

CMPT 413/713: Natural Language Processing

# Scaling laws for LLMs (scaling LMs up and down)

Spring 2025 2025-03-31

Slides adapted from Anoop Sarkar

## New capabilities emerge at scale

**QUESTION ANSWERING** 

ARITHMETIC

https://ai.googleblog.com/2022/04/pathways-language-model-palm-scaling-to.html





#### Larger and larger language models GPTv4 (1.7T?, OpenAl) Megatron-Turing NLG (530B, MS+NVidia) PaLM (540B, Google) GPT-3 (175B) Bloom (176B, HuggingFace+BigScience) LaMDA (137B, Google) Megatron-LM (8.3B) Chinchilla (70B, Llama (65B, FAIR) Turing-NLG (17.2B) DeepMind) T5 (11B) GPT-2 (1.5B) BERT-Large (340M) ELMo (94M) 2022 2023 2019 2020 2021



### https://huggingface.co/blog/large-language-models





### 2023-2024 OPTIMAL LANGUAGE MODELS 2023



### https://lifearchitect.ai/models





Language Models are Few-Shot Learners (Brown et al. 2020)

### Total Compute Used During Training



# How does LLM performance scale as we increase model and data size?

### **Scaling Laws for Neural Language Models**

#### Jared Kaplan \*

Johns Hopkins University, OpenAI

jaredk@jhu.edu

#### **Tom Henighan**

#### Tom B. Brown

OpenAI

OpenAI

henighan@openai.com

#### -

#### **Scott Gray**

OpenAI

scott@openai.com

Alec Radford

OpenAI

alec@openai.com

https://arxiv.org/abs/2001.08361

Sam McCandlish\*

OpenAI sam@openai.com

n Benjamin Chess

**Rewon Child** 

OpenAI

tom@openai.com bchess@openai.com

OpenAI

rewon@openai.com

#### Jeffrey Wu

OpenAI jeffwu@openai.com

#### **Dario Amodei**

OpenAI

damodei@openai.com

Jan 2020

## **Scaling Laws for LLMs Power laws**

For LLMs, we are interested in how the test performance scales with relation to

- Model size: number of model parameters N (excluding subword embeddings)
- Data size: number of tokens trained on **D**
- Amount of compute (MFLOPs) C (1 PetaFLOP-day (PF-day) is  $8.64 \times 10^{19}$  FLOPS)

### Findings

- Model performance scales as **power law** of model size and data size
- Power law: relation between two quantities where one quantity increases as a power of another
  - $f(x) = (a/x)^k$  e.g. model performance vs. model size
- N, D, C are dominant. Other choices in hyperparameters like width vs. depth are less relevant



https://openai.com/research/ai-and-compute



### Model size: computing the number of parameters

Operation	Parameters	FLOPs per Token
Embed	$(n_{ m vocab}+n_{ m ctx})d_{ m model}$	$4d_{ m model}$
Attention: QKV	$n_{ m layer}d_{ m model}3d_{ m attn}$	$2n_{ m layer}d_{ m model}3d_{ m attn}$
Attention: Mask		$2n_{ m layer}n_{ m ctx}d_{ m attn}$
Attention: Project	$n_{ m layer}d_{ m attn}d_{ m model}$	$2n_{ m layer}d_{ m attn}d_{ m embd}$
Feedforward	$n_{ m layer}2d_{ m model}d_{ m ff}$	$2n_{ m layer}2d_{ m model}d_{ m ff}$
De-embed		$2d_{ m model}n_{ m vocab}$
Total (Non-Embedding)	$N = 2d_{\text{model}}n_{\text{layer}} \left(2d_{\text{attn}} + d_{\text{ff}}\right)$	$C_{\text{forward}} = 2N + 2n_{\text{layer}}n_{\text{ctx}}d_{\text{attr}}$

**Table 1** Parameter counts and compute (forward pass) estimates for a Transformer model. Sub-leading terms such as nonlinearities, biases, and layer normalization are omitted.



