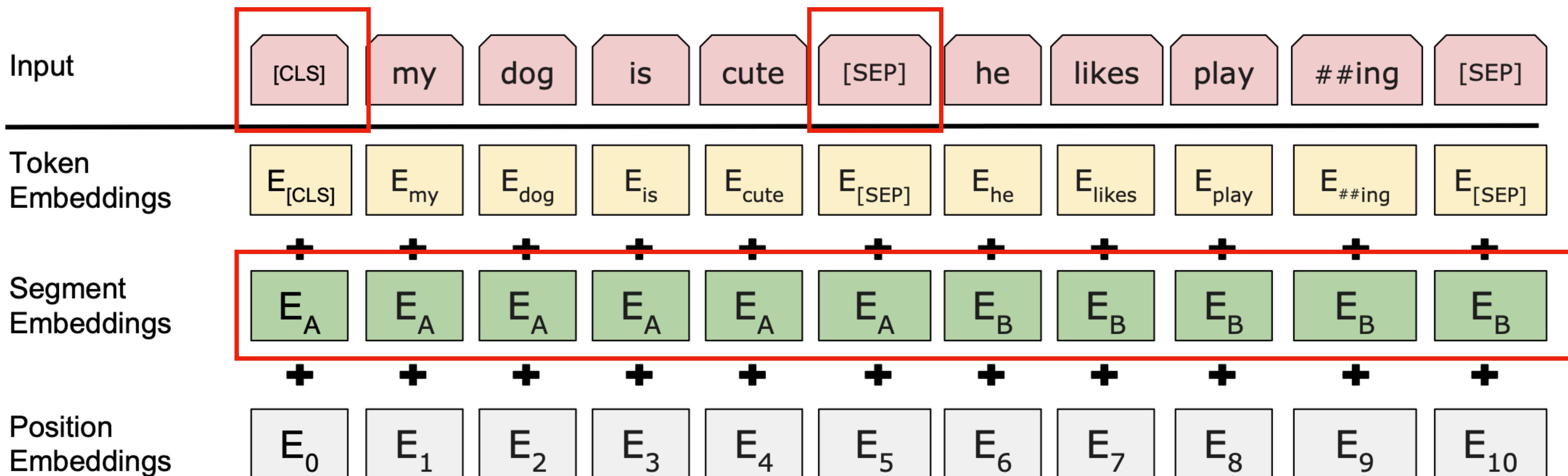
# Encoder: BERT

**Bidirectional Encoder Representations from Transformers**

[Devlin et al., 2018]

Special token added to the beginning of each input sequence

Special token to separate sentence A/B



Learned embedding to every token indicating whether it belongs to sentence A or sentence B

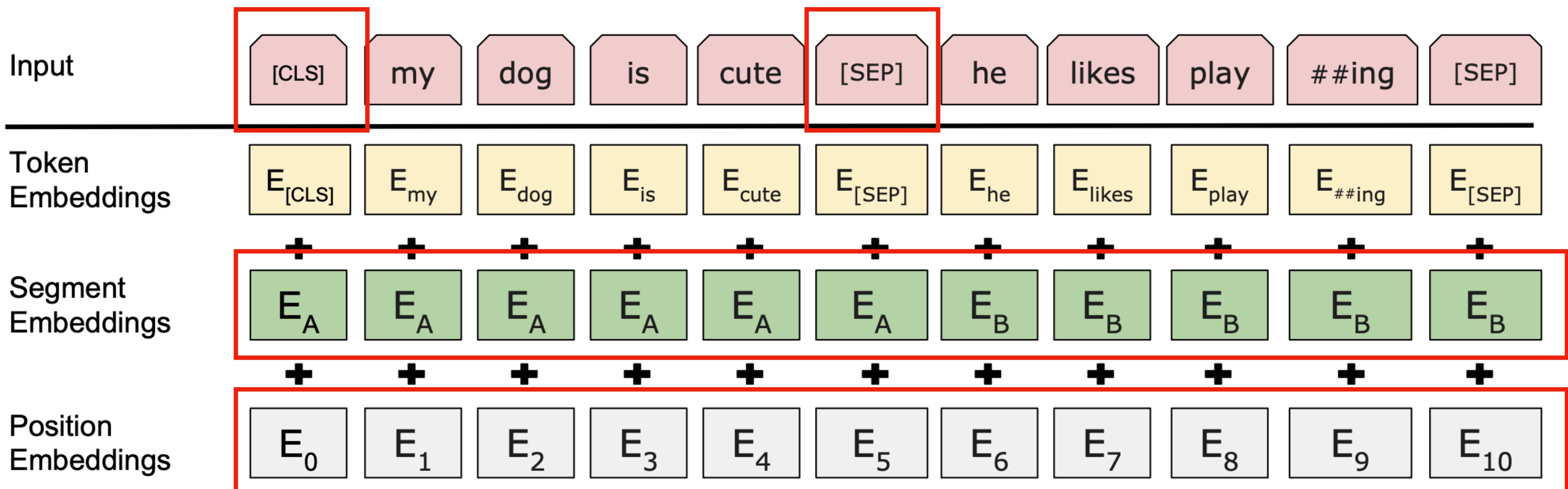
# Encoder: BERT

**Bidirectional Encoder  
Representations from Transformers**

[Devlin et al., 2018]

Special token added to the beginning of each input sequence

Special token to separate sentence A/B



Learned embedding to every token indicating whether it belongs to sentence A or sentence B

Position of the token in the entire sequence

# Encoder: BERT

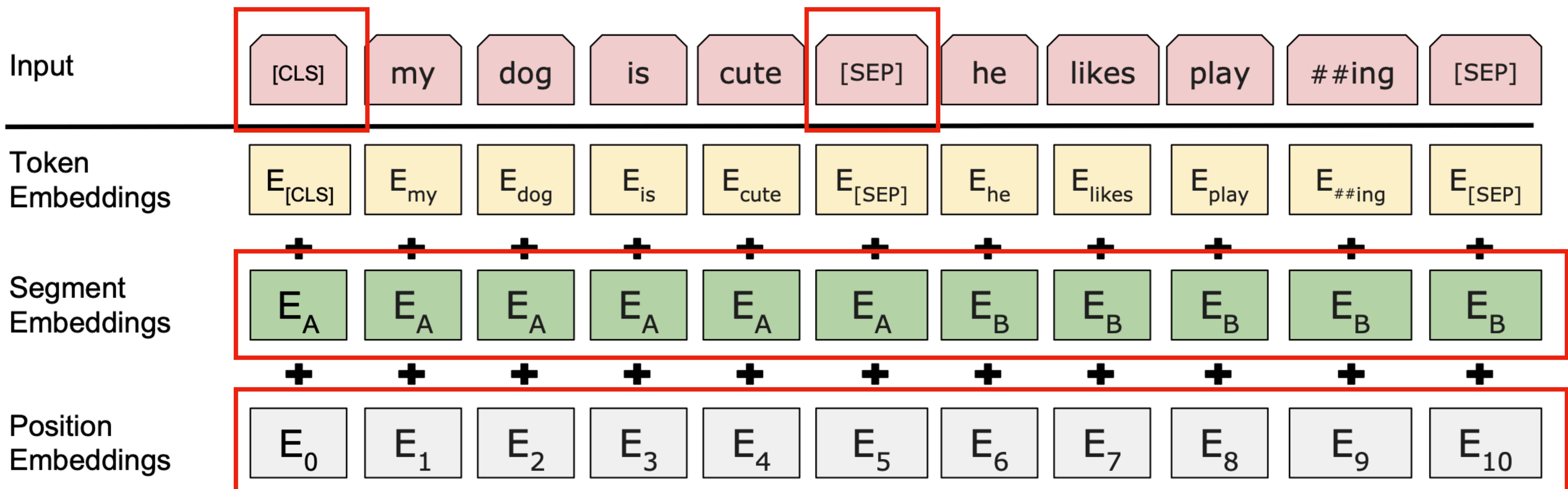
**Bidirectional Encoder Representations from Transformers**

[Devlin et al., 2018]

Special token added to the beginning of each input sequence

Special token to separate sentence A/B

Final embedding is the sum of all three!



Learned embedding to every token indicating whether it belongs to sentence A or sentence B

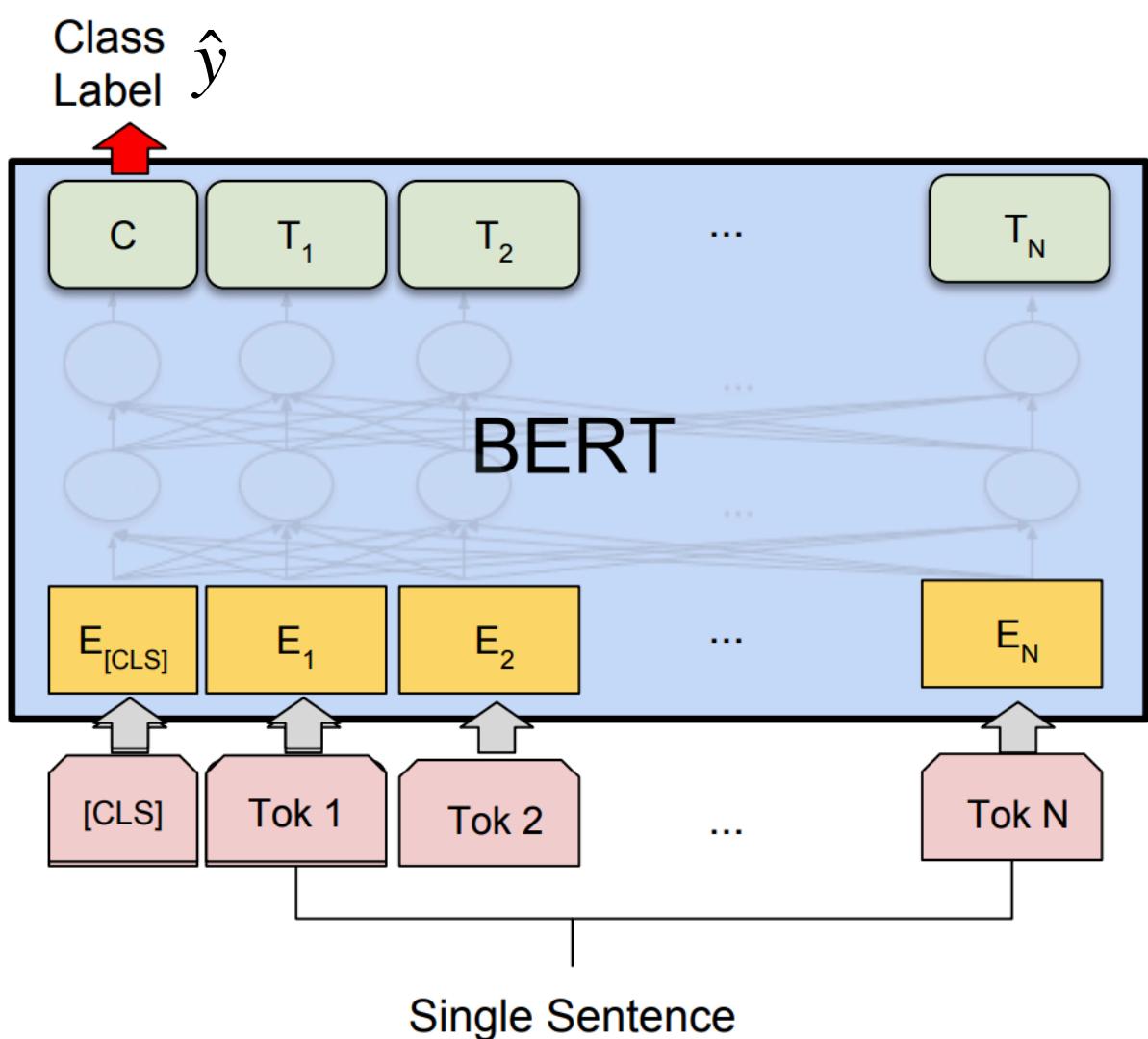
Position of the token in the entire sequence

# Encoder: BERT (Fine-tuning)

# Encoder: BERT (Fine-tuning)

Single-Sentence Tasks like  
SST-2 (Sentiment Analysis)

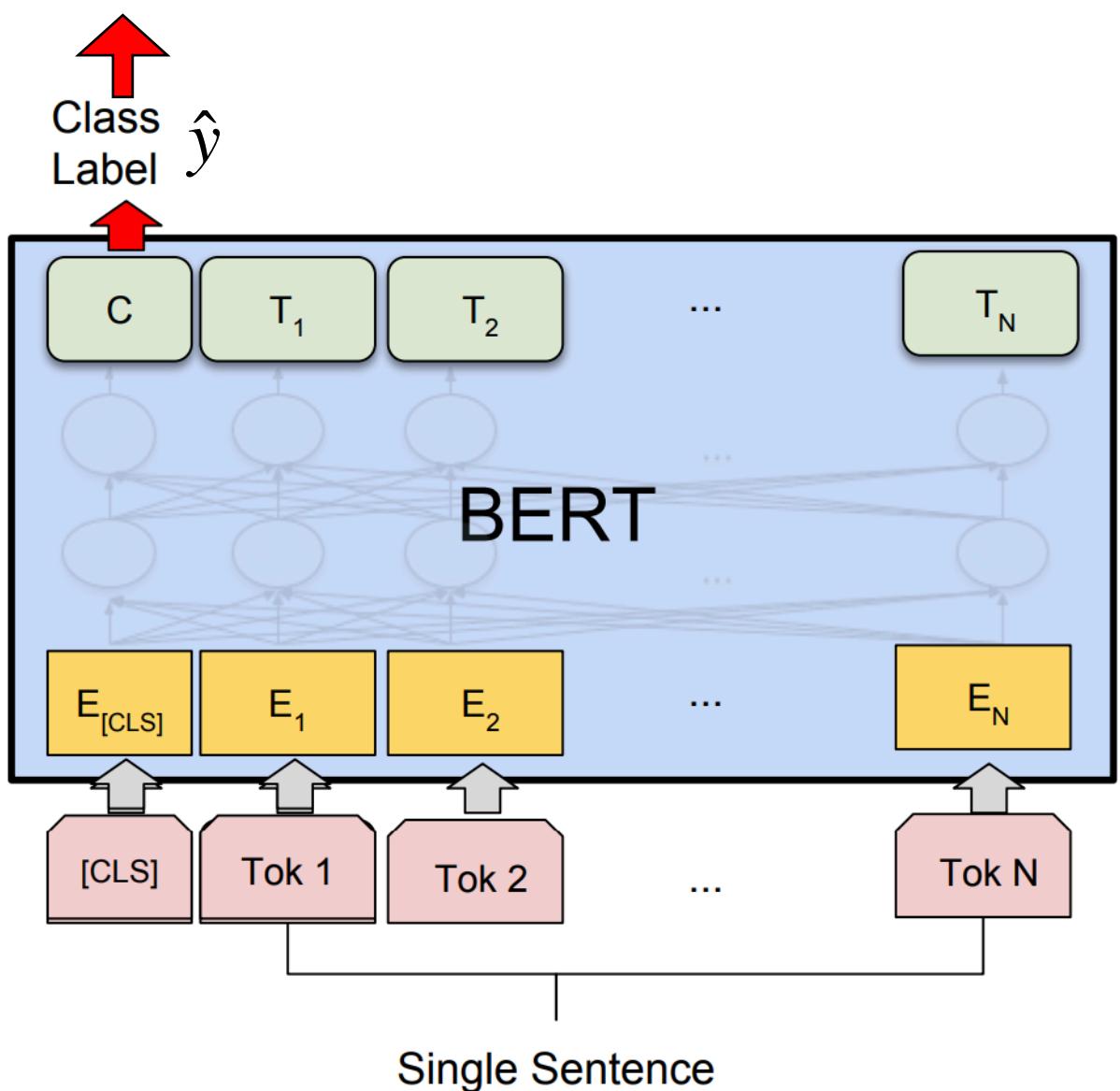
# Encoder: BERT (Fine-tuning)



Single-Sentence Tasks like  
SST-2 (Sentiment Analysis)

