# A Comprehensive Overview of Large Language Models (LLM): Insights from a Machine Learning System Perspective

**Shengyuan Ye**
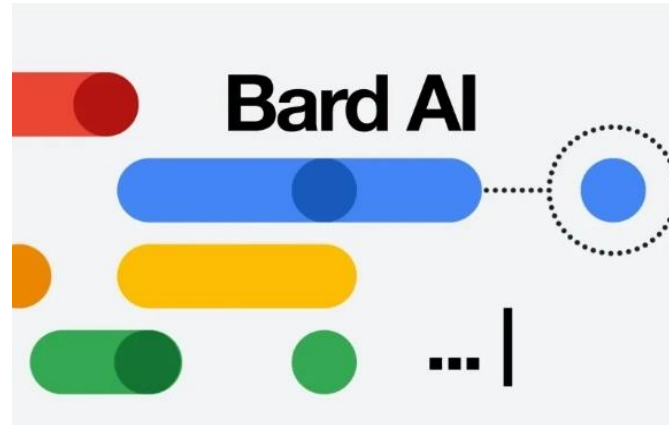
School of Computer Science and Engineering
Sun Yat-sen University
Contact: yeshy8@mail2.sysu.edu.cn

# LLM and ChatGPT

- **Chat-based LLM is walking into our daily life!**



**OpenAI - ChatGPT**

**Google -Bard**

**Microsoft-Bing**

Credit: Google images.

# LLM and ChatGPT

- **LLMs are taking our jobs!**

GPTs are GPTs: An Early Look at the Labor Market Impact
Potential of Large Language Models

Tyna Eloundou[1], Sam Manning[1,2], Pamela Mishkin[*1], and Daniel Rock[3]

[1]OpenAI
[2]OpenResearch
[3]University of Pennsylvania

March 20, 2023

If human logic and creativity can be replaced.
What jobs do you think will be left?

> Lior ⚡ ✔
> @AlphaSignalAI
>
> For the last 10 years I believed AI will free humanity from brainless tasks and push the world towards a more creative future.
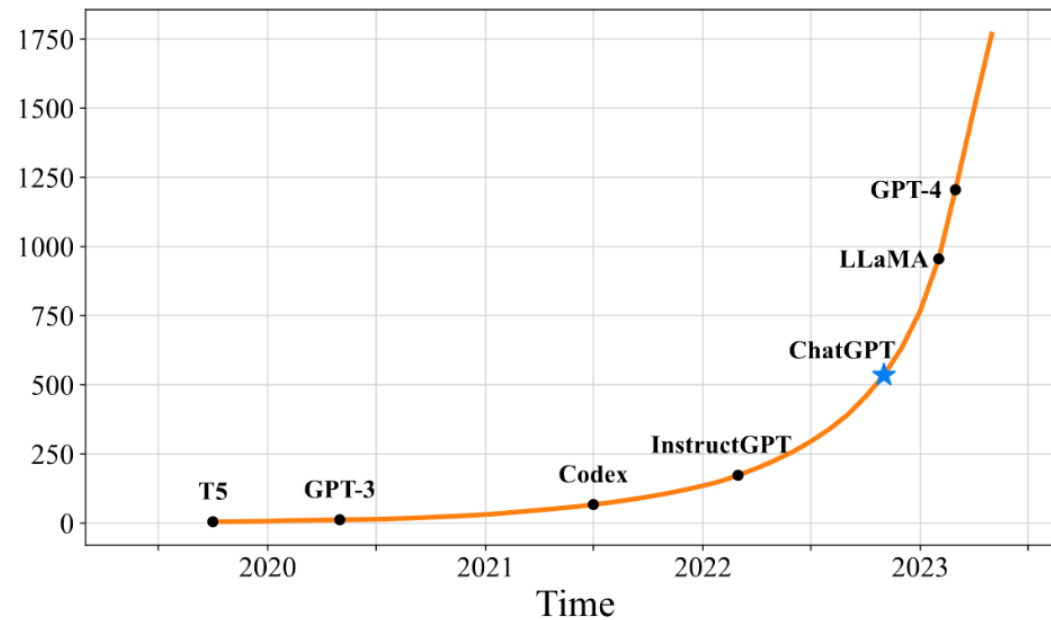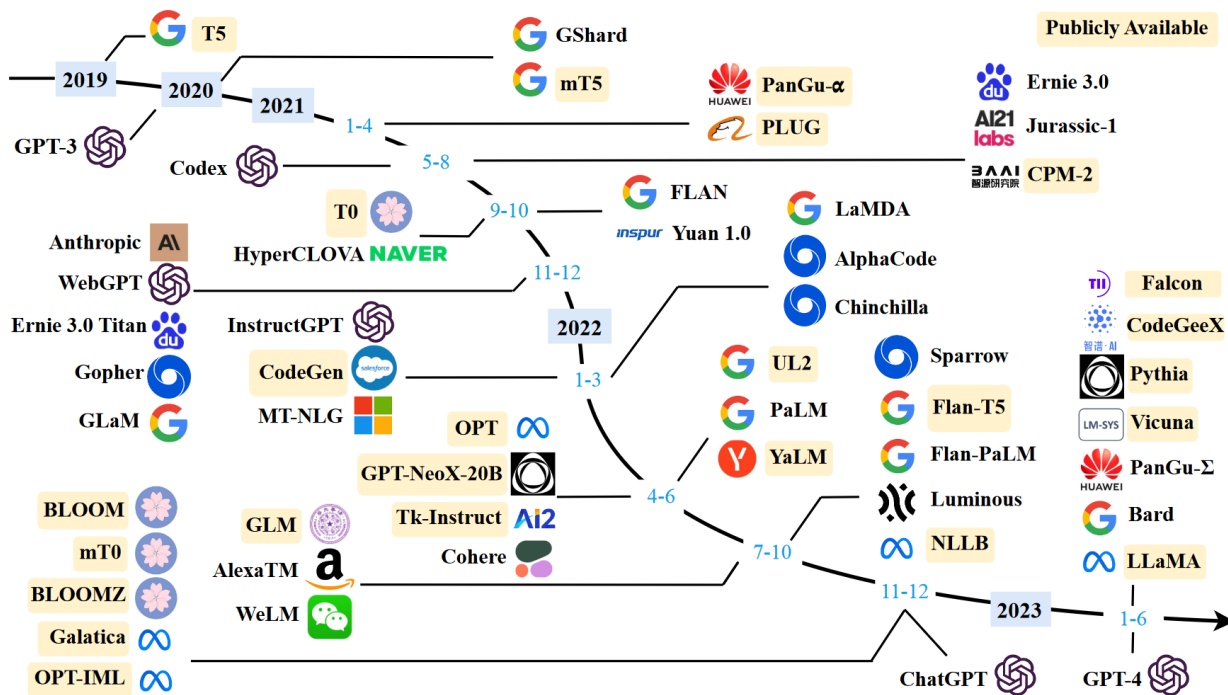>
> However, with models like Stable Diffusion, artists are also being pushed out.

## 没有GPT暴露风险的职业

| | |
|---|---|
| 农业设备操作员 | 油漆匠、泥水匠、瓦匠助手 |
| 运动员与体育竞赛者 | 管道工助手 |
| 汽车玻璃的安装和维修工 | 屋顶工人助手 |
| 公交与卡车机械师与柴油发动机专家 | 鱼肉禽类的切割工 |
| 水泥泥瓦匠和混凝土修整工 | 摩托车机械师 |
| 厨师 | 铺路及夯实设备操作员 |
| 手动切割与裁剪的工作者 | 打桩机操作员 |
| 石油与天然气钻探业的井架操作员 | 金属浇注机操作员 |
| 餐厅与自助餐厅的服务员和调酒师助手 | 铁路轨道铺设和维护设备操作员 |
| 洗碗工 | 耐火材料维修商 |
| 疏浚操作员 | 采矿顶板锚杆机操作员 |
| 电力管线安装与维修工 | 石油和天然气钻探业的体力劳动者 |
| 地表矿业挖掘、装载、拉铲机械操作员 | 屠宰工和肉类包装机操作员 |
| 地板工 | 石匠 |
| 铸造模具和制芯师 | 密封石膏板或其他墙板的接缝工 |
| 砖匠、石匠、瓷砖匠安装工助手 | 轮胎修理工与更换工 |
| 木匠助手 | 井口泵送机操作员 |

Credit: Google images.

# Trend of LLM in Research Fields

- LLM is rapidly emerging as the hottest direction in research fields



(b) Query="Large Language Model"

A sharp increase occurs after the release of ChatGPT: the average number of published arXiv papers that contain "large language model" in title or abstract goes **from 0.40 per day to 8.58 per day**.

