

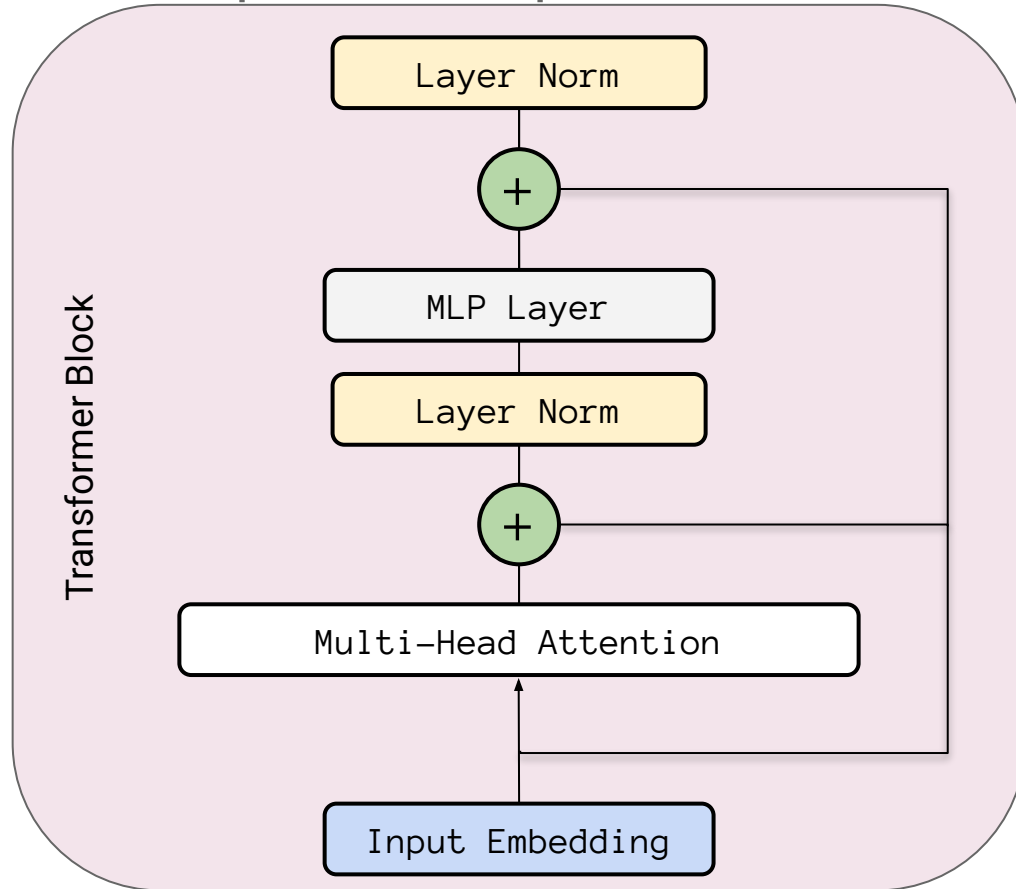
Training Large Language Models

CS 4644 / 7643: Deep Learning

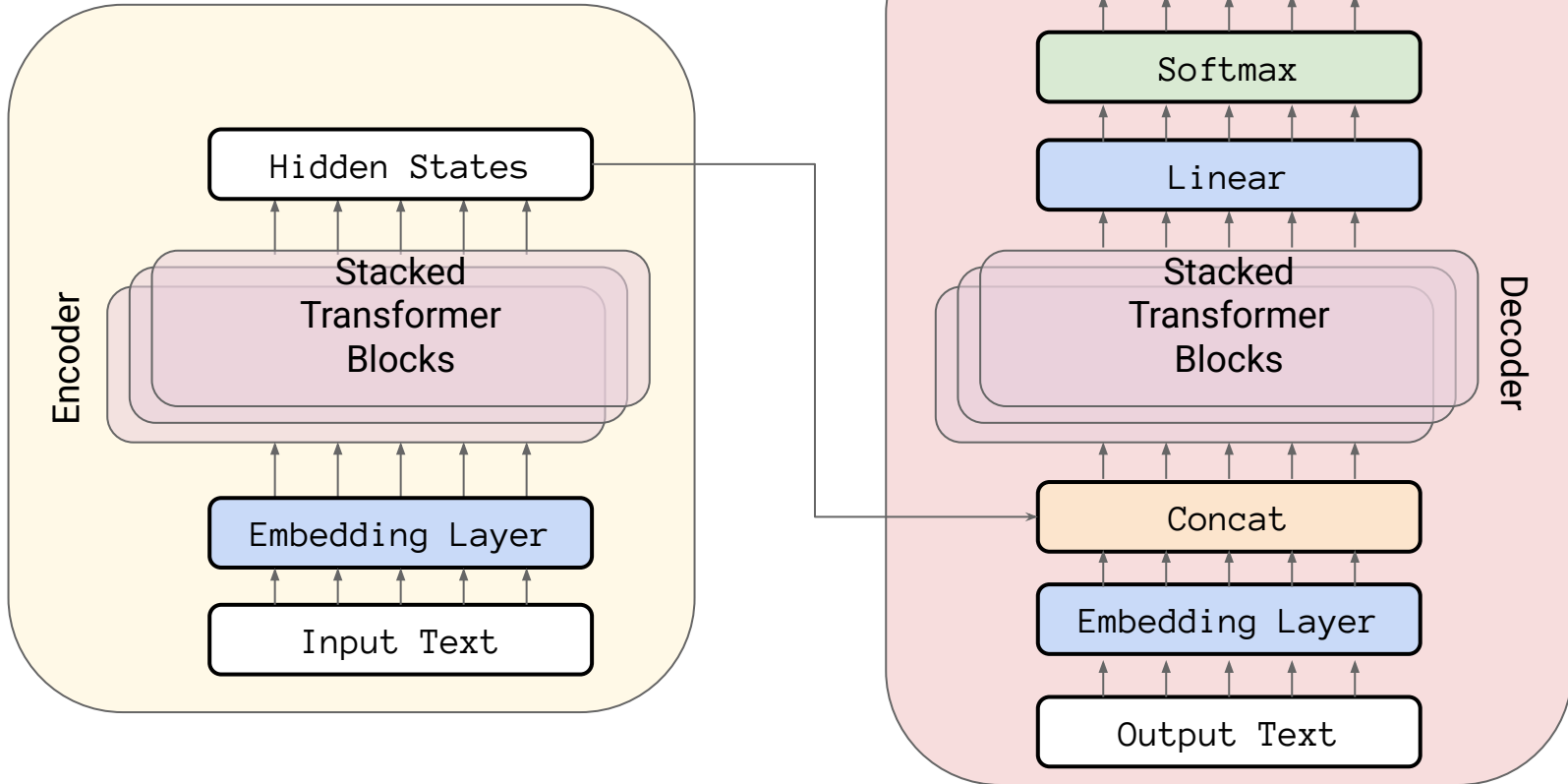
William Held

School of Interactive Computing
Georgia Institute of Technology

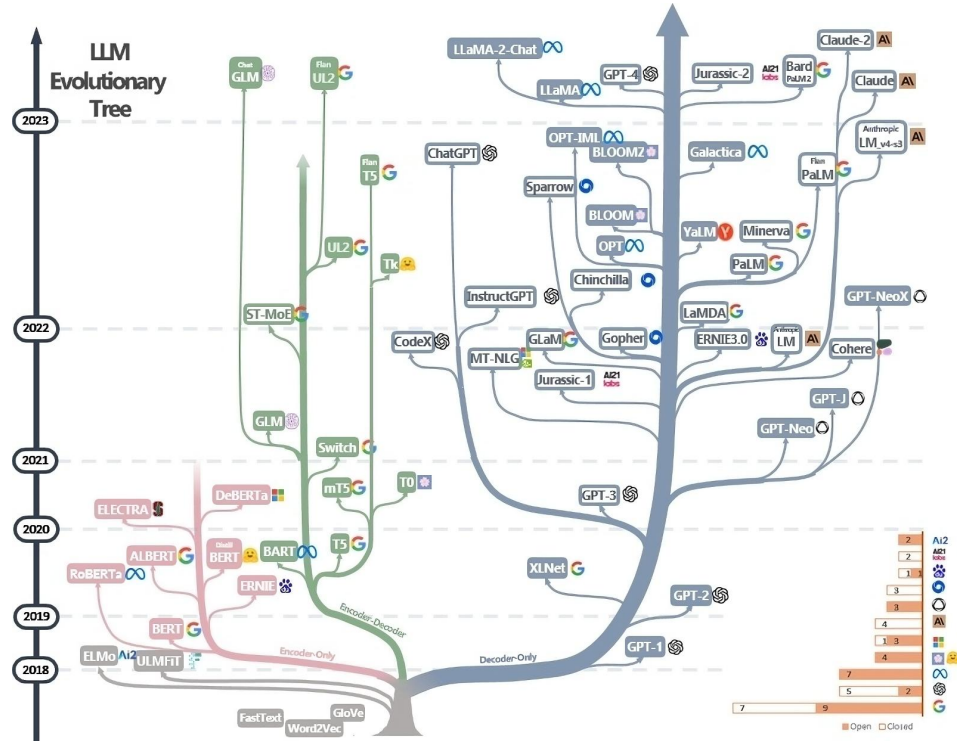
Transformer Lecture Speed Recap: The Transformer Block



Attention is “All” You Need



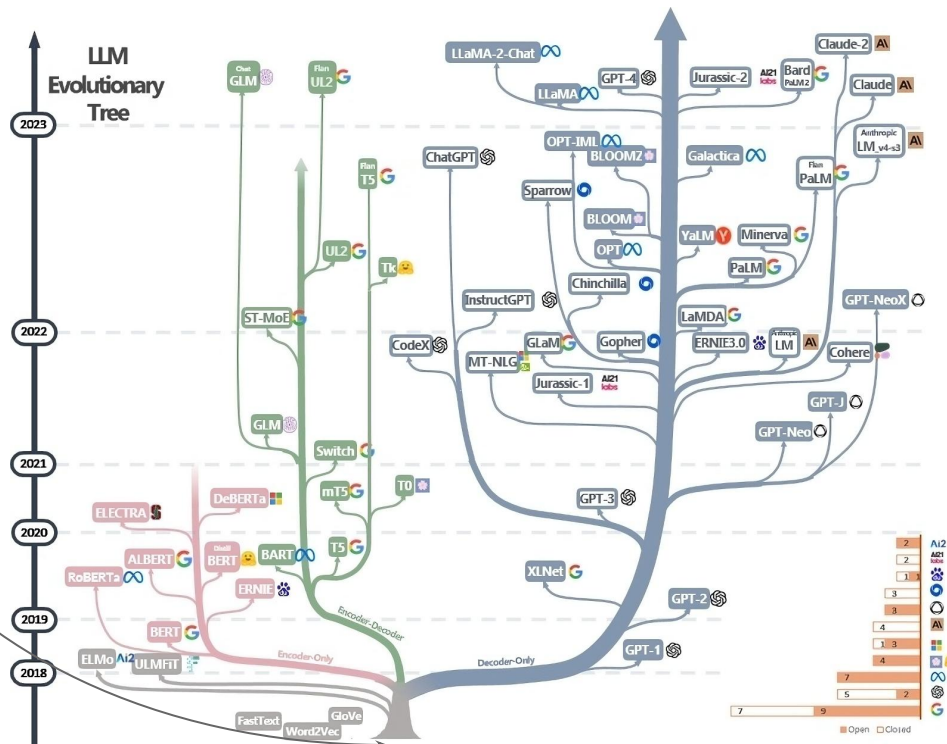
How do we go from purpose driven models to LLMs?



How do we go from purpose driven models to LLMs?

Self-Supervised Learning

How do we most effectively turn raw text into meaningful loss?



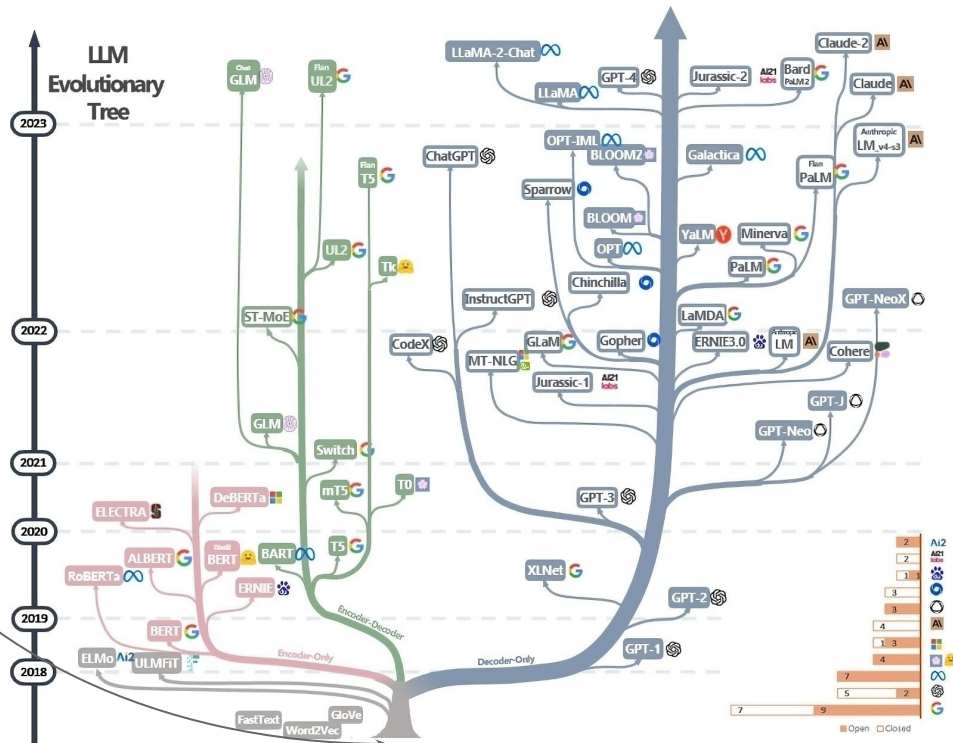
How do we go from purpose driven models to LLMs?

