

# Transformers, Large Language Models (LLMs), and Foundation Models

Ju Sun

Computer Science & Engineering

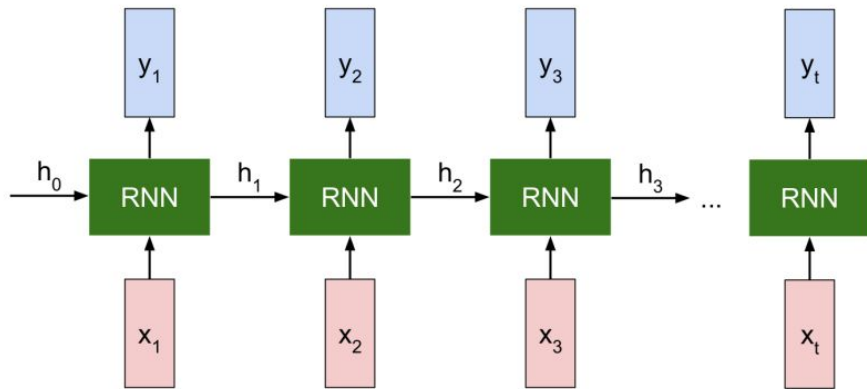
Apr 07/09, 2026



UNIVERSITY OF MINNESOTA  
Driven to Discover<sup>SM</sup>

# Quick recap

## RNN: model sequences



(Credit: Stanford CS231N)

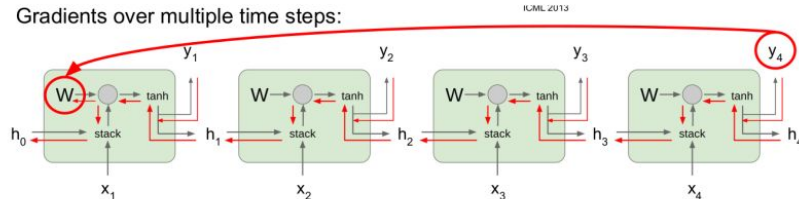
$$\mathbf{h}_t = \tanh(\mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{W}_x \mathbf{x}_t)$$

$$\mathbf{y}_t = \mathbf{V}_y \mathbf{h}_t$$

$\mathbf{W}_h$ ,  $\mathbf{W}_x$  and  $\mathbf{V}_y$  are shared across the sequence

## Vanishing/exploding gradient issue

Gradients over multiple time steps:



(Credit: Stanford CS231N)

$$\begin{aligned} \frac{\partial L_t}{\partial \mathbf{W}} &= \frac{\partial L_t}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \dots \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}} = \frac{\partial L_t}{\partial \mathbf{h}_t} \left( \prod_{k=2}^t \frac{\partial \mathbf{h}_k}{\partial \mathbf{h}_{k-1}} \right) \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}} \\ &= \frac{\partial L_t}{\partial \mathbf{h}_t} \left( \prod_{k=2}^t \text{diag}(\tanh'(\mathbf{W}_h \mathbf{h}_{k-1} + \mathbf{W}_x \mathbf{x}_k)) \mathbf{W}_h \right) \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}} \end{aligned}$$

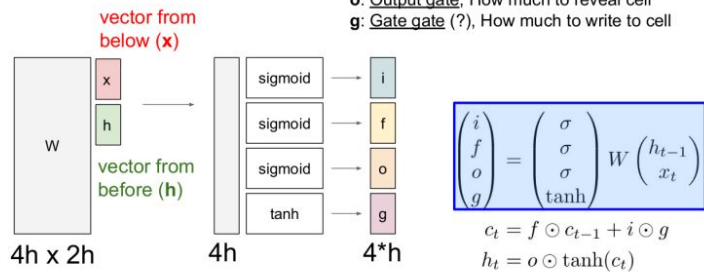
- \* when  $\|\mathbf{W}_h\| > 1$ , gradient **explodes** if  $t$  large
- \* when  $\|\mathbf{W}_h\| < 1$ , gradient **vanishes** if  $t$  large

$$\begin{aligned} &\left\| \prod_{k=2}^t \text{diag}(\tanh'(\mathbf{W}_h \mathbf{h}_{k-1} + \mathbf{W}_x \mathbf{x}_k)) \mathbf{W}_h \right\| \\ &\leq \prod_{k=2}^t \left\| \text{diag}(\tanh'(\mathbf{W}_h \mathbf{h}_{k-1} + \mathbf{W}_x \mathbf{x}_k)) \right\| \|\mathbf{W}_h\| \\ &\leq \prod_{k=2}^t \left\| \text{diag}(\tanh'(\mathbf{W}_h \mathbf{h}_{k-1} + \mathbf{W}_x \mathbf{x}_k)) \right\| \|\mathbf{W}_h\|^{t-1} \end{aligned}$$

# Gated RNNs

## Long Short Term Memory (LSTM)

[Hochreiter et al., 1997]



(Credit: Stanford CS231N)

$u$ : **update gate**, control state update

$r$ : **reset gate**, control how previous state affects new content

$g$ : new content

**Gated recurrent unit (GRU)**

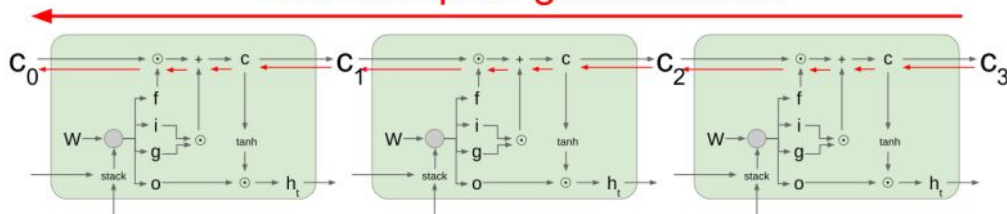
$$\begin{bmatrix} u \\ r \\ g \end{bmatrix} = \begin{bmatrix} \sigma \\ \sigma \end{bmatrix} \left( W \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix} \right)$$

$$g = \tanh(W_h (r \odot h_{t-1}) + W_x x_t + b_g)$$

$$h_t = u \odot h_{t-1} + (1 - u) \odot g$$

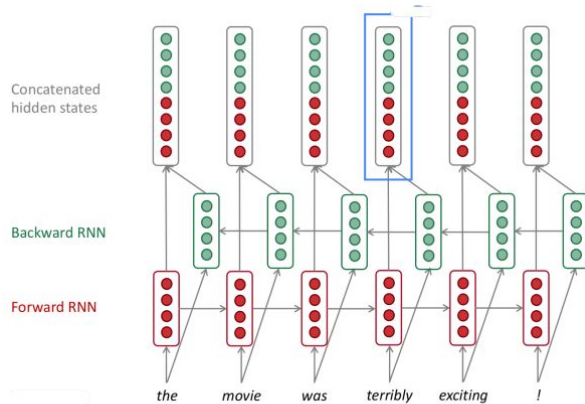
$f, i, o$  are merged

**Uninterrupted gradient flow!**



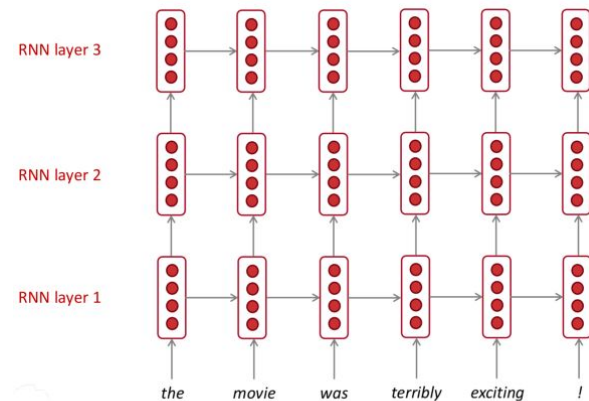
(Credit: Stanford CS231N)

# Modern RNNs



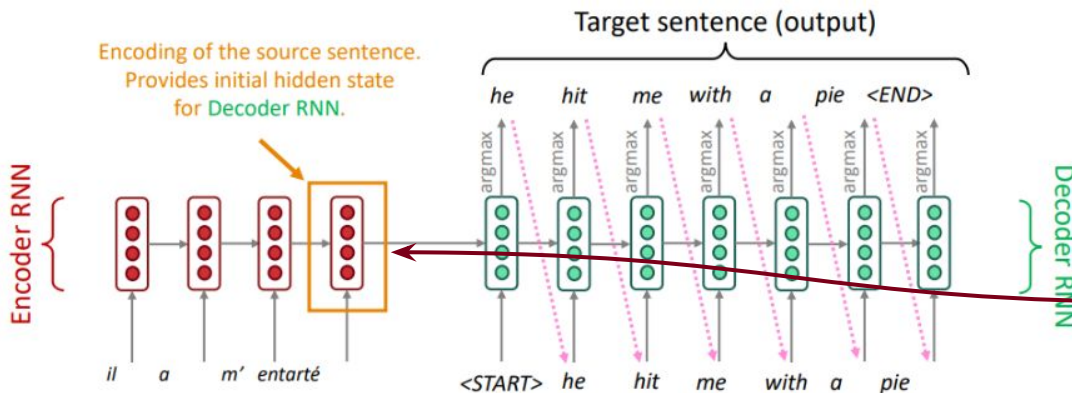
(Credit: Stanford CS224N)

## Bidirectional RNN



(Credit: Stanford CS231N)

## Deep RNN

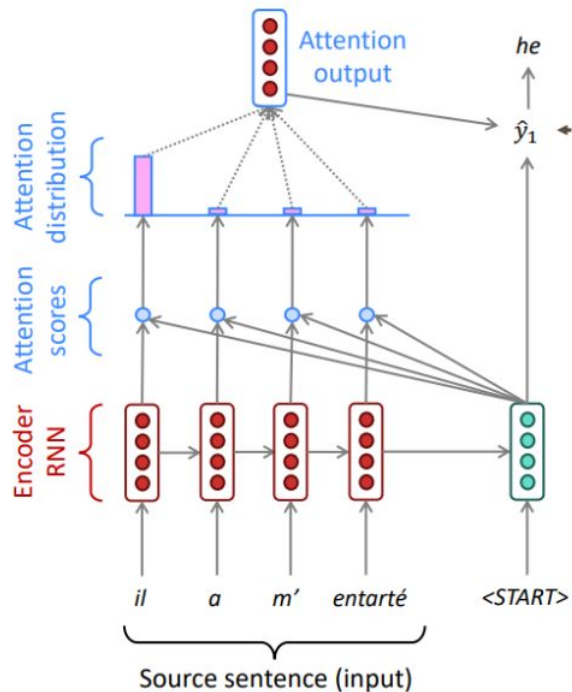


(Credit: Stanford CS231N)

## Seq2Seq model

**Bottleneck problem**

# Attention mechanism



(Credit: Stanford CS231N)

**Input:** source vectors  $s_1, \dots, s_N \in \mathbb{R}^h$ , and target vector  $t$

**Output:** weighted summation

$$\sum_{j=1}^N w_j s_j \quad \text{where } w_j = \text{similarity}(s_j, t)$$

Attention scores

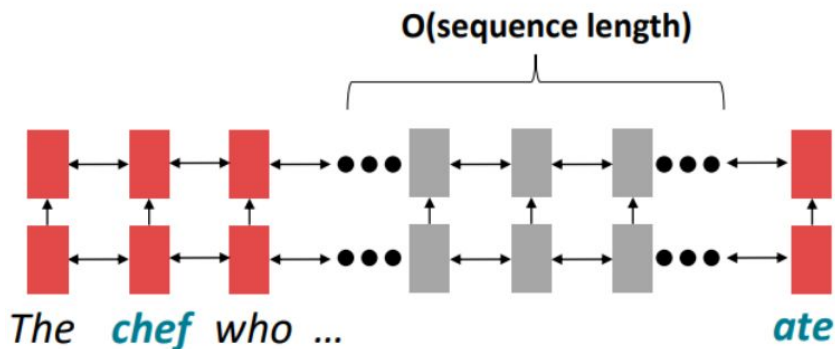
Many possibilities:

- dot-product attention:  $\hat{w}_j = \langle s_j, t \rangle$  (Is it better to normalize this or rescale it by the dimension factor? )
- multiplicative attention:  $\hat{w}_j = \langle s_j, Wt \rangle$
- "additive attention":  $\hat{w}_j = v^T \sigma(W_1 s_j + W_2 t)$

The actual weights are attention scores passed through **softmax**

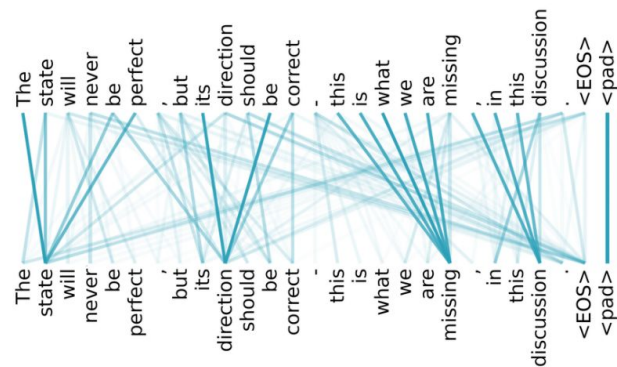
$$w_j = \frac{\exp(\hat{w}_j)}{\sum_k \exp(\hat{w}_k)}$$

