# Contextualized Sequence Representations

# Sequence Encoding - Contextualization

Option 1: Bi-directional LSTM:
(e.g., ELMO)

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

Contextualized
Sequence Encoding

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

I    do    not    like    it

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

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

I    do    not    like    it

Sequential Computation

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

But harder to parallelize…

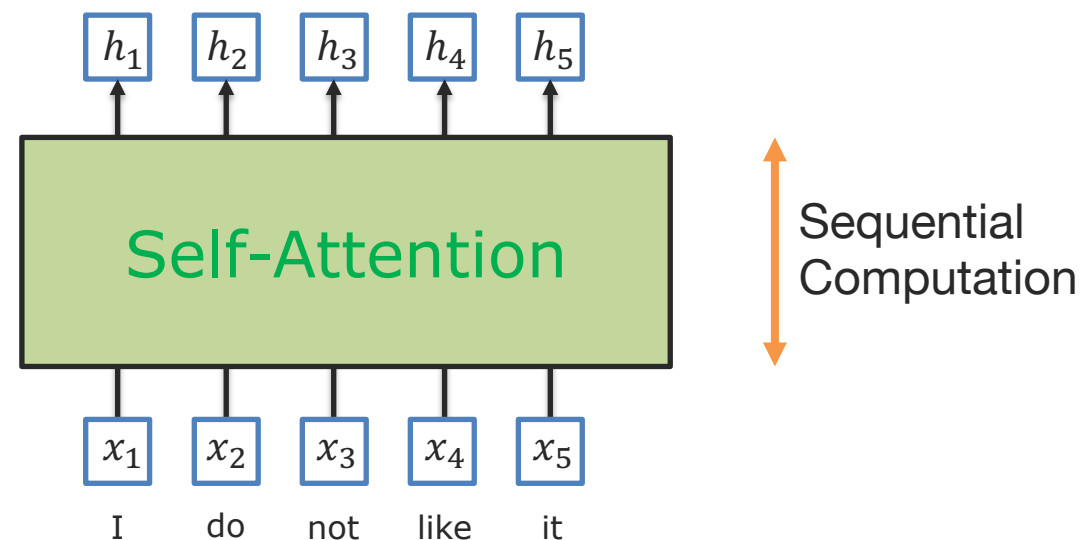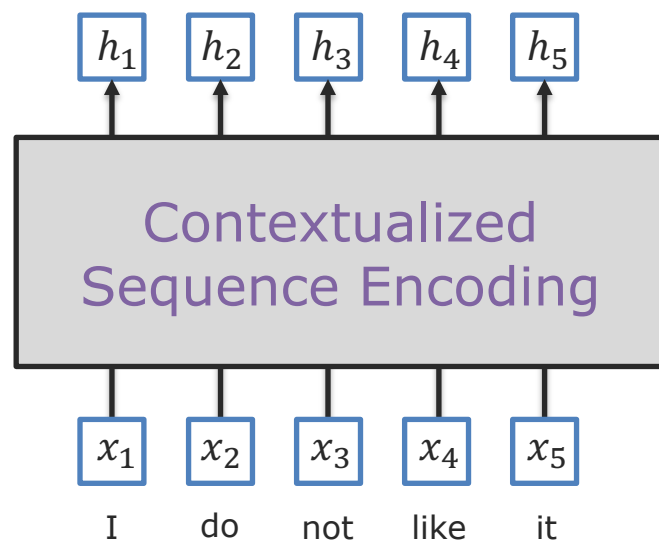# Sequence Encoding - Contextualization

Option 2: Convolutions



Sequential Computation

Can be parallelized!

But modeling long-range dependencies requires multiple layers.

And convolutional kernels are static.

