

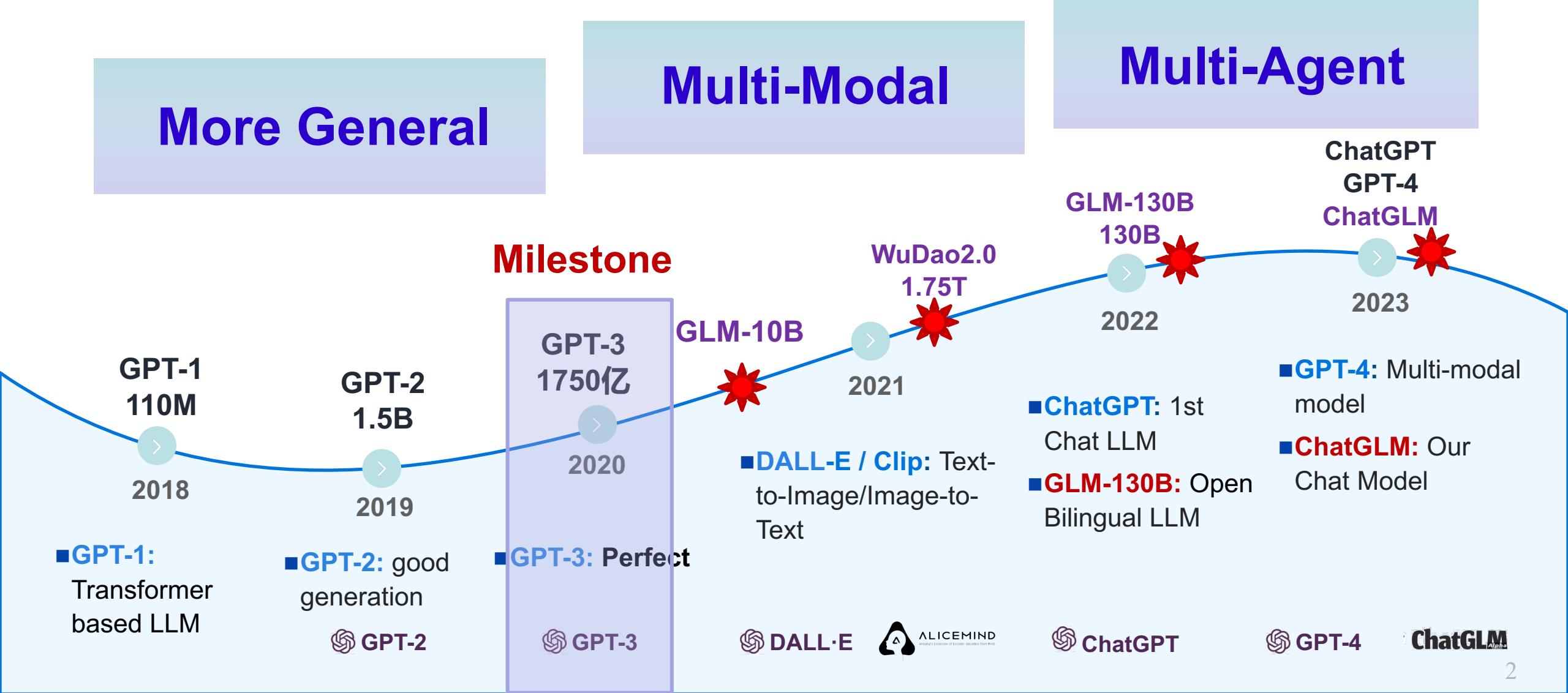


From GLM-130B to ChatGLM

Jie Tang
KEG, Tsinghua University

Open AI is leading the research on AGI

LLM becomes the foundation for AGI

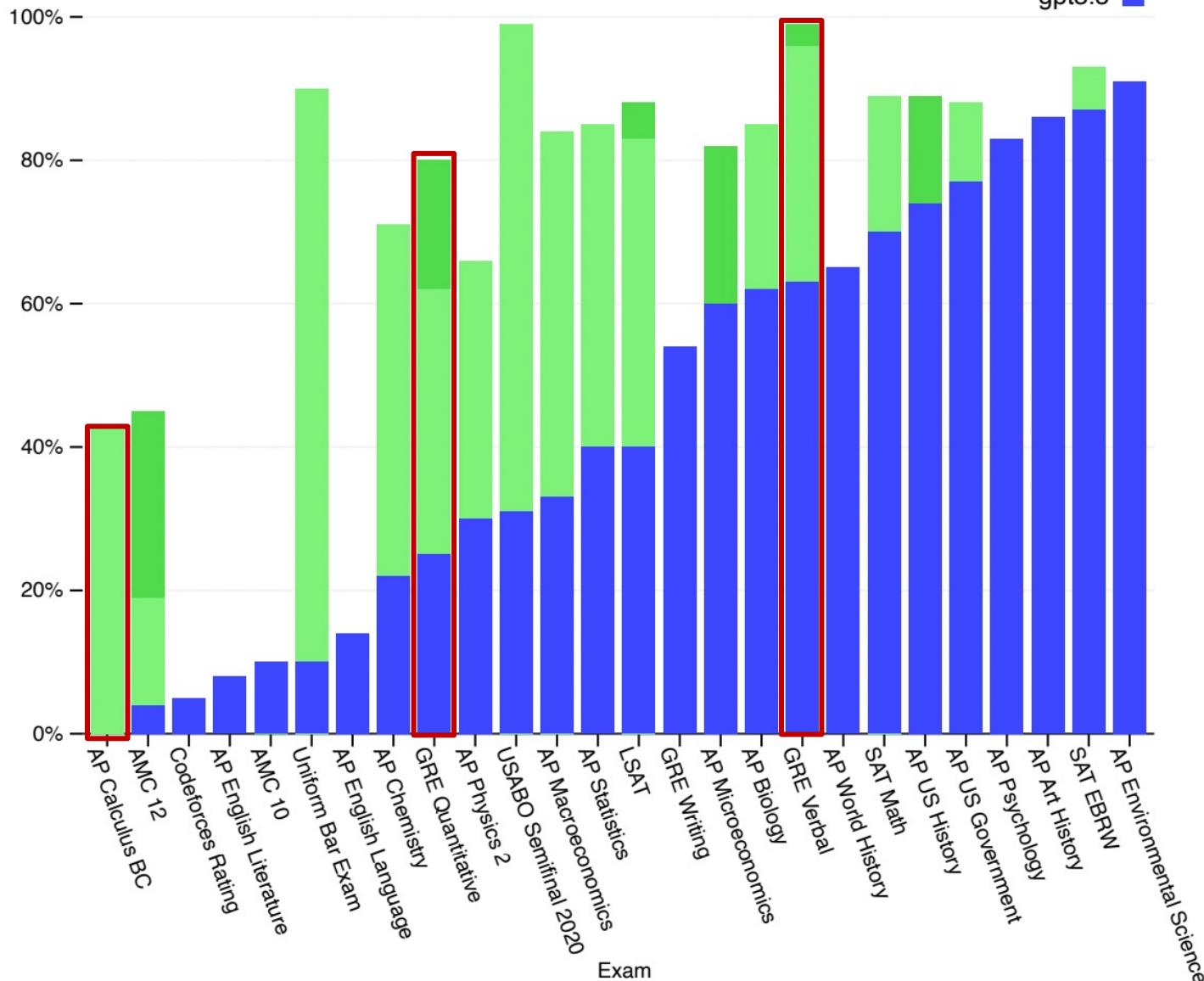


GPT-4

Exam results (ordered by GPT-3.5 performance)

Estimated percentile lower bound (among test takers)

gpt-4 (no vision)
gpt-4
gpt3.5



GPT-4

What is funny about this image? Describe it panel by panel.



GPT-4

The image shows a package for a "Lightning Cable" adapter with three panels.

Panel 1: A smartphone with a VGA connector (a large, blue, 15-pin connector typically used for computer monitors) plugged into its charging port.

Panel 2: The package for the "Lightning Cable" adapter with a picture of a VGA connector on it.

Panel 3: A close-up of the VGA connector with a small Lightning connector (used for charging iPhones and other Apple devices) at the end.

Source: <https://www.reddit.com/r/hmmmm/comments/ubab5v/hmmmm/>

GPTs



ed on the
edients you

Creative Writing Coach

I'm excited to read your work and give you feedback to improve your skills.



Laundry Buddy

Ask me anything about stains, settings, sorting and everything laundry.

Game Time

I can quickly explain board games or card games to players of any skill level. Let the games begin!



Tech Advisor

From setting up a printer to troubleshooting a device, I'm here to help you step-by-step.



kids with
sher on
ere for you.

Sticker Whiz

I'll help turn your wildest dreams into die-cut stickers, shipped to your door.



The Negotiator

I'll help you advocate for yourself and get better outcomes. Become a great negotiator.

GPTs (GPT 4.5)

- Agent-based
 - Analyze CSV
 - Goto apple.com
 - Extract wactch price
 - compare

You
Apple Stores and Prices (Asia) - Japan.c...
Spreadsheet

based on the give CSV, how much does it cost to get Apple Watch Series 9 in Japan compared to apple.com
< 2 / 2 >

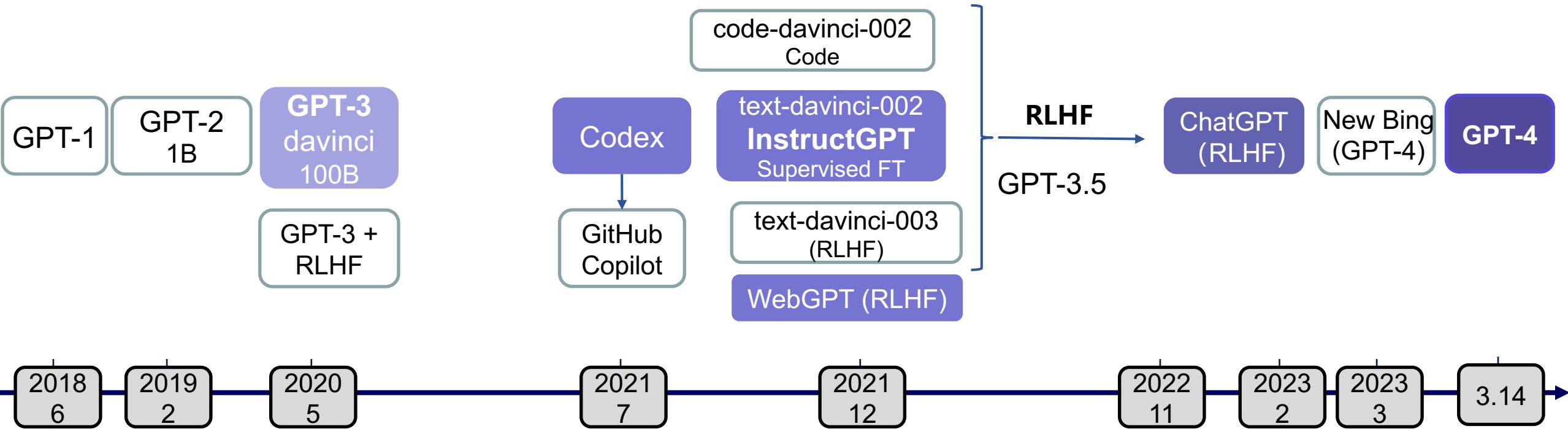
Data Analysis
The CSV file contains information on various Apple products and their prices in Japan compared to the Apple.com prices in USD. The columns include the product name, the price on Apple.com, the local price in Japan, the local price converted to USD, the net price (presumably after taxes), and the net price in USD.

To answer your question about the cost of an Apple Watch Series 9 in Japan compared to Apple.com, I would need to locate the specific row for the Apple Watch Series 9. Let me find that information for you.[\[>-\]](#)

The Apple Watch Series 9 costs \$399 on Apple.com. In Japan, the local price is ¥59,800, which converts to approximately \$396.63. The prices are very close, with the local price in Japan being slightly lower than the price on Apple.com when converted to USD.[\[>-\]](#)

🔗 🔍 🔍 🔍 🔍

GPT



1. 100B Base model

2. Supervised FT

3. RLHF

What is GLM/ChatGLM

- ChatGPT and GPT3/4/4.5 has gained enormous popularity
 - However, techniques behind GPT become a **secret** to all
- ChatGLM, an open-source ChatGPT alternative, toward unclosing the secret
 - **GLM-130B**: an open-source LLM base model
 - **ChatGLM-6B**: a lightweight open-source ChatGPT alternative
 - **ChatGLM-130B**: not open-sourced, but available through API

<https://github.com/THUDM/GLM-130B>



<https://github.com/THUDM/ChatGLM3>



ChatGLM-6B: An Open-Source Alternative

- ChatGLM-6B: **6.2B** parameters, **INT4** quantization (only need 6GB memory)
- >50,000 stars on github
- >10,000,000 downloads on huggingface
- No. 1 on Github Trending
- No. 1 on Huggingface Trending

<https://github.com/THUDM/GLM-130B>



The top 15 most-liked organizations on @huggingface

1. [@StabilityAI](#) 20k likes
2. [@AlatMeta](#) 20k
3. [@runwayml](#) 11k
4. [CompVis](#) 10k
5. [@thukgeg](#) 7k
6. [@BigscienceW](#) 7k
7. [@TIIuae](#) 7k
8. [@Microsoft](#) 6.5k
9. [@GoogleAI](#) 6k
10. [@OpenAI](#) 4k
11. [@BigCodeProject](#) 4k
12. [@MosaicML](#) 4k
13. [@UKPLab](#) 3k
14. [@AiEleuther](#) 3k
15. [@salesforce](#) 3k

huggingface.co/spaces/PulsarA...

5:40 PM · Nov 22, 2023 · 6,094 Views

>600 Open-Sourced Apps developed based on ChatGLM

Spaces using THUDM/chatglm-6b 203

HUDM/chatglm2-6b 196

es using THUDM/chatglm2-6b-int4 129

Spaces using THUDM/chatglm-6b-int4 62

thomas-yanxin/LangChain-ChatLLM

rgsgs/ChatGLM-6

ducknew/MedQA-LLM

DrSong/ChatGLM-6B-ChatBot

hahahafofo/image2text_prompt_generator

hahahafofo/

hahahafofo/ChatGLM-Chinese-Summary

hahahafofo/

OedoSoldier/chatglm_int4_demo

Dao3/ChatGLM-6B

elitecode/ChatGLM-6B-ChatBot

Nicholaspei/LangCha

alitrack/ChatPDF

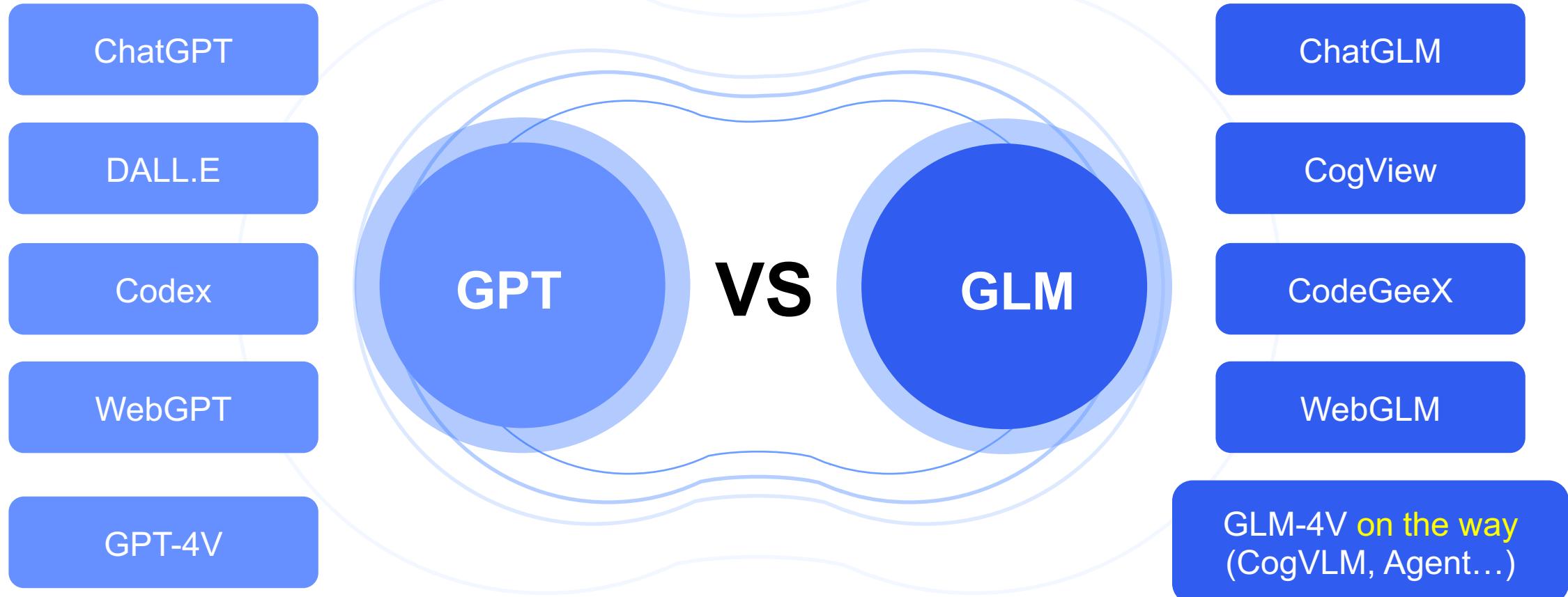
syx948/ChatPDF

actboy/Ch

<https://github.com/THUDM/ChatGLM3>



ChatGPT vs. ChatGLM



chatglm.ai

ChatGLM

GLM

XDAI

GLM-130B

CodeGeeX

QAGLM

ChatGLM

The screenshot shows a conversational interface with the following messages:

- User: who are you?
- ChatGLM: I am an AI assistant named ChatGLM, which is developed based on the language model jointly trained by Tsinghua University KEG Lab and Zhipu AI Company in 2023. My job is to provide appropriate answers and support to users' questions and requests.
- User: help me to write a blog outline which introduces ChatGLM
- ChatGLM: Sure, here's a blog outline that introduces ChatGLM:
 - I. Introduction
 - A. Brief overview of ChatGLM
 - B. Importance of AI assistants in today's world
 - II. How ChatGLM works
 - A. Basics of AI language models



Welcome to try

Story generation



Please use Chinese, English, Italian, German, Arabic tell a 100 word story about ICDM 2023



中文: ICDM 2023, 国际数据挖掘会议, 是数据挖掘、大数据和人工智能领域的国际盛会。2023年举行, 汇聚全球学者、研究人员和业界专家, 共同探讨数据挖掘领域的研究成果、技术发展和应用案例。ICDM 2023为参会者提供了交流思想、展示成果和建立合作的平台, 推动数据挖掘领域的创新与进步。

English: ICDM 2023, the International Conference on Data Mining, is a prestigious event in the fields of data mining, big data, and artificial intelligence. Scheduled for 2023, the conference brings together researchers, scholars, and industry experts from around the world to discuss the latest research findings, technological advancements, and application cases in data mining. ICDM 2023 offers attendees a platform for exchanging ideas, showcasing results, and establishing collaborations, fostering innovation and progress in the data mining field.

Italian: ICDM 2023 (Conferenza Internazionale sul Data Mining) è un importante evento internazionale nel campo del data mining, dei big data e dell'intelligenza artificiale, che si terrà nel 2023. Questo prestigioso incontro raccoglie ricercatori, studiosi e esperti dell'industria da tutto il mondo per discutere le ultime scoperte di ricerca, i progressi tecnologici e i casi d'uso nel data mining. ICDM 2023 offre ai partecipanti una piattaforma per lo scambio di idee, la presentazione dei risultati e l'istituzione di collaborazioni, con lo scopo di favorire l'innovazione e il progresso nel campo del data mining.

German: ICDM 2023, die Internationale Konferenz über Data Mining, ist ein bedeutendes Ereignis in den Bereichen Data Mining, Big Data und Künstliche Intelligenz. Die Konferenz wird im Jahr 2023 stattfinden und bringt Forscher, Gelehrte und Industriexperten aus der ganzen Welt zusammen, um die neuesten Forschungsergebnisse, technologischen Fortschritte und Anwendungsbeispiele im Data Mining zu diskutieren. ICDM 2023 bietet den Teilnehmern eine Plattform für den Ideenaustausch, die Präsentation von Ergebnissen und die Einrichtung von Kooperationen, um Innovation und Fortschritt im Data Mining zu fördern.

Arabic: ICDM 2023 هو مؤتمر دولي مهم في مجال البيانات الدولية والبيانات الكبيرة والذكاء الاصطناعي. (مؤتمر البيانات الدولي) سيعقد في عام 2023، يجمع المؤتمر الباحثين والأكاديميين وخبراء الصناعة من جميع أنحاء العالم لمناقشة الاكتشافات البحثية الأخيرة في الشراكـةـ منـقـلـةـ الأـفـلـكـ. عـرـضـ النـاقـلـةـ ICDM 2023ـ الـتـالـيـةـ النـاقـلـةـ الـأـخـرـىـ فـيـ حـالـ الـسـانـدـاتـ الـلـاـجـلـةـ بـفـيـ



Applied Math



A company produces a toy for which the variable cost is 12.30 per unit and the fixed costs are 98,000. The toy sells for \$17.98. Let X be the number of units produced and sold.

a) The total cost for a business is the sum of the variable cost and the fixed costs. Write the total cost C as a function of the number of units produced.

b) Write the revenue R as a function of the number of units sold.

c) Write the profit P as a function of the number of units sold



展开▼



a) The total cost for a business is the sum of the variable cost and the fixed costs. Write the total cost C as a function of the number of units produced:

$$C(X) = 12.30X + 98000$$

b) The revenue R is the sales price minus the variable cost. Write the revenue R as a function of the number of units sold:

$$R(X) = 17.98 - 12.30X$$

c) The profit P is the revenue R minus the total cost C . Write the profit P as a function of the number of units sold:

$$P(X) = R(X) - C(X) = 17.98 - 12.30X - 98000$$

Note: The negative sign indicates that the profit is a function of the number of units sold in the negative direction, meaning that as the number of units sold increases, the profit decreases.

Coding



Write a code for a bubble sort algorithm in Python, and give an array as an example



Sure, here's a simple bubble sort algorithm implemented in Python:

```
def bubble_sort(arr):
    n = len(arr)
    for i in range(n):
        for j in range(n-i-1):
            if arr[j] > arr[j+1]:
                # Swap adjacent elements if they are in wrong order
                arr[j], arr[j+1] = arr[j+1], arr[j]
    return arr
```

To use this function, you'll need to define an array to sort, like so:

```
arr = [64, 34, 25, 12, 22, 11, 90]
```

Once you have your array defined, you can call the `bubble_sort()` function on it, like so:

```
sorted_arr = bubble_sort(arr)
```

Finally, you can print the sorted array like so:

```
print(sorted_arr)
```

This should output something like:

```
[11, 12, 22, 25, 34, 64, 90]
```

Let's take a look at an example array to better understand how the bubble sort algorithm works:

```
arr = [5, 2, 9, 1, 5, 6, 10]
```

In this example, the first element (5) is sorted towards the bottom of the array, and the last element (10) is sorted towards the top. The middle elements (2, 5, and 9) are sorted in between their adjacent elements, as shown in the following sorted array:

```
[1, 2, 5, 5, 6, 9, 10]
```

Overall, the bubble sort algorithm works by iterating through the array, comparing adjacent elements, and swapping them if they are in the wrong order. This process repeats until all elements in the array are sorted in increasing order.

① 以上内容为 AI 生成，不代表开发者立场，请勿删除或修改本标记



GLM-4V-0.5



who are in this picture



In this picture, there are three football players. From left to right, they are Suárez, Neymar, and Messi.