# Self-attention



## RNN

- Long interaction distance
- Resistant to parallelization

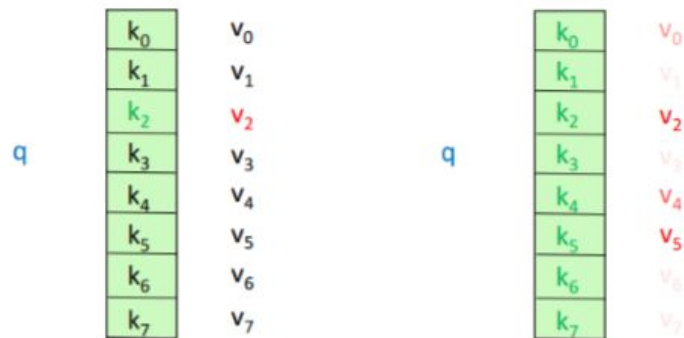


## Self-attention

- $O(1)$  interaction distance
- Highly parallelizable

**Each token gets a selective summary of information from all others**

## Self-attention: a closer look



- Each word now encoded as (query, key, value) triple
- For an input  $x_i$ , we have:

$$\mathbf{q}_i = (\mathbf{W}^Q)^\top \mathbf{x}_i, \quad \mathbf{k}_i = (\mathbf{W}^K)^\top \mathbf{x}_i, \quad \mathbf{v}_i = (\mathbf{W}^V)^\top \mathbf{x}_i$$

- Calculate attention scores between query and all keys:  $e_{ij} = \langle \mathbf{q}_i, \mathbf{k}_j \rangle$
- softmax normalization  $w_{ij} = \exp(e_{ij}) / \sum_k \exp(e_{ik})$
- output the weighted sum of values  $\sum_j w_{ij} \mathbf{v}_j$

## Self-attention in matrix notation

Assume  $X$  collects all input words, each one a row:

- Compute queries, keys, and values

$$Q = XW^Q, \quad K = XW^K, \quad V = XW^V$$

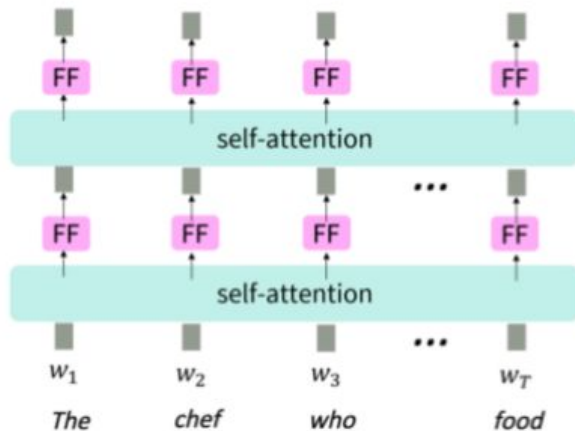
- Calculate attention scores between query and all keys:  $E = QK^T$
- softmax normalization to each row:  $A = \text{softmax}(E)$
- output the weighted sum of values  $AV$

$$\text{output} = \text{softmax}(QK^T)V$$

# Adding in nonlinearity

Equation for Feed Forward Layer

$$m_i = MLP(\text{output}_i) \\ = W_2 * \text{ReLU}(W_1 \times \text{output}_i + b_1) + b_2$$



Encoder



First step toward the Transformer!

# A block/layer of transformer!

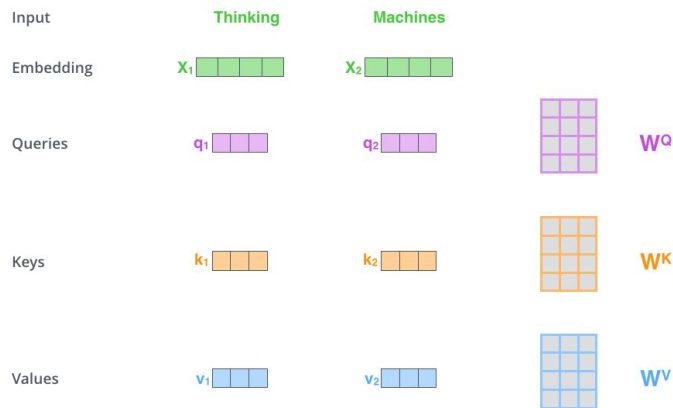


Image credit: <https://jalammargithub.io/illustrated-transformer/>

- Each word now encoded as (query, key, value) triple
- For an input  $x_i$ , we have:

$$q_i = (W^Q)^T x_i, \quad k_i = (W^K)^T x_i, \quad v_i = (W^V)^T x_i$$

- Calculate attention scores between query and all keys:  $e_{ij} = \langle q_i, k_j \rangle$
- softmax normalization  $w_{ij} = \exp(e_{ij}) / \sum_k \exp(e_{ik})$
- output the weighted sum of values  $\sum_j w_{ij} v_j$

In matrix notation

- Compute queries, keys, and values

$$Q = XW^Q, \quad K = XW^K, \quad V = XW^V$$

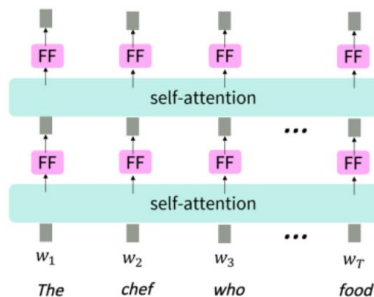
- Calculate attention scores between query and all keys:  $E = QK^T$
- softmax normalization  $A = \text{softmax}(E)$
- output the weighted sum of values  $AV$

$$\text{output} = \text{softmax}(QK^T)V$$

**Question:** why we need both query and key?

Equation for Feed Forward Layer

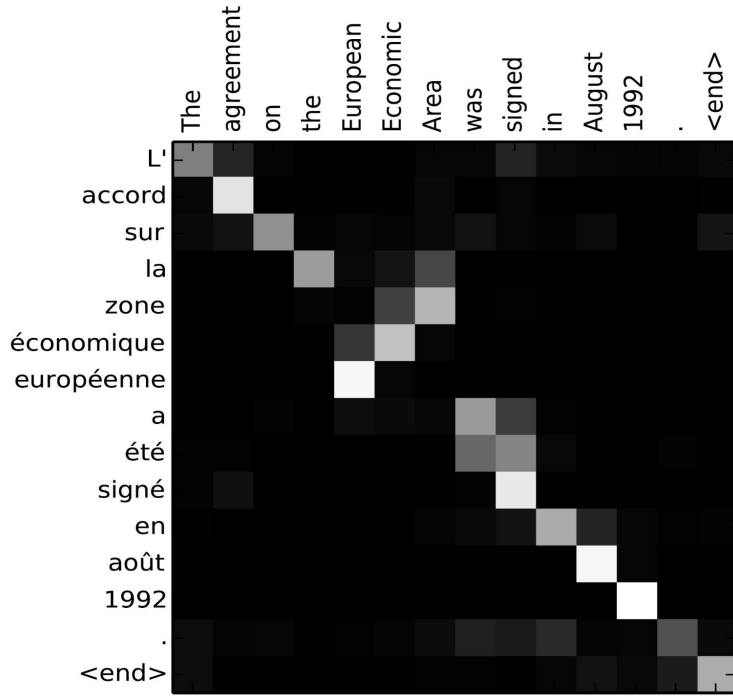
$$m_i = \text{MLP}(\text{output}_i) \\ = W_2 * \text{ReLU}(W_1 \times \text{output}_i + b_1) + b_2$$



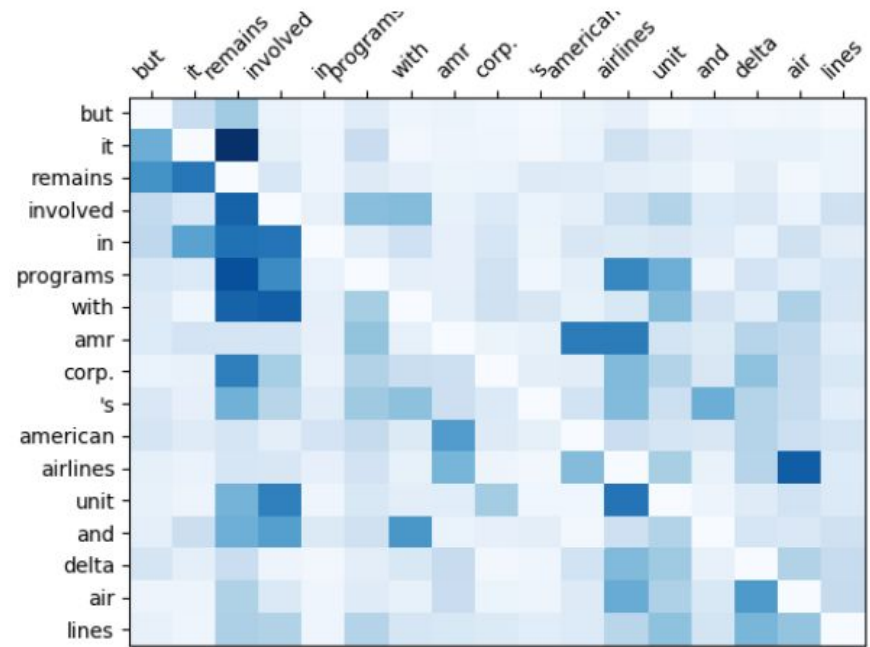
Adding in nonlinearity!

First step toward Transformers!

# Attention matrices—visualizing correlations



General attention



Self-attention

# Transformers

# Transformers

## 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  
lukaszkaier@google.com

**Illia Polosukhin\* ‡**  
illia.polosukhin@gmail.com

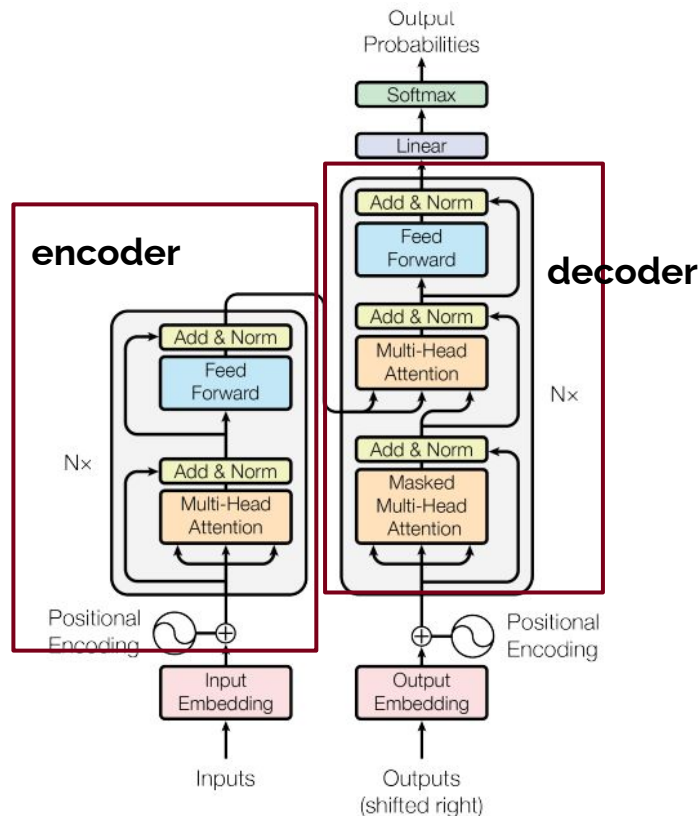
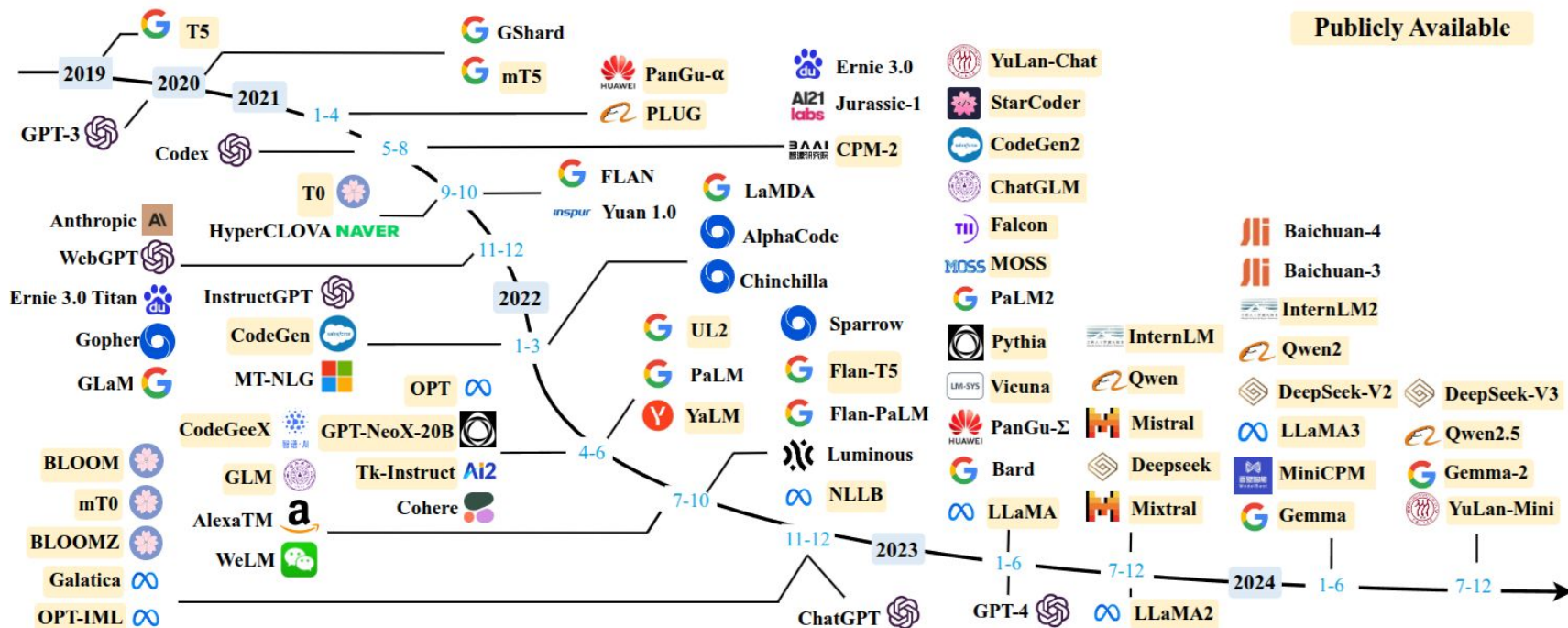