# Encoder: BERT (Fine-tuning)

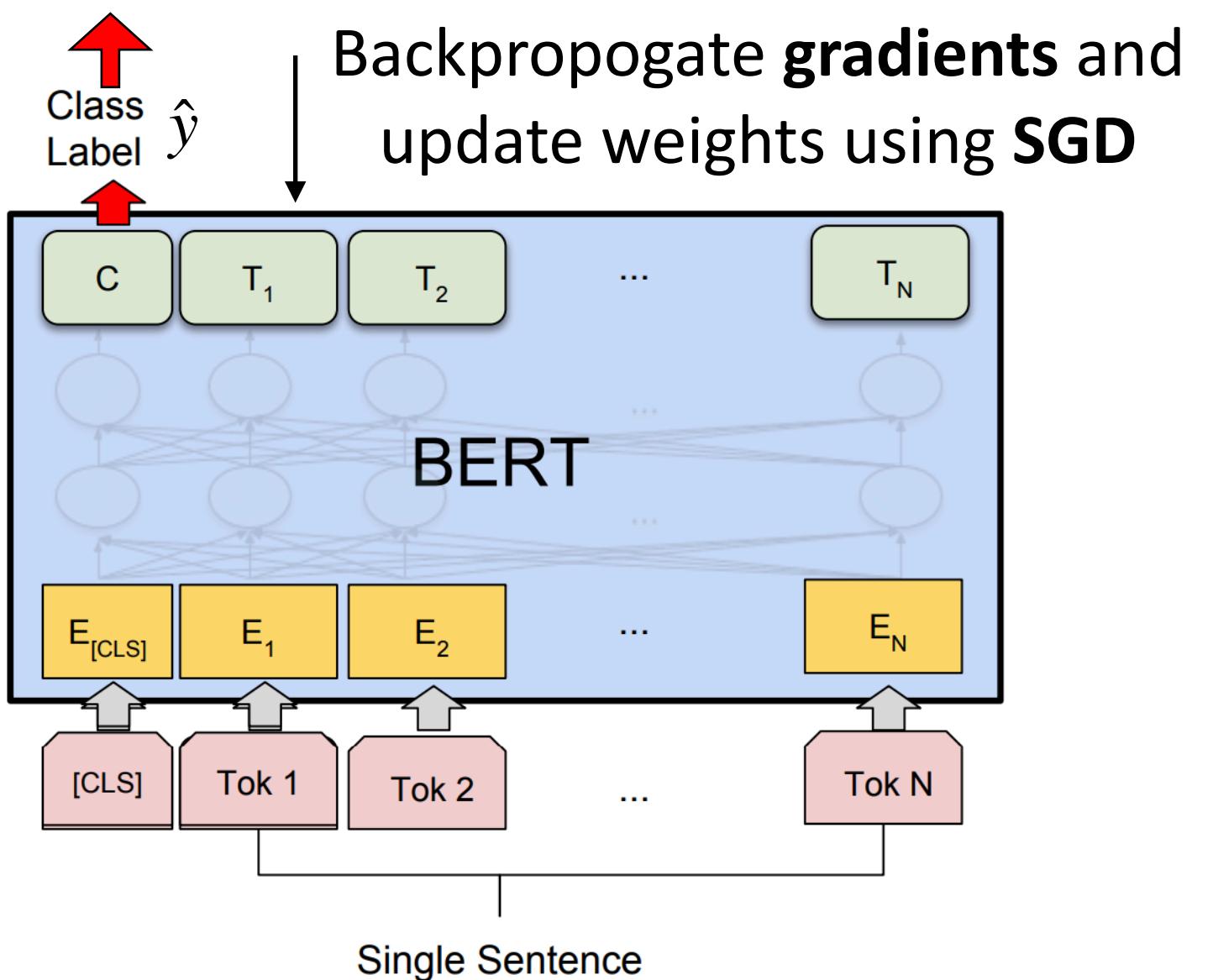
Cross-Entropy Loss  $L(\hat{y}, y)$



Single-Sentence Tasks like  
SST-2 (Sentiment Analysis)

# Encoder: BERT (Fine-tuning)

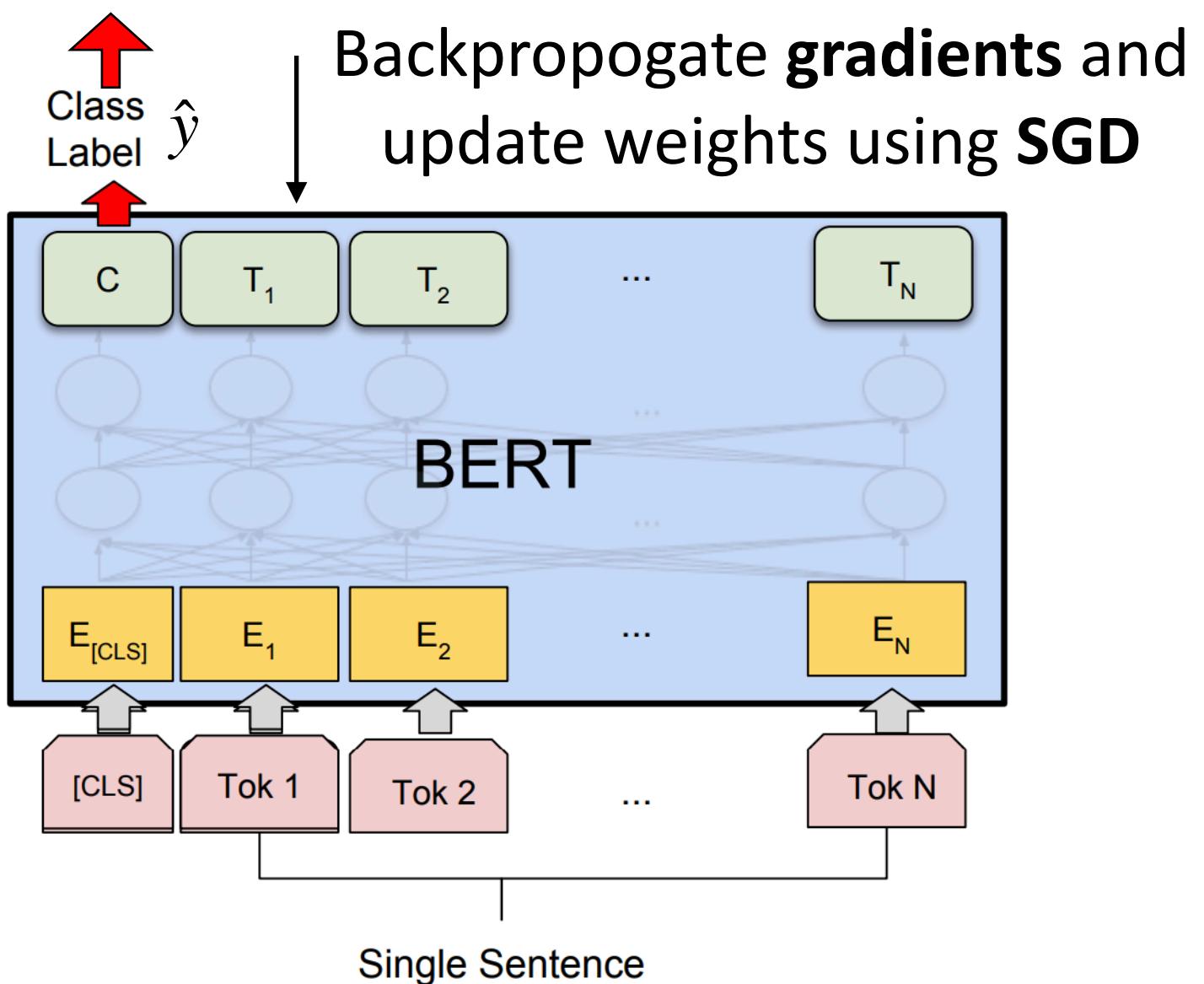
Cross-Entropy Loss  $L(\hat{y}, y)$



Single-Sentence Tasks like  
SST-2 (Sentiment Analysis)

# Encoder: BERT (Fine-tuning)

Cross-Entropy Loss  $L(\hat{y}, y)$

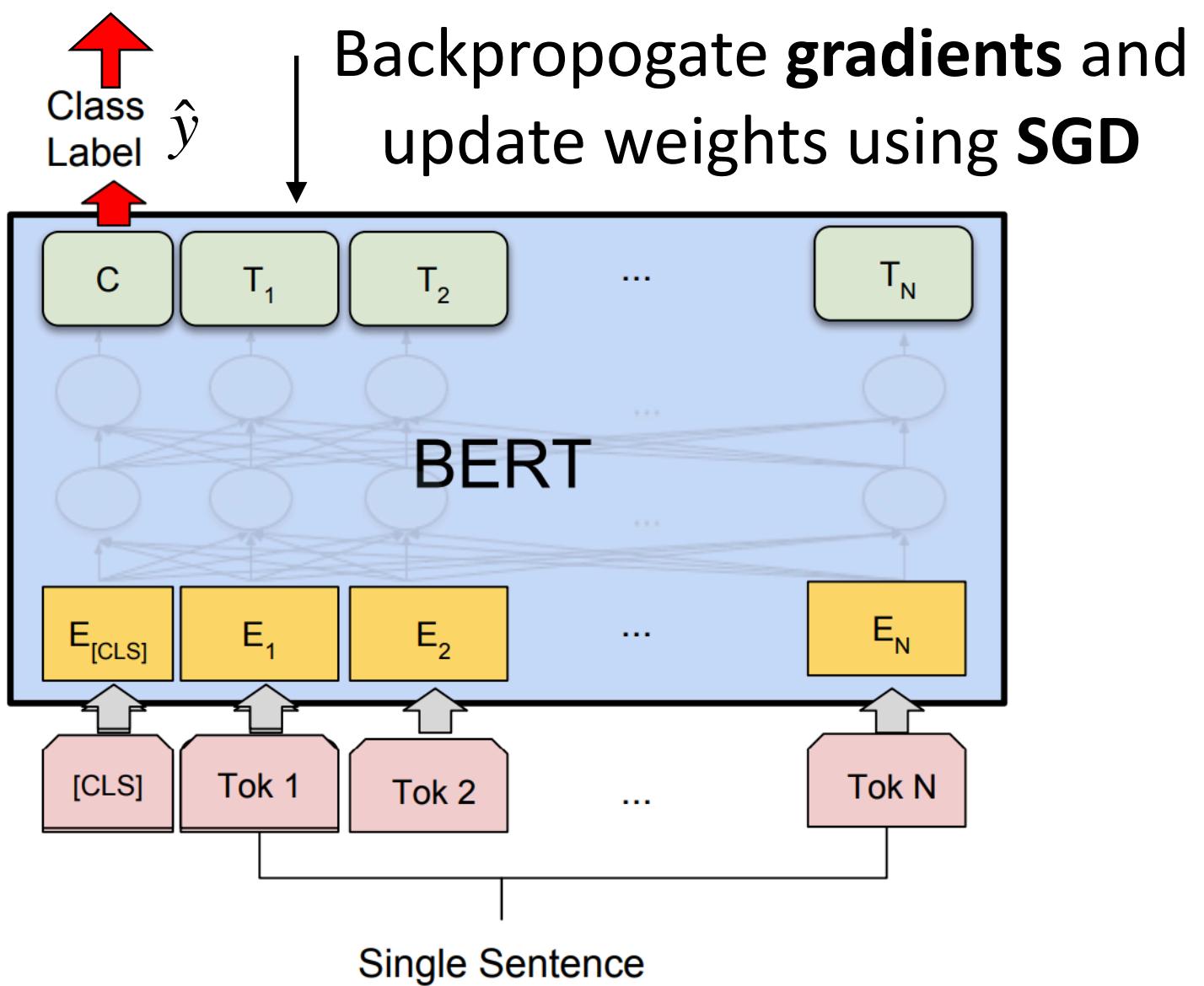


Single-Sentence Tasks like  
SST-2 (Sentiment Analysis)

Sentence Pair Classification  
Tasks like Natural Language  
Inference

# Encoder: BERT (Fine-tuning)

**Cross-Entropy Loss**  $L(\hat{y}, y)$



Single-Sentence Tasks like  
SST-2 (Sentiment Analysis)

Input:

Premise: A soccer game with multiple males playing

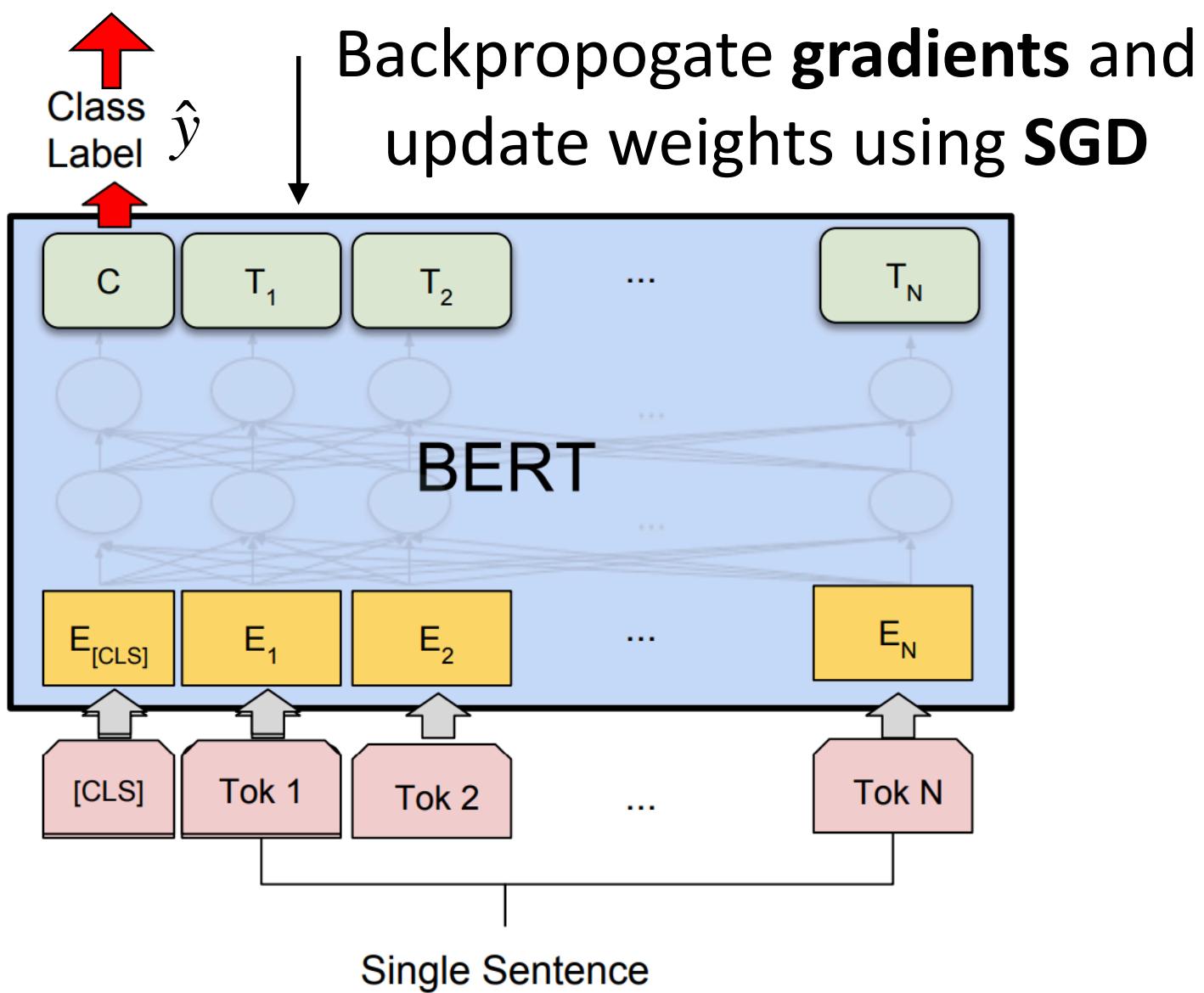
Hypothesis: Some men are playing a sport

Label: Entailment / Neutral / Contadiction

Sentence Pair Classification  
Tasks like Natural Language  
Inference

# Encoder: BERT (Fine-tuning)

**Cross-Entropy Loss**  $L(\hat{y}, y)$



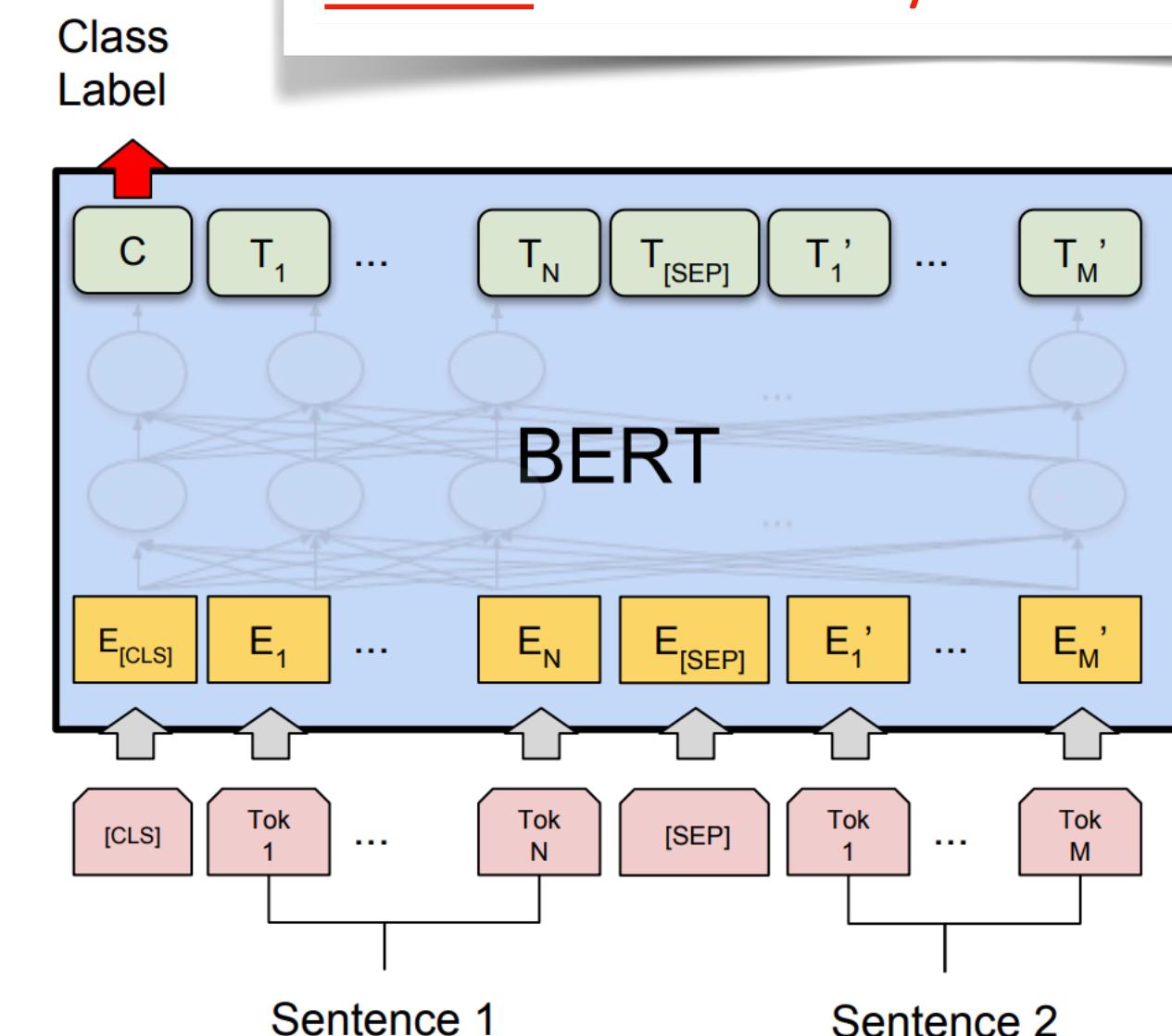
Single-Sentence Tasks like  
SST-2 (Sentiment Analysis)

