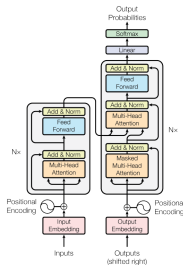


# Transformer

## A universal deep learning architecture



### Learning goals

- Understand the initial use of the Transformer
- Grasp the other application fields since then

# MACHINE TRANSLATION

- Sequence-to-sequence task
- Already served as a motivation for introducing the "ordinary" Attention-mechanism ▶ Bahdanau et al., 2015
- Crucial, that the decoder has access to the whole input sequence  
→ This is very well solved by cross-attention
- Good contextualization in the encoder improves translation quality
  - (Bidirectional) RNNs/LSTMs are only (concatenated) unidirectional architectures
  - Transformer-Encoder layers are bidirectional by construction
  - Stacking them on top of each other makes this bidirectional contextualization even "deeper"

# WMT 2014 EN-TO-DE AND EN-TO-FR

## Parallel training data:

Parallel data:

File	Size	CS-EN	DE-EN	HI-EN	FR-EN	RU-EN	Notes
<a href="#">Europarl v7</a>	628MB	✓	✓		✓		same as previous year, <a href="#">corpus home page</a>
<a href="#">Common Crawl corpus</a>	876MB	✓	✓		✓	✓	same as previous year
<a href="#">UN corpus</a>	2.3GB				✓		same as previous year, <a href="#">corpus home page</a>
<a href="#">News Commentary</a>	77MB	✓	✓		✓	✓	updated, <a href="#">data with document boundaries</a>
<a href="#">10<sup>2</sup>French-English corpus</a>	2.3 GB				✓		same as previous year [ <a href="#">md5</a> ] <a href="#">sha1</a>
CzEng 1.0	115MB	✓					same as previous year, <a href="#">corpus home page</a> (avoid sections 98 and 99)
Yandex 1M corpus	121MB					✓	<a href="#">corpus home page</a> ; v1.3 now in original case
<a href="#">Wiki Headlines</a>	7.8MB			✓		✓	Provided by CMU. The ru-en is unchanged from last year.
<a href="#">HindEnCorp</a>	25MB			✓			Collected by Charles University
The JHU Corpus				✗			This is fully contained in HindEnCorp, so not made available here.

# THE BLEU SCORE

- Based on n-gram overlap from candidate and reference sentence
- Precision (for each n-gram) calculated as follows:

$$P_n = \frac{\sum_{C \in \{\text{Candidates}\}} \sum_{n\text{-gram} \in C} \text{Count}_{\text{clip}}(n\text{-gram})}{\sum_{C' \in \{\text{Candidates}\}} \sum_{n\text{-gram}' \in C'} \text{Count}(n\text{-gram}')}$$

- Finally, the BLEU score can be computed as

$$\text{BLEU} = BP \cdot \exp \left( \sum_{n=1}^N w_n \cdot \log(P_n) \right),$$

- where  $BP$  is a "brevity penalty" to penalize short generations,  $N$  is the number of n-grams &  $w_n$  the weight for each  $P_n$  (usually  $\frac{1}{N}$ )

# TRANSFORMER FOR MT

The Transformer ..

- .. outperforms the previous SOTA models
- .. at a lower number of required training FLOPs

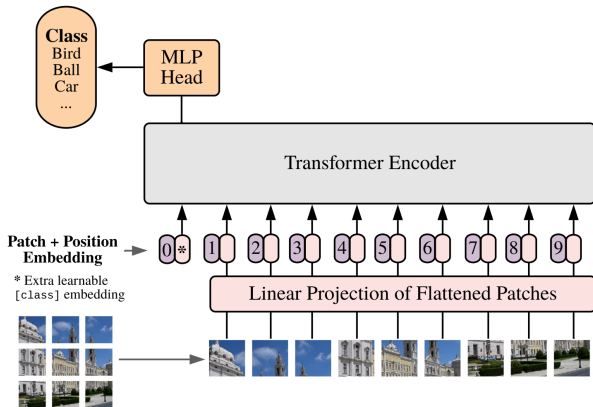
Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	<b>41.29</b>	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1	<b><math>3.3 \cdot 10^{18}</math></b>	
Transformer (big)	<b>28.4</b>	<b>41.8</b>	$2.3 \cdot 10^{19}$	