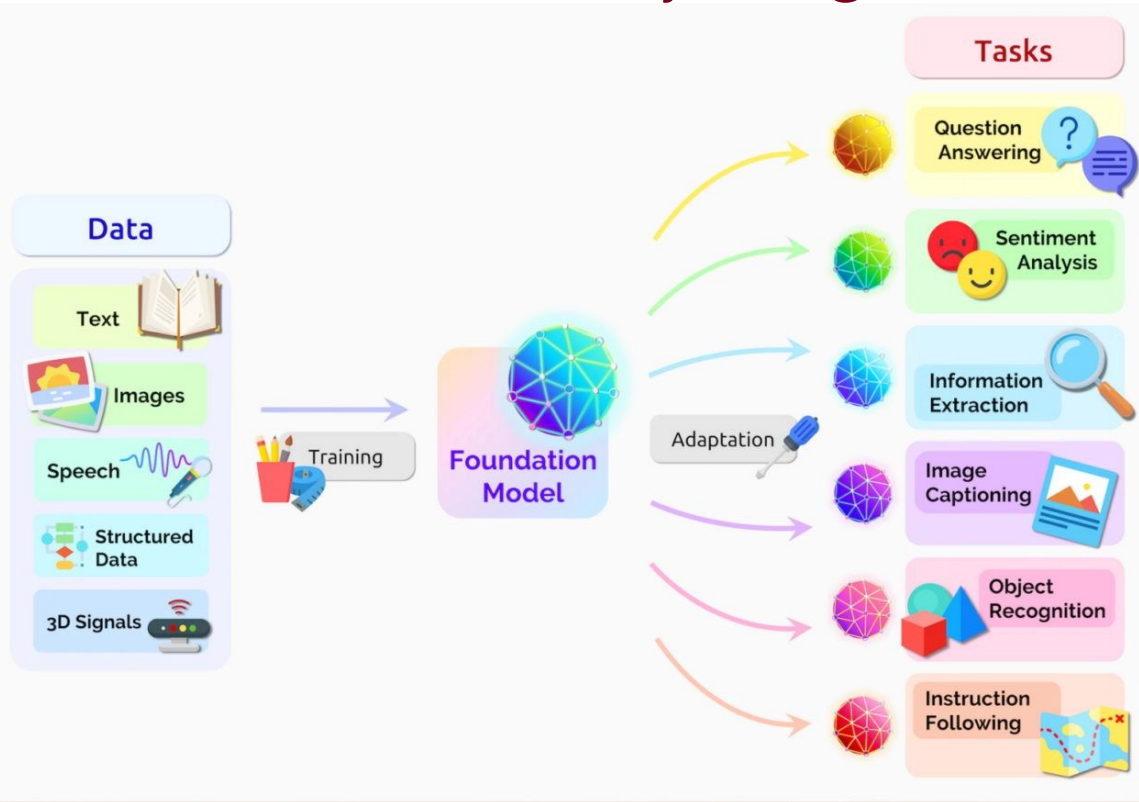
Figure 1: The Transformer - model architecture.

NIPS 2017; <https://arxiv.org/abs/1706.03762>

# Transformers reign in NLP!



# Transformers for everything!

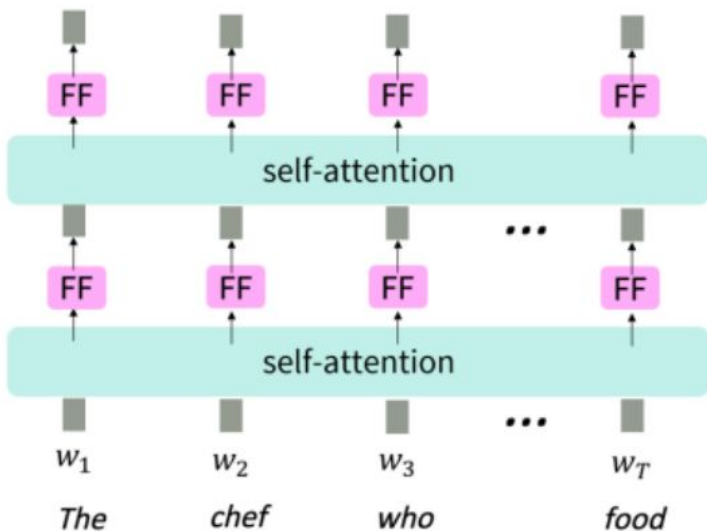


- Transformers have been modified to deal with **almost all** kinds of structured and **unstructured** data
- Enable multimodal data integration and interaction

# Starting from self-attention

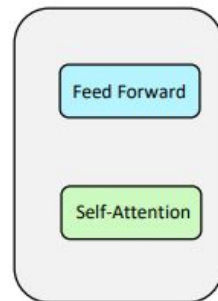
## Equation for Feed Forward Layer

$$m_i = MLP(\text{output}_i) \\ = W_2 * \text{ReLU}(W_1 \times \text{output}_i + b_1) + b_2$$



(Credit: Stanford CS231N)

Encoder



Input  
Embedding

Inputs



Decoder

Output  
Embedding

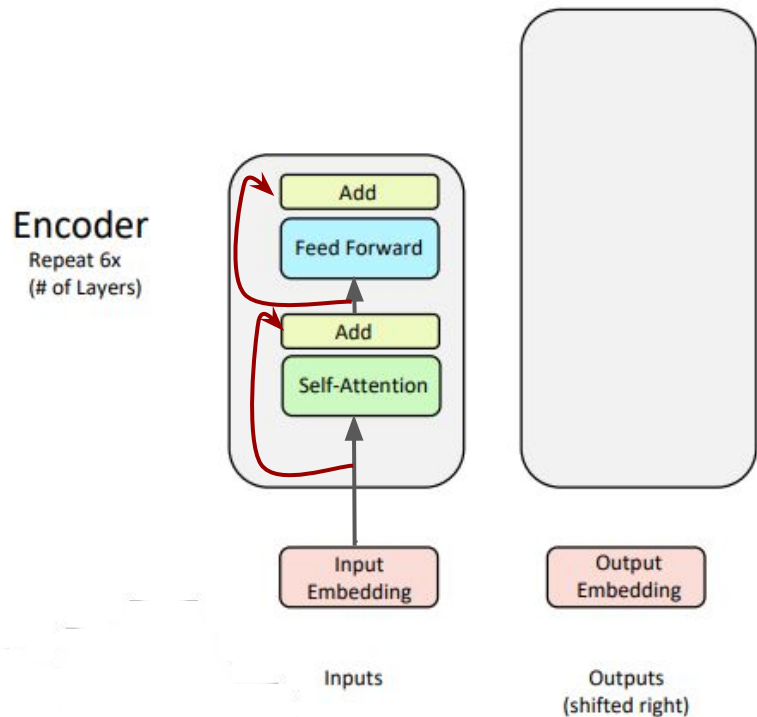
Outputs  
(shifted right)

(Credit: Stanford CS231N)

## Three tricks to build in depth:

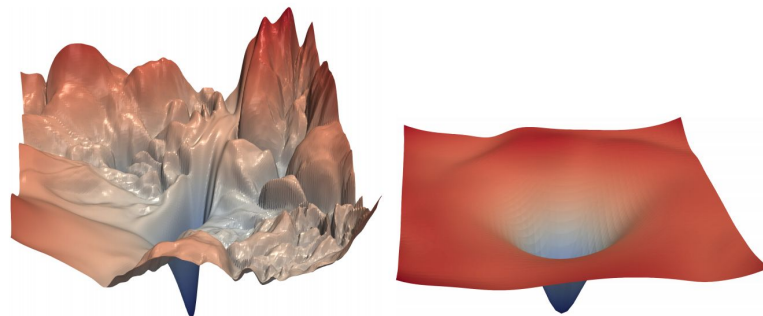
- Residual connection
- Layer normalization
- Scaled inner product attention

# Trick 1: Residual connection



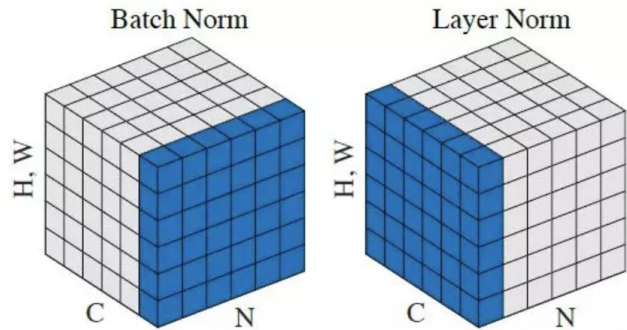
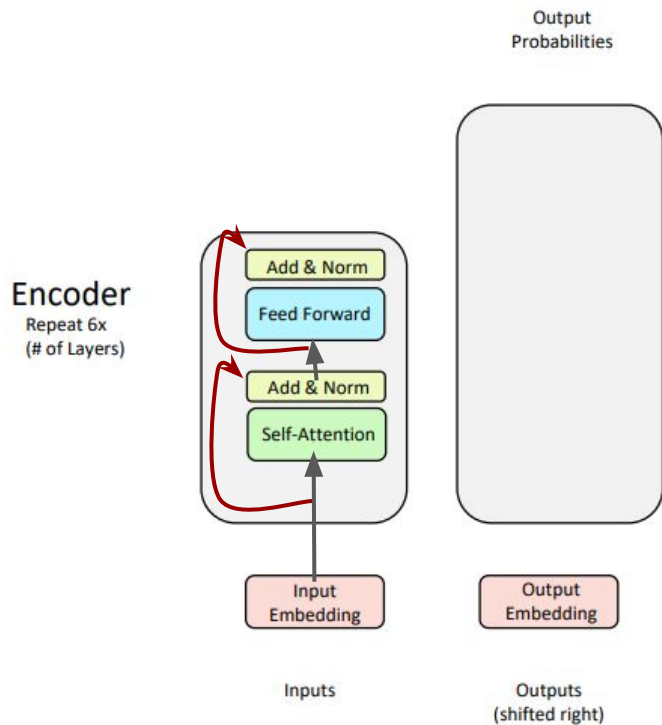
$$\mathbf{x}_k = F(\mathbf{x}_{k-1}) + \mathbf{x}_{k-1}$$

- Automatic architecture pruning
- Mitigating vanishing gradient
- Smoothing out landscape



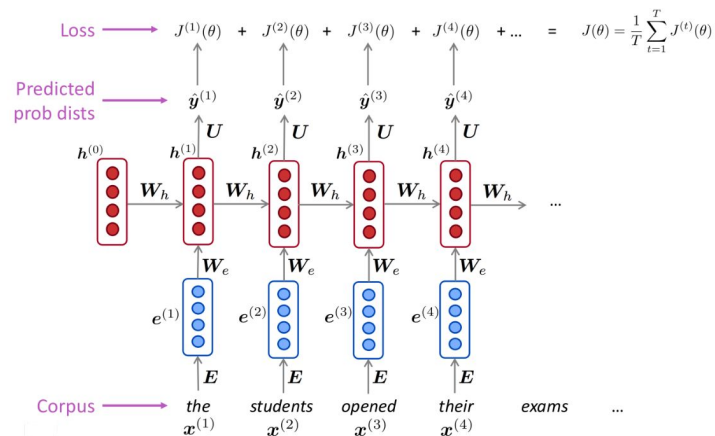
<https://arxiv.org/abs/1712.09913>

# Trick 2: Layer normalization

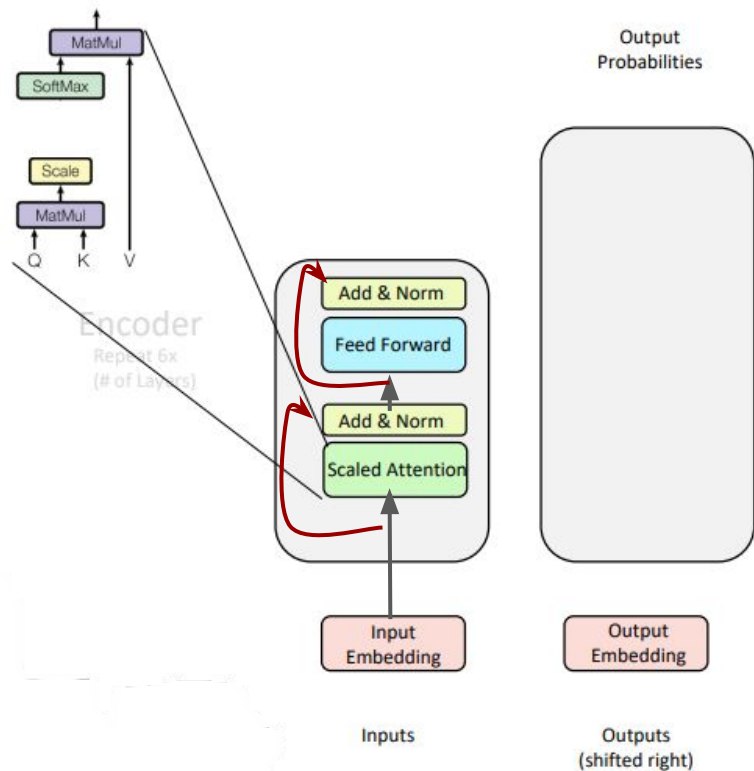


$$x^{\ell'} = \frac{x^{\ell} - \mu^{\ell}}{\sigma^{\ell} + \epsilon}$$

Why not batchnorm?



