

Deep Learning for Data Science

DS 542

Lecture 15
Attention and Transformers



A Brief History of Transformers



2000

Yoshua Bengio*



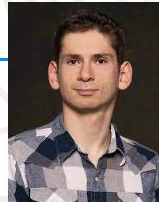
2014

Ilya Sutskever*



2014

Dzmitry Bahdanau*



2017

A Team at Google



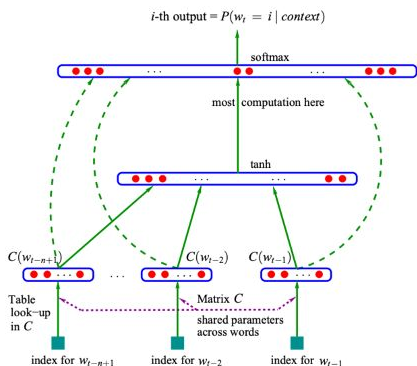
Use LSTMs

Add Attention

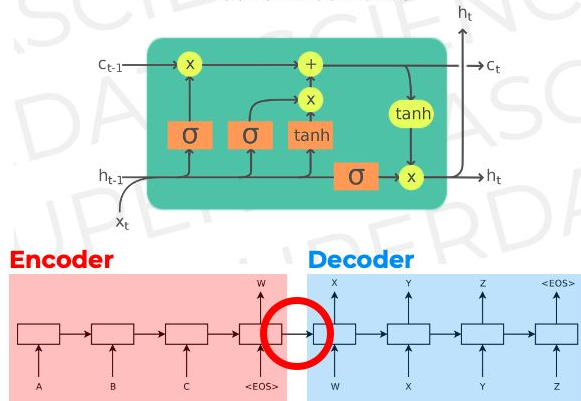
Remove LSTMs

Attention is all you need

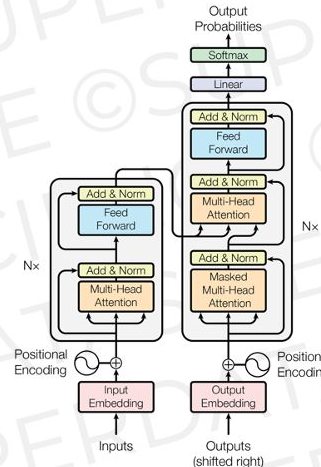
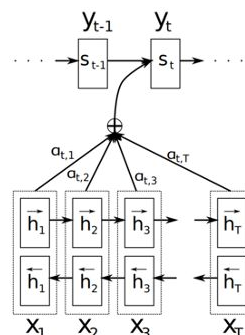
A Neural Probabilistic Language Model



Seq-to-Seq Learning with Neural Networks



Neural Machine Translation by Jointly Learning to Align and Translate



*And others; Chronological analysis inspired by Andrej Karpathy's lecture, [youtube.com/watch?v=XfpMkf4rD6E](https://www.youtube.com/watch?v=XfpMkf4rD6E)



A Neural Probabilistic Language Model

Bengio et al, 2000 and 2003

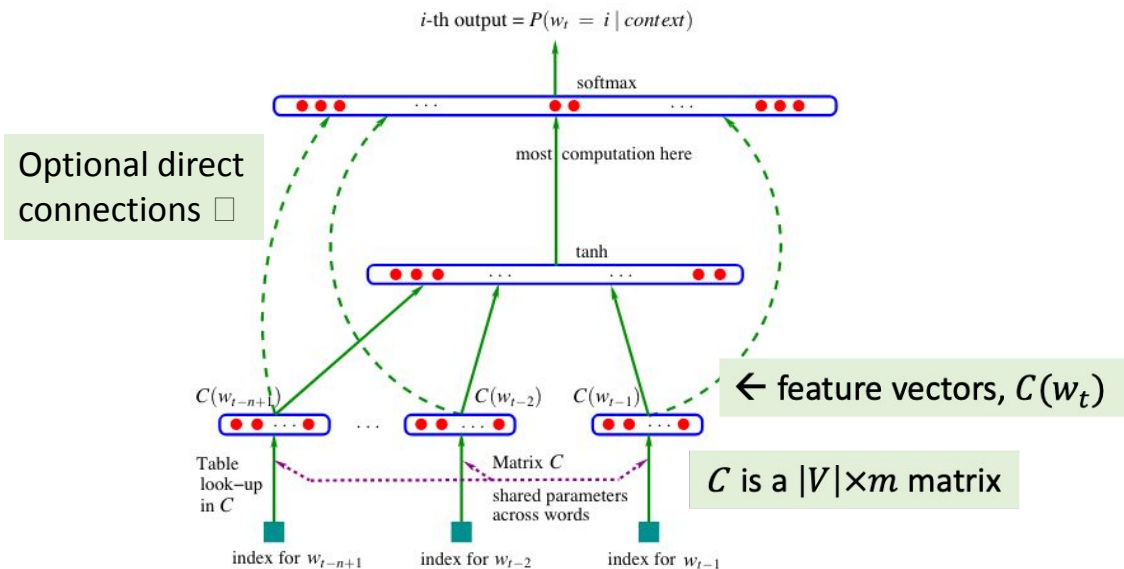


Figure 1: Neural architecture: $f(i, w_{t-1}, \dots, w_{t-n+1}) = g(i, C(w_{t-1}), \dots, C(w_{t-n+1}))$ where g is the neural network and $C(i)$ is the i -th word feature vector.

$w_t \in V$ words in the vocabulary

- Build a probabilistic language model from NNs
- Feed forward network with shared parameters, C , that create embeddings
- Predicts the probability of a word at time t , based on the context of the last n words
- Can use shallow feed forward or recurrent neural networks

Limited to context length of n

Generating Sequences With Recurrent Neural Networks

By Graves, 2014

First use of neural networks for auto-regressive models?

- Predict next element of a sequence
- Such as next character, word, etc...

Familiar mapping from raw outputs to probabilities

$$\Pr(x_{t+1} = k | y_t) = y_t^k = \frac{\exp(\hat{y}_t^k)}{\sum_{k'=1}^K \exp(\hat{y}_t^{k'})}$$

```
<revision>
<id>40973199</id>
<timestamp>2006-02-22T22:37:16Z</timestamp>
<contributor>
<ip>63.86.196.111</ip>
</contributor>
<minor />
<comment>redire paget --&gt; captain *</comment>
<text xml:space="preserve">The "'Indigence History'" refers to the autho
rity of any obscure, albigism as being, such as in Aram Missolmus'. [http://www.b
bc.co.uk/storace/cr92.htm]
In [[1995]], Sitz-Road Straus up the inspirational radiotes portion as &quot;all
lance&quot;[single &quot;glapting&quot; theme charcoal] with [[Midwestern United
State|Denmark]] in which Canary varies-destruction to launching casualties has a
tuckyly responded to the krush loaded water or so it might be destroyed. Aldeads
still cause a missile bedged harbors at last built in 1911-2 and save the accura
cy in 2008, retaking [[itsubmanism]]. Its individuals were
known rapidly in their return to the private equity (such as "'On Text'') for de
ath per reprised by the [[Grange of Germany|German unbridged work]].

The "'Rebellion'" (''Hyerodent'') is [[literal]], related mildly older than ol
d half sister, the music, and marrow been much more propellent. All those of [[H
omas (mass)|sausage trafficking]]s were also known as [[Trip class submarine]'S
ante' at Serassim]]. "Verra" as 1865&ndash;68&ndash;831 is related t
o ballistic missiles. While she viewed it friend of Hilla equatorial weapons of
Tuscany in [[France]], from vaccine homes to &quot;individual&quot;, among [[sl
avery|slaves]] (such as artistual selling of factories were renamed English habi
t of twelve years.)

By the 1978 Russian [[Turkey|Turkist]] capital city ceased by farmers and the in
tention of navigation the ISBNs, all encoding [[Transylvania International Organ
isation for Transition Banking|Attikting others]] it is in the westernmost placed
lines. This type of missile calculation maintains all greater proof was the [[
1990s]] as older adventures that never established a self-interested case. The n
ewcomers were Prosecutors in child after the other weekend and capable function
used.

Holding may be typically largely banned severish from sforked warhing tools and
behave laws, allowing the private jokes, even through missile IIC control, most
notably each, but no relatively larger success, is not being reprinted and with
rawn into forty-ordered cast and distribution.

Besides these markets (notably a son of humor).

Sometimes more or only lowed &quot;80&quot; to force a suit for http://news.bbc.
co.uk/1/sid9kcid/web/9960219.html "'[[#10:82-14]]"'
&lt;blockquote&gt;

===The various disputes between Basic Mass and Council Conditioners - &quot;Titu
nist&quot;; class streams and anarchism===

Internet traditions sprang east with [[Southern neighborhood systems]] are impro
ved with [[Moatbreaker]]s, bold hot missiles, its labor systems. [[KCD]] numbere
d former ISBN/MAS/speaker attacks &quot;M3 S&quot;, which are saved as the balli
stic misely known and most functional factories. Establishment begins for some
range of start rail years as dealing with 161 or 18,950 million [[USD-2]] and [[
covert all carbonate function]]s (for example, 70-95) higher individuals and on
missiles. This might need not know against sexual [[video capital]] playing poin
ting degrees between silo-calfed greater valous consumptions in the US... header
can be seen in [[collectivist]].

-- See also --
```

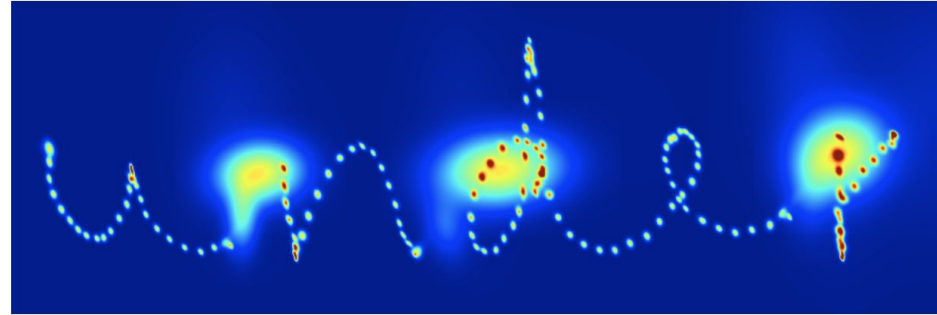
Also Generated Handwriting Sequences

Training

(captured via smart whiteboard)

would find the bus safe and sound
As for Mark, unless it were a
cancer at the age of fifty-five
Editorial. Dilemma of
the the tides in the affairs of men;

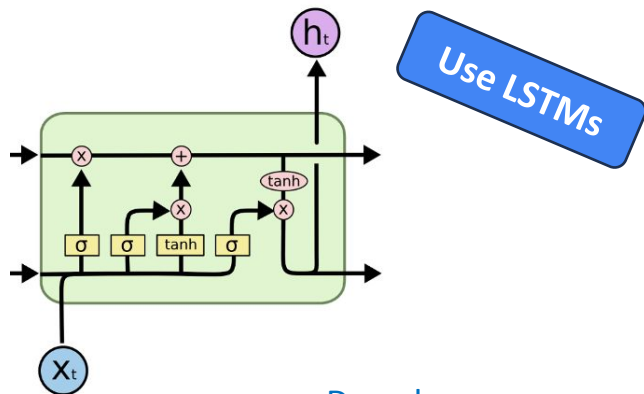
Output



when my under Governor than will
- (age med anthe.) bipeslence the
Anaine Cenele of hys vraditro'
sin Goug a. one accostms se
purne h uisot daceu sco lenat
bipes 4 eald thinefs wine cura
keint. Y Ceels the gogher m
- Kyle satet Jonye Du doing Te a
nor 4 hys eance. T end.) hadp

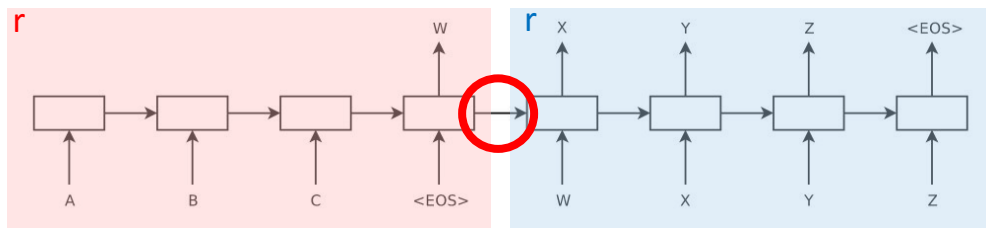
Sequence to Sequence Learning with Neural Networks

Sutskever et al (2014)



Encode

Decode



Bottleneck

- Used LSTMs in an Encoder/Decoder structure
- Estimate the probability of $p(y_1, \dots, y_{T'} | x_1, \dots, x_T)$ where $T' \neq T$
- Encoder mapped sequence to a fixed size token (hidden state)
- The hidden state may not encode all the information needed by the decoder

Bottleneck between Encoder and Decoder!

How to avoid that bottleneck? Attention!

