

# DEEP NETWORK DEVELOPMENT

Imre Molnár PhD student, ELTE, Al Department ⊠ imremolnar@inf.elte.hu ⊕ curiouspercibal.github.io Tamás Takács PhD student, ELTE, AI Department ⊠ <u>tamastheactual@inf.elte.hu</u> ⊕ tamastheactual.github.io **Deep Network Development** 

# Lecture 8.



# Transformer Networks

Budapest, 4<sup>th</sup> April 2025





**3** Transformer Architecture



### Recap

#### **Sequential data**

Text – sequence of words / characters



Speech – sequence of signals / acoustic features



Video – sequence of images (frames)





- ··· GTGCATCTGACTCCTGAGGAGAAG ··· DNA
- $\begin{array}{c} & & \\$

#### **Previous Lecture**

## Recap

### Key properties of sequential data

- Order
- Temporal information
- Context







Key model properties:

- Handle variable length sequences
- Track long term dependencies
- Maintain the order of the input
- Share parameters across sequences







## **Recurrent Neural Network (review)**

### **RNN (unrolled version)**



# **Embeddings: Word representation (review)**

### Featurized representation: Word Embedding

If we subtract man and woman, main difference is gender We can compute word similarities

We can compute word analogies: man is to woman as king is to \_\_\_\_\_





Male-Female

Verb tense

	Man (5391)	Woman (9853)	King (4914)	Queen (7151)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.7	0.68	0.03	-0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97



# **Embeddings: Word representation (review)**

### Featurized representation: Word Embedding

The word embeddings are learned with training.

Therefore, in practice, the features aren't that understandable.

We can visualize lower representations of the embeddings with techniques such as T-SNE





## **Embeddings: Word representation (review)**

RNN (unrolled version)





## Long Short-Term Memory (LSTM) (review)

### LSTM – a more complex network





# Long Short-Term Memory (LSTM) (review)

### LSTM – a more complex network

Activations

- a = h
- $a_t = h_t$

Cell state (memory)

- C
- C~



### Long Short-Term Memory (LSTM) (review) LSTM – Forget Gate





# Long Short-Term Memory (LSTM) (review)

LSTM – Input / Ignore Gate





### Long Short-Term Memory (LSTM) (review) LSTM – Output Gate



#### Previous Lecture









Copy

ncatenate

#### Gates:

$$i_t = \sigma(W_{hi}h_{t-1} + W_{xi}x_t + b_i)$$
  

$$f_t = \sigma(W_{hf}h_{t-1} + W_{xf}x_t + b_f)$$
  

$$o_t = \sigma(W_{ho}h_{t-1} + W_{xo}x_t + b_o)$$
  

$$\widetilde{C}_t = \tanh(W_{hc}h_{t-1} + W_{xc}x_t + b_c)$$

Vector

Transfer

#### Outputs:

$$C_t = f_t \circ C_{t-1} + i_t \circ \widetilde{C}_t$$
$$h_t = o_t \circ \tanh(\widetilde{C}_t)$$

where ° is element-wise multiplication operation



## Long Short-Term Memory (LSTM) (review)

### LSTM – a more complex network



Vanilla RNN:

![](_page_14_Figure_6.jpeg)

![](_page_15_Picture_1.jpeg)

## **Gated Recurrent Unit (GRU) (review)**

![](_page_15_Figure_3.jpeg)

![](_page_15_Figure_4.jpeg)

![](_page_16_Picture_1.jpeg)

## **Bidirectional RNN (review)**

Idea: use 2 independent recurrent models together.

- Input is fed in the proper time order to the first one, and in reverse time order to the second one.
- Outputs are combined at each time step using concatenation or summation.

