



AI EXPERT NETWORK

AI Large Language Models Models & Healthcare

Jon Chun

@jonchun2000

AI for Humanity

Kenyon College

21 Apr 2023

About

Description:

The rapid rise of Large Language Models (LLMs) such as ChatGPT and GPT-4 in recent months has sparked an unprecedented surge in new research, models, frameworks, and industry applications. This talk will provide a concise overview of the historical development of LLMs and examine current models, frameworks, and applications. We will conclude with a forward-looking discussion of future applications, including LLM tool integration (e.g., LangChain), multimodal processing (e.g., JARVIS, HuggingGPT), automation frameworks (e.g., AutoGPT), and potential integration with Distributed Autonomous Organizations (DAOs)/Blockchain. Special emphasis will be placed on the unique challenges and opportunities associated with incorporating LLMs and related technologies within the healthcare and HealthTech startups.

Bio:









Jon Chun is a co-creator of the world's first AI for the Humanities curriculum at Kenyon College. He has mentored hundreds of original Data Science, Machine Learning, and AI research projects, which have been downloaded over 22,000 times by more than 1,700 institutions worldwide, including Stanford, Berkeley, CMU, MIT, Princeton, and Oxford. Before entering academia, Jon was a Silicon Valley entrepreneur, co-founding and leading the world's largest privacy and anonymity web service, as well as developing the first web-based VPN appliance. Apart from startups, he has served as a Director of Development at Symantec, the world's largest security company, and as CTO for a major third-party disability insurance and absence management firm in Silicon Valley. Jon has published patents, research, and presented at national conferences on diverse specialties including network security, privacy and anonymity, medical informatics, gene therapy, Multimodal Affective AI, Narratology, FATE/XAI, and GPT-2, GPT-3, and GPT-4 (upcoming).

Overview

A Snapshot in Time
(21 April 2023)

HOW LONG IT TOOK TOP APPS TO HIT 100M MONTHLY USERS

ChatGPT is estimated to have hit 100M users in January, 2 months after it's launch. Here's how long it took other top apps to reach that:

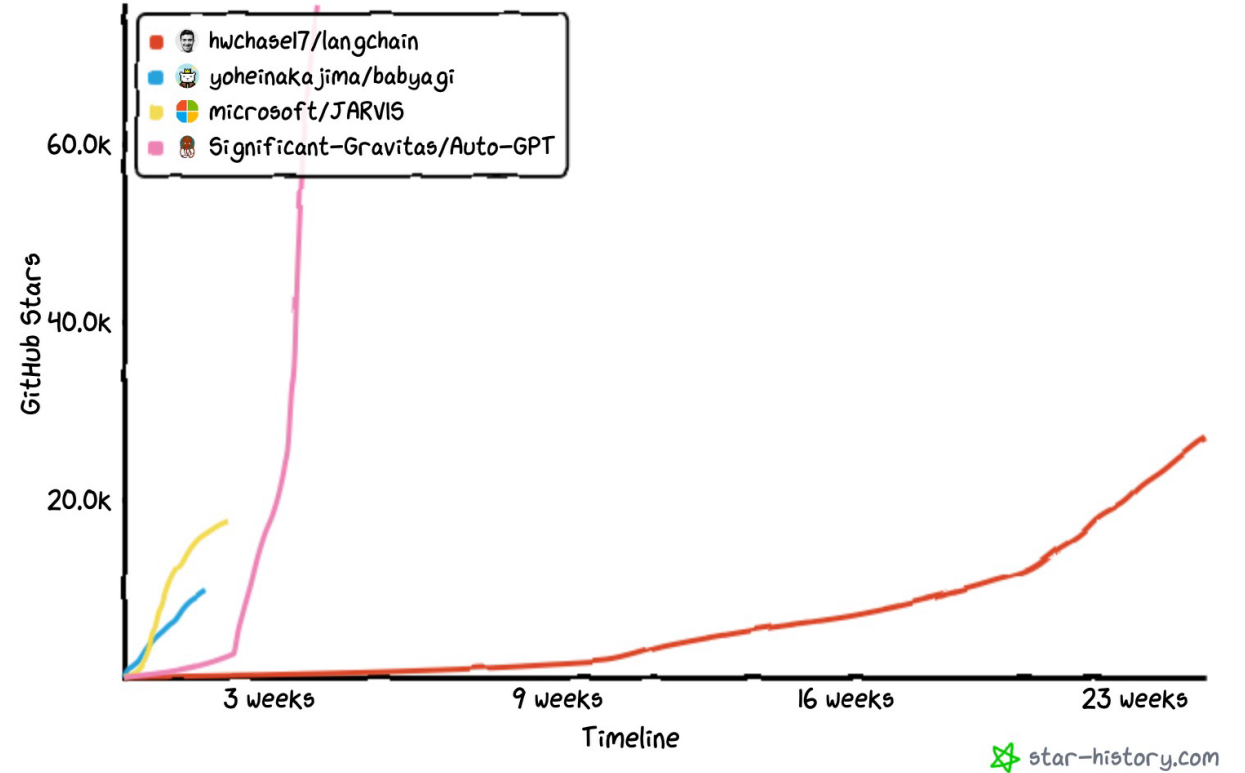
APP	MONTHS TO REACH 100M GLOBAL MAUS
 CHATGPT	2
 TIKTOK	9
 INSTAGRAM	30
 PINTEREST	41
 SPOTIFY	55
 TELEGRAM	61
 UBER	70
 GOOGLE TRANSLATE	78

