

Introduction and Logistic

COMP4551

Binhang Yuan



Amazing Progress of AIGC

 DeepSeek-R1 is now live and open source, rivaling OpenAI's Model o1. Available on web, app, and API. Click for details.

deepseek

Into the unknown

Start Now
Free access to DeepSeek-V3.
Experience the intelligent model.

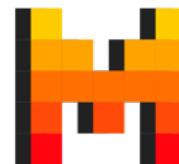
Get DeepSeek App
Chat on the go with DeepSeek-V3
Your free all-in-one AI tool

November 12, 2025 Product Release

GPT-5.1: A smarter, more conversational ChatGPT

We're upgrading GPT-5 while making it easier to customize ChatGPT.
Starting to roll out today to everyone, beginning with paid users.

[Try in ChatGPT ↗](#)





$$\min_x \mathbb{E}_\xi f(\xi, x)$$



$$\min_x \mathbb{E}_\xi f(\xi, x)$$

Data

- (ImageNet) 1.3M Images (est. 160+ GB)
- (Llama-3.1) 15 Tillion Tokens (est. 100+ TB)

Model

- (GPT-2) 1.3 Billion Parameters (2.6 GB fp16)
- (Llama-3.1) 405 Billion Parameters (810GB fp16)

Compute

- (GPT-2) est. 2.5 GFLOPS/token
- (Llama-3.1) est. 1.2 TFLOPS/token

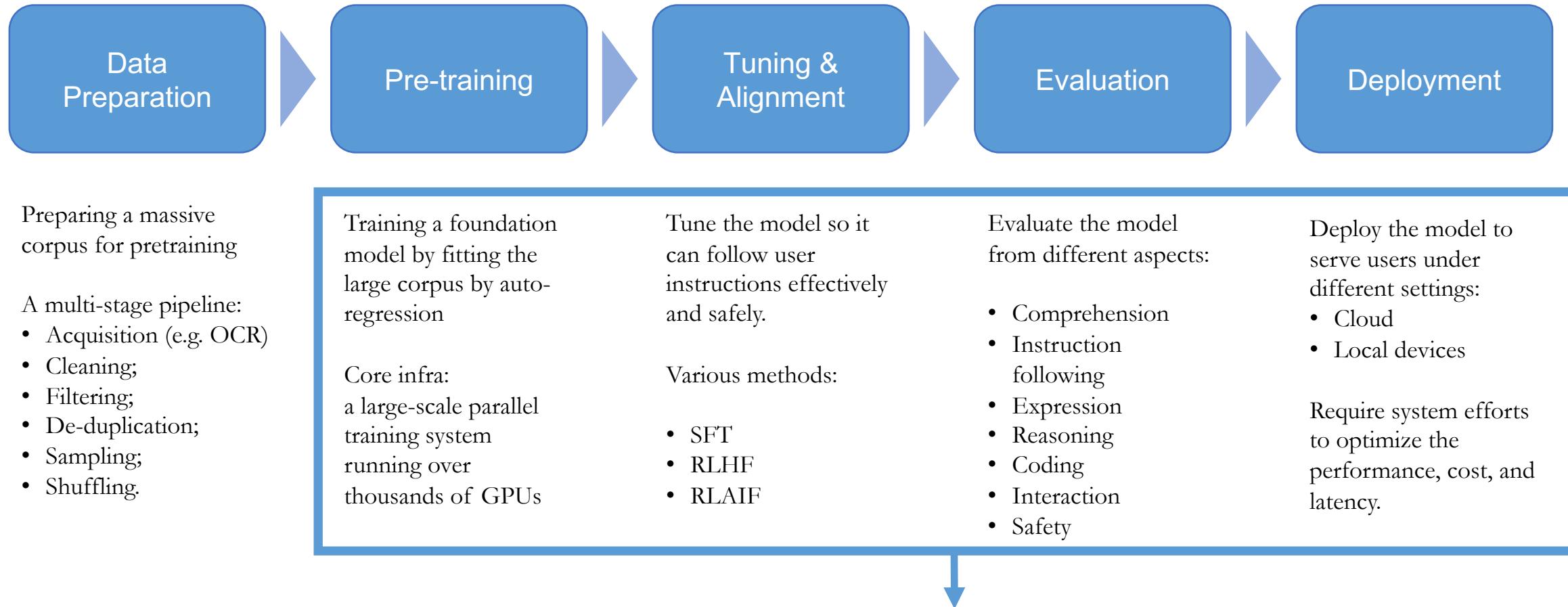


The goal of this course:

Unravel the secrets of such foundation models from the both the algorithm and system perspective!



The Path Towards a Foundation Model



Covered by this course!



RELAXED
SYSTEM LAB

Course Overview



Syllabus

RELAXED
SYSTEM LAB

Date	Topic
W1 - 02/03, 02/05	Introduction and Logistics & ML Preliminary
W2 - 02/10, 02/12	Stochastic Gradient Descent & Automatic Differentiation
W3 - 02/17, 02/19	Spring Festival
W4 - 02/24, 02/26	Language Model Architecture & Large Scale Pretrain Overview
W5 - 03/03, 03/05	Nvidia GPU Performance & Collective Communication Library
W6 - 03/10, 03/12	Data-, Pipeline- Parallel Training & Tensor Model-, Optimizer- Parallel Training
W7 - 03/17, 03/19	Sequence-, MoE- parallelism & Mid-Term Review
W8 - 03/24, 03/26	Mid-Term Exam & Generative Inference



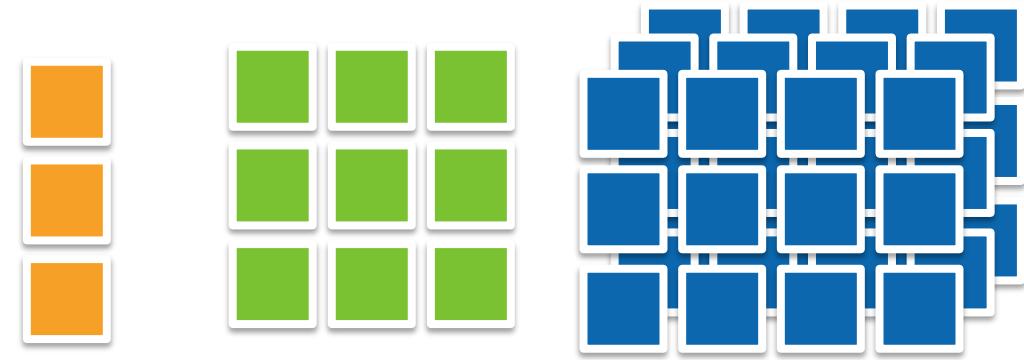
Syllabus

Date	Topic
W9 - 03/31, 04/02	Inference Alogirhtm Optimizations & Inference System Optimizations
W10 - 04/07, 04/09	Spring Break & Prompt Engineering
W11 - 04/14, 04/16	Inference Scaling & Retrieval Augmented Generation
W12 - 04/21, 04/23	LLM Agent & Parameter Efficient Fine-Tuning
W13 - 04/28, 04/30	RL Alignment & LLM Evaluation
W14 - 05/05, 05/07	Guest Speech (TBD) & Final Review



Machine Learning Preliminary

- Linear Algebra:
 - Vector, matrix, tensor.
- PyTorch Tensors.





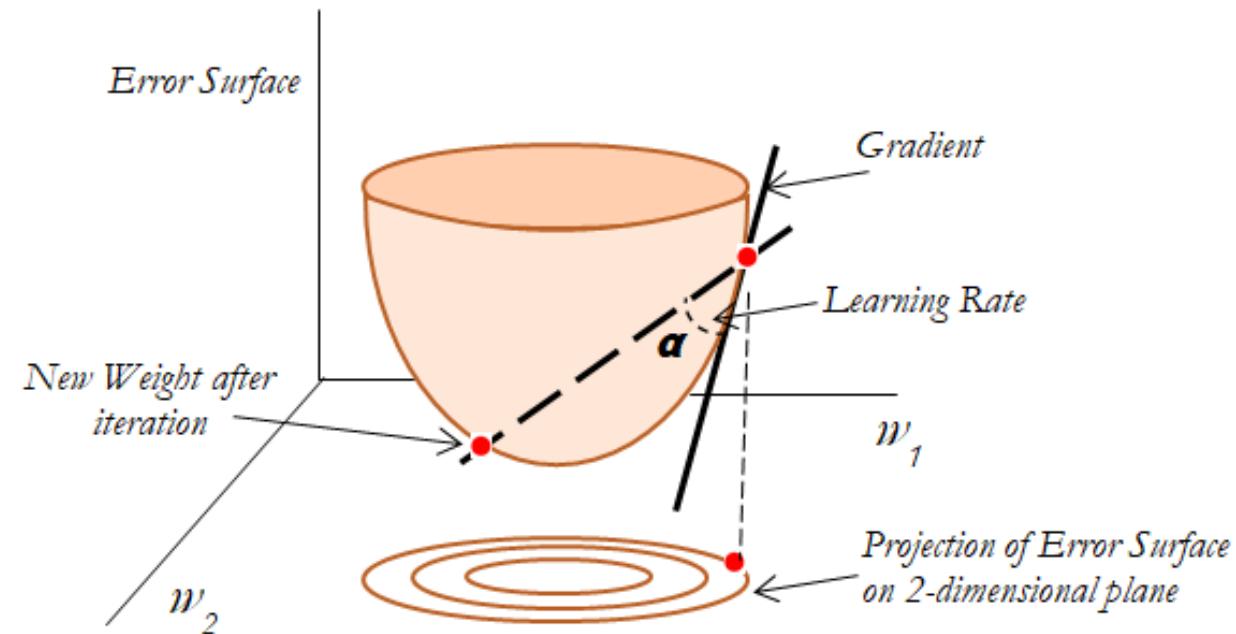
Stochastic Gradient Descent

- Then, suppose we have:
 - $f: \mathbb{R}^d \rightarrow \mathbb{R}$;
- Definition of a derivative/gradient :

- $\nabla f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix} \in \mathbb{R}^d$

