

CS5841/EE5841 Machine Learning

Lecture 14: Transformers

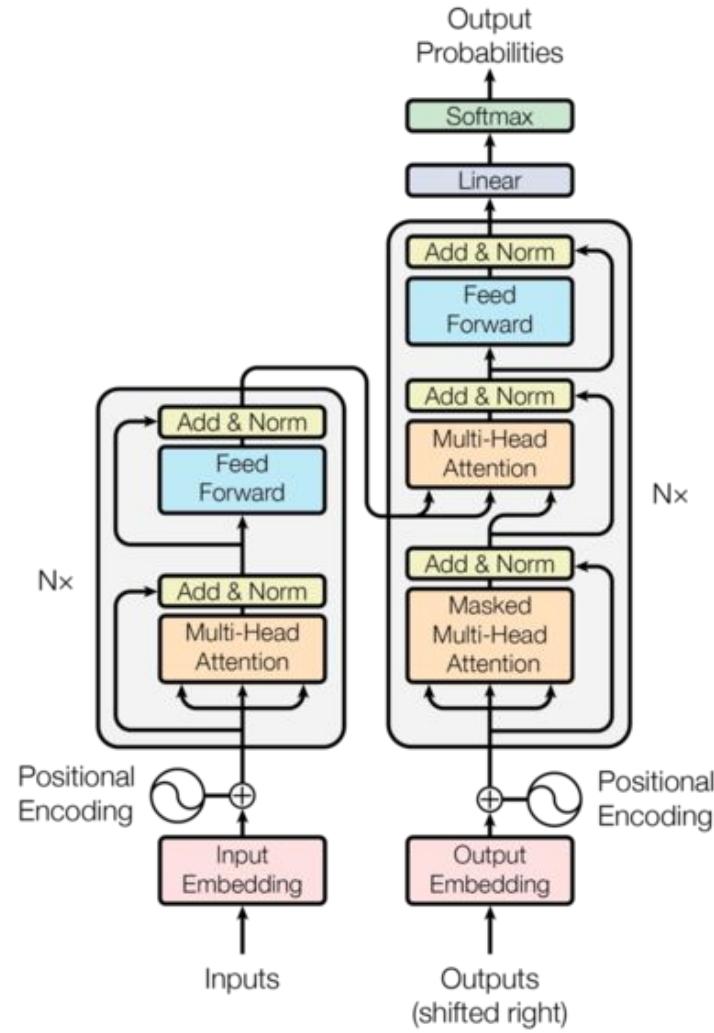
Evan Lucas



Michigan Tech

Transformer Overview

- Built to operate on sequential data
- Uses attention without convolutional or recurrent layers
 - Mimics cognitive attention - focuses on what aspects of input are relevant
- Input to Transformer includes a positional encoding, either absolute or relative
- Can be bidirectional or unidirectional
- Often used in an encoder-decoder format, such as in the image on the right
- Introduced in 2017 by Google and affiliates [1]



Original Transformer
architecture from [1]

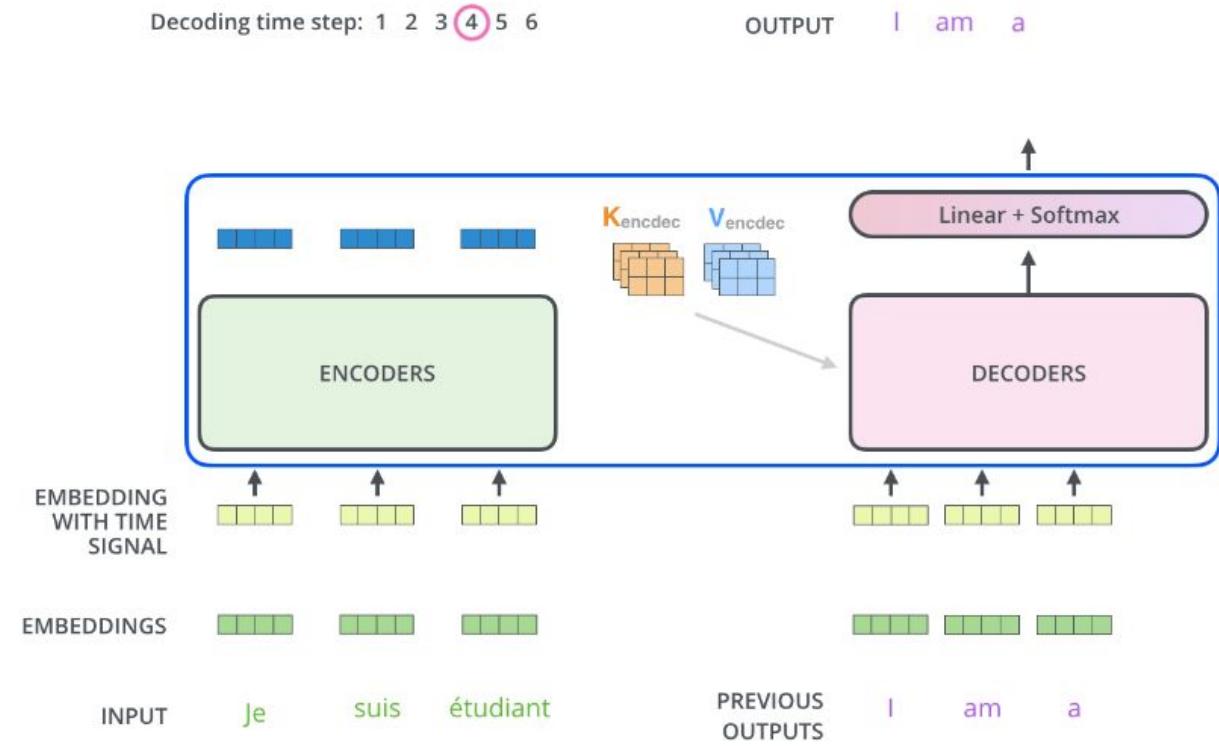
1. Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems*. 2017.



Transformer Architecture Discussion

For summarization, we use encoder-decoder architectures

The decoder uses encoded input plus previous outputs to determine next output

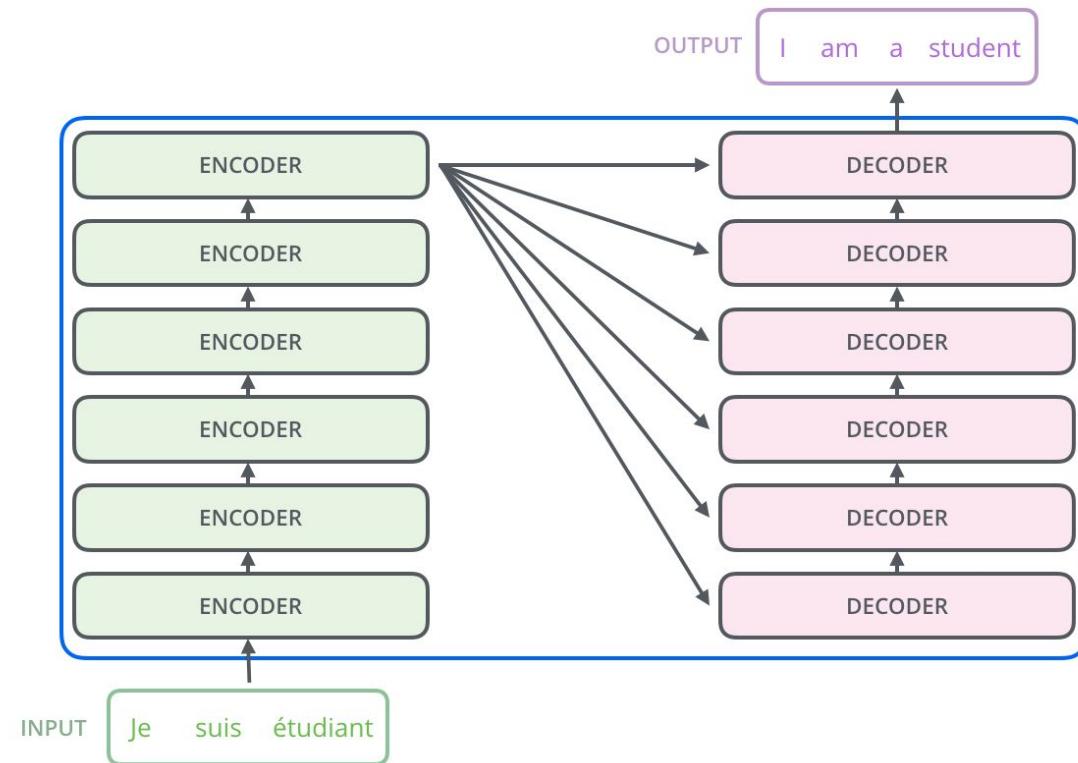


Other transformer architectures

- Encoder only (ex: BERT)
 - Trained by making model predict missing/corrupted tokens (Autoencoding pretraining)
 - Better suited for classification type tasks
- Decoder only (ex: GPT2/3)
 - Trained by having model predict next token (Autoregressive pretraining)
 - Well suited for text generation
- Encoder-Decoder (ex: Longformer Encoder-Decoder)
 - Multiple training methods, including previous as well as sequence to sequence generation
 - Used for summarization, translation, and question answering

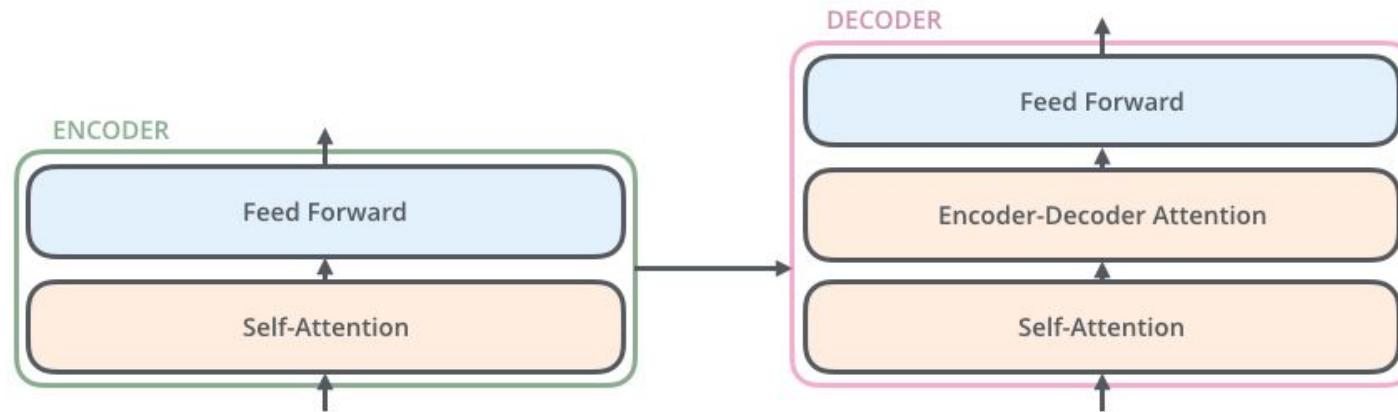
Transformers Step-by-Step

- Input is fed through several layers of encoders
- Final encoder output is fed into each decoder layer



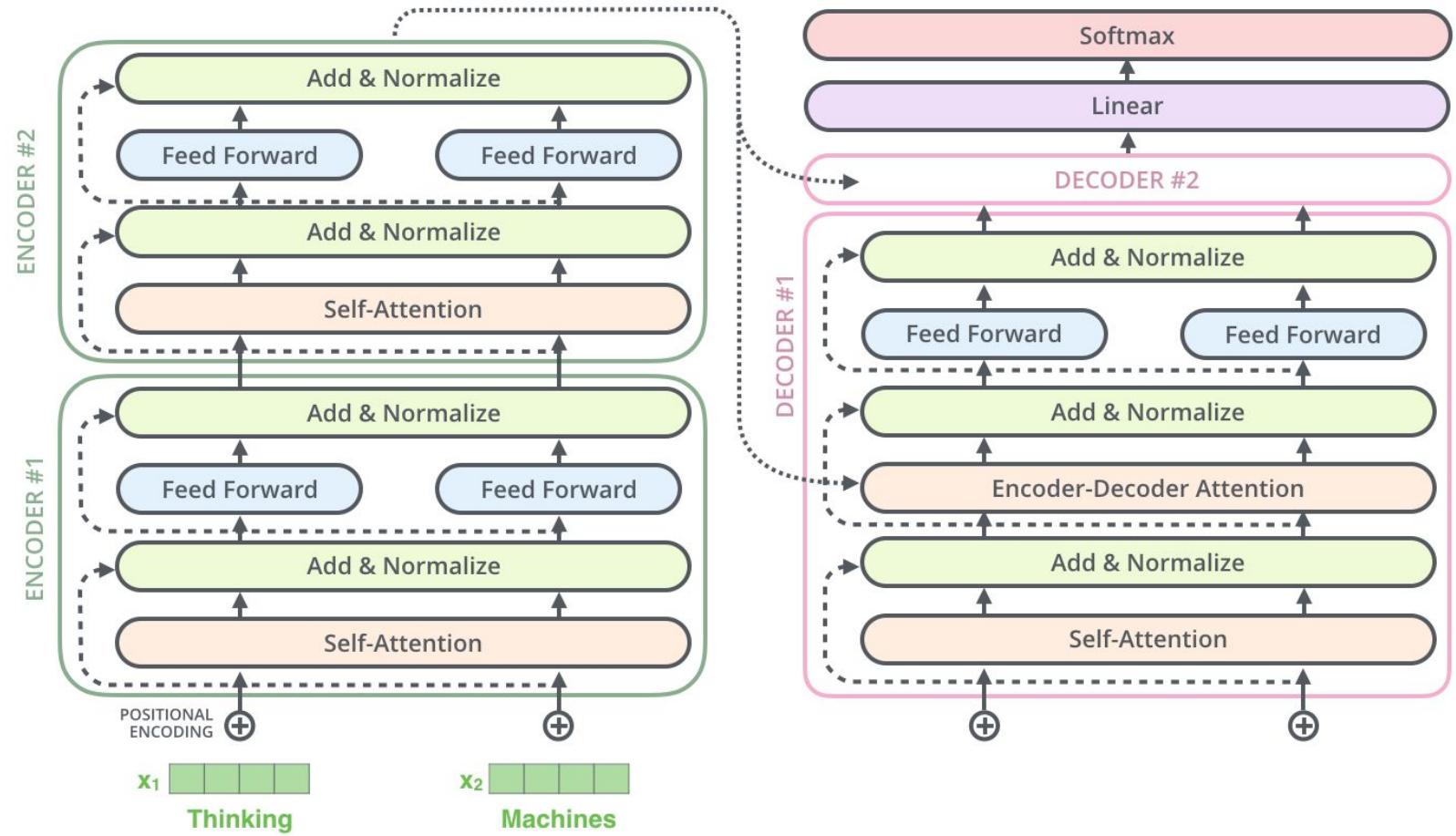
Transformers Step-by-Step

- Each encoder layer includes two main stages - multi-headed self attention on input and a feedforward layer
- Each decoder layer includes self-attention, cross-attention, and a feedforward layer
- Residual skip connections not shown for clarity



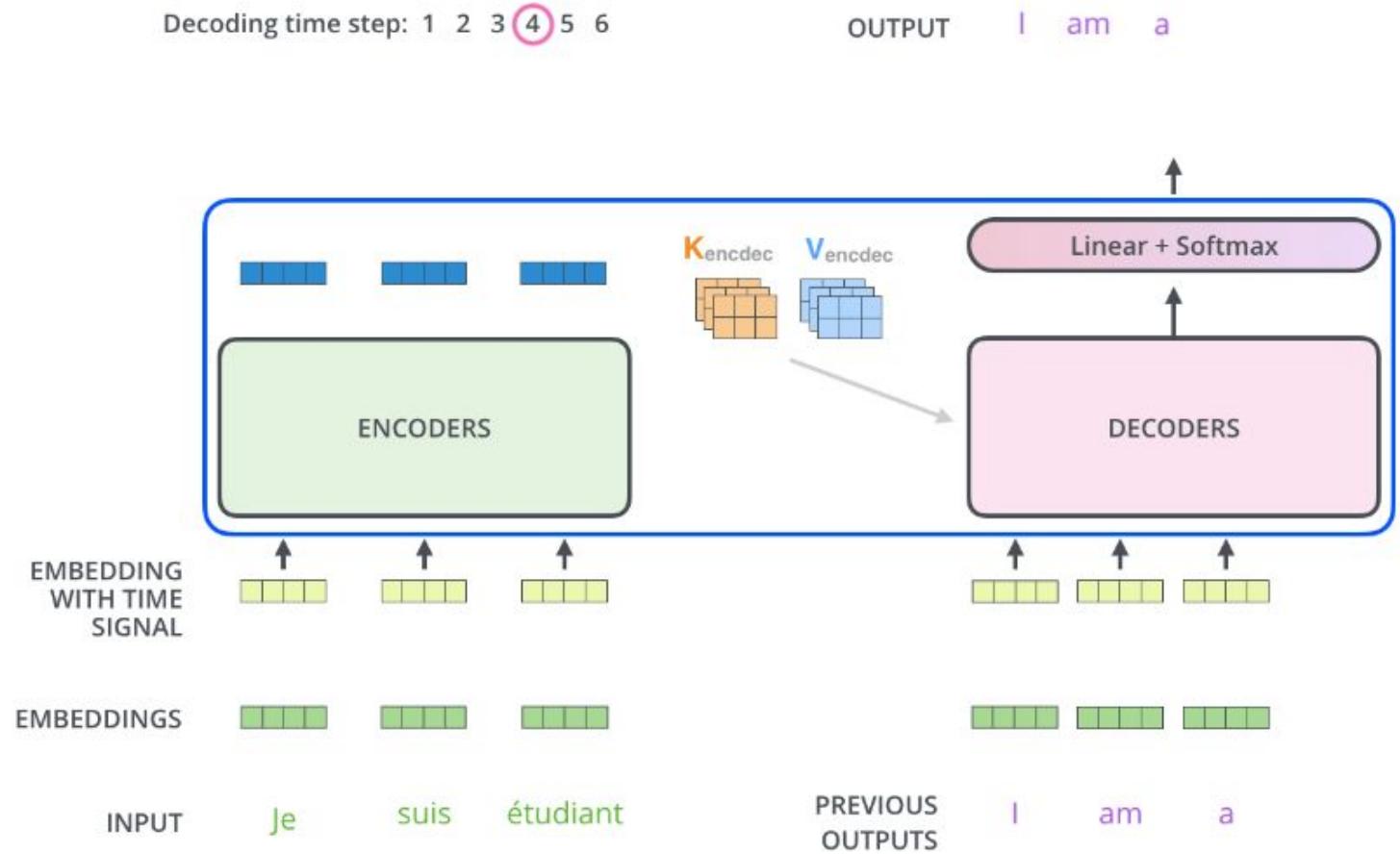
Transformers Step-by-Step

- A more detailed version of the previous slide
- Final output uses Softmax to select output token



Transformers Step-by-Step

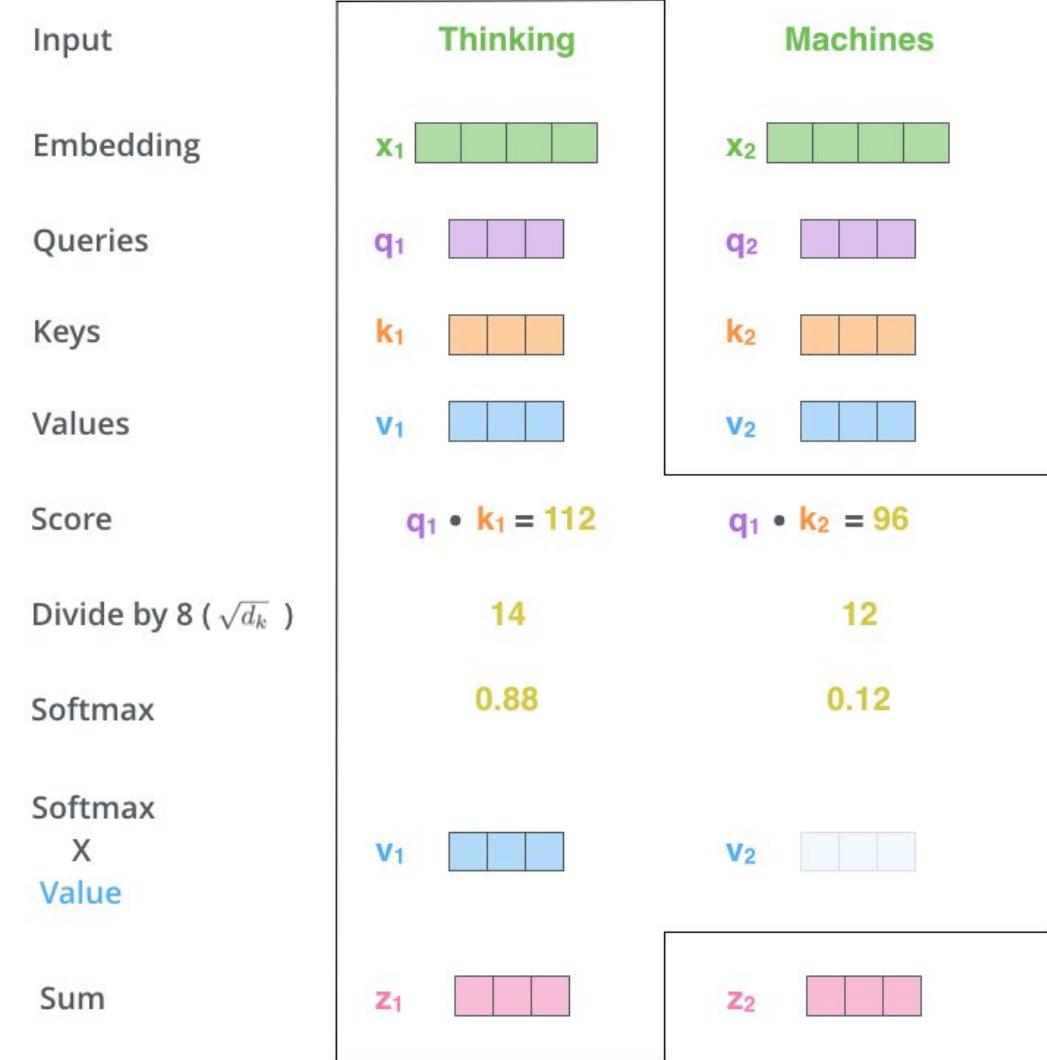
- Decoder outputs one token at a time
- Previous outputs are one of the inputs to the decoder



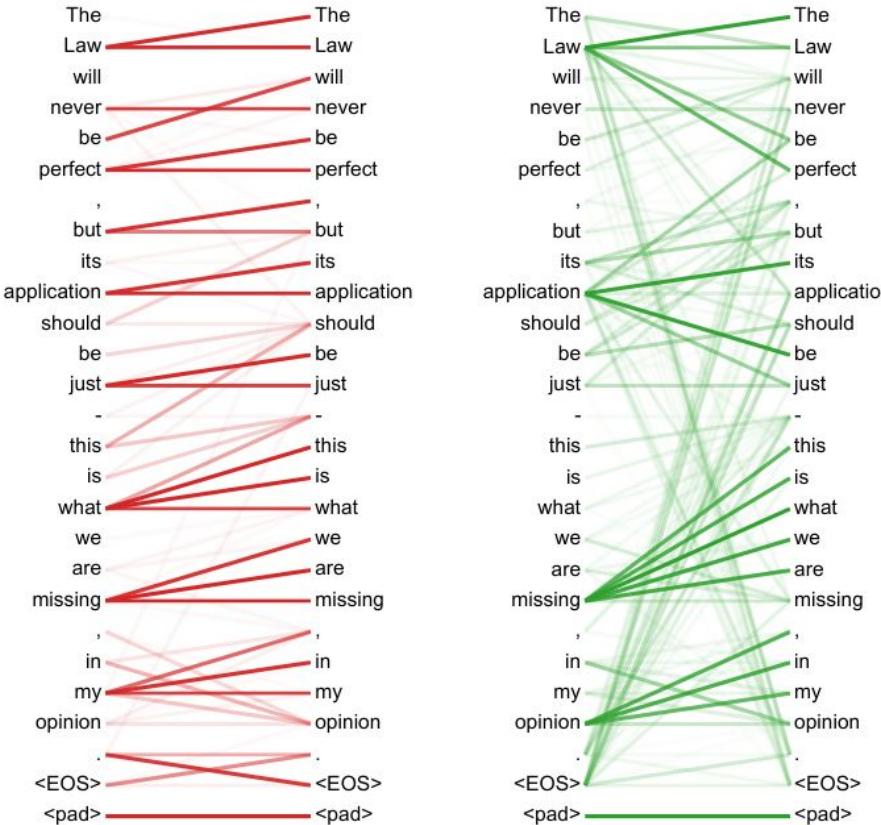
Attention

- Each input embedding is multiplied by a learned weight to create “queries”, “keys”, and “values”
 - These can be considered representations of the tokens
- Attention is calculated using multiple heads that learn different aspects

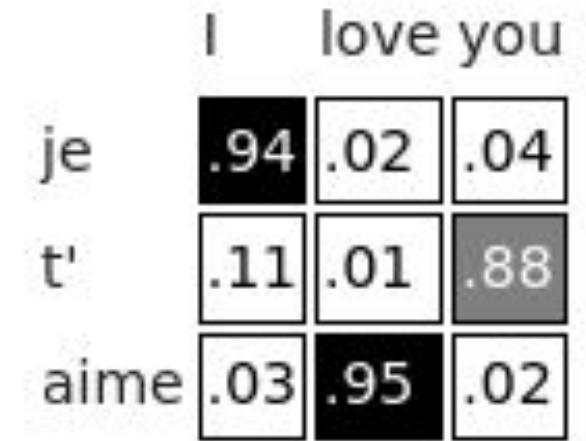
$$\text{softmax} \left(\frac{\begin{matrix} Q \\ \times \\ K^T \end{matrix}}{\sqrt{d_k}} \right) V = Z$$



Attention



Example of two different
self-attention heads
<https://arxiv.org/abs/1706.03762>



Example of single head
cross-attention weight matrix for
translation
[https://en.wikipedia.org/wiki/Attention_\(machine_learning\)](https://en.wikipedia.org/wiki/Attention_(machine_learning))



Michigan Tech

Tokenizing

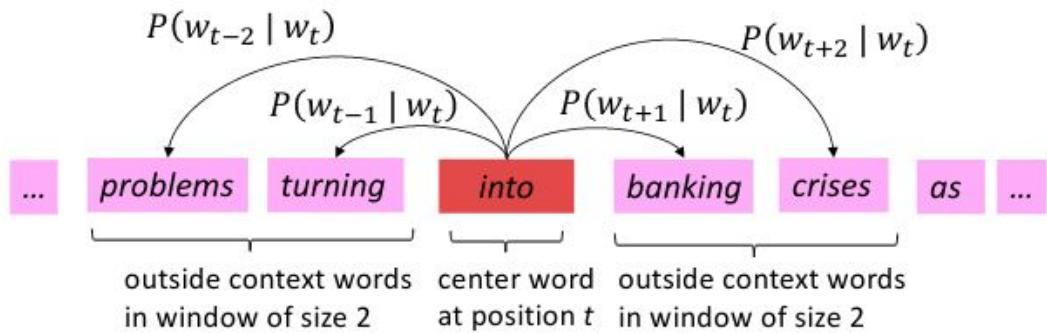
- Tokenizing - preprocessing data to represent text as finite dictionary of words
- Common words get tokens, uncommon words are built from other tokens, non-word tokens are used for structure and model inputs
- Typical dictionary sizes are in the 10,000-50,000 range
- Different strategies - BytePair Encoding (BPE), WordPiece, SentencePiece



Michigan Tech

Word Embeddings

- An n-dimensional representation of the word's meaning
- Historically derived from simple probabilistic models that predict a word given context
- One example: Word2Vec (2013)
- Figures are from a low dimensional example



$$\text{expect} = \begin{pmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 0.109 \\ -0.542 \\ 0.349 \\ 0.271 \\ 0.487 \end{pmatrix}$$



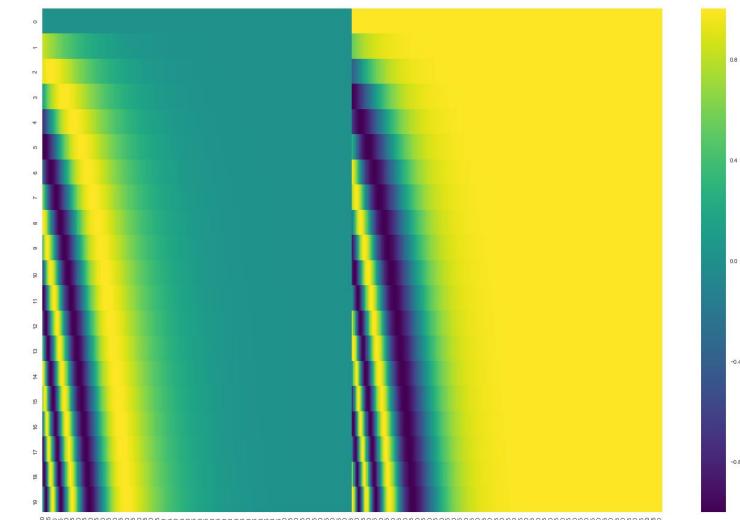
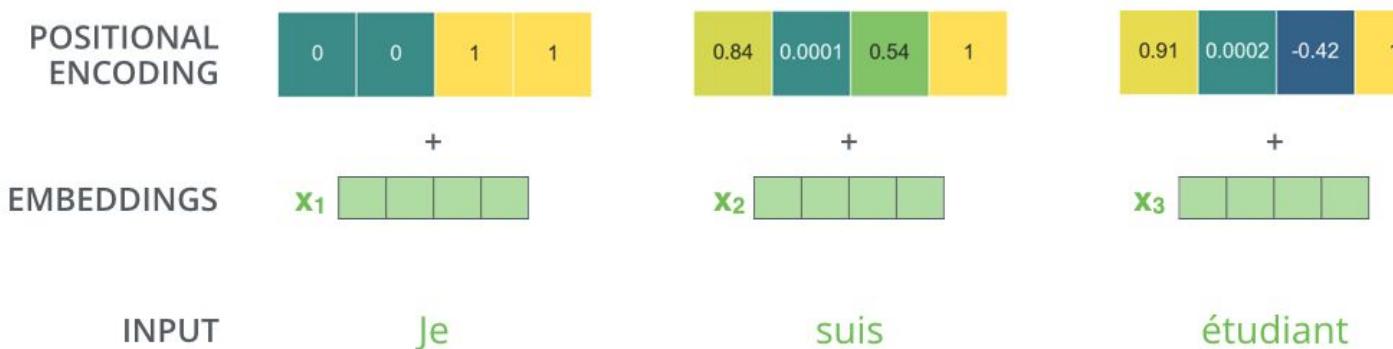
<https://projector.tensorflow.org/>