Carnegie Mellon University

# Sequence Encoding - Contextualization

Option 3: Self-attention



Sequential Computation

Can be parallelized!
Long-range dependencies
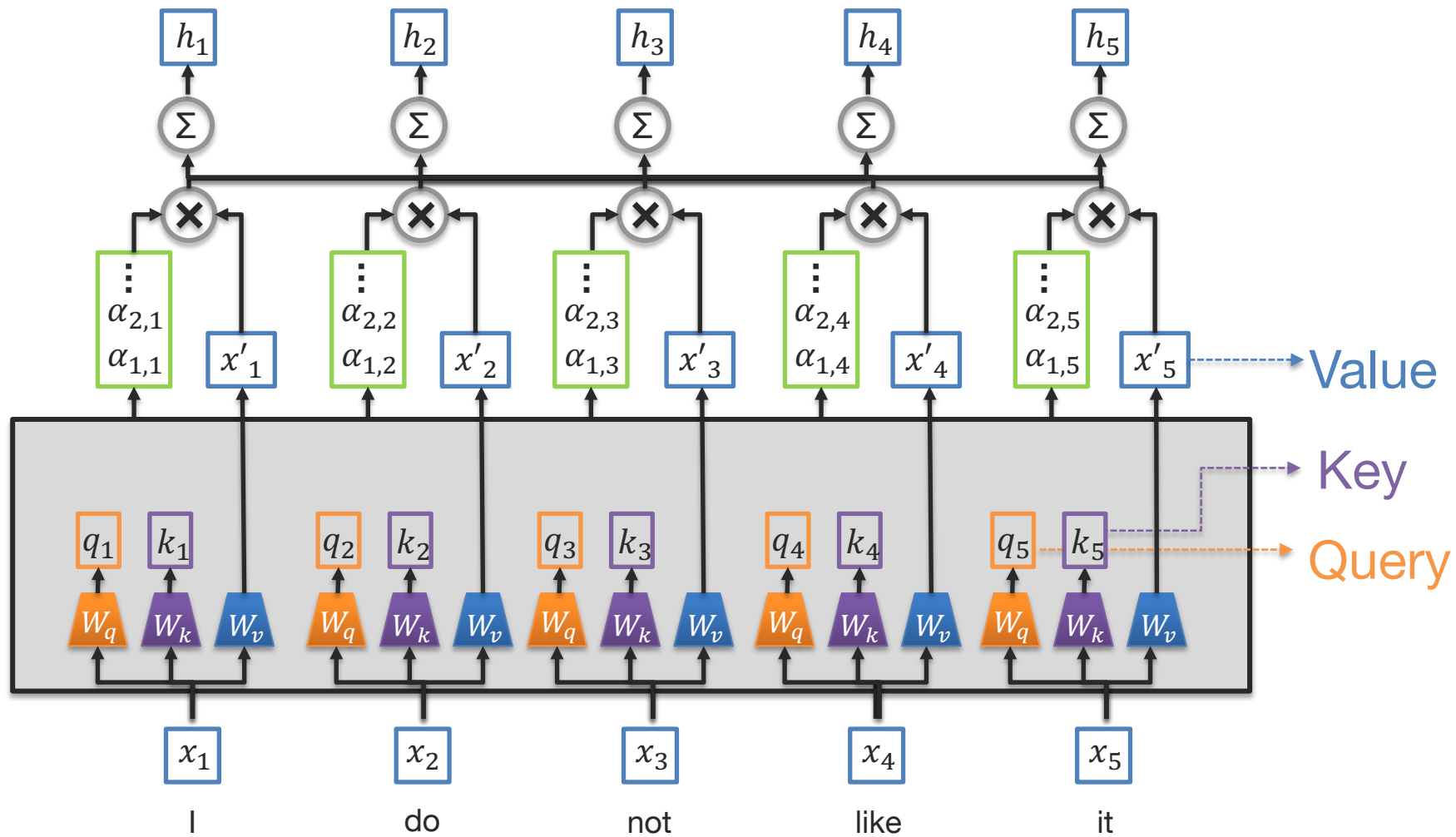Dynamic attention weights

