# AI SUMMER CAMP
## TRANSFORMER & NLP

# Content

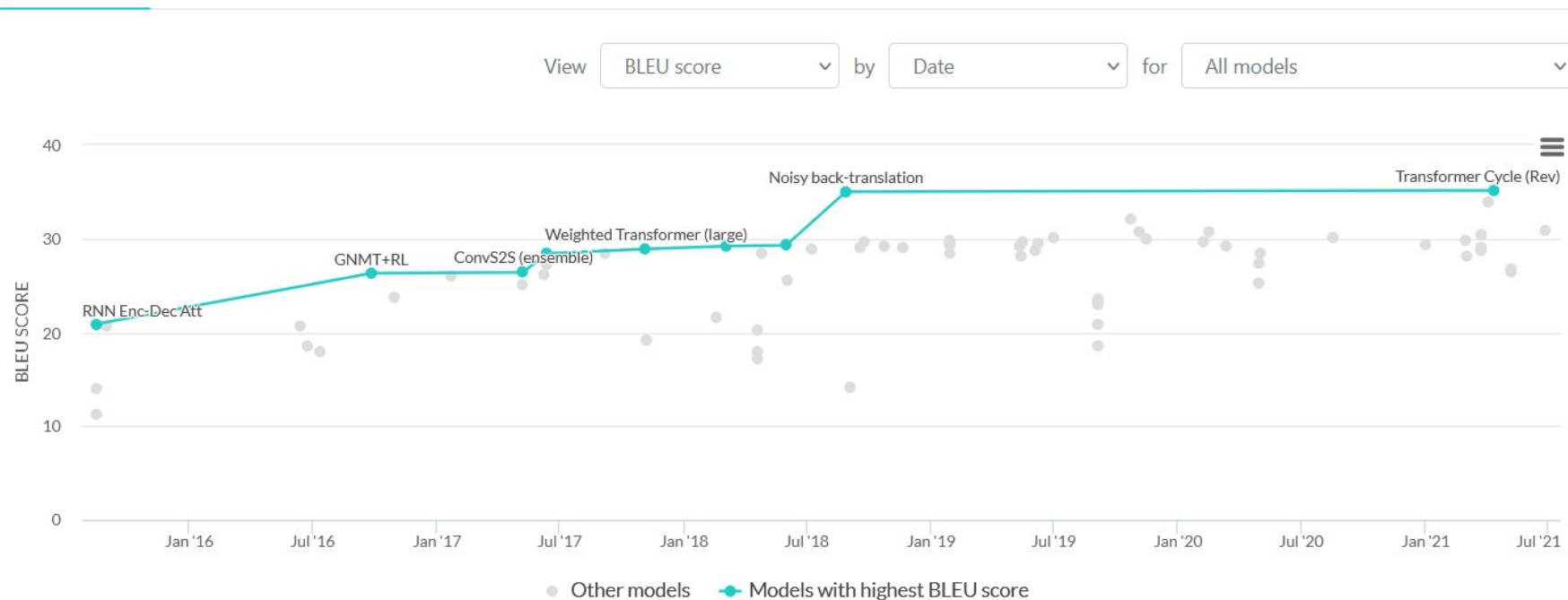# Natural Language Processing

Common tasks include:

- text classification

- translation

- summarization

- named entity recognition

- dialogue (chatbots)

- question answering

For more info visit https://paperswithcode.com/area/natural-language-processing

# Machine Translation on WMT2014 English-German

Leaderboard    Dataset

View  [BLEU score ▾]  by  [Date ▾]  for  [All models ▾]

- Noisy back-translation
- Transformer Cycle (Rev)
- Weighted Transformer (large)
- GNMT+RL
- ConvS2S (ensemble)
- RNN Enc-Dec Att