# Trick 3: Scaled inner product attention



$$\text{output} = \text{softmax}(\mathbf{Q}\mathbf{K}^T)\mathbf{V}$$

- Suppose that entries of  $\mathbf{Q}$  and  $\mathbf{K}$  behaves like IID zero-mean, unit variance
- $\mathbb{E}\langle \mathbf{q}^i, \mathbf{k}^j \rangle = 0$  but

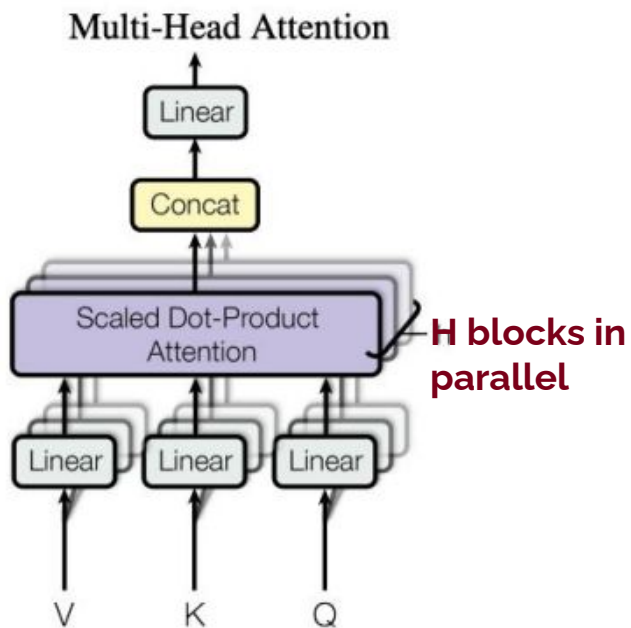
$$\text{Var}\langle \mathbf{q}^i, \mathbf{k}^j \rangle = d_k$$

This can blow up exp computation in the softmax normalization for large  $d_k$ !

Solution: normalize by standard deviation

$$\text{output} = \text{softmax}(\mathbf{Q}\mathbf{K}^T / \sqrt{d_k})\mathbf{V}$$

# Multi-head attention



[Vaswani et al. 2017]

**Multiple, independent** self-attention blocks in parallel

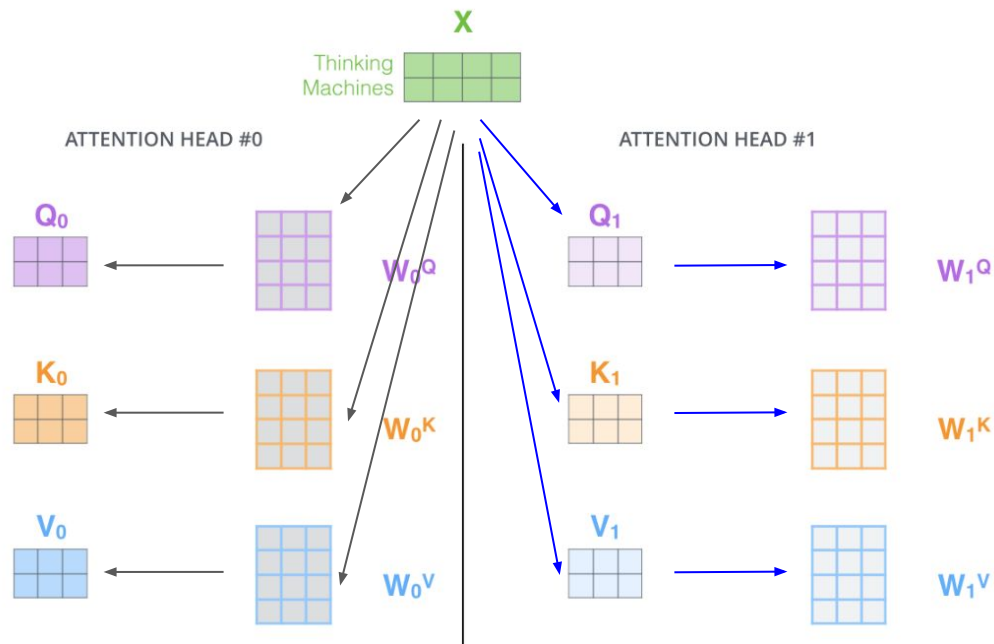
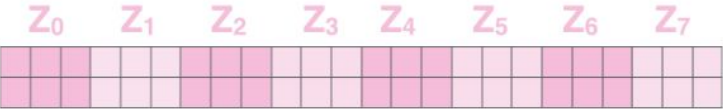


Image credit: <https://ialammargithub.io/illustrated-transformer/>

**Intuition:** allow the flexibility of capturing different kinds of "relevance"/correlations

# Multi-head attention

1) Concatenate all the attention heads



**Concatenate**

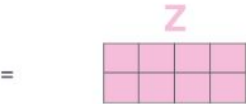


2) Multiply with a weight matrix  $W^O$  that was trained jointly with the model

X

**Multiply**

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



**Output**



Image credit: <https://jalammar.github.io/illustrated-transformer/>

# Multi-head attention in one picture

1) This is our input sentence\*

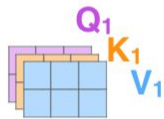
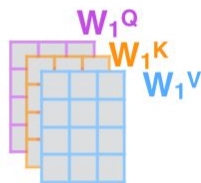
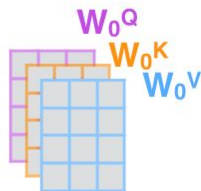
2) We embed each word\*

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

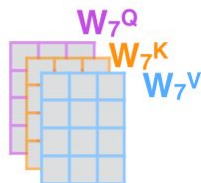
Thinking Machines



...

...

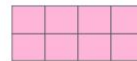
...



$W^O$

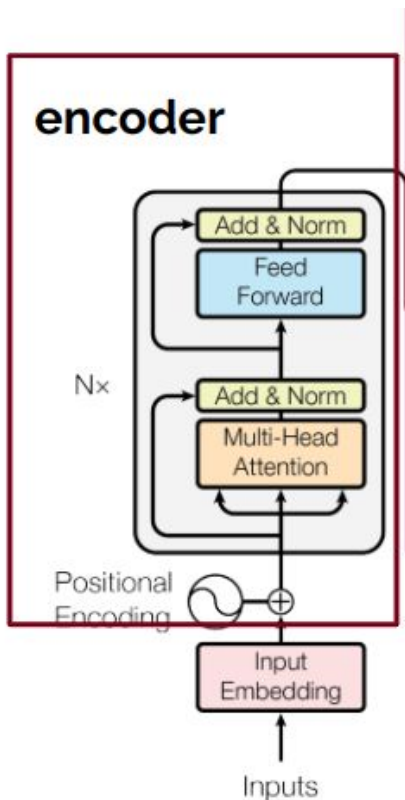


$Z$



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

# Positional encoding



Does the input order matter or not?

$$Q = XW^Q, \quad K = XW^K, \quad V = XW^V$$

$$\text{output} = \text{softmax}(QK^T / \sqrt{d_k})V$$

**Positional encoding** to break the **order invariance**

- Idea: a positional vector to (hopefully) encode the position information

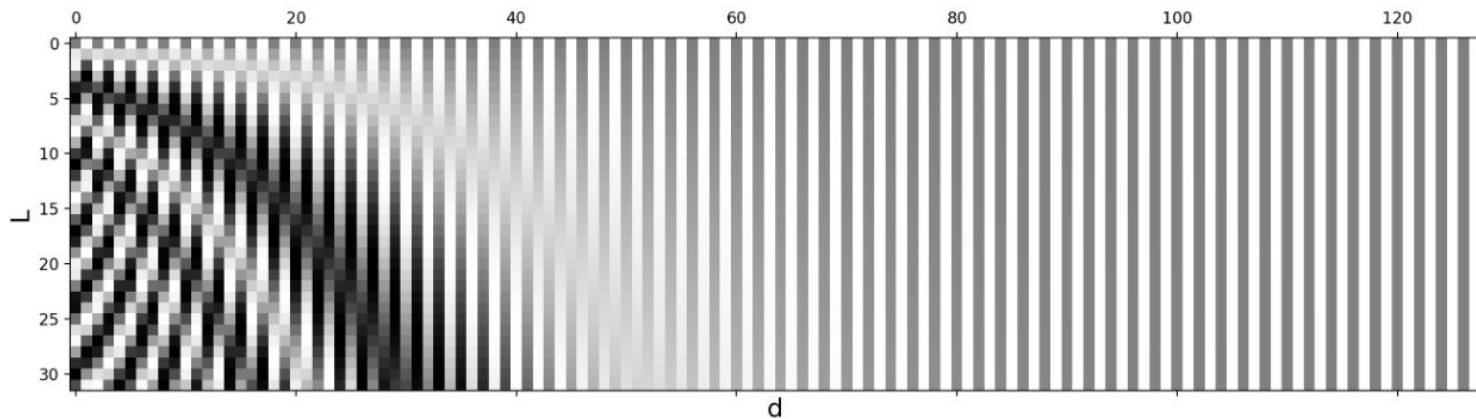
$$\text{E.g., } \mathbf{X}_p = \mathbf{X} + \mathbf{P}, \text{ or } \mathbf{X}_p = [\mathbf{X}, \mathbf{P}]$$

- $\mathbf{P}$  can be pre-defined, or made learnable

# Sinusoidal positional encoding

$$\text{PE}(i, \delta) = \begin{cases} \sin\left(\frac{i}{10000^{2\delta'/d}}\right) & \text{if } \delta = 2\delta' \\ \cos\left(\frac{i}{10000^{2\delta'/d}}\right) & \text{if } \delta = 2\delta' + 1 \end{cases}$$

L: sequence length  
d: embedding dimension



# Decoder

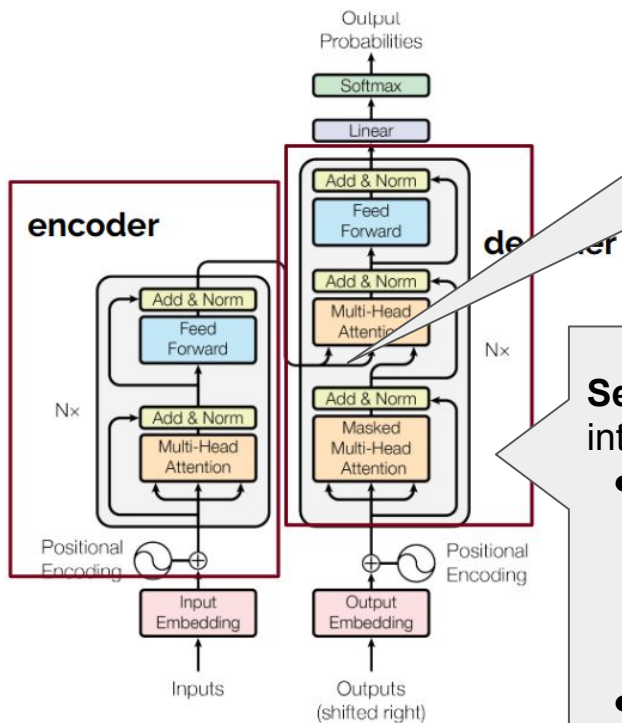
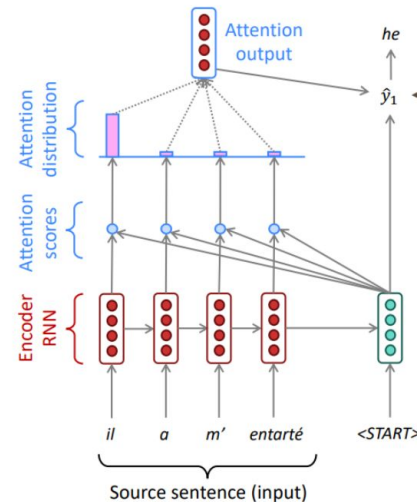


Figure 1: The Transformer - model architecture.

