

Attention & Transformers for Text Classification

(Vaswani et. al. 2017)

(Slides/Figures by Jay Alammam,
original papers)

Attention

Sentence Representation

You can't cram the meaning of a whole %&!\$# sentence into a single \$&!#* vector!



But what if we could use multiple vectors, based on the length of the sentence.

this is an example → 

this is an example → 

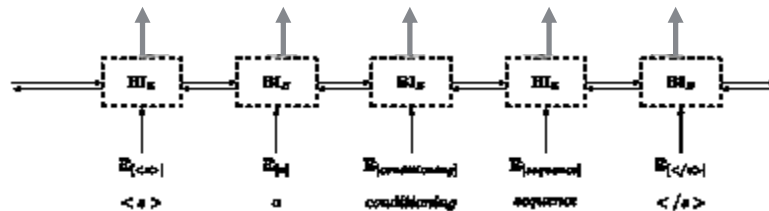
Intuition

- Encoding a single vector is too restrictive
Instead of the encoder producing a single vector for the sentence, it will produce one vector **for each word**.
- But, we still need 1 vector. Multiple vectors \rightarrow Single vector
Sum/Avg operators are good. But give equal importance to each input
- We dynamically decide which input is more/less important for a task. Create a weighted sum to reflect this variation
- Attention:
 - query (q): decides how much importance to give each input
 - attention weights (α_i): importance of each input i (normalized to 1)
 - unnormalized attention weights ($\bar{\alpha}_i$): intermediate step to compute α_i
 - attended summary: weighted average of input with α as weights

Encoder for One Vector/Word

$$\mathbf{c}_{1:n} = \text{ENC}(\mathbf{x}_{1:n}) = \text{biRNN}^*(\mathbf{x}_{1:n})$$

Encoder

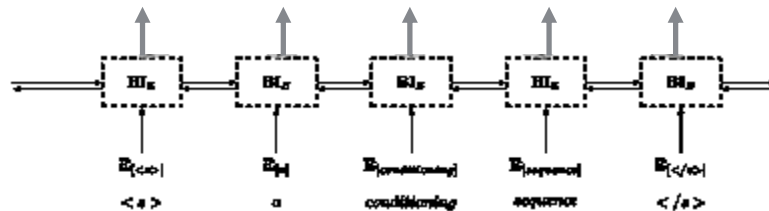


Encoder for One Vector/Word

$$c_{1:n} = \text{ENC}(x_{1:n}) = \text{biRNN}^*(x_{1:n})$$

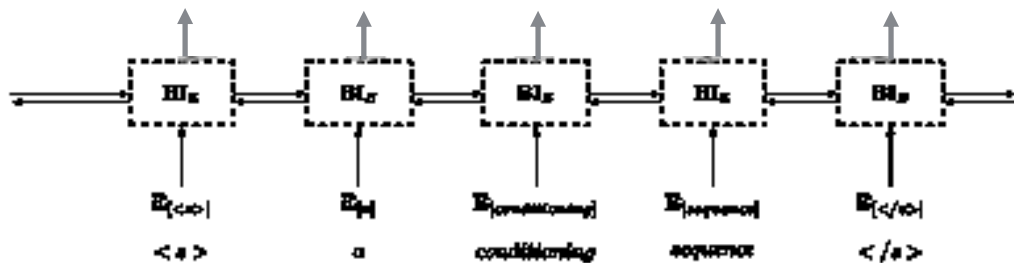
but how do we feed
this sequence
to the MLP?

Encoder



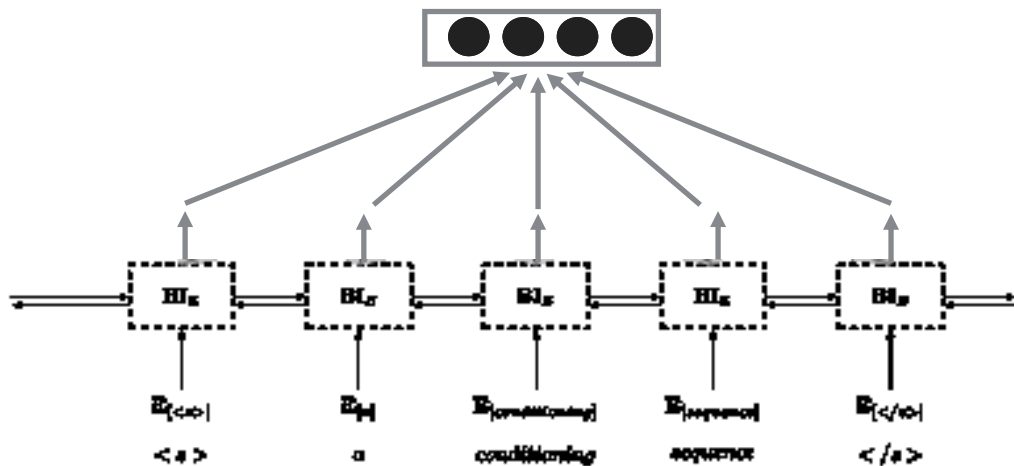
Encoder for One Vector/Word

we can combine the different outputs
into a single vector (attended summary)



Multiple Vectors \rightarrow Single Vector

we can combine the different outputs
into a single vector (attended summary)



Attention



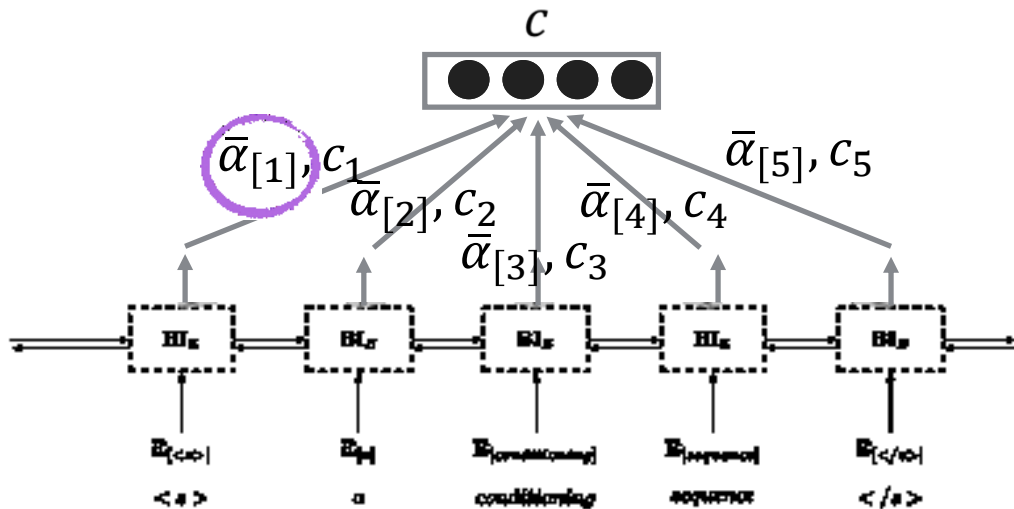
$$c = \sum_{i=1}^n \alpha_{[i]} \cdot c_i$$

Attention

$$\bigcirc = \text{softmax}(\bar{\alpha}_{[1]}, \dots, \bar{\alpha}_{[n]})$$

$$c = \sum_{i=1}^n \bigcirc c_i$$

Encoder



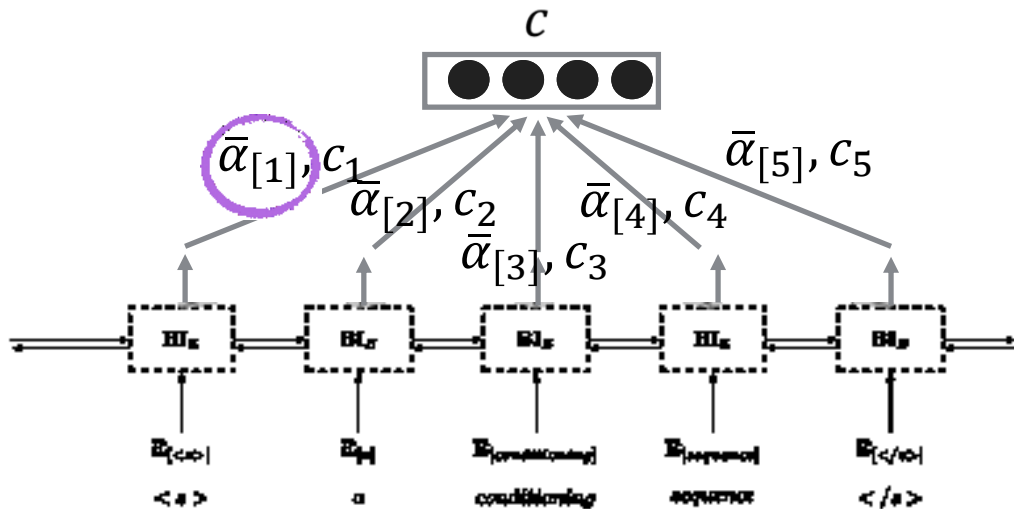
Attention

$$\alpha = \text{softmax}(\bar{\alpha}_{[1]}, \dots, \bar{\alpha}_{[n]})$$

$$\bigcirc = \text{MLP}^{\text{att}}(q, c_i)$$

$$c = \sum_{i=1}^n \bigcirc c_i$$

query



Encoder

Attention

$$y = \text{softmax}(MLP^{out}(c))$$

$$c = \sum_{i=1}^n \alpha_{[i]} \cdot c_i$$

$$c_{1:n} = biLSTM_{enc}(x_{1:n})$$

$$\alpha = \text{softmax}(\bar{\alpha}_{[1]}, \dots, \bar{\alpha}_{[n]})$$

$$\bar{\alpha}_{[i]} = MLP^{att}(q, c_i)$$

Two things missing. What are they?

