

# Attention & Transformers for NLP Tasks



By: John Cisler and Ben Paulson



# BIG QUESTION

Why are transformers, and attention, used in some of the most effective deep-learning models today for NLP?

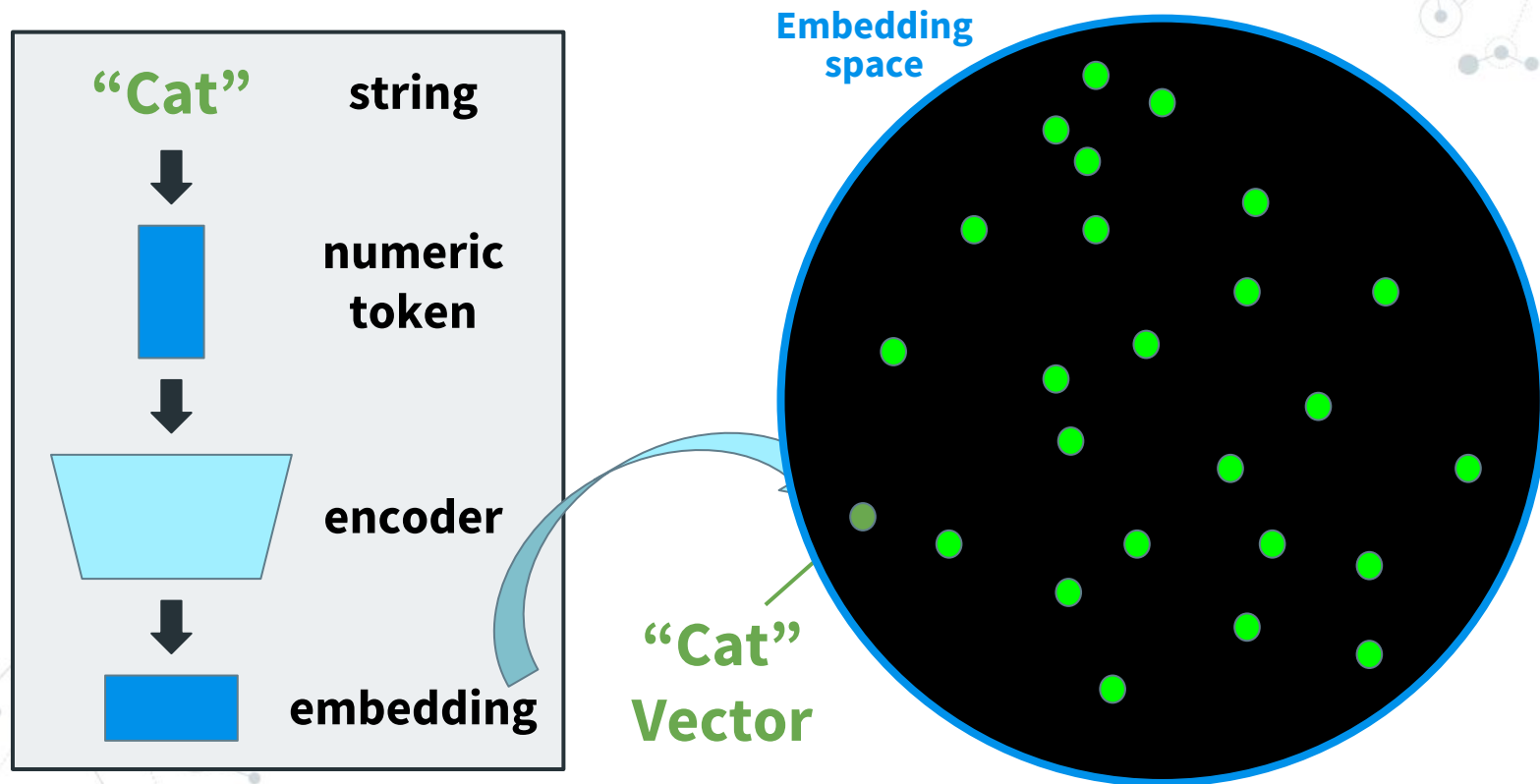
**“Attention Is All You Need” & Language Translation**

# Paper Key Takeaways

- ◎ Language Translation
  - English to French/German
- ◎ Transformer
  - Seq2Seq Model
  - High Throughput
  - Easily parallelizable
  - Mitigated long-range dependency issue
  - State-of-the-art results

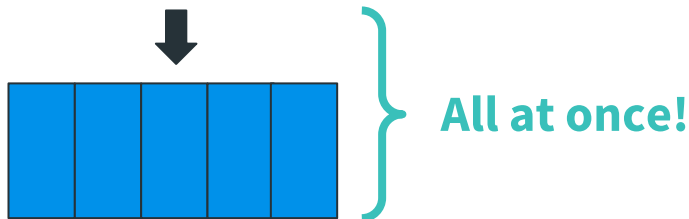
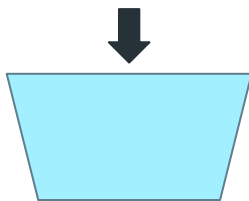


