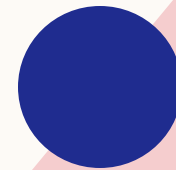# Introduction to training LLMs for AI agents

**LLM Agents MOOC series**

Yann Dubois |Sep. 15th 2025

All views are my own, and the methods based on information found online (esp Kimi, LlaMA, DeepSeek)
(nothing is related to OpenAI, unless explicitly stated)

# LLMs

- LLMs & chatbots took over the world



- How do train those?



All numbers are approximate from different open-source projects, especially LLaMA and DeepSeek and Kimi.

# General LLM training pipeline

Only for reasoning models

## Pretraining
Predict next word on internet

Data: > 10T tokens
Time: months
Compute cost: > $10M
Bottleneck: data & comp.

Eg LLaMA 3

## Reasoning RL
Think on questions with objective answers

Data: ~1M problems
Time: weeks
Compute cost: > $1M
Bottleneck: RL env & hacks

Eg DeepSeek R1

## Classic post-training / RLHF
Max user utility & prefs

Data: ~100k problems
Time: days
Compute cost: > $100K
Bottleneck: data & evals

Eg LLaMA-instruct

Still called post-training

All numbers are approximate from different open-source projects, especially LLaMA and DeepSeek.

# LLM training pipeline

- Architecture

- Training algorithm/loss

What people have been focusing on until 2023

- Data & RL env

- Evaluation

- Systems and infra to scale

What matters in practice

# LLM specializing pipeline

## Prompting

Art of asking the model what you want

Data: 0
Time: hours
Compute cost: 0
Bottleneck: evals

## Finetuning

Second stage of posttraining to domain specific data

Data: ~10k-100k problems
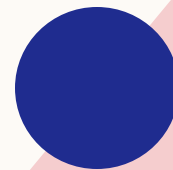Time: days
Compute cost: ~$10k-$100K
Bottleneck: data & evals

# Pretraining

**Method**

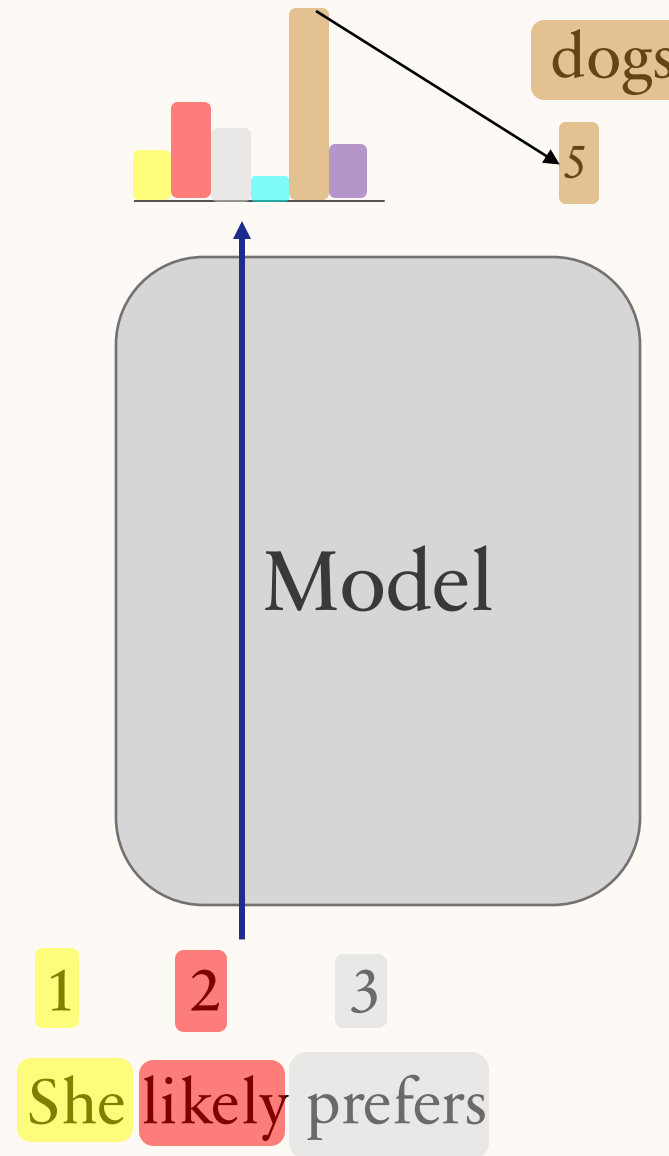Data

Compute

# Pretraining

- Goal: teach the model everything in the world

- Task: predict the next word

- Data: any reasonable data on internet

  - > 10T tokens (20-40T for llama 4, 15T for DSv3)

  - > 20B unique web pages

- Key since GPT-2 (2019)

# AR Language Models

- Task: predict the next word

- Steps:

1. tokenize

2. forward

3. predict probability of next token
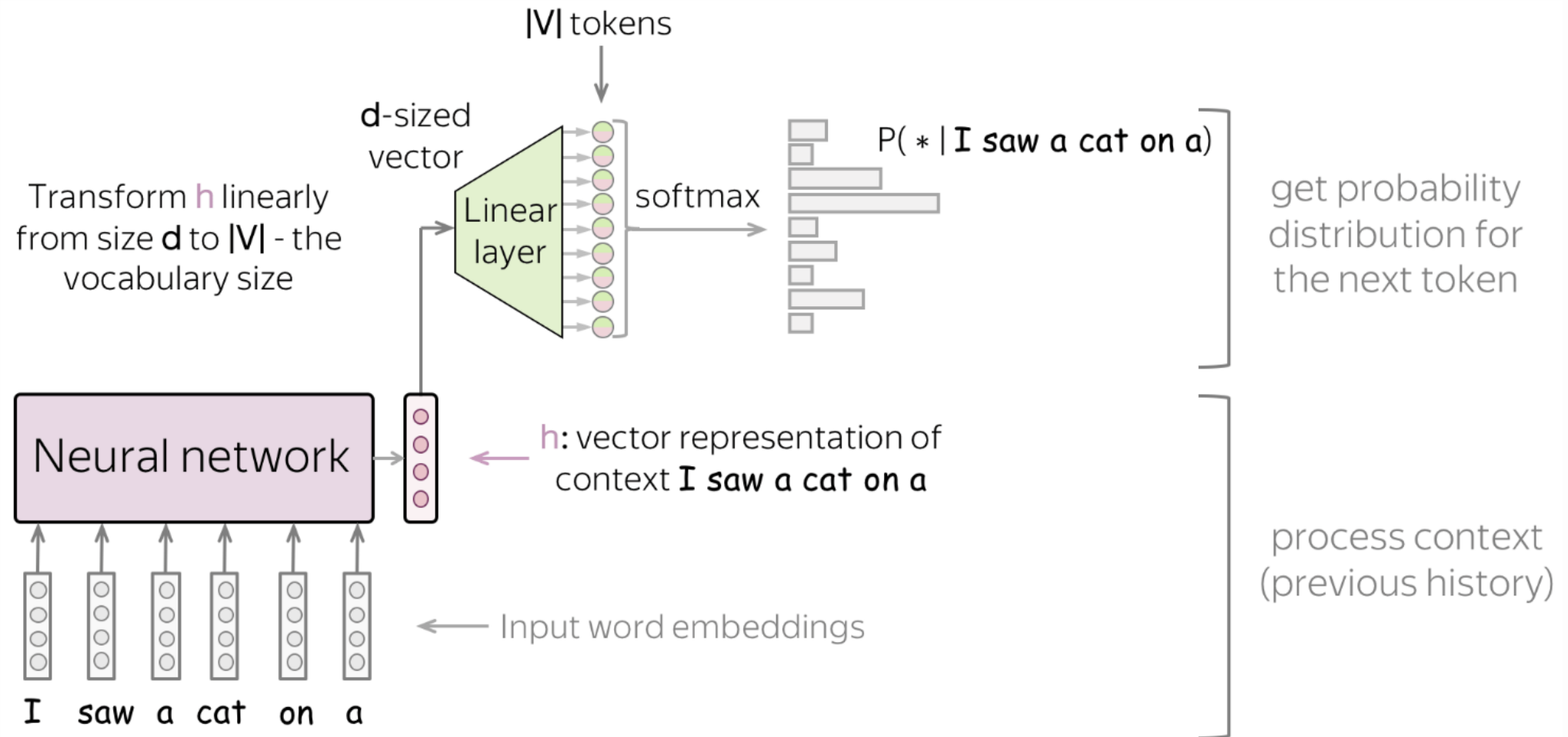
4. sample

5. detokenize

Inference only

Model

dogs

5

1    2    3

She  likely  prefers

# A Simple Language Model: N-grams

- How can you learn what to predict?