Input:

Premise: A soccer game with multiple males playing

Hypothesis: Some men are playing a sport

Label: Entailment / Neutral / Contadiction



Sentence Pair Classification  
Tasks like Natural Language  
Inference

# Encoder: BERT

**Bidirectional Encoder Representations from Transformers**

[Devlin et al., 2018]

- **SOTA at the time on a wide range of tasks after fine-tuning!**

System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average -
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT <sub>BASE</sub>	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>92.7</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>82.1</b>

- **QQP:** Quora Question Pairs (detect paraphrase questions)
- **QNLI:** natural language inference over question answering data
- **SST-2:** sentiment analysis
- **CoLA:** corpus of linguistic acceptability (detect whether sentences are grammatical.)
- **STS-B:** semantic textual similarity
- **MRPC:** microsoft paraphrase corpus
- **RTE:** a small natural language inference corpus

# Encoder: BERT

Bidirectional Encoder  
Representations from Transformers

[Devlin et al., 2018]

- **SOTA at the time on a wide range of tasks after fine-tuning!**

System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT <sub>BASE</sub>	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>92.7</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>82.1</b>

- **QQP:** Quora Question Pairs (detect paraphrase questions)
- **QNLI:** natural language inference over question answering data
- **SST-2:** sentiment analysis
- **CoLA:** corpus of linguistic acceptability (detect whether sentences are grammatical.)
- **STS-B:** semantic textual similarity
- **MRPC:** microsoft paraphrase corpus
- **RTE:** a small natural language inference corpus

# Encoder: BERT

**Bidirectional Encoder**

[Devlin et al., 2018]

**Representations from Transformers**

# Encoder: BERT

SWAG  
(Commonsense  
inference task)

System	Dev	Test
ESIM+GloVe	51.9	52.7
ESIM+ELMo	59.1	59.2
OpenAI GPT	-	78.0
BERT <sub>BASE</sub>	81.6	-
BERT <sub>LARGE</sub>	<b>86.6</b>	<b>86.3</b>
Human (expert) <sup>†</sup>	-	85.0
Human (5 annotations) <sup>†</sup>	-	88.0

## Bidirectional Encoder Representations from Transformers

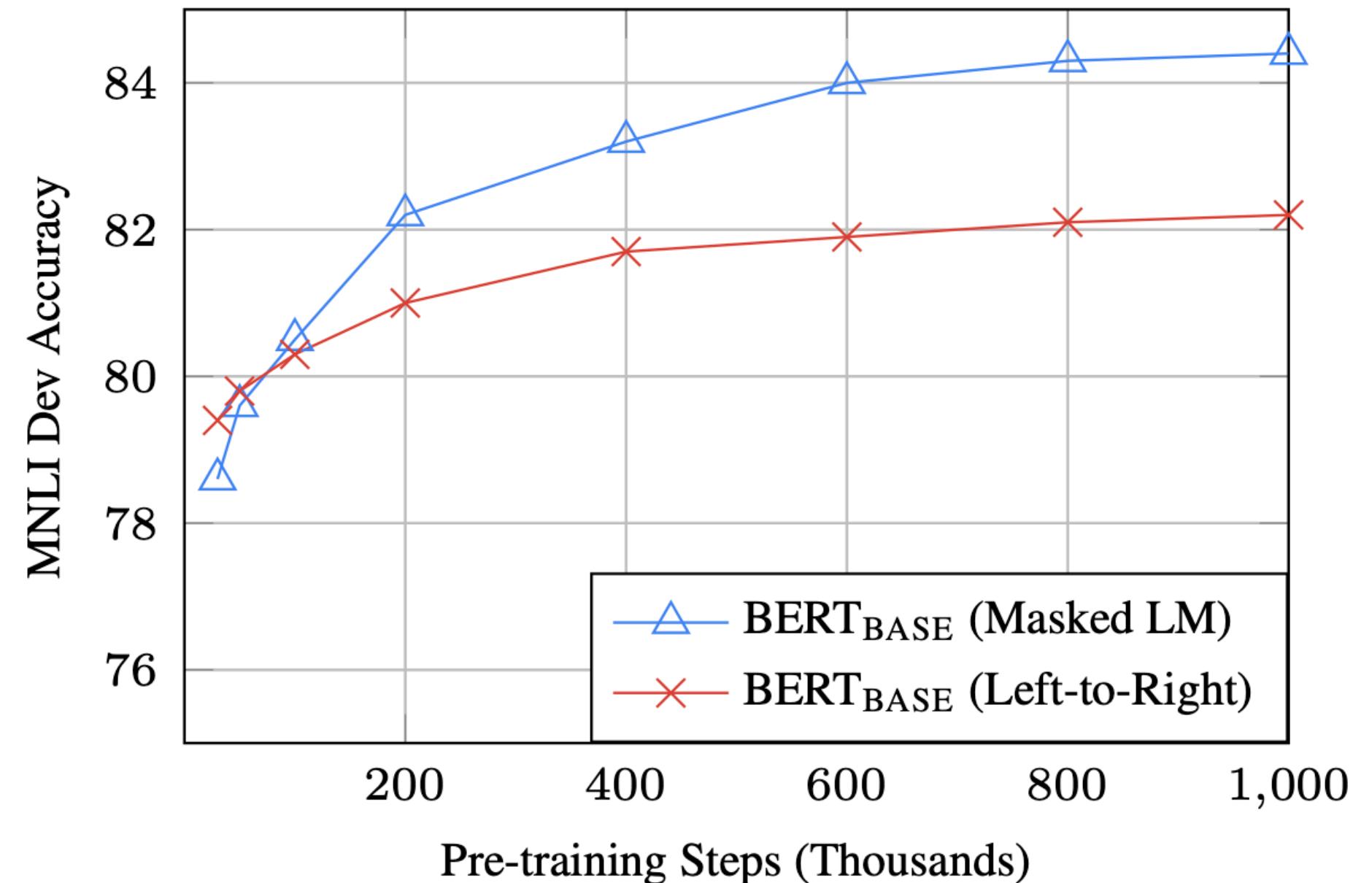
[Devlin et al., 2018]

- **Two Sizes of Models**
  - **Base:** 110M, 4 Cloud TPUs, 4 days
  - **Large:** 340M, 16 Cloud TPUs, 4 days
  - Both models can be fine-tuned with single GPU
  - The larger the better!

# Encoder: BERT

SWAG  
(Commonsense  
inference task)

System	Dev	Test
ESIM+GloVe	51.9	52.7
ESIM+ELMo	59.1	59.2
OpenAI GPT	-	78.0
$\text{BERT}_{\text{BASE}}$	81.6	-
$\text{BERT}_{\text{LARGE}}$	<b>86.6</b>	<b>86.3</b>
Human (expert) <sup>†</sup>	-	85.0
Human (5 annotations) <sup>†</sup>	-	88.0



## Bidirectional Encoder Representations from Transformers

[Devlin et al., 2018]

- **Two Sizes of Models**
  - **Base:** 110M, 4 Cloud TPUs, 4 days
  - **Large:** 340M, 16 Cloud TPUs, 4 days
  - Both models can be fine-tuned with single GPU
  - The larger the better!
- MLM converges slower than Left-to-Right at the beginning, but outperforms it eventually

# Encoder: RoBERTa

[Liu et al., 2019]

- **Original BERT is significantly undertrained!**
- More data (16G => 160G)
- Pre-train for longer
- Bigger batches
- Removing the next sentence prediction (NSP) objective
- Training on longer sequences
- Dynamic masking, randomly masking out different tokens
- A larger byte-level BPE vocabulary containing 50K sub-word units

# Encoder: RoBERTa

[Liu et al., 2019]

- **Original BERT is significantly undertrained!**
- More data (16G => 160G)
- Pre-train for longer
- Bigger batches
- Removing the next sentence prediction (NSP) objective
- Training on longer sequences
- Dynamic masking, randomly masking out different tokens
- A larger byte-level BPE vocabulary containing 50K sub-word units



All around better than BERT!

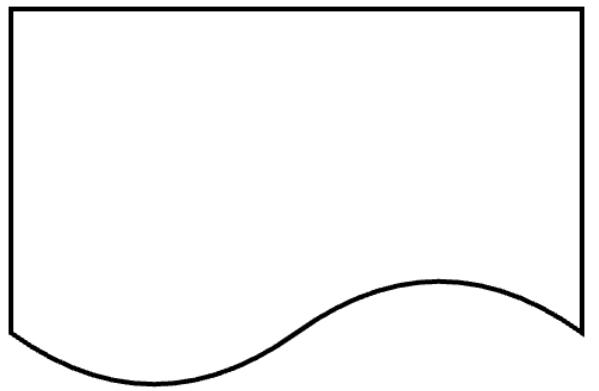
# Encoders for Information Retrieval

# Encoders for Information Retrieval

**Retrieve the set of relevant documents given a query**

# Encoders for Information Retrieval

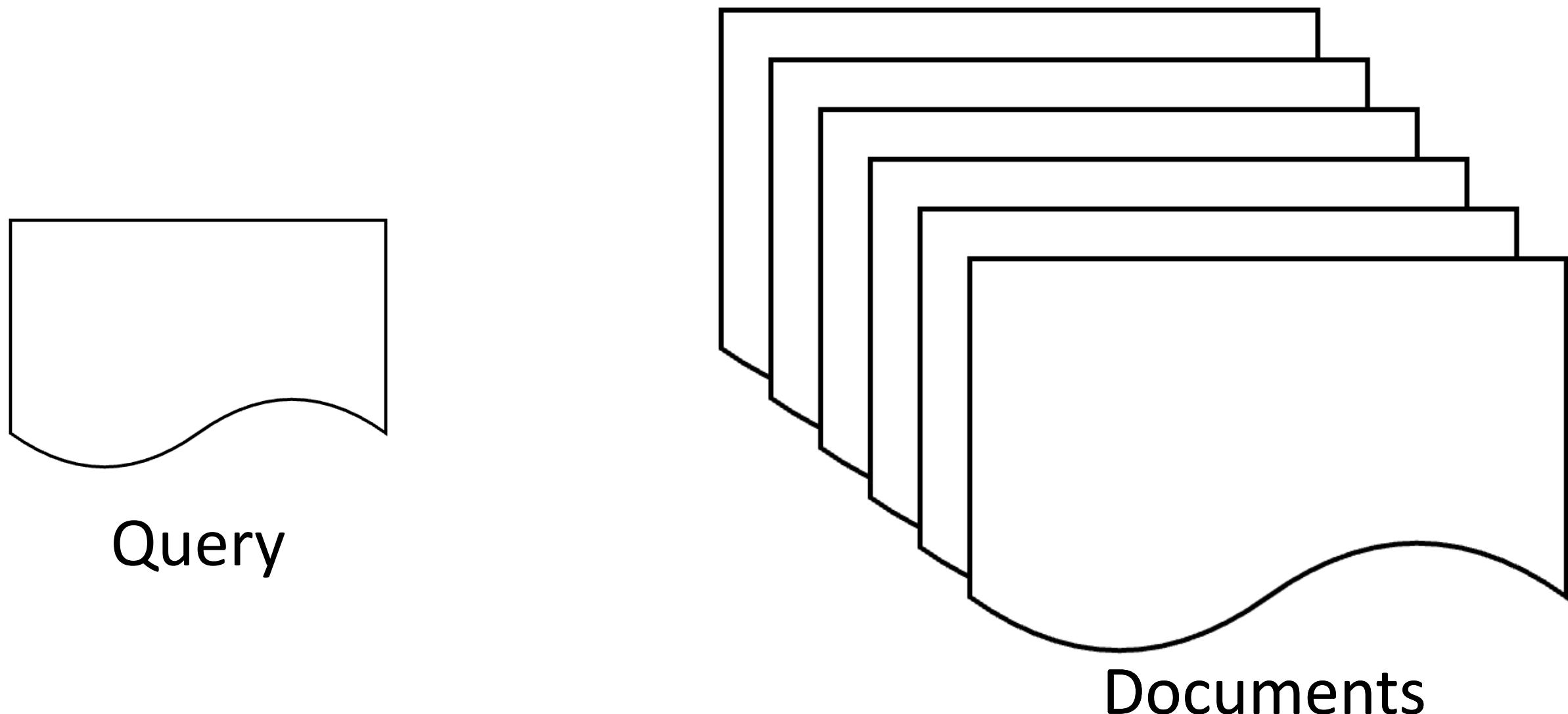
**Retrieve the set of relevant  
documents given a query**