Self-Supervised Learning

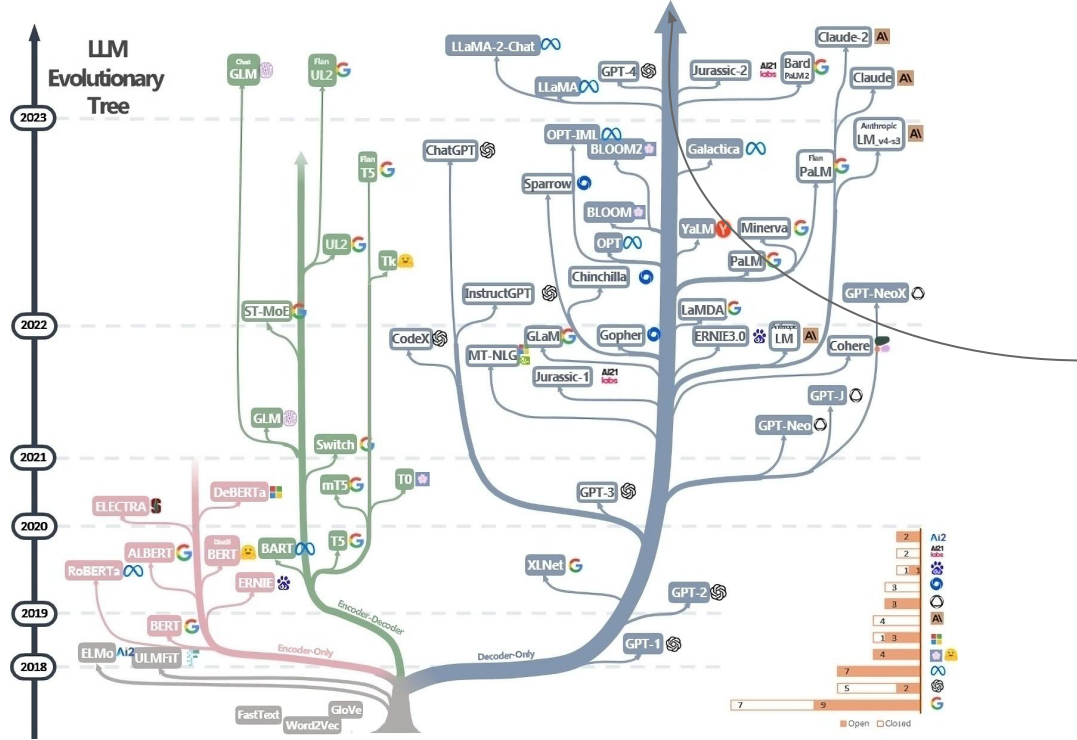
How do we most effectively turn raw text into meaningful loss?

Covered Today

- Encoder Only
- Decoder Only

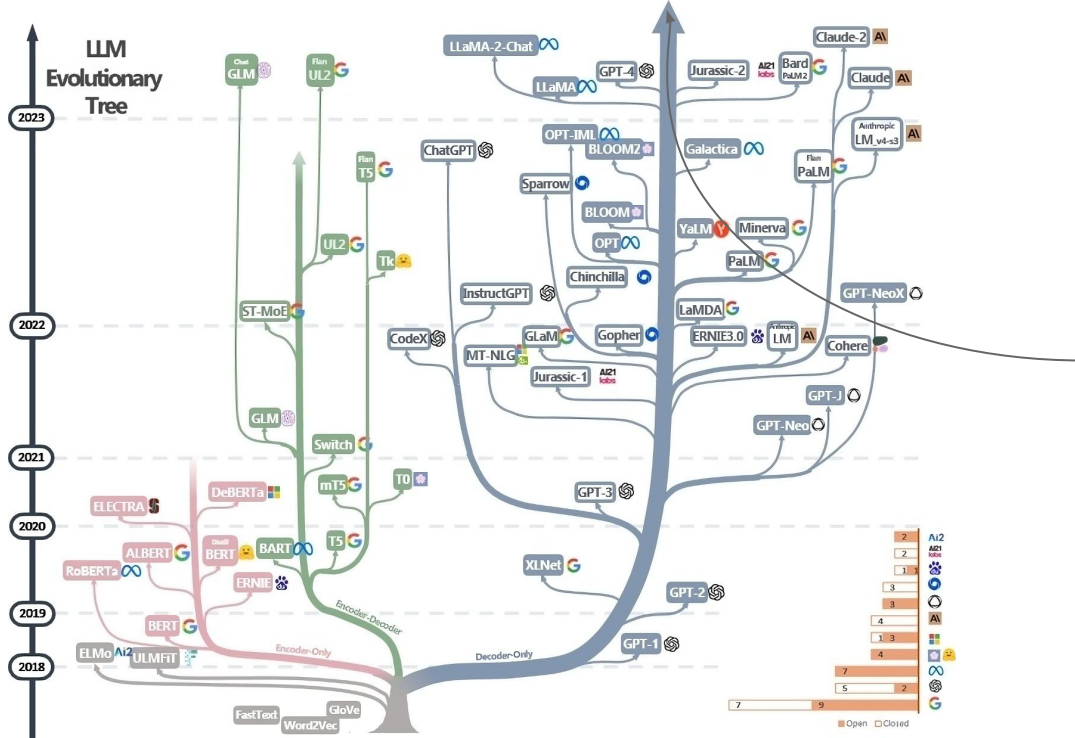


How do we go from purpose driven models to LLMs?



Scaling Laws
How do we train large models on large amounts of quality data?

How do we go from purpose driven models to LLMs?



Scaling Laws

How do we train large models on large amounts of quality data?

Covered Today

- Data Curation
- Distributed Training
- “Alignment”

LLM Advancements have been driven primarily by these two

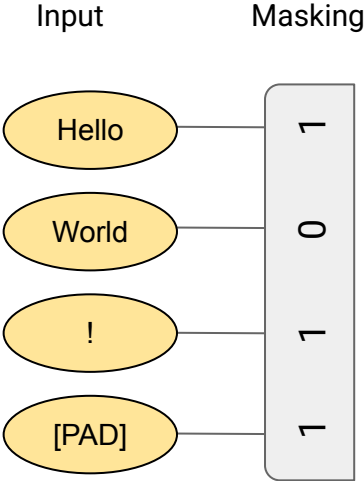
Self-Supervised Learning

How do we most effectively turn raw text into meaningful loss?

Scaling Laws

How do we train large models on large amounts of quality data?

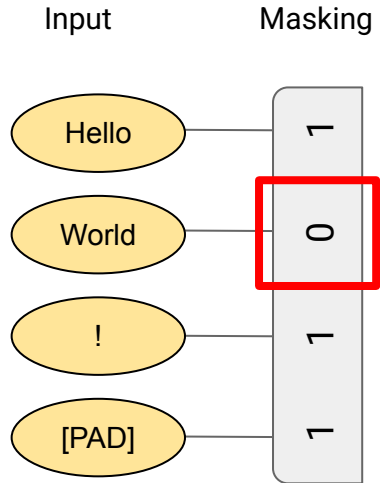
SSL | From raw text to loss!



Masked Language Model

[Devlin et al. 2018 \(BERT\)](#)

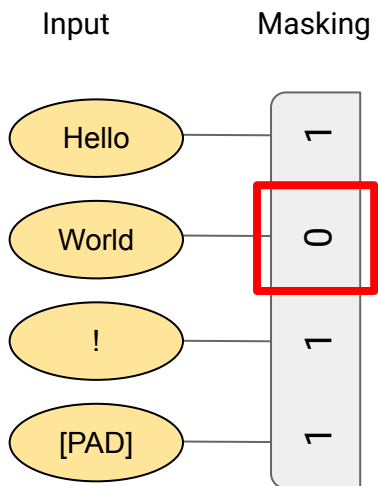
SSL | What is the “Mask” in a Masked Language Model?



Masked Language Model

[Devlin et al. 2018 \(BERT\)](#)

SSL | What is the “Mask” in a Masked Language Model?



Masked Attention

Similarities: $E = (QX^T / \sqrt{DQ}) * \text{MASK}$

Attention Matrix: $A = \text{softmax}(E, \text{dim}=1)$

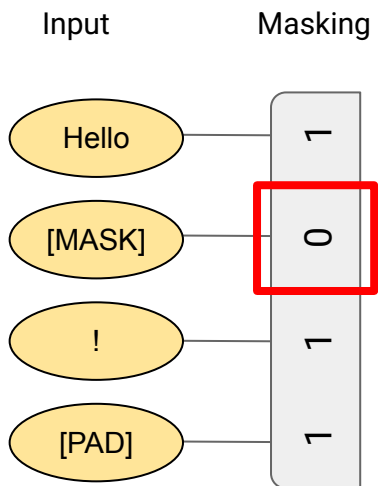
Output vectors: $Y = AX$

$$Y_i = \sum_j A_{i,j} X_j$$

Masked Language Model

[Devlin et al. 2018 \(BERT\)](#)

SSL | What is the “Mask” in a Masked Language Model?



Intuition

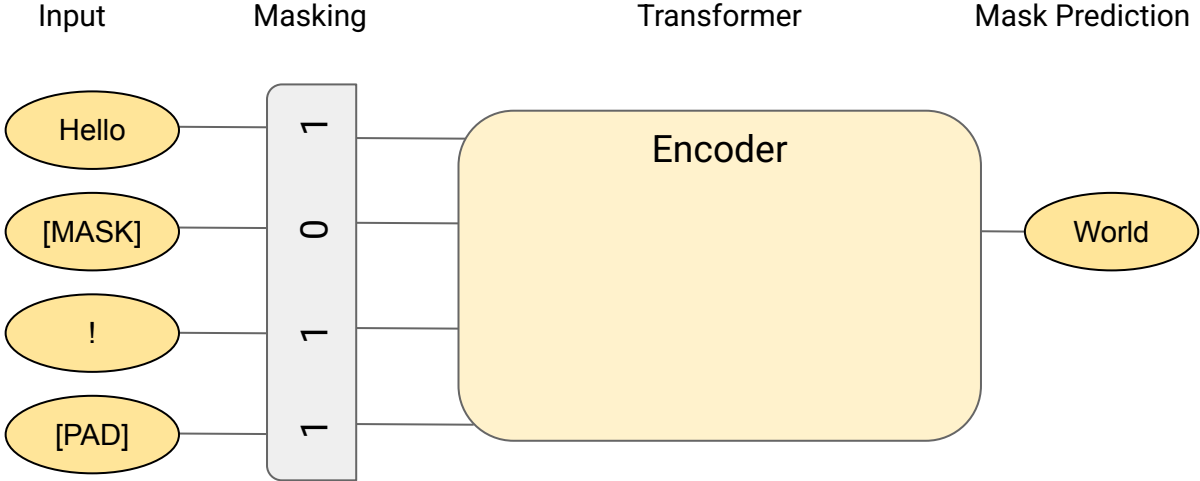
If $MASK_i = 0$, then $Y_i = \sum_{j, j \neq i} A_{i,j} X$

a.k.a the representation of the masked token is created purely from context

Masked Language Model

[Devlin et al. 2018 \(BERT\)](#)

SSL | Masked Token Prediction



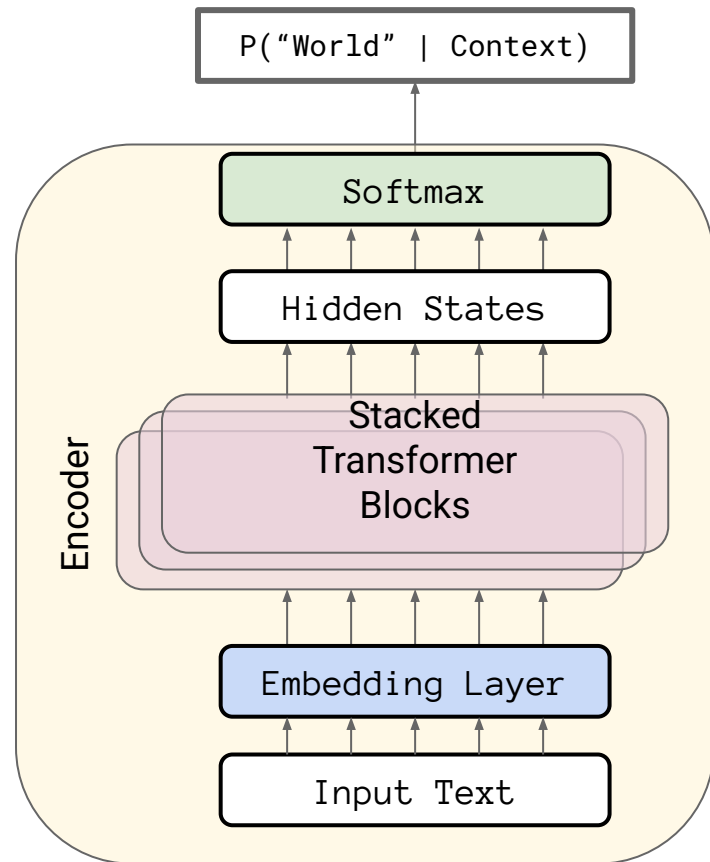
Masked Language Model

[Devlin et al. 2018 \(BERT\)](#)

SSL | Masked Token Prediction

Optimize Negative Log Likelihood

$$\text{loss} = -\log(P(\text{"World"} \mid \text{Context}))$$

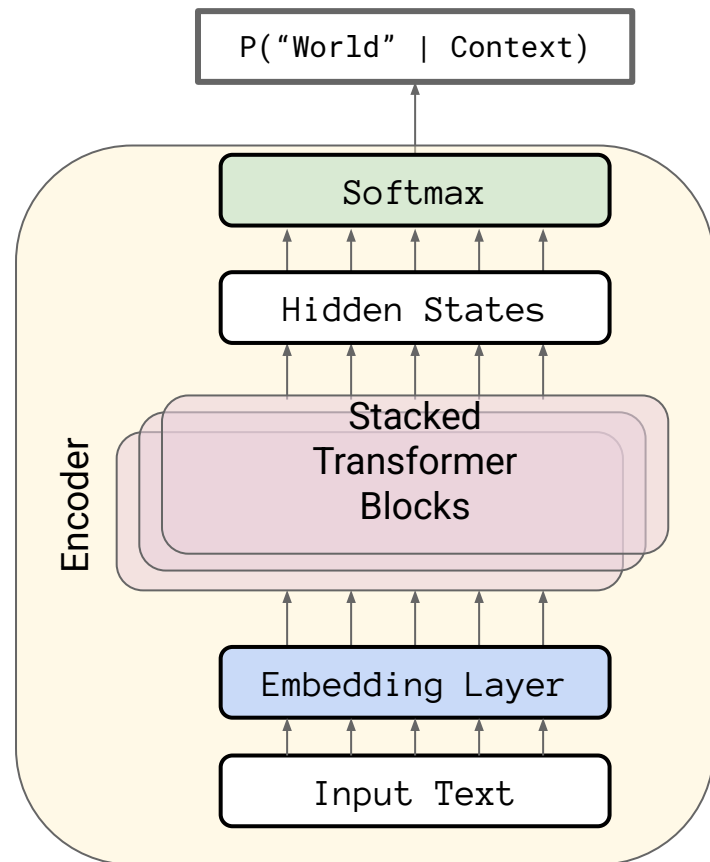


SSL | Masked Token Prediction

Optimize Negative Log Likelihood

$$\text{loss} = -\log(P(\text{"World"} \mid \text{Context}))$$

Equivalent to the Cross-Entropy
Loss term from Lecture 3!



Data | BERT used existing curation!

