Large Language Models

Jian Tang

HEC Montreal

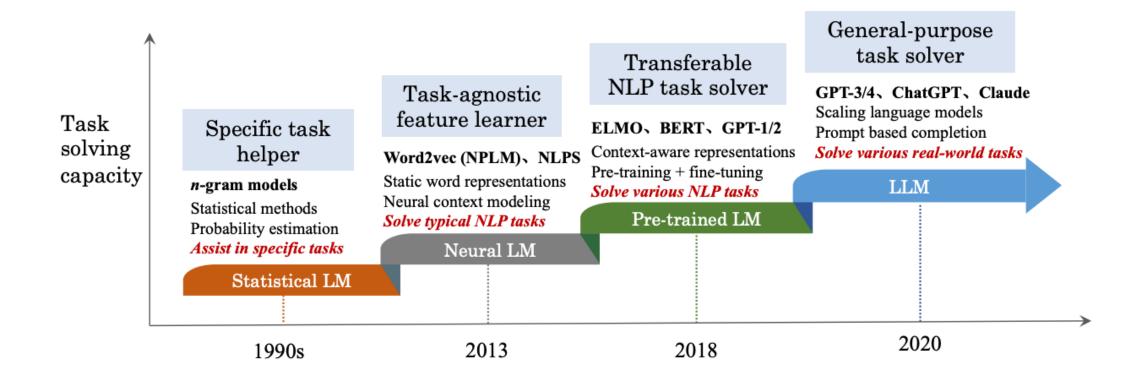
Mila-Quebec AI Institute

Email: jian.tang@hec.ca





History of NLP Techniques



Natural Language Modeling

• Given a text sequence $X = x_1 x_2 \dots x_T$, we want to model the joint distribution of

• X can be a sentence or a document

Two ways of Modeling P(X)

- Masked Language Modeling: predicting the missing word(s) conditioning on the rest of the words, i.e. $P(x_i|X_{-i})$
 - Filling in the blank

$$x_1 \quad x_2 \quad x_3 \dots x_{i-1}$$
 ? $x_{i+1} \dots x_T$

Bidirectional modeling

Two ways of Modeling P(X)

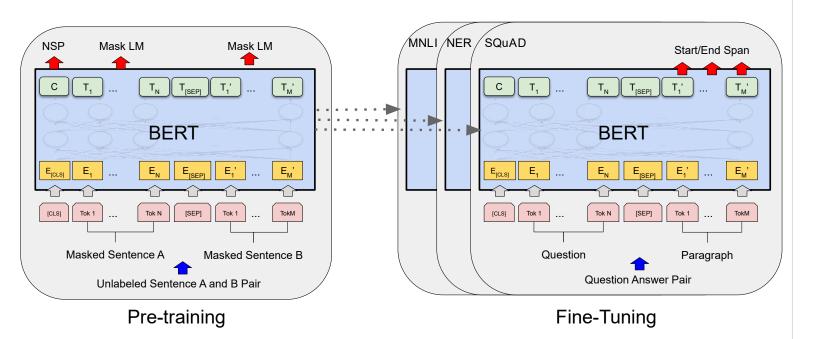
• Next Token Prediction: predicting the next token/word conditioning on the preceding tokens, i.e. $P(x_i|X_{< i})$

$$x_1 \quad x_2 \quad x_3 \dots x_{i-1}$$
 ?

Unidirectional modeling

BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

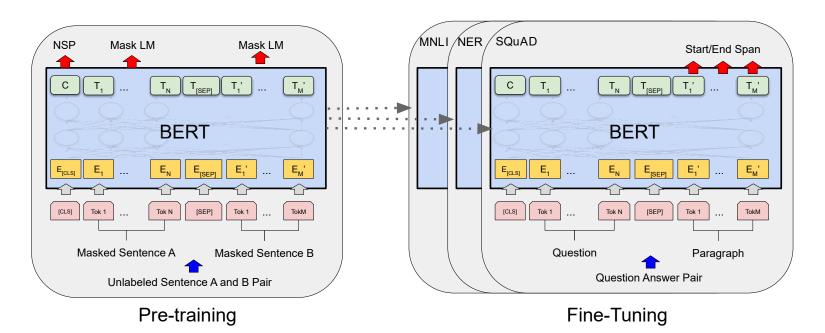
BERT: Pre-training and Fine-Tuning Framework



