# Multimodal Machine Learning

## Lecture 7.1: Alignment and Representations

Louis-Philippe Morency

# Objectives of today's class

- Contextualized sentence embedding
- Transformer networks
    - Self-attention
    - Multi-head attention
    - Position embeddings
    - Sequence-to-sequence modeling
- Multimodal contextualized embeddings
- Language pre-training
    - BERT pre-training and fine-tuning
- Multimodal pre-training

# Administrative Stuff

# Upcoming Schedule

- First project assignment:
  - Proposal presentation (10/1 and 10/3)
  - First project report (Sunday 10/6)
- Midterm project assignment
  - Midterm presentations (11/5 and 11/7)
  - Midterm report (Sunday 11/10)
- Final project assignment
  - Final presentation (12/3 & 12/5)
  - Final report (Sunday 12/8)

# Midterm Project Report Instructions

- **Goal:** Evaluate state-of-the-art models on your dataset and identify key issues through a detailed error analysis
  - It will inform the design of your new research ideas
- **Report format:** 8 pages, 2 column (ICML template)
  - The report should follow a similar structure to a research paper
- **Number of SOTA models**
  - Teams of 3 should have at least two baseline models
  - Teams of 4 or 5 should have at least three baseline models
- **Error analysis**
  - This is one of the most important part of this report. You need to understand where previous models can be improved.

# Midterm Project Report Instructions

Main report sections:

- Abstract
- Introduction
- Related work
- Problem statement
- Multimodal baseline models
- Experimental methodology
- Results and discussion
- New research ideas

# Midterm Presentations

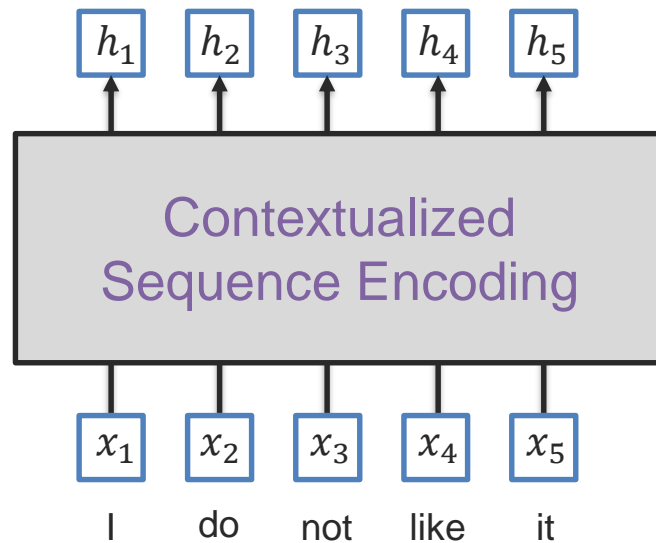Oral Presentations          **OR**          Poster Presentations

Language Technologies Institute          Carnegie Mellon University

# Contextualized Sequence Encoding

# Sequence Encoding - Contextualization



**Option 1: Bi-directional LSTM:**

(e.g., ELMO)

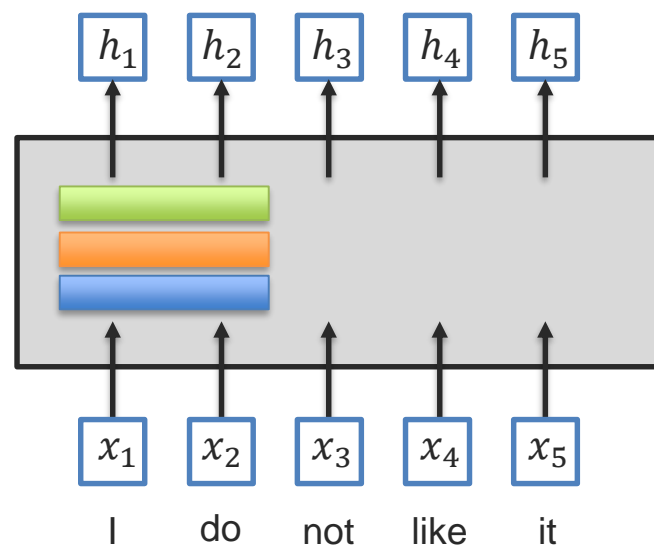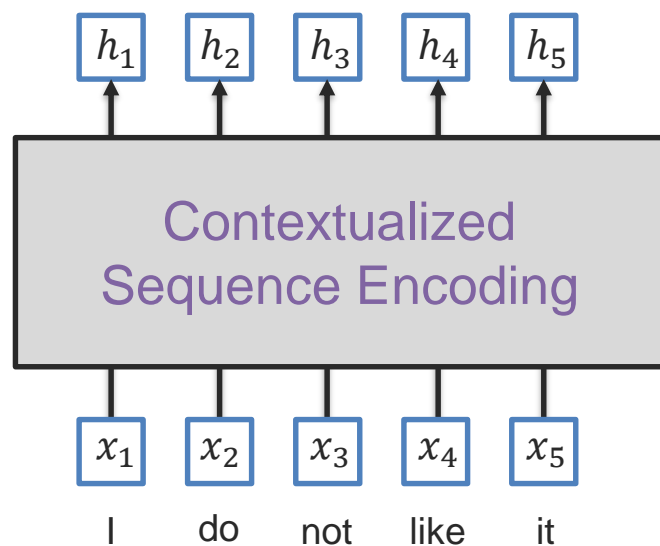How to encode this sequence while modeling the interaction between elements (e.g., words)?

But harder to parallelize…

# Sequence Encoding - Contextualization
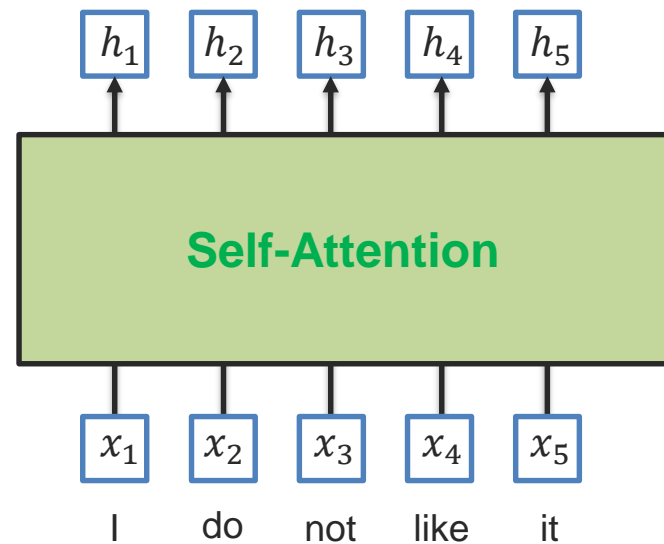
**Option 2: Convolutions**



Can be parallelized!

