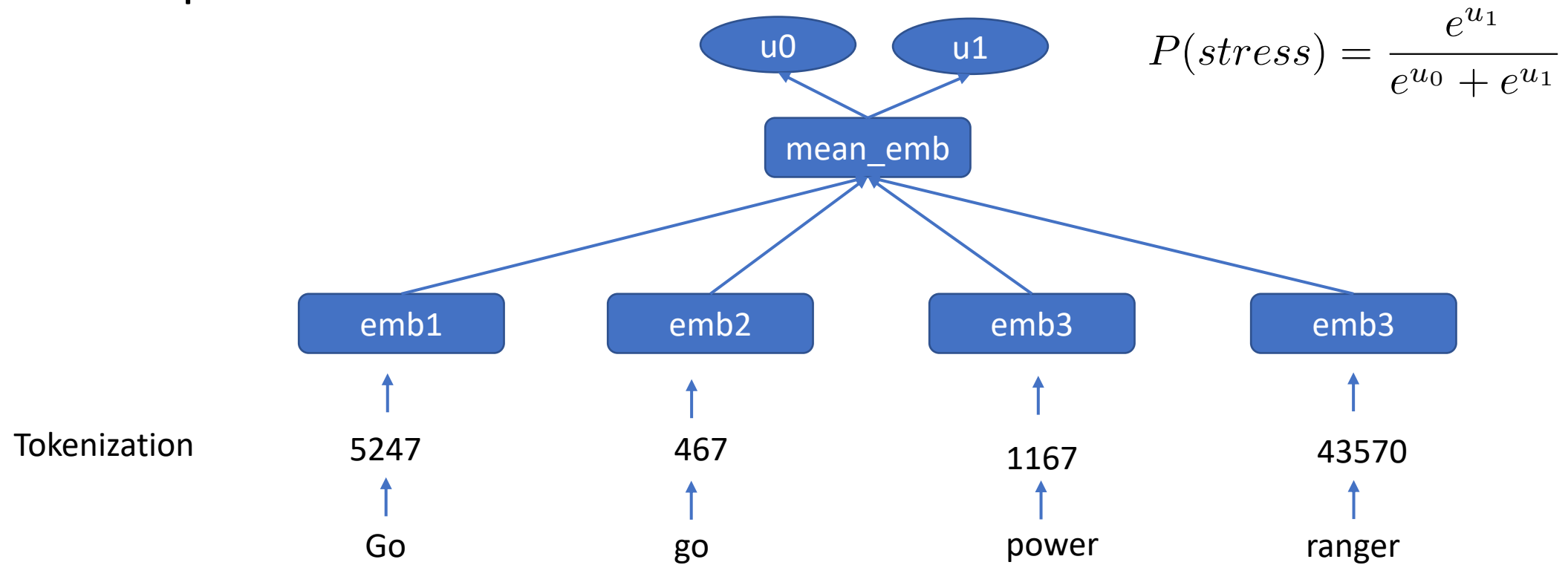# Gen-AI: Technical and Social

Lecture 03: Transformer and Text Generation

# Recap: Text Classification

- A simple model



$$P(stress) = \frac{e^{u_1}}{e^{u_0} + e^{u_1}}$$

# Recap: Text Classification

- A simple model

```python
class EmbeddingModel(nn.Module):
    def __init__(self, vocab_size, embedding_dim, hidden_dim, num_classes):
        super().__init__()
        self.embedding = nn.Embedding(vocab_size, embedding_dim)
        self.fc1 = nn.Linear(embedding_dim, num_classes)

    def forward(self, inputs):
        output = self.embedding(inputs)
        output = torch.mean(output, dim=1)   # Mean pooling
        output = self.fc1(output)
        return output

model = EmbeddingModel(tokenizer.vocab_size, 128, 128, 2)
```

# Issues of Word Embeddings?

- 1. Does not take the surrounding contexts into account
- 2. Does not take care of the order of the words

# Context matters

| Word | Example Contexts |
|------|------------------|
| it | The animal didn't cross the street because it was too tired<br>The animal didn't cross the street because it was too wide |
| station | The train left the station on time<br>The radio station was playing 60s hits<br>I was stationed on a remote island in Polynesia |

# Order matter

She only eats pizza  -> Only she eats pizza
The dog chased the cat -> The cat chased the dog.
….

# Transformer

By failing to prepare, you are preparing to fail.                    ---Benjamin Franklin
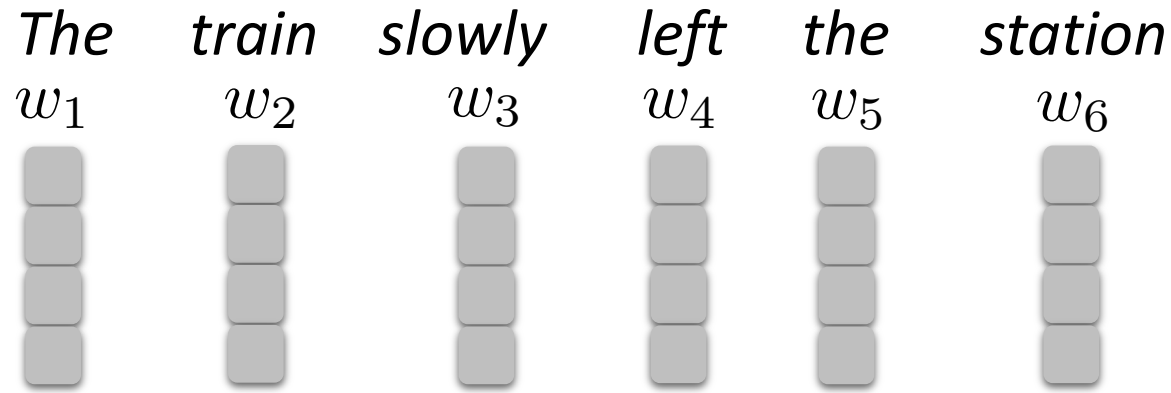
Let's prepare a right architecture…

# We will focus on this first

**How to take the surrounding context of each word into account**

How to take the order of the words into account

# From embeddings to *contextual embeddings*
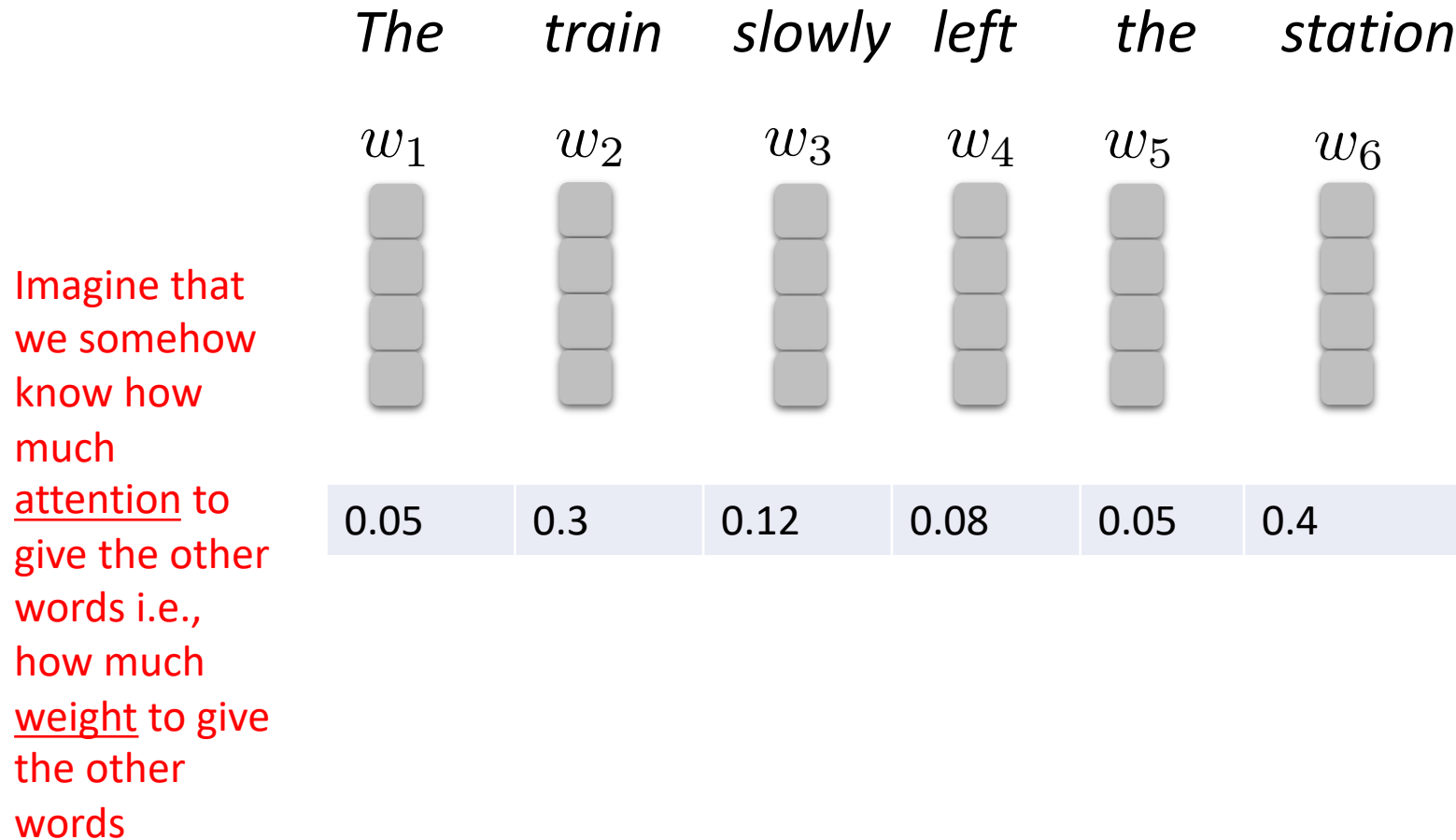
| *The* | *train* | *slowly* | *left* | *the* | *station* |
|-------|---------|----------|--------|-------|-----------|
| $w_1$ | $w_2$ | $w_3$ | $w_4$ | $w_5$ | $w_6$ |