Attention and/vs Interpretation

Dialogue Act

(A) Ground truth : Statement-opinion
Predict : Statement-opinion

And if you try to do anything, uh, like, uh, not identify yourself to the government, they know who you are.

(B) Ground truth : Statement-non-opinion
Predict : Statement-non-opinion

I, uh, ride bicycles, uh, fifteen, twenty miles, I don't know, maybe three times, maybe four times a week.

(C) Ground truth : ios, facebook
5-best predict : ios, facebook-graph-api, facebook, objective-c, iphone

I have an iOS application that already using some methods of Facebook Graph API, but I need to implement sending private message to friend by Facebook from my application. As I know, there is no way to sending private messages by Graph API, but it maybe possible by help Facebook Chat API. I already read documentation but it do not help me. If anybody has some kind of example or tutorial, how to implement Facebook Chat API in iOS application, how sending requests or something, it will be very helpfull. Thanks.

(D) Ground truth : python, numpy, matrix
5-best predict : python, numpy, arrays, matrix, indexing

I have a huge matrix that I saved with savetxt with numpy library. Now I want to read a single cell from that matrix, e.g.,
`cell = getCell(i, j); print cell`
return the value 10 for example.
I tried this:
`x = np.loadtxt("fname.m", dtype="int", usecols=([i]))`
`cell = x[j]`
but it is really slow because I loop over many index. Is there a way to do that without reading useless lines?

Key Term

Attention Functions

v: attended vec, **q**: query vec

$\text{MLP}^{\text{att}}(\mathbf{q};\mathbf{v})=$

- Additive Attention: $\text{ug}(\mathbf{W}^1\mathbf{v} + \mathbf{W}^2\mathbf{q})$
- Dot Product: $\mathbf{v} \cdot \mathbf{q}$
- Multiplicative Attention: $\mathbf{v}^T\mathbf{W}\mathbf{q}$

Additive vs Multiplicative

While the two are similar in theoretical complexity, dot-product attention is much faster and more space-efficient in practice, since it can be implemented using highly optimized matrix multiplication code.

While for small values of d_k the two mechanisms perform similarly, additive attention outperforms dot product attention without scaling for larger values of d_k [3]. We suspect that for large values of d_k , the dot products grow large in magnitude, pushing the softmax function into regions where it has extremely small gradients⁴. To counteract this effect, we scale the dot products by $\frac{1}{\sqrt{d_k}}$.

d_k is the dimensionality of q and v

$$\frac{\mathbf{v} \cdot \mathbf{q}}{\sqrt{d_k}}$$



Scaled dot product attention

Paper's Justification:

To illustrate why the dot products get large, assume that the components of q and k are independent random variables with mean 0 and variance \rightarrow Then their dot product, $q \cdot k$ has mean 0 and variance d_k

Multi-head Key-Value Self Attention

Key-Value Attention

- Split an input vector x_i into two vectors $x_i=[k_i;v_i]$
k: key vector
v: value vector


- Use key vector for computing attention

$$\text{MLP}^{\text{att}}(q;x_i)=\frac{k_i \cdot q}{\sqrt{d}} \quad //\text{scaled multiplicative}$$

- Use value vector for computing attended summary

$$c = \sum_{i=1}^n \alpha_{[i]} \cdot v_i$$

Key-Value Attention (alternative)

-  an input vector x_i into two vectors
k: key vector: $k_i = W^K x_i$
v: value vector: $v_i = W^V x_i$

- Use key vector for computing attention

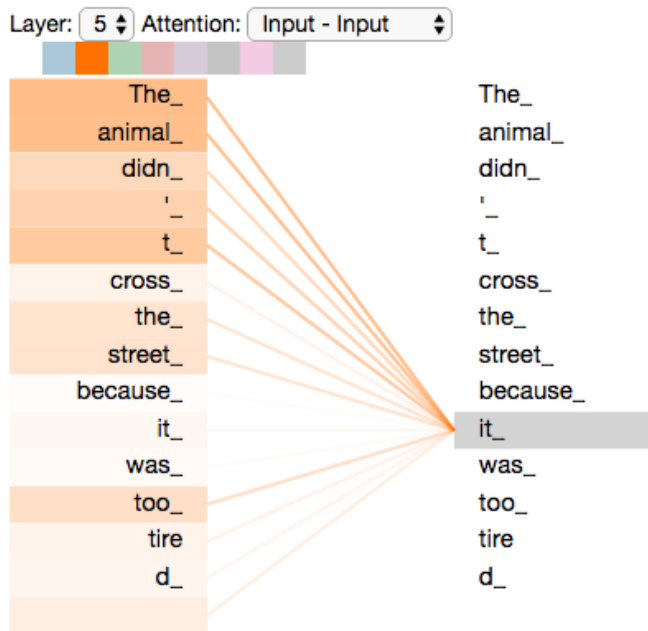
$$\text{MLP}^{\text{att}}(q; x_i) = \frac{k_i \cdot q}{\sqrt{d}} \quad // \text{scaled multiplicative}$$

- Use value vector for computing attended summary

$$c = \sum_{i=1}^n \alpha_{[i]} \cdot v_i$$

Self-attention (single-head, high-level)

"The animal didn't cross the street because it was too tired"





There is no external query q .

The input is also the query.


Many approaches:

<https://ruder.io/deep-learning-nlp-best-practices/>

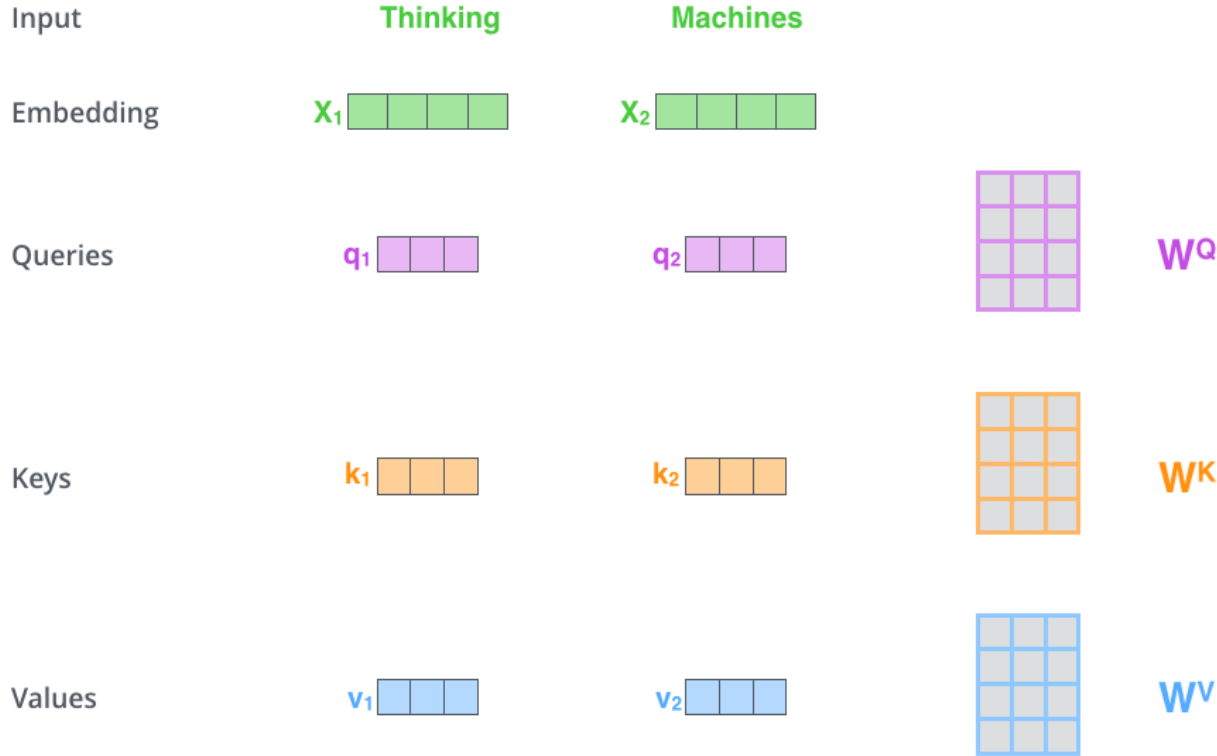
Key-Value Single-Head Self Attention

- Project an input vector x_i into  vectors
k: key vector: $k_i = W^K x_i$
v: value vector: $v_i = W^V x_i$
q:  vector: $q_i = W^Q x_i$
- Use key and query vectors for computing attention of i^{th} word at word j

$$\text{MLP}^{\text{att}}(x_j; x_i) = \frac{k_i \cdot q_j}{\sqrt{d}} \quad // \text{scaled multiplicative}$$

- Use value vector for computing attended summary  = $\sum_{i=1}^n \alpha_{[i]} \cdot v_i$