SOURCE: UBS



98k (4/20)
3 weeks

Star History



Large Language Models and Healthcare

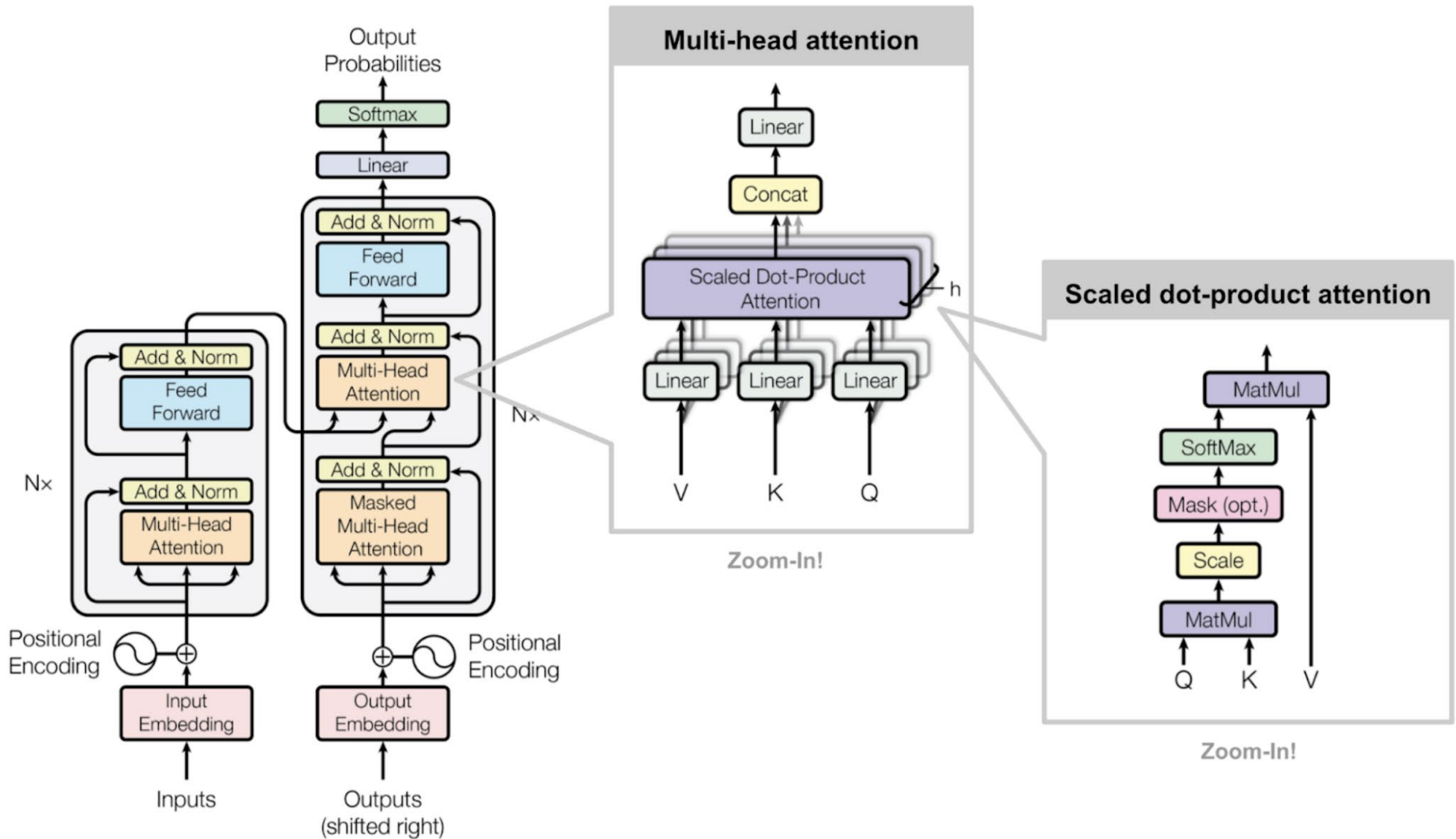
- Evolution of LLMs
- Survey of LLMs
- Extending LLMs
- Medical Research
- MedTech Startups

