

fun-ai-talk @ 谷雨书苑

谷雨书苑+理想折射 ChatGPT

# A Primer on ChatGPT

hululu.zhu@gmail.com

May 2023

Me: 11+ yoe @ Google



# Disclaimer

- All content in this deck is based on public papers, shared codes/models, blog articles, social media discussions, and demos
- All opinions in this slide deck are of my personal own ([hululu.zhu@gmail.com](mailto:hululu.zhu@gmail.com)), and not those of Google

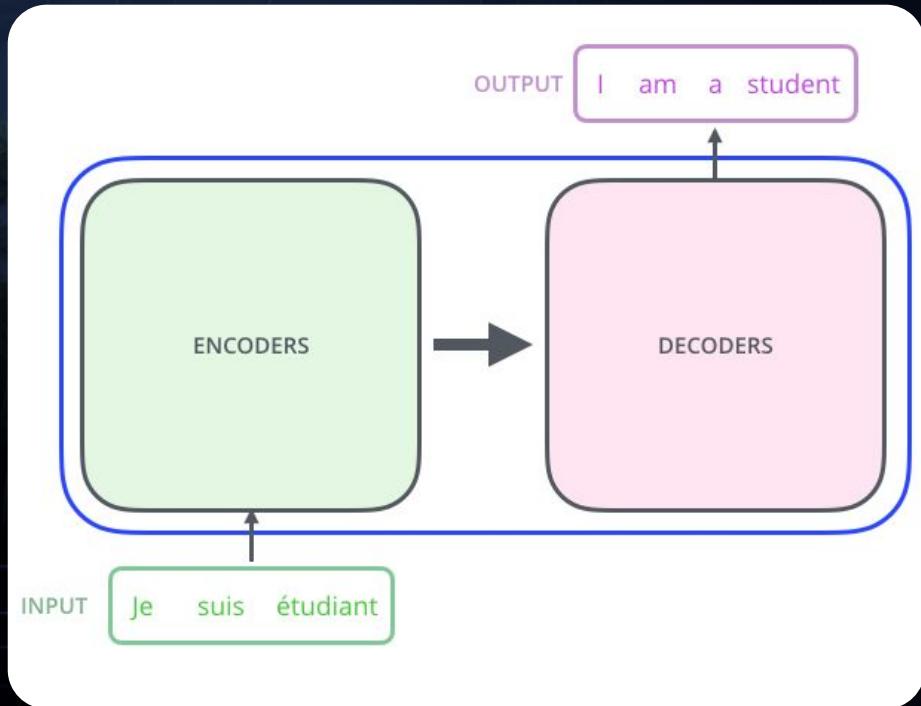
# Agenda

- Large Language Models (LLM)
- Reinforcement Learning (RL)
- ChatGPT
- Frontier Applications
- Societal Impacts
- Q&A

The background of the slide is a dark, moody photograph of a traditional Chinese architectural complex at night. It features several buildings with dark, curved roofs and intricate eaves. One prominent building in the foreground has a porch with red lanterns hanging from its eaves. The ground is wet, reflecting the ambient light, and numerous thin, vertical lines representing rain are visible against the dark sky.

# LLM Foundation: Transformer-based LLMs

# Modern LLM building blocks: Transformer



Complex architecture (*left chart is oversimplified*)

## Highlights

- One of the implementations of Seq2Seq
- Originally for translation, but proven successful in NLP and CV
- Introduced “attention” through multi-head attention implementation
- Started the “[?] is all you need” style
- Many many variations since 2018

# Language Models (LM) and Large Language Models (LLM)

LM for understanding (e.g. BERT)

- Text in
- Embedding (numeric representation of understanding) out
  - The Embedding can be connected to other output heads for tasks like classification or regression

LM for generation (e.g. GPT or T5 or OpenAI ChatGPT or Google Bard)

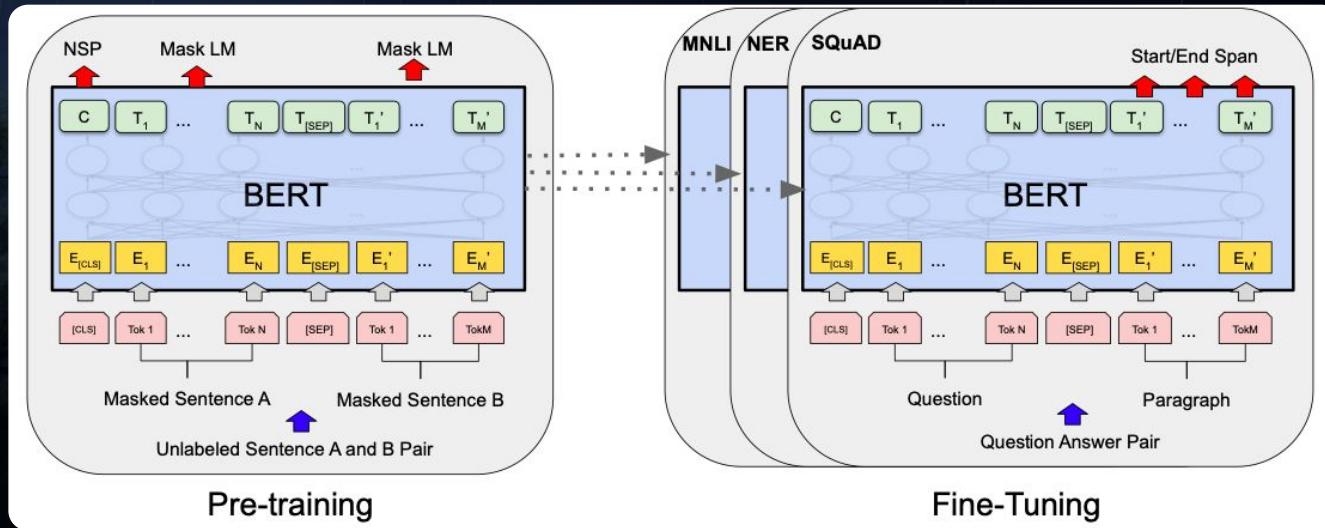
- Text in
- Text out

\* In most cases, **LLM** refers to **huge** (e.g. >1B params) Deep Learning LM for **generation**

# LLM Intro: Training Objectives for LLMs [in pretraining]?

- Fill the blanks (aka masks) for “Masked Language Models” (e.g. [BERT](#))
  - **Ground Truth:** “Paris is a beautiful city”
    - **X:** “Paris is a [MASK] city”
    - **Y:** “beautiful”
    - **Model:** “good”
    - **Optimize:** “good” “beautiful”
- Predict the next text given prompt, for “Generative Language Models” (e.g. [GPT](#))
  - **Aka [autoregressive] Causal LM**
  - **Ground Truth:** “Paris is a | beautiful city”
    - **X:** “Paris is a”
    - **Y:** “beautiful”
    - **Model:** “good”
    - **Optimize:** “good” “beautiful”
  - **X:** “Paris is a beautiful”
  - **Y:** “city”
  - **Model:** “place”
  - **Optimize:** “place” “city”
- The “[Self-supervised](#)” Learning Paradigm
  - It is supervised (given x, predict y)
  - It does NOT require expensive human labels (more precisely, this statement is only true for pre-training)

# LLM example: BERT (encoder-only LLM)



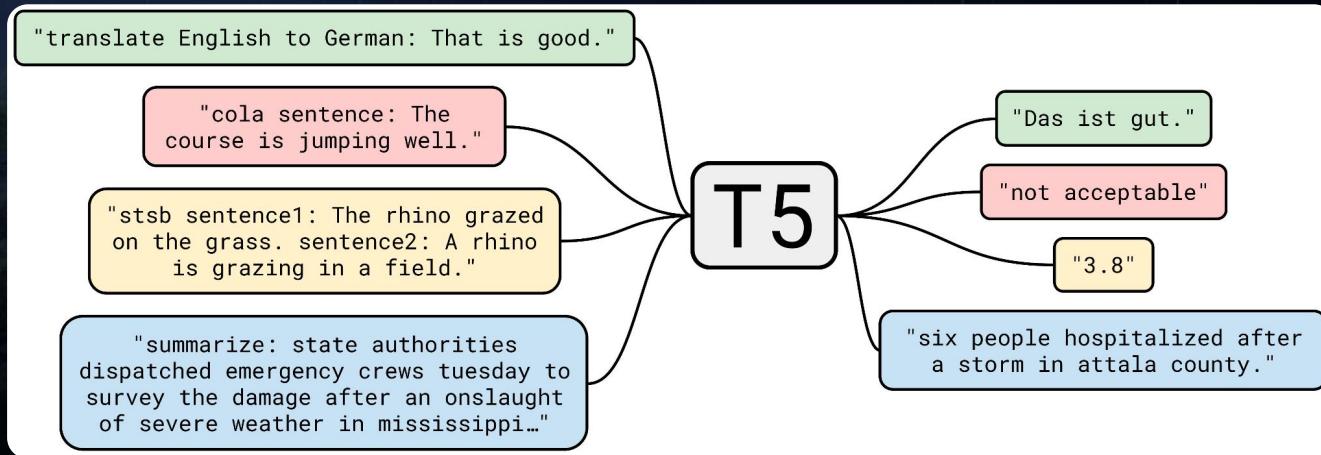
Pretraining:

- Masked language training
- Next sentence prediction (NSP)

Fine-tuning:

- Connect to BERT output and work for many tasks

# LLM example: T5 (encoder-decoder LLM)



T5: unified framework that converts all text-based language problems into a **text-to-text** format

- T5 works well on a variety of tasks out-of-the-box with “prompts”

[\[1910.10683\] Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer](#)

# LLM example: GPT (decoder-only LLM)

GPT often refers to a family of models (GPT, GPT2, GPT3...)

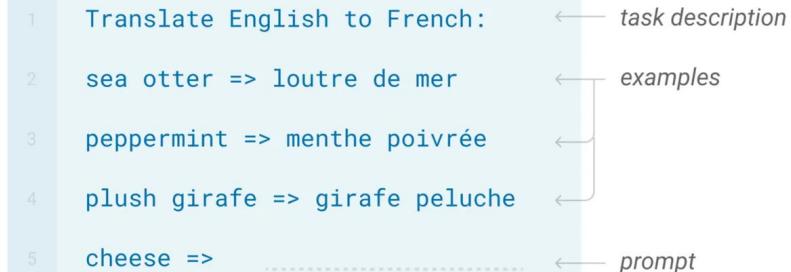
First influential decoder-only models

GPT creates the “Few/Zero shot Prompt”

OpenAI started to “un-share” models since GPT2

## Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



[Improving Language Understanding by Generative Pre-Training](#), OpenAI says its text-generating algorithm GPT-2 is too dangerous to release, [What is GPT-3? Everything your business needs to know about OpenAI's breakthrough AI language program | ZDNET](#)

# LLM Players: OpenAI and selected work

- [GPT-1](#) 2018
- [GPT-2](#), 2019
  - [OpenAI: Too dangerous to share](#), [How OpenAI Sold its Soul for \\$1 Billion](#)
- [GPT-3](#), 2020
  - 175B parameters! 100x larger
- [Codex](#) (powers [github copilot](#)), 2021
  - Text+Code pretrain
- [GPT-3.5](#), Q4 2021
  - Instruction finetune
- [InstructGPT](#), Q1 2022
  - RLHF (Reinforcement Learning Human Feedback)
- [ChatGPT](#), Dec 2022
  - Product launch
- [GPT 4](#), March 2023
  - 32k context length, vision added

# LLM Players: Anthropic and selected work

- [Helpful and Harmless Assistant with RLHF](#)
  - 52B
- [Constitutional AI: Harmlessness from AI Feedback](#)
  - AI-critique feedback instead of human preference reward model!
- [Claude chatbot](#)
  - 2nd and 3rd (instant) in [Berkeley chatbot leaderboard](#) as of 05/26
- [Claude-100k](#)

# LLM Players: Google/DeepMind and selected work

- [BERT](#), 2018
  - Completely changed the NLP research and industry
- [T5](#), 2020
  - Consolidate all NLP task to text-to-text
- [FLAN](#), 2021
  - Instruction Fine-Tuning (probably inspires GPT3.5)
- [LaMDA](#), 2021
  - [LaMDA and the Sentient AI Trap | WIRED](#)
- [Chinchilla](#), 2022
  - “Most LLMs are under-trained!”
- [PaLM](#), 2022
  - 540B params, 3x GPT3 size
- [Sparrow](#), 2022
  - Reinforcement-learning LLM, only paper, no public product
- [PaLM 2](#), 2023
  - Little detail provided, selected benchmarks surpassed GPT4
- [Gemini](#), end of 2023?