Michigan Tech

Positional Embeddings

- Values to include representation of location in input are added
 - Current standard is sine based
- A current topic of research.
 - Alternate positional embedding methods
 - Are they needed at all?
 - ALiBi - a modification of positional embedding

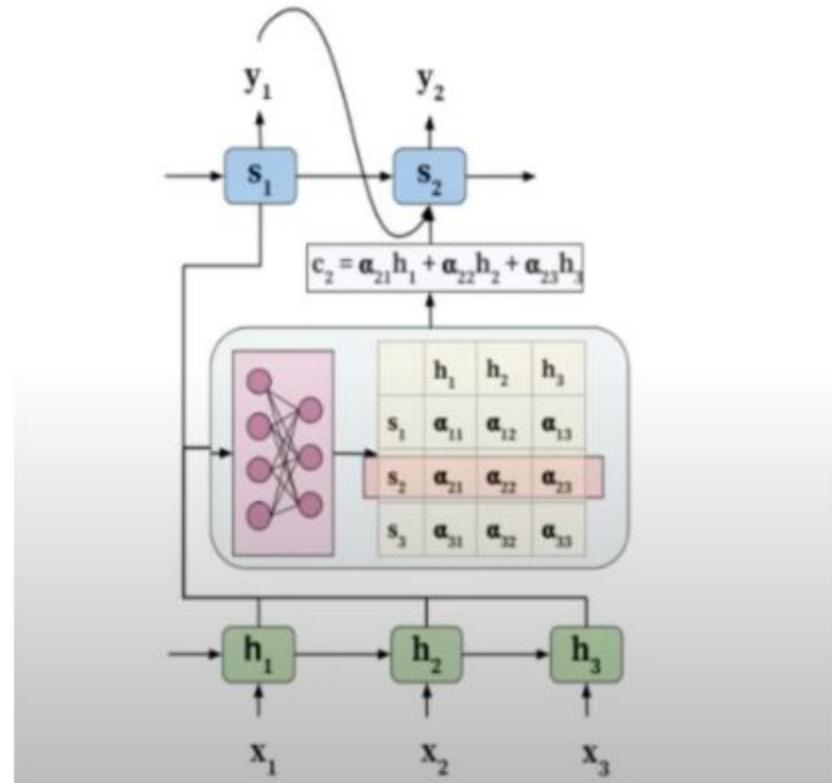


A real example for 20 tokens (rows)
with a 512 embedding size (columns)
using sine and cosine

Previous Work

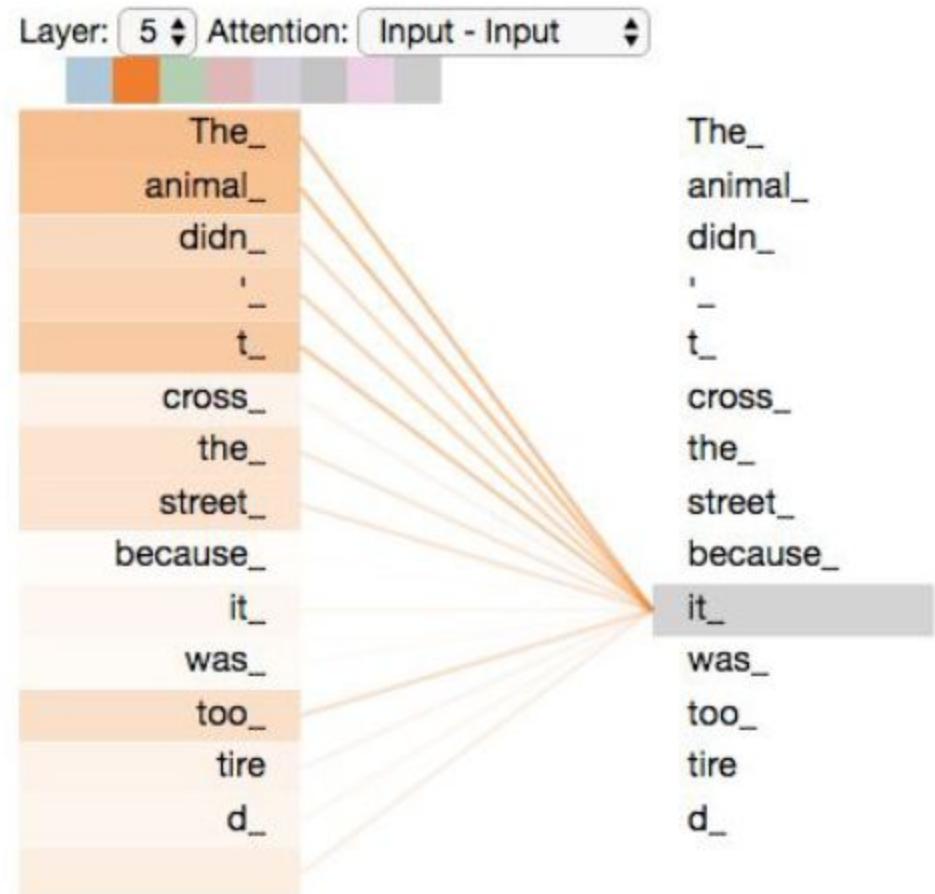
- Pre-2014 - RNN, LSTM, GRU
 - RNNs struggled with long sequence data
 - Could not handle long sequence lengths (forgetful)
 - Example on next slide
- Neural Machine Translation by Jointly Learning to Align and Translate (Bahdanau, 2014)
 - Landmark Paper
 - Proposed New Architecture:
 - bidirectional RNN as an encoder and an RNN decoder, with an attention mechanism in between
 - Proposed initial idea of attention
 - Model gave attention to particular hidden states when decoding each word
 - Learned which words to give attention to during translation for each word of decoder
 - Still required on RNN

Encoder-decoder model with Attention



Self-Attention

- Each element attends to every other element
- It is an attention mechanism relating different positions of a single sequence in order to compute a representation of the same sequence
- Each element becomes query, key, and value from the input embeddings by multiplying by a weight matrix



<https://jalammar.github.io/illustrated-transformer/>



Michigan Tech

Idea Behind Attention

Goal: Learn (differentiable) how to pick relevant information from input data

1. Create three vectors from each of the encoder's input values (**query**, **key**, **value**)
2. Calculate a score for how much to focus on each part of the input when we encode words at specific positions
3. How?
 - o Actual math on the next slide, but the idea is to select a **value** (referenced by a **key**) relevant to a **query** (what trying to pull from input)

- ❖ From Paper - “A mapping a query and a set of key-value pairs to an output, where the query, keys, values, and output are all vectors. The output is computed as a weighted sum of the values, where the weight assigned to each value is computed by a compatibility function of the query with the corresponding key” - (Vaswani, 2017)

key	value
name	Quinn
position	quarterback
Handedness	right

Keys and values are actually word embeddings, not words



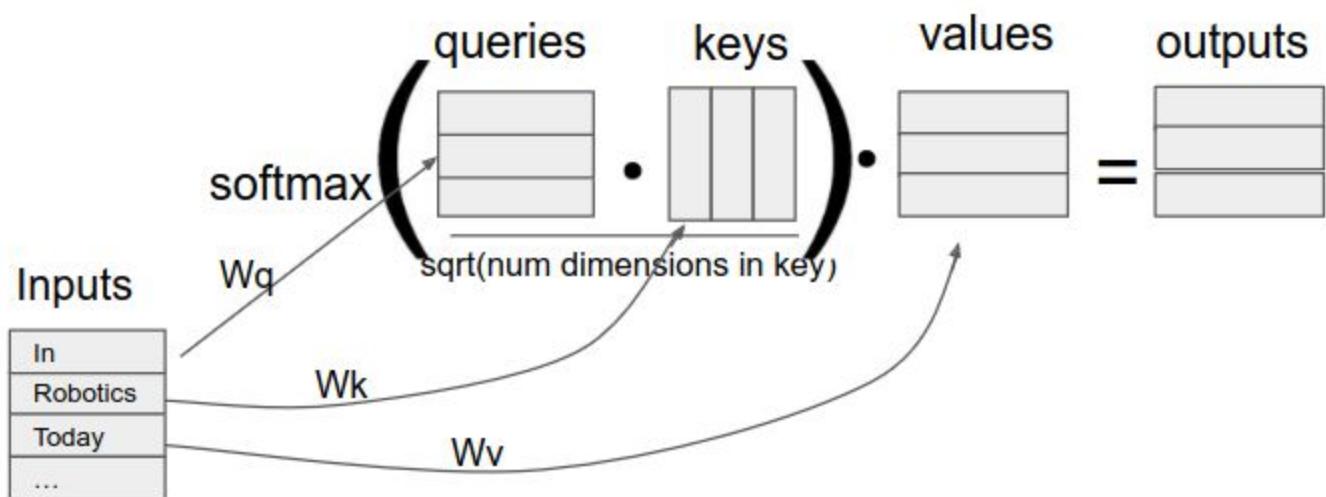
Michigan Tech

Scaled Dot Product Attention

- Key algorithm in transformers
- Attention weights: how likely each query matches key (used for weighted sum of values)
- Dividing by d_k makes algorithm easier to train
 - prevents binary prediction of one 1 and a whole bunch of zeros
 - helps the gradients flow
- then can take loss, backpropagate, update the weights and values
- NOT stepping through a sequence like with a RNN

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

- input consists of queries and keys of dimension d_k , and values of dimension d
- Queries packed into matrix Q
- The keys and values are also packed together into matrices K and V .



Multi-Head Attention

- Idea:
 - a. Stack linear layers (weight matrices without biases) that are independent each for keys, queries, values.
 - b. Concatenate output of attention heads to form (plus non-linearity) output layer
- Why?
 - a. Allows for model to focus on different positions
 - b. Gives attention layer multiple “representation subspaces”
 - c. No longer need to oversaturate one attention mechanism

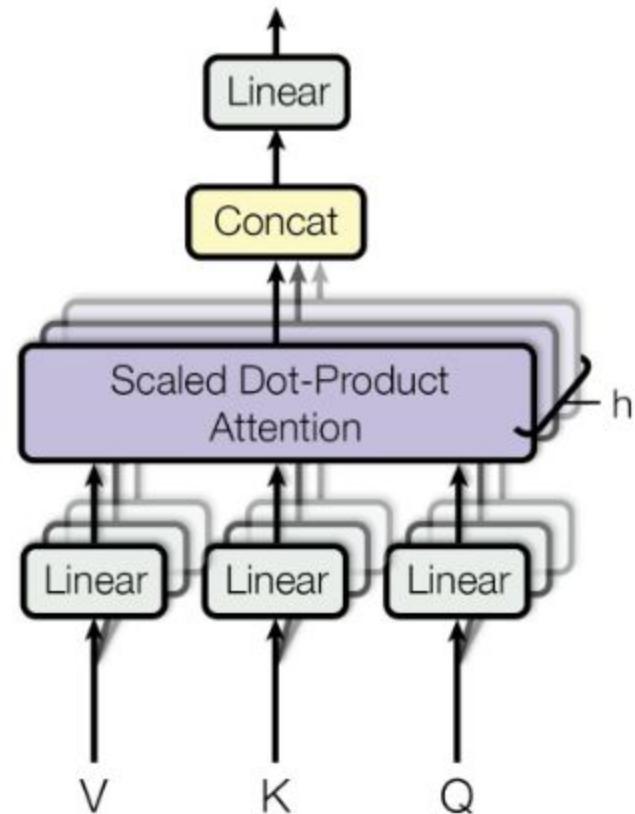
$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

$h = 8$ parallel attention layers, or heads.

Learnable parameter matrices

Multi-Head Attention



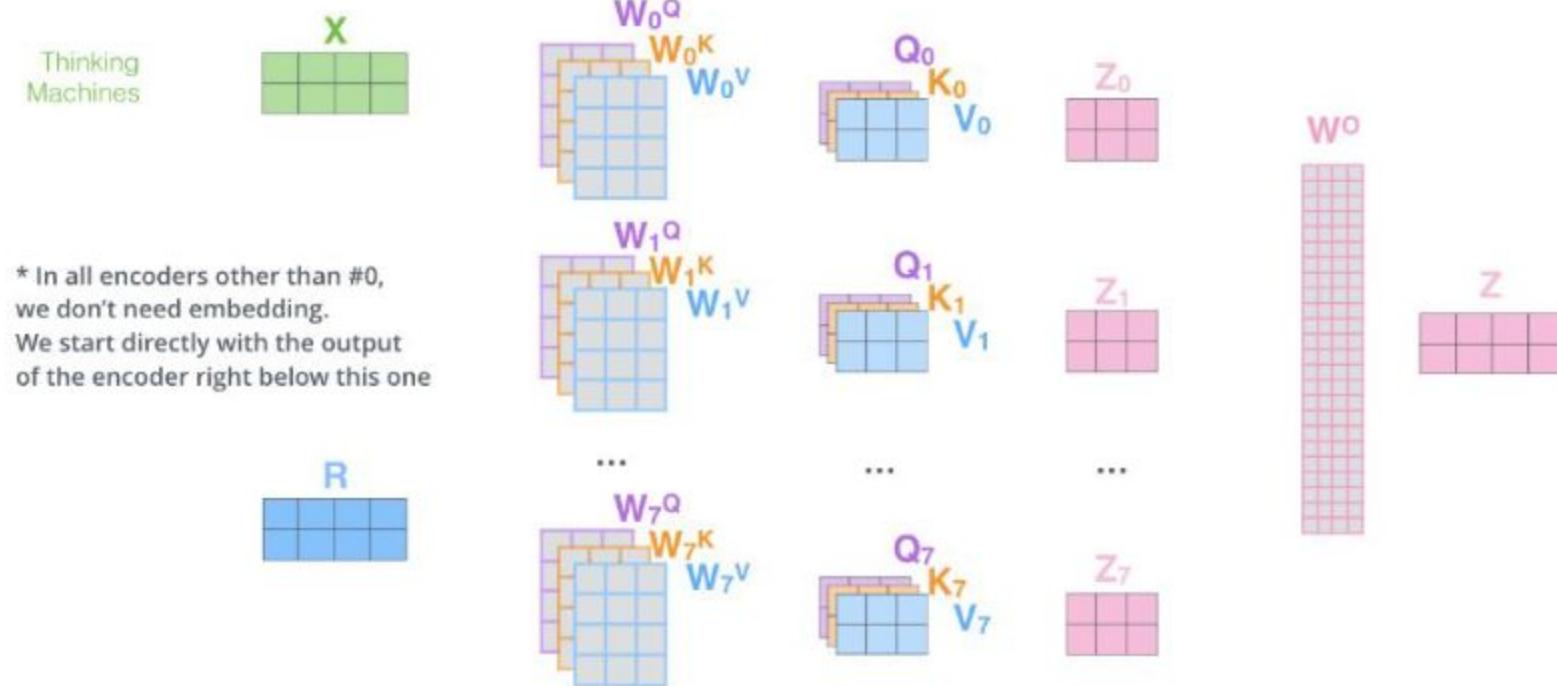
Michigan Tech
1885

Multi-Head Attention Continued

Key points:

- We calculate 8 different attention heads, but we need to combine them
- Attention heads are independent of each other

- 1) This is our input sentence* each word*
- 2) We embed each word*
- 3) Split into 8 heads. We multiply X or R with weight matrices
- 4) Calculate attention using the resulting $Q/K/V$ matrices
- 5) Concatenate the resulting Z matrices, then multiply with weight matrix W^o to produce the output of the layer



<https://jalammar.github.io/illustrated-transformer/>



Michigan Tech

Positional Encoding

- Need for information about the position and order of tokens in a sequence
- Positional Embedding: Vector that represents position of each token
 - Element wise addition the position embedding to the word embedding vector
 - Positional embedding be fixed or learned

"We chose this function because we hypothesized it would allow the model to easily learn to attend by relative positions, since for any fixed offset k , $\text{PE}_{\text{pos}+k}$ can be represented as a linear function of PE_{pos} ."

For every sine-cosine pair corresponding to frequency ω_k , there is a linear transformation $M \in \mathbb{R}^{2 \times 2}$ (independent of t) where the following equation holds:

$$M \cdot \begin{bmatrix} \sin(\omega_k \cdot t) \\ \cos(\omega_k \cdot t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k \cdot (t + \phi)) \\ \cos(\omega_k \cdot (t + \phi)) \end{bmatrix}$$
$$M_{\phi,k} = \begin{bmatrix} \cos(\omega_k \cdot \phi) & \sin(\omega_k \cdot \phi) \\ -\sin(\omega_k \cdot \phi) & \cos(\omega_k \cdot \phi) \end{bmatrix}$$

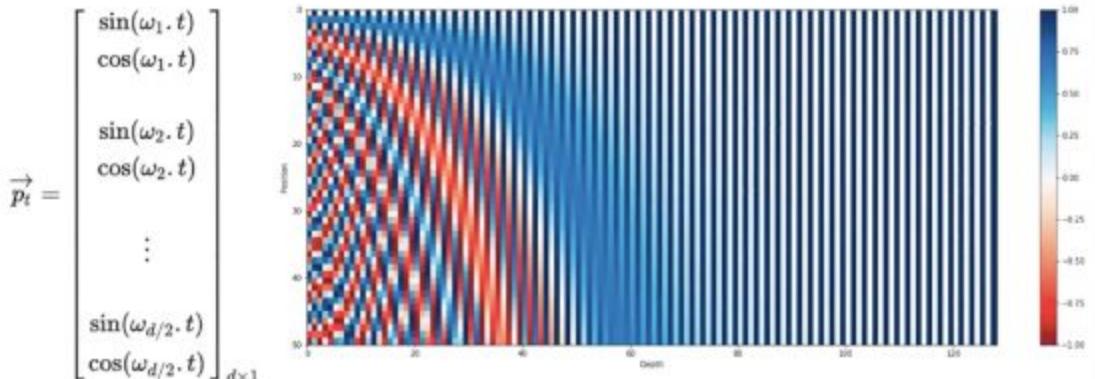
Encoding used in paper

$$\text{PE}_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$\text{PE}_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

pos is the position and i is the dimension.

* Sine method:

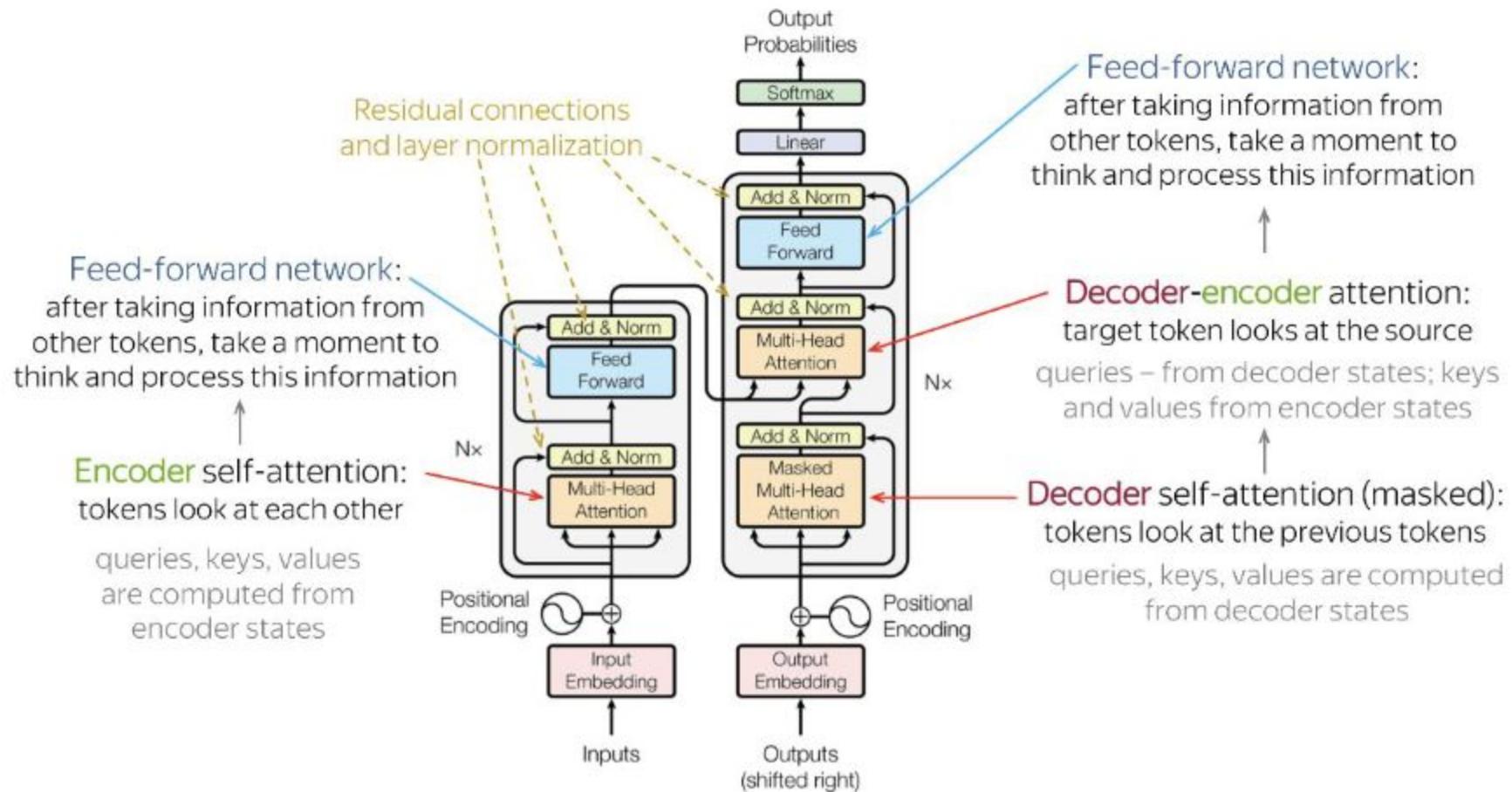


Note: positional embeddings are now learned



Michigan Tech

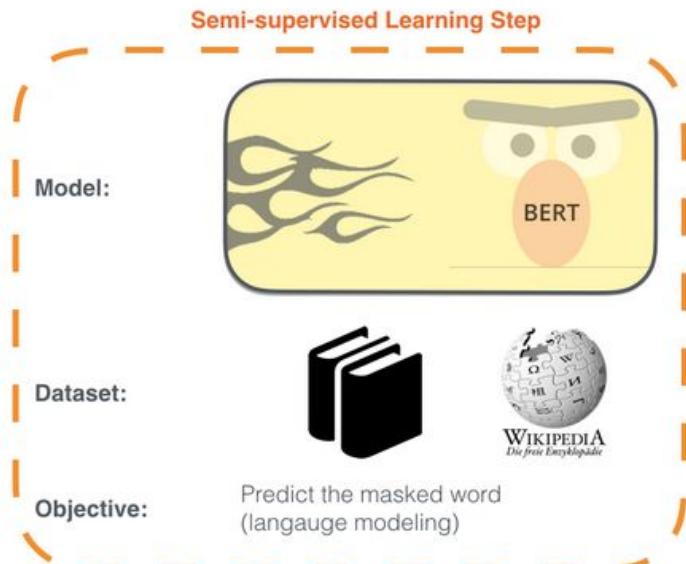
Full-Architecture



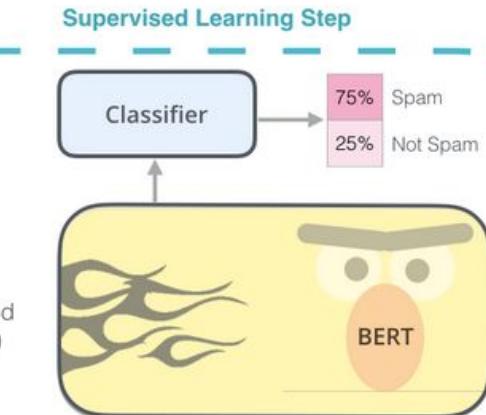
BERT

1 - Semi-supervised training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.



2 - Supervised training on a specific task with a labeled dataset.



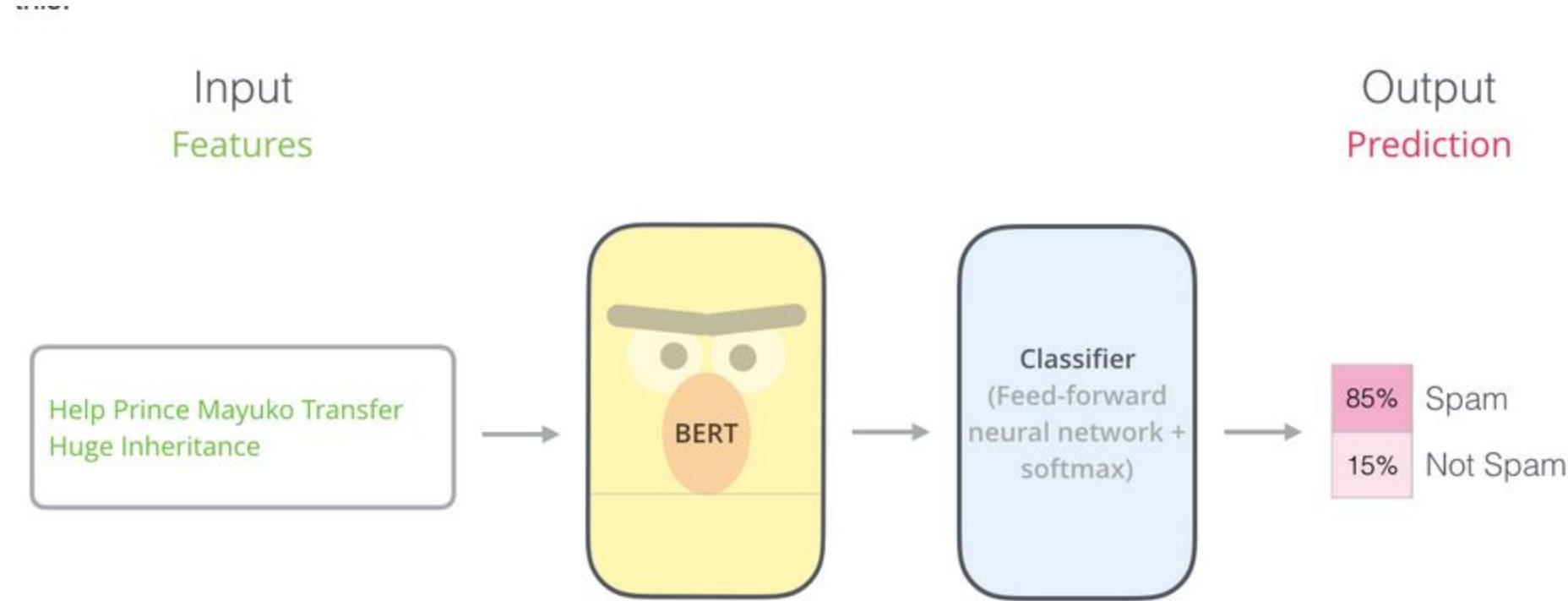
Email message	Class
Buy these pills	Spam
Win cash prizes	Spam
Dear Mr. Atreides, please find attached...	Not Spam

The two steps of how BERT is developed. You can download the model pre-trained in step 1 (trained on un-annotated data), and only worry about fine-tuning it for step 2. [Source for book icon].