Key-Value Single-Head Self Attention

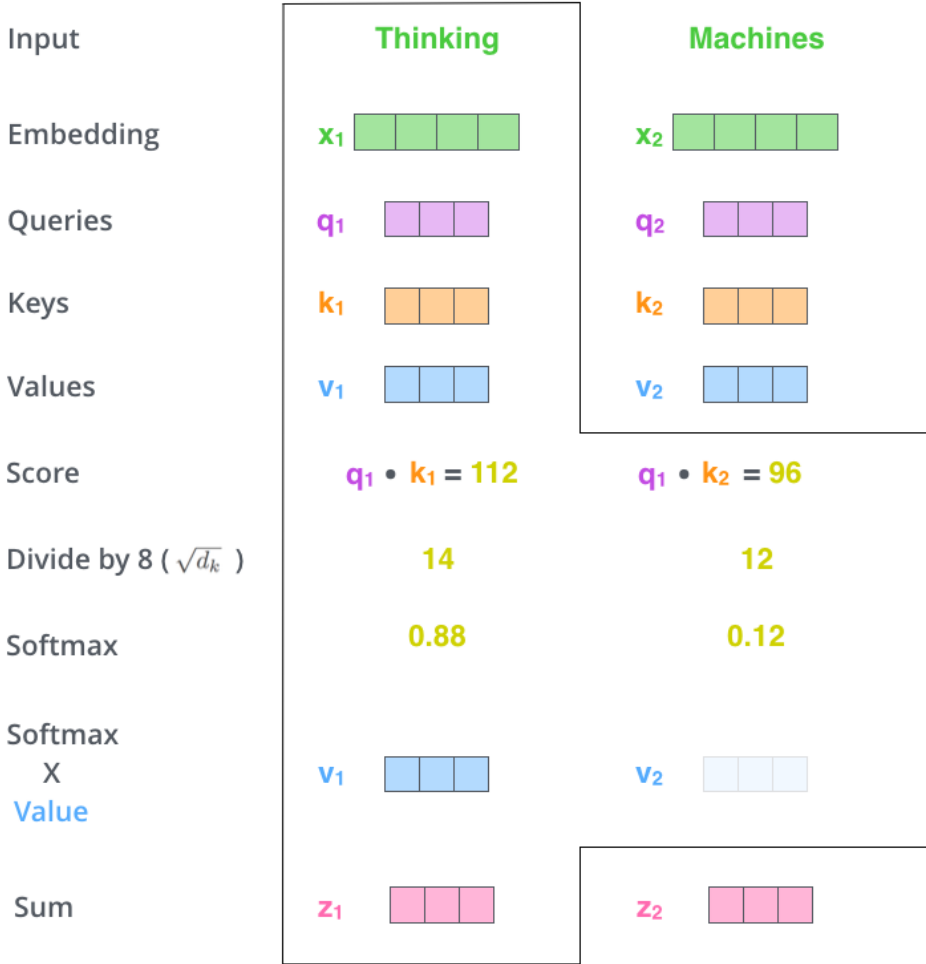


Creation of query, key and value vectors by multiplying by weight matrices

Separation of Value and Key

Matrix multiplications are quite efficient and can be done in aggregated manner

Key-Value Single-Head Self Attention

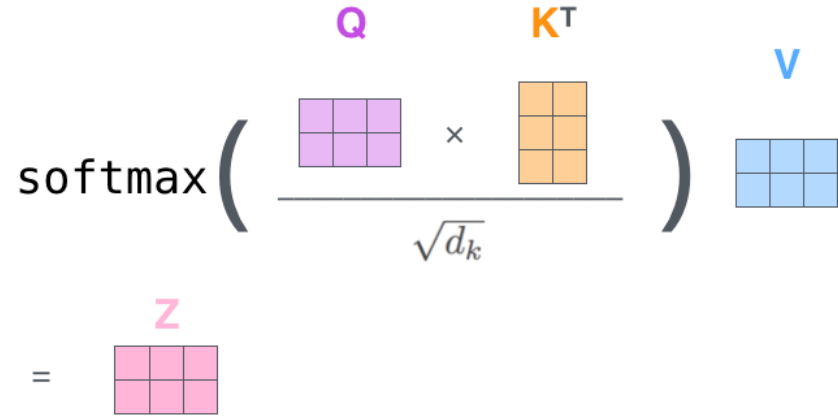


Key-Value Single-Head Self Attention

$$X \times W^Q = Q$$


$$X \times W^K = K$$

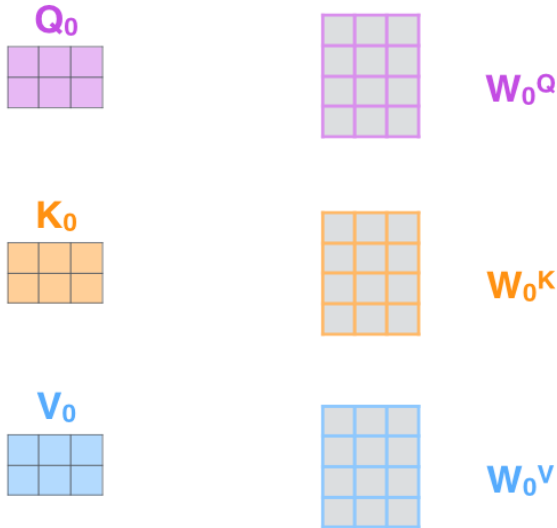

$$X \times W^V = V$$


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


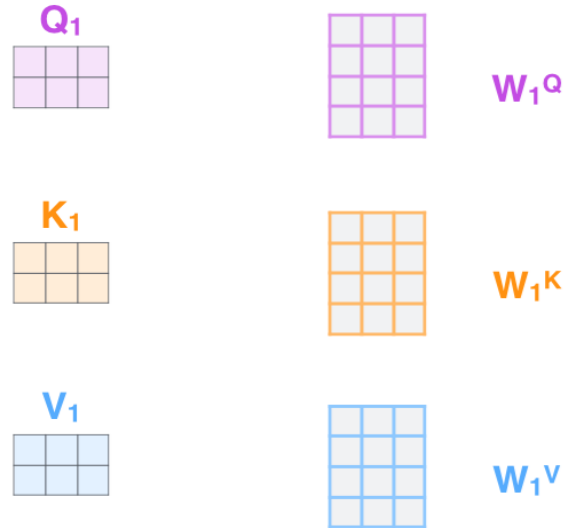