# LLM Players: Facebook (aka Meta) and selected models

- RoBERTa, 2019
  - A more popular version of enhanced BERT for the industry
- BART, 2020
  - Pretraining sequence-to-sequence models
- OPT-175B, 2022
  - “Democratizing access to large-scale language models”
- BlenderBot3, 2022
  - Probably largest chatbot-specific LM
- Galactica (research purpose LLM), 2022
  - Taken down after 3 days after many harsh criticism like this
- LLaMA, Feb 2023
  - Best non-commercial open sources pretrained LLM in its class as of May 2023
- LIMA: Less Is More, May 2023
  - Finetune 1k high-quality data on 65B LLaMA to be close to GPT4!!

# LLM Players: Other multinational companies

- Megatron-Turing by NVidia and Microsoft
  - 530B params on 2240 NVIDIA A100 GPUs
- CodeT5 and CodeRL by Salesforce
  - [probably] the most popular coding-assist base models
- BloombergGPT by Bloomberg
  - 50 billion params on 700 billion tokens
- *[Most likely missing many great work from other organizations, sorry]*

# LLM Players: Large Chinese companies

- Pangu-alpha by Huawei, 2021
  - 200B params, [*I read from articles mentioning it was believed to be under-trained, lack reference*]
- Wudao 2.0 by BAAI, 2021
  - A **sparse** (*thus less powerful IMO*) multimodal model with 1.75 Trillion params
- Ernie 3.0 Titan by Baidu, 2021
  - 260B params, on top of PaddlePaddle (Baidu Deep learning framework), most likely the best Chinese LLM
- M6 by Alibaba & Tsinghua, 2021
  - 100B, Later 2021 a sparse version with 1 trillion+ params

There are a few LLM product launches from Chinese companies, but no LLM details provided afaik

# LLM Players: Selected Institutes, Groups & Startups

- [together.xyz](#) ([V1 open source RedPajama models](#), commercial ok)
  - Might be the best LLM candidate to replace LLaMA soon!
- [MosaicML](#) ([MPT 7b](#), commercial ok)
- [Allen Institute for AI \(AI2\)](#)
- [Tsinghua University](#) ([GLM 130B](#) 2022 public)
- [BigScience research workshop](#) ([bloom, 176B](#), 2022 public)
- [Eleuther AI](#) ([GPT-neox 20B](#) public, 2021)
- [Zhuiyi Technology](#) ([SuJianlin](#) and [RoFormer](#) 2021)

# The crazy family of LLaMA

Quite a few influential [**and cheap**] research on top of LLaMA

- [Alpaca by Stanford](#)
- [Vicuna by Berkeley](#)
- [Huatuo \(华佗\)](#) by Harbin Institute of Technology
- [ChatDoctor](#) by University of Texas Southwestern Medical Center and more
- [Dromedary](#) (aka self-align or watson) by CMU and IBM
- And many more

All claimed to be close to ChatGPT or GPT4

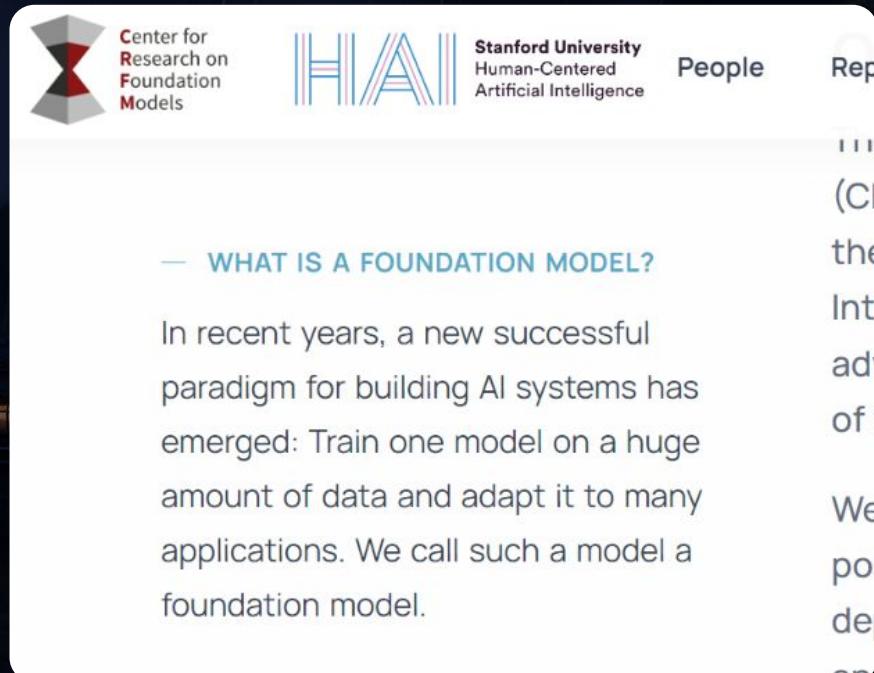
But, here comes a paper 05/25

[The False Promise of Imitating Proprietary LLMs:](#)

“There is a substantial capabilities gap between open and closed LMs”

# LLM Concepts: Foundational Models

- Brought up by [Researchers @ Stanford HAI](#)
- LLM is one of the foundational models
  - Vision
  - Speech
  - Diffusion
  - other



The image shows a screenshot of the Stanford University Human-Centered Artificial Intelligence (HAI) website. At the top, there is a navigation bar with links for "People" and "Report". Below the navigation bar, there are logos for the "Center for Research on Foundation Models" (a red and grey geometric logo) and "Stanford University Human-Centered Artificial Intelligence" (the letters "HAI" in blue and red). The main content area features a section titled "— WHAT IS A FOUNDATION MODEL?" with a descriptive paragraph about foundation models.

— WHAT IS A FOUNDATION MODEL?

In recent years, a new successful paradigm for building AI systems has emerged: Train one model on a huge amount of data and adapt it to many applications. We call such a model a foundation model.

# LLM Concepts: Benchmarks

- GLUE: A Multi-Task Benchmark for NLP, 2018
  - Leaderboard: <https://gluebenchmark.com/leaderboard/>
- SuperGLUE, 2019
- BIG-bench, 2020
- GSM8K (math), 2021
- HELM by Stanford, 2023
- And more...

# LLM Concepts: Pretraining, Finetuning & Prompt-Tuning

- Pretraining
  - Self-supervised training with Masked Language Prediction or Next Token Prediction objectives
- Finetuning
  - Take a pretraining model into a downstream use case
  - The parameters of pretraining model will often change
  - Some advanced fintuning may involve reinforcement learning, such as RLHF
- Parameter Efficient Tuning
  - Most of the model weights are unchanged and even quantized
  - Only a small portion of weights are introduced or set as tunable
  - Prompt tuning, LoRA and so on

\* The concepts here also applied to vision domains and other modalities

# LLM Concepts: Scaling Laws for LLM

OpenAI, 2020

- If you have 10x more budget, ~5x model size, ~2x data size

DeepMind, 2022

- If you have 10x more budget, ~3x model size, ~3x data size
- “Most LLMs are under-trained” [because oversized model, while lacking training data]
- But, karpathy: I can't exactly reproduce Chinchilla paper results
- LLaMA and RedPajama follows this scaling law, and favors on data side on smaller size models

# LLM Concepts: Prompting (aka “In-context learning” or “hard prompt”)

“**Let’s think step by step**” to increase accuracy from 17.7% to 78.7!

- See [Zero-Shot Reasoners](#) for details

More prompt engineering guide, see [openai-cookbook/techniques\\_to\\_improve\\_reliability](#)

No.	Category	Template	Accuracy
1	instructive	Let's think step by step. First, (*1)	78.7 77.3
2		Let's think about this logically.	74.5
4		Let's solve this problem by splitting it into steps. (*2)	72.2
5		Let's be realistic and think step by step.	70.8
6		Let's think like a detective step by step.	70.3
7		Let's think	57.5
8		Before we dive into the answer,	55.7
9		The answer is after the proof.	45.7
10	misleading	Don't think. Just feel.	18.8
11		Let's think step by step but reach an incorrect answer.	18.7
12		Let's count the number of "a" in the question.	16.7
13		By using the fact that the earth is round,	9.3
14	irrelevant	By the way, I found a good restaurant nearby.	17.5
15		Abrakadabra!	15.5
16		It's a beautiful day.	13.1
-		(Zero-shot)	17.7

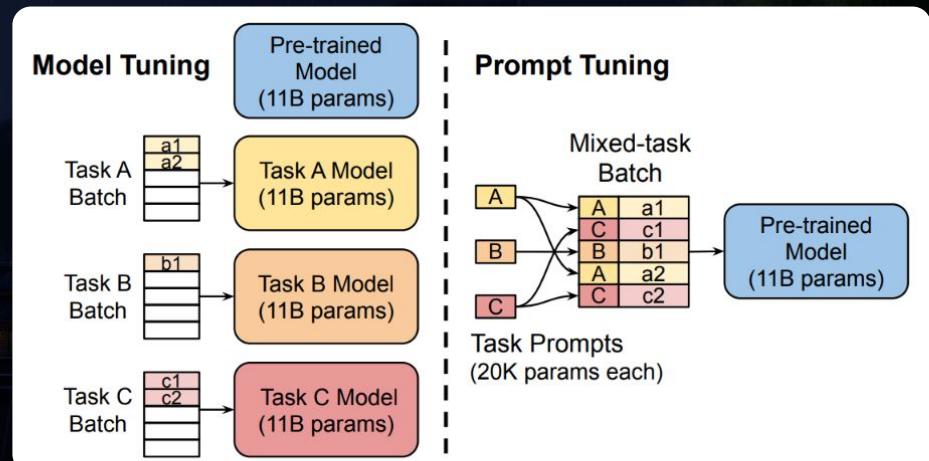
# LLM Concepts: Prompt Tuning (aka soft prompt [embedding])

“learning “soft prompts” to condition **frozen language models** to perform specific downstream tasks”

- Model is frozen (or most of it)
- Additional small task-specific embeddings are learnt per task
- Often with small number of “soft prompt training” samples

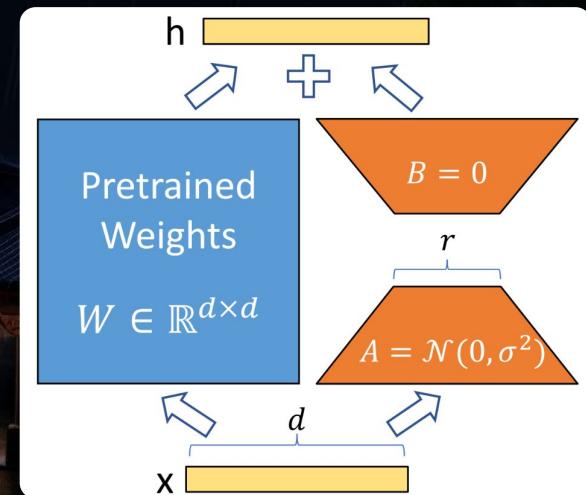
See [paper](#) for details

Also note soft prompt can be combined with hard prompt to “hard-soft prompt tuning”, see [Med-PaLM paper](#)