- Eg how can you know what comes after the grass is

- Idea: statistics!

  - Take all occurrences of the grass is on Wikipedia

  - Predicted probability for X is

    P(X| the grass is) = Count(X| the grass is)/Count(the grass is)

- Problem:

  - You need to keep count of all occurrences for each n-gram

  - Most sentences are unique: this can't generalize

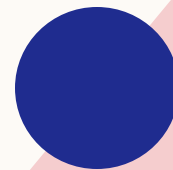- Solution: neural networks

# Neural Language Models



Transform **h** linearly from size **d** to **|V|** - the vocabulary size

**|V|** tokens

**d**-sized vector

Linear layer

softmax

$P( * \mid$ **I saw a cat on a**$)$

get probability distribution for the next token

Neural network

**h**: vector representation of context **I saw a cat on a**

process context (previous history)

Input word embeddings

I saw a cat on a

https://lena-voita.github.io/nlp_course/language_modeling.html#intro

# Pretraining

Method

**Data**

Compute

# Pretraining Data

- Idea: use ~~all of the web~~...

- Note: int...

  1. Dow... C file

  2. Text...

  3. Filte...

  4. Ded... ...ays same

  5. Heu... s, dirty toks)

  6. Mod...

  7. Dat... ing scaling

     laws...

- Also: lr ... t

```
<!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Transitional//EN" "http://www.w3.org/TR/xhtml1/DTD/xhtml1-transitional.dtd"> <html xmlns="http://www.w3.org/1999/xhtml"> <head> <meta http-equiv="Content-Type" content="text/html; charset=utf-8" /> <title>000 084 - Downloads Free 000 084 - Download 000 084 Software</title> <meta type="description" content="000 084 at Smart Code for free download. 000 084 freeware and shareware free downloads." /> <meta type="keywords" content="000 084, downloads, freeware, software, free, 000-084 Test Prep Training, Pass4sure IBM 000-084, TopCerts 000-084 Questions and Answers, HP0-084 Free Practice Exam Questions, Pass4sure ADOBE 9A0-084" /> <link rel="shortcut icon" href="/design/favicon.ico" type="image/x-icon" /> <link href="/design/look.css" rel="stylesheet" type="text/css" /> <!--[if lte ie 6]> <link href="/design/iexplorer6.css" rel="stylesheet" type="text/css" /> <![endif]--> <script type="text/javascript">link = "http://www.smartcode.com";</script> <script type="text/javascript" src="/design/bmark.js"></script> <link href="/design/generic/ui.css" rel="stylesheet" type="text/css" /> <script src="/cms/generic.2/scripts/jquery.js" type="text/javascript" language="javascript"></script> <script src="/cms/generic.2/scripts/ui.js" type="text/javascript" language="javascript"></script> </head> <body> <div id="wrapper"> <div id="header"> <a class="logo_type" href="http://www.smartcode.com/"><!--<img src="/design/sitelogo.jpg" alt="Premier website for Windows Shareware and Freeware Downloads">--></a> </div> <div id="menu"> <script type="text/javascript" src="/design/jcoding.js"></script> <form action="http://www.smartcode.com/downloads" id="search" name="search" method="get" onsubmit="return do_search(false);"> <input class="inp_text" type="text" name="query" id="query" value=""/> <a href="#" onclick="do_search (false);return false;" class="input_href"><img src="/design/search_button.gif" class="input_href" /></a> </form> <ul> <li><a href="http://www.smartcode.com/">Home</a></li> <li><a href="http://www.smartcode.com/db/allrootandsubcats.php">Categories</a></li> <li><a href="http://www.smartcode.com/db/new.php">New</a></li> <li><a href="http://www.smartcode.com/db/top.php">Popular</a></li> <li><a href="http://www.smartcode.com/submit/">Submit</a></li> <li><a href="http://www.smartcode.com/main/rss/">RSS</a></li> <li><a href="http://www.smartcode.com/main/contact.html">Contact</a></li> </ul> </div> <div id="content"> <div id="content_right"> <style> h1 { padding-bottom: 15px; } h1 strong { float: left; } div.pager { font-size: 11px; float: right; padding-top: 5px; } </style> <div style="display:none;"> false<br> 5<br> 5<br> 5 </div> <table border="0" cellspacing="0" cellpadding="0" id="table1"> <tr> <td width="4"> </td> <td> <script src="/design/ccoding.js" type="text/javascript"></script> </td> <td width="15"> </td> <td> <script src="/design/ccoding_im.js" type="text/javascript"></script> </td> </tr> </table> <br> <div class="hr" style="width:100%; margin-bottom:10px;"></div> <h1> <strong>000 084</strong> <div class="pager"> Pages:  1 <a href="/p2.html"> 2 </a> <a href="/p3 .html"> 3 </a>  <a href="http://000-084.smartcode.com/freeware.html">Freeware</a> <a href="http://000-084.smartcode.com/macosx .html">Mac</a> </div> </h1> <br> <p ct="1=Testkingworld.com is your ultimatel source for the System x High Performance Servers..."> Testkingworld.com is your ultimatel source for the System x High Performance Servers... <a class="details-link" href="http://000-084-test-prep-training.smartcode.com/info.html">Details</a></p><br> <p ct="1=One of the best and most rewarding features of the 000-084 training materials are that..."> One of the best and most rewarding features of the 000-084 training materials are that... <a class="details-link" href="http://pass4sure-ibm-000-084.smartcode.com/info.html">Details</a></p><br> <p ct="1=Download free 000-084 questions and answers. 000-084 exam questions are ultimate..."> Download free 000-084 questions and answers. 000-084 exam questions are ultimate... <a class="details-link" href="http://topcerts-000-084-questions-and-answers.smartcode.com/info .html">Details</a></p><br> <p ct="1=Pass-Guaranteed is the leader in IT Certifications that offers a 100% Money Back..."> Pass-Guaranteed is the leader in IT Certifications that offers a 100% Money Back... <a class="details-link"
```
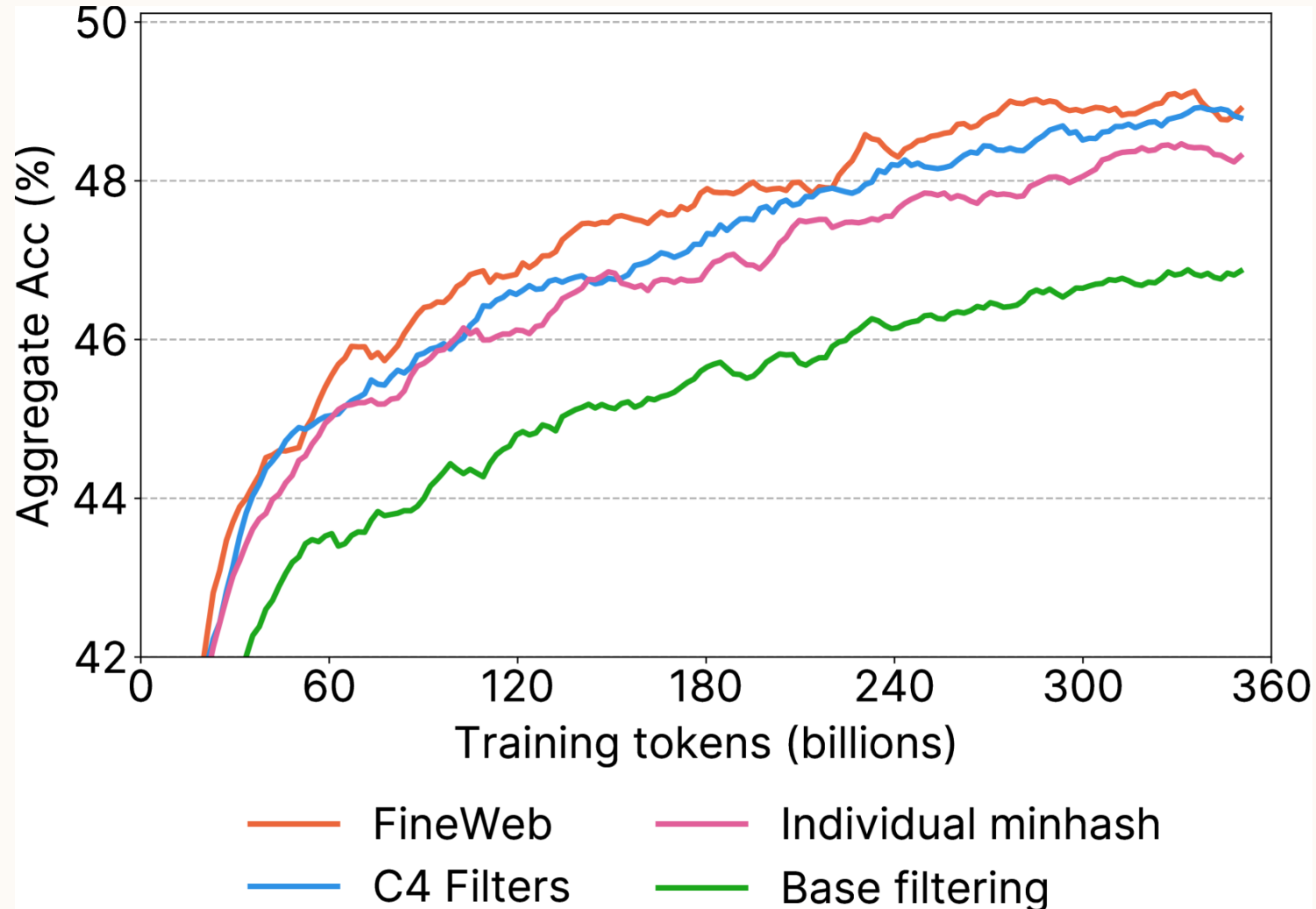
# Pretraining Data