BLEU SCORE

○ Other models    ●—● Models with highest BLEU score

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 1 | **Transformer Cycle**<br>(Rev) | 35.14 | 33.54 | | ✓ | Lessons on Parameter Sharing across Layers in Transformers | ⌂ | →] | 2021 | Transformer |
| 2 | **Noisy back-translation** | 35.0 | 33.8 | | ✓ | Understanding Back-Translation at Scale | ⌂ | →] | 2018 | |
| 3 | **Transformer+Rep**<br>(Uni) | 33.89 | 32.35 | | ✓ | Rethinking Perturbations in Encoder-Decoders for Fast Training | ⌂ | →] | 2021 | Transformer |
| 4 | **T5-11B** | 32.1 | | | ✓ | Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer | ⌂ | →] | 2019 | Transformer |
| 5 | **Transformer + R-Drop** | 30.91 | | | ✗ | R-Drop: Regularized Dropout for Neural Networks | ⌂ | →] | 2021 | Transformer |
| 6 | **BERT-fused NMT** | 30.75 | | | ✗ | Incorporating BERT into Neural Machine Translation | ⌂ | →] | 2020 | Transformer |
| 7 | **Data Diversification - Transformer** | 30.7 | | | ✗ | Data Diversification: A Simple Strategy For Neural Machine Translation | ⌂ | →] | 2019 | Transformer |
| 8 | **Mask Attention Network**<br>(big) | 30.4 | | 215M | ✗ | Mask Attention Networks: Rethinking and Strengthen Transformer | | →] | 2021 | |
| 9 | **Transformer**<br>(ADMIN init) | 30.1 | 29.5 | | ✗ | Very Deep Transformers for Neural Machine Translation | ⌂ | →] | 2020 | Transformer |

# Transformer

## Attention Is All You Need

**Ashish Vaswani***
Google Brain
avaswani@google.com

**Noam Shazeer***
Google Brain
noam@google.com

**Niki Parmar***
Google Research
nikip@google.com

**Jakob Uszkoreit***
Google Research
usz@google.com

**Llion Jones***
Google Research
llion@google.com

**Aidan N. Gomez*** [†]
University of Toronto
aidan@cs.toronto.edu

**Łukasz Kaiser***
Google Brain
lukaszkaiser@google.com

**Illia Polosukhin*** [‡]
illia.polosukhin@gmail.com

### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

# Recurrent Neural Network (RNN)

## The Vanilla RNN Model

First time-step ($t = 1$):

$$\mathbf{h}_1 = tanh(W^{xh} \cdot \mathbf{x}_1 + W^{hh} \cdot \mathbf{h}_0)$$

$$\hat{\mathbf{y}}_1 = softmax(W^{hy} \cdot \mathbf{h}_1)$$

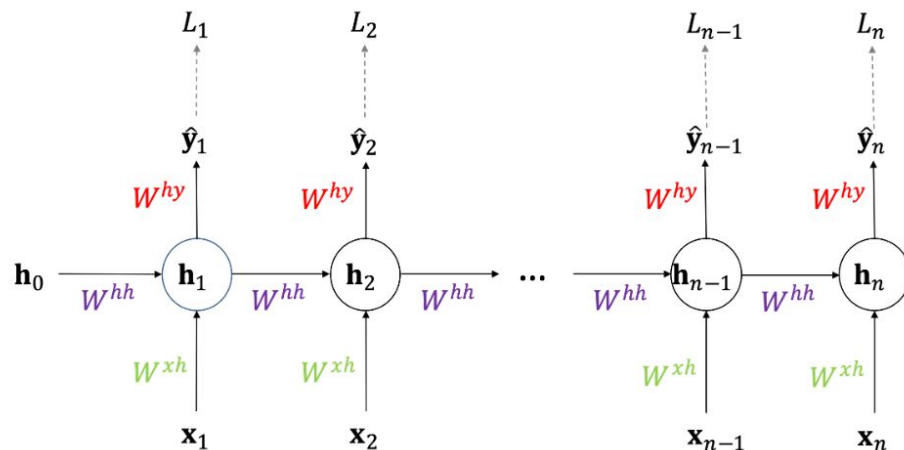$$L_1 = CE(\hat{\mathbf{y}}_1, \mathbf{y}_1)$$