- We can easily get  stand-alone embeddings for all the words

# From embeddings to *contextual embeddings*

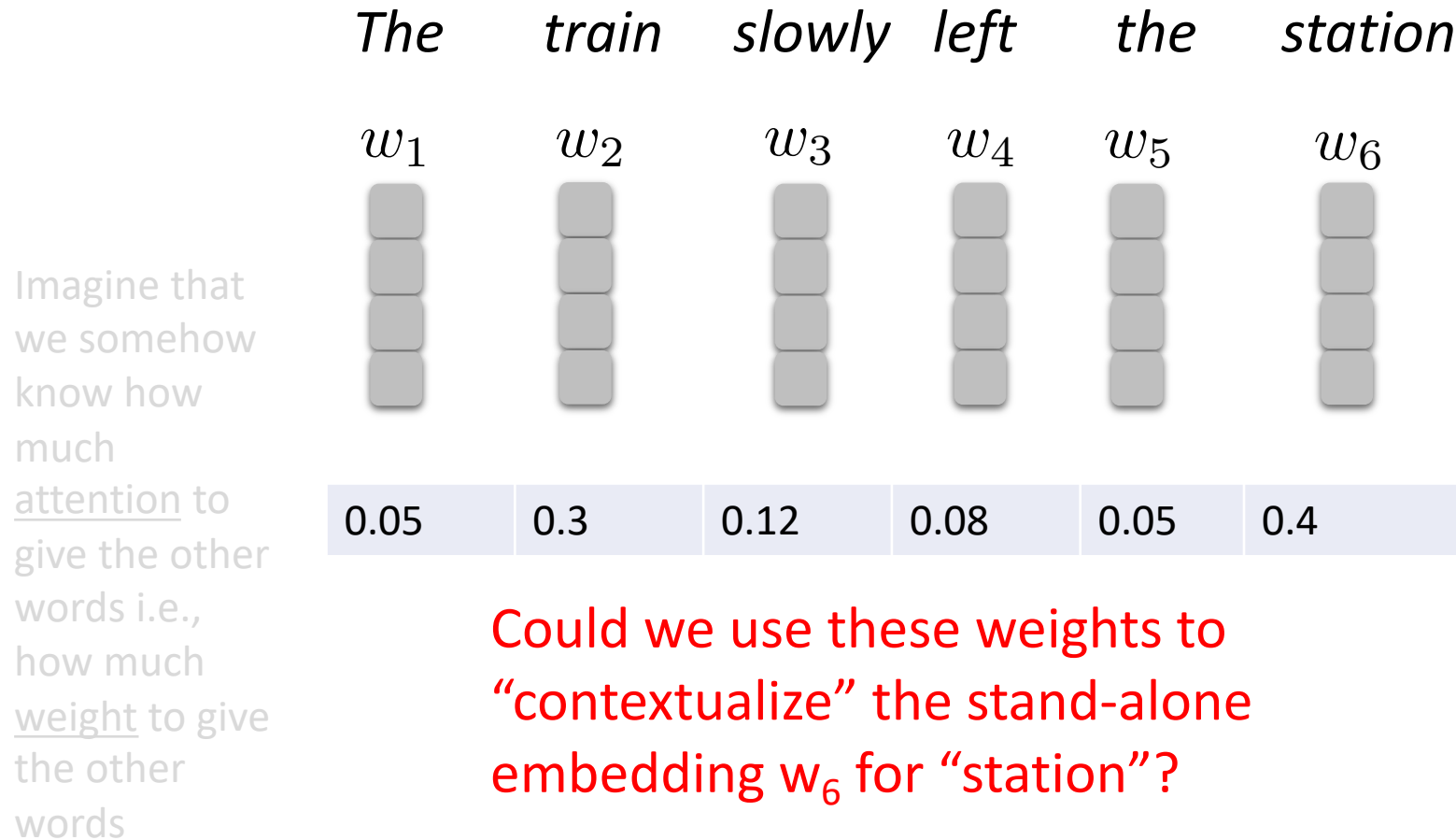| *The* | *train* | *slowly* | *left* | *the* | *station* |
|-------|---------|----------|--------|-------|-----------|
| $w_1$ | $w_2$ | $w_3$ | $w_4$ | $w_5$ | $w_6$ |

- We can easily get stand-alone embeddings for all the words

- **How can we modify station's embedding so that it incorporates the other words?**

# From embeddings to *contextual embeddings*

| The | train | slowly | left | the | station |
|-----|-------|--------|------|-----|---------|
| $w_1$ | $w_2$ | $w_3$ | $w_4$ | $w_5$ | $w_6$ |

Imagine that we somehow know how much <u>attention</u> to give the other words i.e., how much <u>weight</u> to give the other words

| 0.05 | 0.3 | 0.12 | 0.08 | 0.05 | 0.4 |

# From embeddings to *contextual embeddings*

*The*     *train*    *slowly*   *left*    *the*     *station*

$w_1$      $w_2$      $w_3$      $w_4$      $w_5$      $w_6$

Imagine that we somehow know how much attention to give the other words i.e., how much weight to give the other words

| 0.05 | 0.3 | 0.12 | 0.08 | 0.05 | 0.4 |

Could we use these weights to "contextualize" the stand-alone embedding $w_6$ for "station"?

# From embeddings to *contextual embeddings*

The        train      slowly     left        the        station

$w_1$        $w_2$        $w_3$        $w_4$        $w_5$        $w_6$

$0.05$ + $0.3$ + $0.12$ + $0.08$ + $0.05$ +$0.4$

**Idea**: For each word, we can calculate a <u>weighted average</u> of the stand-alone embeddings of *all* the words in the sentence

| 0.05 | 0.3 | 0.12 | 0.08 | 0.05 | 0.4 |
|------|-----|------|------|------|-----|

# From embeddings to *contextual embeddings*

*The        train    slowly   left      the      station*

$w_1$        $w_2$        $w_3$        $w_4$        $w_5$        $w_6$

0.05 ⬚ + 0.3 ⬚ + 0.12 ⬚ + 0.08 ⬚ + 0.05 ⬚ +0.4 ⬚

Idea: For each word, we can calculate a weighted average of the stand-alone embeddings of *all* the words in the sentence

# From embeddings to *contextual embeddings*

The       train     slowly   left     the     station

$w_1$      $w_2$      $w_3$      $w_4$      $w_5$    $w_6$

$0.05 \quad + 0.3 \quad + 0.12 \quad + 0.08 \quad + 0.05 \quad + 0.4$

Stand-alone embedding

Contextual version of $w_6$

# Same thing but more abstract

The      train     slowly   left      the      station

$w_1$      $w_2$      $w_3$    $w_4$      $w_5$      $w_6$

$s_{1,6}$      $s_{2,6}$      $s_{3,6}$    $s_{4,6}$      $s_{5,6}$      $s_{6,6}$

$\hat{w}_6$

$$\hat{w}_6 = s_{1,6}w_1 + s_{2,6}w_2 + s_{3,6}w_3 + s_{4,6}w_4 + s_{5,6}w_5 + s_{6,6}w_6$$

For a given word (e.g., 'station'), how should the weights of the other words be chosen?

# For a given word (e.g., 'station'), how should the weights of the other words be chosen?

## Intuition

- The weight of a word should be proportional to how related it is to the word "station"
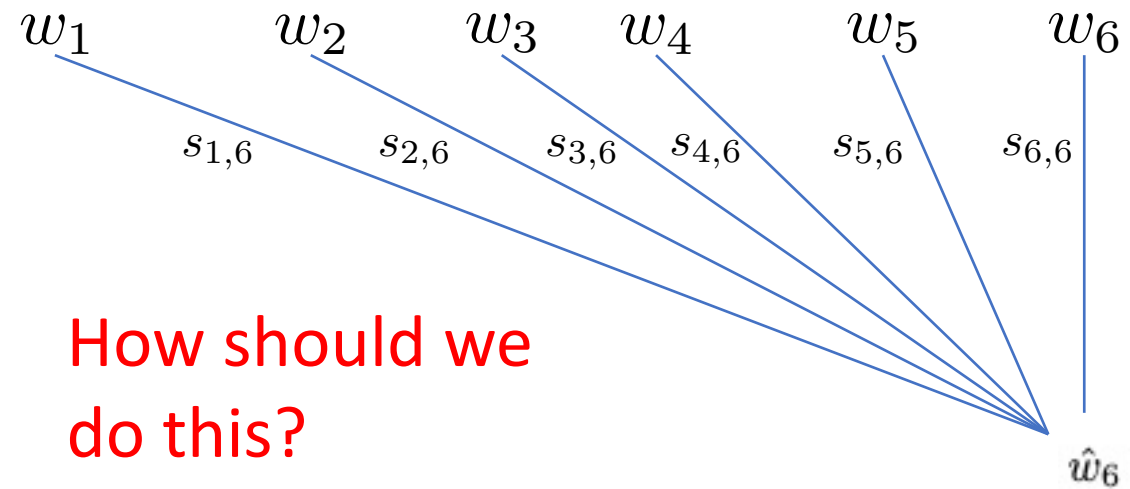
# For a given word (e.g., 'station'), how should the weights of the other words be chosen?

## Intuition

- The weight of a word should be proportional to how related it is to the word "station"

- One way to quantify how "related" two words are: the *dot-product* of their stand-alone embeddings
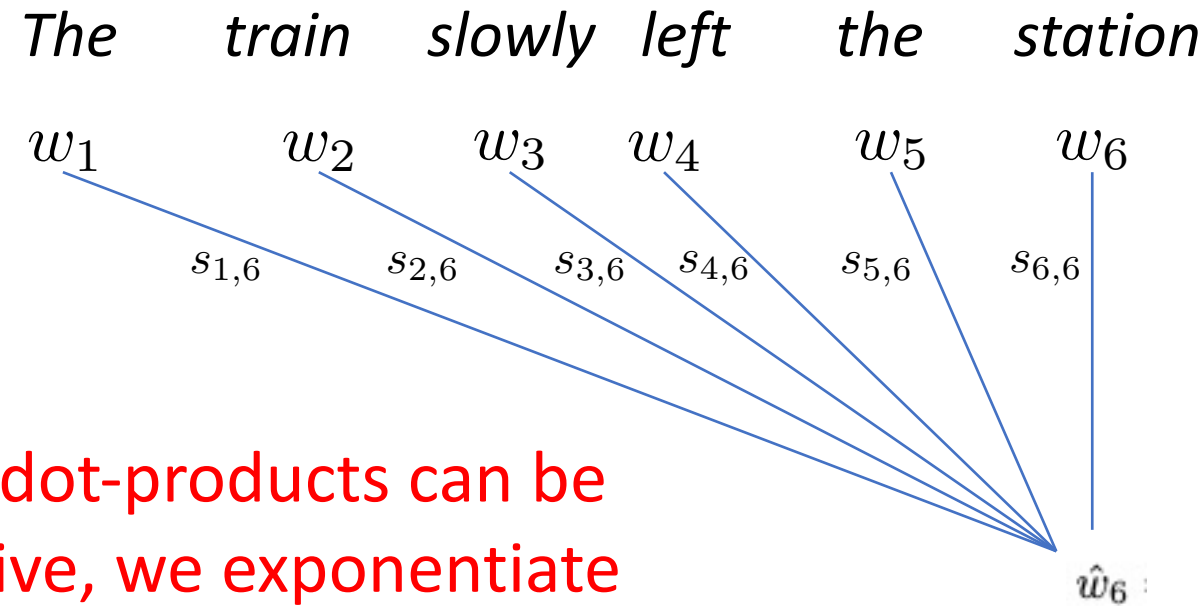
# We need to convert dot-products to proper weights*

The      train    slowly  left    the    station

$$w_1 \qquad w_2 \qquad w_3 \qquad w_4 \qquad w_5 \qquad w_6$$

$$s_{1,6} \qquad s_{2,6} \qquad s_{3,6} \quad s_{4,6} \qquad s_{5,6} \qquad s_{6,6}$$