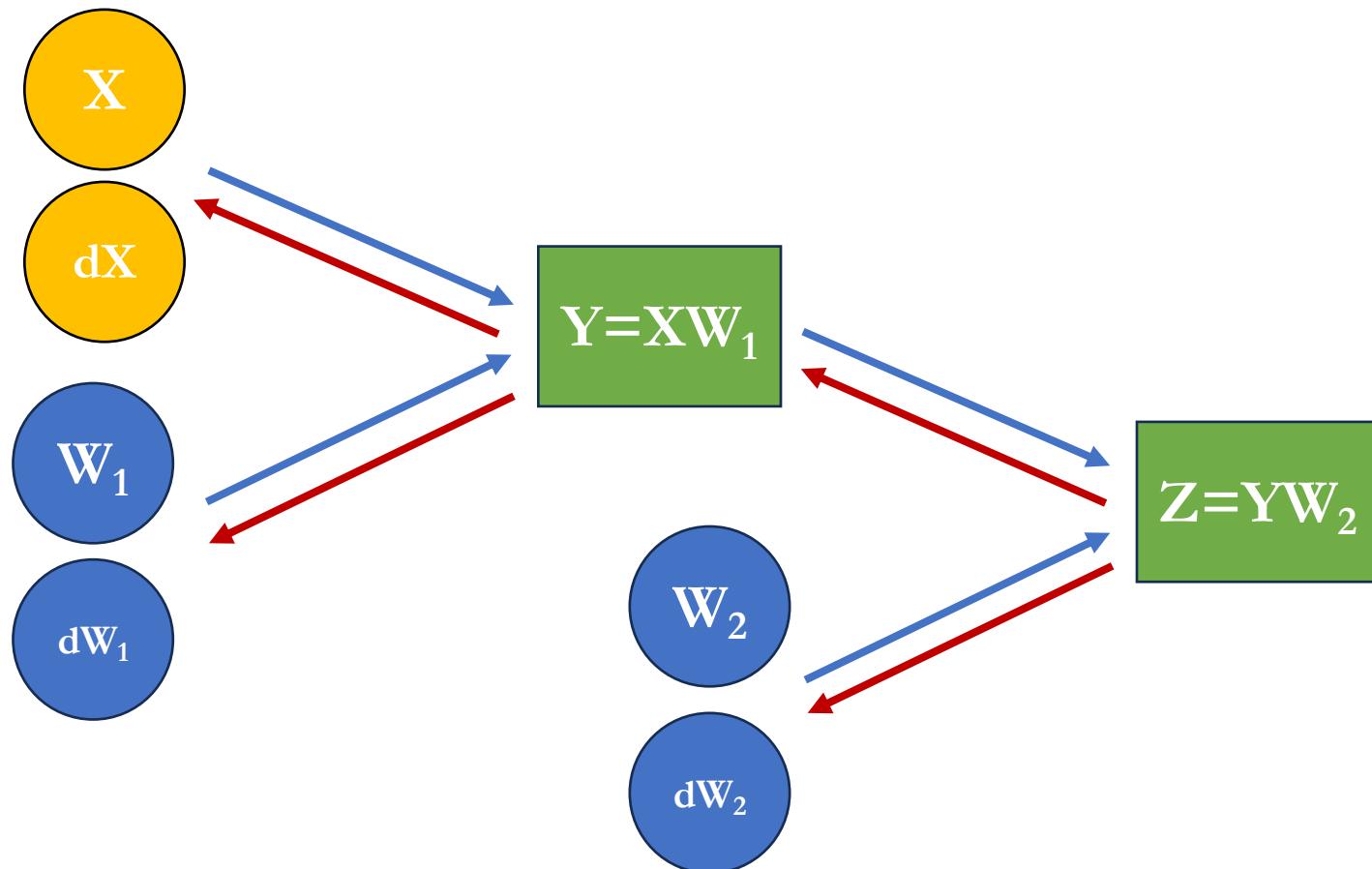
- Where:

- $$\frac{\partial f}{\partial x_i} = \lim_{\epsilon \rightarrow 0} \frac{f(x_1, x_2, \dots, x_i + \epsilon, x_{i+1}, \dots, x_d) - f(x_1, x_2, \dots, x_i, x_{i+1}, \dots, x_d)}{\epsilon}$$
$$= \lim_{\epsilon \rightarrow 0} \frac{f(x + \epsilon e_i) - f(x)}{\epsilon}$$





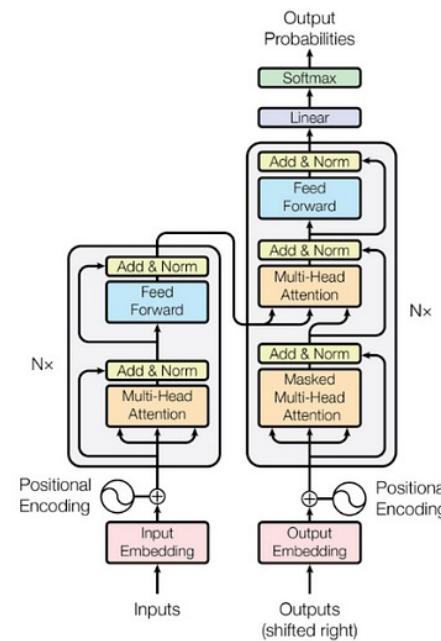
Auto-Differentiation & PyTorch Autograd





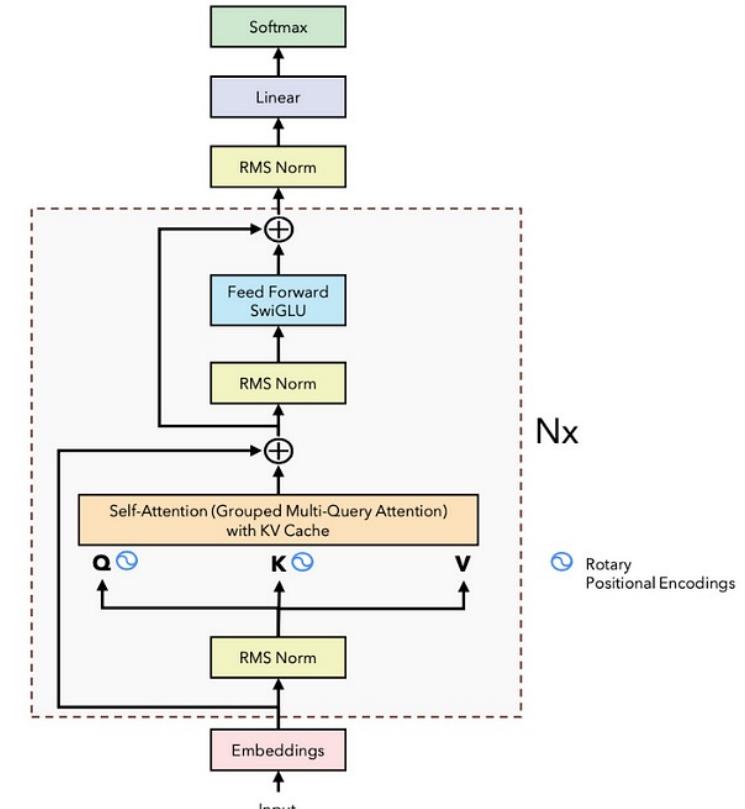
Transformer Architecture

Transformer vs LLaMA



Transformer

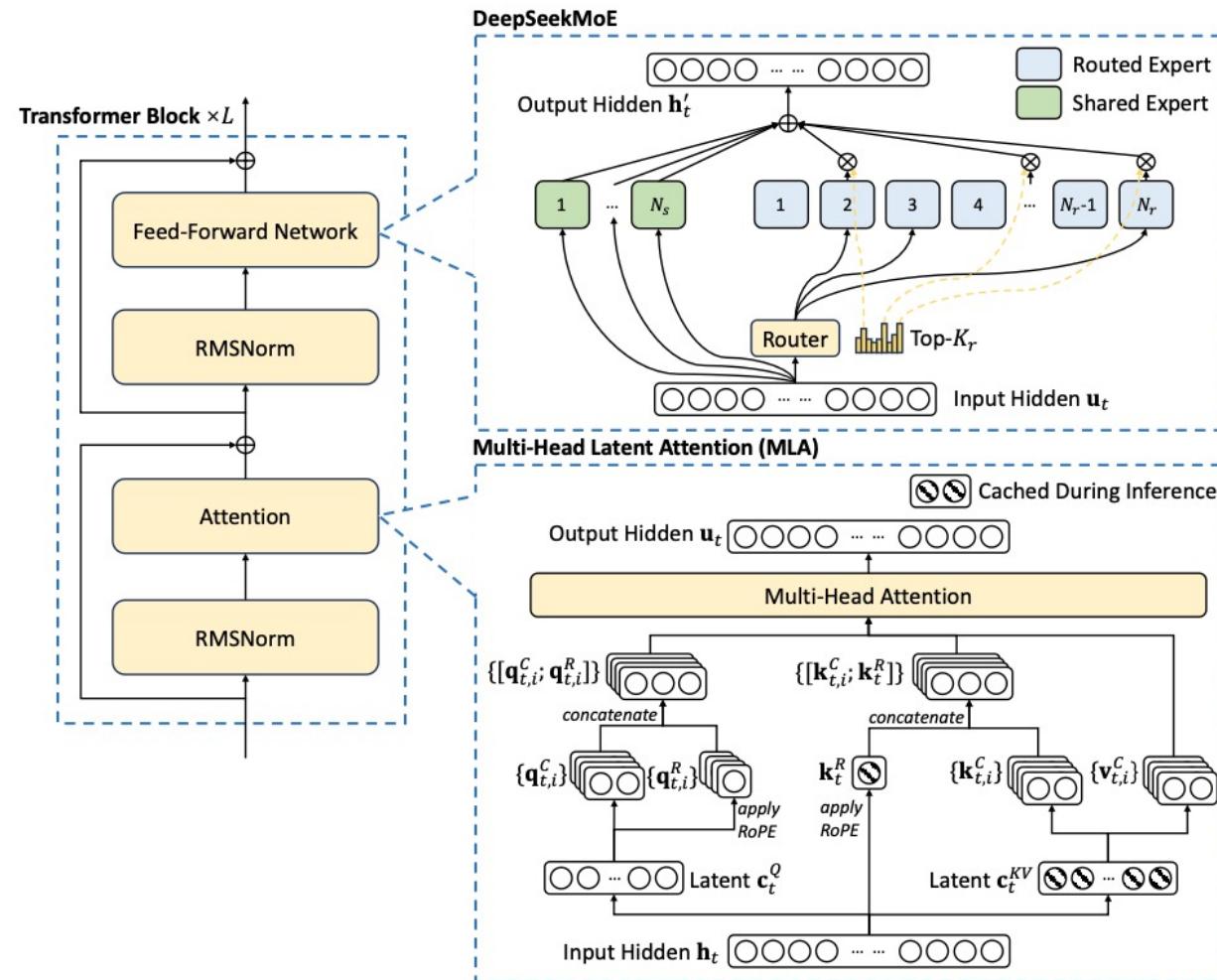
("Attention is all you need")