Key-Value Multi-Head Self Attention



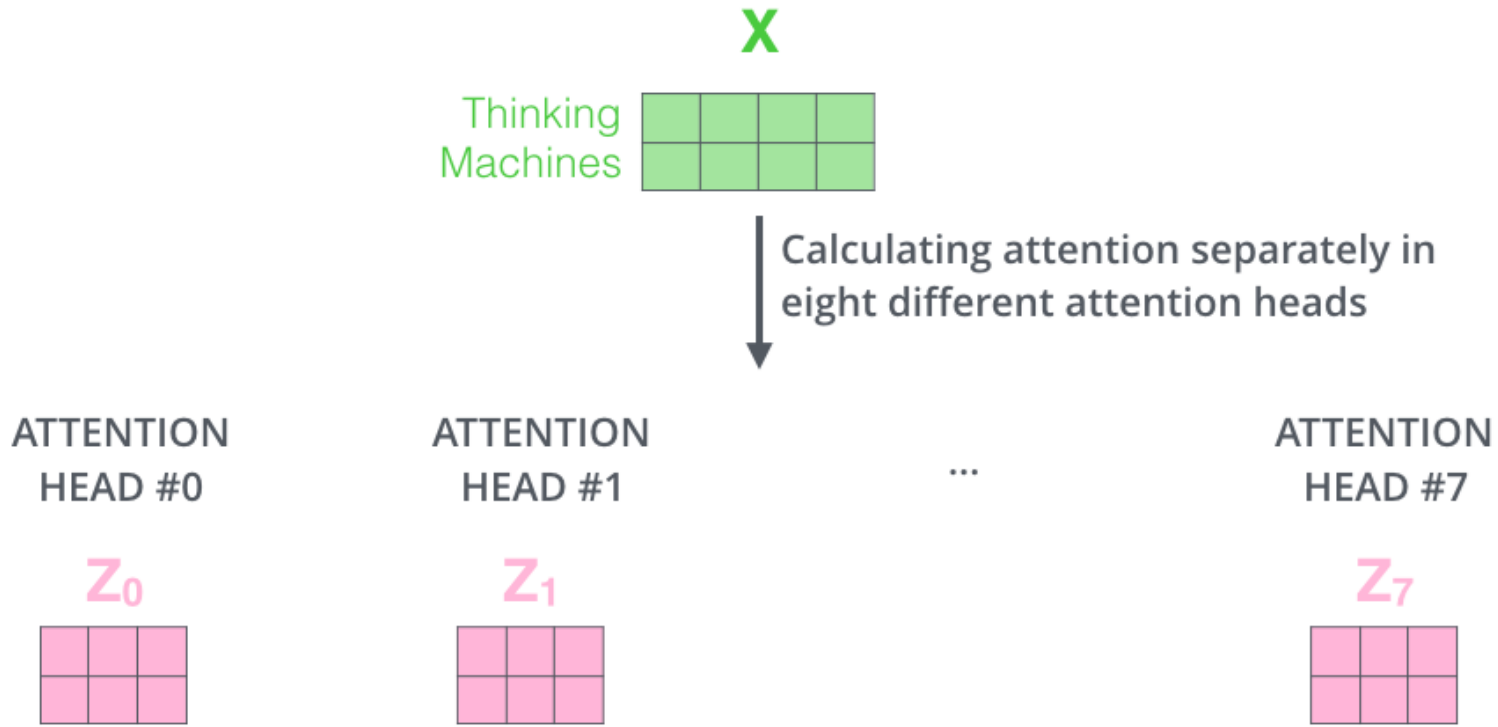
ATTENTION HEAD #0



ATTENTION HEAD #1

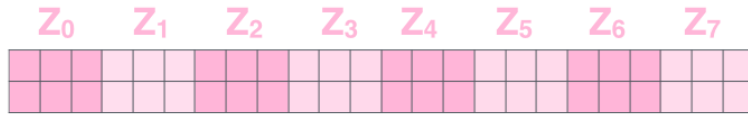


Multi-Head Attention



Multi-Head Attended Vector \rightarrow Output

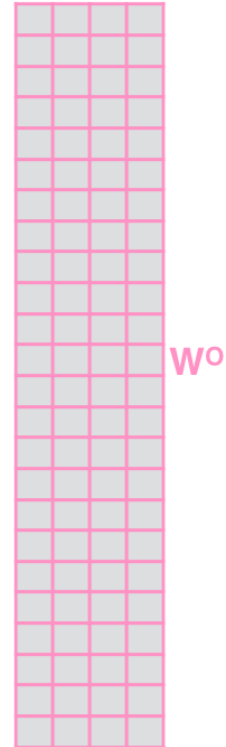
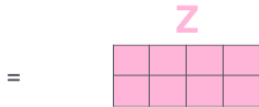
1) Concatenate all the attention heads



2) Multiply with a weight matrix W^O that was trained jointly with the model

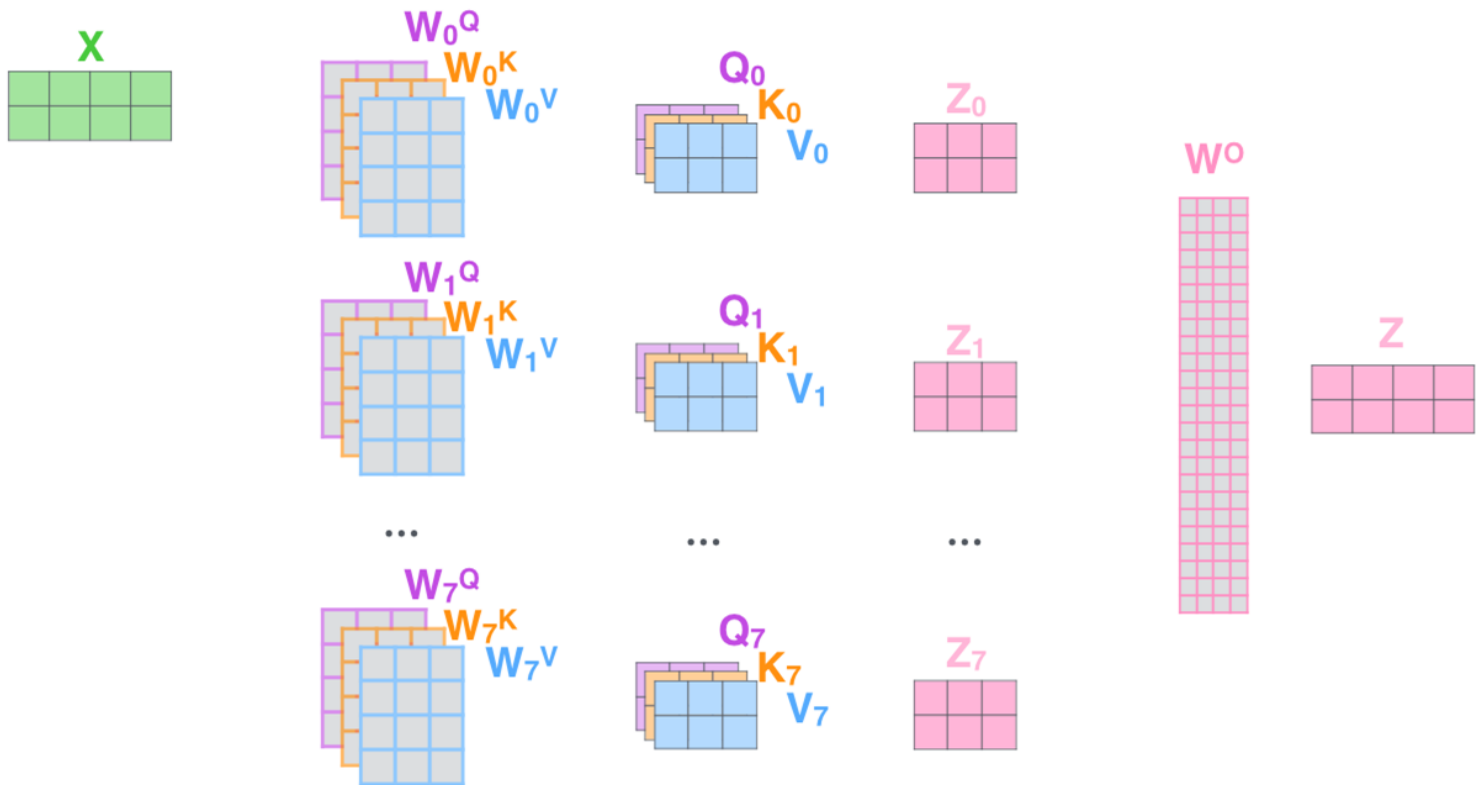
\times

3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN

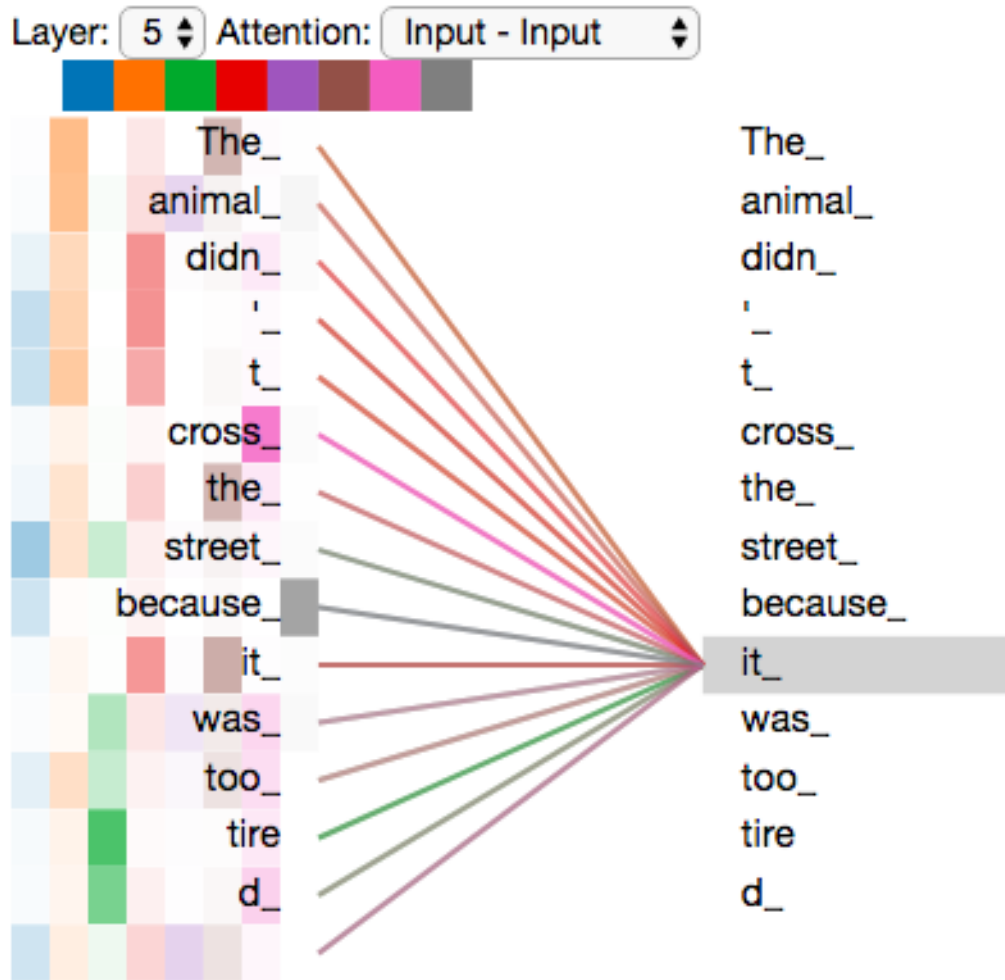


Key-Value Multi-Head Self Attention (summary)