**Cross-attention** (to model the interaction between the encoder key-values and the current decoder query)

**Self-attention** (to model the interaction within itself)

- Respect the sequential nature (e.g., language modeling, assuming access to the future is cheating! )
- Masked out future tokens



(Credit: Stanford CS231N)

we can look at these (not greyed out) words

For encoding these words

	[START]	The	chef	who
[START]	-∞	-∞	-∞	-∞
The		-∞	-∞	-∞
chef			-∞	-∞
who				-∞

Attention matrix

# Strong performance in machine translation

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>	$2.3 \cdot 10^{19}$	

# Computation

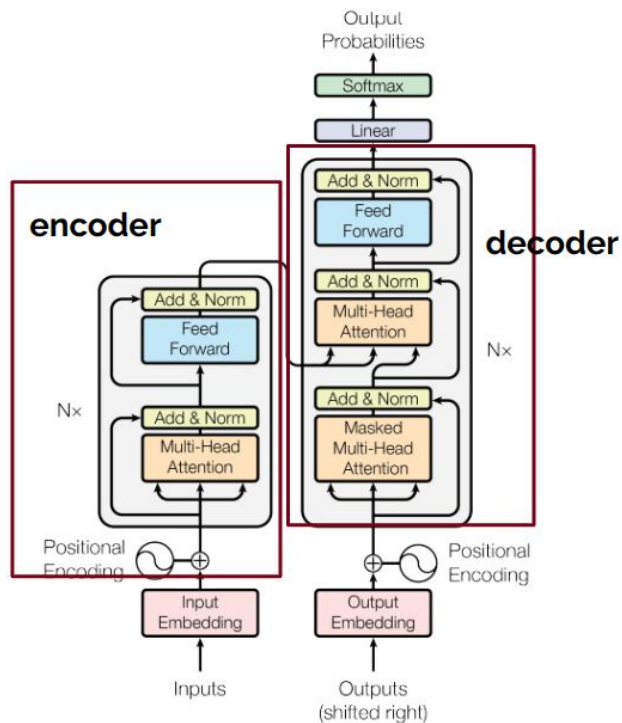


Figure 1: The Transformer - model architecture.

What's the total computation?

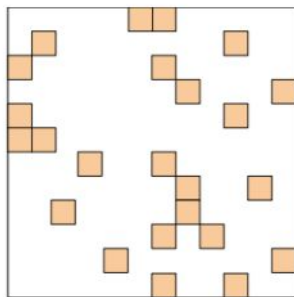
$$Q = XW^Q, \quad K = XW^K, \quad V = XW^V$$

$$\text{output} = \text{softmax}(QK^T / \sqrt{d_k})V$$

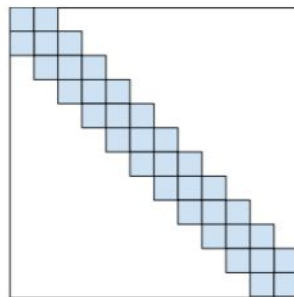
$$O(T^2 d)$$

Quadratic computation vs. linear computation in RNNs ( $T$  is the length of each input sequence,  $d$  is the embedding dimension)

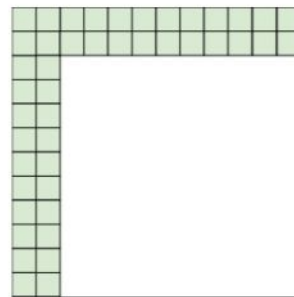
# Computation



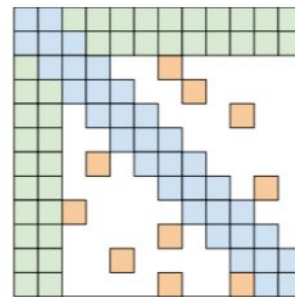
(a) Random attention



(b) Window attention



(c) Global Attention



(d) BIGBIRD

**Idea; building in sparsity** <https://arxiv.org/abs/2007.14062>

## Do Transformer Modifications Transfer Across Implementations and Applications?

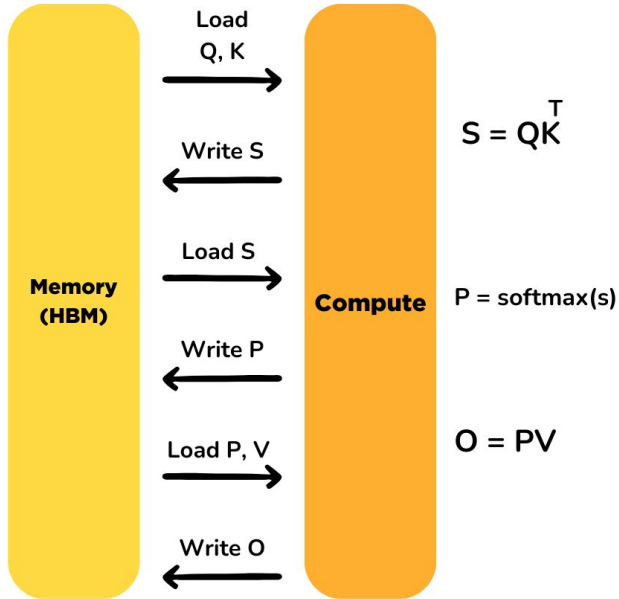
Sharan Narang*	Hyung Won Chung	Yi Tay	William Fedus
Thibault Fevry†	Michael Matena†	Karishma Malkan†	Noah Fiedel
Noam Shazeer	Zhenzhong Lan†	Yanqi Zhou	Wei Li
Nan Ding	Jake Marcus	Adam Roberts	Colin Raffel†

**But not much consistent improvement so far**

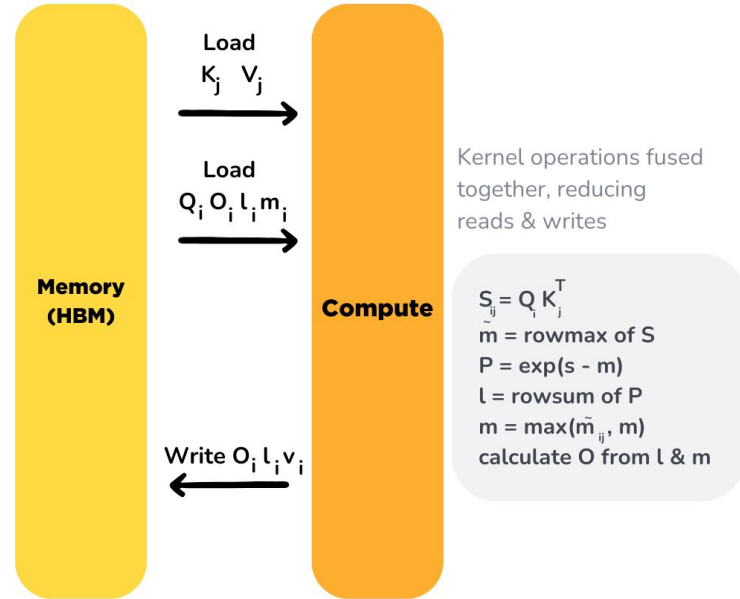
<https://arxiv.org/abs/2102.11972>

# Hardware-level acceleration: Flash Attention

Standard Attention Implementation



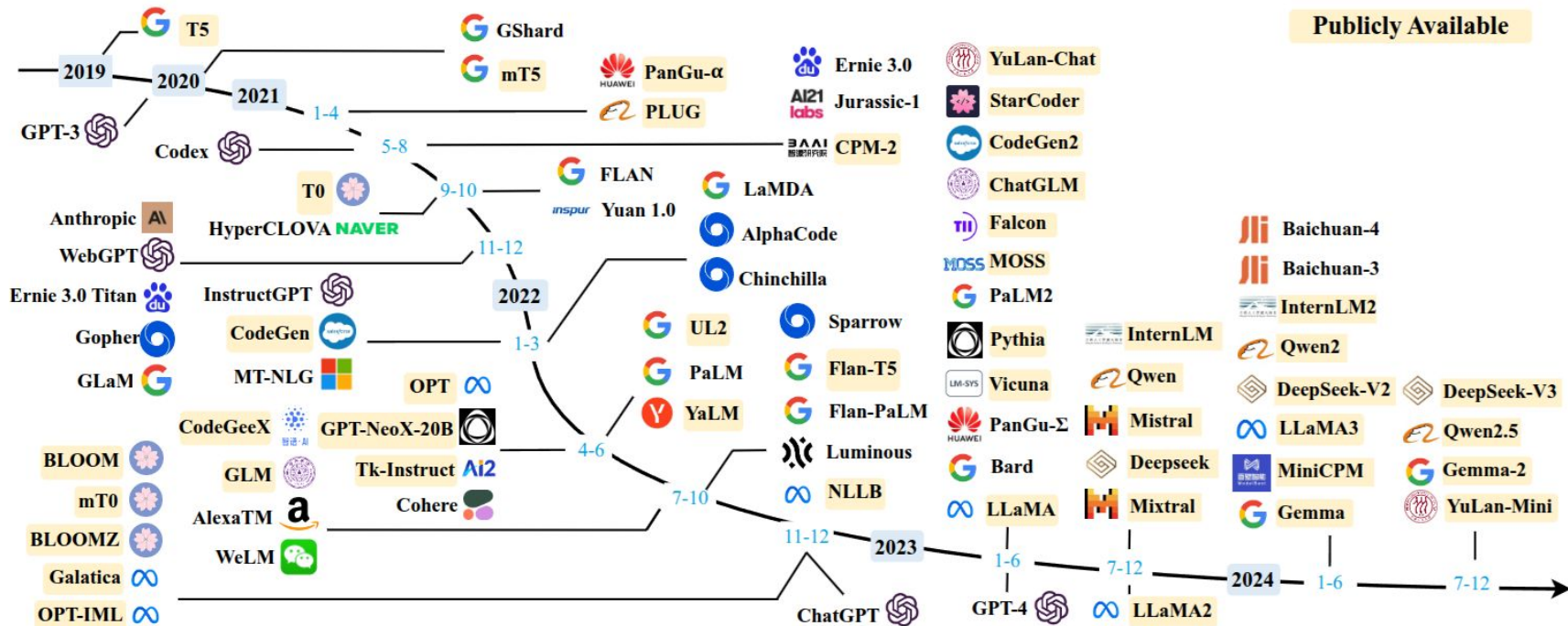
Flash Attention



Initialize O, l and m matrices with zeroes. m and l are used to calculate cumulative softmax. Divide Q, K, V into blocks (due to SRAM's memory limits) and iterate over them, for i is row & j is column.