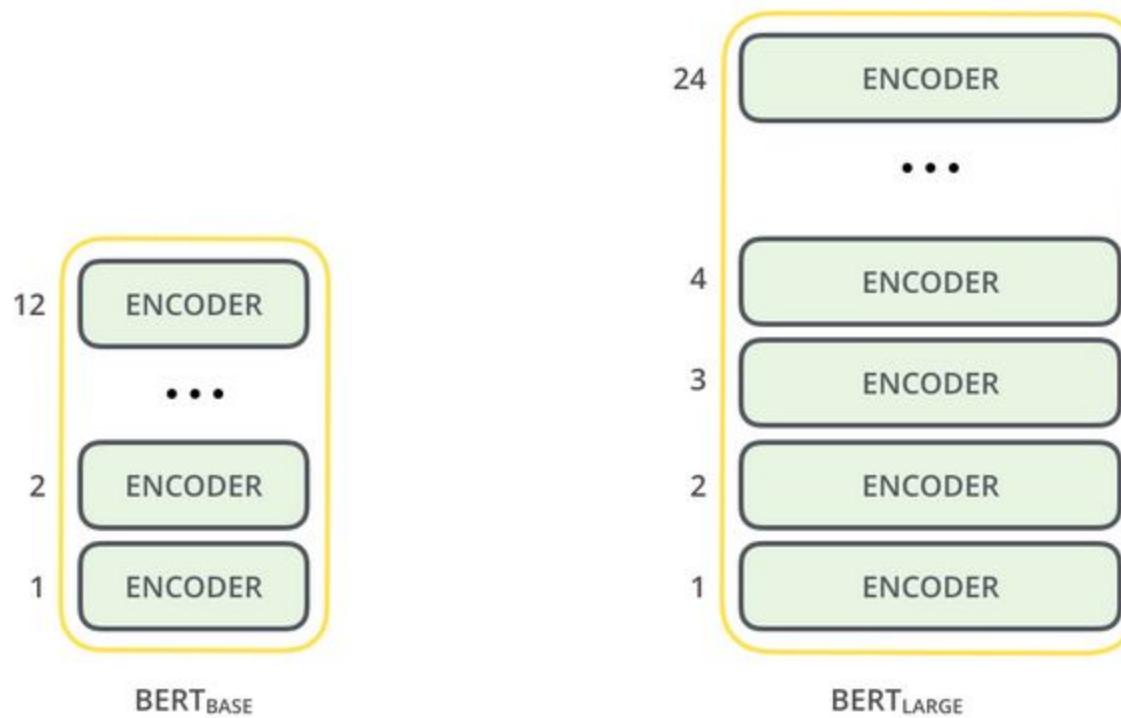
Michigan Tech

Sentence classification example



Michigan Tech

Different model sizes



Michigan Tech

You can download all 24 from [here](#), or individually from the table below.

	H=128	H=256	H=512	H=768
L=2	<u>2/128 (BERT-Tiny)</u>	<u>2/256</u>	<u>2/512</u>	<u>2/768</u>
L=4	<u>4/128</u>	<u>4/256 (BERT-Mini)</u>	<u>4/512 (BERT-Small)</u>	<u>4/768</u>
L=6	<u>6/128</u>	<u>6/256</u>	<u>6/512</u>	<u>6/768</u>
L=8	<u>8/128</u>	<u>8/256</u>	<u>8/512 (BERT-Medium)</u>	<u>8/768</u>
L=10	<u>10/128</u>	<u>10/256</u>	<u>10/512</u>	<u>10/768</u>
L=12	<u>12/128</u>	<u>12/256</u>	<u>12/512</u>	<u>12/768 (BERT-Base)</u>



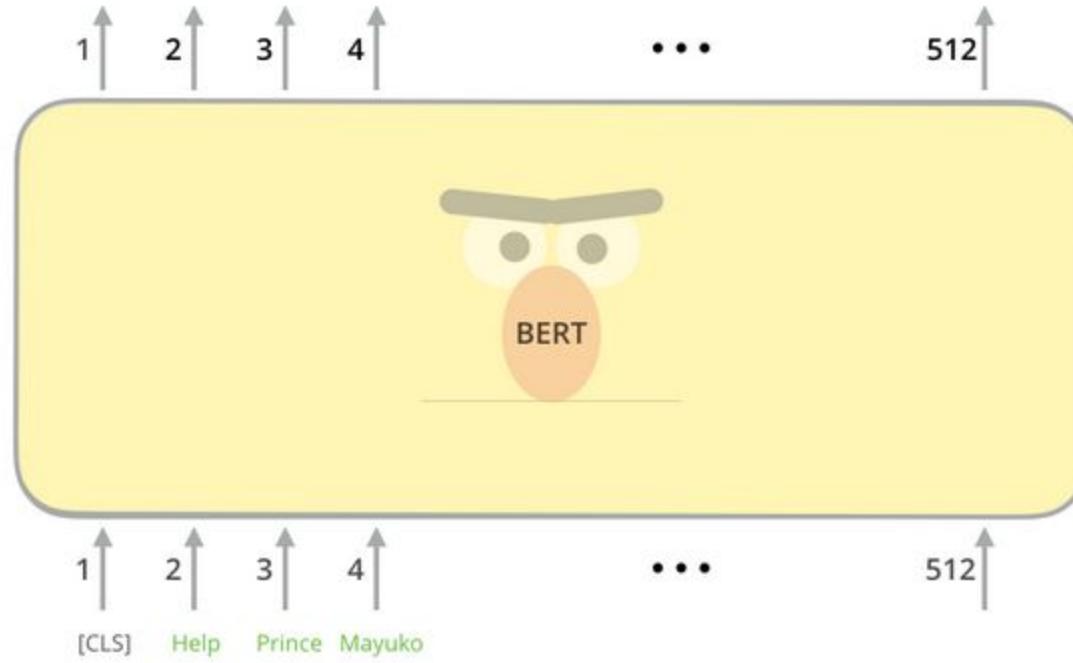
Michigan Tech

Model	Score	CoLA	SST-2	MRPC	STS-B	QQP	MNLI-m	MNLI-mm	QNLI(v2)
BERT-Tiny	64.2	0.0	83.2	81.1/71.1	74.3/73.6	62.2/83.4	70.2	70.3	81.5
BERT-Mini	65.8	0.0	85.9	81.1/71.8	75.4/73.3	66.4/86.2	74.8	74.3	84.1
BERT-Small	71.2	27.8	89.7	83.4/76.2	78.8/77.0	68.1/87.0	77.6	77.0	86.4
BERT-Medium	73.5	38.0	89.6	86.6/81.6	80.4/78.4	69.6/87.9	80.0	79.1	87.7