Thinking
Machines



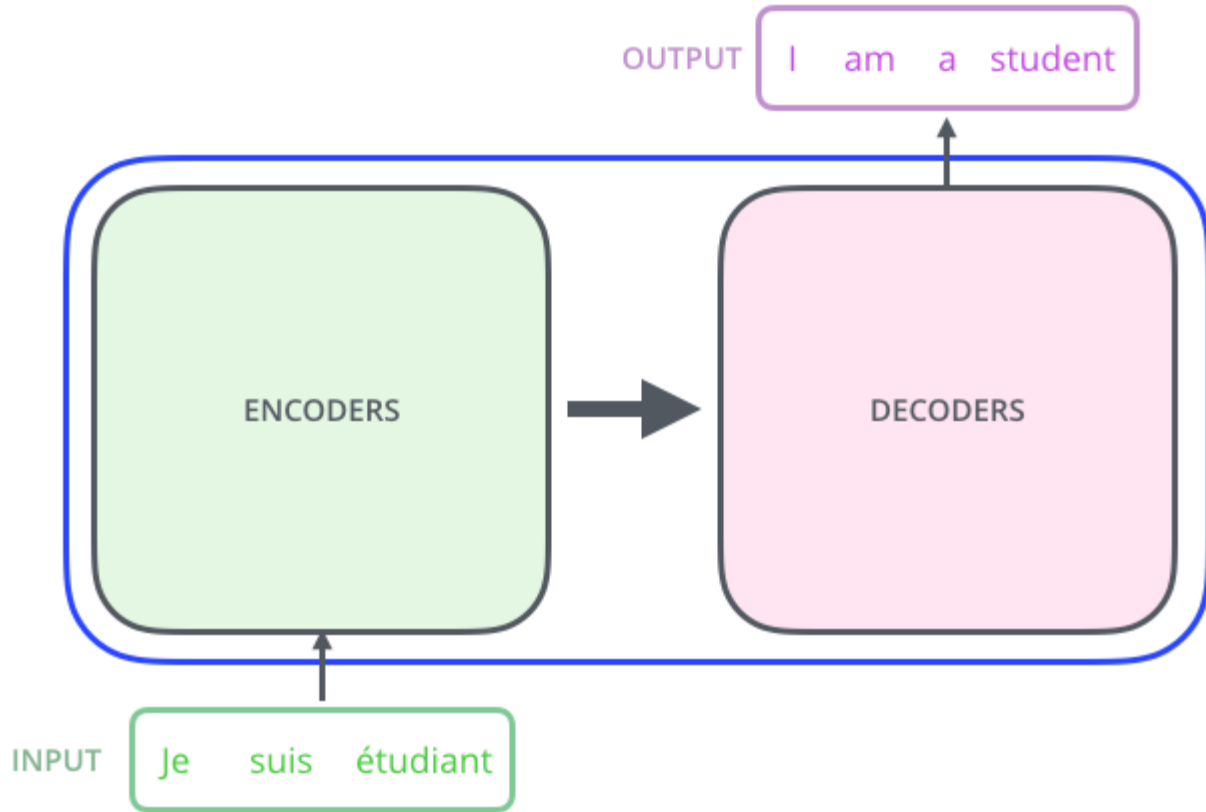
Self attention visualisation (Interpretable?!)



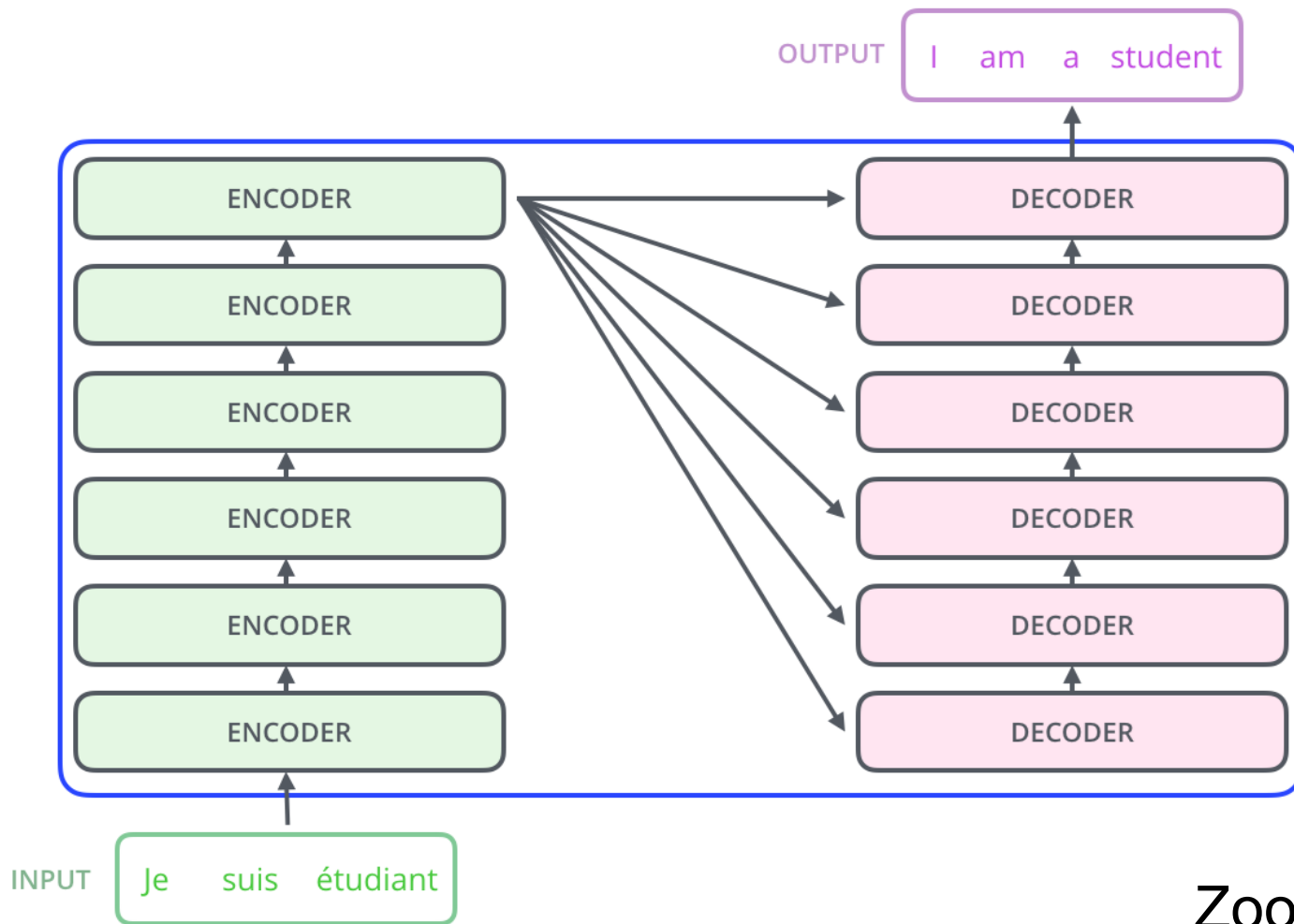
Transformer Architecture (for text classification)

Motivation

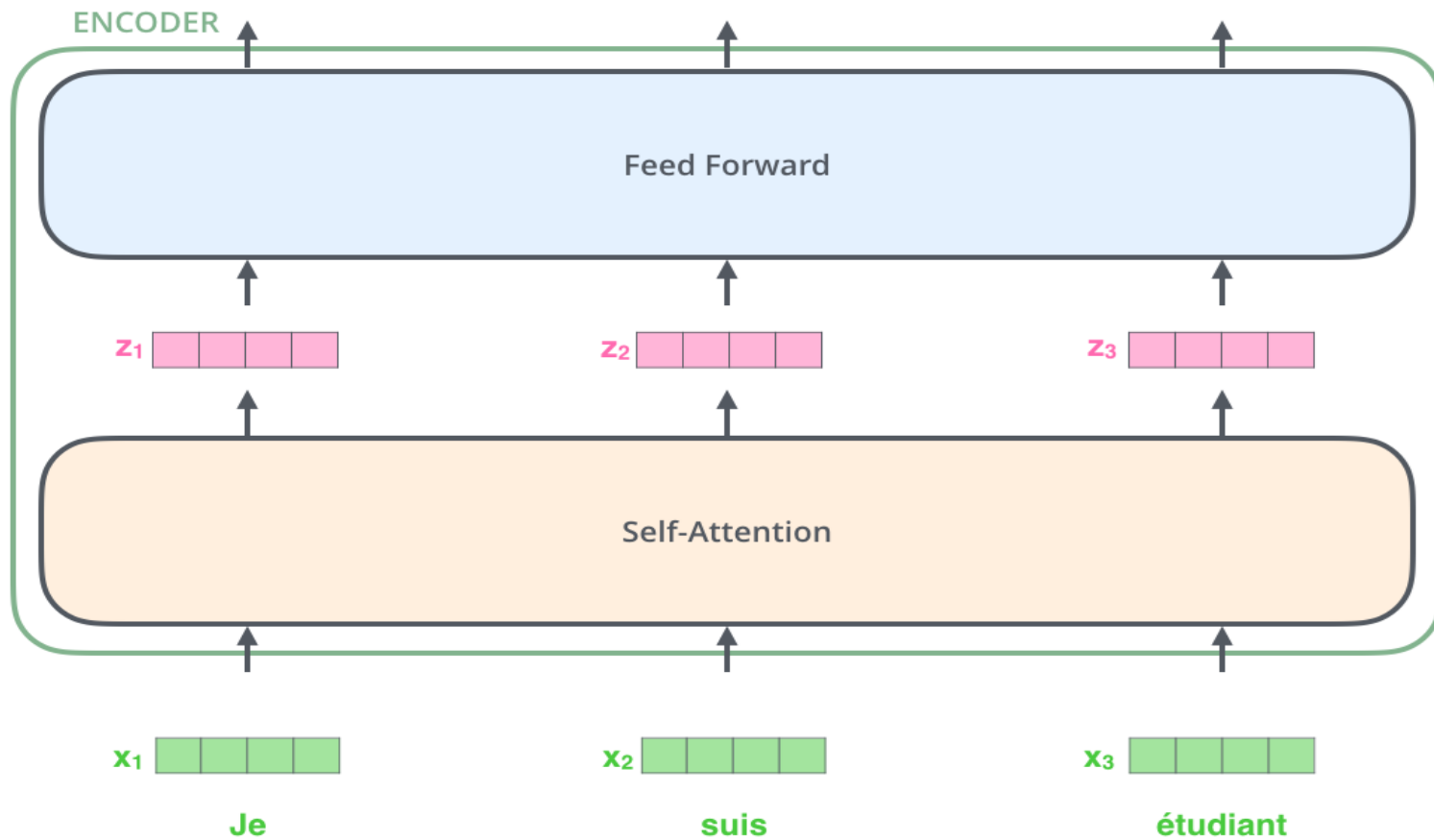
- Recurrence is powerful
 - Issues with learnability: vanishing gradients
 - Issues with remembering long sentences
 - Issues with scalability: backpropagation time high due to sequentiality in sentence length
 - Issues with scalability: can't be parallelized even at test time – $O(\text{sentence length})$
- Remove recurrence: only use attention
“Attention is All You Need”