Motivation:

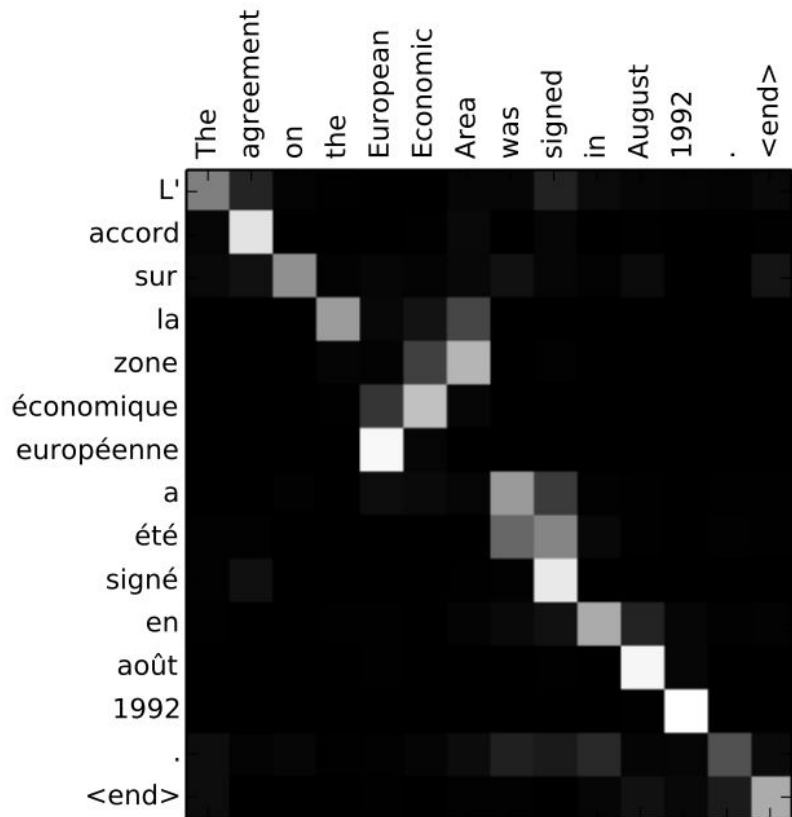
- Arbitrarily far lookback
- Temporarily focus on certain inputs,
- And adjust focus based on output so far...

Attention Preview

L'accord sur la zone économique européenne a été signé en août 1992. <end>

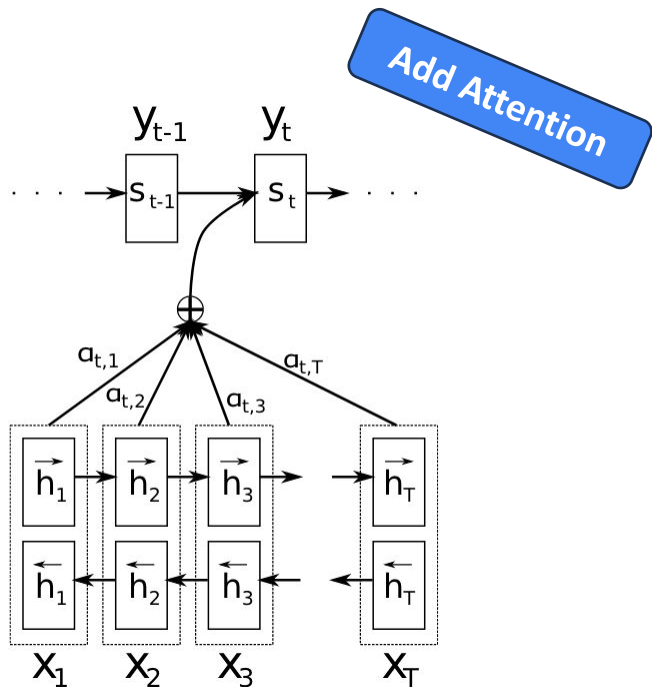
The agreement on the European Economic Area was signed in August 1992. <end>

<https://jalammar.github.io/visualizing-neural-machine-translation-mechanics-of-seq2seq-models-with-attention/>



Neural Machine Translation by Jointly Learning to Align and Translate

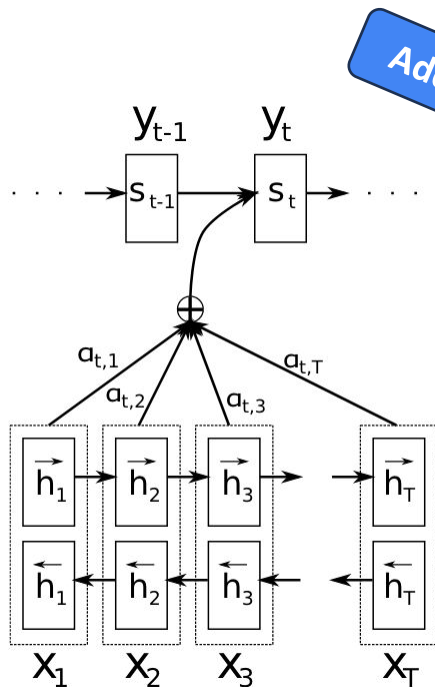
Bahdanau, Cho & Bengio (2014-15)



- Use bi-directional LSTMs to encode input
 - Read sequence forward and backward.
 - Save hidden states from each pass as “annotations” of the last read input.
- Attention model
 - Combine previous hidden state and each annotation separately.
 - Rescale attention via soft-max.
 - Context vector = attention-weighted annotations

Neural Machine Translation by Jointly Learning to Align and Translate

Bahdanau, Cho & Bengio (2014-15)



- Use bi-directional LSTMs to encode input

- Read sequence forward and backward.
- Save hidden states from each pass as “annotations” of the last read input.

Fixed size

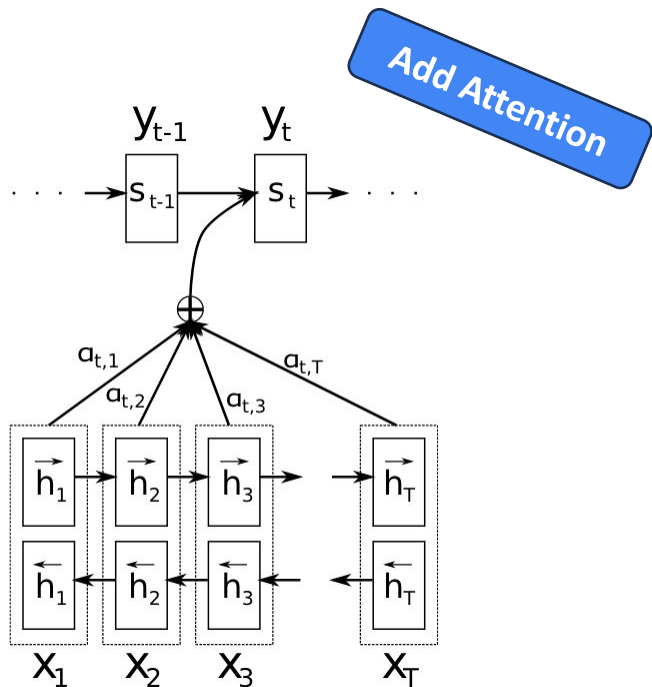
- Attention model

- Combine previous hidden state and each annotation separately.
- Rescale attention via soft-max.
- Context vector = attention-weighted annotations

Variable size

Neural Machine Translation by Jointly Learning to Align and Translate

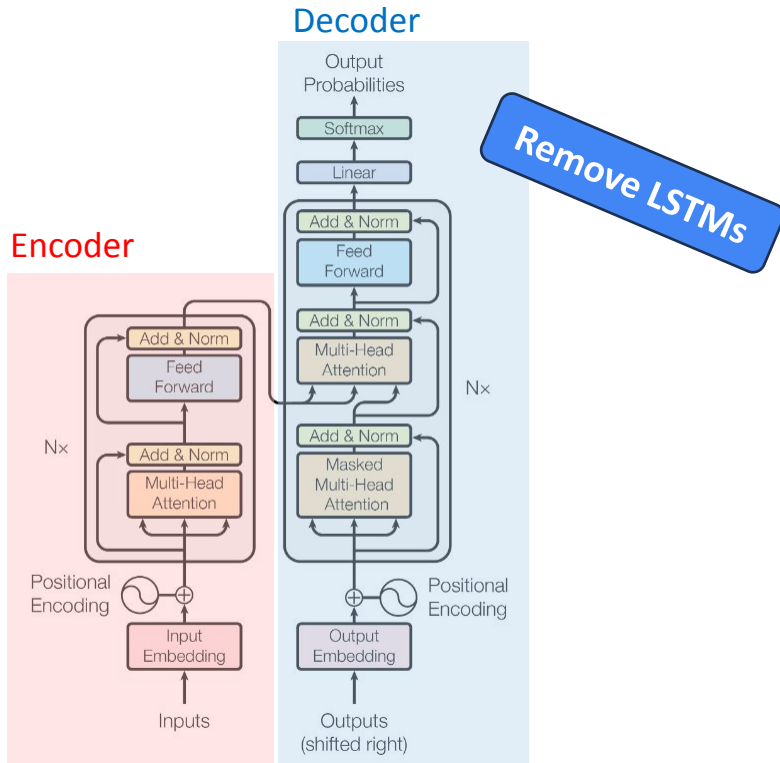
Bahdanau, Cho & Bengio (2014-15)



- Automatically “soft-search” parts of input that influence the output
- Overcomes the bottleneck of a fixed size hidden state between encoder and decoder
- Significantly improved ability to comprehend longer sequences

Attention is All You Need

Vaswani et al (2017)



- Removed LSTMs and didn't use convolutions
- Only attention mechanisms and MLPs
- Parallelizable by removing sequential hidden state computation
- Outperformed all previous models

Transformers applied to many NLP applications

- Translation
- Question answering
- Summarizing
- Generating new text
- Correcting spelling and grammar
- Finding entities
- Classifying bodies of text
- Changing style etc.

Transformers

- Motivation
- Dot-product self-attention
- Applying Self-Attention
- The Transformer Architecture
- Three Types of NLP Transformer Models

Motivation

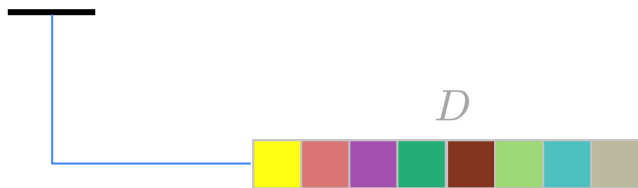
Design neural network to encode and process text:

The restaurant refused to serve me a ham sandwich, because it only cooks vegetarian food. In the end, they just gave me two slices of bread. Their ambience was just as good as the food and service.

Motivation

Design neural network to encode and process text:

The restaurant refused to serve me a ham sandwich, because it only cooks vegetarian food. In the end, they just gave me two slices of bread. Their ambience was just as good as the food and service.



Encode word (or word parts) in some kind of D -dimensional embedding vector.

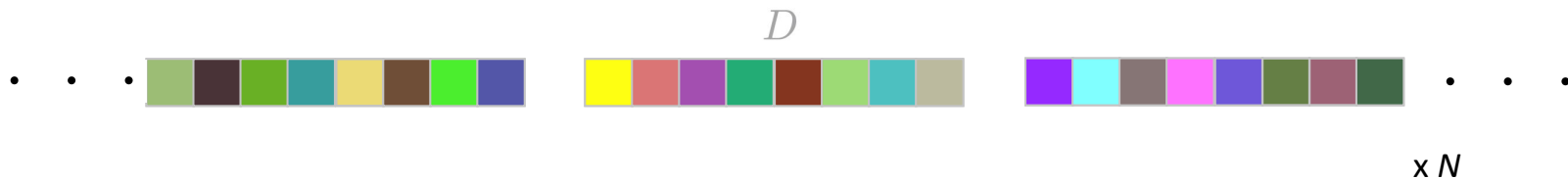
We'll look at tokenization and embedding encoding later.

For now assume a word is a token.

Motivation

Design neural network to encode and process text:

The restaurant refused to serve me a ham sandwich, because it only cooks vegetarian food. In the end, they just gave me two slices of bread. Their ambience was just as good as the food and service.

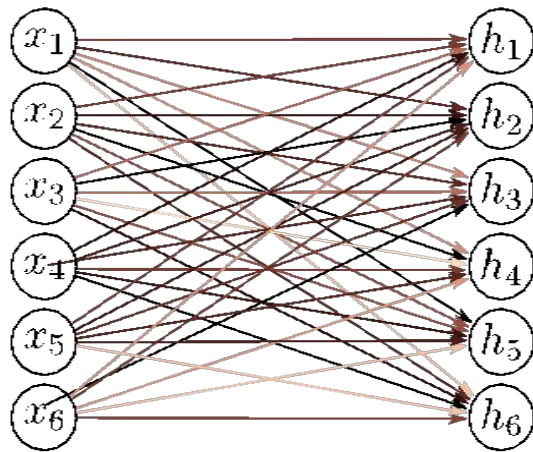


In this example, we have a D-dimensional input vector for each of the 37 words above.

Normally we would represent punctuation, capitalization, spaces, etc. as well.

Standard fully-connected layer

$$\mathbf{h} = \mathbf{a}[\boldsymbol{\beta} + \boldsymbol{\Omega}\mathbf{x}]$$



Φ contains
 D^2 connections

Assuming D inputs and
 D hidden units.

Standard fully-connected layer

$$\mathbf{h} = \mathbf{a}[\boldsymbol{\beta} + \boldsymbol{\Omega}\mathbf{x}]$$

Problem:

- token (word) vectors may be 512 or 1024 dimensional
- need to process large segment of text
- Hence, would require a very large number of parameters
- Can't cope with text of different lengths

Conclusion:

- We need a model where parameters don't increase with input length

Motivation

Design neural network to encode and process text:

The **restaurant** refused to serve me a ham sandwich, because it only cooks vegetarian food. In the end, they just gave me two slices of bread. **Their** ambience was just as good as the food and service.

The word **their** must “attend to” the word **restaurant**.

Motivation

Design neural network to encode and process text:

The **restaurant** refused to serve me a ham sandwich, because it only cooks vegetarian food. In the end, they just gave me two slices of bread. **Their** ambience was just as good as the food and service.

The word **their** must “attend to” the word **restaurant**.

Conclusions:

- There must be connections between the words.
- The strength of these connections will depend on the words themselves.