BERT Corpus

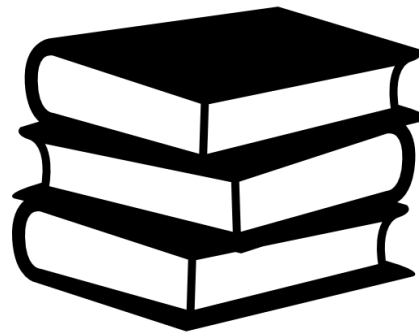
English Wikipedia + BooksCorpus

Size

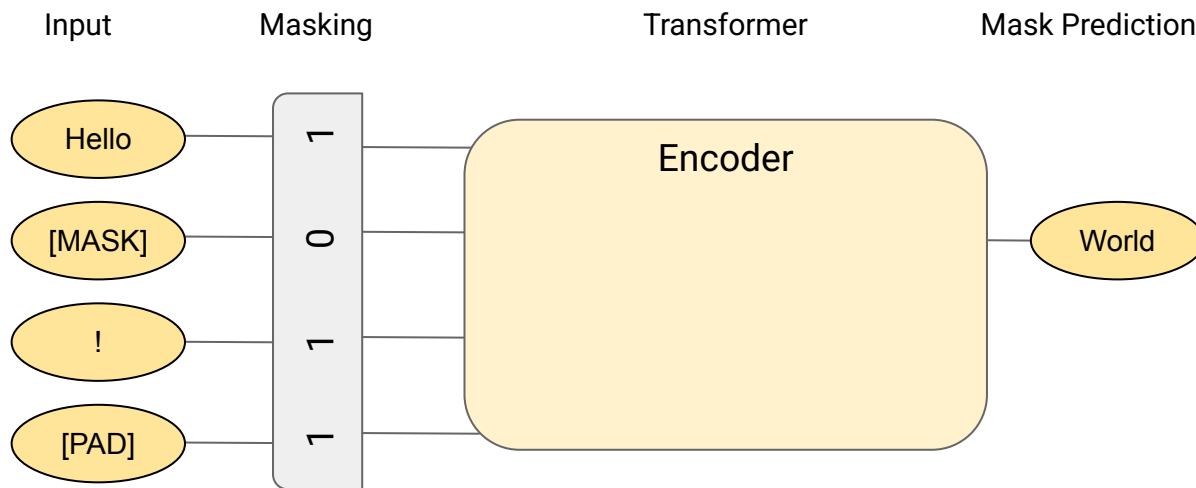
~3 Billion Tokens

Quality

High quality text,
Broad “Academic” Knowledge,
Limited Diversity

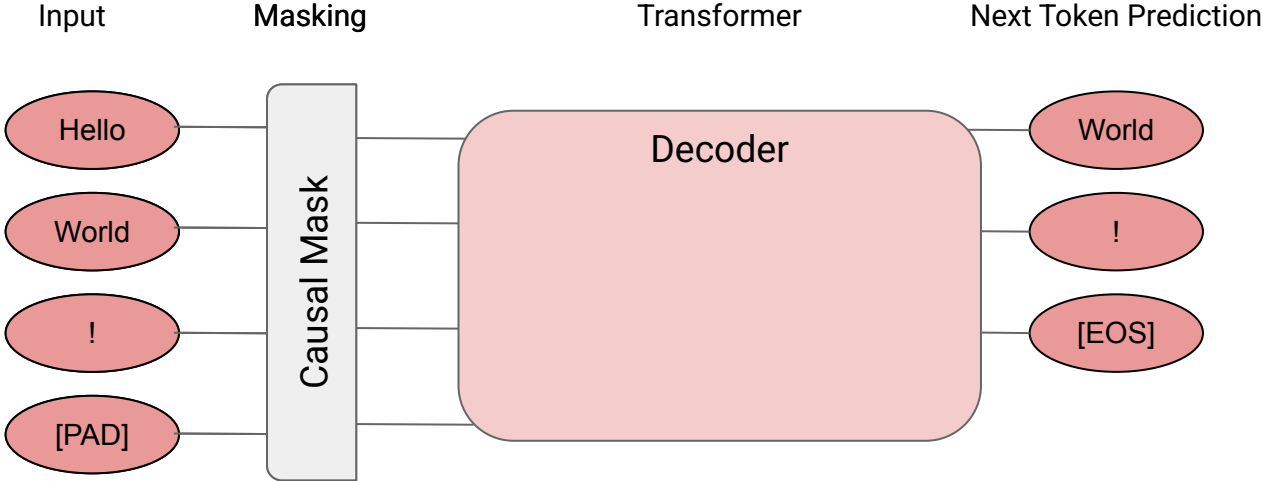


Questions?



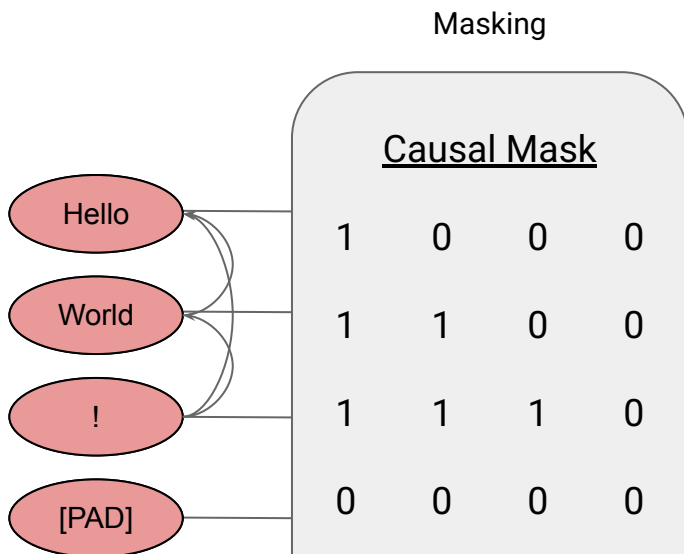
Masked Language Model

SSL | “How does GPT work?”



[Radford et al. 2019 \(GPT-2\)](#)

SSL | Autoregressive Language Modeling



Masked Attention Again!

Similarities: $E = (QXT / \sqrt{DQ}) * \text{MASK}$

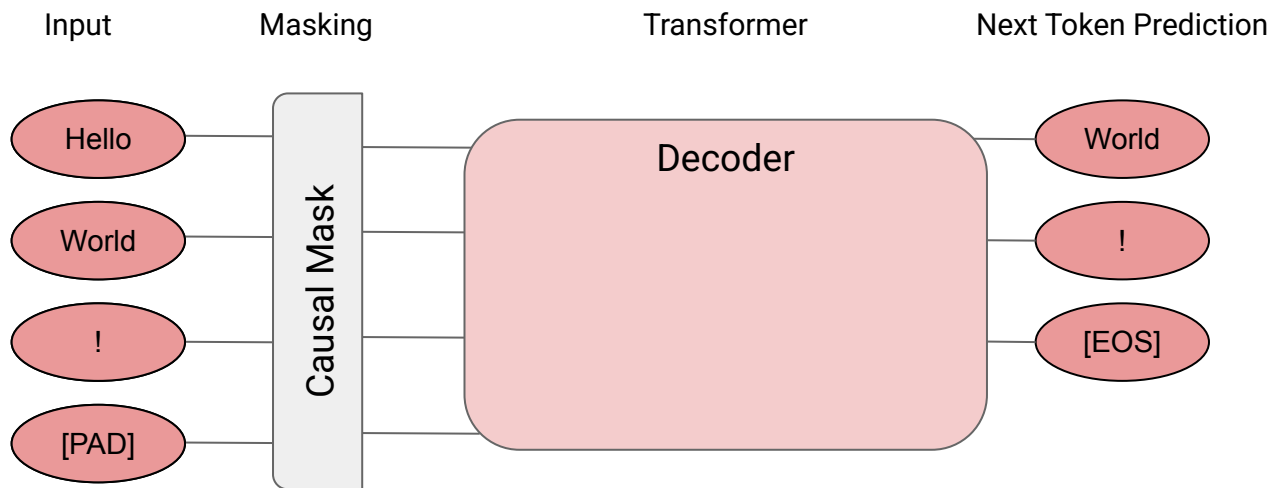
Attention Matrix: $A = \text{softmax}(E, \text{dim}=1)$

Output vectors: $Y = AX$

$$Y_i = \sum_j A_{i,j} X_j$$

Tokens only affected by preceding tokens

SSL | Purely Autoregressive



Optimize Negative Log Likelihood of Whole Sequence

$$\text{loss} = -(\log(P(\text{"World"} \mid \text{"Hello"})) + \log(P(\text{"!"} \mid \text{"Hello World"})) + \log(P(\text{"[EOS]"} \mid \text{"Hello World!"})))$$

[Radford et al. 2019 \(GPT-2\)](#)

Data | Increasing Token Count via Human Curation Heuristics

GPT-2 Corpus

All Reddit Outbound links with at least 3 karma

Size

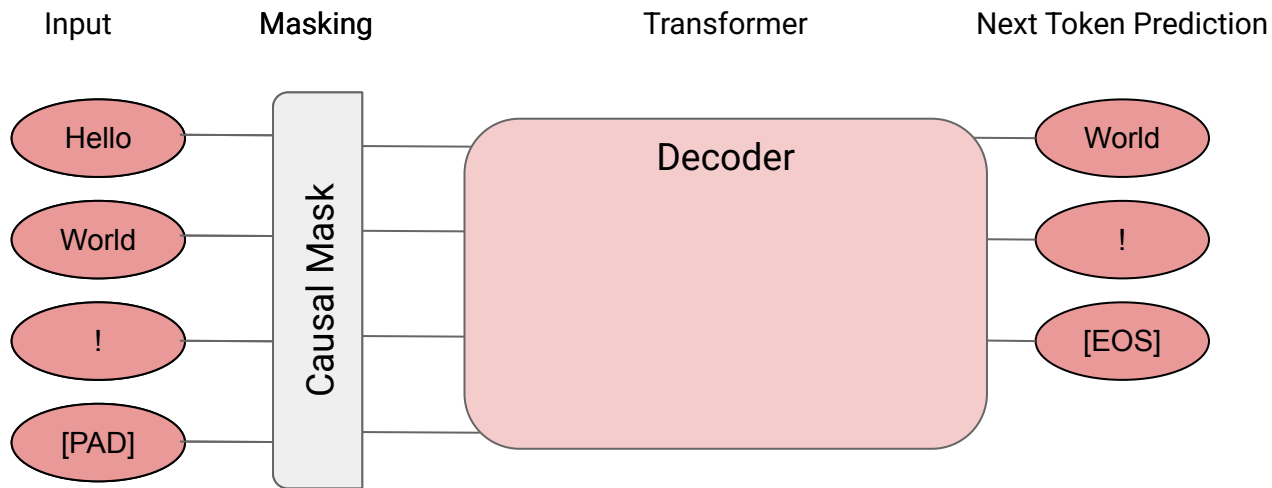
~10 Billion Tokens

Quality

High quality text,
Broad Knowledge,
Improved Diversity

URL Domain	# Docs	% of Total Docs
bbc.co.uk	116K	1.50%
theguardian.com	115K	1.50%
washingtonpost.com	89K	1.20%
nytimes.com	88K	1.10%
reuters.com	79K	1.10%
huffingtonpost.com	72K	0.96%
cnn.com	70K	0.93%
cbc.ca	67K	0.89%
dailymail.co.uk	58K	0.77%
go.com	48K	0.63%

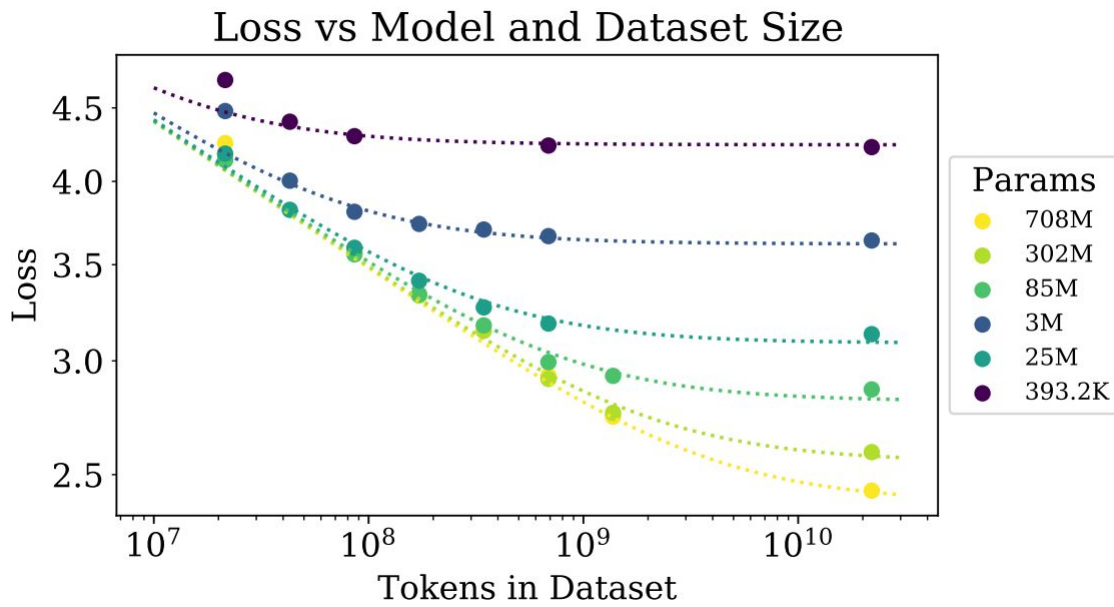
Questions?



Autoregressive Language Model

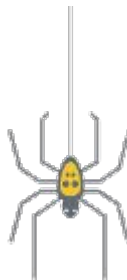
Data & Parameter Scaling | Moving to Large Language Models

Today's LLMs are driven by data and model scaling



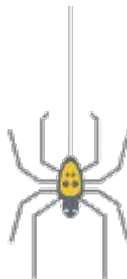
[Kaplan et al. 2020](#)

Data Scaling | Collecting High-Quality Self-Supervision at Scale



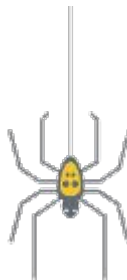
We could get a lot more data from CommonCrawl!

Data Scaling | Collecting High-Quality Self-Supervision at Scale



We could get a lot more data from CommonCrawl!
A lot of it is spam though...

Data Scaling | Collecting High-Quality Self-Supervision at Scale



We could get a lot more data from CommonCrawl!
A lot of it is spam though...
How do we get “useful” data?

T5 - Encoder-Decoder with Common Crawl Scale Data

T5 Corpus (AKA C4)

All Common Crawl Text Which
Meets Heuristics

Size

~350 Billion Tokens

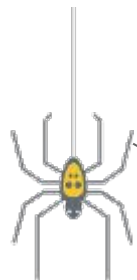
Quality

Varying quality text,
Broad Knowledge,
Improved Diversity

- We only retained lines that ended in a terminal punctuation mark (i.e. a period, exclamation mark, question mark, or end quotation mark).
- We discarded any page with fewer than 3 sentences and only retained lines that contained at least 5 words.
- We removed any page that contained any word on the “List of Dirty, Naughty, Obscene or Otherwise Bad Words”.⁶
- Many of the scraped pages contained warnings stating that Javascript should be enabled so we removed any line with the word Javascript.
- Some pages had placeholder “lorem ipsum” text; we removed any page where the phrase “lorem ipsum” appeared.
- Some pages inadvertently contained code. Since the curly bracket “{” appears in many programming languages (such as Javascript, widely used on the web) but not in natural text, we removed any pages that contained a curly bracket.
- Since some of the scraped pages were sourced from Wikipedia and had citation markers (e.g. [1], [citation needed], etc.), we removed any such markers.
- Many pages had boilerplate policy notices, so we removed any lines containing the strings “terms of use”, “privacy policy”, “cookie policy”, “uses cookies”, “use of cookies”, or “use cookies”.
- To deduplicate the data set, we discarded all but one of any three-sentence span occurring more than once in the data set.