Query

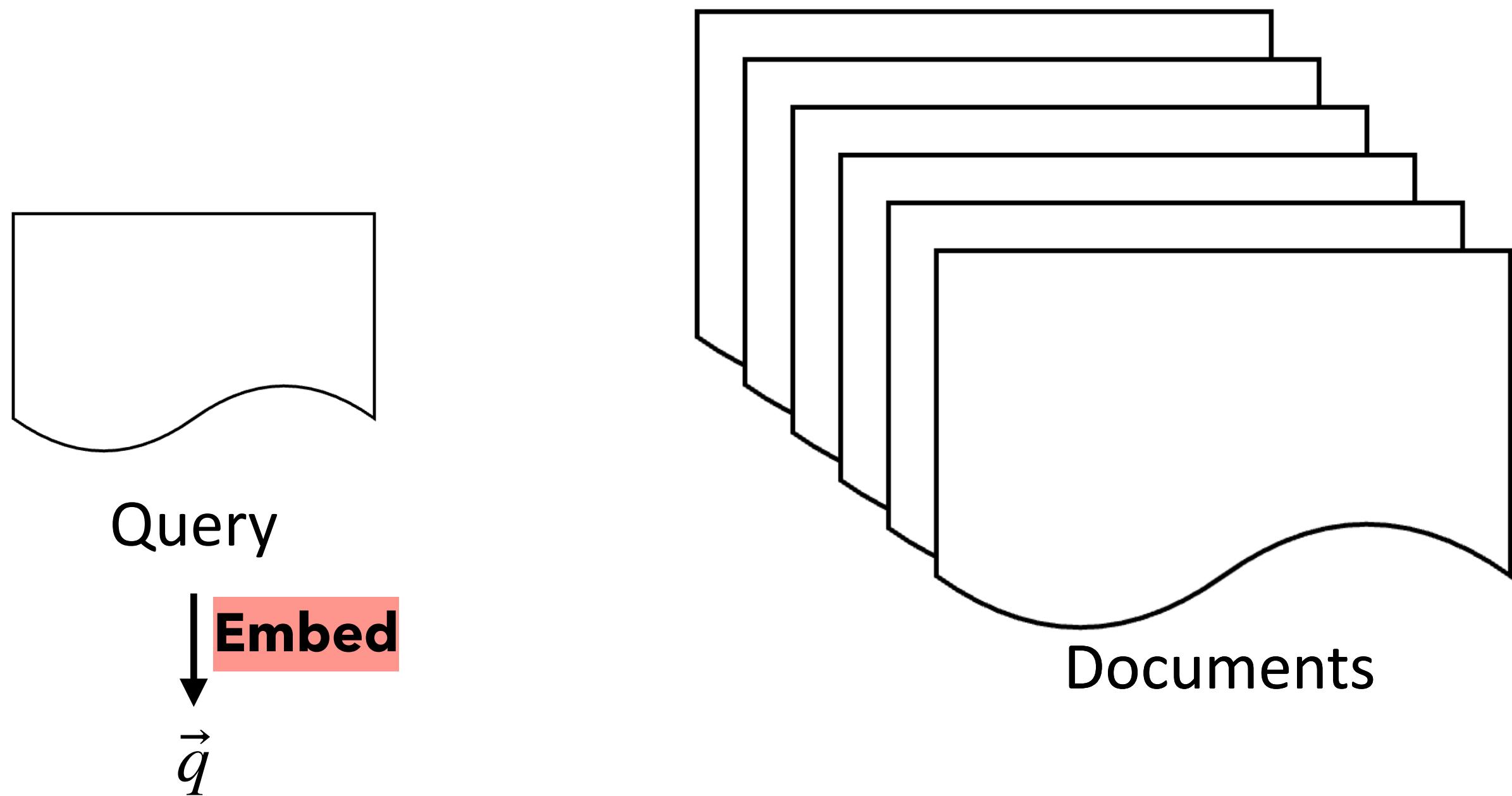
# Encoders for Information Retrieval

**Retrieve the set of relevant documents given a query**



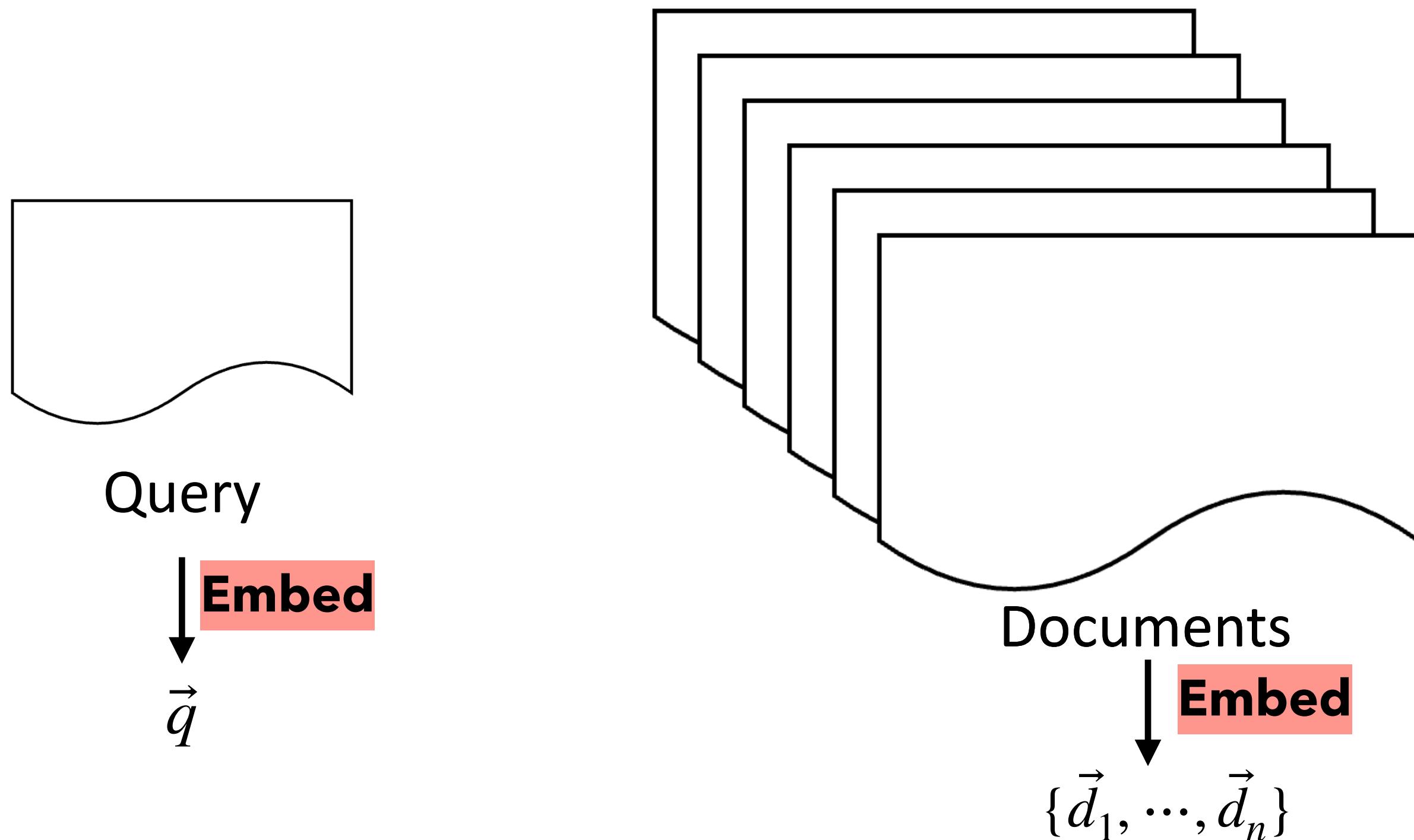
# Encoders for Information Retrieval

**Retrieve the set of relevant documents given a query**



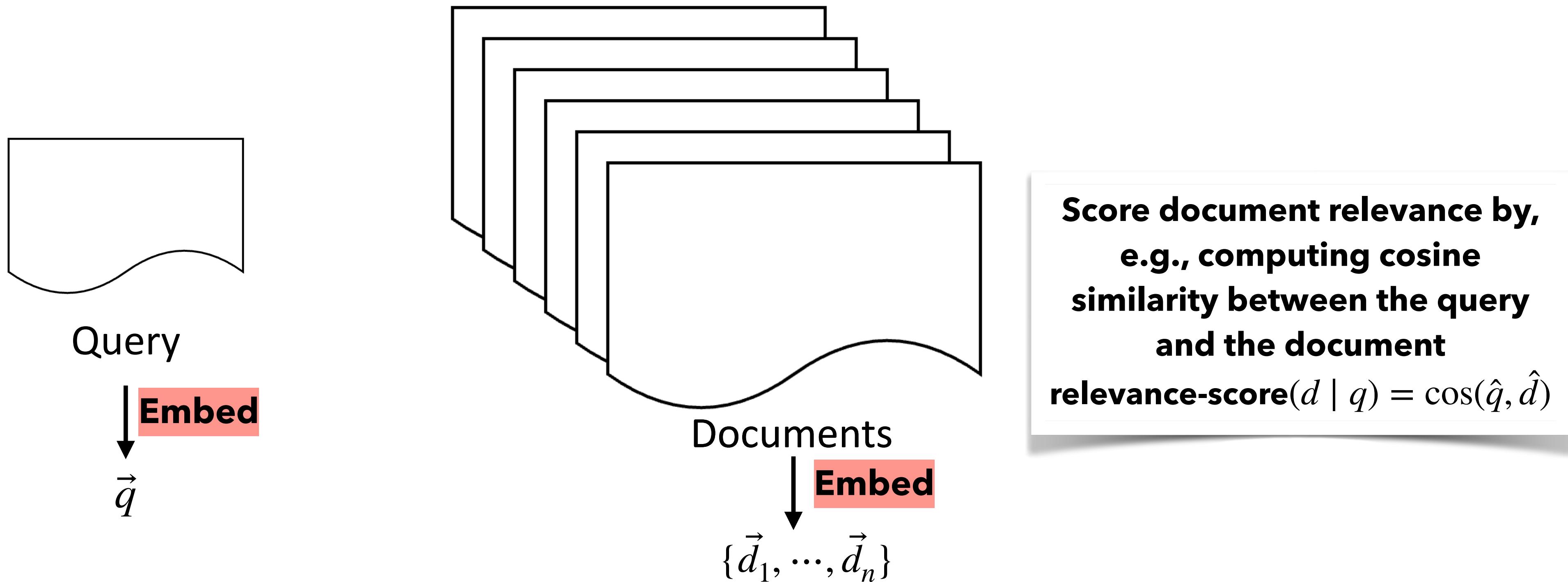
# Encoders for Information Retrieval

**Retrieve the set of relevant documents given a query**



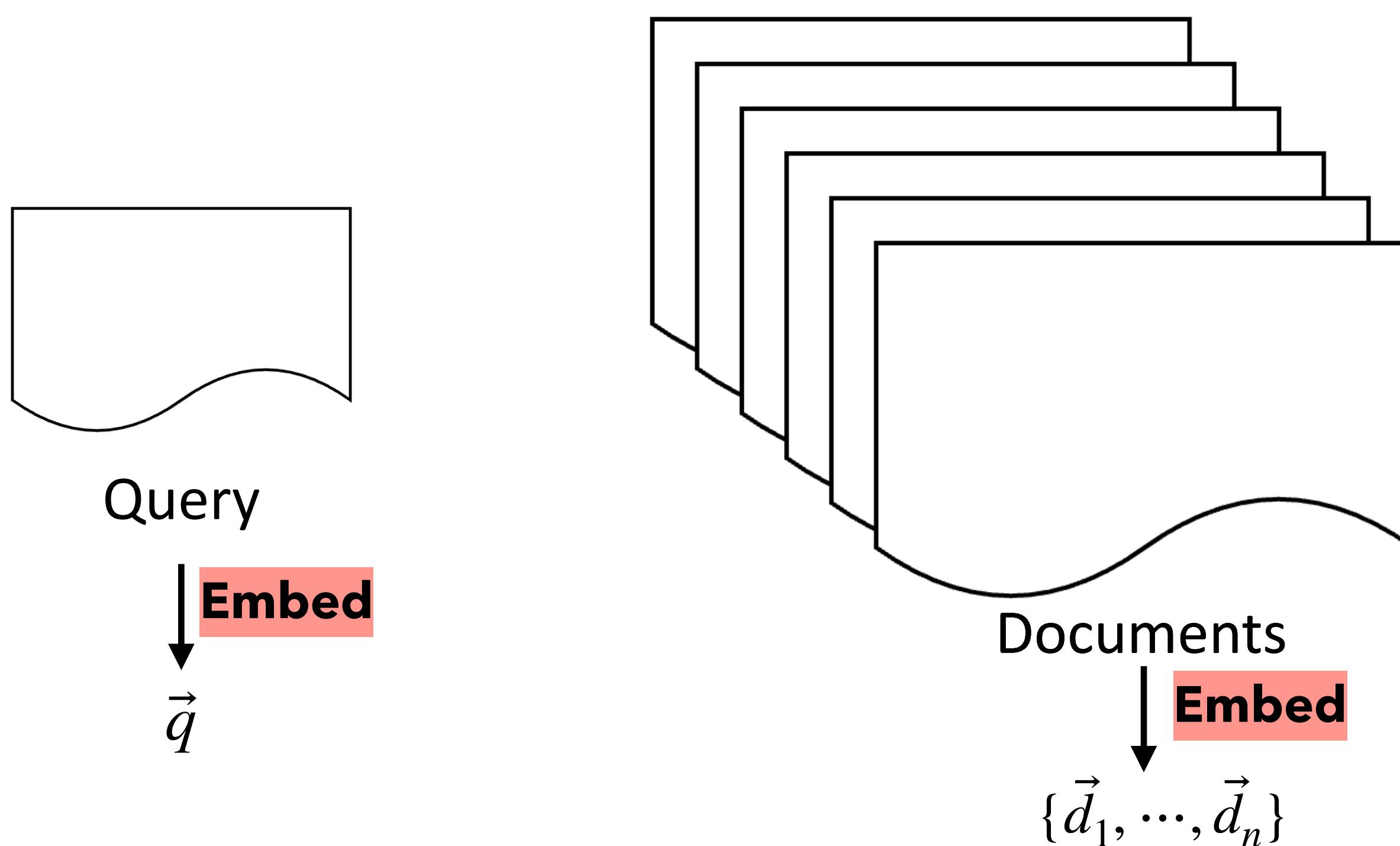
# Encoders for Information Retrieval

**Retrieve the set of relevant documents given a query**



# Encoders for Information Retrieval

**Retrieve the set of relevant documents given a query**



**Score document relevance by,  
e.g., computing cosine  
similarity between the query  
and the document**

$$\text{relevance-score}(d | q) = \cos(\hat{q}, \hat{d})$$

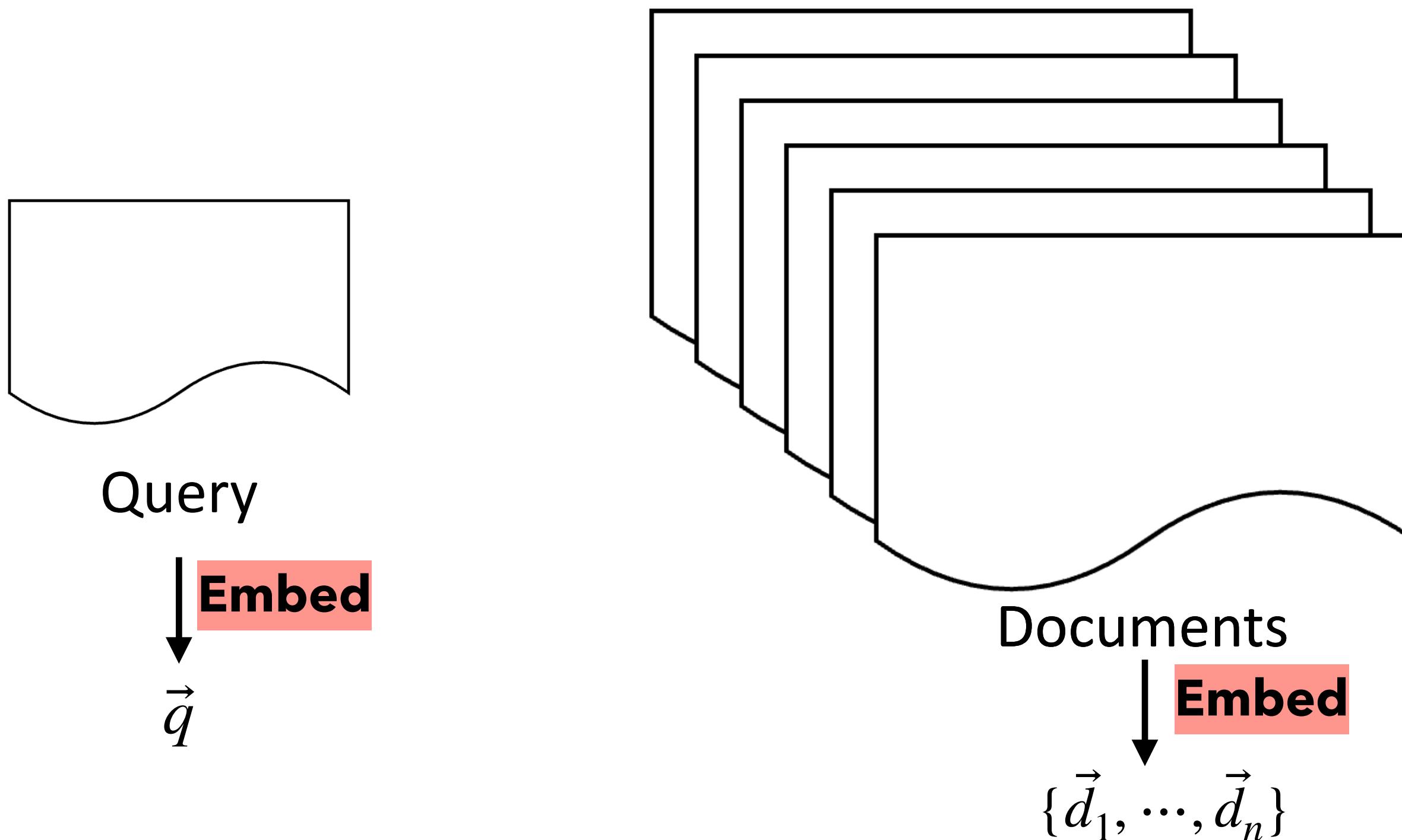
HW2!

# Encoders for Information Retrieval

**Retrieve the set of relevant documents given a query**

## Applications:

- Search Engines (This is how google works!)
- Retrieval Augmented Language Models



**Score document relevance by,  
e.g., computing cosine  
similarity between the query  
and the document**

$$\text{relevance-score}(d | q) = \cos(\hat{q}, \hat{d})$$

HW2!

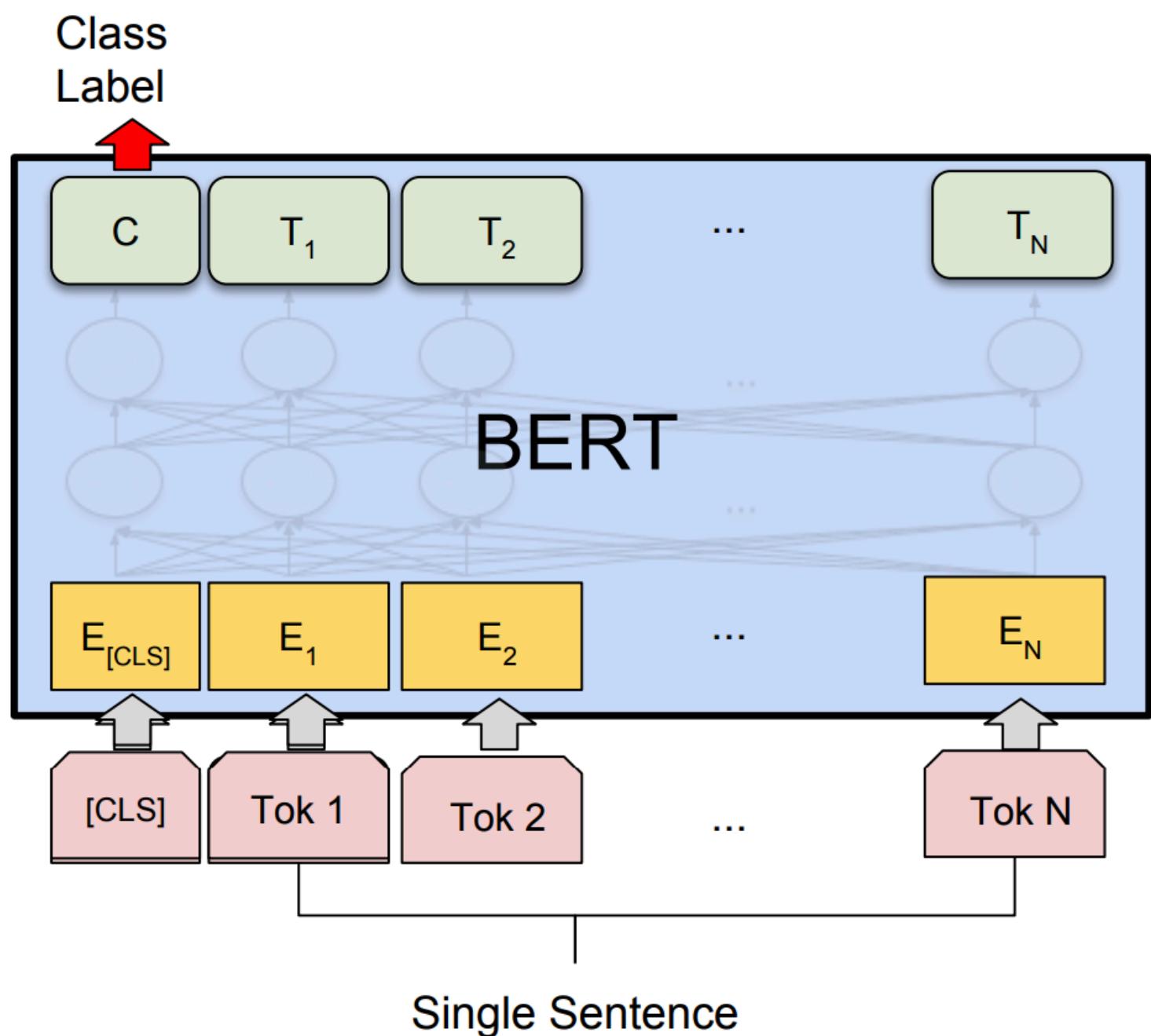
# Encoders for Information Retrieval

# Encoders for Information Retrieval

How do we get sentence embeddings from an encoder-based model like BERT?