We focus only on encoder... (decoder is for sequence generation – will study later)



Zooming in...

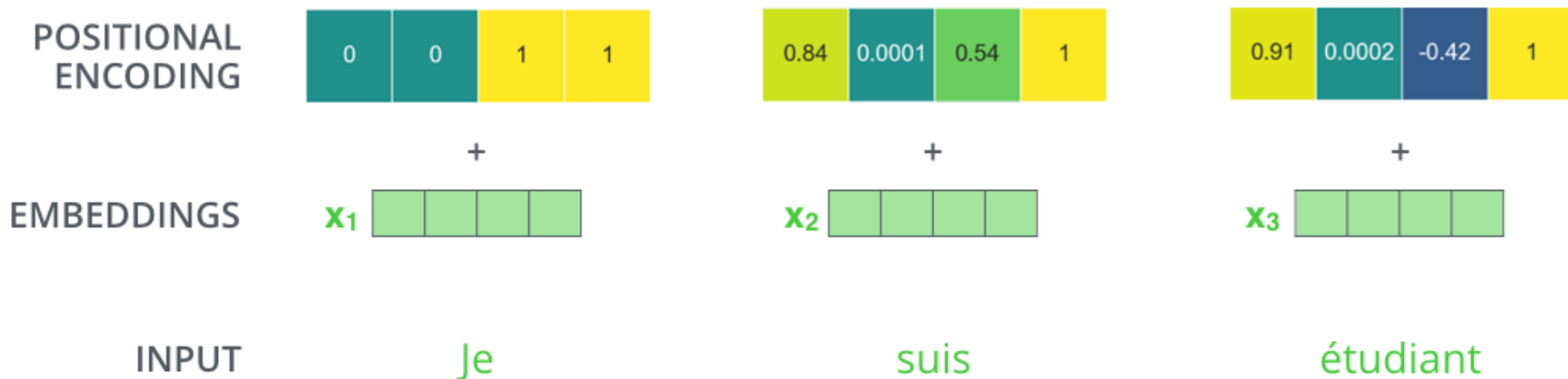


Can you see a fundamental limitation?

Encoders have same architecture but different weights...

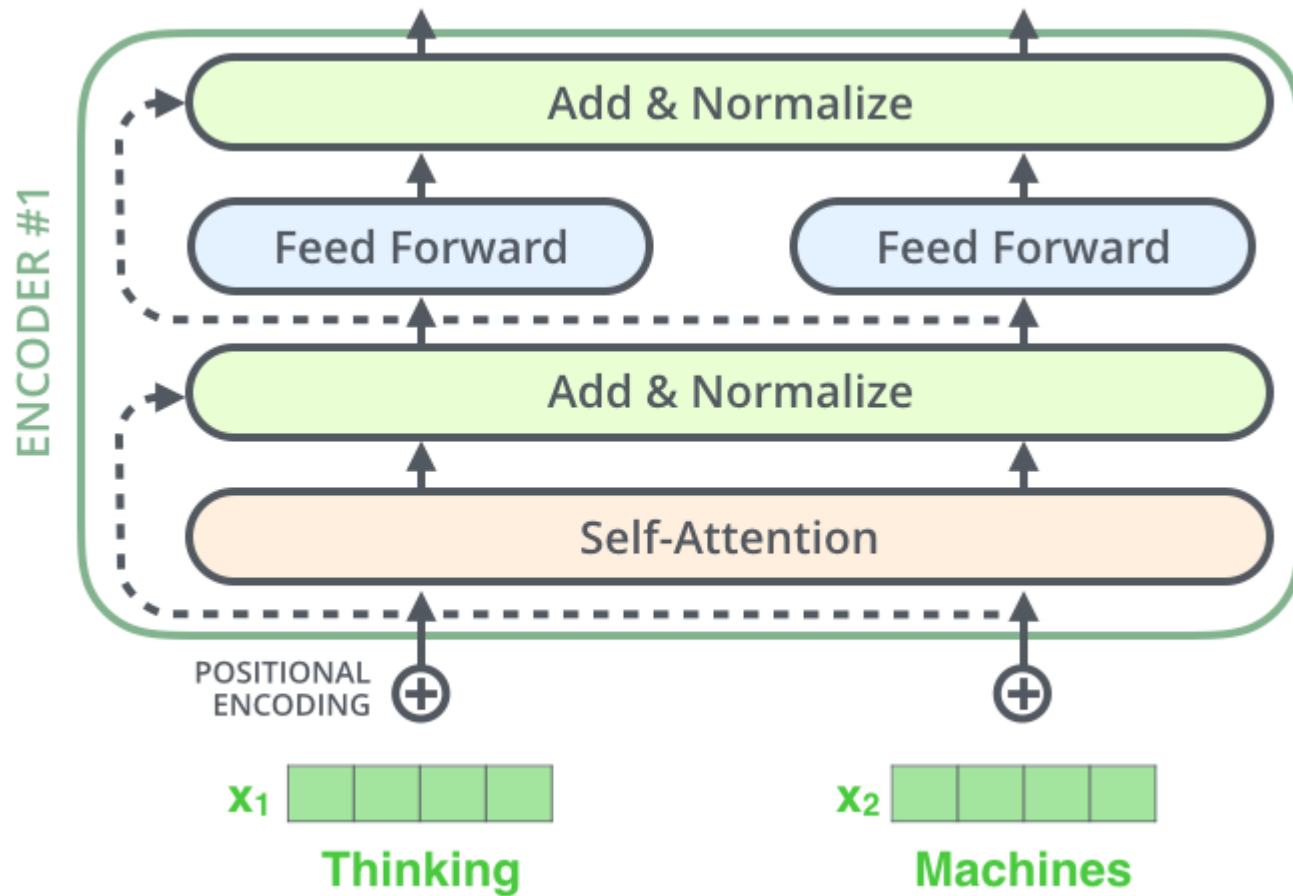
Zooming in further...

