# Advanced Topics in Machine Learning

Alejo Nevado-Holgado

## Lecture 16 (NLP 8) - Transformers
V 0.1 (4 Mar 2020 - final version)

# Course structure

➢ **Introduction:** What is NLP. Why it is hard. Why NNs work well ← Lecture 9 (NLP 1)

➢ **Word representation:** How to represent the meaning of individual words
- Old technology: One-hot representations, synsets ← Lecture 9 (NLP 1)
- Embeddings: First trick that boosted the performance of NNs in NLP ← Lecture 9 (NLP 1)
  - Word2vec: Single layer NN. CBOW and skip-gram ← Lecture 10 (NLP 2)
  - Co-occurrence matrices: Basic counts and SVD improvement ← Lecture 10 (NLP 2)
  - Glove: Combining word2vec and co-occurrence matrices idea ← Lecture 10 (NLP 2)
  - Evaluating performance of embeddings ← Lecture 10 (NLP 2)

➢ **Named Entity Recognition (NER):** How to find words of specific meaning within text
- Multilayer NNs: Margin loss. Forward- and back-propagation ← Lecture 11 (NLP 3)
- Better loss functions: margin loss, regularisation ← Lecture 11 (NLP 3)
- Better initializations: uniform, xavier ← Lecture 11 (NLP 3)
- Better optimizers: Adagrad, RMSprop, Adam... ← Lecture 11 (NLP 3)

# Course structure

➢ **Language modelling:** How to represent the meaning of full pieces of text
- Old technology: N-grams ← Lecture 12 (NLP 4)
- Recursive NNs language models (RNNs) ← Lecture 12 (NLP 4)
- Evaluating performance of language models ← Lecture 12 (NLP 4)
- Vanishing gradients: Problem. Gradient clipping ← Lecture 13 (NLP 5)
- Improved RNNs: LSTM, GRU, Bidirectional... ← Lecture 13 (NLP 5)

➢ **Machine translation:** How to translate text
- Old technology: Georgetown−IBM experiment and ALPAC report ← Lecture 14 (NLP 6)
- Seq2seq: Greedy decoding, encoder-decoder, beam search ← Lecture 14 (NLP 6)
- Attention: Simple attention, transformers, reformers ← Lecture 14 (NLP 6)
- Evaluating performance: BLEU ← Lecture 14 (NLP 6)

# Course structure

➢ **Question Answering:** X

- Task definition, datasets, cloze-style tasks, Attentive Reader ← Lecture 15 (NLP 7)

➢ **Conference Resolution:** X

- Task definition, pairs method, clustering method, language models ← Lecture 15 (NLP 7)

➢ **Convolutional Neural Networks:** X

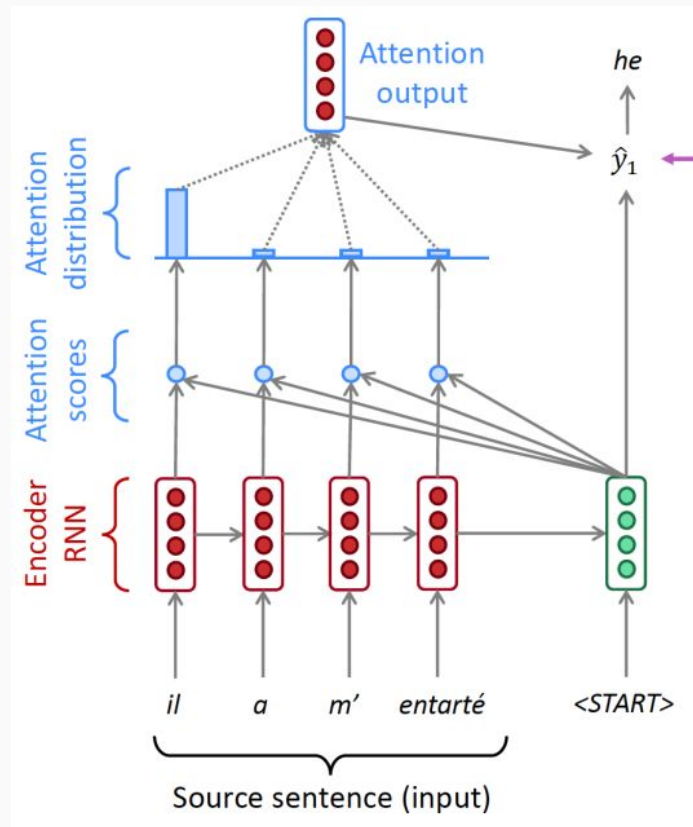- CNNs in vision, CNNs in language, example ← Lecture 15 (NLP 7)

➢ **Transformers:** X

- Architecture: encoder, self-attention, encoding position, decoder ← Lecture 16 (NLP 8)
- Existing systems. Ranking ← Lecture 16 (NLP 8)

# Transformers: Why a new architecture?

**The problem:** Recurrences are very slow to train, and their computations cannot be parallelized. Although LSTMs and GRUs capture long terms relationships much better than vanilla RNNs, they still don't do it well enough.

**The solution:** We saw in previous lectures that attention can give any time step access to any other time step, no matter the length of the input. The whole purpose of recurrence in RNN architectures was accessing previous time steps no matter the length of the input. Why don't we simply use pure attention to access all time steps? It is parallelizable, and maybe it captures long term relationships better than recurrence.
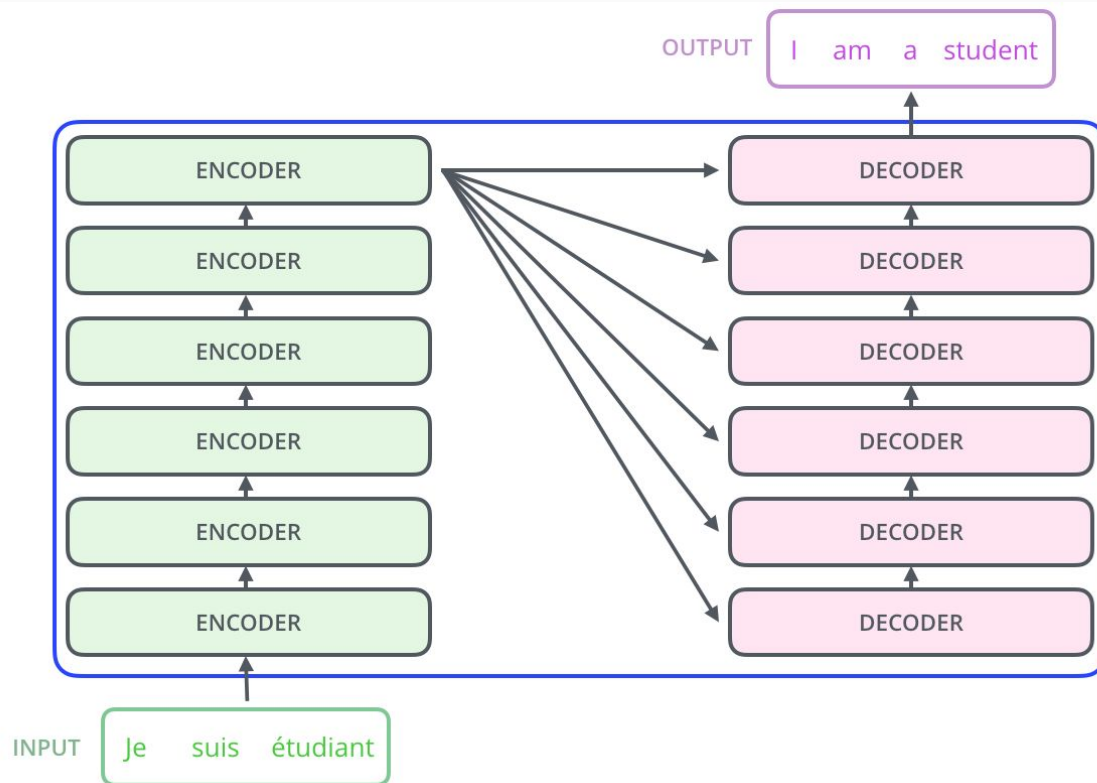
# Transformers: Architecture