[Raffel et al. 2019](#)

GPT-3 - Increased Scaling Via Automated Data Curation

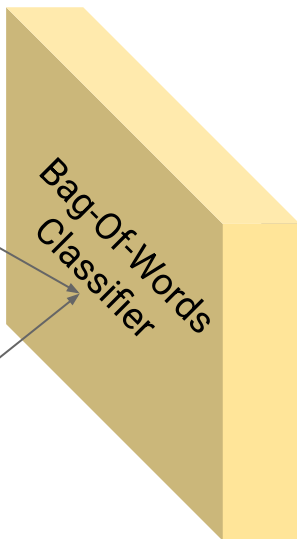


Low-Quality, High Volume

URL Domain	# Docs	% of Total Docs
bbc.co.uk	116K	1.50%
theguardian.com	115K	1.50%
washingtonpost.com	89K	1.20%
nytimes.com	88K	1.10%
reuters.com	79K	1.10%
huffingtonpost.com	72K	0.96%
cnn.com	70K	0.93%
cbc.ca	67K	0.89%
dailymail.co.uk	58K	0.77%
go.com	48K	0.63%

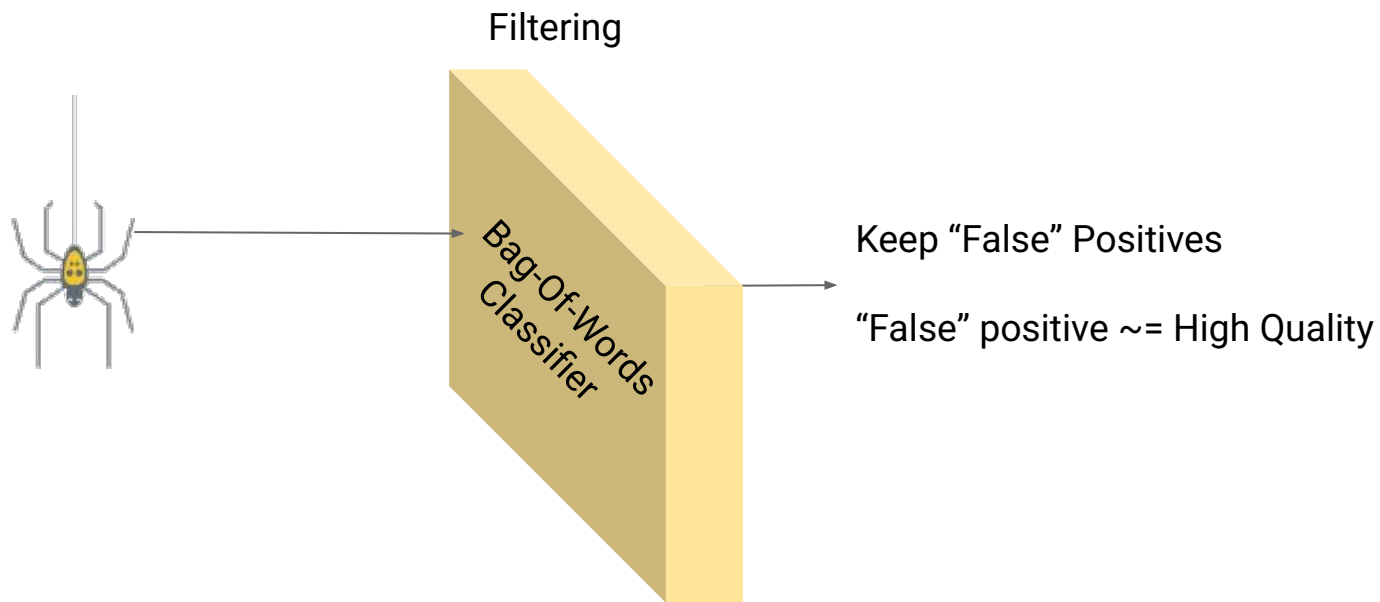
High Quality, Medium Volume

Training



Distinguish High and Low Quality

GPT-3 - Increased Scaling Via Automated Data Curation



Data | GPT-2 to Original GPT-3 was mostly data scaling

GPT-3 Corpus

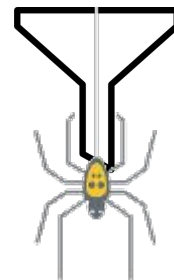
Common-Crawl Filtered using
GPT-2 Training Data

Size

~400 Billion Tokens

Quality

High-ish quality text,
Broad Knowledge,
Web-scale Diversity



Data | Recent Open Source models focus heavily on data scaling

Llama 1 Corpus

Size

~1.4 Trillion Tokens

Quality

Varying quality text,
Broad Knowledge,
Web-scale Diversity,
Includes Code!

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

Data | Data Mixture has become the biggest “secret”

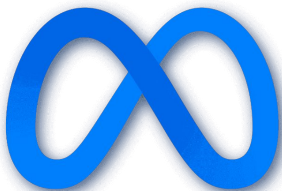
Llama 3 Corpus

Size

15 Trillion Tokens

Quality

Minimal details known



[Dubey et al. 2024](#)

Gemini Corpus

Size

Unknown

Quality

No details known



[Anil et al. 2023](#)

GPT-4 Corpus

Size

Unknown (Est. >11T Tokens)

Quality

No details known



[OpenAI 2023](#) 643 Deep Learning - William Held

Questions?

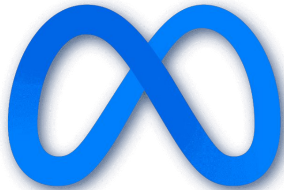
Llama 3 Corpus

Size

15 Trillion Tokens

Quality

Minimal details known



[Dubey et al. 2024](#)

Gemini Corpus

Size

Unknown

Quality

No details known



[Anil et al. 2023](#)

GPT-4 Corpus

Size

Unknown (Est. >11T Tokens)