Evolution of Large Language Models (LLMs)

From DNNs to Transformers to LLM

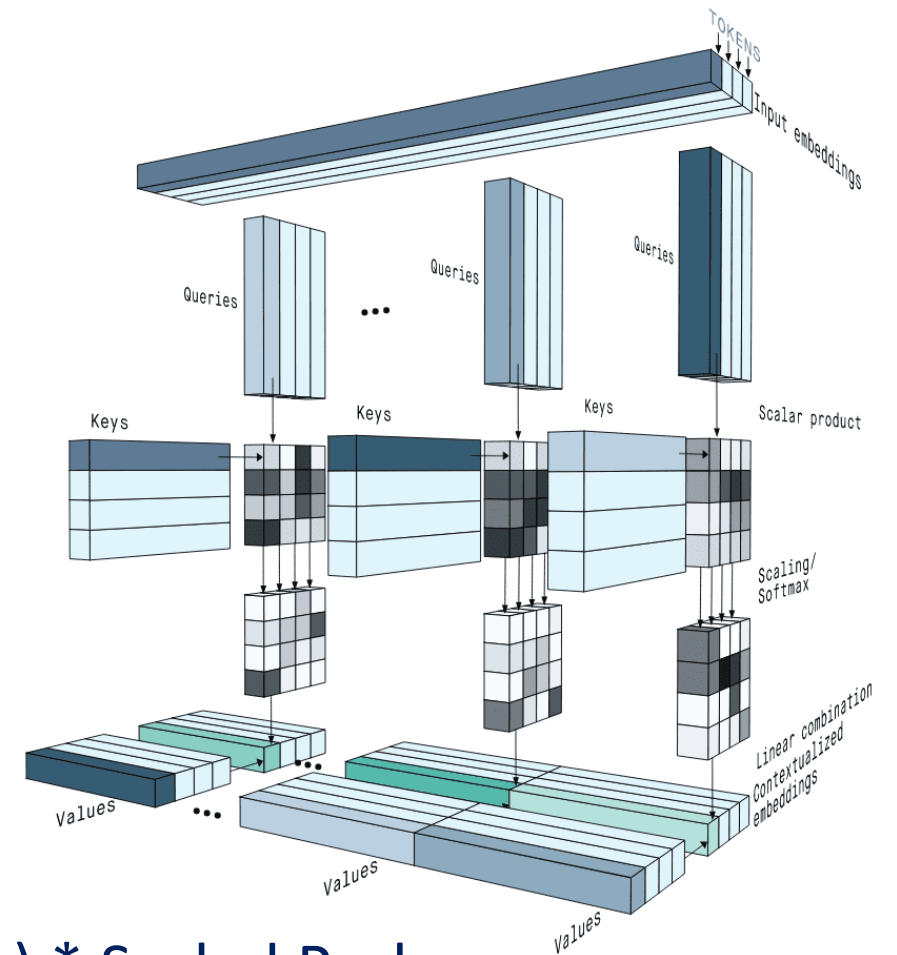
Milestones

- 2012 Oct 13: ImageNet Results
- 2017 June 12: Attention is All You Need *
- 2019 Aug 20: 774M GPT2
- 2022 Nov 30: ChatGPT *
- 2022 Mar 14: GPT4