The transformer follows an encoder-decoder architecture, with all decoders attending to the last state of the encoder. This is the same as we studied a few lectures ago for translation (Lecture NLP 6)

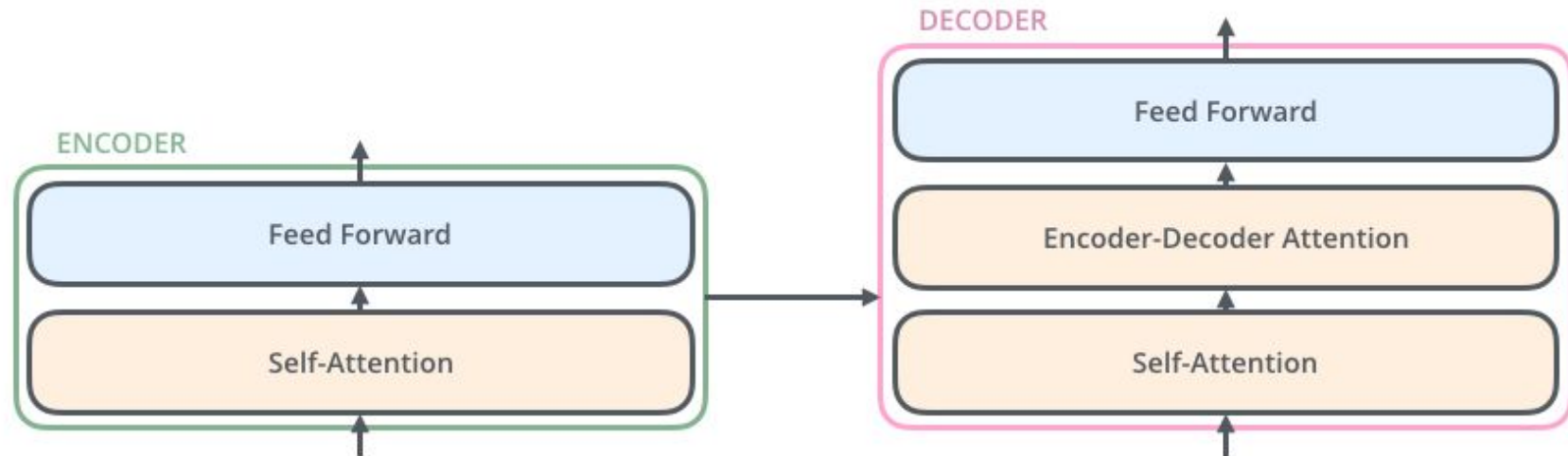Original paper: https://arxiv.org/abs/1706.03762

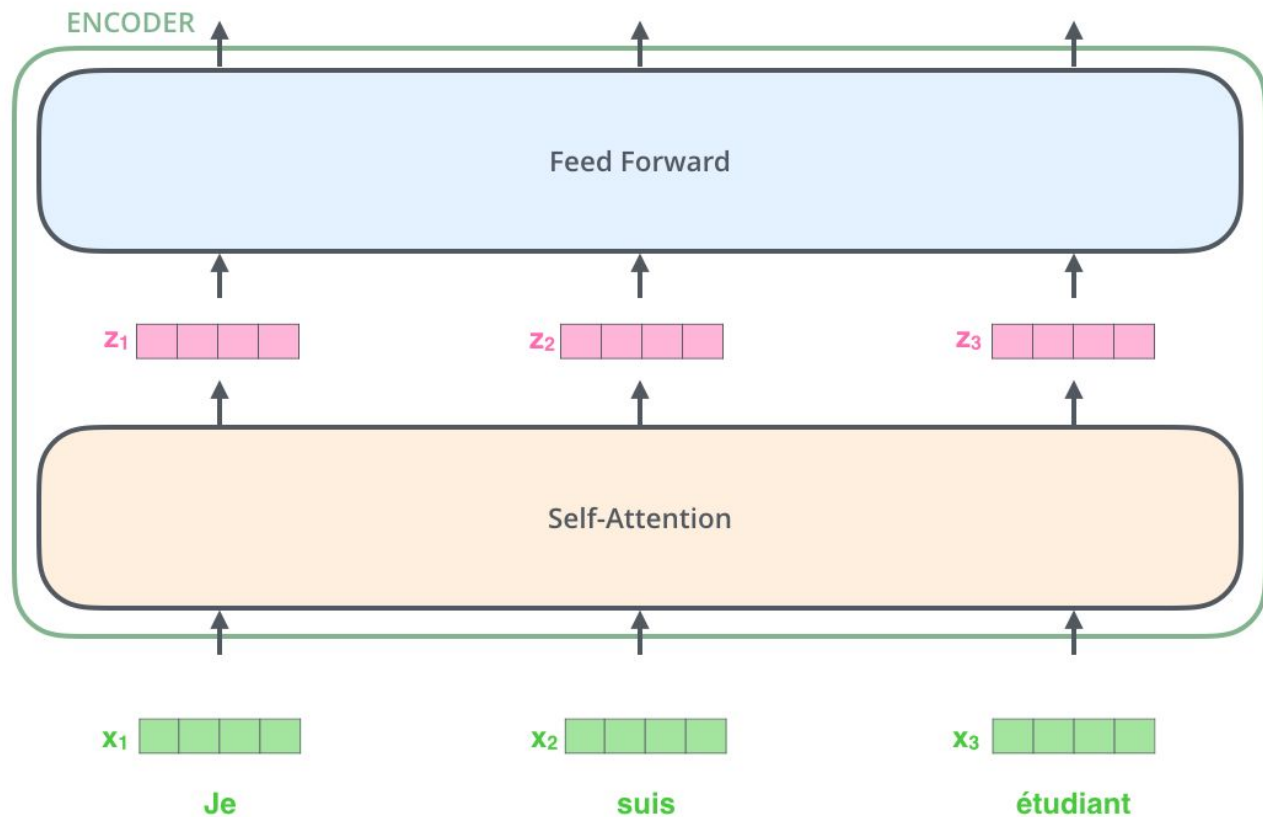Best description out there: http://jalammar.github.io/illustrated-transformer/