In general:

$$\mathbf{h}_t = tanh(W^{xh} \cdot \mathbf{x}_t + W^{hh} \cdot \mathbf{h}_{t-1})$$

$$\hat{\mathbf{y}}_t = softmax(W^{hy} \cdot \mathbf{h}_t)$$

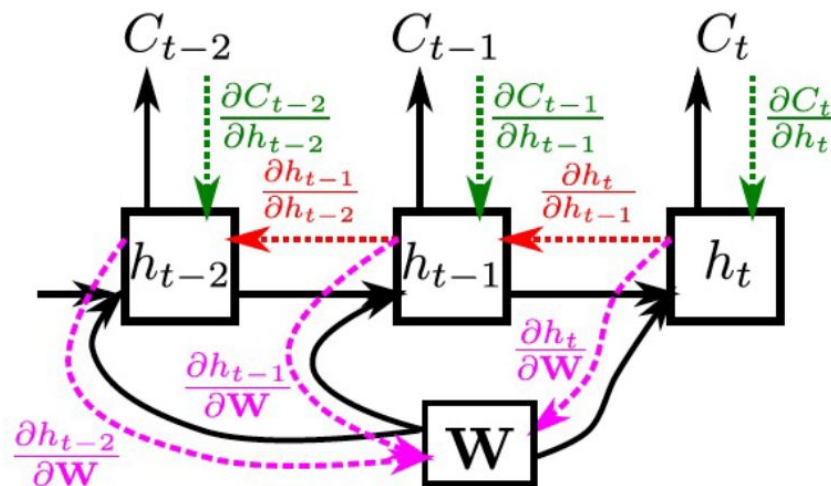$$L_t = CE(\hat{\mathbf{y}}_t, \mathbf{y}_t)$$
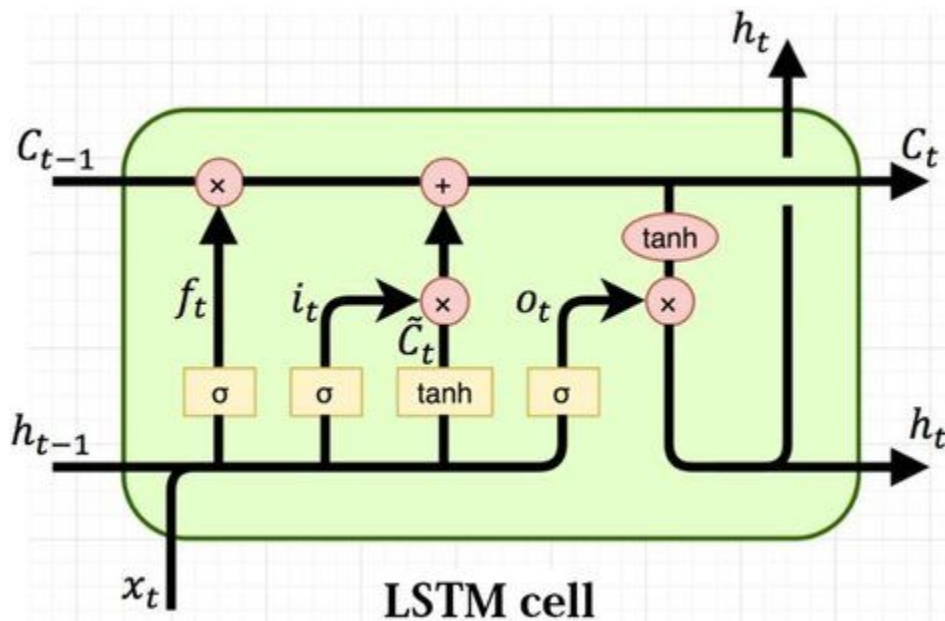
In total:

$$L = \sum_t L_t$$



45

# Backpropagation through time



$$\frac{\partial C_t}{\partial \mathbf{W}} = \sum_{t'=1}^{t} \frac{\partial C_t}{\partial h_t} \frac{\partial h_t}{\partial h_{t'}} \frac{\partial h_{t'}}{\partial \mathbf{W}}, \text{ where } \frac{\partial h_t}{\partial h_{t'}} = \prod_{k=t'+1}^{t} \frac{\partial h_k}{\partial h_{k-1}}$$