Transformer Architecture

- Tokenization
- Embedding
- Positional Embeddings
- $QK = \text{Queries}(\text{loc}) * \text{Keys}(\text{preposition})$
- $\text{Scaled Prob} = \text{SoftMax}(QK) / \text{Normalizer}$
- $\text{Contextualized Embeddings} = \text{Values}(\text{e.g. place}) * \text{Scaled Prob}$
- Multi-head Attention
- Multi-Layers (BERT 16)



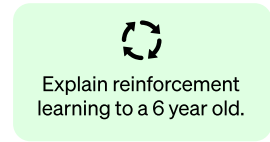
ChatGPT

- Transformer Architecture
- InstructGPT
- RLHF
- Prompt Engineering

Step 1

Collect demonstration data and train a supervised policy.

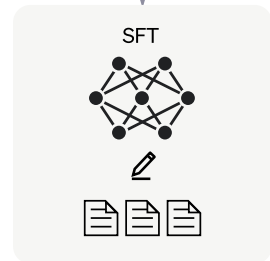
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



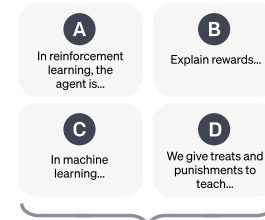
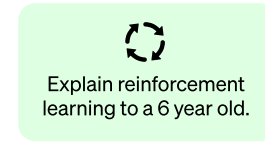
This data is used to fine-tune GPT-3.5 with supervised learning.



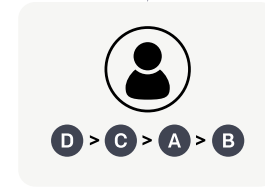
Step 2

Collect comparison data and train a reward model.

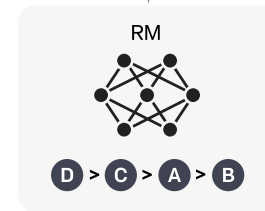
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



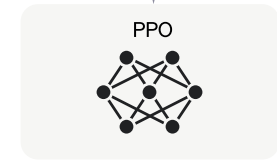
Step 3

Optimize a policy against the reward model using the PPO reinforcement learning algorithm.

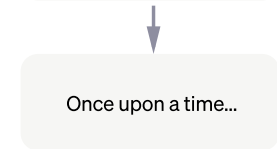
A new prompt is sampled from the dataset.



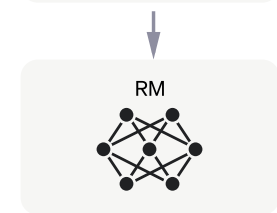
The PPO model is initialized from the supervised policy.



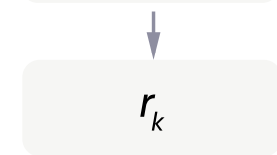
The policy generates an output.



The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.