Carnegie Mellon University

# Self-Attention

# Self-Attention

# Self-Attention

# Transformer Self-Attention

Carnegie Mellon University

# Transformer Self-Attention



Scale dot-product attention weights

Value

Key

Query

I    do    not    like    it

# Transformer Self-Attention



Scale dot-product attention weights

Value

Key

Query

I     do     not     like     it

Carnegie Mellon University

# Transformer Self-Attention

# Transformer Self-Attention

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

Carnegie Mellon University

# Transformer Multi-Head Self-Attention



Linear projection

Transformer's Self-Attention Layer

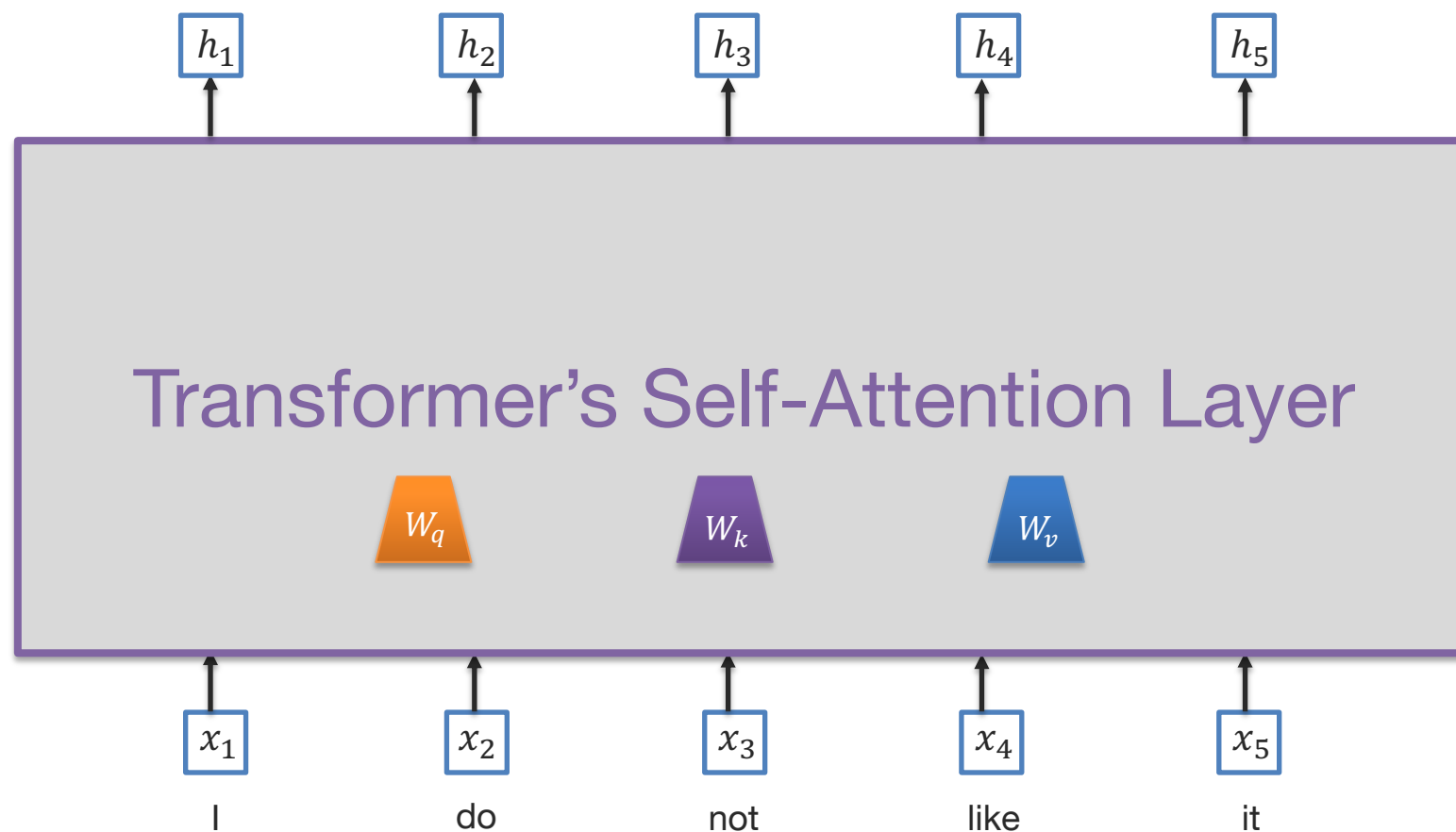$W_q^1$   $W_k^1$   $W_v^1$

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

I   do   not   like   it

Carnegie Mellon University

# Transformer Multi-Head Self-Attention

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
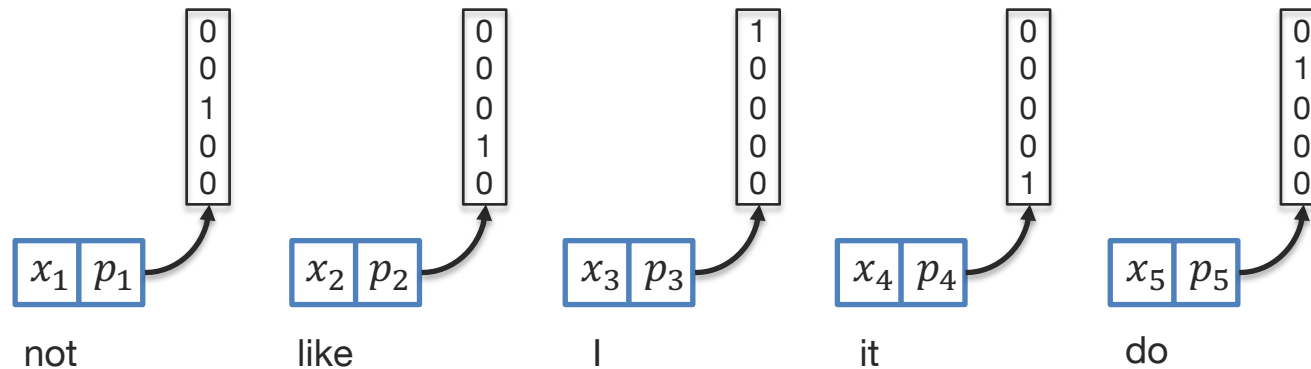
<span style="color:red">How can we encode position information?</span>