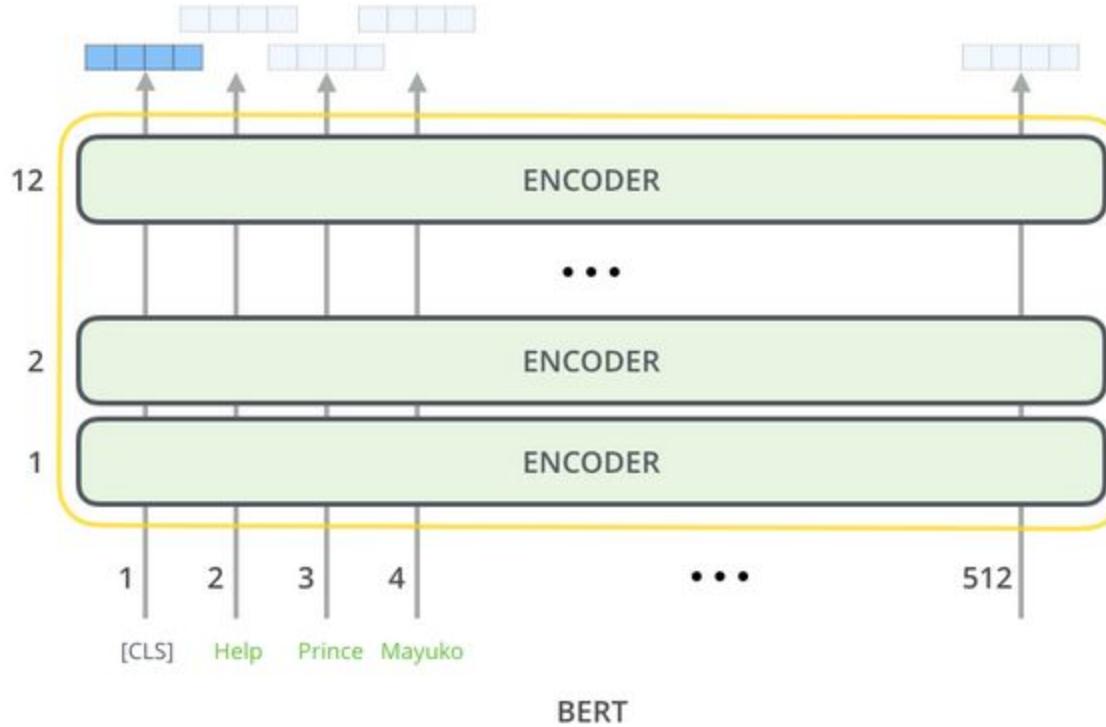
Michigan Tech

Data passthrough for BERT



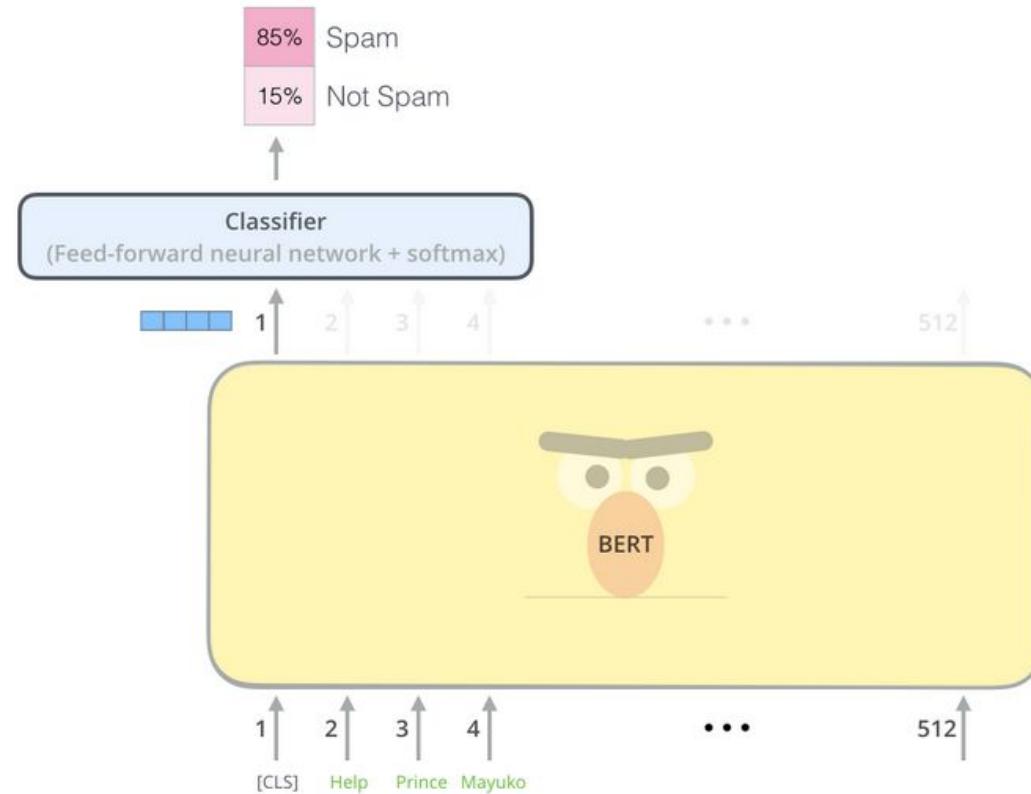
Michigan Tech

Data passthrough for BERT



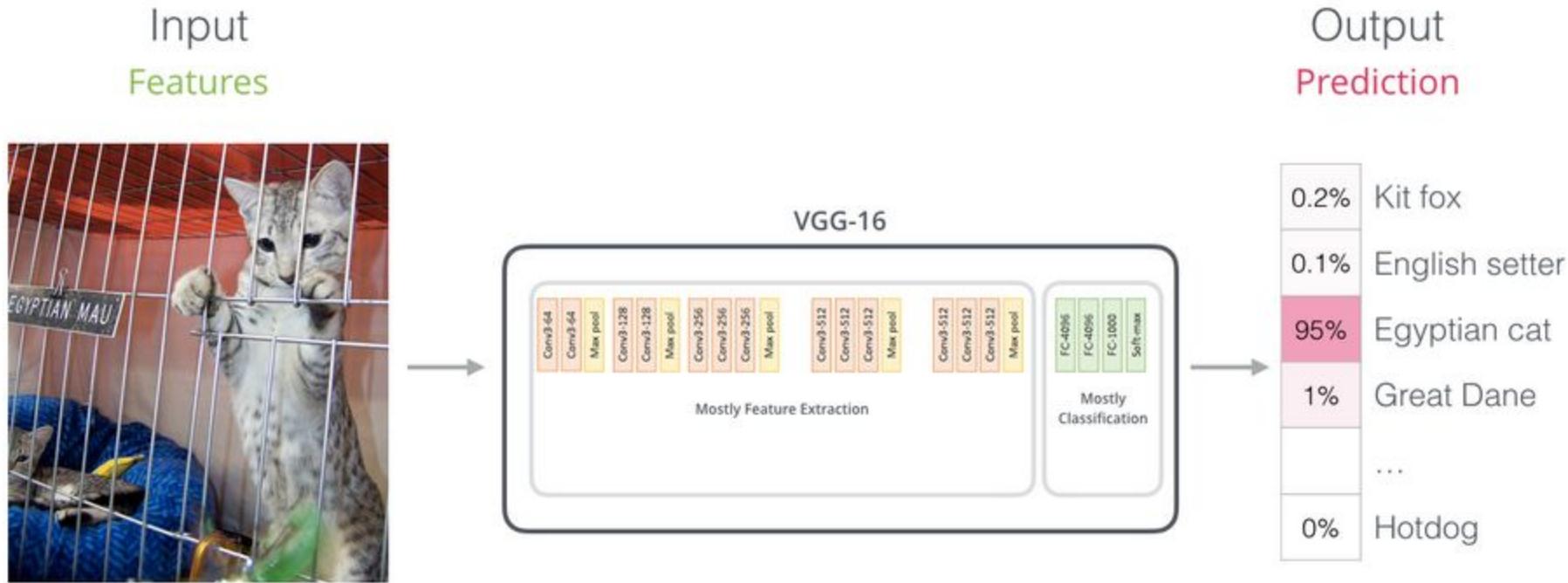
Michigan Tech

Transfer learning with BERT



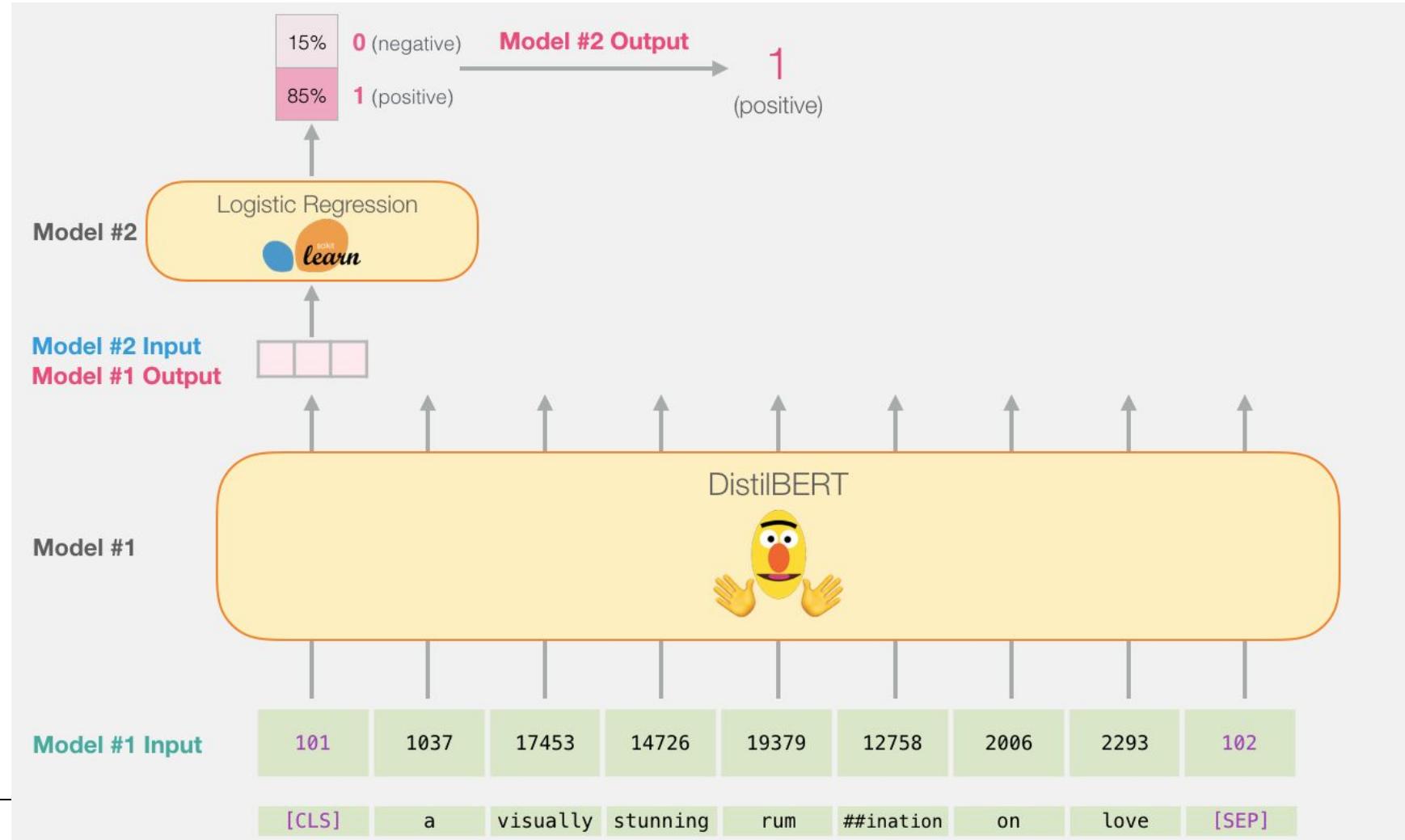
Michigan Tech
1885

Similarity to transfer learning with CNN

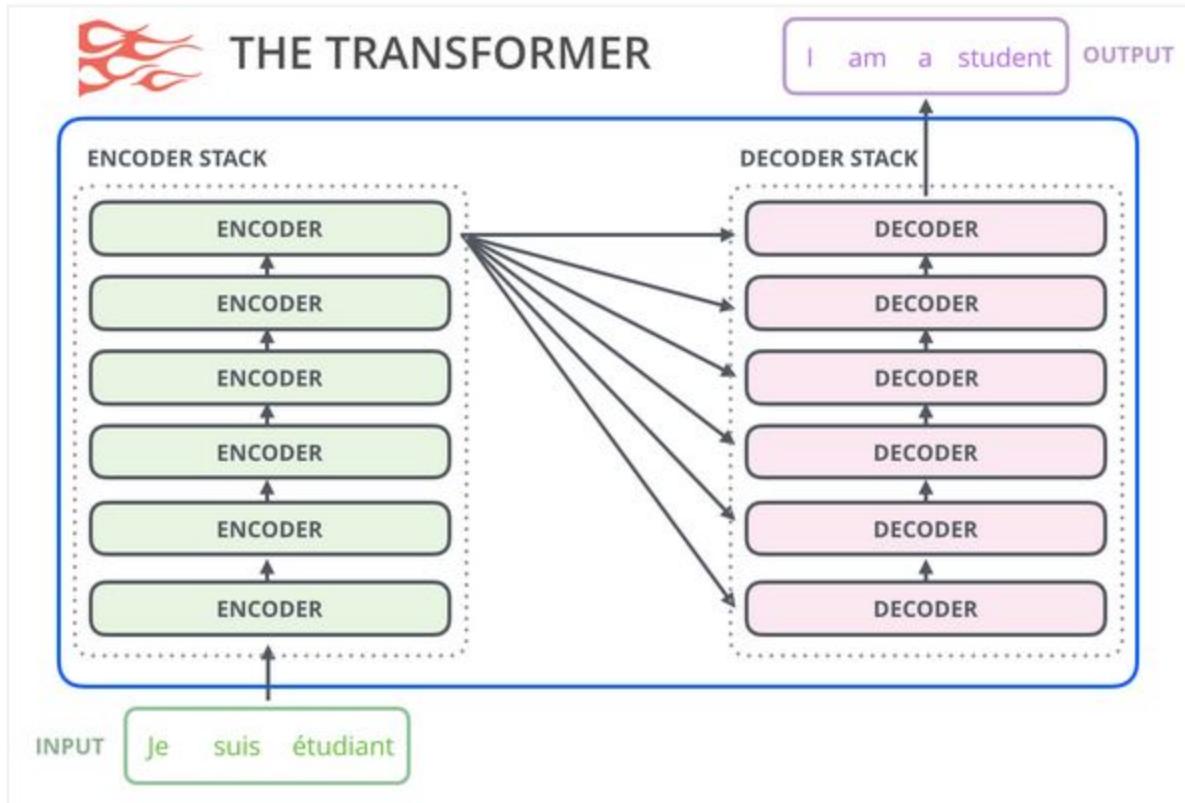


Michigan Tech

BERT as featurizer for other model



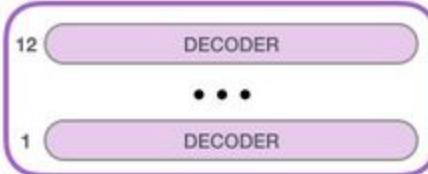
Encoder-Decoder



Michigan Tech



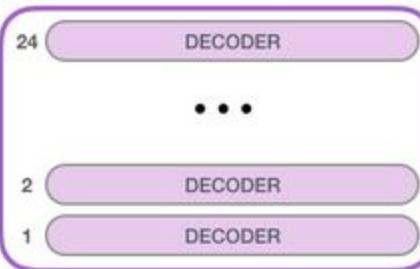
GPT-2
SMALL



Model Dimensionality: 768



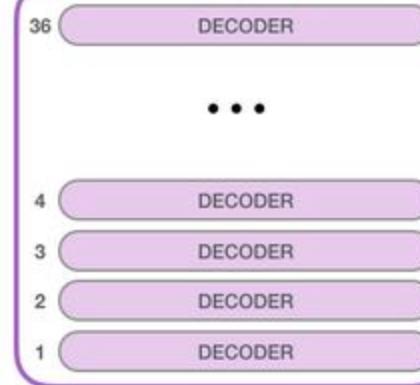
GPT-2
MEDIUM



Model Dimensionality: 1024



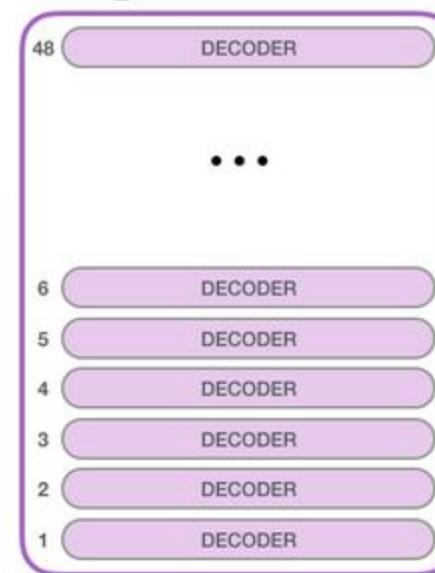
GPT-2
LARGE



Model Dimensionality: 1280



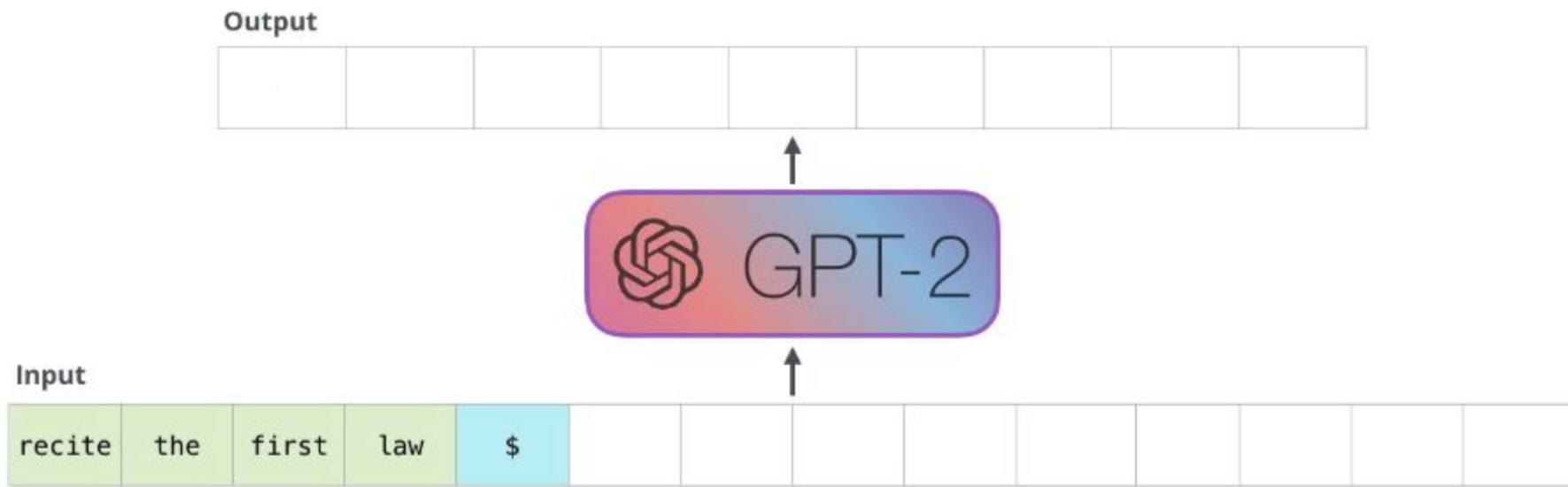
GPT-2
EXTRA
LARGE



Model Dimensionality: 1600

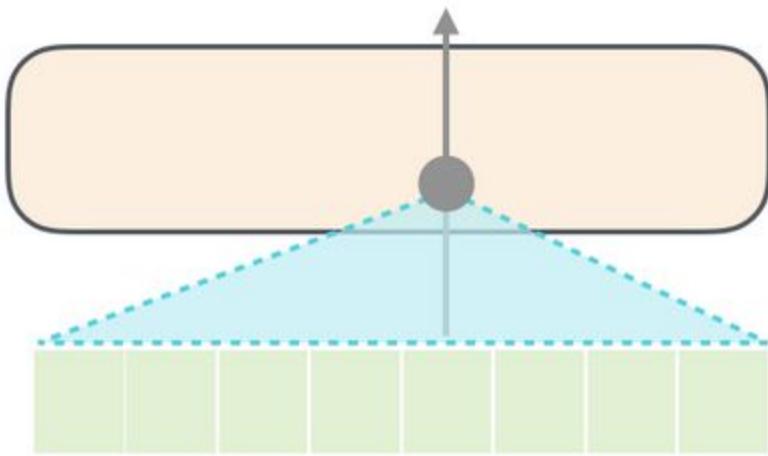


Michigan Tech

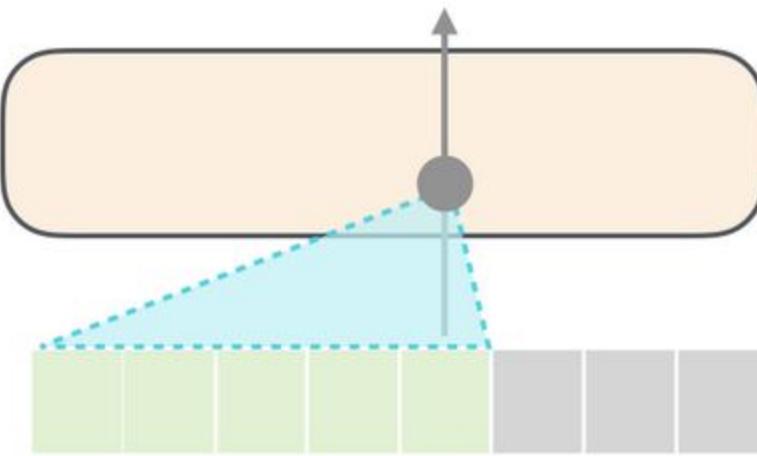


Michigan Tech

Self-Attention

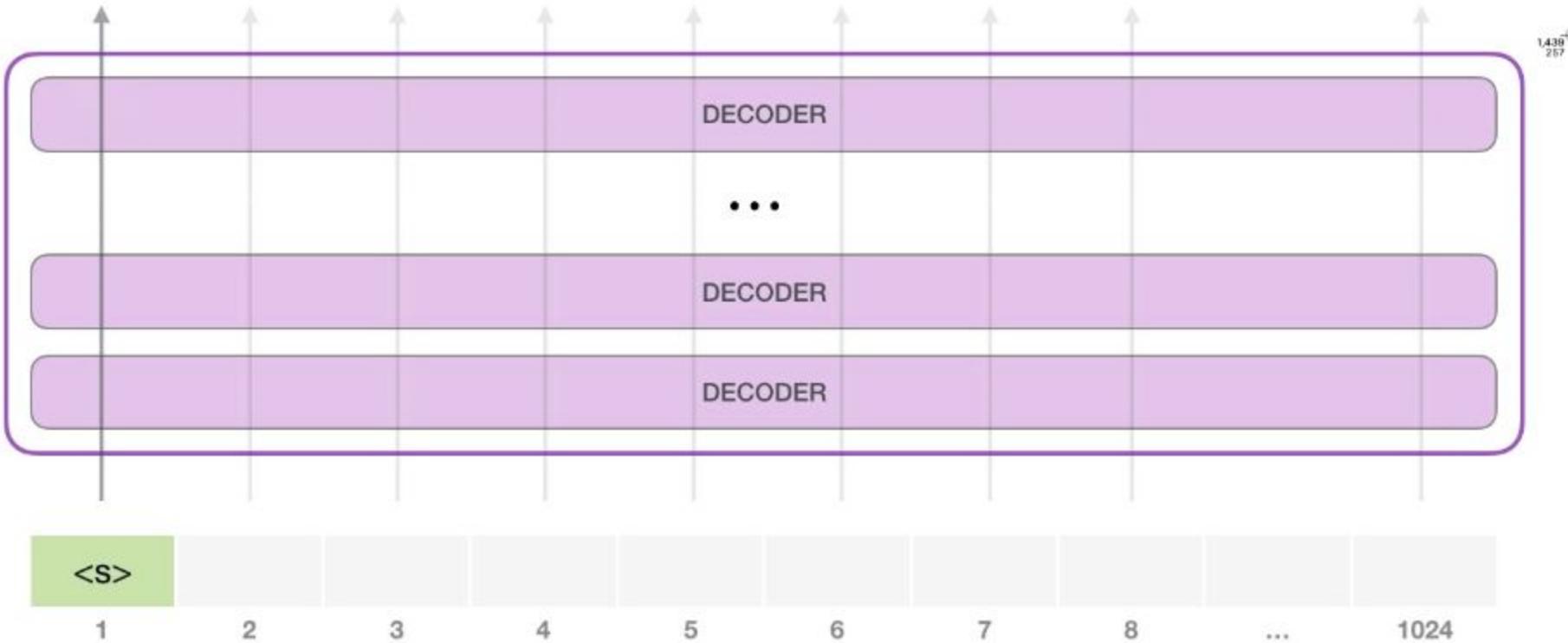


Masked Self-Attention

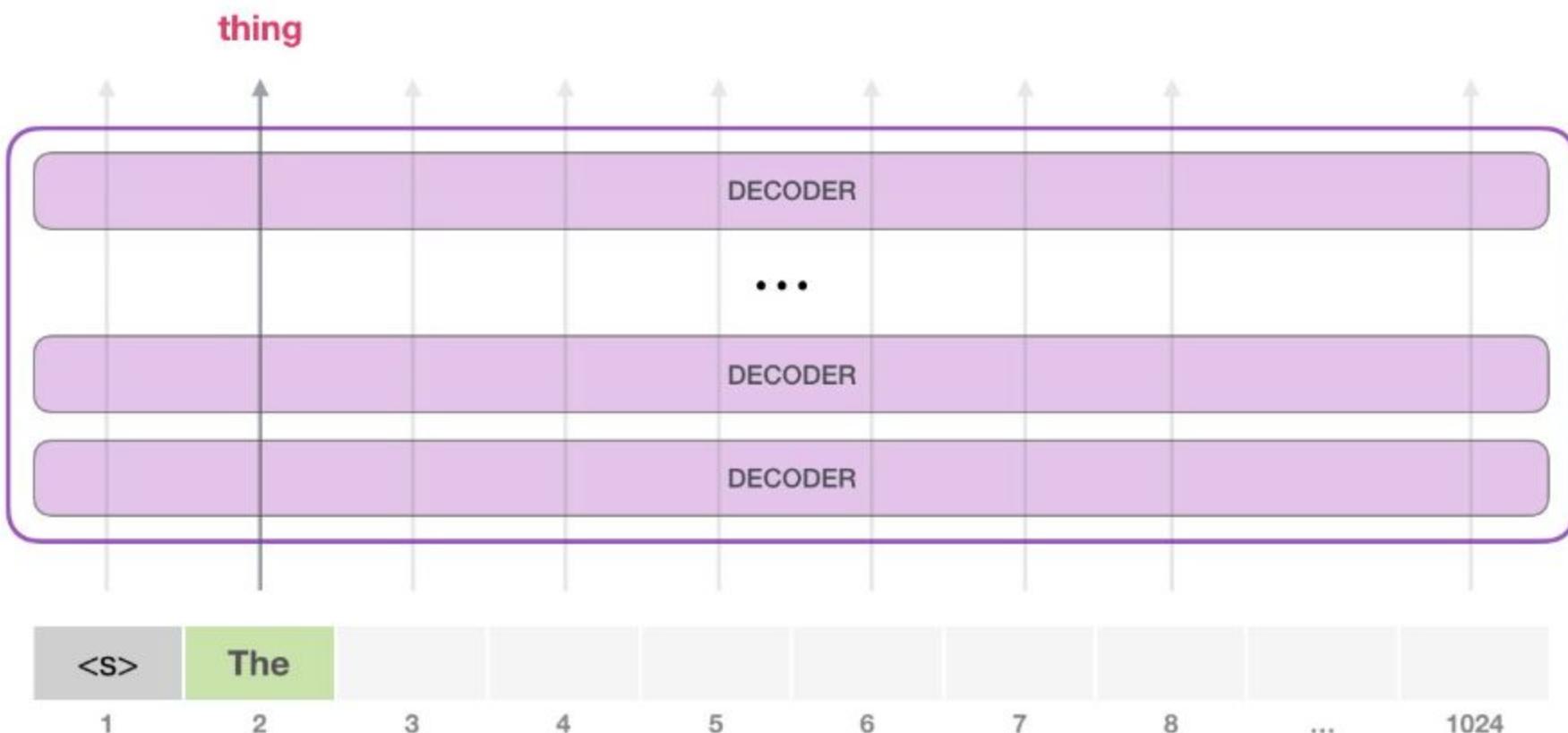


Michigan Tech

The



Michigan Tech

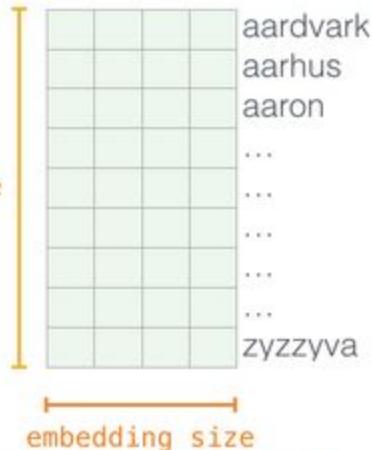


Michigan Tech

as part of a unique model.

Token Embeddings (wte)

model vocabulary size
50,257



768 (small) / **1024** (medium) / **1280** (large) / **1600** (extra large)

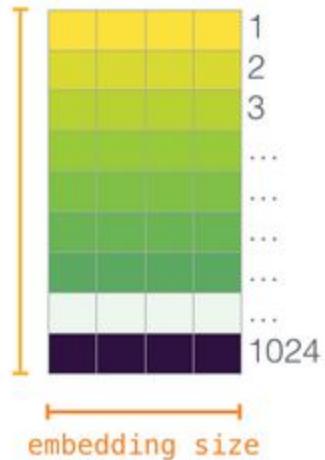
Each row is a word embedding: a list of numbers representing a word and capturing some of its meaning. The size of that list is different in different GPT2 model sizes. The smallest model uses an embedding size of 768 per word/token.



Michigan Tech

Positional Encodings (wpe)

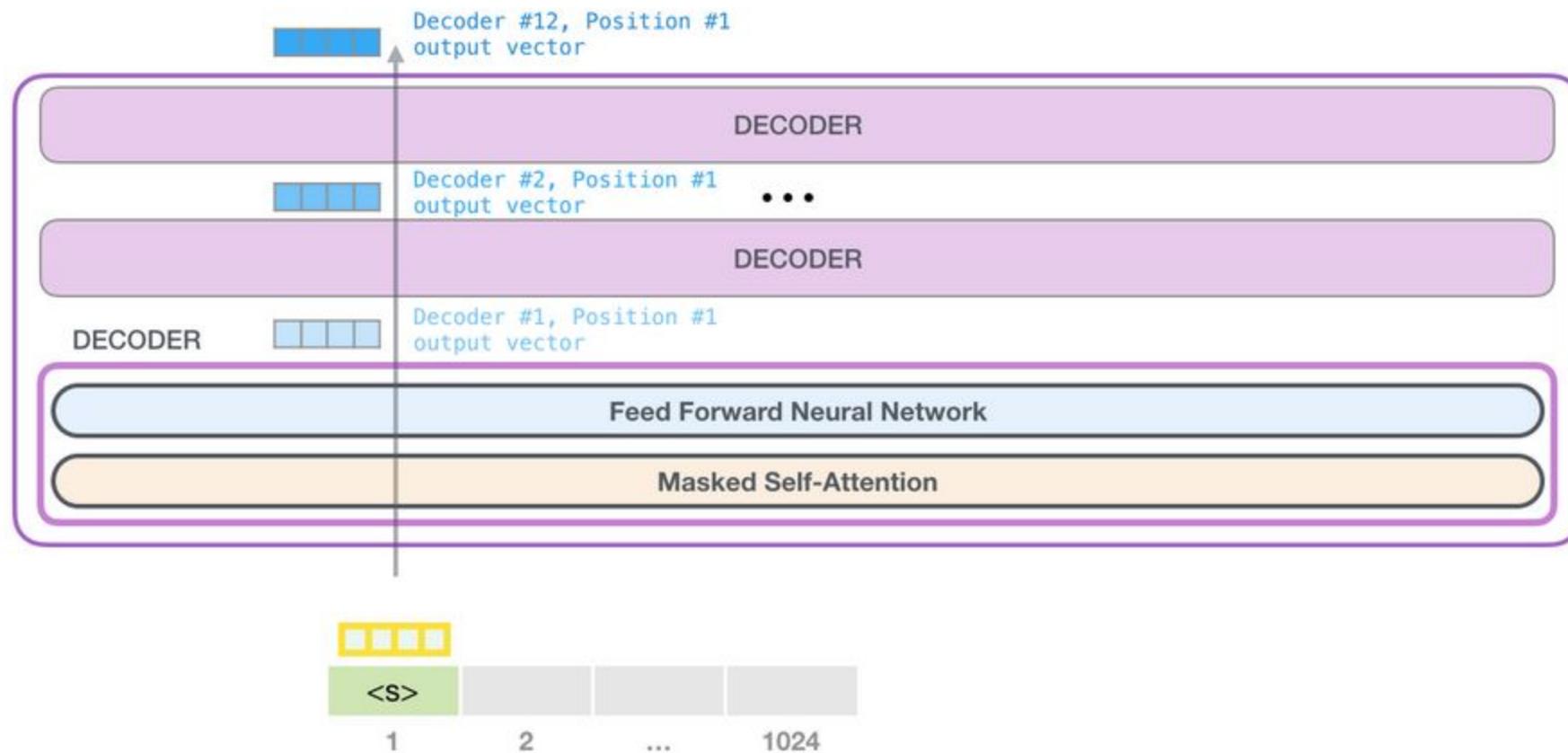
Context size
1024



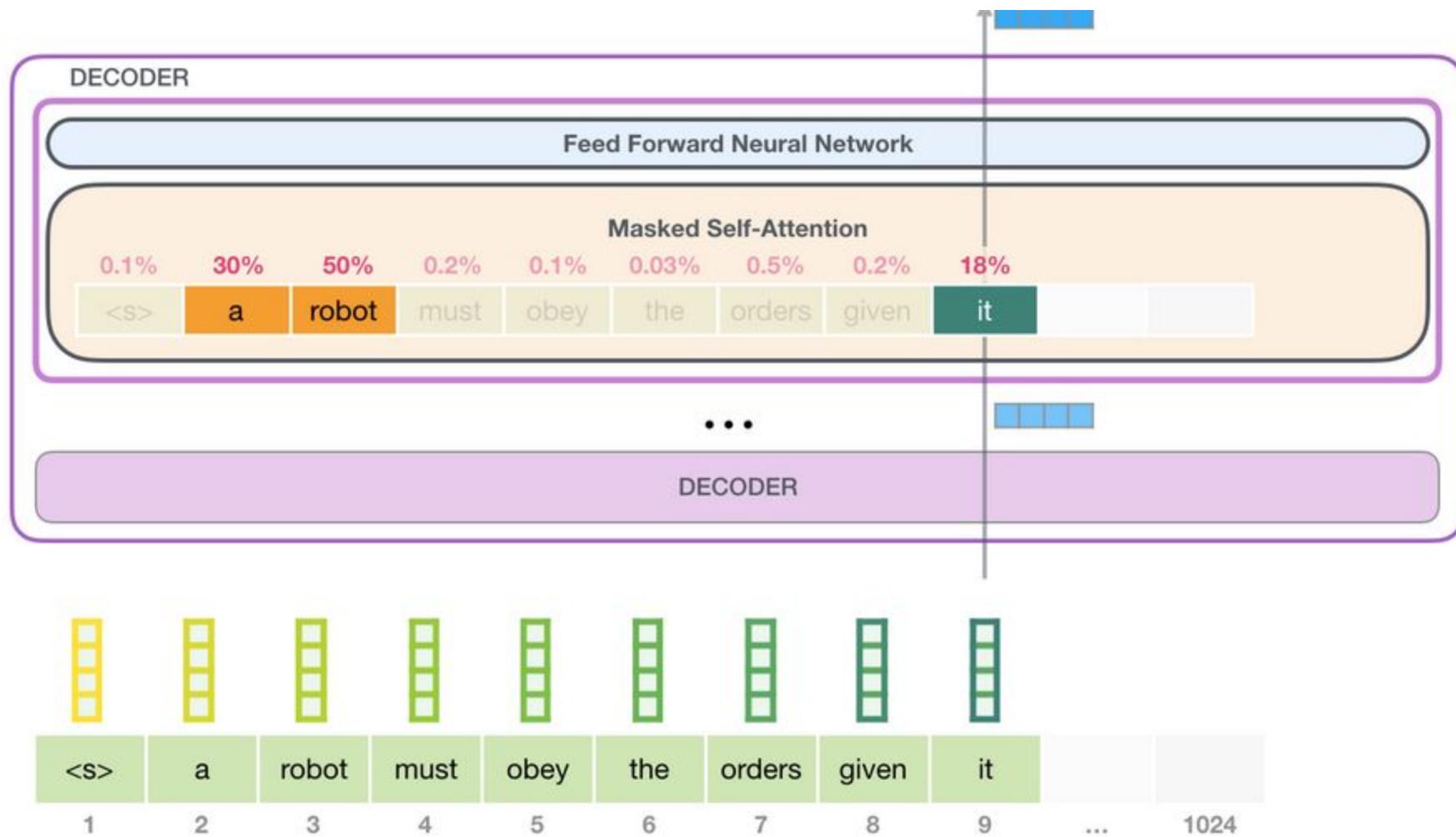
Michigan Tech



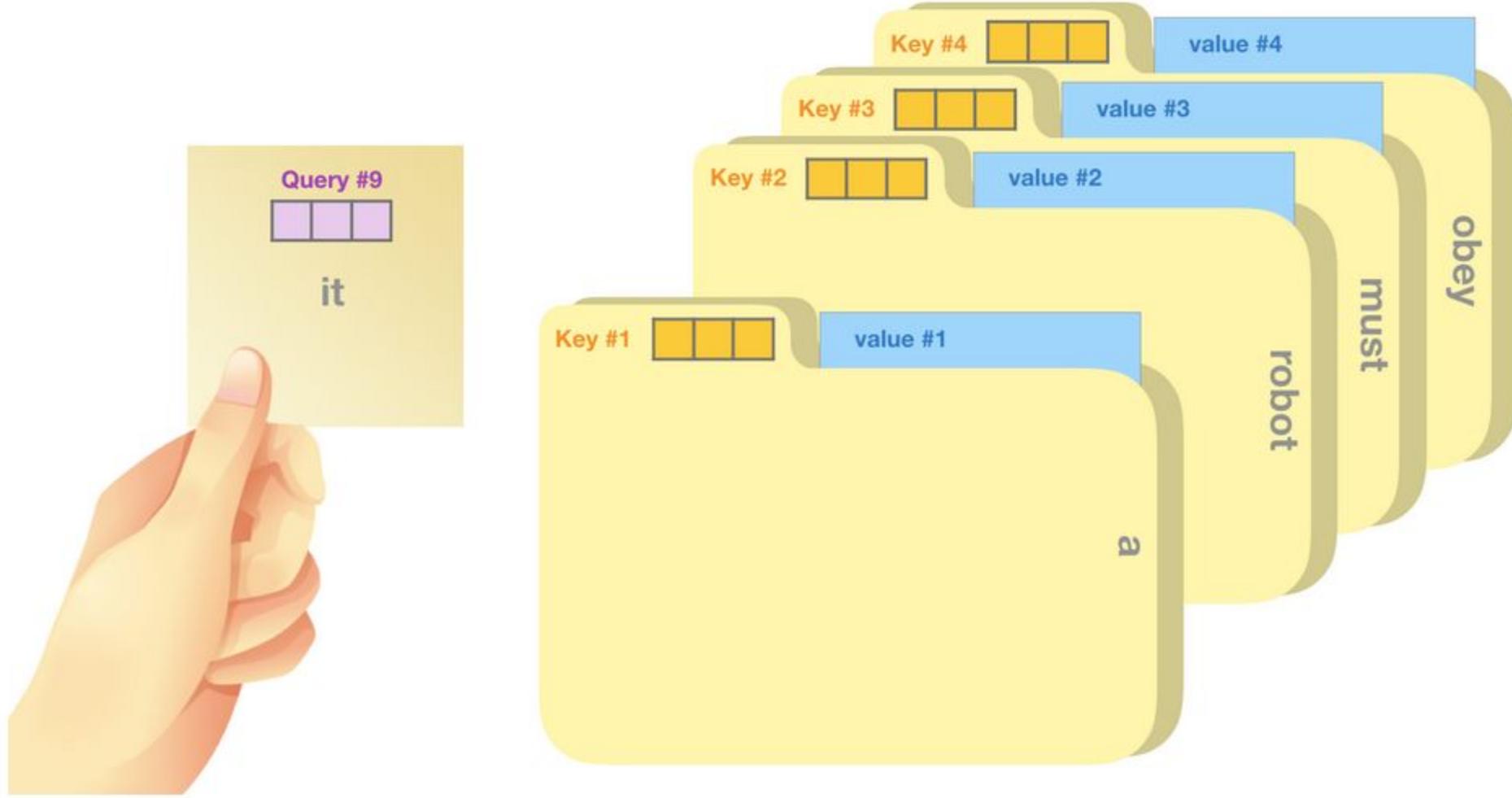
Michigan Tech



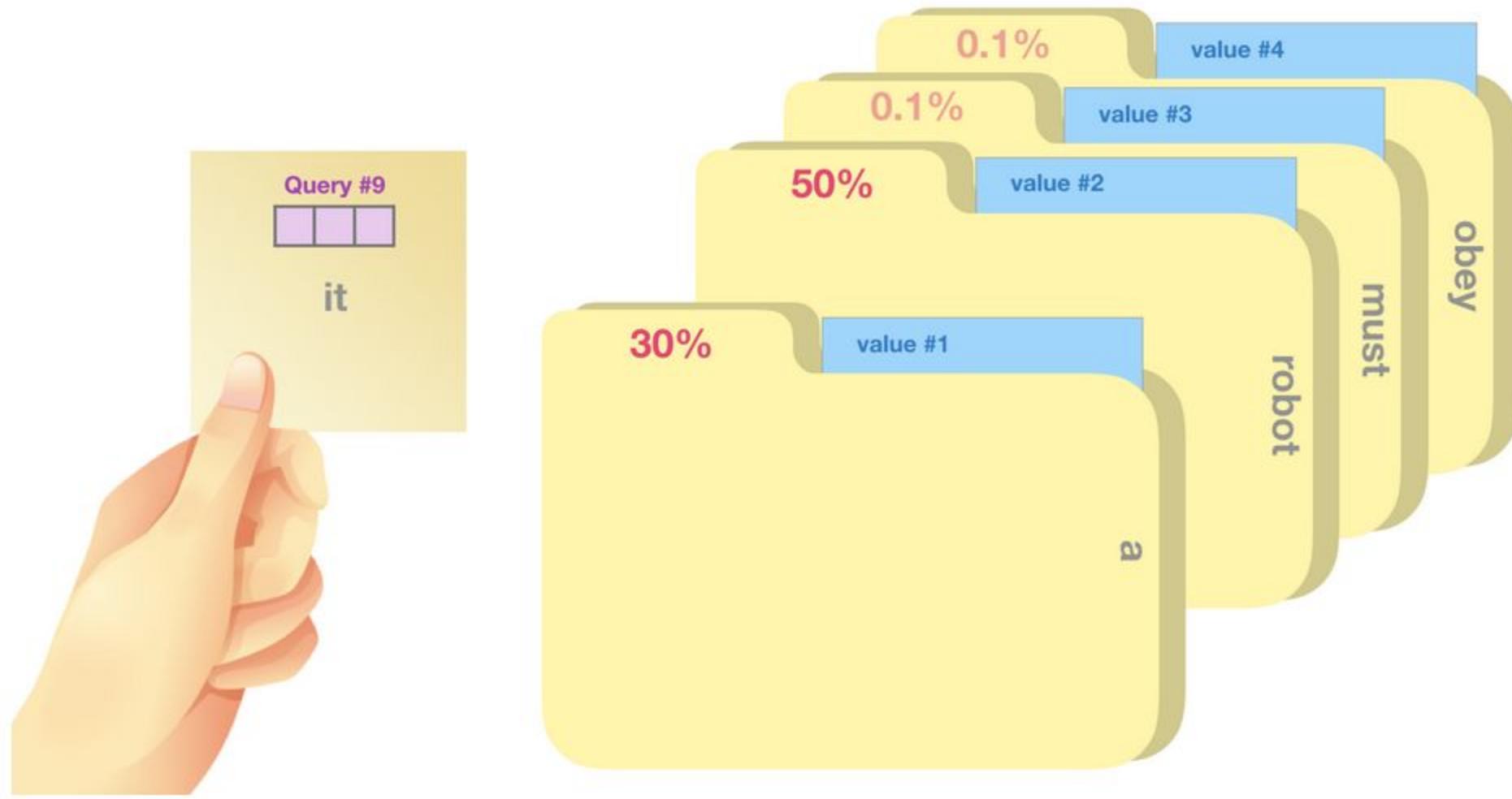
Michigan Tech



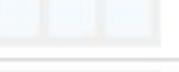
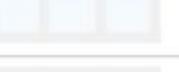
Michigan Tech



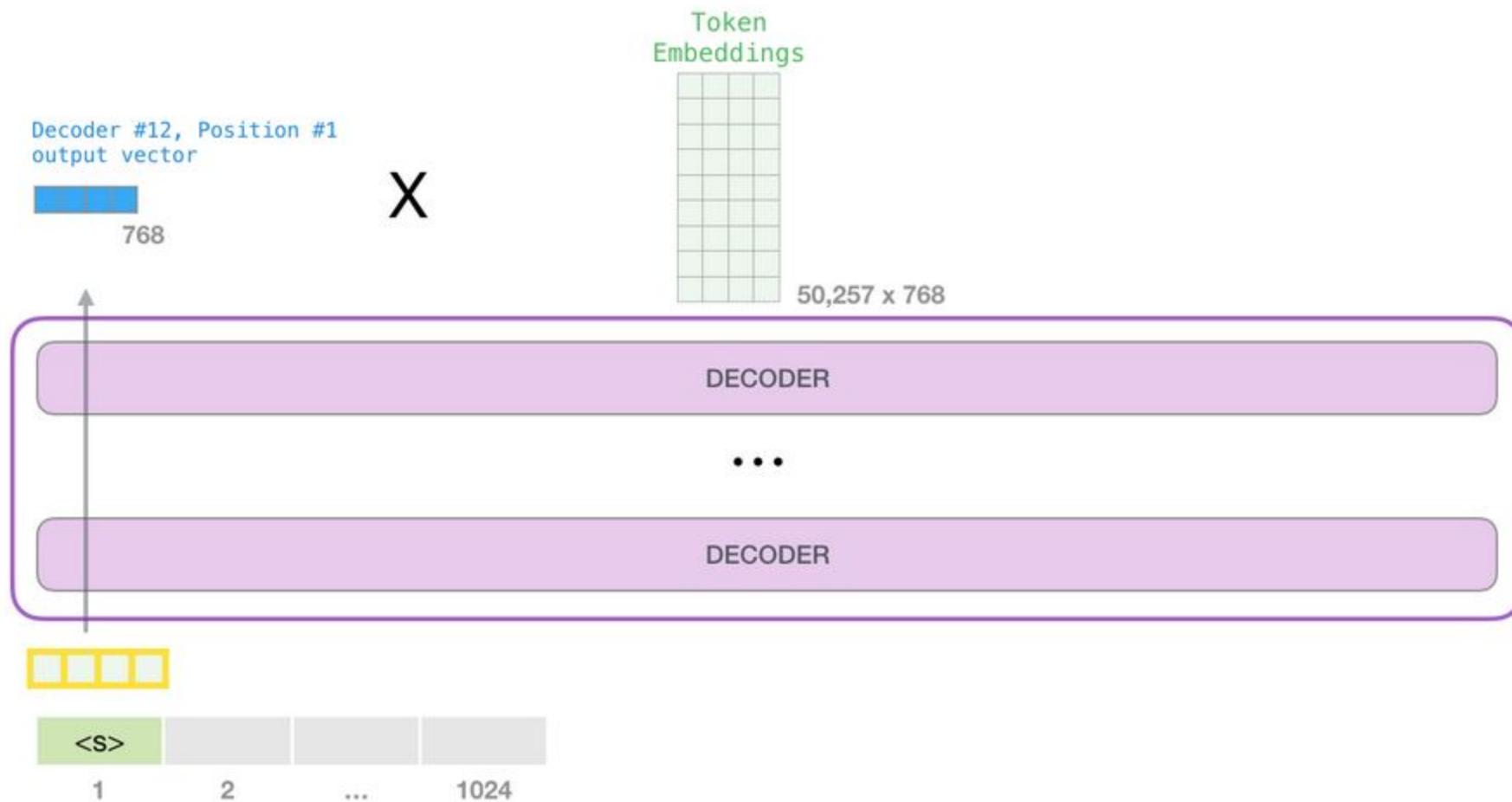
Michigan Tech



Michigan Tech

Word	Value vector	Score	Value X Score
<S>		0.001	
a		0.3	
robot		0.5	
must		0.002	
obey		0.001	
the		0.0003	
orders		0.005	
given		0.002	
it		0.19	
Sum:			





Michigan Tech

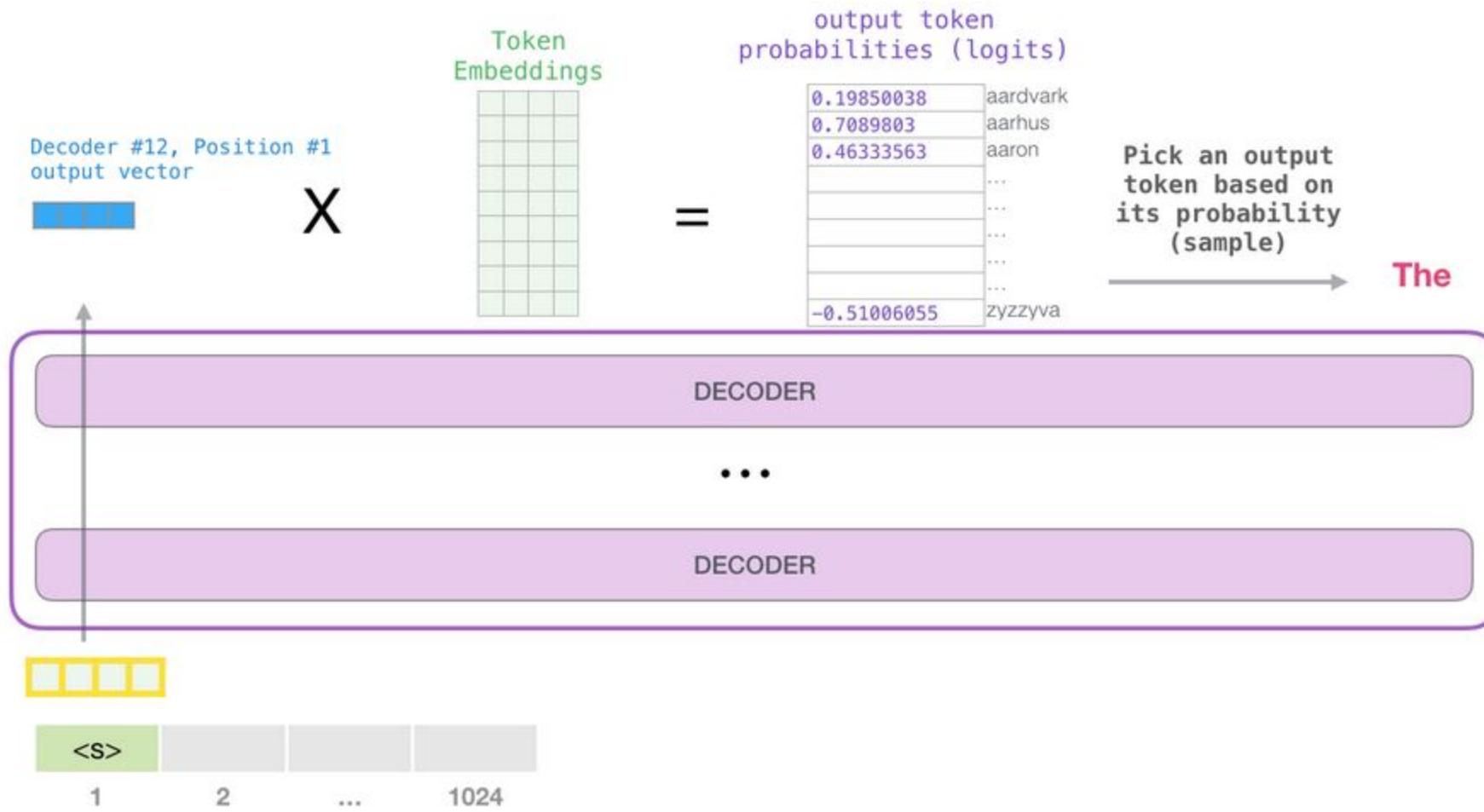
output token probabilities (logits)

model vocabulary size
50,257

0.19850038	aardvark
0.7089803	aarhus
0.46333563	aaron
...	...
...	...
...	...
...	...
...	...
-0.51006055	zyzzyva