- Pre-train the model on unlabeled data over pre-training tasks
- Initializing the model with the pre-trained parameters, and fine-tune all the parameters using labeled data from downstream tasks
- Unified architecture across different tasks

Model Architecture

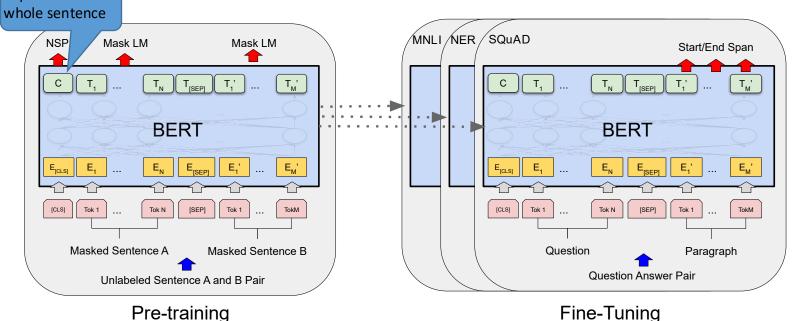
A multi-layer bidirectional Transformer encoder



- **BERT_BASE**: Number of layers: 12, Hidden Size: 768, the number of self-attention heads: 12, total number of parameters: 110M
- **BERT_LARGE**: Number of layers: 23, Hidden Size: 1024, the number of self-attention heads: 16, total number of parameters: 340M

Model Architecture

Input/Output Representations Representation of

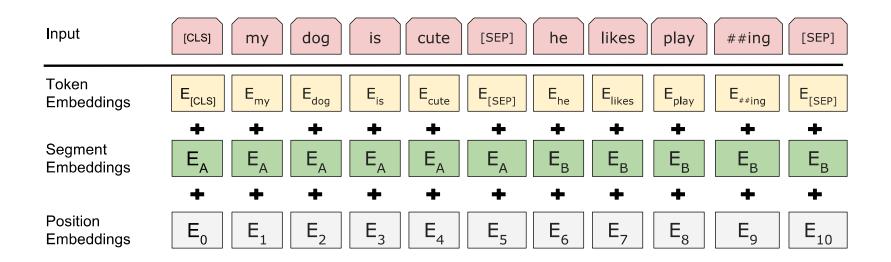


Fine-Tuning

- Both a single sentence and a pair of sentences are represented in one token sequence
 - Start with a special token [CLS]
 - Separate sentences with a special token [SEP]
 - Add a learned embeddings to every token indicating whether it belongs to sentence A or sentence B

Input Representation

 Each token embedding is the summation of token embedding, segment embedding, and position embeddings.



Pre-training BERT

- Task #1: Masked LM
- Mask some percentages of the input tokens at random, and then predict those masked tokens.
 - 15% of the tokens in each sequence are masked out at random

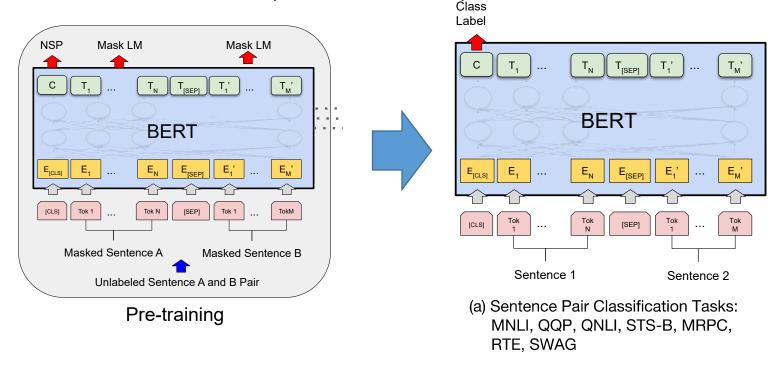
Pre-training BERT

- Task #2: Next Sentence Prediction (NSP)
 - Important in many downstream tasks such as question answering (QA) and natural language inference (NLI)
- 50% sentence pairs (A,B) are positive
 - Actually sentence B that follows A
- 50% sentence pairs (A, B) are negative
 - Randomly select a sentence B from the training corpus