LLaMA



Transformer with MoE

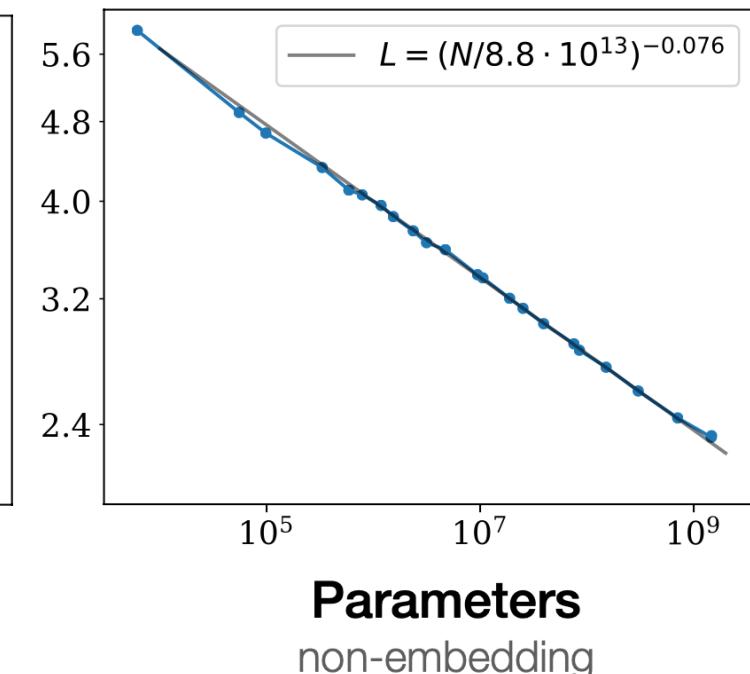
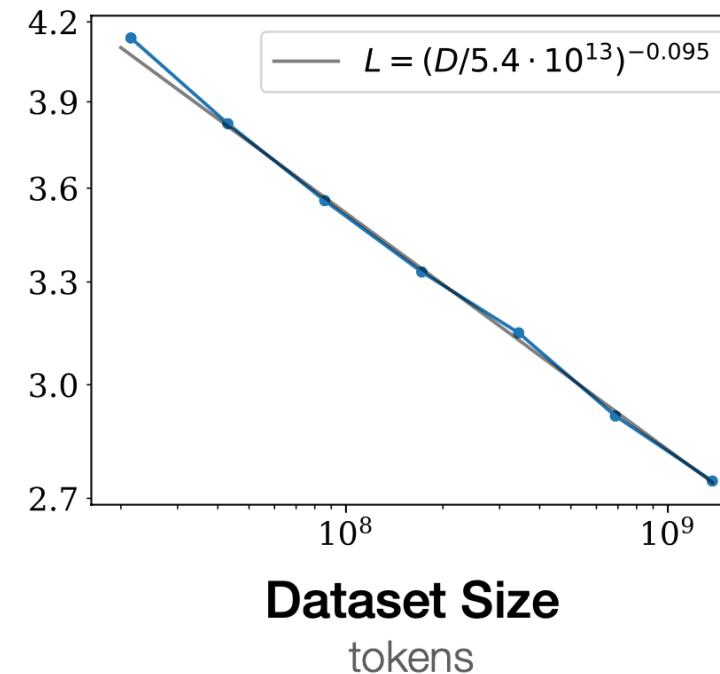
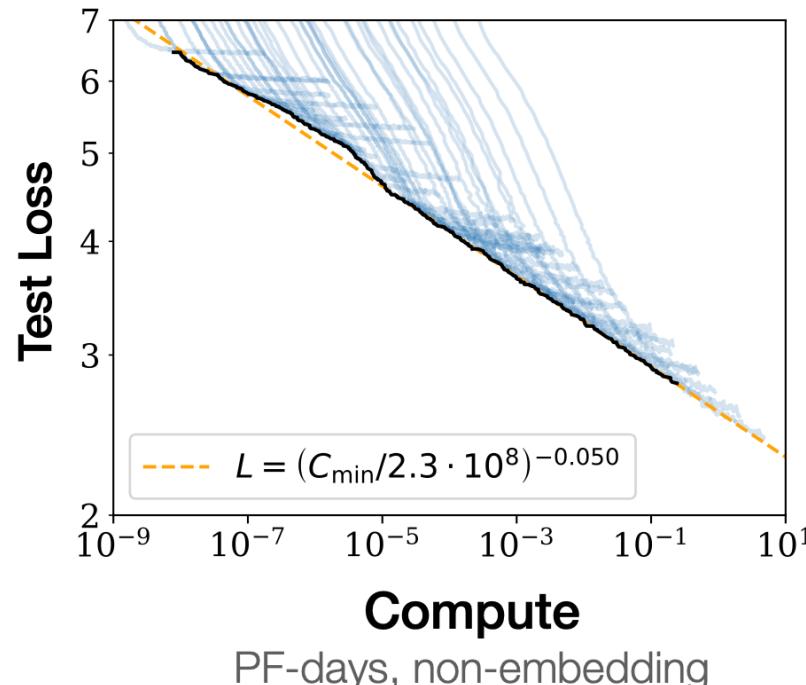


Deepseek V3



Large Scale Pretrain Overview

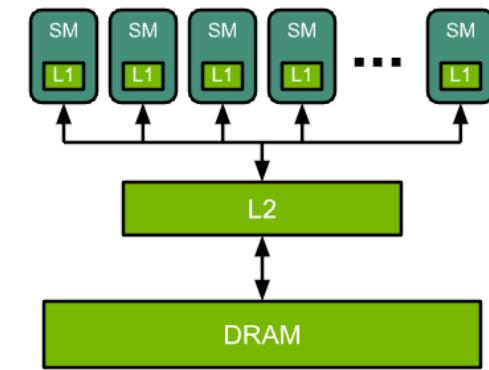
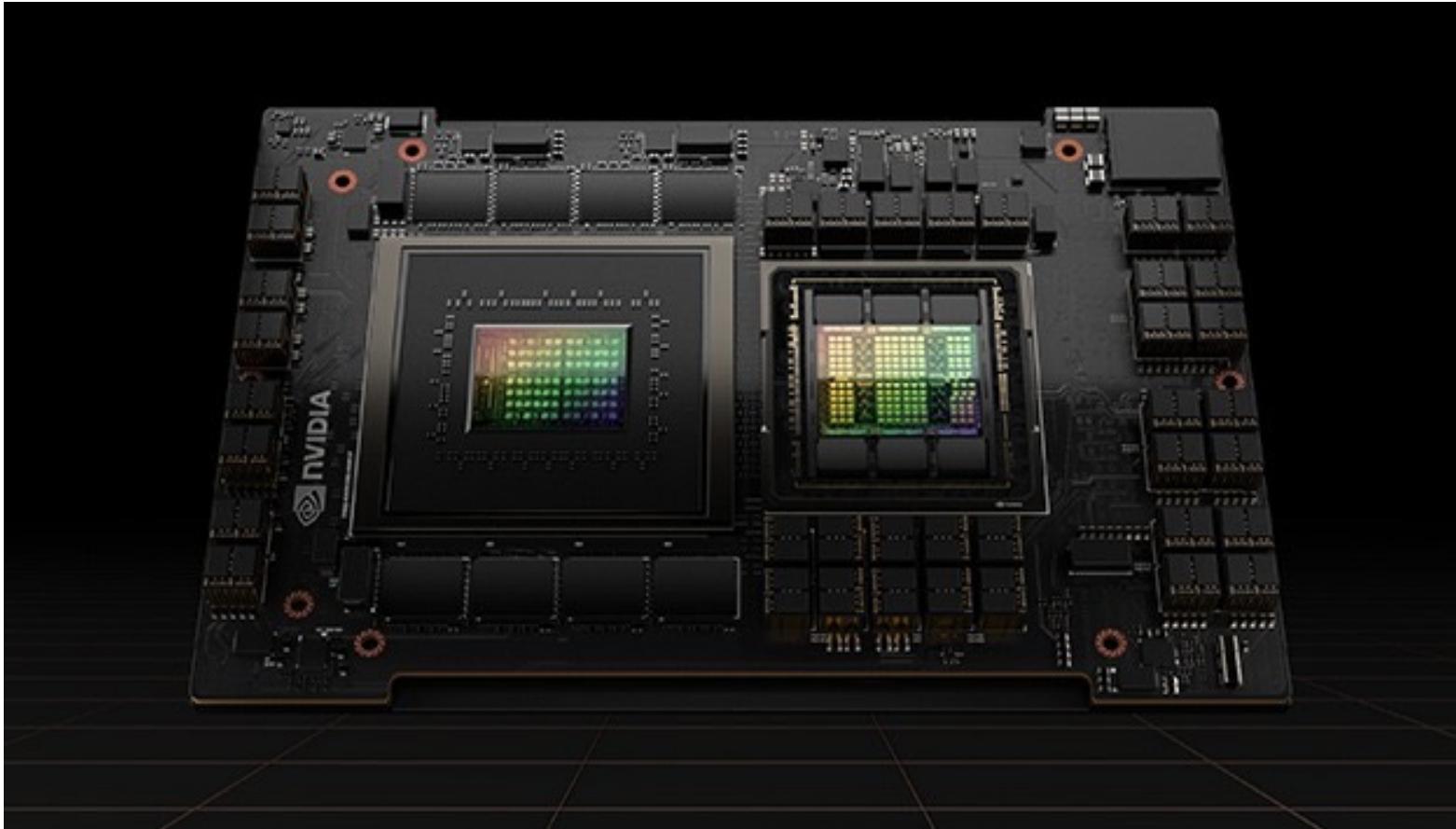
Scaling Laws for Neural Language Models