# Long Short Term Memory (LSTM)



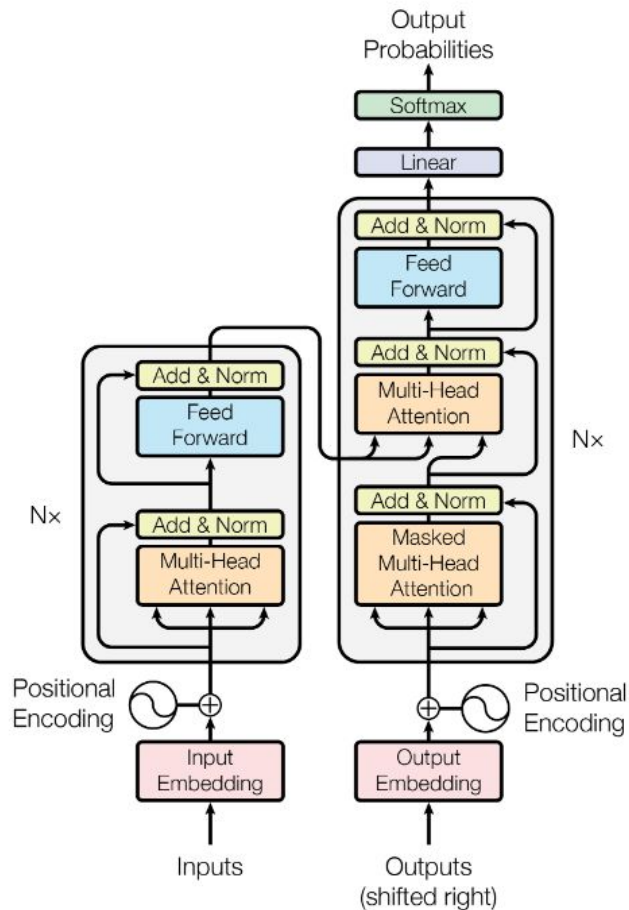$$i_t = \sigma\left(x_t U^i + h_{t-1} W^i\right)$$