### Test performance



Figure 1 Language modeling performance improves smoothly as we increase the model size, datasetset size, and amount of compute<sup>2</sup> used for training. For optimal performance all three factors must be scaled up in tandem. Empirical performance has a power-law relationship with each individual factor when not bottlenecked by the other two.



### Excluding embeddings from parameter count

• Power law relationship not so clear when embeddings are included



# **Power laws for test loss**

- compute C
- For models with limited number of parameters:  $L(N) = (N_c/N)^{\alpha_N}$   $\alpha_N \approx 0.076$ ,  $N_c \approx 8.8 \times 10^{13}$  (non-embd params)
- For models with limited dataset size:  $L(D) = (D_c/D)^{\alpha_D}$   $\alpha_D \approx 0.095$ ,  $D_c \approx 5.4 \times 10^{13}$  (tokens)
- For models trained with limited compute:  $L(C) = (C_c^{min}/C_{min})^{\alpha_c^{min}} \quad \alpha_c^{min} \approx 0.050, \quad C_c^{min} \approx 3.1 \times 10^8 \text{ (PF-days)}$

• Let  $L(\cdot)$  represent the test loss dependent on either parameters N, or dataset size D or

## Scaling laws for LLMs

### Test loss *L* as function of model size N and dataset size D

$$L(N,D) = \left[ \left(\frac{N_c}{N}\right)^{\frac{\alpha_N}{\alpha_D}} + \frac{D_c}{D} \right]^{\alpha_D}$$



N is number of model parameters (not including vocabulary and positional embeddings) D is the number of tokens

Test loss *L* after transient period as function of model size N and number of update steps S

$$L(N,S) = \left(\frac{N_c}{N}\right)^{\alpha_N} + \left(\frac{S_c}{S_{\min}(S)}\right)^{\alpha_S}$$

# Scaling laws for LLMs

### • Keeping model size N fixed, architecture shape doesn't matter that much



Scaling Laws for Neural Language Models, Kaplan et al, OpenAl, 2020



# Comparing LSTM vs Transformers

### • LSTM cannot take advantage of long context (>100 tokens)



# Given a compute budget, what size model and amount of data should we train on?

### Large models are more sample-efficient than small models

Larger models require **fewer samples** to reach the same performance



**Figure 2** We show a series of language model training runs, with models ranging in size from  $10^3$  to  $10^9$ parameters (excluding embeddings).

The optimal model size grows smoothly with the loss target and compute budget





## How to allocate increasing compute?

For compute-efficient training

- Increase **model size** more than data (increase data sublinearly).
- Increase **batch size** as data size increases



Data requirements grow relatively slowly

Optimal model size increases very quickly

# **Optimal Allocation of Compute Budget**

Training at fixed batch size (should increase batch size with more data)



Models larger than the optimal-size can train faster (with less steps)

# **Critical batch size**

 $E_{\min}$ For compute efficient training, train with  $B_{crit} =$ 

- Larger B: more stable gradient, less training steps
- Critical batch size: above which scaling efficiency decreases significantly















# **Critical batch size as function of test loss**



The critical batch size  $B_{\rm crit}$  follows a power law in the loss as performance increase, and does Figure 10 not depend directly on the model size. We find that the critical batch size approximately doubles for every 13% decrease in loss.  $B_{\rm crit}$  is measured empirically from the data shown in Figure 18, but it is also roughly predicted by the gradient noise scale, as in [MKAT18]. arXiv:1812.06162



# Lessons from scaling LLMs

- Number of model parameters
- Size of dataset D
- Amount of compute (MFLOPs) C
- Performance depends strongly on scale, weakly on model shape
- N, D, C when not bottlenecked by the other two
- independent of the model size

Performance has a power-law relationship with each of the three scale factors

Performance improves predictably as long as we scale up N and D in tandem

Training curves follow predictable power-laws whose parameters are roughly

# Lessons from scaling LLMs

- Transfer to a different distribution incurs a constant penalty but otherwise improves roughly in line with performance on the training set.
- Large models are more sample-efficient than small models, reaching the same level of performance with fewer optimization steps and using fewer data points
- The ideal batch size for training these models is roughly a power of the loss only, and continues to be determinable by measuring the gradient noise scale
- If no constraints on data and model size, with given compute budget C

$$N \propto C^{lpha_C^{\min} / lpha_N} \qquad B \propto C^{lpha_C^{\min} / lpha_B}$$