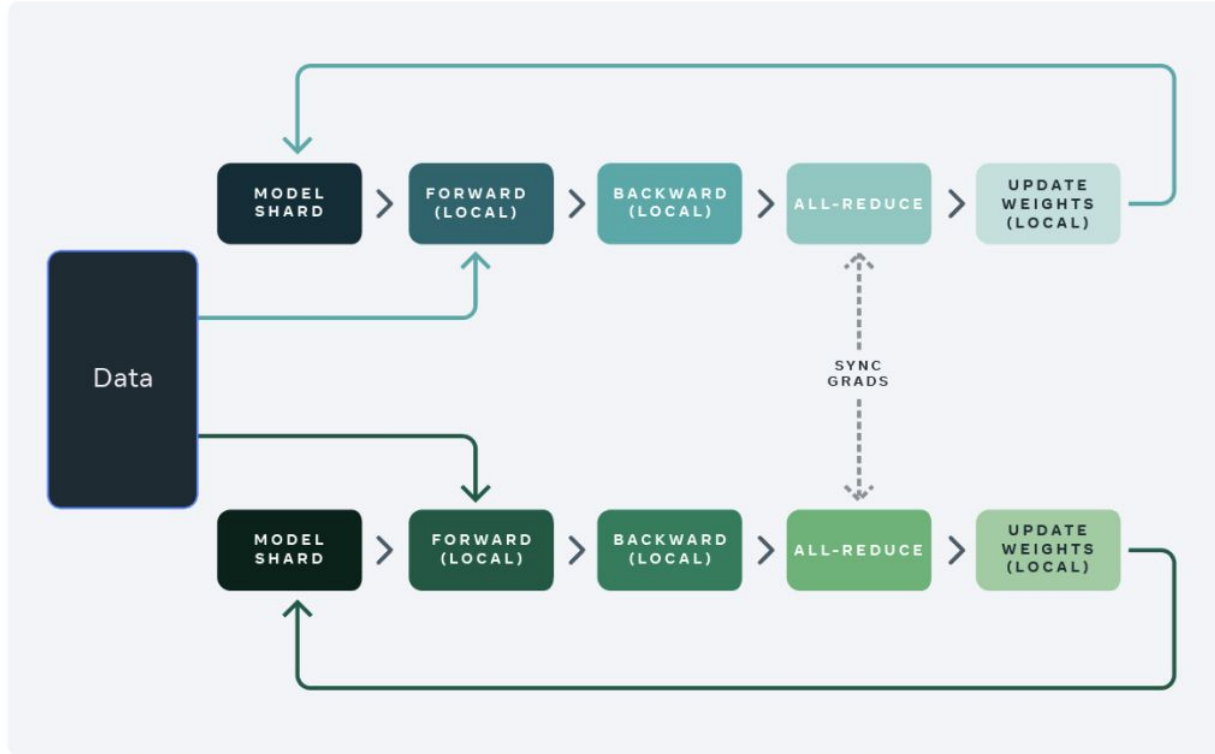
Quality

No details known

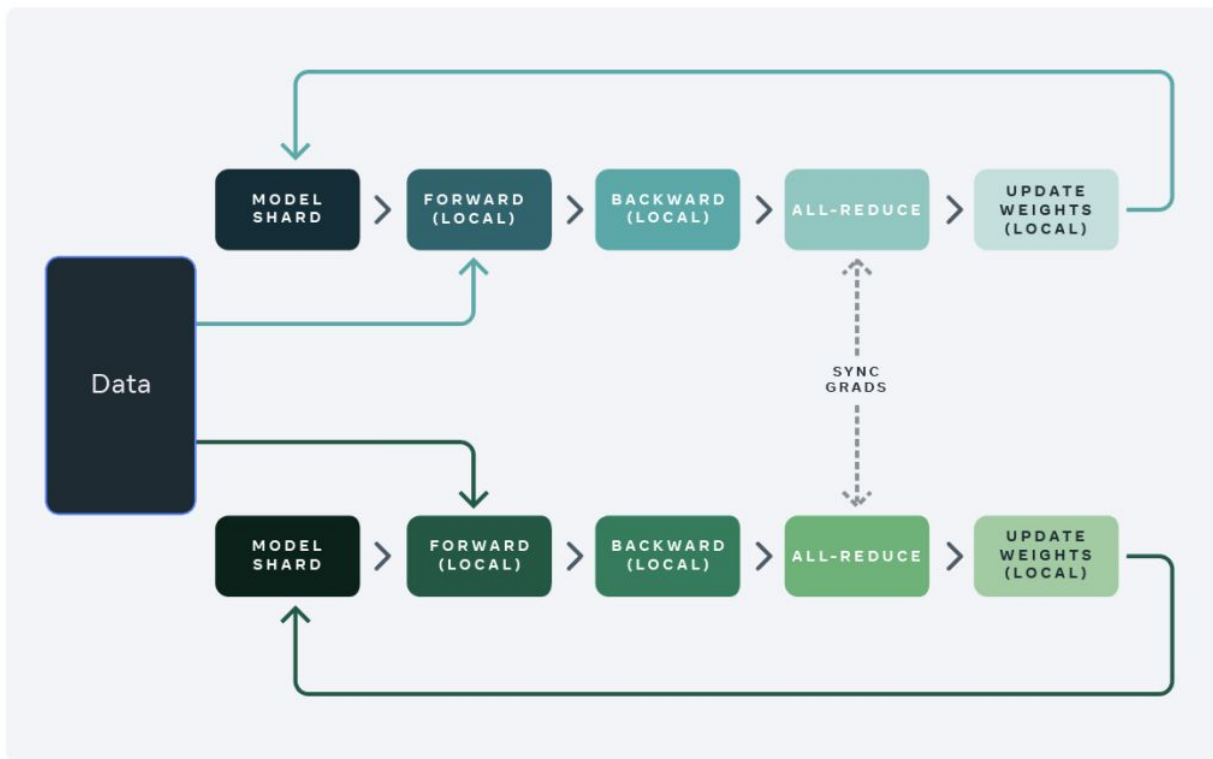


[OpenAI 2023](#) 643 Deep Learning - William Held

Scaling Parameters | Data Parallel Training

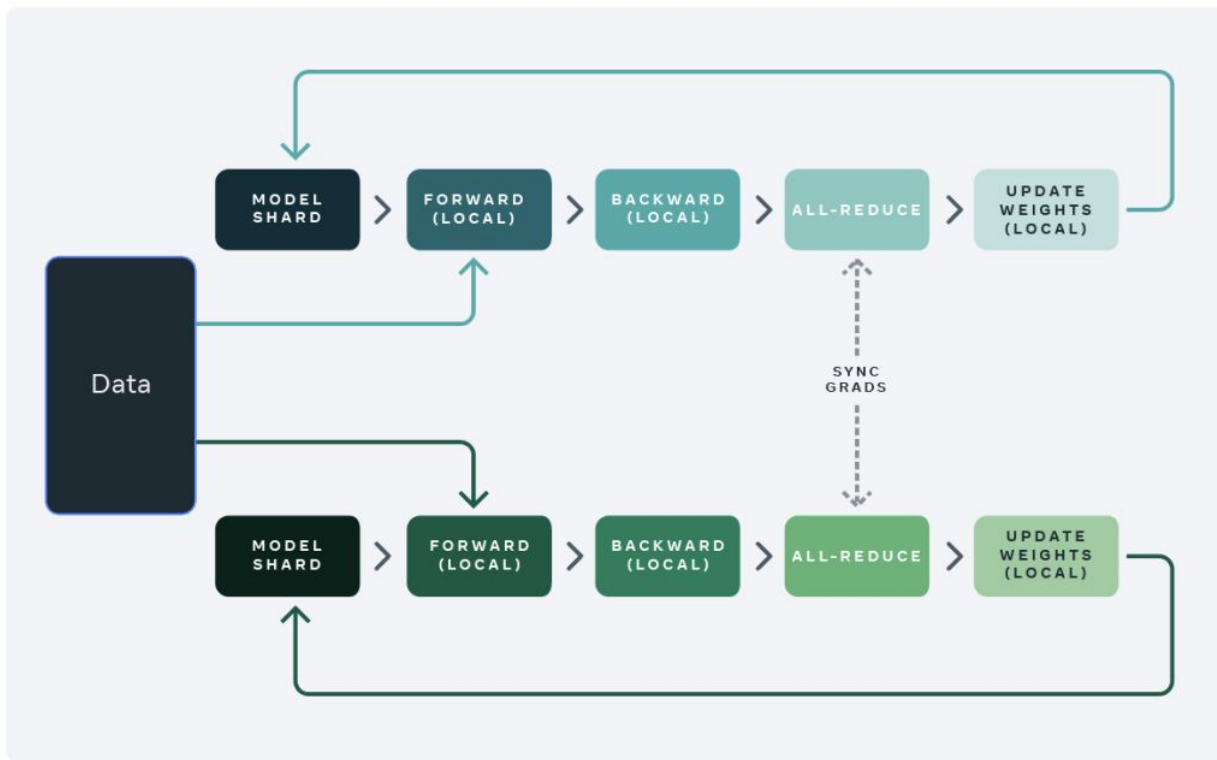


Scaling Parameters | Data Parallel Training



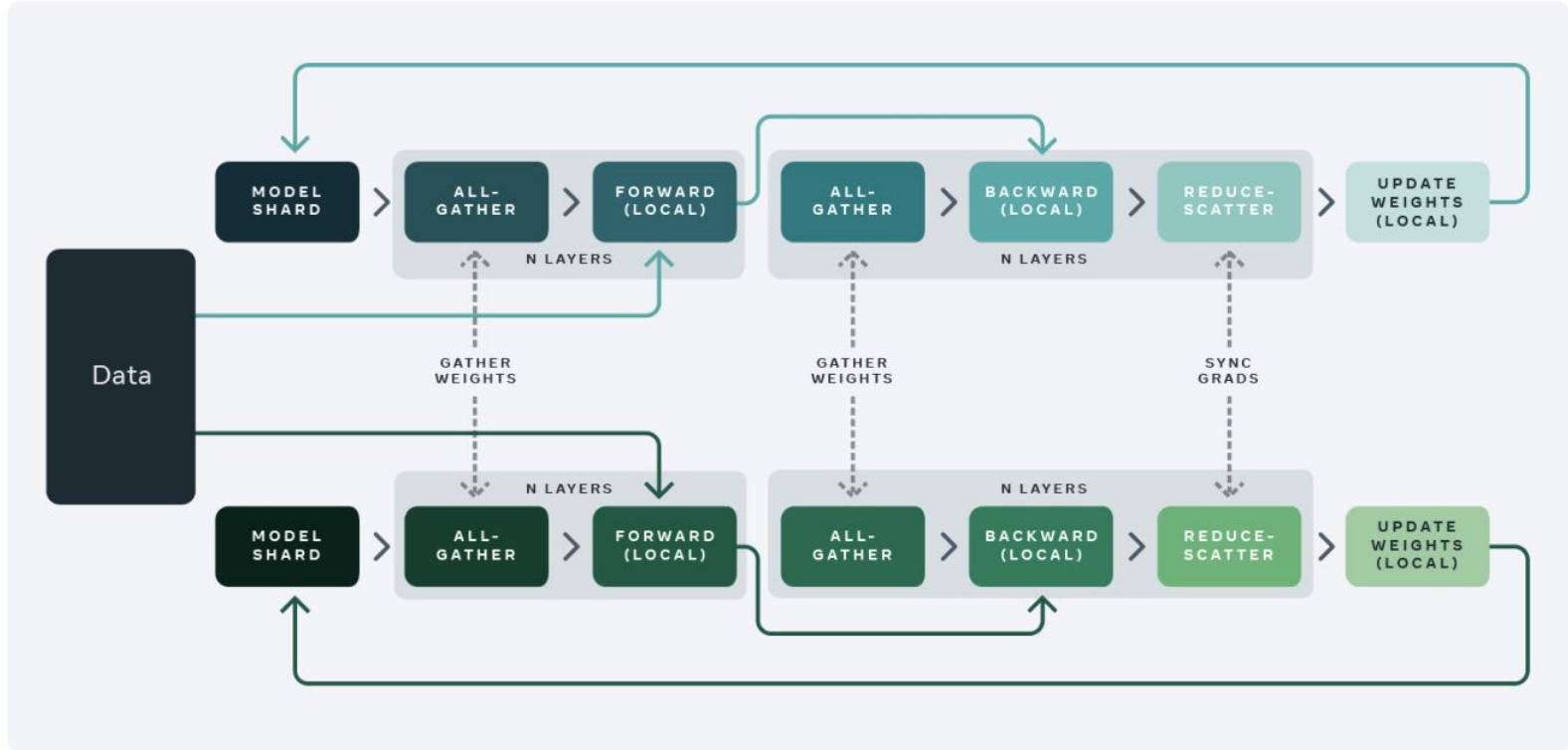
Total memory increases linearly with shards

Scaling Parameters | Data Parallel Training

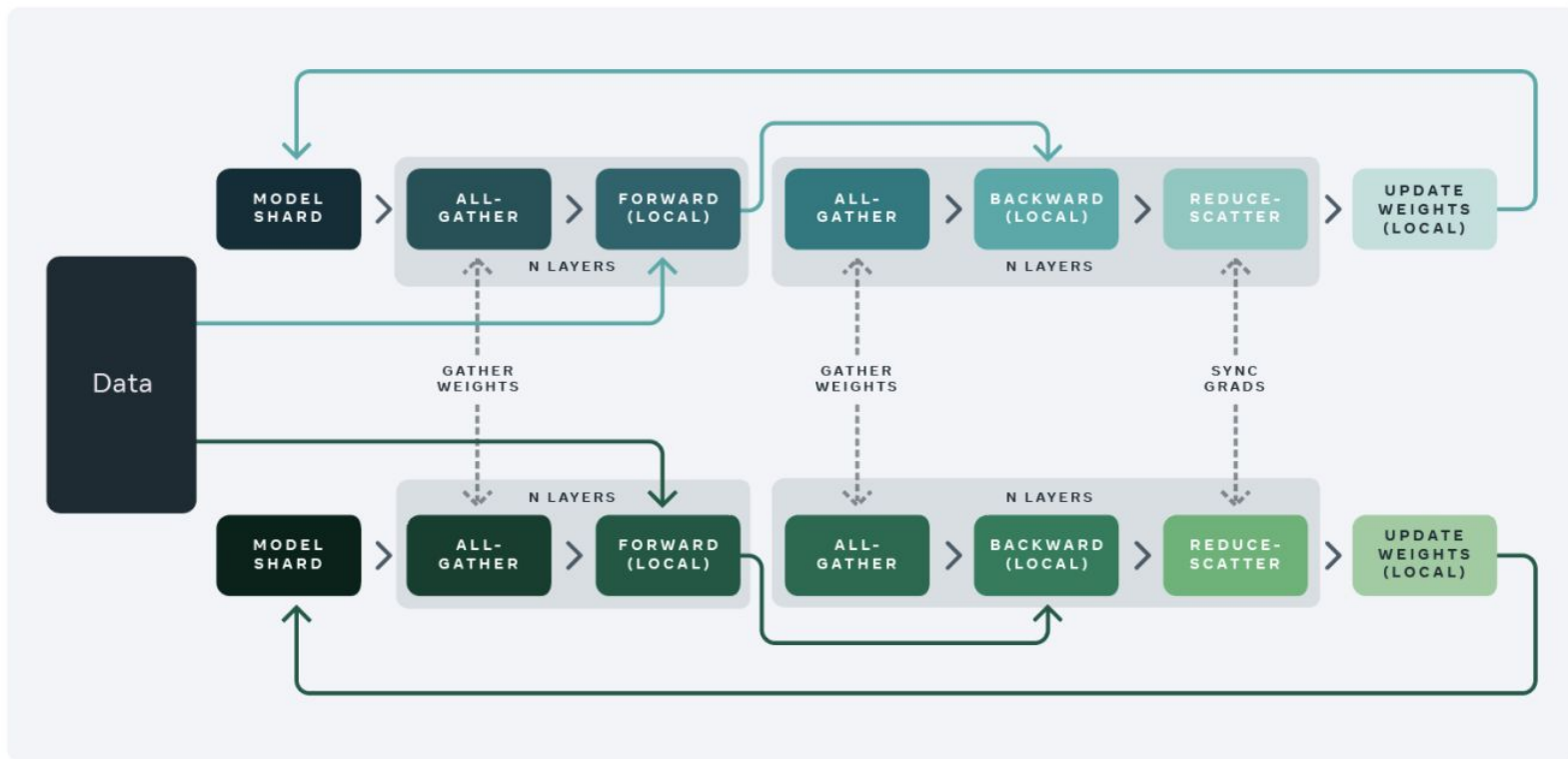


Max memory constrains model size

Scaling Parameters | *Fully* Sharded Data Parallel Training

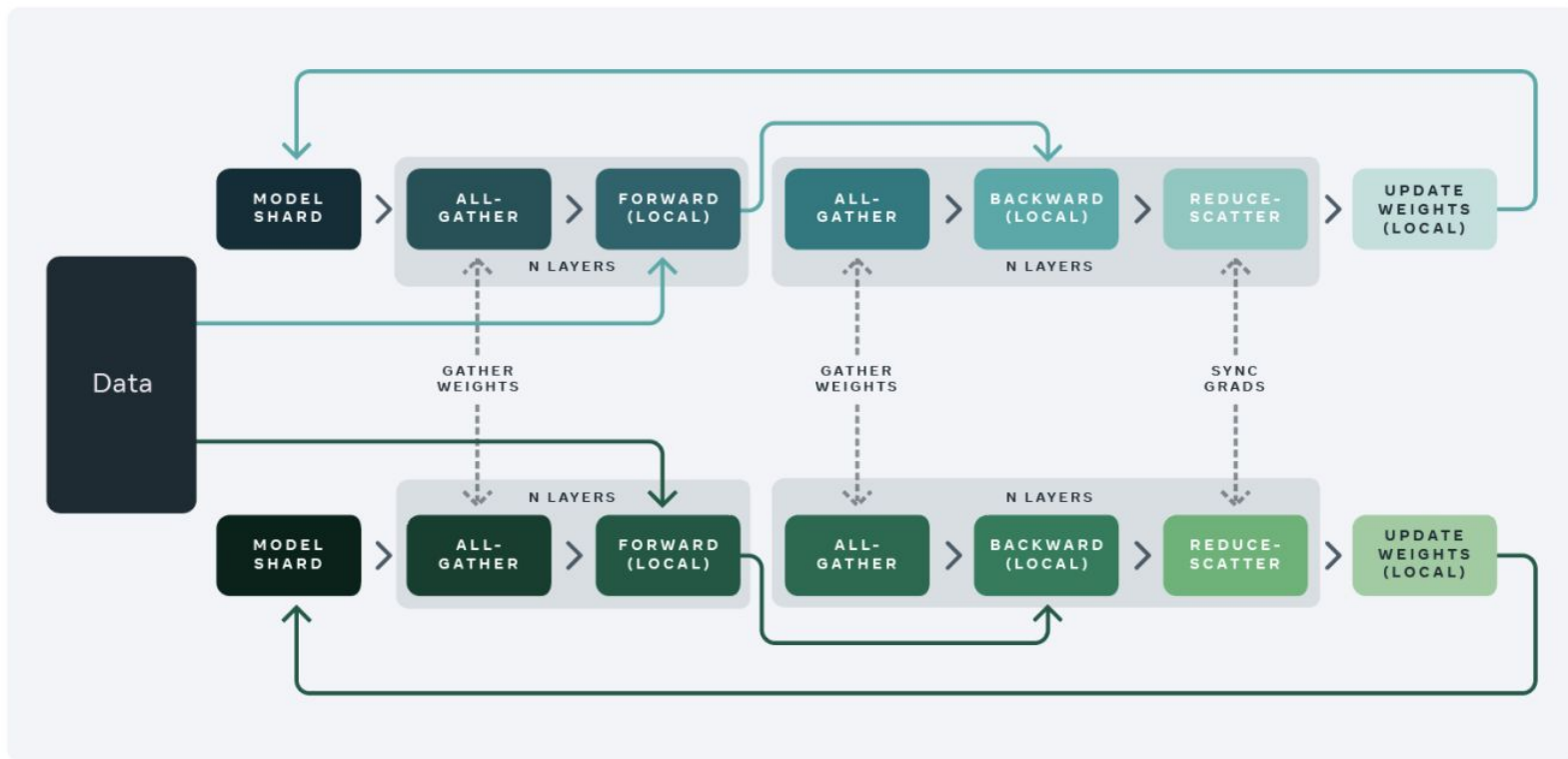


Scaling Parameters | *Fully* Sharded Data Parallel Training



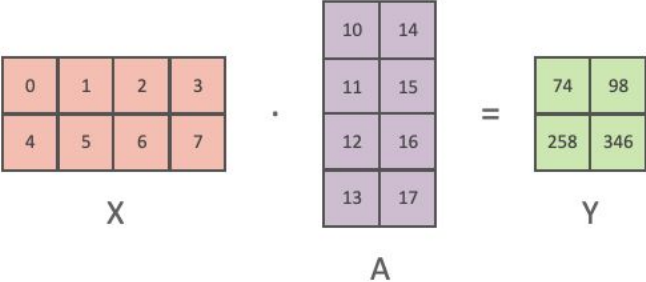
Total memory is constant

Scaling Parameters | *Fully* Sharded Data Parallel Training

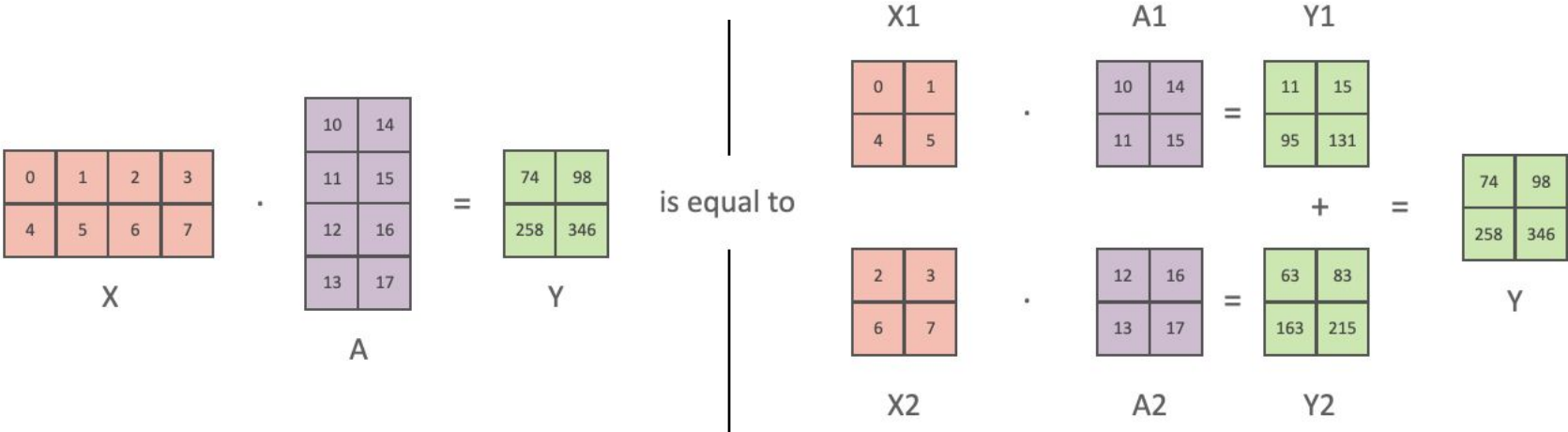


Max single GPU memory constrains layer size

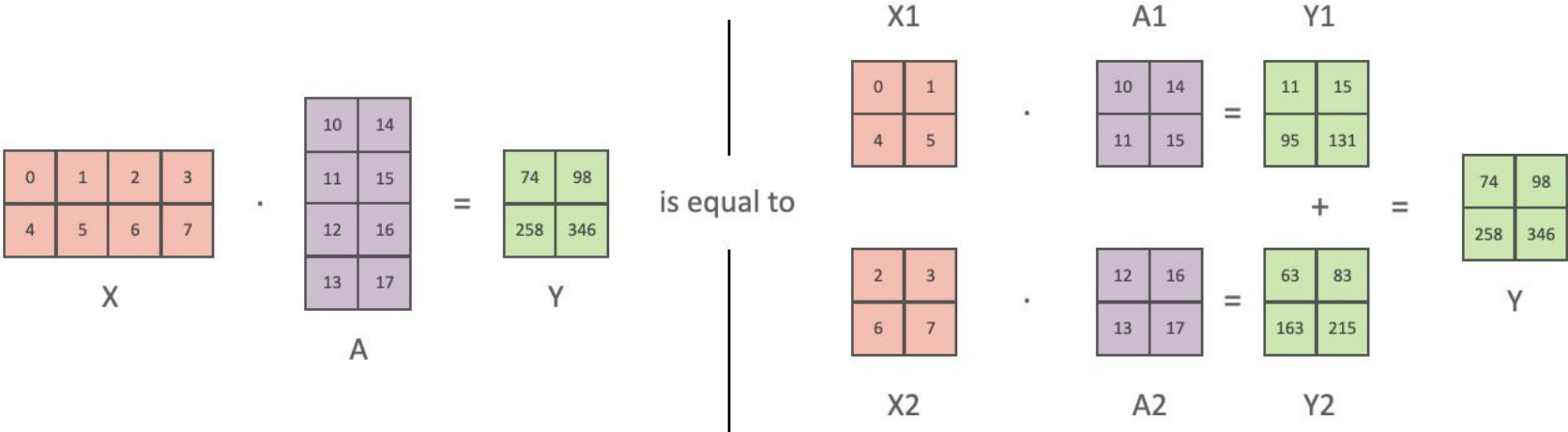
Scaling Parameters | Tensor Parallel Training