Credit: https://github.com/RUCAIBox/LLMSurvey.

# The Cost Barrier of LLM



- The cost of training GPT-3 is estimated to be around $1.4 million, and for some larger LLM models, the training costs range between $2 million to $12 million.

- The cost of operating OpenAI's ChatGPT could potentially reach $0.7 million per day.

# The Cost Barrier of LLM

- How much does it cost when using ChatGPT to finish a writing task?

## Prices of GPT4

| Model | Input | Output |
|---|---|---|
| 8K context | $0.03 / 1K tokens | $0.06 / 1K tokens |
| 32K context | $0.06 / 1K tokens | $0.12 / 1K tokens |

## Prices of GPT3.5-Turbo

| Model | Input | Output |
|---|---|---|
| 4K context | $0.0015 / 1K tokens | $0.002 / 1K tokens |
| 16K context | $0.003 / 1K tokens | $0.004 / 1K tokens |

全国甲卷

阅读下面的材料，根据要求写作。（60分）

人们因技术发展得以更好地掌控时间，但也有人因此成了时间的仆人。

这句话引发了你怎样的联想与思考？请写一篇文章。

要求：选准角度，确定立意，明确文体，自拟标题；不要套作，不得抄袭；不得泄露个人信息；不少于800字。

- Input Tokens: 100 tokens
- Output Tokens: 800 tokens

## Price of using ChatGPT4:

- Input Price: **100 tokens × $0.03/1K = $ 0.003**
- Output Price: **800 tokens × $0.06/1k = $ 0.048**
- **Total Price: $ 0.003 + $ 0.048 = $ 0.051 = ¥ 0.371**

# The Cost Barrier of LLM



ChatGPT and Comparisons, Worldwide
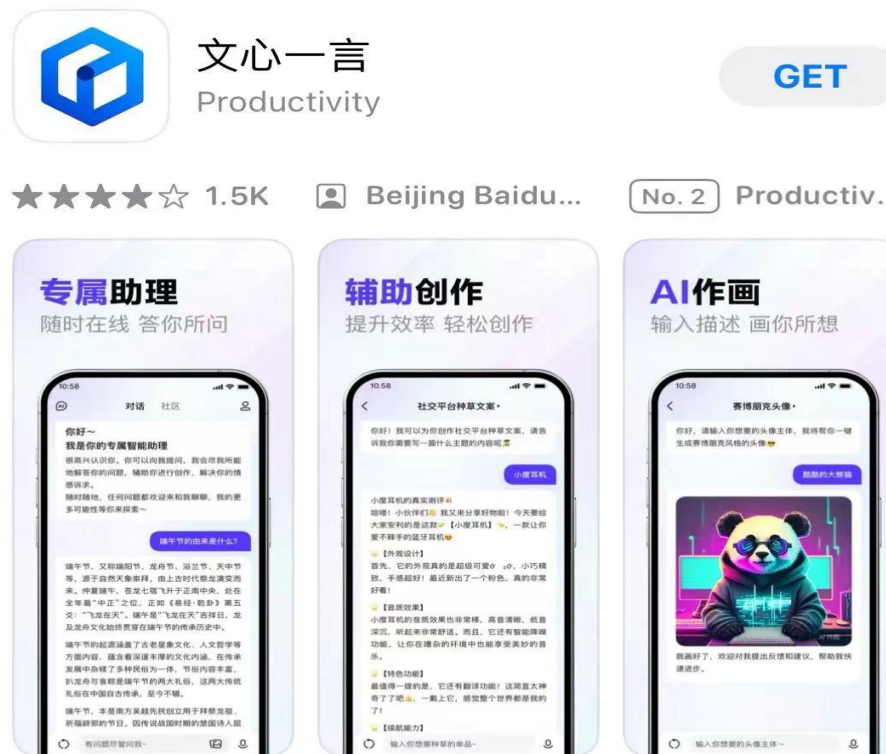Monthly Visits Desktop & Mobile Web Worldwide

* Preliminary Estimate

- chat.openai.com
- bing.com
- character.ai
- bard.google.com

- According to data released in May of this year, the ChatGPT website has surpassed 1.5 billion monthly active users.
- Due to immense cost pressures, companies that fail to capture market share will ultimately be eliminated.



百度、商汤等大模型产品获批，今日起全面开放上线

机器之心  2023-08-31 13:21  发表于北京

文心一言
Productivity

GET

★★★★☆ 1.5K  Beijing Baidu...  No. 2  Productiv...

专属助理
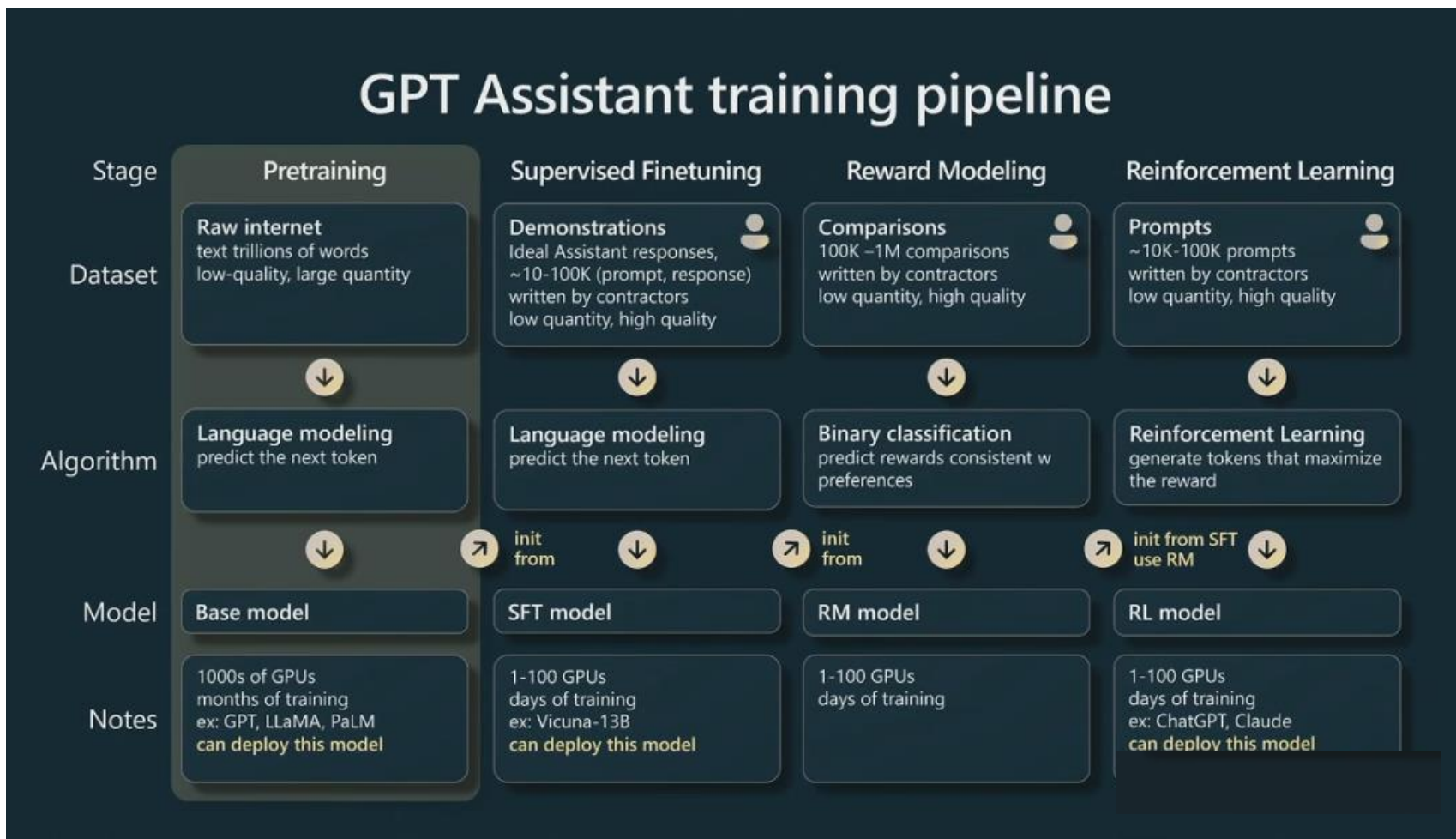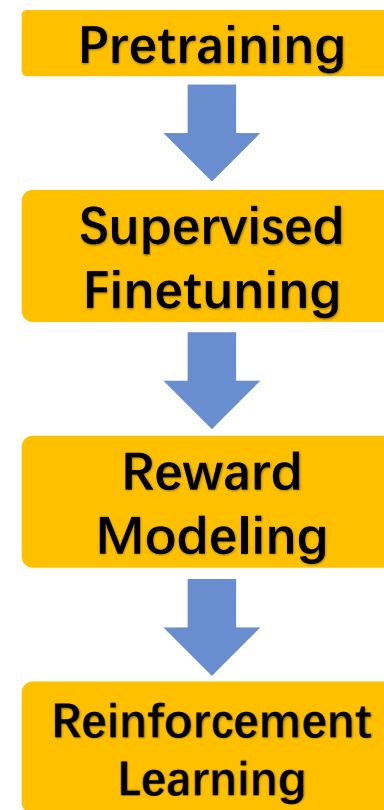随时在线 答你所问

辅助创作
提升效率 轻松创作

AI作画
输入描述 画你所想

Credit:Google.