But modeling long-range dependencies require multiple layers

And convolutional kernels are static

Language Technologies Institute

Carnegie Mellon University

# Sequence Encoding - Contextualization

**Option 3: Self-attention**

| $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ |
|-------|-------|-------|-------|-------|

Contextualized
Sequence Encoding

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |
|-------|-------|-------|-------|-------|

I     do     not     like     it

| $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ |
|-------|-------|-------|-------|-------|

**Self-Attention**

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |
|-------|-------|-------|-------|-------|

I     do     not     like     it
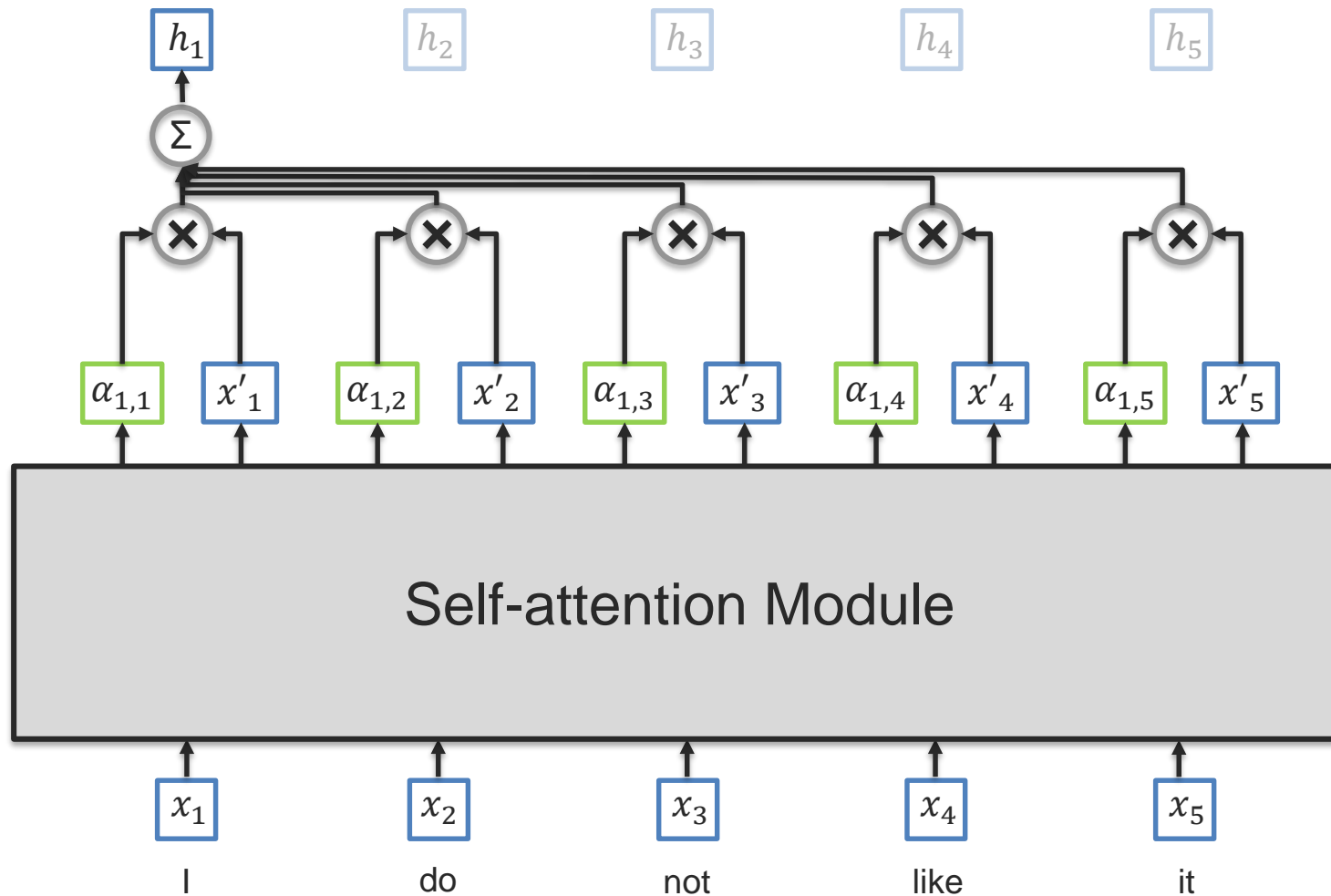
Can be parallelized!