![](_page_16_Figure_6.jpeg)

![](_page_17_Picture_1.jpeg)

### **Bidirectional RNN (review) BIDIRECTIONAL RNN (BRNN)** $h_1$ $h_0$ ... **A BRNN layer:** $\overrightarrow{\boldsymbol{h}_{d,t}^{(n)}} = \overrightarrow{g_d} \left( \overrightarrow{\boldsymbol{W}_d^h} \overrightarrow{\boldsymbol{h}_{d,t-1}^{(n)}} + \overrightarrow{\boldsymbol{W}_d^x} \overrightarrow{\boldsymbol{h}_{d-1,t}^{(n)}} + \overrightarrow{\boldsymbol{b}_d} \right)$ $\overleftarrow{\boldsymbol{h}_{d,t}^{(n)}} = \overleftarrow{g_d} \left( \overleftarrow{\boldsymbol{W}_d^h} \overleftarrow{\boldsymbol{h}_{d,t+1}^{(n)}} + \overleftarrow{\boldsymbol{W}_d^x} \overleftarrow{\boldsymbol{h}_{d-1,t}^{(n)}} + \overleftarrow{\boldsymbol{b}_d} \right)$ $oldsymbol{h}_{d,t}^{(n)} = g_d \left( \overrightarrow{oldsymbol{W}_d} \overrightarrow{oldsymbol{h}_{d,t}^{(n)}} + \overleftarrow{oldsymbol{W}_d} \overleftarrow{oldsymbol{h}_{d,t}^{(n)}} + oldsymbol{b}_d ight)$ $x_0$ $x_1$ . . . Where

• " $\rightarrow$ " means normal time order, while " $\leftarrow$ " is associated with reverse time order

$$t = 1, ..., T_n; \theta_d = \left\{ \overrightarrow{W_d^h}, \overrightarrow{W_d^x}, \overrightarrow{b_d}, \overleftarrow{W_d^h}, \overleftarrow{W_d^x}, \overleftarrow{b_d}, \overrightarrow{W_d}, \overleftarrow{W_d}, \overleftarrow{b_d}, \overrightarrow{W_d}, \overleftarrow{b_d} \right\}$$

**Deep Network Development** 

# Lecture 8.

![](_page_18_Picture_2.jpeg)

# Transformer Networks

Budapest, 4<sup>th</sup> April 2025

![](_page_18_Picture_5.jpeg)

![](_page_18_Picture_6.jpeg)

**3** Transformer Architecture

### 1. Natural Language Understanding

## **Natural Language Understanding**

![](_page_19_Picture_2.jpeg)

![](_page_19_Picture_3.jpeg)

-

## **Machine Translation**

- Most of the proposed neural machine translation models belog to the a family of encode-decoders
- From probabilistic perspective:

 $argmax_{y}p(y \mid x)$ 

• Which means finding a target **y** that maximizes the conditional probability of **y** given a source sentence **x** 

![](_page_20_Picture_6.jpeg)

![](_page_20_Picture_7.jpeg)

Machine Translation

## **Machine Translation**

- In the **encoder-decoder** framework an **encoder** reads and reduces the input sentence (a sequence of vectors  $\mathbf{x} = (x_1, \dots, x_T)$ ) into a vector c
- The approach for RNNs:

 $h_t = f(x_t, h_{t-1})$  $c = q(\{h_1, \dots, h_{T_x}\})$ 

- Where q is a nonlinear function,  $h_t$  is the hidden state at time t
- If we choose f as an LSTM then  $q(\{h_1, ..., h_{T_x}\}) = h_{T_x}$
- The decoder then predicts the next word  $y_t$  given the context vector c and all the previously predicted words  $\{y_1, \dots, y_{t-1}\}$

![](_page_21_Picture_8.jpeg)

![](_page_21_Picture_9.jpeg)

Machine Translation

## **Machine Translation**

- The **decoder** then predicts the next word  $y_t$  given the context vector c and all the previously predicted words  $\{y_1, \dots, y_{t-1}\}$
- The predicted translation y can be expressed as the joint probability of:

$$p(\mathbf{y}) = \prod_{t=1}^{T} p(y_t \mid \{y_1, \dots, y_{t-1}\}, c)$$

- where  $\mathbf{y} = \{y_1, ..., y_{T_y}\}$
- Whit and RNN this is modelled as:

 $p(y_t \mid \{y_1, \dots, y_{t-1}\}, c) = g(y_{t-1}, s_t, c)$ 

• Where g nonlinear multi layered function and  $s_t$  is the hidden state

![](_page_22_Picture_9.jpeg)

![](_page_22_Picture_10.jpeg)

Machine Translation

### 1. Natural Language Understanding

![](_page_23_Picture_1.jpeg)

## Encoder – Decoder (Seq2Seq)

![](_page_23_Figure_3.jpeg)

![](_page_24_Picture_1.jpeg)

### Encoder – Decoder (Seq2Seq)

• The encoder and decoder are nothing more than stacked RNN layers, such as LSTM's. The encoder processes the input and produces one compact representation, called z, from all the input timesteps. It can be regarded as a compressed format of the input.