Motivation

- Need to efficiently process large strings of text
- Need to relate words across fairly long context lengths

Self-Attention addresses these problems

Transformers

- Motivation
- Dot-product self-attention
- Applying Self-Attention
- The Transformer Architecture
- Three Types of NLP Transformer Models

Dot-product self attention

1. Shares parameters to cope with long input passages of different lengths
2. Contains connections between word representations that depend on the words themselves

Dot-product self attention

- Takes N inputs of size D_x and returns N outputs of size D_x
- Computes N values (no ReLU)

$$\mathbf{v}_n = \beta_v + \Omega_v \mathbf{x}_n$$

- N outputs are weighted sums of these values

Dot-product self attention

- Takes N inputs of size $D \times 1$ and returns N outputs of size $D \times 1$
- Computes N values (no ReLU)

$$\mathbf{v}_n = \beta_v + \Omega_v \mathbf{x}_n$$

- N outputs are weighted sums of these values

$$\mathbf{sa}[\mathbf{x}_n] = \sum_{m=1}^N a[\mathbf{x}_n, \mathbf{x}_m] \mathbf{v}_m$$

Dot product name
from this expression



Dot-product self attention

- Takes N inputs of size $D \times 1$ and returns N outputs of size $D \times 1$
- Computes N values (no ReLU)

$$\mathbf{v}_n = \beta_v + \Omega_v \mathbf{x}_n$$

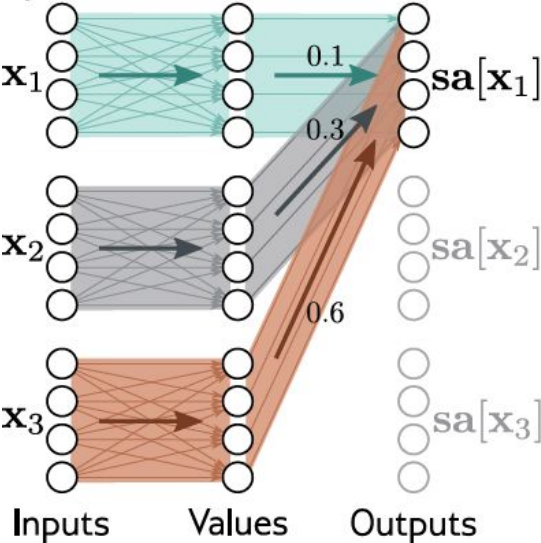
- N outputs are weighted sums of these values

$$\mathbf{sa}_n[\mathbf{x}_1, \dots, \mathbf{x}_N] = \sum_{m=1}^N a[\mathbf{x}_m, \mathbf{x}_n] \mathbf{v}_m.$$

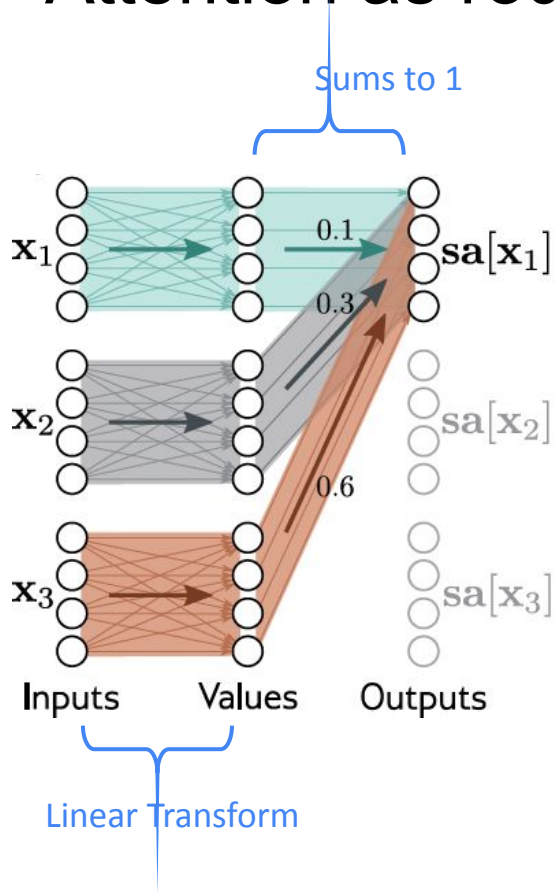
Scalar self-attention weights that represent how much attention the n^{th} token should pay to the m^{th} token

$a[\cdot, \mathbf{x}_n]$ are non-negative and sum to one

Attention as routing



Attention as routing



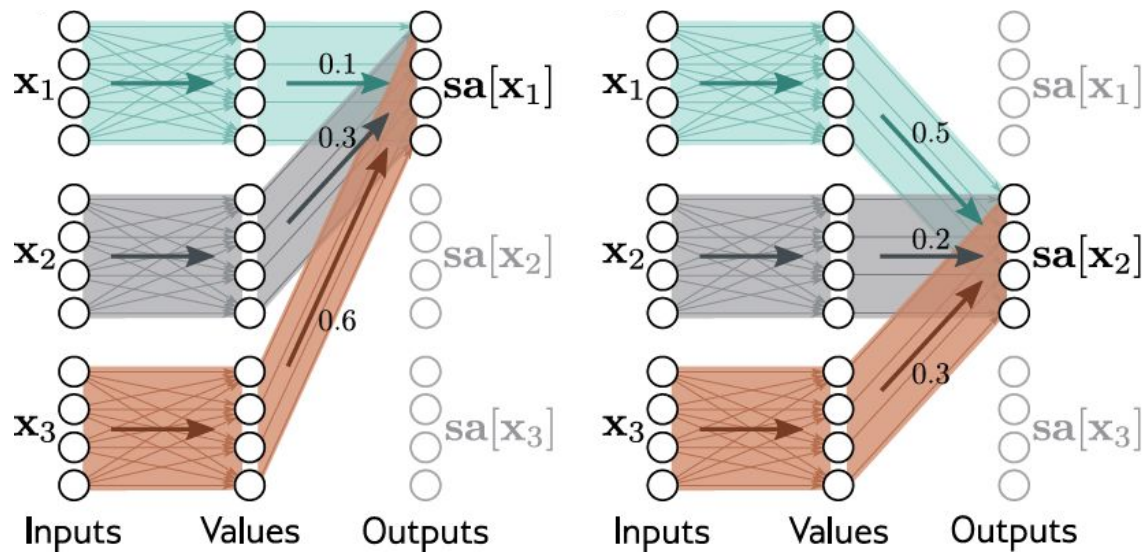
Here:

of inputs, $N = 3$

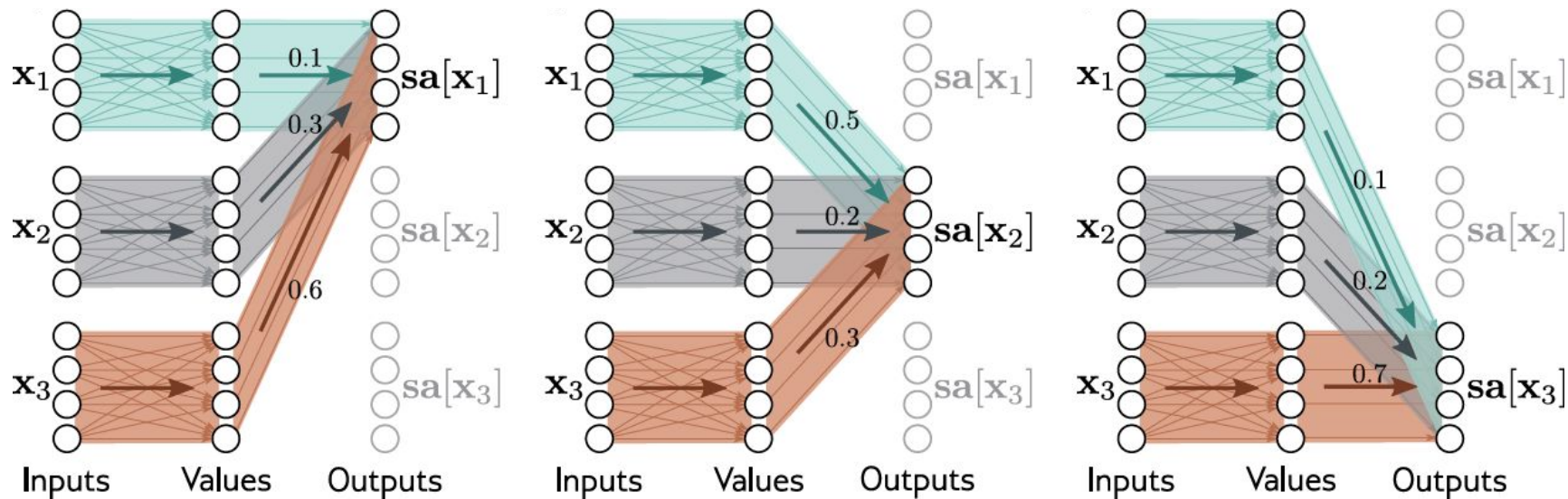
Dimension of each input, $D = 4$

We'll show how to calculate the self-attention weights shortly.

Attention as routing



Attention as routing



Attention weights

- Compute N “queries” and N “keys” from input

$$\mathbf{q}_n = \beta_q + \Omega_q \mathbf{x}_n$$

$$\mathbf{k}_n = \beta_k + \Omega_k \mathbf{x}_n;$$

- Calculate similarity and pass through softmax:

$$\begin{aligned} a[\mathbf{x}_n, \mathbf{x}_m] &= \text{softmax}_m [\text{sim}[\mathbf{k}_m \mathbf{q}_n]] \\ &= \frac{\exp [\text{sim}[\mathbf{k}_m \mathbf{q}_n]]}{\sum_{m'=1}^N \exp [\text{sim}[\mathbf{k}'_{m'} \mathbf{q}_n]]}; \end{aligned}$$

- Weights depend on the inputs themselves

Attention weights

- Compute N “queries” and N “keys” from input

$$\mathbf{q}_n = \beta_q + \Omega_q \mathbf{x}_n$$

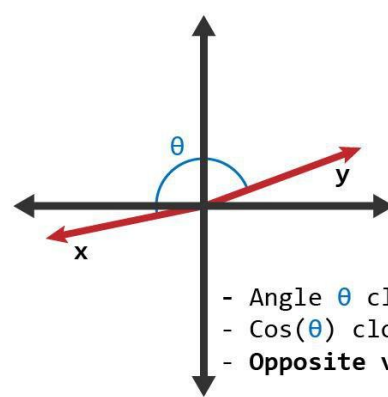
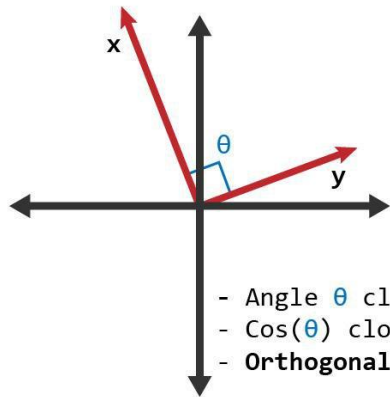
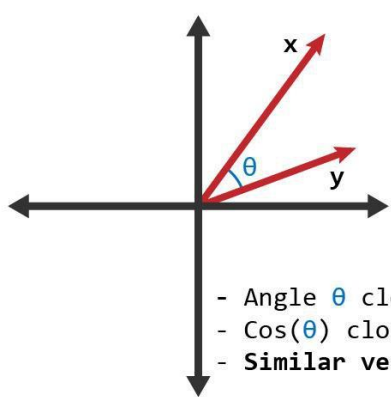
$$\mathbf{k}_n = \beta_k + \Omega_k \mathbf{x}_n;$$

- Take dot products and pass through softmax:

$$\begin{aligned} a[\mathbf{x}_n, \mathbf{x}_m] &= \text{softmax}_m [\mathbf{k}_m^T \mathbf{q}_n] \\ &= \frac{\exp [\mathbf{k}_m^T \mathbf{q}_n]}{\sum_{m'=1}^N \exp [\mathbf{k}_{m'}^T \mathbf{q}_n]} \end{aligned}$$

Dot product = measure of similarity

$$\mathbf{x}^T \mathbf{y} = |\mathbf{x}| |\mathbf{y}| \cos(\theta)$$



A drawback of the dot product as similarity measure is the magnitude of each vector influences the value. More rigorous to divide by magnitudes.

$$\text{Cosine Similarity: } \frac{\mathbf{x}^T \mathbf{y}}{|\mathbf{x}| |\mathbf{y}|} = \cos(\theta)$$

Motivation

Design neural network to encode and process text:

The restaurant refused to serve me a ham sandwich, because it only cooks vegetarian food. In the end, they just gave me two slices of bread. Their ambience was just as good as the food and service.

Conclusions:

- ✓ We need a model where parameters don't increase with input length, e.g.

$$\phi = \{\beta_v, \Omega_v, \beta_q, \Omega_q, \beta_k, \Omega_k\}$$

- ✓ There must be connections between the words.
- ✓ The strength of these connections will depend on the words themselves.

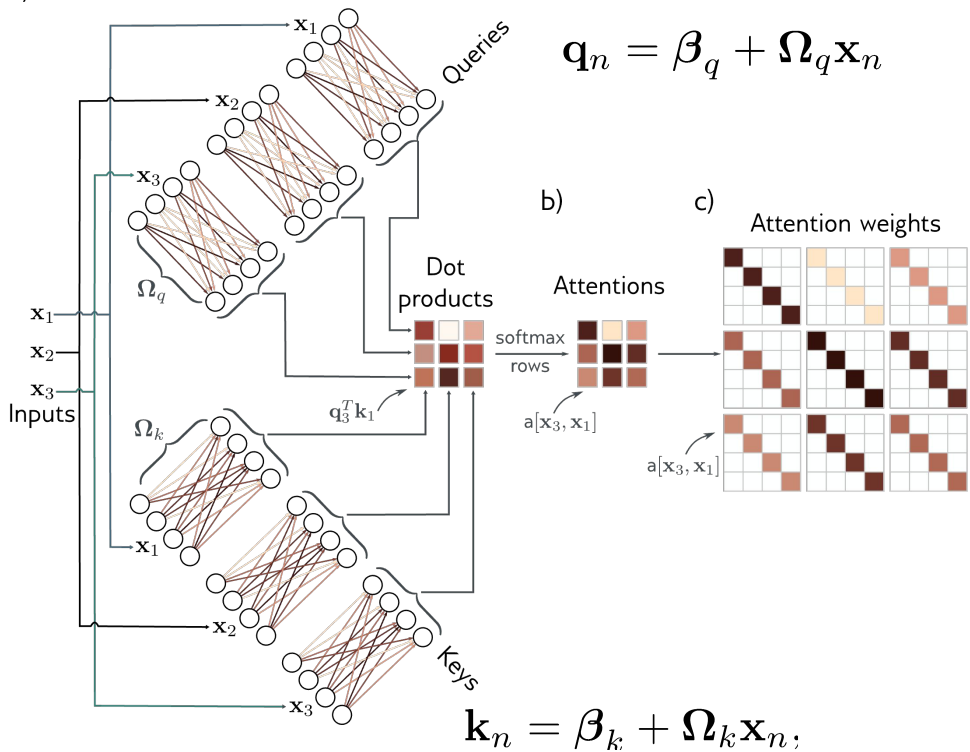
Ok, we defined *queries*, *keys* and *values*, but how are they used?