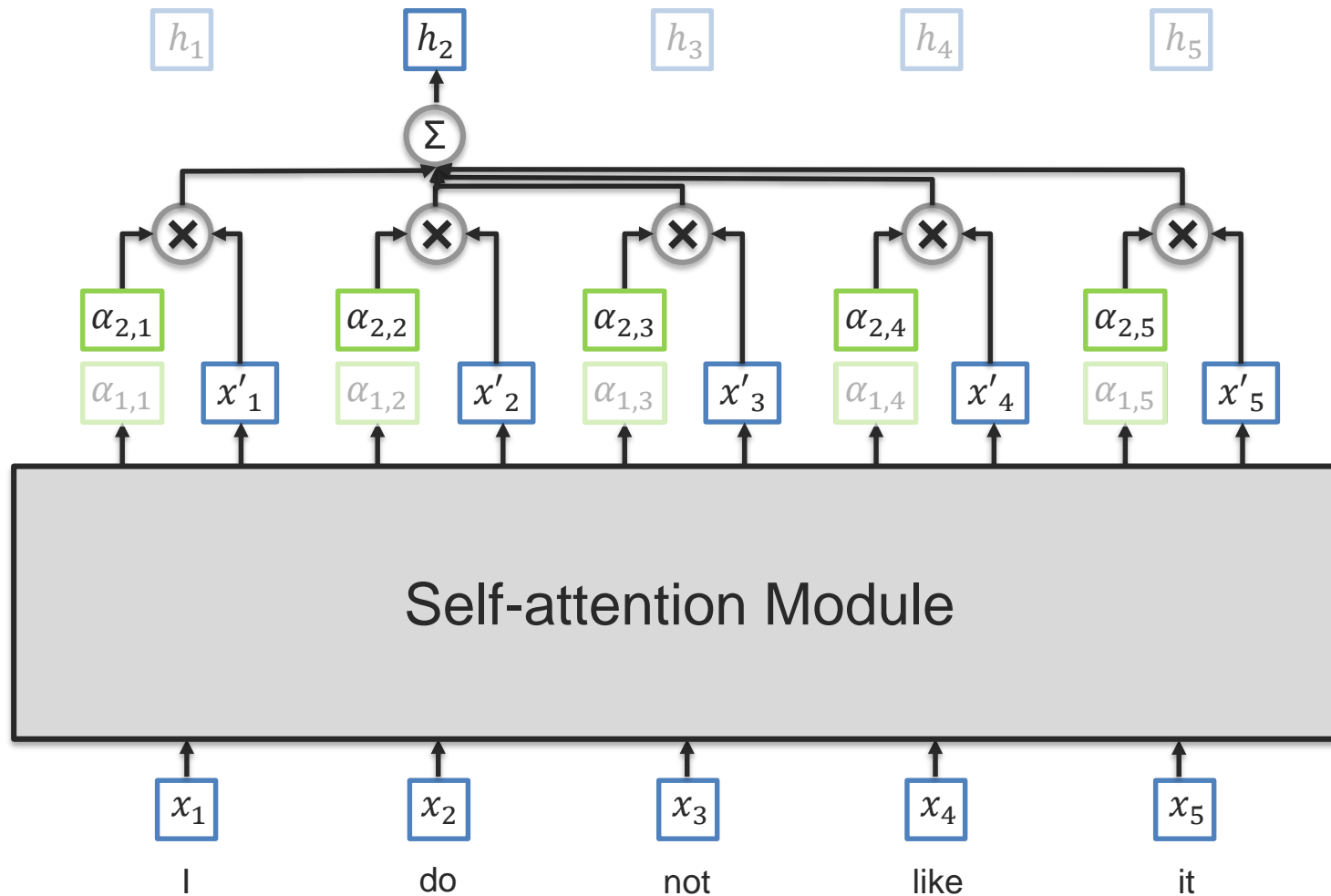Long-range dependencies

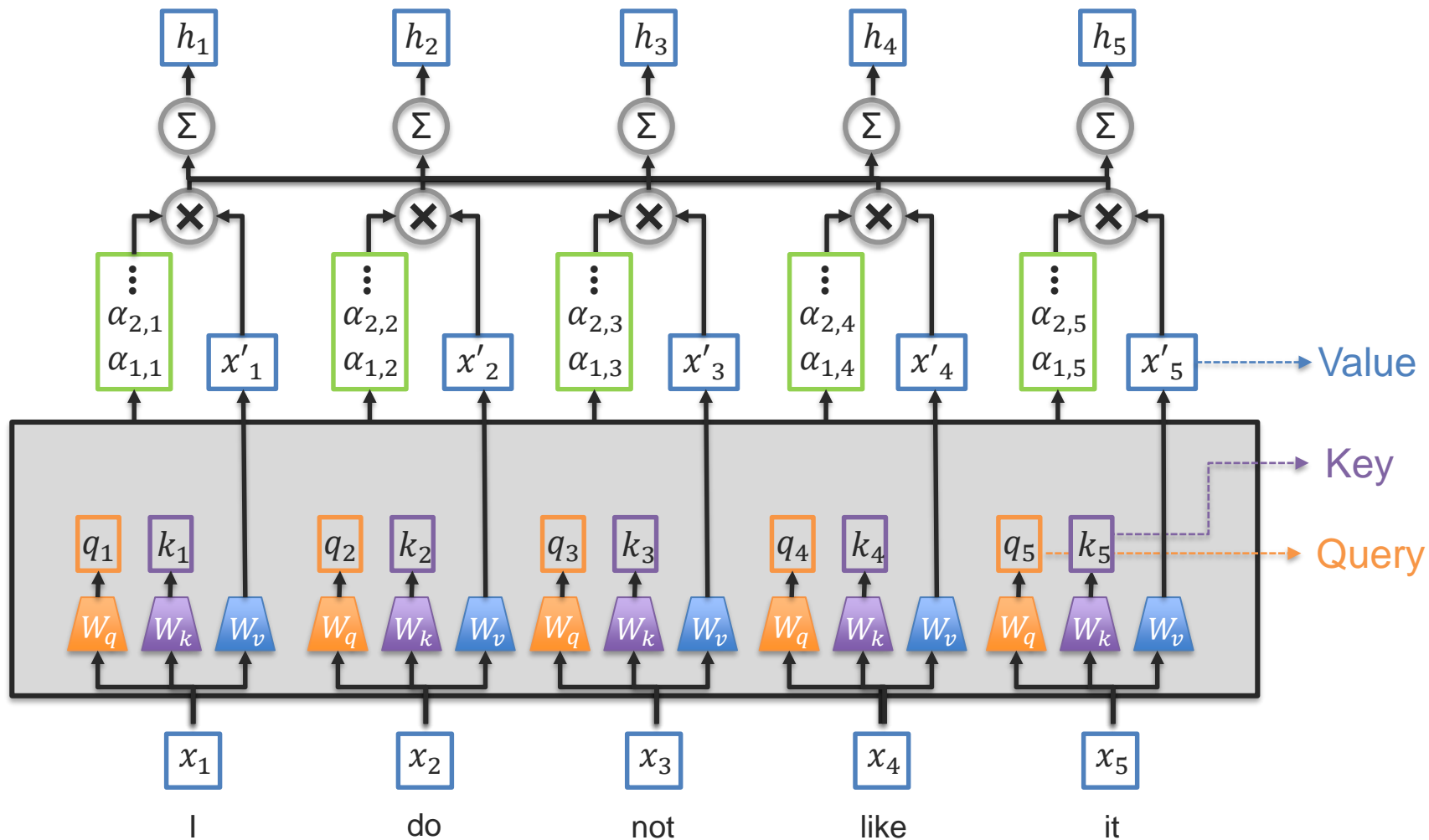Dynamic attention weights

# Self-Attention

# Self-Attention

# Self-Attention

# Transformer Self-Attention

# Transformer Self-Attention



Scale dot-product attention weights

Value

Key

Query

I    do    not    like    it

Language Technologies Institute

Carnegie Mellon University

# Transformer Self-Attention

# Transformer Self-Attention

Language Technologies Institute

Carnegie Mellon University

# Transformer Self-Attention

What if we want to attend simultaneously to multiple subspaces of $x$?

# Transformer Multi-Head Self-Attention



Linear projection

$h_1$ $h_2$ $h_3$ $h_4$ $h_5$

$h_1^1$ $h_1^2$ $h_1^3$ $h_2^1$ $h_2^2$ $h_2^3$ $h_3^1$ $h_3^2$ $h_3^3$ $h_4^1$ $h_4^2$ $h_4^3$ $h_5^1$ $h_5^2$ $h_5^3$

## Transformer's Self-Attention Layer

$W_q^1$ $W_k^1$ $W_v^1$

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$

I           do          not          like          it

Language Technologies Institute

Carnegie Mellon University

# Transformer Multi-Head Self-Attention

Language Technologies Institute

Carnegie Mellon University

# Transformer Multi-Head Self-Attention

$h_1$ $h_2$ $h_3$ $h_4$ $h_5$

Transformer's Multi-Head Self-Attention Layer

$W_q^3$ $W_k^3$ $W_v^3$

$W_q^2$ $W_k^2$ $W_v^2$

$W_q^1$ $W_k^1$ $W_v^1$

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$

not     like     I     it     do

What happens if the words are shuffled?

# Position embeddings

❑ Position information is not encoded in a self-attention module

How can we encode position information?

**Simple approach:** one-hot encoding

Language Technologies Institute

Carnegie Mellon University

# Position embeddings

❑ Position information is not encoded in a self-attention module

How can we encode position information?