 $S \propto C^{\alpha_C^{\min}/\alpha_S} \qquad D = B \cdot S$ 

### Is larger models always better? Can we train high-performance smaller models with more data?

# Is bigger always better?





# **Training Compute-Optimal Large Language Models**

Jordan Hoffmann<sup>\*</sup>, Sebastian Borgeaud<sup>\*</sup>, Arthur Mensch<sup>\*</sup>, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals and Laurent Sifre\*

https://arxiv.org/abs/2203.15556



## **Train longer on more tokens** Lessons from training Chinchilla

- From GPT3: large models should not be trained to lowest possible loss to be compute optimal
- Question: Given a fixed FLOPs budget how should one trade off model size and number of training tokens?
- Pre-training loss L(N, D) for N parameters and D training tokens. Find the optimal N and D values for a given compute budget.
- Empirical study on training 400 models from 70M to 16B parameters, trained on 5B to 400B tokens.
- Answer: Train smaller models for (a lot) more training steps.

• Better to scale model size and number of tokens linearly!



- Approach 1
- Approach 2
- Approach 3
- Kaplan et al (2020)



Chinchilla (70B) Gopher (280B) GPT-3 (175B) Megatron-Turing NLG (530B)

### • For different model sizes, choose number of training tokens to keep FLOPs constant



for an *optimal* model trained with the compute budget of *Gopher*.

Figure 3 | IsoFLOP curves. For various model sizes, we choose the number of training tokens such that the final FLOPs is a constant. The cosine cycle length is set to match the target FLOP count. We find a clear valley in loss, meaning that for a given FLOP budget there is an optimal model to train (left). Using the location of these valleys, we project optimal model size and number of tokens for larger models (center and right). In green, we show the estimated number of parameters and tokens

### Model

LaMDA (Thoppilan et al., 2022) GPT-3 (Brown et al., 2020) Jurassic (Lieber et al., 2021) *Gopher* (Rae et al., 2021) MT-NLG 530B (Smith et al., 2022)

Chinchilla

Size (# Parameters)	Training Toker
137 Billion	168 Billion
175 Billion	300 Billion
178 Billion	300 Billion
280 Billion	<b>300</b> Billion
530 Billion	270 Billion
70 Billion	1.4 Trillion



# Is larger models always better? Can we go to even smaller models?

# Small language models



- regime
- Mostly open-weight

Slide Credit: Mengzhou Xia https://princeton-cos597r.github.io/lectures/lec16.pdf

### • SLM: language models <10B • Fierce competition in small model

# Small language models are getting better and better





Slide Credit: Mengzhou Xia https://princeton-cos597r.github.io/lectures/lec16.pdf



# How to obtain such small models?

- Train on high quality data
  - content, we can learn the task better, even with a smaller model."
  - Phi-2: The surprising power of small language models
- Model pruning and distillation
  - Pruning: <u>Sheared-Llama</u>, <u>Llama 3.2</u>, <u>Minitron</u>
  - Distillation: <u>Gemma2</u>, <u>Llama 3.2</u>, <u>Minitron</u>

# "If we have a small dataset is focused on text-book quality educational



# **Types of pruning**

- Unstructured pruning
  - Magnitude pruning
  - Weights and activations
- Structured pruning
  - Prune entire blocks / components

	U	Т	$\nearrow$	Params	MNLI
BERT <sub>base</sub> (teacher)	×	×	$1.0 \times$	85M	84.8
Distillation					
DistillBERT <sub>6</sub>	1	×	2.0  imes	43M	82.2
TinyBERT <sub>6</sub>	1	1	2.0  imes	43M	84.0
MobileBERT <sup>‡</sup>	1	×	2.3  imes	20M	83.9
DynaBERT	×	1	6.3  imes	11M	76.3
AutoTinyBERT <sup>‡</sup>	1	1	9.1  imes	3.3M	78.2
$TinyBERT_4$	1	1	$11.4 \times$	4.7M	78.8
Pruning					
Movement Pruning	×	1	$1.0 \times$	9M	81.2
Block Pruning	×	1	$2.7 \times$	25M	83.7
CoFi Pruning (ours)	X	1	$2.7 \times$	26M	84.9
CoFi Pruning (ours)	X	1	$12.1 \times$	4.4M	80.6