Michigan Tech



Michigan Tech

Let's talk about sampling



Michigan Tech

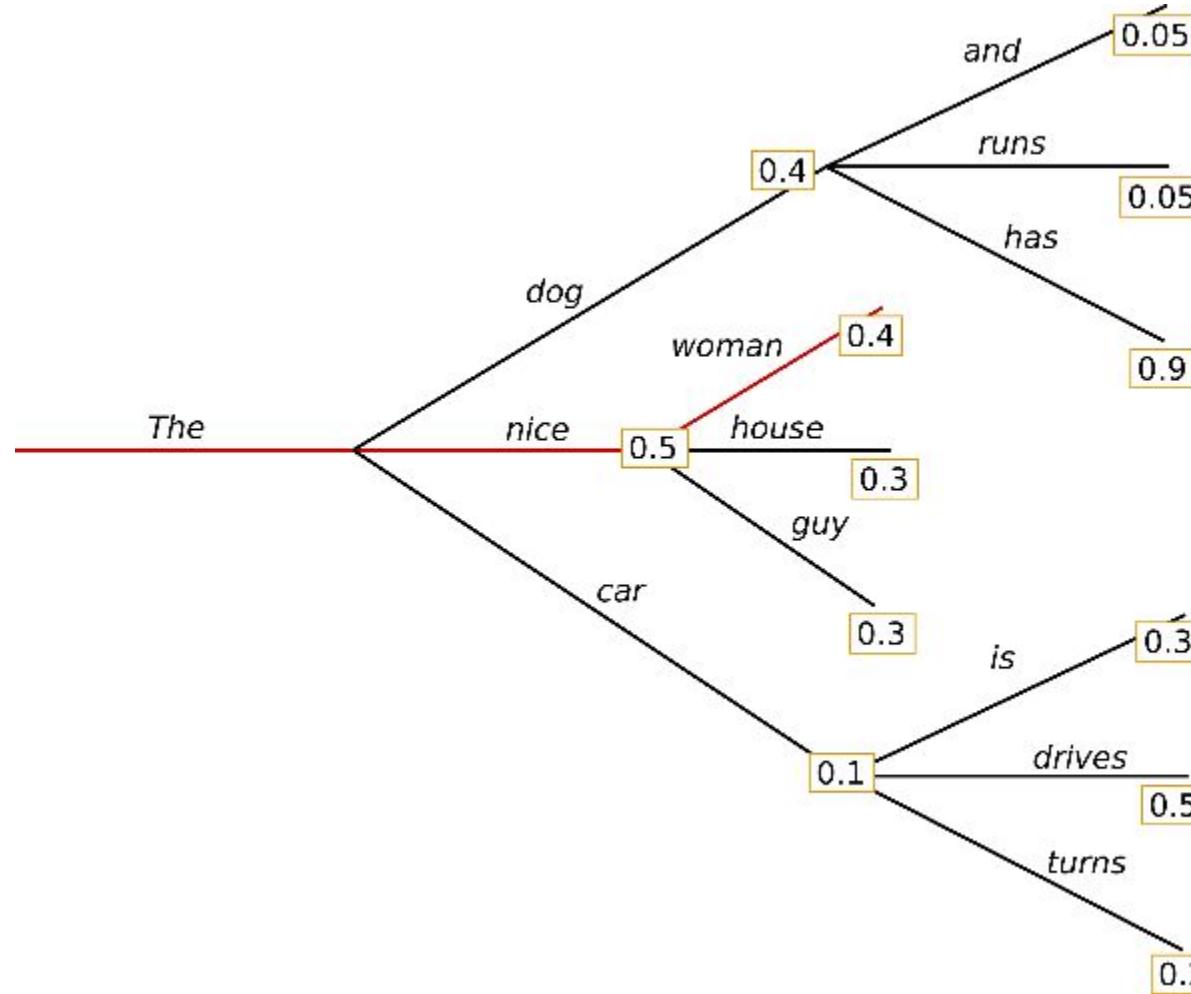
Conditional token probabilities

$$P(w_{1:T}|W_0) = \prod_{t=1}^T P(w_t|w_{1:t-1}, W_0), \text{with } w_{1:0} = \emptyset,$$



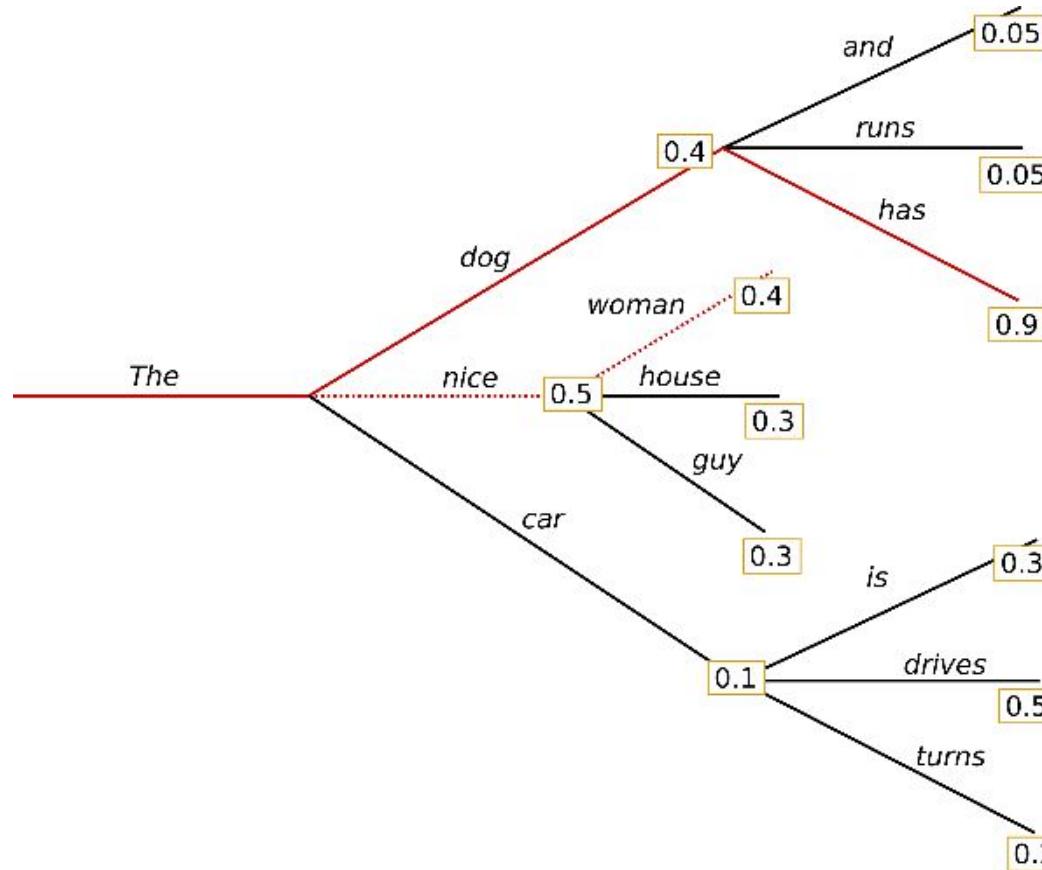
Michigan Tech

The easiest method - greedy decoding



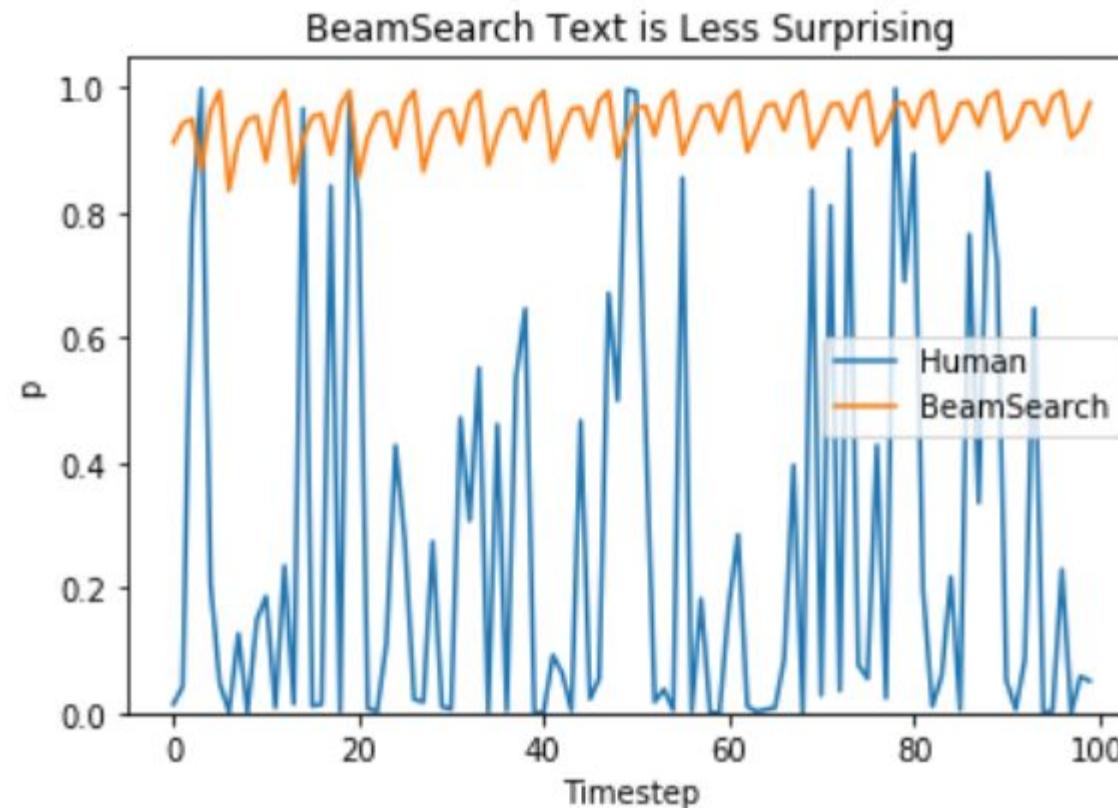
Michigan Tech

More advanced deterministic approach - beam search

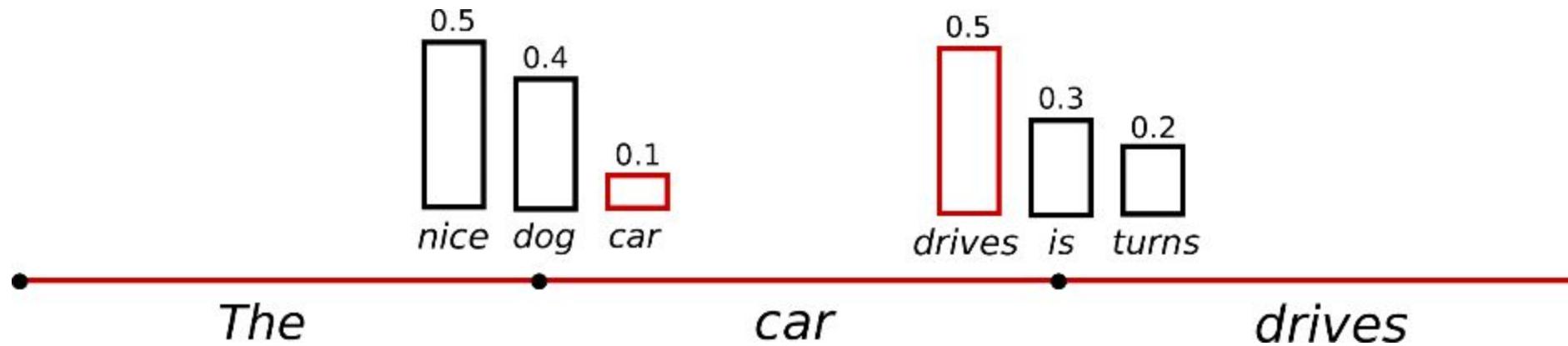


Michigan Tech

Is this good?



Sampling based approaches



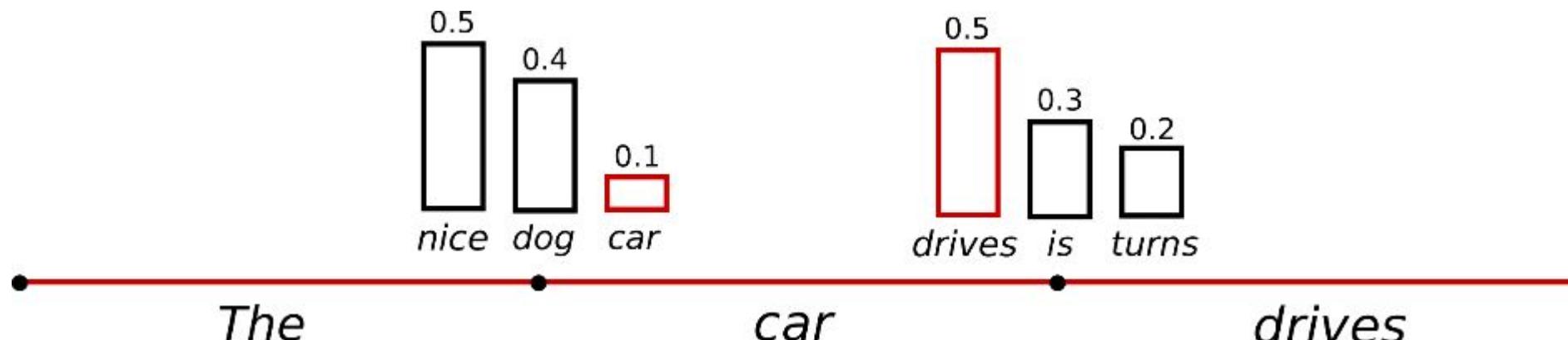
Michigan Tech

Too much randomness is bad

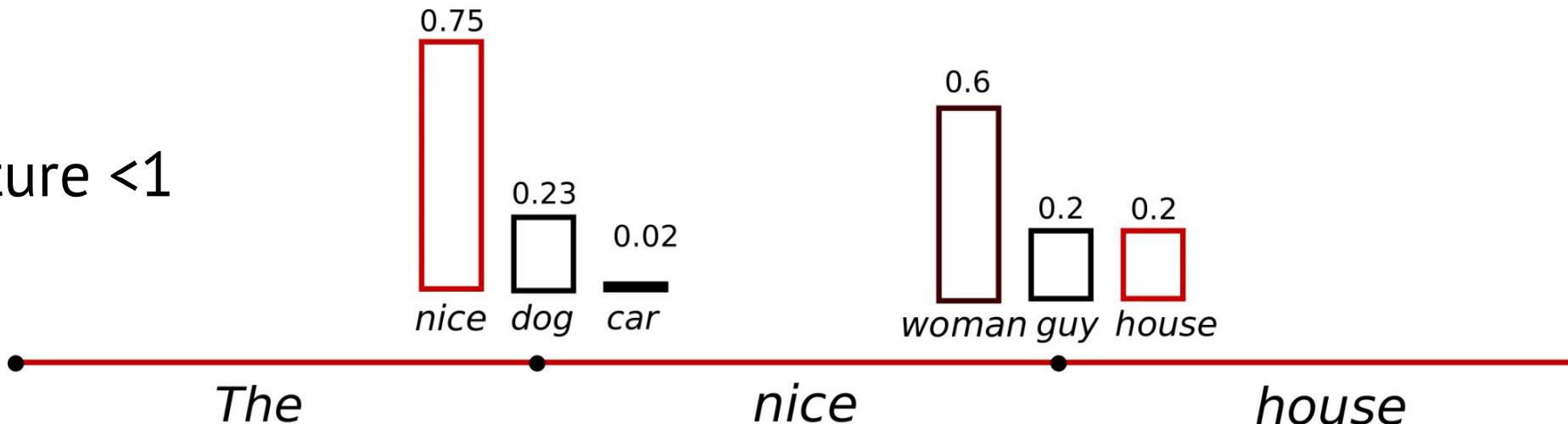


Michigan Tech

“Cooling down” distribution

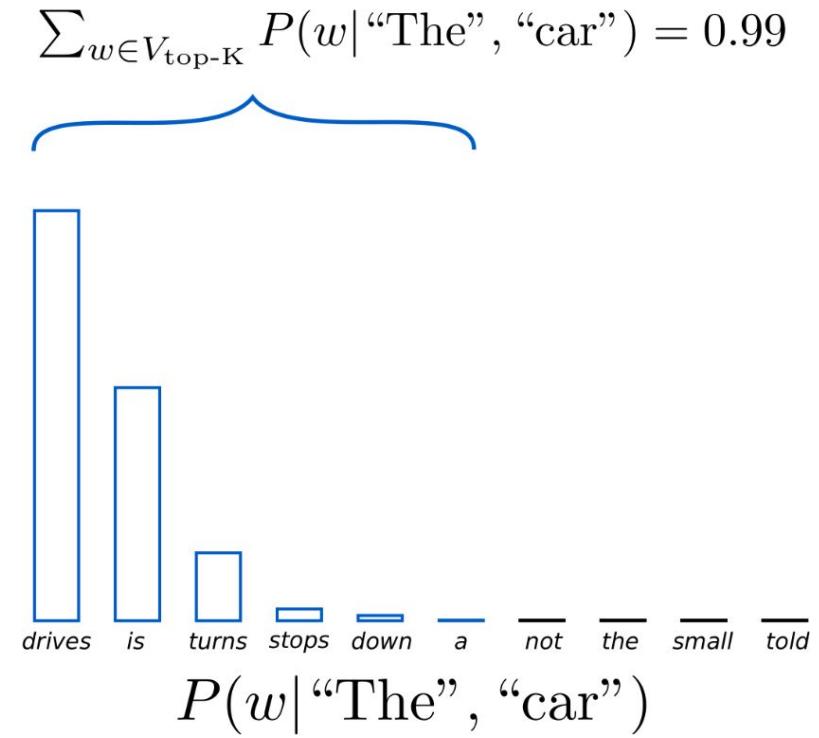
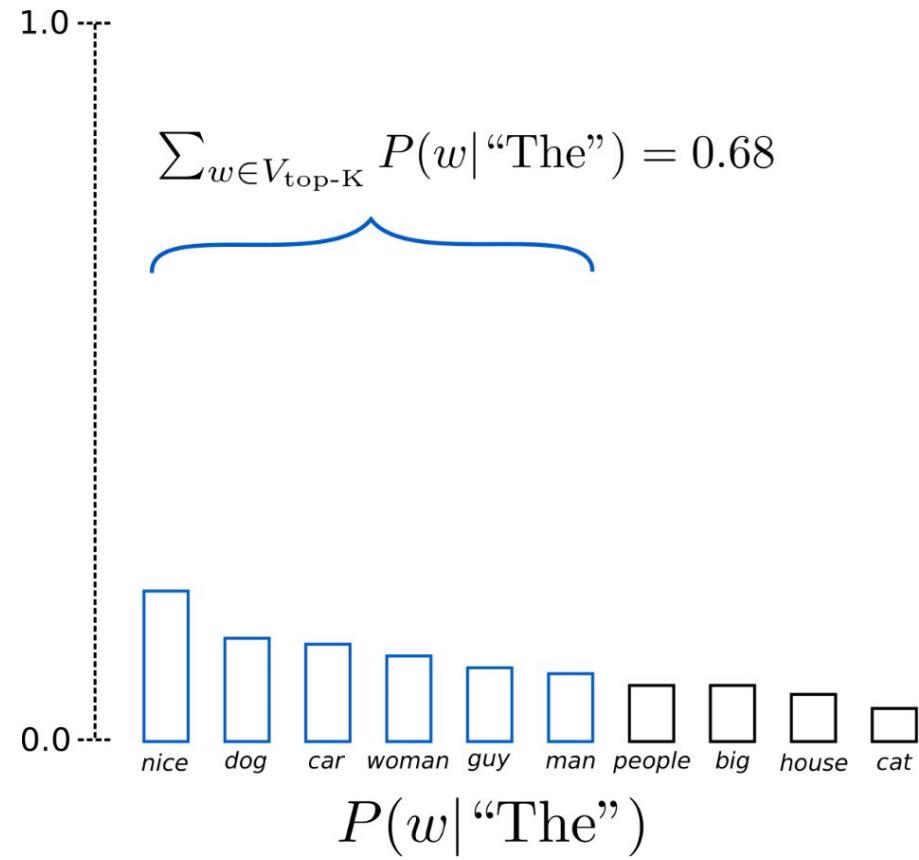


With
temperature < 1



Michigan Tech

Top-k sampling



Michigan Tech
1885

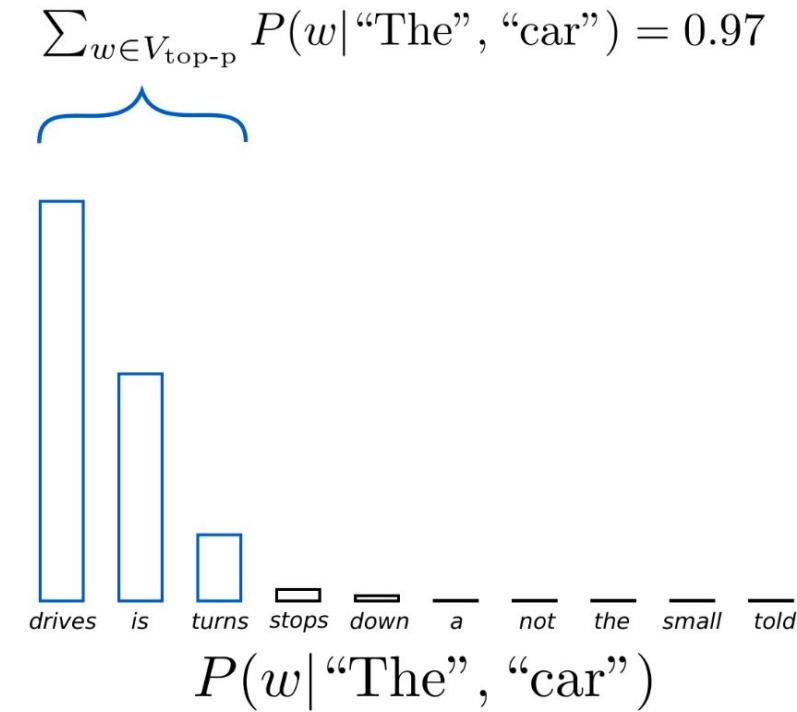
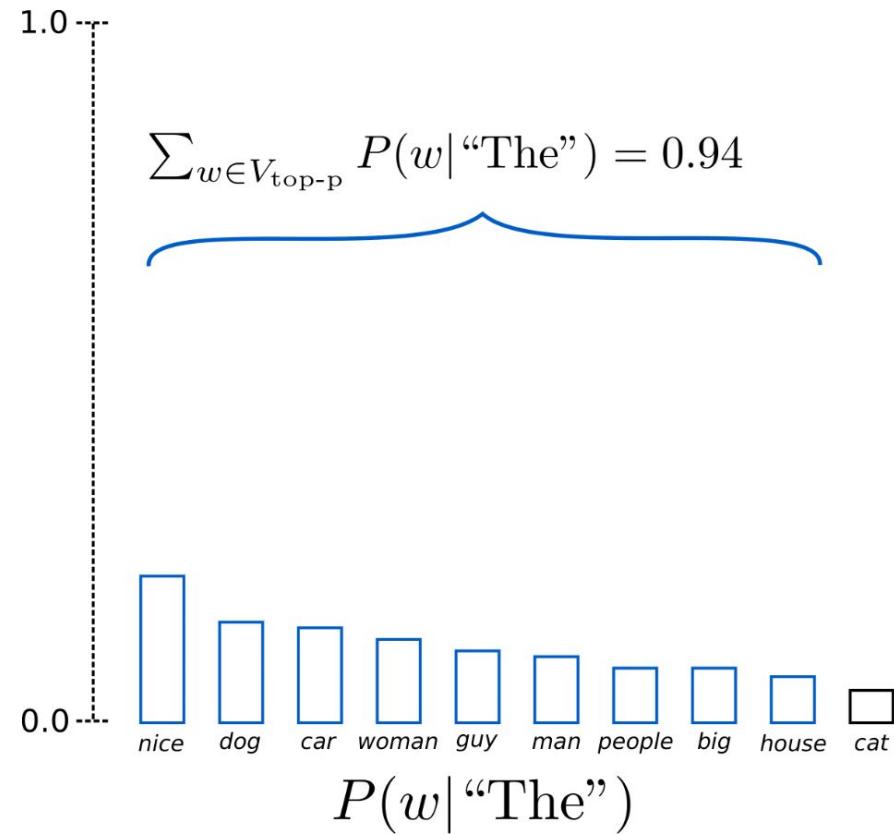
Top-k

- Better than temperature
- More human sounding than beam-search
- Doesn't dynamically adapt to different token prediction distributions



Michigan Tech

Top-p (nucleus) sampling



Top-p

- Generally works pretty well
- Can be used in conjunction with temperature and top-k



Michigan Tech

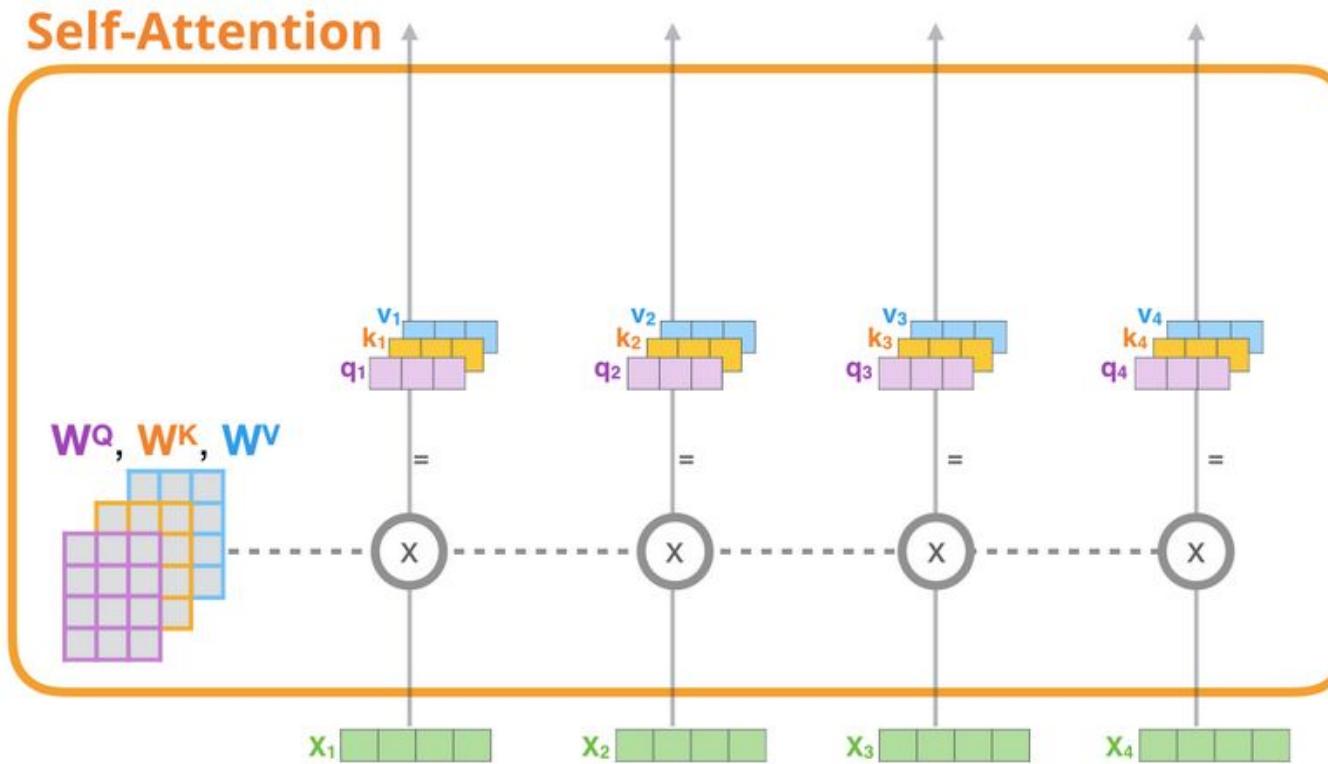
Back to self-attention



Michigan Tech

Self attention (without masking)