![](_page_24_Figure_4.jpeg)

Encoder

![](_page_25_Picture_1.jpeg)

## Encoder – Decoder (Seq2Seq)

**Problems with this approach** 

- The intermediate representation z cannot encode information from all the input timesteps. This is commonly known as the bottleneck problem. The vector z needs to capture all the information about the source sentence.
- In practice, how far we can see in the past (the so-called reference window) is finite. RNN's tend to **forget information** from timesteps that are far behind.

![](_page_25_Figure_6.jpeg)

**Vanishing Gradient:** where the contribution from the earlier steps becomes insignificant in the gradient for the vanilla RNN unit.

![](_page_26_Picture_1.jpeg)

# **BLEU (bilingual evaluation understudy)** [1]

- **BLEU** is an algorithm for evaluating the **quality of text** which has been machine-translated from one natural language to another.
- It's a metric, that measures the precision of n-grams (such as bigrams or trigrams) between the predicted translation and the reference translations

![](_page_26_Picture_5.jpeg)

The children are sitting with their shoes off .

A man wearing a long white robe and a white cap , stands over a group of seated people .

A man in a white robe stands and talks to people who are sitting .

A man in a long white robe and white head cap is pointing to a group of people sitting on a blanket .

A man dressed in all white is talking to the people seated on the dirt floor .

![](_page_27_Picture_0.jpeg)

#### 1. Natural Language Understanding

## N-gram

- It is just a fancy way of describing "a set of 'n' consecutive words in a sentence".
- For instance, in the sentence "The ball is blue", we could have ngrams such as:
  - 1-gram (unigram): "The", "ball", "is", "blue"
  - 2-gram (bigram): "The ball", "ball is", "is blue"
  - 3-gram (trigram): "The ball is", "ball is blue"
  - 4-gram: "The ball is blue"

![](_page_27_Picture_9.jpeg)

The children are sitting with their shoes off .

A man wearing a long white robe and a white cap , stands over a group of seated people .

A man in a white robe stands and talks to people who are sitting .

A man in a long white robe and white head cap is pointing to a group of people sitting on a blanket .

A man dressed in all white is talking to the people seated on the dirt floor .

# Precision

- Let's say, that we have:
  - Target Sentence: He eats an apple
  - **Predicted Sentence:** He ate an apple
- We would normally compute the Precision using the formula:

 $Precision = \frac{Number of correct predicted words}{Number of total predicted words}$ 

- Here **eats** and **ate** count as different words
- Precision = 3 / 4

![](_page_28_Picture_10.jpeg)

![](_page_28_Picture_12.jpeg)

### <u>1. Natural Language</u> Understanding

## Precision

- The issue is that only using precision allow us to cheat:
  - Target Sentence: He eats an apple
  - Predicted Sentence: He he he
- Precision = 3 / 3

![](_page_29_Picture_6.jpeg)

# **Clipped Precision**

- Target Sentence 1: He eats a sweet apple
- Target Sentence 2: He is eating a tasty apple
- Predicted Sentence: He He He eats tasty fruit
- 1. We compare each word from the predicted sentence with all of the target sentences.
- 2. If the word matches any target sentence, it is considered to be correct.
- 3. We limit the count for each correct word to the maximum number of times that that word occurs in the Target Sentence.