6

# Transformers: Architecture

The difference between the transformer and the encoder-decoder of Lecture NLP 6 is on the internal architecture of each encoder and each decoder. Rather than simple hidden states, each encoder and decoder is a mini-NN of its own. This is sometimes called a 'module', 'block' or even 'layer', and it is very common in modern NNs (e.g. VGG16, ResNet, ByteNet…)
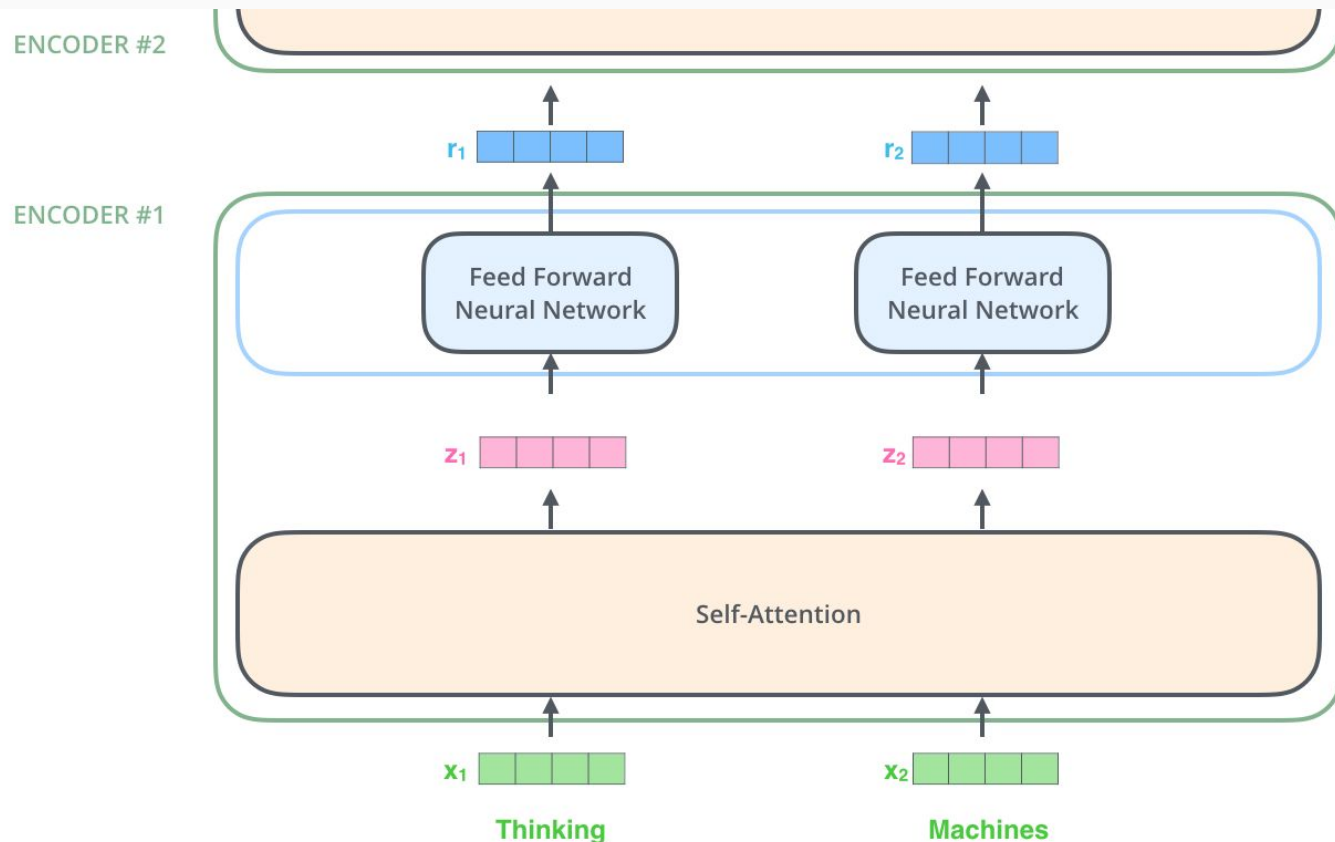
# Transformers: Encoder

The first encoder receives an embedding $[\mathbf{x}^{(t)}]_X$ per word (X = num embedding dimensions). The self-attention layer mixes information across words, producing a new presentation per word $[\mathbf{z}^{(t)}]_Z$, like the attention mechanism of a seq2seq.
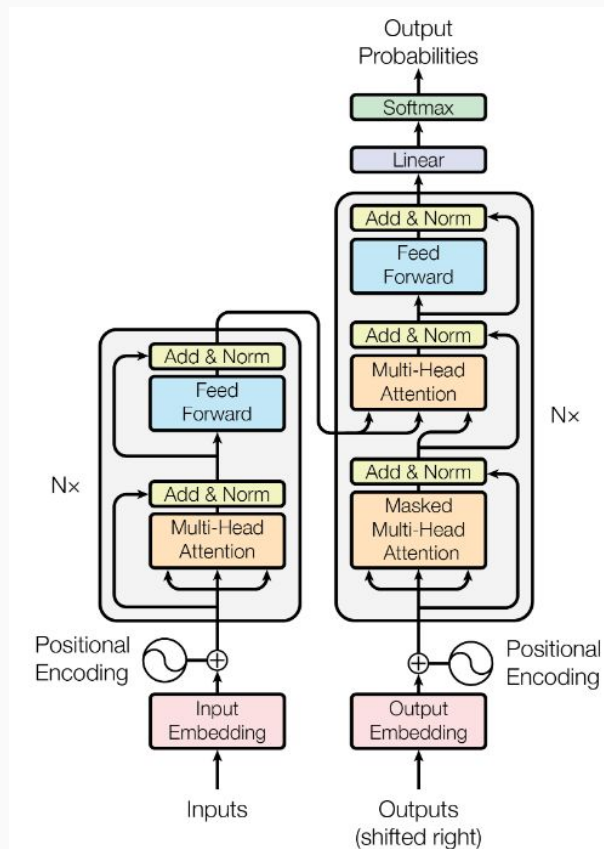
# Transformers: Encoder