**How should we do this?**

$\hat{w}_6$

* non-negative, and summing to 1.0

# We need to convert dot-products to proper weights*

*The*     *train*    *slowly*   *left*     *the*     *station*

$w_1$      $w_2$      $w_3$    $w_4$      $w_5$     $w_6$

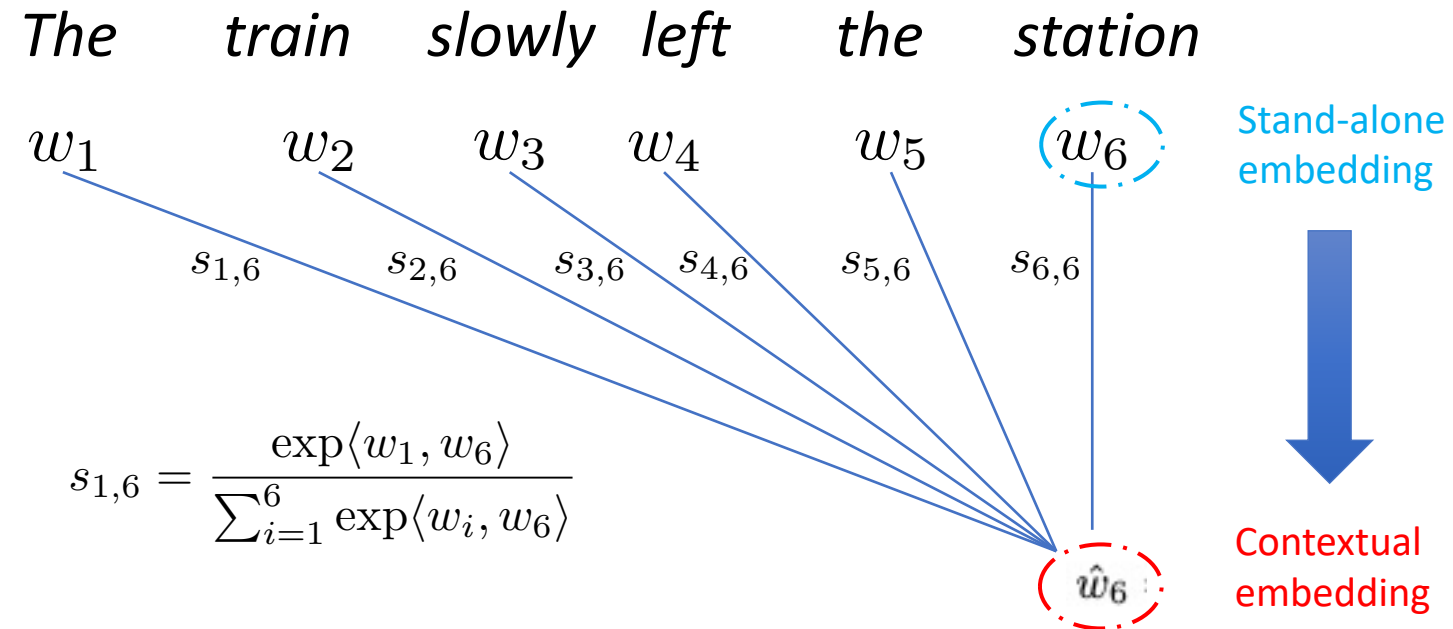$s_{1,6}$     $s_{2,6}$     $s_{3,6}$   $s_{4,6}$     $s_{5,6}$     $s_{6,6}$

Since dot-products can be negative, we exponentiate them and then normalize (remember softmax?)

$\hat{w}_6$

\* non-negative, and summing to 1.0

# Normalized attention weights



The     train    slowly  left    the    station

$w_1$     $w_2$     $w_3$    $w_4$     $w_5$     $w_6$

$s_{1,6}$     $s_{2,6}$     $s_{3,6}$    $s_{4,6}$     $s_{5,6}$     $s_{6,6}$

Since dot-products can be negative, we exponentiate them and then normalize

$$s_{1,6} = \frac{\exp\langle w_1, w_6 \rangle}{\sum_{i=1}^{6} \exp\langle w_i, w_6 \rangle}$$

$\hat{w}_6$

# From embedding to *contextual* embedding!

*The*     *train*     *slowly*   *left*     *the*     *station*

$w_1$      $w_2$     $w_3$    $w_4$     $w_5$    $w_6$

Stand-alone embedding

$s_{1,6}$    $s_{2,6}$    $s_{3,6}$   $s_{4,6}$    $s_{5,6}$    $s_{6,6}$

$$s_{1,6} = \frac{\exp\langle w_1, w_6 \rangle}{\sum_{i=1}^{6} \exp\langle w_i, w_6 \rangle}$$

$\hat{w}_6$

Contextual embedding

$$\hat{w}_6 = s_{1,6}w_1 + s_{2,6}w_2 + s_{3,6}w_3 + s_{4,6}w_4 + s_{5,6}w_5 + s_{6,6}w_6$$

By choosing weights in this manner, the embedding of a word moves closer to the embeddings of the other words in the current context, in proportion to how related they are

train                    station                    radio

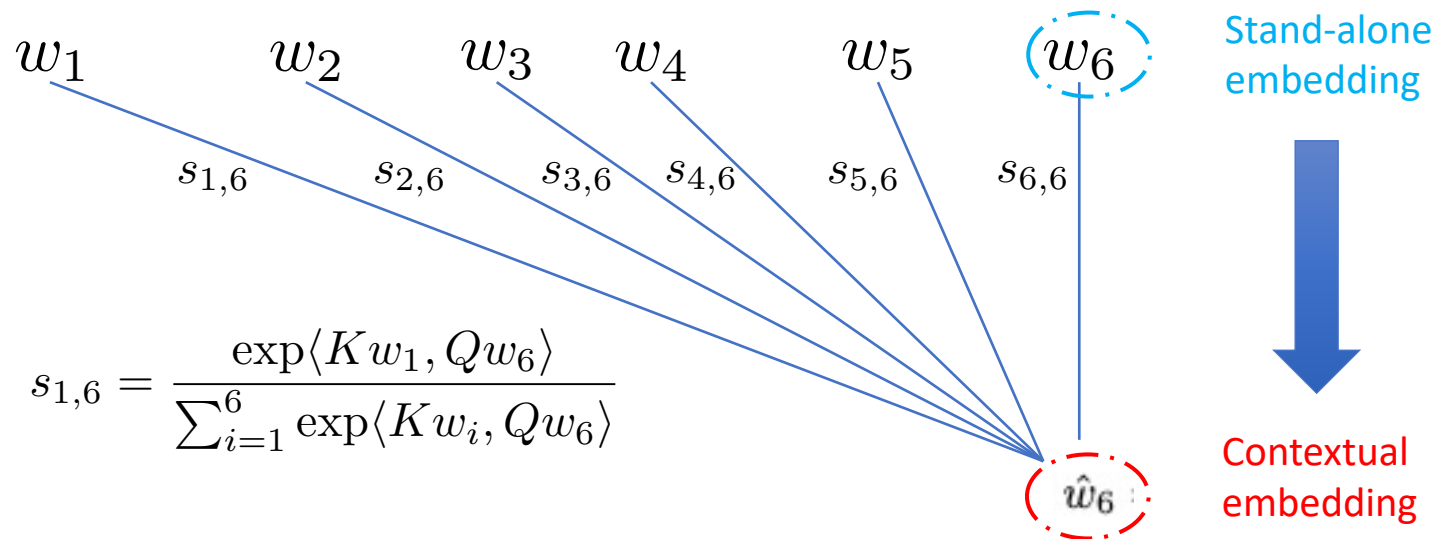- The word 'station' has many contexts.

By choosing weights in this manner, the embedding of a word moves closer to the embeddings of the other words in the current context, in proportion to how related they are

train             station             radio

- The word 'station' has many contexts.
  - In the current context, 'train' is closely related to 'station' and therefore exerts a strong "pull" on it

By choosing weights in this manner, the embedding of a word moves closer to the embeddings of the other words in the current context, in proportion to how related they are



- The word 'station' has many contexts.
  - In the current context, 'train' is closely related to 'station' and therefore exerts a strong "pull" on it
  - 'radio' is also related to 'station' but doesn't appear in the current context so (automatically) has zero weight

By choosing weights in this manner, the embedding of a word moves closer to the embeddings of the other words in the current context, in proportion to how related they are

train ⟵--------------------------- station radio

- The word 'station' has many contexts.
  - In the current context, 'train' is closely related to 'station' and therefore exerts a strong "pull" on it
  - 'radio' is also related to 'station' but doesn't appear in the current context so (automatically) has zero weight

- By moving station closer to train (equivalently – paying more "attention" to train), we are contextualizing station's embedding to the context of trains, platforms, departures, etc.

# From embedding to *contextual* embedding!
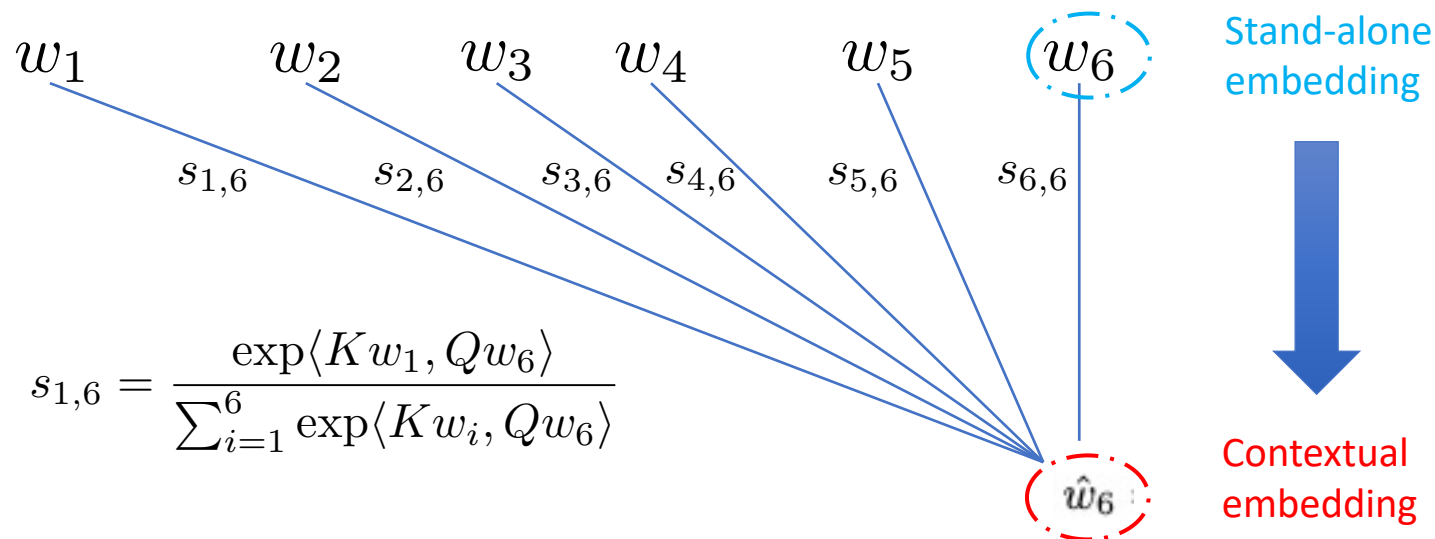
The  train  slowly  left  the  station

$w_1$  $w_2$  $w_3$  $w_4$  $w_5$  $w_6$

<span style="color:cyan">Stand-alone embedding</span>

Make it more flexible

$s_{1,6}$  $s_{2,6}$  $s_{3,6}$  $s_{4,6}$  $s_{5,6}$  $s_{6,6}$

$$s_{1,6} = \frac{\exp\langle Kw_1, Qw_6 \rangle}{\sum_{i=1}^{6} \exp\langle Kw_i, Qw_6 \rangle}$$

$\hat{w}_6$

<span style="color:red">Contextual embedding</span>

$$\hat{w}_6 = s_{1,6}w_1 + s_{2,6}w_2 + s_{3,6}w_3 + s_{4,6}w_4 + s_{5,6}w_5 + s_{6,6}w_6$$