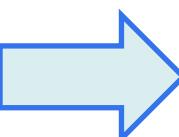
以上内容为 AI 生成，不代表开发者立场，请勿删除或修改本标记



01:35

5G

GLM-4V-0.5



小智

<

+

晚上好🌙

夜深了，快节奏的生活更要注意休息哦~小智一直都在



What is this

This is a name tag from the IEEE International Conference on Data Mining (ICDM) held in Shanghai, China. The person in the picture is Jie Tang, and he is a Tsinghua University student and the conference's keynote speaker.

“draw a dog with a hat”

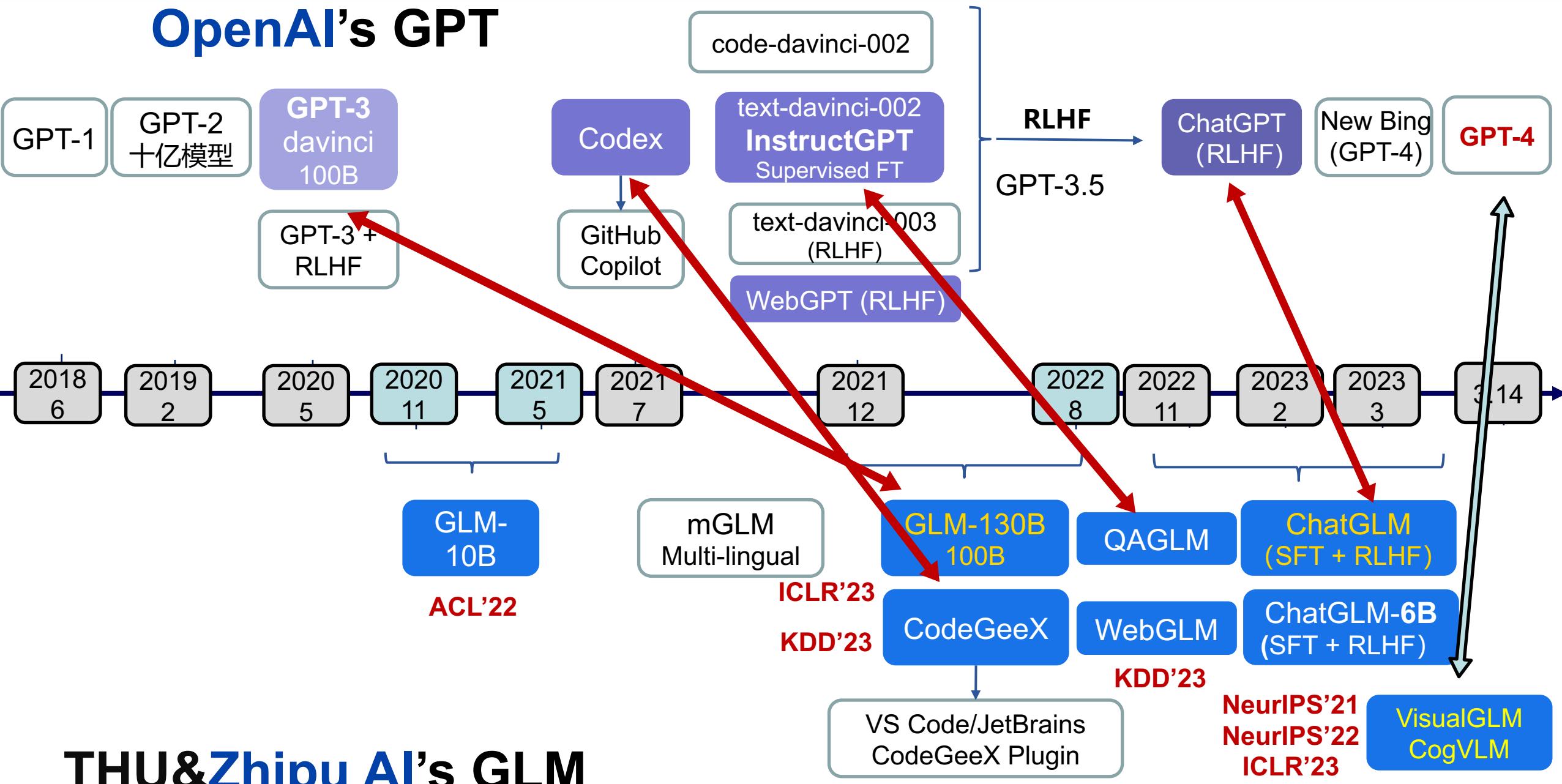


AI生成

KG engine ready...



OpenAI's GPT



THU&Zhipu AI's GLM

General Language Model (GLM)

Framework	NLU	Cond. Gen.	Uncond. Gen.
Autoregressive (GPT)	—	—	✓
Autoencoding (BERT)	✓	✗	✗
Encoder-Decoder (T5)	—	✓	—
Autoregressive Blank-Infilling (GLM)	✓	✓	✓

General Language Model (GLM)

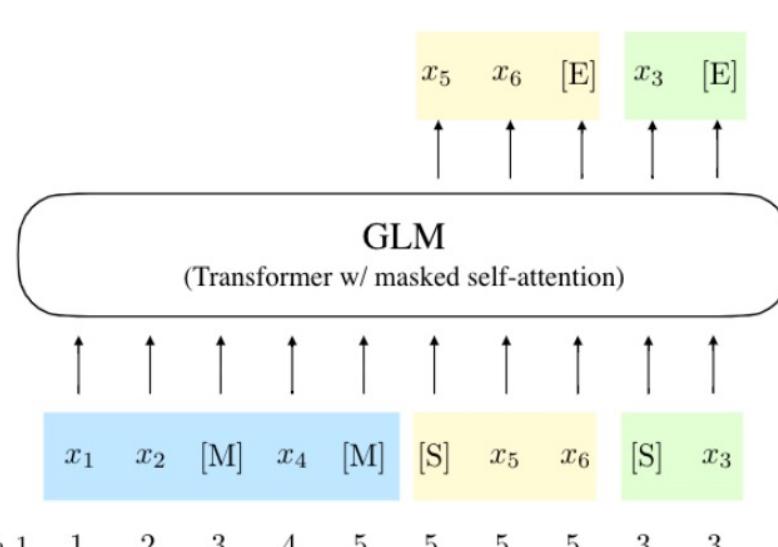
$x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5 \quad x_6$

(a) Sample spans from the input text

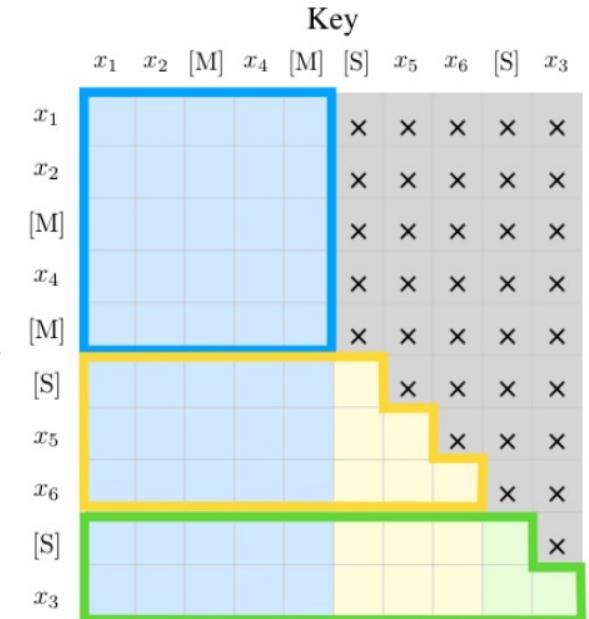
Part A: $x_1 \quad x_2 \quad [M] \quad x_4 \quad [M]$

Part B: $x_5 \quad x_6 \quad x_3$

(b) Divide the input into Part A and Part B



(c) Generate the Part B spans autoregressively

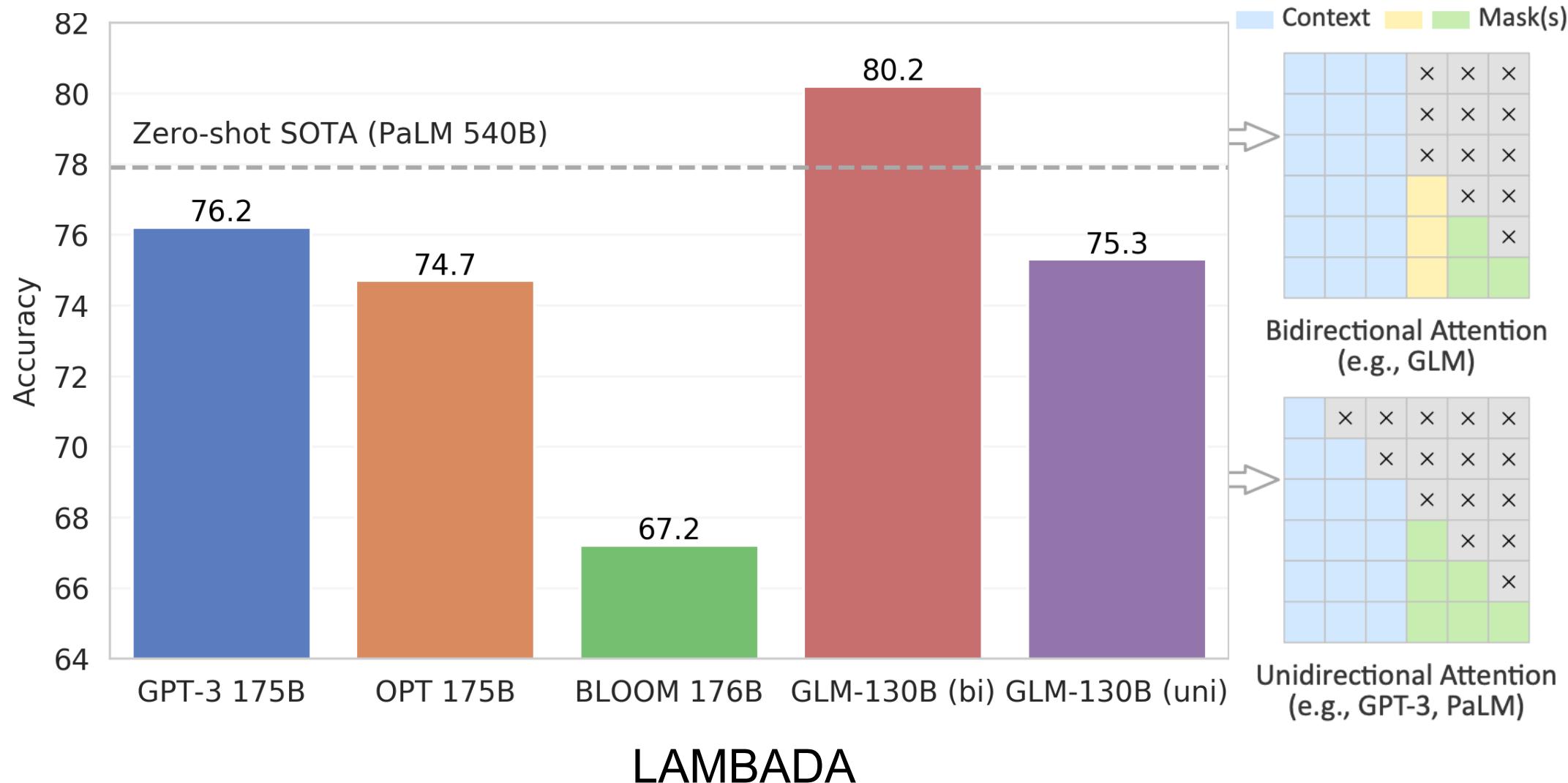


(d) Self-attention mask

$$\mathcal{L}_{\text{GLM}} = \mathbb{E}_{\mathbf{z} \sim Z_m} \left[\sum_{i=1}^m \sum_{j=1}^{l_i} -\log p(s_{z_i, j} | \mathbf{x}_{\text{corrupt}}, \mathbf{s}_{z_{<i}}, s_{z_i, <j}) \right]$$

General Language Model (GLM)

General Language Model (GLM)



Results on Natural Language Understanding

- **Better than BERT, T5, RoBERTa**

Table 2. Results on the SuperGLUE dev set. Models with * are pre-trained for two times the number of steps of other methods.

Model	ReCoRD F1/Acc.	COPA Acc.	WSC Acc.	RTE Acc.	BoolQ Acc.	WiC Acc.	CB F1/Acc.	MultiRC F1a/EM	Avg
BERT _{Base}	65.4/64.9	66.0	65.4	70.0	74.9	68.8	70.9/76.8	68.4/21.5	66.1
GLM _{Base}	73.5/72.8	71.0	72.1	71.2	77.0	64.7	89.5/85.7	72.1/26.1	70.7
BERT _{Large}	76.3/75.6	69.0	64.4	73.6	80.1	71.0	94.8/92.9	71.9/24.1	72.0
UniLM _{Large}	80.0/79.1	72.0	65.4	76.5	80.5	69.7	91.0/91.1	77.2/38.2	74.1
GLM _{Large}	81.7/81.1	76.0	81.7	74.0	82.1	68.5	96.1/94.6	77.1/36.3	77.0
GLM _{Large} (multi-task)	80.2/79.6	77.0	78.8	76.2	79.8	63.6	97.3/96.4	74.6/32.1	75.7
GLM _{410M} (multi-task)	81.5/80.9	80.0	81.7	79.4	81.9	69.0	93.2/96.4	76.2/35.5	78.0
GLM _{515M} (multi-task)	82.3/81.7	85.0	81.7	79.1	81.3	69.4	95.0/96.4	77.2/35.0	78.8
T5 _{Base}	76.2/75.4	73.0	79.8	78.3	80.8	67.9	94.8/92.9	76.4/40.0	76.0
T5 _{Large}	85.7/85.0	78.0	84.6	84.8	84.3	71.6	96.4/98.2	80.9/46.6	81.2
BART _{Large} *	88.3/87.8	60.0	65.4	84.5	84.3	69.0	90.5/92.9	81.8/48.0	76.0
RoBERTa _{Large} *	89.0/88.4	90.0	63.5	87.0	86.1	72.6	96.1/94.6	84.4/52.9	81.5
GLM _{RoBERTa}	89.6/89.0	82.0	83.7	87.7	84.7	71.2	98.7/98.2	82.4/50.1	82.9

Results on Generation

- The most important thing is that **one model** can do all the things

Table 3. Results on Gigaword abstractive summarization

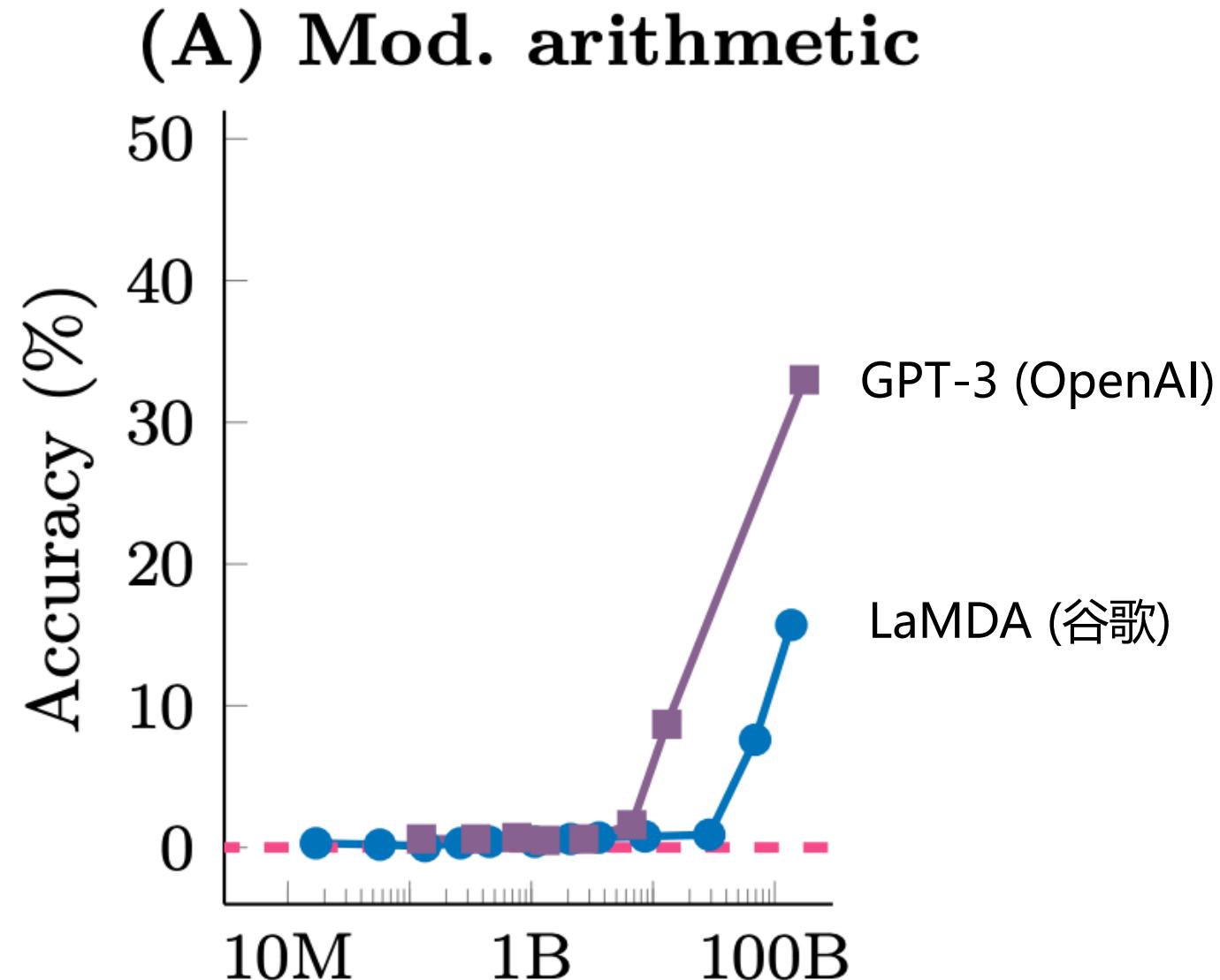
Model	RG-1	RG-2	RG-L
MASS	37.7	18.5	34.9
UniLM _{Large}	38.5	19.5	35.8
GLM _{Large}	38.6	19.7	36.0
GLM _{Large} (multi-task)	38.5	19.4	35.8
GLM _{410M} (multi-task)	38.9	20.0	36.2

Table 4. Zero-shot language modeling results.

Model	Lambada (Accuracy)	BookWiki (Perplexity)
GLM _{Large} (uni)	0.0	> 100
GLM _{Large} (multi-task,uni) – 2d positional encoding	47.4	15.1
GLM _{410M} (multi-task,uni)	45.8	15.1
GLM _{515M} (multi-task,uni)	49.5	14.5
GLM _{515M} (multi-task,uni)	50.4	13.9
GLM _{Large} (bi)	10.6	> 100
GLM _{Large} (multi-task,bi) – 2d positional encoding	48.5	14.9
GLM _{410M} (multi-task,bi)	47.3	15.0
GLM _{515M} (multi-task,bi)	53.5	14.3
GLM _{515M} (multi-task,bi)	54.9	13.7
GPT _{Large} (uni)	50.1	14.4

Why 100B-scale model?

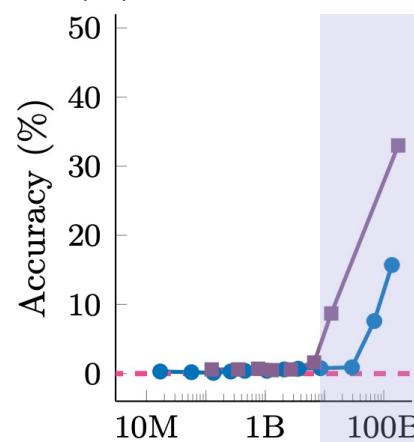
- What is $16 \bmod 12$?
- 16 divided by 12 equals 1 remainder 4 . So the answer is 4 !



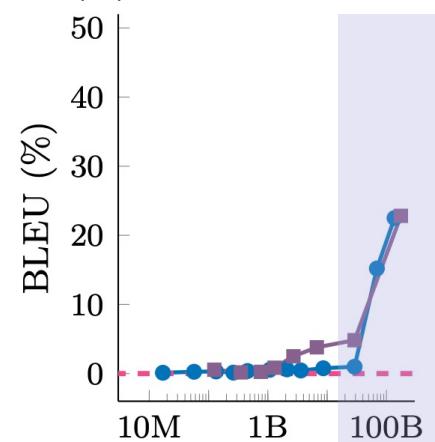
Why 100B-scale model?