Additional filters
20B -> 15B

JS, len, lorem ipsum, {
36B -> 20B

Dedup for <100 docs
36B -> 20B

NSFW blocklist, mostly
English, simple
document filtering
(repetition, length, etc)
200B -> 36B

The FineWeb Datasets: Decanting the Web for the Finest Text Data at Scale

# Midtraining data

- Continued pretraining to **adapt the model** to desired properties / higher quality data (<1T toks)

- Data mix changes shifts: eg more scientific, coding, multilingual data

- Longer context extension: bump (eg 4 → 128k for DSv3)

- Desired formatting/instruction following

- Higher quality data

- Reasoning data

- …

# Pre/mid training Data

- Collec...
- Lot of
  - He
  - He
- A lot o
  - C
- Comm
  - C
  - T
- Closed



**Composition of the Pile by Category**

■ Academic ■ Internet ■ Prose ■ Dialogue ■ Misc

PubMed Central · ArXiv · Pile-CC · Bibliotik · PG-19 · BC2 · StackExchange · Github · DM Math · Subtitles · FreeLaw · USPTO · PMA · Phil · NIH · OpenWebText2 · Wikipedia · IRC · EP · HN · YT

# Pretraining

Method

Data

**Compute**

# Pretraining compute

- Empirically: for any type of data and model, the most important is how much compute you spend on training (data & size)

- You can even predict performance with compute with scaling laws!



You can now do research at low scale and then predict how well it would hold at larger scale

Scaling laws
[Kaplan+ 2020]

# Scaling laws: tuning

- You have 10K GPUs for a month, what model do you train?

- Old pipeline:

  - Tune hyperparameters on big models (e.g. 30 models)

  - Pick the best => final model is trained for as much as each filtered out ones (e.g. 1 day)

- New pipeline:

  - Find scaling recipes (eg lr decrease with size)

  - Tune hyperparameters on small models of different sizes (e.g. for <3 days)

  - Extrapolate using scaling laws to larger ones

  - Train the final huge model (e.g. >27 days)

# Scaling laws for development

- Q: Should we use transformers or LSTM?



A: Transformers have a better constant and scaling rate (slope)

# Scaling laws: eg Chinchilla

- Q: How do we optimally allocate training* resources (size vs data)?



A: Use 20:1 tokens for each parameter (20:1)

*doesn't consider inference cost => in practice use larger (> 150:1)

Chinchilla
[Hoffmann+ 2022]

# Bitter lesson

- **Bitter lesson:** models improve with scale & Moore's Law

  => "only thing that matters in the long run is the leveraging of computation."

  Bitter [Sutton 2019]  http://www.incompleteideas.net/IncIdeas/BitterLesson.html

- Don't spend time over complicating: do the simple things and scale them!



nooooo you can't just scale up pure connectionist models on Internet data without inductive biases and modularization and expect them to learn real-world knowledge and grammar from form, or arithmetic and logical reasoning and causal inference—that's just memorization and superficial pattern-matching  like Eliza,you need grounding in real-world communication with intent and social dynamics and multimodal robotic embodiment which can foster disentangled learning from guided exploration and self-directed goals expressed in Bayesian programs and probabilistic graphical models which are interpretable and pin down a unique semantics which can be debiased and expressed with uncertainty, and learned efficiently on tiny academic budgets, the cost only shows how this is a dead-end, we need to stop chasing SOTAs and model the complexity of the brain and consider the social context to decolonize AI's structural biases for Third World researchers...

haha gpus go bitterrr

Image from Gwern

# Training a SOTA model

- Example of current SOTA: LLaMA 3 400B                    ~40 tok/param => train compute optimal

    Data: 15.6T tokens                              Parameters: 405B

- FLOPs: 6NP = 6 * 15.6e12 * 405e9 = 3.8 e25 FLOPs        ~2x less than executive order

- Compute: 16K H100 with average throughput of 400 TFLOPS

- Time: 3.8e25 / (400e12 * 3600) = 26M GPU hour / (16e3 * 24) = 70 days    From paper: ~30M h

- Cost: rented compute + salary=~$2/h*26Mh = ~$52M        $50-80M

- Carbon emitted= 26Mh*0.7kW*0.24kg/kWh = 4400 tCO2eq          ~2k return tickets JFK-LHR

- Next model? ~10x more FLOPs

# Pretraining
Predict next word on internet

Data: > 10T words
Time: months
Compute cost: > $10M
Bottleneck: data & comp.