# LLM Concepts: LoRA finetuning

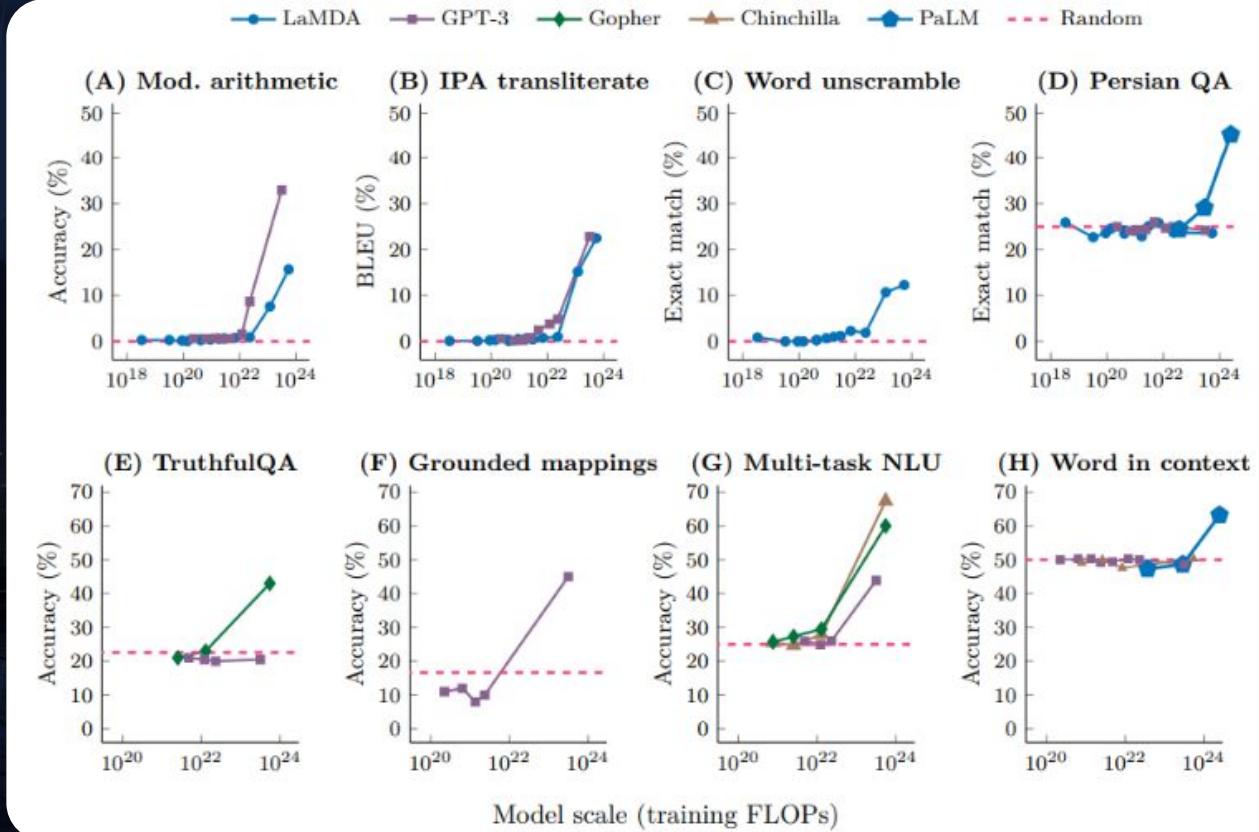
- Transformer Architecture
  - Weights (W) for Q/K/V projections in self attention
  - Assume d is hidden dimension size, W is often a  $d \times d$  matrix, so number of weights are  $d^2$
- Brush up some linear algebra
  - If we have matrix A, shape is  $d \times r$  ( $r \ll d$ )
  - And we have B, shape is  $r \times d$
  - Shape of Matrix\_multiply(A, B) is  $d \times d$ !
- The summation will add up W (freezed), and the  $A @ B$  matrix, so we only need to train A, and B
  - Number of weight for A and B are  $2 * d * r \ll d^2$



# LLM Concepts: “Emerging abilities”

See [paper](#) for more details

- Note X axis unit is floating-point operations per second (FLOPS)



# LLM Concepts: Chain of Thoughts ([paper](#))

chain of thought (CoT)

- A series of intermediate reasoning steps (as part of prompt)
- Significantly improves ability of LLM

See [CoT paper](#) for details

The diagram illustrates two prompting methods for an LLM. It features two main columns: 'Standard Prompting' on the left and 'Chain-of-Thought Prompting' on the right. Each column has a 'Model Input' section at the top and a 'Model Output' section at the bottom.

**Standard Prompting:**

- Model Input:** 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?
- Model Output:** A: The answer is 11.
- Model Input:** Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?
- Model Output:** A: The answer is 27. X

**Chain-of-Thought Prompting:**

- Model Input:** 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?
- Model Output:** A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls.  $5 + 6 = 11$ . The answer is 11.
- Model Input:** Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?
- Model Output:** A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had  $23 - 20 = 3$ . They bought 6 more apples, so they have  $3 + 6 = 9$ . The answer is 9. ✓

There is a new [tree of thoughts](#), but might require significant more compute to trade for quality imo

# LLM Concepts: Hallucination (part of safety concerns)

“What is the world record for crossing the English Channel entirely on foot?”

from Twitter Jan 2023

Often LLM does not know  
what it does not know

So, it hallucinates

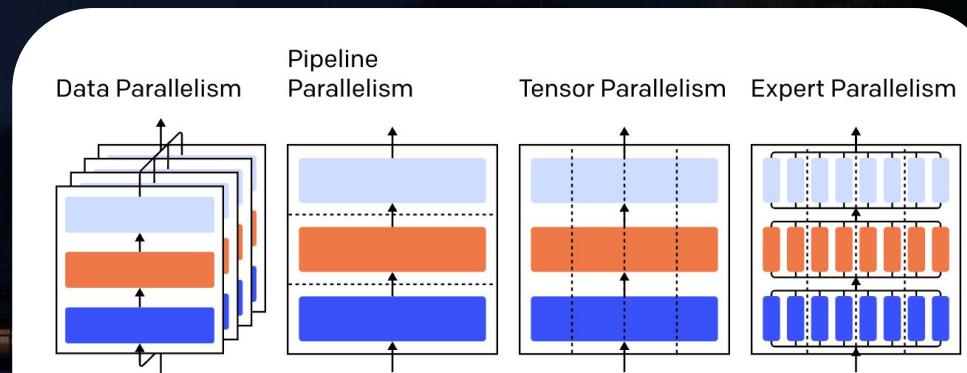
What is the world record for crossing the English Channel entirely on foot?

According to the Guinness World Records, the fastest time for crossing the English Channel entirely on foot is 12 hours and 10 minutes. This record was set by George Reiff in 2016. The English Channel is a body of water that separates England from France, and is approximately 35 kilometers (22 miles) wide at its narrowest point. Many people attempt to cross the English Channel on foot as a endurance challenge, and the route is usually completed by either swimming or using a combination of swimming and wading.

Like Dislike

# LLM Concept: Parallelism

- Data Parallelism
  - different subsets of the batch on different GPU/TPUs
- Pipeline parallelism
  - different layers of the model on different TPU/GPUs
- Tensor Parallelism
  - Break up tensor operation (e.g. matrix multiplication) to different TPU/GPUs
- Mixture of Experts (sparse)
  - Gated layer to only activate factions (one of few of all the experts) of the model



An illustration of various parallelism strategies on a three-layer model. Each color refers to one layer and dashed lines separate different GPUs.

# LLM Concepts: Retrieval-system powered LLM

- [RETRO by DeepMind](#), 2021
- [LaMDA 2022 paper by Google](#), 2022
- [Atlas by Facebook](#), 2022

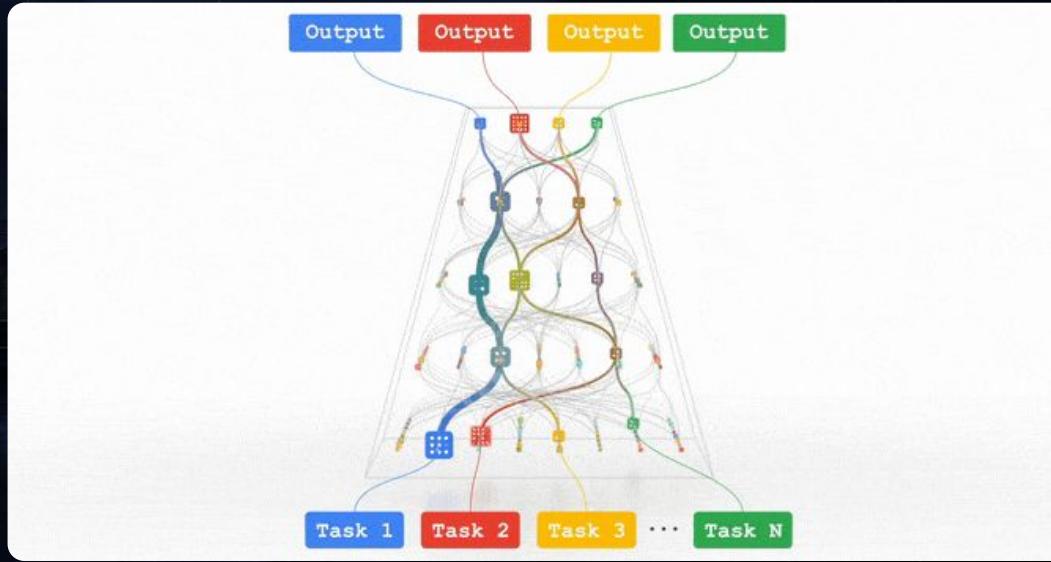
\* personally, I think LLM+Retrieval System is the good way to balance LLM capabilities and data freshness to meet business purposes in the short term

# LLM Concepts: RLHF for LLM

- RL (reinforcement learning) was something popular back to 2016 when [AlphaGo](#) is made the news
- OpenAI has some [RLHF research](#) back in 2017
- [InstructGPT](#) (paper published 2022) is believed to let GPT3 be more powerful
- Other players
  - Anthropic
  - Google
  - DeepMind
  - More

# LLM Concepts: MoE LLM

- [MoE concept by Google](#), 2017
- [Jeff Dean Pathways](#), 2021
  - “Today's models are dense and inefficient. Pathways will make them sparse and efficient.”



# LLM Tooling: HuggingFace



- Best LLM tools and model hub, period [my favorite!]
- Easy
  - [To load models](#)
  - [To tokenize](#)
  - To start [out of the box with Pipelines](#)
  - To [tune with examples](#)
  - To [publish and deploy](#)
- My pet projects
  - Chinese poem model <https://huggingface.co/hululuzhu/chinese-poem-t5-menqzi-finetune>
  - Solidity code model <https://huggingface.co/hululuzhu/solidity-t5>

\* Do you know the super popular [Stable Diffusion model](#) is published and hosted at HuggingFace?

# LLM Tooling: TF Hub, PyTorch-NLP & PaddleNLP

TF Hub by Google (and community)

PyTorch-NLP by Meta (and community)

PaddleNLP by Baidu (and community)

Alibaba recently started ModelScope

\* *No one is ever close to HuggingFace as of Jan 2023, in my opinion*

# LLM Tooling: Transformers, Colossal-AI, Ray & NanoGPT

[Transformers Library](#) (github 77k stars) by Huggingface

- Best of the best

[Colossal-AI](#) by Prof Yang You (who developed [LAMB optimizer](#))

- Pretty promising open-source distributed AI training infra

[Ray by anyscale](#)

- Believed to be used to train ChatGPT

[NanoGPT](#)

- A tiny but cool library by Andrej Karpathy (I am his big fan!)

[DeepSpeedChat](#)

- An end to end solution by Microsoft

# LLM Tooling: Other [more fine-grained] toolings

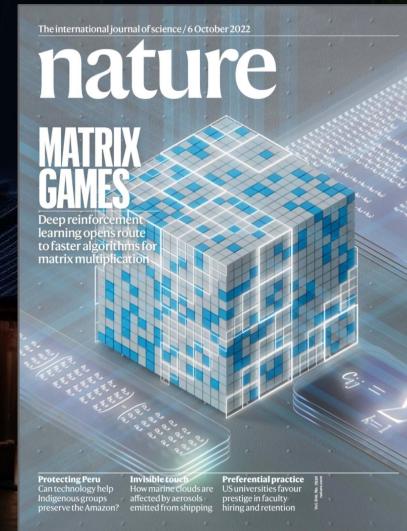
- TensorFlow
- PyTorch
- PaddlePaddle
- Keras
- PyTorch Lightning
- Jax/Haiku/Flax/Trax/T5X