A note on Positional embeddings

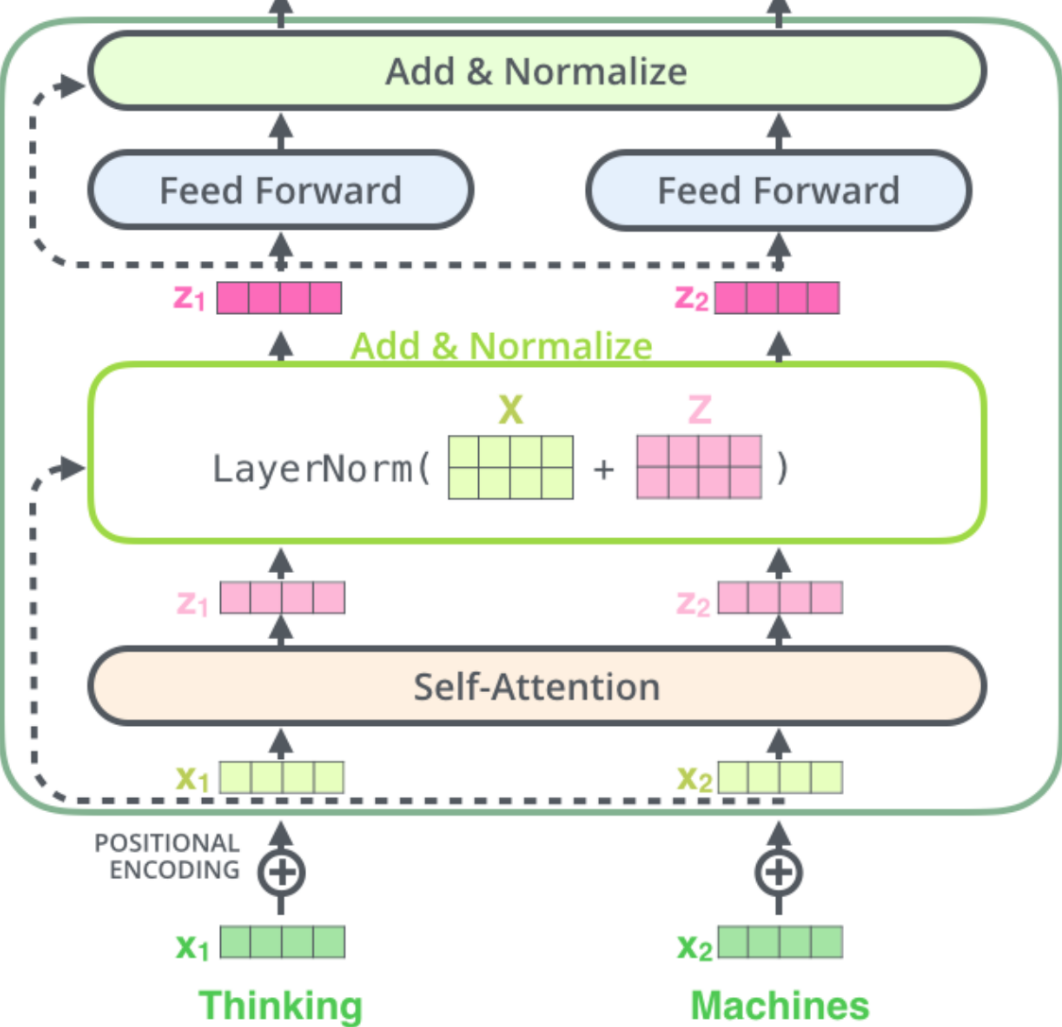


Positional embeddings can be extended to any sentence length but if any test input is longer than all training inputs then we will face issues.

Solution: use a functional form (as in Transformer paper – sinusoidal encoding)



Adding residual connections...

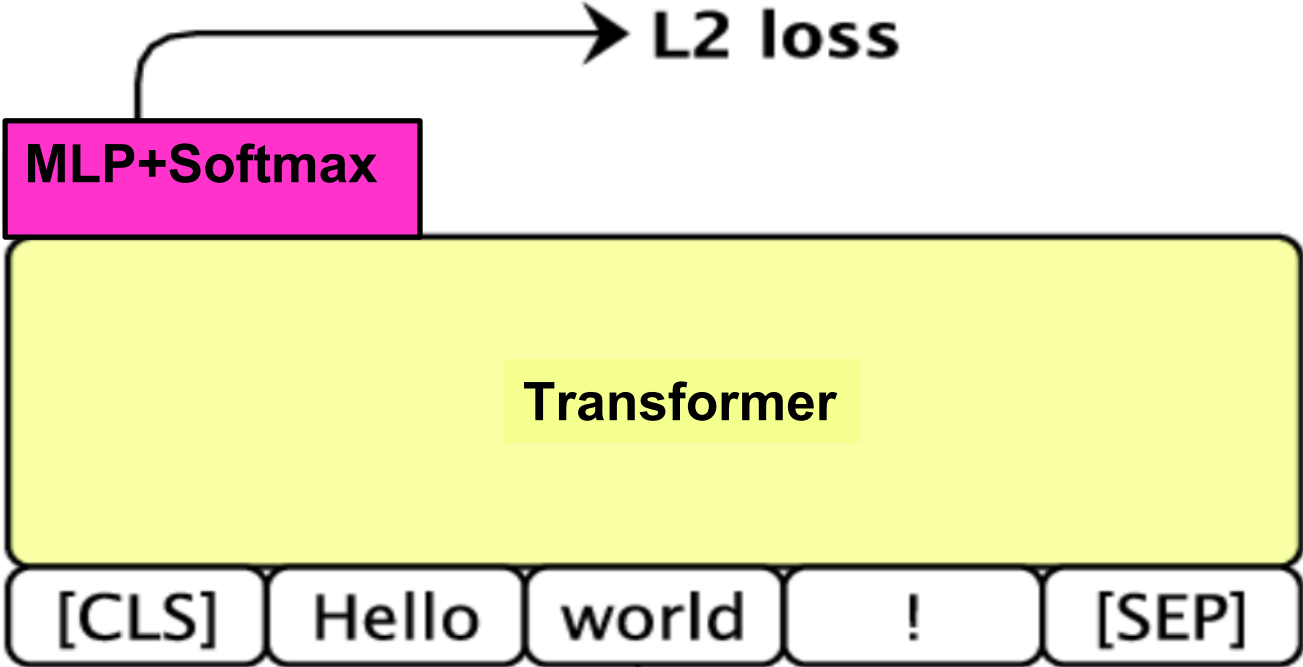


The residual connections help the network train, by allowing gradients to flow through the networks directly.

The layer normalizations stabilize the network -- substantially reducing the training time necessary.

The pointwise feedforward layer is used to project the attention outputs potentially giving it a richer representation.

Use of [CLS] for Text Classification



Pros

- Current state-of-the-art in machine translation and text simplification.
- Enables deep architectures
- Intuition of model well explained
- Easier learning of long-range dependencies
- Can be efficiently parallelized
- Gradients don't suffer from vanishing gradients

Cons

Huge number of parameters so-

- Very data hungry
- Takes a long time to train
- No study of memory utilisation

Other issues

- Keeping sentence length limited
- How to ensure multi-head attention has diverse perspectives.

Reformer & Longformer

The Efficient Transformers

Kitaev et. al. (January 2020, ICLR)

Beltagy et. al. (April 2020, Arxiv)

Concerns about the transformer

“Transformer models are also used on increasingly long sequences. Up to 11 thousand tokens of text in a single example were processed in (Liu et al., 2018) ... These large-scale long-sequence models yield great results but **strain resources to the point where some argue that this trend is breaking NLP research**”

“Many large Transformer models **can only realistically be trained in large industrial research laboratories** and such models trained with model parallelism cannot even be fine-tuned on a single GPU as their **memory requirements demand a multi-accelerator hardware setup**”

Memory requirement estimate (per layer)

Largest transformer layer ever: 0.5B parameters = 2GB

Activations for 64K tokens for embedding size 1K and batch size 8

$$= 64K * 1K * 8 = 2GB$$

Training data used in BERT = 17GB

Why can't we fit everything in one GPU? 32GB GPUs are common today.

Caveats follow ->>>>>

Caveats

1. There are N layers in a transformer, whose activations need to be stored for backpropagation
2. We have been ignoring the feed-forward networks upto now, whose depth even exceeds the attention mechanism so contributes to significant fraction of memory use.
3. Dot product attention is $O(L^2)$ in space complexity where L is length of text input.

Solutions

1. Reversible layers, first introduced in Gomez et al. (2017), enable storing only a single copy of activations in the whole model, so the N factor disappears.
2. Splitting activations inside feed-forward layers and processing them in chunks saves memory inside feed-forward layers.
3. Approximate attention computation based on locality-sensitive hashing replaces the $O(L^2)$ factor in attention layers with $O(L \log L)$ and so allows operating on long sequences.

Locality Sensitive Hashing

Hypothesis: Attending on all vectors is approximately same as attending to the 32/64 closest vectors to query in key projection space.

To find such vectors easily we require:

- Key and Query to be in same space
- Locality sensitive hashing i.e. if distance between key and query is less then distance between their hash values is less.

Locality sensitive hashing scheme taken from Andoni et al., 2015

For simplicity, a bucketing scheme chosen: attend on everything in your bucket

Locality sensitive hashing

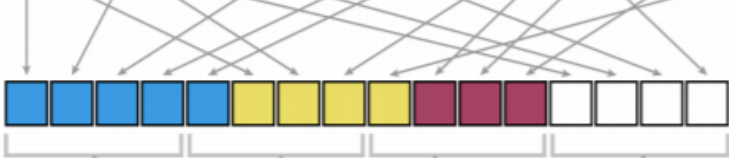
Sequence of queries=keys



LSH bucketing



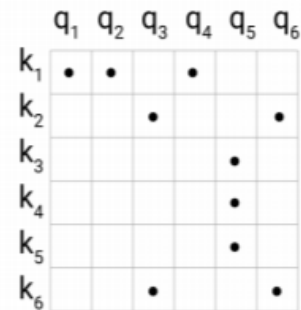
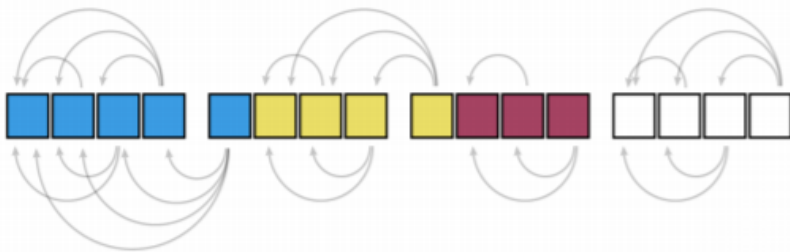
Sort by LSH bucket