Scaling Parameters | Tensor Parallel Training

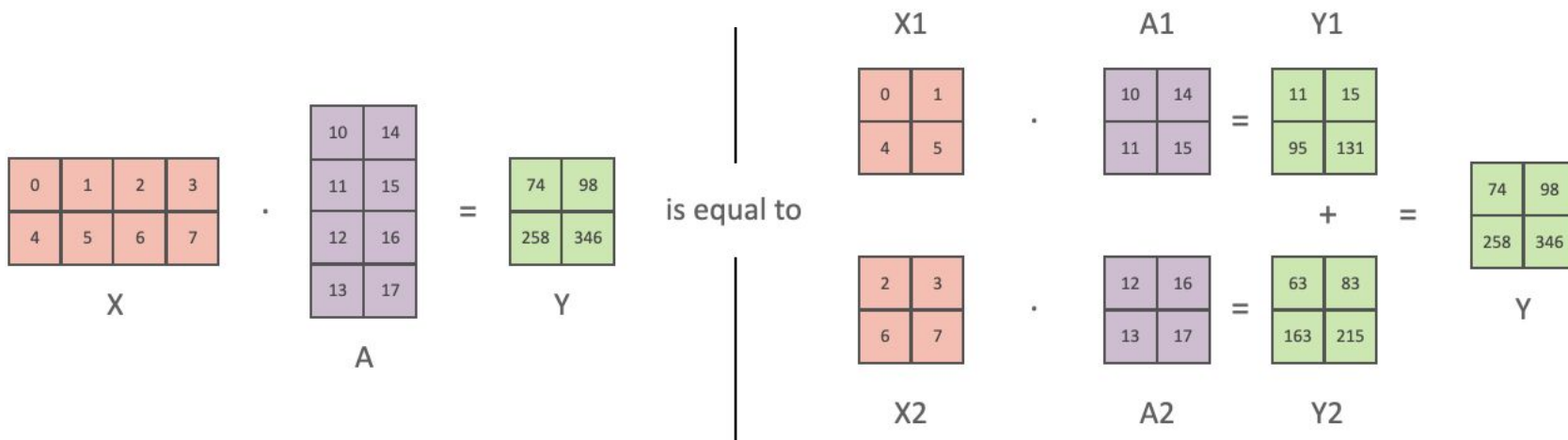


Scaling Parameters | Tensor Parallel Training



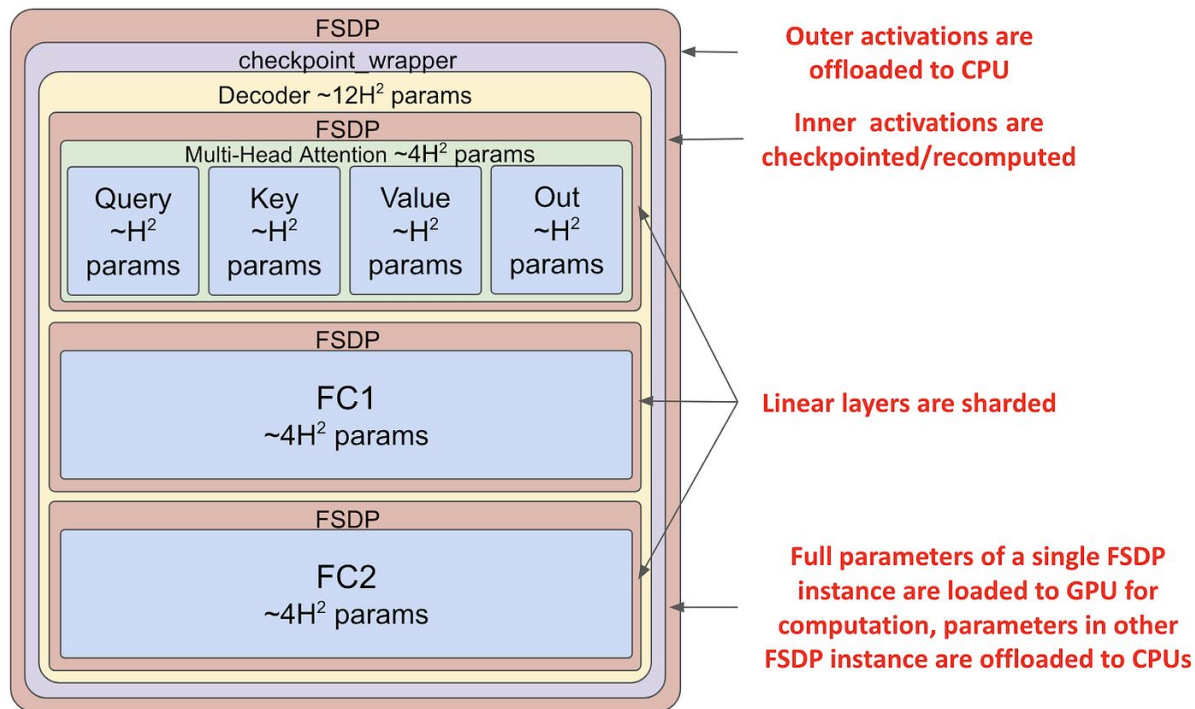
Don't need to sync gradients!

Scaling Parameters | Tensor Parallel Training

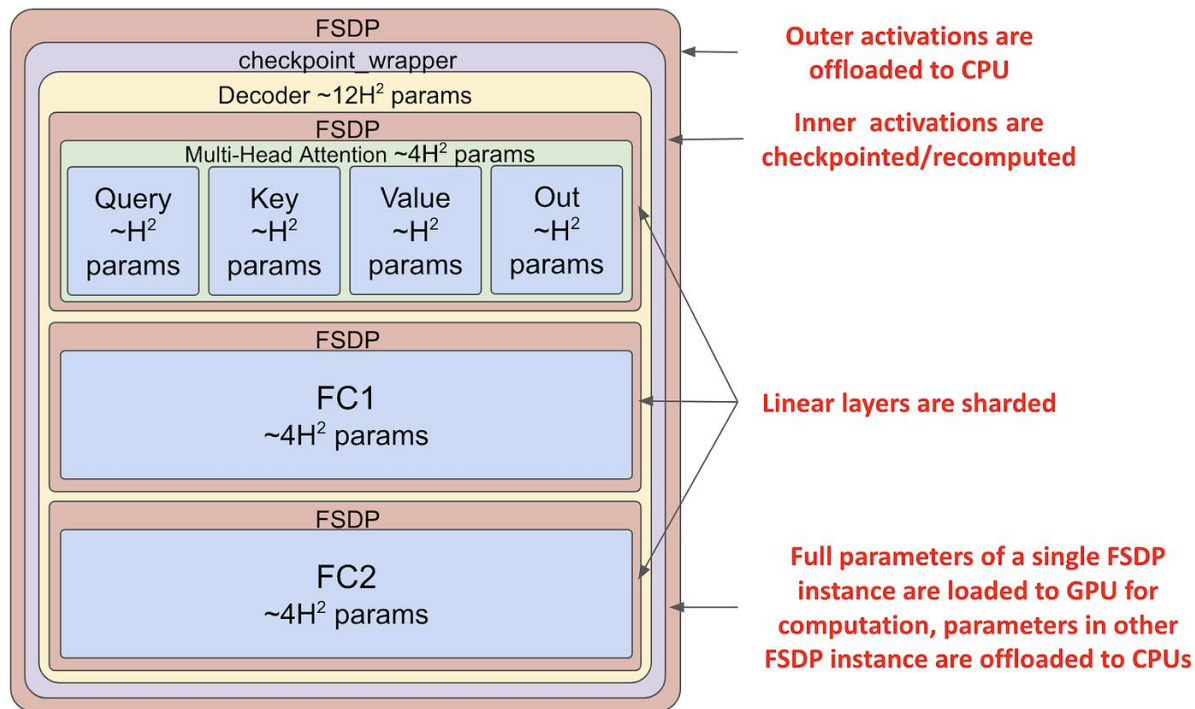


Don't need to sync gradients!
Max GPU memory constrains layer shard size

Scaling Parameters | FSDP + TP = ~Limitless Scaling

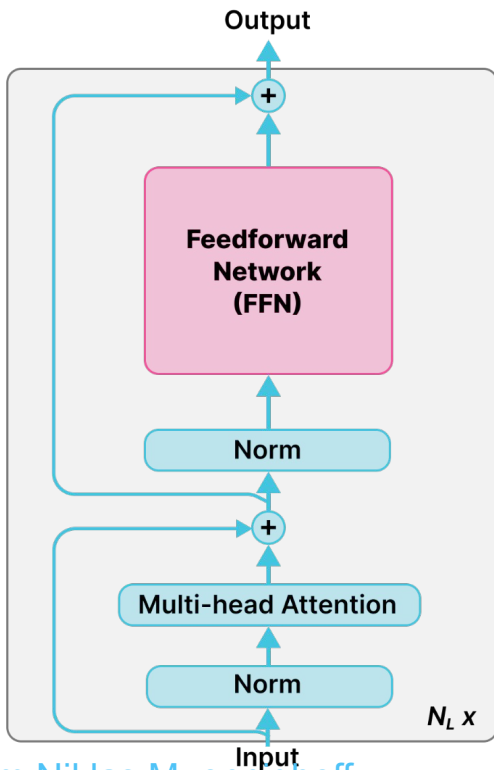


Scaling Parameters | FSDP + TP = ~Limitless Scaling

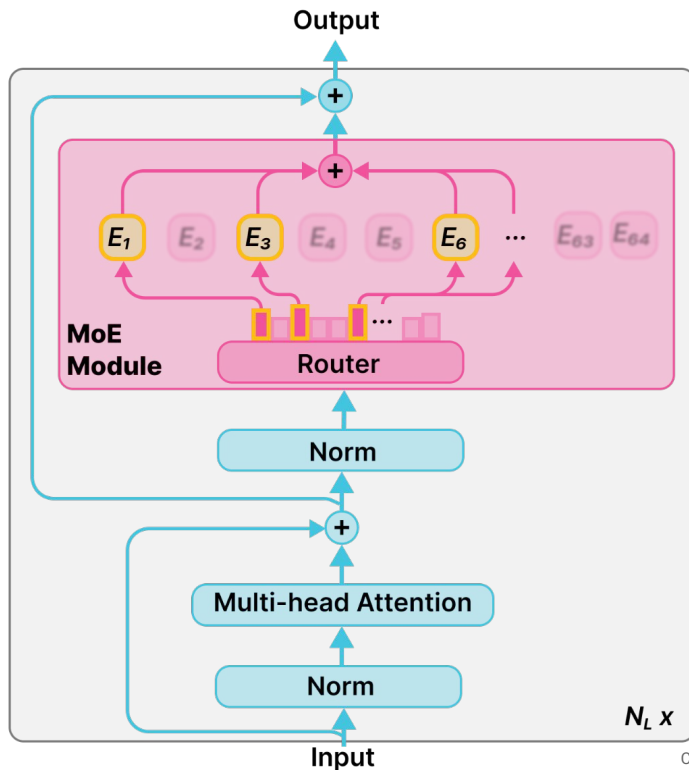


Scaling Parameters | Mixture Of Experts

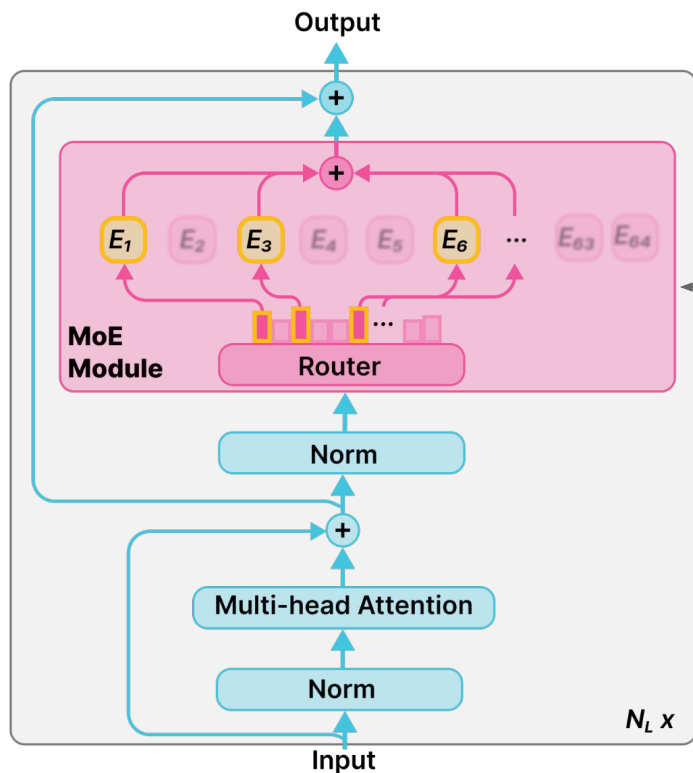
Dense LMs (OLMo, Llama...)