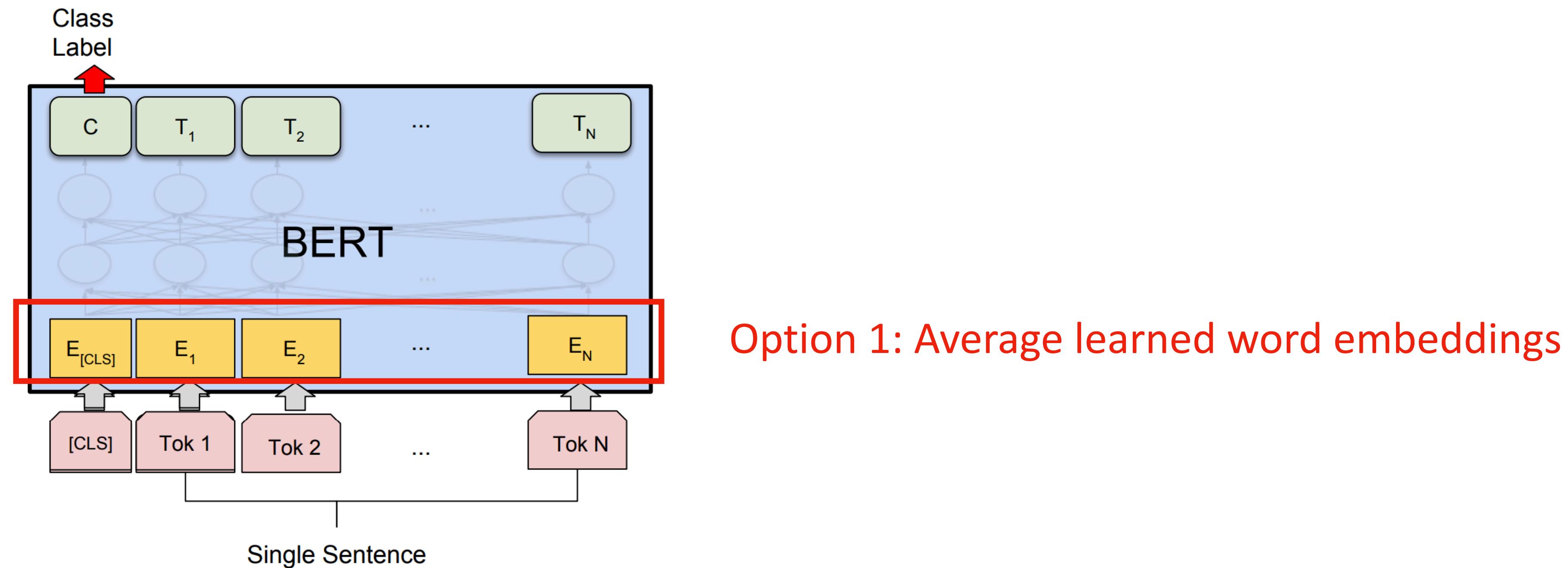
# Encoders for Information Retrieval

How do we get sentence embeddings from an encoder-based model like BERT?



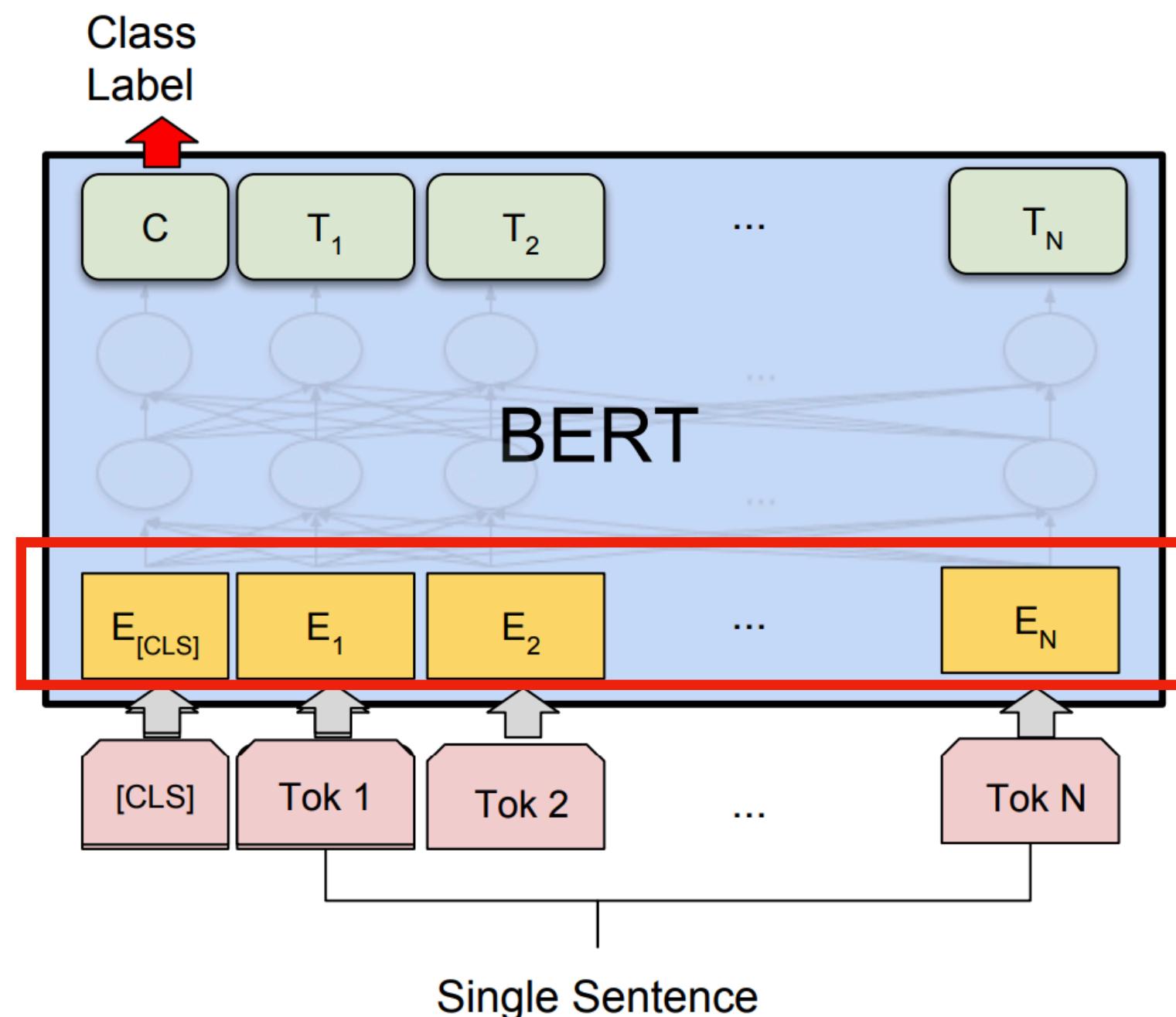
# Encoders for Information Retrieval

How do we get sentence embeddings from an encoder-based model like BERT?



# Encoders for Information Retrieval

How do we get sentence embeddings from an encoder-based model like BERT?

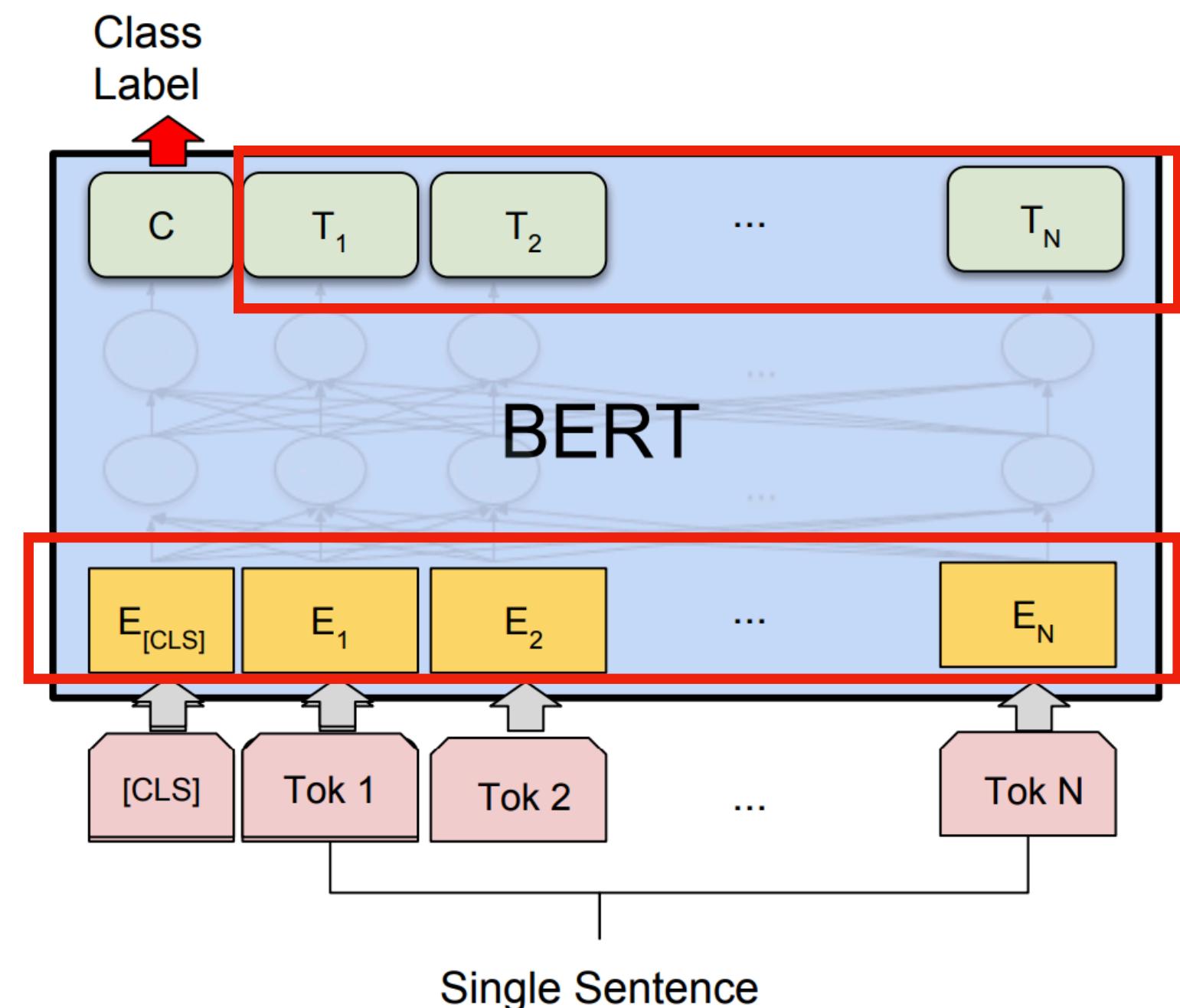


Option 1: Average learned word embeddings

**Problem:**  
Representations not contextual!  
Equivalent to using GloVe vectors

# Encoders for Information Retrieval

How do we get sentence embeddings from an encoder-based model like BERT?



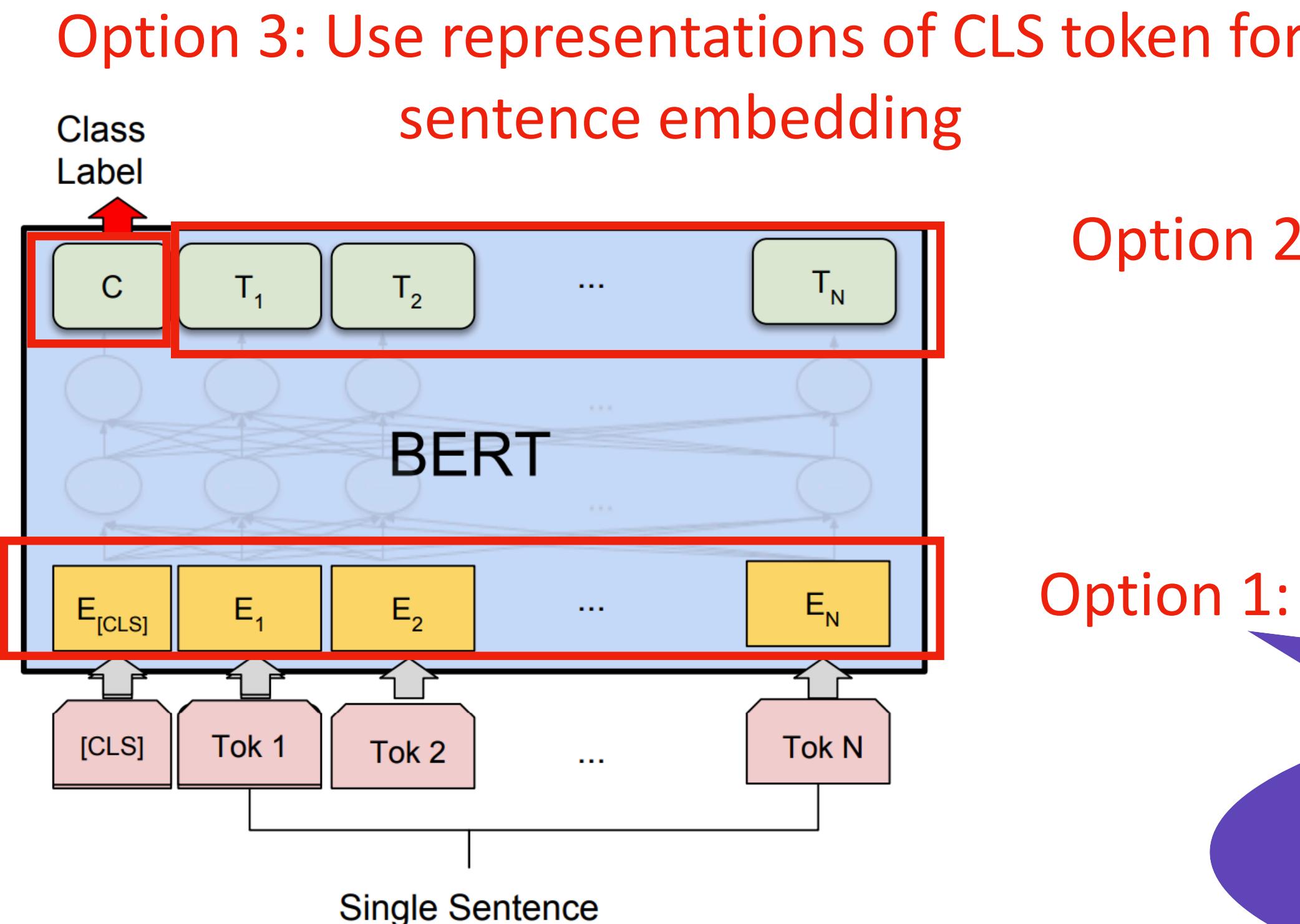
Option 2: Average learned **contextual** word embeddings

Option 1: Average learned word embeddings

**Problem:**  
Representations not contextual!  
Equivalent to using GloVe vectors

# Encoders for Information Retrieval

How do we get sentence embeddings from an encoder-based model like BERT?