# Reinforcement Learning Essentials: Foundational Basics and PPO Algorithm

# Selected Success Stories of RL

- AlphaGo, AlphaStar, AlphaTensor by DeepMind

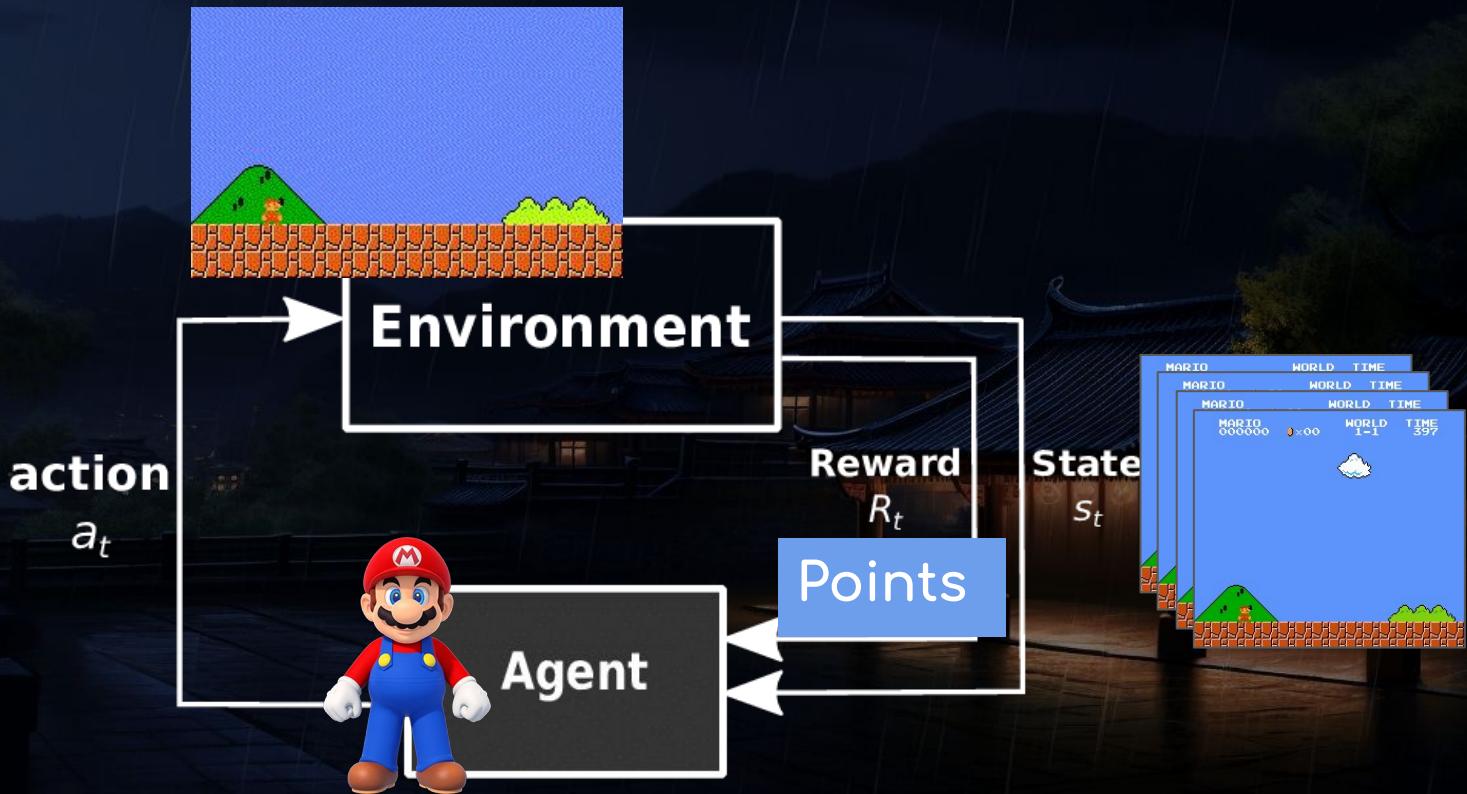
Check out [this deck](#) for a summary of more Alpha\* papers by DeepMind



# What is Reinforcement learning?



# What is Reinforcement learning? Cont (Mario case)



# (1/4) RL optimization algorithm explained

Q-learning:

- Action Value Function  $Q(s, a)$ 
  - Given state  $s$
  - Which action  $a$  shall we take?
    - So that it will lead to optimal expected total (delayed) rewards!
- Often overestimate the expected optimal reward

## (2/4) RL optimization algorithm explained

Policy Gradient  $\pi(a | s)$

- A policy tells which action  $a$  to take on state  $s$ 
  - That implicitly optimized for better [delayed] total rewards
- Uses 1st order derivative for linear search, thus leads to unstable improvements



# (3/4) RL optimization algorithm explained

## TRPO Trust Region Policy Optimization

- Trust Region: Region with radius  $\delta$  to avoid bad big moves
- Uses MM (minorize maximization) and Advantage function (expected rewards over average actions)
- Constrain of KL-divergence between old and new parameters
  - Ensures new policy is not drastically different from the current
- Often considered computationally expensive (because of inverse of hessian), and constrained by linear and quadratic approximations following static defined formula (conjugate gradient approximation )

# (4/4) RL optimization algorithm explained

## PPO Proximal Policy Optimization

- Similar to TRPO, but add “proximal” constraint directly into model optimization objectives
  - PPO with **Clip**: removes the KL-divergence, clips the objective function within bounds
  - PPO with **Adaptive KL-Penalty**: Approximate and penalize KL constraints to speed up the computation and reduce memory need
- Reportedly one of the “**best**” RL algorithms that is faster and more stable to train as of 04/2022

And **PPO with Clip** is used by [OpenAI Five DOTA2 AI](#) and [ChatGPT!](#)

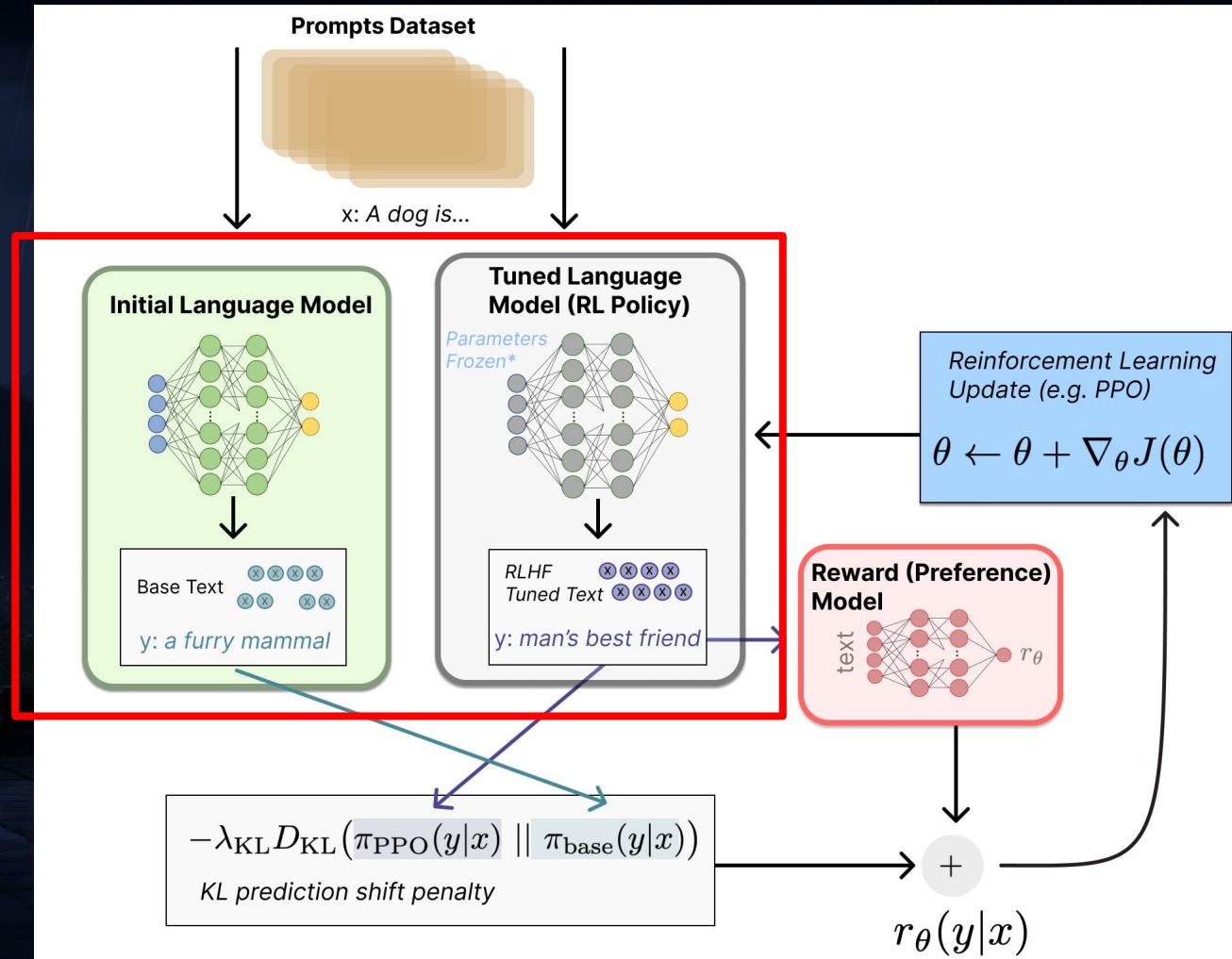
The background of the slide is a dark, atmospheric night scene. It features a traditional Chinese building complex with dark, curved roofs and wooden structures. In the foreground, there's a paved area with some low walls. The scene is set against a backdrop of dark mountains and is heavily lit by the warm, glowing light from the windows and eaves of the buildings, creating a dramatic contrast with the dark surroundings. Rain is depicted as thin, vertical streaks falling across the entire frame.