<https://arxiv.org/pdf/2001.08361.pdf>



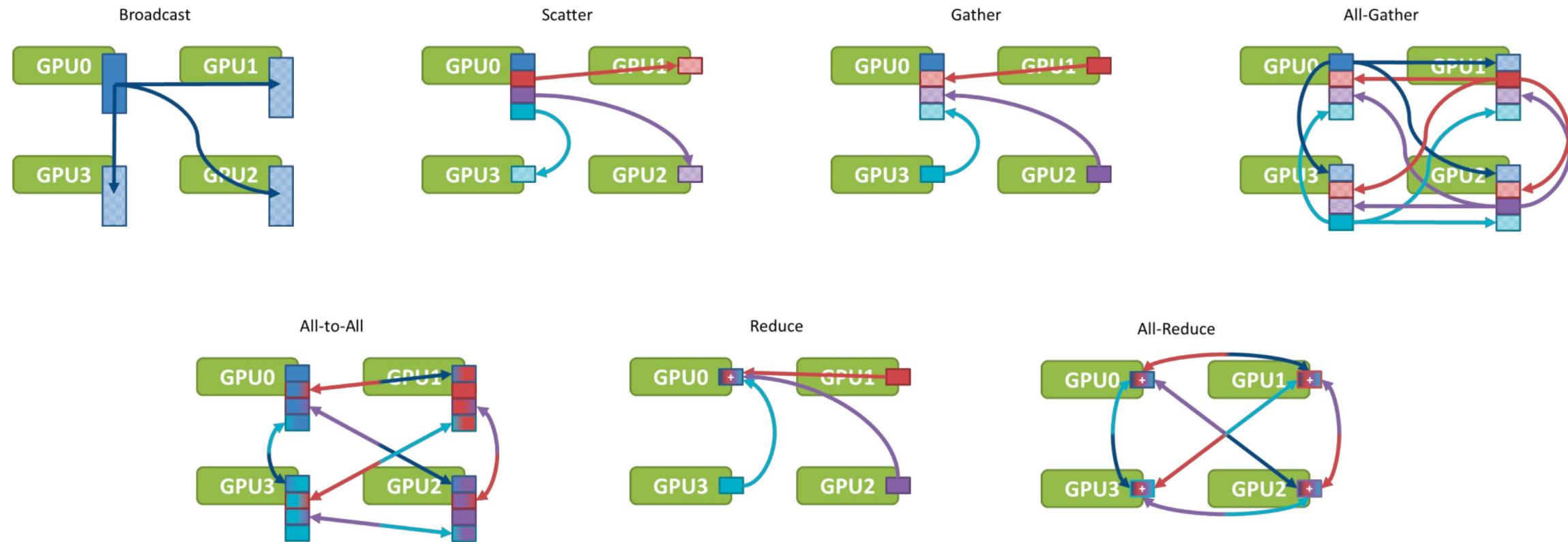
Nvidia GPU Performance



<https://www.nvidia.com/en-us/data-center/h100/>

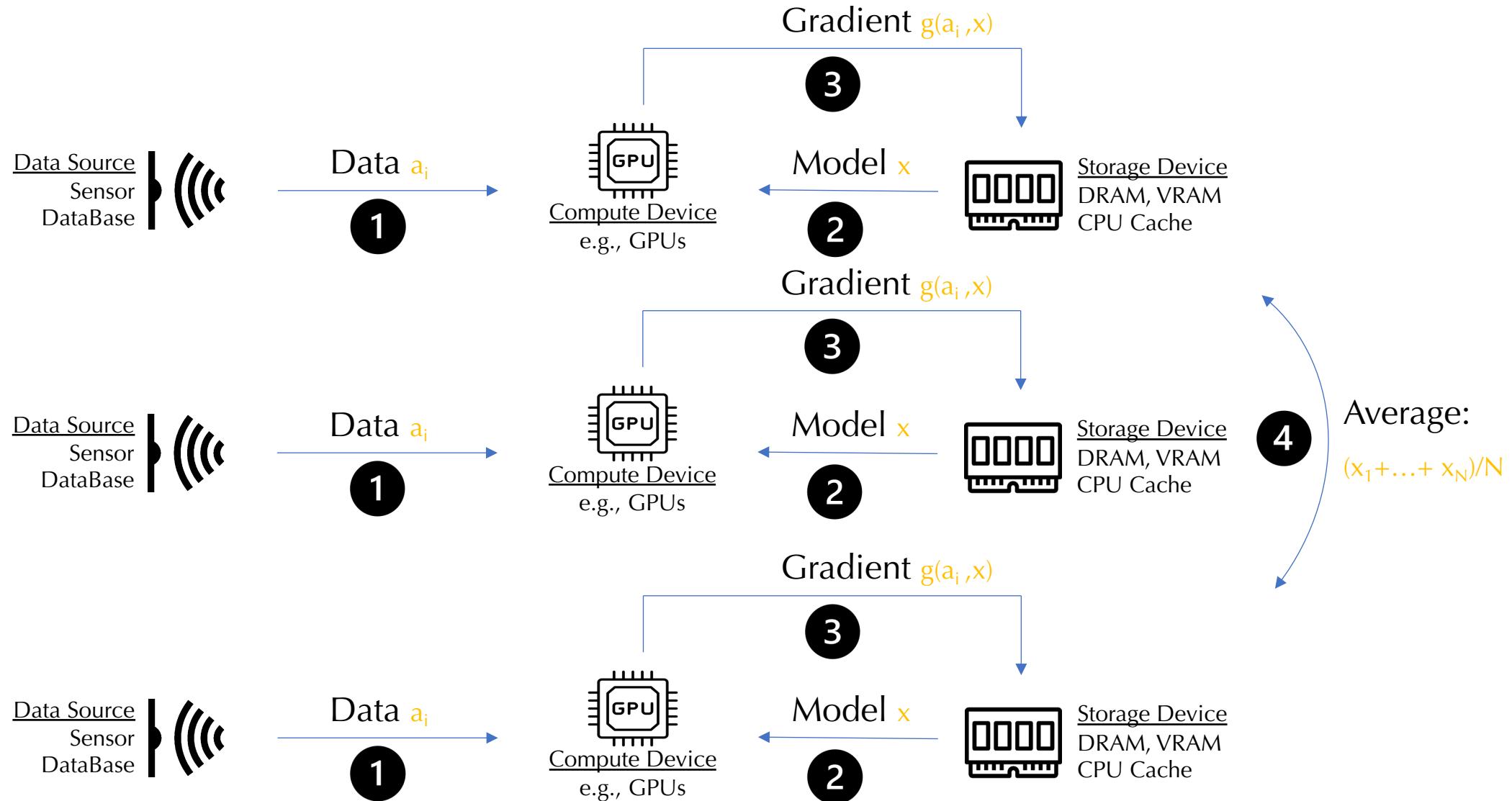


Nvidia Collective Communication Library

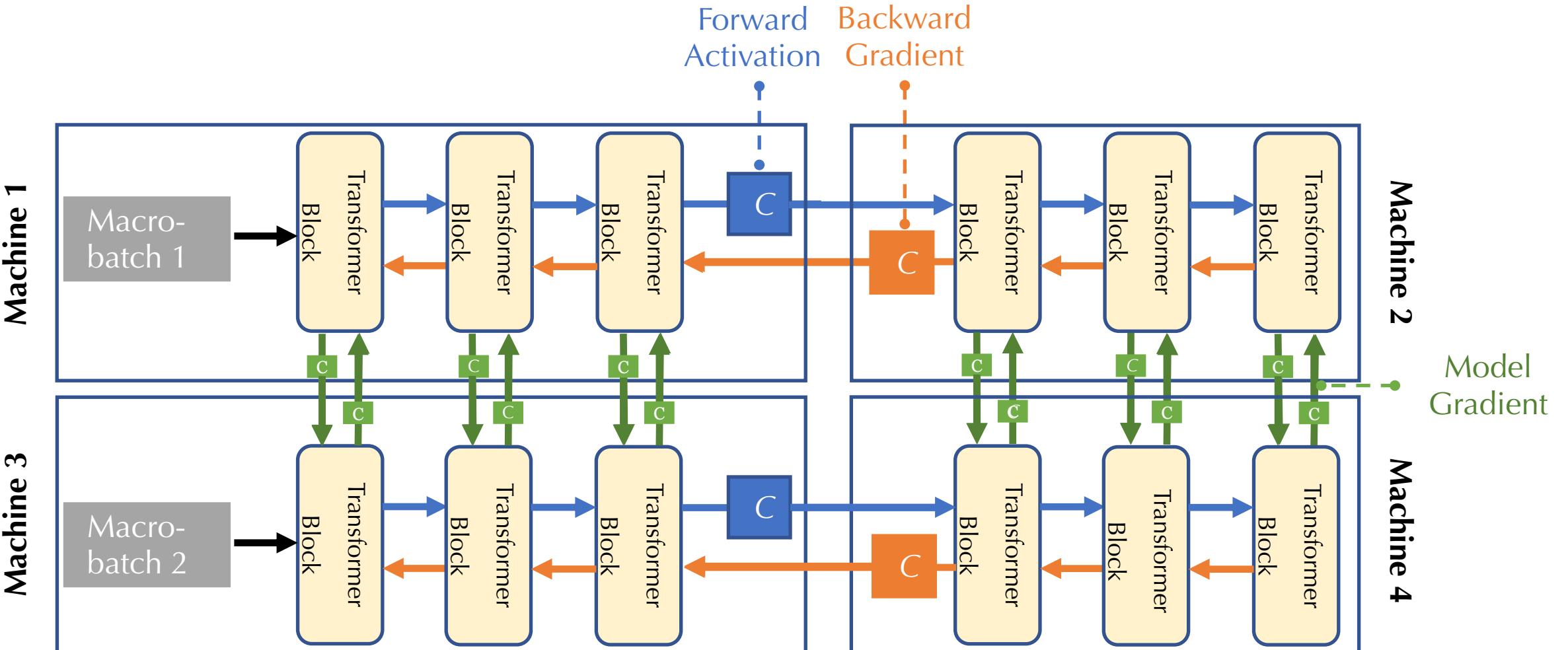




Data Parallelism

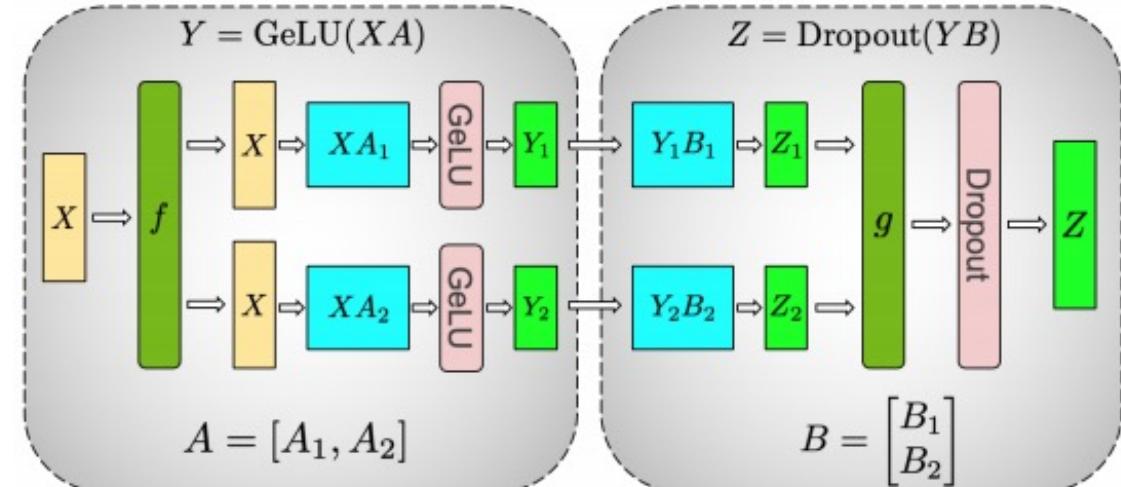
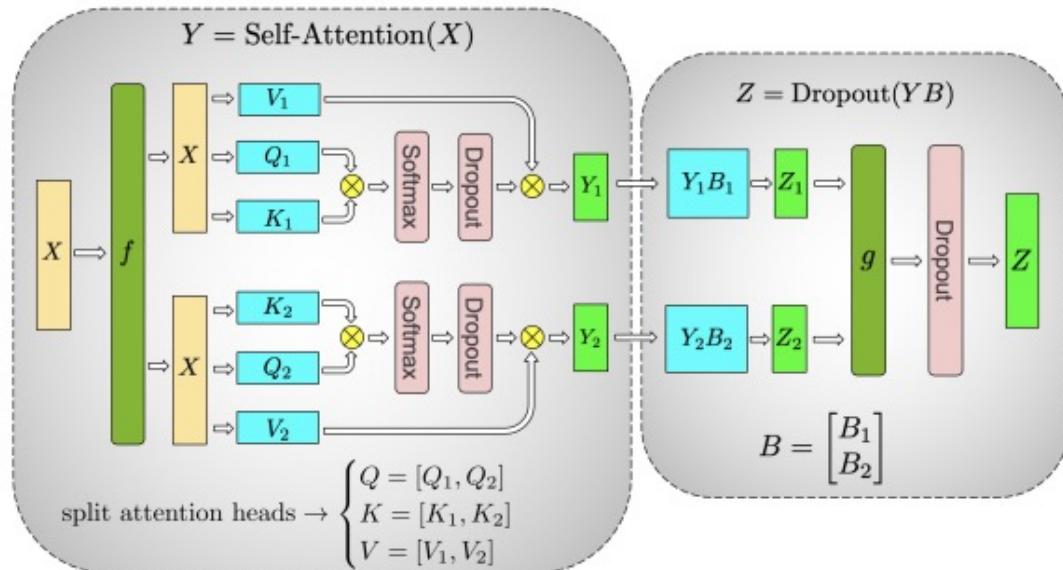