# Large language models (LLMs)

# Transformers reign in NLP!



# LLMs: large models trained on large datasets

	Model	Release Time	Size (B)
	T5 [73]	Oct-2019	11
	mT5 [74]	Oct-2020	13
	PanGu- $\alpha$ [75]	Apr-2021	13*
	CPM-2 [76]	Jun-2021	198
	T0 [28]	Oct-2021	11
	CodeGen [77]	Mar-2022	16
	GPT-NeoX-20B [78]	Apr-2022	20
	Tk-Instruct [79]	Apr-2022	11
	UL2 [80]	May-2022	20
	OPT [81]	May-2022	175
	NLLB [82]	Jul-2022	54.5
	CodeGeeX [83]	Sep-2022	13
	GLM [84]	Oct-2022	130
	Flan-T5 [64]	Oct-2022	11
Publicly Available	BLOOM [69]	Nov-2022	176
	mT0 [85]	Nov-2022	13
	Galactica [35]	Nov-2022	120
	BLOOMZ [85]	Nov-2022	176
	OPT-IML [86]	Dec-2022	175
	LLaMA [57]	Feb-2023	65
	Pythia [87]	Apr-2023	12
	CodeGen2 [88]	May-2023	16
	StarCoder [89]	May-2023	15.5
	LLaMA2 [90]	Jul-2023	70

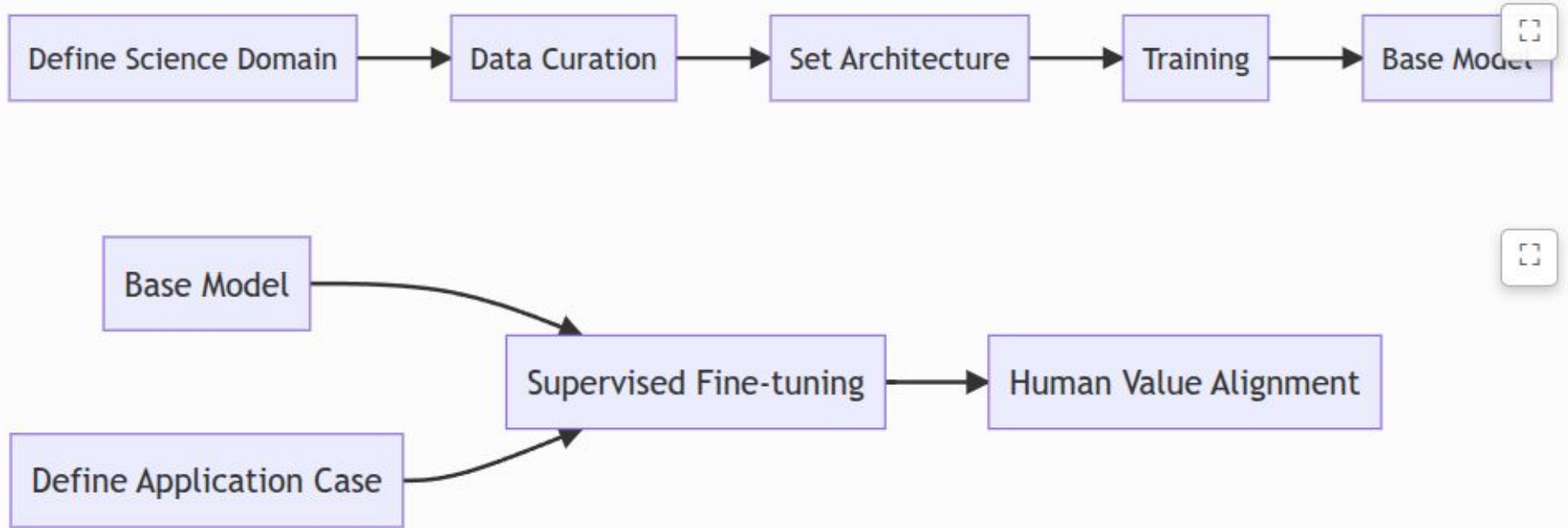
	GPT-3 [55]	May-2020	175
	GShard [91]	Jun-2020	600
	Codex [92]	Jul-2021	12
	ERNIE 3.0 [93]	Jul-2021	10
	Jurassic-1 [94]	Aug-2021	178
	HyperCLOVA [95]	Sep-2021	82
	FLAN [62]	Sep-2021	137
	Yuan 1.0 [96]	Oct-2021	245
	Anthropic [97]	Dec-2021	52
	WebGPT [72]	Dec-2021	175
	Gopher [59]	Dec-2021	280
	ERNIE 3.0 Titan [98]	Dec-2021	260
	GLaM [99]	Dec-2021	1200
	LaMDA [63]	Jan-2022	137
	MT-NLG [100]	Jan-2022	530
Closed Source	AlphaCode [101]	Feb-2022	41
	InstructGPT [61]	Mar-2022	175
	Chinchilla [34]	Mar-2022	70
	PaLM [56]	Apr-2022	540
	AlexaTM [102]	Aug-2022	20
	Sparrow [103]	Sep-2022	70
	WeLM [104]	Sep-2022	10
	U-PaLM [105]	Oct-2022	540
	Flan-PaLM [64]	Oct-2022	540
	Flan-U-PaLM [64]	Oct-2022	540
	GPT-4 [46]	Mar-2023	-
	PanGu- $\Sigma$ [106]	Mar-2023	1085
	PaLM2 [107]	Mar-2023	16

# LLMs: large models trained on large datasets

TABLE 2: Statistics of commonly-used data sources.

Corpora	Size	Source	Latest Update Time
BookCorpus [138]	5GB	Books	Dec-2015
Gutenberg [139]	-	Books	Dec-2021
C4 [73]	800GB	CommonCrawl	Apr-2019
CC-Stories-R [140]	31GB	CommonCrawl	Sep-2019
CC-NEWS [27]	78GB	CommonCrawl	Feb-2019
REALNEWS [141]	120GB	CommonCrawl	Apr-2019
OpenWebText [142]	38GB	Reddit links	Mar-2023
Pushift.io [143]	2TB	Reddit links	Mar-2023
Wikipedia [144]	21GB	Wikipedia	Mar-2023
BigQuery [145]	-	Codes	Mar-2023
the Pile [146]	800GB	Other	Dec-2020
ROOTS [147]	1.6TB	Other	Jun-2022

# Developing LLMs in practice



# Three crucial technical steps toward LLMs

- **Pretraining**
- **Supervised finetuning**
- **Alignment**

# Pretraining: data collection

14

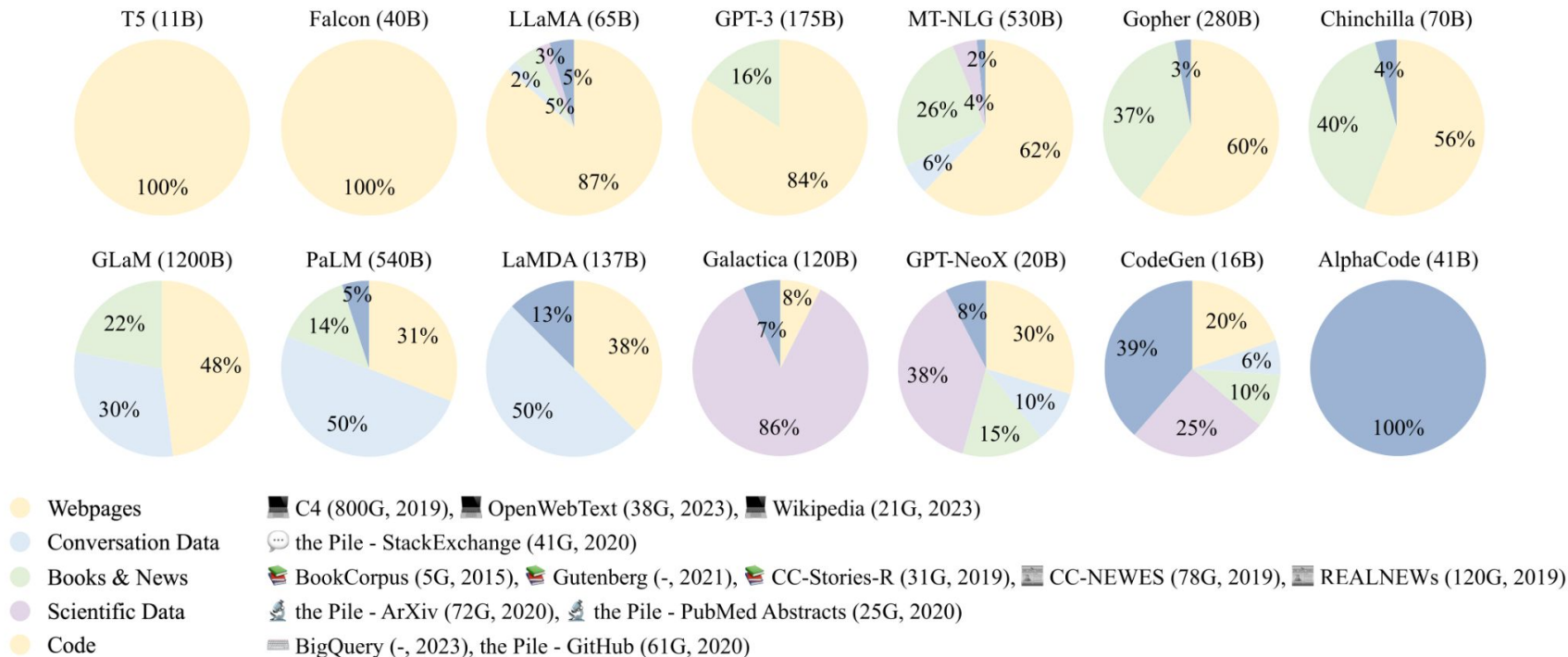


Image credit: A Survey of Large Language Models <https://arxiv.org/abs/2303.18223>

37

# Pretraining: data pre-processing

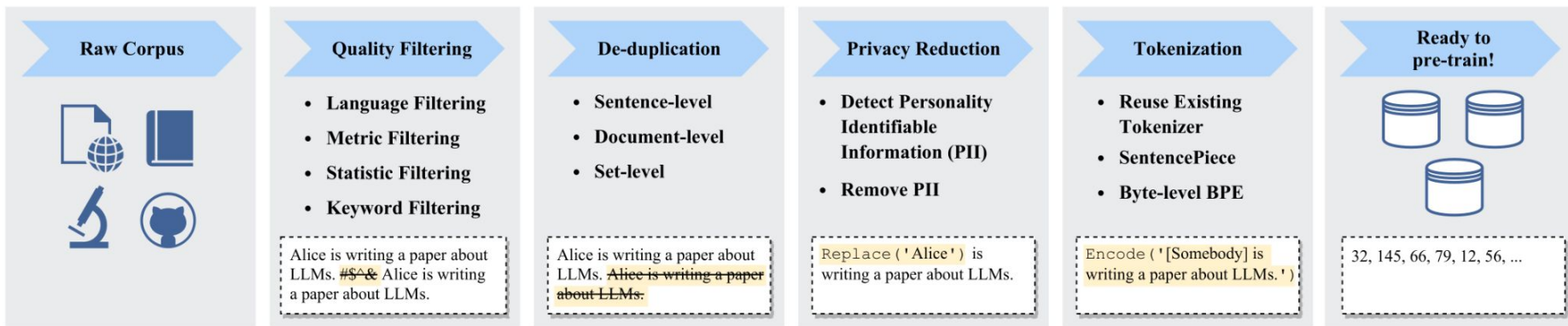
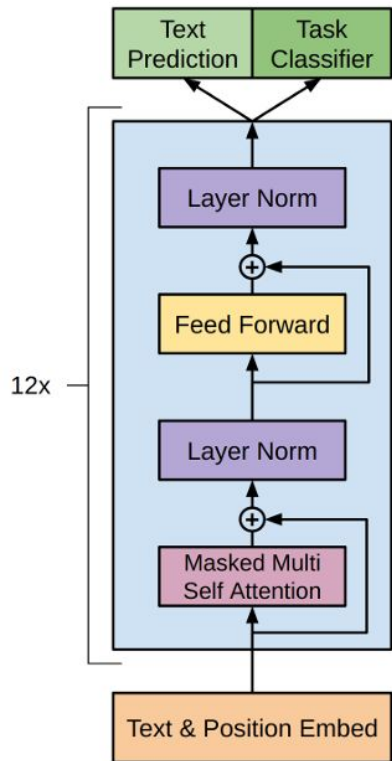


Fig. 6: An illustration of a typical data preprocessing pipeline for pre-training large language models.

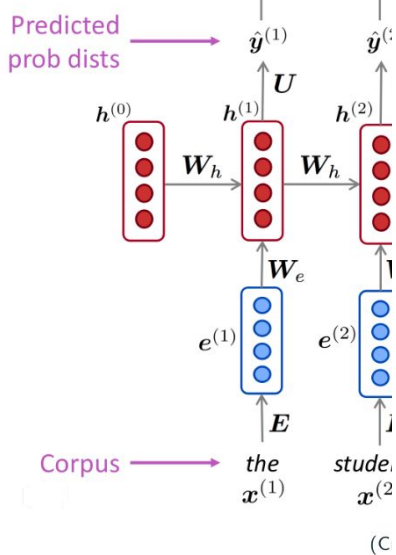
# Pretraining: architecture & task



Most popular: (transformer-based) **decoder-only** architectures pretrained on **language modeling**, i.e. model

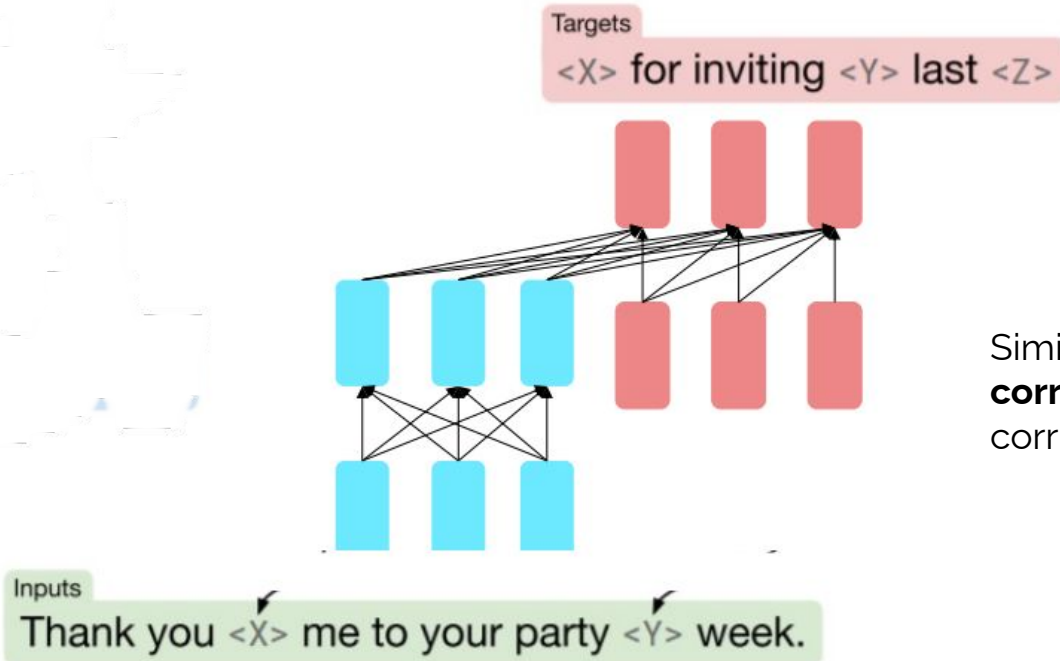
$$\mathbb{P} [x^{(t+1)} | x^{(t)}, \dots, x^{(1)}]$$

$$\text{Loss} \rightarrow J^{(1)}(\theta) + J^{(2)}(\theta) + J^{(3)}(\theta) + J^{(4)}(\theta) + \dots = J(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)}(\theta)$$