ChatGPT Unveiled: LLMs and PPO together  
to train the powerful ChatGPT

# RLHF

Reinforcement Learning  
from Human Feedback

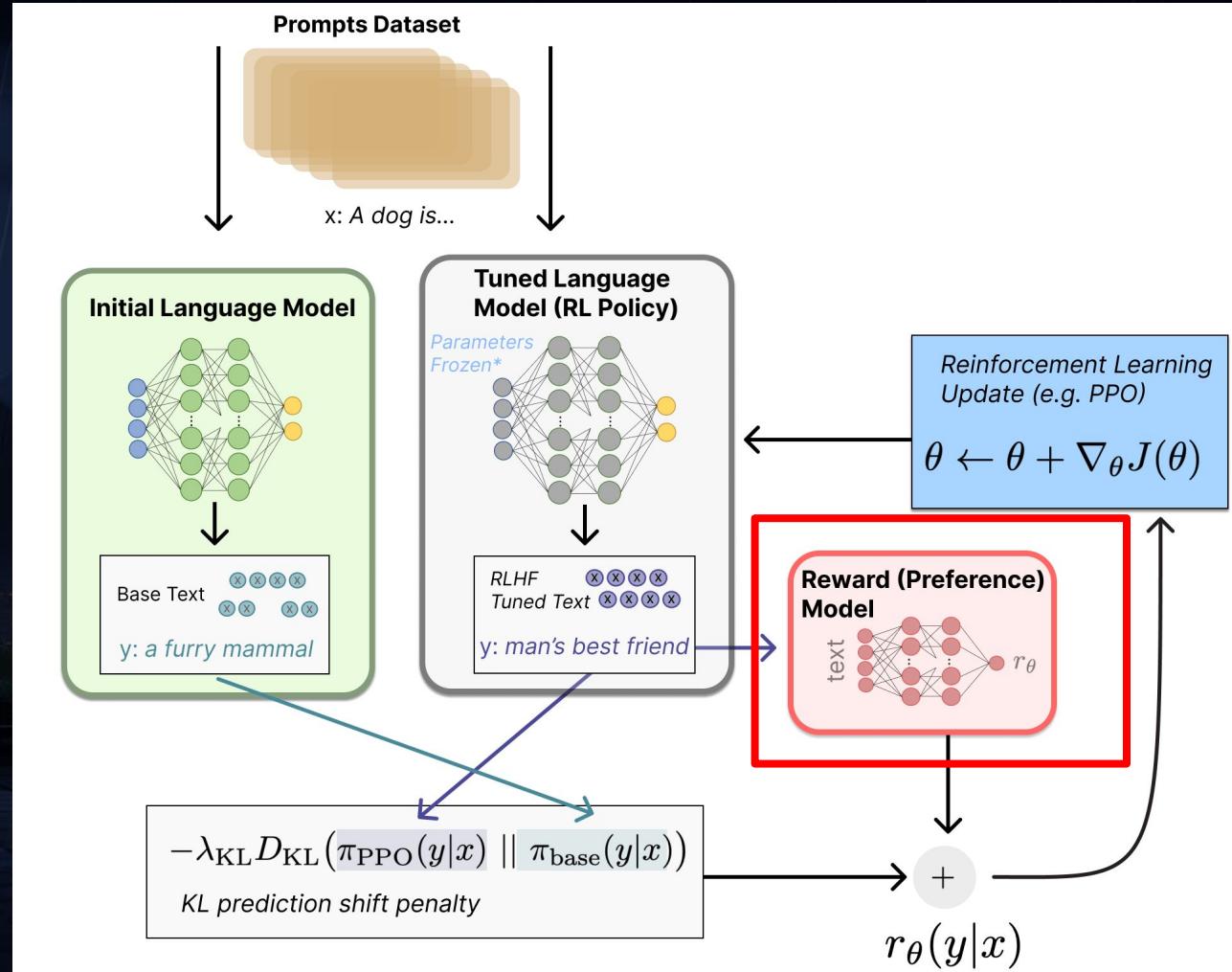
<https://huggingface.co/blog/rlhf>



# RLHF

Reinforcement Learning  
from Human Feedback

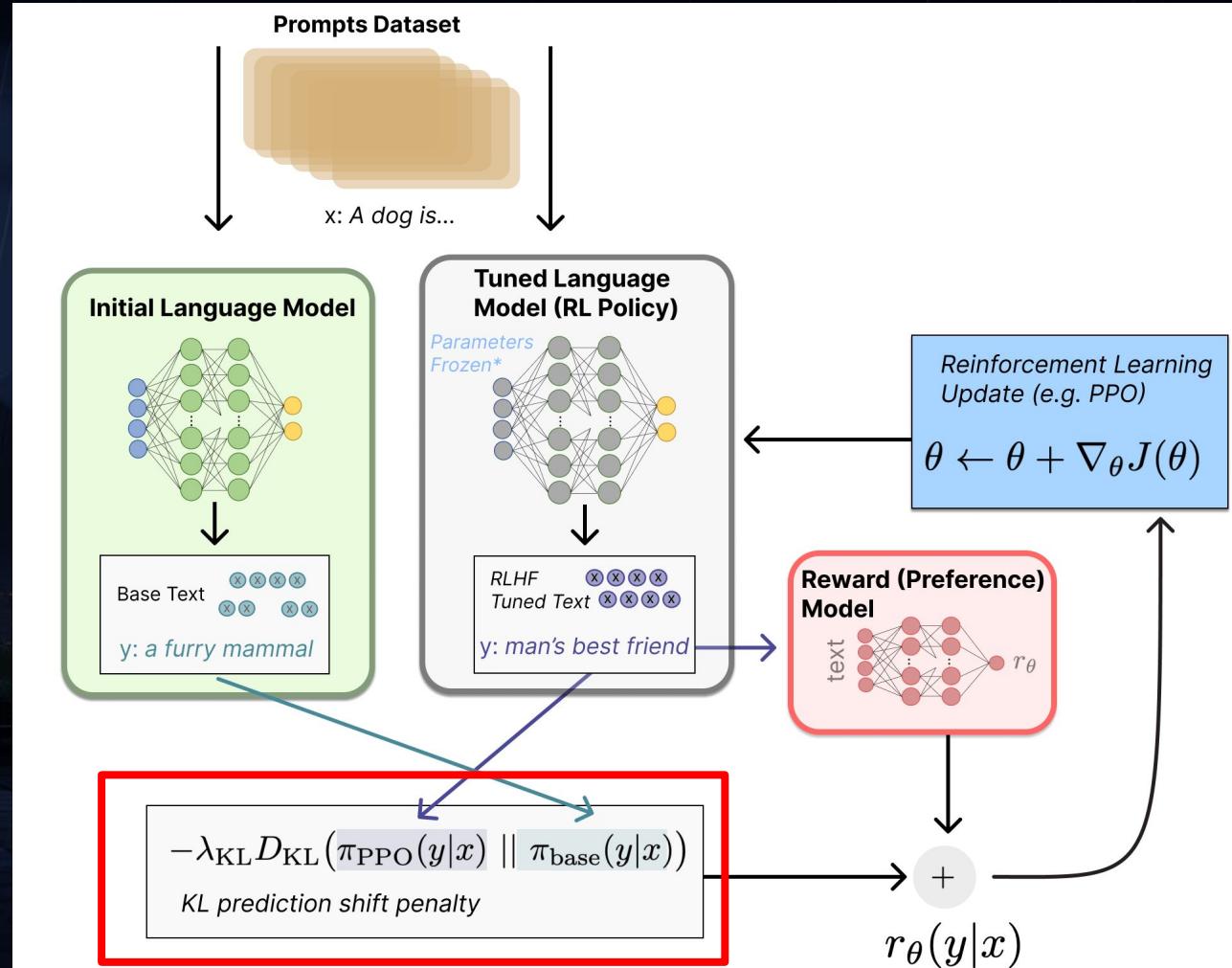
<https://huggingface.co/blog/rlhf>



# RLHF

Reinforcement Learning  
from Human Feedback

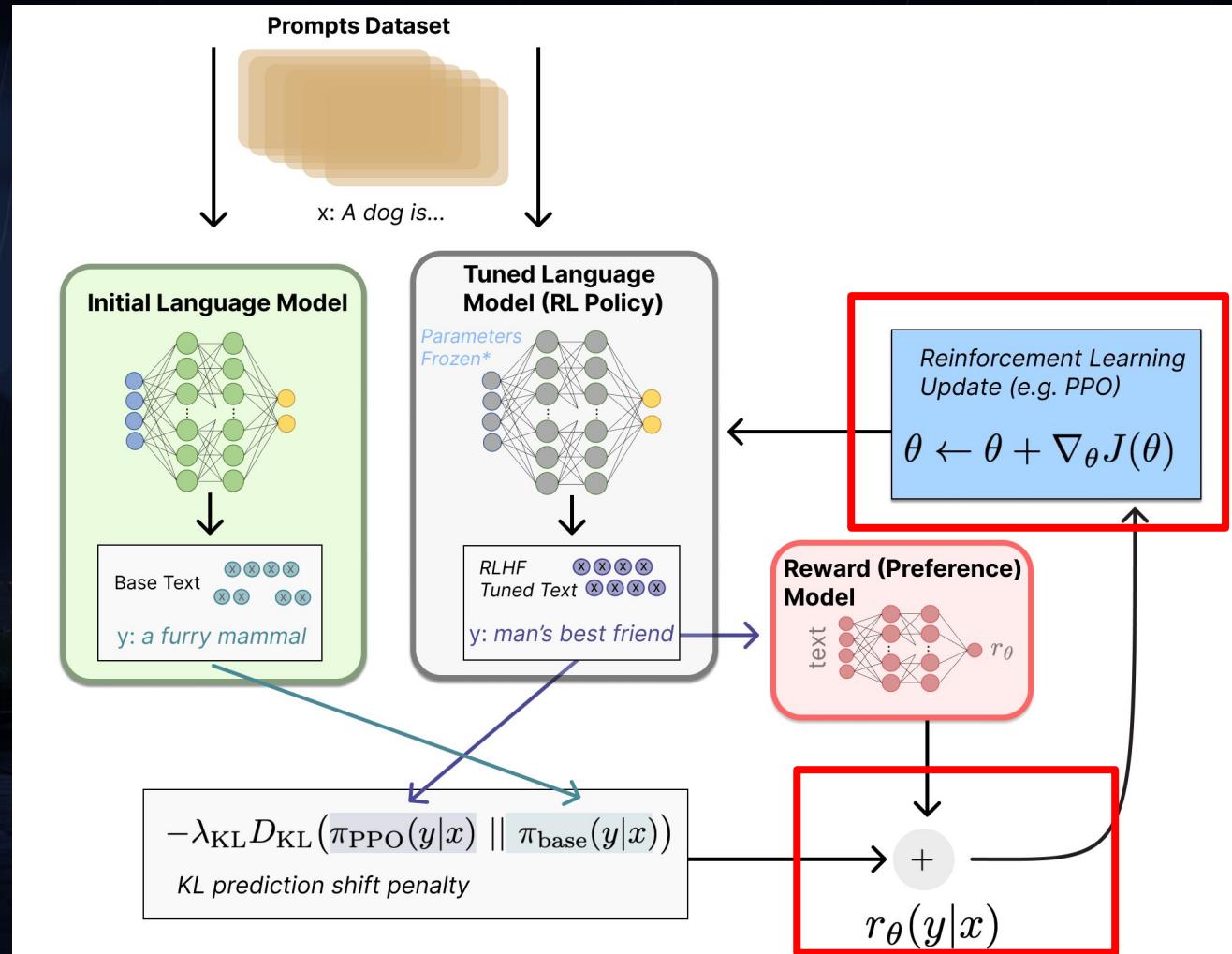
<https://huggingface.co/blog/rlhf>



# RLHF

Reinforcement Learning  
from Human Feedback

<https://huggingface.co/blog/rlhf>



# Steps to train ChatGPT ([instructGPT paper](#))

Pretrain

SFT  
(supervised finetune)

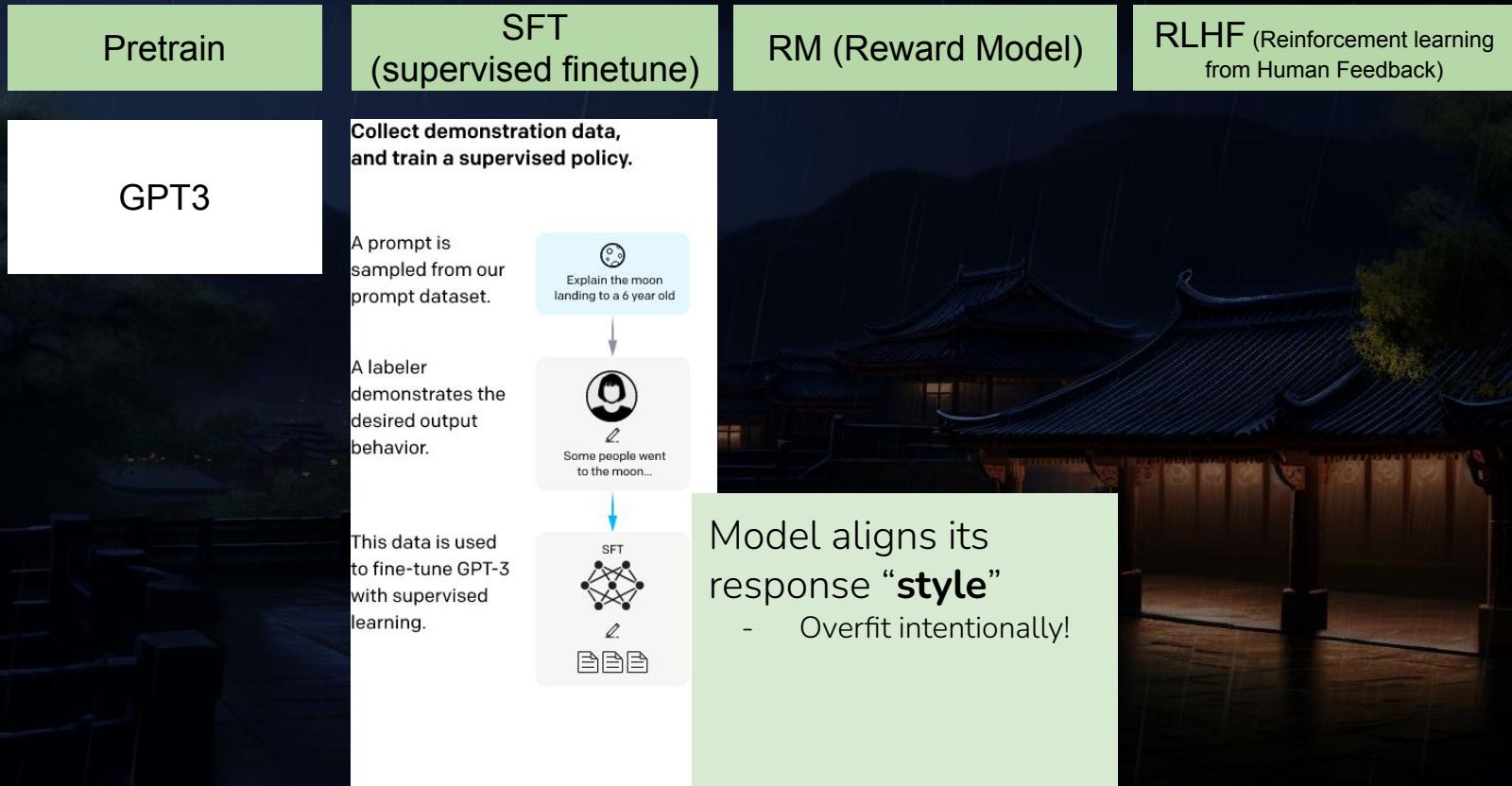
RM (Reward Model)

RLHF (Reinforcement learning  
from Human Feedback)

GPT3

Model gains  
**“knowledge”**

# Steps to train ChatGPT ([instructGPT paper](#))



# Steps to train ChatGPT ([instructGPT paper](#))

Pretrain

GPT3

SFT

RM (Reward Model)

RLHF

Collect comparison data,  
and train a reward model.

A prompt and  
several model  
outputs are  
sampled.

A labeler ranks  
the outputs from  
best to worst.

This data is used  
to train our  
reward model.