Examples:

- DeepSeek v3

- LLaMA 4

# Post-training

Method

Data

Compute

# Language Modeling ≠ assisting users

- Problem: language modeling is not what we want

| PROMPT | *Explain the moon landing to a 6 year old in a few sentences.* |
|---|---|
| COMPLETION | GPT-3 |

```
Explain the theory of gravity to a 6 year old.

Explain the theory of relativity to a 6 year old in a few sentences.

Explain the big bang theory to a 6 year old.

Explain evolution to a 6 year old.
```

# Classic PT / alignment / IF

- Goal: steer the model to be useful on real-world tasks

- Task: maximize answer preferences of humans

- Data: 5k-500k problems

- Important since ChatGPT (2022)

# Reasoning

- Goal: teach the model to reason

- Task: answer correctly

- Data: any hard task with verifiable answer

- Important since o1(2024)

# Reasoning

- Goal: test-time scaling



o1 AIME accuracy
during training

o1 AIME accuracy
at test time

# Post-training

Methods:

SFT

RL

Data

Compute

# Eg task: "alignment"

- We want LLM follows user instructions and designer's desires (eg moderation)



- Background:

  - data of desired behaviors is what we want **but** scarce and expensive

  - pretraining data scales **but** is not what we want

- Idea: finetune pretrained LLM on a little desired data => "post-"training

# Supervised finetuning (SFT)

- Idea: finetune the LLM with language modeling of the desired answers

Next word prediction          "supervised"

- How do we collect the data? Ask humans

| Can you write a short introduction about the relevance of the term "monopsony" in economics? Please use examples related to potential monopsonies in the labour market and cite relevant research. | "Monopsony" refers to a market structure where there is only one buyer for a particular good or service. In economics, this term is particularly relevant in the labor market, where a monopsony employer has significant power over the wages and working conditions of their employees. […]. Overall, the concept of monopsony is essential to understanding the dynamics of labor markets and the impact of market power on workers. Further research is needed to understand the extent and impact of monopsonies on the economy and to develop policies to address this issue. References: Bivens, J., & Mishel, L. (2013). The Pay of Corporate Executives and Financial Professionals as Evidence of Rents in Top 1 Percent Incomes. Journal of Economic Perspectives, 27(3), 57-78. |
| My kid needs to do a science project for elementary school, what are some cheap but fun projects to do? | There are many cheap and fun science projects that elementary school kids can do! Here are a few ideas: 1. Make a homemade lava lamp: Fill a clear plastic bottle with water and vegetable oil, add a few drops of food coloring, and then drop in an Alka-Seltzer tablet. The reaction will cause bubbles to rise and fall like a lava lamp. […]7. Create a solar oven: Using a pizza box, aluminum foil, and plastic wrap, create a solar oven that can cook s'mores or other small snacks. 8. Explore magnets: Using magnets and a variety of objects, have your child experiment with magnetism to see what objects are attracted to magnets and which ones are not. These are just a few ideas to get you started! There are many more cheap and fun science projects that your child can do with household items and a little creativity. |

OpenAssistant
[Kopf+ 2023]

This was the ~key to GPT3 -> ChatGPT model!

# Scalable data for SFT: eg Alpaca

- Problem: human data is slow to collect and expensive

- Idea: use LLMs to scale data collection

| | |
|---|---|
| Give three tips for staying healthy. | 1.Eat a balanced diet and make sure to include plenty of fruits and vegetables. 2. Exercise regularly to keep your body active and strong. 3. Get enough sleep and maintain a consistent sleep schedule. |
| What does "algorithm" mean? | An algorithm is a step-by-step set of instructions used to solve a problem or achieve a goal. It is a process of breaking down complex tasks into smaller, more manageable steps. Algorithms can be used to find the shortest path from one place to another, sort data, or render data from an immense database in a few seconds. |
| Find the average number in a list. | def avg_list(nums): total = 0 for n in nums: total += n return total / len(nums) print(avg_list([4, 3, 6, 10, 8])) # Output 6.4 |

Alpaca
[Taori+ 2023]

Started for academic replication of ChatGPT but "synthetic data generation" is now hot topic!

# Scalable data for SFT: eg DeepSeek R1

- Problem: LLM-generated data ~assumes that you have a smarter LLM

- Idea: use rejection sampling based on verifiers

  DeepSeek-R1
  [DeepSeek-AI 2025]

  1. Temporary LLM generates many answers

  2. Keep answer if it's correct (eg, passes test case), or preferred over others



Image from Jay Alammar

# What to learn during SFT

- Instruction following

- Desired format or style

- Tool use [eg Kimi 2 or xLAM]

- Early reasoning [eg DeepSeek R1]

- … anything where you can get good input, output pairs…


- SFT is either seen as a final stage or as a preparation for RL

# SFT pipelines can be complex: eg Kimi k2

- SFT pipelines can use complex system to build useful SFT data

- Eg Kimi K2 uses LLM simulated user & tools, and rubric based rejection sampling to build data for agentic tool use



(a) Synthesizing tool specs, agents and tasks

(b) Generating agent trajectories

Kimi K2
[Kimi team 2025]

# Scalable data for SFT: quantity?

- You learn quickly from SFT! ~ 10k are sufficient for learning style & IF



LIMA
[Zhou+ 2023]

- Just learns the format of desired answers (length, bullet points, …)

  - The knowledge is already in the pretrained LLM => Specializes to one "type of user"

- DeepSeek-R1 uses ~800K examples for learning reasoning using R1-zero

# Post-training

**Methods:**

  SFT

  **RL**

Data

Compute

# Reinforcement learning

- Problem: SFT is **behavior cloning** of humans

  1.  **Bound by human abilities**: humans may prefer things that they are not able to generate

  2.  **Hallucination**: cloning *correct* answer teaches LLM to hallucinate if it didn't know about it!

  Can you write a short introduction about the relevance of the term "monopsony" in economics? Please use examples related to potential monopsonies in the labour market and cite relevant research.

  "Monopsony" refers to a market structure where there is only one buyer for a particular good or service. In economics, this term is particularly relevant in the labor market, where a monopsony employer has significant power over the wages and working conditions of their employees. […]. Overall, the concept of monopsony is essential to understanding the dynamics of labor markets and the impact of market power on workers. Further research is needed to understand the extent and impact of monopsonies on the economy and to develop policies to address this issue. **References: Bivens, J., & Mishel, L. (2013). The Pay of Corporate Executives and Financial Professionals as Evidence of Rents in Top 1 Percent Incomes. Journal of Economic Perspectives, 27(3), 57-78.**

  If LLM doesn't know [Bivens 2013] => teaches the model to make up plausibly sounding references

  3.  **Price:** collecting ideal answers can be expensive

# Reinforcement learning

- Idea: maximize desired behavior rather than clone it

- Read: DeepSeek R1 & Kimi K2

Kimi K2
[Kimi team 2025]

DeepSeek-R1
[DeepSeek-AI 2025]

- Key: where the reward comes from:

  - Rule-based rewards (eg string match for close-ended QA, or passing coding function)

  - Reward model trained to predict human preferences (RLHF)

  - LLM as a judge

  - …



Image from Jay Alammar

# Reinforcement learning: DeepSeek R1

# Reinforcement learning: GRPO

DeepSeek R1 optimizes GRPO which uses MC estimates for computing the advantage

Similar loss for Kimi K1.5 & K2



DeepSeek-R1
[DeepSeek-AI 2025]

# RL: Infra is key

- Sampling is an important bottleneck since you sample multiple outputs per problem

- Infra is key, especially for agents:

  - Long rollouts: Kimi pauses long tail rollouts

  - Environment feedback can be slow: Kimi uses concurrent rollouts & dedicated services for envs



Engines are collocated to avoid communication overhead

Kimi K2
[Kimi team 2025]

# RL from Human Feedback (RLHF)

- Idea: maximize human preference rather than clone their behavior

- Made ChatGPT

- Pipeline:

  1. For each instruction: generate 2 answers from a pretty good model (SFT)

  2. Ask labelers to select their preferred answers

  3. Finetune the model to generate more preferred answers (PPO or DPO)

**How?**

Instruction

# RLHF: gains



| Method | Simulated Win-rate (%) |
| --- | --- |
| GPT-4*[†] | $79.0 \pm 1.4$ |
| ChatGPT*[†] | $61.4 \pm 1.7$ |
| PPO | $46.8 \pm 1.8$ |
| DPO | $46.8 \pm 1.7$ |
| Best-of-1024 | $45.0 \pm 1.7$ |
| Expert Iteration | $41.9 \pm 1.7$ |
| SFT 52k | $39.2 \pm 1.7$ |
| SFT 10k | $36.7 \pm 1.7$ |
| Binary FeedME | $36.6 \pm 1.7$ |
| Quark | $35.6 \pm 1.7$ |
| Binary Reward Conditioning | $32.4 \pm 1.6$ |
| Davinci001* | $24.4 \pm 1.5$ |
| LLaMA 7B* | $11.3 \pm 1.1$ |

Learn to summarize
[Stiennon+ 2020]

AlpacaFarm
[Dubois+ 2023]

# RLHF: human data

- Data: human crowdsourcing



guidelines



example

# RLHF: challenges of human data

- Slow & expensive

- Hard to focus on correctness rather than form (eg length)

- Annotator distribution shifts its behavior

- Crowdsourcing ethics

LLM Opinion
[Santurkar+ 202...

Pretrain    Posttrain

Long way to go
[Singhal+ 2024]

Question: *Why don't adults roll off the bed?*

☆ **SFT (Before); 59 tokens**

*Adults typically do not roll off of the bed because they have developed the muscle memory to ... d maintaining proper posture.*

**h longer / more details**

developed muscle control
eping. Additionally, most
in their sleep, so rolling off
also be safety concerns
ne were to lose balance or

TECHNOLOGY

# AMERICA ALREADY HAS AN AI UNDERCLASS

Search engines, ChatGPT, and other AI tools wouldn't function without an army of contractors. Now those workers say they're underpaid and mistreated.
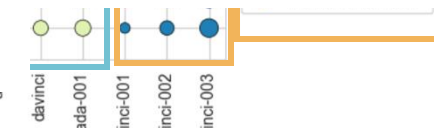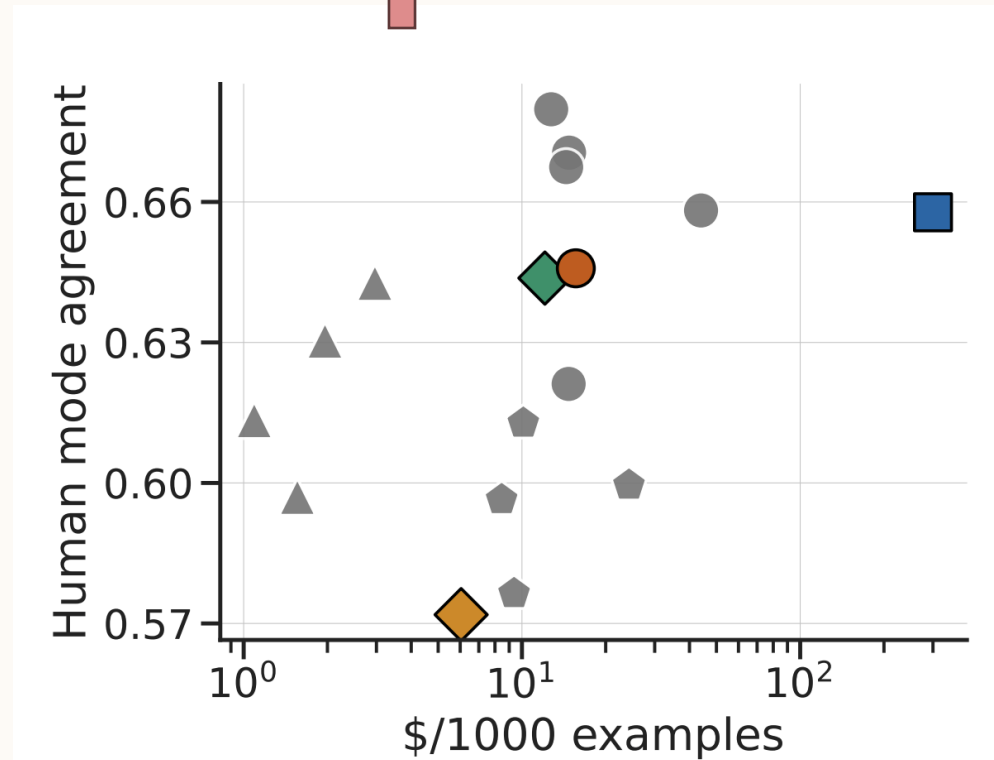
By Matteo Wong

# RLHF: LLM data

- Idea: replace human preferences with LLM preferences



Works surprisingly well!

AlpacaFarm
[Dubois+ 2023]

# Evaluation

**Close-ended**

Open-ended

# Importance of evaluation in AI

Quantify progress towards desired task to:

Identify improvements

Select models

Decide if production ready

# Close-ended Evaluation

- Idea: few possible answers =>automate verification

- Eg MMLU



**Astronomy**

**What is true for a type-Ia supernova?**
    A. This type occurs in binary systems.
    B. This type occurs in young galaxies.
    C. This type produces gamma-ray bursts.
    D. This type produces high amounts of X-rays.
    Answer: A

**High School Biology**

**In a population of giraffes, an environmental change occurs that favors individuals that are tallest. As a result, more of the taller individuals are able to obtain nutrients and survive to pass along their genetic information. This is an example of**
    A. directional selection.
    B. stabilizing selection.
    C. sexual selection.
    D. disruptive selection
    Answer: A

MMLU
[Hendrycks+ 2020]

# Evaluation: challenges

- Sensitivity to prompting/inconsistencies

| | MMLU (HELM) | MMLU (Harness) | MMLU (Original) |
|---|---|---|---|
| llama-65b | **0.637** | 0.488 | **0.636** |
| tiiuae/falcon-40b | 0.571 | **0.527** | 0.558 |
| llama-30b | 0.583 | 0.457 | 0.584 |
| EleutherAI/gpt-neox-20b | 0.256 | 0.333 | 0.262 |
| llama-13b | 0.471 | 0.377 | 0.47 |
| llama-7b | 0.339 | 0.342 | 0.351 |
| tiiuae/falcon-7b | 0.278 | 0.35 | 0.254 |

# Evaluation: challenges

- Sensitivity to prompting/inconsistencies

- Train & test contamination (~not important for development)

# Evaluation

Close-ended

**Open-ended**

# Evaluation: aligned LLM

- How do we evaluate something like ChatGPT?

- Challenges:

  - Large diversity

  - Open-ended tasks => hard to automate

- Idea: ask for annotator preference between answers

Table 1: Distribution of use case categories from our API prompt dataset.

| Use-case | (%) |
| --- | --- |
| Generation | 45.6% |
| Open QA | 12.4% |
| Brainstorming | 11.2% |
| Chat | 8.4% |
| Rewrite | 6.6% |
| Summarization | 4.2% |
| Classification | 3.5% |
| Other | 3.5% |
| Closed QA | 2.6% |
| Extract | 1.9% |

InstructGPT
[Ouyang+ 2022]

# Human evaluation: eg ChatBot Arena

- Idea: have users interact (blinded) with two chatbots, rate which is better.