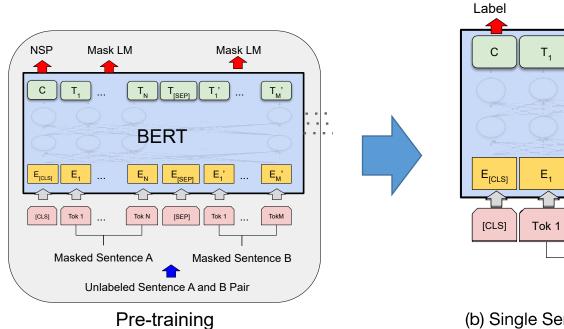
Pre-training Data

- BookCorpus (800M words)
- English Wikipedia(2,500M words)

- Sentence Pair Classification Tasks, e.g., nature language inference
 - Given a pair of sentences, the goal is to predict whether the second sentence is an *entailment*, *contradiction*, *or neutral* w.r.t. the first one



Sentence Classification



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

Single Sentence

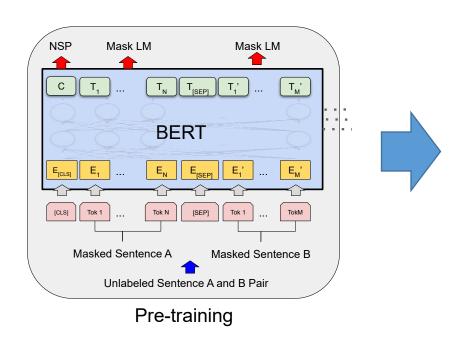
Tok 2

BERT

Tok N

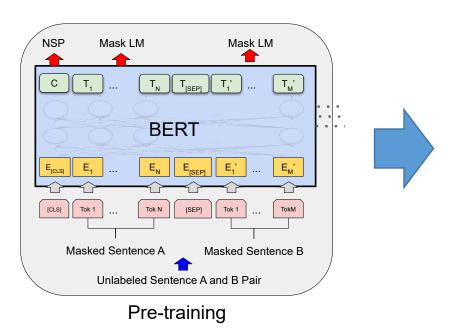
Class

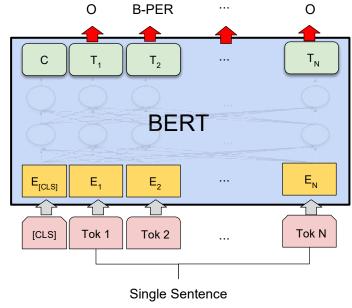
Question Answering



(c) Question Answering Tasks: SQuAD v1.1

Sentence Tagging





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

Experimental Results

• GLUE: General Language Understanding Evaluation Benchmark

System	MNLI-(m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Average
	392k	363k	108k	67k	8.5k	5.7k	3.5k	2.5k	-
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT _{BASE}	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
$BERT_{LARGE}$	86.7/85.9	72.1	92.7	94.9	60.5	86.5	89.3	70.1	82.1

Ablation Study

• LTR: left-to-right, unidirectional language modeling

]	Dev Set		
Tasks	MNLI-m	QNLI	MRPC	SST-2	SQuAD
	(Acc)	(Acc)	(Acc)	(Acc)	(F1)
$\overline{\mathrm{BERT}_{\mathrm{BASE}}}$	84.4	88.4	86.7	92.7	88.5
No NSP	83.9	84.9	86.5	92.6	87.9
LTR & No NSP	82.1	84.3	77.5	92.1	77.8
+ BiLSTM	82.1	84.1	75.7	91.6	84.9

Effect of Model Size

• L: number of layers

• H: hidden dimension

• A: number of self-attention heads

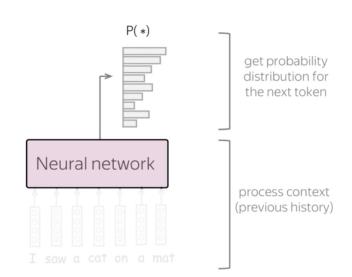
Ну	perpar	ams		Dev Set Accuracy				
#L	#H	#A	LM (ppl)	MNLI-m	MRPC	SST-2		
3	768	12	5.84	77.9	79.8	88.4		
6	768 768	3 12	5.24 4.68	80.6 81.9	82.2 84.8	90.7 91.3		
12	768	12	3.99	84.4	86.7	92.9		
12 24	1024 1024	16 16	3.54 3.23	85.7 86.6	86.9 87.8	93.3 93.7		