► Source: Vaswani et al., 2017

# TRANSFORMER FOR COMPUTER VISION

## Vision Transformer

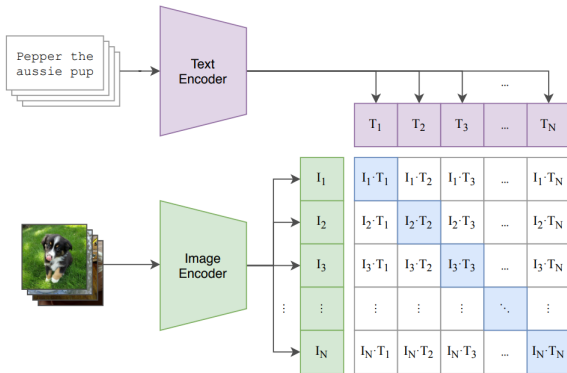


► Source: Dosovitskiy et al., 2020

# TRANSFORMER FOR MULTIMODAL LEARNING

## CLIP

(1) Contrastive pre-training

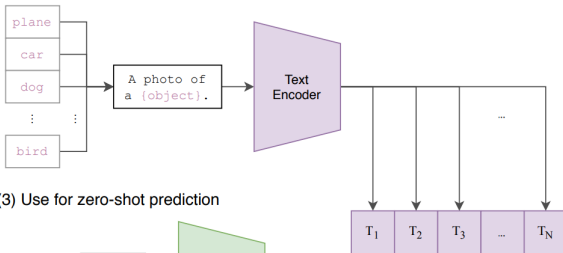


► Source: Radford et al., 2021

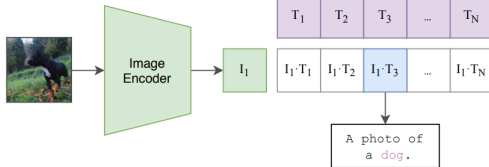
# TRANSFORMER FOR MULTIMODAL LEARNING

## CLIP

(2) Create dataset classifier from label text



(3) Use for zero-shot prediction

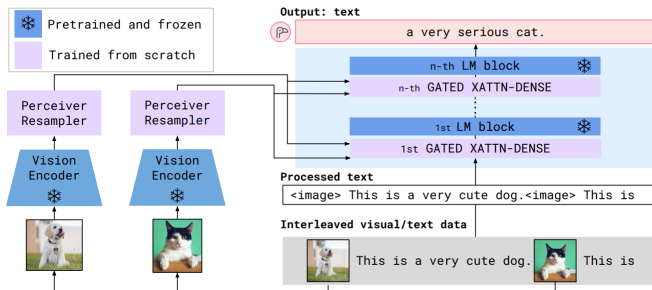


► Source: Radford et al., 2021



# TRANSFORMER FOR MULTIMODAL LEARNING

## FLAMINGO



► Source: Alayrac et al., 2022

# TRANSFORMER FOR MULTIMODAL LEARNING

## GPT-4

---

### Example of GPT-4 visual input:

---

User      What is funny about this image? Describe it panel by panel.



Source: <https://www.reddit.com/r/t/1amm/comments/abab5v/1amm/>

GPT-4      The image shows a package for a "Lightning Cable" adapter with three panels.

Panel 1: A smartphone with a VGA connector (a large, blue, 15-pin connector typically used for computer monitors) plugged into its charging port.

Panel 2: The package for the "Lightning Cable" adapter with a picture of a VGA connector on it.

Panel 3: A close-up of the VGA connector with a small Lightning connector (used for charging iPhones and other Apple devices) at the end.

The humor in this image comes from the absurdity of plugging a large, outdated VGA connector into a small, modern smartphone charging port.

---

**Table 3.** Example prompt demonstrating GPT-4's visual input capability. The prompt consists of a question about an image with multiple panels which GPT-4 is able to answer.

► Source: OpenAI, 2023