# How to Train your ChatGPT Assistant

- **ChatGPT Training Pipeline**



**Four Stage Pipeline**

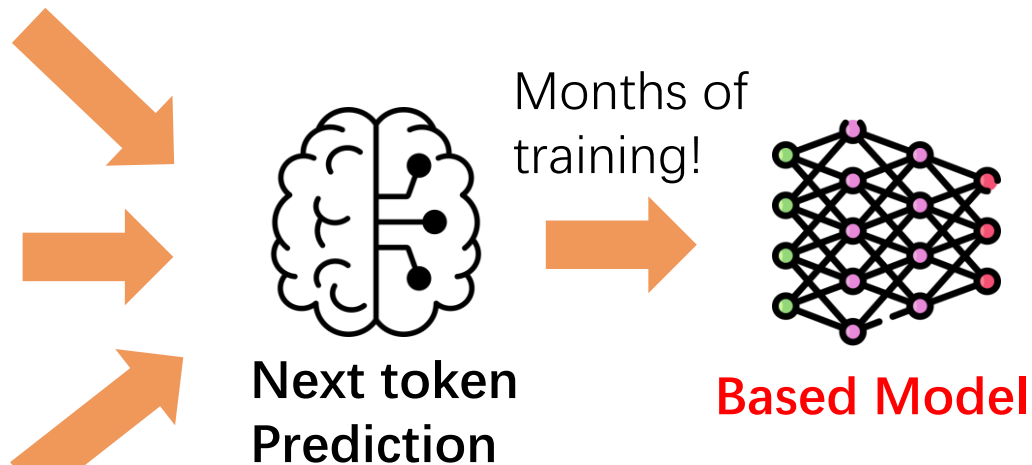Pretraining → Supervised Finetuning → Reward Modeling → Reinforcement Learning

Credit: Andrej Karpathy @ OpenAI.
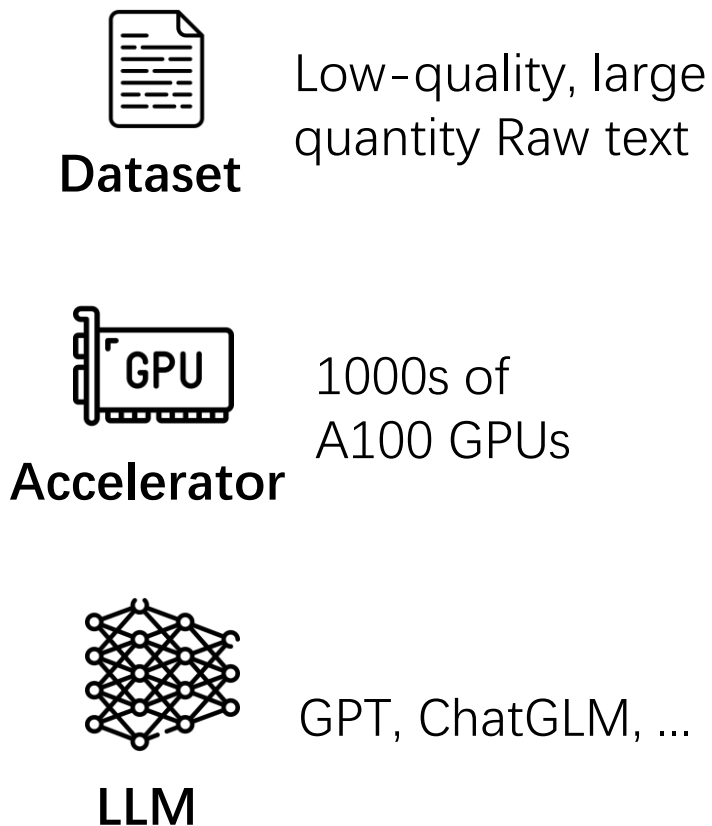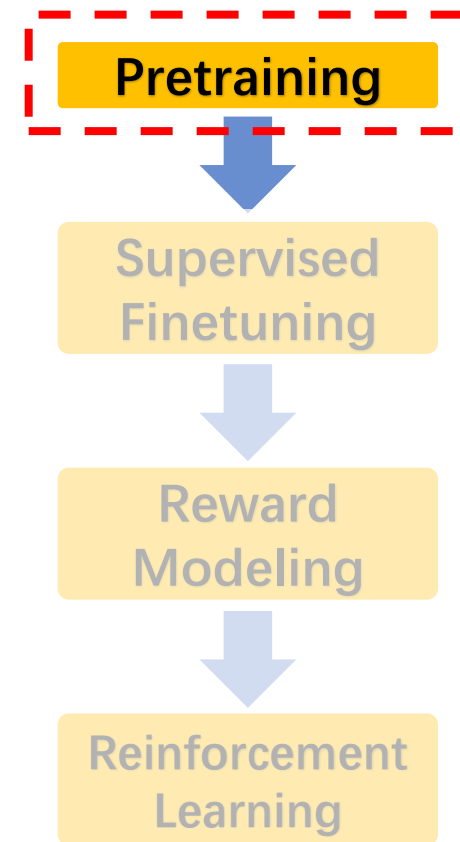
# How to Train a ChatGPT Assistant?

The LLM pre-training takes up **99%** of the entire training pipeline's time and typically requires **thousands of GPUs** for training over **several months**.

**Dataset**
Low-quality, large quantity Raw text

**Accelerator**
1000s of A100 GPUs

**LLM**
GPT, ChatGLM, ...

**Next token Prediction**

Months of training!

**Based Model**

**Four Stage Pipeline**

**Pretraining**

**Supervised Finetuning**

**Reward Modeling**

**Reinforcement Learning**

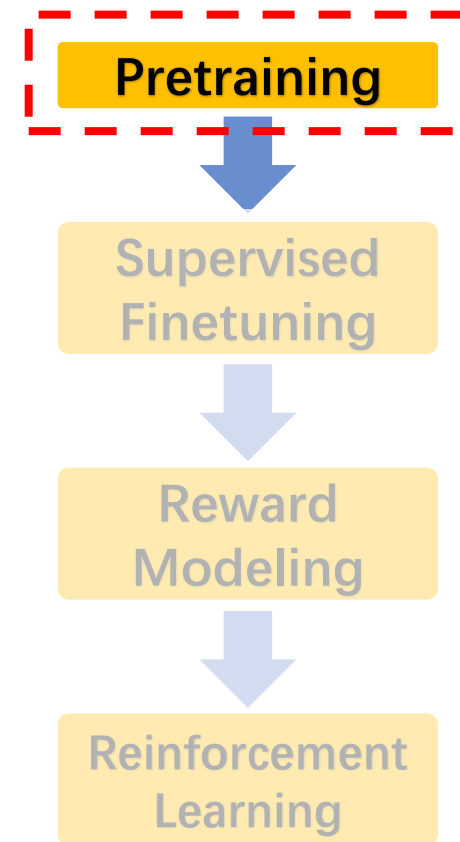Credit: Andrej Karpathy @ OpenAI.

# How to Train a ChatGPT Assistant?

Smaller technology enterprises and educational research labs often can't afford the cost of pre-training LLMs.



**Four Stage Pipeline**

**Pretraining**

⬇

**Supervised Finetuning**

⬇

**Reward Modeling**

⬇

**Reinforcement Learning**

Credit: Andrej Karpathy @ OpenAI.

# How to Train a ChatGPT Assistant?

- **Based model are not Chat Assistant!**

Write a poem about bread!

Write a poem about someone.
Write a poem about angel.
Write a poem about basketball.