Model	Category	Size
GPT3 [55]	Causal decoder	175B
PanGU- $\alpha$ [75]	Causal decoder	207B
OPT [81]	Causal decoder	175B
PaLM [56]	Causal decoder	540B
BLOOM [69]	Causal decoder	176B
MT-NLG [100]	Causal decoder	530B
Gopher [59]	Causal decoder	280B
Chinchilla [34]	Causal decoder	70B
Galactica [35]	Causal decoder	120B
LaMDA [63]	Causal decoder	137B
Jurassic-1 [94]	Causal decoder	178B
LLaMA [57]	Causal decoder	65B
LLaMA 2 [90]	Causal decoder	70B
Falcon [127]	Causal decoder	40B
GLM-130B [84]	Prefix decoder	130B
T5 [73]	Encoder-decoder	11B

# Pretraining: architecture & task – alternative



Similar to pretraining encoder,  
**corruption removal!** (called span corruption)

# Pretraining: architecture details

Configuration	Method	Equation
Normalization position	Post Norm [22] Pre Norm [26] Sandwich Norm [201]	$\text{Norm}(\mathbf{x} + \text{Sublayer}(\mathbf{x}))$ $\mathbf{x} + \text{Sublayer}(\text{Norm}(\mathbf{x}))$ $\mathbf{x} + \text{Norm}(\text{Sublayer}(\text{Norm}(\mathbf{x})))$
Normalization method	LayerNorm [202] RMSNorm [203] DeepNorm [204]	$\frac{\mathbf{x} - \mu}{\sqrt{\sigma}} \cdot \gamma + \beta, \quad \mu = \frac{1}{d} \sum_{i=1}^d x_i, \quad \sigma = \sqrt{\frac{1}{d} \sum_{i=1}^d (x_i - \mu)^2}$ $\frac{\mathbf{x}}{\text{RMS}(\mathbf{x})} \cdot \gamma, \quad \text{RMS}(\mathbf{x}) = \sqrt{\frac{1}{d} \sum_{i=1}^d x_i^2}$ LayerNorm( $\alpha \cdot \mathbf{x} + \text{Sublayer}(\mathbf{x})$ )
Activation function	ReLU [205] GeLU [206] Swish [207] SwiGLU [208] GeGLU [208]	$\text{ReLU}(\mathbf{x}) = \max(\mathbf{x}, \mathbf{0})$ $\text{GeLU}(\mathbf{x}) = 0.5\mathbf{x} \otimes [1 + \text{erf}(\mathbf{x}/\sqrt{2})], \quad \text{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$ $\text{Swish}(\mathbf{x}) = \mathbf{x} \otimes \text{sigmoid}(\mathbf{x})$ $\text{SwiGLU}(\mathbf{x}_1, \mathbf{x}_2) = \text{Swish}(\mathbf{x}_1) \otimes \mathbf{x}_2$ $\text{GeGLU}(\mathbf{x}_1, \mathbf{x}_2) = \text{GeLU}(\mathbf{x}_1) \otimes \mathbf{x}_2$
Position embedding	Absolute [22] Relative [73] RoPE [209] Alibi [210]	$\mathbf{x}_i = \mathbf{x}_i + \mathbf{p}_i$ $A_{ij} = \mathbf{W}_q \mathbf{x}_i \mathbf{x}_j^T \mathbf{W}_k^T + r_{i-j}$ $A_{ij} = \mathbf{W}_q \mathbf{x}_i \mathbf{R}_{\theta, i-j} \mathbf{x}_j^T \mathbf{W}_k^T$ $A_{ij} = \mathbf{W}_q \mathbf{x}_i \mathbf{R}_{\theta, i-j} \mathbf{x}_j^T \mathbf{W}_k^T \quad A_{ij} = \mathbf{W}_q \mathbf{x}_i \mathbf{x}_j^T \mathbf{W}_k^T - m(i - j)$

# Pretraining: optimization details

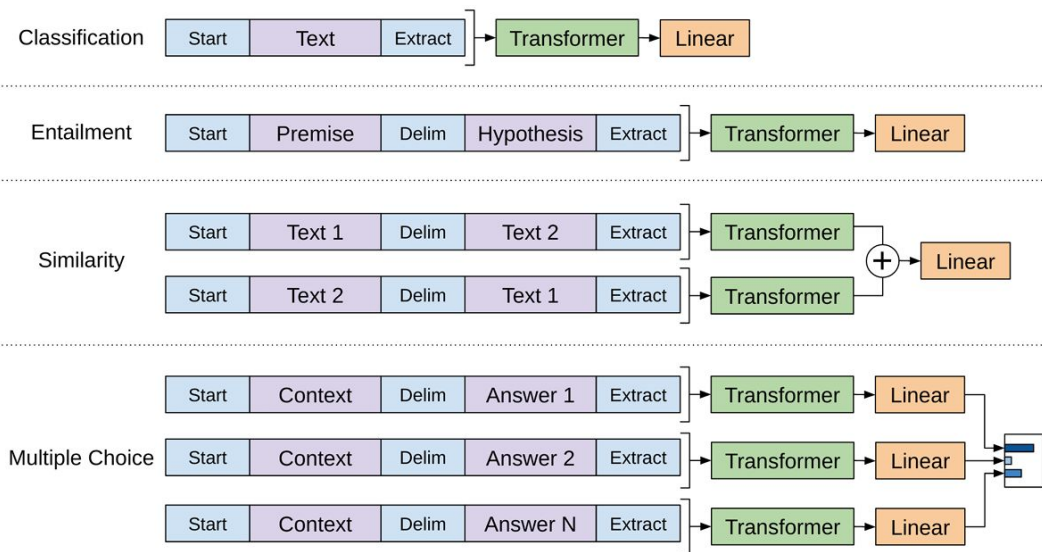
TABLE 5: Detailed optimization settings of several existing LLMs.

Model	Batch Size (#tokens)	Learning Rate	Warmup	Decay Method	Optimizer	Precision Type	Weight Decay	Grad Clip	Dropout
GPT3 (175B)	32K→3.2M	$6 \times 10^{-5}$	yes	cosine decay to 10%	Adam	FP16	0.1	1.0	-
PanGu- $\alpha$ (200B)	-	$2 \times 10^{-5}$	-	-	Adam	-	0.1	-	-
OPT (175B)	2M	$1.2 \times 10^{-4}$	yes	manual decay	AdamW	FP16	0.1	-	0.1
PaLM (540B)	1M→4M	$1 \times 10^{-2}$	no	inverse square root	Adafactor	BF16	$lr^2$	1.0	0.1
BLOOM (176B)	4M	$6 \times 10^{-5}$	yes	cosine decay to 10%	Adam	BF16	0.1	1.0	0.0
MT-NLG (530B)	64 K→3.75M	$5 \times 10^{-5}$	yes	cosine decay to 10%	Adam	BF16	0.1	1.0	-
Gopher (280B)	3M→6M	$4 \times 10^{-5}$	yes	cosine decay to 10%	Adam	BF16	-	1.0	-
Chinchilla (70B)	1.5M→3M	$1 \times 10^{-4}$	yes	cosine decay to 10%	AdamW	BF16	-	-	-
Galactica (120B)	2M	$7 \times 10^{-6}$	yes	linear decay to 10%	AdamW	-	0.1	1.0	0.1
LaMDA (137B)	256K	-	-	-	-	BF16	-	-	-
Jurassic-1 (178B)	32 K→3.2M	$6 \times 10^{-5}$	yes	-	-	-	-	-	-
LLaMA (65B)	4M	$1.5 \times 10^{-4}$	yes	cosine decay to 10%	AdamW	-	0.1	1.0	-
LLaMA 2 (70B)	4M	$1.5 \times 10^{-4}$	yes	cosine decay to 10%	AdamW	-	0.1	1.0	-
Falcon (40B)	2M	$1.85 \times 10^{-4}$	yes	cosine decay to 10%	AdamW	BF16	0.1	-	-
GLM (130B)	0.4M→8.25M	$8 \times 10^{-5}$	yes	cosine decay to 10%	AdamW	FP16	0.1	1.0	0.1
T5 (11B)	64K	$1 \times 10^{-2}$	no	inverse square root	AdaFactor	-	-	-	0.1
ERNIE 3.0 Titan (260B)	-	$1 \times 10^{-4}$	-	-	Adam	FP16	0.1	1.0	-
PanGu- $\Sigma$ (1.085T)	0.5M	$2 \times 10^{-5}$	yes	-	Adam	FP16	-	-	-

# Supervised finetuning

TABLE 6: A detailed list of available collections for instruction tuning.

Categories	Collections	Time	#Examples
Task	Nat. Inst. <a href="#">264</a>	Apr-2021	193K
	FLAN <a href="#">62</a>	Sep-2021	4.4M
	P3 <a href="#">265</a>	Oct-2021	12.1M
	Super Nat. Inst. <a href="#">79</a>	Apr-2022	5M
	MVPCorpus <a href="#">266</a>	Jun-2022	41M
	xP3 <a href="#">85</a>	Nov-2022	81M
	OIC <a href="#">22</a>	Mar-2023	43M
Chat	HH-RLHF <a href="#">267</a>	Apr-2022	160K
	HC3 <a href="#">268</a>	Jan-2023	87K
	ShareGPT <a href="#">23</a>	Mar-2023	90K
	Dolly <a href="#">24</a>	Apr-2023	15K
	OpenAssistant <a href="#">269</a>	Apr-2023	161K
Synthetic	Self-Instruct <a href="#">129</a>	Dec-2022	82K
	Alpaca <a href="#">123</a>	Mar-2023	52K
	Guanaco <a href="#">25</a>	Mar-2023	535K
	Baize <a href="#">270</a>	Apr-2023	158K
	BELLE <a href="#">271</a>	Apr-2023	1.5M