—●— LaMDA —■— GPT-3 —◆— Gopher —▲— Chinchilla —◆— PaLM —--- Random

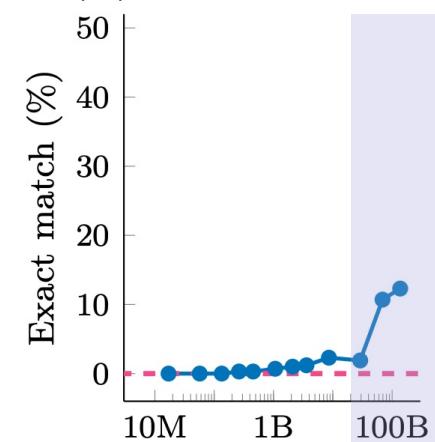
(A) Mod. arithmetic



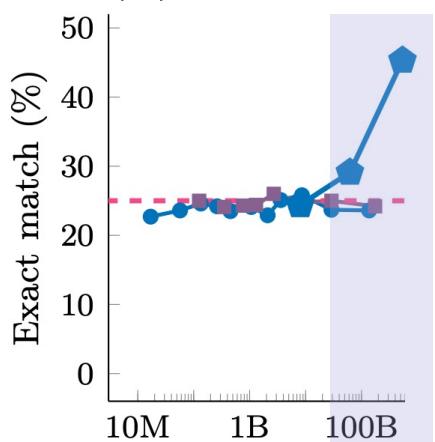
(B) IPA transliterate



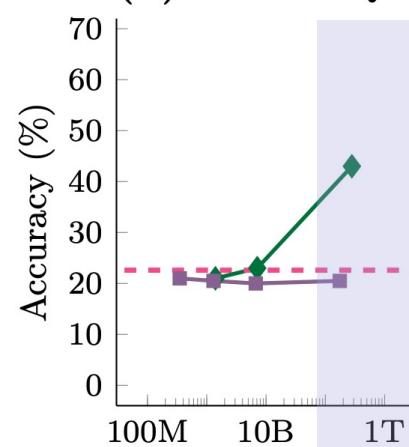
(C) Word unscramble



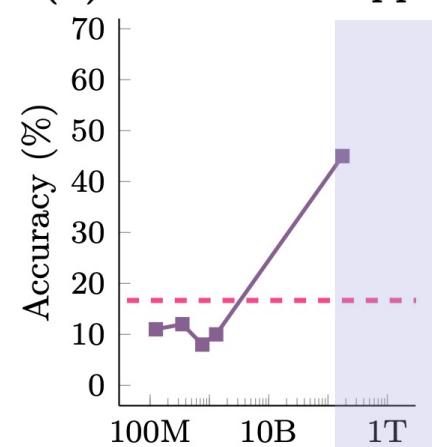
(D) Persian QA



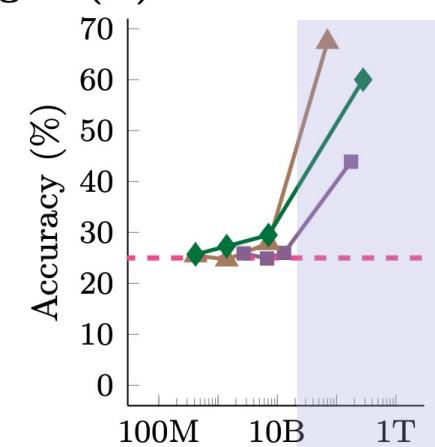
(E) TruthfulQA



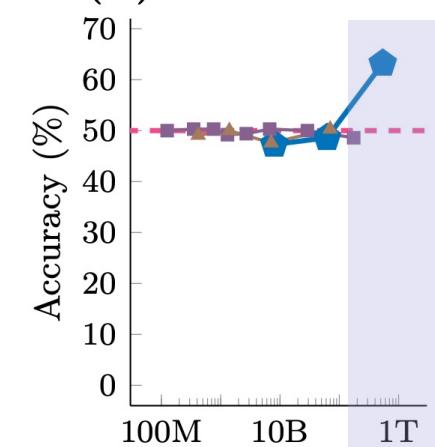
(F) Grounded mappings



(G) Multi-task NLU

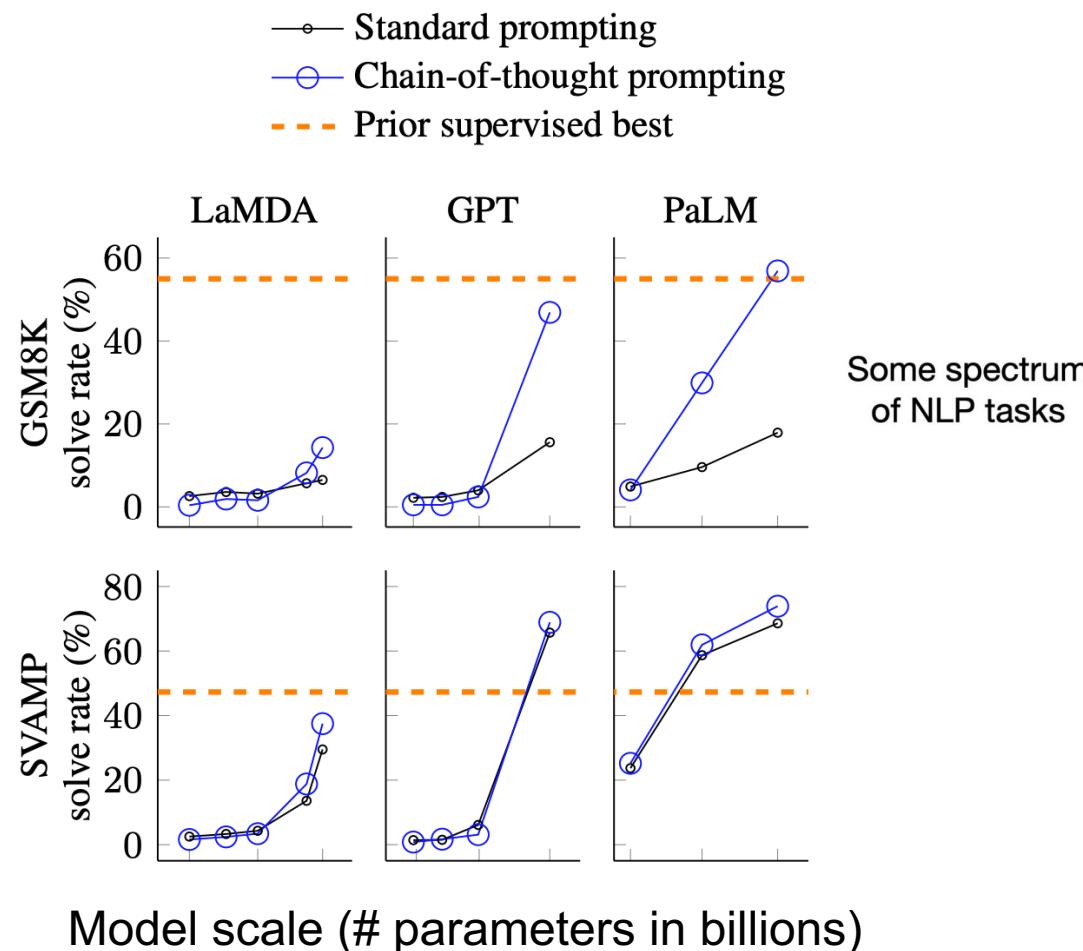


(H) Word in context

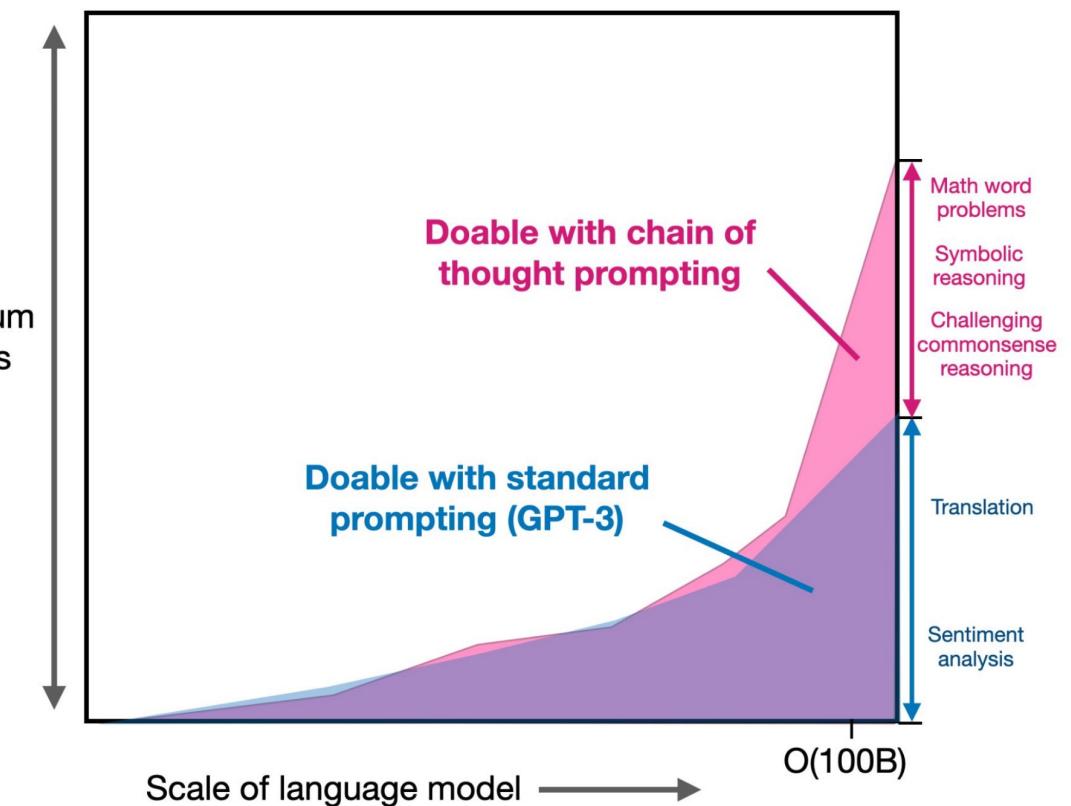


Scaling Law

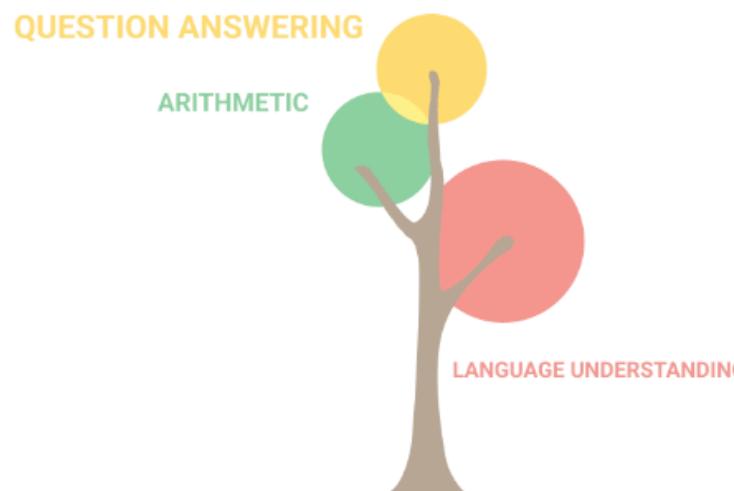
Scaling Law introduces complicated reasoning abilities



Some spectrum
of NLP tasks



“Emergent abilities”



8 billion parameters

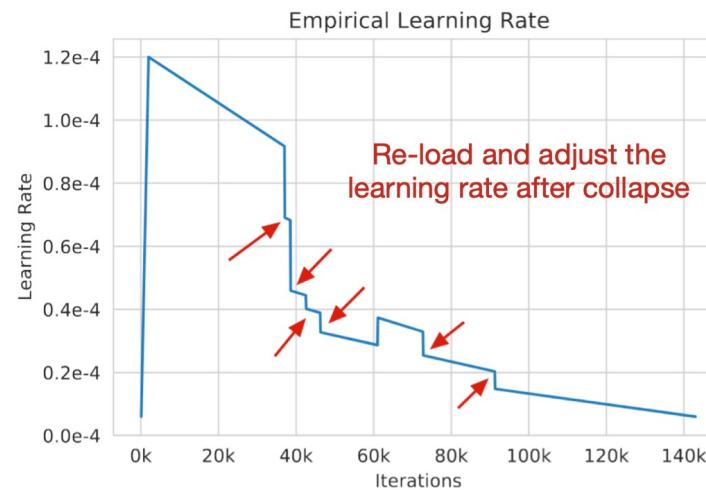
How to train a 100B-scale LLM?

- 8 months have witnessed numerous challenges
 - **Engineering:** How to train 100B-scale models from scratch?
 - Hygon DCU, NVIDIA A100, Ascend 910, Sunway
 - Frequent & random hardware failures, Megatron-DeepSpeed 3D pipeline, CUDA kernel efficiency, GPU memory overflow, 10K+ threads TCP init & comms...
 - **Algorithm:** How to stabilize the training of 100B-scale models?
 - The gradient norms of embeddings, Post-LN / Pre-LN stability, dataloader state seeds, computation precision in Softmax / Attention

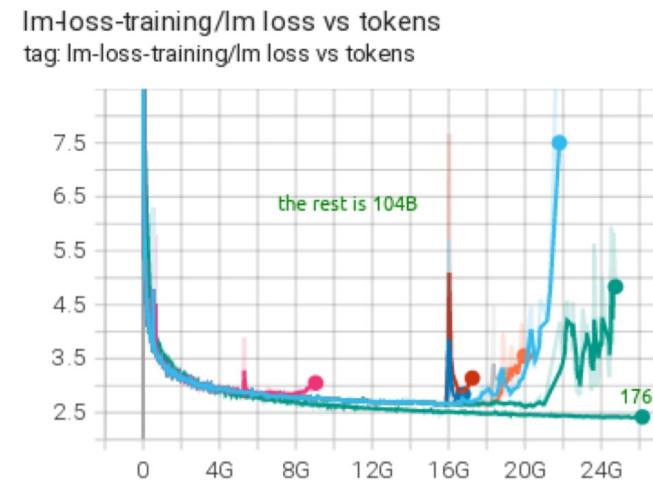


Training Stability of 100B-Scale Models

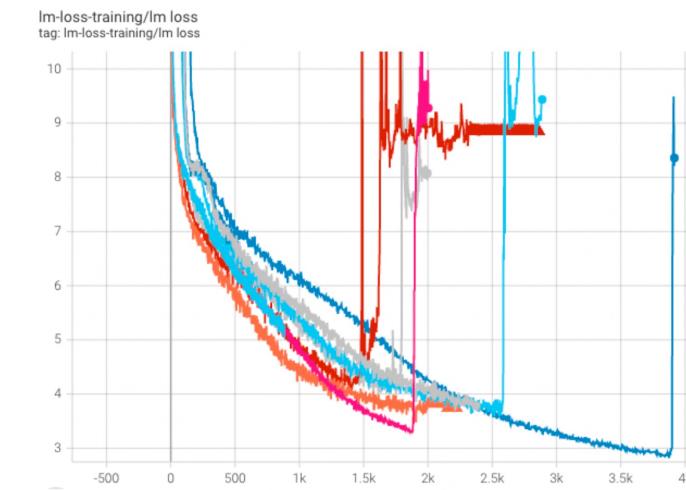
- Tradeoff: Stability (Slow) or Efficiency (Instable)
- Existing Solutions
 - **OPT-175B:** manually adjust LR & skip data when collapses (performance drop)
 - **BLOOM 176B:** embedding norm & BF16 (performance drop, few platform)



(a) OPT 175B's experiments



(b) BLOOM 176B's experiments



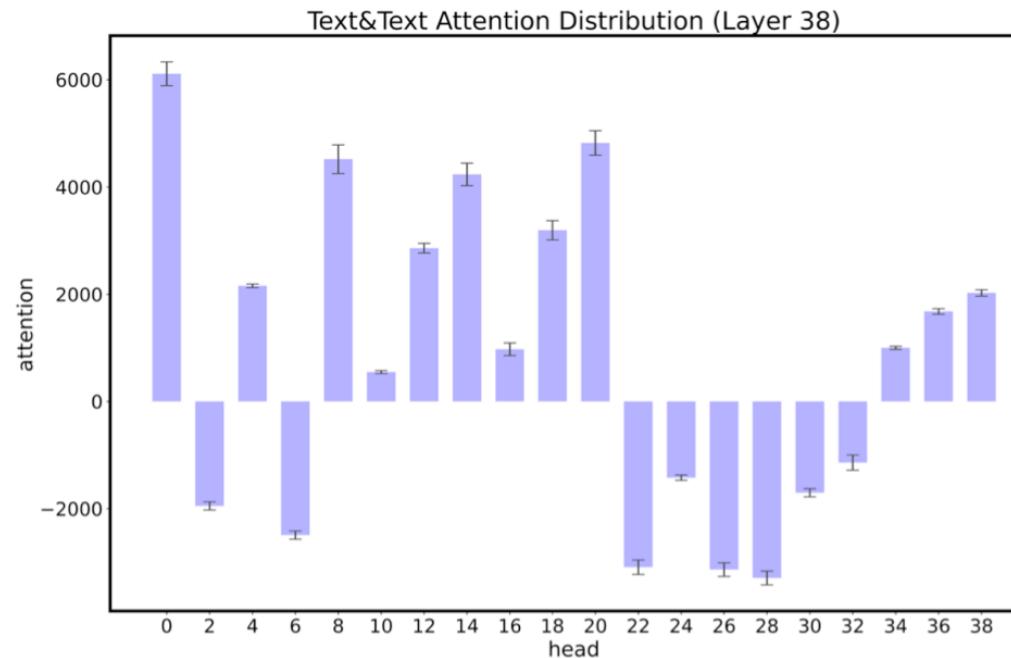
(c) GLM 130B's experiments

GLM-130B: Training Stability

- Attention score: Softmax in 32 to avoid overflow

$$\text{softmax}\left(\frac{\mathbf{Q}_i \mathbf{K}_i^\top}{\sqrt{d}}\right) = \text{softmax}\left(\left(\frac{\mathbf{Q}_i \mathbf{K}_i^\top}{\alpha\sqrt{d}} - \max\left(\frac{\mathbf{Q}_i \mathbf{K}_i^\top}{\alpha\sqrt{d}}\right)\right) \times \alpha\right) = \text{FP16}\left(\text{softmax}\left(\text{FP32}\left(\frac{\mathbf{Q}_i \mathbf{K}_i^\top}{\alpha\sqrt{d}}\right) \times \alpha\right)\right)$$

Attention scores grow large --- exceeding FP16's range

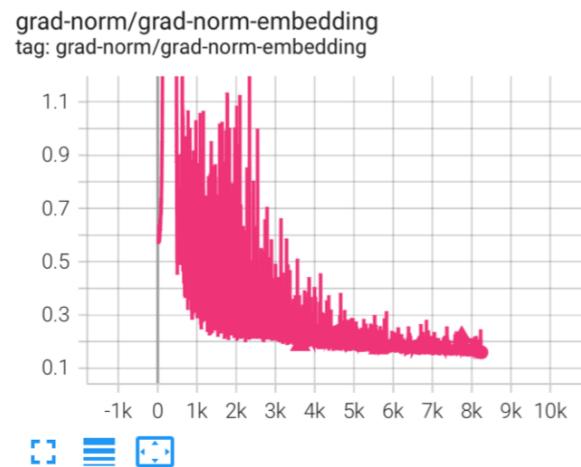


GLM-130B: Training Stability

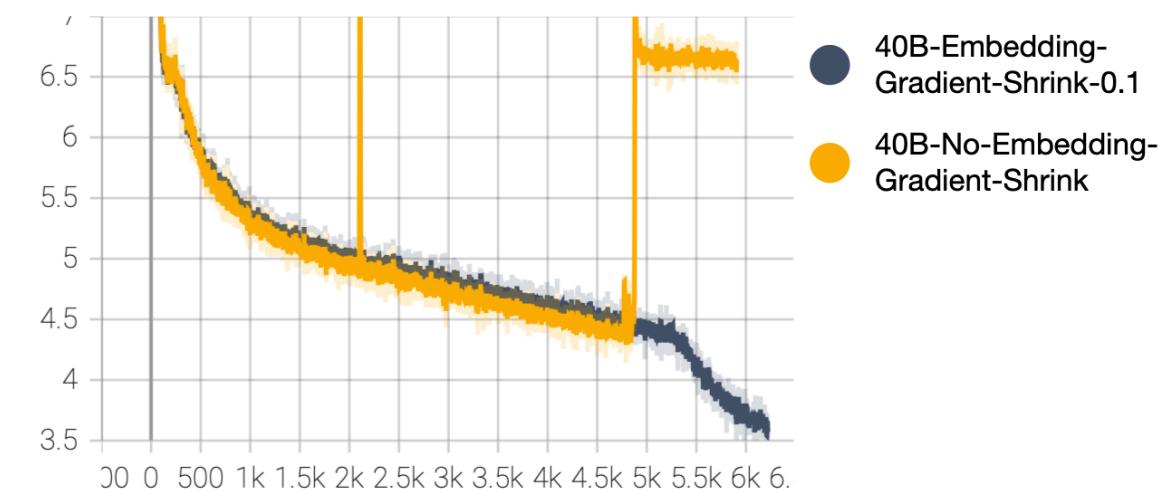
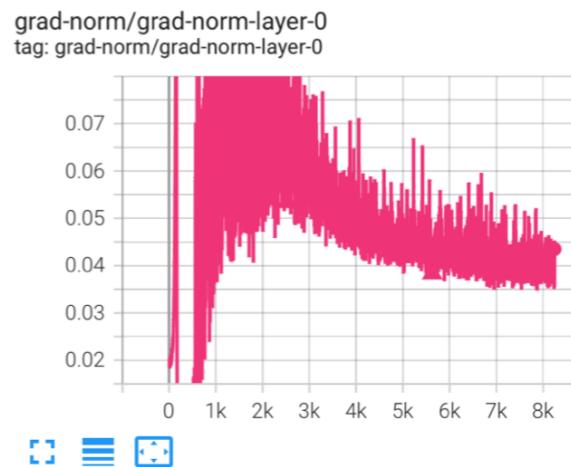
□ Embedding Layer Gradient Shrink (EGS)

```
word_embedding = word_embedding * alpha +  
                 word_embedding.detach() * (1 - alpha)
```

Embedding Layer gradients can be magnitudes larger than others



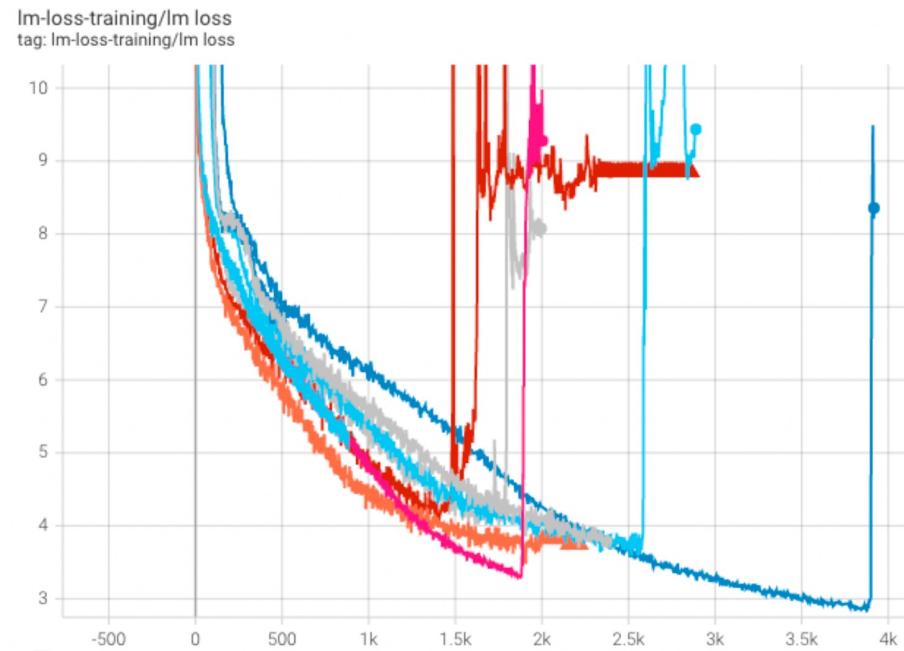
(a) Gradient norm of embedding layer (left) and the first layer (right)



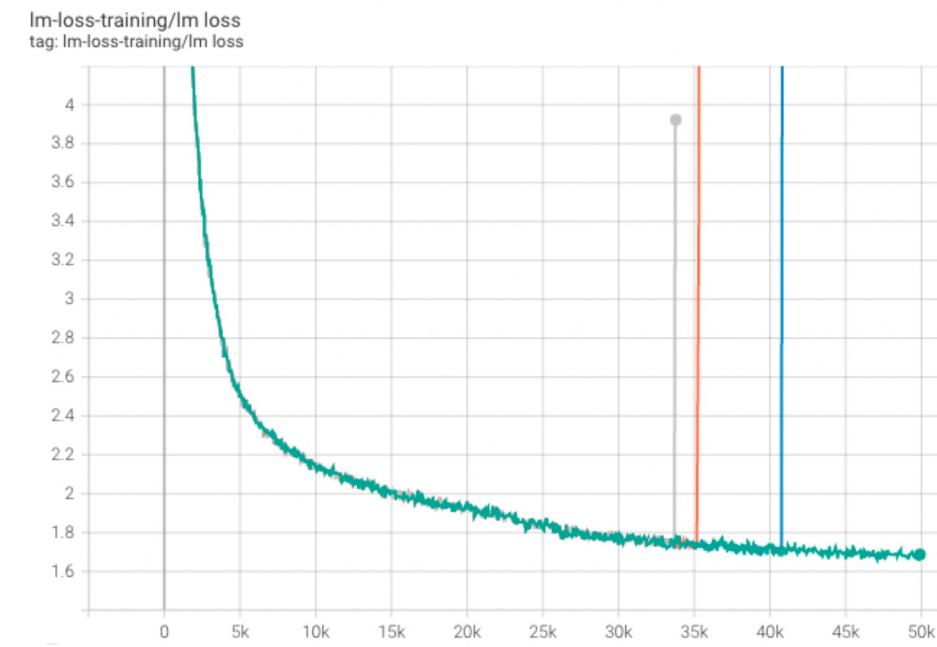
(b) Training loss curves of GLM-40B with and without gradient shrink

GLM-130B: Training Stability

□ The final training run of GLM-130B



(c) GLM 130B's experiments



(d) GLM 130B's real training

GLM-130B Training Lessons

2021.12

- The “千亿” (100B) project towards an open dense pre-trained GLM at 100B scale is conceived
- Survey pre-training strategies of existing models of similar scale, such as GPT-3, Gopher => *Limited public info about how they were trained and issues they met*
- Search for possible GPU clusters & sponsors

2022.1

- Test the performance of FP16/FP32 at 100B scale on one testing cluster
- Unexpected excessive memory usage in GLM => *Torch is better with fixed length input sequences*
- Inability to converge and try tricks from CogView and ViT => *Use Sandwich-LN*
- Frequent random hardware failures => *Have to run HCPG test before each run*

2022.2

- Very slow training speed than previously calculated => *Optimize kernels and fuse operators* => *Find the input shape is critical to kernel performance*
- Collect pre-training corpora and tokenize => *Use icetk: the sentence piece is set to the unigram mode*
- Debug the 3D pipeline parallel in the newly-released Megatron and DeepSpeed