The training task of pre-trained language model is to predict the next token, rather than engaging in QA (Question-Answering) dialogues. It may use more questions to answer a question.
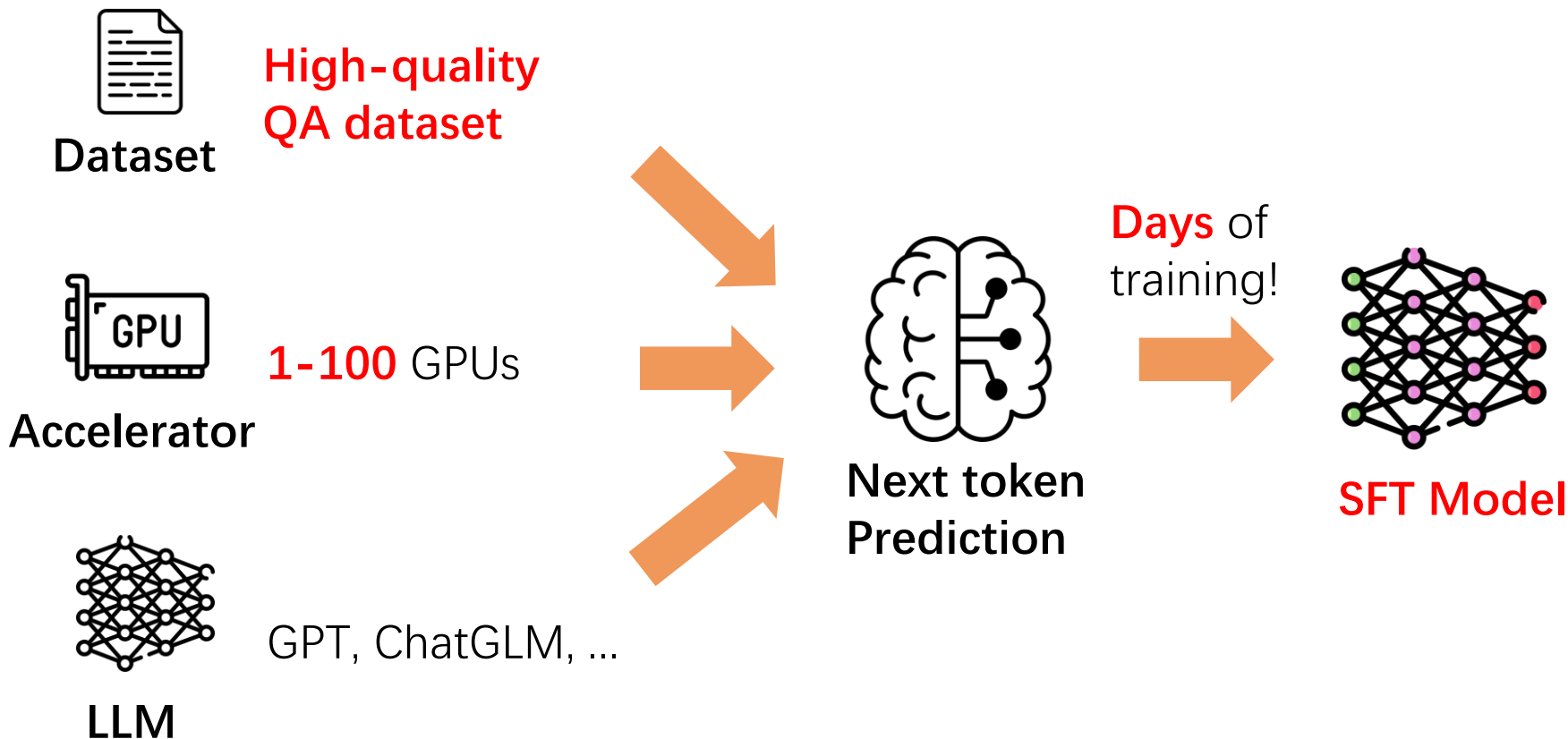
**Four Stage Pipeline**

**Pretraining**

**Supervised Finetuning**

**Reward Modeling**

**Reinforcement Learning**
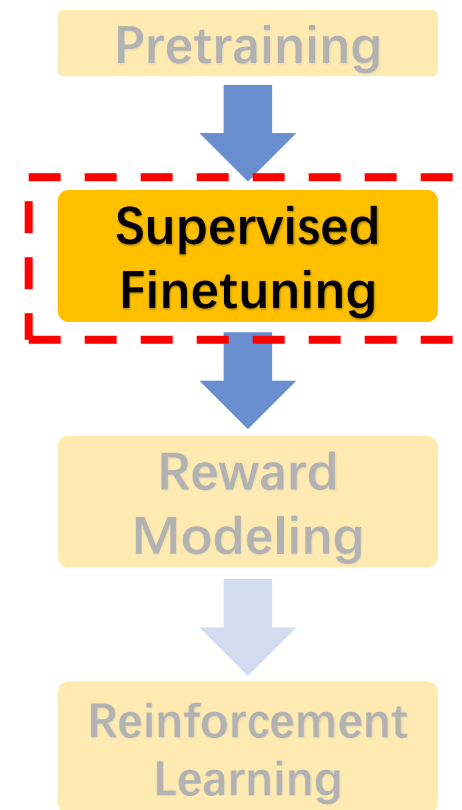
Credit: Andrej Karpathy @ OpenAI.

# How to Train a ChatGPT Assistant?

During the Supervised Finetuning of LLM, workers are hired to gather **high-quality QA data** and the **based-model is fine-tuned**, requiring fewer GPUs and just days to train.

**Dataset**

**High-quality QA dataset**

**Accelerator**

**1-100** GPUs

**LLM**

GPT, ChatGLM, …

**Next token Prediction**

**Days** of training!

**SFT Model**

**Four Stage Pipeline**

Pretraining

⬇

**Supervised Finetuning**

⬇

Reward Modeling

⬇

Reinforcement Learning

Credit: Andrej Karpathy @ OpenAI.

# How to Train a ChatGPT Assistant?

Multiple answers are generated from the same prompt based on SFT model. Workers rank these answers to compile a large ranking dataset, which is then used to train a Transformer-based Reward Model.

**Four Stage Pipeline**



Pretraining

Supervised Finetuning

Reward Modeling

Reinforcement Learning

Credit: Andrej Karpathy @ OpenAI.

# How to Train a ChatGPT Assistant?

In reinforcement learning, policy gradient algorithm are used to amplify the generation probabilities of "favorable" responses and minimize those of "toxic" ones.

**Four Stage Pipeline**

- **Action Space: all the vocabulary**
- **State: currently generated token sequence**
- **Reward: Provided by Reward Model**

Pretraining

Supervised Finetuning

Reward Modeling

Reinforcement Learning



RL Fine-tuning

Prompts → LM Outputs → Reward Model → 😊/😞 Reward → Training with RL algorithm (PPO) → Aligned LM

Credit: Andrej Karpathy @ OpenAI.