Table 1: A comparison of state-of-the-art distillation and pruning methods. U and T denote whether Unlabeled and Task-specific are used for distillation or pruning. The inference speedups  $(\nearrow)$  are reported against a BERT<sub>base</sub> model and we evaluate all the models on an NVIDIA V100 GPU (§4.1). The models labeled as <sup>‡</sup> use a different teacher model and are not a direct comparison. Models are one order of magnitude faster.<sup>3</sup>

Structured Pruning Learns Compact and Accurate Models [Xia et al. ACL 2022]



# **Different ways to distill**

- Response-based
  - Try to match output (either predictions) or logits)
- Feature-based  $L_{FeaD}(f_t(x), f_s(x)) = \mathcal{L}_F(\Phi_t(f_t(x)), \Phi_s(f_s(x)))$ 
  - Try to match feature activation
  - Can be at different hidden layers
- **Relation-based** 
  - Relationship between different features  $L_{RelD}(f_t, f_s) = \mathcal{L}_{R^1}(\Psi_t(\hat{f}_t, \check{f}_t), \Psi_s(\hat{f}_s, \check{f}_s))$



### Phi: Transformation functions to map features to the same shape



Distillation losses can be L\_2, CE, MMD, KL 



H
ea
che
er
$\leq$
ğ
2



### **Distillation of DeepSeek-R1 Training recipe** Generated / Filtered using DeepSeek-R1\* Reasoning (600K) General (200K) Self-curated data Reasoning LLM (800k)(DeepSeek-R1\*) Small LLM SFT Qwen/Llama)

Small Reasoning  $\mathsf{I} \mathsf{I} \mathsf{M}$ Qwen/Llama) SFT



# Distillation

Model	AIME 2024		MATH-500	GPQA Diamond	LiveCode Bench	CodeForces
	pass@1	cons@64	pass@1	pass@1	pass@1	rating
GPT-4o-0513	9.3	13.4	74.6	49.9	32.9	759
Claude-3.5-Sonnet-1022	16.0	26.7	78.3	65.0	38.9	717
OpenAI-o1-mini	63.6	80.0	90.0	60.0	53.8	1820
QwQ-32B-Preview	50.0	60.0	90.6	54.5	41.9	1316
DeepSeek-R1-Distill-Qwen-1.5B	28.9	52.7	83.9	33.8	16.9	954
DeepSeek-R1-Distill-Qwen-7B	55.5	83.3	92.8	49.1	37.6	1189
DeepSeek-R1-Distill-Qwen-14B	69.7	80.0	93.9	59.1	53.1	1481
DeepSeek-R1-Distill-Qwen-32B	72.6	83.3	94.3	62.1	57.2	1691
DeepSeek-R1-Distill-Llama-8B	50.4	80.0	89.1	49.0	39.6	1205
DeepSeek-R1-Distill-Llama-70B	70.0	86.7	94.5	65.2	57.5	1633

Table 5 | Comparison of DeepSeek-R1 distilled models and other comparable models on reasoning-related benchmarks.

# **Combined pruning and distillation**



Slide Credit: Mengzhou Xia https://princeton-cos597r.github.io/lectures/lec16.pdf

### 1B & 3B Pruning & Distillation

https://ai.meta.com/blog/llama-3-2-connect-2024-vision-edge-mobile-devices/



# **Combined pruning and distillation**



Slide Credit: Mengzhou Xia https://princeton-cos597r.github.io/lectures/lec16.pdf

https://ai.meta.com/blog/llama-3-2-connect-2024-vision-edge-mobile-devices/



# Test time scaling

- Scaling laws for training
- What about scaling laws for inference?
  - Tradeoff model size vs inference time compute
    - Large models can only be deployed on machines with large memory
    - Smaller models can be deployed on smaller devices
- Generate multiple output chains and do search to select / aggregate to form • a final output