![](_page_30_Picture_8.jpeg)

## Clipped Precision = $\frac{\text{Clipped # of correct predicted words}}{\text{# of total predicted words}}$

Word	Matching Sentence	Matched Predicted Count	Clipped Count
He	Both	3	1
eats	Target 1	1	1
tasty	Target 2	1	1
fruit	None	0	0
Total		5	3

Clipped Precision = **3** / **6** 

### 1. Natural Language Understanding

## How is BLEU score calculated?

- Target Sentence: The guard arrived late because it was raining
- **Predicted Sentence:** The guard arrived late because of the rain

 $p_1 = \frac{5}{6}$  $p_2 = \frac{4}{7}$ **Precision 2-gram: Precision 1-gram:** Target Sentence: The guard arrived late because it was raining The guard arrived late because it was raining Target Sentence: Predicted Sentence: The guard arrived late because of the rain Predicted Sentence: The guard arrived late because of the rain  $p_3 = \frac{3}{6}$ **Precision 3-gram:**  $p_2 = \frac{2}{r}$ **Precision 4-gram:** The guard arrived late because it was raining Target Sentence: The guard arrived late because it was raining Target Sentence: Predicted Sentence: The guard arrived late because of the rain Predicted Sentence: The guard arrived late because of the rain

![](_page_31_Picture_9.jpeg)

## How is BLEU score calculated?

### **Geometric Average Precision (GAP)**

• We combine these precision scores using the following formula

$$\mathbf{GAP}(N) = \prod_{n=1}^{N} p_n^w$$

- We usually use N = 4 and uniform weights  $w = \frac{1}{N}$
- In case of N = 4 the Geometric Average Precision:

 $\mathbf{GAP}(N) = (p_1)^{\frac{1}{4}} \cdot (p_2)^{\frac{1}{4}} \cdot (p_3)^{\frac{1}{4}} \cdot (p_1)^{\frac{1}{4}}$ 

![](_page_32_Picture_9.jpeg)

![](_page_32_Picture_11.jpeg)

![](_page_32_Picture_13.jpeg)

# How is BLEU score calculated?

### **Brevity Penalty**

• The Brevity Penalty penalizes sentences that are too short.

**BrevityPenalty** =  $\begin{cases} 1, & \text{if } c > r \\ e^{(1-\frac{r}{c})}, & \text{if } c \le r \end{cases}$ 

- c is the predicted length -> number of words in the predicted sentence
- r is the target length -> number of words in the target sentence (length closes to prediction length)
- In the previous example c = 8 and r = 8

![](_page_33_Picture_11.jpeg)

![](_page_33_Picture_12.jpeg)

## How is BLEU score calculated?

 To get the BLEU score we multiply the Brevity Penalty with the Geometric Average of the Precision Scores

 $BLEU(N) = BrevityPenalty \cdot GAP(N)$ 

- **BLEU** Score can be computed for different values of N.
  - BLEU-1 uses the unigram Precision score
  - **BLEU-2** uses the geometric average of **unigram** and **bigram** precision
  - BLEU-3 uses the geometric average of **unigram**, **bigram**, and **trigram** precision
  - and so on.

![](_page_34_Picture_10.jpeg)

**Deep Network Development** 

# Lecture 8.

![](_page_35_Picture_2.jpeg)

# **Attention Mechanism**

Budapest, 4<sup>th</sup> April 2025

1 Natural Language Understanding

**2** Attention Mechanism

**3** Transformer Architecture


#### 2. Attention Mechanism

#### **Introduction to Attention Mechanism**

- All the information from the encoder is represented in the z vector (context)
- However, as seen previously, information from earlier timestamps is not preserved
- Can we create a better context vector?



#### 2. Attention Mechanism



#### **Introduction to Attention Mechanism**

• In the new approach, we define each conditional probability as:

$$p(y_i | \{y_1, \dots, y_{i-1}\}, c) = g(y_{i-1}, s_i, c_i)$$

• Where  $s_i$  is an RNN hidden state for time *i*, computed by

 $s_i = f(s_{i-1}, y_{i-1}, c_i)$ 

• Note that here we produce a distinct context vector  $c_i$  for each target word  $y_i$ 