$$f_t = \sigma\left(x_t U^f + h_{t-1} W^f\right)$$

$$o_t = \sigma\left(x_t U^o + h_{t-1} W^o\right)$$

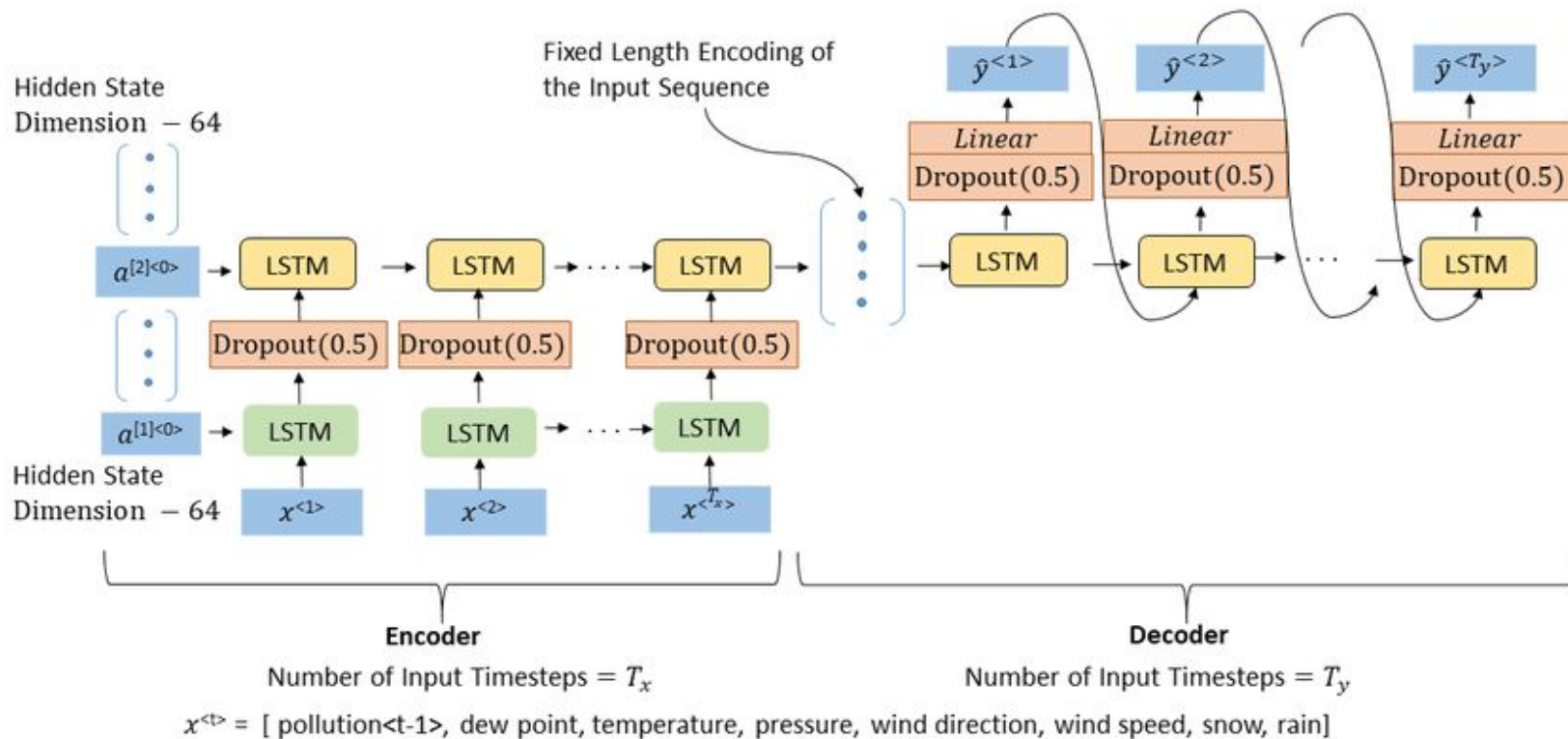$$\tilde{C}_t = \tanh\left(x_t U^g + h_{t-1} W^g\right)$$

$$C_t = \sigma\left(f_t * C_{t-1} + i_t * \tilde{C}_t\right)$$
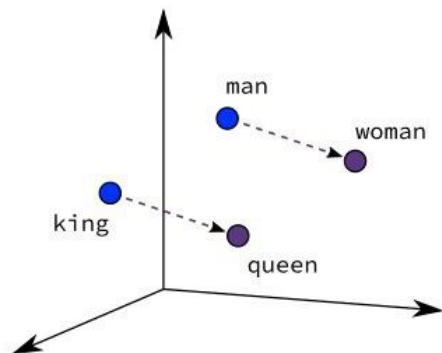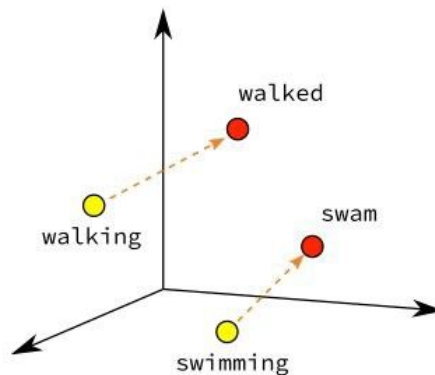
$$h_t = \tanh(C_t) * o_t$$

# Transformer



Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Nx

Nx

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

# Encoder Decoder



**Hidden State Dimension — 64**

$a^{[2]<0>}$

**Hidden State Dimension — 64**

$x^{<1>}$ $x^{<2>}$ $x^{<T_x>}$

$a^{[1]<0>}$

Fixed Length Encoding of the Input Sequence

$\hat{y}^{<1>}$ $\hat{y}^{<2>}$ $\hat{y}^{<T_y>}$

Linear Dropout(0.5)  Linear Dropout(0.5)  Linear Dropout(0.5)