# From embedding to *contextual* embedding!

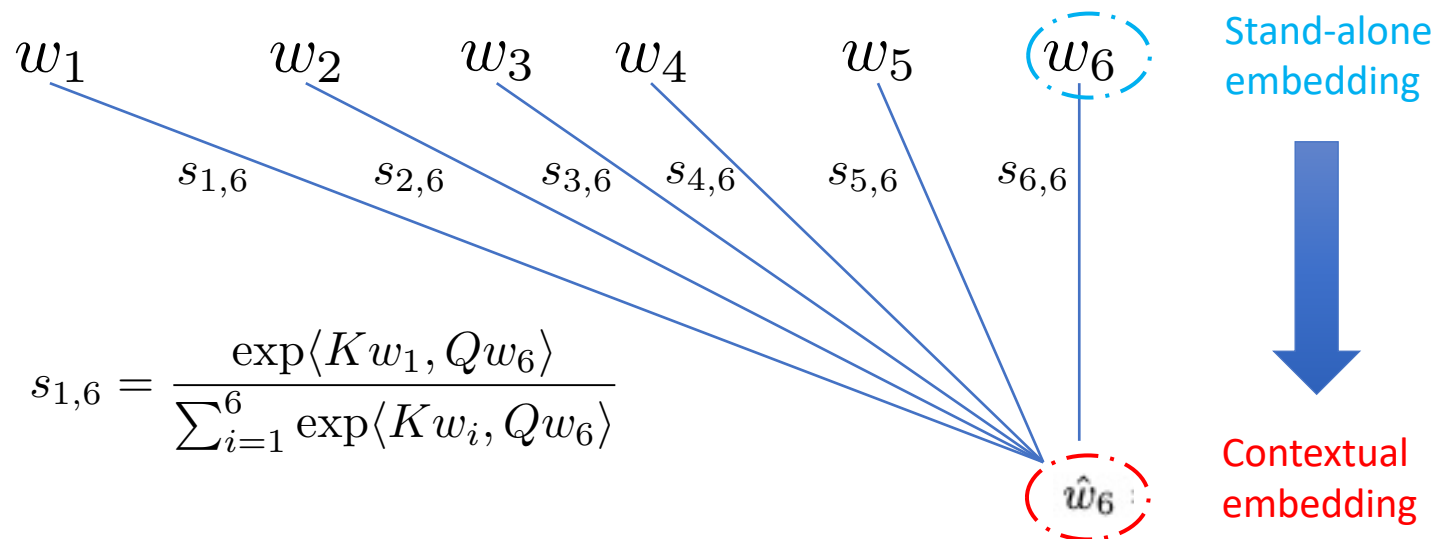*The    train    slowly  left    the    station*

$w_1$        $w_2$        $w_3$    $w_4$        $w_5$    $(w_6)$    Stand-alone embedding

$s_{1,6}$    $s_{2,6}$    $s_{3,6}$  $s_{4,6}$    $s_{5,6}$    $s_{6,6}$

Make it more flexible:
- e.g., in K=transportation subspace, Q=location subspace, how closely two words are related?

$$s_{1,6} = \frac{\exp\langle Kw_1, Qw_6\rangle}{\sum_{i=1}^{6} \exp\langle Kw_i, Qw_6\rangle}$$

$(\hat{w}_6)$    Contextual embedding

$$\hat{w}_6 = s_{1,6}w_1 + s_{2,6}w_2 + s_{3,6}w_3 + s_{4,6}w_4 + s_{5,6}w_5 + s_{6,6}w_6$$

# From embedding to *contextual* embedding!

*The        train     slowly   left      the      station*

$w_1$          $w_2$          $w_3$      $w_4$          $w_5$          $w_6$   <span style="color:#3daae0">Stand-alone embedding</span>

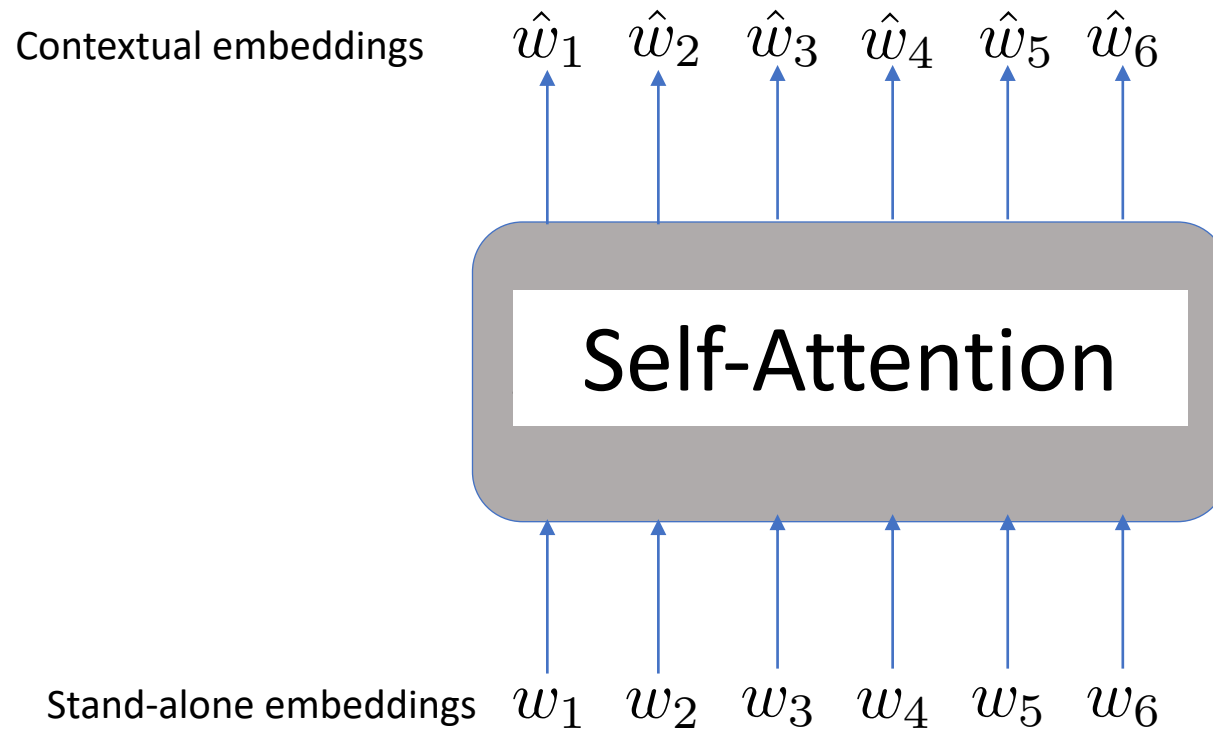$s_{1,6}$        $s_{2,6}$        $s_{3,6}$   $s_{4,6}$      $s_{5,6}$        $s_{6,6}$

Make it more flexible:
- e.g., in K=transportation subspace, Q=location subspace, how closely two words are related?

$$s_{1,6} = \frac{\exp\langle Kw_1, Qw_6\rangle}{\sum_{i=1}^{6}\exp\langle Kw_i, Qw_6\rangle}$$

$\hat{w}_6$

<span style="color:#c00000">Contextual embedding</span>

$$\hat{w}_6 = s_{1,6}Vw_1 + s_{2,6}Vw_2 + s_{3,6}Vw_3 + s_{4,6}Vw_4 + s_{5,6}Vw_5 + s_{6,6}Vw_6$$

# From embedding to *contextual* embedding!

*The*       *train*       *slowly*   *left*       *the*       *station*

$w_1$       $w_2$       $w_3$       $w_4$       $w_5$       $w_6$

Stand-alone embedding

$s_{1,6}$       $s_{2,6}$       $s_{3,6}$   $s_{4,6}$       $s_{5,6}$       $s_{6,6}$

Make it more flexible:
- e.g., in K=transportation subspace, Q=location subspace, how closely two words are related?

$$s_{1,6} = \frac{\exp\langle Kw_1, Qw_6\rangle}{\sum_{i=1}^{6}\exp\langle Kw_i, Qw_6\rangle}$$

$\hat{w}_6$

Contextual embedding

$$\hat{w}_6 = s_{1,6}Vw_1 + s_{2,6}Vw_2 + s_{3,6}Vw_3 + s_{4,6}Vw_4 + s_{5,6}Vw_5 + s_{6,6}Vw_6$$

Key (K), Query (Q), Value (V) projection are all tunable parameters: $K, Q, V \in R^{d \times d_{model}}$

This operation is referred to as a 'Self Attention' layer and can be done very efficiently with matrix operations
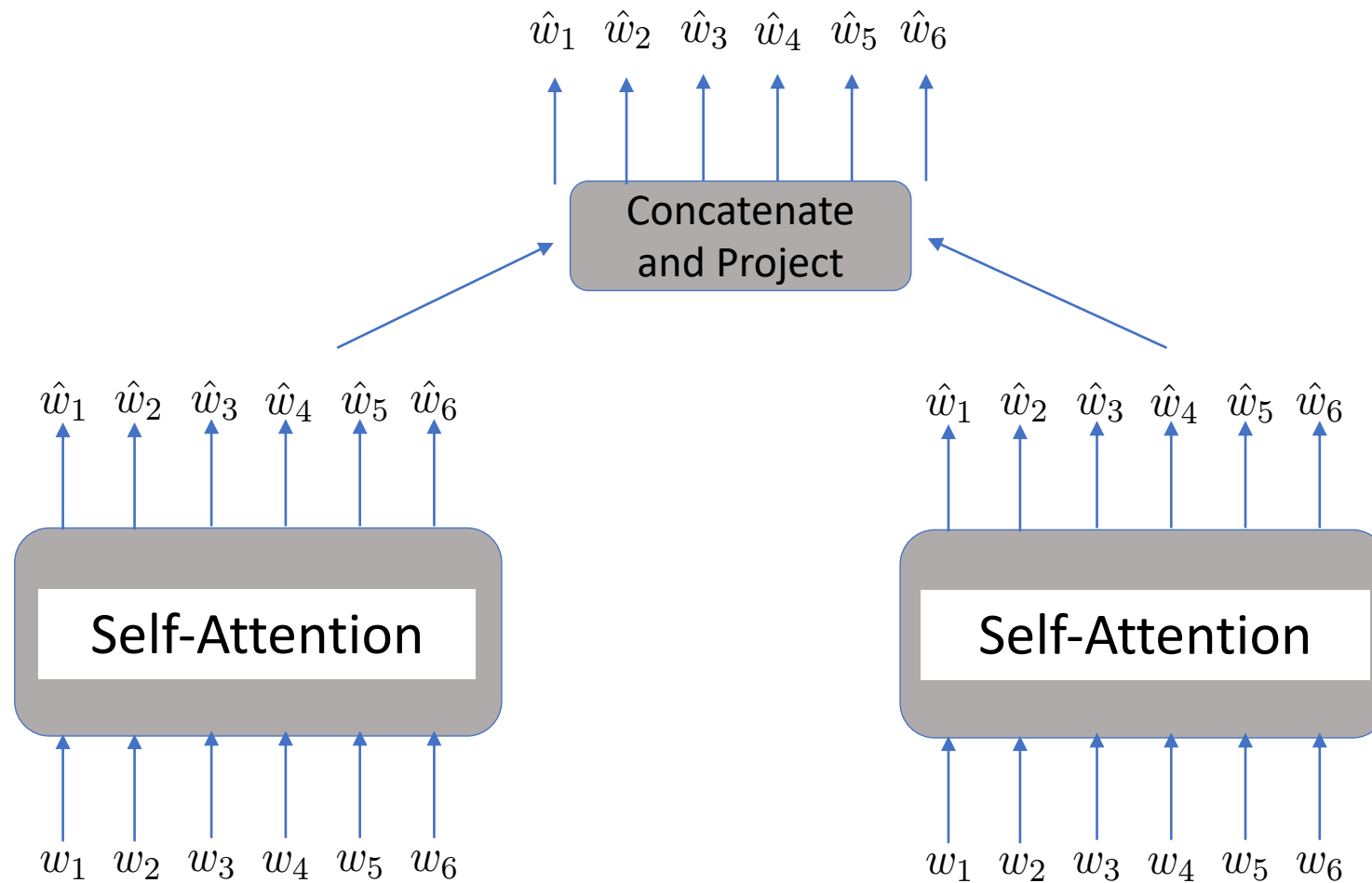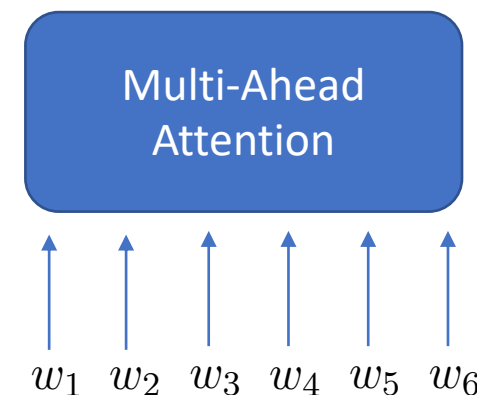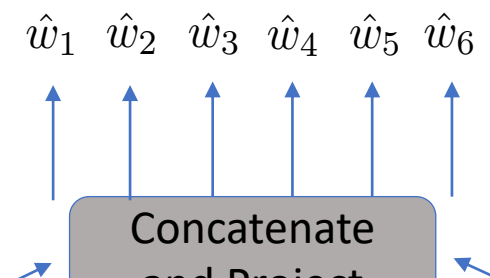
Contextual embeddings $\hat{w}_1$ $\hat{w}_2$ $\hat{w}_3$ $\hat{w}_4$ $\hat{w}_5$ $\hat{w}_6$