# How to Train a ChatGPT Assistant?

Supervised
Finetuning

Reinforcement
Learning Alignment

**Based Model** → **SFT Model** → **Aligned Model**

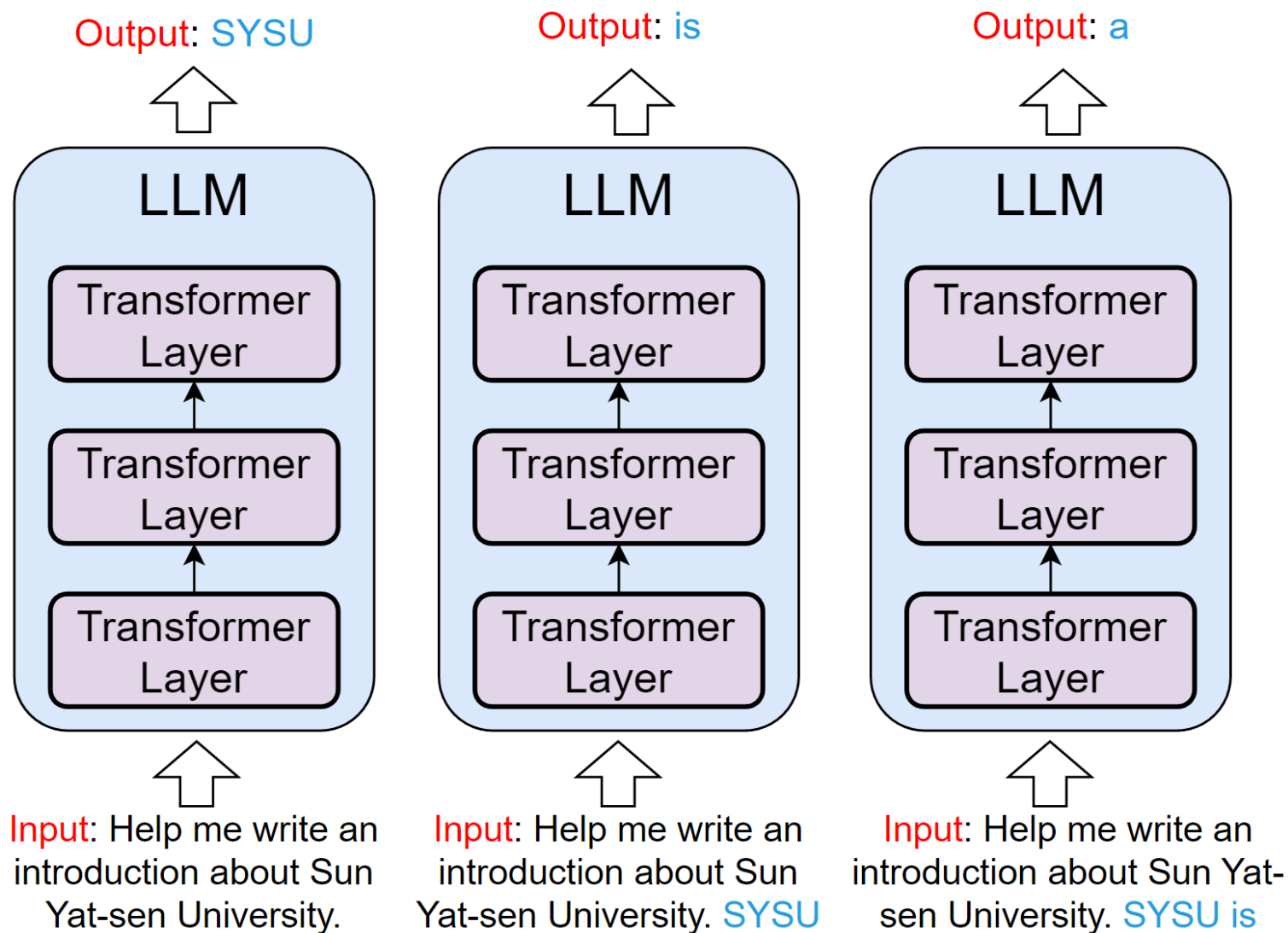| Rank | Model | Elo Rating | Description | License |
|---|---|---|---|---|
| 1 | GPT-4 | 1274 | ChatGPT-4 by OpenAI | Proprietary |
| 2 | Claude-v1 | 1224 | Claude by Anthropic | Proprietary |
| 3 | GPT-3.5-turbo | 1155 | ChatGPT-3.5 by OpenAI | Proprietary |
| 4 | Vicuna-13B | 1083 | a chat assistant fine-tuned from LLaMA on user-shared conversations by LMSYS | Weights available; Non-commercial |
| 5 | Koala-13B | 1022 | a dialogue model for academic research by BAIR | Weights available; Non-commercial |
| 6 | RWKV-4-Raven-14B | 989 | an RNN with transformer-level LLM performance | Apache 2.0 |
| 7 | Oasst-Pythia-12B | 928 | an Open Assistant for everyone by LAION | Apache 2.0 |
| 8 | ChatGLM-6B | 918 | an open bilingual dialogue language model by Tsinghua University | Weights available; Non-commercial |
| 9 | StableLM-Tuned-Alpha-7B | 906 | Stability AI language models | CC-BY-NC-SA-4.0 |
| 10 | Alpaca-13B | 904 | a model fine-tuned from LLaMA on instruction-following demonstrations by Stanford | Weights available; Non-commercial |
| 11 | FastChat-T5-3B | 902 | a chat assistant fine-tuned from FLAN-T5 by LMSYS | Apache 2.0 |
| 12 | Dolly-V2-12B | 863 | an instruction-tuned open large language model by Databricks | MIT |
| 13 | LLaMA-13B | 826 | open and efficient foundation language models by Meta | Weights available; Non-commercial |

**Aligned Model** (ranks 1–3)

**SFT Model** (ranks 4–5)

- The SFT model significantly outperforms the based pre-trained model in QA tasks.

- The Aligned Model can further filter out expressions from the SFT model's output that are harmful or not in line with human norms.

Credit: Andrej Karpathy @ OpenAI.

# What happens when LLM generates an answer?

Output: SYSU

Output: is

Output: a

LLM

Transformer Layer

Transformer Layer

Transformer Layer

LLM

Transformer Layer

Transformer Layer

Transformer Layer

LLM

Transformer Layer

Transformer Layer

Transformer Layer

Input: Help me write an introduction about Sun Yat-sen University.

Input: Help me write an introduction about Sun Yat-sen University. SYSU

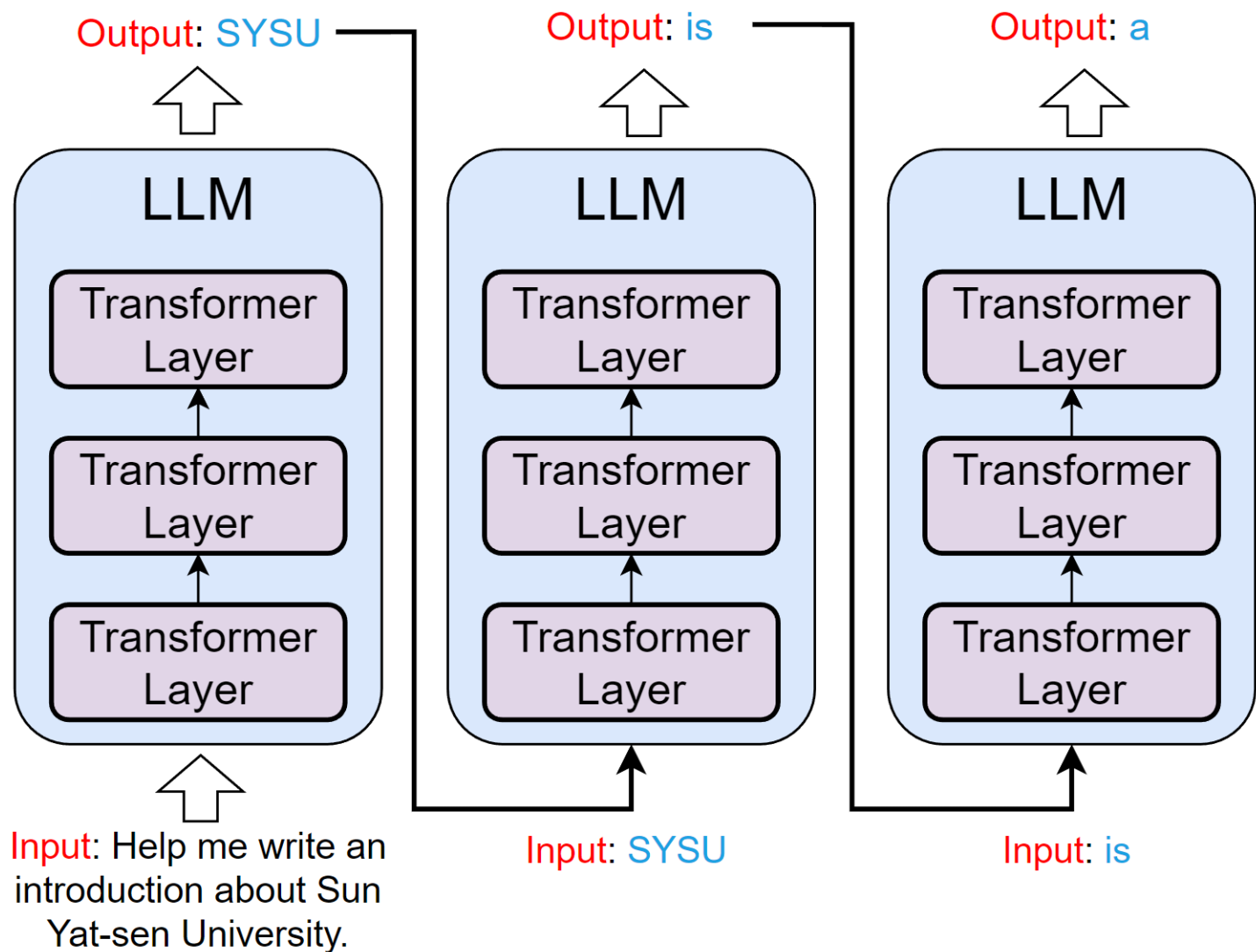Input: Help me write an introduction about Sun Yat-sen University. SYSU is

**Autoregressive**

- When a LLM generates response, it uses its own previous outputs as inputs for future predictions, forming a chain of dependencies.