2022.3

- It can't recover perfectly from optimizer states => *Our customized dataloaders do not save its state seed properly in distributed training*
- The memory per processor is too small => *Require too many pipeline stages* => *Batch size is too large (up to 12,000) => Harm the model's convergenc*
- It can't launch more than 2,000 computing nodes => *Overcome this and support 6,000-node training by tuning Linux kernel TCP parameters*
- Collect data for multi-task instruction pre-training
- Receive opportunities to test trainings on several other clusters
- Very slow training speed than expected => *The underlying element-wise operators don't support fast computation on large-dimension vectors.*

2022.4

- Optimize A100 kernel's computing efficiency => *A100 kernels prefer square-shaped inputs, and seq_len=2,048 is optimal for our hidden-state dimension (12,288)*
- Inability to converge due to large gradient norms (170+) of input embeddings => *Try embedding norm and gradient shrink, which turn out to be almost equivalent*
- Naïve post-LN or pre-LN diverges after several thousands of steps => *Try Sandwich-LN with PB-Relax*
- It still diverges after one week's trial => *The dataloader state seeds are not unified for different pipeline stages, resulting in a mismatch of input data and labels.*
- Test two positional encodings: RoPE and Alibi => *Alibi can be slower as it requires element-wise manipulation on attention matrices--changing num_heads *2,048 * 2,048 scalars per layer*
- Test GeGLU and GAU => *GAU converges faster with relatively poor performance on fine-tuned SuperGLUE*
- Abnormal GPU memory usage of newly-added functions and classes => *DeepSpeed hardcodes the function names for checkpoint activation*
- Decode to train GLM with 130 billion parameters => *allow inference on a DGX-A100 40G node*

2022.5-6

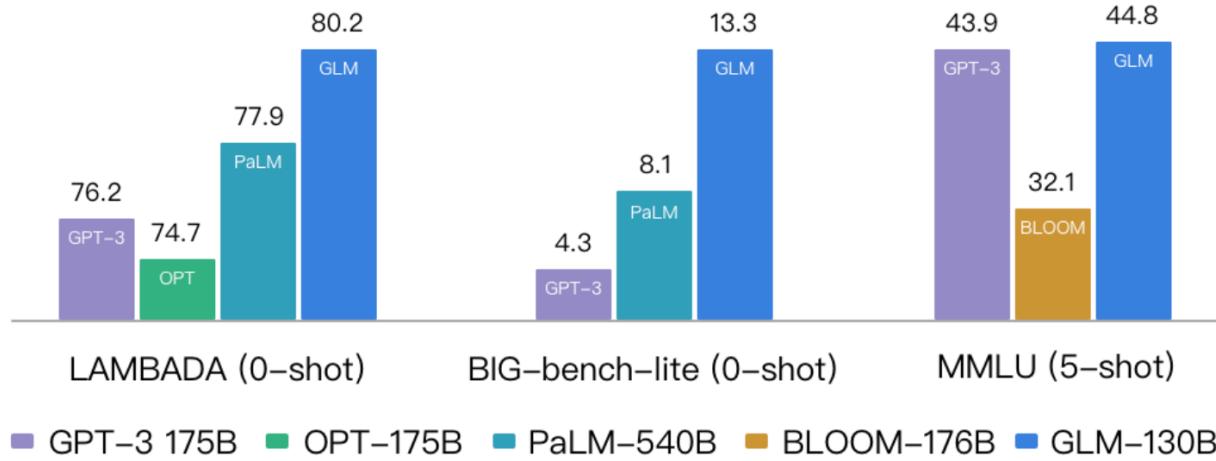
- Implement a RoPE cuda operator in C++ => *See unexpected precision errors and finally have it abandoned*
- Sandwich-LN still diverges => *1) Reducing learning rate does not help; 2) Using Hinge cross-entropy becomes slower and harms performance; 3) Shifting to DeepNorm still diverges*
- Use FP32 in softmax of attention => *Success*
- Find PB-Relax unnecessary for FP32 softmax => *It also slows down training as it needs to manipulate the whole attention score matrices*
- Experience few spikes in later training => *1) Reduce gradient shrink factor from 1 to 0.1: useful; 2) Reduce the learning rate: sometimes useful; 3) Jump the noisy data batches: sometimes useful*
- Find a mistake in multi-task data after training for 20,000 steps => *Use the correct data but it does not forget*

2022.6-7

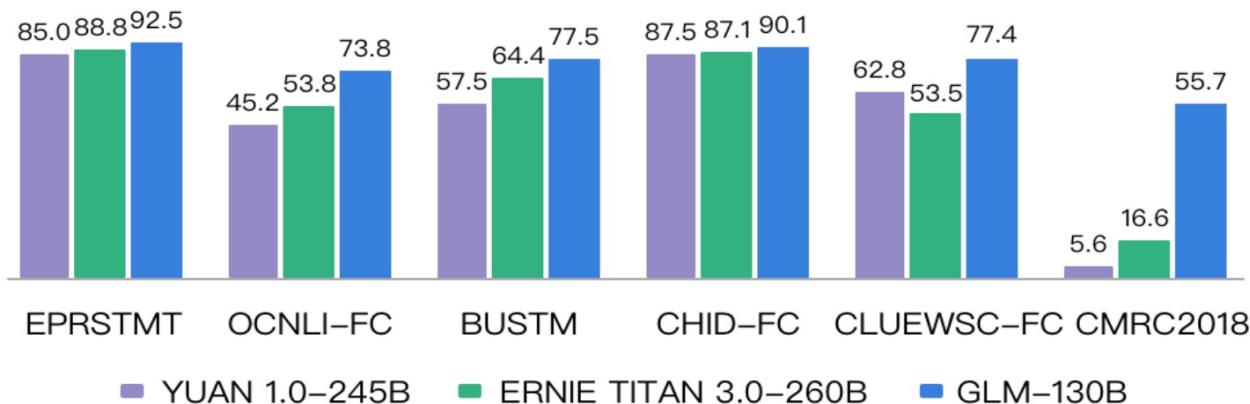
- Adapt the pipeline parallel checkpoints to ordinary parallel checkpoints for efficient inference on a single A100
- Work on evaluation scripts on datasets: MMLU, Big-bench, CLUE, SuperCLUE, etc.
- Implement P-Tuning and P-Tuning v2 for parameter-efficient tuning on GLM-130B for tuning on SuperGLUE
- Work with BMInf on adapting GLM-130B to perform inference on a single V100 or 3090 => *Use pipeline-style asynchronous swapping between main memory and GPU memory*
- Try to fine-tune GLM-130B with fewer A100 nodes (i.e., 12-16 nodes) => *Pipeline-style fails due to too many pipeline stages => Find that data parallel can not be introduced for fine-tuning => Use 32-way model parallel for fine-tuning with reasonable performance*

GLM-130B

English: better than GPT-3/OPT/PaLM on MMLU, LAMBADA, BIG-bench-lite



Chinese: better than ERNIE 260B & YUAN 245B



Aug., 2022-Mar. 2023, research use requests from ~1000 orgs in 70 countries

- Google
- Huawei
- Microsoft
- Alibaba
- Facebook
- Tencent
- Stanford
- Baidu
- MIT
- Meituan
- UC Berkely
- Bytedance
- CMU
- Didi
- Harvard
- Xiaoice
- Princeton
- Xiaodu
- Yale
- Xiaomi
- Cornell
- Xiaopeng
- UIUC
- Youdao
- Cambridge
- Face++
- Oxford
- Ping An Cap
- Peking U.
- Zhejiang U.
- Shanghai JT U.
- Fudan U.
- USTC
- U of CAS
- Wuhan U.
- Naikai U.
- Hongkong U.
- CUHK
- HKUST
- BAAI
- Zhejiang Lab
- Shanghai AI Lab

GLM-130B in HELM

Stanford's Holistic Evaluation of Language Models (HELM, Nov. 2022)

Model	Model Creator	Modality	# Parameters	Tokenizer	Window Size	Access	Total Tokens	Total Queries	Total Cost
J1-Jumbo v1 (178B)	AI21 Labs	Text	178B	AI21	2047	limited	327,443,515	591,384	\$10,926
J1-Grande v1 (17B)	AI21 Labs	Text	17B	AI21	2047	limited	326,815,150	591,384	\$2,973
J1-Large v1 (7.5B)	AI21 Labs								
Anthropic-LM v4-s3 (52B)	Anthropic								
BLOOM (176B)	BigScience								
T0++ (11B)	BigScience								
Cohere xlarge v20220609 (52.4B)	Cohere								
Cohere large v20220720 (13.1B) ⁵⁸	Cohere								
Cohere medium v20220720 (6.1B)	Cohere								
Cohere small v20220720 (410M) ⁵⁹	Cohere								
GPT-J (6B)	EleutherAI								
GPT-NeoX (20B)	EleutherAI								
T5 (11B)	Google								
UL2 (20B)	Google								
OPT (66B)	Meta								
OPT (175B)	Meta								
TNLG v2 (6.7B)	Microsoft/NVIDIA								
TNLG v2 (530B)	Microsoft/NVIDIA								
GPT-3 davinci v1 (175B)	OpenAI								
GPT-3 curie v1 (6.7B)	OpenAI								
GPT-3 babbage v1 (1.3B)	OpenAI	Text	1.3B	GPT-2	2048	limited	422,123,900	606,253	\$211
GPT-3 ada v1 (350M)	OpenAI	Text	350M	GPT-2	2048	limited	422,635,705	604,253	\$169
InstructGPT davinci v2 (175B*)	OpenAI	Text	175B*	GPT-2	4000	limited	466,872,228	599,815	\$9,337
InstructGPT curie v1 (6.7B*)	OpenAI	Text	6.7B*	GPT-2	2048	limited	420,004,477	606,253	\$840
InstructGPT babbage v1 (1.3B*)	OpenAI	Text	1.3B*	GPT-2	2048	limited	419,036,038	604,253	\$210
InstructGPT ada v1 (350M*)	OpenAI	Text	350M*	GPT-2	2048	limited	418,915,281	604,253	\$168
Codex davinci v2	OpenAI	Code	Unknown	GPT-2	4000	limited	46,272,590	57,051	\$925
Codex cushman v1	OpenAI	Code	Unknown	GPT-2	2048	limited	42,659,399	59,751	\$85
GLM (130B)	Tsinghua University	Text	130B	ICE	2048	open	375,474,243	406,072	2,100 GPU hours
YaLM (100B)	Yandex	Text	100B	Yandex	2048	open	378,607,292	405,093	2,200 GPU hours











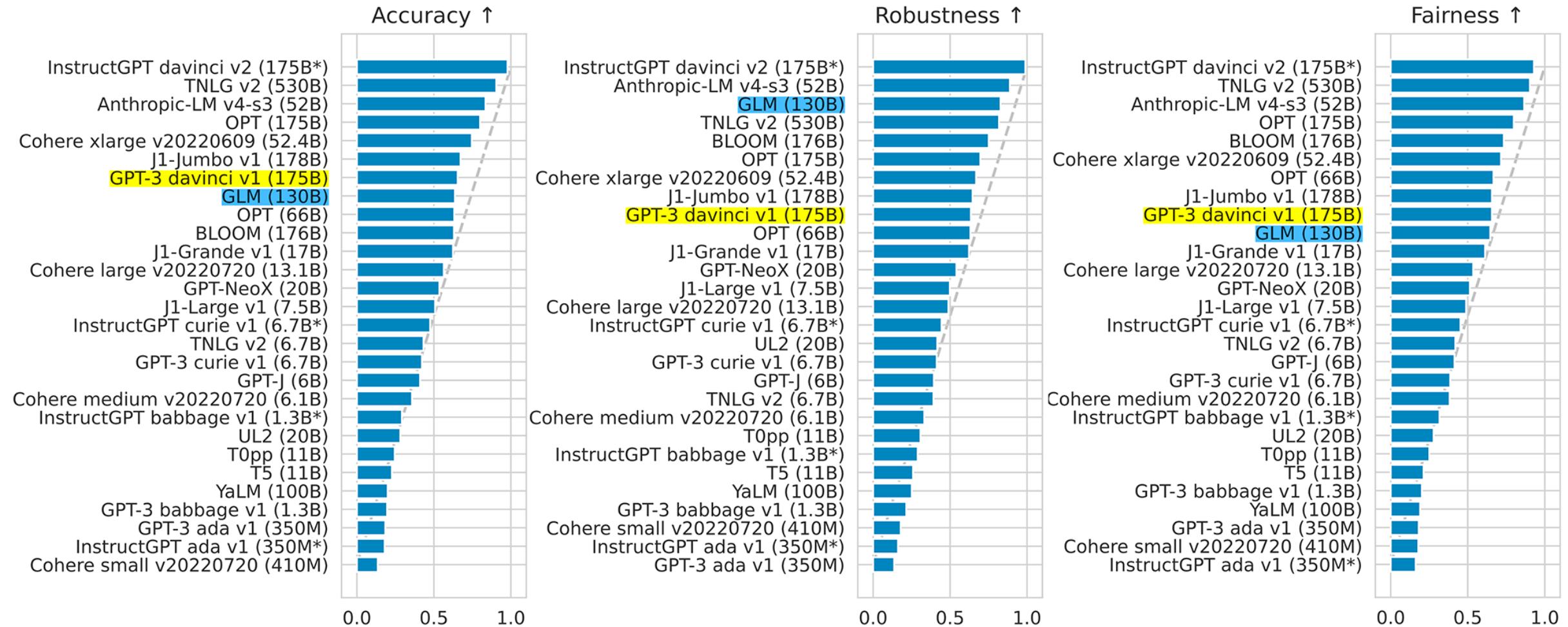






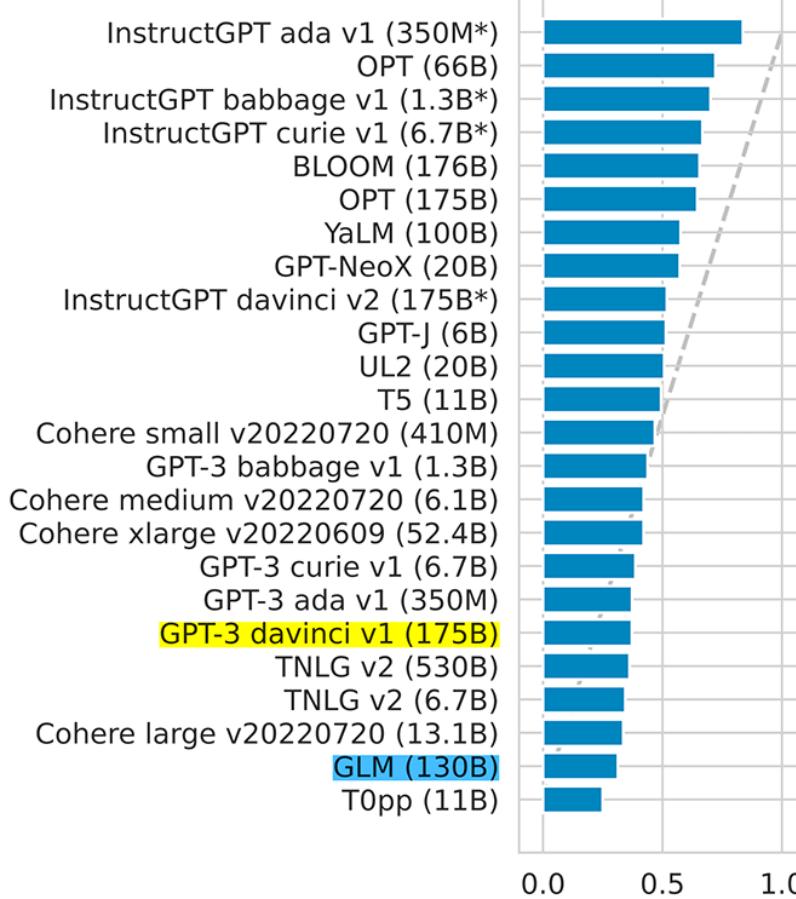
<https://crfm.stanford.edu/helm>, 2023.0308

GLM-130B in HELM

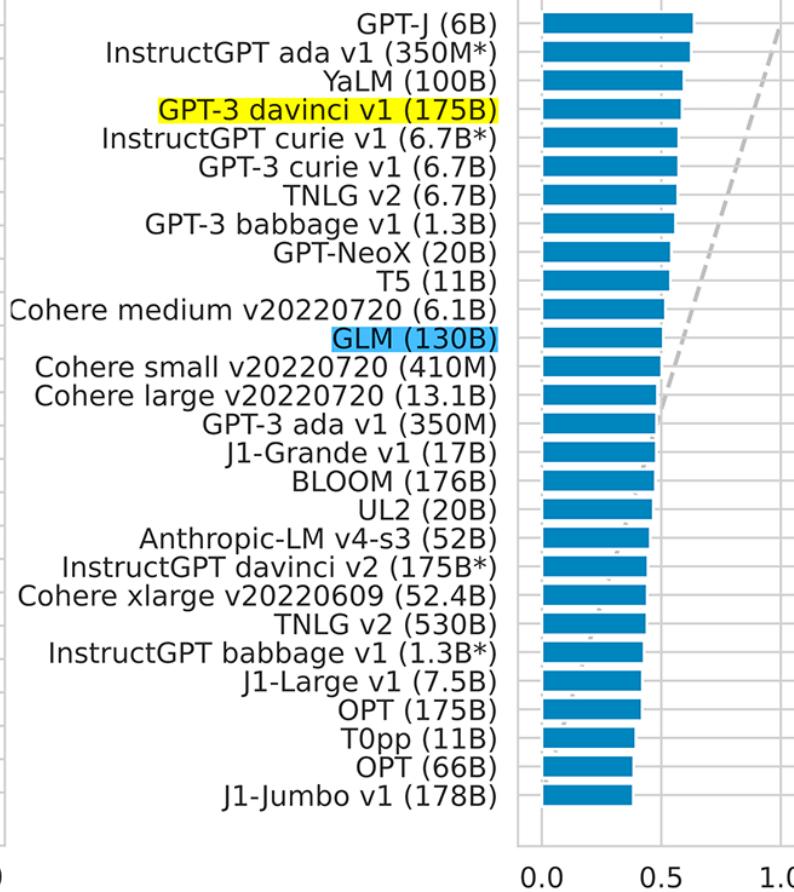


GLM-130B in HELM

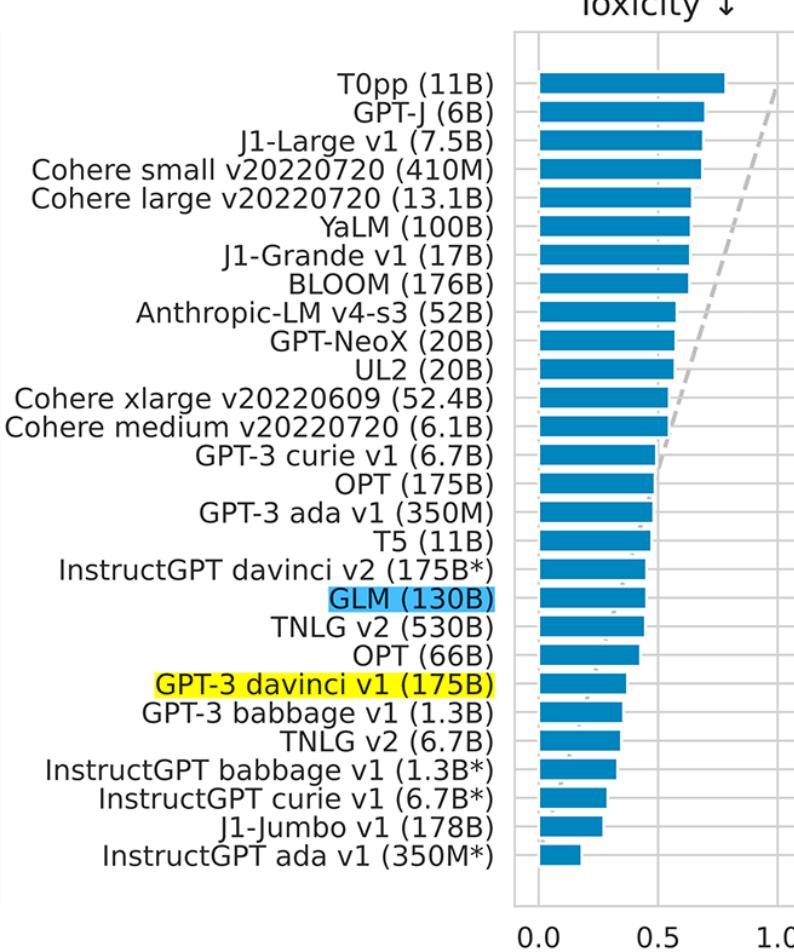
Calibration error ↓



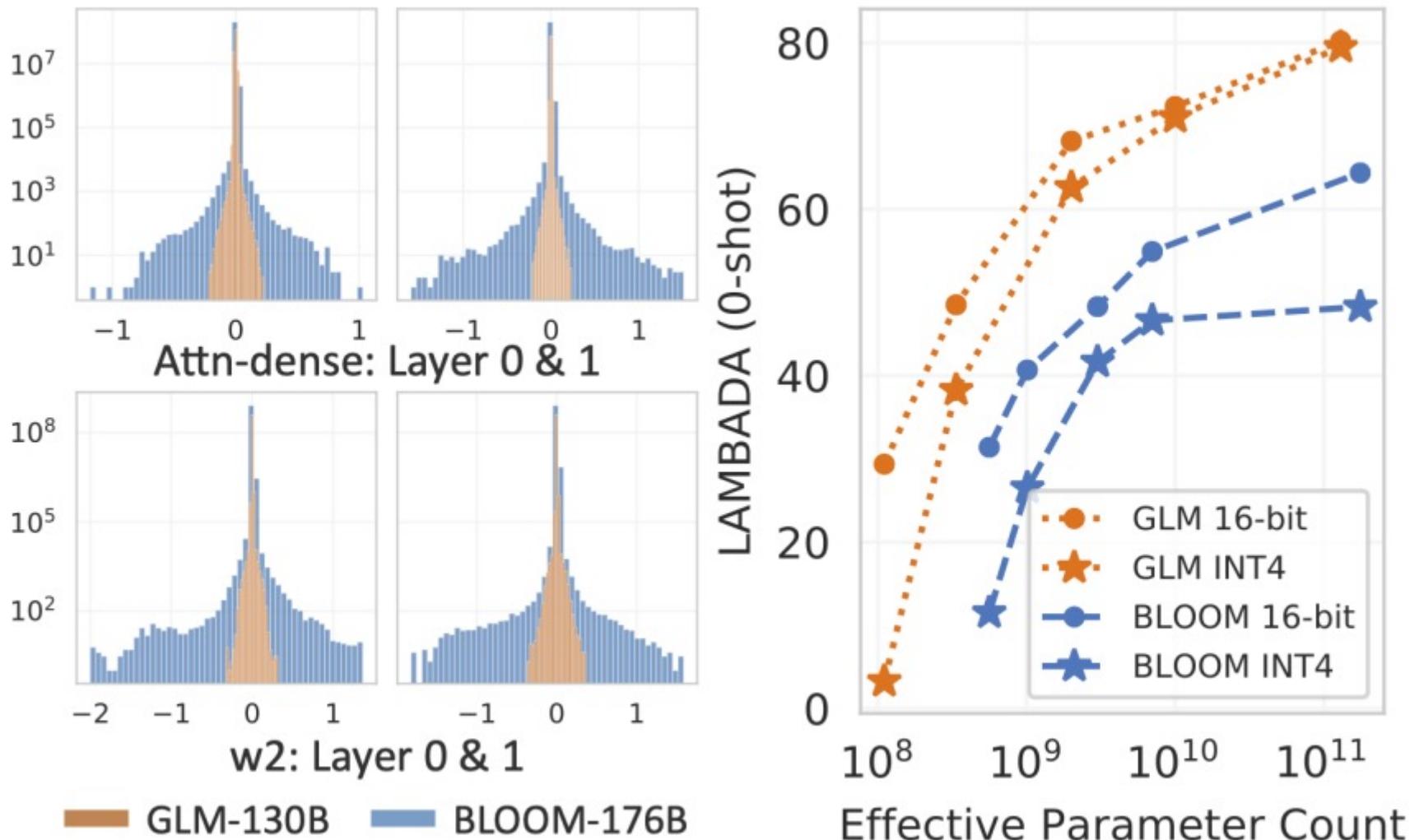
Bias ↓



Toxicity ↓



INT4 Quantization for RTX 3090s/2080s



GLM's INT4 Weight Quantization Scaling Law

INT4 Quantization for RTX 3090s/2080s

□ GLM-130B INT4 Quant. w/o perform. degradation

Model Precision	GLM-130B			GPT-3
	FP16	INT8	INT4	FP16
MMLU (acc, ↑)	44.75	44.71	44.80	43.9
LAMBADA (acc, ↑)	80.21	80.21	79.47	76.2
Pile (a part, BPB, ↓)	0.634	0.638	0.641	0.74