**Simple approach:** one-hot encoding + linear embeddings + $\begin{cases} \text{Sum} \\ \text{- or -} \\ \text{concat} \end{cases}$

Language Technologies Institute

Carnegie Mellon University

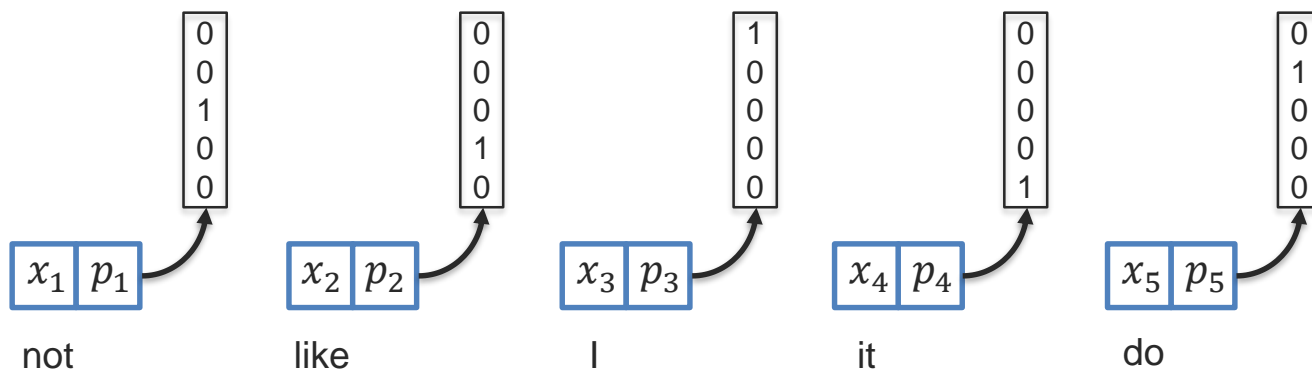# Transformer Multi-Head Self-Attention

Language Technologies Institute

Carnegie Mellon University

# Transformer Multi-Head Self-Attention

In vector format…



Transformer's Multi-Head Self-Attention Layer

$h$

$W_q^3$  $W_k^3$  $W_v^3$

$W_q^2$  $W_k^2$  $W_v^2$

$W_q^1$  $W_k^1$  $W_v^1$

$x$

$p$

# Transformer Multi-Head Attention

Language Technologies Institute

Carnegie Mellon University

# Sequence-to-Sequence Using Transformer

# Sequence-to-Sequence Modeling

Je    n'    aime    pas    cela

$\hat{y}_1$ $\hat{y}_2$ $\hat{y}_3$ $\hat{y}_4$ $\hat{y}_5$

## How can we perform seq2seq translation with transformer attention?

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$

I    do    not    like    it

Language Technologies Institute

Carnegie Mellon University

# Seq2Seq with Transformer Attentions

Je    n'    aime    pas    cela
$\hat{y}_1$  $\hat{y}_2$  $\hat{y}_3$  $\hat{y}_4$  $\hat{y}_5$

$h_1$  $h_2$  $h_3$  $h_4$  $h_5$

self-attention

$x_1$  $x_2$  $x_3$  $x_4$  $x_5$

I    do    not    like    it

Language Technologies Institute

Carnegie Mellon University

# Seq2Seq with Transformer Attentions

# Seq2Seq with Transformer Attentions

How should we connect the encoder and decoder self-attention to the transformer attention?

Je    n'    aime    pas    cela

$\hat{y}_1$   $\hat{y}_2$   $\hat{y}_3$   $\hat{y}_4$   $\hat{y}_5$

Transformer attention

$W_q$    $W_k$    $W_v$

**Query**    **Key**    **Value**

$h$

self-attention

$x_1$    $x_2$    $x_3$    $x_4$    $x_5$

I    do    not    like    it

$g$

"masked" self-attention

$y_0$    $y_1$    $y_2$    $y_3$    $y_4$

START    Je    n'    aime    pas

# Seq2Seq with Transformer Attentions

Language Technologies Institute

Carnegie Mellon University

# Contextualized Multimodal Embedding

# Multimodal Embeddings

| $h^L$ | $h^V$ | $h^A$ |
|-------|-------|-------|

How to learn contextualized
representations from multiple modalities?

I really liked it this time

**Language**

**Visual**

**Acoustic**

Language Technologies Institute

Carnegie Mellon University

# Contextualized Multimodal Embeddings

$h^L$     $h^V$     $h^A$

Transformer self-attention

I really liked it this time

Language     Visual     Acoustic

Any other approach?

# Multimodal Transformer



Tsai et al., Multimodal Transformer for Unaligned Multimodal Language Sequences, ACL 2019

Language Technologies Institute

Carnegie Mellon University

# Cross-Modal Transformer

$$\text{softmax}\left(\frac{Q_\alpha K_\beta^\top}{\sqrt{d_k}}\right)V_\beta \in \mathbb{R}^{T_\alpha \times d_v}$$

$\text{CM}_{\beta \to \alpha}(X_\alpha, X_\beta)$

$$\text{softmax}\left(\frac{Q_\alpha K_\beta^\top}{\sqrt{d_k}}\right)$$

$Q_\alpha \in \mathbb{R}^{T_\alpha \times d_k}$

$K_\beta \in \mathbb{R}^{T_\beta \times d_k}$

$V_\beta \in \mathbb{R}^{T_\beta \times d_v}$

$W_{Q_\alpha}$

$W_{K_\beta}$

$W_{V_\beta}$

$X_\alpha \in \mathbb{R}^{T_\alpha \times d_\alpha}$

$X_\beta \in \mathbb{R}^{T_\beta \times d_\alpha}$

Modality $\alpha$

Modality $\beta$

Tsai et al., Multimodal Transformer for Unaligned Multimodal Language Sequences, ACL 2019

Language Technologies Institute

Carnegie Mellon University

# Language Pre-training

# Token-level and Sentence-level Embeddings

Token-level embeddings

Sentence-level embedding



Which tasks?

Which tasks?

Language Technologies Institute

Carnegie Mellon University

# Pre-Training and Fine-Tuning



**Pre-training**

(e.g., language model)

**Fine-Tuning**

# BERT:
# Bidirectional Encoder Representations from Transformers

**Advantages:**

1. Jointly learn representation for token-level and sentence level

2. Same network architecture for pre-training and fine-tuning

Carnegie Mellon University

# BERT:
# Bidirectional Encoder Representations from Transformers

**Advantages:**

1    Jointly learn representation for token-level and sentence level

2    Same network architecture for pre-training and fine-tuning

3    Can be used learn relationship between sentences

4    Models bidirectional and long-range interactions between tokens

How can we do all this?

| $h_s$ | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ | $h_{sep}$ | $h'_1$ | $h'_2$ | $h'_3$ | $h'_4$ | $h'_5$ |

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | sep | $x'_1$ | $x'_2$ | $x'_3$ | $x'_4$ | $x'_5$ |

I    do    not    like    it       I    enjoy    my    time    here

Language Technologies Institute
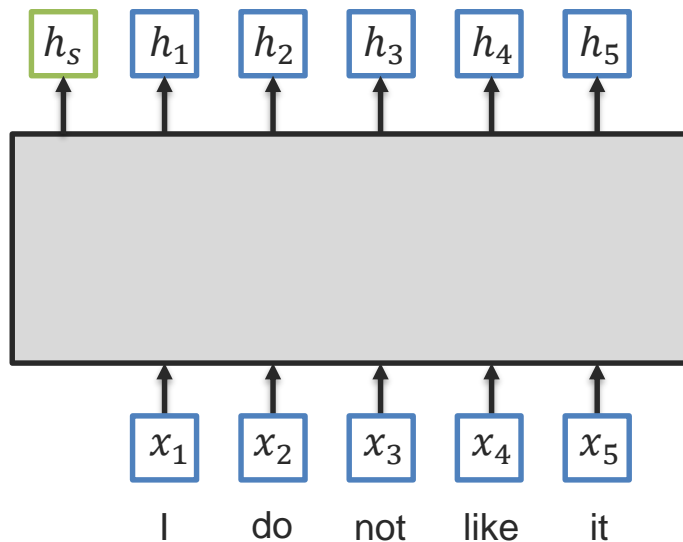
Carnegie Mellon University