The new presentation per word $[z^{(t)}]_Z$ that emerges from the self-attention layer, is then transformed with a feed forward fully connected NN into $[r^{(t)}]_R$.
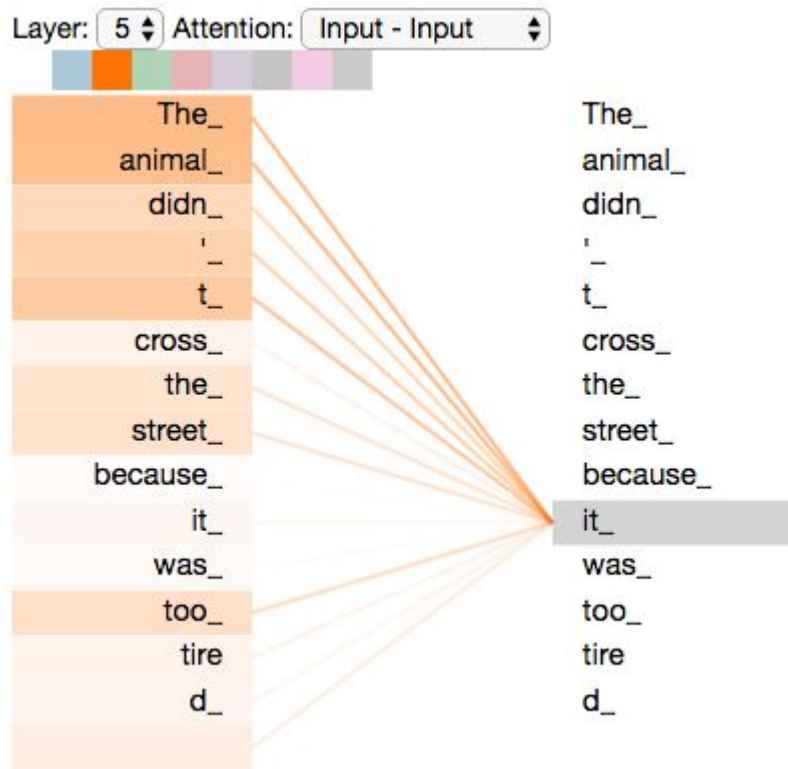
# Transformers: Encoder

Be careful! In the original paper they use a residual connection in the self-attention and feed-forwards layers. The description by Jalammar does not emphasize this

**Encoder:** The encoder is composed of a stack of $N = 6$ identical layers. Each layer has two sub-layers. The first is a multi-head self-attention mechanism, and the second is a simple, position-wise fully connected feed-forward network. We employ a residual connection [11] around each of the two sub-layers, followed by layer normalization [1]. That is, the output of each sub-layer is LayerNorm($x$ + Sublayer($x$)), where Sublayer($x$) is the function implemented by the sub-layer itself. To facilitate these residual connections, all sub-layers in the model, as well as the embedding layers, produce outputs of dimension $d_{\text{model}} = 512$.

# Transformers: Self-attention

The effect of attention is that $[z^{(t)}]_Z$ (for a given time step 't') becomes a mixture of the original embeddings $[x^{(t)}]_X$. The idea is that each time step 'borrows' information from other time steps that it is related to. For instance, a pronoun may borrow information from the complement of the name that it refers to. Layer, the feed-forwards NN further transforms $[z^{(t)}]_Z$ into $[r^{(t)}]_R$, but this time without mixing information across time steps.
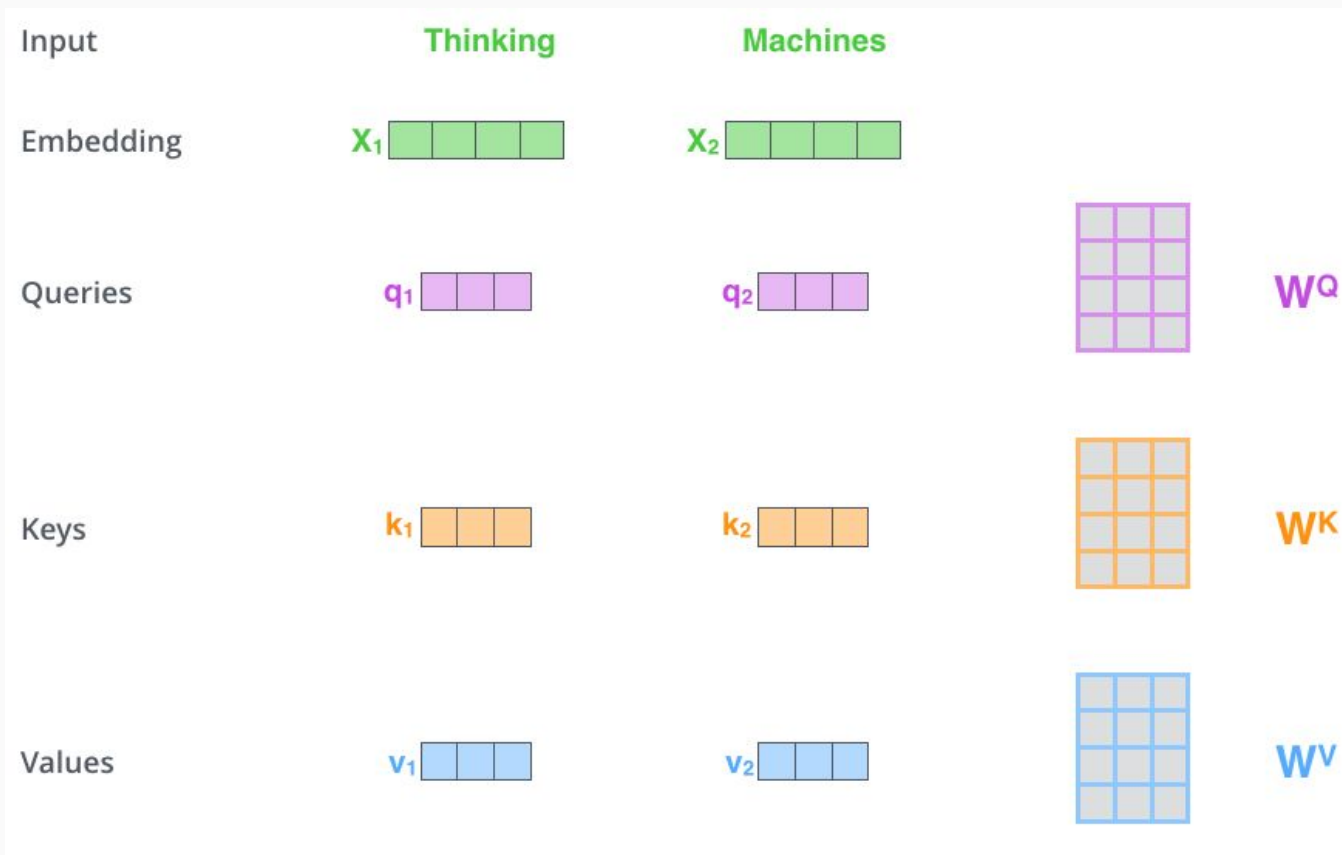
# Transformers: Self-attention