Large Language Models

• Given a sequence $X = x_1 x_2 \dots x_T$, Next Token Prediction: predicting the next token/word conditioning on the preceding tokens, i.e. $P(x_i|X_{< i})$

$$O = 1/T \sum_{i=1}^{T} \log P(x_i | X_{< i})$$

• In practice, we will specify a maximum context window



Next Token Prediction

- Essentially learning a mapping function f: context -> word
 - a classification problem

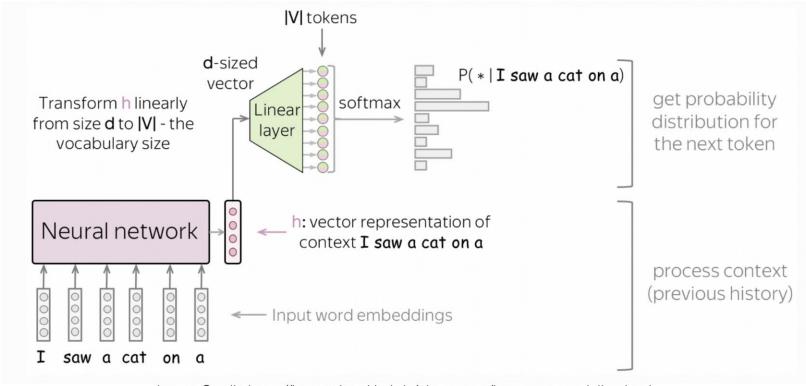


Image Credit: https://lena-voita.github.io/nlp_course/language_modeling.html

History of Generative LLMs

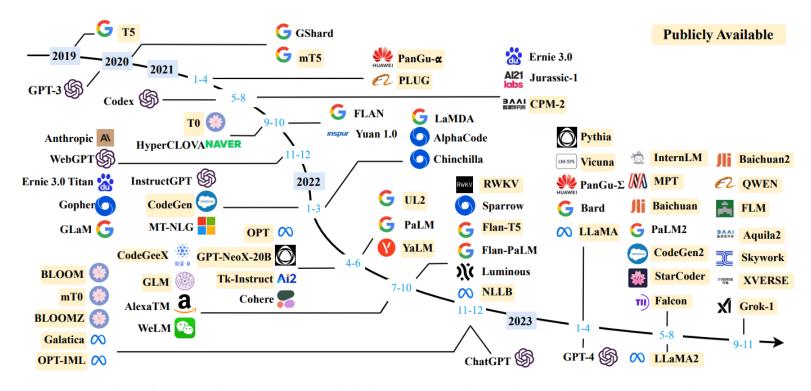
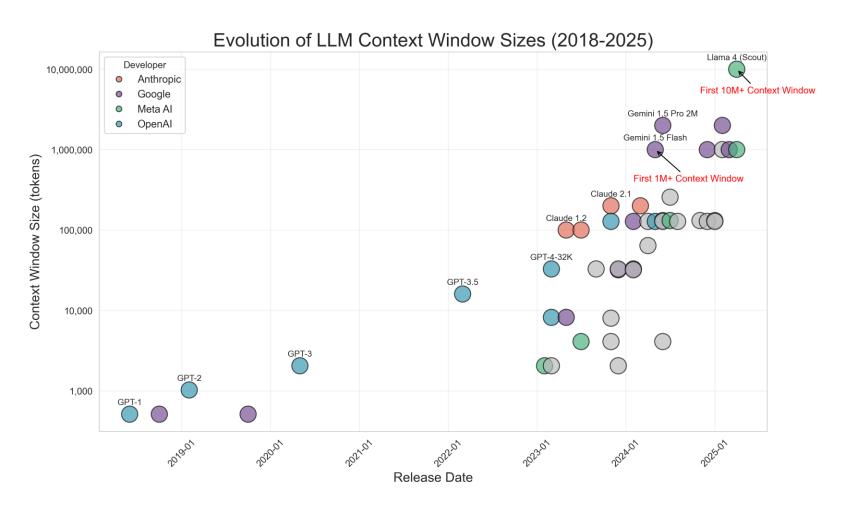