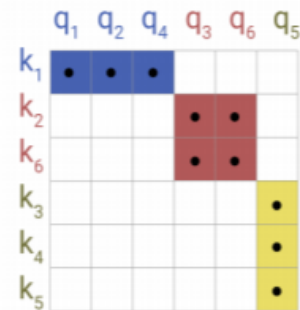
Chunk sorted sequence to parallelize



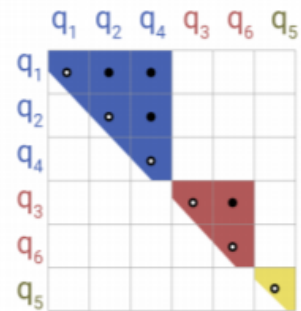
Attend within same bucket in own chunk and previous chunk



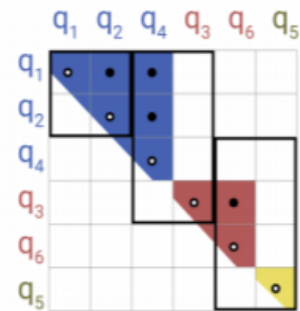
(a) Normal



(b) Bucketed



(c) Q = K



(d) Chunked

Solutions

1. Reversible layers, first introduced in Gomez et al. (2017), enable storing only a single copy of activations in the whole model, so the N factor disappears.
2. Splitting activations inside feed-forward layers and processing them in chunks saves memory inside feed-forward layers.
3. Approximate attention computation based on locality-sensitive hashing replaces the $O(L^2)$ factor in attention layers with $O(L \log L)$ and so allows operating on long sequences.

RevNets

Reversible residual layers were introduced in Gomez et. al. 2017

Idea: Activations of previous layer can be recovered from activations of subsequent layers, using model parameters.

Normal residual layer: $y = x + F(x)$

Reversible layer:

$$y_1 = x_1 + F(x_2)$$

$$x_2 = y_2 - G(y_1)$$

$$y_2 = x_2 + G(y_1)$$

$$x_1 = y_1 - F(x_2)$$

So, for transformer:

$$Y_1 = X_1 + \text{Attention}(X_2)$$

$$Y_2 = X_2 + \text{FeedForward}(Y_1)$$

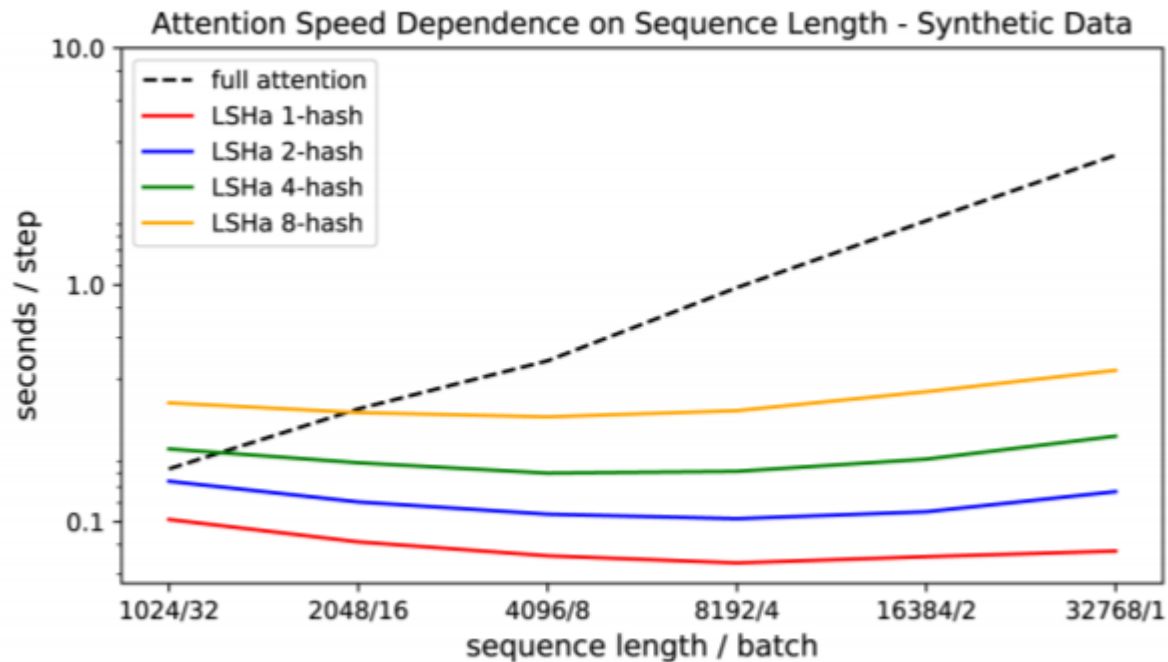
Chunking

$$Y_2 = [Y_2^{(1)}; \dots; Y_2^{(c)}] = [X_2^{(1)} + \text{FeedForward}(Y_1^{(1)}); \dots; X_2^{(c)} + \text{FeedForward}(Y_1^{(c)})]$$

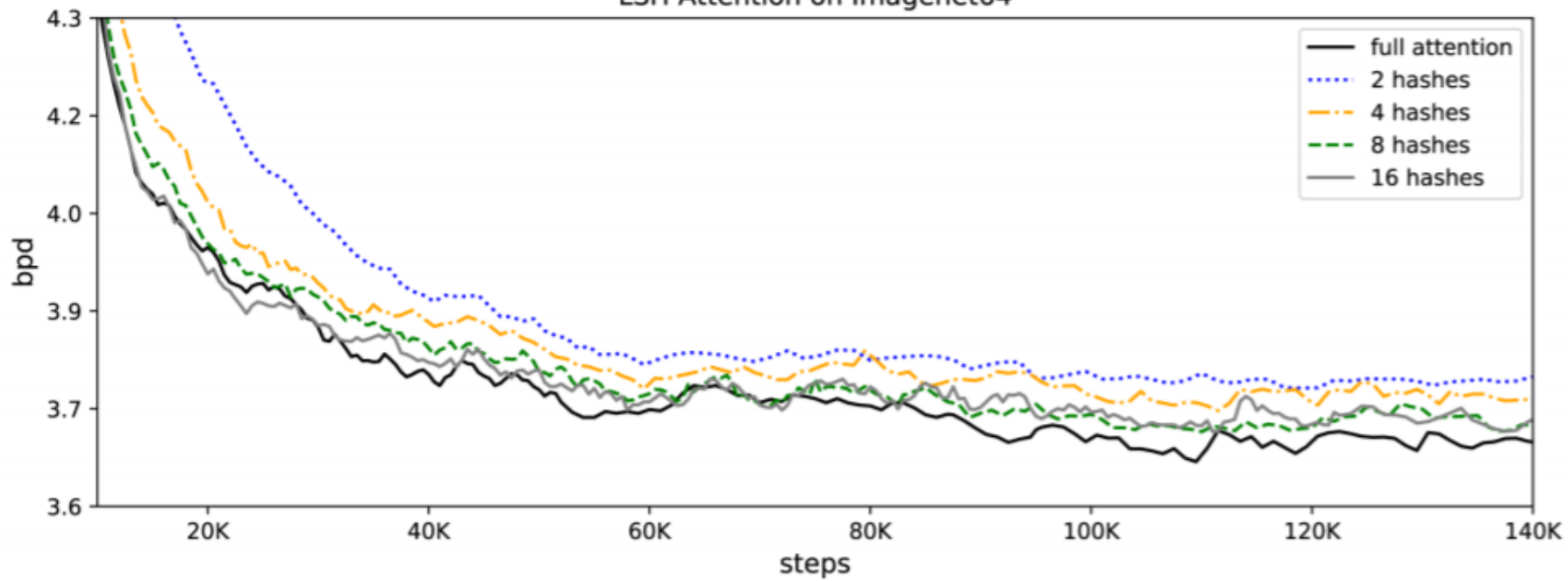
Operations done a chunk at a time:

- Forward pass of Feed-forward network
- Reversing the activations during backpropagation
- For large vocabularies, chunk the log probabilities

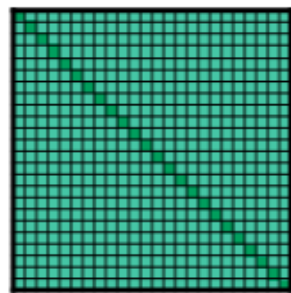
Experiments



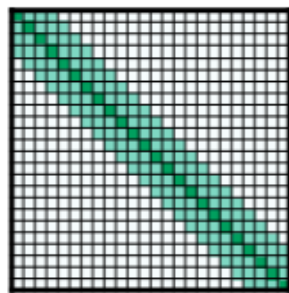
LSH Attention on Imagenet64



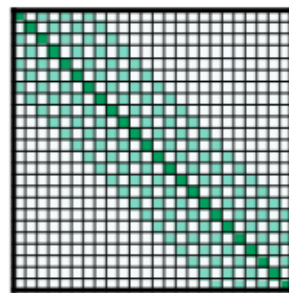
Longformer



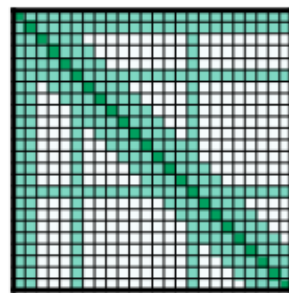
(a) Full n^2 attention



(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.