The attention mechanism is multiplicative. We first create a query, a key and a value vector per time step by linearly transforming each embedding:

$$[\mathbf{q}^{(t)}]_Q = [\mathbf{W^Q}]_{QX} [\mathbf{x}^{(t)}]_X$$

$$[\mathbf{k}^{(t)}]_Q = [\mathbf{W^K}]_{KQ} [\mathbf{x}^{(t)}]_X$$
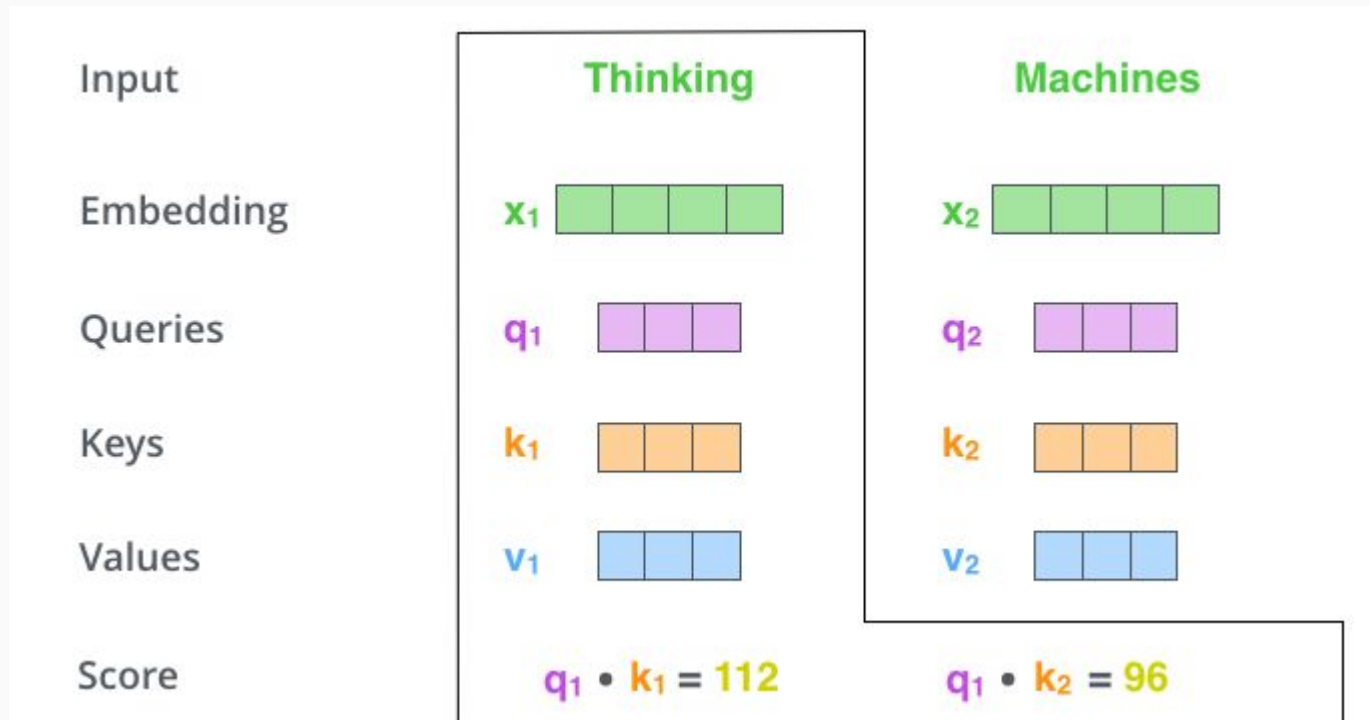
$$[\mathbf{v}^{(t)}]_V = [\mathbf{W^V}]_{VX} [\mathbf{x}^{(t)}]_X$$
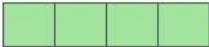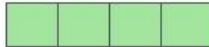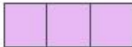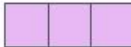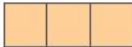
# Transformers: Self-attention

Then we dot-multiplicate each query with each key:

$$[q^{(t)}]_Q \cdot [k^{(\tau)}]_Q$$



| Input | **Thinking** | **Machines** |
|---|---|---|
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \bullet k_1 = 112$ | $q_1 \bullet k_2 = 96$ |

# Transformers: Self-attention

Then we rescale dividing by sqrt($Q$) and apply softmax. $Q$ is the number of dimensions of the query $[\mathbf{q}^{(t)}]_Q$ and the key $[\mathbf{k}^{(\tau)}]_Q$ (both vectors need to have the same size to allow for the dot product)

| Input | Thinking | Machines |
|---|---|---|
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ |
| Divide by 8 ( $\sqrt{d_k}$ ) | 14 | 12 |
| Softmax | 0.88 | 0.12 |

4

# Transformers: Self-attention

The result is a weight given to each value to form the new hidden state:

$$[\mathbf{z}^{(t)}]_Z = \sum_\tau [\, sm_\tau (\, [\mathbf{q}^{(t)}]_Q \cdot [\mathbf{k}^{(\tau)}]_Q \,) / \sqrt{Q} \,]_1 [\mathbf{v}^{(\tau)}]_Z$$



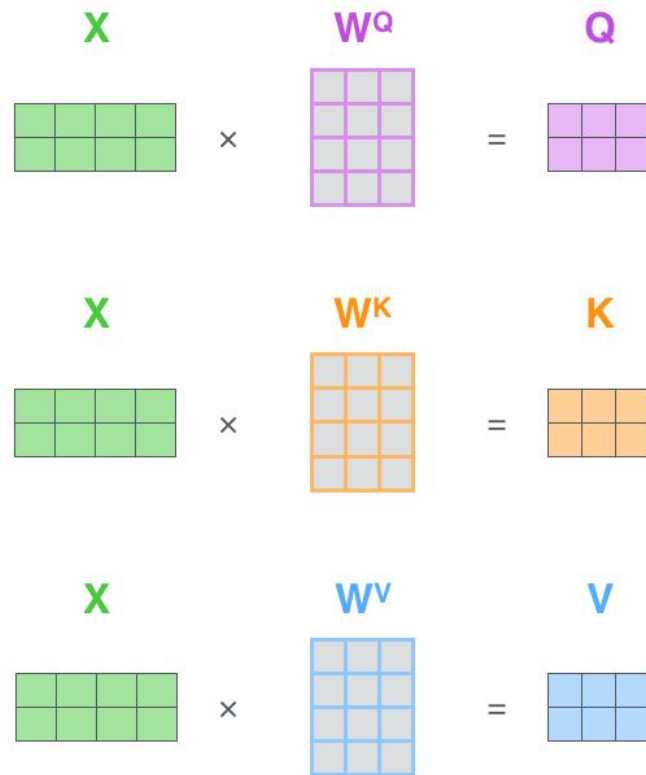| | Thinking | Machines |
|---|---|---|
| Input | | |
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \bullet k_1 = 112$ | $q_1 \bullet k_2 = 96$ |
| Divide by 8 ($\sqrt{d_k}$) | 14 | 12 |
| Softmax | 0.88 | 0.12 |
| Softmax X Value | $v_1$ | $v_2$ |
| Sum | $z_1$ | $z_2$ |