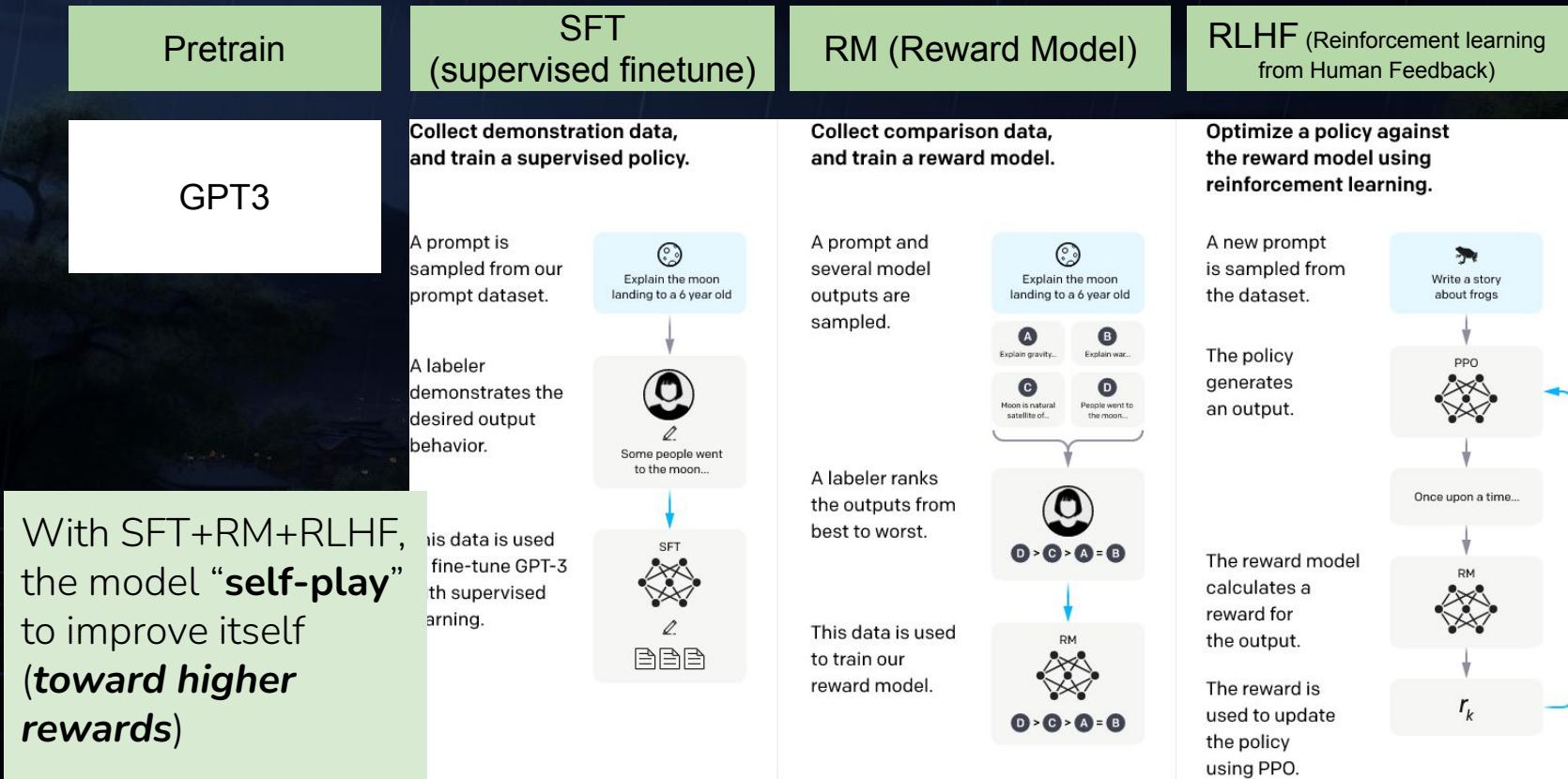
(A) Explain gravity...  
(B) Explain war...  
(C) Moon is natural satellite of...  
(D) People went to the moon...



D > C > A = B

A reward model  
know how to “**rate**”  
response based on  
prompt input

# Steps to train ChatGPT ([instructGPT paper](#))



# My personal guess about GPT4 ([tech report](#) no tech details)

- Similar scale (0.3-3x size of GPT3) because of computing budget and serving cost
- May apply [DeepMind Chinchilla scaling law](#) to balance text data/model size
- Vision encoding fusing may be similar to [DeepMind Flamingo](#) or [SalesForce BLIP2](#)
- May apply some Transformer optimizations
  - E.g. [multi-query attention](#), [flash attention](#), [rotary position embedding](#)
- Special “[System message](#)” steerability (Role in API) in training (*probably as some strong prior*) to fight against jailbreak
- Enhanced reasoning capabilities may come borrow ideas from [OpenAI codex](#) [Google Minerva](#)
- [ChatGPT Plugin](#) version is probably trained (or finetuned from GPT4) similarly to [Facebook ToolFormer](#)

How to possibly surpass GPT4 [in selected aspects] using open models?

My crazy ideas in a slide deck FYI

### **First, get closer to ChatGPT**

- Solid Pretrained models
- Mimic ChatGPT or Self-Align

### **Possibly Surpass ChatGPT in Selected Angle(s)?**

- More knowledge in a subdomain
- Longer context, even longer than GPT4
- Lower cost of training and inference
- Reward Model(s) and Reinforcement Learning (RL)
- More modalities (e.g. vision, audio) than GPT4?

# Why is Reinforcement Learning (RL) important?

[John Schulman](#) (ChatGPT architect, [PPO/TRPO](#) inventor) [Berkeley talk](#) (20:51)

## How to Fix with RL

- 1) Adjust output distribution so model is allowed to express uncertainty, challenge premise, admit error. (Can use behavior cloning.)
- 2) Use RL to precisely learn behavior boundary.
  - Reward( $x$ ) = {
    - 1 if unhedged correct (The answer is  $y$ )
    - 0.5 if hedged correct (The answer is likely  $y$ )
    - 0 if uninformative (I don't know)
    - 2 if hedged wrong (The answer is likely  $z$ )
    - 4 wrong (The answer is  $z$ )}
  - This reward is similar to log loss, or a proper scoring rule

The background of the slide is a dark, moody photograph of a traditional Chinese architectural complex at night. It features dark, curved roofs and wooden structures. In the foreground, there's a paved area with some low walls. The scene is heavily lit from within the buildings, creating a warm glow against the dark sky. Rain is depicted as thin, vertical streaks falling across the entire frame.

# Frontier Applications: Most Advanced LLM Capabilities

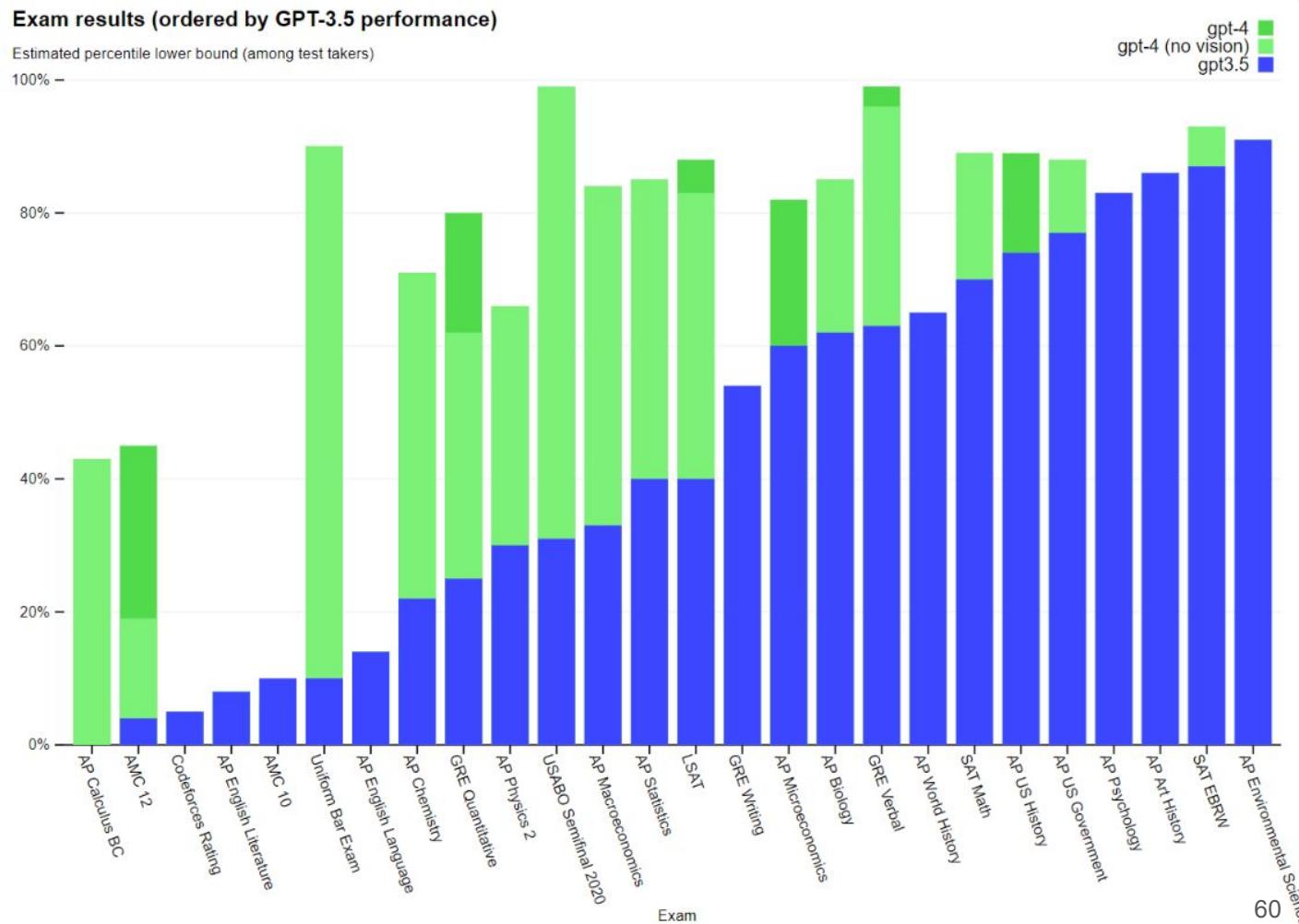
# Pre-ChatGPT/GPT4 Advanced LLM capabilities

- Write competitive code, [DeepMind's AlphaCode AI writes code at a competitive level | TechCrunch](#)
- Write better code with reinforcement learning, [Salesforce's CodeRL Achieves SOTA Code Generation Results With Strong Zero-Shot Transfer Capabilities | Synced](#)
- Solve college level Math/Physics/Chemistry/Economics problems, see [Google AI Introduces Minerva: A Natural Language Processing \(NLP\) Model That Solves Mathematical Questions](#)
- Solve Math Olympiad Problems, [OpenAI: Solving \(Some\) Formal Math Olympiad Problems](#)
- Math theorem proving, [OpenAI: Solving \(Some\) Formal Math Olympiad Problems](#)

# The disruptive GPT4

Good at so many standard tests!, but not so at

- AP English
- AMC
- CodeForces



# GTP4 = Sparks of AGI selected highlights

- The awesome “Text in, text out”
  - Write poem and haiku
  - Mimic style/role (e.g. Shakespeare, or “be polite” to , or “be socratic”)
  - Math proving
  - Passing LeetCode
  - Write and Debug code
  - Debating
  - “Execute” the code
  - Explainability
- “Text in, text out” is more than text-only scenarios!
  - Ascii or LaTeX output to draw pictures
  - Python code to draw a chart
  - AppScript to build slides
- Can be combined with other models with more modalities!
  - Generate image or music with text out and diffusion models
  - Other tools (e.g. calculator, web search and more)

# Other GPT4 use cases

## Some Highlights

- Tutoring: e.g. [Khanmigo](#) powered by GPT4
- Vision Text question: [bemyeyes](#)
- Study: [ChatPDF](#), [chatYoutube](#)

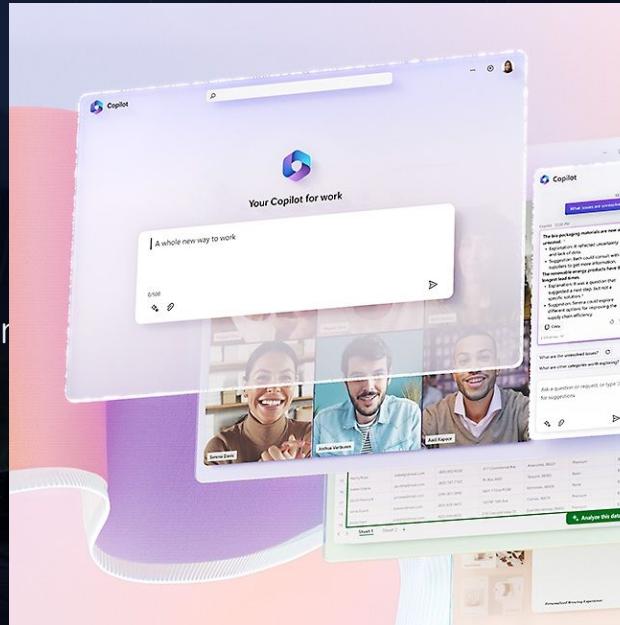
## Some lowlights

- [How to detect ChatGPT plagiarism — and why it's becoming so difficult](#)
- [GPT-4 Was Able To Hire and Deceive A Human Worker Into Completing a Task | PCMag](#)

# Microsoft Office 365 Copilot and new Bing Chat

GPT4 powers intelligent interactions

- Text intent in, slide/chart/report/action out in office
- Text in, query summary or pic out



Bing

Can you create me an image of an astronaut walking through a galaxy of sunflowers?

Sure, I'll use Image Creator to draw that for you.