Dropout(0.5)  Dropout(0.5)  Dropout(0.5)

LSTM

**Encoder**
Number of Input Timesteps $= T_x$

**Decoder**
Number of Input Timesteps $= T_y$

$x^{<t>}$ = [ pollution<t-1>, dew point, temperature, pressure, wind direction, wind speed, snow, rain]

# Word Embedding



Male-Female

Verb Tense

Country-Capital

https://projector.tensorflow.org/

# Attention

# Scaled Dot Product Attention

$$attention(Q,K,V) = softmax\left(\frac{QK^T}{\sqrt{d}}\right) V \qquad (4)$$

- Q -> query
- K -> key
- V -> value

yemek
| 2 | -1 |

●

koşarken
| -3 |
| 1 |

$= 2*(-3) + (-1)*1$
$= -7$

yemek
| 2 | -1 |

●

yemek
| 2 |
| -1 |

$= 5$

yemek
| 2 | -1 |

●

yedik
| 2 |
| 1 |

$= 3$

(1)

$$softmax(-7) = 10^{-6}$$

$$softmax(5) = 0.88$$

$$softmax(3) = 0.12$$

(2)

$$10^{-6} \; \boxed{\begin{array}{c|c} \text{koşarken} \\ \hline -3 & 1 \end{array}} \; + \; 0.88 \; \boxed{\begin{array}{c|c} \text{yemek} \\ \hline 2 & -1 \end{array}} \; + \; 0.12 \; \boxed{\begin{array}{c|c} \text{yedik} \\ \hline 2 & 1 \end{array}} \quad\quad (3)$$

$$= \; \boxed{\begin{array}{c|c} \text{yemek'} \\ \hline 2 & -0.76 \end{array}}$$

**Query**

| | |
|---|---|
| -3 | 1 |
| 2 | -1 |
| 2 | 1 |

koşarken, yemek, yedik

**Key$^T$**

| | | |
|---|---|---|
| -3 | 2 | 2 |
| 1 | -1 | 1 |

koşarken, yemek, yedik

✖

=

| | | |
|---|---|---|
| 10 | -7 | -5 |
| -7 | 5 | 3 |
| -5 | 3 | 5 |

koşarken, yemek, yedik

koşarken, yemek, yedik

(5)

$$softmax\left( \frac{\begin{array}{|c|c|c|} \hline 10 & -7 & -5 \\ \hline -7 & 5 & 3 \\ \hline -5 & 3 & 5 \\ \hline \end{array}}{\sqrt{2}} \right) = \begin{array}{|c|c|c|} \hline 0.99 & 10^{-6} & 10^{-5} \\ \hline 10^{-4} & 0.80 & 0.19 \\ \hline 10^{-4} & 0.19 & 0.80 \\ \hline \end{array} \quad (6)$$

Word embedding

| | | |
|---|---|---|
| koşarken | -3 | 1 |
| yemek | 2 | -1 |
| yedik | 2 | 1 |

$\text{Weight}_Q$ $\times$ ... $=$ $Q$

$\text{Weight}_K$ $\times$ ... $=$ $K$ (8)

$\text{Weight}_V$ $\times$ ... $=$ $V$

$Y_1$

$Y_2$

$Y_3$

$Y$

$\times$

$=$

(9)

# Positional Encoding

Solution #1

# Positional Encoding

Solution #2

# Positional Encoding

Solution #3

$$PE_{(pos,2i)} = sin(pos/10000^{2i/d_{model}})$$
$$PE_{(pos,2i+1)} = cos(pos/10000^{2i/d_{model}})$$

p₀

i=0
i=1
i=2
i=3
i=4

d=5

pos = 0

https://www.youtube.com/watch?v=dichIcUZfOw
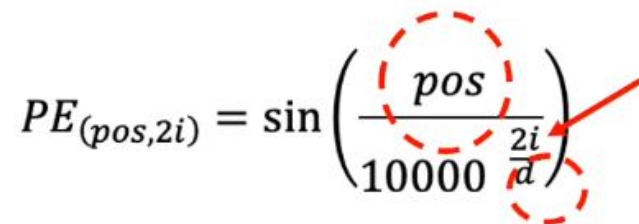
$$PE_{(pos,2i)} = sin\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$

# Dropout



(a) Standard Neural Net           (b) After applying dropout.