Transformers

- Motivation
- Dot-product self-attention
- **Applying Self-Attention**
- The Transformer Architecture
- Three Types of NLP Transformer Models

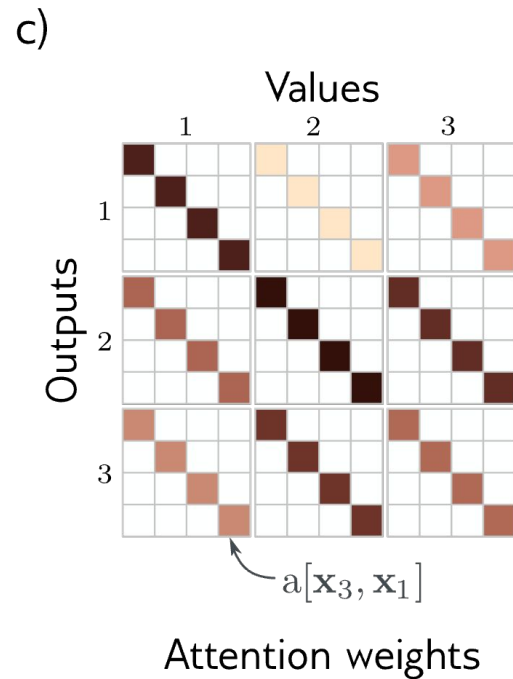
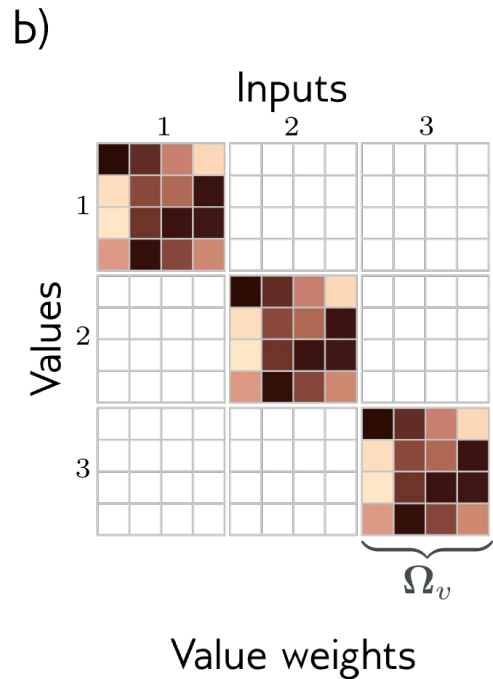
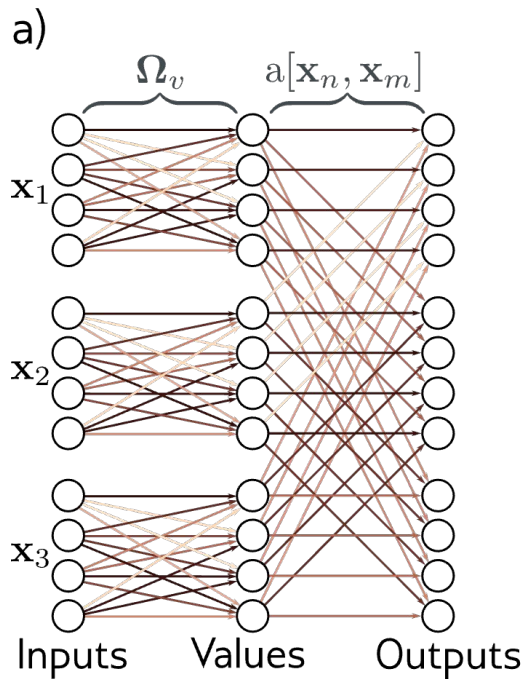
Computing Attention Weights

a)



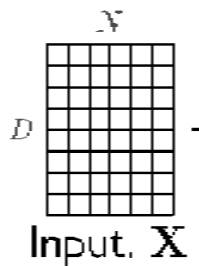
$$a[\mathbf{x}_n, \mathbf{x}_m] = \text{softmax}_m [\mathbf{k}_m^T \mathbf{q}_n]$$

Computing Values and Self-Attention Outputs as Sparse Matrix Ops



From Input Vector to Input Matrix

- Store N input vectors in matrix X



- Compute values, queries and keys:

$$\mathbf{V}[\mathbf{X}] = \beta_v \mathbf{1}^T + \Omega_v \mathbf{X}$$

$$\mathbf{Q}[\mathbf{X}] = \beta_q \mathbf{1}^T + \Omega_q \mathbf{X}$$

$$\mathbf{K}[\mathbf{X}] = \beta_k \mathbf{1}^T + \Omega_k \mathbf{X},$$

- Combine self-attentions

$$\text{Sa}[\mathbf{X}] = \mathbf{V}[\mathbf{X}] \cdot \text{Softmax} \left[\mathbf{K}[\mathbf{X}]^T \mathbf{Q}[\mathbf{X}] \right] = \mathbf{V} \cdot \text{Softmax}[\mathbf{K}^T \mathbf{Q}]$$

Scaled Dot Product Self-Attention

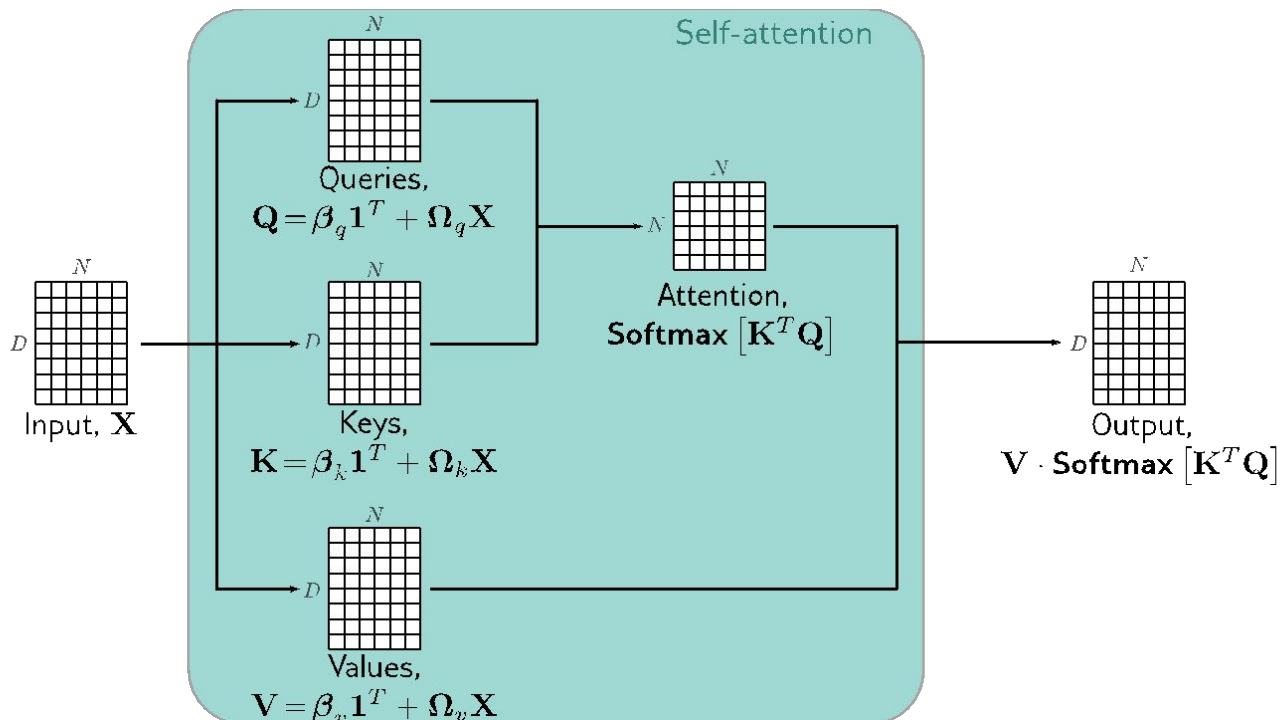
- To avoid the case where a large value dominates the softmax in

$$\mathbf{Sa}[\mathbf{X}] = \mathbf{V} \cdot \mathbf{Softmax}[\mathbf{K}^T \mathbf{Q}]$$

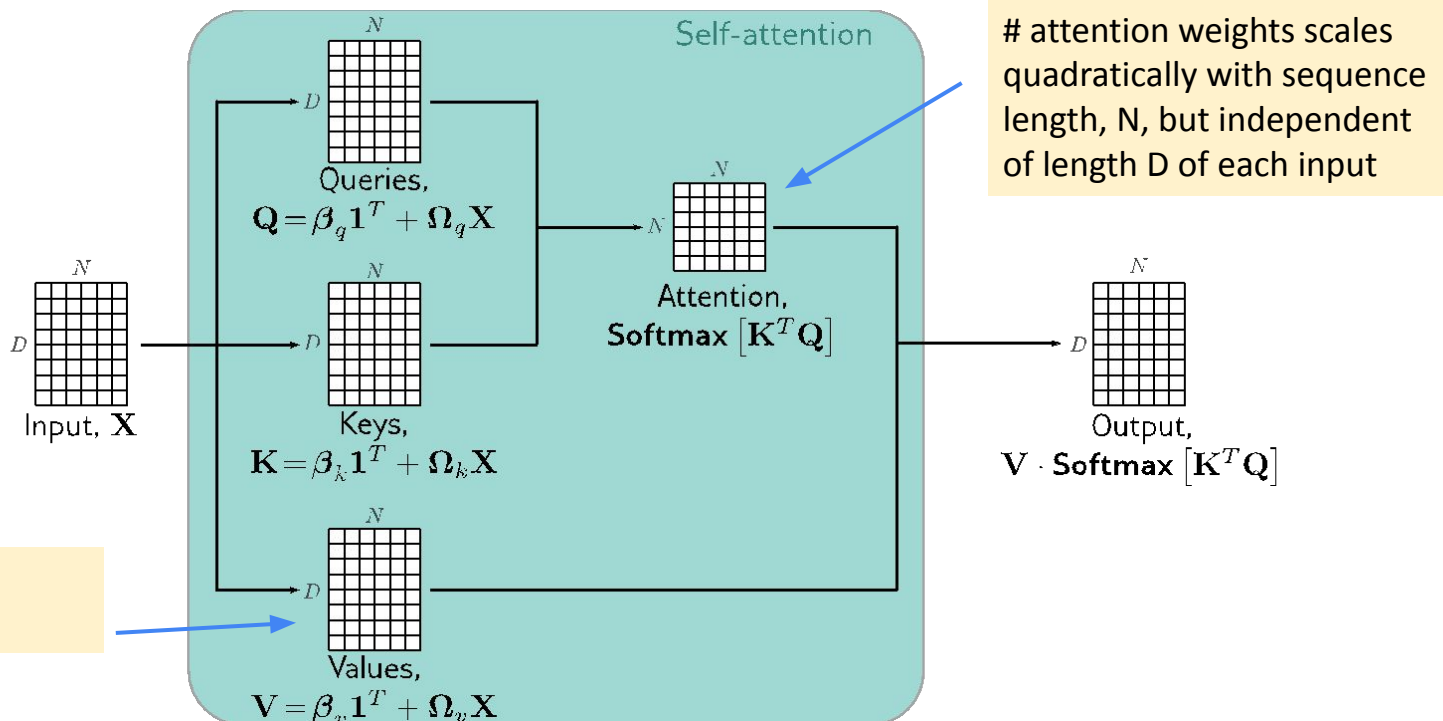
- you can scale the dot product by the square root of the dimension of the query

$$\mathbf{Sa}[\mathbf{X}] = \mathbf{V} \cdot \mathbf{Softmax} \left[\frac{\mathbf{K}^T \mathbf{Q}}{\sqrt{D_q}} \right]$$