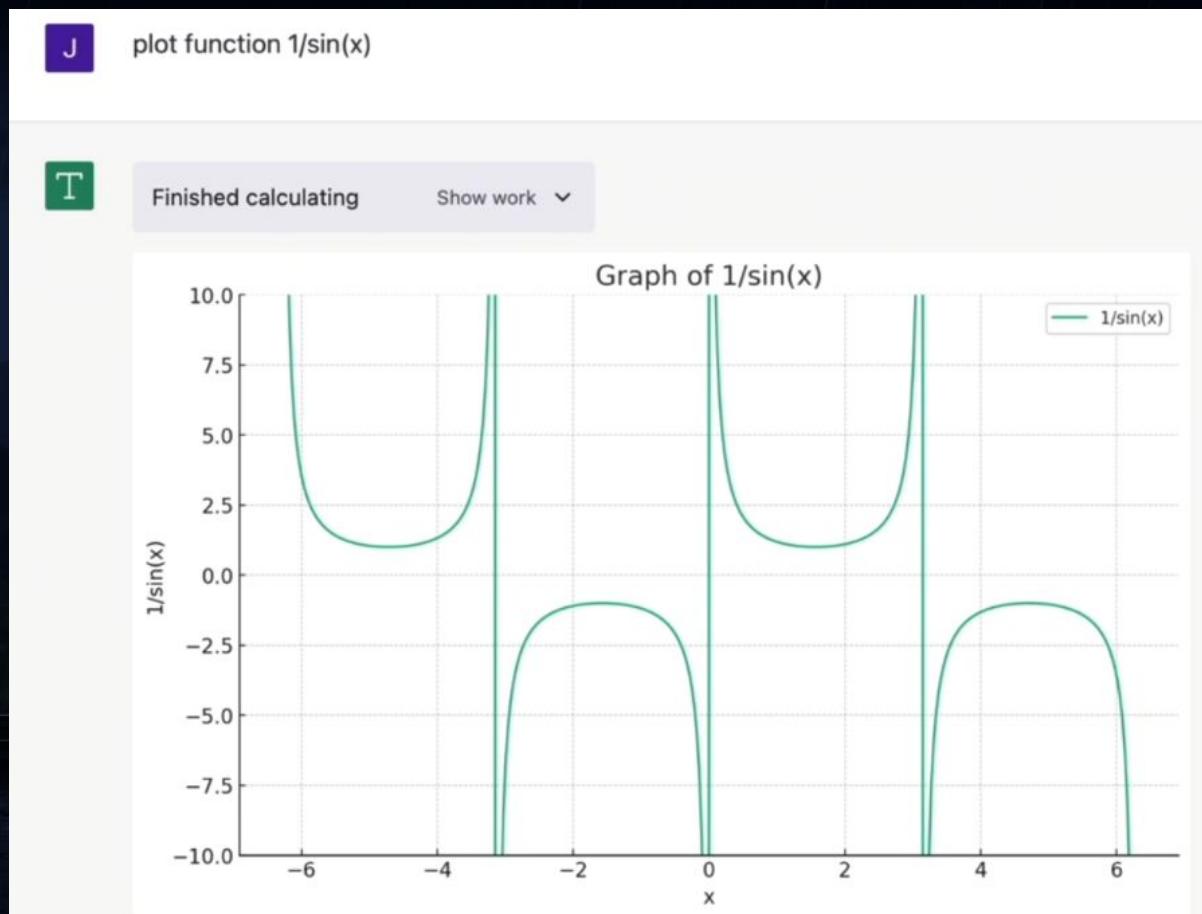
Made with Image Creator

Change the astronaut to a cat Change the sunflowers to roses Add a moon in the background

Type message

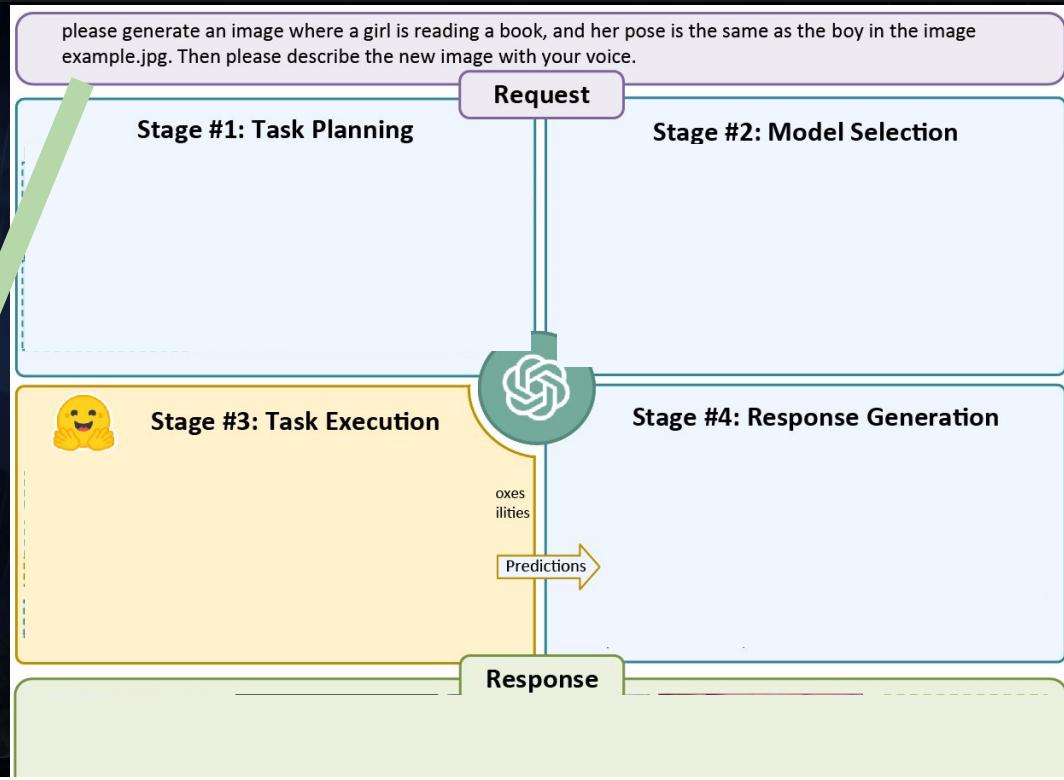
# ChatGPT Plugins

Web browsing, code interpreter,  
[Expedia](#), [FiscalNote](#), [Instacart](#),  
[KAYAK](#), [Klarna](#), [Milo](#), [OpenTable](#),  
[Shopify](#), [Slack](#), [Speak](#), [Wolfram](#),  
and [Zapier](#).



Here is the graph of the function  $1/\sin(x)$ . Please note that the graph has vertical

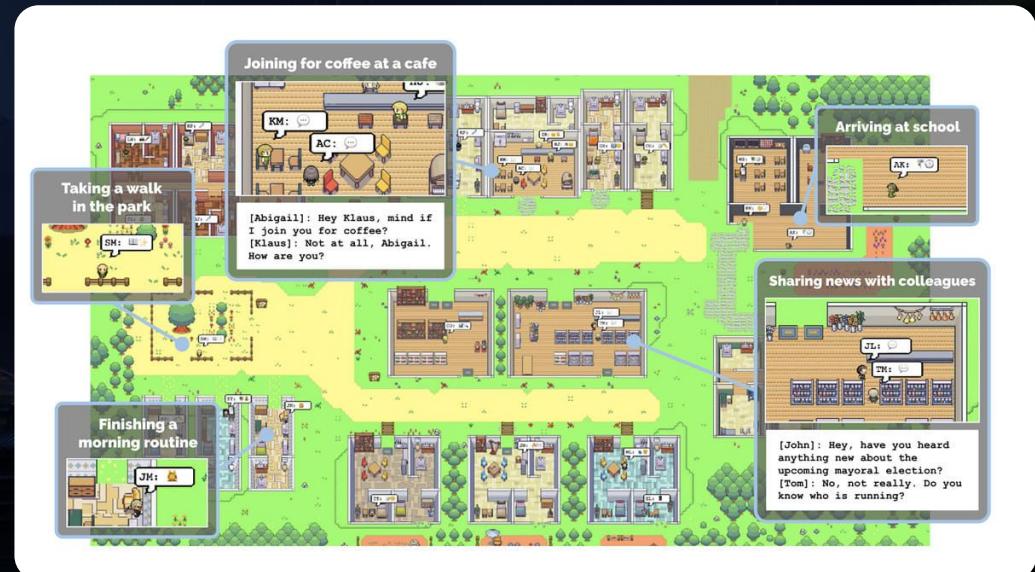
HuggingGPT (more powerful “plugin”), aka Microsoft Jarvis at Github



please generate an image where a girl is reading a book, and her pose is the same as the boy in the image example.jpg. Then please describe the new image with your voice.

# Interactive Simulacra of Human Behavior by Stanford

“fusing large language models with computational, interactive agents, this work introduces architectural and interaction patterns for enabling believable simulations of human behavior”



The background of the slide is a dark, moody photograph of a traditional Chinese architectural complex at night. It features several buildings with dark, curved roofs and intricate wooden structures. One prominent building in the foreground has a porch with red lanterns hanging from its eaves. The ground is wet, reflecting the light from the buildings, and it is raining heavily, with numerous thin, bright lines representing rain streaks against the dark background.

# Societal Impacts: Imminent Effects of ChatGPT-like AI

# Impact Assess to US Job Market ([OpenAI report](#))

“The projected [LLM] effects span all wage levels, with **higher-income jobs potentially facing greater exposure** to LLM capabilities and LLM-powered software...”

“...with access to an LLM, about 15% of all worker tasks in the US could be completed significantly faster at the same level of quality. When incorporating software and tooling built on top of LLMs, this share increases to between 47 and 56% of all tasks”

- My person take: *The higher your income is, statistically more impacted by LLM*
- My person take: Positive: “assistive AI to help humans”, Negative: “automation AI to replace humans”

## Impact Assess to US Job Market ([OpenAI report](#)) cont

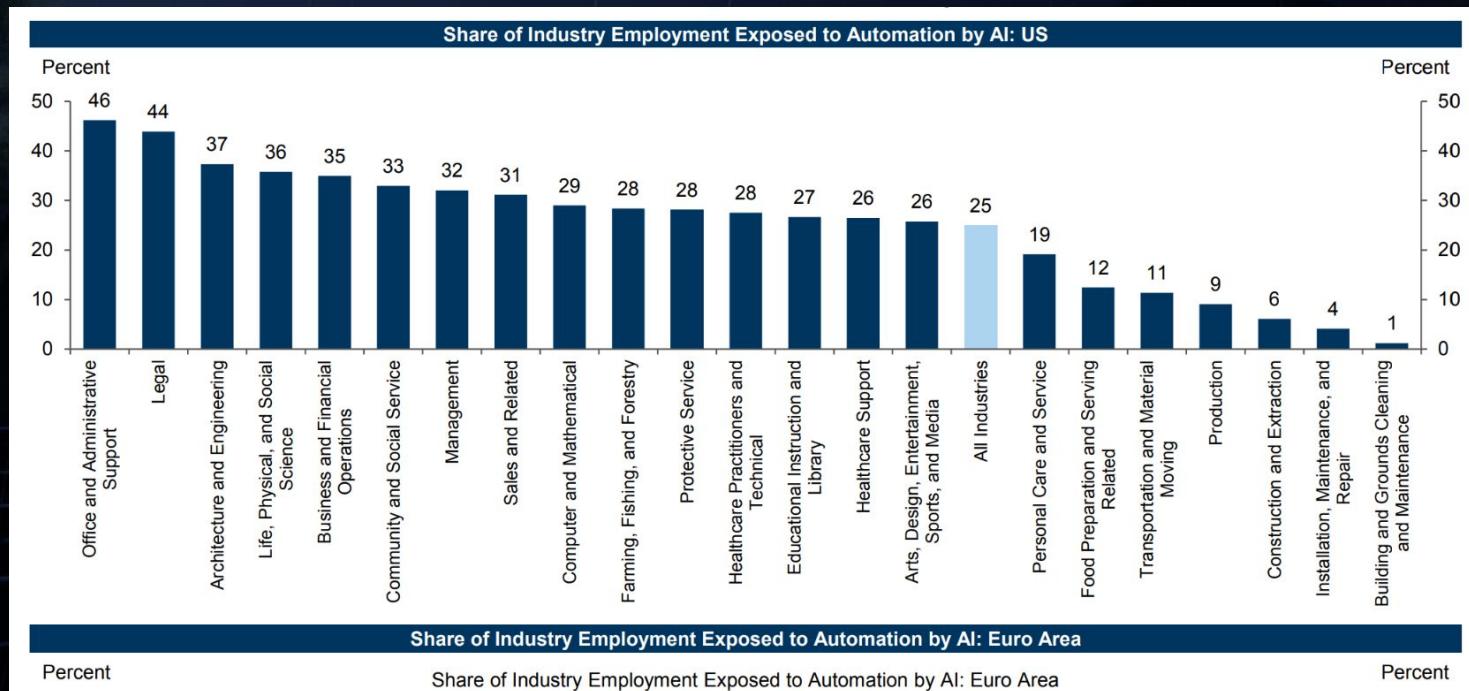
“Our findings reveal that around 80% of the U.S. workforce could have at least 10% of their work tasks **affected** by the introduction of LLMs, while approximately 19% of workers may see at least 50% of their tasks **impacted.**”