## Self-Attention

Stand-alone embeddings $w_1$ $w_2$ $w_3$ $w_4$ $w_5$ $w_6$

# Key Tweak: Multi*-Head Attention

$\hat{w}_1 \quad \hat{w}_2 \quad \hat{w}_3 \quad \hat{w}_4 \quad \hat{w}_5 \quad \hat{w}_6$

Concatenate and Project

$\hat{w}_1 \quad \hat{w}_2 \quad \hat{w}_3 \quad \hat{w}_4 \quad \hat{w}_5 \quad \hat{w}_6$

Self-Attention

$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6$

$\hat{w}_1 \quad \hat{w}_2 \quad \hat{w}_3 \quad \hat{w}_4 \quad \hat{w}_5 \quad \hat{w}_6$

Self-Attention

$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6$

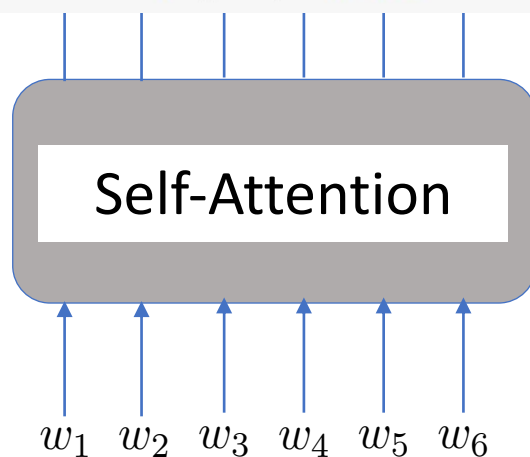# Key Tweak: Different attention 'heads' learn different patterns

$$\hat{w}_1 \quad \hat{w}_2 \quad \hat{w}_3 \quad \hat{w}_4 \quad \hat{w}_5 \quad \hat{w}_6$$

Concatenate and Project

$$\hat{w}_1 \quad \hat{w}_2 \quad \hat{w}_3 \quad \hat{w}_4 \quad \hat{w}_5 \quad \hat{w}_6$$

Self-Attention

$$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6$$

$$\hat{w}_1 \quad \hat{w}_2 \quad \hat{w}_3 \quad \hat{w}_4 \quad \hat{w}_5 \quad \hat{w}_6$$

Self-Attention

$$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6$$

# Key Tweak: Different attention 'heads' learn different patterns

Combine together..



$$\hat{w}_1 \quad \hat{w}_2 \quad \hat{w}_3 \quad \hat{w}_4 \quad \hat{w}_5 \quad \hat{w}_6$$

Multi-Ahead Attention

$$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6$$
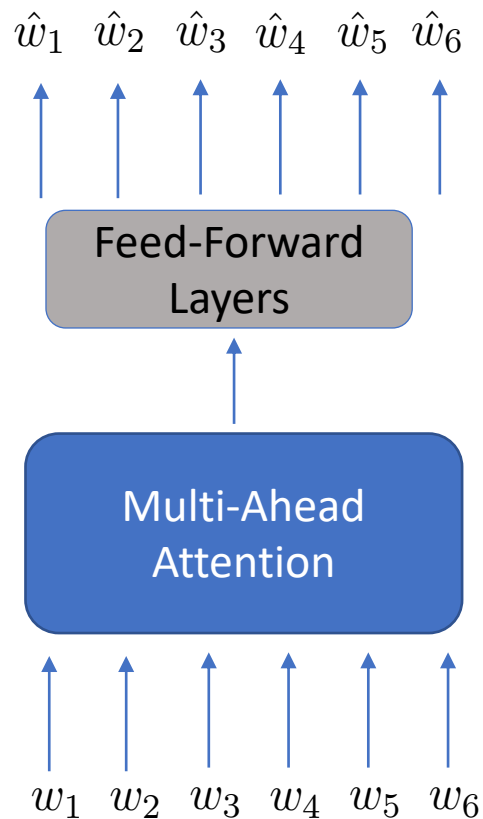
Concatenate and Project

```python
def forward(self, x):
    ### dimension of x is [batch_size, sequence_length, embedding_dim]
    self.attention = nn.MultiheadAttention(embedding_dim, num_heads, dropout=dropout, batch_first=True)
    return self.attention(x, x, x)[0]
```

Self-Attention

Self-Attention

$$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6$$

$$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6$$

# Key Tweak: Inject some non-linearity with feed-forward layers at the end

$\hat{w}_1 \quad \hat{w}_2 \quad \hat{w}_3 \quad \hat{w}_4 \quad \hat{w}_5 \quad \hat{w}_6$

Feed-Forward Layers

Multi-Ahead Attention

$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6$

# Key Tweak: Inject some non-linearity with feed-forward layers at the end

$\hat{w}_1$  $\hat{w}_2$  $\hat{w}_3$  $\hat{w}_4$  $\hat{w}_5$  $\hat{w}_6$

Feed-Forward Layers

Multi-Ahead Attention

$w_1$  $w_2$  $w_3$  $w_4$  $w_5$  $w_6$

```python
self.feed_forward = nn.Sequential(
    nn.Linear(embedding_dim, 4 * embedding_dim),
    nn.GELU(),
    nn.Linear(4 * embedding_dim, embedding_dim),
)
```

$$GELU(x) = xP(X \leq x), X \sim N(0,1)$$

# The story so far

*End with contextual embeddings*

$$\hat{w}_1 \quad \hat{w}_2 \quad \hat{w}_3 \quad \hat{w}_4 \quad \hat{w}_5 \quad \hat{w}_6$$

Feed-Forward Layers

Multi-Ahead Attention

$$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6$$

*Start with random embeddings*

# The story so far

*End with contextual embeddings*

$\hat{w}_1 \quad \hat{w}_2 \quad \hat{w}_3 \quad \hat{w}_4 \quad \hat{w}_5 \quad \hat{w}_6$

Feed-Forward Layers

Multi-Ahead Attention

$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6$

*Start with random embeddings*

Transformer Layer

Add & Norm

Feed Forward

N×

Add & Norm

Multi-Head Attention

Attention is all you need

# The story so far

*End with contextual embeddings*

*Start with random embeddings*

$$\hat{w}_1 \quad \hat{w}_2 \quad \hat{w}_3 \quad \hat{w}_4 \quad \hat{w}_5 \quad \hat{w}_6$$

Feed-Forward Layers

Multi-Ahead Attention

$$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6$$

Transformer Layer

dropout

dropout

Add & Norm

Feed Forward

Nx

Add & Norm

Multi-Head Attention

Attention is all you need

# Three regularization techniques for Deep Learning

- Dropout

- Residual connections (Add)

- Layer normalization (Norm)