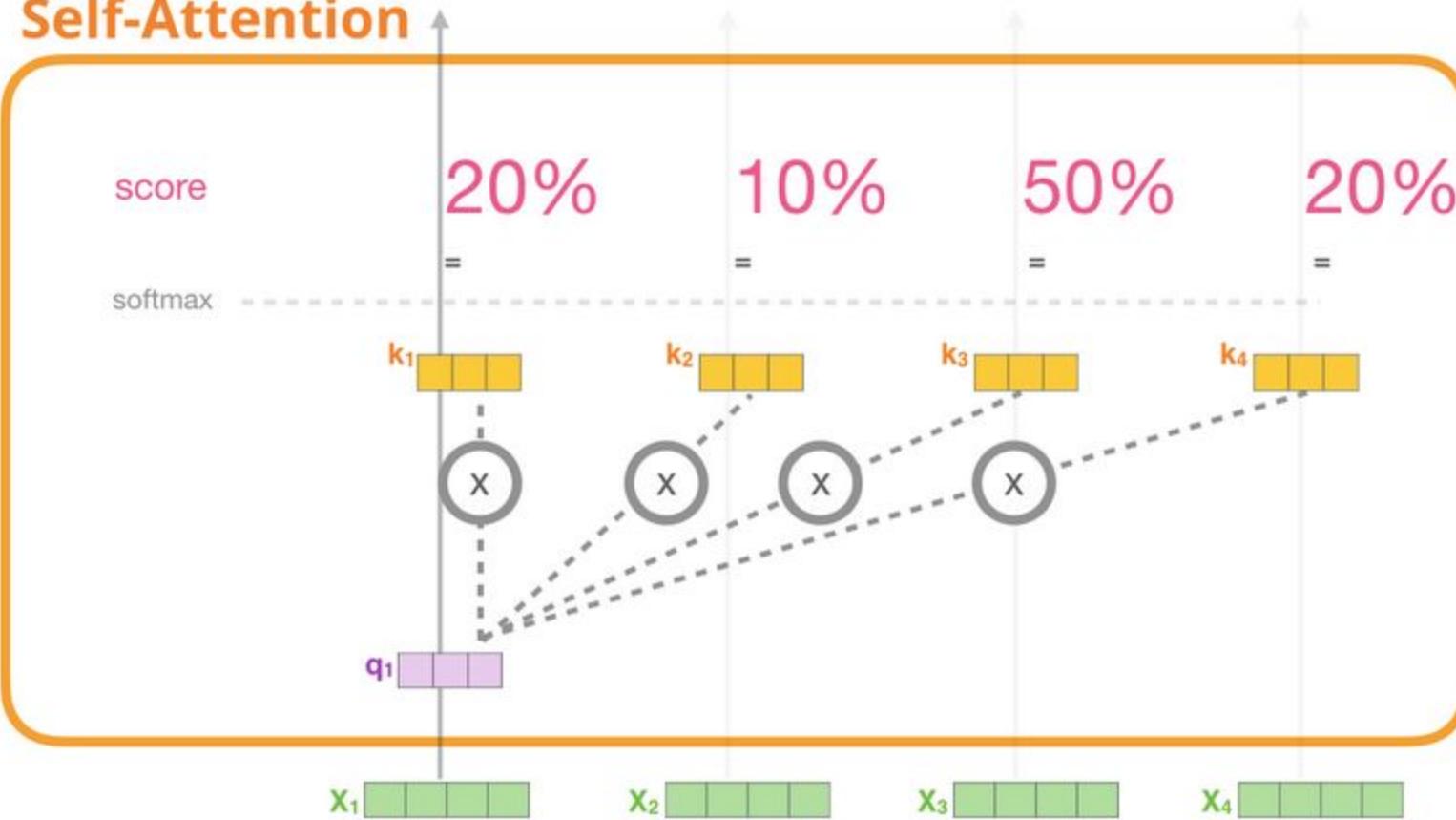
1) For each input token, create a **query vector**, a **key vector**, and a **value vector** by multiplying by weight Matrices \mathbf{W}^Q , \mathbf{W}^K , \mathbf{W}^V



Michigan Tech
1885

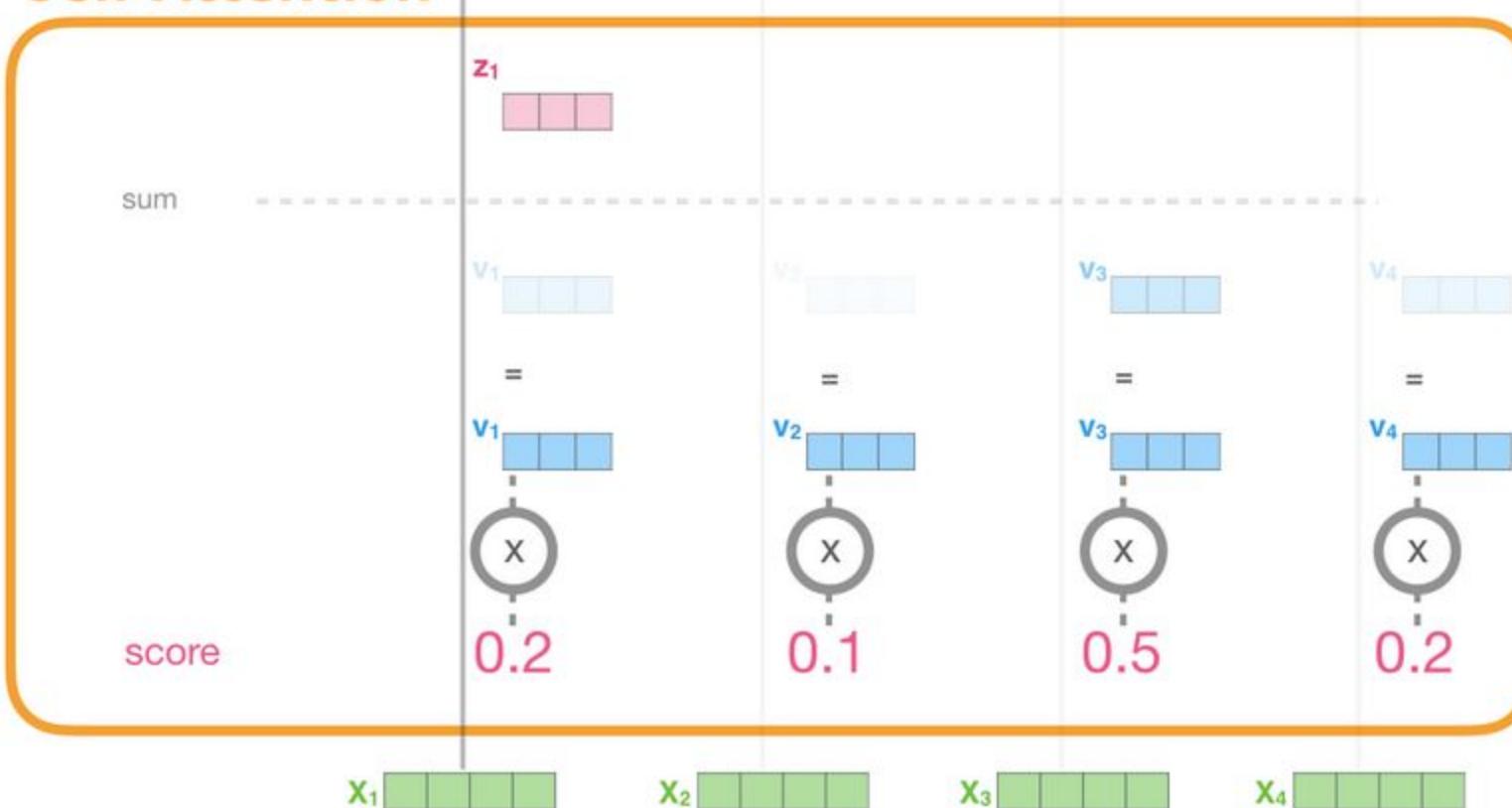
2) Multiply (dot product) the current **query vector**, by all the **key vectors**, to get a score of how well they match

Self-Attention



3) Multiply the value vectors by the scores, then sum up

Self-Attention

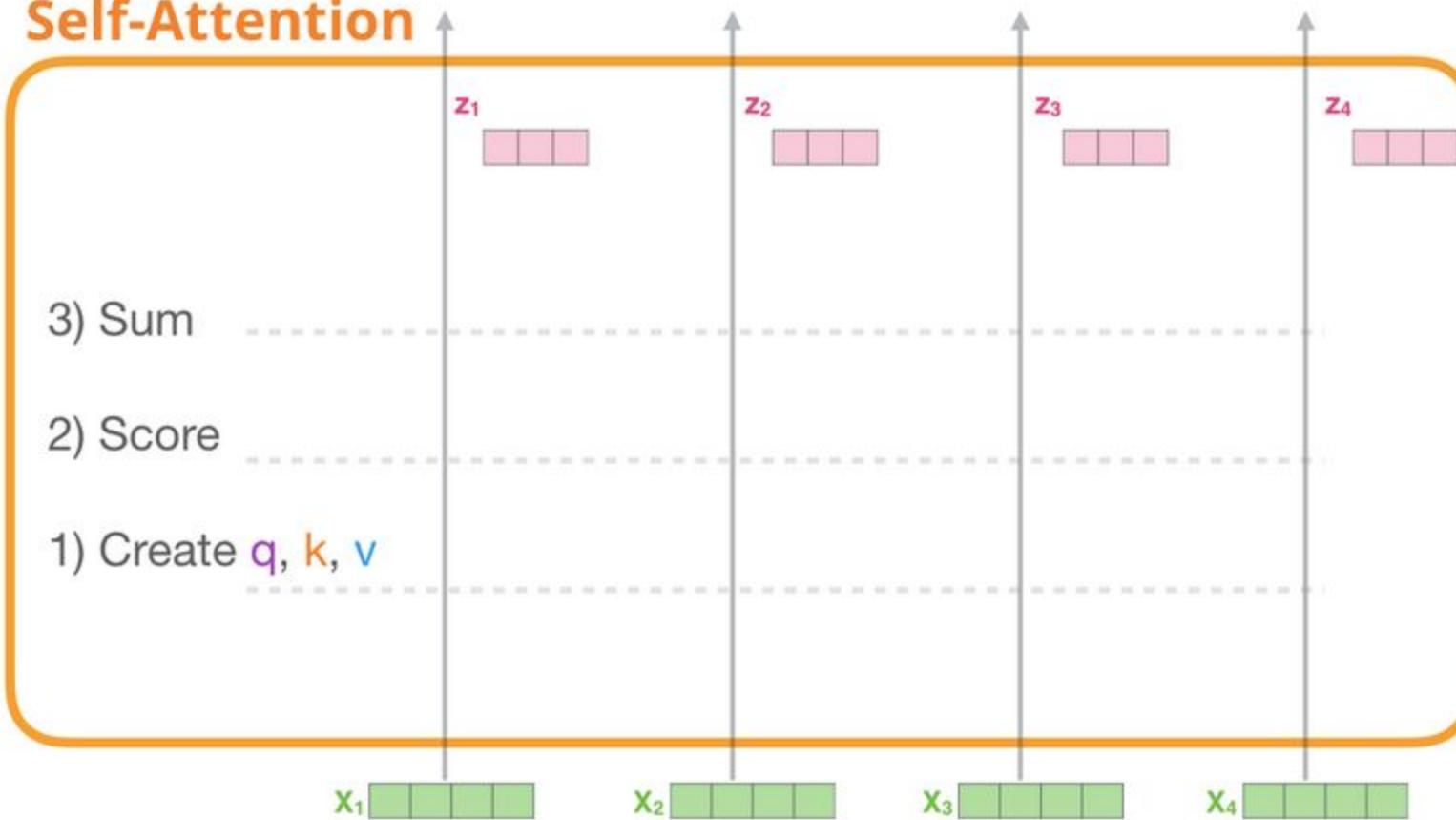


The lower the score, the more transparent we're showing the value vector. That's to indicate how multiplying by a small number dilutes the values of the vector:



Michigan Tech

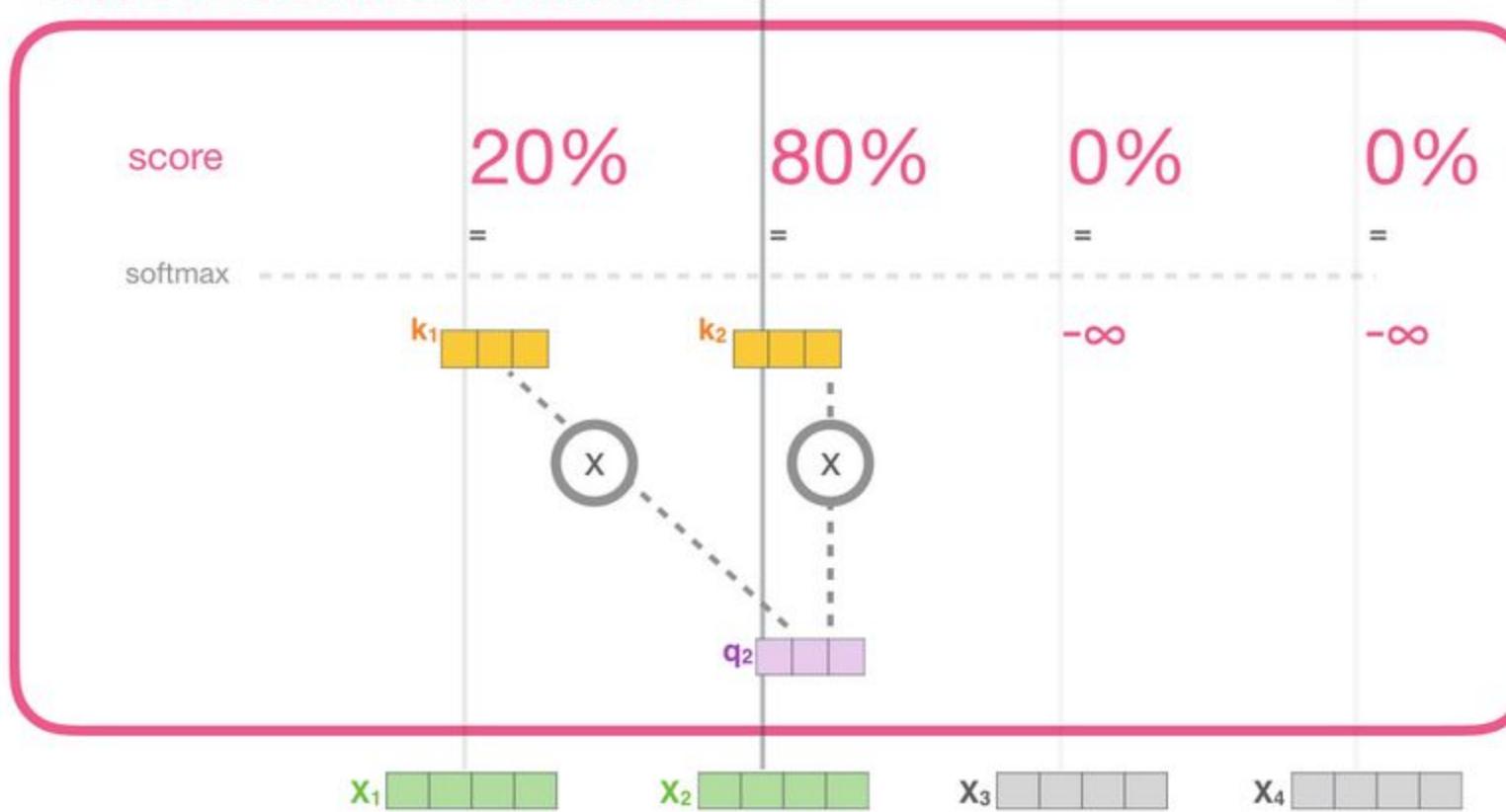
Self-Attention



Michigan Tech

Masked self-attention

Masked Self-Attention



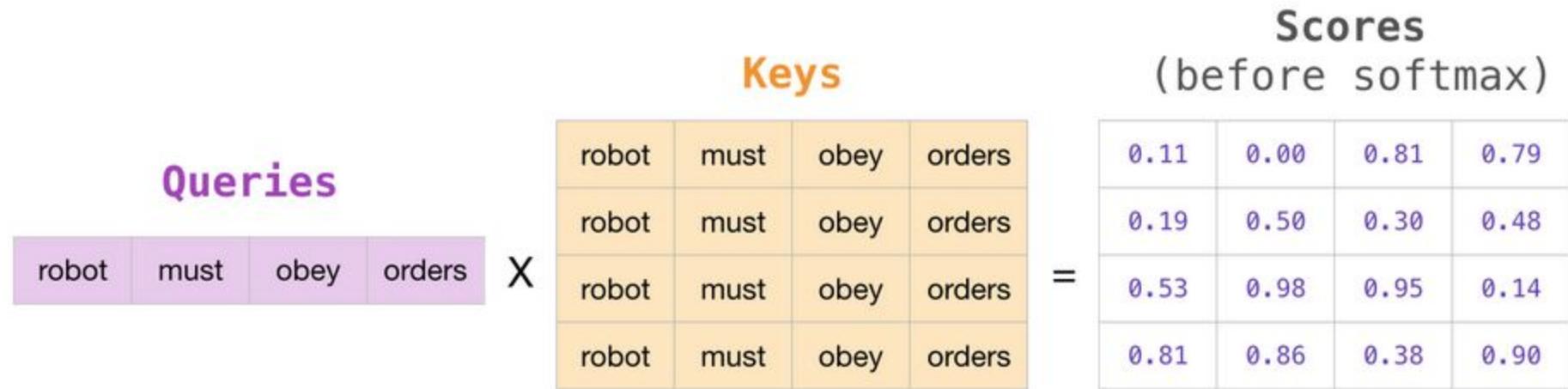
Michigan Tech

Features				Labels	
position:	1	2	3	4	
Example:	1	robot	must	obey	orders
	2	robot	must	obey	orders
	3	robot	must	obey	orders
	4	robot	must	obey	orders



Michigan Tech

Linear algebra tricks



Michigan Tech

Scores
(before softmax)

0.11	0.00	0.81	0.79
0.19	0.50	0.30	0.48
0.53	0.98	0.95	0.14
0.81	0.86	0.38	0.90

Apply Attention
Mask



Masked Scores
(before softmax)

0.11	-inf	-inf	-inf
0.19	0.50	-inf	-inf
0.53	0.98	0.95	-inf
0.81	0.86	0.38	0.90



Masked Scores
(before softmax)

0.11	-inf	-inf	-inf
0.19	0.50	-inf	-inf
0.53	0.98	0.95	-inf
0.81	0.86	0.38	0.90

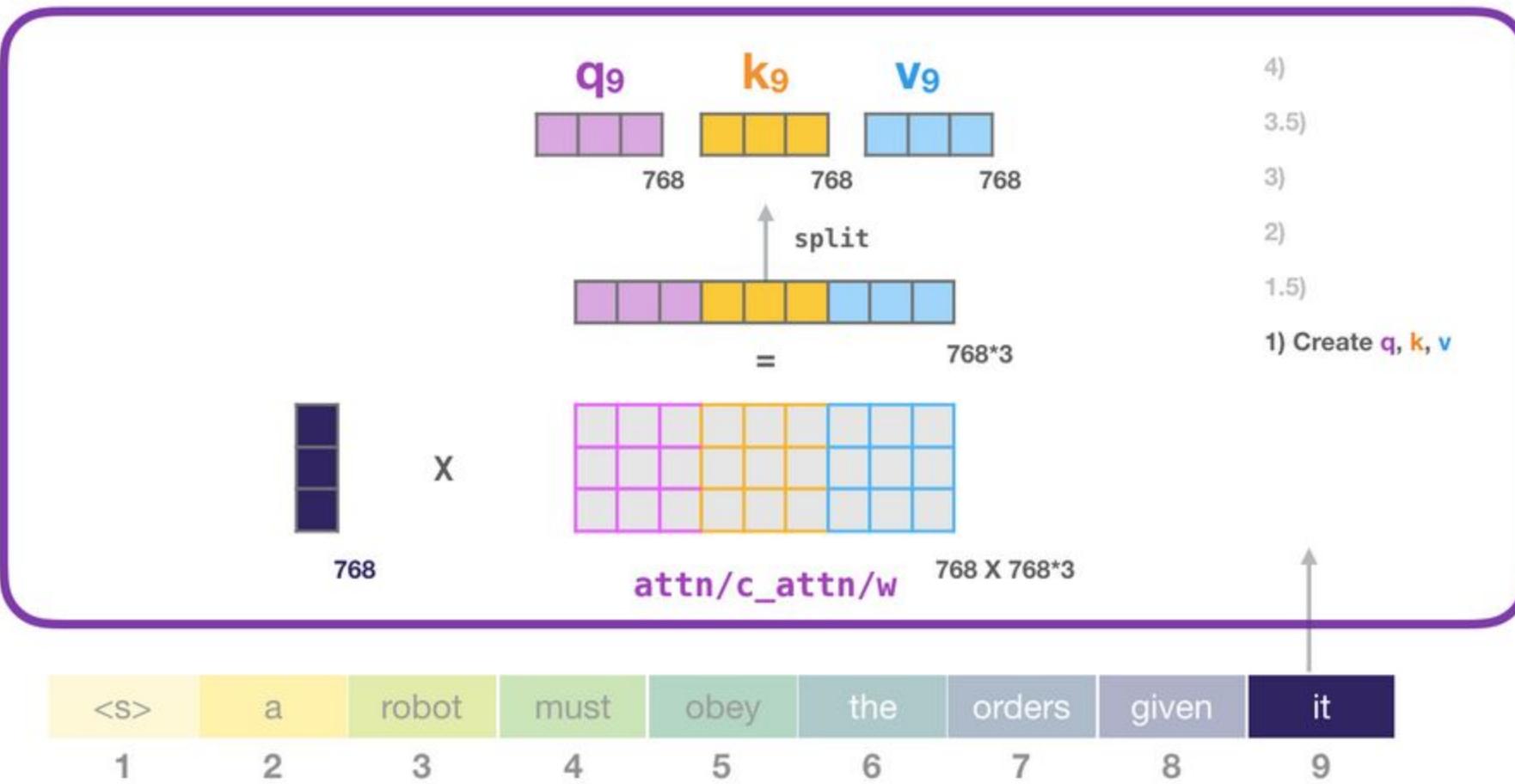
Softmax
(along rows) 

Scores

1	0	0	0
0.48	0.52	0	0
0.31	0.35	0.34	0
0.25	0.26	0.23	0.26



GPT2 Self-Attention

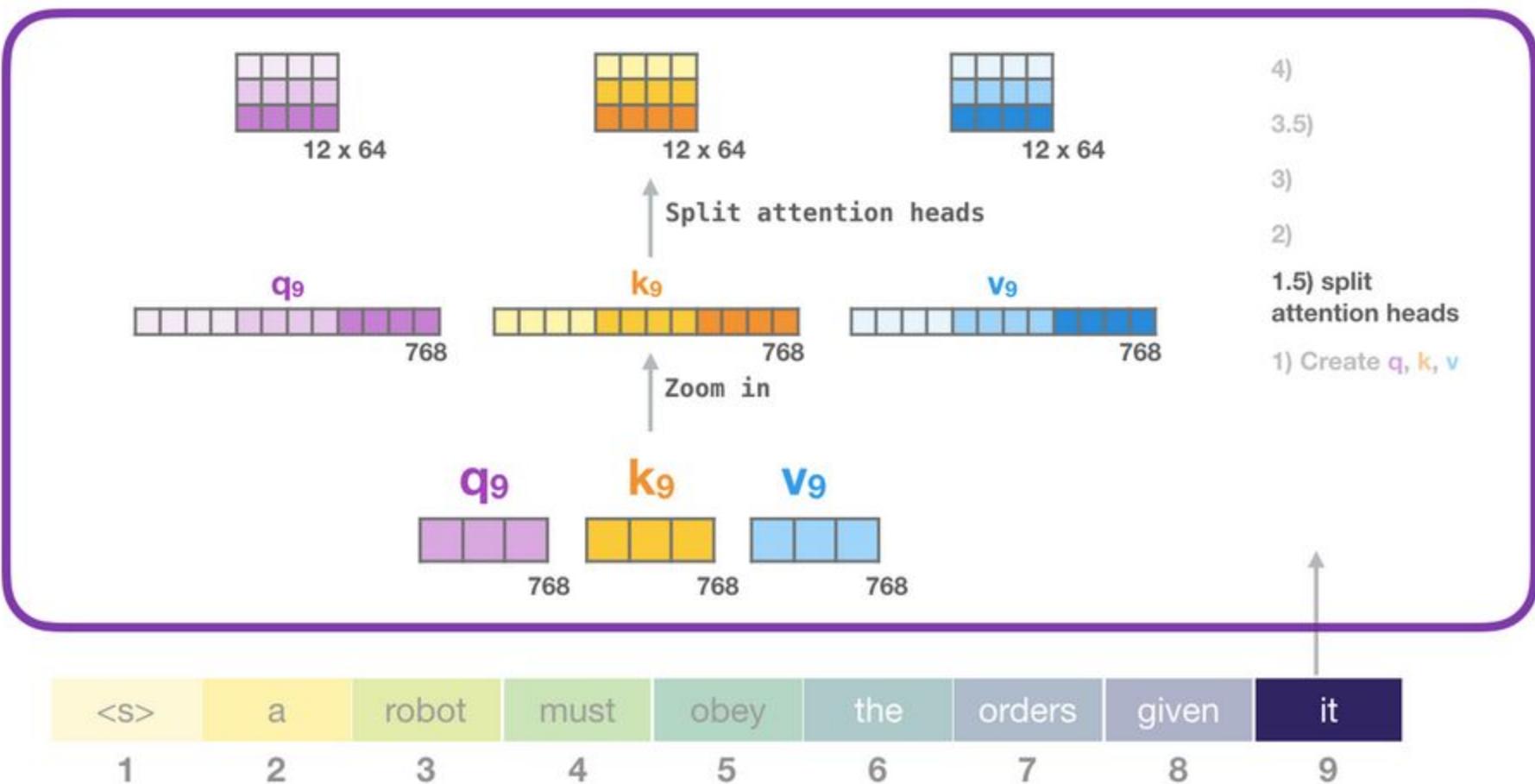


Multiplying the input vector by the attention weights vector (and adding a bias vector afterwards) results in the key, value, and query vectors for this token.