Fig. 3: A timeline of existing large language models (having a size larger than 10B) in recent years. The timeline was established mainly according to the release date (e.g., the submission date to arXiv) of the technical paper for a model. If there was not a corresponding paper, we set the date of a model as the earliest time of its public release or announcement. We mark the LLMs with publicly available model checkpoints in yellow color. Due to the space limit of the figure, we only include the LLMs with publicly reported evaluation results.

Mode Size over Time



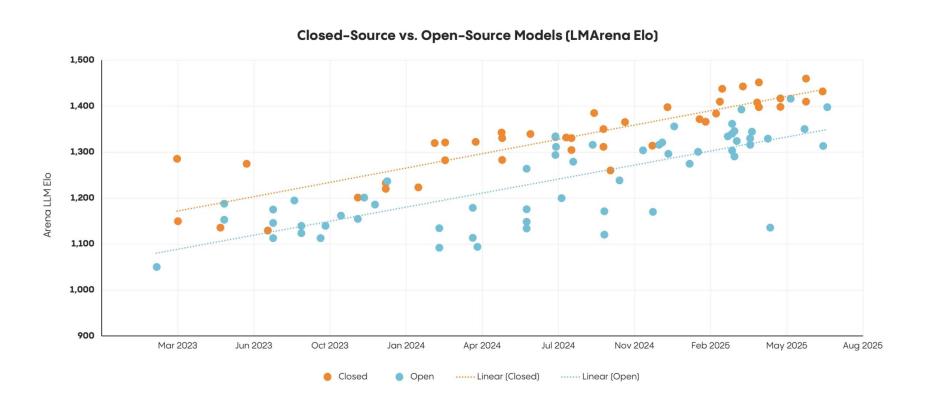
Context Size over Time



https://www.meibel.ai/post/understanding-the-impact-of-increasing-llm-context-windows

Open vs Closed LLMs

Closed-Source vs. Open-Source Models



© 2025 Menlo Ventures

Closed-Source LLMs

• GPT

• Gemini

• Claude





ANTHROP\C

View rankings across various LLMs on their versatility, linguistic precision, and cultural context across text.

Last Updated

Oct 16, 2025

Tot

Total Votes 4,278,480 Total Models 258

Y Overall	✓ Q Search by model name				Defau	lt
Rank (UB) ↑	Model ↑↓	Score ↑↓	95% CI (±) ↑↓	Votes ↑↓	Organization ↑↓	License ↑↓
1	G gemini-2.5-pro	1451	±4	54,087	Google	Proprietary
1	A\ claude-opus-4-1-20250805-thinking-16k	1447	±5	21,306	Anthropic	Proprietary
1	A\ claude-sonnet-4-5-20250929-thinking-32k	1445	±8	6,287	Anthropic	Proprietary
1	\$ gpt-4.5-preview-2025-02-27	1441	±6	14,644	OpenAl	Proprietary
2	\$\text{\$\text{\$\text{\$o\$}}} \text{ chatgpt-4o-latest-20250326}	1440	±4	40,013	OpenAl	Proprietary
2		1440	±4	51,293	OpenAl	Proprietary
2	A\ claude-sonnet-4-5-20250929	1438	±8	6,144	Anthropic	Proprietary
2		1437	±5	23,580	OpenAl	Proprietary
2	A\ claude-opus-4-1-20250805	1437	±5	33,298	Anthropic	Proprietary
3		1434	±6	18,078	Alibaba	Proprietary
10		1425	±5	21,630	OpenAl	Proprietary
10	\$\times qwen3-max-2025-09-23	1423	±7	6,919	Alibaba	Proprietary
10	ℤ glm-4.6	1422	±9	4,401	Z.ai	MIT
11	x grok-4-fast	1420	±8	7,104	xAI	Proprietary
11	A\ claude-opus-4-20250514-thinking-16k	1419	±5	35,522	Anthropic	Proprietary
11	♥ deepseek-v3.2-exp-thinking	1419	±9	4,320	DeepSeek Al	MIT
11	🏂 qwen3-vl-235b-a22b-instruct	1418	±8	6,312	Alibaba	Apache 2.0