Pipeline Parallelism





Tensor Model Parallelism





Optimizer Parallelism

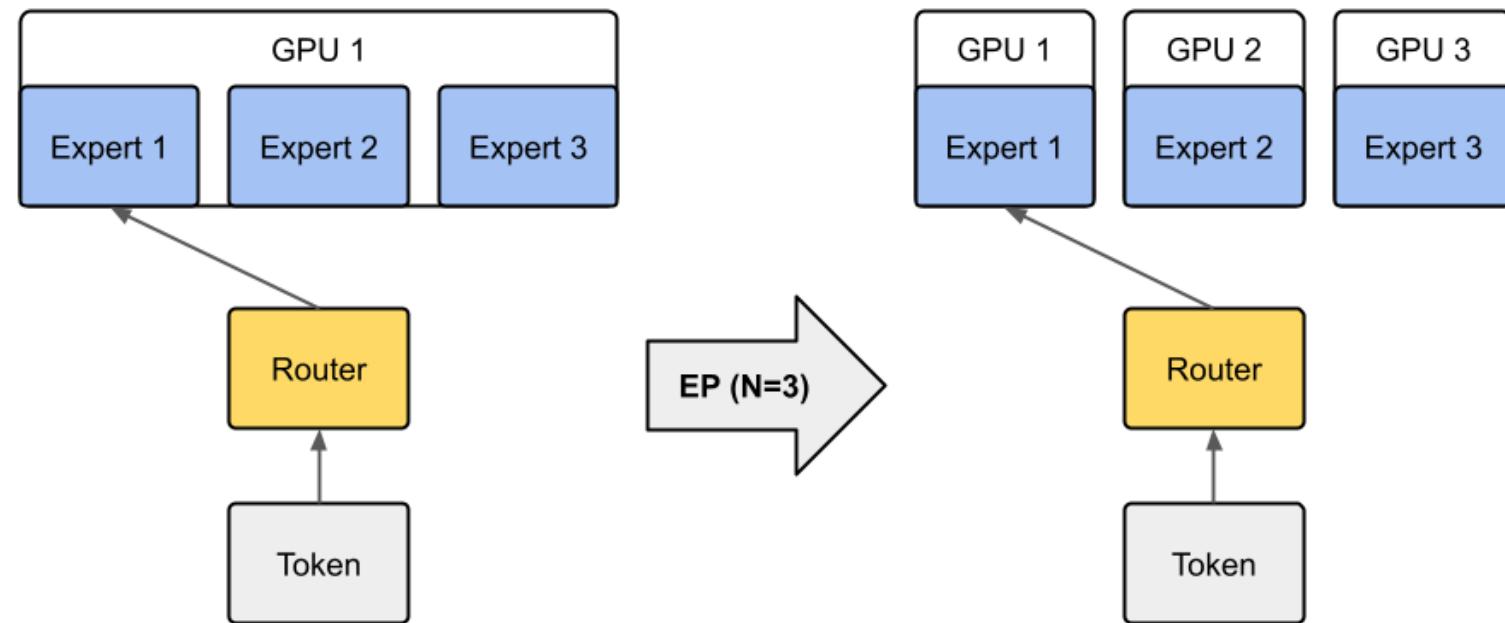
	gpu ₀	...	gpu _i	...	gpu _{N-1}	Memory Consumed
Baseline			$(2 + 2 + K) * \Psi$
P _{os}			$2\Psi + 2\Psi + \frac{K * \Psi}{N_d}$
P _{os+g}			$2\Psi + \frac{(2+K)*\Psi}{N_d}$
P _{os+g+p}			$\frac{(2+2+K)*\Psi}{N_d}$

Zero Redundancy Optimizer (ZeRO)

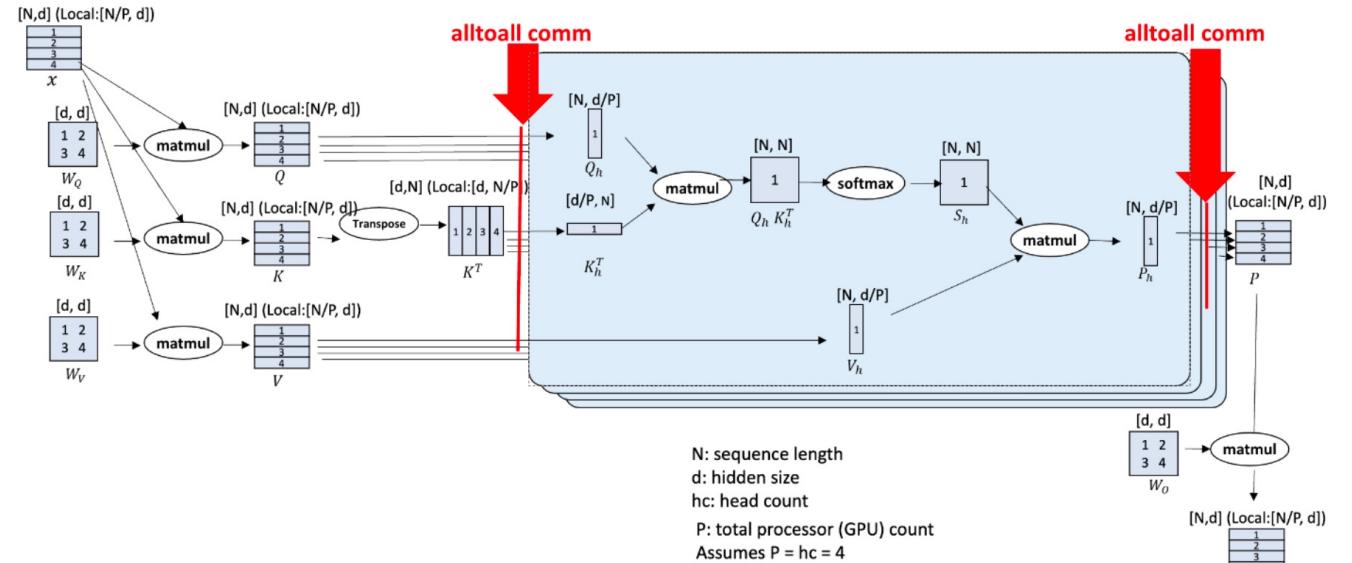
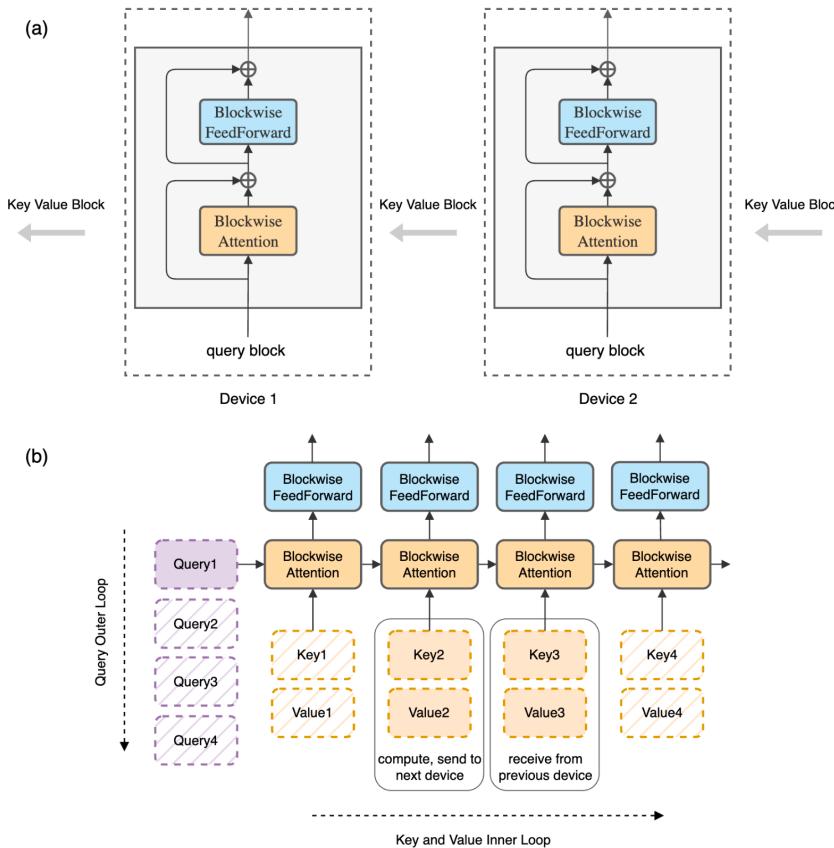


Expert Parallelism

Expert Parallelism applied on Mixture-of-Experts.

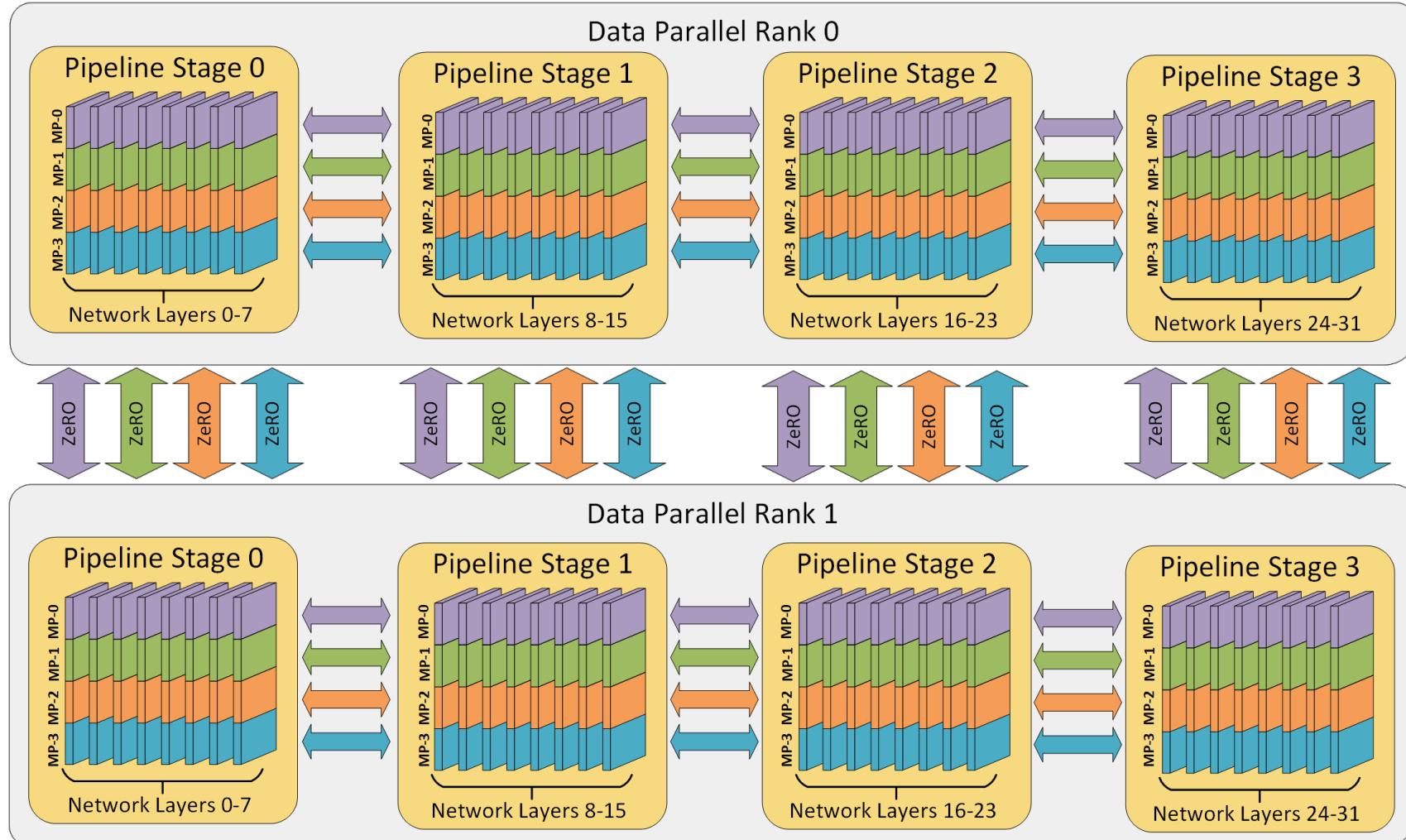


Long Sequence Parallelism





Data-, Pipeline-, Tensor Model-, Optimizer- Parallelisms



<https://www.deepspeed.ai/getting-started/>



Generative Inference & Hugging Face

```
from transformers import AutoTokenizer
import transformers
import torch

model = "meta-llama/Llama-2-7b-chat-hf"

tokenizer = AutoTokenizer.from_pretrained(model)
pipeline = transformers.pipeline(
    "text-generation",
    model=model,
    torch_dtype=torch.float16,
    device_map="auto",
)

sequences = pipeline(
    'I liked "Breaking Bad" and "Band of Brothers". Do you have any recommendations of other shows I',
    do_sample=True,
    top_k=10,
    num_return_sequences=1,
    eos_token_id=tokenizer.eos_token_id,
    max_length=200,
)
for seq in sequences:
    print(f"Result: {seq['generated_text']}")
```

The screenshot shows the Hugging Face Model Hub interface. At the top, there's a search bar with placeholder text "Search models, datasets, users...". Below the search bar, a navigation menu includes "Models", "Datasets", "Spaces", "Posts", "Docs", "Solutions", and "Pricing". A "Full-text search" button and a "Sort: Most downloads" dropdown are also present.