# Dropout

In Training, randomly zero out the output from some of the nodes (typically 10% of the nodes) in a hidden layer (implemented as a "dropout layer" in PyTorch)
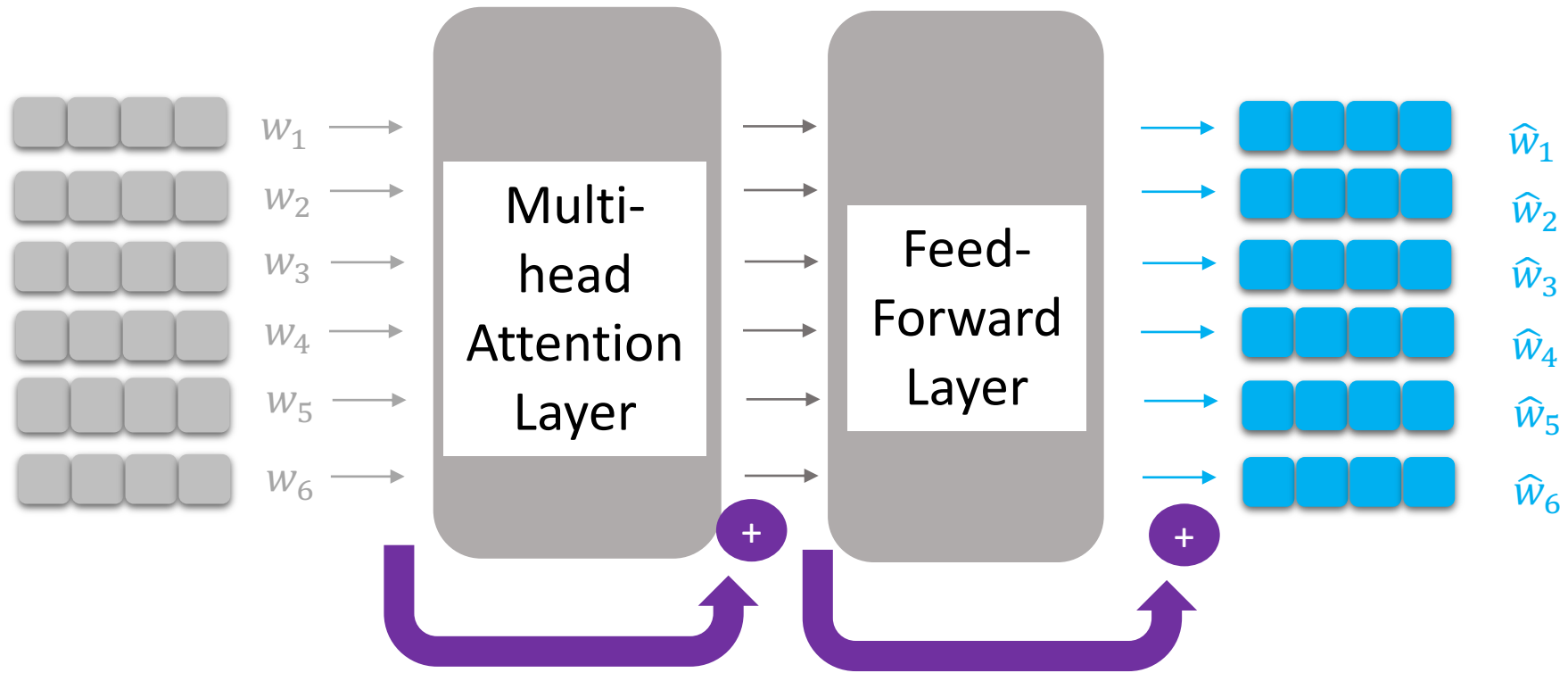


```
self.dropout = nn.Dropout(0.1)
x = self.dropout(x)
```
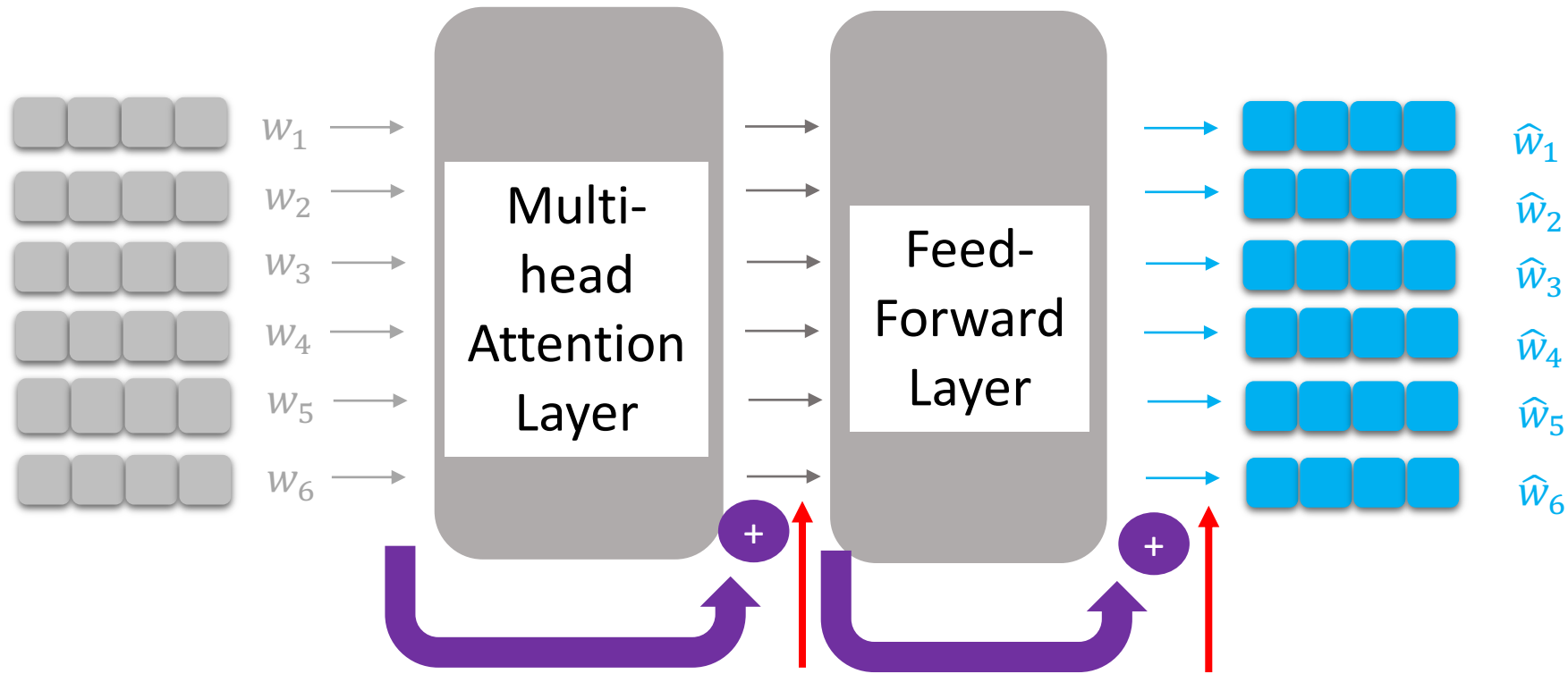
# Residual Connection

$$y = \text{Layer}(x) + x$$

We sum the input embedding to the output embedding of the Attention / Feed-Forward Layers. This helps gradients flow better during backpropagation.

# Layer Normalization

$$y = \frac{x - \mathrm{E}[x]}{\sqrt{\mathrm{Var}[x] + \epsilon}} * \gamma + \beta$$
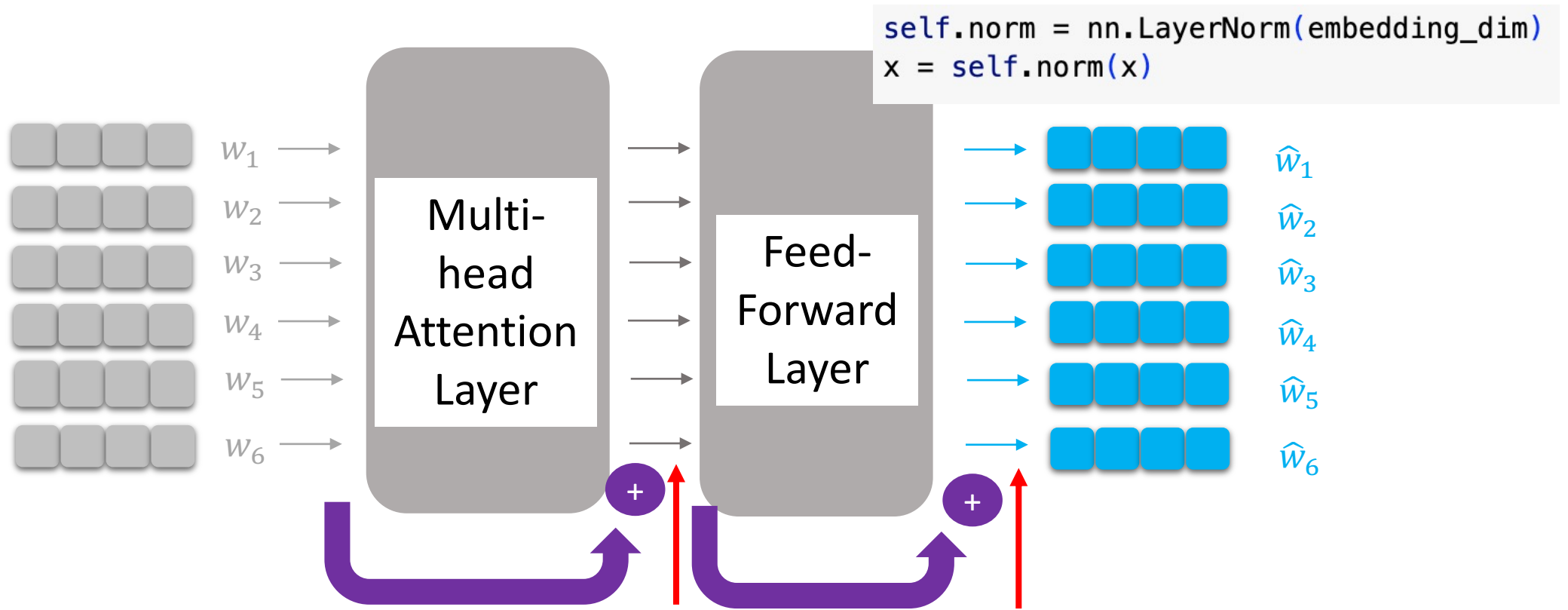
After the Attention / Feed-Forward Layers, we standardize (i.e., subtract mean and divide by std) each embedding. This ensures that the weights stay small.

# Layer Normalization

$$y = \frac{x - \mathrm{E}[x]}{\sqrt{\mathrm{Var}[x] + \epsilon}} * \gamma + \beta$$

After the Attention / Feed-Forward Layers, we standardize (i.e., subtract mean and divide by std) each embedding. This ensures that the weights stay small.

```python
self.norm = nn.LayerNorm(embedding_dim)
x = self.norm(x)
```
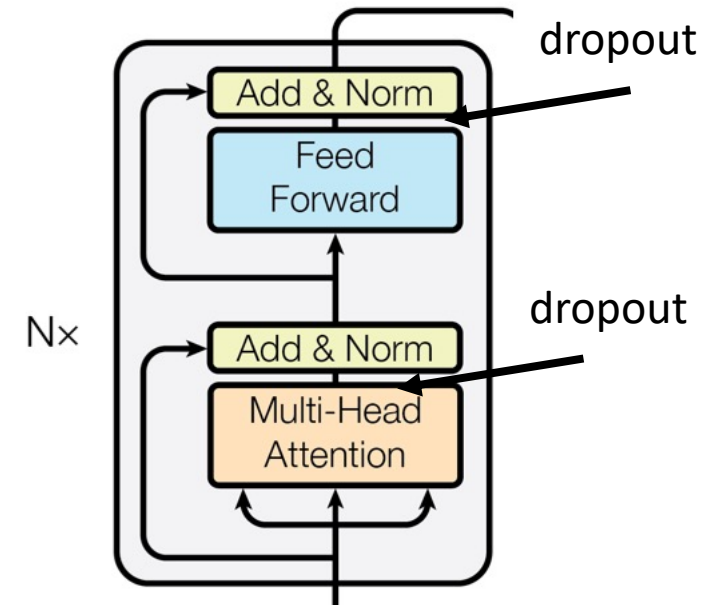


45

# Taken Together: A Transformer Layer

```
y = MultiheadAttention(x) ## Attention Layer
y = Dropout(y) ## Dropout layer
x = x + y ## Add Residule
x = LayerNorm(x) ## Layer Norm

y = FeedForward(x) ## Feed Forward Layer
y = Dropout(y) ## Dropout layer
x = x + y ## Add Residule
x = LayerNorm(x) ## Layer Norm
```

Pseudo code

dropout

dropout

Nx

Add & Norm

Feed Forward

Add & Norm

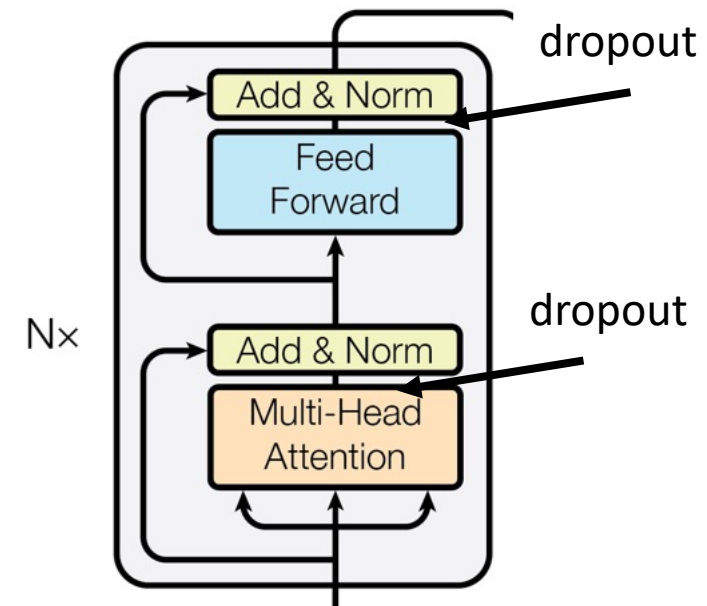Multi-Head Attention

Attention is all you need

# Taken Together: A Transformer Layer

What is dimension of x?

```
y = MultiheadAttention(x) ## Attention Layer
y = Dropout(y) ## Dropout layer
x = x + y ## Add Residule
x = LayerNorm(x) ## Layer Norm

y = FeedForward(x) ## Feed Forward Layer
y = Dropout(y) ## Dropout layer
x = x + y ## Add Residule
x = LayerNorm(x) ## Layer Norm
```