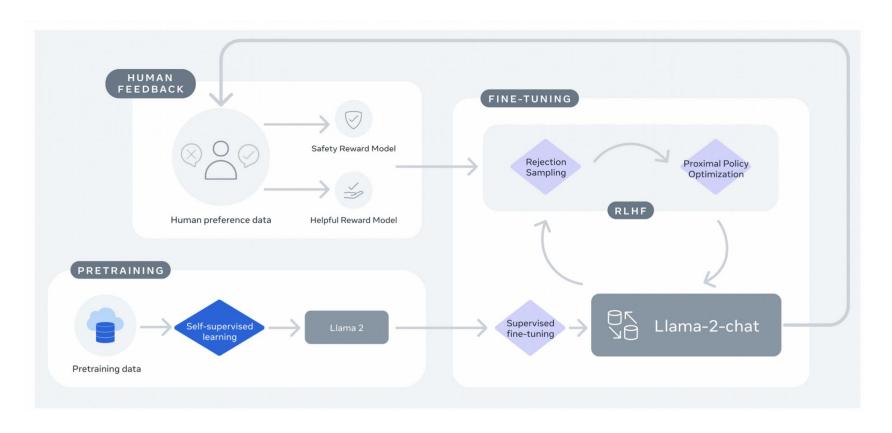
https://lmarena.ai/leaderboard/text

10	Z glm-4.6	1422	±9	4,401	Z.ai	MIT
11	X grok-4-fast	1420	±8	7,104	xAI	Proprietary
11	A\ claude-opus-4-20250514-thinking-16k	1419	±5	35,522	Anthropic	Proprietary
11	♥ deepseek-v3.2-exp-thinking	1419	±9	4,320	DeepSeek Al	MIT
11	🌣 qwen3-v1-235b-a22b-instruct	1418	±8	6,312	Alibaba	Apache 2.0
11	🍃 qwen3-235b-a22b-instruct-2507	1418	±5	29,343	Alibaba	Apache 2.0
11	❤ deepseek-r1-0528	1417	±6	19,284	DeepSeek	MIT
11	kimi-k2-0905-preview	1417	±7	10,772	Moonshot	Modified MIT
11	♥ deepseek-v3.1	1416	±6	15,380	DeepSeek	MIT
11	❤ deepseek-v3.1-thinking	1415	±7	12,098	DeepSeek	MIT
11	kimi-k2-0711-preview	1415	±5	28,321	Moonshot	Modified MIT
11	♥ deepseek-v3.1-terminus	1414	±10	3,775	DeepSeek Al	MIT
11	♥ deepseek-v3.1-terminus-thinking	1413	±10	3,541	DeepSeek Al	MIT
12	X grok-4-0709	1413	±5	29,264	xAI	Proprietary
12	A\ claude-opus-4-20250514	1411	±4	43,310	Anthropic	Proprietary
12	♥ deepseek-v3.2-exp	1408	±9	4,684	DeepSeek Al	MIT
13	⑤ gpt-4.1-2025-04-14	1411	±4	41,918	OpenAl	Proprietary
14	x grok-3-preview-02-24	1409	±4	34,154	xAI	Proprietary
18	№ mistral-medium-2508	1406	±5	23,844	Mistral	Proprietary
18	Z glm-4.5	1406	±5	22,612	Z.ai	MIT

Open-source LLMs

Overview of LLMs Training

 Pretraining -> Supervised Fine-tuning (SFT) -> Reinforcement Learning Human Feedback (RLHF)



Pre-training Data



Colossal Clean Crawled Corpus (C4)