The main area displays a list of 473,734 models. Each model entry includes the owner's profile picture, the model name, its purpose (e.g., "Feature Extraction", "Text-to-Image"), the last update date, file size, and the number of downloads. Some entries also show the number of stars or reviews. The models are categorized into sections: Multimodal, Computer Vision, Natural Language Processing, Audio, Tabular, and Reinforcement Learning.

For example, the first few models listed are:

- pysentimiento/robertuito-sentiment-analysis (Updated Feb 25, 2023, 67.1M downloads)
- Supabase/gte-small (Updated Sep 21, 2023, 47.7M downloads)
- openai/clip-vit-large-patch14 (Zero-Shot Image Classification, Updated Sep 15, 2023, 38M downloads)
- roberta-base (Fill-Mask, Updated Mar 6, 2023, 22.8M downloads)
- gpt2 (Text Generation, Updated Jun 30, 2023, 20.8M downloads)

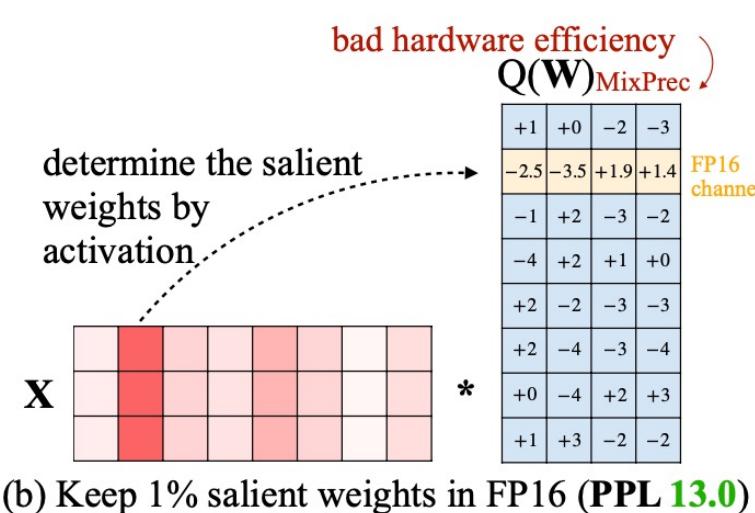


Generative Inference Algorithm Optimization

Quantized LLM Inference

\mathbf{W}_{FP16}	$\mathbf{Q}(\mathbf{W})_{\text{INT3}}$
+1.2 -0.2 -2.4 -3.4	+1 +0 -2 -3
-2.5 -3.5 +1.9 +1.4	-3 -4 +2 +1
-0.9 +1.6 -2.5 -1.9	-1 +2 -3 -2
-3.5 +1.5 +0.5 -0.1	-4 +2 +1 +0
+1.8 -1.6 -3.2 -3.4	+2 -2 -3 -3
+2.4 -3.5 -2.8 -3.9	+2 -4 -3 -4
+0.1 -3.8 +2.4 +3.4	+0 -4 +2 +3
+0.9 +3.3 -1.9 -2.3	+1 +3 -2 -2

(a) RTN quantization (**PPL 43.2**)



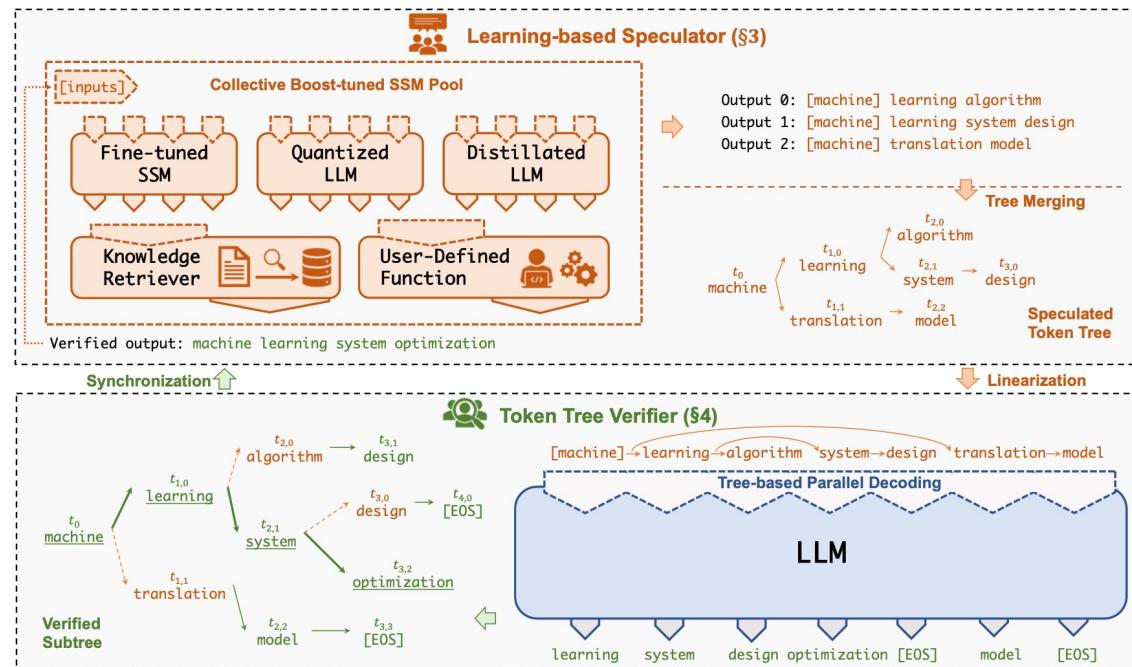
(b) Keep 1% salient weights in FP16 (**PPL 13.0**)

\mathbf{X}	$\mathbf{Q}(\mathbf{W})_{\text{INT3}}$
scale before quantize	a
average mag.	
*	

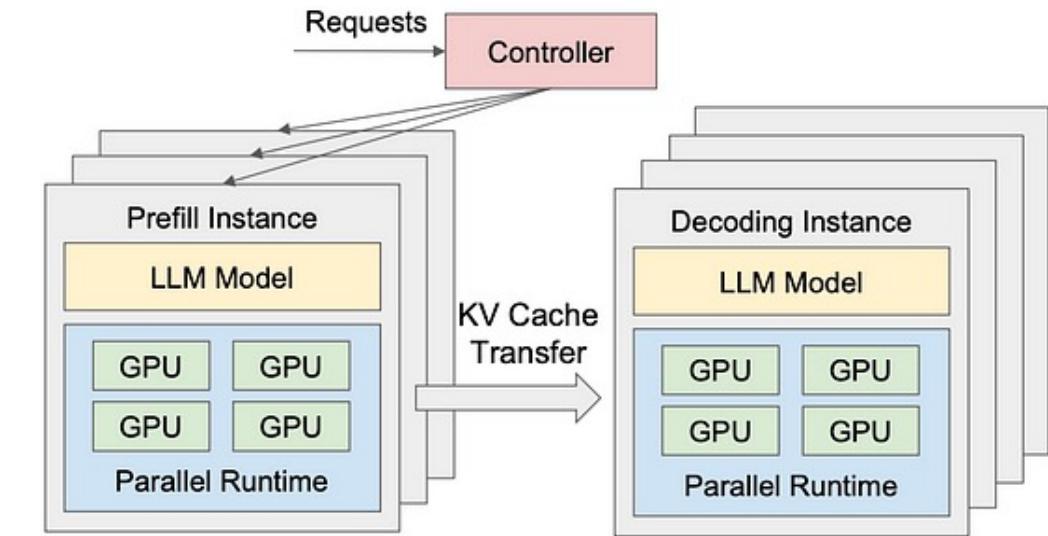
(c) Scale the weights before quantization (**PPL 13.0**)