# Constructing the instruction sets

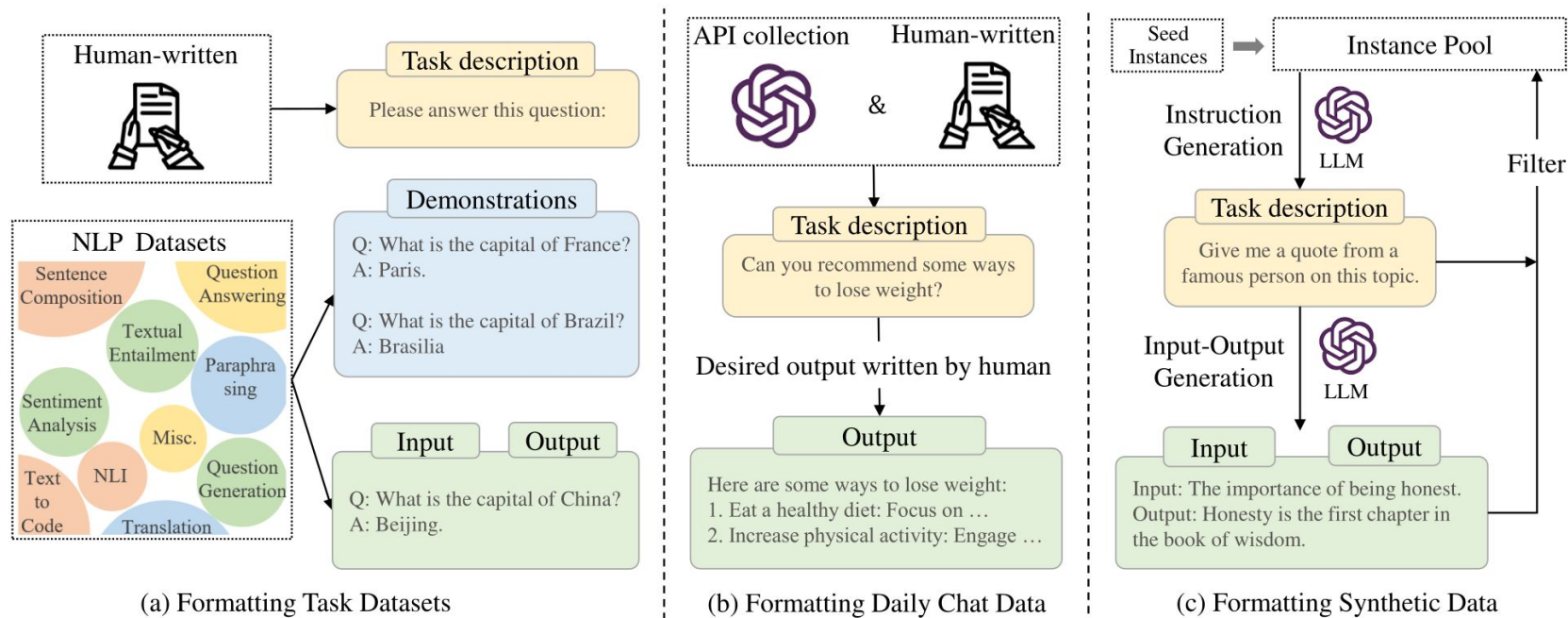
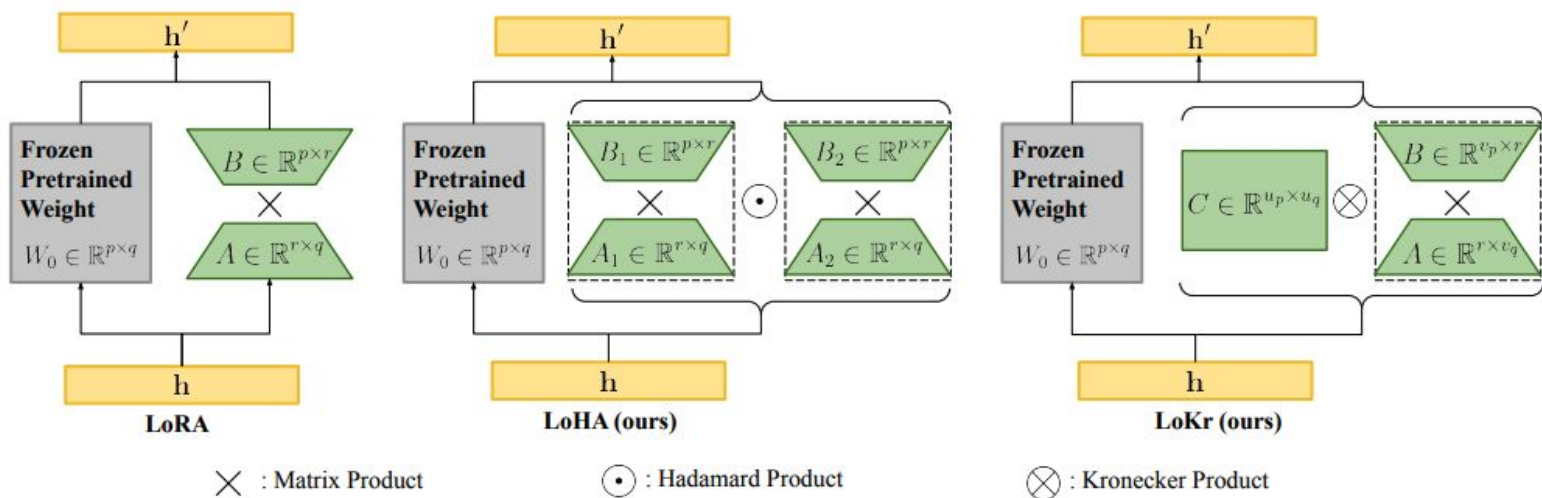


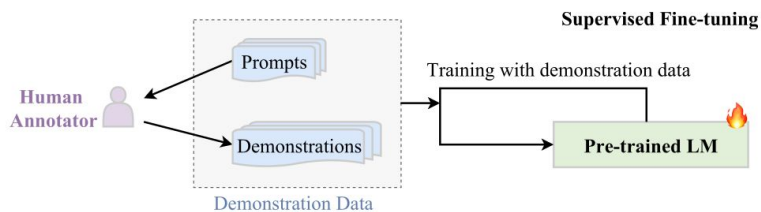
Fig. 9: An illustration of instance formatting and three different methods for constructing the instruction-formatted instances.

# Parameter-efficient finetuning

## Low-rank adaptation (LoRA)

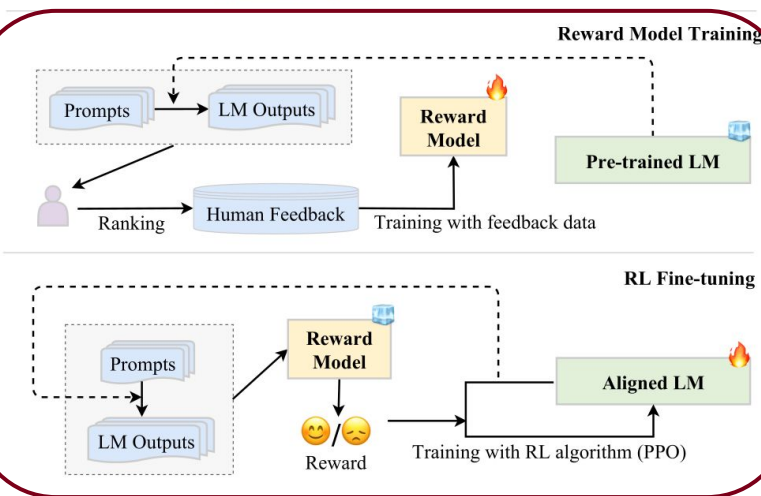


# Alignment



Make sure the output is aligned with human values/scientific truth/etc and not harmful

## Reinforcement learning with human feedback (RLHF)



OpenAI

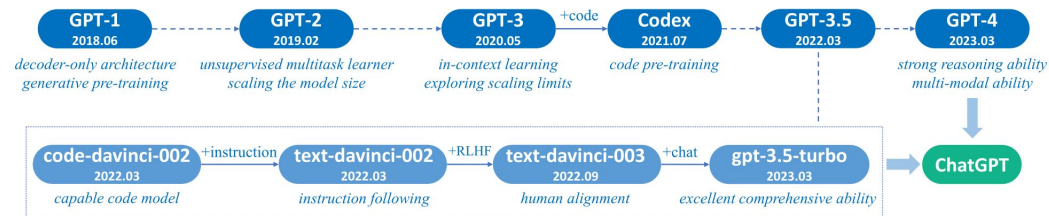


Fig. 3: A brief illustration for the technical evolution of GPT-series models. We plot this figure mainly based on the papers, blog articles and official APIs from OpenAI. Here, *solid lines* denote that there exists an explicit evidence (e.g., the official statement that a new model is developed based on a base model) on the evolution path between two models, while *dashed lines* denote a relatively weaker evolution relation.

# Alignment: details

$$\max_{p \in \mathbb{P}} \mathbb{E}_{x \sim \mathcal{D}, y \sim p(\cdot | x)} \left\{ R(x, y) - \beta \cdot \text{KL}(p(\cdot | x) \| p_0(\cdot | x)) \right\}$$

reward

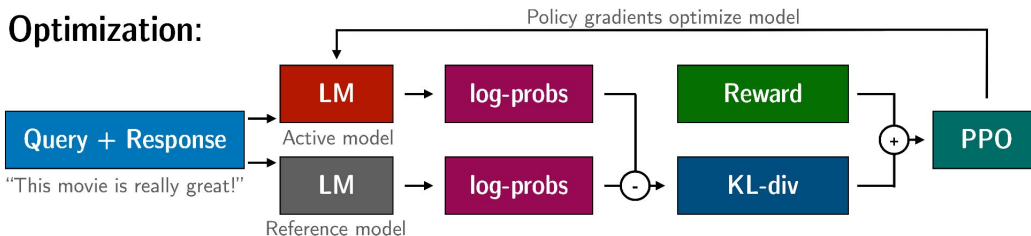
## Rollout:



## Evaluation:

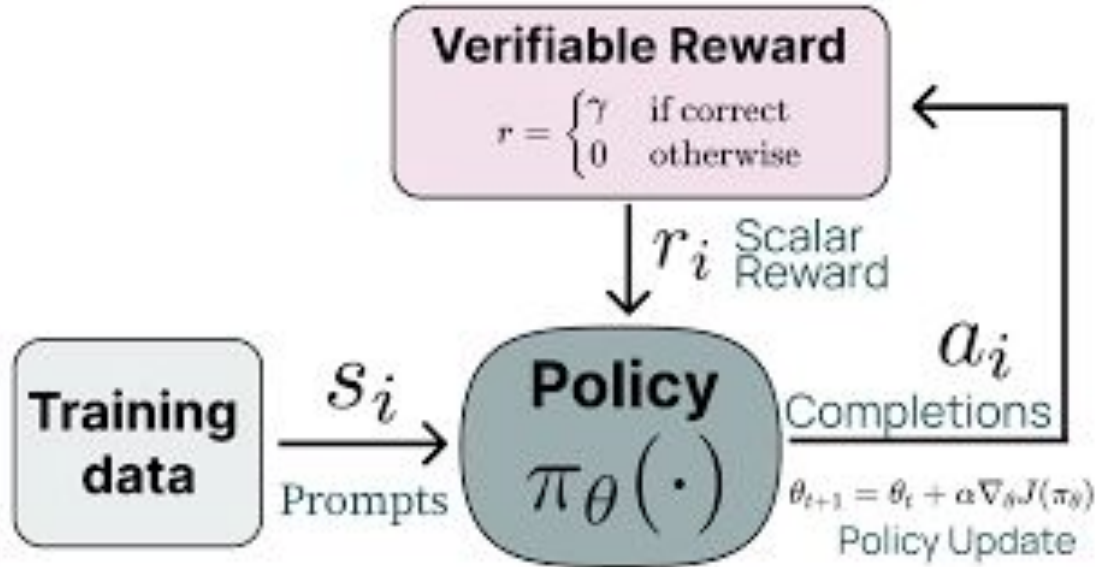


## Optimization:



Regularization to penalize over-alignment

# Alignment: RL with verifiable rewards (RLVR)



Incentivizes chain-of-thoughts reasoning on math/coding/etc tasks

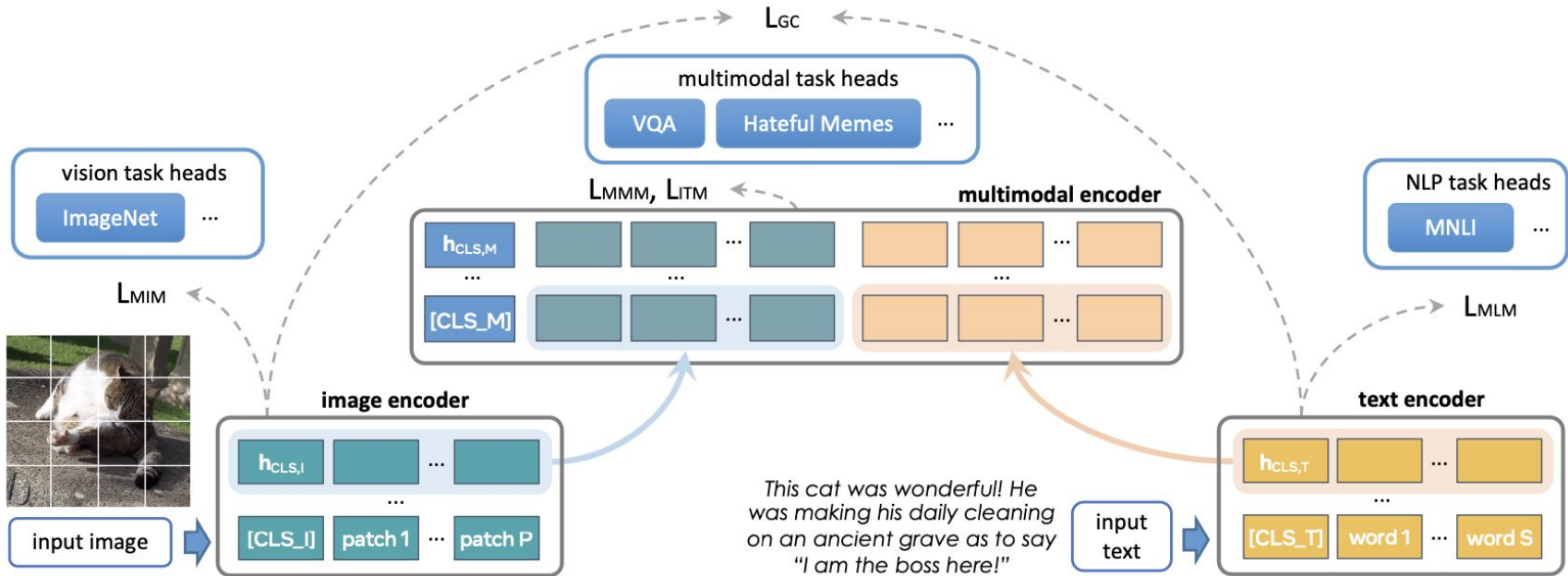
# Transformers for other domains

# Vision Transformers

*Transformers | Davide Cocomini | 2021*

[https://en.wikipedia.org/wiki/Vision\\_transformer](https://en.wikipedia.org/wiki/Vision_transformer)

# Multimodal foundation models



<https://pytorch.org/blog/scaling-multimodal-foundation-models-in-torchmultimodal-with-pytorch-distributed/>