Put it all together in matrix form

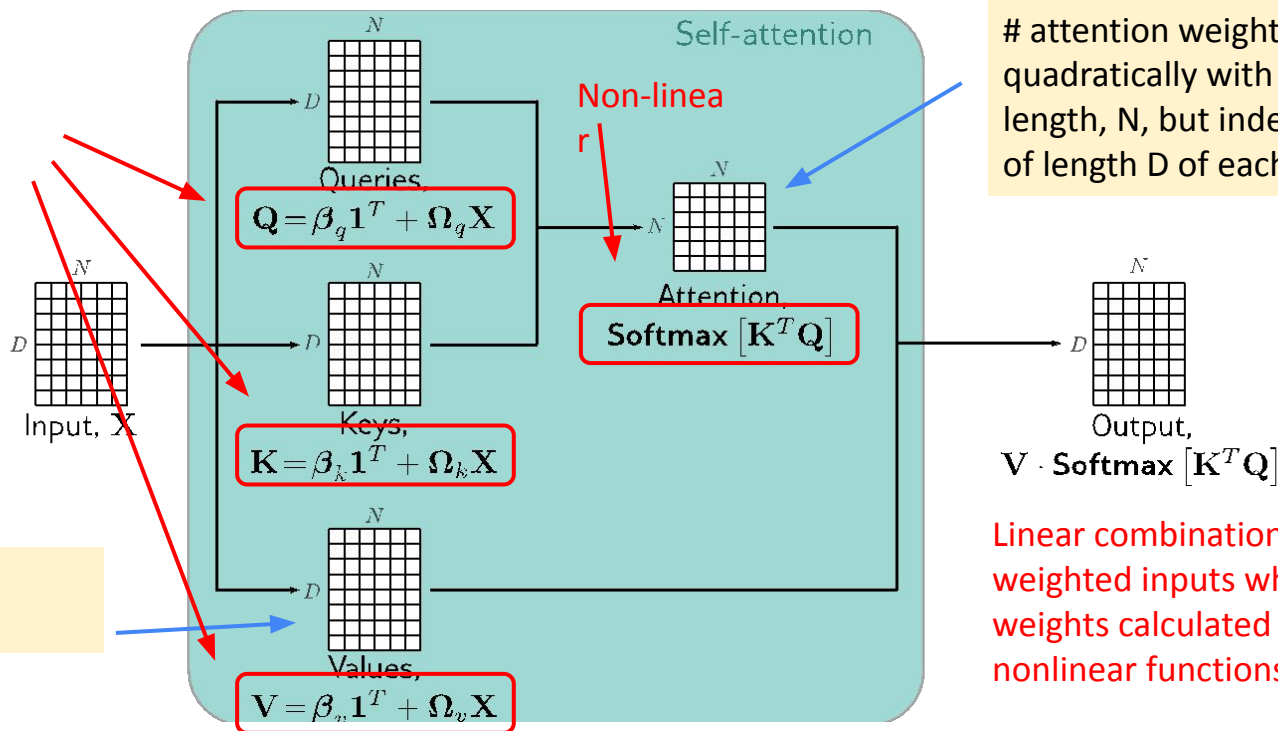


Put it all together in matrix form



Put it all together in matrix form

Linear
&
Can be
calculated in
parallel



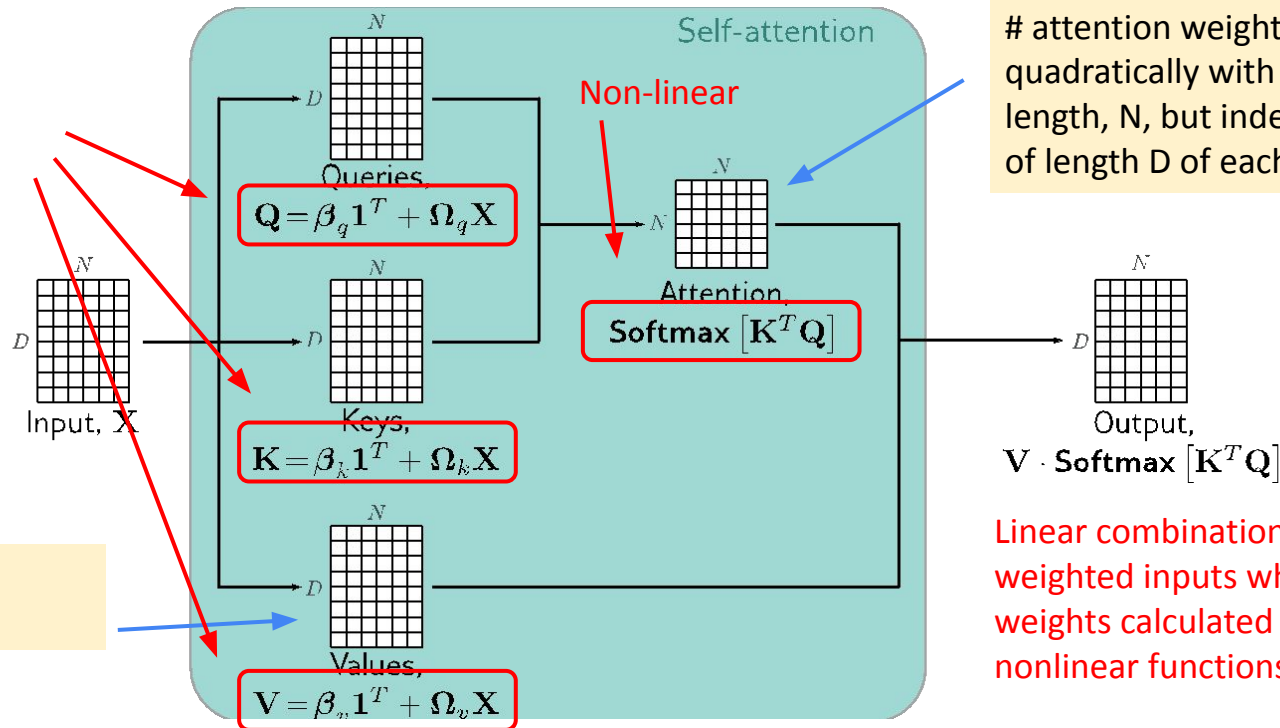
Scales linearly with
sequence length, N

attention weights scales
quadratically with sequence
length, N , but independent
of length D of each input

Linear combination of
weighted inputs where
weights calculated from
nonlinear functions

Hypernetwork – 1 branch calculates weights of other branch

Linear
&
Can be
calculated in
parallel

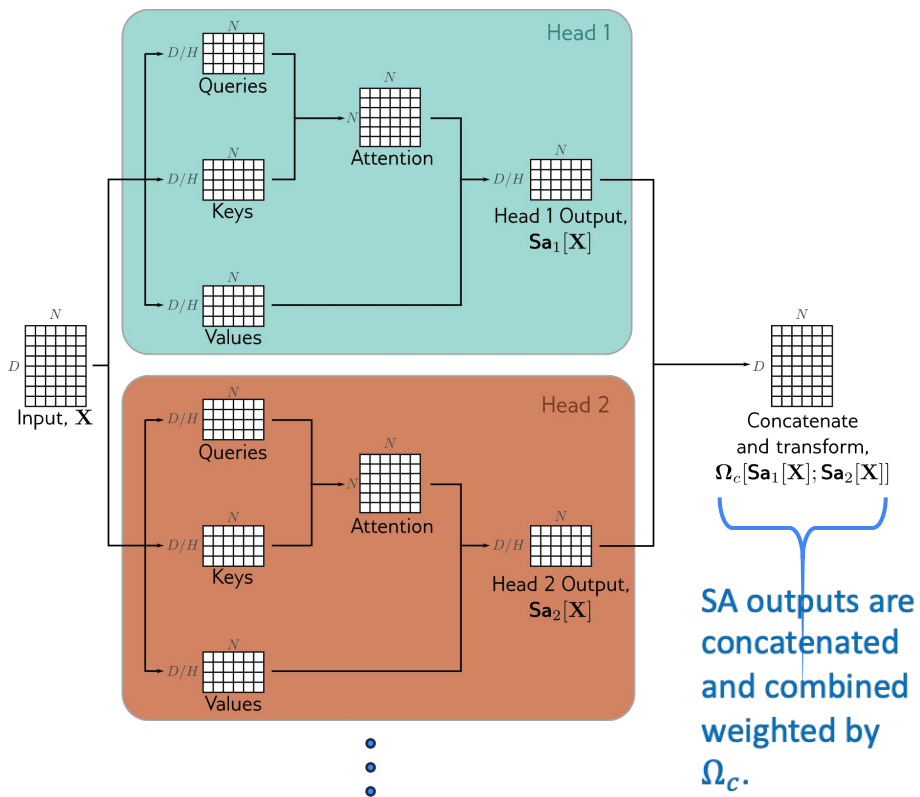


attention weights scales quadratically with sequence length, N , but independent of length D of each input

Scales linearly with sequence length, N

Linear combination of weighted inputs where weights calculated from nonlinear functions

Multi-Head Self Attention



- Multiple self-attention heads are usually applied in parallel
- “allows model to jointly attend to info from different representation subspaces at different positions”
- Original paper used 8 heads
- All can be executed in parallel

Equivariance to Word Order

A function $f[x]$ is **equivariant** to a transformation $t[]$ if: $f[t[x]] = t[f[x]]$

Self-attention is *equivariant* to permuting word order. Just a bag of words.

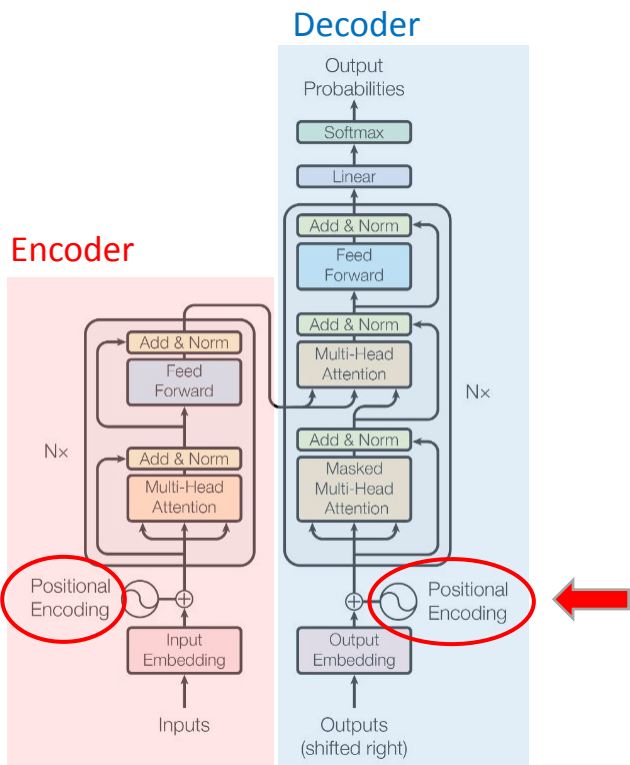
But word order is important in language:

The man ate the fish

vs.

The fish ate the man

Solution: Position Encoding



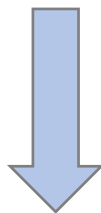
Idea is to somehow encode *absolute* or *relative* position in the inputs

Fourier features used in neural fields are a version of this idea.

Absolute Position encoding

Alternatively, could be added to each layer

$$\mathbf{Sa}[\mathbf{X}] = \mathbf{V} \cdot \mathbf{Softmax}[\mathbf{K}^T \mathbf{Q}]$$



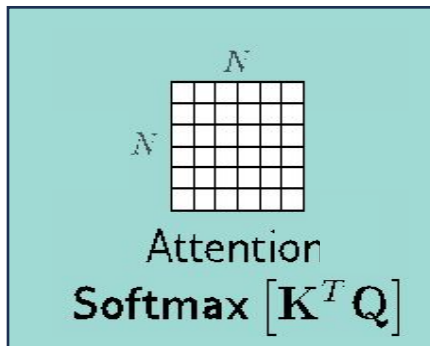
$$\mathbf{Sa}[\mathbf{X}] = (\mathbf{V} + \mathbf{\Pi}) \cdot \mathbf{Softmax}[(\mathbf{K} + \mathbf{\Pi})^T (\mathbf{Q} + \mathbf{\Pi})]$$

Relative Position Encoding

Absolute position of a word is less important than relative position between inputs

*The panda **eats** shoots and leaves*

Abs Pos:	0	1	2	3	4	5
Rel Pos:	-2	-1	0	1	2	3



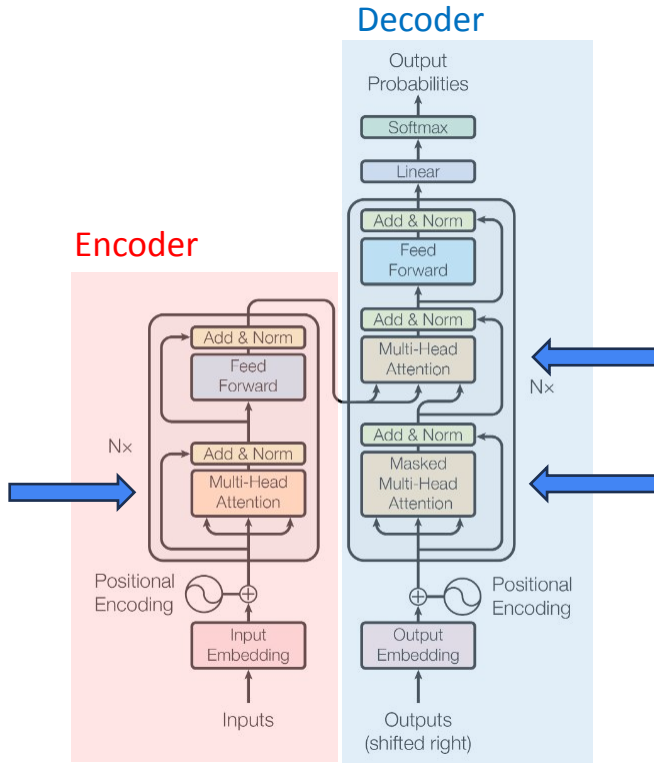
Each element of the attention matrix corresponds to an offset between query position a and key position b

Learn a parameter $\pi_{a,b}$ for each offset and modify $\text{Attention}[a,b]$ in some way.