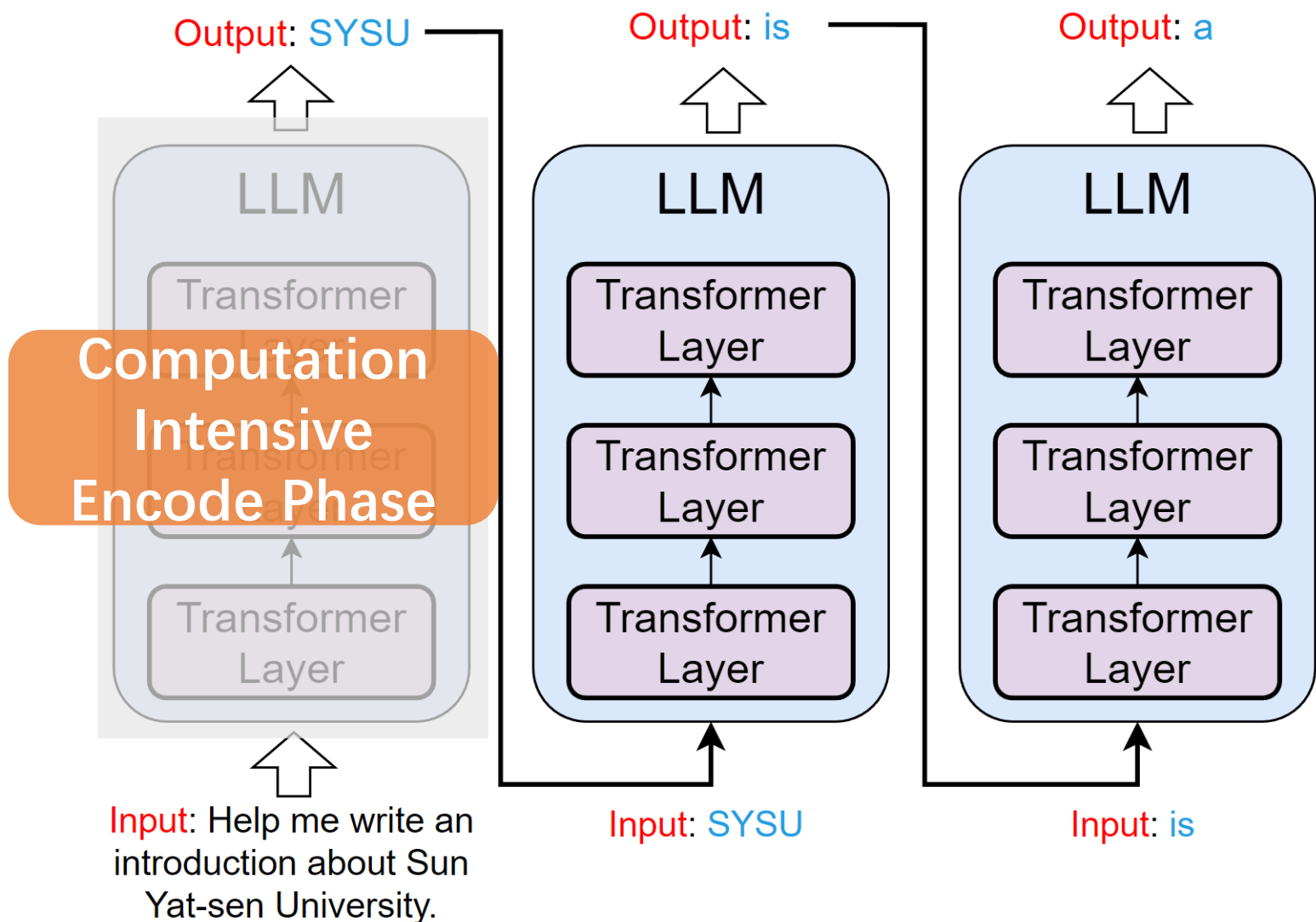- This autoregressive behavior allows the model to generate coherent and contextually relevant text.

Credit: Andrej Karpathy @ OpenAI.

# What happens when LLM generates an answer?

Output: SYSU

Output: is

Output: a

LLM

Transformer Layer

Transformer Layer

Transformer Layer

Input: Help me write an introduction about Sun Yat-sen University.

LLM

Transformer Layer

Transformer Layer

Transformer Layer

Input: SYSU

LLM

Transformer Layer

Transformer Layer

Transformer Layer

Input: is

## Intermediate Cache

- LLM will cache previous computational results (such as the calculations from the black parts of Self-Attention) in memory to avoid redundant calculations.

- The inference process of LLM can be divided into two distinct phases: Encode and Decode
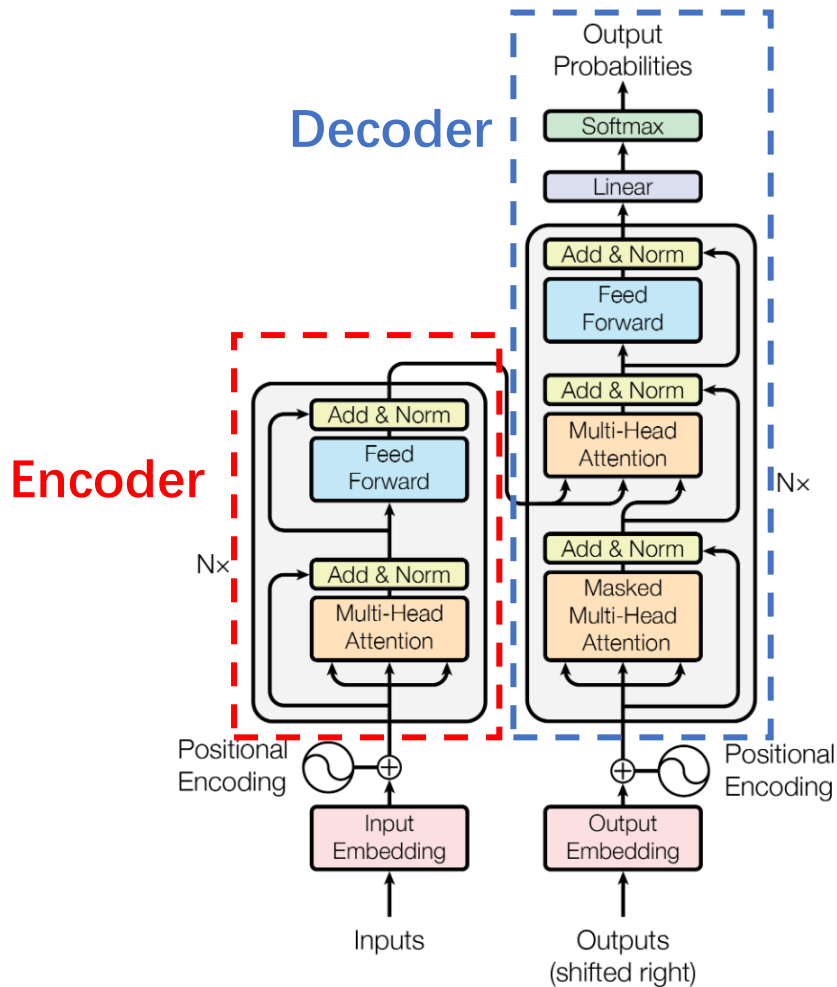
Credit: Andrej Karpathy @ OpenAI.

# What happens when LLM generates an answer?



- **Encode Phase**: Inference on LLM's input involves potentially hundreds or even thousands of tokens, making it computationally intensive.

Credit: Andrej Karpathy @ OpenAI.

# What happens when LLM generates an answer?

Output: SYSU     Output: is     Output: a

**LLM**

Transformer Layer

Transformer Layer

Transformer Layer

LLM

Transformer Layer

Transformer Layer

Transformer Layer

LLM

Transformer Layer

Transformer Layer

Transformer Layer

**Memory Intensive Decode Phase**

Input: Help me write an introduction about Sun Yat-sen University.

Input: SYSU

Input: is

- **Encode Phase**: Inference on LLM's input involves potentially hundreds or even thousands of tokens, making it computationally intensive.

- **Decode Phase**: each previously predicted token is inputted one at a time, requiring frequent retrieval of intermediate cache from storage.

Credit: Andrej Karpathy @ OpenAI.