### **Introduction to Attention Mechanism (review)**

- The core idea is that the context vector c<sub>i</sub> should have access to all parts of the input sequence instead of just the last one.
- In other words, we need to form a **direct connection** with each timestamp. We can look at all the different words at the same time and learn to "pay attention" to the correct ones depending on the task at hand.
- In the encoder-decoder:
  - Given the hidden states of the encoder at each time step  $h = \{h_1, \dots, h_{T_x}\}$
  - Given the previous state in the decoder  $y_{i-1}$
  - We define an attention network that gives the attention scores for the current state of the decoder

$$e_{ij} = \operatorname{attention}_{\operatorname{net}}(y_{i-1}, \mathbf{h}) \in \mathbb{R}^n$$

• We convert this scores into probabilities

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$$

• Finally, we get our new context vector z:

 $c_i = \sum_{i=1}^{T_x} \alpha_{ij} h_j$ 



### **Introduction to Attention Mechanism (review)**

- Attention mechanism helps creating a better context vector
- It learns which information from the encoder is relevant for the decoder



### **Introduction to Attention Mechanism (review)**

- Attention mechanism helps creating a better context vector
- It learns which information from the encoder is relevant for the decoder





Important





#### **Attention Mechanism (review)**



Deep Network Development

2. Attention Mechanism



### **Attention Mechanism (review)**



 $\alpha_{ii}$  = amount of "attention"  $y_i$  should pay to  $h_*$ 



### **Attention Mechanism (review)**



Deep Network Development

#### 2. Attention Mechanism

#### **How does Attention work?**









#### **Attention Mechanism**



#### Several ways to calculate the scores

Name	Alignment score function	Citation
Content-base attention	$ ext{score}(oldsymbol{s}_t,oldsymbol{h}_i) =  ext{cosine}[oldsymbol{s}_t,oldsymbol{h}_i]$	Graves2014
Additive(*)	$ ext{score}(oldsymbol{s}_t,oldsymbol{h}_i) = \mathbf{v}_a^ op  ext{tanh}(\mathbf{W}_a[oldsymbol{s}_t;oldsymbol{h}_i])$	Bahdanau2015
Location- Base	$lpha_{t,i} =  ext{softmax}(\mathbf{W}_a \mathbf{s}_t)$ Note: This simplifies the softmax alignment to only depend on the target position.	Luong2015
General	$\operatorname{score}(\boldsymbol{s}_t, \boldsymbol{h}_i) = \boldsymbol{s}_t^\top \mathbf{W}_a \boldsymbol{h}_i$ where $\mathbf{W}_a$ is a trainable weight matrix in the attention layer.	Luong2015
Dot-Product	$ ext{score}(oldsymbol{s}_t,oldsymbol{h}_i)=oldsymbol{s}_t^{ op}oldsymbol{h}_i$	Luong2015
Scaled Dot- Product(^)	$\operatorname{score}(s_t, h_i) = \frac{s_t^{T} h_i}{\sqrt{n}}$ Note: very similar to the dot-product attention except for a scaling factor; where n is the dimension of the source hidden state.	Vaswani2017

**Deep Network Development** 

## Lecture 8.



# Transformer Architecture

Budapest, 4<sup>th</sup> April 2025

1 Natural Language Understanding





#### 3. Transformer Architecture

### **Transformers** [1]

- Multi-Head Attention
- Self-Attention





[1] Vaswani, Ashish et al. "Attention is All you Need." (2017)



#### **Transformer Attention**

- Feature-based attention: The Key, Value, and Query
- Key-value-query concepts come from information retrieval systems.
- Example of searching for a video on YouTube:
  - When you search (query) for a particular video, the search engine will map your query against a set of keys (video title, description, etc.) associated with possible stored videos. Then the algorithm will present you the best-matched videos (values). This is the foundation of content/feature-based lookup.
- Bringing this idea closer to the transformer's attention we have something like this:



# **Transformer Attention**





### **Self-Attention**



We can also define the attention of the same sequence, called <u>self-attention</u>. Instead of looking for an input-output sequence association/alignment, we are now looking for scores between the elements of the sequence

Hello love you Hello 0.8 0.1 0.05 0.05 0.6 0.2 0.1 0.1 0.05 0.2 0.65 0.1 ove 0.2 0.6 you 0.1 0.1