OLMoE



Scaling Parameters | Mixture Of Experts

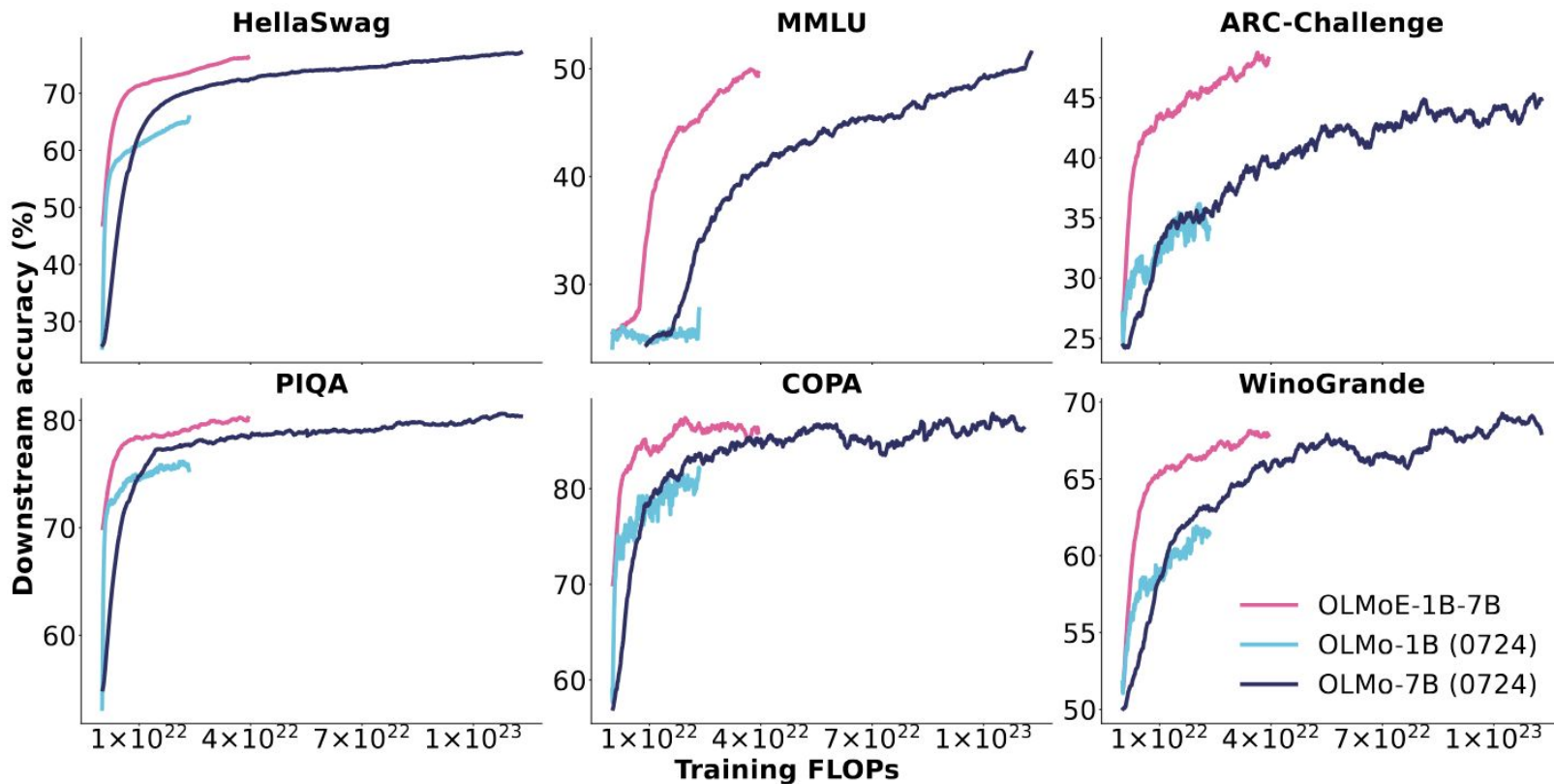


Router activation is sparse!

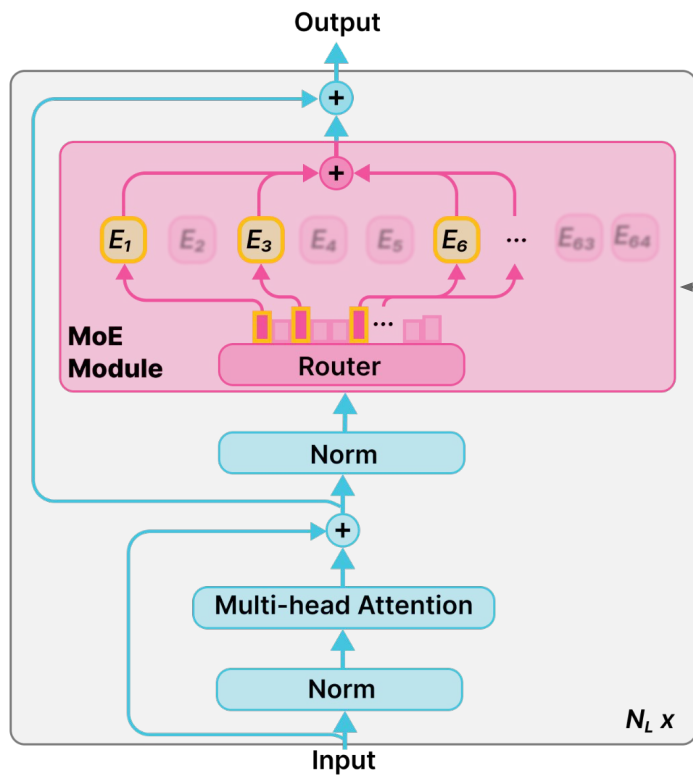
Pros of Mixture of Experts

- + Cheaper, Large Scale Training
- + Lower Inference Requirements

Scaling Parameters | Mixture Of Experts



Scaling Parameters | Mixture Of Experts

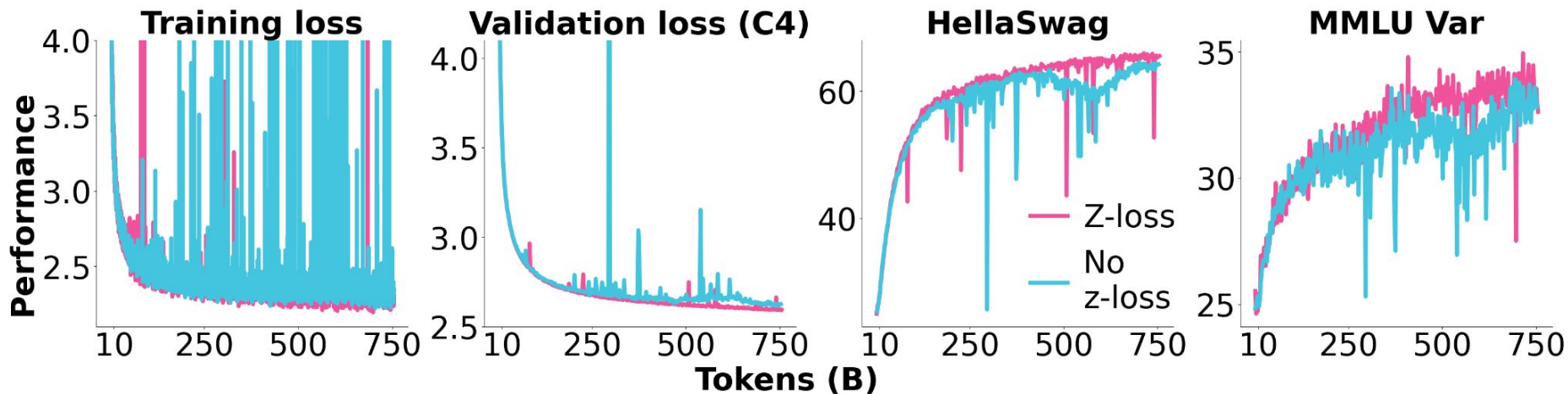


Router activation is sparse!

Cons of Mixture of Experts

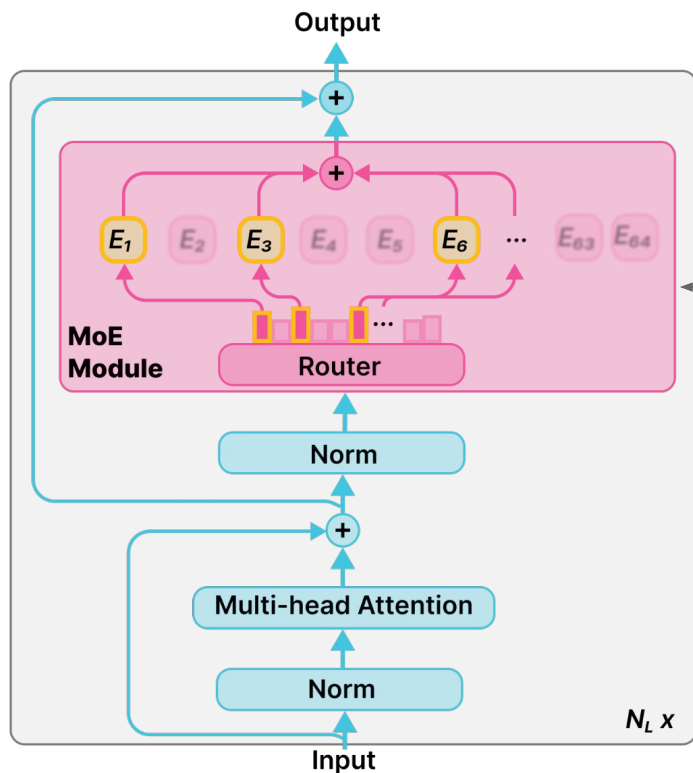
- More Unstable Training Runs

Scaling Parameters | Mixture Of Experts



Z-Loss decreases the magnitude of the logits into the router to stabilize gradients

Scaling Parameters | Mixture Of Experts

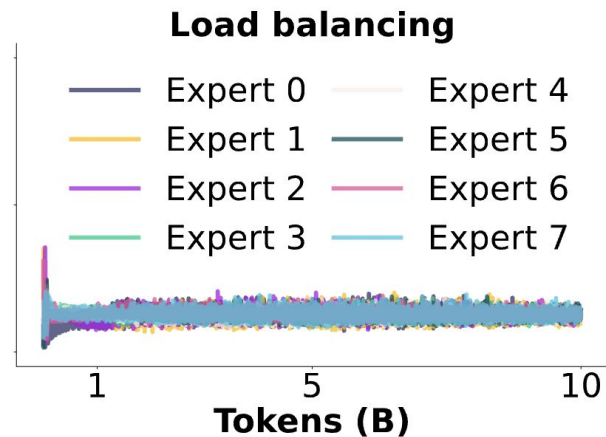
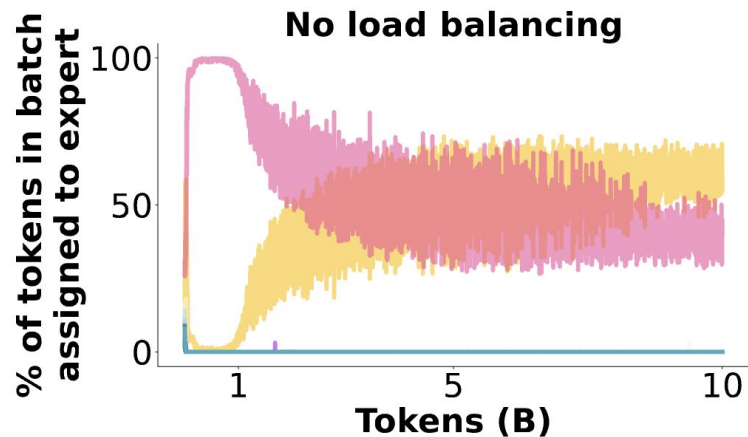


Router activation is sparse!

Cons of Mixture of Experts

- More Unstable Training Runs
- Experts frequently collapse

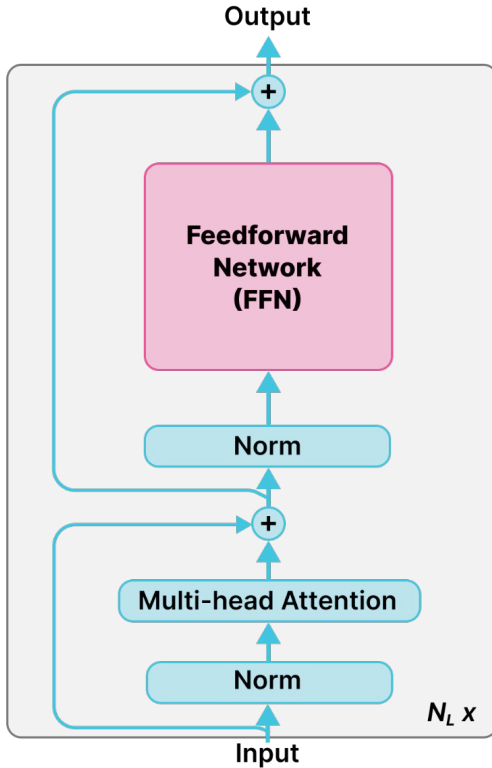
Scaling Parameters | Mixture Of Experts



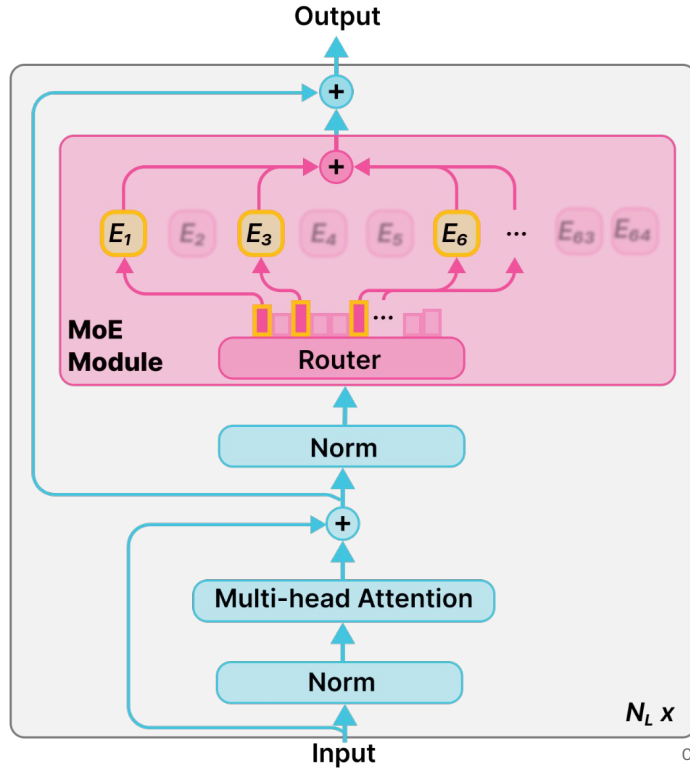
Load Balancing loss to update all experts similarly

Questions?

Dense LMs (OLMo, Llama...)



OLMoE



LLMs aren't very good for “assistance” by default



Suggest 3 things to do in San Mateo.



San Mateo is a city in San Mateo County, California, in the high-tech enclave of Silicon Valley in the San Francisco Bay Area. With a population of 97,207 at the 2010 census, the city is located on the San Francisco Peninsula halfway between San Francisco and San Jose, about 20 miles (32 km) from the Pacific Ocean.

Massively Multitask Supervised “Instruction” Finetuning

“Reverse this string:
‘OLLEH’.”

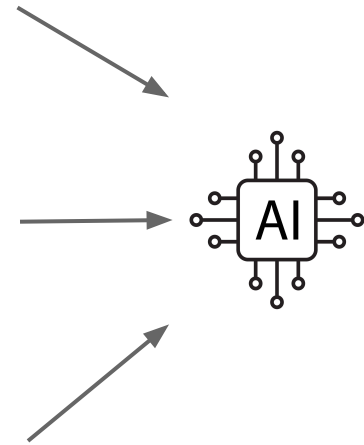
+ ““HELLO””

“Suggest 3 things to
do in San Mateo.”