# The Backbone Architecture in LLM

- The vast majority of LLMs are based on the **Transformer architecture**.
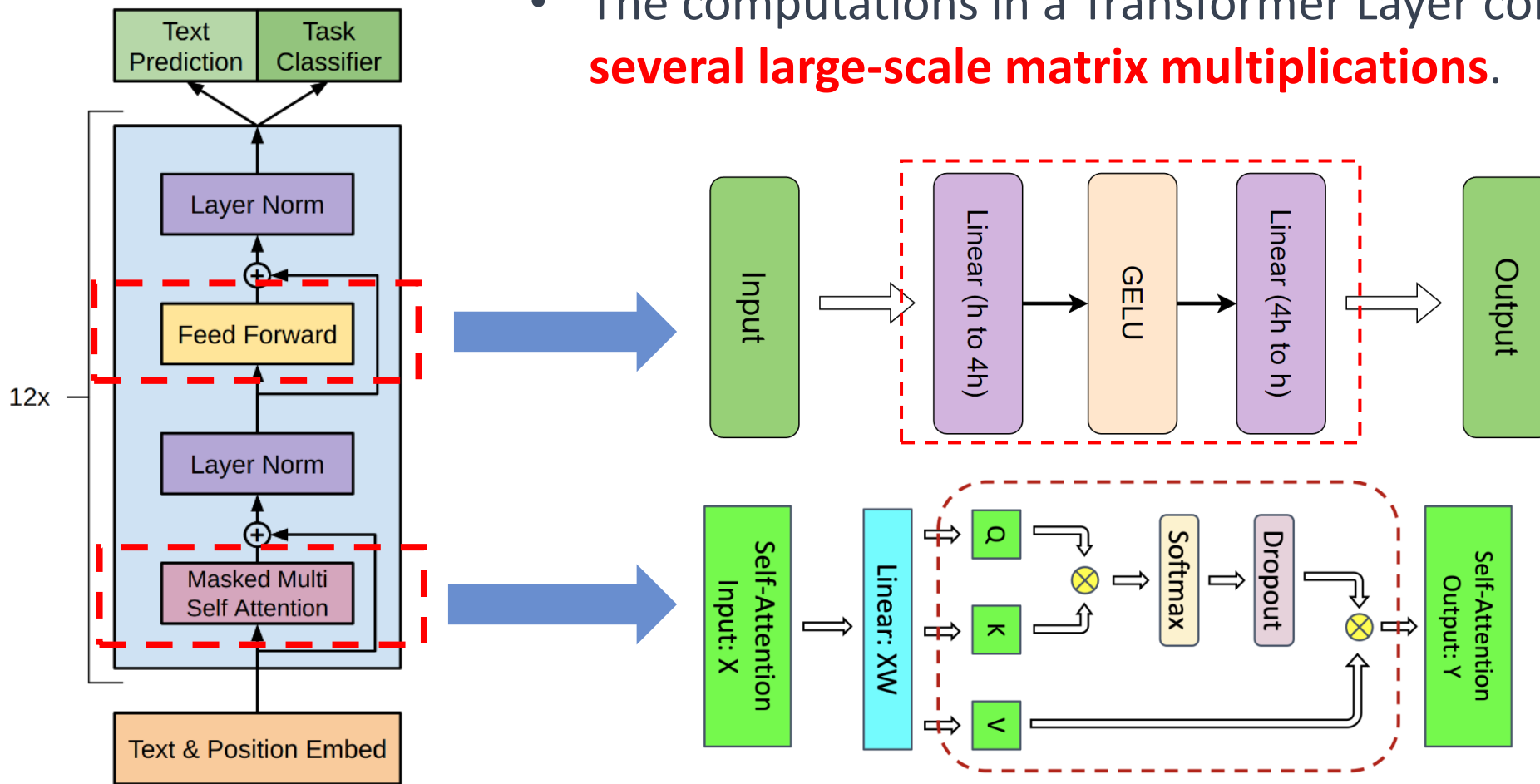


**Attention is all you need! (NIPS 2017)**

**GPT Series (2019-now)**

# The Backbone Architecture in LLM

- The vast majority of LLMs are based on the **Transformer architecture**.

  - The computations in a Transformer Layer consist of **several large-scale matrix multiplications**.

Credit:Reducing Activation Recomputation in Large Transformer Models.

# Parameter Size and Memory Footprint in LLM

- A language model is generally considered "large" if it has <span style="color:red">tens of millions to billions of parameters</span>.

| 模型名称 | 发布时间 | 发布机构 | 参数规模 |
|---|---|---|---|
| T5 | 2019-10 | Google | 13B |
| GPT-3 | 2020-05 | OpenAI | 175B |
| LaMDA | 2021-05 | Google | 137B |
| Jurassic | 2021-08 | AI21 | 178B |
| MT-NLG | 2021-10 | Microsoft、NVIDIA | 530B |
| ERNIE 3.0 Titan | 2021-12 | Baidu | 260B |
| Gopher | 2021-12 | DeepMind | 280B |
| Chinchilla | 2022-04 | DeepMind | 70B |
| PaLM | 2022-04 | Google | 540B |
| OPT | 2022-05 | Meta | 125M-175B |
| BLOOM | 2022-07 | BigScience | 176B |
| GLM-130B | 2022-08 | Tsinghua | 130B |
| LLaMA | 2023-02 | Meta | 7B-65B |

# Parameter Size and Memory Footprint in LLM

- A language model is generally considered "large" if it has tens of millions to billions of parameters.

| 模型名称 | 发布时间 | 发布机构 | 参数规模 |
|---|---|---|---|
| T5 | 2019-10 | Google | 13B |
| GPT-3 | 2020-05 | OpenAI | 175B |
| LaMDA | 2021-05 | Google | 137B |
| Jurassic | 2021-08 | AI21 | 178B |
| MT-NLG | 2021-10 | Microsoft、NVIDIA | 530B |
| ERNIE 3.0 Titan | 2021-12 | Baidu | 260B |
| Gopher | 2021-12 | DeepMind | 280B |
| Chinchilla | 2022-04 | DeepMind | 70B |
| PaLM | 2022-04 | Google | 540B |
| OPT | 2022-05 | Meta | 125M-175B |
| BLOOM | 2022-07 | BigScience | 176B |
| GLM-130B | 2022-08 | Tsinghua | 130B |
| LLaMA | 2023-02 | Meta | 7B-65B |

- During training, model parameters use Float64 and require 4 bytes each, while during inference, they use Float32 and require 2 bytes each.

The peak memory footprint for accommodating the ChatGLM-130B model is:

- Float64: 130B * 4 /1024/1024 = **480GB!!!**
- Float32: 130B * 2 /1024/1024 = **240GB!!!**

# How to Break the Resource Wall of a Single GPU?

- The memory of a single device is insufficient to accommodate an entire LLM.

**Megatron-LM: Training Multi-Billion Parameter Language Models Using Model Parallelism**
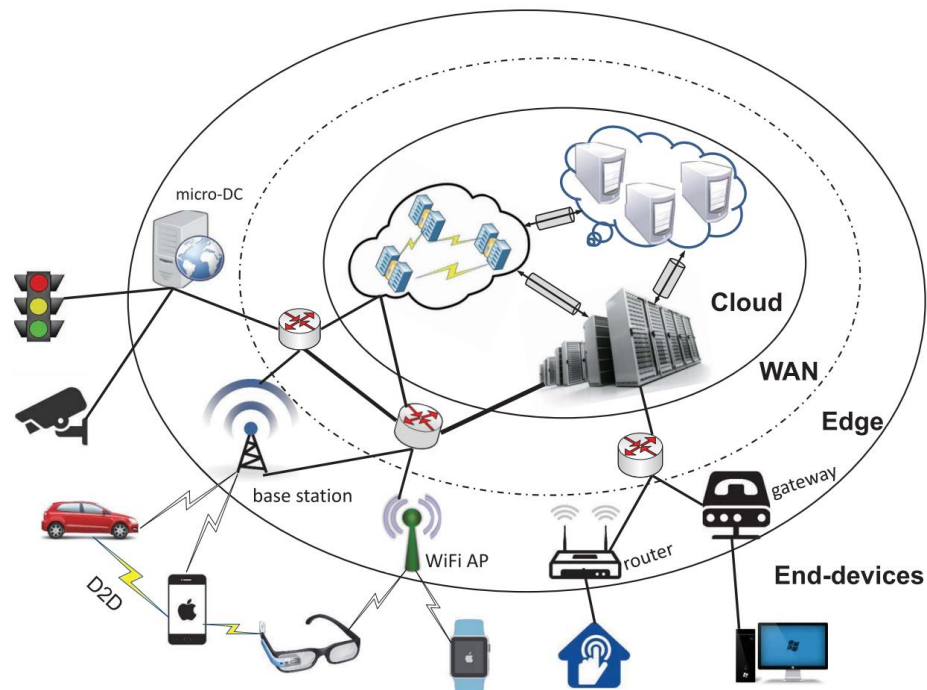
Mohammad Shoeybi[1,2]  Mostofa Patwary[1,2]  Raul Puri[1,2]  Patrick LeGresley[2]  Jared Casper[2]
Bryan Catanzaro[2]

- Leveraging matrix decomposition techniques, a Tensor-Parallel (TP) distributed algorithm was developed to partition the model across multiple GPUs, with each GPU storing only a fraction of the model's weights.

Credit:Megatron-LM: Training Multi-Billion Parameter Language Models Using Model Parallelism.