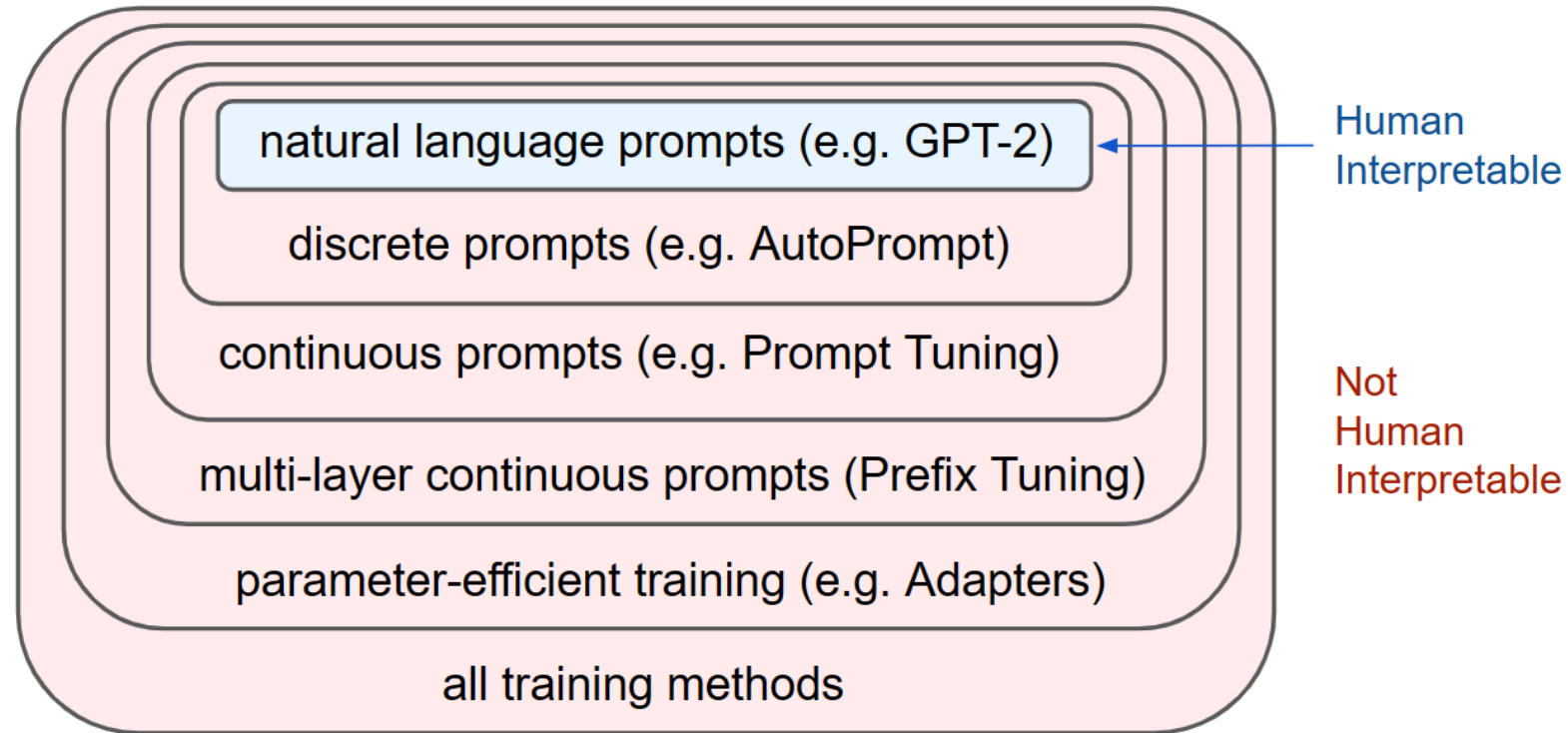
Prompt Engineering

- Zero, One & n-Shot Learning
- Context Window
- Ensemble
- Chain of Thought
- Reflection / Self-Critique
- Jailbreak

A Taxonomy of Prompting Methods

By Graham Neubig (10/15/2022)

See [CMU ANLP Prompting Lecture](#), [A Unified View of Parameter-Efficient Transfer Learning](#)



GPT-2: <https://openai.com/blog/better-language-models/>

AutoPrompt: <https://arxiv.org/abs/2010.15980>

Prefix Tuning: <https://arxiv.org/abs/2101.00190>

Prompt Tuning: <https://arxiv.org/abs/2104.08691>

Adapters: <https://arxiv.org/abs/2010.15980>

Survey of LLMs

As of April 2023

(Don't Blink)

Large Language Model Metaverse

- SOTA: GPT4 (PaLM, LLaMA)
- Limitations
- Characteristics
 - Sequential
 - Language
 - Auto-Regressive / Self-Supervised Learning
 - General > Fine-Tuning
 - Emergent Functionality
- Variations
 - Commercial: GPT4
 - Open-Source
 - Distilled
 - Uncensored

GPT4

- “OpenAI” Technical Report
- Scale, Dataset and Training
- “Lot’s of little things”
- Context Window
- Multimodal
- Performance
- Emergence

GPT-4 Technical Report

OpenAI*

Abstract

We report the development of GPT-4, a large-scale, multimodal model which can accept image and text inputs and produce text outputs. While less capable than humans in many real-world scenarios, GPT-4 exhibits human-level performance on various professional and academic benchmarks, including passing a simulated bar exam with a score around the top 10% of test takers. GPT-4 is a Transformer-based model pre-trained to predict the next token in a document. The post-training alignment process results in improved performance on measures of factuality and adherence to desired behavior. A core component of this project was developing infrastructure and optimization methods that behave predictably across a wide