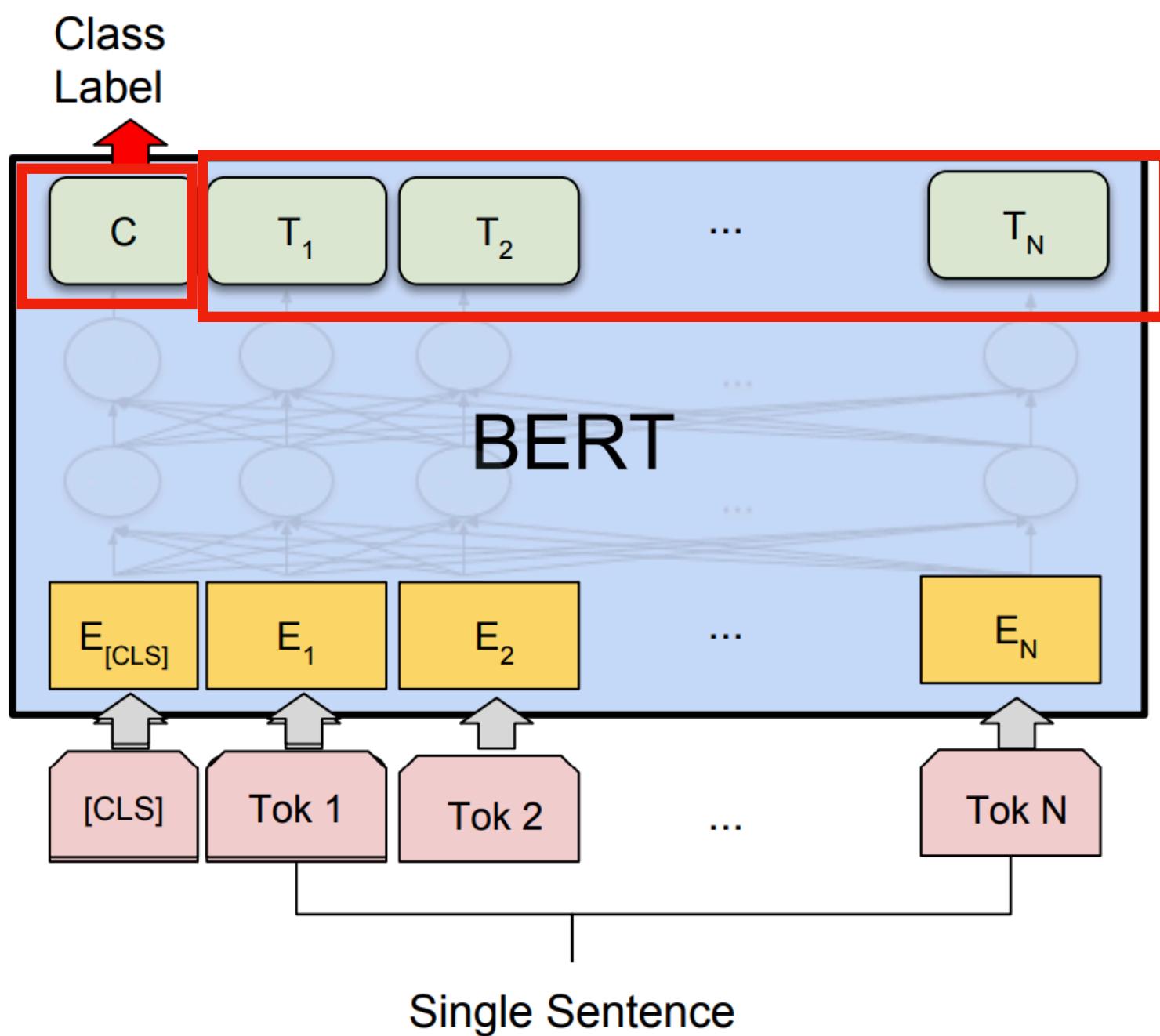


Option 2: Average learned **contextual** word embeddings

Option 1: Average learned word embeddings

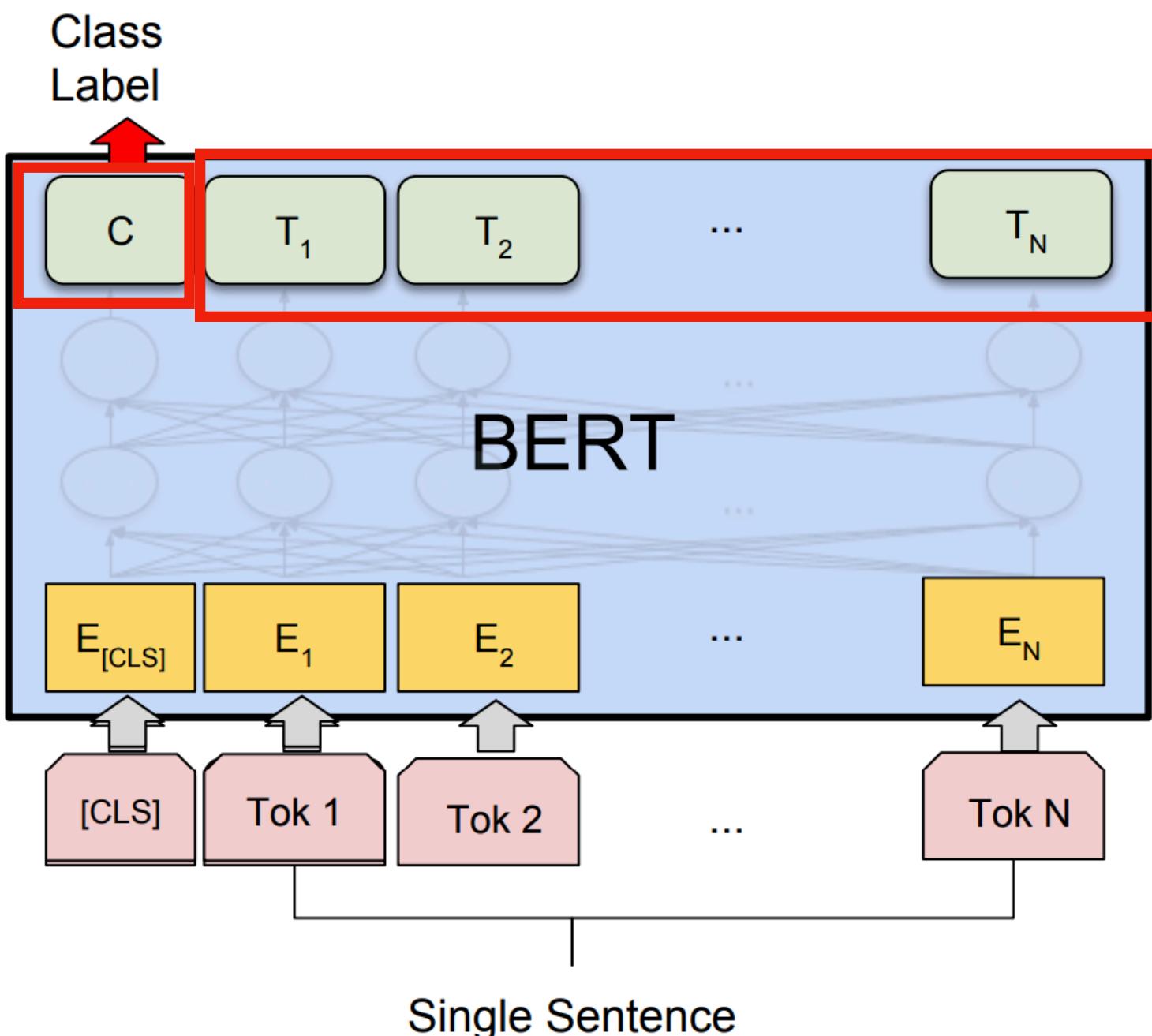
**Problem:**  
Representations not contextual!  
Equivalent to using GloVe vectors

# Encoders for Information Retrieval



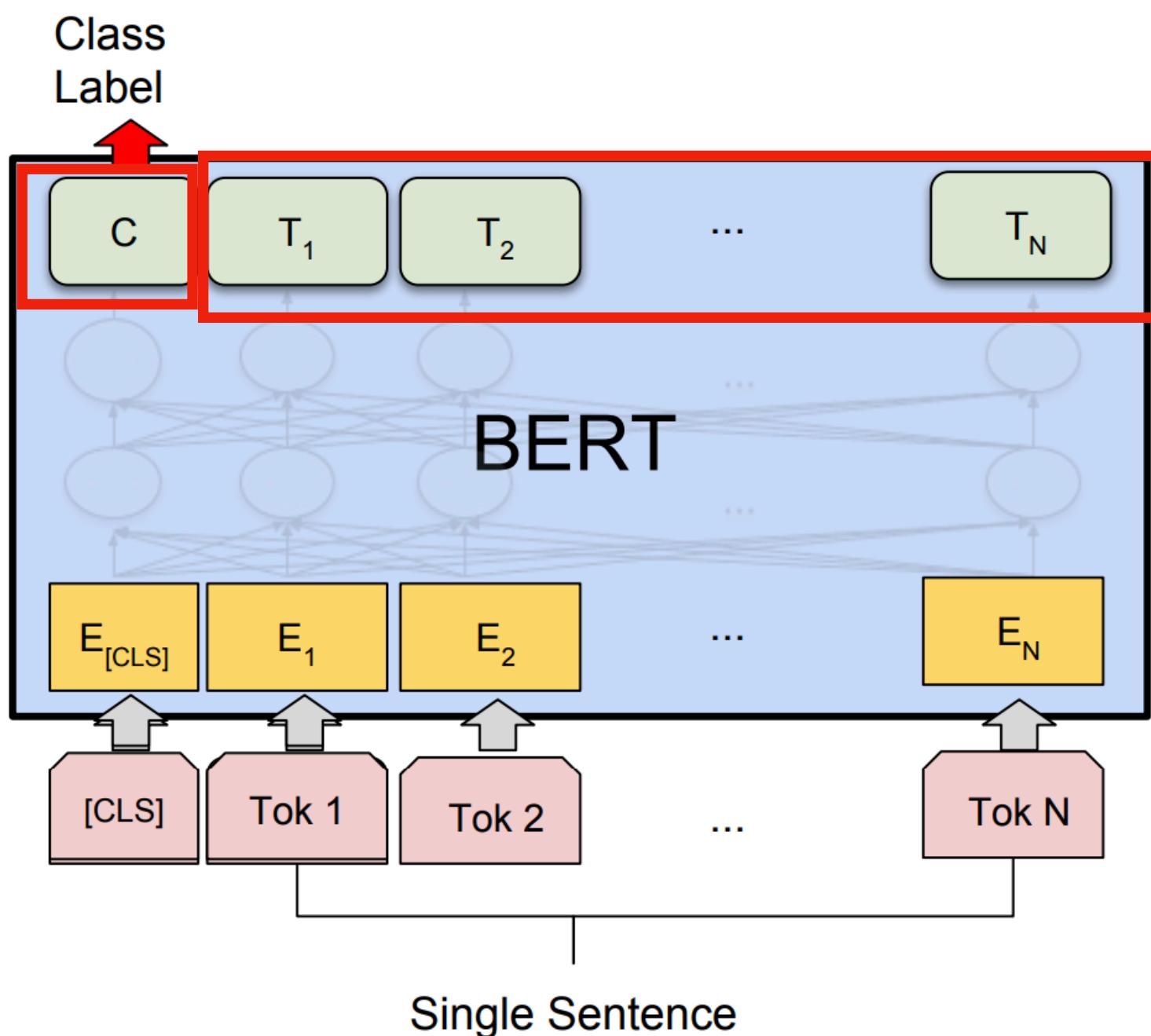
# Encoders for Information Retrieval

Out of the box even contextual representations are not very good for retrieval!



# Encoders for Information Retrieval

Out of the box even contextual representations are not very good for retrieval!

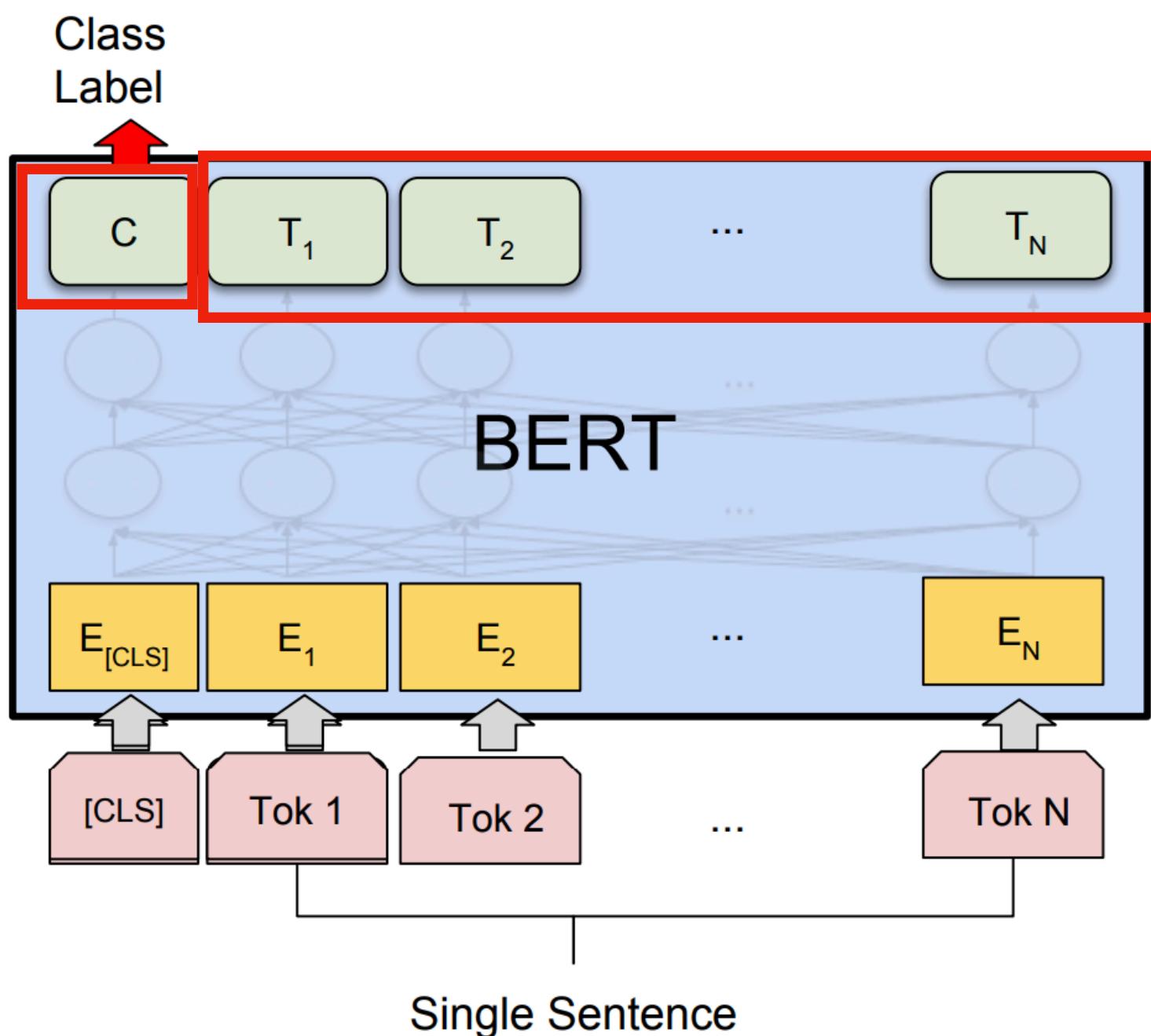


Model	STS12	STS13	STS14	STS15	STS16	STSb	SICK-R	Avg.
Avg. GloVe embeddings	55.14	70.66	59.73	68.25	63.66	58.02	53.76	61.32
Avg. BERT embeddings	38.78	57.98	57.98	63.15	61.06	46.35	58.40	54.81
BERT CLS-vector	20.16	30.01	20.09	36.88	38.08	16.50	42.63	29.19
InferSent - Glove	52.86	66.75	62.15	72.77	66.87	68.03	65.65	65.01
Universal Sentence Encoder	64.49	67.80	64.61	76.83	73.18	74.92	<b>76.69</b>	71.22

Spearman correlations for Textual Similarity (STS) tasks  
(higher is better)

# Encoders for Information Retrieval

Out of the box even contextual representations are not very good for retrieval!



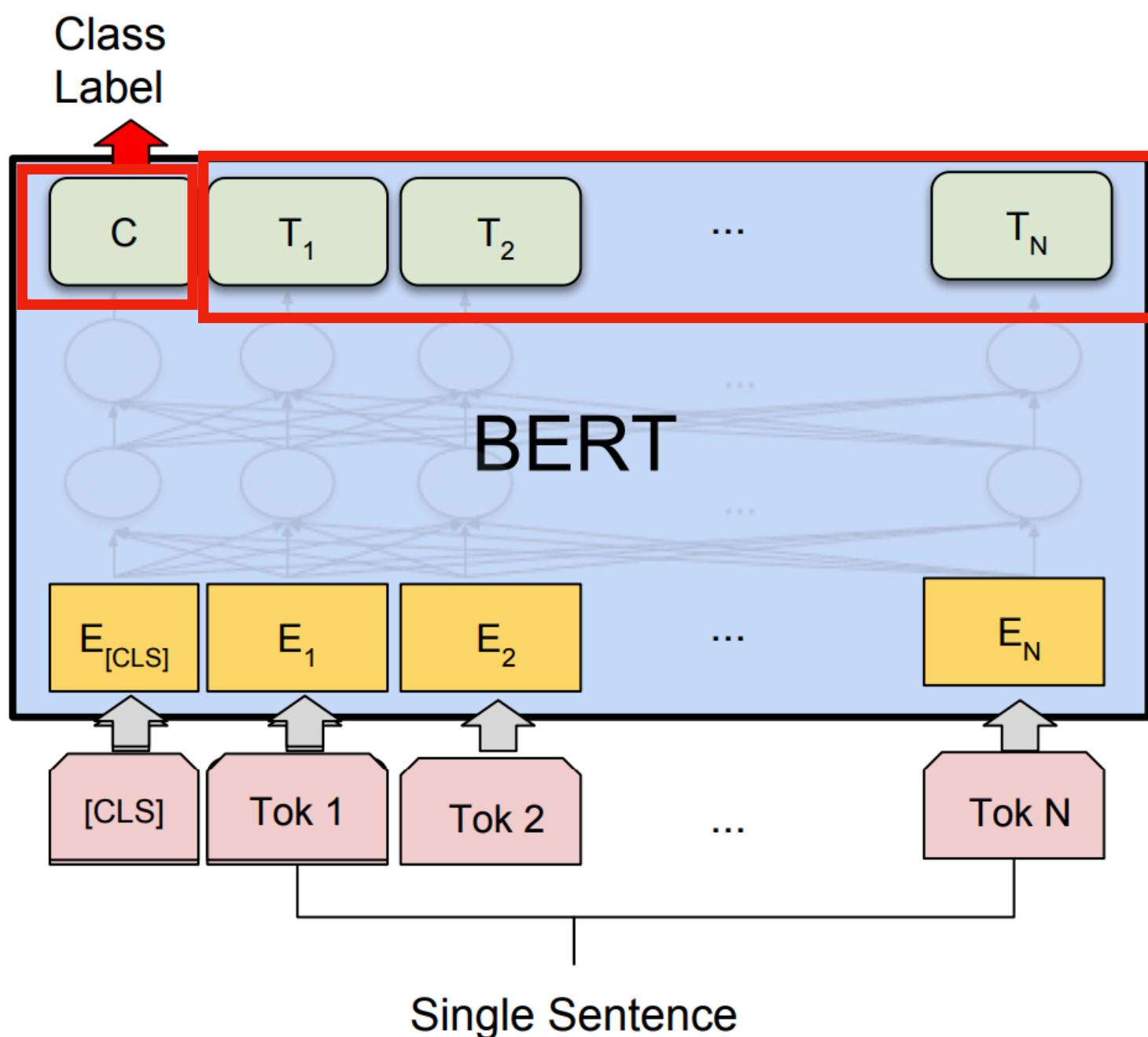
Option 2  
Option 3

Model	STS12	STS13	STS14	STS15	STS16	STSb	SICK-R	Avg.
Avg. GloVe embeddings	55.14	70.66	59.73	68.25	63.66	58.02	53.76	61.32
Avg. BERT embeddings	38.78	57.98	57.98	63.15	61.06	46.35	58.40	54.81
BERT CLS-vector	20.16	30.01	20.09	36.88	38.08	16.50	42.63	29.19
InferSent - Glove	52.86	66.75	62.15	72.77	66.87	68.03	65.65	65.01
Universal Sentence Encoder	64.49	67.80	64.61	76.83	73.18	74.92	<b>76.69</b>	71.22

Spearman correlations for Textual Similarity (STS) tasks  
(higher is better)

# Encoders for Information Retrieval

Out of the box even contextual representations are not very good for retrieval!