# **Tree of thought**

- Decision making by considering multiple paths of reasoning
- Consider different "thoughts" expressed in language
- Use to solve different types of problems



Tree of Thoughts: Deliberate Problem Solving with Large Language Models [Yao et al, 2023] Large Language Model Guided Tree-of-Thought [Long 2023]





Figure 2 | **Comparing different PRM search methods. Left:** Best-of-N samples N full answers and then selects the best answer according to the PRM final score. Center: Beam search samples N candidates at each step, and selects the top M according to the PRM to continue the search from. Right: lookahead-search extends each step in beam-search to utilize a k-step lookahead while assessing which steps to retain and continue the search from. Thus lookahead-search needs more compute.

Scaling LLM Test-Time Compute Optimally can be More Effective than Scaling Model Parameters [Snell et al. 2024] - https://arxiv.org/pdf/2408.03314 ue equibare eprinair) eau periore Energine ruan eequilê mederi arannerere femer er an





# Closer look at some recent model architectures

### **PaLM: Scaling Language Modeling with Pathways**

Aakanksha Chowdhery<sup>\*</sup> Sharan Narang<sup>\*</sup> Jacob Devlin<sup>\*</sup> Maarten Bosma Gaurav Mishra Adam Roberts Paul Barham Hyung Won Chung Charles Sutton Sebastian Gehrmann Parker Schuh Kensen Shi Sasha Tsvyashchenko Joshua Maynez Abhishek Rao<sup>†</sup> Parker Barnes Yi Tay Noam Shazeer<sup>‡</sup> Vinodkumar Prabhakaran Emily Reif Nan Du Ben Hutchinson Reiner Pope James Bradbury Jacob Austin Michael Isard Guy Gur-Ari Pengcheng Yin Toju Duke Anselm Levskaya Sanjay Ghemawat Sunipa Dev Henryk Michalewski Xavier Garcia Vedant Misra Kevin Robinson Liam Fedus Denny Zhou Daphne Ippolito David Luan<sup>‡</sup> Hyeontaek Lim Barret Zoph Alexander Spiridonov Ryan Sepassi David Dohan Shivani Agrawal Mark Omernick Andrew M. Dai Thanumalayan Sankaranarayana Pillai Marie Pellat Aitor Lewkowycz Erica Moreira Rewon Child Oleksandr Polozov<sup>†</sup> Katherine Lee Zongwei Zhou Xuezhi Wang Brennan Saeta Mark Diaz Orhan Firat Michele Catasta<sup>†</sup> Jason Wei Kathy Meier-Hellstern Douglas Eck Jeff Dean Slav Petrov Noah Fiedel

https://arxiv.org/abs/2204.02311

Google Research

# PaLM

Architecture

- SwiGLU activation:  $Swish(xW) \otimes xV$
- Parallel layers
  - Serial: y = x + MLP(LayerNorm(x + Attention(LayerNorm(x))))
  - Parallel: y = x + MLP(LayerNorm(x)) + Attention(LayerNorm(x))
- 15% faster training speed (degradation for small models 8B, but no degradation at 62B) Attention: Shared key-value across heads, query is still separately projected per head
- RoPE (rotary position) embeddings
- Shared input-output embeddings
- No biases: increased training stability
- Vocabulary: SentencePiece with 256k tokens

Training data

• 780 billion tokens of natural language + source code from github

PaLM: Scaling Language Modeling with Pathways, Chowdhery et al, Google, 2022

# PaLM: model architecture

- written as:

Whereas the parallel formulation can be written as:

y = x + MLP(LayerNorm(x)) + Attention(LayerNorm(x))

The parallel formulation results in roughly 15% faster training speed at large scales, since the MLP and Attention input matrix multiplications can be fused. Ablation experiments showed a small quality degradation at 8B scale but no quality degradation at 62B scale, so we extrapolated that the effect of parallel layers should be quality neutral at the 540B scale.