<https://arxiv.org/pdf/2306.00978.pdf>

Generative Inference System Optimization



Speculative Decoding

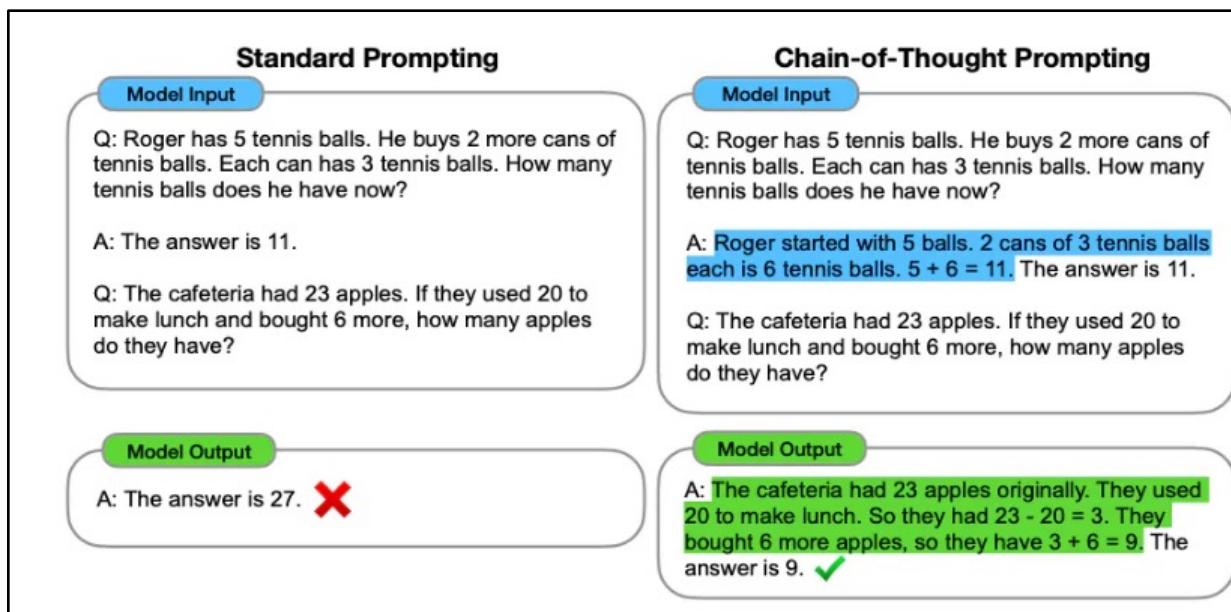


Prefill-decode Disaggregated Inference



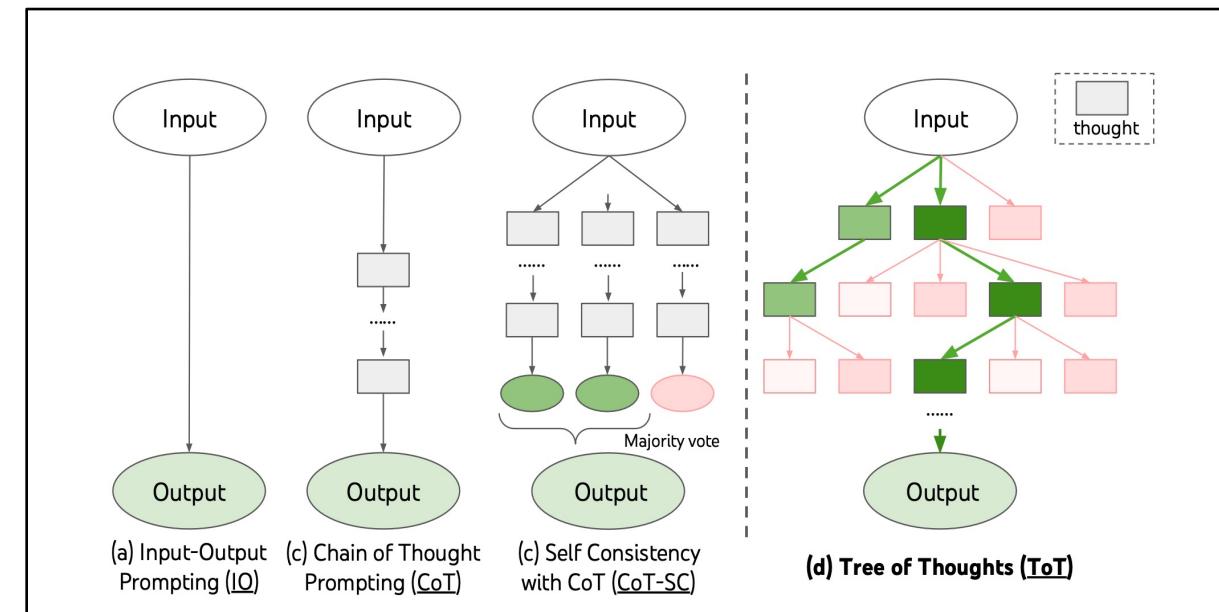
Prompt Engineering Overview & Practices

Chain of Thoughts



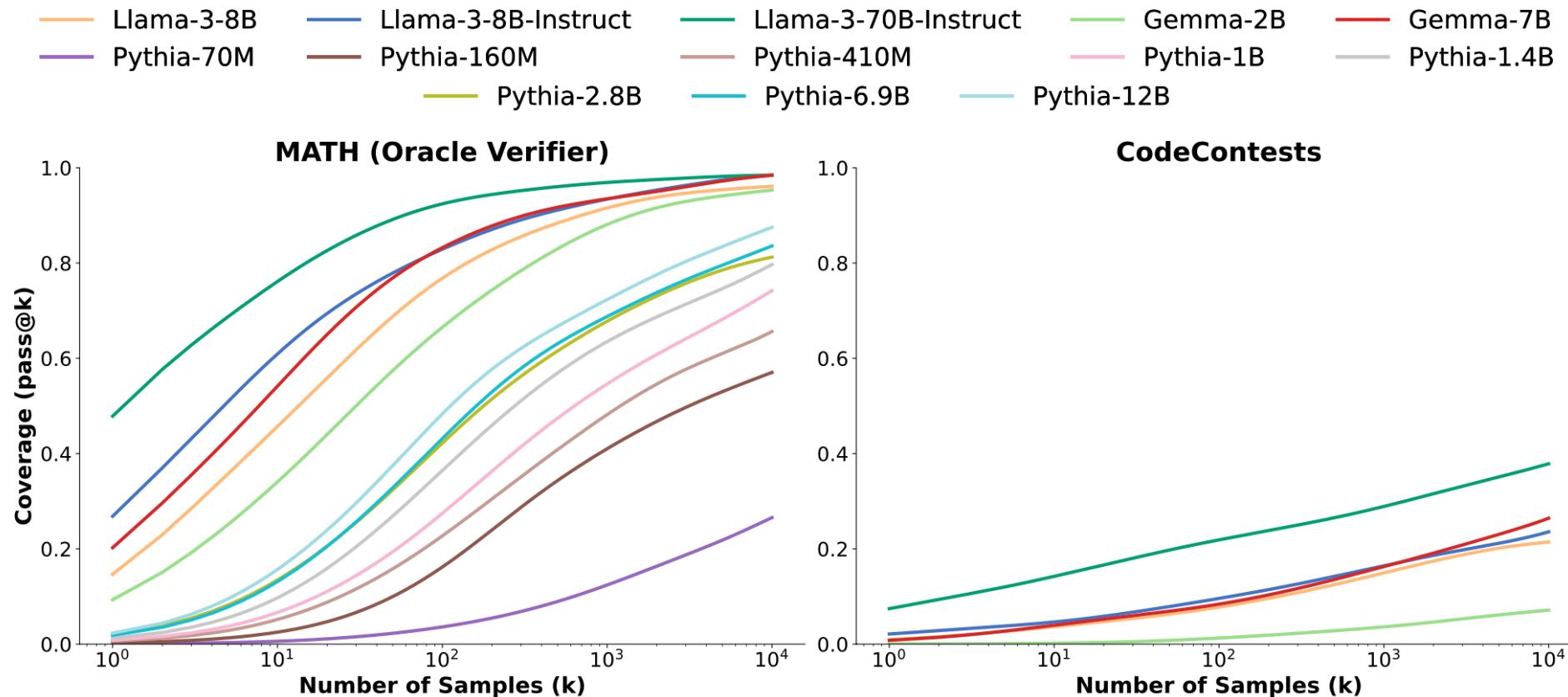
<https://arxiv.org/pdf/2201.11903.pdf>

Tree of Thoughts



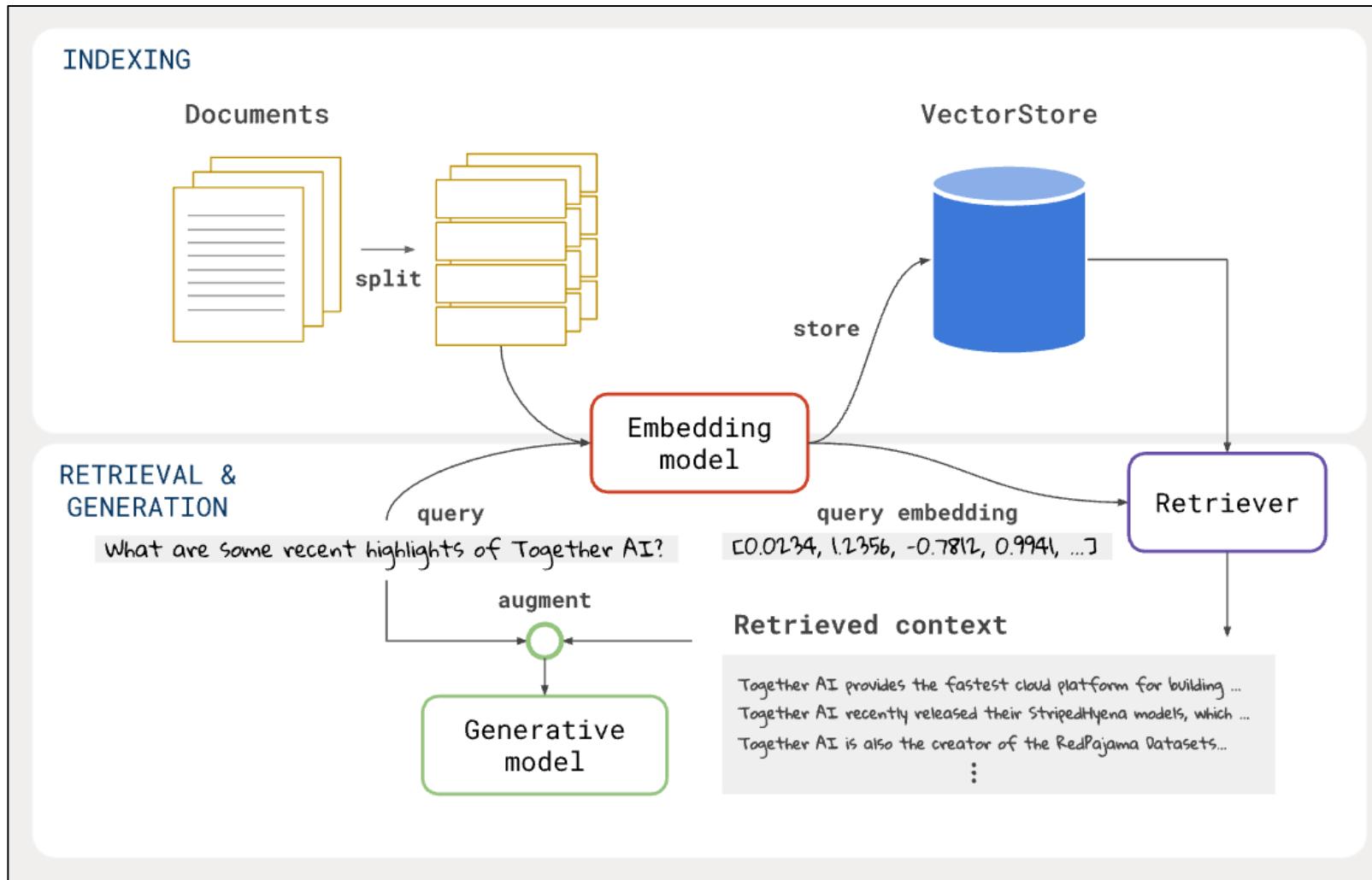
<https://arxiv.org/pdf/2305.10601.pdf>

Inference Time Scaling



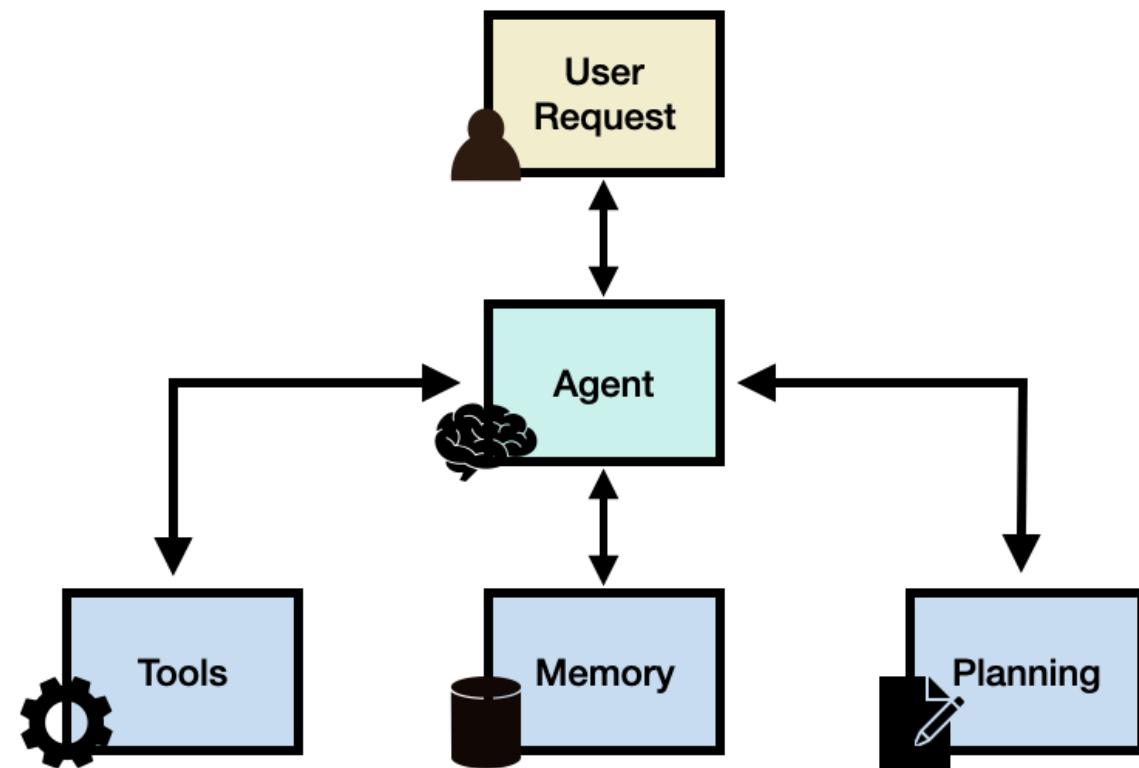


Retrieval Augmented Generation





LLM Agent

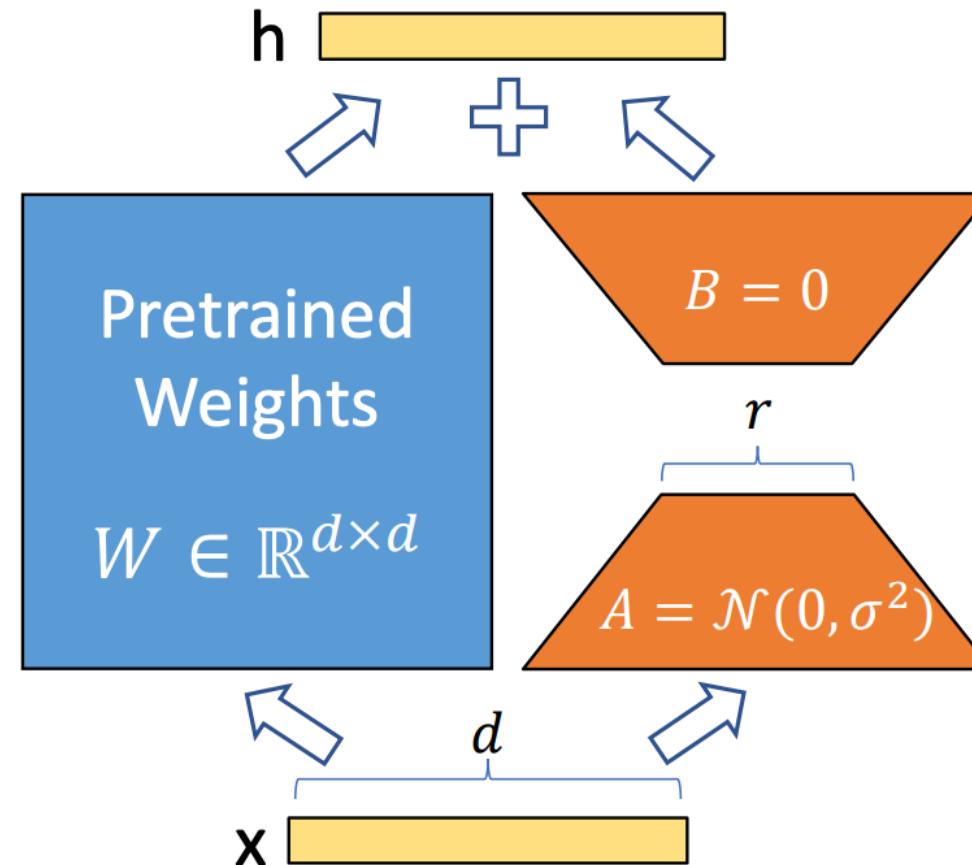




Parameter Efficient Fine-tuning (LoRA)

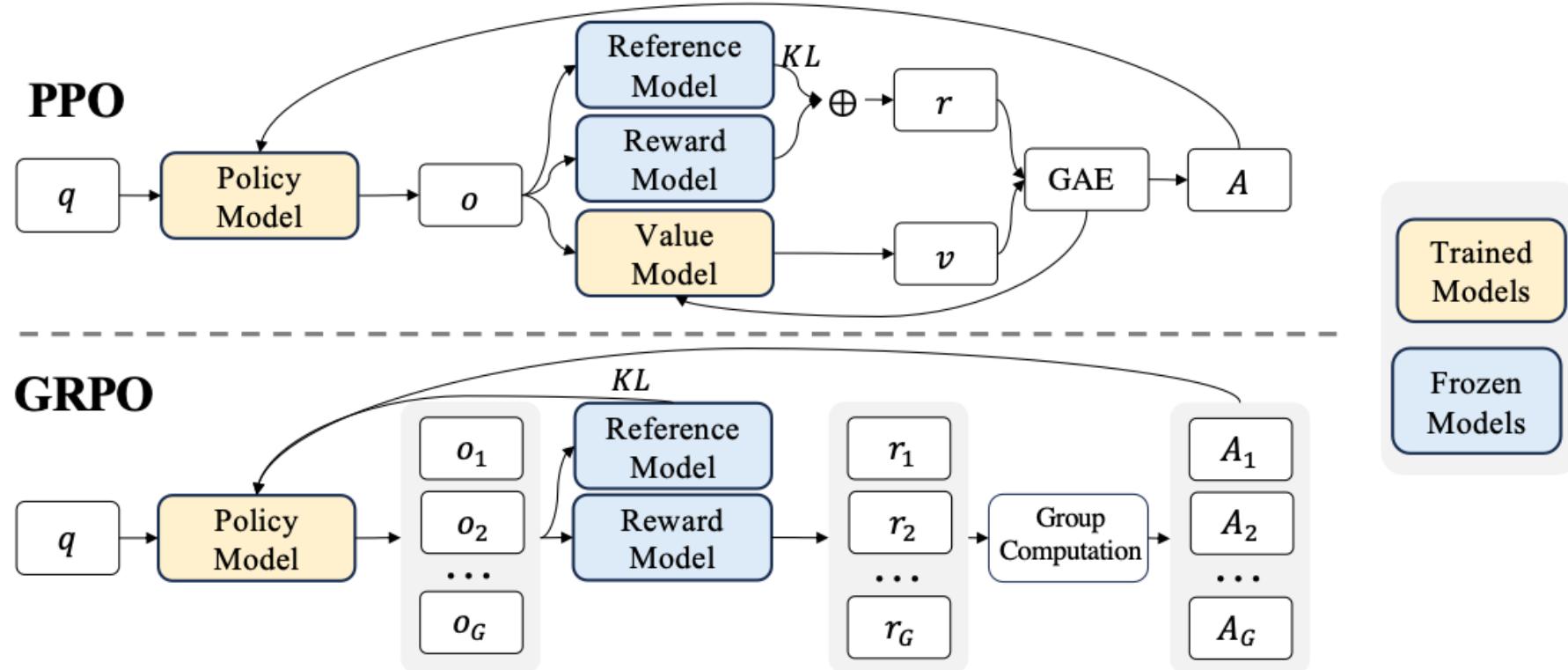
Low-Rank Adaptation (LoRA)

<https://arxiv.org/pdf/2106.09685.pdf>





RL Alignment





LLM Evaluation

Arena (battle) Arena (side-by-side) Direct Chat Leaderboard Arena Explorer About Us

Chatbot Arena LLM Leaderboard: Community-driven Evaluation for Best LLM and AI chatbots

[小红书](#) | [Twitter](#) | [Discord](#) | [Blog](#) | [GitHub](#) | [Paper](#) | [Dataset](#) | [Kaggle Competition](#)

Chatbot Arena is an open platform for crowdsourced AI benchmarking, developed by researchers at UC Berkeley [SkyLab](#) and [LMArena](#). With over 1,000,000 user votes, the platform ranks best LLM and AI chatbots using the Bradley-Terry model to generate live leaderboards. For technical details, check out our [paper](#).

Chatbot Arena thrives on community engagement — cast your vote to help improve AI evaluation!

New Launch! WebDev Arena: [web.lmarena.ai](#) - AI Battle to build the best website!

Language Overview Vision Text-to-Image Copilot Arena WebDev Arena Arena-Hard-Auto

Total #models: 197. Total #votes: 2,604,203. Last updated: 2025-02-02.

Code to recreate leaderboard tables and plots in this [notebook](#). You can contribute your vote at [lmarena.ai](#)!

Category	Overall	Apply filter	Style Control	Show Deprecated	Overall Questions	#models: 197 (100%)	#votes: 2,604,203 (100%)