Pseudo code



dropout

dropout

Nx

Attention is all you need

# Taken Together: A Transformer Layer

What is dimension of x? [batch_size, seq_Len, embedding_dim]

```
y = MultiheadAttention(x) ## Attention Layer
y = Dropout(y) ## Dropout layer
x = x + y ## Add Residule
x = LayerNorm(x) ## Layer Norm

y = FeedForward(x) ## Feed Forward Layer
y = Dropout(y) ## Dropout layer
x = x + y ## Add Residule
x = LayerNorm(x) ## Layer Norm
```

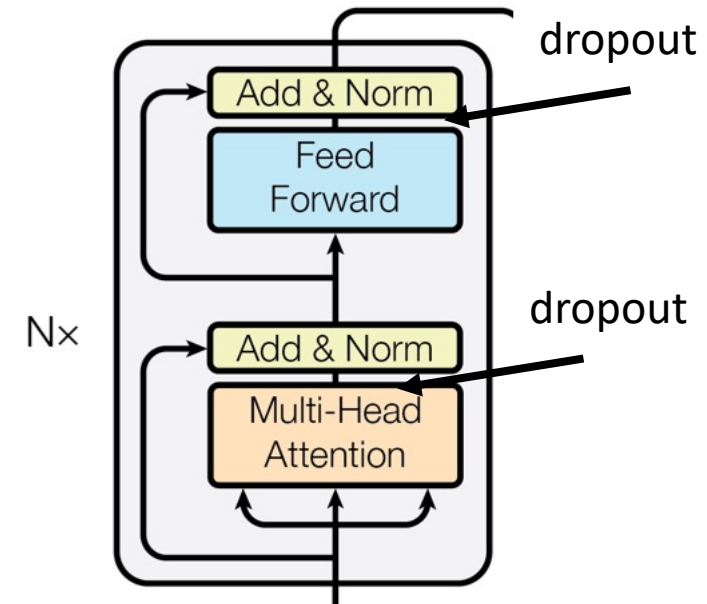Pseudo code



dropout

dropout

Nx

Attention is all you need

# Taken Together: A Transformer Layer

What is dimension of x? [batch_size, seq_Len, embedding_dim]

```
y = MultiheadAttention(x) ## Attention Layer
y = Dropout(y) ## Dropout layer
x = x + y ## Add Residule
x = LayerNorm(x) ## Layer Norm

y = FeedForward(x) ## Feed Forward Layer
y = Dropout(y) ## Dropout layer
x = x + y ## Add Residule
x = LayerNorm(x) ## Layer Norm
```
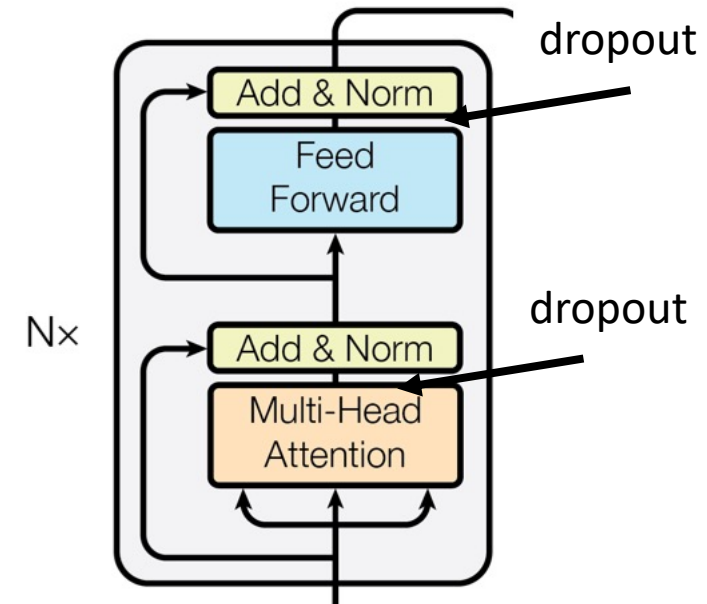
Pseudo code



dropout

dropout

Nx

Attention is all you need

We can have multiple Transformer Layers stacked

# Does everything look good?



Tokenization

emb1     emb2     emb3     emb4

5247     467     1167     43570

Go     go     power     ranger

# Does everything look good?

# Does everything look good?



Transformer Layer

Transformer Layer

Let's directly take positions as inputs and learn a position embedding!

Tokenization

| emb1 | emb2 | emb3 | emb4 |

5247    467    1167    43570

Go    go    power    ranger

# Does everything look good?

Transformer Layer

Transformer Layer

Let's directly take positions as inputs and learn a position embedding!

Tokenization

emb1   emb2   emb3   emb4

+   +   +   +

P1_emb   P2_emb   P3_emb   P4_emb

5247   467   1167   43570

Go   go   power   ranger

# This is called a Transformer (Encoder)!

What is dimension of x? [batch_size, seq_Len, embedding_dim]



Nx

Positional Encoding

Input Layer

Inputs

Transformer Layer

```
y = MultiheadAttention(x) ## Attention Layer
y = Dropout(y) ## Dropout layer
x = x + y ## Add Residule
x = LayerNorm(x) ## Layer Norm

y = FeedForward(x) ## Feed Forward Layer
y = Dropout(y) ## Dropout layer
x = x + y ## Add Residule
x = LayerNorm(x) ## Layer Norm
```
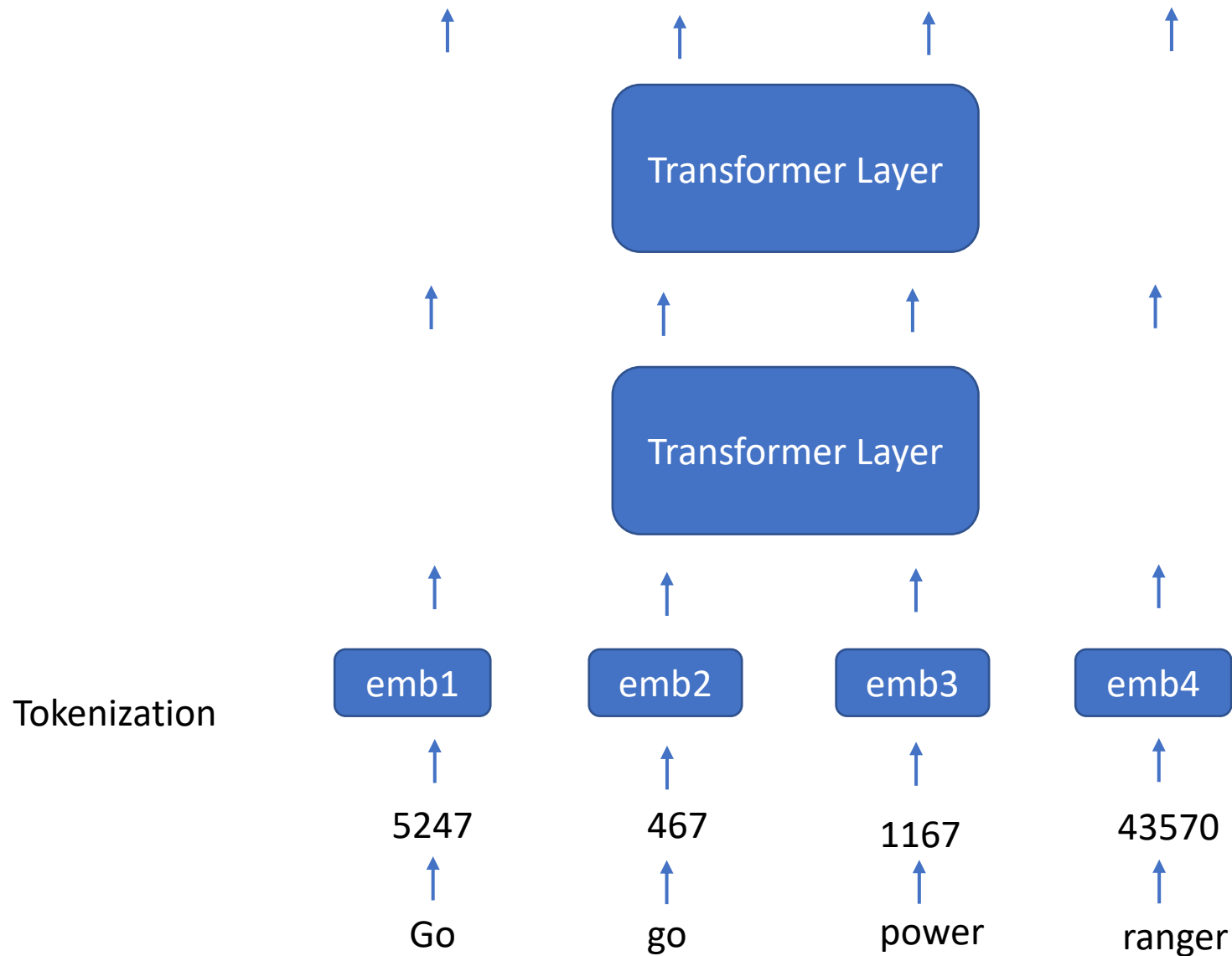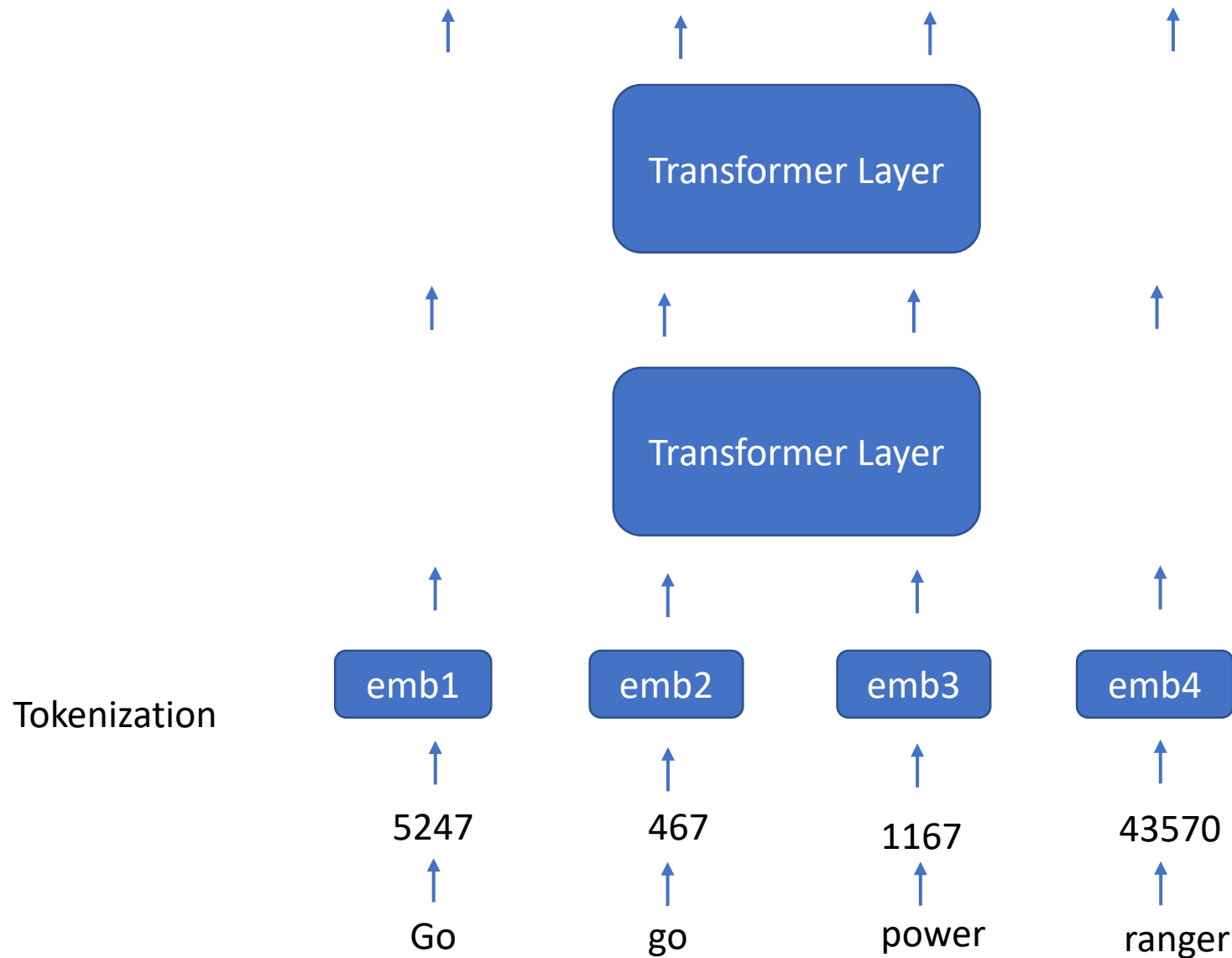
```
x = Embedding(inputs)
x = x + Embedding(torch.arange(seq_len))
x = LayerNorm(x)
x = Dropout(x)
```