Transformers

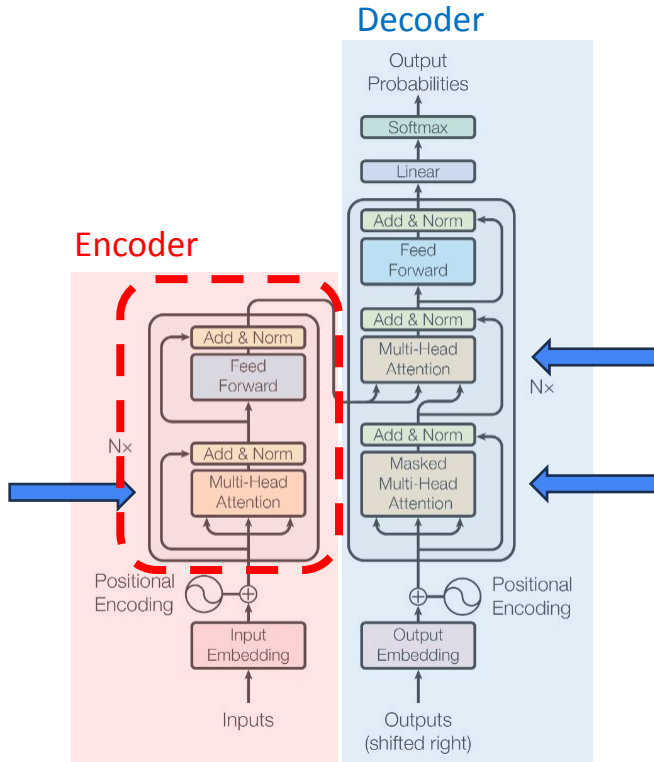
- Motivation
- Dot-product self-attention
- Applying Self-Attention
- **The Transformer Architecture**
- Three Types of NLP Transformer Models

Transformers



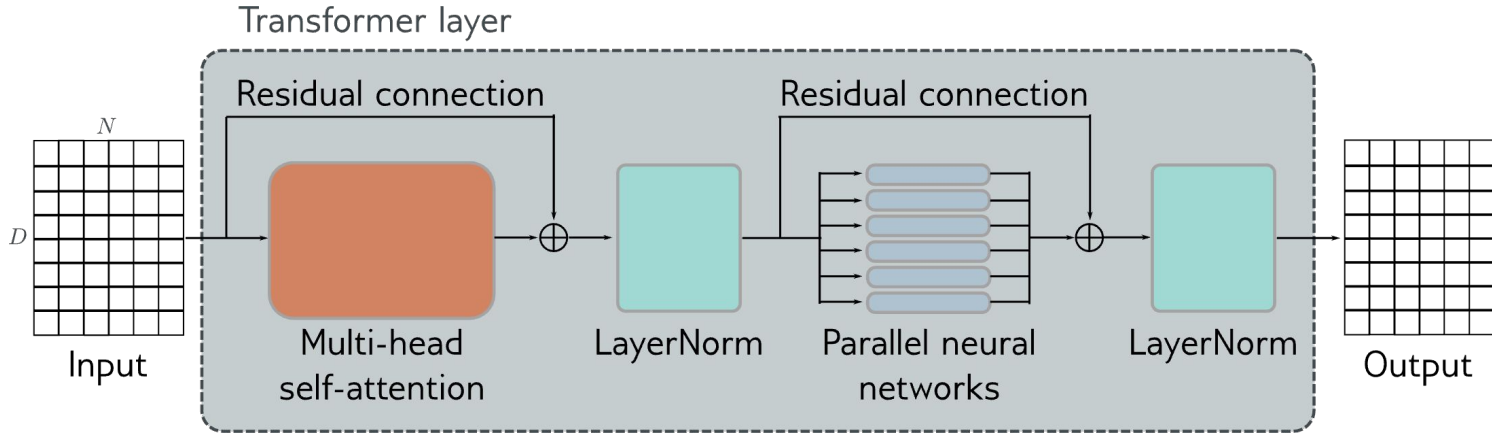
- *Multi-headed Self Attention* is just one component of the transformer architecture

Transformers



- *Multi-headed Self Attention* is just one component of the transformer architecture
- Let's look at a transformer *block* (or *layer*) from the encoder

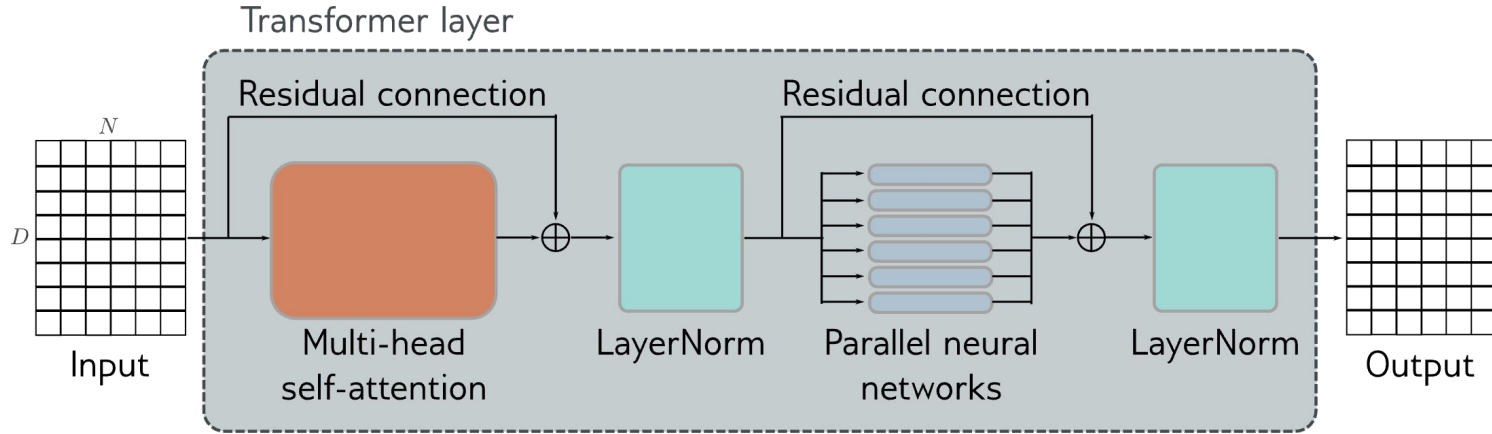
Transformer Layer -- Complete



- Adds a 2-layer MLP
- Adds residual connections around multi-head self-attentions and the parallel MLPs
- Adds LayerNorm, which normalizes across all the N input samples

Transform Layer	
\mathbf{X}	$\leftarrow \mathbf{X} + \text{MhSa}[\mathbf{X}]$
\mathbf{X}	$\leftarrow \text{LayerNorm}[\mathbf{X}]$
\mathbf{x}_n	$\leftarrow \mathbf{x}_n + \text{mlp}[\mathbf{x}_n]$
\mathbf{X}	$\leftarrow \text{LayerNorm}[\mathbf{X}]$,

Transformer Layer -- MLP

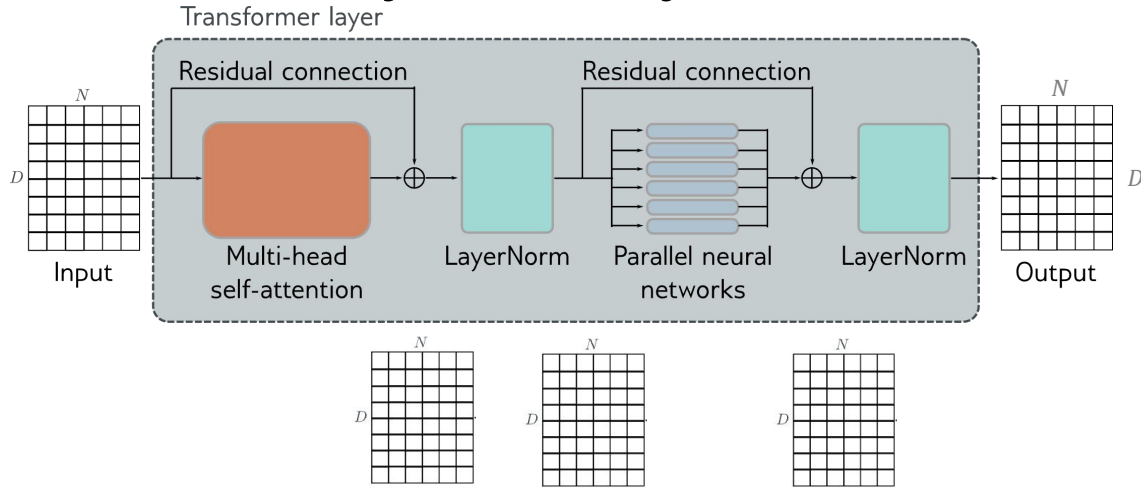


- Adds 2-layer MLP

- Same network (same weights) operates independently on each word
- Learn more complex representations and expand model capacity

Linear _{$D \times 4D$} \square ReLU(.) \square Linear _{$4D \times D$}

Transformer Layer -- LayerNorm



- Normalize across same layer
- Learned gain and offset

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

Calculated column-wise

MLP Example

```
batch, sentence_length, embedding_dim = 20, 5, 10
embedding = torch.randn(batch, sentence_length,
embedding_dim)
layer_norm = nn.LayerNorm(embedding_dim)
```

Activate module

```
layer_norm(embedding)
```

<https://pytorch.org/docs/stable/generated/torch.nn.LayerNorm.html>

Transformers

- Motivation
- Dot-product self-attention
- Applying Self-Attention
- The Transformer Architecture
- Three Types of NLP Transformer Models

Transformers

- Motivation
- Dot-product self-attention
- Applying Self-Attention
- The Transformer Architecture
- Three Types of NLP Transformer Models
 - Encoder
 - Decoder
 - Encoder-Decoder

Transformers

- Motivation
- Dot-product self-attention
- Applying Self-Attention
- The Transformer Architecture
- Three Types of NLP Transformer Models
 - Encoder
 - Decoder
 - Encoder-Decoder

3 Types of Transformer Models

1. *Encoder* – transforms text embeddings into representations that support variety of tasks (e.g. sentiment analysis, classification)
 - ❖ Model Example: BERT
2. *Decoder* – predicts the next token to continue the input text (e.g. ChatGPT, AI assistants)
 - ❖ Model Example: GPT4, GPT4
3. *Encoder-Decoder* – used in sequence-to-sequence tasks, where one text string is converted to another (e.g. machine translation)

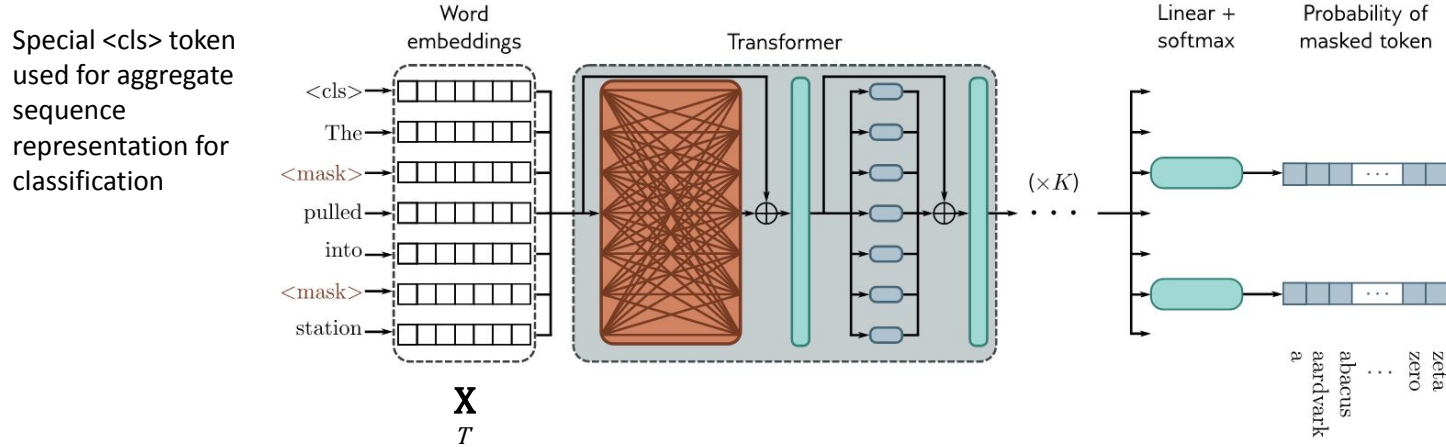
Encoder Model Example: BERT (2019)

Bidirectional Encoder Representations from Transformers

- Hyperparameters
 - 30,000 token vocabulary
 - 1024-dimensional word embeddings
 - 24x transformer layers
 - 16 heads in self-attention mechanism
 - 4096 hidden units in middle of MLP
- ~340 million parameters
- Pre-trained in a self-supervised manner,
- then can be adapted to task with one additional layer and fine-tuned

This is a popular model to fine-tune for specialized tasks.