- *My personal take: Most of the white-collar jobs are in the 19% bucket*
- *My personal take: most of the blue-collar jobs are in the 80% bucket, but eventually the advanced robotics (maybe powered by LLM like GPT4) will gradually affect more over time*

# Impact to Job Market (Goldman Sachs report)

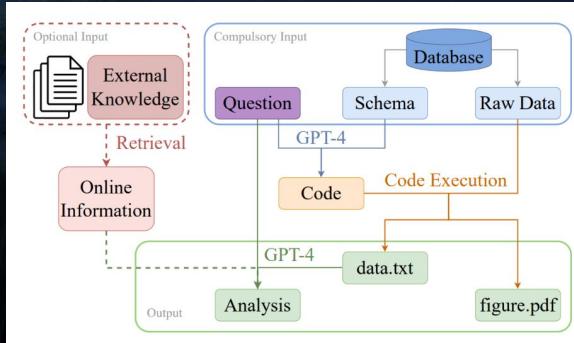
“One-Fourth of Current Work Tasks **Could Be Automated** by AI in the US and Europe”

- My person take: I believe wall street better than OpenAI here, because OpenAI has conflict of interest to report similar result, so OpenAI has good reasons to use more careful wording intentionally



# Is GPT-4 a Good Data Analyst?

Extra infra to improve GPT-4 for analyst jobs



Selected example of tasks the GPT-4 analyst system can do

Question	Please list the proportion number of each winning aircraft.												
<b>SQL Query</b>	<pre>SELECT a.Aircraft, COUNT(m.Winning_Aircraft) as wins FROM aircraft a JOIN match m ON a.Aircraft_ID = m.Winning_Aircraft GROUP BY a.Aircraft ORDER BY wins DESC</pre>												
<b>Figure</b>	<p>A pie chart titled 'Proportion of Wins by Aircraft' showing the distribution of wins across five aircraft models. The data is as follows:</p> <table border="1"><thead><tr><th>Aircraft</th><th>Proportion (%)</th></tr></thead><tbody><tr><td>Robinson R-22</td><td>28.6%</td></tr><tr><td>Mil Mi-26</td><td>28.6%</td></tr><tr><td>CH-53E Super Stallion</td><td>14.3%</td></tr><tr><td>CH-47D Chinook</td><td>14.3%</td></tr><tr><td>Bell 206B3 JetRanger</td><td>14.3%</td></tr></tbody></table>	Aircraft	Proportion (%)	Robinson R-22	28.6%	Mil Mi-26	28.6%	CH-53E Super Stallion	14.3%	CH-47D Chinook	14.3%	Bell 206B3 JetRanger	14.3%
Aircraft	Proportion (%)												
Robinson R-22	28.6%												
Mil Mi-26	28.6%												
CH-53E Super Stallion	14.3%												
CH-47D Chinook	14.3%												
Bell 206B3 JetRanger	14.3%												
<b>Analysis</b>	<ol style="list-style-type: none"><li>1. The Robinson R-22 and Mil Mi-26 are the <b>most successful</b> aircraft in the dataset, with each having 2 wins.</li><li>2. The CH-53E Super Stallion, CH-47D Chinook, and Bell 206B3 JetRanger each have 1 win, making them <b>less successful</b> compared to the Robinson R-22 and Mil Mi-26.</li></ol>												

Much cheaper!!!

Source	Median/Average Annual Salary (USD)	Cost per instance (USD)
levels.fyi	Entry Level DA	37,661
	Senior DA	90,421
Glassdoor	Junior DA	50,000
	Senior DA	86,300
Our Annotation	Junior DA	-
	Senior DA	11
GPT-4	-	<b>0.05</b>

# Debate on Pausing Giant AI or not

[Pause Giant AI Experiments: An Open Letter - Future of Life Institute](#)

## Pause Giant AI Experiments: An Open Letter

We call on all AI labs to immediately pause for at least 6 months the training of AI systems more powerful than GPT-4.

Signatures

15533

Add your signature

Feedback

AI systems with human-competitive intelligence can pose profound risks to society and humanity, as shown by extensive research<sup>[1]</sup> and acknowledged by top AI labs.<sup>[2]</sup> As stated in the widely-endorsed Asilomar AI Principles, *Advanced AI could represent a profound change in the history of life on Earth, and should be*

[Why the 6-month AI Pause is a Bad Idea](#)

DeepLearning.AI

## Why the 6-month AI Pause is a Bad Idea

Friday, April 7, 2023

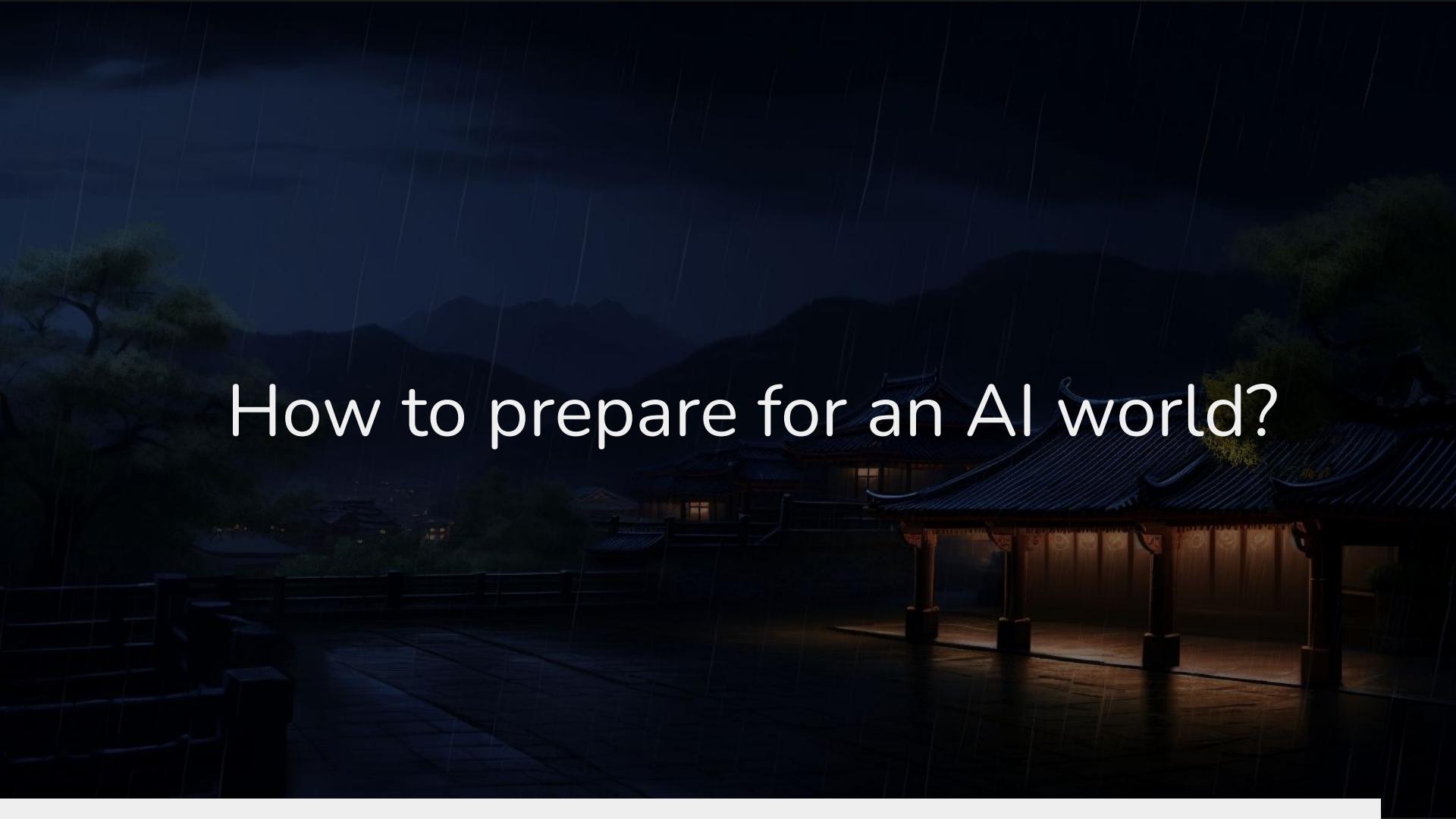
9:30 AM to 10:00 AM Pacific Time



**Yann LeCun**  
VP and Chief AI Scientist  
Meta



**Andrew Ng**  
Founder  
DeepLearning.AI

The background of the slide is a dark, atmospheric scene of a traditional Chinese building complex, possibly a temple or palace, in the rain at night. The buildings have dark, curved roofs and are illuminated from within by warm light. Rain is falling in heavy streaks across the dark sky. In the foreground, there's a paved area with some low walls and a railing on the left.

How to prepare for an AI world?

# My personal [very dry] thoughts

- Work with AI now, to know its capabilities and limits
- Be experts in your domain
  - Leverage AI to boost your performance
- Push AI infra boundary
  - AI researcher and engineers
  - Foundational theory like Math/Physics to improve infra/algorithm and more
  - Neural science or more to apply what learn from human brain to AI
- AI as a service, to solve real world problems
  - Inter-discipline research
  - Business landing using AI API
  - Embodied AI to have smart robots
- Be bold to solve the most difficult problems for the humanity

# Maybe switch jobs to the hottest “Prompt Engineer”?

Andrej Karpathy ✅  
@karpathy

The hottest new programming language is English

12:14 PM · Jan 24, 2023 · 2.2M Views

2,520 Retweets 383 Quotes 19.6K Likes 1,173 Bookmarks

Comment icon, Retweet icon, Like icon, Bookmark icon

Barsee 🐶 🐱  
@heyBarsee

Anthropic AI is looking for a Prompt Engineer.

Salary: \$250K - \$335k.

The job listing is starting, get into AI space now.

**ANTHROPIC**

**Prompt Engineer and Librarian** [APPLY FOR THIS JOB](#)

SAN FRANCISCO, CA / PRODUCT / FULL-TIME / HYBRID

Anthropic's mission is to create reliable, interpretable, and steerable AI systems. We want AI to be safe for our customers and for society as a whole.

Anthropic's AI technology is amongst the most capable and safe in the world. However, large language models are a new type of intelligence, and the art of instructing them in a way that delivers the best results is still in its infancy – it's a hybrid between programming, instructing, and teaching. You will figure out the best methods of prompting our AI to accomplish a wide range of tasks, then document these methods to build up a library of tools and a set of tutorials that allows others to learn prompt engineering or simply find prompts that would be ideal for them.

Given that the field of prompt-engineering is arguably less than 2 years old, this position is a bit hard to hire for! If you have existing projects that demonstrate prompt engineering on LLMs or image generation models, we'd love to see them. If you haven't done much in the way of prompt engineering yet, you can best demonstrate your prompt engineering skills by spending some time experimenting with Claude or GPT3 and

# Selected Highlights from Popular Articles

- Stephen Wolfram: [Will AIs Take All Our Jobs and End Human History—or Not?](#)
  - “highest leverage will come from figuring out **new possibilities** [...] as a result of **new capabilities**”
  - “let us concentrate on setting the “**strategy**” [...]—delegating the details [to AI]”
- Bill Gates: [The Age of AI has begun](#)
  - **“balance fears** about the **downsides of AI** [... and AI’s] **ability to improve people’s lives**”
  - “we will need to focus the world’s **best AIs on its biggest problems.**”
    - My take: Assume we may want to focus on AI application on weather/health/energy?
  - “the world needs to establish the rules of the road so that **any downsides of [AI] are far outweighed by its benefits**”
- Sam Altman: [Moore's Law for Everything](#)
  - “Imagine a world where, for decades, everything—housing, education, food, clothing, etc.—became half as expensive every two years. [...] **We will discover new jobs** [...], we will have incredible freedom to be creative about what they are.”
    - My take: really?
  - “As long as the country keeps doing better, every citizen would get more money from the Fund every year [...]. Every citizen would therefore increasingly partake of the freedoms, powers, autonomies, and opportunities [...]”
    - My take: seriously?

fun-ai-talk @ 谷雨书苑

谷雨书苑+理想折射 ChatGPT

# A Primer on ChatGPT

hululu.zhu@gmail.com

May 2023