- Problem: cost & speed!

ChatBot Arena
[Chiang+ 2024]

# LLM evaluation: eg AlpacaEval

- Idea: use LLM instead of human

- Steps:

  - For each instruction: generate output by baseline and model to eval

  - Ask GPT-4 which output is better

  - Average win-probability => win rate

- Benefits:

  - 98% correlation with ChatBot Arena

  - < 3 min and < $10

- Challenge: **spurious correlation**

Chat Arena Spearman correlation

| 0.35 | 0.51 | 0.59 | 0.63 | 0.66 | 0.69 | 0.83 | 0.87 | 0.94 | 0.98 |
|------|------|------|------|------|------|------|------|------|------|

Output Length · TruthfulQA · HellaSwag · GSM-8K · Open LLM · WinoGrande · ARC-C · MMLU · MT-bench · LC AlpacaEval 2.0

AlpacaEval
[Li+ 2023]

# Systems

# Systems

- Problem: everyone is bottlenecked by compute!

- Why not buy more GPUs?

  - GPUs are expensive and scarce!

  - Physical limitations (eg communication between GPUs)

- => importance of resource allocation (scaling laws) and optimized pipelines



haha gpus go bitterrr

# Systems 101: GPUs

- Massively parallel: same instruction applied on all thread but different inputs.

    => Optimized for throughput!



SM
Streaming
Multiprocessors

# Systems 101: GPUs

- Massively parallel

- Fast matrix multiplication: special cores >10x faster than other floating point ops



Matmul vs. non-matmul FLOPS across GPUs

# Systems 101: GPUs

- Massively parallel

- Fast matrix multiplication

- Compute > memory & communication:

  - Hard to keep processors fed with data



Scaling of Peak hardware FLOPS, and Memory/Interconnect Bandwidth

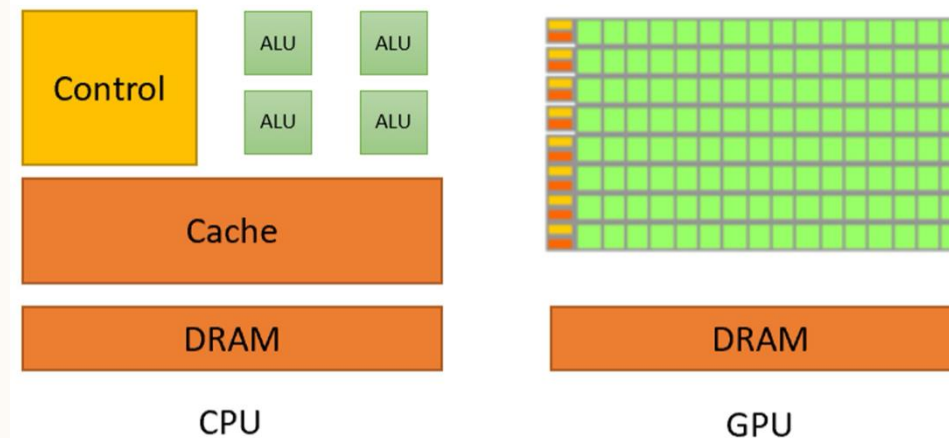HW FLOPS: 60000x / 20 yrs (3.0x/2yrs)
DRAM BW: 100x / 20 yrs (1.6x/2yrs)
Interconnect BW: 30x / 20 yrs (1.4x/2yrs)

BERT transformer

Table 1. Proportions for operator classes in PyTorch.

|  | Operator class | % flop | % Runtime |
|---|---|---|---|
| Matmul | △ Tensor contraction | 99.80 | 61.0 |
|  | □ Stat. normalization | 0.17 | 25.5 |
| Activation | ○ Element-wise | 0.03 | 13.5 |

DataMovement
[Ivanov+ 2020]

# Systems 101: GPUs

- Massively parallel

- Fast matrix multiplication

- Compute > memory & communication

- Memory hierarchy:

  - Closer to cores => faster but less memory

  - Further from cores => more memory but slower

**TABLE IV**
THE MEMORY ACCESSES LATENCIES

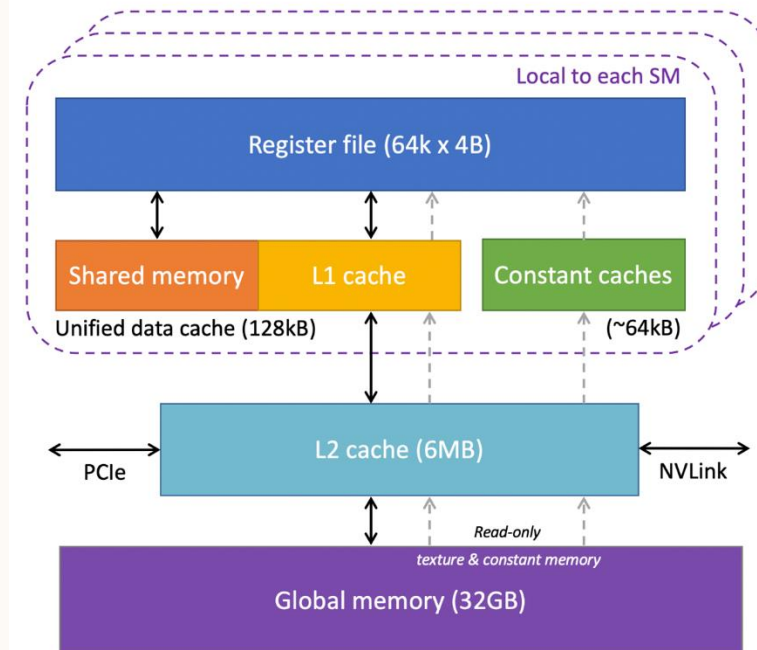| Memory type | CPI (cycles) |
|---|---|
| Global memory | 290 |
| L2 cache | 200 |
| L1 cache | 33 |
| Shared Memory (ld/st) | (23/19) |

# Systems 101: GPUs

- Massively parallel

- Fast matrix multiplication

- Compute > memory & communication

- Memory hierarchy


- Metric: **Model Flop Utilization (MFU)**

  - Ratio: observed throughput / theoretical best for that GPU

  - 50% is great!

# Systems: low precision

- Fewer bits => faster communication & lower memory consumption

- For deep learning: decimal precision ~doesn't matter except exp & updates

  - Matrix multiplications can use bf16 instead of fp32

- For training: **Automatic Mixed Precision** (**AMP**)

  - Weights stored in fp32, but before computation convert to bf16

  - Activation in bf16 => main memory gains

  - (Only) matrix multiplication in bf16 => speed gains

  - Gradients in bf16 => memory gains

  - Master weights updated fp32 => full precision

# Systems: operator fusion

- Problem:
  - communication is slow

  - every new PyTorch line moves variables to global memory

- Idea: communicate once

- torch.compile

```
x1 = x.cos() # Read from x in global memory, write to x1
x2 = x1.cos() # Read from x1 in global memory, write to x2
```



Naïve (non-fused)          Fused kernel

# Systems: tiling

- Idea: group and order threads to minimize global memory access (slow)

- Eg matrix multiplication

- Compute matrix multiplications in subphases to reuse memory

  1. Load M_00 and N_00 tiles into SM

  2. Compute partial sums for P

  3. Load M_00 and N_20 into SM

  4. …

- => reuse reads (~cache)

- T reduction of global reads

E.g. assume that thread can
Then have to reread a

# Systems: eg FlashAttention

- Idea: kernel fusion, tiling, recomputation for attention!