• SwiGLU Activation – We use SwiGLU activations  $(Swish(xW) \cdot xV)$  for the MLP intermediate activations because they have been shown to significantly increase quality compared to standard ReLU, GeLU, or Swish activations (Shazeer, 2020). Note that this does require three matrix multiplications in the MLP rather than two, but Shazeer (2020) demonstrated an improvement in quality in computeequivalent experiments (i.e., where the standard ReLU variant had proportionally larger dimensions).

• Parallel Layers – We use a "parallel" formulation in each Transformer block (Wang & Komatsuzaki, 2021), rather than the standard "serialized" formulation. Specifically, the standard formulation can be

```
y = x + MLP(LayerNorm(x + Attention(LayerNorm(x))))
```



# PaLM: model architecture

- not shared between examples, and only a single token is decoded at a time.
- sequence lengths.
- is done frequently (but not universally) in past work.

• Multi-Query Attention – The standard Transformer formulation uses k attention heads, where the input vector for each timestep is linearly projected into "query", "key", and "value" tensors of shape [k, h], where h is the attention head size. Here, the key/value projections are shared for each head, i.e. "key" and "value" are projected to [1, h], but "query" is still projected to shape [k, h]. We have found that this has a neutral effect on model quality and training speed (Shazeer, 2019), but results in a significant cost savings at autoregressive decoding time. This is because standard multi-headed attention has low efficiency on accelerator hardware during auto-regressive decoding, because the key/value tensors are

• **RoPE Embeddings** – We use RoPE embeddings (Su et al., 2021) rather than absolute or relative position embeddings, since RoPE embeddings have been shown to have better performance on long

• Shared Input-Output Embeddings – We share the input and output embedding matrices, which



# PaLM: model architecture

in increased training stability for large models.

digit tokens (e.g., " $123.5 \rightarrow 1\ 2\ 3\ .\ 5$ ").

• No Biases – No biases were used in any of the dense kernels or layer norms. We found this to result

• Vocabulary – We use a SentencePiece (Kudo & Richardson, 2018a) vocabulary with 256k tokens, which was chosen to support the large number of languages in the training corpus without excess tokenization. The vocabulary was generated from the training data, which we found improves training efficiency. The vocabulary is completely lossless and reversible, which means that whitespace is completely preserved in the vocabulary (especially important for code) and out-of-vocabulary Unicode characters are split into UTF-8 bytes, with a vocabulary token for each byte. Numbers are always split into individual

# PaLM: model hyperparameters

Model	Layers	# of Heads	$d_{ m model}$	# of Parameters (in billions)	Batch Size
PaLM 8B	32	16	4096	8.63	$256 \rightarrow 512$
PaLM 62B	64	32	8192	62.50	$512 \rightarrow 1024$
PaLM 540B	118	48	18432	540.35	$512 \rightarrow 1024 \rightarrow 2048$

Table 1: Model architecture details. We list the number of layers,  $d_{\text{model}}$ , the number of attention heads and attention head size. The feed-forward size  $d_{\rm ff}$  is always  $4 \times d_{\rm model}$  and attention head size is always 256.



# PaLM: training data

Total dataset s

Data source

Social media conversations Filtered webpages (multilin Books (English) GitHub (code) Wikipedia (multilingual) News (English)

Table 2: Proportion of data from each source in the training dataset. The multilingual corpus contains text from over 100 languages, with the distribution given in Appendix Table 29.

ize = 780 billion tokens					
	Proportion of data				
(multilingual)	50% 27% 13% 5% 4% 1%				



# PaLM: Pathways data parallelism

- Trained on two TPU v4 pods
  - Each pod had 3072 TPU chips attached to 768 hosts (total 6144 chips)
  - Each pod had full copy of model parameters
- Model + data parallelism, no pipeline parallelism
  - 12-way model parallelism, 256-way data sharing







Figure 2: The Pathways system (Barham et al., 2022) scales training across two TPU v4 pods using two-way data parallelism at the pod level.