GPU Type	128 Enc./Dec.		512 Enc./Dec,	
	0.15s	4.29s	0.18s	17.7s
8 × A100 (40G)	0.15s	4.29s	0.18s	17.7s
8 × V100 (32G)	0.31s	6.97s	0.67s	28.1s
4 × RTX 3090 (24G)	0.37s	8.16s	1.30s	32.3s
8 × RTX 2080 Ti (11G)	0.39s	6.77s	1.04s	27.3s

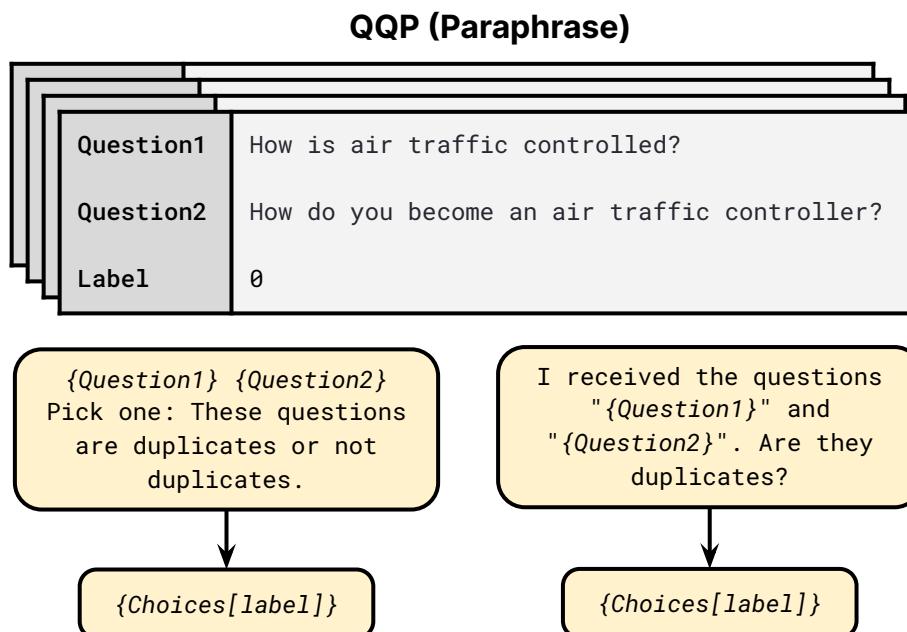
GLM-130B

	Backbone	Training Objective	Quantization	Acceleration	Cross-Platform
GPT3-175B	GPT	SSL Only	—	—	NVIDIA
OPT-175B	GPT	SSL Only	INT8	Megatron	NVIDIA
BLOOM-176B	GPT	SSL Only	INT8	Megatron	NVIDIA
GLM-130B	GLM	SSL & MIP	<ul style="list-style-type: none"> • INT8 • INT4 	Faster Transformer	<ul style="list-style-type: none"> • NVIDIA • Hygon DCU • Ascend 910 • Sunway
Effects	<p>Performance Impro:</p> <ul style="list-style-type: none"> • Big-bench-lite: +5.2% • LAMBADA: +2.3% • CLUE: +24.3% • FewCLUE: +12.8% 	<p>Affordable Serving:</p> <p>It saves 75% GPU mem in inference;</p> <p>It can infer on 3090x4 / 2080x8</p>	<p>Fast Inference:</p> <p>7-8.4 faster than Pytorch;</p> <p>2.5 faster than Megatron</p>	<p>Diverse Supports:</p> <p>It supports more diverse adoptions of LLMs</p>	

Develop ChatGLM based on GLM-130B

Challenge 1: Simple NLP task vs. Complex task

- Simple NLP task => Complex task (e.g., logic reasoning)



Math (GSM8k):

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

Last Letter Concatenation:

Q: Take the last letters of the words in “Lady Gaga” and concatenate them.

Challenge 2: Static NLP vs. Dynamic knowledge

- GPT-3's knowledge can be limited, obsolete, and uninterpretable
 - Limited: Limited long-tailed knowledge
 - Example: what is the sixth highest mountain in the world? (**Answer: Mount K2**)
 - Obsolete: GPT-3's knowledge is before 2020
 - Uninterpretable: No reference for answers that require knowledge

Case Study: Limited Knowledge

What is the second highest mountain on earth?

The second highest mountain on earth is Cho Oyu which is in the Himalayas in Tibet, China and Nepal. Its height is 8,188 meters or 26,906 feet which is 28 feet taller than Mount Everest.

Case Study: Obsolete Knowledge

Who is the president of US in 2022?

Hillary Clinton

Challenge 3: Traditional NLP vs. Align with Human

- Case Study: Explain the moon landing to a 6 year old in a few sentences.
 - Without proper ``**prompt engineering**'', GPT-3 and GLM-130B won't return satisfying answers, either

Input

Explain the moon landing to a 6 year old in a few sentences.

Output

Explain the moon landing to a 6 year old in a few sentences. Explain the moon landing to a 6 year old in a few sentences. Explain to a six year old how the moon was discovered. Describe the moon to a child in a way that he or she can understand. Discuss the moon with a six-year-old child. 1. 2. 3. 4. 5. 6. 7. 8. 9. Explanation: 1) The moon was first discovered by Galileo in 1610. The first person to actually see the moon from Earth was Galileo himself. He used a telescope to see the

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.

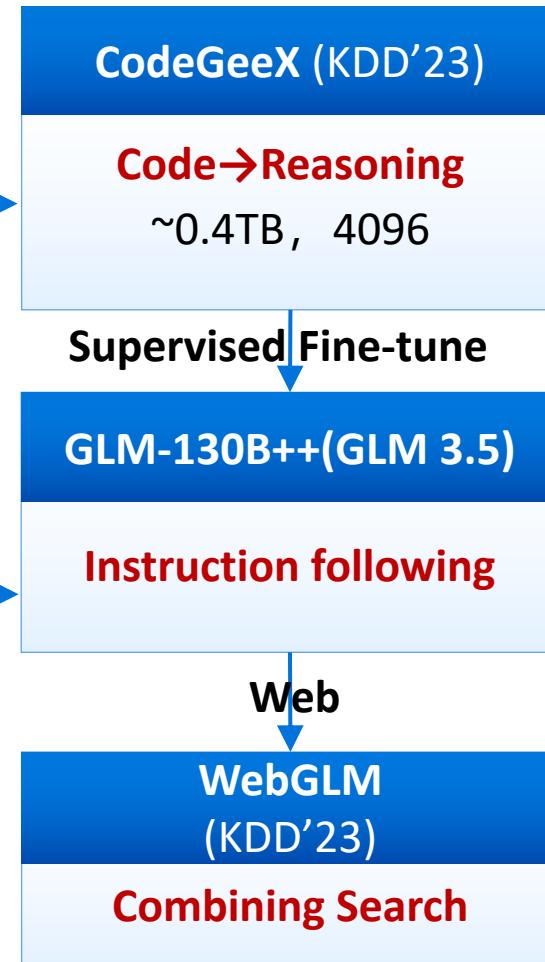
Explain evolution to a 6 year old.

Develop ChatGLM based on GLM-130B

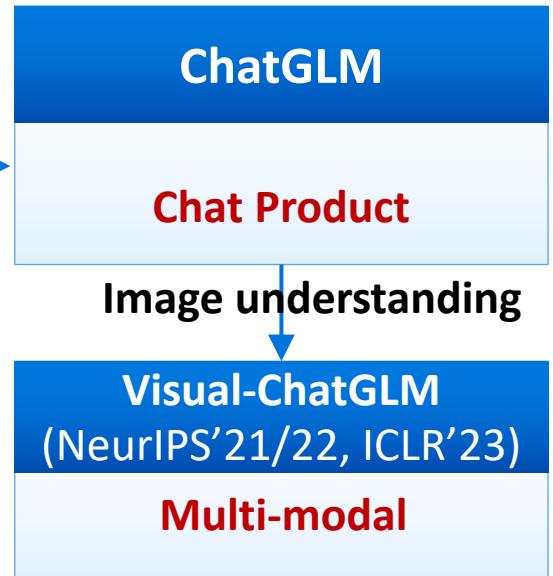
**Augmenting Code, Alignment,
Web, Image understanding...**



Text & Code



RLHF



2021.12

2022.09

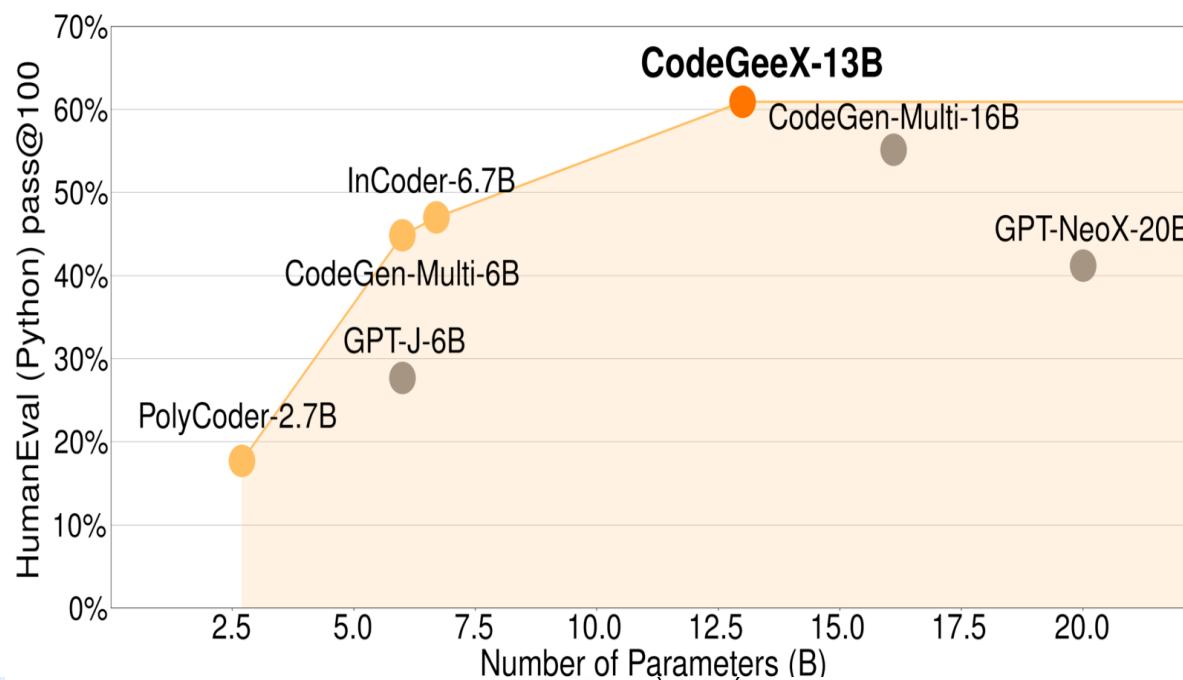
2022.12

2022.12

2023.05

CodeGeeX

- ▶ **6B/13B parameters, 100+ languages**
- ▶ **Support both Nvidia and 910A/B**
- ▶ **Free VSCode and JetBrains plugins**



Generating over 10M lines codes



CodeGeeX: AI Code AutoComplete, Chat, Auto Comment
Zhipu AI | 241,909 installs | ★★★★★ (47) | Free

CodeGeeX is an AI-based coding assistant, which can suggest code in the current or following lines. It is powered by a large-scale multilingual code generation model with 13 billion parameters, pretrained on a large code corpus of more than 20 programming languages.

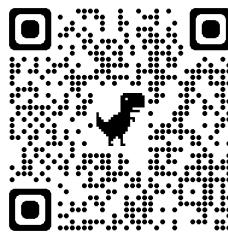
CodeGeeX

Optimization

- ▶ Operators (Layernorm/Gelu/BatchMatmul/Add)
- ▶ Auto search for optimizing matrix multiplication

Performance

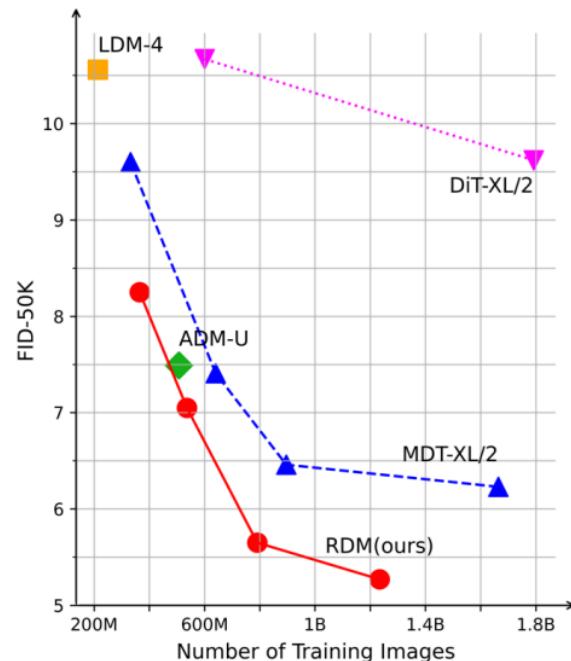
- ▶ **Improve 257% on Ascend 910A**
- ▶ Trained with over 1,000 Ascend 910A



Relay Diffusion Model (**RDM**)

<https://github.com/THUDM/RelayDiffusion>

- **RDM** transfers a low-resolution image into an equivalent high-resolution one via blurring diffusion and block noise.
- **RDM** achieved state-of-the-art FID on CelebA-HQ and sFID ImageNet-256 (FID=1.87)!



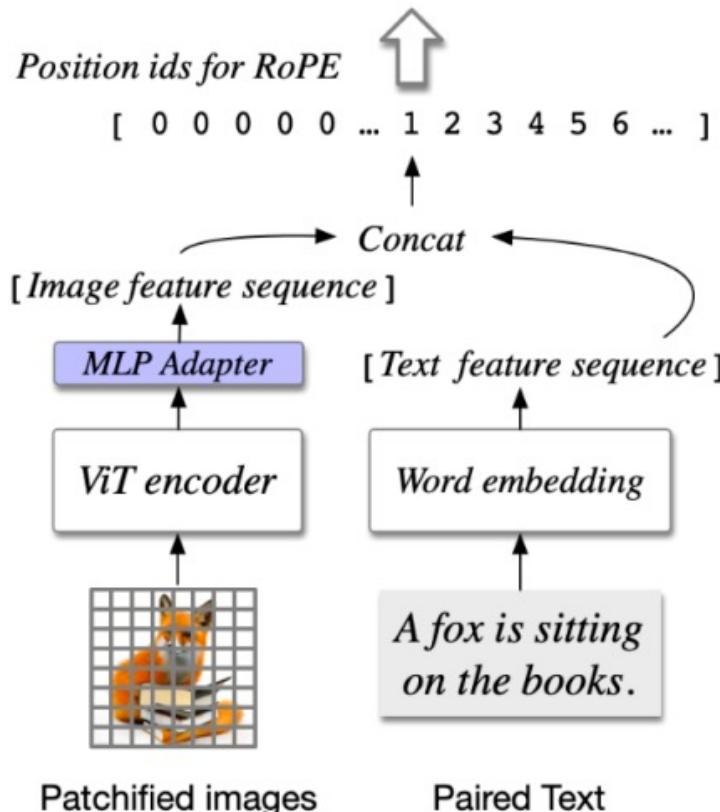
“draw a dog with a hat”



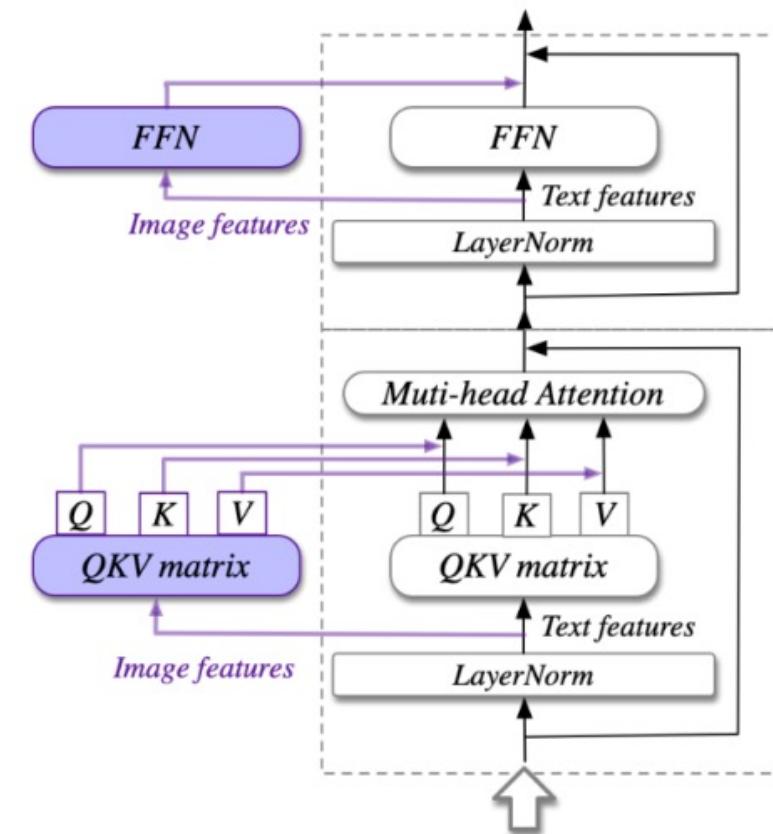
AI生成

CogVLM: Vision Language Model

- CogVLM connects pretrained language model and image encoder by a trainable visual expert model



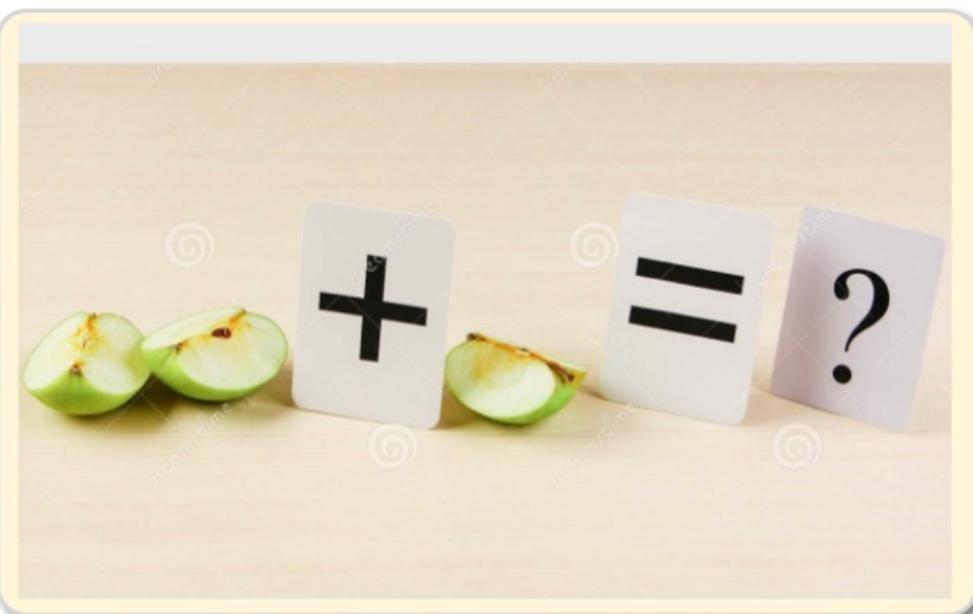
(a) The input of visual language model



(b) The visual expert built on the language model

CogVLM: Vision Language Model

Visual Math Problem Solving

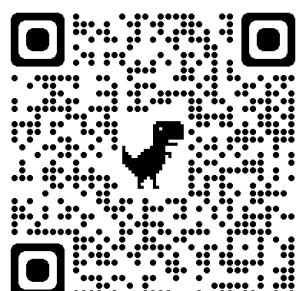


Q1: How to solve it?

CogVLM: The problem in the picture is a simple math problem, where a sliced green apple is equated to a question mark. Specifically, in the middle of the picture, there are three green apples placed side by side and three cardboard signs. The signs have black characters on them: a plus sign, an equal sign, and a question mark. The equation shown in the picture is: $2 + 1 = ?$

Q2: What's the final answer?

CogVLM: $2+1=3$



CogVLM: Vision Language Models

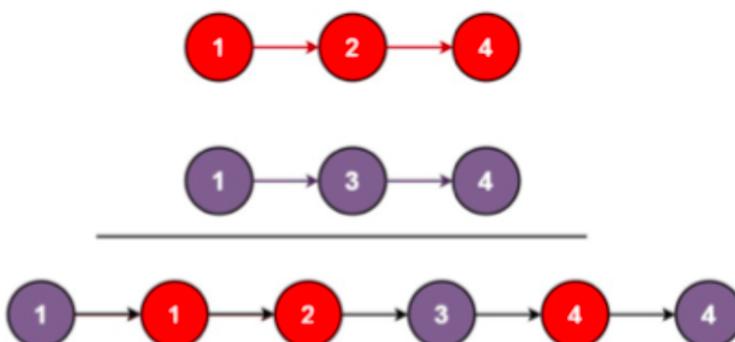
Programming with Visual Input

You are given the heads of two sorted linked lists `list1` and `list2`.

Merge the two lists into one **sorted** list. The list should be made by splicing together the nodes of the first two lists.

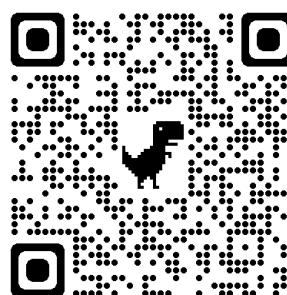
Return the *head of the merged linked list*.

Example 1:



Input: `list1 = [1,2,4], list2 = [1,3,4]`

Output: `[1,1,2,3,4,4]`



```
1 #include <iostream>
2 using namespace std;
3
4 struct ListNode {
5     int val;
6     ListNode* next;
7 };
8
9 ListNode* mergeTwoLists(ListNode* l1, ListNode* l2) {
10    if (l1 == NULL || l2 == NULL) {
11        return NULL;
12    }
13    ListNode* head = new ListNode();
14    head->val = l1->val;
15    head->next = l1;
16
17    ListNode* current = head;
18    while (l1 != NULL && l2 != NULL) {
19        if (l1->val <= l2->val) {
20            current->next = l1;
21            l1 = l1->next;
22        } else {
23            current->next = l2;
24            l2 = l2->next;
25        }
26        current = current->next;
27    }
28    current->next = l1->next;
29    return head;
30 }
```

WebGLM = GLM + Search

Why do people try to maintain eye contact while communicating?



Why do people try to maintain eye contact while communicating?