# Residual Connection



$\mathbf{x}$

weight layer

$\mathcal{F}(\mathbf{x})$

relu

weight layer

$\mathbf{x}$
identity

$\mathcal{F}(\mathbf{x}) + \mathbf{x}$

relu

# Complexity

Table 1: Maximum path lengths, per-layer complexity and minimum number of sequential operations for different layer types. $n$ is the sequence length, $d$ is the representation dimension, $k$ is the kernel size of convolutions and $r$ the size of the neighborhood in restricted self-attention.
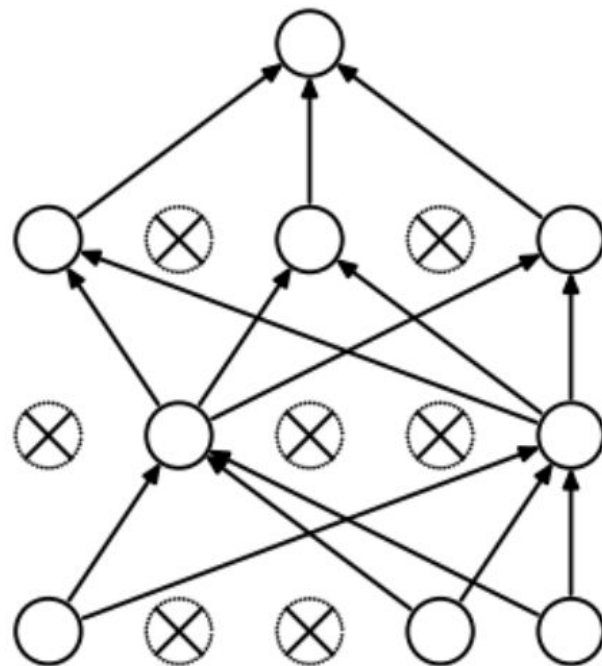
| Layer Type | Complexity per Layer | Sequential Operations | Maximum Path Length |
|---|---|---|---|
| Self-Attention | $O(n^2 \cdot d)$ | $O(1)$ | $O(1)$ |
| Recurrent | $O(n \cdot d^2)$ | $O(n)$ | $O(n)$ |
| Convolutional | $O(k \cdot n \cdot d^2)$ | $O(1)$ | $O(log_k(n))$ |
| Self-Attention (restricted) | $O(r \cdot n \cdot d)$ | $O(1)$ | $O(n/r)$ |

# Advantages

parallelizable (thus faster)

computationally less complex (most of the time)

better capture longer dependencies

more interpretable

# Improved Transformers

- **BERT**
- GPT
- T5
- BART
- Pegasus
- XLM
- Reformer
- Longformer
- ELECTRA
- RoBERTa
- ...

# Vision Transformers

AN IMAGE IS WORTH 16x16 WORDS:TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

https://arxiv.org/pdf/2010.11929.pdf

# Courses & Resources

- Stanford University NLP w/ DL - http://web.stanford.edu/class/cs224n/
- Huggingface - https://huggingface.co/course/chapter1
- Deeplearning.ai NLP -
  https://www.deeplearning.ai/program/natural-language-processing-specialization/

# Bibliography

- https://user.ceng.metu.edu.tr/~skalkan/DL/week_13.pdf
- https://kazemnejad.com/blog/transformer_architecture_positional_encoding/
- https://www.youtube.com/watch?v=dichIcUZfOw
- https://arxiv.org/pdf/1706.03762.pdf