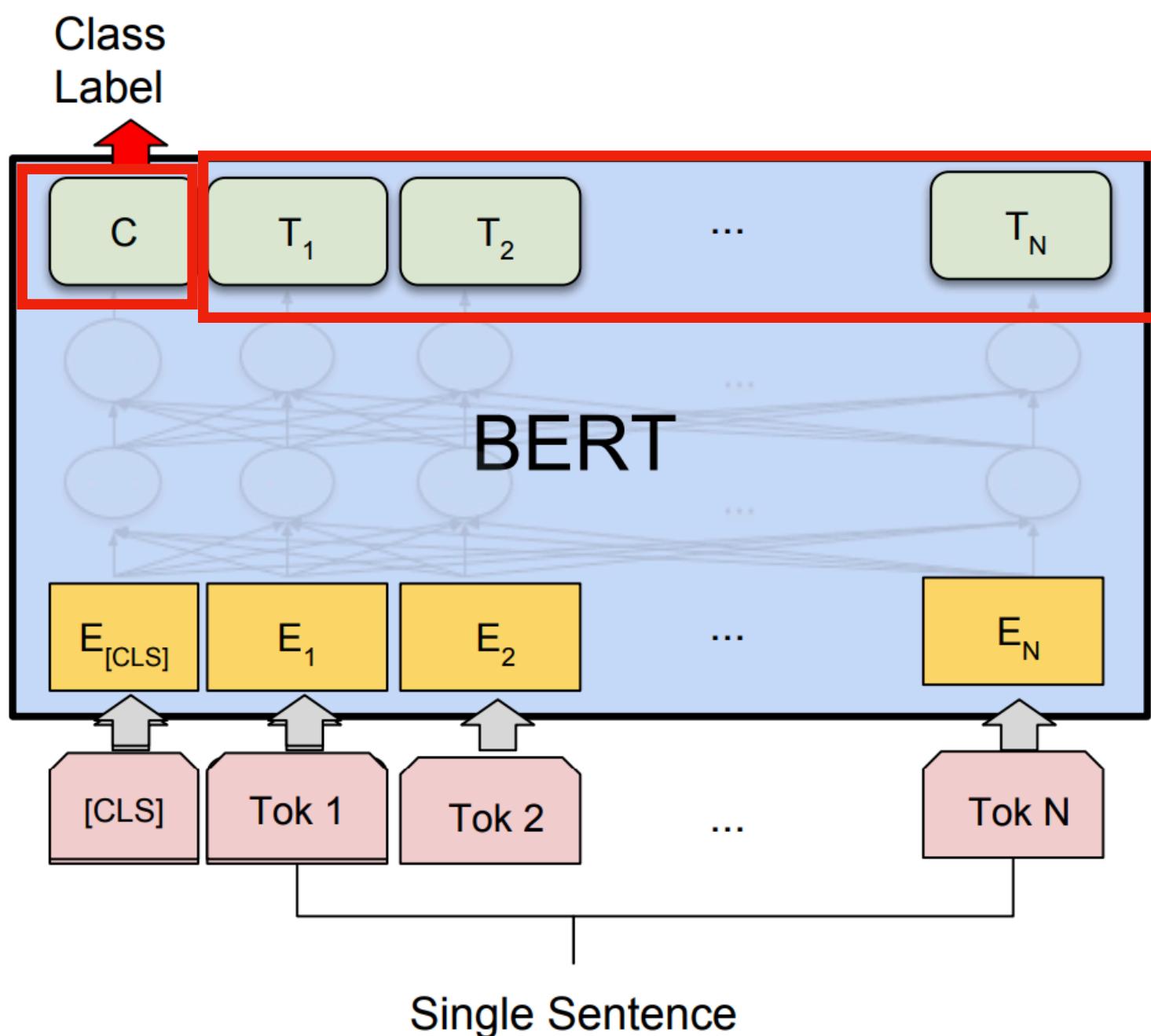
Option 2  
Option 3

Model	STS12	STS13	STS14	STS15	STS16	STSb	SICK-R	Avg.
Avg. GloVe embeddings	55.14	70.66	59.73	68.25	63.66	58.02	53.76	61.32
Avg. BERT embeddings	38.78	57.98	57.98	63.15	61.06	46.35	58.40	54.81
BERT CLS-vector	20.16	30.01	20.09	36.88	38.08	16.50	42.63	29.19
InferSent - Glove	52.86	66.75	62.15	72.77	66.87	68.03	65.65	65.01
Universal Sentence Encoder	64.49	67.80	64.61	76.83	73.18	74.92	76.69	71.22

Spearman correlations for Textual Similarity (STS) tasks  
(higher is better)

# Encoders for Information Retrieval

Out of the box even contextual representations are not very good for retrieval!



Option 1  
Option 2  
Option 3

Model	STS12	STS13	STS14	STS15	STS16	STSb	SICKER	Avg.
Avg. GloVe embeddings	55.14	70.66	59.73	68.25	63.66	58.02	53.76	61.32
Avg. BERT embeddings	38.78	57.98	57.98	63.15	61.06	46.35	58.40	54.81
BERT CLS-vector	20.16	30.01	20.09	36.88	38.08	16.50	42.63	29.19
InferSent - Glove	52.86	66.75	62.15	72.77	66.87	68.03	65.65	65.01
Universal Sentence Encoder	64.49	67.80	64.61	76.83	73.18	74.92	76.69	71.22

Spearman correlations for Textual Similarity (STS) tasks  
(higher is better)

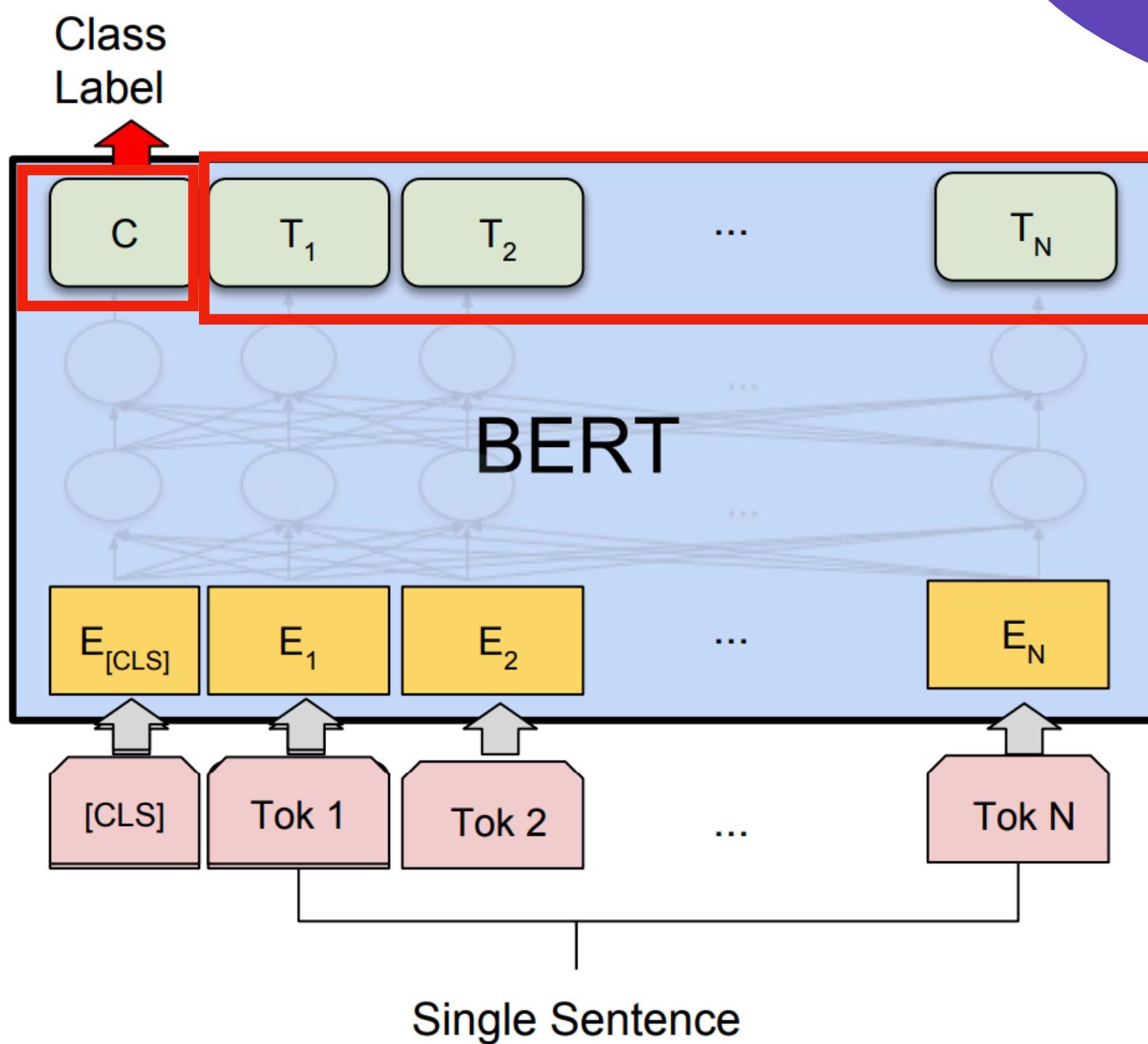
Performance is even  
worse than averaging word  
embeddings!

# Encoders for Information Retrieval

Out of the box even contextual representations are not very good for retrieval!

Why?

Performance is even worse than averaging word embeddings!



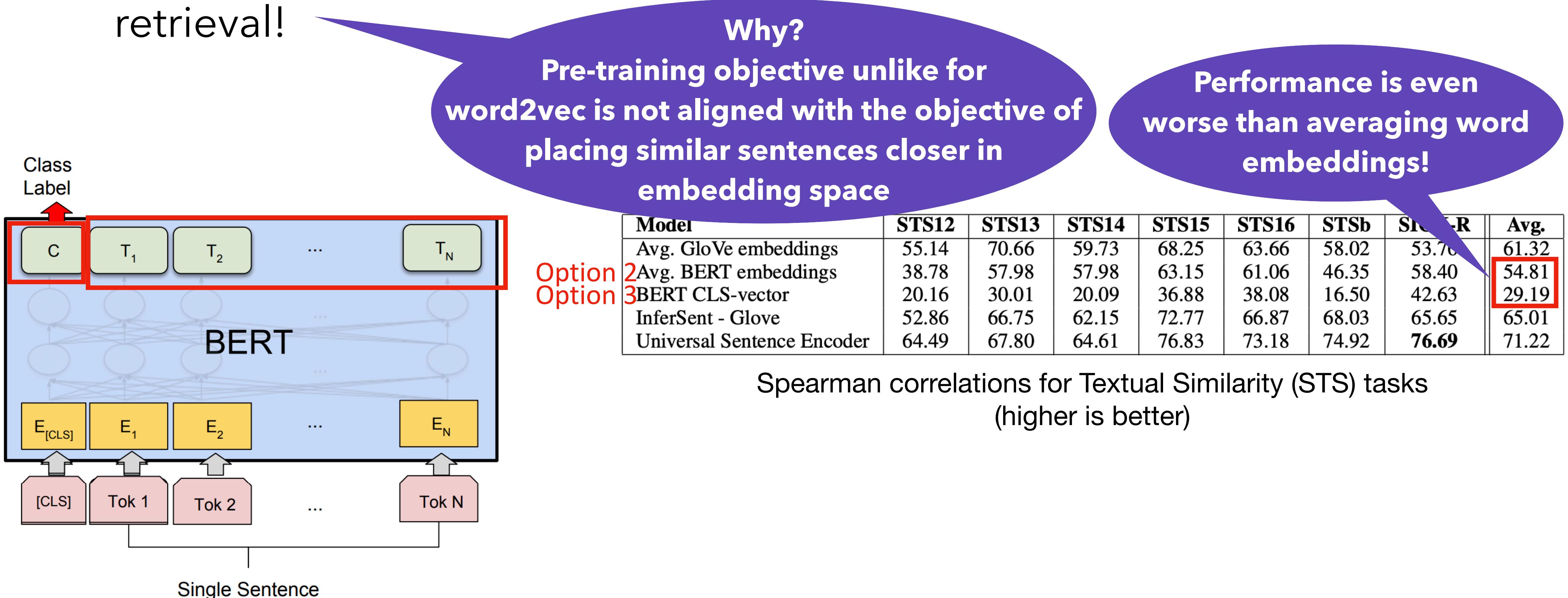
Option 1  
Option 2  
Option 3

Model	STS12	STS13	STS14	STS15	STS16	STSb	SICKER	Avg.
Avg. GloVe embeddings	55.14	70.66	59.73	68.25	63.66	58.02	53.76	61.32
Avg. BERT embeddings	38.78	57.98	57.98	63.15	61.06	46.35	58.40	54.81
BERT CLS-vector	20.16	30.01	20.09	36.88	38.08	16.50	42.63	29.19
InferSent - Glove	52.86	66.75	62.15	72.77	66.87	68.03	65.65	65.01
Universal Sentence Encoder	64.49	67.80	64.61	76.83	73.18	74.92	76.69	71.22

Spearman correlations for Textual Similarity (STS) tasks  
(higher is better)

# Encoders for Information Retrieval

Out of the box even contextual representations are not very good for retrieval!



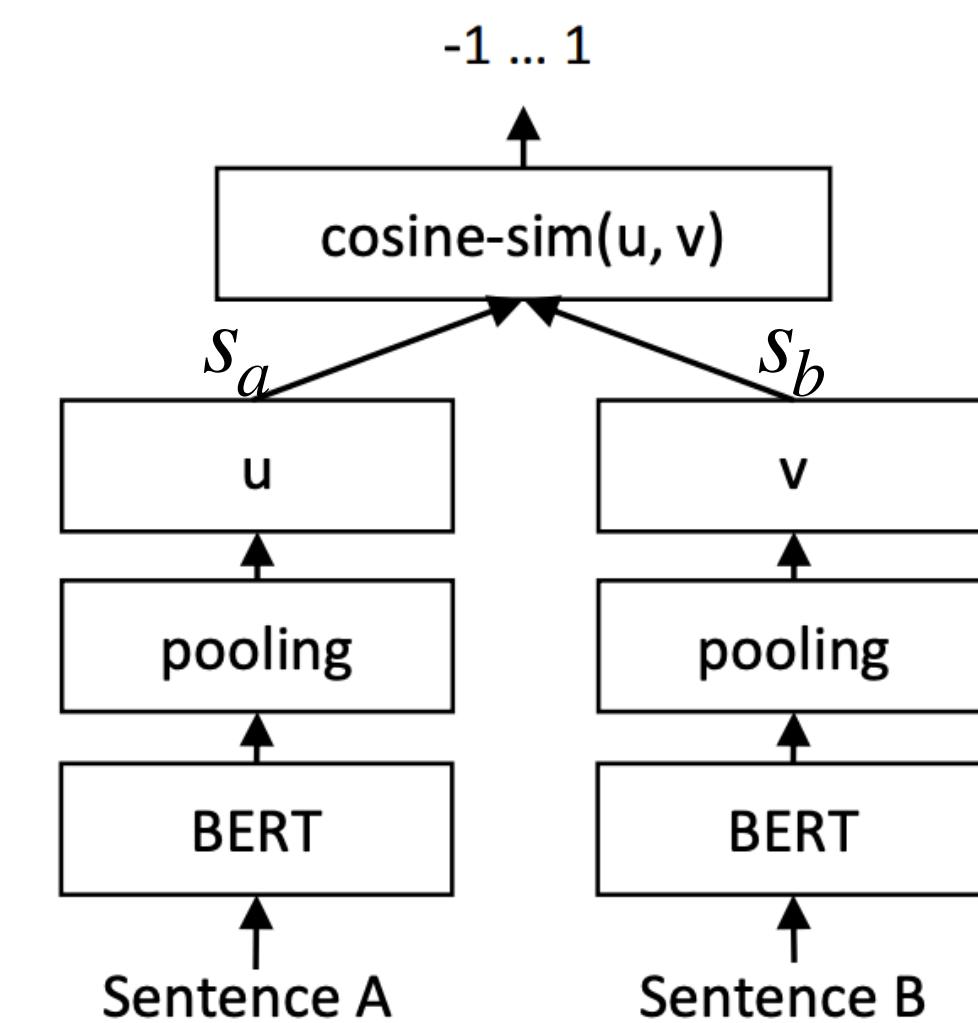
# Encoders for Information Retrieval: Sentence BERT (S-BERT)

# Encoders for Information Retrieval: Sentence BERT (S-BERT)

- Finetune BERT / RoBERTa to learn sentence/document level representations such that similar sentences are located closer in the embedding space