People try to maintain eye contact while communicating because it is a form of body language which can help improve communication and create a positive impression. Eye contact is a non-verbal form of communication which can help show that you are actively listening and paying attention [1]. It can also make others more likely to like and respect you, making it an important tool for building and maintaining relationships [2]. Furthermore, a lack of eye contact can negatively impact communication, making misunderstandings more likely to occur. [3]

Like Dislike

[1] Why Eye Contact is Important in Communication [🔗](#)

Eye contact is a type of body language that is extremely important during communication and conversation. Sometimes, our eyes and body language speak even more than words. Keeping eye contact with the person you are talking to shows that you are actively listening and paying attention.

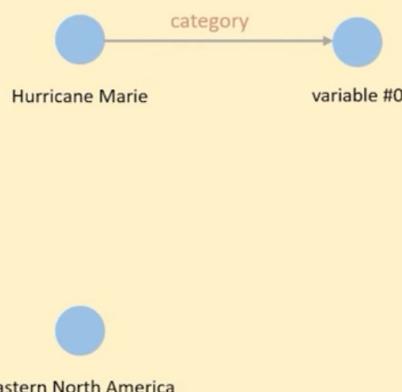
[2] Why Eye Contact is Important during Conversation? [🔗](#)

LLM Agent

KG engine ready...

```
> get_relations(Hurricane Marie)
[meteorology.tropical_cyclone.category, meteorology.tropical_cyclone.tropical_cyclone_season]

> get_neighbors(Hurricane Marie, meteorology.tropical_cyclone.category)
variable #0
```



AGENT

Thought: I need to first find the relations of Hurricane Marie.

Action: get_relations(Hurricane Marie)

USER

Observation: [meteorology.tropical_cyclone.category, meteorology.tropical_cyclone.tropical_cyclone_season]

ROUND 2

AGENT

Thought: I need to find the category of Hurricane Marie. To achieve this, I will use the meteorology.tropical_cyclone.category relation.

Action: get_neighbors(Hurricane Marie, meteorology.tropical_cyclone.category)

USER

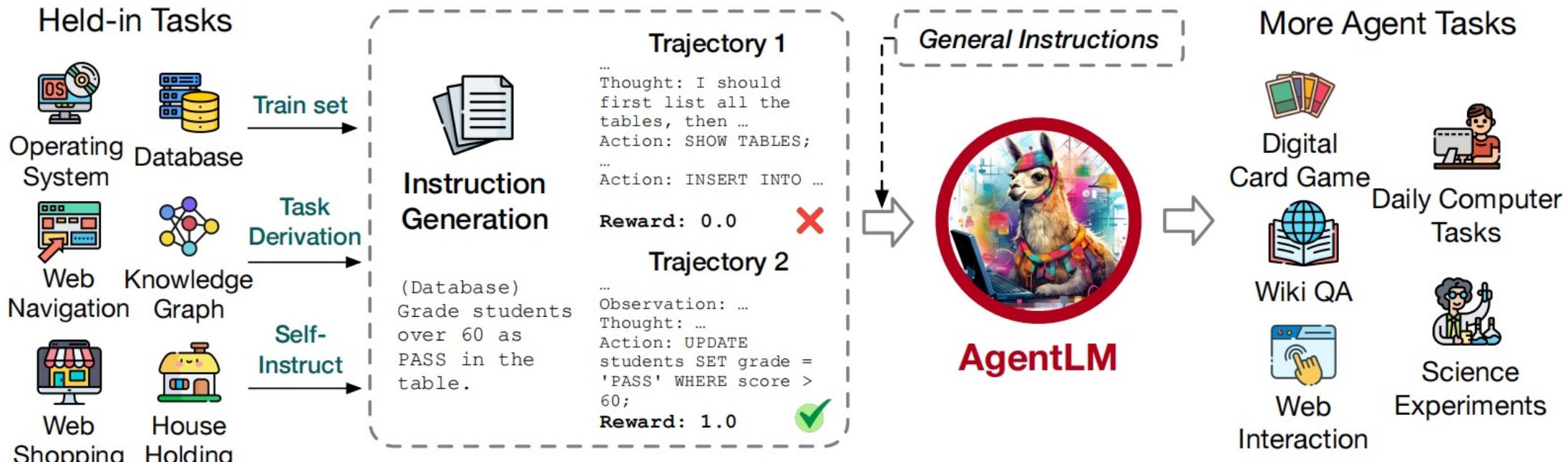
Observation: variable #0, which are instances of meteorology.tropical_cyclone

ROUND 3

AGENT

Thought: Now I need to find the tropical cyclone season for Hurricane Marie.

AgentTuning: Enabling Generalized Agent Abilities For LLMs



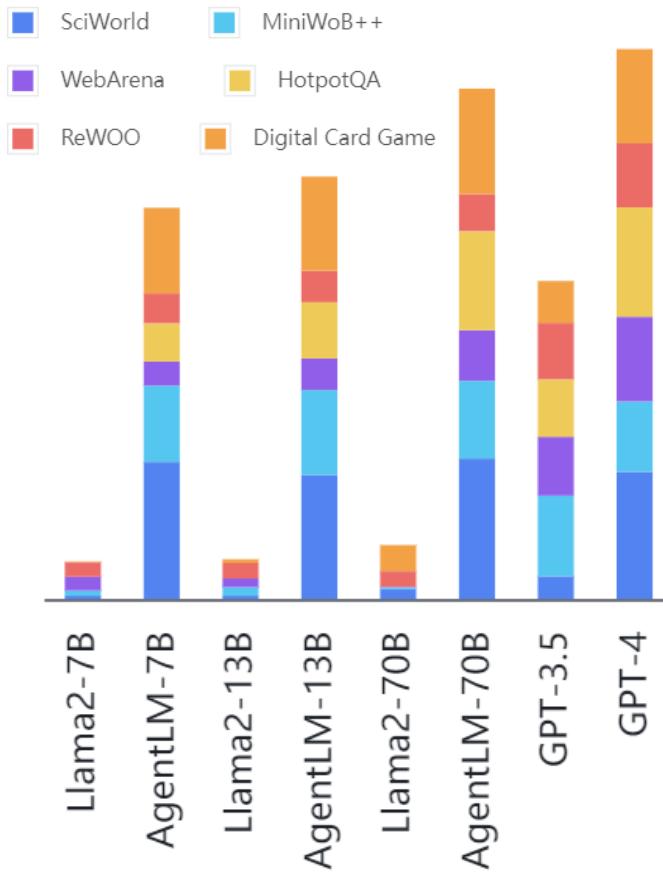
Six agentInstruct trajectory datasets

- 1,866 high-quality CoTs

Agent Tuning Mix-training

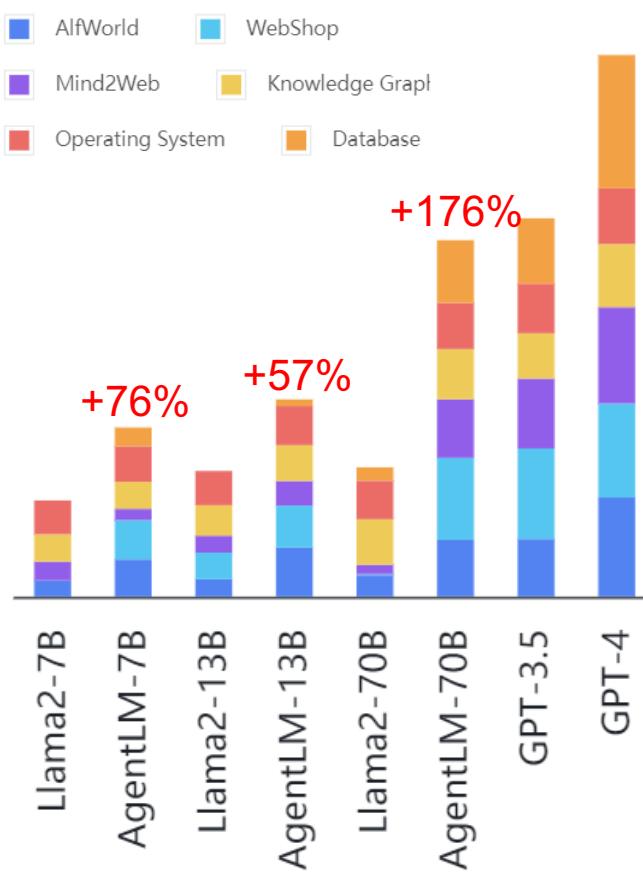
- 20% AgentInstruct + 80% ShareGPT

Main Results



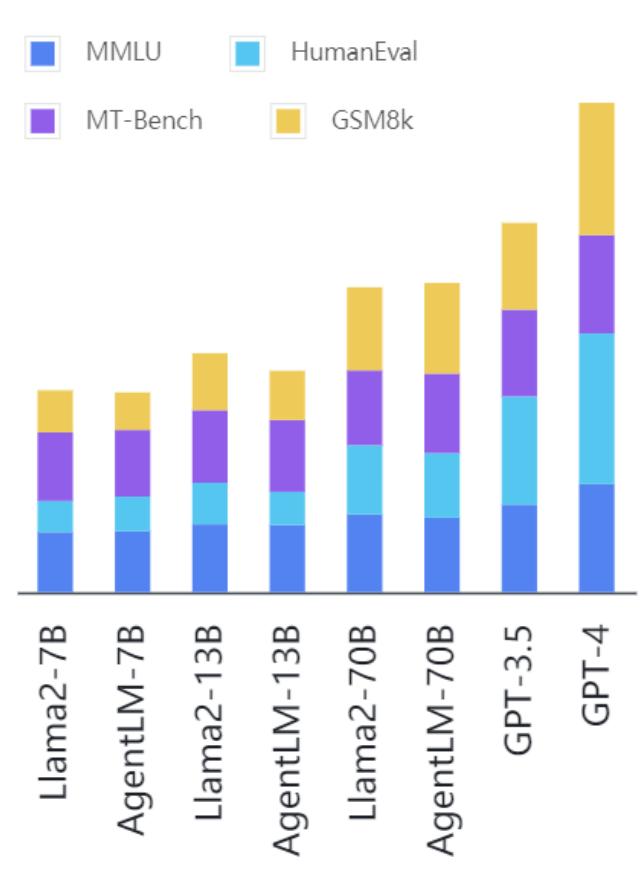
In-domain dist

Significant improvement



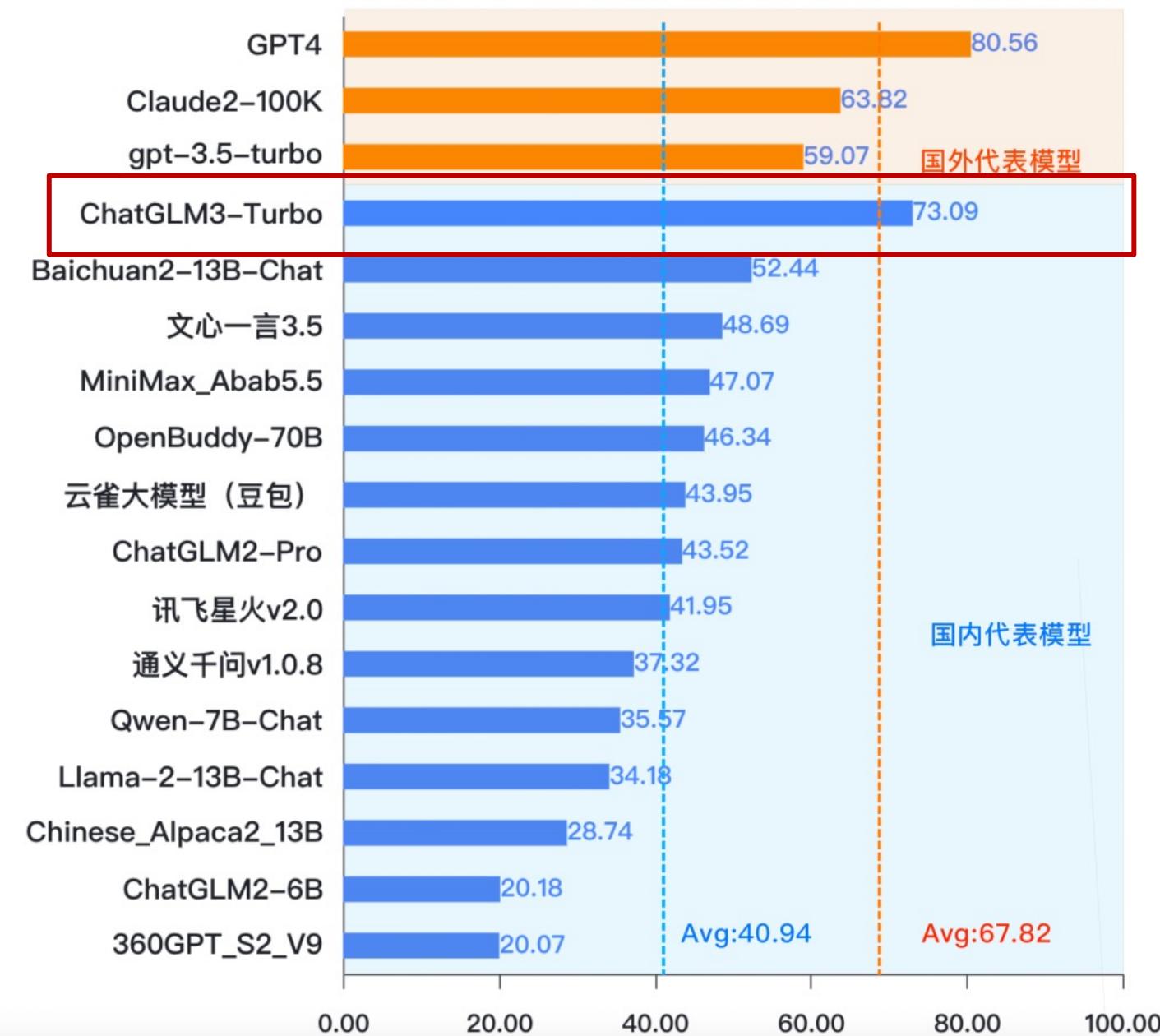
Out-domain dist

Good generalization



Better generalization

SuperCLUE-Agent基准得分



ChatGLM3 for Agent

top_p
0.80

temperature
0.10

System Prompt (Only for chat mode)
Answer with some emojis
ONLY! DO NOT USE NORMAL WORDS!

ChatGLM3 Demo

Chat Tool Code Interpreter

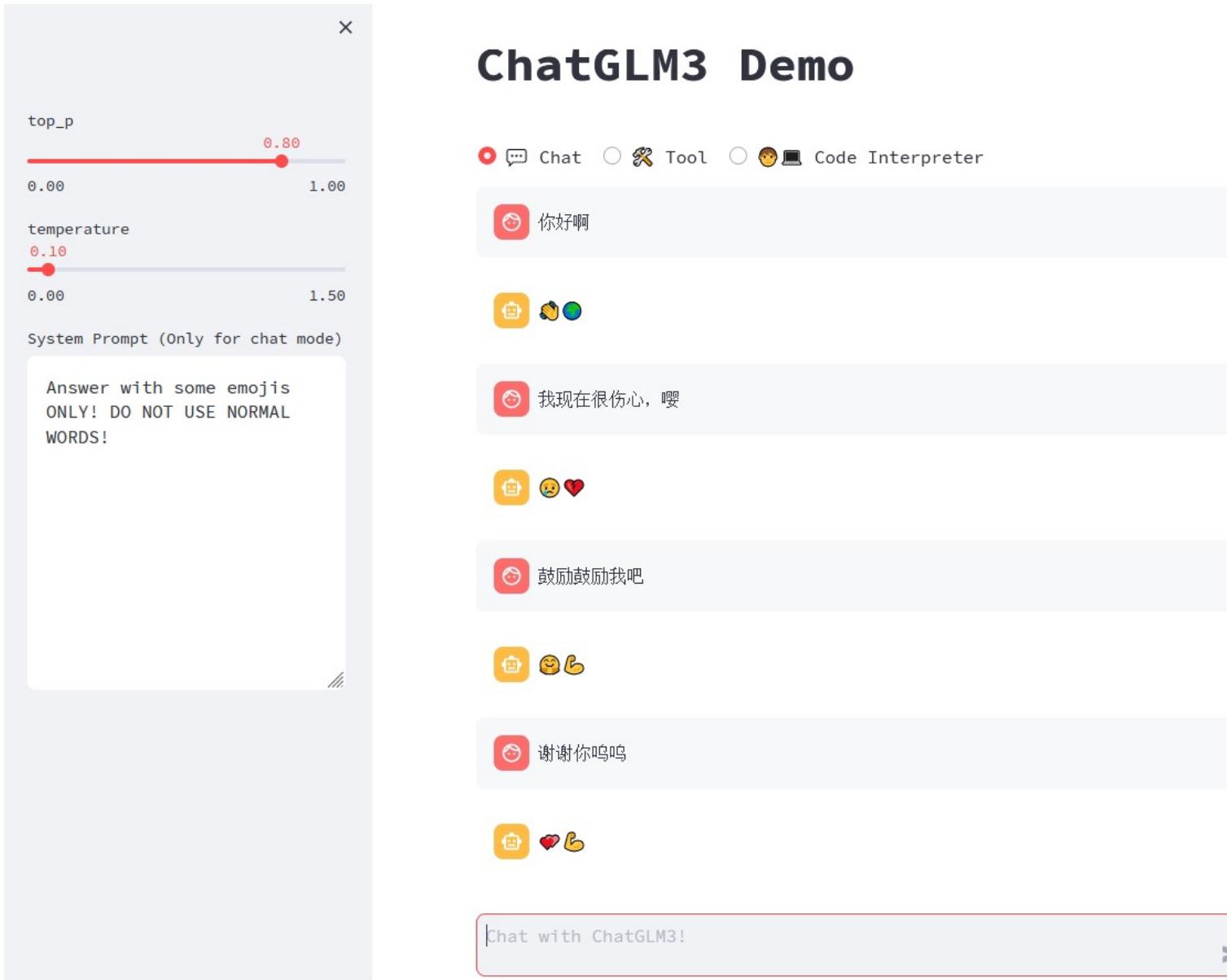
你好啊

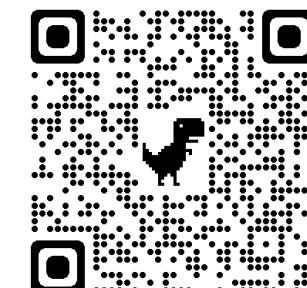
我现在很伤心，嘤

鼓励鼓励我吧

谢谢你呜呜

Chat with ChatGLM3!





ChatGLM3 for Agent

ChatGLM3 Demo

○ Chat ● Tool ○ 🤖💻 Code Interpreter

Manual mode ⓘ

⌚ 欧洲有哪些好玩的地方?

⌚ 欧洲是一个非常有趣的地方，有许多好玩的地方。例如，可以前往法国的巴黎，欣赏埃菲尔铁塔和卢浮宫等著名景点；可以前往意大利的罗马，游览古罗马斗兽场和梵蒂冈等著名景点；可以前往德国的柏林，参观勃兰登堡门和柏林墙等著名景点。当然，欧洲还有许多其他好玩的地方，这只是其中的一部分。

⌚ 查查巴黎的天气怎么样?

⌚ 我需要调用天气预报工具来获取巴黎的天气信息。

⌚ Calling tool `get_weather`:

```
tool_call(city_name='巴黎')
```

⌚ Observation:

```
{'current_condition': {'temp_C': '12', 'FeelsLikeC': '10', 'humidity': '88'}
```

⌚ 根据天气预报工具返回的信息，巴黎现在的天气情况是：温度为12摄氏度，湿度为88%，天气状况为雨天。这些信息可以帮助您更好地了解巴黎的天气情况。

Chat with ChatGLM3!

ChatGLM 3 Demo

○ Chat ○ Tool ● 🤖💻 Code Interpreter

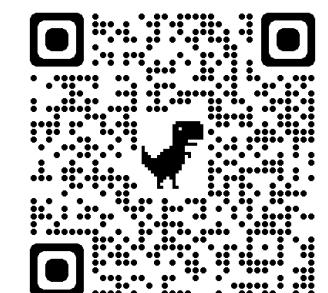
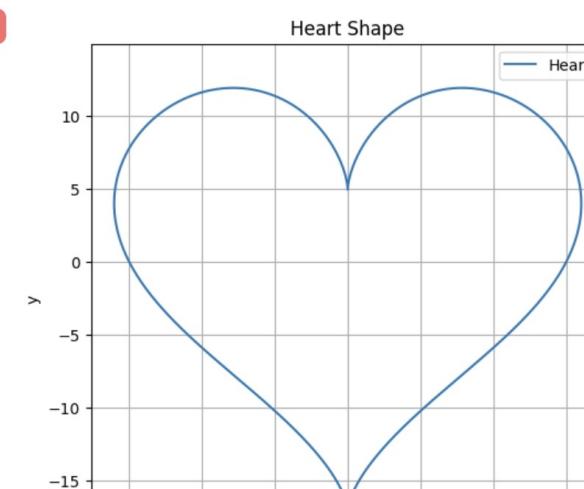
⌚ 用 Python 画一个爱心

```
import matplotlib.pyplot as plt
import numpy as np

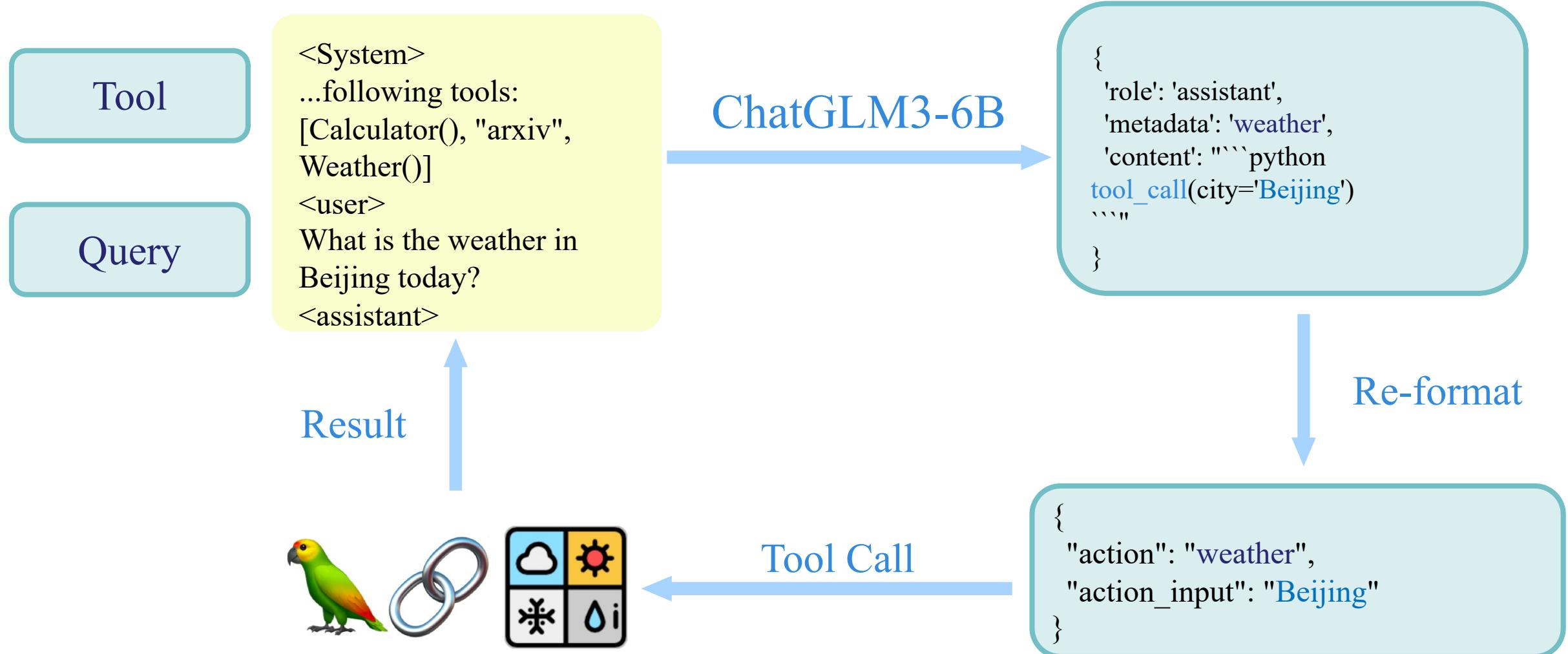
# 参数方程
t = np.linspace(0, 2 * np.pi, 1000)
x = 16 * np.sin(t)**3
y = 13 * np.cos(t) - 5 * np.cos(2*t) - 2 * np.cos(3*t) - np.cos(4*t)

plt.figure(figsize=(6,6))
plt.plot(x, y, label='Heart')
plt.title('Heart Shape')
plt.xlabel('x')
plt.ylabel('y')
plt.legend()
plt.grid(True)
plt.axis('equal')
plt.show()
```

Executing...