**Simple approach:** one-hot encoding



not  like  I  it  do

Carnegie Mellon University

# Position embeddings

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

How can we encode position information?
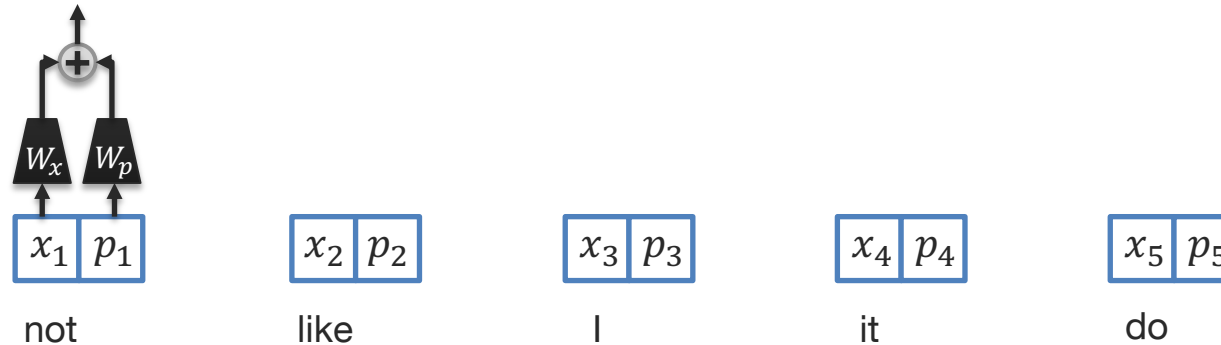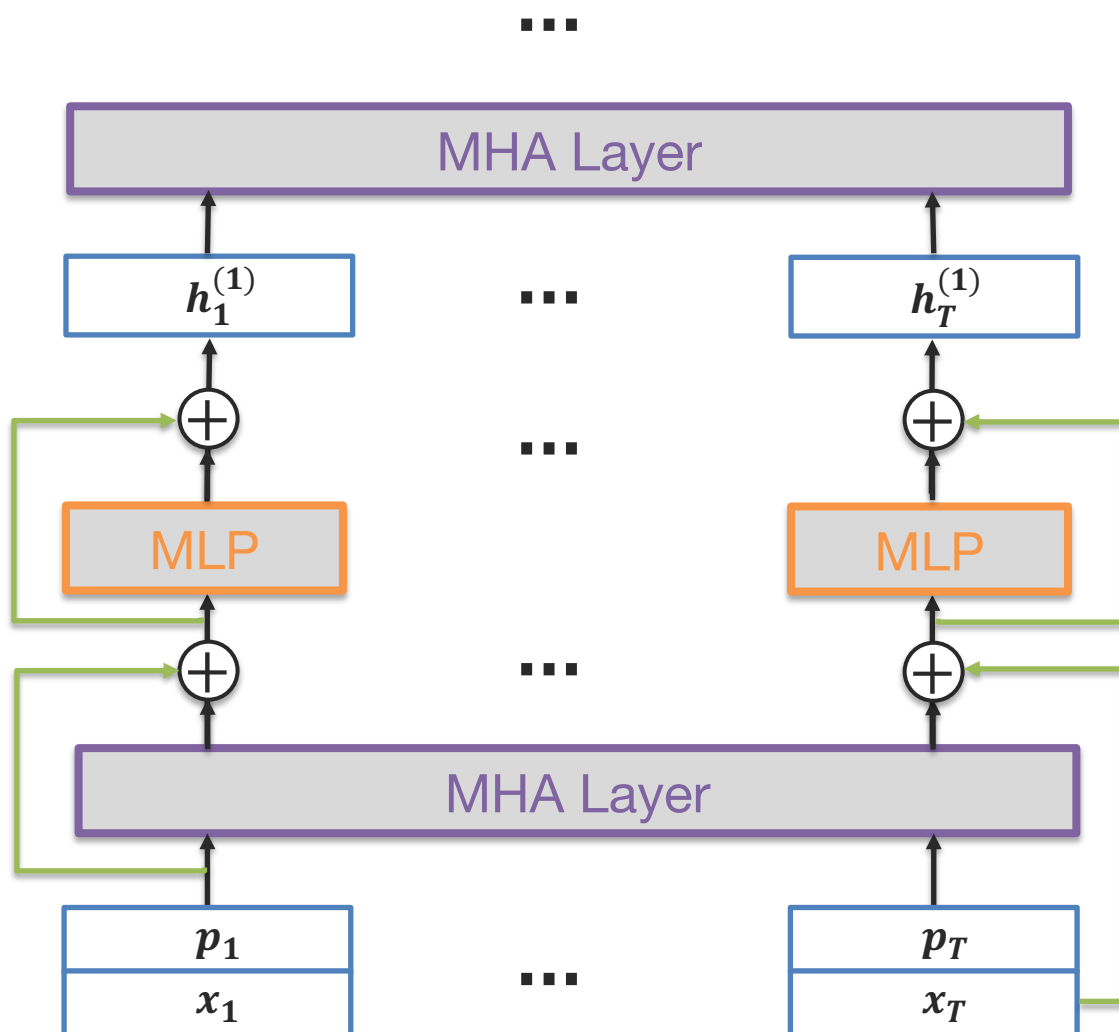
**Simple approach:** one-hot encoding + linear embeddings $+\begin{cases} \text{Sum} \\ \text{- or -} \\ \text{concat} \end{cases}$

**Carnegie Mellon University**

# Transformer – Layers

...



**Multi-head attention** is basically linear

=> need **MLPs**, and many layers (12 in BERT, 96 in GPT-3)

**Residual connections** allow each layer to update only a small subspace of the hidden space

=> more than just contextualized word vectors

# Transformer Key Ideas

1.  Learn to route information, then apply local computation to compose it.

2.  No fixed-weight similarity kernels; attention is input-dependent similarity

3.  Compose multiple simple learned operations

4.  Highly parallelizable; allows data to trump inductive biases

For more intuition: https://transformer-circuits.pub/2021/framework/index.html

**Carnegie Mellon University**