# Transformers: Self-attention

Computationally, we do all these calculations in parallel by using matrices rather than vectors. A matrix represents all the time steps in one go:

$$[Z]_{TZ} = [\text{ sm}( [Q]_{TQ} [K^*]_{QT} ) / \sqrt{Q} ]_{TT} [V]_{TZ}$$

This is extremely efficient, because modern computers (and specially GPUs) have hardware optimized to perform these operations very fast.
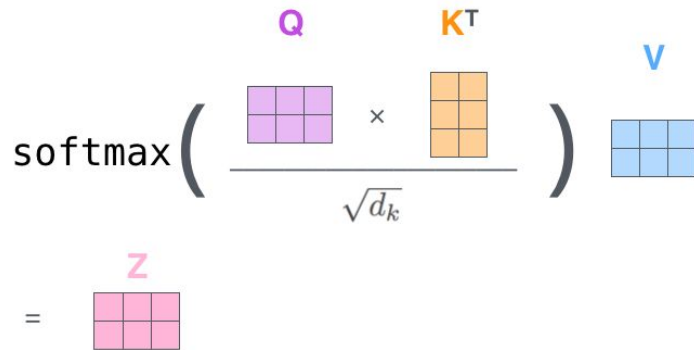
# Transformers: Self-attention

Computationally, we do all these calculations in parallel by using matrices rather than vectors. A matrix represents all the time steps in one go:
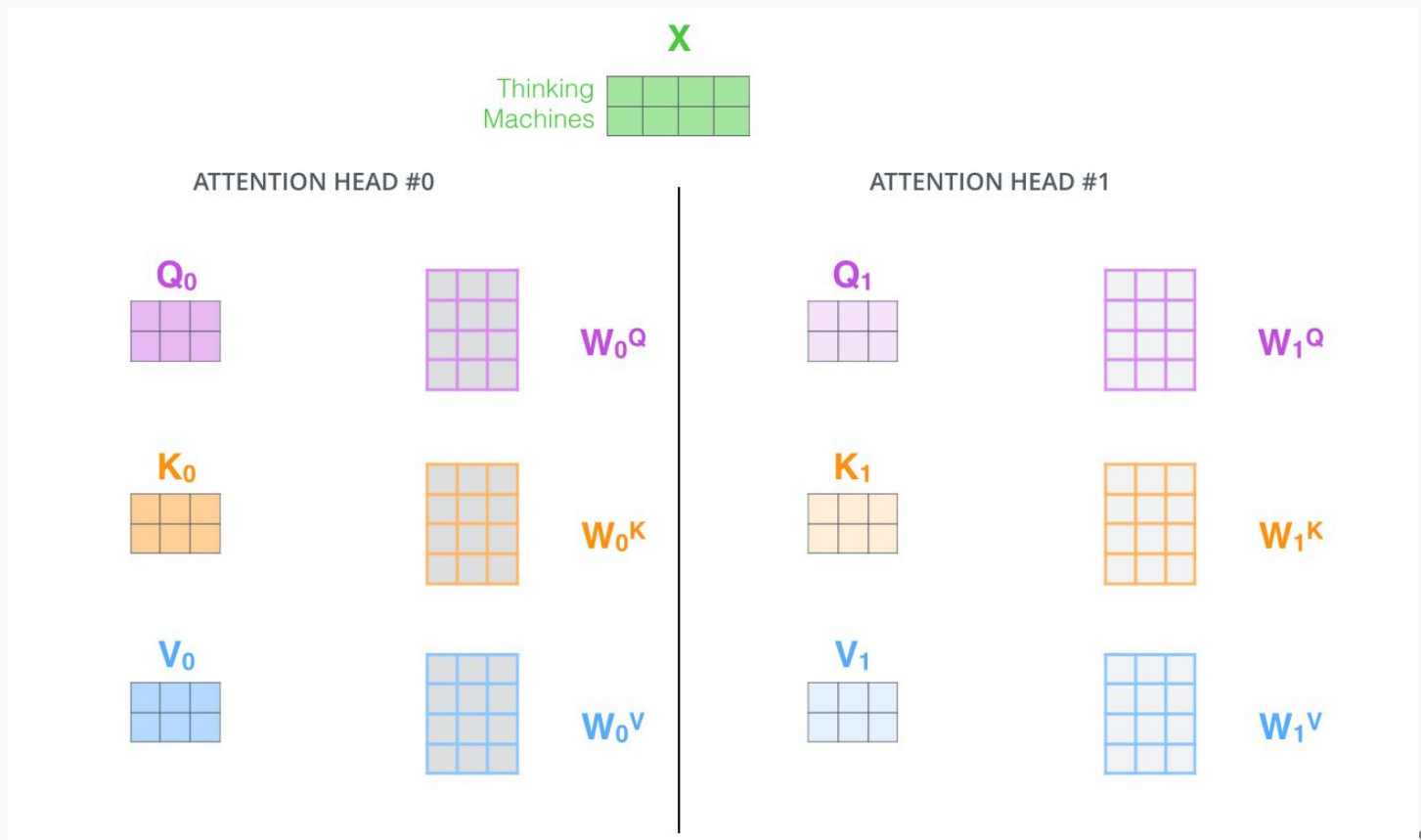
$$[Z]_{TZ} = [\, \text{sm}(\, [Q]_{TQ}\, [K^*]_{QT}\, )\, /\sqrt{Q}\, ]_{TT}\, [V]_{TZ}$$