# BERT:
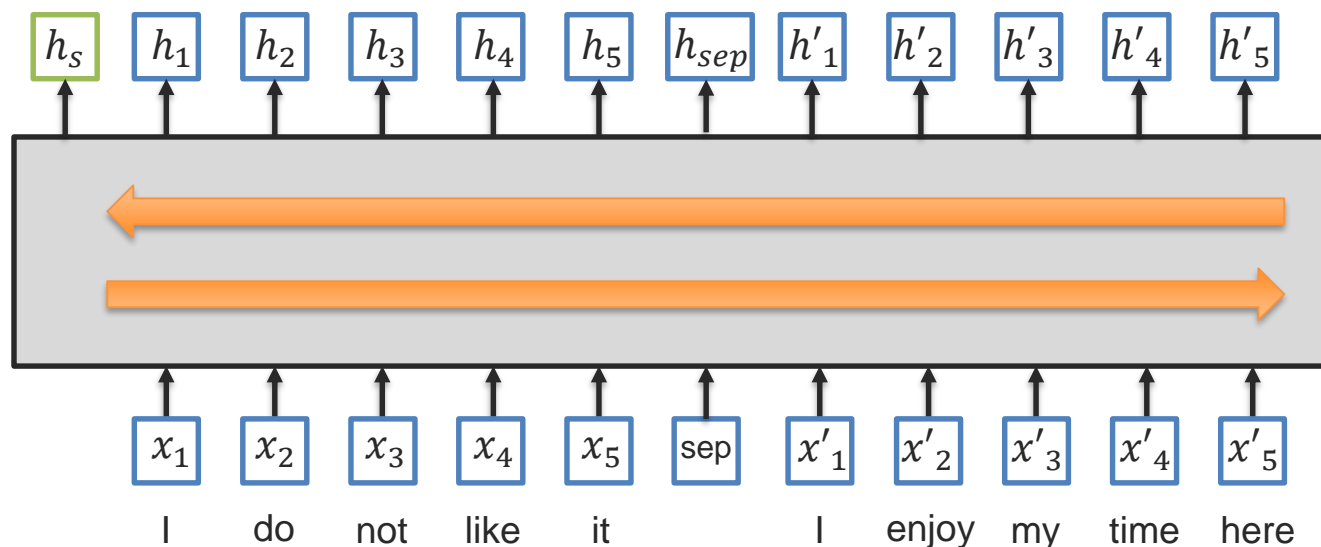# Bidirectional Encoder Representations from Transformers

**Advantages:**

**1** Jointly learn representation for token-level and sentence level

**2** Same network architecture for pre-training and fine-tuning

**3** Can be used learn relationship between sentences

**4** Models bidirectional interactions between tokens

| $h_s$ | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ | $h_{sep}$ | $h'_1$ | $h'_2$ | $h'_3$ | $h'_4$ | $h'_5$ |

## Transformer Self-Attention

Special sentence-level token

But how to train unsupervised?

| cls | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | sep | $x'_1$ | $x'_2$ | $x'_3$ | $x'_4$ | $x'_5$ |

I do not like it    I enjoy my time here

# Pre-training BERT Model

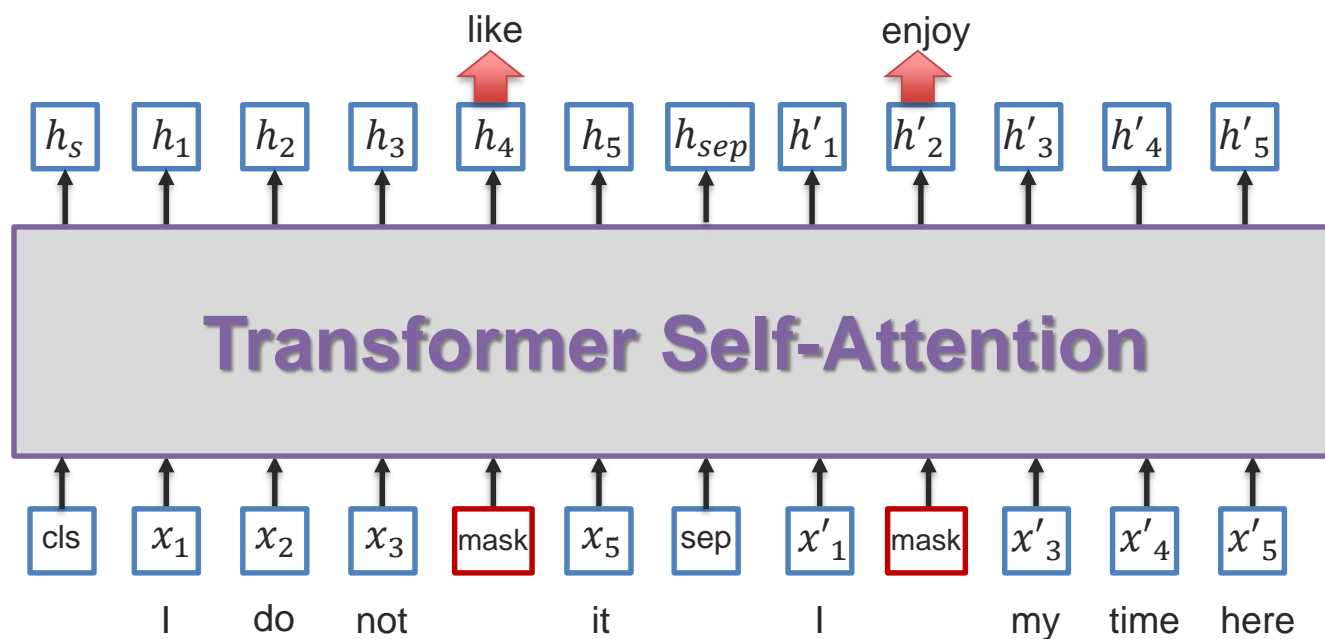**1** **Masked Language Model**

Randomly mask input tokens and then try to predict them

What is the loss function?

# Pre-training BERT Model

**(2)  Next Sentence Prediction**

Given two sentences, predict if this is the next one or not

What is the loss function?

Where can we find training data?

How can BERT know the difference between both sentences?