27 Mar 2023

OpenAI codebase next word prediction

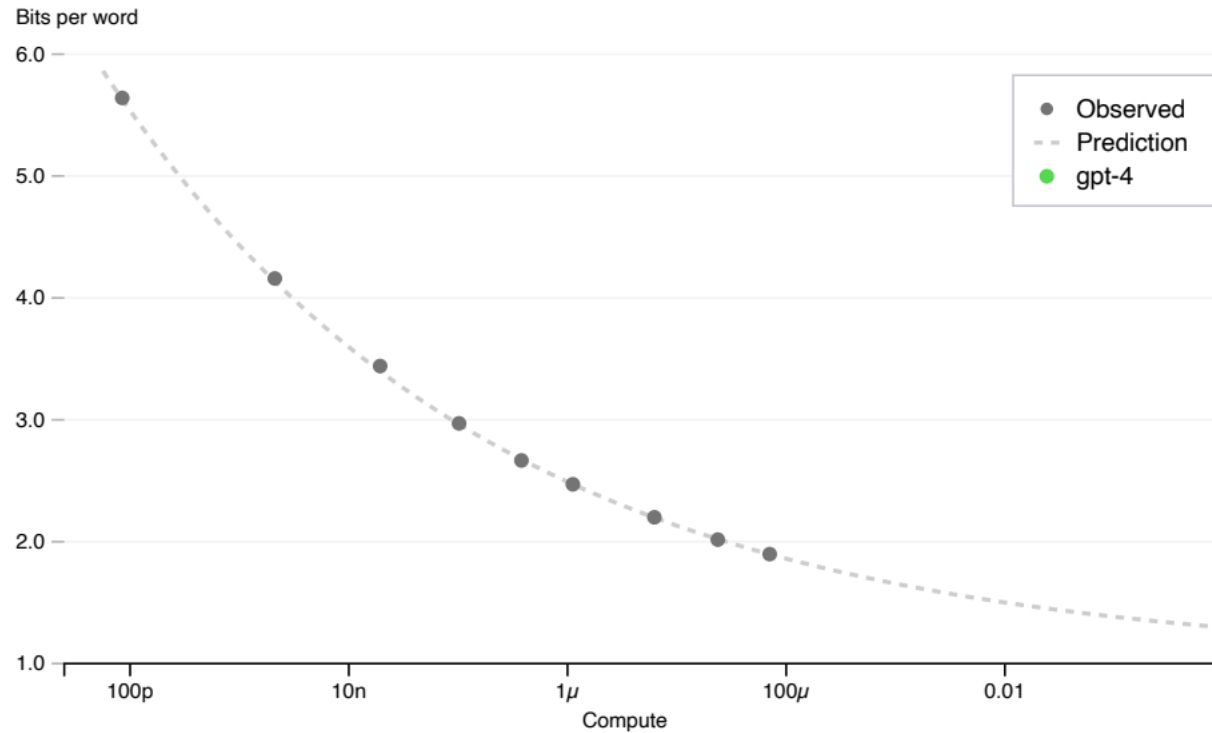


Figure 1. Performance of GPT-4 and smaller models. The metric is final loss on a dataset derived from our internal codebase. This is a convenient, large dataset of code tokens which is not contained in the training set. We chose to look at loss because it tends to be less noisy than other measures across different amounts of training compute. A power law fit to the smaller models (excluding GPT-4) is shown as the dotted line; this fit accurately predicts GPT-4’s final loss. The x-axis is training compute normalized so that GPT-4 is 1.

Exam results (ordered by GPT-3.5 performance)

Estimated percentile lower bound (among test takers)