# LLM on Mobile Edge Devices

- **Mobile Computing + LLM** has emerged as a new paradigm
  - Popularization of mobile devices in both magnitude and variety
  - Proliferation of mobile data in both scale and modality

Credit: https://coruzant.com/opinion/the-future-is-edge-computing/

# LLM on Mobile Edge Devices

- **Model lightweighting** and **edge deployment** will become new research focuses in LLMs.



雷军：小米研发大模型的方向是轻量化和本地部署

南方都市报APP·湾财AI快报
综合 2023-08-25 14:19



交互升级 全新对话升级

全新唤醒页　　全新对话框架

Credit: Google Image

# Challenges of LLM on Mobile

**How to apply?**

*Constrained Capability*

*Mobile Devices*

*Heterogeneous Hardware*

*Dynamic Resources*

**Huge Gap**

- LLM Computing is extremely computation-intensive and resource-demanding

- Mobile devices are resource-constrained and heterogeneous

# Break the Memory Wall of Mobile Devices

- Utilizing the concept of paging from operating systems, Transformer layers not in use are offloaded to auxiliary storage like SD cards to expand the available memory on mobile devices.

## POET: Training Neural Networks on Tiny Devices with Integrated Rematerialization and Paging

Shishir G. Patil[1]  Paras Jain[1]  Prabal Dutta[1]  Ion Stoica[1]  Joseph E. Gonzalez[1]

Credit: POET: Training Neural Networks on Tiny Devices with Integrated Rematerialization and Paging.

# Break the Memory Wall of Mobile Devices

- Propose a framework support memory-efficient on-device LLM training



**Memory-efficient DNN Training on Mobile Devices**

In Gim and JeongGil Ko
School of Integrated Technology
College of Engineering
Yonsei University, Seoul, Korea
{hyunjun.kim,jeonggil.ko}@yonsei.ac.kr

# Collaborative Execution on Mobile Cluster

- Federated Few-shot Learning on Mobile Cluster
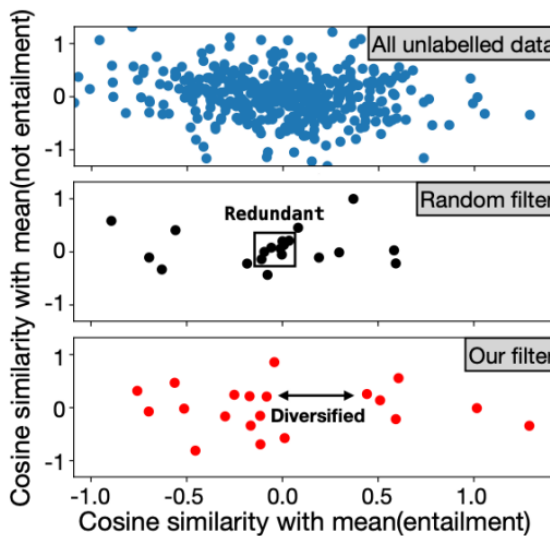
## Federated Few-Shot Learning for Mobile NLP

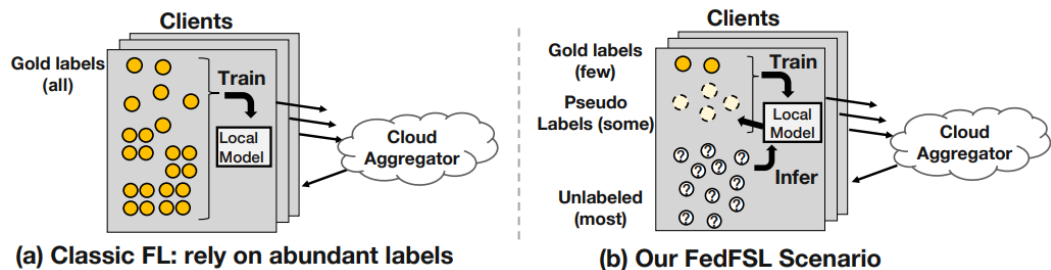**Dongqi Cai**
Beiyou Shenzhen Institute

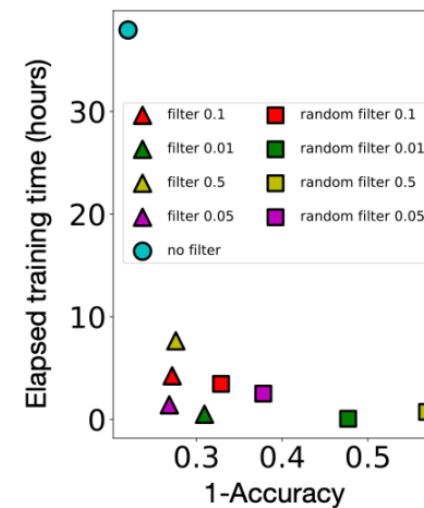**Shangguang Wang**
Beiyou Shenzhen Institute

**Yaozong Wu**
Beiyou Shenzhen Institute

**Felix Xiaozhu Lin**
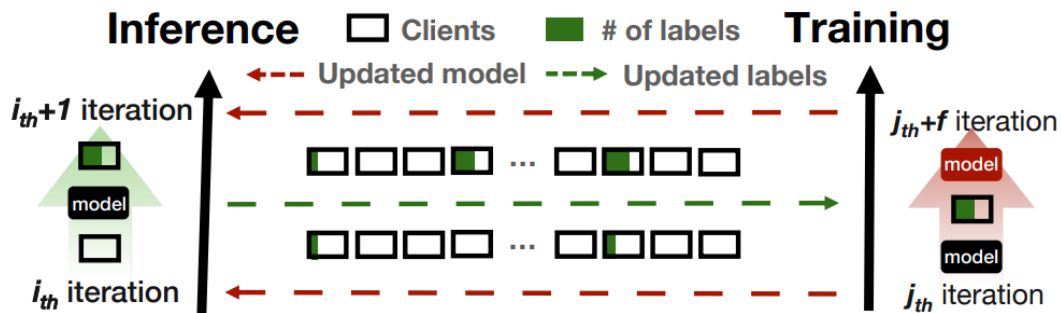University of Virginia

**Mengwei Xu**
Beiyou Shenzhen Institute

(a) Representative diversity

(b) End-to-end performance

# Summary

- **LLM+Mobile** is the new frontier, teeming with open questions that are ripe for exploration—let's pioneer the unknown!

- Awesome On-device-AI
https://github.com/ysyisyourbrother/awesome-on-device-AI
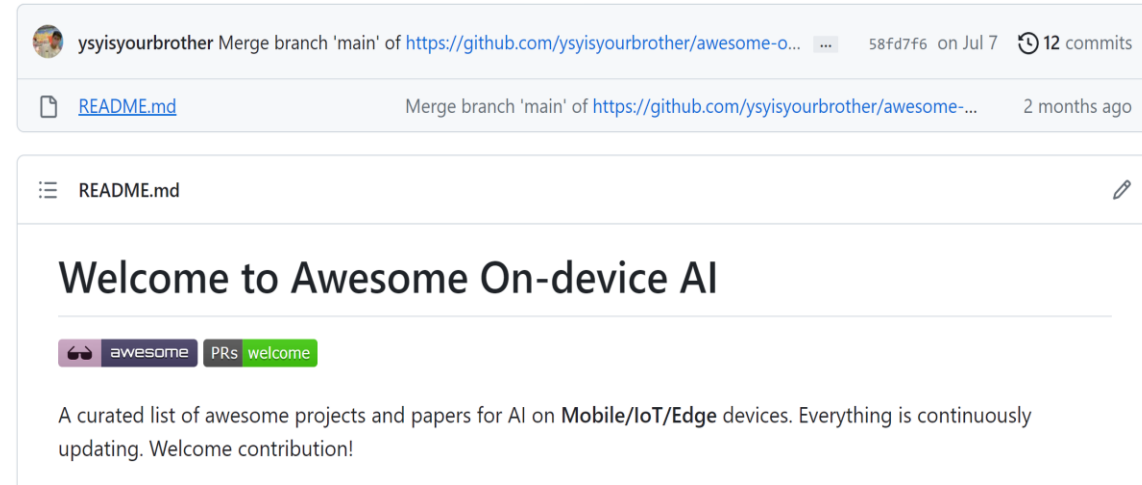
- A Reading List for Machine Learning Systems:
https://jeongseob.github.io/readings_mlsys.html

- Edge AI Paper List:
https://github.com/xumengwei/Edge-AI-Paper-List

- Resource Efficient Large Language Model
https://github.com/UbiquitousLearning/Paper-list-resource-efficient-large-language-model

# Thanks

**Shengyuan Ye**

School of Computer Science and Engineering
Sun Yat-sen University
Contact: yeshy8@mail2.sysu.edu.cn