IsNext
**- or -**
NotNext

| $h_s$ | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ | $h_{sep}$ | $h'_1$ | $h'_2$ | $h'_3$ | $h'_4$ | $h'_5$ |

## Transformer Self-Attention

| cls | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | sep | $x'_1$ | $x'_2$ | $x'_3$ | $x'_4$ | $x'_5$ |

I    do   not   like   it        I   enjoy   my   time   here

# Three Embeddings: Token + Position + Sentence

| Input | [CLS] | my | dog | is | cute | [SEP] | he | likes | play | ##ing | [SEP] |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Token Embeddings | $E_{[CLS]}$ | $E_{my}$ | $E_{dog}$ | $E_{is}$ | $E_{cute}$ | $E_{[SEP]}$ | $E_{he}$ | $E_{likes}$ | $E_{play}$ | $E_{\#\#ing}$ | $E_{[SEP]}$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Segment Embeddings | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Position Embeddings | $E_0$ | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ | $E_6$ | $E_7$ | $E_8$ | $E_9$ | $E_{10}$ |

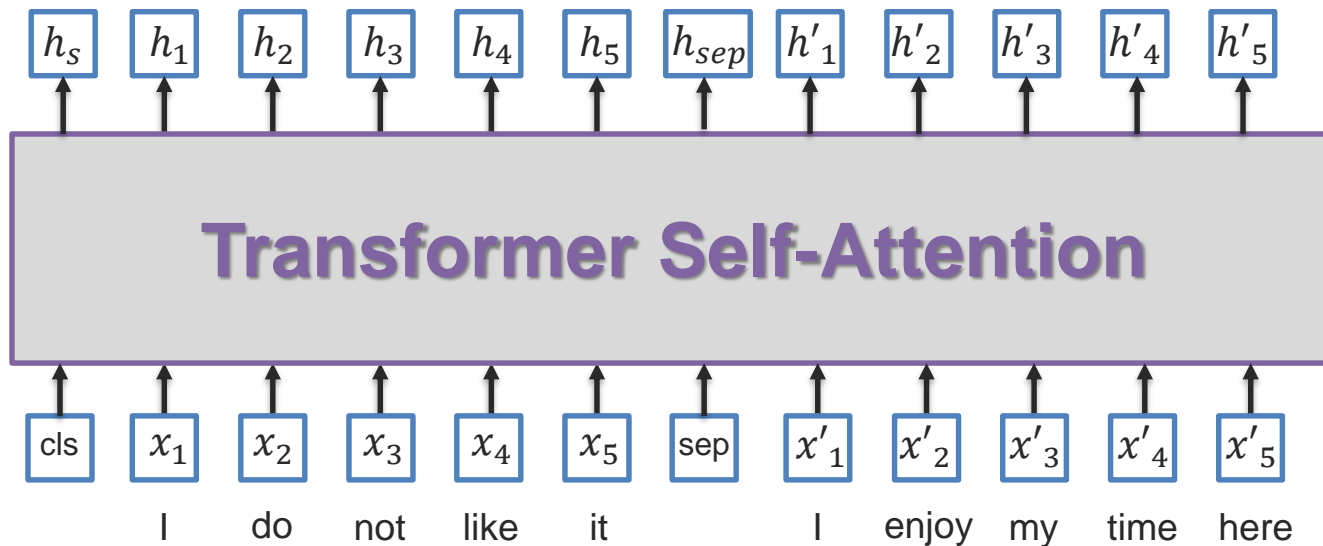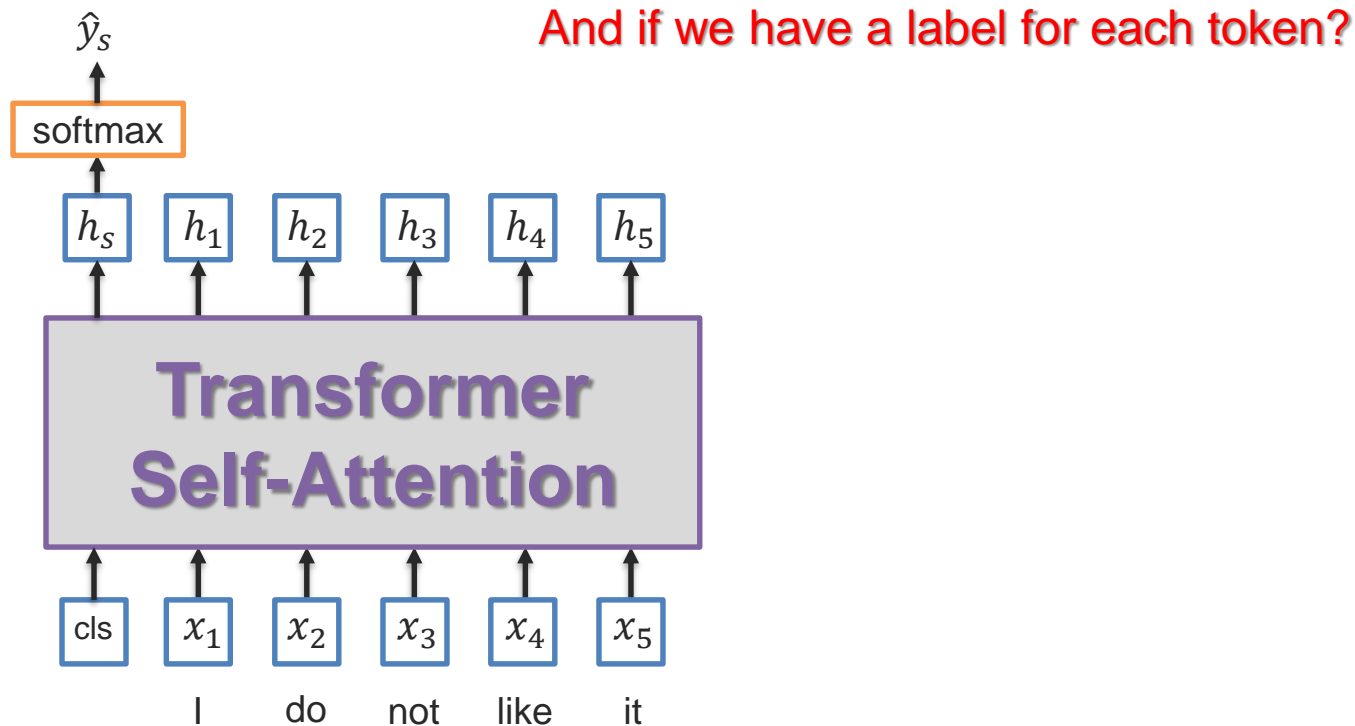Language Technologies Institute

Carnegie Mellon University

# Fine-Tuning BERT

① Sentence-level classification for only one sentence

Examples: sentiment analysis, document classification

How?