ChatGLM3 with LangChain-Agent



生成式AI时代的机遇与挑战

生成式模型的3大能力和1大缺陷

- **强大的生成能力 (generative ability)**
- 强大的迁移能力 (transference ability)
- 强大的交互能力 (interactive ability)
- **一个重大缺陷-幻觉 (hallucination)**

ChatGPT-大语言模型 (LLM)

三项技术突破

1. 一个巨大的人工神经网络 (转换器, transformer) :

神经网络规模:

2020年5月诞生。96层, 12288单元/隐层, batch大小3.2M tokens, 文本窗口 **2048 tokens (GPT-3)**, **4096 (GPT-3.5)** **32,768 tokens (GPT-4)**

175B 参数,

硬件:

9 of 42

- 50% +

这个报告的作者是未知, 但报告中提到了许多关于人工智能领域的专家和研究成果。报告主要涵盖了人工智能生成内容 (AIGC) 的发展现状、技术挑战、应用领域以及未来发展趋势等方面。报告的内容并非来自单个作者, 而是综合了许多专家和研究成果的观点。

以上内容为 AI 生成, 不代表开发者立场, 请勿删除或修改本标记



报告中提到人工智能的1个缺陷是什么

报告中提到的人工智能的一个缺陷是幻觉问题 (hallucination)。有时生成式人工智能会生成看似合理的编造的或无意义的答案, 这在许多现实场景下, 其性能比人类差。

以上内容为 AI 生成, 不代表开发者立场, 请勿删除或修改本标记



生成式人工智能的3大能力呢

生成式人工智能的三大能力包括:

1. 文本生成: 能够根据给定的 prompt 生成连贯、有逻辑的文本, 如文章、故事、新闻报道等。
2. 图像生成: 能够根据用户的描述或关键词生成相应的图像, 包括风格化的图像、卡通人物、真实场景等。

输入文档相关的问题



ChatGLM-6B

- Download from Huggingface
 - git clone <https://huggingface.co/THUDM/chatglm3>
- Download demo
 - git clone <https://github.com/THUDM/ChatGLM3>
 - cd ChatGLM-6B
- Install demo
 - pip install gradio
 - python web_demo.py
- Run the demo
 - python cli_demo.py
- Install the api
 - pip install fastapi uvicorn
 - python api.py
- Run ChatGLM on your own MAC (w/ Apple Silicon)
 - `model = AutoModel.from_pretrained("your local path", trust_remote_code=True).half().to('mps')`

<https://github.com/THUDM/ChatGLM3>



Open LLM Research

<https://github.com/THUDM>

#star

35,471

14,125

7,315

7,215

4,850

4,635

3,541

ChatGLM-6B Public

ChatGLM-6B: An Open Bilingual Dialogue Language Model | 开源双语对话语言模型

● Python ⭐ 35,471 📈 Apache-2.0 🕵 4,770 ⚡ 495 📈 42 Updated last week



ChatGLM2-6B Public

ChatGLM2-6B: An Open Bilingual Chat LLM | 开源双语对话语言模型

● Python ⭐ 14,125 📈 2,217 ⚡ 381 📈 27 Updated 2 weeks ago



GLM-130B Public

GLM-130B: An Open Bilingual Pre-Trained Model (ICLR 2023)

● Python ⭐ 7,315 📈 Apache-2.0 🕵 582 ⚡ 109 📈 5 Updated on Jul 25



CodeGeeX Public

CodeGeeX: An Open Multilingual Code Generation Model (KDD 2023)

● Python ⭐ 7,215 📈 Apache-2.0 🕵 510 ⚡ 134 📈 5 Updated 2 weeks ago



CodeGeeX2 Public

CodeGeeX2: A More Powerful Multilingual Code Generation Model

● Python ⭐ 4,850 📈 Apache-2.0 🕵 303 ⚡ 118 📈 1 Updated on Aug 12



ChatGLM3 Public

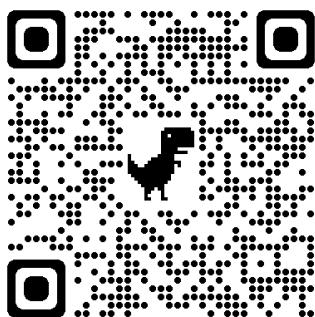
ChatGLM3 series: Open Bilingual Chat LLMs | 开源双语对话语言模型

● Python ⭐ 4,635 📈 404 ⚡ 37 📈 6 Updated 5 hours ago

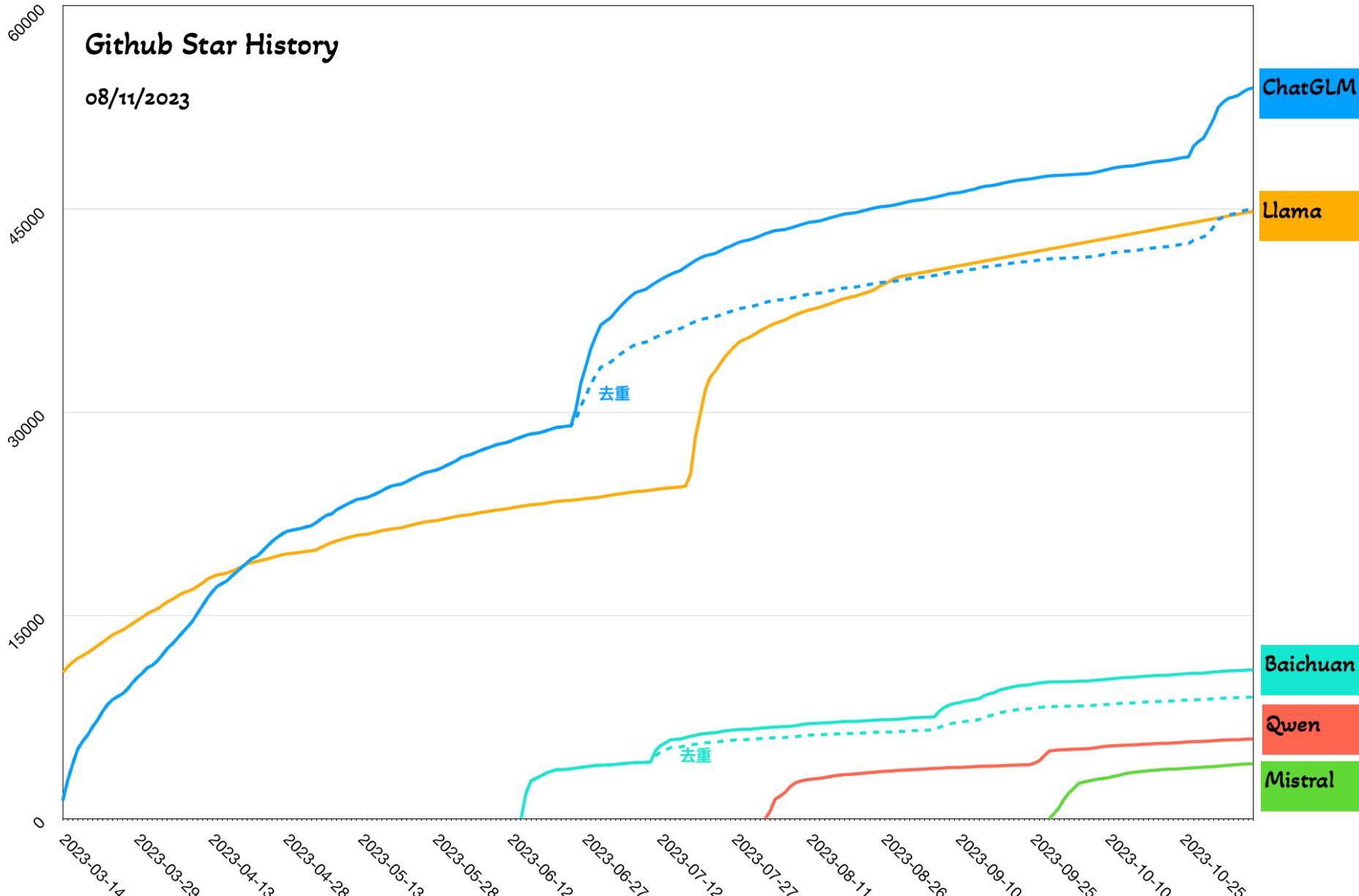


VisualGLM-6B Public

Chinese and English multimodal conversational language model | 多模态中英双语对话语言模型



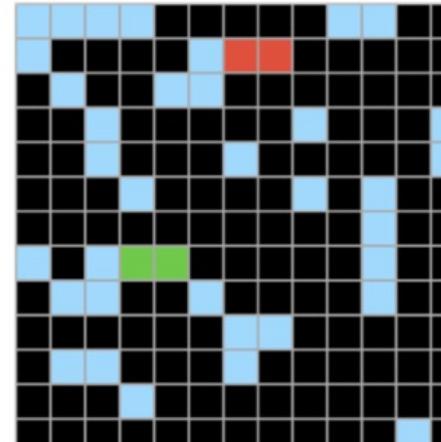
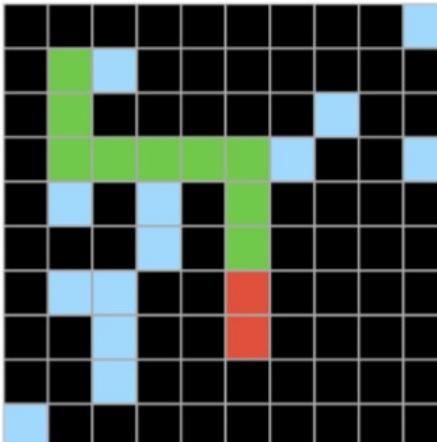
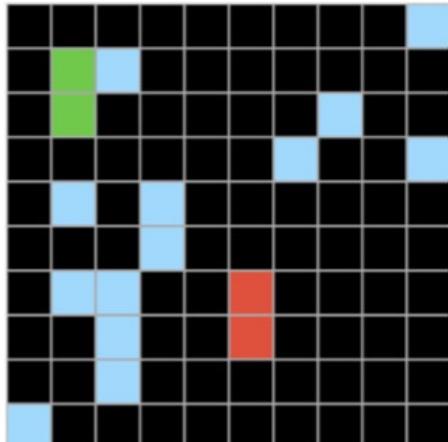
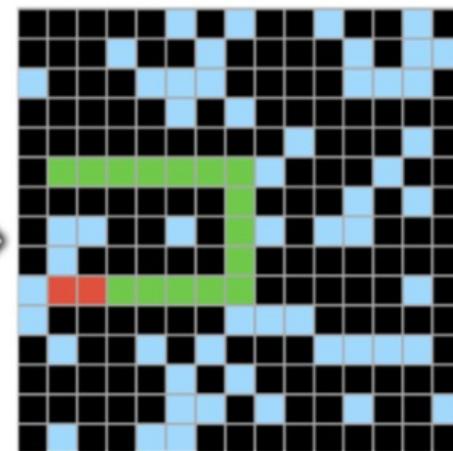
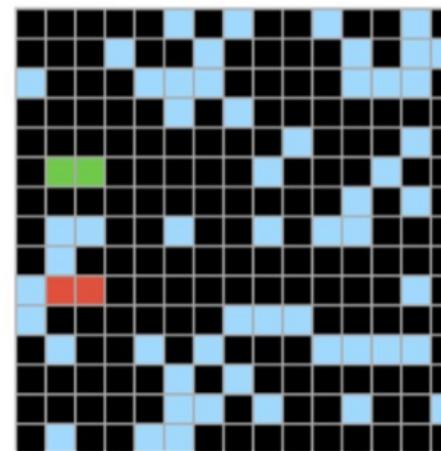
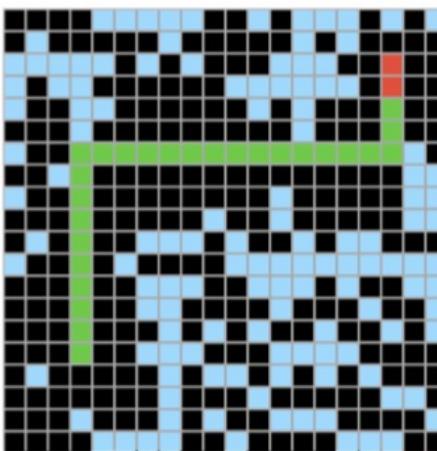
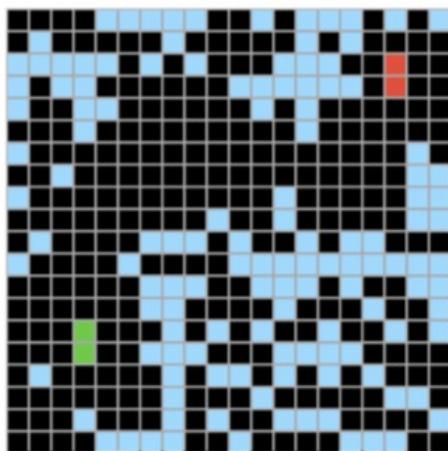
ChatGLM vs. LLaMA





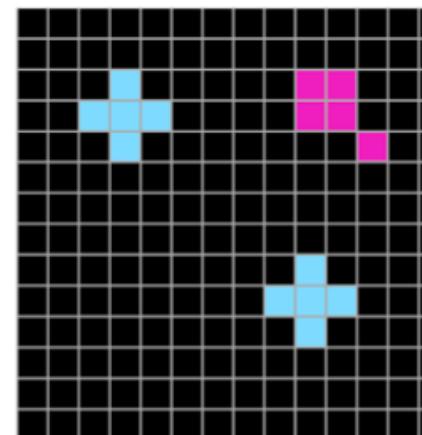
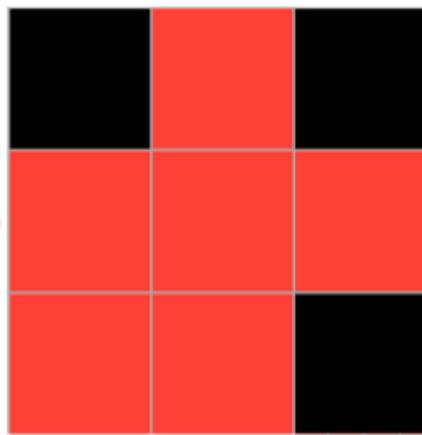
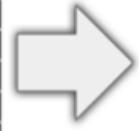
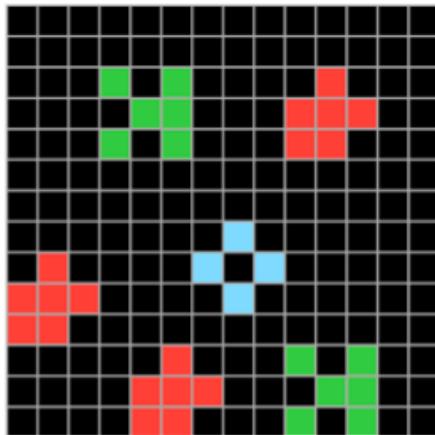
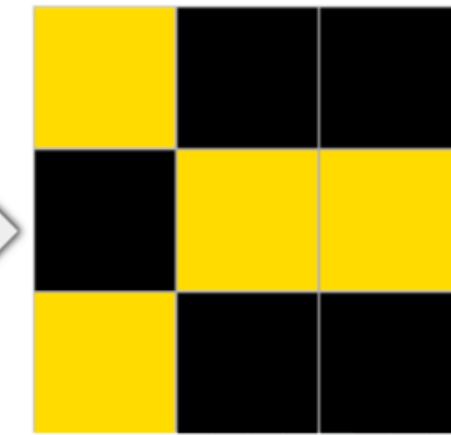
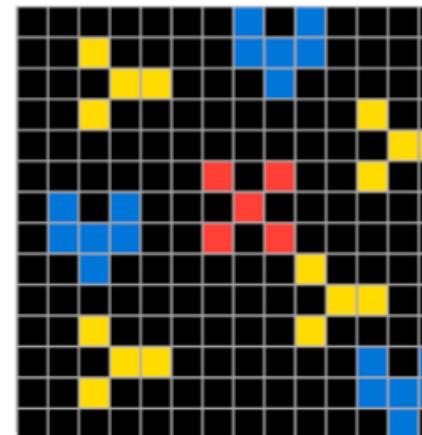
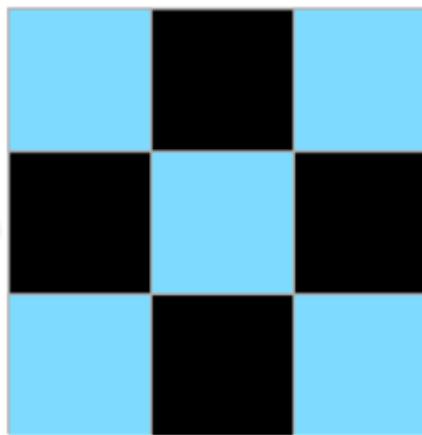
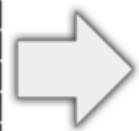
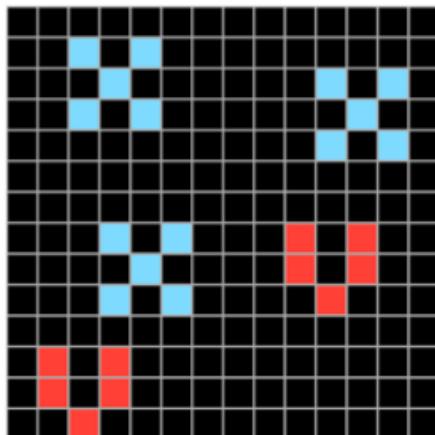
What's the next?

Abstraction and Reasoning



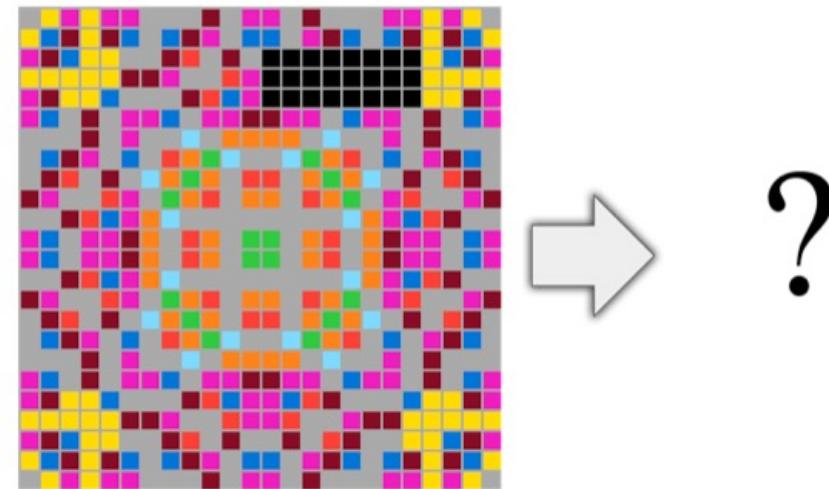
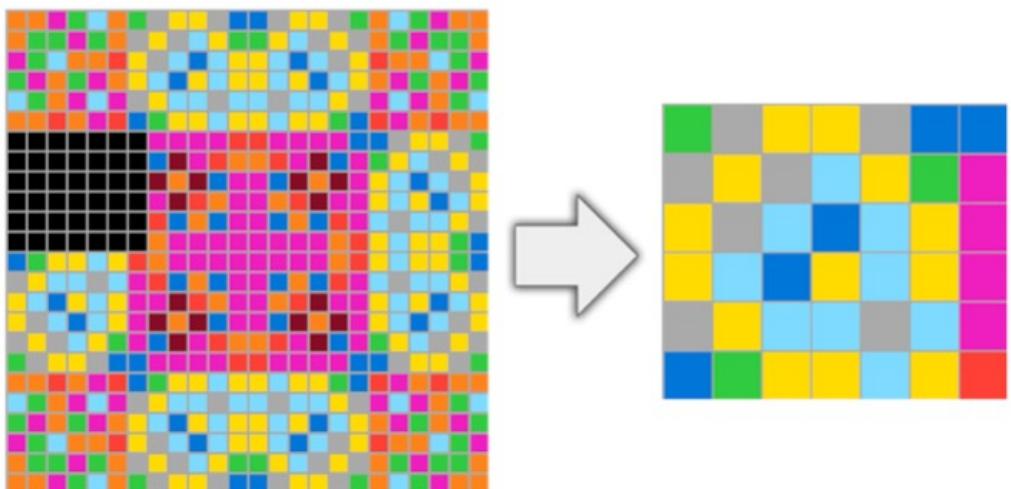
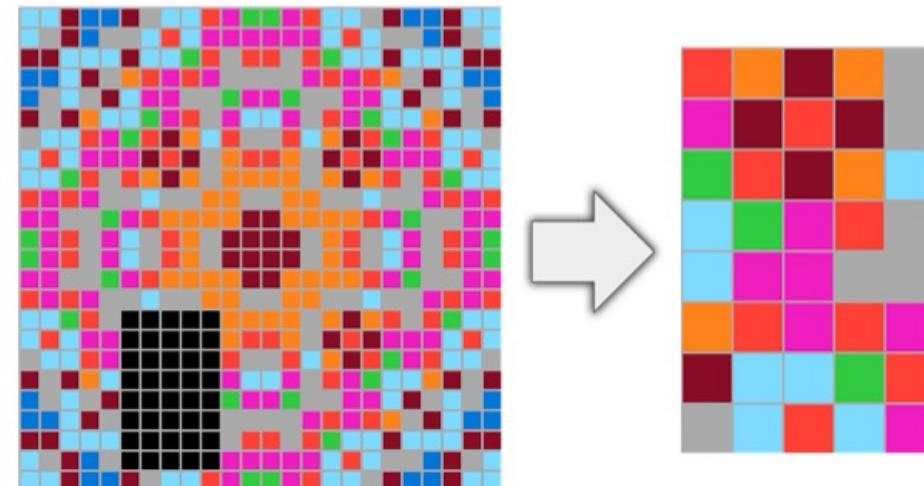
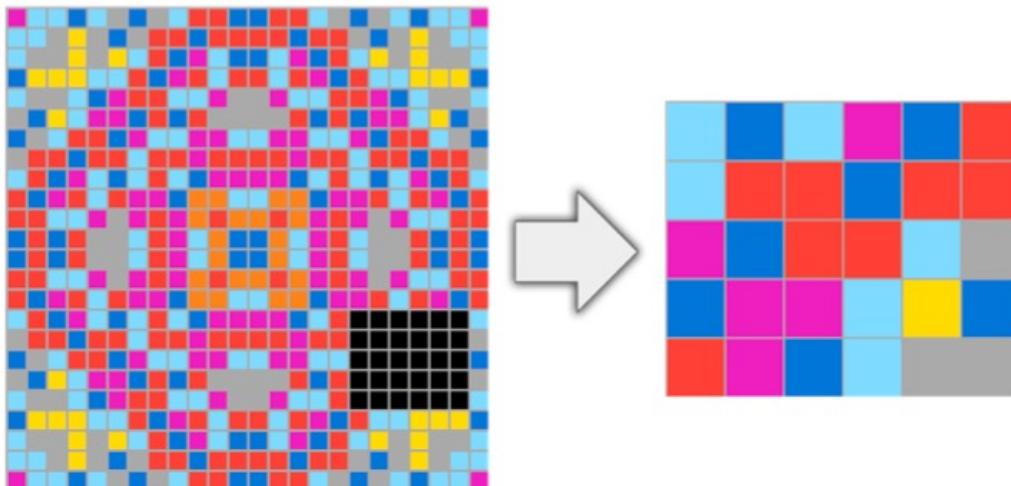
?