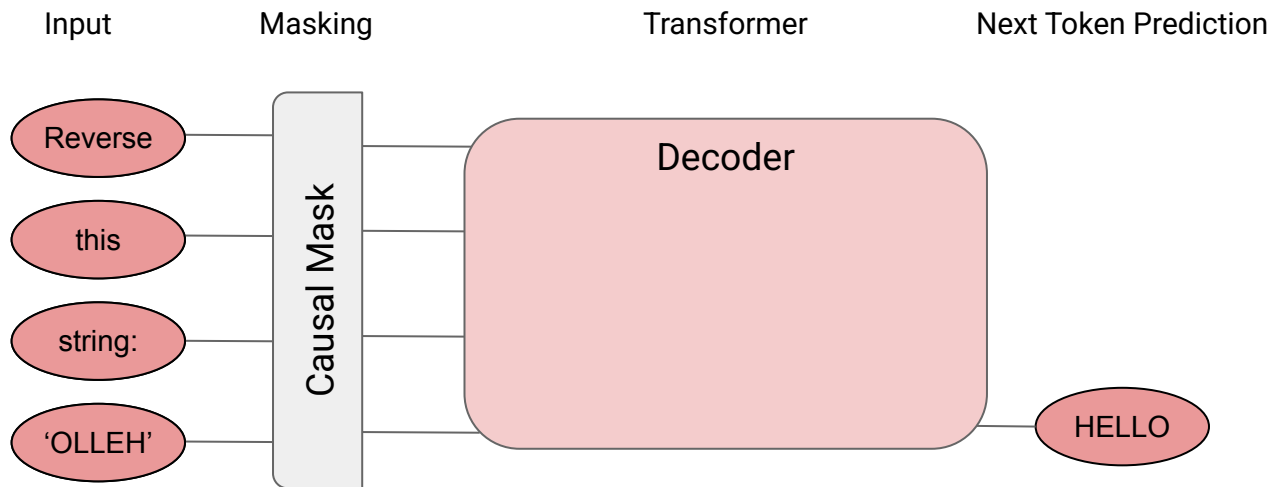
+ “1. Visit San Mateo,
Central Park...”

“What is the capital
of Qatar?”

+ “Doha”



Instruction Tuning | Just keep training!









Optimize Negative Log Likelihood of The Response

$$\text{loss} = -\log(P(\text{RESPONSE} \mid \text{INSTRUCTION}))$$

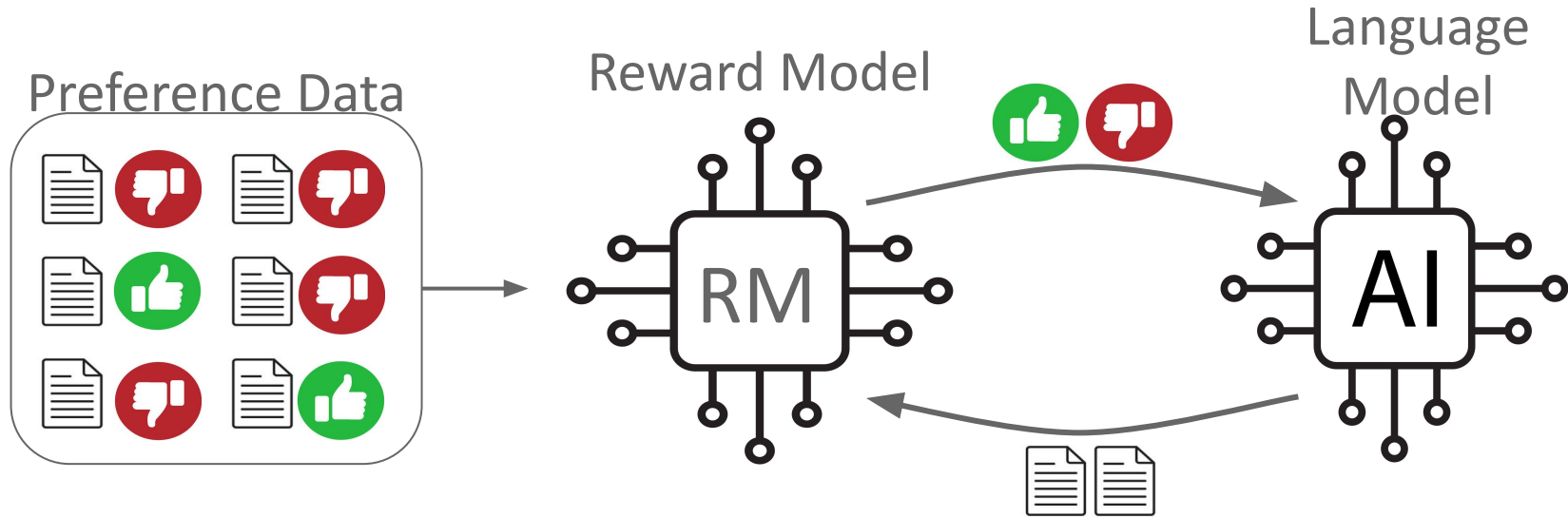
Further Refinement from Sparse Reward (RLHF)

“Suggest 3 things to do in San Mateo.”

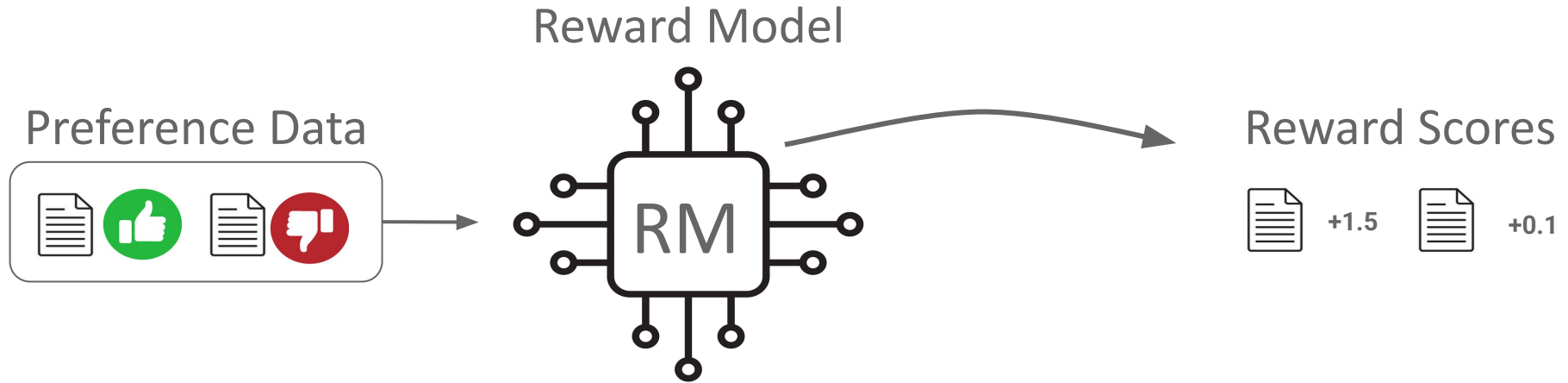


- “San Mateo is a city in San Mateo County...”  
- “1. Visit San Mateo, Central Park...”  
- “I’m sorry I can’t help.”  

Reinforcement Learning From Human Feedback (RLHF)



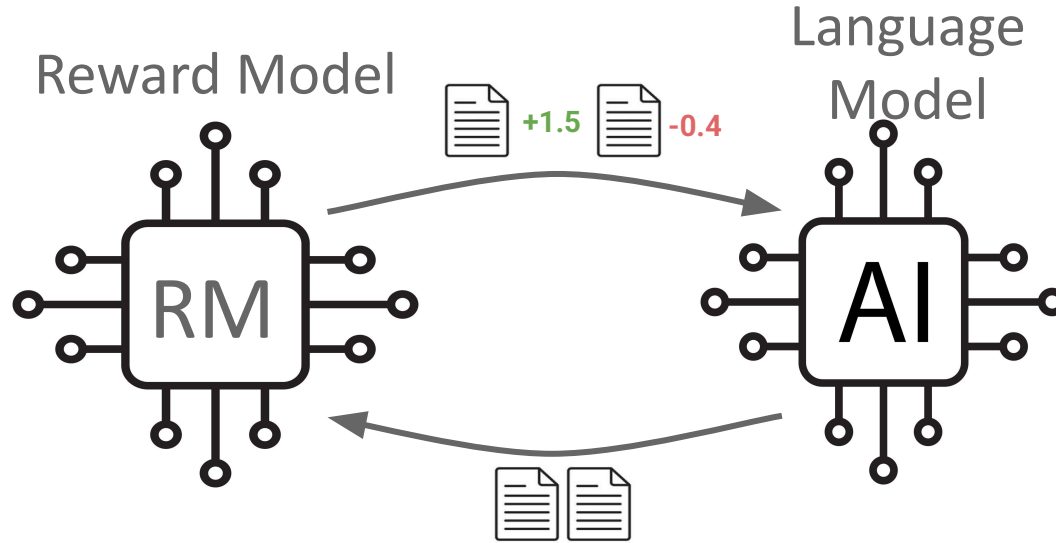
Reinforcement Learning From Human Feedback (RLHF)



Optimize Reward Margin between Preferences

$$\text{loss} = -\log(\sigma(\text{RM}(\text{POSITIVE}) - \text{RM}(\text{NEGATIVE})))$$

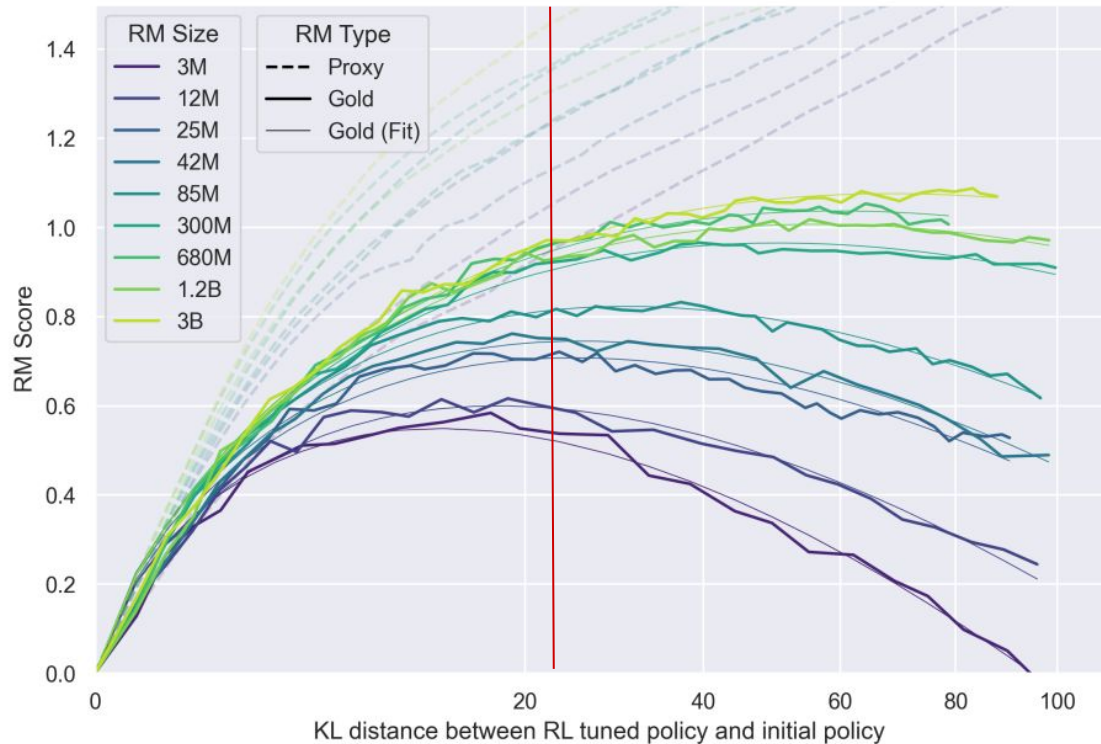
Reinforcement Learning From Human Feedback (RLHF)



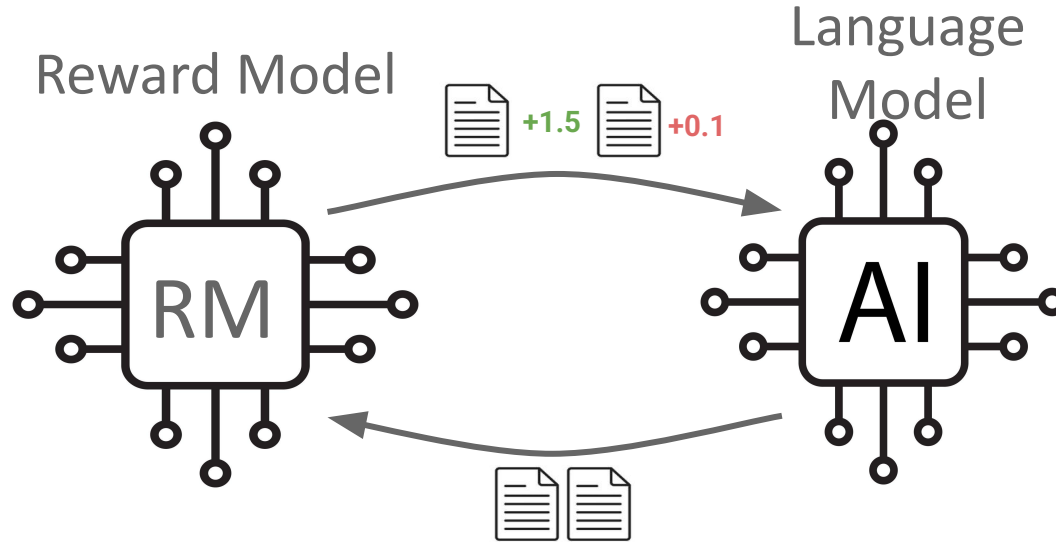
Optimize Reward Margin between Preferences

$$\text{loss}_{\text{RM}} = -\text{RM}(\text{GENERATED_EXAMPLES})$$

Models Quickly Overfit to Naively Optimized Reward



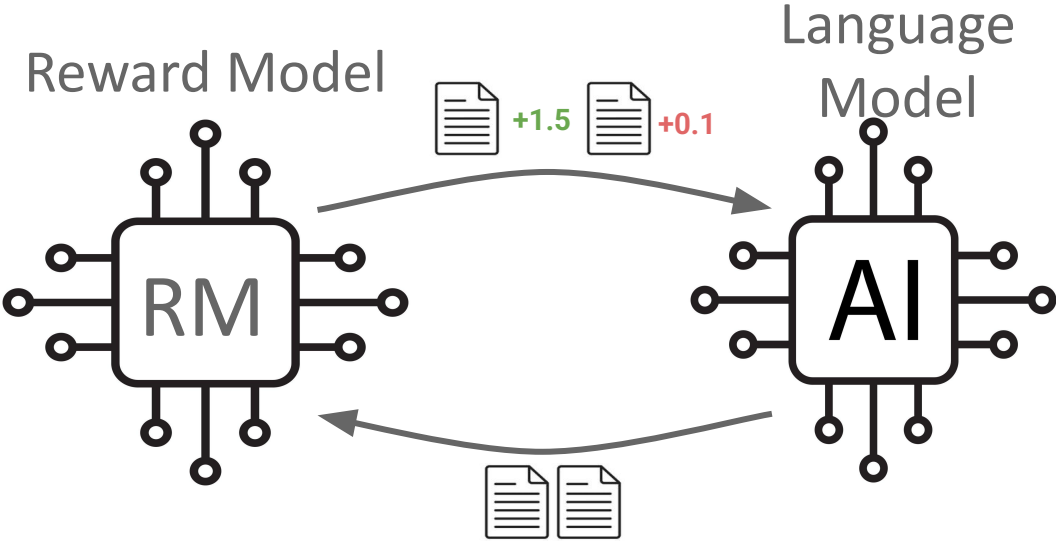
Reinforcement Learning From Human Feedback (RLHF)



Optimize Reward Without Drifting Too Far from SFT

$$\text{loss}_{\text{RLHF}} = \text{loss}_{\text{RM}} + \text{KL}(\text{LM}_{\text{RLHF}}, \text{LM}_{\text{SFT}})$$

Questions?



Final Questions?

Fill out my anonymous feedback form