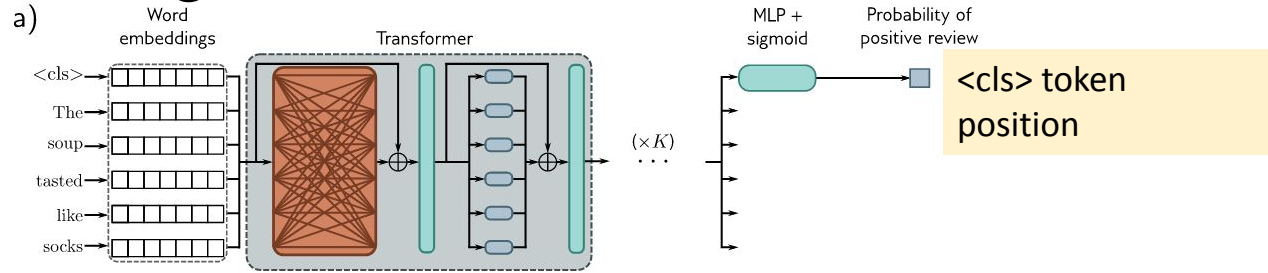
Encoder Pre-Training



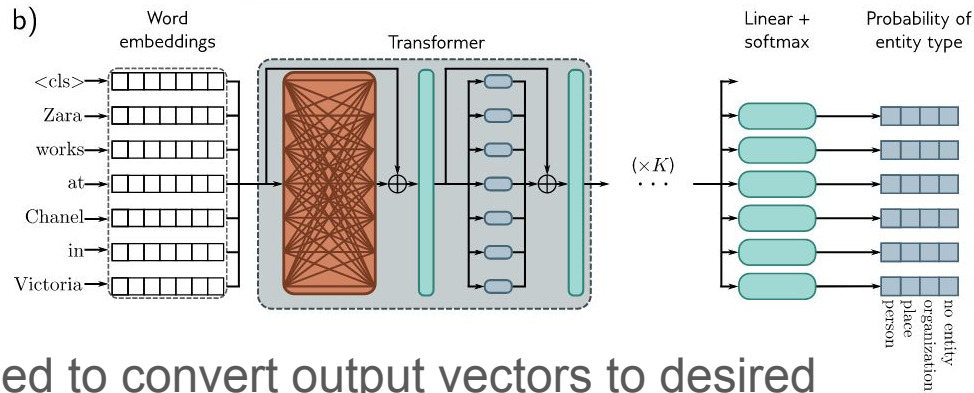
- A small percentage of input embedding replaced with a generic <mask> token
- Predict missing token from output embeddings
- Added linear layer and softmax to generate probabilities over vocabulary
- Trained on BooksCorpus (800M words) and English Wikipedia (2.5B words)

Encoder Fine-Tuning

Sentiment Analysis



Named Entity Recognition (NER)



- Extra layer(s) appended to convert output vectors to desired output format
- 3rd Example: Text span prediction -- predict start and end location of answer to a question in passage of Wikipedia, see <https://rajpurkar.github.io/SQuAD-explorer/>

Transformers

- Motivation
- Dot-product self-attention
- Applying Self-Attention
- The Transformer Architecture
- Three Types of NLP Transformer Models
 - Encoder
 - Decoder
 - Encoder-Decoder

Decoder Model Example: GPT3 (2020)

Generative Pre-trained Transformer

- One purpose: generate the next token in a sequence
- By constructing an autoregressive model

We saw this interface before,
but better internals now.

Decoder Model Example: GPT3 (2020)

Generative Pre-trained Transformer

- One purpose: *generate the next token in a sequence*
- By constructing an autoregressive model
- Factors the probability of the sentence:

$$\begin{aligned} \Pr(\textit{Learning deep learning is fun}) = & \\ & \Pr(\textit{Learning}) \times \Pr(\textit{deep} \mid \textit{learning}) \times \\ & \Pr(\textit{learning} \mid \textit{Learning deep}) \times \\ & \Pr(\textit{is} \mid \textit{Learning deep learning}) \times \\ & \Pr(\textit{fun} \mid \textit{Learning deep learning is}) \end{aligned}$$

Decoder Model Example: GPT3 (2020)

Generative Pre-trained Transformer

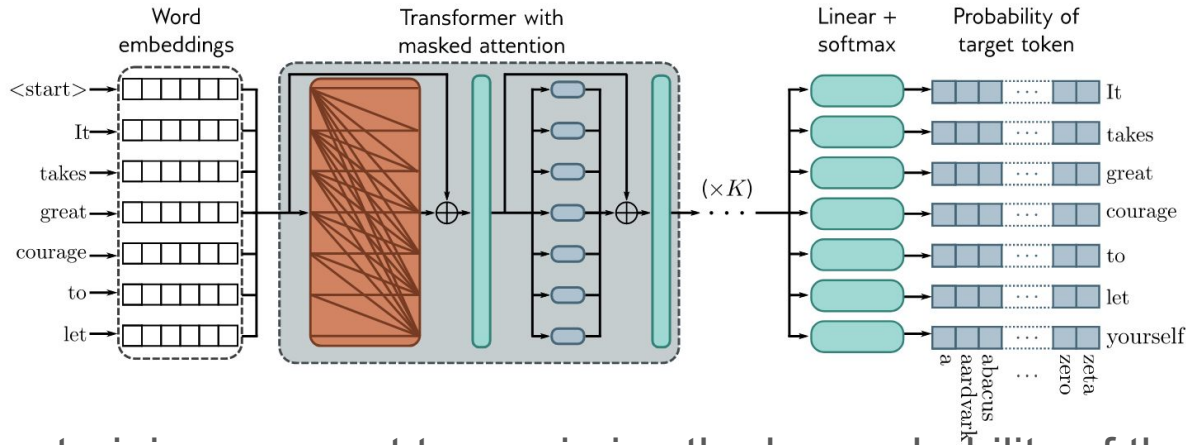
- One purpose: *generate the next token in a sequence*
- By constructing an autoregressive model
- Factors the probability of the sentence:

$$\begin{aligned} \Pr(\textit{Learning deep learning is fun}) = & \\ & \Pr(\textit{Learning}) \times \Pr(\textit{deep} \mid \textit{learning}) \times \\ & \Pr(\textit{learning} \mid \textit{Learning deep}) \times \\ & \Pr(\textit{is} \mid \textit{Learning deep learning}) \times \\ & \Pr(\textit{fun} \mid \textit{Learning deep learning is}) \end{aligned}$$

- More formally: Autoregressive model_N

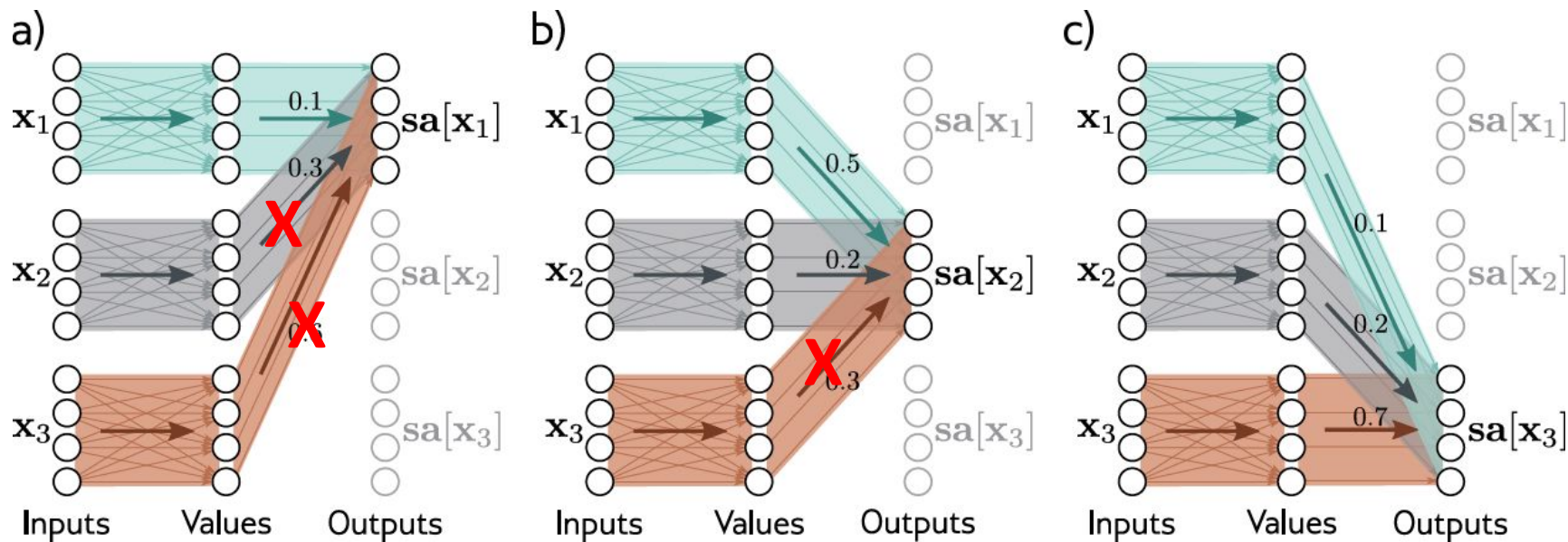
$$\Pr(t_1, t_2, \dots, t_N) = \Pr(t_1) \prod_{n=2}^N \Pr(t_n \mid t_1, t_2, \dots, t_{n-1})$$

Decoder: *Masked Self-Attention*



- During training we want to maximize the log probability of the input text under the autoregressive model
- We want to **make sure the model doesn't "cheat"** during training by looking ahead at the next token
- Hence we mask the self attention weights corresponding to current and right context to negative infinity

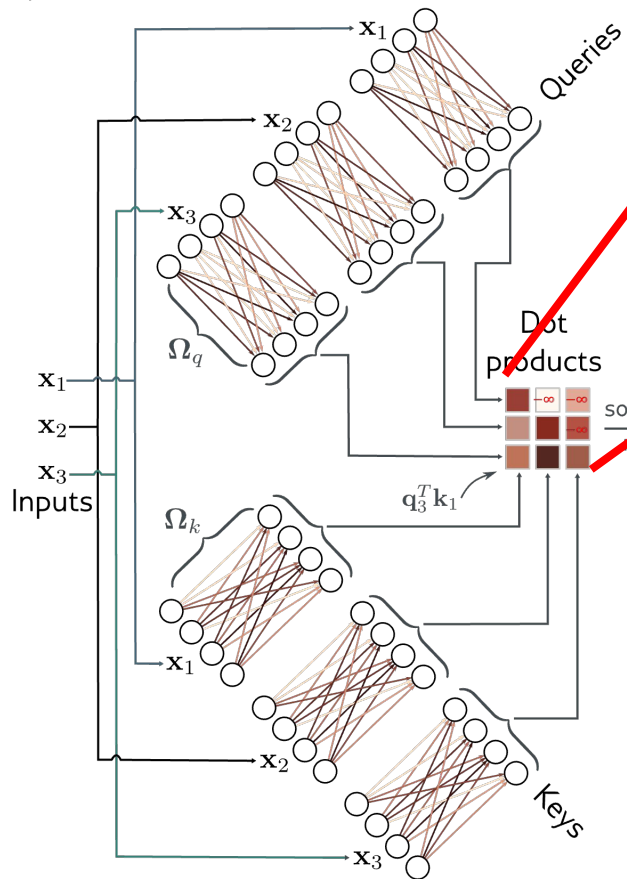
Masked Self-Attention



Mask right context self-attention weights to zero

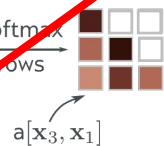
Masked Self-Attention

a)



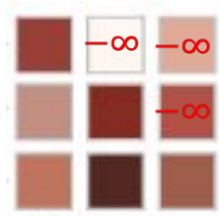
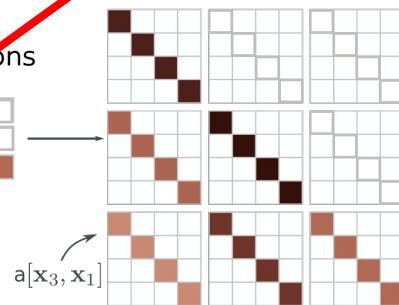
b)

Attentions

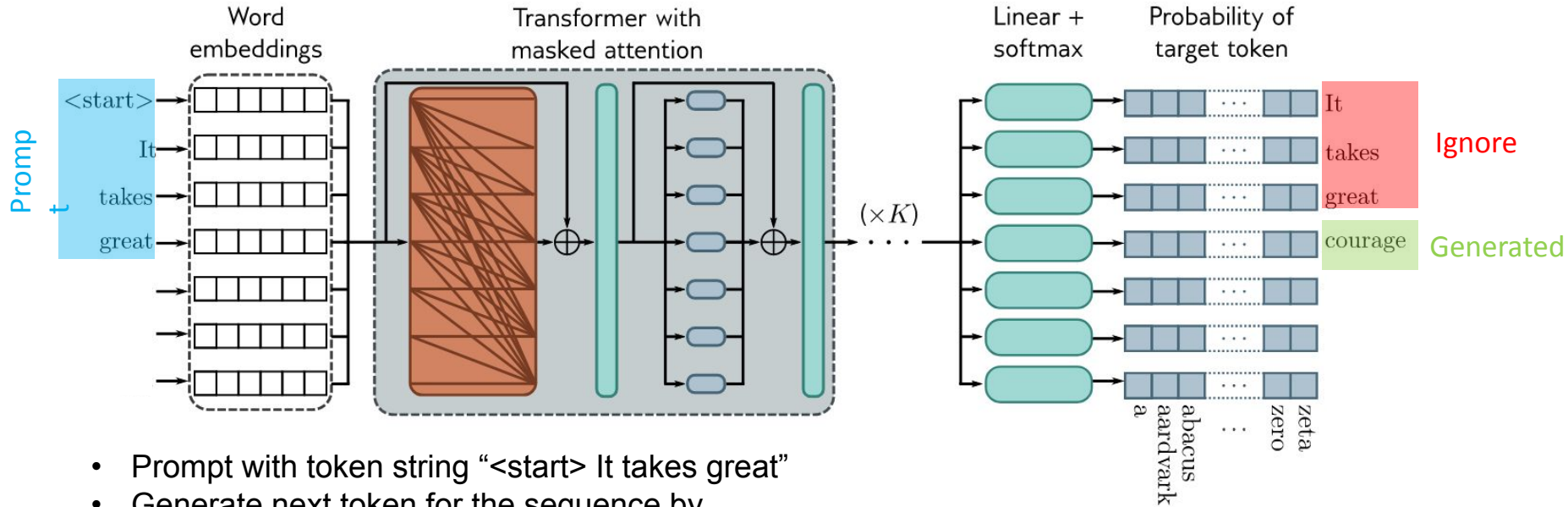


c)