# Fine-Tuning BERT

① Sentence-level classification for only one sentence

Examples: sentiment analysis, document classification

**And if we have a label for each token?**

$\hat{y}_s$

| softmax |

| $h_s$ | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ |

**Transformer Self-Attention**

| cls | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |

I    do    not    like    it

# Fine-Tuning BERT

② Token-level classification for only one sentence

Examples: part-of-speech tagging, slot filling
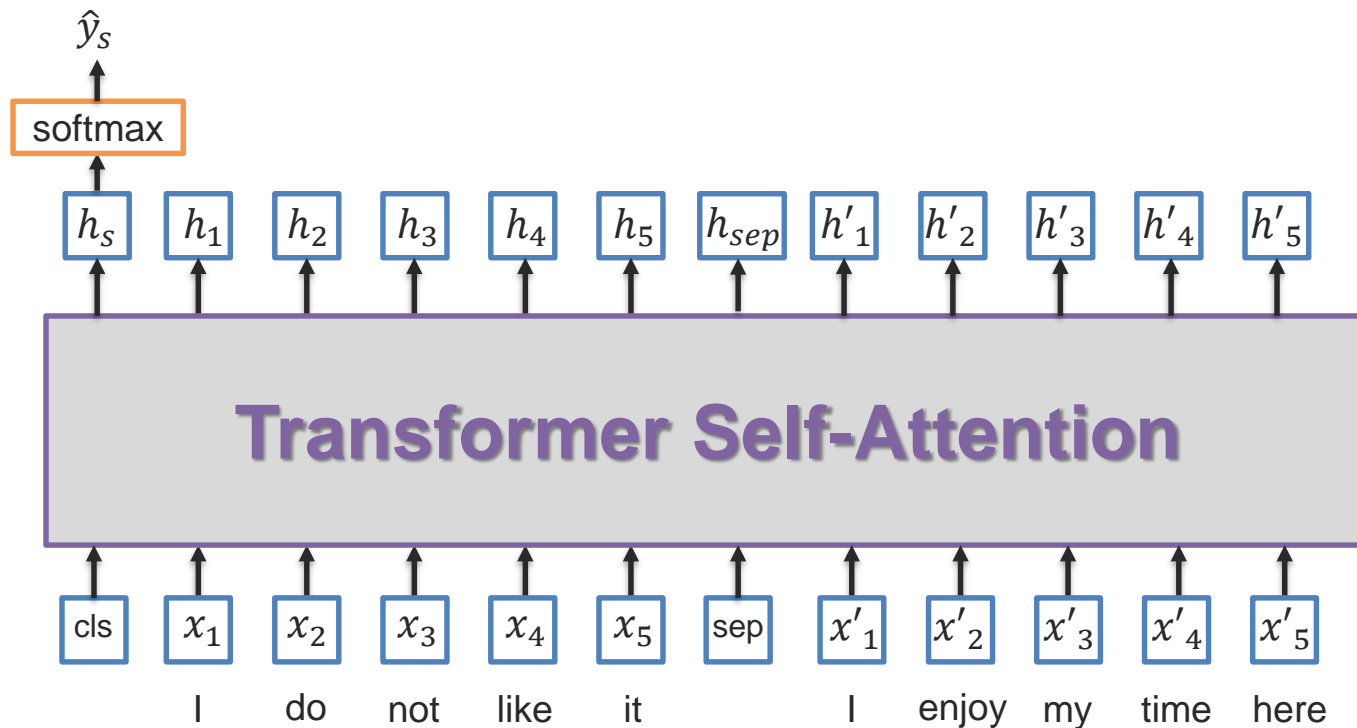


How to compare two sentences?

Language Technologies Institute

Carnegie Mellon University

# Fine-Tuning BERT

③ Sentence-level classification for two sentences

Examples: natural language inference

$\hat{y}_s$

| softmax |

| $h_s$ | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ | $h_{sep}$ | $h'_1$ | $h'_2$ | $h'_3$ | $h'_4$ | $h'_5$ |

## Transformer Self-Attention

| cls | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | sep | $x'_1$ | $x'_2$ | $x'_3$ | $x'_4$ | $x'_5$ |

I do not like it I enjoy my time here

# Fine-Tuning BERT

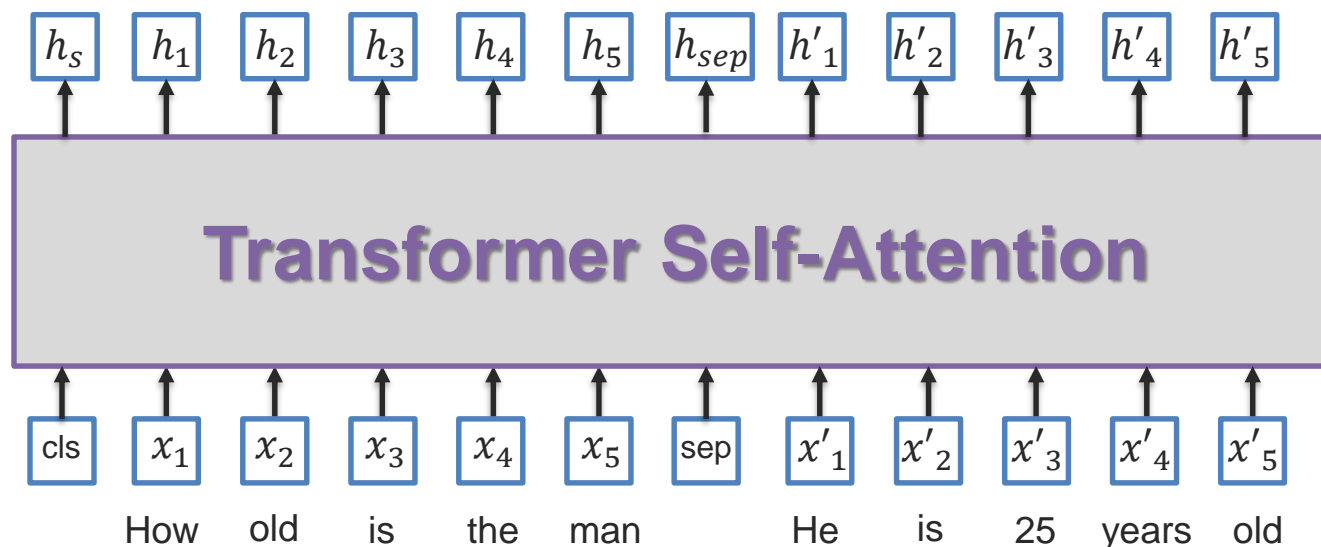**(4)** Question-answering: find start/end of the answer in the document

**Paragraph:** " ... *Other legislation followed, including the Migratory Bird Conservation Act of 1929, a 1937 treaty prohibiting the hunting of right and gray whales, and the Bald Eagle Protection Act of 1940. These later laws had a low cost to society—the species were relatively rare—and little opposition was raised.*"