### Chinchilla: 70B GPT-3: 175B Gopher: 280B PaLM: 540B







PaLM: Scaling Language Modeling with Pathways, Chowdhery et al, Google, 2022

## PaLM









### PaLM

PaLM: Scaling Language Modeling with Pathways, Chowdhery et al, Google, 2022





PaLM

PaLM: Scaling Language Modeling with Pathways, Chowdhery et al, Google, 2022

![](_page_55_Picture_5.jpeg)

# Toward multimodal agents

#### Mobile Manipulation

![](_page_56_Picture_2.jpeg)

![](_page_56_Picture_3.jpeg)

Human: Bring me the rice chips from the drawer. Robot: 1. Go to the drawers, 2. Open top drawer. I see <img>. 3. Pick the green rice chip bag from the drawer and place it on the counter.

![](_page_56_Picture_6.jpeg)

#### Visual Q&A, Captioning ...

![](_page_56_Picture_8.jpeg)

Given **<img>**. Q: What's in the image? Answer in emojis. A: 🍏 🍌 🍻 為 🍑 🐃 🚣.

![](_page_56_Picture_10.jpeg)

hurdle at a dog show.

Haiku about embodied LLMs. A: Embodied language. Models learn to understand. The world around them.

#### <u>An Embodied Multimodal Language Model</u> [Dreiss et al, Google, 2023] https://palm-e.github.io/ 57

![](_page_56_Picture_16.jpeg)

### LLaMA: Open and Efficient Foundation Language Models

- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet Marie-Anne Lachaux, Timothee Lacroix, Baptiste Rozière, Naman Goyal Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin Edouard Grave, Guillaume Lample\*
  - Meta AI

https://arxiv.org/abs/2302.13971

![](_page_57_Picture_7.jpeg)

# LLaMa

![](_page_58_Figure_1.jpeg)

### • 65B model trained on 1.4T tokens for ~21 days on 2048 A100 GPU with 80GB RAM.

params	dimension	n heads	n layers	learning rate	batch size	n tokens
6.7B	4096	32	32	$3.0e^{-4}$	<b>4M</b>	1.0T
13.0B	5120	40	40	$3.0e^{-4}$	<b>4M</b>	1.0T
32.5B	6656	52	60	$1.5e^{-4}$	<b>4M</b>	1.4T
65.2B	8192	64	80	$1.5e^{-4}$	<b>4M</b>	1.4T

LLaMA: Open and Efficient Foundation Language Models [Touvron et al. FAIR, 2023]

### LLaMA

![](_page_59_Picture_5.jpeg)

### Architecture

**Pre-normalization [GPT3].** To improve the training stability, we normalize the input of each transformer sub-layer, instead of normalizing the output. We use the RMSNorm normalizing function, introduced by Zhang and Sennrich (2019).

SwiGLU activation function [PaLM]. We replace the ReLU non-linearity by the SwiGLU activation function, introduced by Shazeer (2020) to improve the performance. We use a dimension of  $\frac{2}{3}4d$  instead of 4d as in PaLM.

**Rotary Embeddings [GPTNeo].** We remove the absolute positional embeddings, and instead, add rotary positional embeddings (RoPE), introduced by **Su et al.** (2021), at each layer of the network.

Data

Com C4 Gith Wiki Bool

ArX

Stac

![](_page_60_Picture_8.jpeg)

## Training data

aset	Sampling prop.	Epochs	Disk siz
nmonCrawl	67.0%	1.10	3.3 TI
	15.0%	1.06	783 GI
ub	4.5%	0.64	328 GI
ipedia	4.5%	2.45	83 GI
ks	4.5%	2.23	85 GI
liv	2.5%	1.06	92 GI
kExchange	2.0%	1.03	78 GI

LLaMA: Open and Efficient Foundation Language Models [Touvron et al. FAIR, 2023]

![](_page_60_Picture_12.jpeg)

# Grouped multi-query attention

- Reduce number of heads used for keys and values
- Shared values and keys across heads

![](_page_61_Figure_3.jpeg)

- GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints [Ainslie et al. 2023]

![](_page_62_Figure_0.jpeg)