#### 

### **Visualizing Self-Attention**



#### **Transformers Attention**

$$A(q, K, V) = \sum_{i} \frac{\exp(q \cdot k^{\langle i \rangle})}{\sum_{j} \exp(q \cdot k^{\langle j \rangle})} v^{\langle i \rangle}$$

Embedded	tensor	in Q space

Self Attention	walk	by	river	bank	
walk	0.88	0.1	0.21	0.12	
by	0.08	0.8	0.12	0.14	
river	0.03	0.14	0.83	0.7	
bank	0.01	0.01	0.4	0.77	

### **Self-Attention**



Layer: 5 🖨 Attent	ion: Input - Input 🔶	
The_		The_
animal_		animal_
didn_		didn_
<u></u>		<u></u>
t_		t_
cross_		cross_
the_		the_
street_		street_
because_		because_
it_		it_
was_		was_
too_		too_
tire		tire
d_		d_

4/28/2025

**Summary** 

3. Transformer Architecture

- How does the input elements relate to the output (what part of the input is relevant for the output)
  - Helps the model in focusing on the most relevant information



### Summary

#### 2. Self-Attention

- How do the different elements of the input relate to each other (what part of the input is relevant in understanding the other parts of the input)
  - Helps in capturing dependencies and contextual information within the input





#### 3. Transformer Architecture

### Summary

#### **3. Multi-Head Attention**

- Repeating self-attention multiple times to better understand the input (multiple ways of deciding the relationship between the parts of the input)
  - Looks at the input from different perspectives





#### 1st Head Attention

Linear

Concat

Scaled Dot-Product Attention

Linear

Κ

Linear

Q

Linear

V



 $\begin{array}{l} <\!\!\mathrm{SOS}\!\!>\!x^{<\!1\!\!>}\;x^{<\!2\!\!>}\;\ldots\;x^{<\!T_x\!-\!1\!\!>}\;x^{<\!T_x\!\!>}<\!\!\mathrm{EOS}\!\!>\\ \mathrm{Jane\;visite\;l'Afrique\;en\;septembre} \end{array}$ 





























<SOS>



<SOS>





<SOS>



















#### 3. Transformer Architecture

- When you convert a sequence into a set (tokenization), you lose the notion of order.
- **Positional encoding** is a set of small constants, which are added to the word embedding vector before the first self-attention layer.
- So if the same word appears in a different position, the actual representation will be slightly different, depending on where it appears in the input sentence.
- They use a sinusoidal function for the positional encoding. The sine function tells the model to pay attention to a particular wavelength.










#### **Transformer**







#### **Transformer**





### **Transformer Encoder**

#### Summary: the Transformer encoder

To process a sentence, we need these 3 steps:

- 1. Word embeddings of the input sentence are computed simultaneously.
- 2. Positional encodings are then applied to each embedding resulting in word vectors that also include positional information.
- 3. The word vectors are passed to the first encoder block.

Each block consists of the following layers in the same order:

- 1. A multi-head self-attention layer to find correlations between each word
- 2. A normalization layer
- 3. A residual connection around the previous two sublayers
- 4. A linear layer
- 5. A second normalization layer
- 6. A second residual connection

Note that the above block can be replicated several times to form the Encoder. In the original paper, the encoder composed of 6 identical blocks.





#### 3. Transformer Architecture

# **Transformer Decoder**

#### Transformer decoder: what is different?

The decoder consists of all the aforementioned components plus two novel ones. As before:

- 1. The output sequence is fed in its entirety and word embeddings are computed
- 2. Positional encoding are again applied
- 3. And the vectors are passed to the first Decoder block

#### Each decoder block includes:

- 1. A Masked multi-head self-attention layer
- 2. A normalization layer followed by a residual connection
- 3. A new multi-head attention layer (known as Encoder-Decoder attention)
- 4. A second normalization layer and a residual connection
- 5. A linear layer and a third residual connection

The decoder block appears again 6 times. The final output is transformed through a final linear layer and the output probabilities are calculated with the standard softmax function.