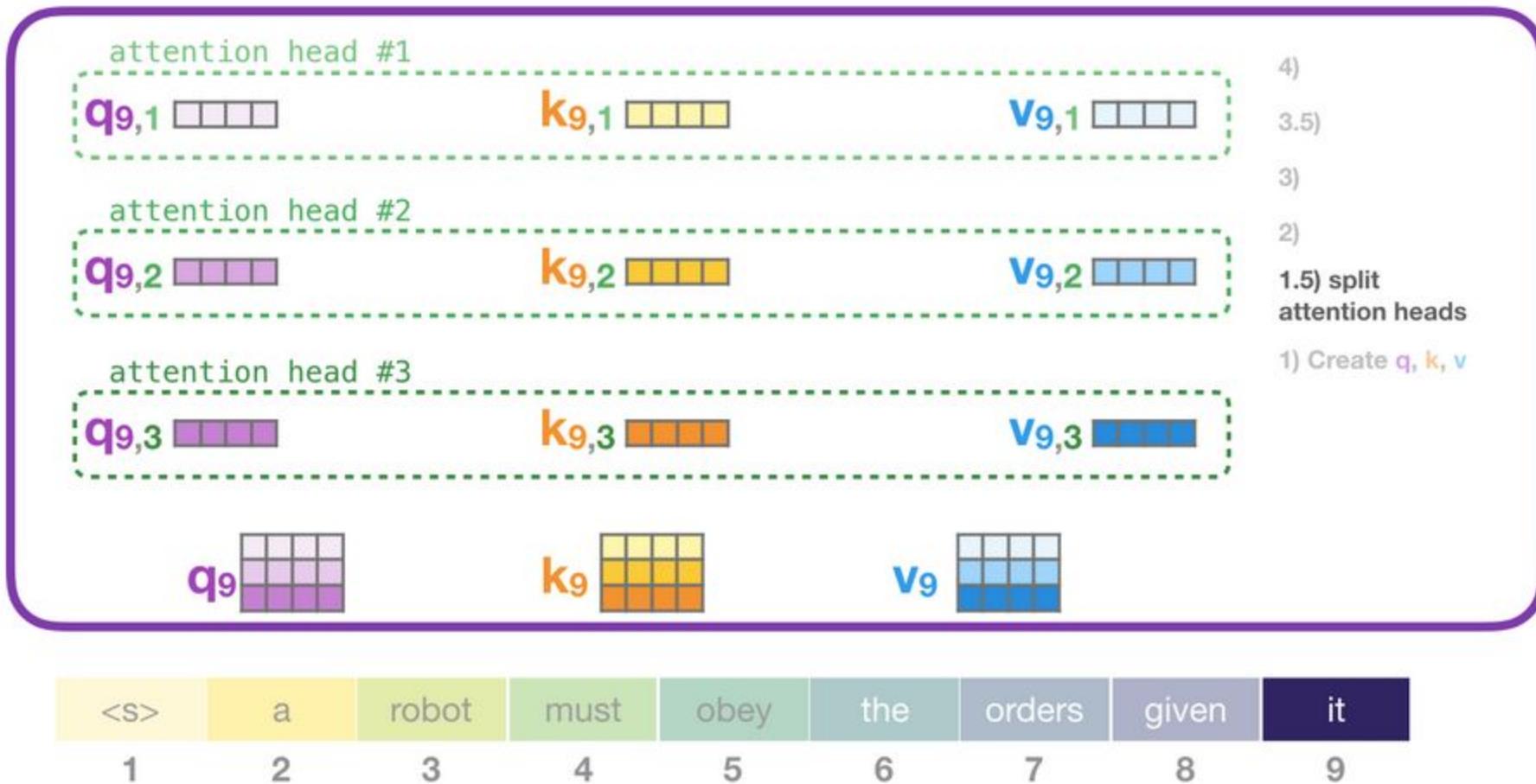
Michigan Tech

GPT2 Self-Attention



Michigan Tech

GPT2 Self-Attention



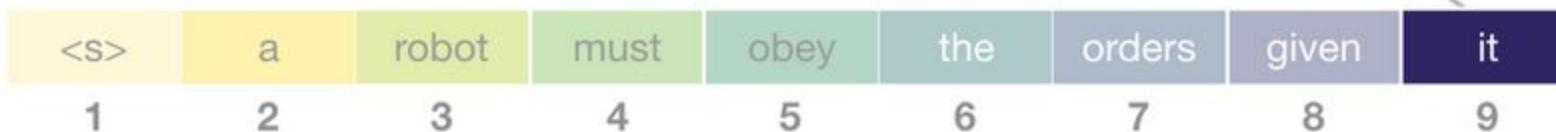
Michigan Tech

GPT2 Self-Attention

attention head #1

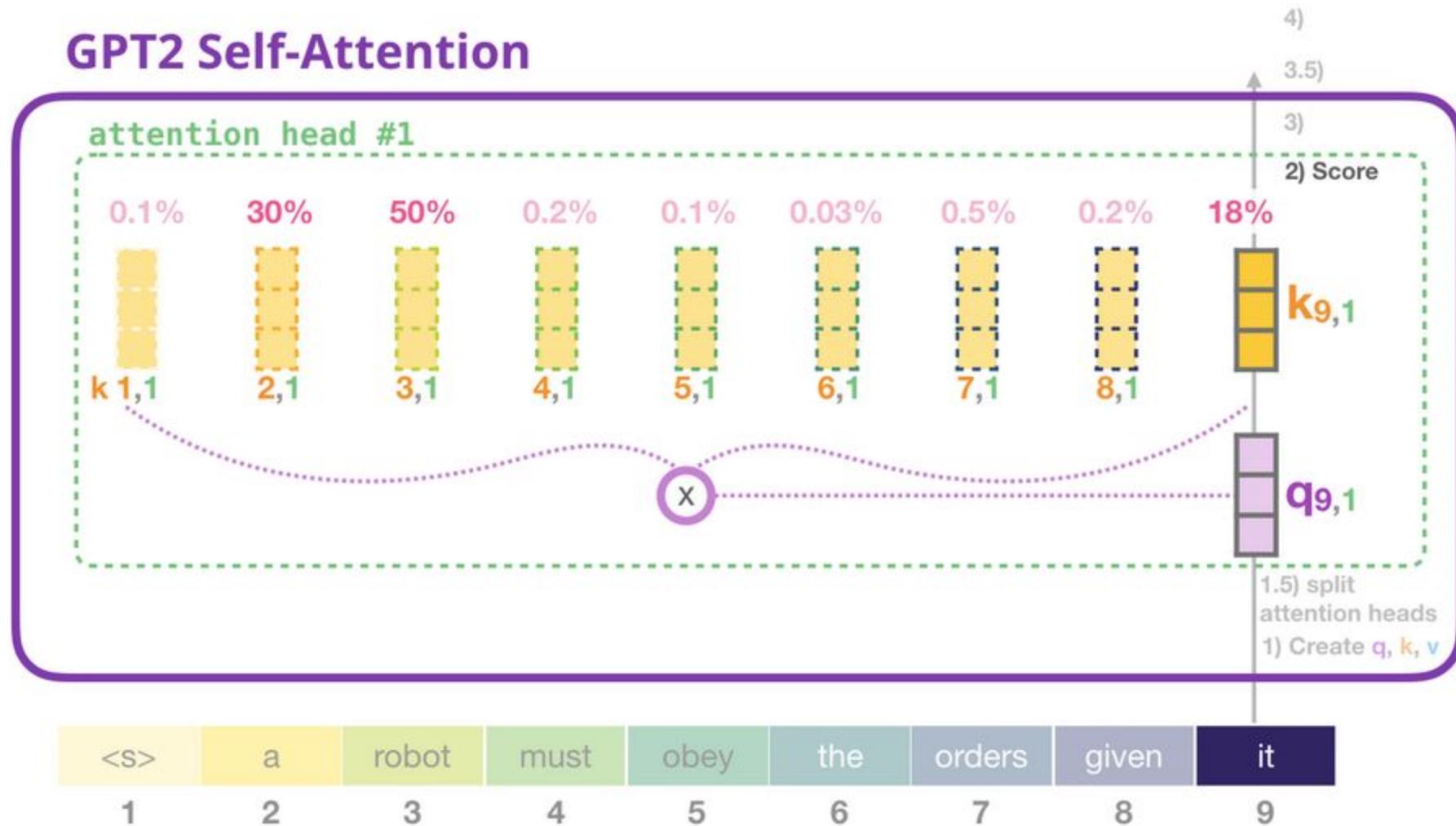


- 4)
- 3.5)
- 3)
- 2) Score
- 1.5) split attention heads
- 1) Create q, k, v



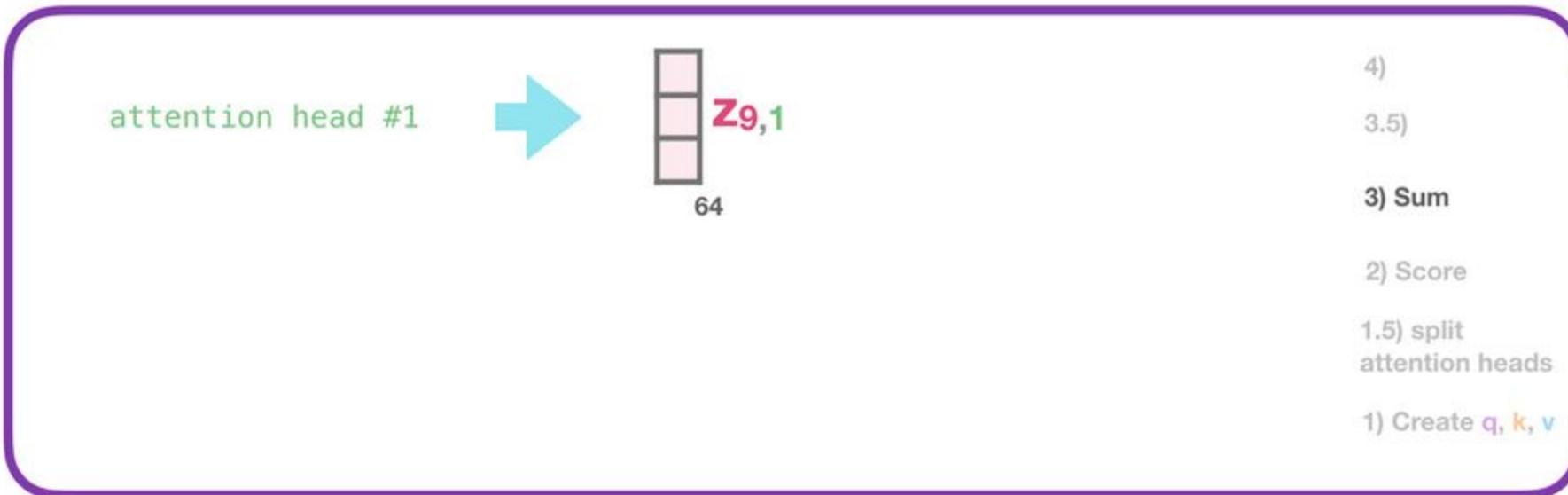
Michigan Tech

GPT2 Self-Attention



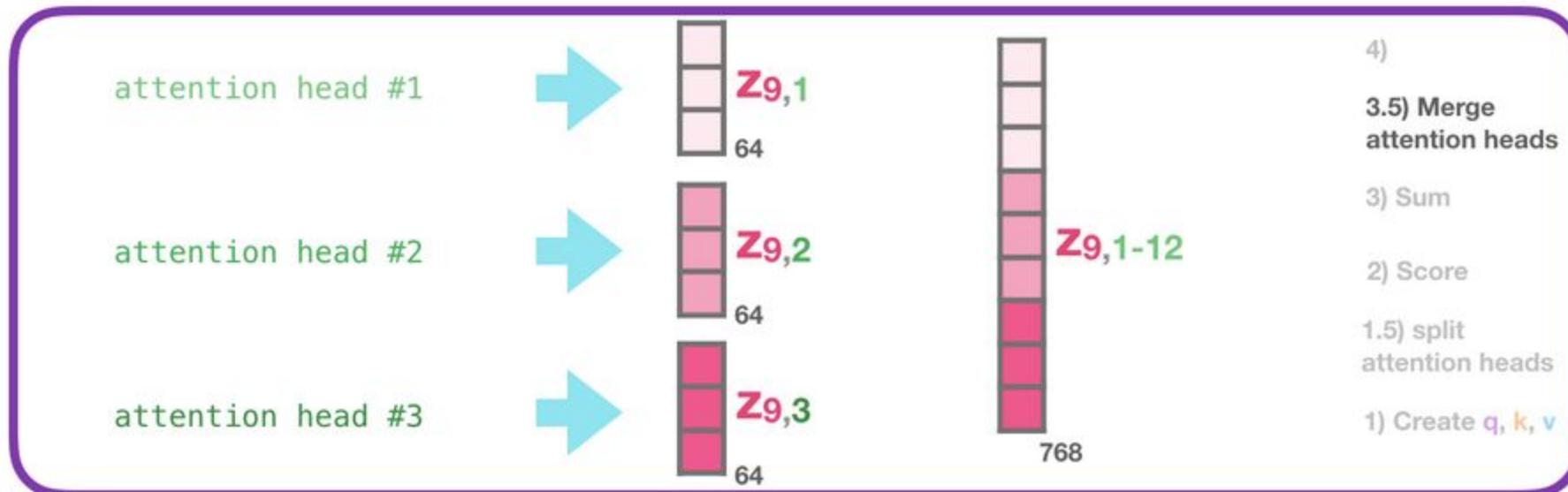
Michigan Tech
1885

GPT2 Self-Attention



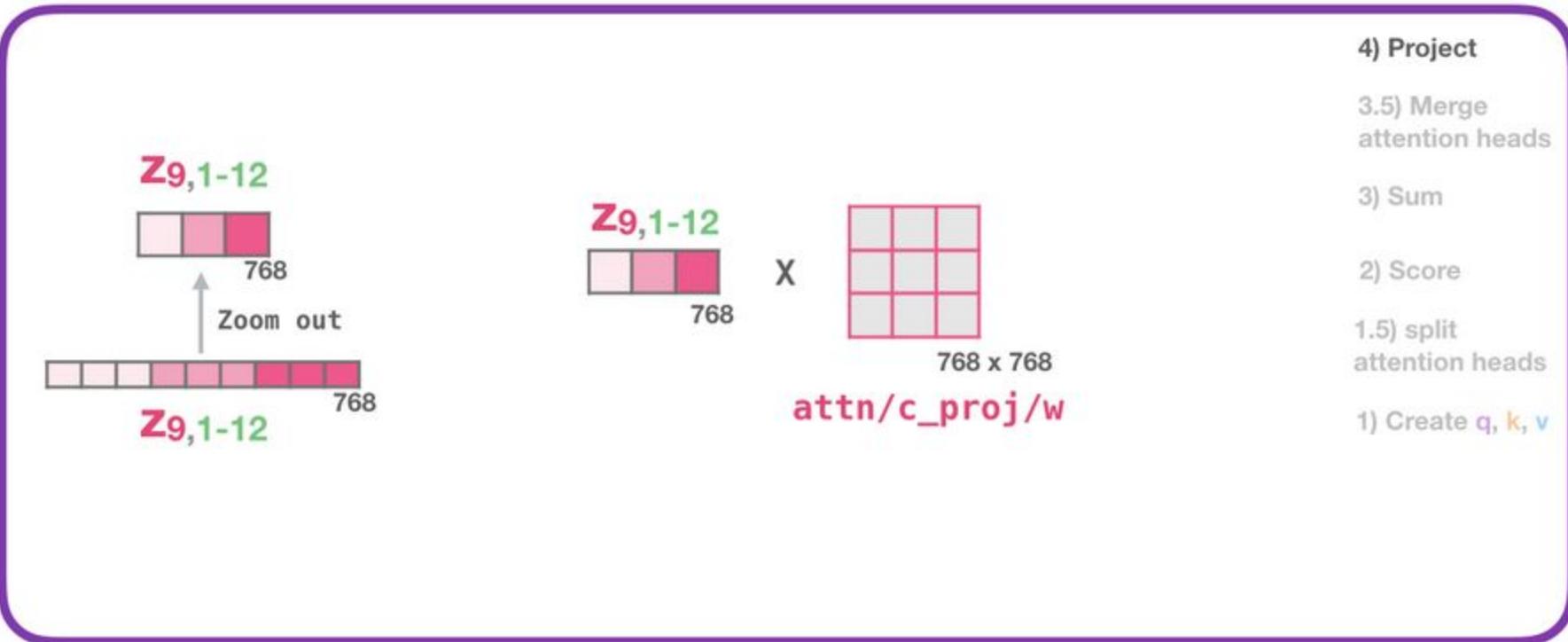
Michigan Tech

GPT2 Self-Attention

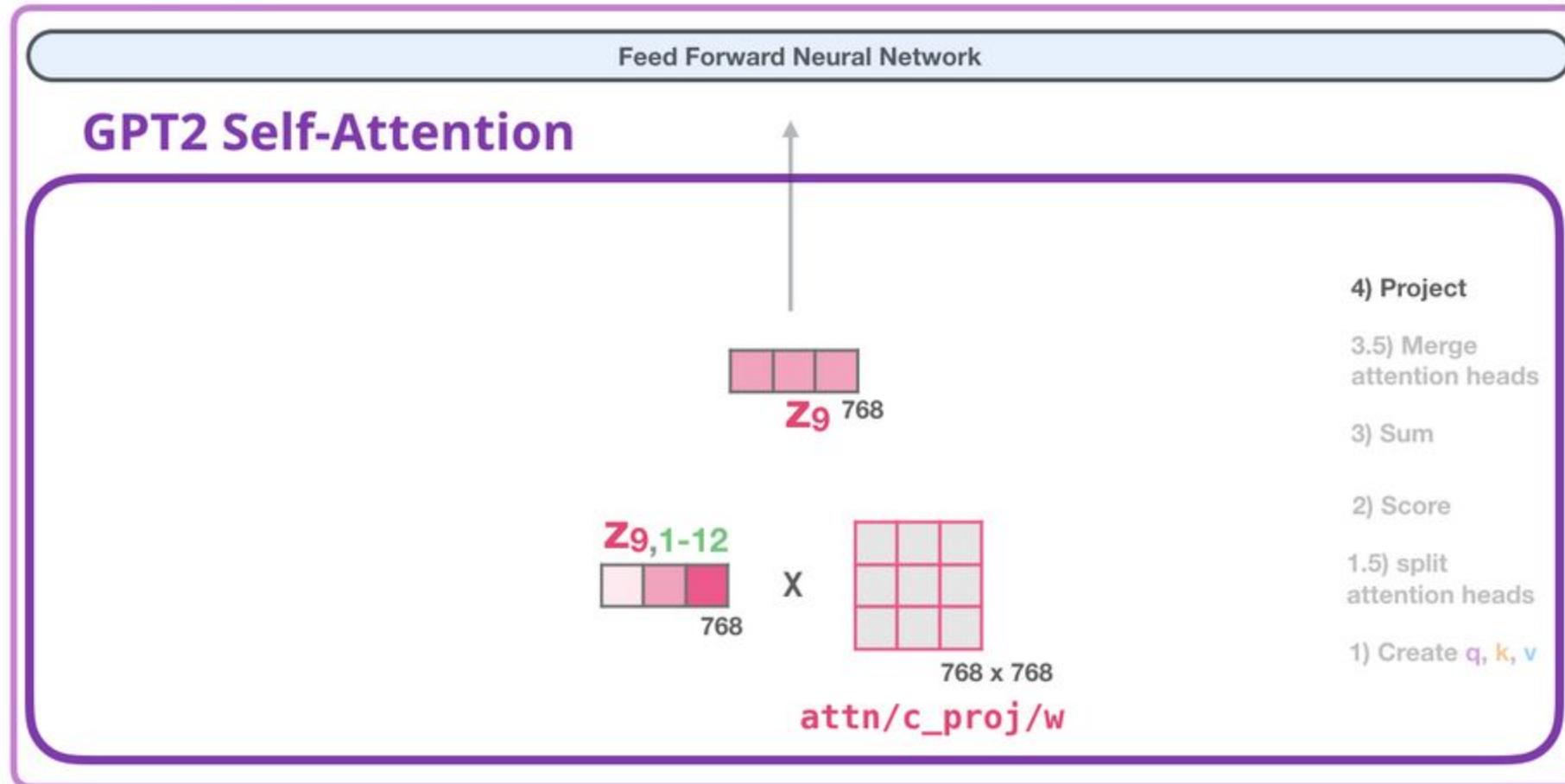


Michigan Tech

GPT2 Self-Attention



Michigan Tech



GPT2 Fully-Connected Neural Network

$$\begin{matrix} \text{Z}_9 \\ \text{768} \end{matrix} \times \begin{matrix} \text{mlp/c_fc/w} \\ \text{768 X 768*4} \end{matrix} = \begin{matrix} \text{768*4} \end{matrix}$$

2)

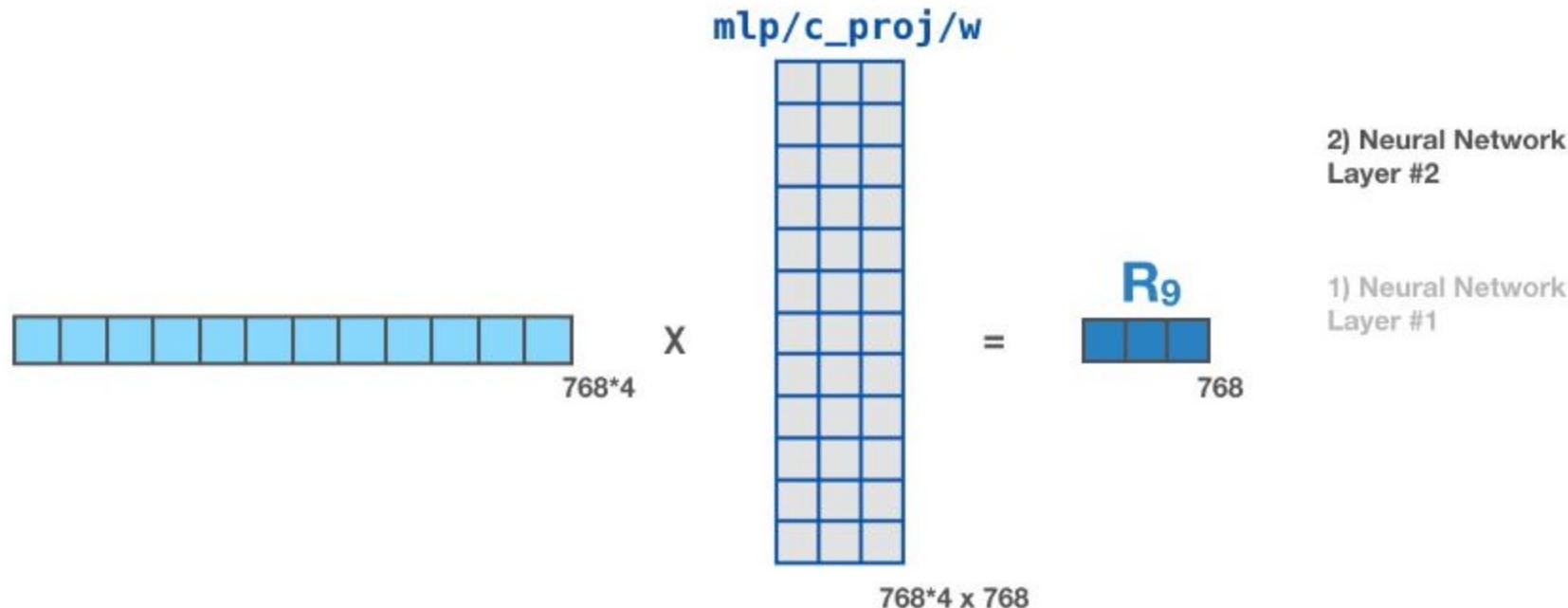
1) Neural Network
Layer #1

(Not shown: A bias vector)



Michigan Tech

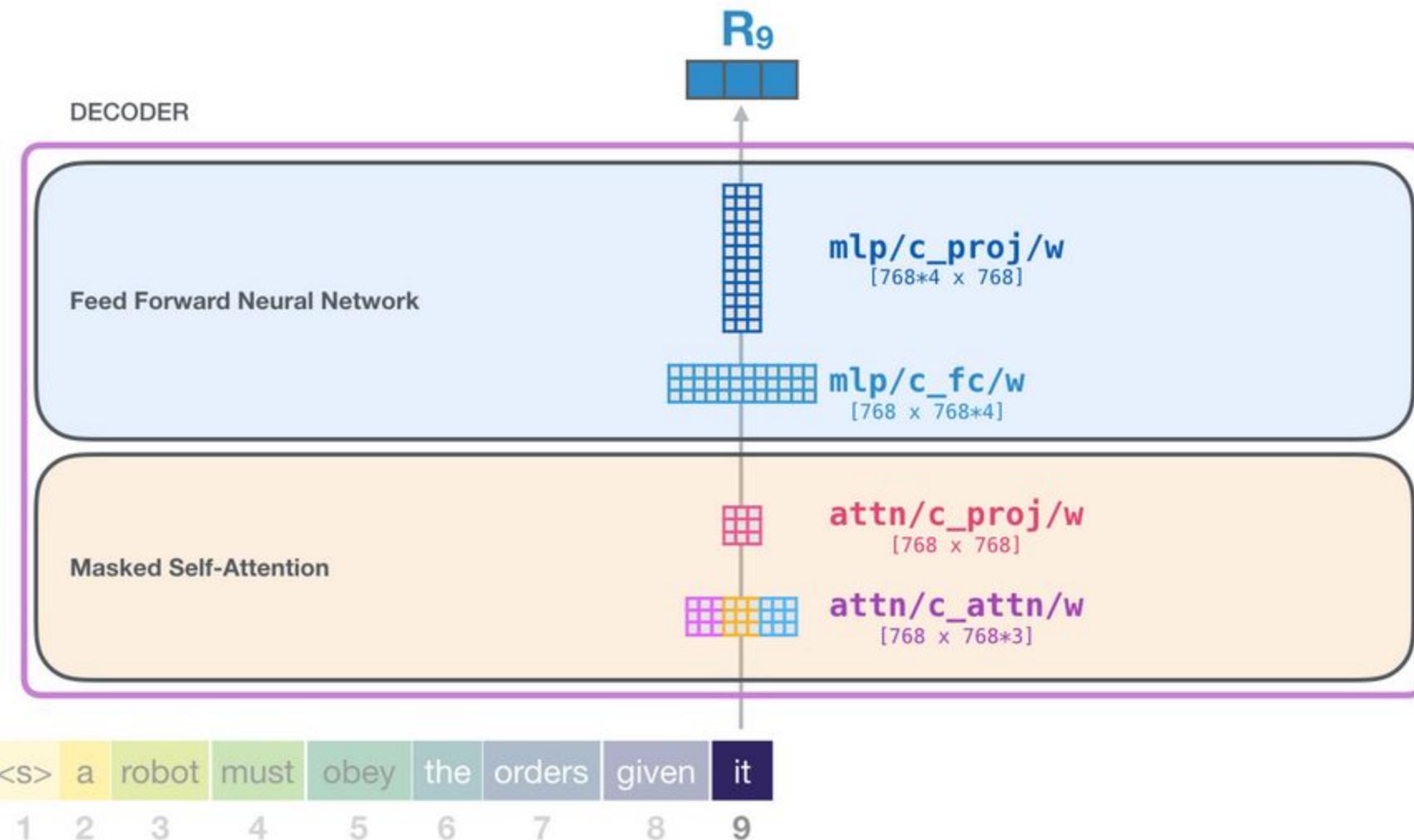
GPT2 Fully-Connected Neural Network



(Not shown: A bias vector)



Michigan Tech



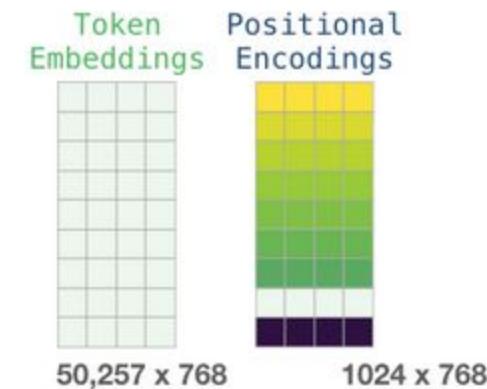
And each block has its own set of those weights. On the other hand, the model has only one token embedding matrix.



Michigan Tech



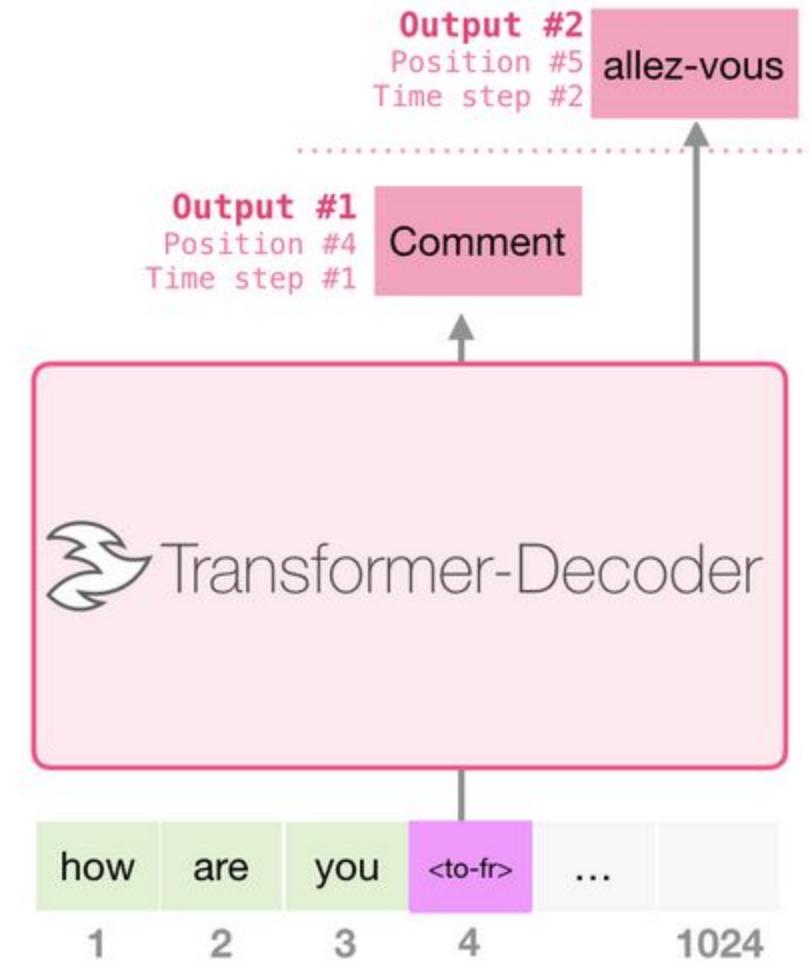
GPT-2 SMALL



Michigan Tech

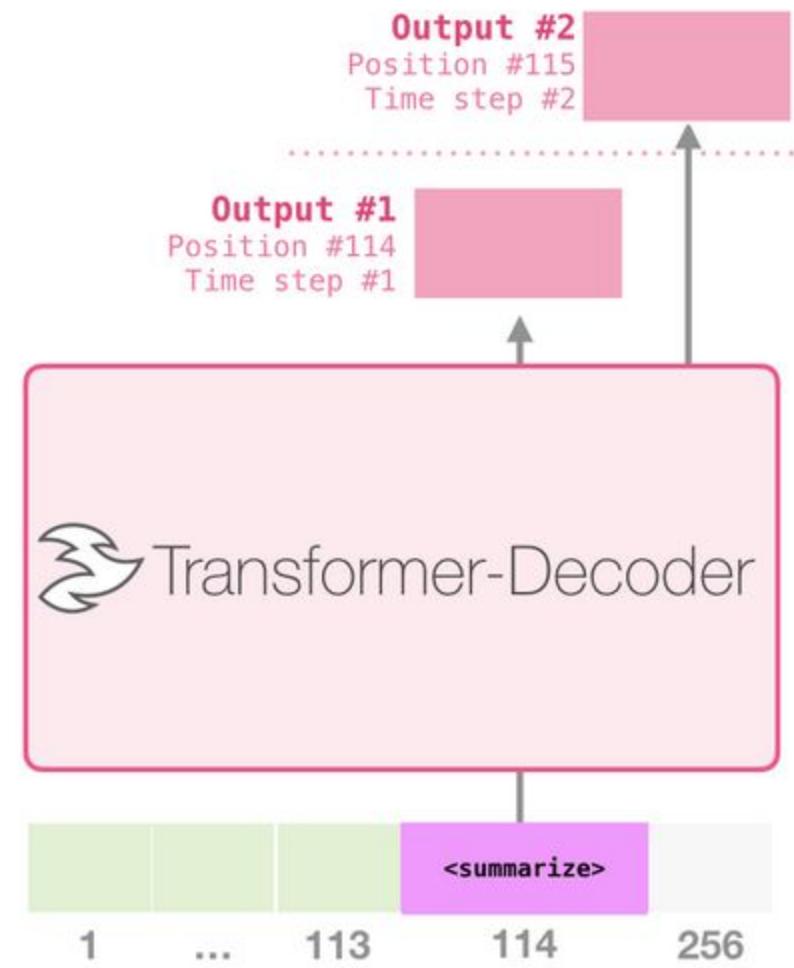
Training Dataset

I	am	a	student	<to-fr>	je	suis	étudiant
let	them	eat	cake	<to-fr>	Qu'ils	mangent	de
good	morning	<to-fr>	Bonjour				



Training Dataset

Article #1 tokens	<summarize>	Article #1 Summary
Article #2 tokens	<summarize>	Article #2 Summary padding
Article #3 tokens	<summarize>	Article #3 Summary



Critique / Limitations / Open Issues

- Baseline transformers use $O(n^2)$ in memory/computation.
 - increasing size of length increases computation quadratically
- Fixed length unlike RNN of arbitrary length
- Inability to process input sequentially (not like how the human brain works)
- Paper addressing limitation of transformers:
<https://arxiv.org/pdf/2002.09402.pdf>
 - Limited access to long memory
 - Limited ability to update state.