![](_page_62_Figure_1.jpeg)

![](_page_62_Picture_2.jpeg)

LLaMA: Open and Efficient Foundation Language Models [Touvron et al. FAIR, 2023] 63

![](_page_62_Picture_4.jpeg)

### LLaMA

		BoolQ	PIQA	SIQA	HellaSwag	WinoGrande	ARC-e	ARC-c	OBQA
GPT-3	175B	60.5	81.0	-	78.9	70.2	68.8	51.4	57.6
Gopher	280B	79.3	81.8	50.6	79.2	70.1	-	-	-
Chinchilla	70B	83.7	81.8	51.3	80.8	74.9	-	-	-
PaLM	62B	84.8	80.5	-	79.7	77.0	75.2	52.5	50.4
PaLM-cont	62B	83.9	81.4	-	80.6	77.0	-	-	-
PaLM	540B	88.0	82.3	-	83.4	81.1	76.6	53.0	53.4
	7B	76.5	79.8	48.9	76.1	70.1	72.8	47.6	57.2
ΤΤοΜΛ	13B	78.1	80.1	50.4	79.2	73.0	74.8	52.7	56.4
LLawiA	33B	83.1	82.3	50.4	82.8	76.0	80.0	<b>57.8</b>	58.6
	65B	85.3	82.8	52.3	84.2	77.0	78.9	56.0	60.2

LLaMA: Open and Efficient Foundation Language Models [Touvron et al. FAIR, 2023] 64

![](_page_63_Picture_3.jpeg)

### **GLaM: Efficient Scaling of Language Models with Mixture-of-Experts**

Nan Du<sup>\*1</sup> Yanping Huang<sup>\*1</sup> Andrew M. Dai<sup>\*1</sup> Simon Tong<sup>1</sup> Dmitry Lepikhin<sup>1</sup> Yuanzhong Xu<sup>1</sup> Maxim Krikun<sup>1</sup> Yanqi Zhou<sup>1</sup> Adams Wei Yu<sup>1</sup> Orhan Firat<sup>1</sup> Barret Zoph<sup>1</sup> Liam Fedus<sup>1</sup> Maarten Bosma<sup>1</sup> Zongwei Zhou<sup>1</sup> Tao Wang<sup>1</sup> Yu Emma Wang<sup>1</sup> Kellie Webster<sup>1</sup> Marie Pellat<sup>1</sup> Kevin Robinson<sup>1</sup> Kathleen Meier-Hellstern<sup>1</sup> Toju Duke<sup>1</sup> Lucas Dixon<sup>1</sup> Kun Zhang<sup>1</sup> Quoc V Le<sup>1</sup> Yonghui Wu<sup>1</sup> Zhifeng Chen<sup>1</sup> Claire Cui<sup>1</sup>

https://arxiv.org/abs/2112.06905

![](_page_64_Picture_3.jpeg)

# Mixture of Experts (MoE) for LLMs

![](_page_65_Figure_1.jpeg)

Interleaved transformer and MoE layers Sparse activation of experts

Weighted average of outputs from selected experts is passed to the transformer layer

Experts: each expert is a FFN

For each input token (e.g. 'roses'), the **Gating** module selects the two most relevant experts out of 64. Two different experts selected for each token.

Larger models with less activated parameters per input token More performant with similar amount of compute

Model Name	Model Type	$n_{ m params}$	$n_{ m act-params}$
BERT	Dense Encoder-only	340M	340M
T5	Dense Encoder-decoder	13B	13B
GPT-3	Dense Decoder-only	175B	175B
Jurassic-1	Dense Decoder-only	178B	178B
Gopher	Dense Decoder-only	280B	280B
Megatron-530B	Dense Decoder-only	530B	530B
GShard-M4	MoE Encoder-decoder	600B	1.5B
Switch-C	MoE Encoder-decoder	1.5T	1.5B
GLaM (64B/64E)	MoE Decoder-only	1.2T	96.6B

## **Mixture of Experts (MoE) for LLMs** Better effective FLOPs per token prediction in causal LMs

![](_page_67_Figure_1.jpeg)