This is extremely efficient, because modern computers (and specially GPUs) have hardware optimized to perform these operations very fast.
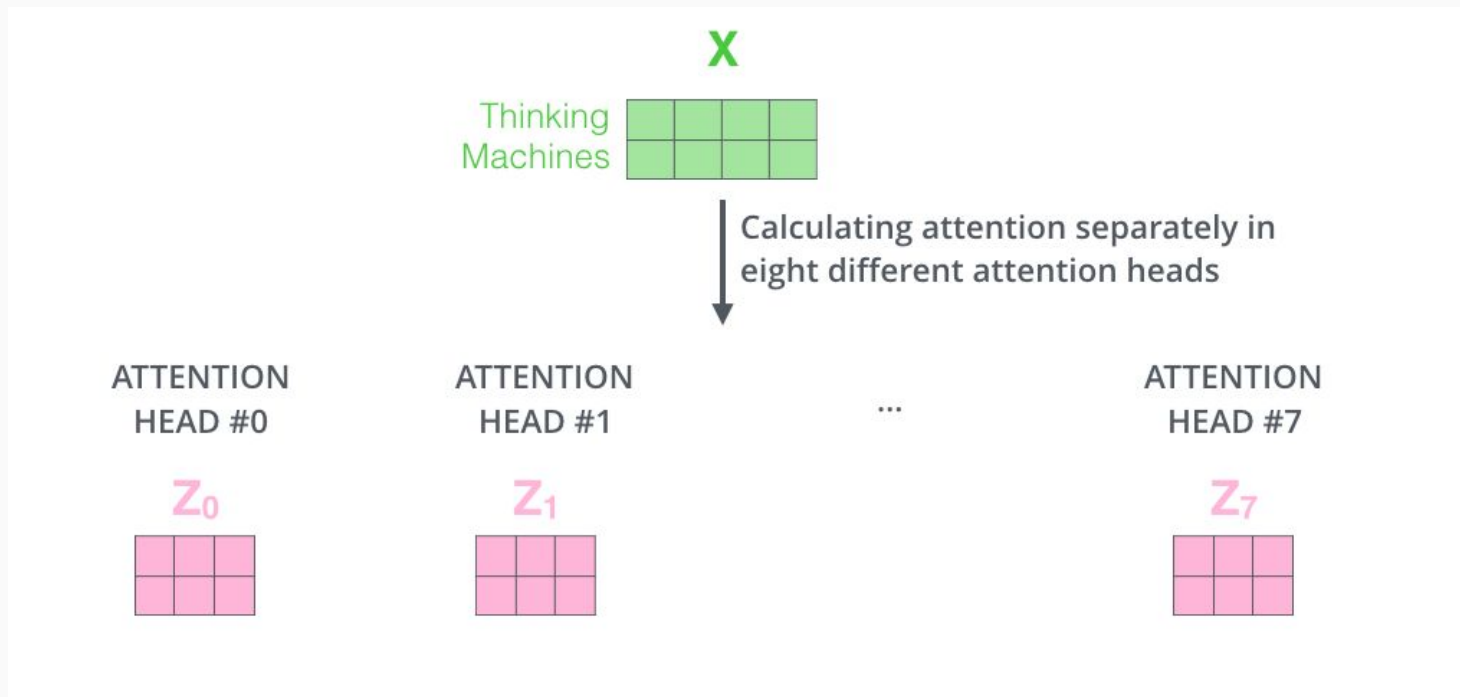
# Transformers: Self-attention

Another novelty of the transformer, is that it uses several attention channels in parallel. They are called 'heads':

# Transformers: Self-attention

Another novelty of the transformer, is that it uses several attention channels in parallel. They are called 'heads', and their results are $[Z^1]_{TZ}$, $[Z^2]_{TZ}$, $[Z^3]_{TZ}$, ..., $[Z^H]_{TZ}$ (H = number of heads):

# Transformers: Self-attention

But we so many heads we end up with too many matrixes. To prevent the hidden states to group exponentially in size, we pool all heads with a linear transformation.

1) Concatenate all the attention heads

$Z_0$  $Z_1$  $Z_2$  $Z_3$  $Z_4$  $Z_5$  $Z_6$  $Z_7$

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

X

$W^O$

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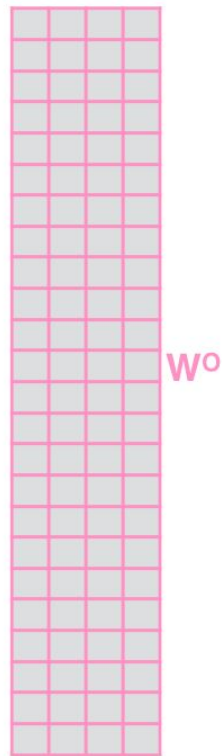
Z

=

$$\text{concat}( [Z^1]_{TZ}, [Z^2]_{TZ}, [Z^3]_{TZ}, ..., [Z^H]_{TZ} ) = [Z]_{T(Z \times H)}$$

$$[Z]_{T(Z \times H)} \times [W^O]_{(Z \times H)Z} = [Z^{\text{next layer}}]_{TZ}$$
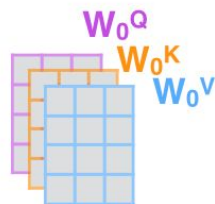
# Transformers: Self-attention
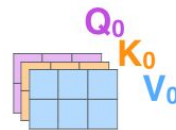


1) This is our input sentence*

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

Thinking Machines

$X$

* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

$R$

$W_0^Q$
$W_0^K$
$W_0^V$

$W_1^Q$
$W_1^K$
$W_1^V$

$W_7^Q$
$W_7^K$
$W_7^V$

$Q_0$
$K_0$
$V_0$

$Q_1$
$K_1$
$V_1$

$Q_7$
$K_7$
$V_7$

$Z_0$

$Z_1$

$Z_7$

$W^O$

$Z$

# Transformers: Self-attention
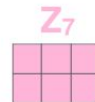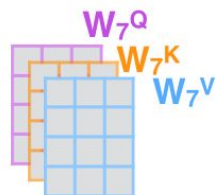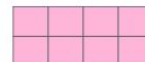
The effect of attention is that $[\mathbf{z}^{(t)}]_Z$ (for a given time step 't') becomes a mixture of the original embeddings $[\mathbf{x}^{(t)}]_X$. The idea is that each time step 'borrows' information from other time steps that it is related to. For instance, a pronoun may borrow information from the complement of the name that it refers to. Layer, the feed-forwards NN further transforms $[\mathbf{z}^{(t)}]_Z$ into $[\mathbf{r}^{(t)}]_R$, but this time without mixing information across time steps.

Multi-heads = It does all of this several times, with different with a different $[\mathbf{W^Q}]_{QX}$, $[\mathbf{W^K}]_{KQ}$, $[\mathbf{W^V}]_{VX}$ per head



22

# Transformers: Self-attention

# Transformers: Encoding position

We also add information about the position of each word. We do this with a positional encoding vector $[\mathbf{t}^{(t)}]_X$ per possible position.

# Transformers: Encoding position

➢ This is a similar idea that we have used before for conference resolution, where we concatenated extra features to the word embeddings.

➢ This is a common trick in NN NLP

➢ However the transformer e-wise multiplies $[\mathbf{t}^{(t)}]_X$ rather than concatenating it to $[\mathbf{x}^{(t)}]_X$.

# Transformers: Encoding position

➢ The positional encoding vectors have pre-specified values

➢ These values follow some sort of wavelet function

➢ This is quite similar to how the hippocampus in the human brain encodes position!