Michigan Tech

The challenge of long document summarization

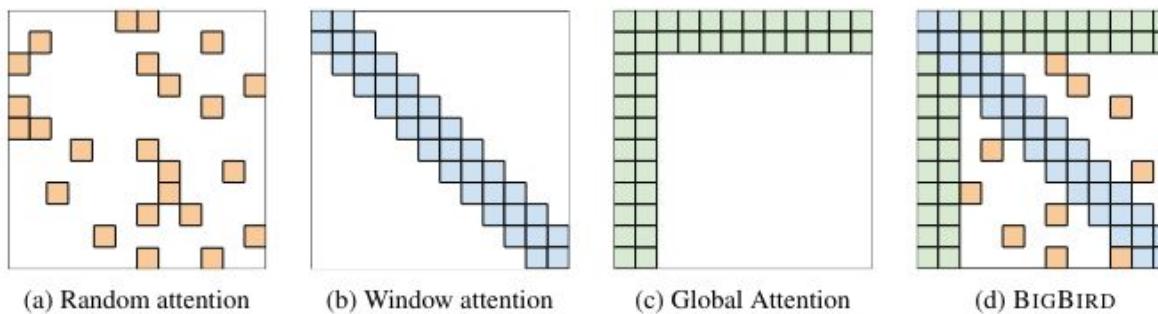
- Most transformer models use self-attention, which requires $O(n^2)$ time and memory complexity for n-length input
- This limits input size, historically as small as 512 tokens (2022)



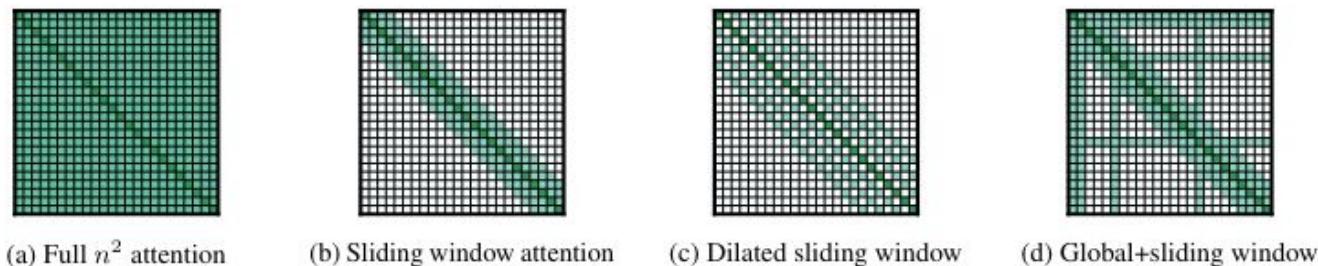
Michigan Tech

The challenge of long document summarization, modifying attention

1. Sparse Attention - Used by BigBird[1] and BigBirdPegasus



2. Sliding Dilated Window + Global Attention used by Longformer [2]

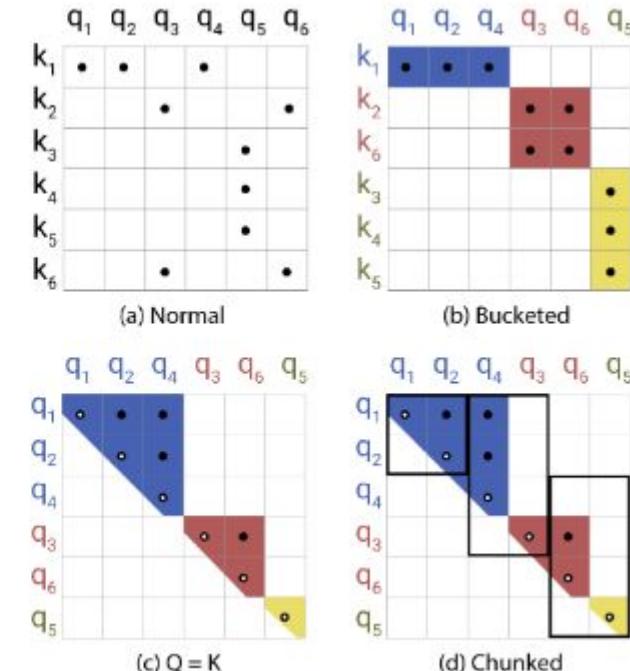
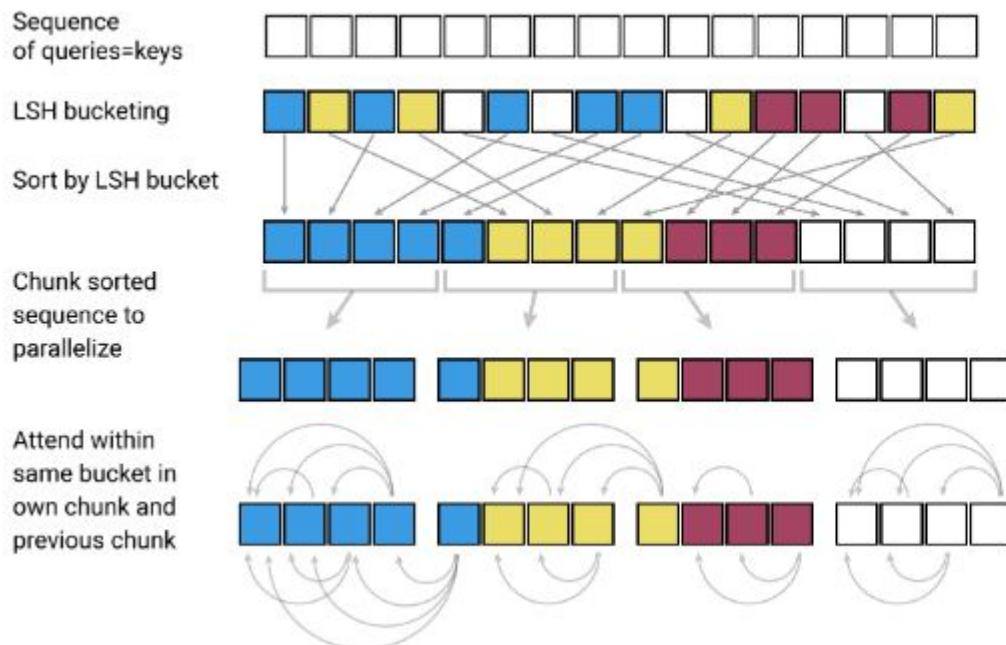


1. Zaheer, Manzil, et al. "Big Bird: Transformers for Longer Sequences." *NeurIPS*. 2020.
2. Beltagy, Iz, Matthew E. Peters, and Arman Cohan. "Longformer: The long-document transformer." *arXiv preprint arXiv:2004.05150* (2020).
3. Kitaev, Nikita, Łukasz Kaiser, and Anselm Levskaya. "Reformer: The efficient transformer." *arXiv preprint arXiv:2001.04451* (2020).



Michigan Tech

Reformer - Locality Sensitive Hashing Attention

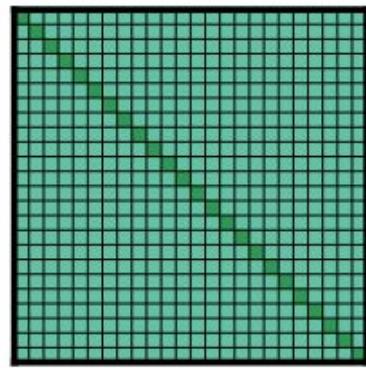


1. Zaheer, Manzil, et al. "Big Bird: Transformers for Longer Sequences." *NeurIPS*. 2020.
2. Beltagy, Iz, Matthew E. Peters, and Arman Cohan. "Longformer: The long-document transformer." *arXiv preprint arXiv:2004.05150* (2020).
3. Kitaev, Nikita, Łukasz Kaiser, and Anselm Levskaya. "Reformer: The efficient transformer." *arXiv preprint arXiv:2001.04451* (2020).

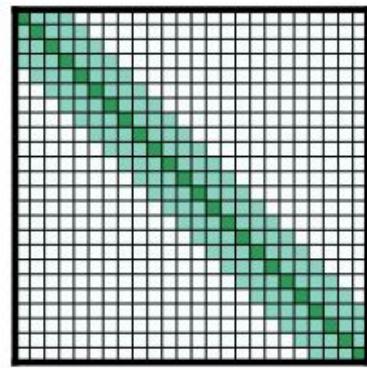


Michigan Tech

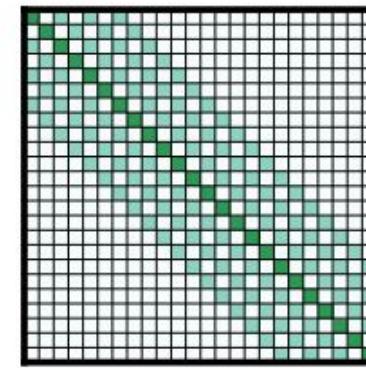
Longformer



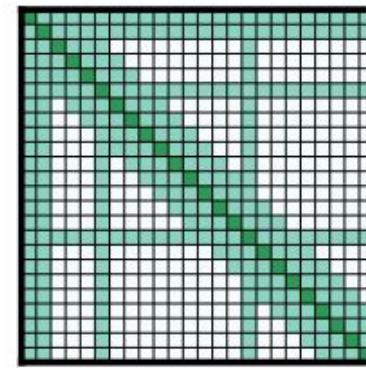
(a) Full n^2 attention



(b) Sliding window attention



(c) Dilated sliding window



(d) Global+sliding window

Figure 2: Comparing the full self-attention pattern and the configuration of attention patterns in our Longformer.



Michigan Tech

Longformer

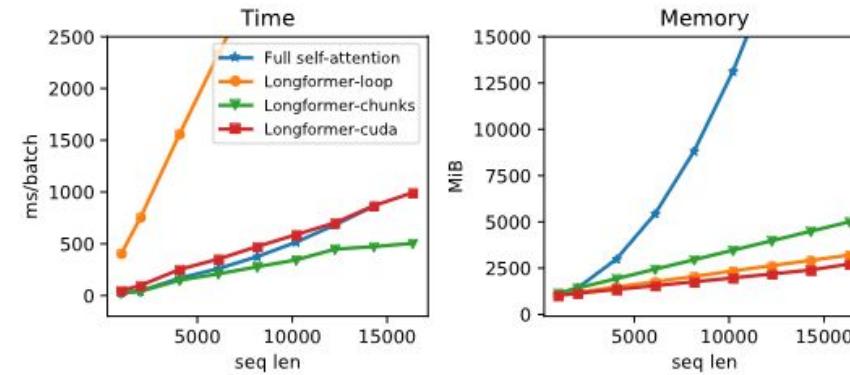


Figure 1: Runtime and memory of full self-attention and different implementations of Longformer’s self-attention; Longformer-loop is non-vectorized, Longformer-chunk is vectorized, and Longformer-cuda is a custom cuda kernel implementations. Longformer’s memory usage scales linearly with the sequence length, unlike the full self-attention mechanism that runs out of memory for long sequences on current GPUs. Different implementations vary in speed, with the vectorized Longformer-chunk being the fastest. More details are in section 3.2.



Big Bird

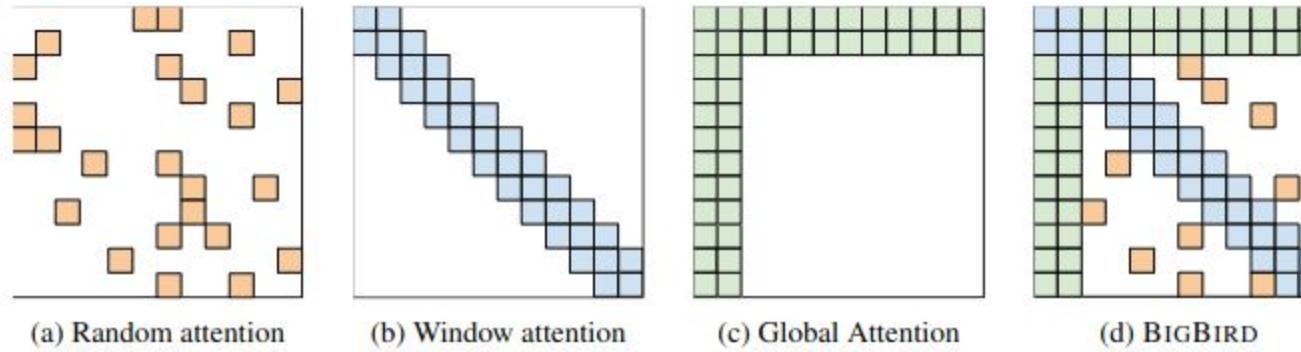
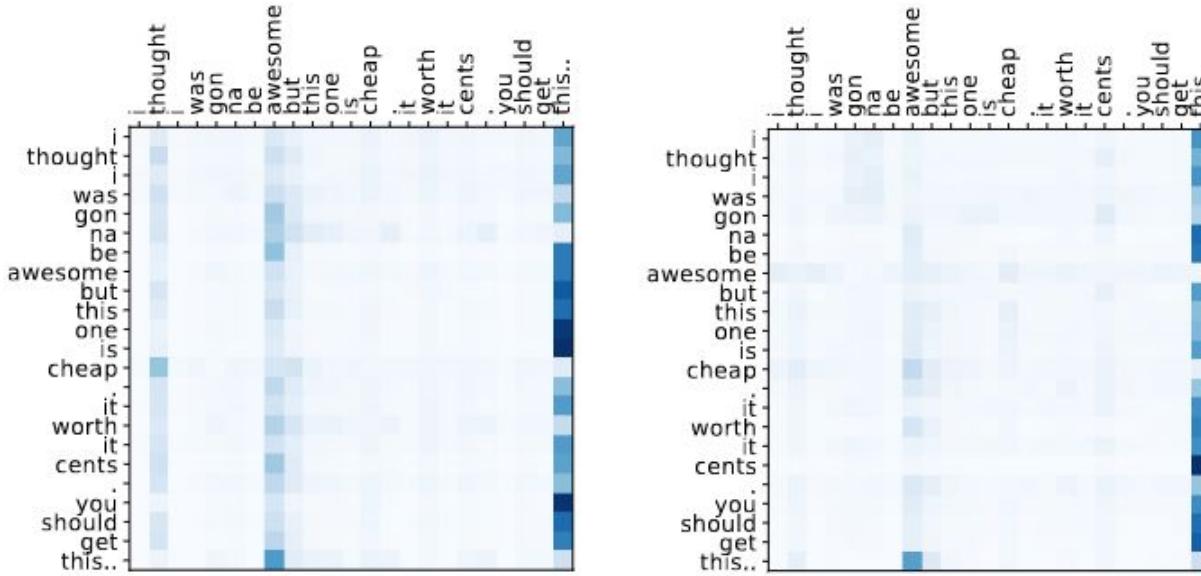


Figure 1: Building blocks of the attention mechanism used in BIGBIRD. White color indicates absence of attention. (a) random attention with $r = 2$, (b) sliding window attention with $w = 3$ (c) global attention with $g = 2$. (d) the combined BIGBIRD model.



Michigan Tech

Smart Bird



(a) Attention heatmap of a 4-dim Transformer. (b) Attention heatmap of a 256-dim Transformer.

Figure 1: The attention heatmaps learned by Transformers with 4 or 256 hidden dimensions.

Method	CNN/DM			PubMed		
	R-1	R-2	R-L	R-1	R-2	R-L
Transformer	38.51	16.08	35.90	34.43	11.84	31.76
Sparse Transformer	38.02	15.50	35.10	37.19	14.63	33.85
Longformer	37.99	15.34	35.28	37.04	14.41	33.82
BigBird	38.57	15.80	35.86	37.77	15.16	34.53
Smart Bird	39.22	16.96	37.04	38.83	16.01	35.49

Table 4: Results of text summarization.



Smart Bird

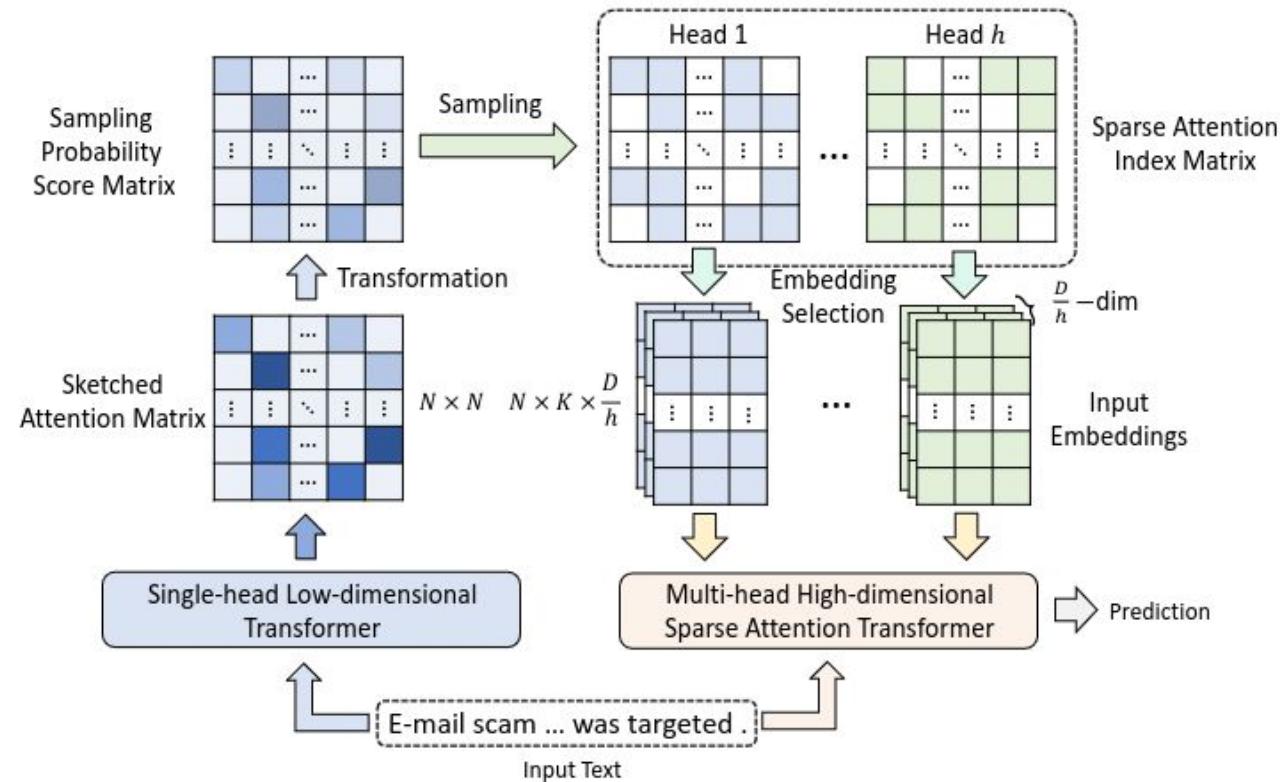


Figure 2: The overall framework of *Smart Bird*.



Michigan Tech

Sparse Transformer

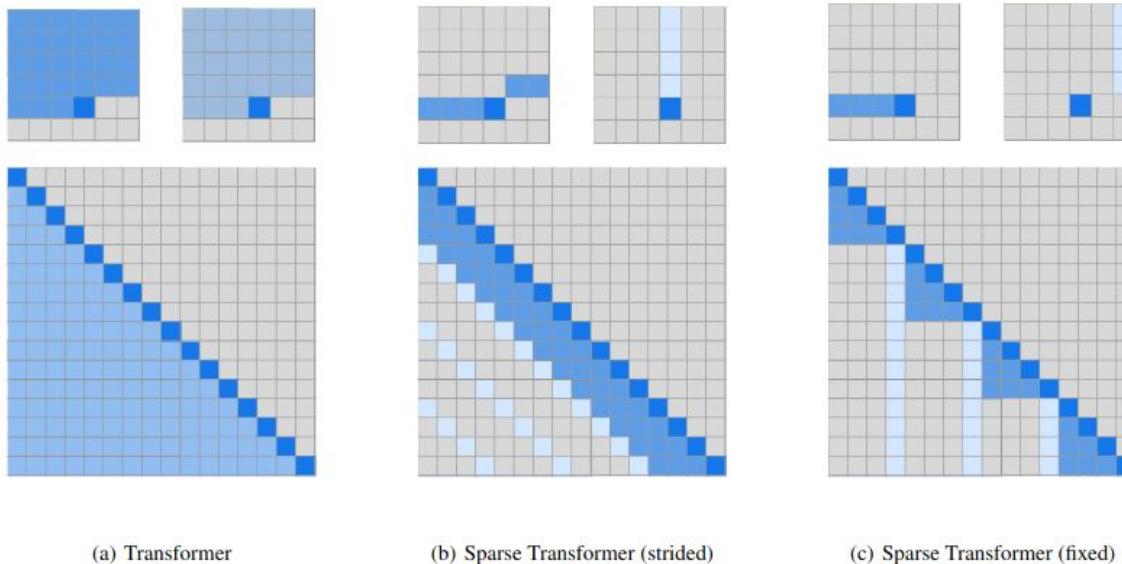


Figure 3. Two 2d factorized attention schemes we evaluated in comparison to the full attention of a standard Transformer (a). The top row indicates, for an example 6x6 image, which positions two attention heads receive as input when computing a given output. The bottom row shows the connectivity matrix (not to scale) between all such outputs (rows) and inputs (columns). Sparsity in the connectivity matrix can lead to significantly faster computation. In (b) and (c), full connectivity between elements is preserved when the two heads are computed sequentially. We tested whether such factorizations could match in performance the rich connectivity patterns of Figure 2.



Michigan Tech

Bertology

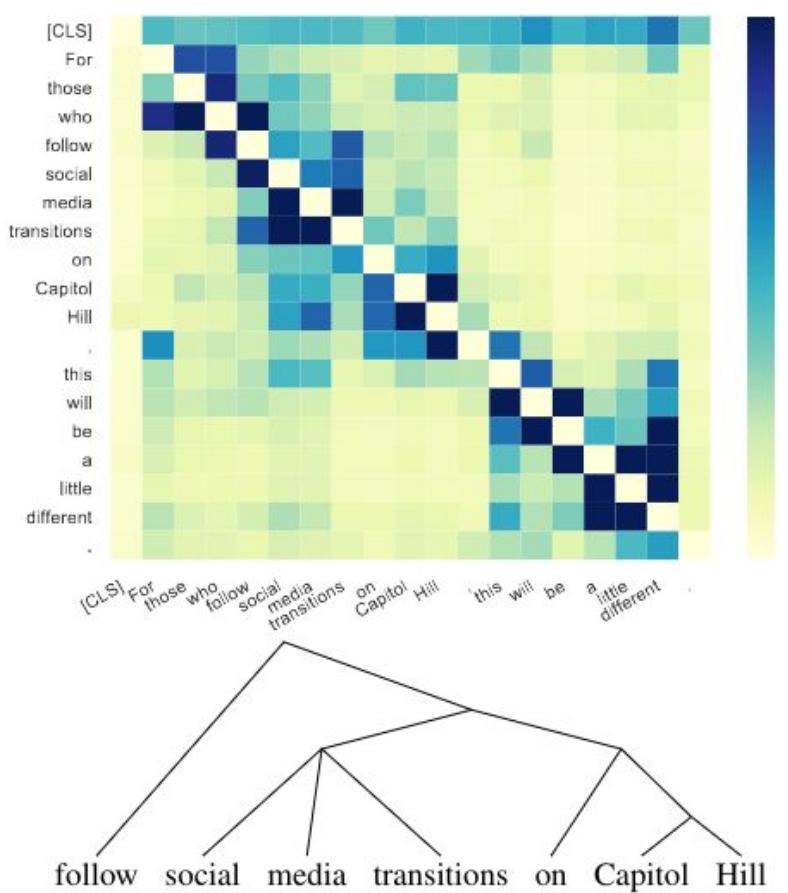


Figure 1: Parameter-free probe for syntactic knowledge: words sharing syntactic subtrees have larger impact on each other in the MLM prediction (Wu et al., 2020)



Michigan Tech

Bertology

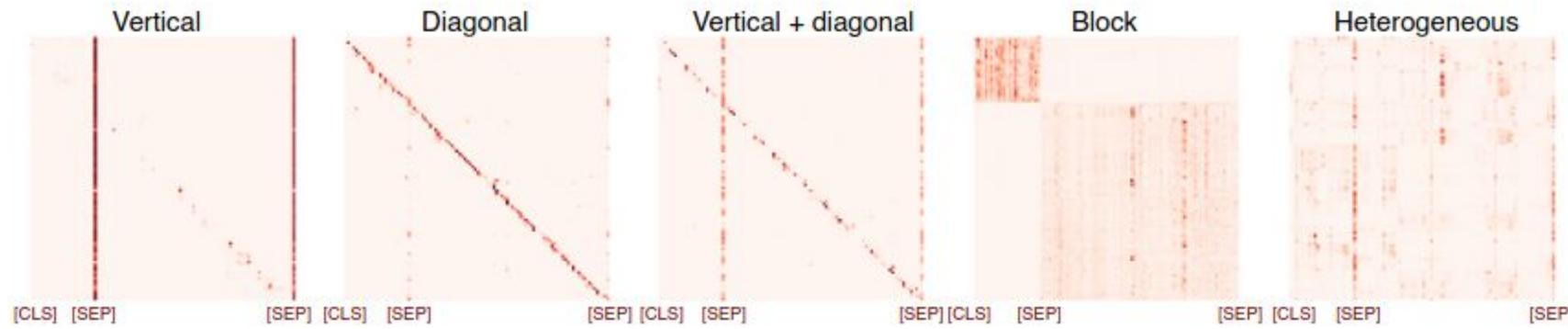


Figure 3: Attention patterns in BERT (Kovaleva et al., 2019)



Michigan Tech

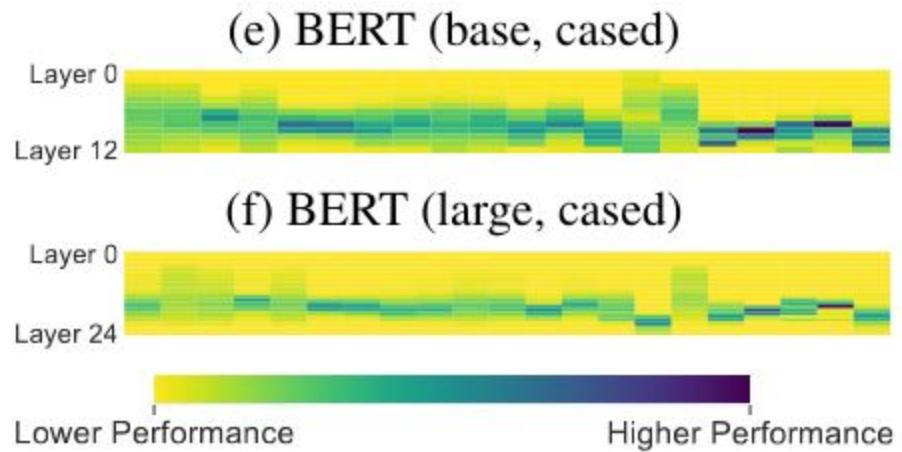


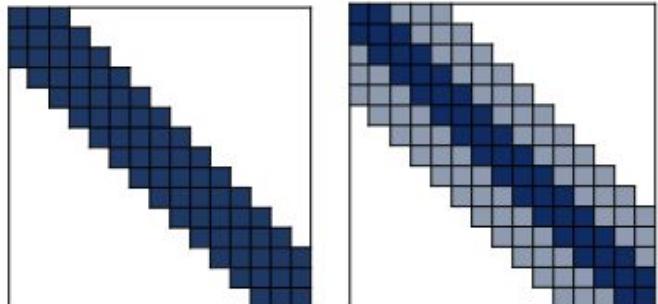
Figure 4: BERT layer transferability (columns correspond to probing tasks, Liu et al. (2019a).



Michigan Tech

Poolingformer

- Paper released in May, 2021 by Microsoft/Sichuan University
- Paper proposes a pooled sliding attention to gain context further away from token of interest



(a) Single-level local attention (b) Two-level pooling attention

Figure 1. (a): The receptive field of single-level local attention (b): The receptive field of our two-level pooling attention.

Table 4. The results on the arXiv test set. We report the results of ROUGE-1 (R-1), ROUGE-2 (R-2), and ROUGE-L (R-L).

Model	R-1	R-2	R-L
Sent-PTR (Pilault et al., 2020)	42.32	15.63	38.06
Extr-Abst-TLM (Pilault et al., 2020)	41.62	14.69	38.03
PEGASUS (Zhang et al., 2020)	44.21	16.95	38.83
Dancer (Gidiotis & Tsoumakas, 2020)	45.01	17.60	40.56
BigBird (Zaheer et al., 2020)	46.63	19.02	41.77
LED _{4k} (Beltagy et al., 2020)	44.40	17.94	39.76
LED _{16k} (Beltagy et al., 2020)	46.63	19.62	41.83
Poolingformer _{4k}	47.86	19.54	42.35
Poolingformer _{16k}	48.47	20.23	42.69



Michigan Tech

Questions + Comments?



Michigan Tech

Resources used

http://cs231n.stanford.edu/slides/2023/lecture_8.pdf

<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

<https://huggingface.co/blog/how-to-generate>

<http://jalammar.github.io/illustrated-gpt2/>



Michigan Tech