# High Throughput – Embeddings



# High Throughput – Embeddings

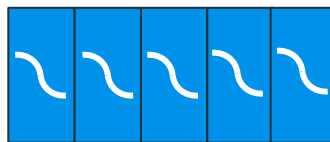
“The cat eats the mouse”



**embeddings**



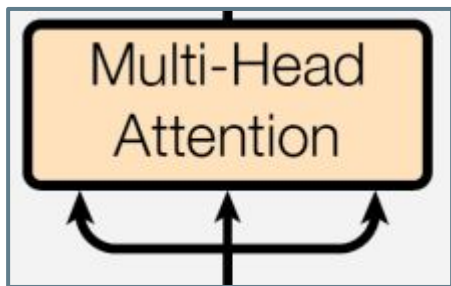
+



**positional**

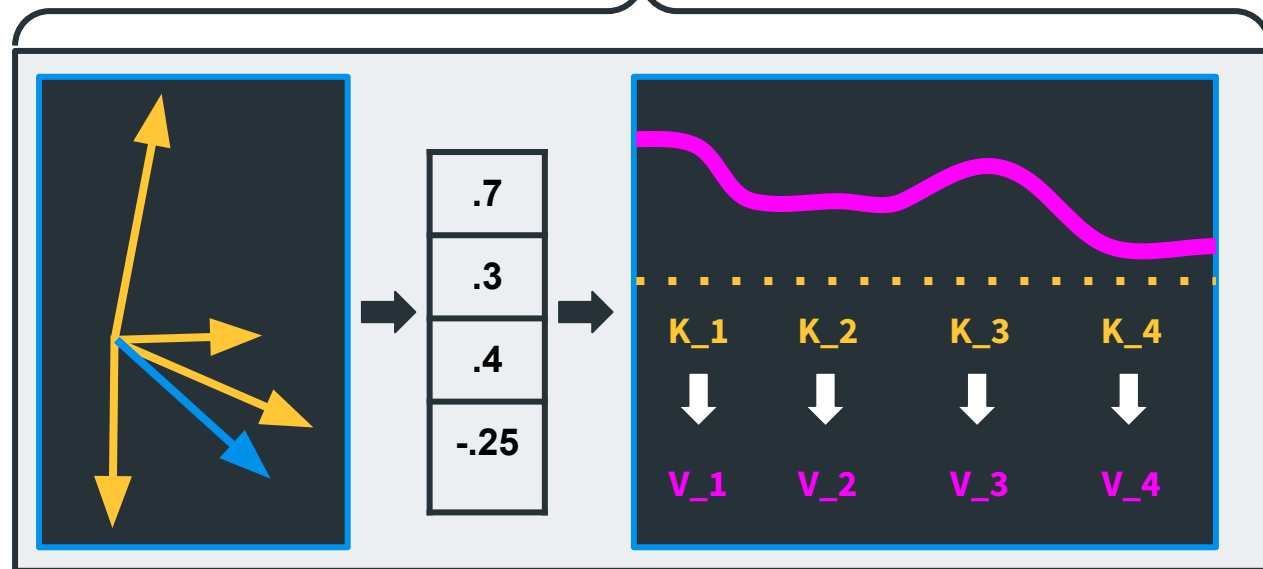
**= order-aware embeddings**

# Long Range – Attention



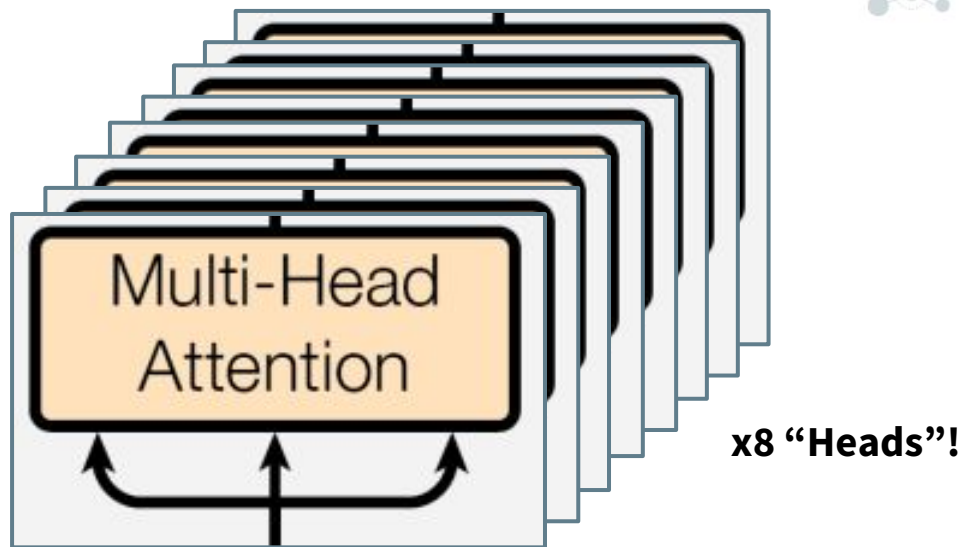
$A(\mathbf{Q}, \mathbf{K}, \mathbf{V})$

$$\Rightarrow A(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_K}}\right) * \mathbf{V}$$