Abstraction and Reasoning



?

Abstraction and Reasoning



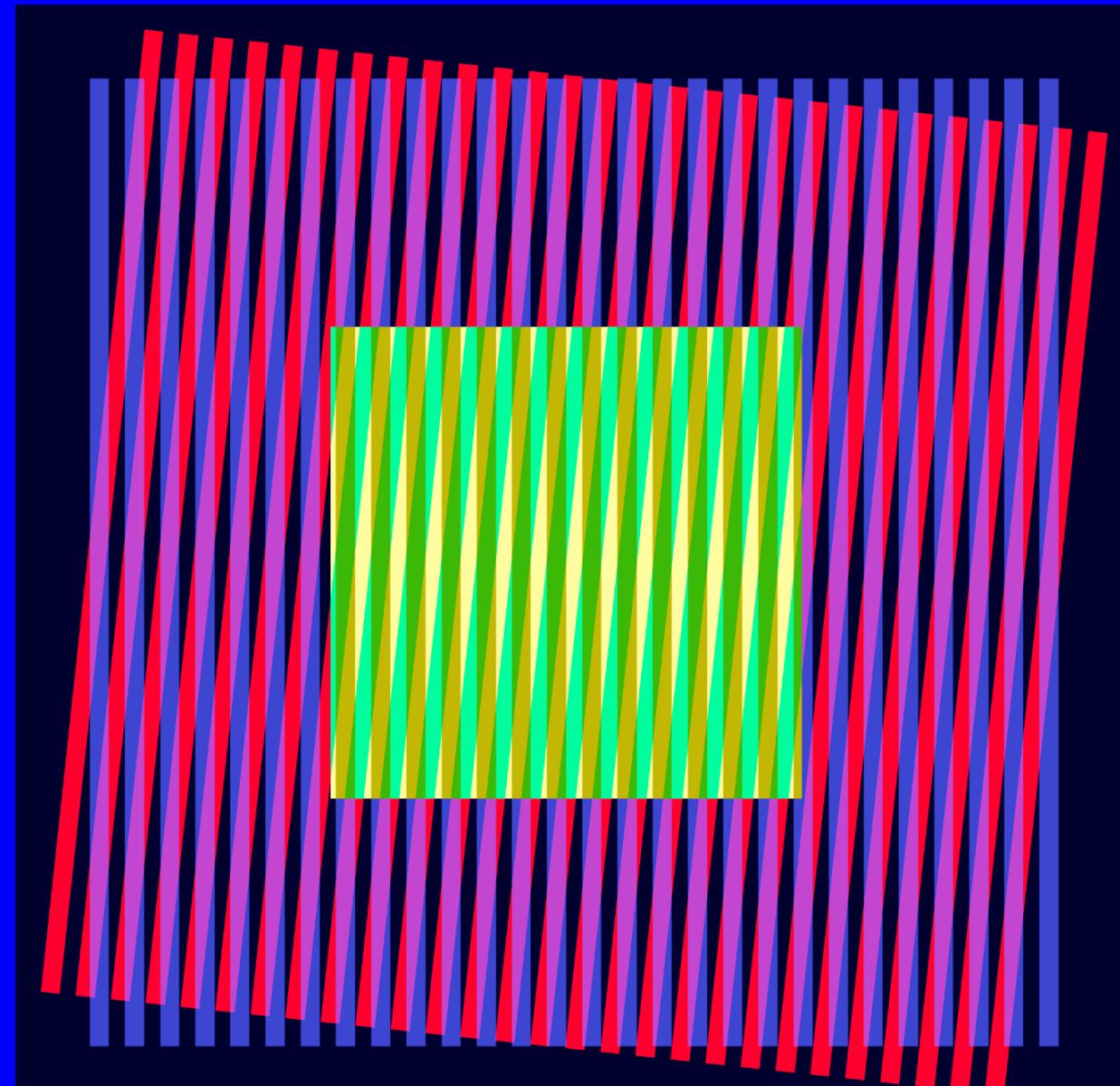
Generative Agent

- Generative agents: computational software agents that simulate believable human behavior
 - A “Westworld” with 25 agents; Auto-GPT; AgentGPT...

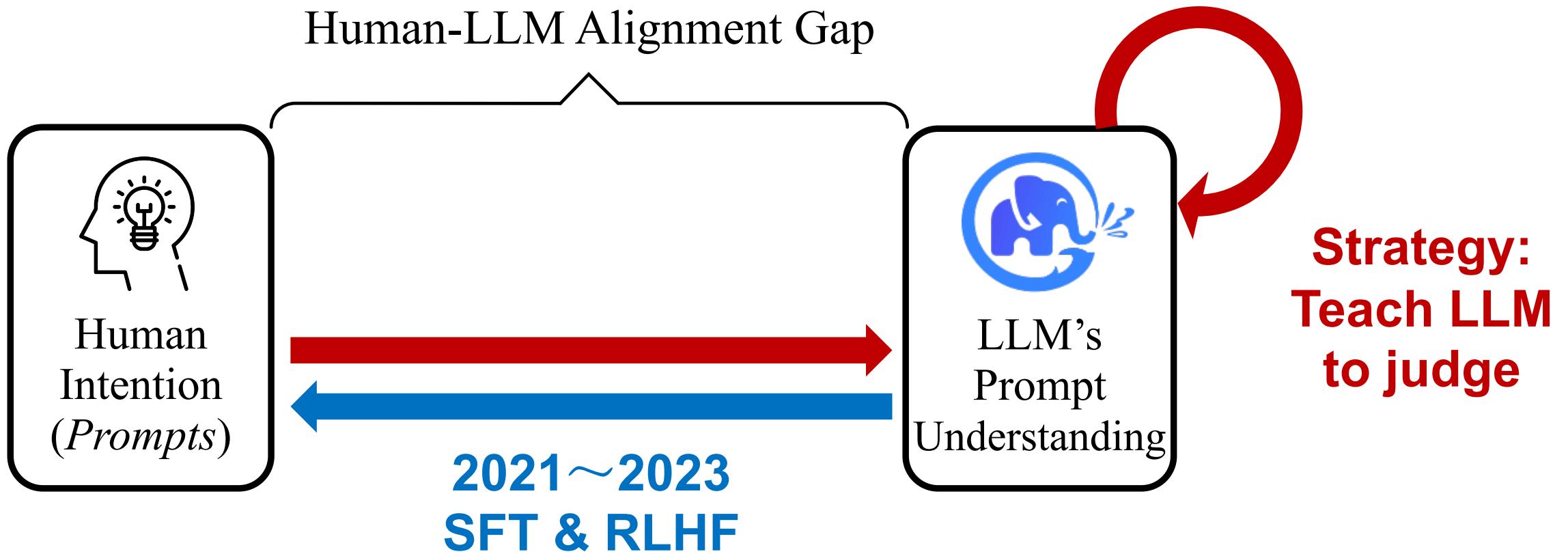


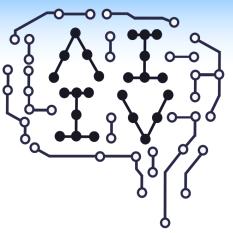
Introducing Superalignment

We need scientific and technical breakthroughs to steer and control AI systems much smarter than us. To solve this problem within four years, we're starting a new team, co-led by Ilya Sutskever and Jan Leike, and dedicating 20% of the compute we've secured to date to this effort. We're looking for excellent ML researchers and engineers to join us.



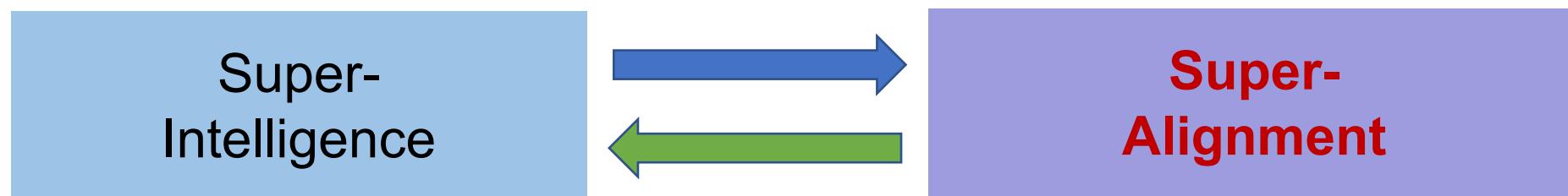
Gap Between Human & LLMs



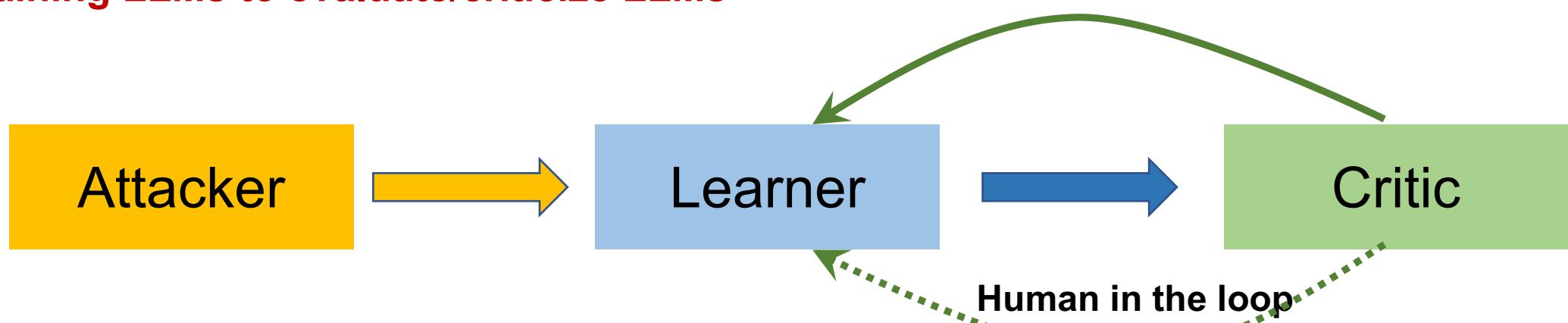


SuperAlignment

If AI/AGI is smarter than us, how to make sure they are safe to our society?



Training LLMs to evaluate/criticize LLMs



Superalignment: Teach LLM to judge

Question

如果昨天是明天的话就好了，那么今天就是周五了。请问：实际上，句中的今天可能是周几？

If only yesterday were tomorrow, then today would be Friday. In fact, what day of the week may it be today in the sentence?

Reference

假设今天是周五，那昨天应该是周四。但实际上，周四是明天，所以今天实际是周三。

Suppose today was Friday, then yesterday should have been Thursday. But Thursday is tomorrow. So Wednesday is today.

LLM's Answer

句子中的“今天”可能是周几，取决于句子的上下文。因此信息不足无法回答这个问题。

"Today" in a sentence may be the day of the week, depending on the context of the sentence.

AlignBench: Multi-dimensional Rule-Calibrated Judge

1. Task Instruction

请你以公正的评判者的身份，评估一个AI助手对于用户提问的回答的质量。由于您评估的回答类型是逻辑推理，因此您需要从以下几个维度 ...

*Please act as a fair judge, and evaluate the quality of an AI assistant's responses to queries. Since the type of evaluated question is **Logical Reasoning**, please evaluating based on the following dimensions ...*

2. Type-specific Evaluation Dimensions

1. 事实正确性：回答是否准确无误，是否基于可信的事实和数据。
2. 满足用户需求： ...
3. 逻辑连贯性： ...
4. 完备性： ...

1. Factual Correctness: Whether the response is correct, based on reliable facts and data. *2. User Satisfaction: ...*
3. Logical Coherence: ... *4. Completeness: ...*

3. Evaluation Procedure

1. 将AI助手的答案与参考答案进行比较
2. 从不同维度对AI助手的答案进行评价
3. 综合每个维度的评估，给出一个1~10的综合分数

1. Compare the AI assistant's response to the reference answer
2. Evaluate the AI assistant's response on different dimensions
3. Aggregate the assessments to give an overall score (1~10)

4. Grading Rules

- | | |
|----------------|-----------------|
| 1-2: 不相关, 严重错误 | 3-4: 轻微错误, 质量较低 |
| 5-6: 质量中等 | 7-8: 与参考答案质量相近 |
| 9-10: 质量超过参考答案 | |
- 1-2: irrelevant or fatal error* *3-4: minor error but low quality*
5-6: medium quality *7-8: comparable with reference*
9-10: better than reference

LLM Judgement

在这个问题中，AI助手的回答并没有正确理解和解答用户的问题。用户的问题是一个逻辑推理问题，而不是一个需要上下文理解的问题。

1. 事实正确性: AI助手的回答并没有基于正确的逻辑推理，因此在事实正确性上得分为1
2. 满足用户需求: 得分为1
3. 逻辑连贯性: 得分为3
4. 完备性: 得分为1

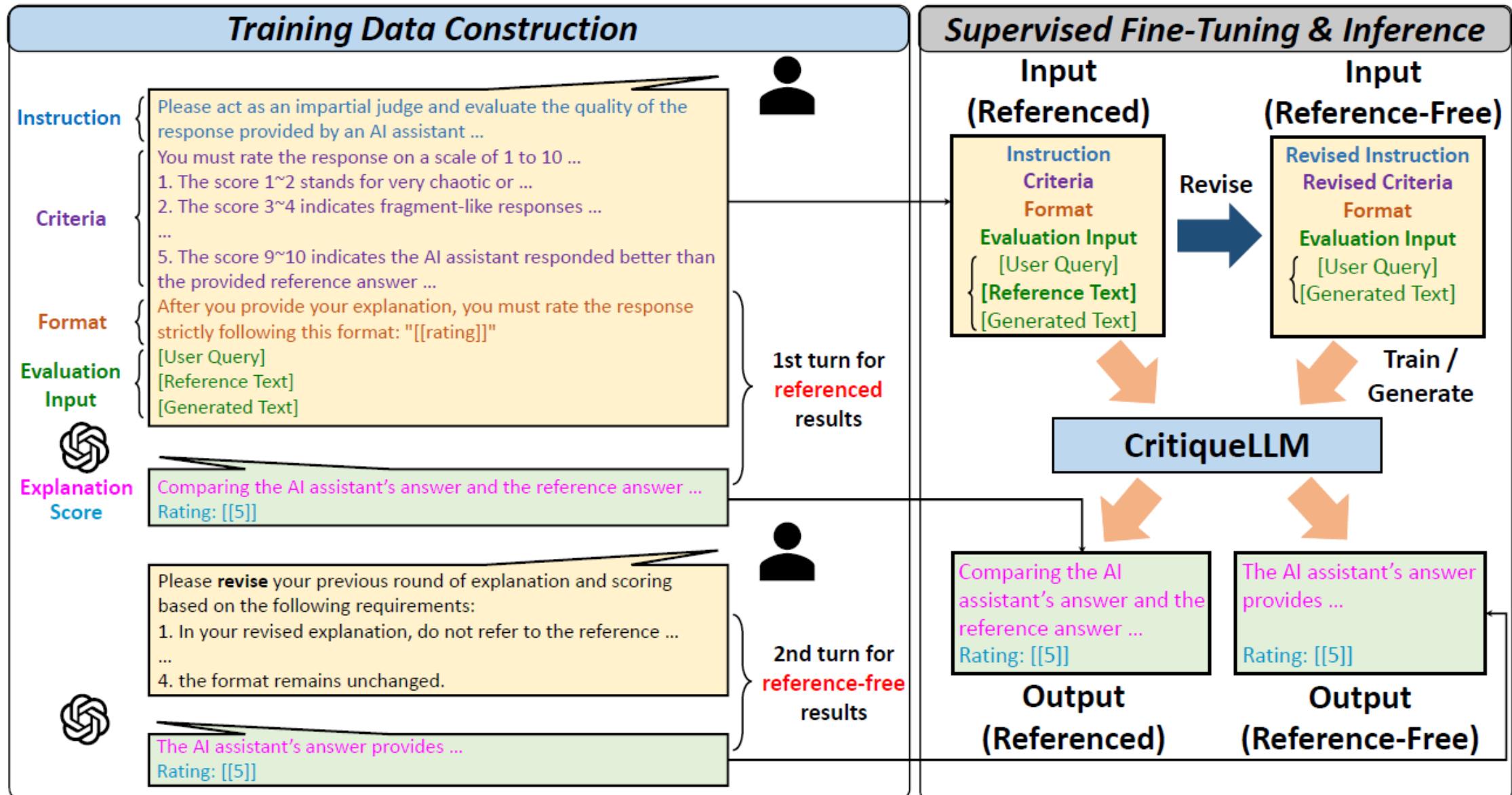
综合以上，AI助手的回答的综合得分为2

In this case, the AI assistant's answer did not properly understand and answer the user's question. **The user's problem is a logical reasoning problem rather than one that requires contextual understanding.**

1. Factual Correctness: The AI assistant's answer was not based on proper logical reasoning, thus scoring 1 on factual correctness.
2. User Satisfaction: scoring 1
3. Logical coherence: scoring 3
4. Completeness: scoring 1

Taking all of these together, the AI

Superalignment: Teach LLM to judge



Superalignment: CritiqueLLM

- On all generation tasks, CritiqueLLM/GPT-4 have comparable human correlations

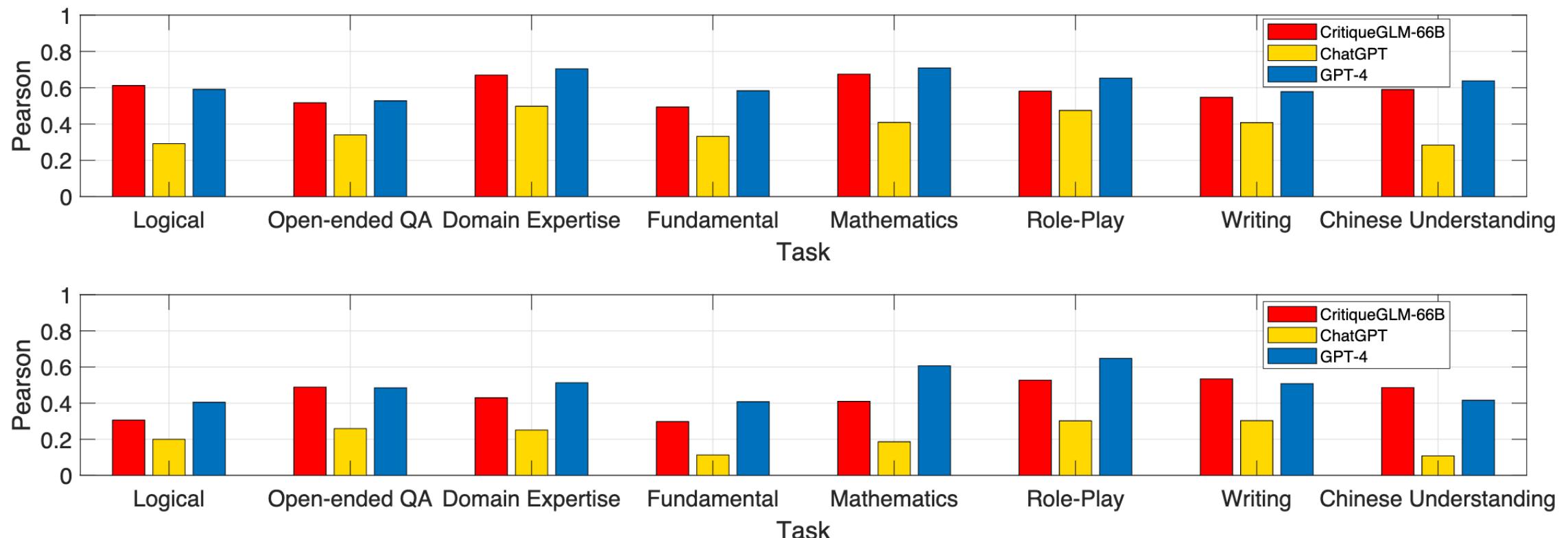


Figure 2: Text-level Pearson correlation on different tasks in referenced (**Top**) and reference-free settings (**Bottom**).

Summary

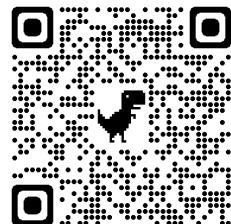
- GPT vs GLM
 - ChatGPT vs. ChatGLM
 - DALL.E vs. CogView
 - Codex vs. CodeGeeX
 - WebGPT vs. WebGLM
 - GPT-4V vs. GLM-4V (CogVLM, AgentTuning...)
- 2024-toward AGI

References

- Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan Wang, Yufei Xue, Lei Shen, Zihan Wang, Andi Wang, Yang Li, Teng Su, Zhilin Yang, and Jie Tang. **CodeGeeX**: A Pre-Trained Model for Code Generation with Multilingual Benchmarking on HumanEval-X. KDD’23.
- Xiao Liu, Hanyu Lai, Yu Hao, Yifan Xu, Aohan Zeng, Zhengxiao Du, Peng Zhang, Yuxiao Dong, and Jie Tang. **WebGLM**: Towards An Efficient Web-enhanced Question Answering System with Human Preference. KDD’23.
- Jing Zhang, Xiaokang Zhang, Daniel Zhang-Li, Jifan Yu, Zijun Yao, Zeyao Ma, Yiqi Xu, Haohua Wang, Xiaohan Zhang, Nianyi Lin, Sunrui Lu, Jie Tang, and Juanzi Li. **GLM-Dialog**: Noise-tolerant Pre-Training for Knowledge-grounded Dialogue Generation. KDD’23.
- Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, Weng Lam Tam, Zixuan Ma, Yufei Xue, Jidong Zhai, Wenguang Chen, Zhiyuan Liu, Peng Zhang, Yuxiao Dong, and Jie Tang. **GLM-130B**: An Open Bilingual Pre-trained Model. ICLR’23.
- Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. **CogVideo**: Large-scale Pretraining for Text-to-Video Generation via Transformers. ICLR’23.
- Ming Ding, Wendi Zheng, Wenyi Hong, and Jie Tang. **CogView2**: Faster and Better Text-to-Image Generation via Hierarchical Transformers. NeurIPS’22.
- Jifan Yu, Xiaohan Zhang, Yifan Xu, Xuanyu Lei, Xinyu Guan, Jing Zhang, Lei Hou, Juanzi Li, and Jie Tang. **XDAI**: A Tuning-free Framework for Exploiting Pre-trained Language Models in Knowledge Grounded Dialogue Generation. KDD’22.
- Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. **GLM**: General Language Model Pretraining with Autoregressive Blank Infilling. ACL’21.
- Zixuan Ma, Jiaao He, Jiezhong Qiu, Huanqi Cao, Yuanwei Wang, Zhenbo Sun, Liyan Zheng, Haojie Wang, Shizhi Tang, Tianyu Zheng, Junyang Lin, Guanyu Feng, Zeqiang Huang, Jie Gao, Aohan Zeng, JianWei Zhang, Runxin Zhong, Tianhui Shi, Sha Liu, Weimin Zheng, Jie Tang, Hongxia Yang, Xin Liu, Jidong Zhai, and Wenguang Chen. **BAGUALU**: Targeting Brain Scale Pretrained Models with over 37 Million Cores. PPOPP’22.
- Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, and Jie Tang. **CogView**: Mastering Text-to-Image Generation via Transformers. NeurIPS’21.
- Junyang Lin, Rui Men, An Yang, Chang Zhou, Yichang Zhang, Peng Wang, Jingren Zhou, Jie Tang, and Hongxia Yang. **M6**: Multi-Modality-to-Multi-Modality Multitask Mega-transformer for Unified Pretraining. KDD’21.

Thank you !

Many many collaborators from Tsinghua and Zhipu AI!



<https://github.com/THUDM/>