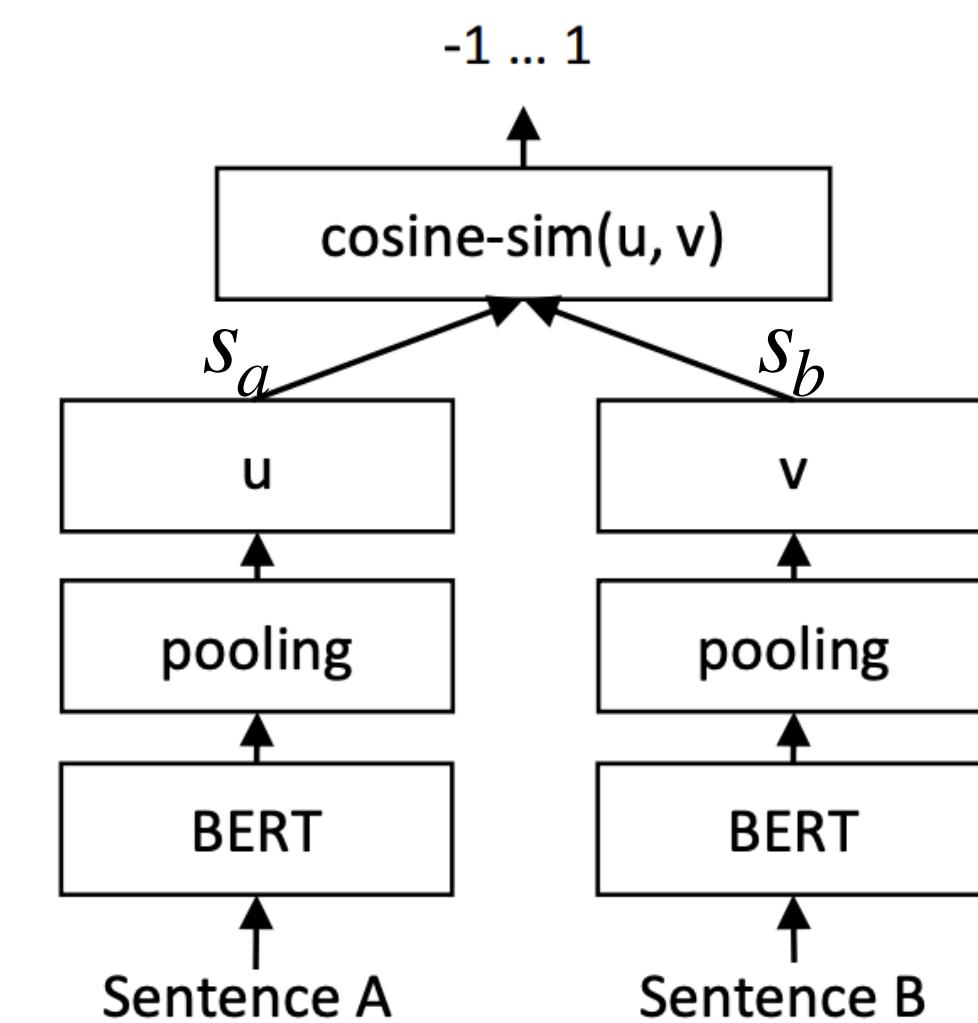
# Encoders for Information Retrieval: Sentence BERT (S-BERT)

- Finetune BERT / RoBERTa to learn sentence/document level representations such that similar sentences are located closer in the embedding space



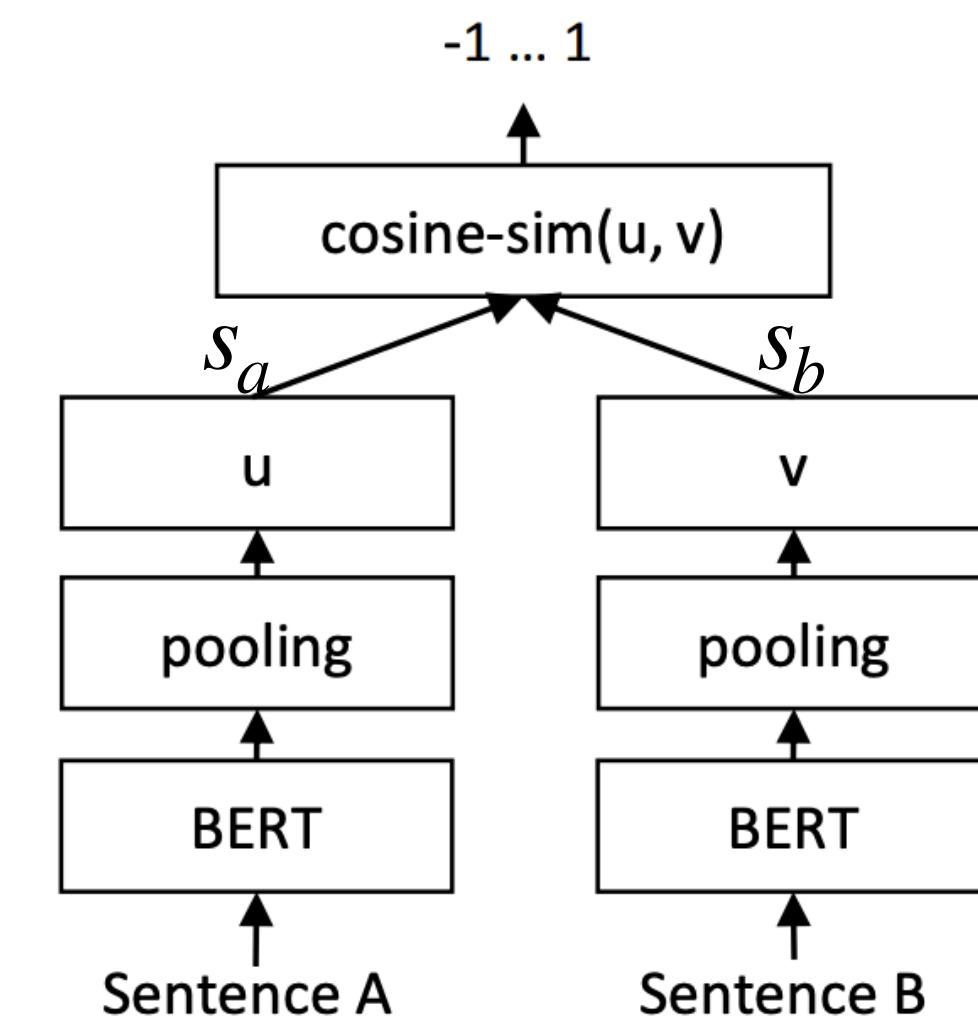
# Encoders for Information Retrieval: Sentence BERT (S-BERT)

- Finetune BERT / RoBERTa to learn sentence/document level representations such that similar sentences are located closer in the embedding space
- Uses a triplet objective function – Given an anchor sentence  $a$ , a positive sentence  $p$ , and a negative sentence  $n$ , triplet loss tunes the network such that the distance between  $a$  and  $p$  is smaller than the distance between  $a$  and  $n$ .



# Encoders for Information Retrieval: Sentence BERT (S-BERT)

- Finetune BERT / RoBERTa to learn sentence/document level representations such that similar sentences are located closer in the embedding space
- Uses a triplet objective function – Given an anchor sentence  $a$ , a positive sentence  $p$ , and a negative sentence  $n$ , triplet loss tunes the network such that the distance between  $a$  and  $p$  is smaller than the distance between  $a$  and  $n$ .



Triplet objective function  
 $\max(\|s_a - s_p\| - \|s_a - s_n\| + \epsilon, 0)$

# Encoders for Information Retrieval: Sentence BERT (S-BERT)

# Encoders for Information Retrieval: Sentence BERT (S-BERT)

Model	STS12	STS13	STS14	STS15	STS16	STSb	SICK-R	Avg.
Avg. GloVe embeddings	55.14	70.66	59.73	68.25	63.66	58.02	53.76	61.32
Avg. BERT embeddings	38.78	57.98	57.98	63.15	61.06	46.35	58.40	54.81
BERT CLS-vector	20.16	30.01	20.09	36.88	38.08	16.50	42.63	29.19
InferSent - Glove	52.86	66.75	62.15	72.77	66.87	68.03	65.65	65.01
Universal Sentence Encoder	64.49	67.80	64.61	76.83	73.18	74.92	<b>76.69</b>	71.22
SBERT-NLI-base	70.97	76.53	73.19	79.09	74.30	77.03	72.91	74.89
SBERT-NLI-large	72.27	<b>78.46</b>	<b>74.90</b>	80.99	76.25	<b>79.23</b>	73.75	76.55
SRoBERTa-NLI-base	71.54	72.49	70.80	78.74	73.69	77.77	74.46	74.21
SRoBERTa-NLI-large	<b>74.53</b>	77.00	73.18	<b>81.85</b>	<b>76.82</b>	79.10	74.29	<b>76.68</b>

Spearman correlations for Textual Similarity (STS) tasks  
(higher is better)

# Encoders for Information Retrieval: Sentence BERT (S-BERT)

Model	STS12	STS13	STS14	STS15	STS16	STSb	SICK-R	Avg.
Avg. GloVe embeddings	55.14	70.66	59.73	68.25	63.66	58.02	53.76	61.32
Avg. BERT embeddings	38.78	57.98	57.98	63.15	61.06	46.35	58.40	54.81
BERT CLS-vector	20.16	30.01	20.09	36.88	38.08	16.50	42.63	29.19
InferSent - Glove	52.86	66.75	62.15	72.77	66.87	68.03	65.65	65.01
Universal Sentence Encoder	64.49	67.80	64.61	76.83	73.18	74.92	<b>76.69</b>	71.22
SBERT-NLI-base	70.97	76.53	73.19	79.09	74.30	77.03	72.91	74.89
SBERT-NLI-large	72.27	<b>78.46</b>	<b>74.90</b>	80.99	76.25	<b>79.23</b>	73.75	76.55
SRoBERTa-NLI-base	71.54	72.49	70.80	78.74	73.69	77.77	74.46	74.21
SRoBERTa-NLI-large	<b>74.53</b>	77.00	73.18	<b>81.85</b>	<b>76.82</b>	79.10	74.29	<b>76.68</b>

Spearman correlations for Textual Similarity (STS) tasks  
(higher is better)

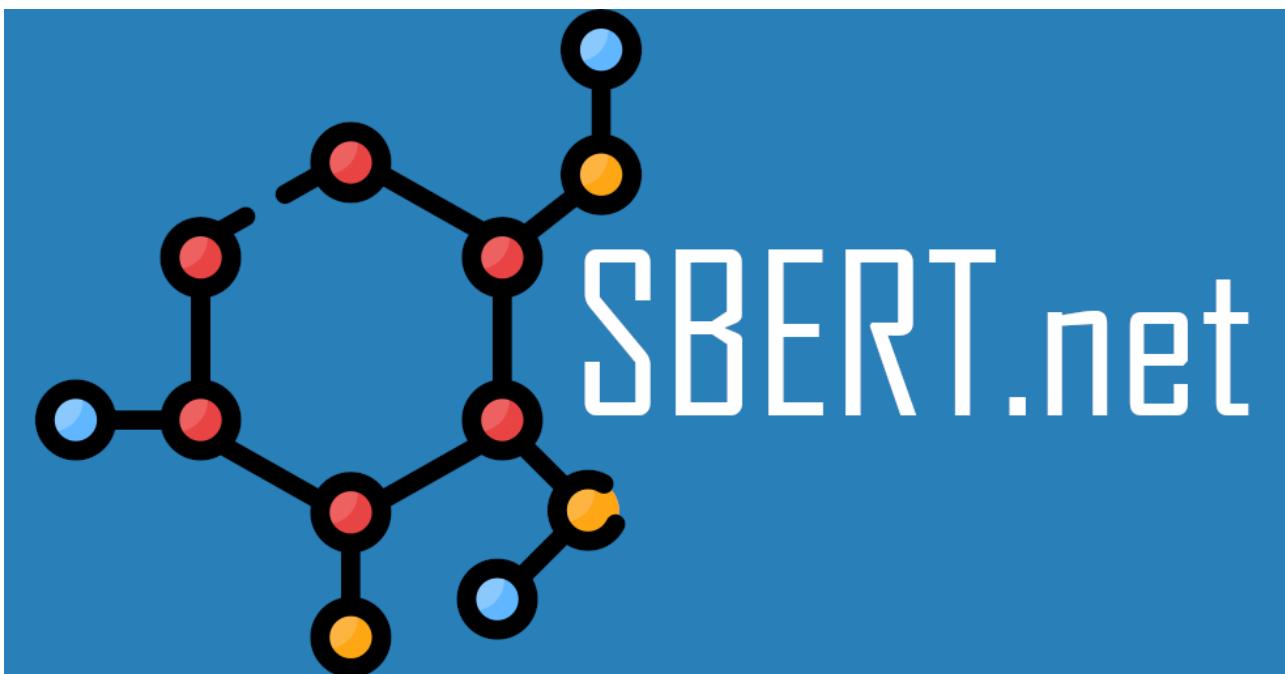
# Encoders for Information Retrieval: Sentence BERT (S-BERT)

Model	STS12	STS13	STS14	STS15	STS16	STSb	SICK-R	Avg.
Avg. GloVe embeddings	55.14	70.66	59.73	68.25	63.66	58.02	53.76	61.32
Avg. BERT embeddings	38.78	57.98	57.98	63.15	61.06	46.35	58.40	54.81
BERT CLS-vector	20.16	30.01	20.09	36.88	38.08	16.50	42.63	29.19
InferSent - Glove	52.86	66.75	62.15	72.77	66.87	68.03	65.65	65.01
Universal Sentence Encoder	64.49	67.80	64.61	76.83	73.18	74.92	<b>76.69</b>	71.22
SBERT-NLI-base	70.97	76.53	73.19	79.09	74.30	77.03	72.91	74.89
SBERT-NLI-large	72.27	<b>78.46</b>	<b>74.90</b>	80.99	76.25	<b>79.23</b>	73.75	76.55
SRoBERTa-NLI-base	71.54	72.49	70.80	78.74	73.69	77.77	74.46	74.21
SRoBERTa-NLI-large	<b>74.53</b>	77.00	73.18	<b>81.85</b>	<b>76.82</b>	79.10	74.29	<b>76.68</b>

Spearman correlations for Textual Similarity (STS) tasks  
(higher is better)

Sentence-BERT /  
RoBERTa performs  
remarkably better than the  
existing approaches!

# Encoders for Information Retrieval: Sentence BERT (S-BERT)



Sentence Transformers Library.  
Very handy for using pre-trained Sentence-BERT-like models

Model	STS12	STS13	STS14	STS15	STS16	STSb	SICK-R	Avg.
Avg. GloVe embeddings	55.14	70.66	59.73	68.25	63.66	58.02	53.76	61.32
Avg. BERT embeddings	38.78	57.98	57.98	63.15	61.06	46.35	58.40	54.81
BERT CLS-vector	20.16	30.01	20.09	36.88	38.08	16.50	42.63	29.19
InferSent - Glove	52.86	66.75	62.15	72.77	66.87	68.03	65.65	65.01
Universal Sentence Encoder	64.49	67.80	64.61	76.83	73.18	74.92	<b>76.69</b>	71.22
SBERT-NLI-base	70.97	76.53	73.19	79.09	74.30	77.03	72.91	74.89
SBERT-NLI-large	72.27	<b>78.46</b>	<b>74.90</b>	80.99	76.25	<b>79.23</b>	73.75	76.55
SRoBERTa-NLI-base	71.54	72.49	70.80	78.74	73.69	77.77	74.46	74.21
SRoBERTa-NLI-large	<b>74.53</b>	77.00	73.18	<b>81.85</b>	<b>76.82</b>	79.10	74.29	<b>76.68</b>

Spearman correlations for Textual Similarity (STS) tasks  
(higher is better)

Sentence-BERT /  
RoBERTa performs  
remarkably better than the  
existing approaches!

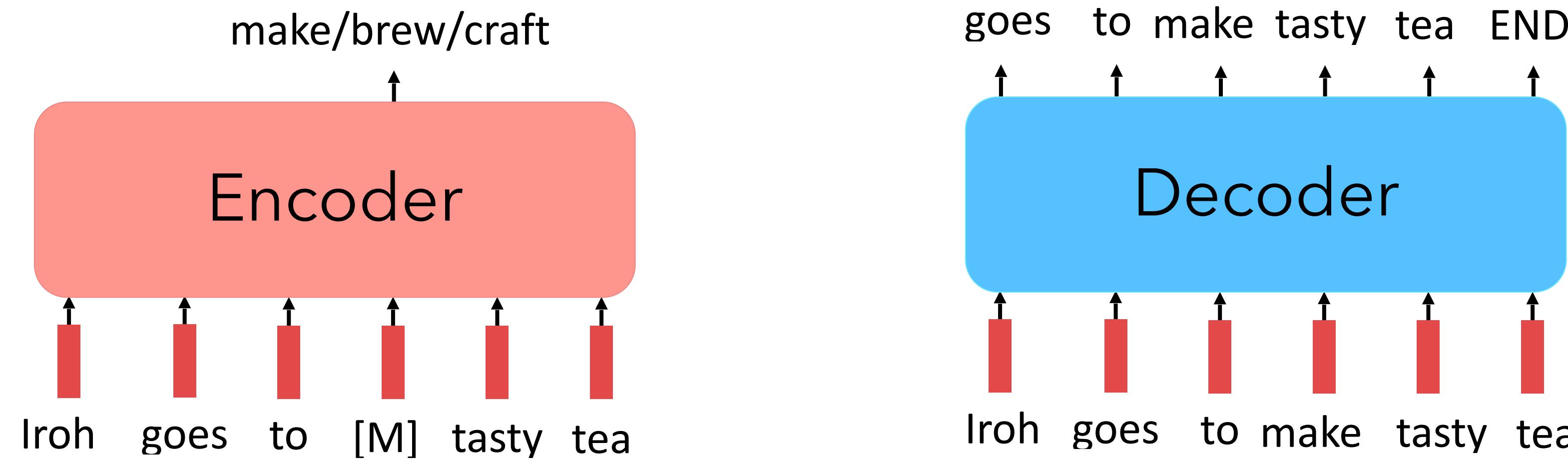
# Encoder: Pros & Cons



- Consider both left and right context
- Capture intricate contextual relationships



- Not good at generating open-text from left-to-right, one token at a time



# Thank you!