

Lecture 4: Basic Intro to Transformers

① Most modern LLMs rely on the transformer architecture → deep neural network architecture introduced in 2017 paper

Attention is All you Need

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

Lukasz Kaiser^{*}
Google Brain
lukasz.kaiser@google.com

Illa 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.

miro

Convolutional Neural Networks

Upgrade

deep

introduced in 2017 paper

Attention is All you Need



Original Transformer: Developed for machine translation

translating English texts to German and French.

2

Encoder returns embedding

example



Encoder returns embedding
vector as input to
decoder

Produces text encodings
used by decoder

vector
embedding

Input text
prepared for
encoder

tokenization

English

Embeddings

Encoder

Preprocessing steps

Input text

"This is an example"

(Input text to be
translated)

"Das est ein Beispiel"

Output layers

Decoder

Preprocessing steps

Input text

"Das est ein"

Partial output text

Model completes one
word at a time

The complete output
(translation)

Generates translated
one word at a time

Input text is
prepared for
decoder

SIMPLIFIED TRANSFORMER
ARCHITECTURE

word at a time

Generates output text from encoded vectors

A note on self-attention mechanism:



Encodes input
text into vectors

Generates output text
from encoded vectors

A note on self-attention mechanism:

- Key part of transformers. Allows model to weigh importance of different words / tokens relative to each other.
- Enables model to capture long range dependencies.
- We will look at this in detail later.



- We will look at this in detail later.

④ later variations of transformer architecture: ^{→ 2017}

BERT

GPT models

★ "Bidirectional encoder representations"

"Generative pretrained"

miro

Convolutional Neural Networks

Upgrade

Present

Share

BERT

GPT models

“Bidirectional encoder representations from transformers”

“Generative pretrained Transformers”

Predicts hidden words

Generates new

miro

Convolutional Neural Networks

Upgrade

Present

Share

“Bidirectional
encoder representations
from transformers”

Predicts hidden words
in a given sentence

This is an example of how
Llm can perform

“Generative pretrained
Transformers”

Generates new
word

This is an example of
Llm can perform

BERT

Fills
missing
words

This is an example of how
LLM can perform

Encoder

Preprocessing steps

input text

This is an ? of how
LLM ? perform

Receives inputs where words are
randomly masked during training

This is an example of how
LLM can perform

Decoder

Preprocessing steps

input text

This is an example of how
LLM can —?



in a given sentence

This is an example of how
LLM can perform

Encoder

Preprocessing steps

input text

This is an ? of how
LLM ? perform

word

This is an example of how
LLM can perform

Decoder

Preprocessing steps

input text

This is an example of how
LLM can —?

learns to
one word
a time.

receives



receives inputs where words are

this is an example of how
LLM can perform

Encoder

Preprocessing steps

input text

This is an ? of how
LLM ? perform

outputs where words are
masked during training

this is an example of how
LLM can perform

Decoder

Preprocessing steps

input text

This is an example of how
LLM can ?

learns to generate
one word at
a time.

receives
incomplete text

⑤

Transformers vs LLMs

- Not all transformers are LLMs
- Transformers can also be used for Computer Vision
- Not all LLMs are transformers
- LLMs can be based on recurrent or ~~convolutional~~ architectures as well