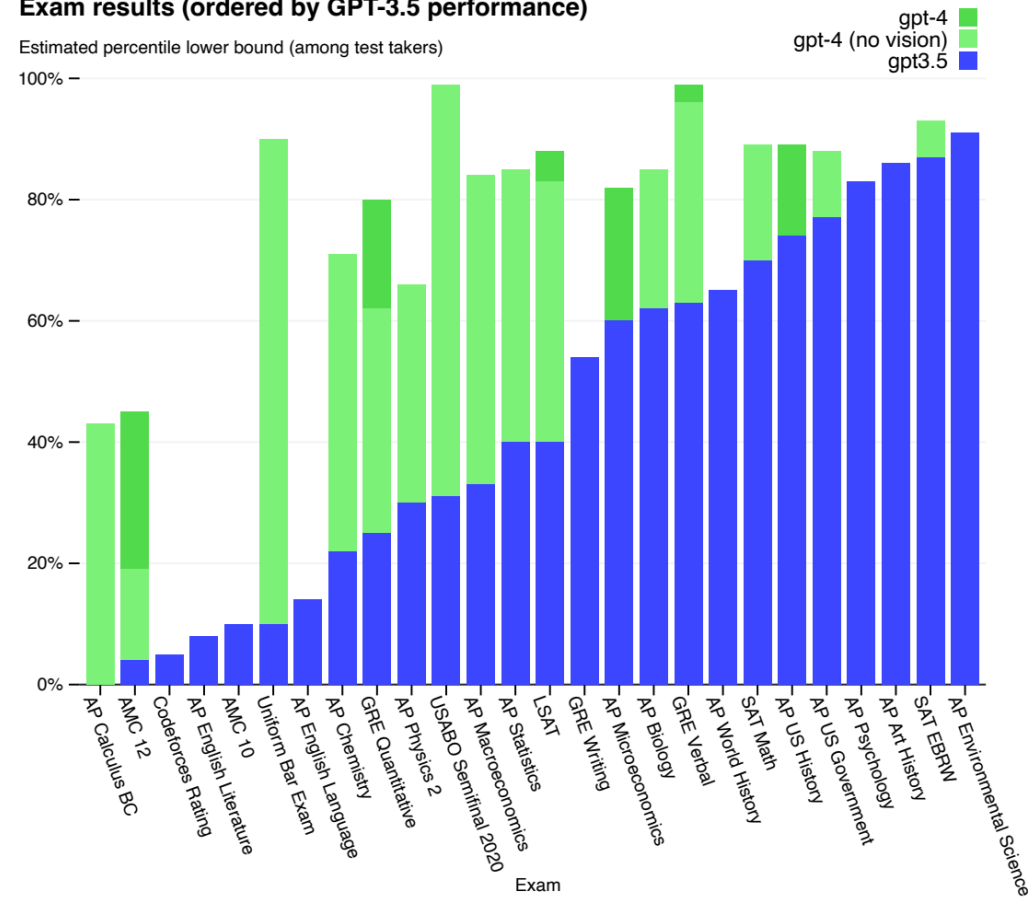


Figure 4. GPT performance on academic and professional exams. In each case, we simulate the conditions and scoring of the real exam. Exams are ordered from low to high based on GPT-3.5 performance. GPT-4 outperforms GPT-3.5 on most exams tested. To be conservative we report the lower end of the range of percentiles, but this creates some artifacts on the AP exams which have very wide scoring bins. For example although GPT-4 attains the highest possible score on AP Biology (5/5), this is only shown in the plot as 85th percentile because 15 percent of test-takers achieve that score.

	GPT-4 Evaluated few-shot	GPT-3.5 Evaluated few-shot	LM SOTA Best external LM evaluated few-shot	SOTA Best external model (incl. benchmark-specific tuning)
MMLU [49] Multiple-choice questions in 57 subjects (professional & academic)	86.4% 5-shot	70.0% 5-shot	70.7% 5-shot U-PaLM [50]	75.2% 5-shot Flan-PaLM [51]
HellaSwag [52] Commonsense reasoning around everyday events	95.3% 10-shot	85.5% 10-shot	84.2% LLaMA (validation set) [28]	85.6 ALUM [53]
AI2 Reasoning Challenge (ARC) [54] Grade-school multiple choice science questions. Challenge-set.	96.3% 25-shot	85.2% 25-shot	85.2% 8-shot PaLM [55]	86.5% ST-MOE [18]
WinoGrande [56] Commonsense reasoning around pronoun resolution	87.5% 5-shot	81.6% 5-shot	85.1% 5-shot PaLM [3]	85.1% 5-shot PaLM [3]
HumanEval [43] Python coding tasks	67.0% 0-shot	48.1% 0-shot	26.2% 0-shot PaLM [3]	65.8% CodeT + GPT-3.5 [57]
DROP [58] (F1 score) Reading comprehension & arithmetic.	80.9 3-shot	64.1 3-shot	70.8 1-shot PaLM [3]	88.4 QDGAT [59]
GSM-8K [60] Grade-school mathematics questions	92.0%* 5-shot chain-of-thought	57.1% 5-shot	58.8% 8-shot Minerva [61]	87.3% Chinchilla + SFT+ORM-RL, ORM reranking [62]