- 1.7x end to end speed up!



FlashAttention
[Dao+ 2022]

| Model implementations | OpenWebText (ppl) | Training time (speedup) |
|---|---|---|
| GPT-2 small - Huggingface [87] | 18.2 | 9.5 days (1.0×) |
| GPT-2 small - Megatron-LM [77] | 18.2 | 4.7 days (2.0×) |
| GPT-2 small - FLASHATTENTION | 18.2 | **2.7 days (3.5×)** |

# Systems: parallelization

- Problem:

  - model very big => can't fit on one GPU

  - Want to use as many GPUs as possible

- Idea: split memory and compute across GPUs

- Background: to naively train a P parameter model you need at least 16P GB of DRAM

  - 4P GB for model weights

  - 2 * 4P GB for optimizer

  - 4P GP for gradients

- E.g. for 7B model you need 112GB!

# Systems: data parallelism

- Goal: use more GPUs

- Naïve data parallelization:

  1. Copy model & optimizer on each GPU

  2. Split data

  3. Communicate and reduce (sum) gradients

- Pro: use parallel GPU

- Con: no memory gains!

# Systems: data parallelism

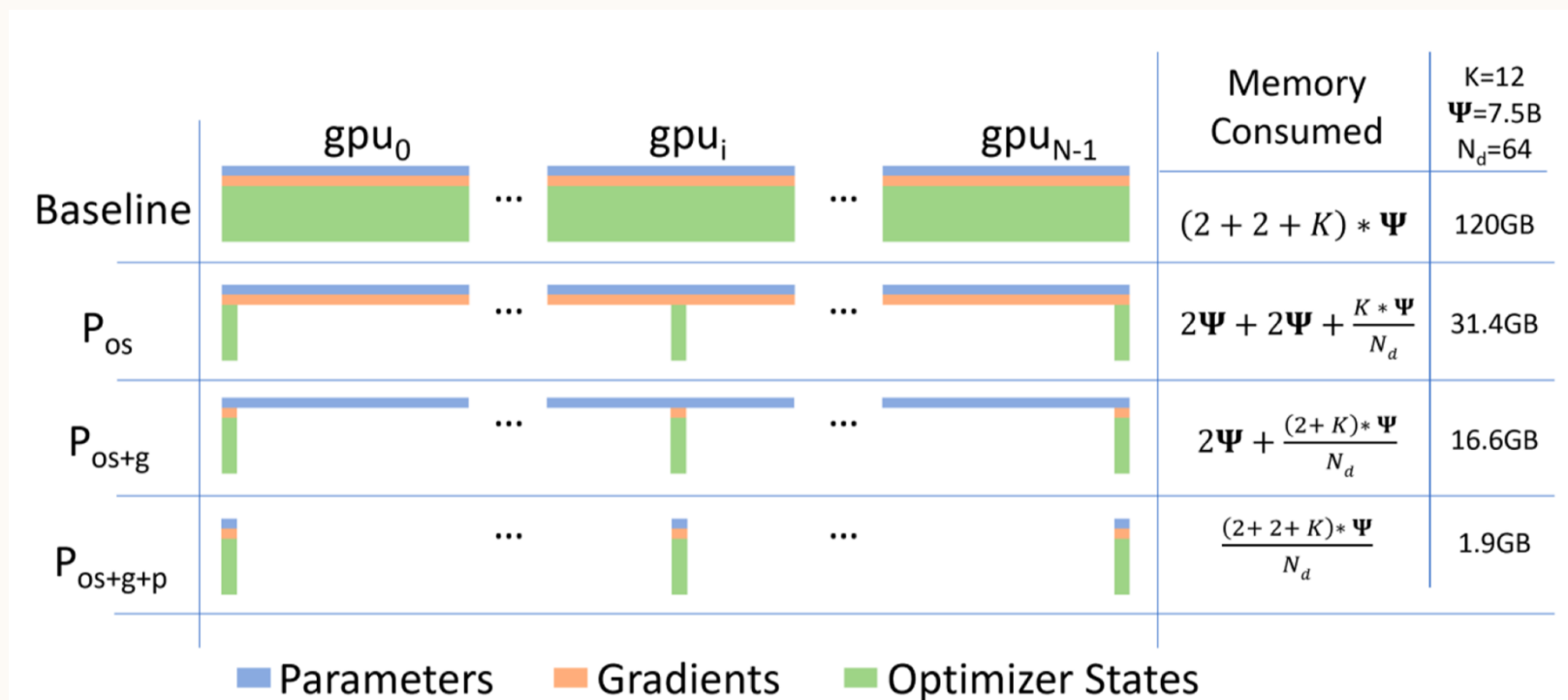- Goal: split up memory

- Idea: each GPU updates subset of weights and them before next step => sharding
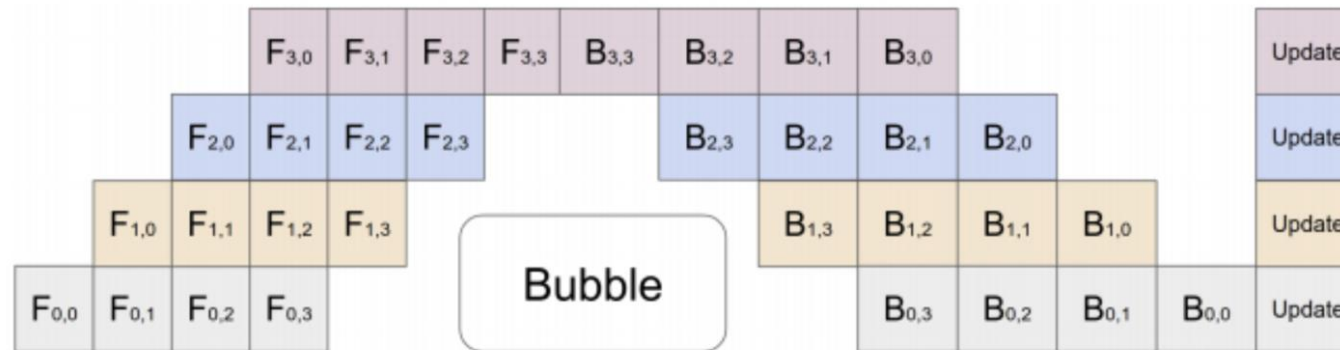


ZeRO
[Rajbhandari+ 2019]

# Systems: model parallelism

- Problem: data parallelism only works if batch size >= # GPUS

- Idea: have every GPU take care of applying specific parameters (rather than updating)

  - Eg **pipeline parallel:** every GPU has different layer



Pipelined Execution

To alleviate this problem, pipeline parallelism splits the input minibatch into multiple microbatches and pipelines the execution of these microbatches across multiple GPUs. This is outlined in the figure below:

The figure represents a model with 4 layers placed on 4 different GPUs (vertical axis). The horizontal axis represents training this model through time demonstrating that the GPUs are utilized much more efficiently. However, there still exists a bubble (as demonstrated in the figure) where certain GPUs are not utilized. (image source).