Rank* (UB)	Rank (StyleCtrl)	Model	Arena Score	95% CI	Votes	Organization	License
1	2	Gemini-2.0-Flash-Thinking-Exp-01-21	1384	+6/-5	9649	Google	Proprietary
2	1	Gemini-Exp-1206	1373	+4/-3	23766	Google	Proprietary
3	1	ChatGPT-4o-latest-(2024-11-20)	1365	+3/-4	37760	OpenAI	Proprietary
3	1	DeepSeek-R1	1361	+7/-8	4195	DeepSeek	MIT
4	6	Gemini-2.0-Flash-Exp	1356	+4/-4	22591	Google	Proprietary
4	1	o1-2024-12-17	1352	+6/-7	11637	OpenAI	Proprietary
7	5	o1-preview	1335	+3/-4	33177	OpenAI	Proprietary
7	7	Owen-Max-2025-01-25	1332	+11/-11	2757	Alibaba	Proprietary
8	9	DeepSeek-V3	1316	+5/-4	16374	DeepSeek	DeepSeek
9	12	GLM-4-Plus-0111	1305	+12/-7	2584	Zhipu	Proprietary
10	13	o1-mini	1305	+3/-3	52364	OpenAI	Proprietary
10	13	Step-2-16K-Exp	1304	+7/-6	5126	StepFun	Proprietary
10	9	Gemini-1.5-Pro-002	1302	+3/-3	49232	Google	Proprietary



RELAXED
SYSTEM LAB

Guest Speech

TBD.



Logistics



Grading Policy

- 4 Homework:
 - $4 \times 5\% = 20\%$
- Mid-term Exam (30%):
 - In-class exam (80 Minutes): **2025/03/24**
- Final Exam (50%):
 - 2 hour
 - Time: **TBD**





Temporary Homework Schedule

- **Homework-1**
 - Release: 2025/02/13
 - DDL: 2025/02/22 23:59
- **Homework-2**
 - Release: 2025/03/04
 - DDL: 2025/03/12 23:59
- **Homework-3**
 - Release: 2025/03/25
 - DDL: 2025/04/10 23:59
- **Homework-4**
 - Release: 2025/04/22
 - DDL: 2025/05/06 23:59



Homework

- Delay penalty (20% each day after the DDL)
- The homework must be finished by yourself:
 - Discussion is welcomed;
 - Copy source code or any other content from other people is strictly forbidden;
 - Using generative AI is cool.





Exams



- Close-book;
- **One A4 page of cheating sheet allowed;**
- The mid-term exam covers the topics during the first half of the semester;
- The final exam covers all the topics.

No Attendance Requirement for the Lecture!

The Laboratory will be hosted with Announcements!



<https://github.com/Relaxed-System-Lab/HKUST-COMP4551-2026spring>