# The real brain: Encoding position



Position of mouse

One place cell's firing pattern

One grid cell's firing pattern

Path of mouse

Firing patterns of multiple place cells

Firing patterns of multiple grid cells

**Mapping One Location**

If a mouse is in one corner of a room, then there is one place cell that fires uniquely at that location. A grid cell that fires at that location also fires at other positions around it in a hexagonal array.

**Mapping a Path**

As the mouse moves, the activity of many place cells records the locations that it visited. Grid cell activity tracks how the mouse moved through overlapping hexagonal coordinate systems that tile the plane.

Place cell

Grid cell

# Transformers: Encoder

To simplify the explanation, we have so far ignored two smaller details of the architecture.
- 1) There is a **residual** connection bypassing each layer
- 2) There is a normalization step after each layer.

# Transformers: Encoder

To simplify the explanation, we have so far ignored two smaller details of the architecture.

- 1) There is a **residual** connection bypassing each layer
- 2) There is a normalization step after each layer.

# Transformers: Encoders → decoders

Besides self-attention, the decoder also uses encoder-decoder attention. This attention is the same as simple self-attention, but it also uses the outputs of the last layer of the encoder

# Transformers: Encoders → decoders

# Transformers: Encoders → decoders

# Transformers: Output

The NN is trained in a language model task. Remember from lecture NLP 4, this consists on predicting the next word.

The output of the neural network tries to find the 1-hot representation of the next word

Once trained, the transformer can be re-used in many other NLP tasks

Which word in our vocabulary is associated with this index?

am

Get the index of the cell with the highest value (argmax)

5

log_probs

0 1 2 3 4 5 ... vocab_size

Softmax

logits

0 1 2 3 4 5 ... vocab_size

Linear

Decoder stack output

# Transformers: Output

# Transformers: Output

After training the model, its outputs will approximate the desired 1-hot representations of words in the vocabulary



**Trained Model Outputs**

| Output Vocabulary: | a | am | I | thanks | student | <eos> |
|---|---|---|---|---|---|---|
| position #1 | 0.01 | 0.02 | 0.93 | 0.01 | 0.03 | 0.01 |
| position #2 | 0.01 | 0.8 | 0.1 | 0.05 | 0.01 | 0.03 |
| position #3 | 0.99 | 0.001 | 0.001 | 0.001 | 0.002 | 0.001 |
| position #4 | 0.001 | 0.002 | 0.001 | 0.02 | 0.94 | 0.01 |
| position #5 | 0.01 | 0.01 | 0.001 | 0.001 | 0.001 | 0.98 |
| | a | am | I | thanks | student | <eos> |

**Target Model Outputs**

| Output Vocabulary: | a | am | I | thanks | student | <eos> |
|---|---|---|---|---|---|---|
| position #1 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 |
| position #2 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| position #3 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| position #4 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 |
| position #5 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 |
| | a | am | I | thanks | student | <eos> |

# Transformers: Existing systems

All of these models are Transformer architecture models

**ULMfit**

Jan 2018

Training:

1 GPU day

**GPT**

June 2018

Training

240 GPU days

**BERT**

Oct 2018

Training

256 TPU days

~320–560 GPU days

**GPT-2**

Feb 2019

Training

~2048 TPU v3 days according to a reddit thread

**fast.ai**

**OpenAI**

**Google AI**

**OpenAI**

# Transformers: Existing systems

| Rank | Model | EM | F1 |
|------|-------|-----|-----|
| | Human Performance<br>*Stanford University*<br>(Rajpurkar & Jia et al. '18) | 86.831 | 89.452 |
| 1<br>Jan 15, 2019 | BERT + MMFT + ADA (ensemble)<br>*Microsoft Research Asia* | 85.082 | 87.615 |
| 2<br>Jan 10, 2019 | BERT + Synthetic Self-Training<br>(ensemble)<br>*Google AI Language*<br>https://github.com/google-research/bert | 84.292 | 86.967 |
| 3<br>Dec 13, 2018 | BERT finetune baseline (ensemble)<br>*Anonymous* | 83.536 | 86.096 |
| 4<br>Dec 16, 2018 | Lunet + Verifier + BERT (ensemble)<br>*Layer 6 AI NLP Team* | 83.469 | 86.043 |
| 4<br>Dec 21, 2018 | PAML+BERT (ensemble model)<br>*PINGAN GammaLab* | 83.457 | 86.122 |
| 5<br>Dec 15, 2018 | Lunet + Verifier + BERT (single<br>model)<br>*Layer 6 AI NLP Team* | 82.995 | 86.035 |

| Rank | Model | EM | F1 |
|------|-------|-----|-----|
| | Human Performance<br>*Stanford University*<br>(Rajpurkar & Jia et al. '18) | 86.831 | 89.452 |
| 1<br>Jan 10, 2020 | Retro-Reader on ALBERT (ensemble)<br>*Shanghai Jiao Tong University*<br>http://arxiv.org/abs/2001.09694 | 90.115 | 92.580 |
| 2<br>Nov 06, 2019 | ALBERT + DAAF + Verifier (ensemble)<br>*PINGAN Omni-Sinitic* | 90.002 | 92.425 |
| 3<br>Sep 18, 2019 | ALBERT (ensemble model)<br>*Google Research & TTIC*<br>https://arxiv.org/abs/1909.11942 | 89.731 | 92.215 |
| 4<br>Jan 23, 2020 | albert+transform+verify (ensemble)<br>*qianxin* | 89.528 | 92.059 |
| 5<br>Dec 08, 2019 | ALBERT+Entailment DA (ensemble)<br>*CloudWalk* | 88.761 | 91.745 |
| 6<br>Feb 20, 2020 | Tuned ALBERT (ensemble model)<br>*Group Data & Analytics Cell | Aditya Birla<br>Group)*<br>https://www.adityabirla.com/About/group-data-and-analytics | 88.637 | 91.230 |

# Course structure

➢ **Question Answering:** X

- Task definition, datasets, cloze-style tasks, Attentive Reader ← Lecture 15 (NLP 7)

➢ **Conference Resolution:** X

- Task definition, pairs method, clustering method, language models ← Lecture 15 (NLP 7)

➢ **Convolutional Neural Networks:** X

- CNNs in vision, CNNs in language, example ← Lecture 15 (NLP 7)

➢ **Transformers:** X

- Architecture: encoder, self-attention, encoding position, decoder ← Lecture 16 (NLP 8)
- Existing systems. Ranking ← Lecture 16 (NLP 8)

# Literature

➢ **Papers** =

- Attention is all you need. https://arxiv.org/abs/1706.03762

- The illustrated transformer. http://jalammar.github.io/illustrated-transformer/

- Language Models are Unsupervised Multitask Learners.
  https://d4mucfpksywv.cloudfront.net/better-language-models/language-models.pdf

- Reformer: The Efficient Transformer. https://arxiv.org/abs/2001.04451

- Illustrating the reformer.
  https://towardsdatascience.com/illustrating-the-reformer-393575ac6ba0