https://arxiv.org/abs/1706.03762

# Coding Exercise

- Text Classification: https://colab.research.google.com/drive/1hWA0Tf4DpeSbhzpb5DP6M_KxfvZOxEyw?usp=sharing

```python
x = Transformer(inputs)
x = torch.mean(x, dim=1)
x = nn.Linear(embedding_dim, 2)(x)
```

Classification Head

# A better training wrapper

- DataLoader
  - The DataLoader in PyTorch efficiently batches and shuffles data, while also enabling parallel processing and flexible collation, making data loading faster and more scalable for training large models.

- Pytorch_lightning
  - streamlines the training process by automating loops, logging, and device management, enabling cleaner, more efficient model development.

# Text Generation

老子

A journey of a thousand miles begins with a single step. ---- Laozi (老子)

# How would we train a model to generate the next token/word?

# How would we train a model to generate the next token/word?

- Multi-class classification problem!

| Next Word* | Probability |
|---|---|
| aardvark | 0.0003 |
| ... | |
| rainy | 0.3 |
| ... | |
| stormy | 0.6 |
| ... | |
| zebra | 0.00009 |

Cross Entropy Loss

# How would we train a model to generate the next token/word?

logits

Linear  (Emb_dim, Vocab_size)

$\hat{w}_1 \quad \hat{w}_2 \quad \hat{w}_3 \quad \hat{w}_4 \quad \hat{w}_5 \quad \hat{w}_6$

Transformer

$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6$

# How would we train a model to generate the next token/word?

logits

Linear (Emb_dim, Vocab_size)

$\hat{w}_1$ $\hat{w}_2$ $\hat{w}_3$ $\hat{w}_4$ $\hat{w}_5$ $\hat{w}_6$

Transformer

$w_1$ $w_2$ $w_3$ $w_4$ $w_5$ $w_6$

Ideally we want to use

$\hat{w}_1 \rightarrow w_2$

$\hat{w}_2 \rightarrow w_3$

$\hat{w}_3 \rightarrow w_4$

$\hat{w}_4 \rightarrow w_5$

# How would we train a model to generate the next token/word?

logits

Linear   (Emb_dim, Vocab_size)

$\hat{w}_1$  $\hat{w}_2$  $\hat{w}_3$  $\hat{w}_4$  $\hat{w}_5$  $\hat{w}_6$

Transformer

$w_1$  $w_2$  $w_3$  $w_4$  $w_5$  $w_6$

$$\hat{w}_1 \rightarrow w_2$$

$$\hat{w}_2 \rightarrow w_3$$

Ideally we want to use

$$\hat{w}_3 \rightarrow w_4$$

$$\hat{w}_4 \rightarrow w_5$$

But $\hat{w}_1$ has information of other words

# How would we train a model to generate the next token/word?

logits

$\uparrow$

Linear    (Emb_dim, Vocab_size)

$\uparrow$

$\hat{w}_1 \quad \hat{w}_2 \quad \hat{w}_3 \quad \hat{w}_4 \quad \hat{w}_5 \quad \hat{w}_6$

$\uparrow \quad \uparrow \quad \uparrow \quad \uparrow \quad \uparrow \quad \uparrow$

Transformer

$\uparrow \quad \uparrow \quad \uparrow \quad \uparrow \quad \uparrow \quad \uparrow$

$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6$

$$\hat{w}_1 \rightarrow w_2$$
$$\hat{w}_2 \rightarrow w_3$$

Ideally we want to use

$$\hat{w}_3 \rightarrow w_4$$
$$\hat{w}_4 \rightarrow w_5$$

But $\hat{w}_1$ has information of other words
This is where Causal Mask comes in

# In transformer, the information spread across tokens by Self-Attention

*The        train      slowly  left      the      station*

$w_1$        $w_2$        $w_3$      $w_4$        $w_5$        $w_6$        Stand-alone embedding

$s_{1,6}$        $s_{2,6}$        $s_{3,6}$      $s_{4,6}$        $s_{5,6}$        $s_{6,6}$

$$s_{1,6} = \frac{\exp\langle Kw_1, Qw_6 \rangle}{\sum_{i=1}^{6} \exp\langle Kw_i, Qw_6 \rangle}$$

$\hat{w}_6$

Contextual embedding

$$\hat{w}_6 = s_{1,6}Vw_1 + s_{2,6}Vw_2 + s_{3,6}Vw_3 + s_{4,6}Vw_4 + s_{5,6}Vw_5 + s_{6,6}Vw_6$$

# Causal Mask requires a token cannot use the information of a future token

The     train    slowly   left     the     station

$w_1$        $w_2$      $w_3$   $w_4$      $w_5$      $w_6$

Stand-alone embedding

$s_{1,4}$     $s_{2,4}$    $s_{3,4}$   $s_{4,4}$

$$s_{1,4} = \frac{\exp\langle Kw_1, Qw_4\rangle}{\sum_{i=1}^{4}\exp\langle Kw_i, Qw_4\rangle}$$

$\hat{w}_4$

Contextual embedding

$$\hat{w}_4 = s_{1,4}Vw_1 + s_{2,4}Vw_4 + s_{3,4}Vw_3 + s_{4,4}Vw_4$$

# Causal Mask requires a token cannot use the information of a future token

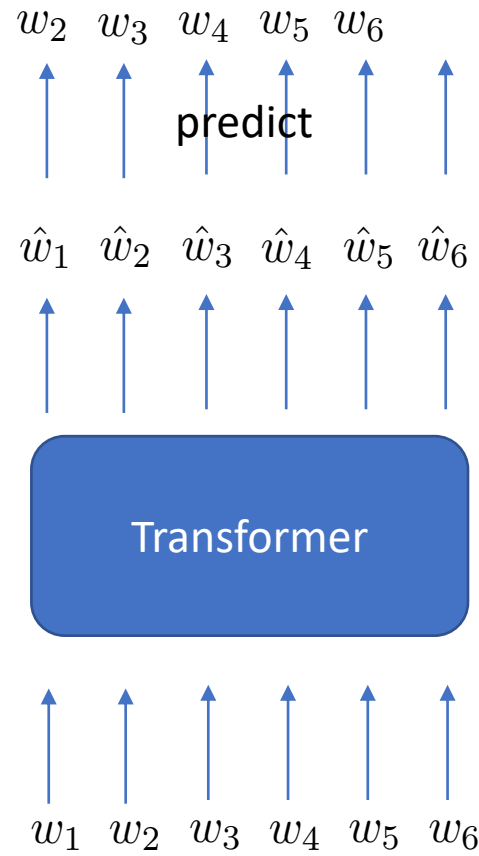$$A_{ij} = \text{mask}_{ij} \cdot \exp \langle Kw_i, Qw_j \rangle$$

$$s_{i,k} = \frac{A_{ij}}{\sum_{i=1}^{n} A_{ij}}$$

|  | the | cat | sat | on | the |
|---|---|---|---|---|---|
| the | 1 |  |  |  |  |
| cat | 1 | 1 |  |  |  |
| sat | 1 | 1 | 1 |  |  |
| on | 1 | 1 | 1 | 1 |  |
| the | 1 | 1 | 1 | 1 | 1 |

mask

Causal Self-Attention, Or Masked Self Attention

# With Causal Mask, *training* is much more efficient

$$w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6$$

predict

$$\hat{w}_1 \quad \hat{w}_2 \quad \hat{w}_3 \quad \hat{w}_4 \quad \hat{w}_5 \quad \hat{w}_6$$

All predictions/loss calculations can be done by one pass

**Transformer**

$$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6$$

# Autoregressive Language Model

With this training, essentially we obtain an oracle: $p(x_k|x_{<k})$, an autoregressive model for language

## This is what we called GPT!

GPT: Generative Pretrained Transformers
1. Transformer-based text-generation model
2. Pre-trained on **Massive** amount data

# Autoregressive Language Model

With this training, essentially we obtain an oracle: $p(x_k|x_{<k})$ , an autoregressive model for language

GPT2 Model Architecture

| Parameters | Layers | $d_{model}$ |
|---|---|---|
| 117M | 12 | 768 |
| 345M | 24 | 1024 |
| 762M | 36 | 1280 |
| 1542M | 48 | 1600 |

# Autoregressive Language Model

With this training, essentially we obtain an oracle: $p(x_k|x_{<k})$, an autoregressive model for language

### GPT2 Model Architecture

| Parameters | Layers | $d_{model}$ |
| --- | --- | --- |
| 117M | 12 | 768 |
| 345M | 24 | 1024 |
| 762M | 36 | 1280 |
| 1542M | 48 | 1600 |

### GPT2 Dataset

How would you collect high quality data from Internet?

# Autoregressive Language Model

With this training, essentially we obtain an oracle: $p(x_k|x_{<k})$, an autoregressive model for language

GPT2 Model Architecture

| Parameters | Layers | $d_{model}$ |
|---|---|---|
| 117M | 12 | 768 |
| 345M | 24 | 1024 |
| 762M | 36 | 1280 |
| 1542M | 48 | 1600 |

GPT2 Dataset

Links from Reddit Posts that have at least 3 karmas

WebText: 45M links -> 8M docs (40GB)

# Autoregressive Language Model

With this training, essentially we obtain an oracle: $p(x_k|x_{<k})$, an autoregressive model for language

GPT2 Model Architecture

GPT2 Dataset

| Parameters | Layers | $d_{model}$ |
|---|---|---|
| 117M | 12 | 768 |
| 345M | 24 | 1024 |
| 762M | 36 | 1280 |
| 1542M | 48 | 1600 |

Links from Reddit Posts that have at least 3 karmas

WebText: 45M links -> 8M docs (40GB)

## Right Data + Right Model

# Recap: Transformer

- Transformer
  - Inputs: Input embedding + Position embedding
  - Transformer Layers
    - MultiHeadAttention
    - Regularizations
      - Dropout
      - LayerNorm
      - Residual Connection
    - Fast Forward Neural Networks
- Text Generation
  - Causal Mask for Transformer

# Assignment