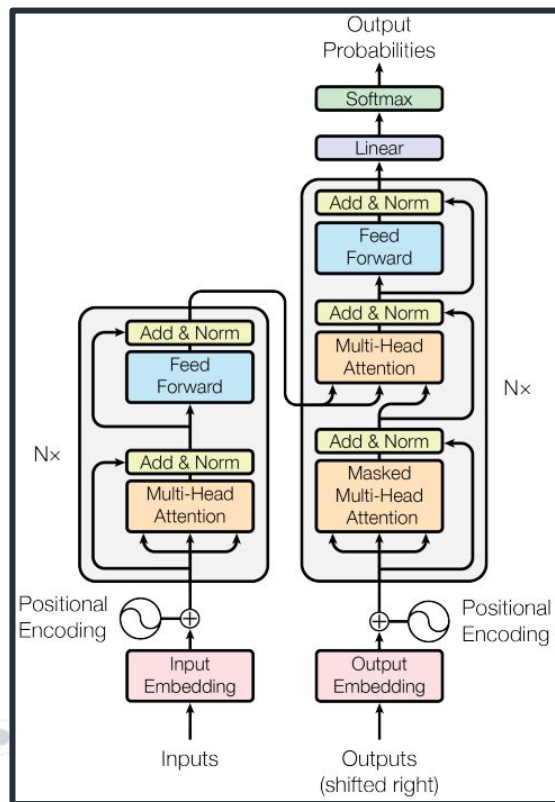
# Parallelization – Multi-Headed

$$A(Q, K, V) =$$



**Multiple heads capture  
more abstract features**

# All Created the Transformer!

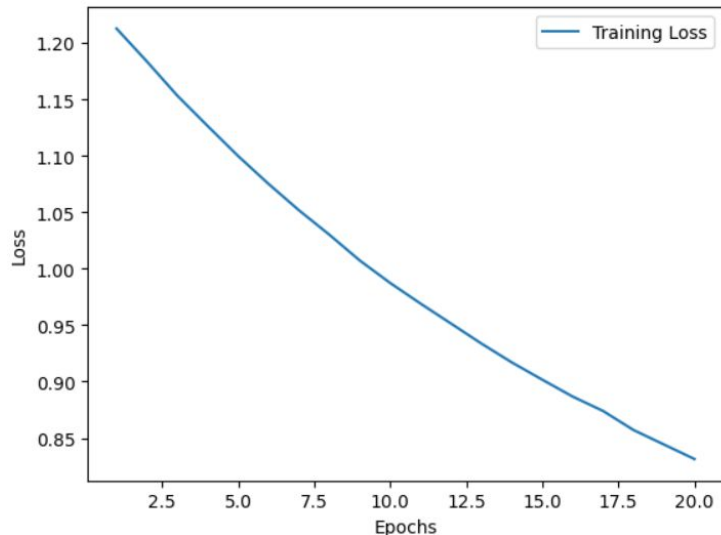


- ⊙ High Throughput
- ⊙ Captures long-range dependencies
- ⊙ Parallelizable

**Next time, we'll walk through inputs as they travel through this transformer**



# Language Translation



```
Epoch 10/20
810/810 [=====] - 181s 224ms/step - loss: 0.9873 - masked_accuracy: 0.7576 - val_loss: 2.1494 -
val_masked_accuracy: 0.6296
Epoch 11/20
810/810 [=====] - 182s 225ms/step - loss: 0.9690 - masked_accuracy: 0.7608 - val_loss: 2.1680 -
val_masked_accuracy: 0.6274
Epoch 12/20
810/810 [=====] - 180s 223ms/step - loss: 0.9513 - masked_accuracy: 0.7642 - val_loss: 2.1657 -
val_masked_accuracy: 0.6332
Epoch 13/20
810/810 [=====] - 182s 224ms/step - loss: 0.9336 - masked_accuracy: 0.7670 - val_loss: 2.1649 -
val_masked_accuracy: 0.6321
Epoch 14/20
810/810 [=====] - 180s 222ms/step - loss: 0.9169 - masked_accuracy: 0.7703 - val_loss: 2.1904 -
val_masked_accuracy: 0.6289
Epoch 15/20
516/810 [=====>.....] - ETA: 1:04 - loss: 0.8831 - masked accuracy: 0.7778
```

```
▶ sentence = 'este é um problema que temos que resolver.'
ground_truth = 'this is a problem we have to solve .'

translated_text, translated_tokens, attention_weights = translator(
    tf.constant(sentence))
print_translation(sentence, translated_text, ground_truth)
```

Input: : este é um problema que temos que resolver.  
Prediction : this is a problem that we have to solve .  
Ground truth : this is a problem we have to solve .

4 Layers vs 6 Layers

30 min vs 20 hours

BLEU Score

# • Roadmap



“Attention is All You Need” Analysis



English -> French implementation in Tensorflow



Trained Model



Effectively explain attention & transformers to a general AI audience



Week 9  
Preso

Trained English to Portuguese Model (w/ metrics)