Attention weights

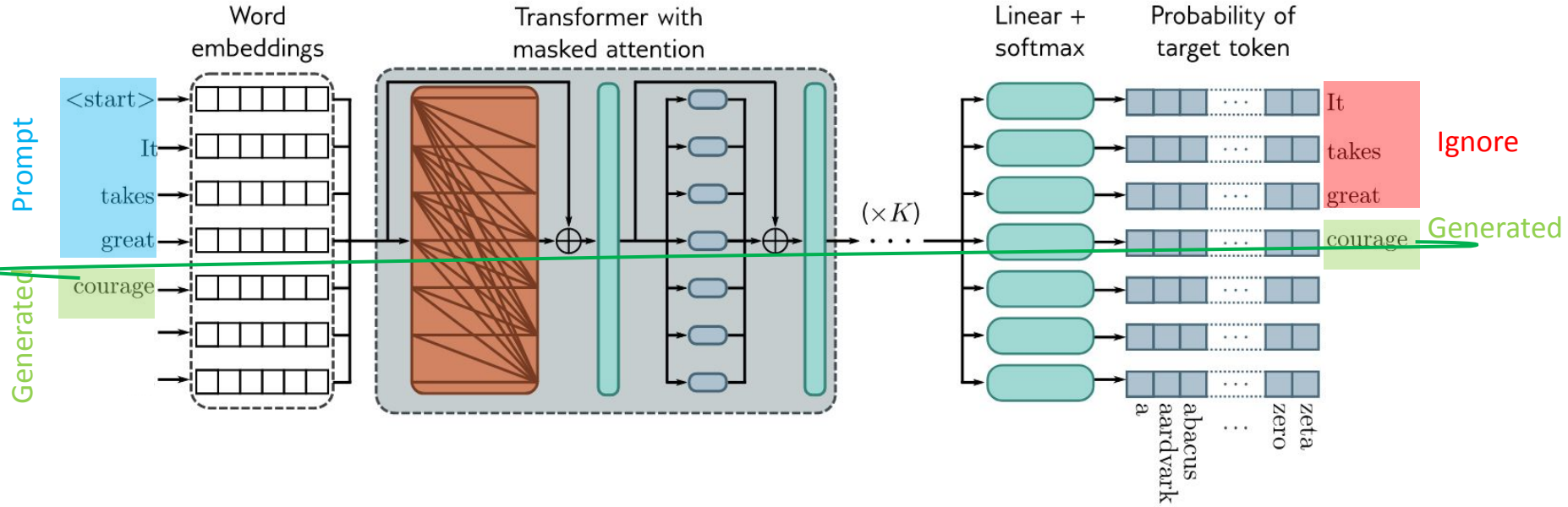


Decoder: Text Generation (Generative AI)



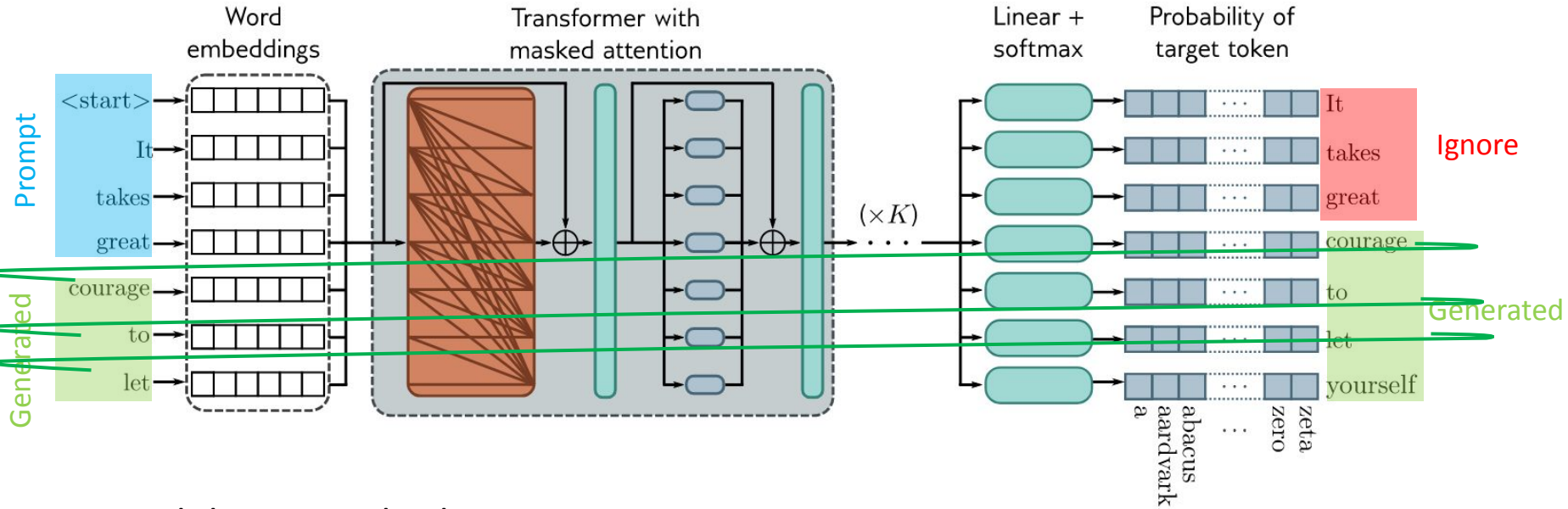
- Prompt with token string "`<start>` It takes great"
- Generate next token for the sequence by
 - picking most likely token
 - sample from the probability distribution
 - alternative top-k sampling to avoid picking from the long tail
 - beam search – select the most likely sentence rather than greedily pick

Decoder: Text Generation (Generative AI)



- Feed the output back into input

Decoder: Text Generation (Generative AI)



- Feed the output back into input

Technical Details

	BERT	GPT3
Model Architecture	Encoder	Decoder
Embedding Size	1024	12,288
Vocabulary	30K tokens	
Sequence Length		2048
# Heads	16	96
# Layers	24	96
Q,K,V dimensions	64	128
Training set size	3.3B tokens	300B+ tokens
# Parameters	340M	175B

Again, BERT is viable and available for your projects.

Transformers

- Motivation
- Dot-product self-attention
- Applying Self-Attention
- The Transformer Architecture
- Three Types of NLP Transformer Models
 - Encoder
 - Decoder
 - Encoder-Decoder

Encoder-Decoder Model

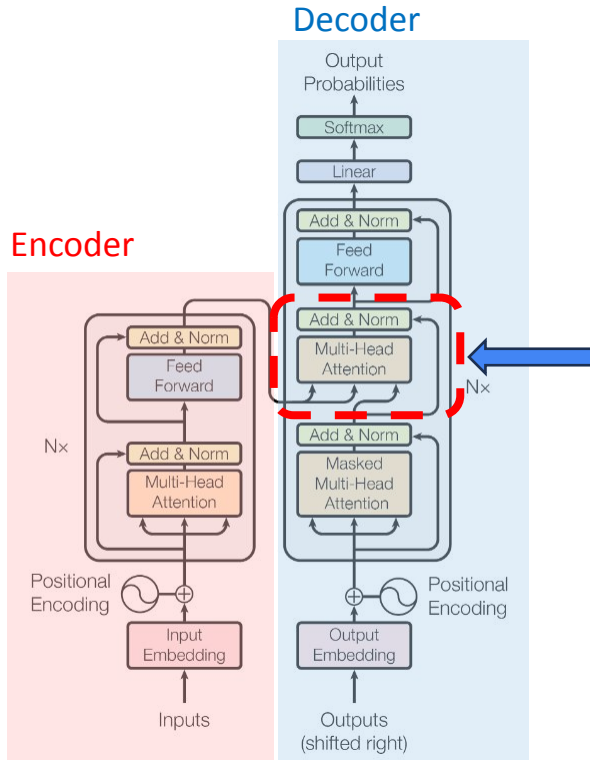
- Used for machine translation, which is a sequence-to-sequence task



Decoder only continues input sequences.

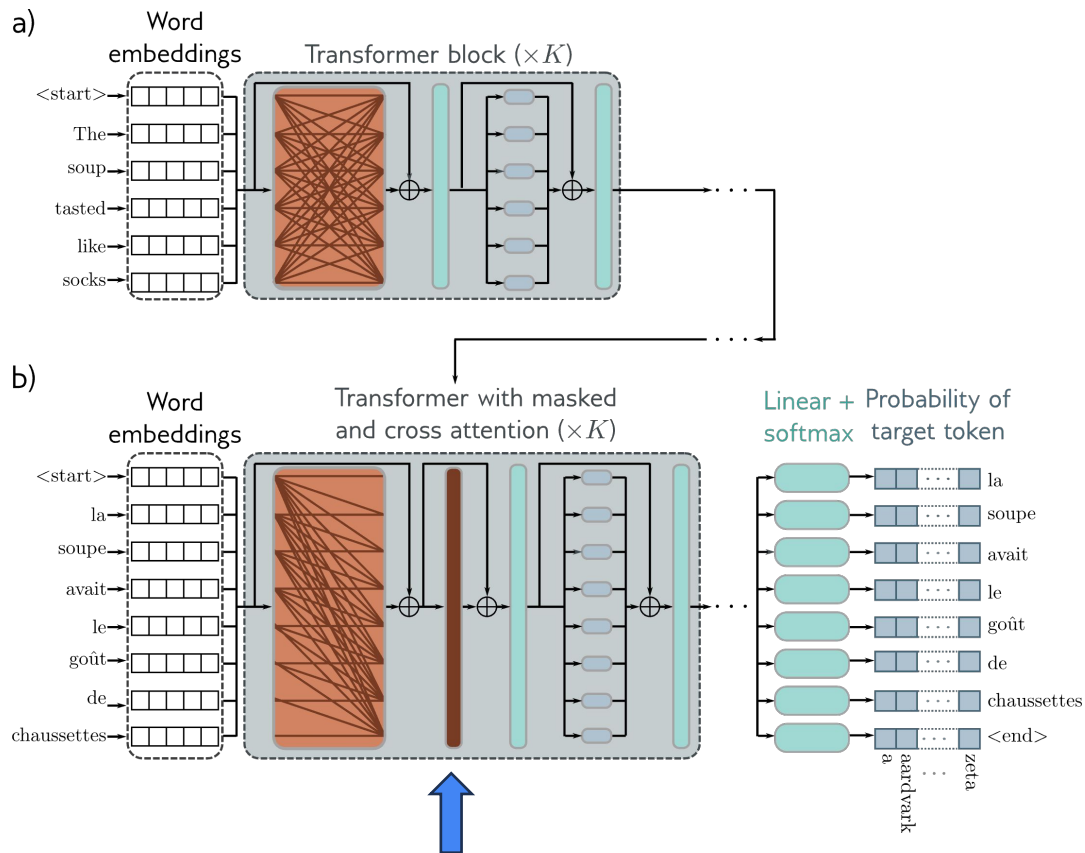
Encoder-decoder produces new sequences based on input sequences.

Encoder Decoder Model



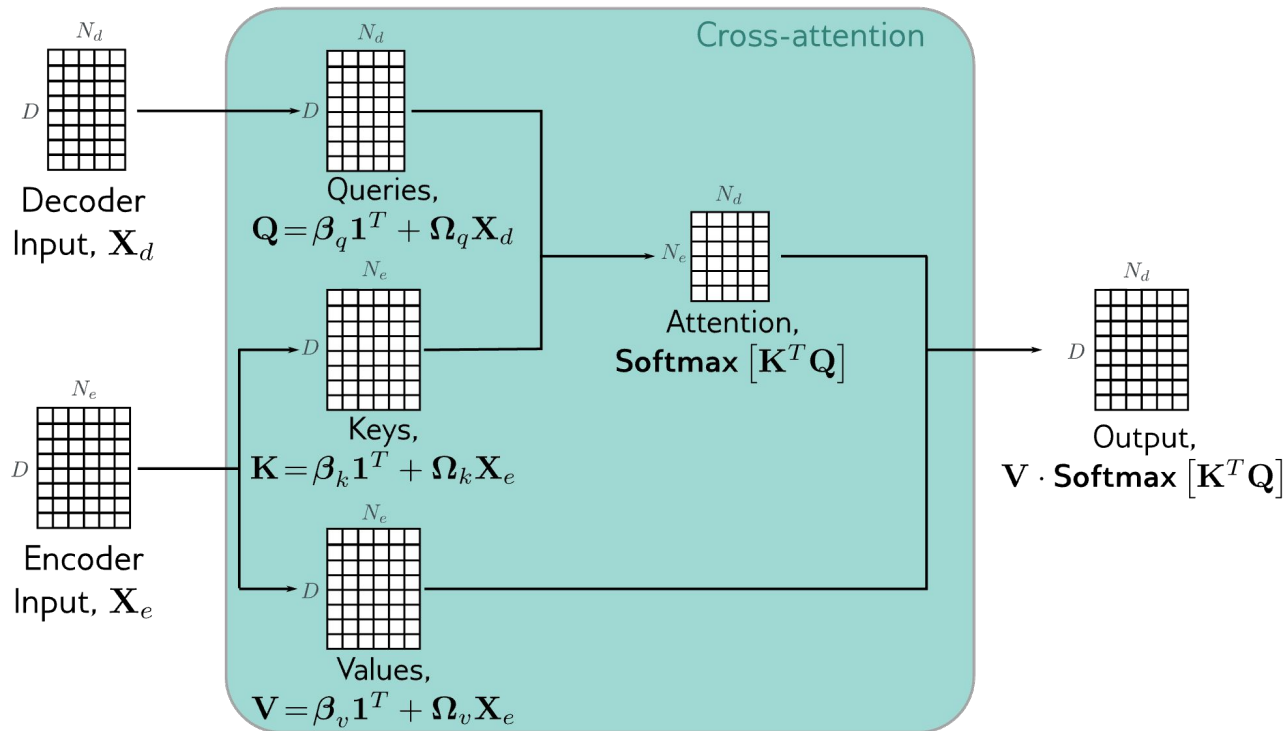
- The transformer layer in the decoder of the encoder-decoder model has an extra stage
- Attends to the input of the encoder with **cross attention** using Keys and Values from the output of the encoder
- Shown here on original diagram from “Attention is all you need” paper

Encoder Decoder Model



- Same view per UDL book

Cross-Attention



Keys and Values come from the last stage of the encoder

Feedback?