GPipe
[Huang+ 2018]

# Systems: model parallelism

- Problem: data parallelism only works if batch size >= # GPUS

- Idea: have every GPU take care of applying specific parameters (rather than updating)

  - Eg **pipeline parallel:** every GPU has different layer

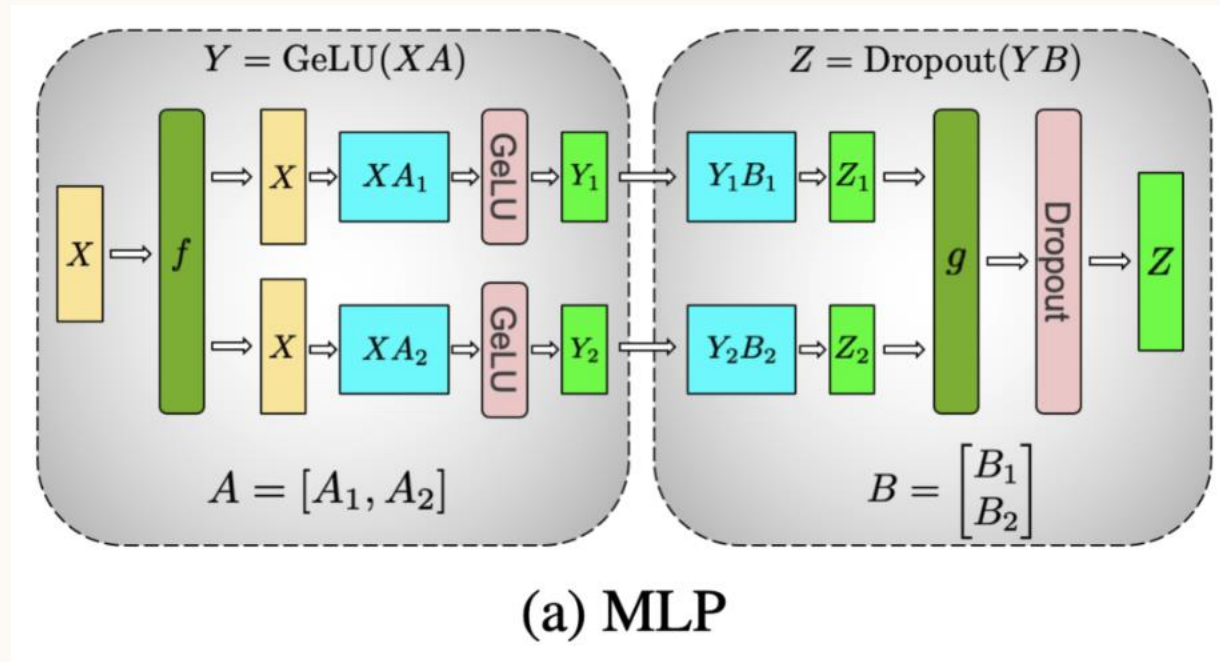  - Eg **tensor parallel:** split single matrix across GPUs and use partial sum



(a) MLP

Megatron-LM:
[Shoeybi+ 2019]

# Systems: architecture sparsity

- Idea: models are huge => not every datapoint needs to go through every parameter

  - Eg **Mixture of Experts:** use a selector layer to have less "active" parameter => same FLOPs



Sparse Expert Models:
[Fedus+ 2012]

# Questions?

# LLM evaluation: Perplexity

- Idea: validation loss

$$PPL(x_{1:L}) = 2^{\frac{1}{L}\mathcal{L}(x_{1:L})} = \prod p(x_i|x_{1:i-1})^{-1/L}$$

- To be more interpretable: use **perplexity**

  - avg per token (~independent of length)

  - Exponentiate => units independent of log base

- Perplexity: between 1 and |Vocab|

  - Intuition: number of tokens that you are hesitating between

# LLM evaluation: Perplexity



Between 2017-2023, models went from "hesitating" between ~70 tokens to <10 tokens

Perplexity not used anymore for academic benchmark but still important for development

# LLM evaluation: spurious correlation

- e.g. LLM prefers longer outputs

- Possible solution: regression analysis / causal inferece to "control" length

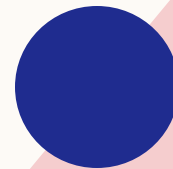Annotator:  ● Human $p_{ref}$   ● Trainer $p_{sim}^{ann}$   ● Evaluator $p_{sim}^{eval}$   ● GPT4 $p_{sim}^{GPT4}$
Model:  ■ Human $p_{ref}$   ◆ Simulated $p_{sim}$   ● GPT4   ▲ ChatGPT   ⬟ Davinci003

| | AlpacaEval | | | Length-controlled AlpacaEval | | |
|---|---|---|---|---|---|---|
| | concise | standard | verbose | concise | standard | verbose |
| gpt4_1106_preview | 22.9 | 50.0 | 64.3 | 41.9 | 50.0 | 51.6 |
| Mixtral-8x7B-Instruct-v0.1 | 13.7 | 18.3 | 24.6 | 23.0 | 23.7 | 23.2 |
| gpt4_0613 | 9.4 | 15.8 | 23.2 | 21.6 | 30.2 | 33.8 |
| claude-2.1 | 9.2 | 15.7 | 24.4 | 18.2 | 25.3 | 30.3 |
| gpt-3.5-turbo-1106 | 7.4 | 9.2 | 12.8 | 15.8 | 19.3 | 22.0 |
| alpaca-7b | 2.0 | 2.6 | 2.9 | 4.5 | 5.9 | 6.8 |

AlpacaEval LC
[Dubois+ 2023]

# Wrap-up

# Outlook

Haven't touched upon:

- Architecture: MoE & SSM
- Decoding & inference
- UI & tools: ChatGPT
- Multimodality

- Misuse
- Context size
- Data wall
- Legality of data collection

Going further:

- CS224N: more of the background and historical context. Some adjacent material.
- CS324: more in-depth reading and lectures.
- CS336: you actually build your LLM. Heavy workload!

# Tokenizer

- Why?

  - More general than words (eg typos)

  - Shorter sequences than with characters

- Idea: tokens as common subsequences (~3 letters)

- Eg: Byte Pair Encoding (BPE). Train steps:

1. Take large corpus of text

tokenizer: text to token index

# Tokenizer

- Why?

  - More general than words (eg typos)

  - Shorter sequences than with characters

- Idea: tokens as common subsequences

- Eg: Byte Pair Encoding (BPE). Train steps:

1. Take large corpus of text

2. Start with one token per character

# Tokenizer

- Why?

  - More general than words (eg typos)

  - Shorter sequences than with characters

- Idea: tokens as common subsequences

- Eg: Byte Pair Encoding (BPE). Train steps:

1. Take large corpus of text

2. Start with one token per character

3. Merge common pairs of tokens into a token

# Tokenizer

- Why?

  - More general than words (eg typos)

  - Shorter sequences than with characters

- Idea: tokens as common subsequences

- Eg: Byte Pair Encoding (BPE). Train steps:

1. Take large corpus of text

2. Start with one token per character

3. Merge common pairs of tokens into a token

4. Repeat until desired vocab size or all merged

# Tokenizer

- Why?

  - More general than words (eg typos)

  - Shorter sequences than with characters

- Idea: tokens as common subsequences

- Eg: Byte Pair Encoding (BPE). Train steps:

1. Take large corpus of text

2. Start with one token per character

3. Merge common pairs of tokens into a token

4. Repeat until desired vocab size or all merged

# Tokenizer

- Why?

  - More general than words (eg typos)

  - Shorter sequences than with characters

- Idea: tokens as common subsequences

- Eg: Byte Pair Encoding (BPE). Train steps:

1. Take large corpus of text

2. Start with one token per character

3. Merge common pairs of tokens into a token

4. Repeat until desired vocab size or all merged

# Tokenizer

- Why?

  - More general than words (eg typos)

  - Shorter sequences than with characters

- Idea: tokens as common subsequences

- Eg: Byte Pair Encoding (BPE). Train steps:

1. Take large corpus of text

2. Start with one token per character

3. Merge common pairs of tokens into a token

4. Repeat until desired vocab size or all merged