**Question 1:** "*Which laws faced significant opposition?*"
**Plausible Answer:** *later laws*

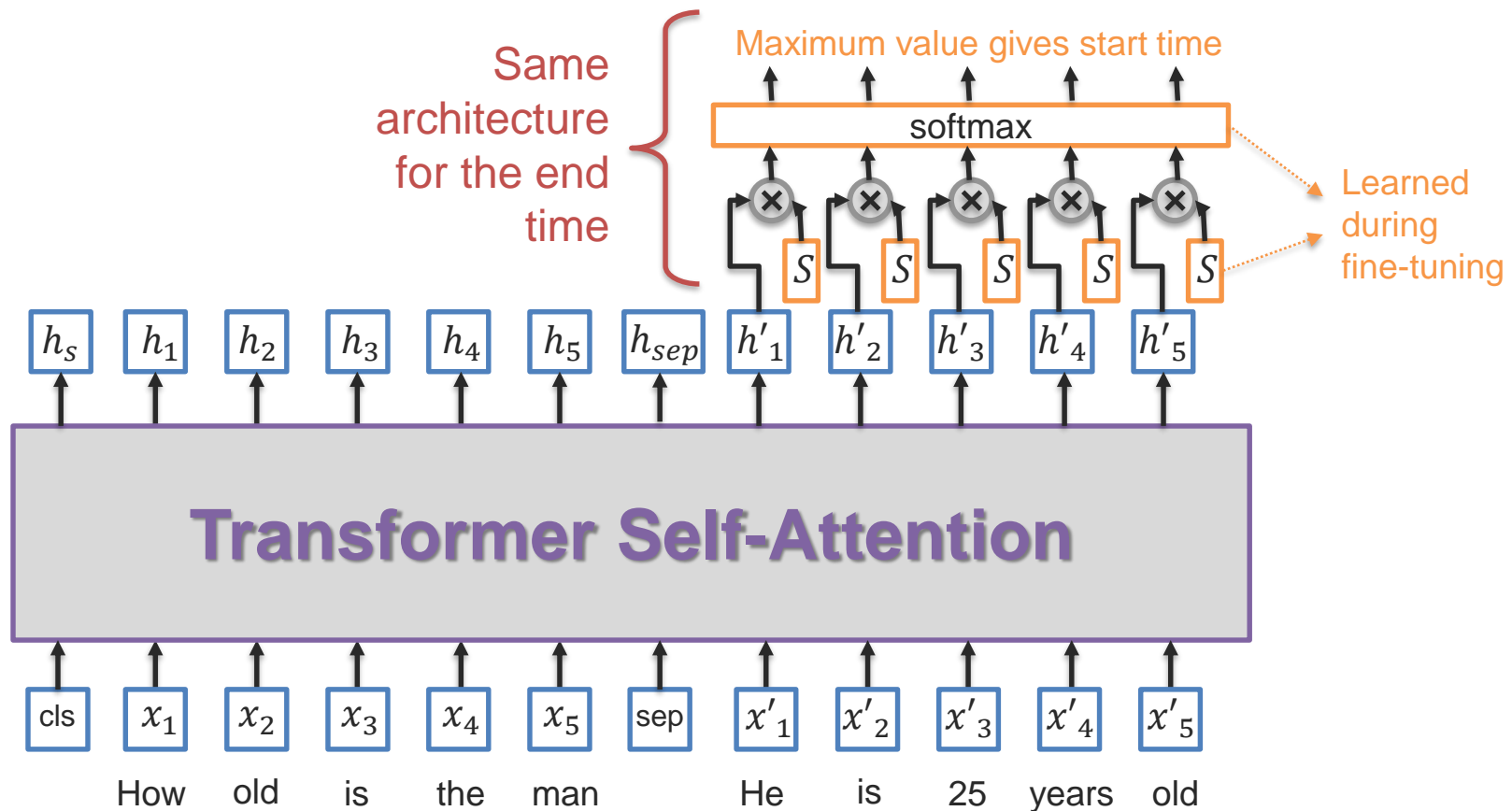**Question 2:** "*What was the name of the 1937 treaty?*"
**Plausible Answer:** *Bald Eagle Protection Act*

| $h_s$ | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ | $h_{sep}$ | $h'_1$ | $h'_2$ | $h'_3$ | $h'_4$ | $h'_5$ |

## Transformer Self-Attention

How?

| cls | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | sep | $x'_1$ | $x'_2$ | $x'_3$ | $x'_4$ | $x'_5$ |

How   old   is   the   man       He   is   25   years   old

# Fine-Tuning BERT

Maximum value gives start time

softmax

Same architecture for the end time

Learned during fine-tuning

$h_s$ $h_1$ $h_2$ $h_3$ $h_4$ $h_5$ $h_{sep}$ $h'_1$ $h'_2$ $h'_3$ $h'_4$ $h'_5$

**Transformer Self-Attention**

cls $x_1$ $x_2$ $x_3$ $x_4$ $x_5$ sep $x'_1$ $x'_2$ $x'_3$ $x'_4$ $x'_5$

How old is the man    He is 25 years old

# Multimodal Pre-training

Language Technologies Institute

Carnegie Mellon University

# Multimodal Pre-Training

How to extend to multimodal modalities?

Language Technologies Institute

Carnegie Mellon University

## How to extend to multimodal modalities?

**Option 1:** Simply concatenate tokens from different modalities



https://arxiv.org/pdf/1908.08530.pdf

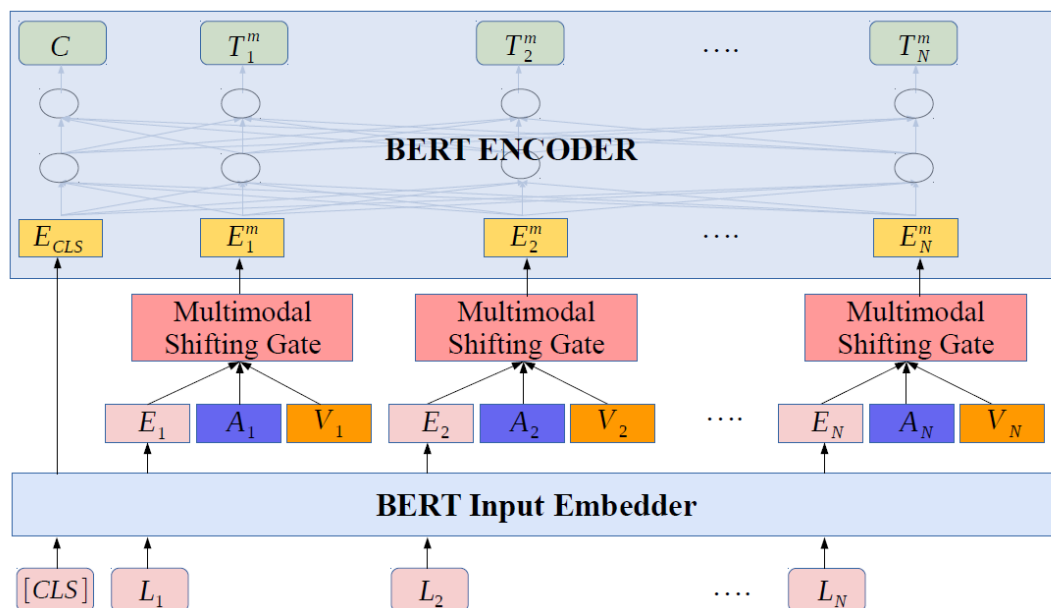Language Technologies Institute
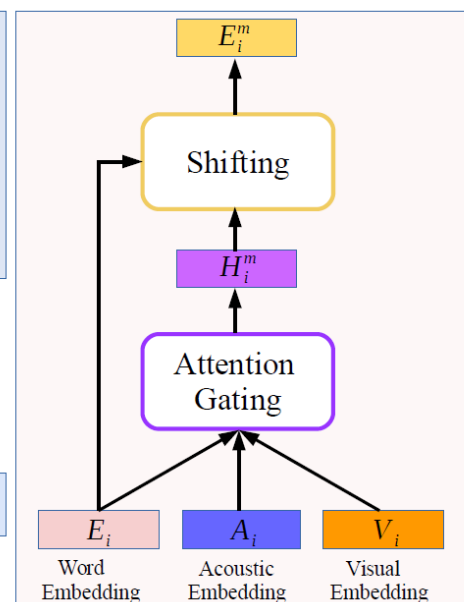
Carnegie Mellon University

# M-BERT

## How to extend to multimodal modalities?

**Option 2:** "Shift" language representation based on the other modalities



(a) Multimodal BERT

(b) Multimodal Shifting Gate

https://arxiv.org/pdf/1908.05787.pdf

Language Technologies Institute

Carnegie Mellon University