#### 3. Transformer Architecture Output Probabilities **Transformer** [1] Softmax Linear Add & Norm Feed Forward Add & Norm Add & Norm Multi-Head Feed Attention Forward N× Add & Norm N× Add & Norm Masked Multi-Head Multi-Head Attention Attention Positional Positional Encoding Encoding Input Output Embedding Embedding Inputs Outputs (shifted right) [1] Vaswani, Ashish et al. "Attention is All you Need." (2017)





#### **Transformer Architectures**





TYPES OF Machine learning PAPER		TYPE	APERS	
Baseline is all you need		HERE'S A NEW TASK WHERE OUR MODELS DON'T SUCCEED JUST YET	NEVER MIND. TURNS OUT WITH SOME CLEVER TRICKS, VE ALREADY GET SUPER- HUMAN PERFORMANCE	WE COMBINED TWO WELL KNOWN TECHNIQUES IN AN UNSURPRISING WAY
We got more compute and it works better We spent \$1M on compute and it looks really cool		TRANSFORMERS ALSO WORK ON THIS TYPE OF DATA	A TASK-SPECIFIC IMPROVEMENT THAT MAY OR MAY NOT WORK ON YOUR DATA	
We figured out how deep learning generalizes this time, I swear		NEURAL NETWORKS ARE LIKE THE BRAIN IN THIS SPECIFIC VAY, VHICH SLIGHTLY HELPS FOR OUR TASK	GOOD LUCK RUNNING A STATISTICAL SIGNIFICANCE TEST ON OUR RESULTS	CHECK OUT THESE NICE CHERRY-PICKED SAMPLES OF OUR MODEL
Results are 0.3% better than that other paper!		WE USED LOTS OF COMPUTE TO TRAIN A SLIGHTLY BETTER LANGUAGE MODEL	SOME THOUGHTS ON WHY THE WAY WE DO THINGS IS URONG	



# **Generative Pre-trained Transformer (GPT) – June 2018**

The original goal is Language Modelling (LM)

Uses Masked Self-Attention to limit the attention to the previous tokens only (left-to-right)

Two stage training:

- 1. Unsupervised pre-training:
  - The goal is to predict the next token based on the previous tokens.
- 2. Supervised fine-tuning:
  - Predict the label (y) based on the input tokens





## **Generative Pre-trained Transformer (GPT) – June 2018**



#### Bidirectional Encoder Representations from Transformers (BERT) – May 2019

- Compared to OpenAI GPT it uses a bidirectional self-attention
- Trained on 2 tasks at the same time during pre-training
  - 1. Masked LM (15% of the tokens are masked)
  - 2. Next Sentence Prediction
- A special [CLS] token is introduced at the beginning of each sequence.







Pre-training



#### Bidirectional Encoder Representations from Transformers (BERT) – May 2019



(a) Sentence Pair Classification Tasks: MNLI, QQP, QNLI, STS-B, MRPC, RTE, SWAG



(b) Single Sentence Classification Tasks: SST-2, CoLA



(d) Single Sentence Tagging Tasks: CoNLL-2003 NER



### **Further Links + Resources**

- <u>https://theaisummer.com/attention/</u>
- <u>https://theaisummer.com/transformer/</u>
- <u>https://jalammar.github.io/illustrated-transformer/</u>
- <u>https://nlp.seas.harvard.edu/annotated-transformer/</u> !!!!
- <u>https://www.youtube.com/watch?v=bCz40MemCcA</u>
- Coursera Deep Learning Specialization

Please read this:

- Beam search: <u>https://towardsdatascience.com/foundations-of-nlp-explained-visually-beam-search-how-it-works-1586b9849a24</u>
- BLEU score: <u>https://cloud.google.com/translate/automl/docs/evaluate#bleu</u>



#### Resources

Books:

- Courville, Goodfellow, Bengio: Deep Learning
  Freely available: <u>https://www.deeplearningbook.org/</u>
- Zhang, Aston and Lipton, Zachary C. and Li, Mu and Smola, Alexander J.: Dive into Deep Learning Freely available: <u>https://d2l.ai/</u>

Courses:

- Deep Learning specialization by Andrew NG
- <u>https://www.coursera.org/specializations/deep-learning</u>



# That's all for today!