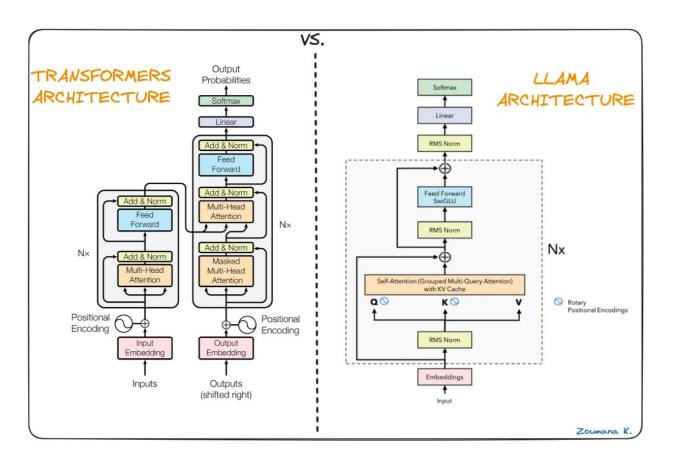
StackExchange

Example: Llama1

• 1.4 Trillion Tokens!!

Dataset	Sampling prop.	Epochs	Disk size
CommonCrawl	67.0%	1.10	3.3 TB
C4	15.0%	1.06	783 GB
Github	4.5%	0.64	328 GB
Wikipedia	4.5%	2.45	83 GB
Books	4.5%	2.23	85 GB
ArXiv	2.5%	1.06	92 GB
StackExchange	2.0%	1.03	78 GB

LlaMa Architecture



https://medium.com/@pranjalkhadka/llama-explained-a70e71e706e9

LLaMA:OpenandEfficient Foundation Language Models (https://arxiv.org/pdf/2302.13971)

Training Loss

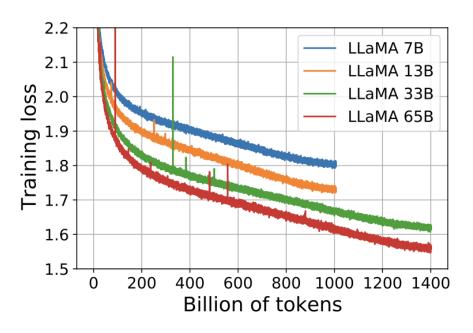


Figure 1: Training loss over train tokens for the 7B, 13B, 33B, and 65 models. LLaMA-33B and LLaMA-65B were trained on 1.4T tokens. The smaller models were trained on 1.0T tokens. All models are trained with a batch size of 4M tokens.





Play with Llama:

https://www.meta.ai/?utm_source=llama_meta_site&utm_medium=web&utm_content=Llama_nav <u>&utm_campaign=July_moment</u>

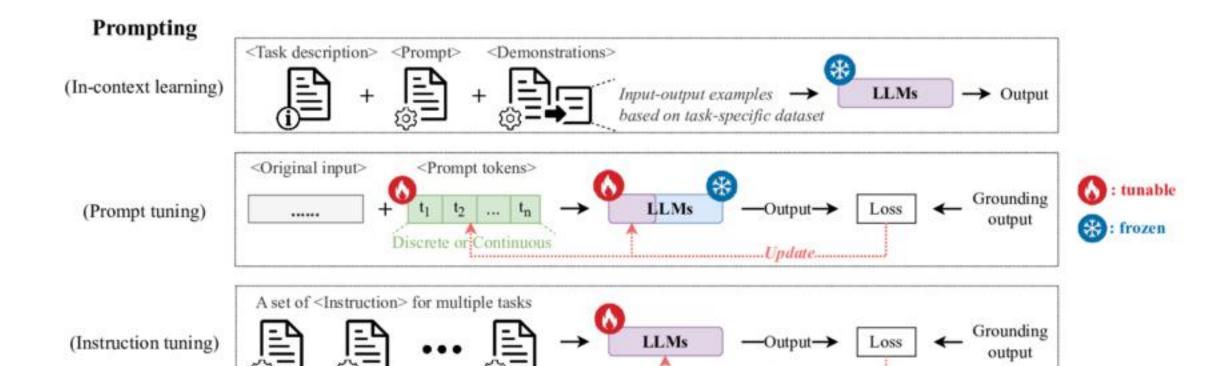
LoRa:

Efficient fine-tuning

Fine-tuning Llama 2 in Google Colab

 https://colab.research.google.com/drive/1wbPpB3fY9YRzebrZq6WkP 5HxMuHMQR72?usp=sharing

Different tasks in Natural Language Understanding



Thanks!