Table 2. Performance of GPT-4 on academic benchmarks. We compare GPT-4 alongside the best SOTA (with benchmark-specific training) and the best SOTA for an LM evaluated few-shot. GPT-4 outperforms existing LMs on all benchmarks, and beats SOTA with benchmark-specific training on all datasets except DROP. For each task we report GPT-4’s performance along with the few-shot method used to evaluate. For GSM-8K, we included part of the training set in the GPT-4 pre-training mix (see Appendix E), and we use chain-of-thought prompting [11] when evaluating. For multiple-choice questions, we present all answers (ABCD) to the model and ask it to choose the letter of the answer, similarly to how a human would solve such a problem.

AI Gold Rush / AI Arms Race

- Scale vs Fine-Tuning
- Dataset
- Training
- MLOPs
- Multimodal
- Embodied
- Considerations
 - Tools & Automation
 - Cognitive Disintermediation
 - Natural Monopoly
 - 5th Industrial Revolution

Risks

- Probabilistic
- Hallucination/False Confidence
- Copyright
- Bias, Offensive & Dangerous (e.g. Euthanasia in JP vs NL)
- FATE/XAI
- Privacy/Security
- Low Resource
- Causality, Models and New Knowledge
- (Supra-)National Regulations, Laws and Liabilities

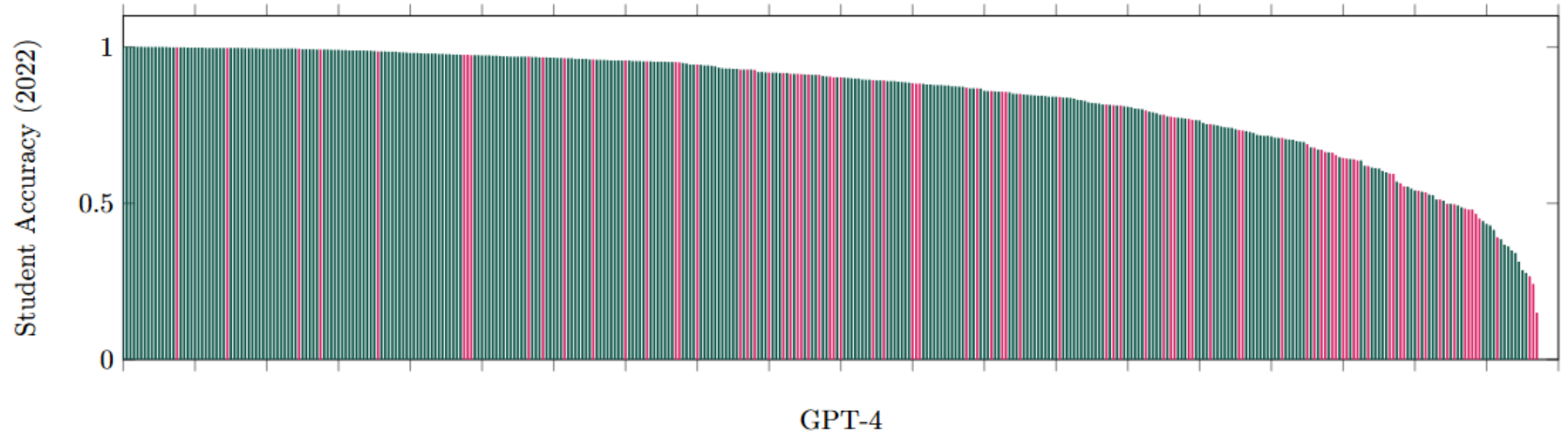


Figure 8: Student (test taker) accuracy vs. GPT-4 results. All problems from 2022 are sorted by the student accuracy, and the bar is green when GPT-4 predicts the correct choice(s) and red otherwise. We see correlation between the student accuracy and the likelihood of the correct prediction. We see similar patterns for other models (Appendix §C).

Extending LLMs

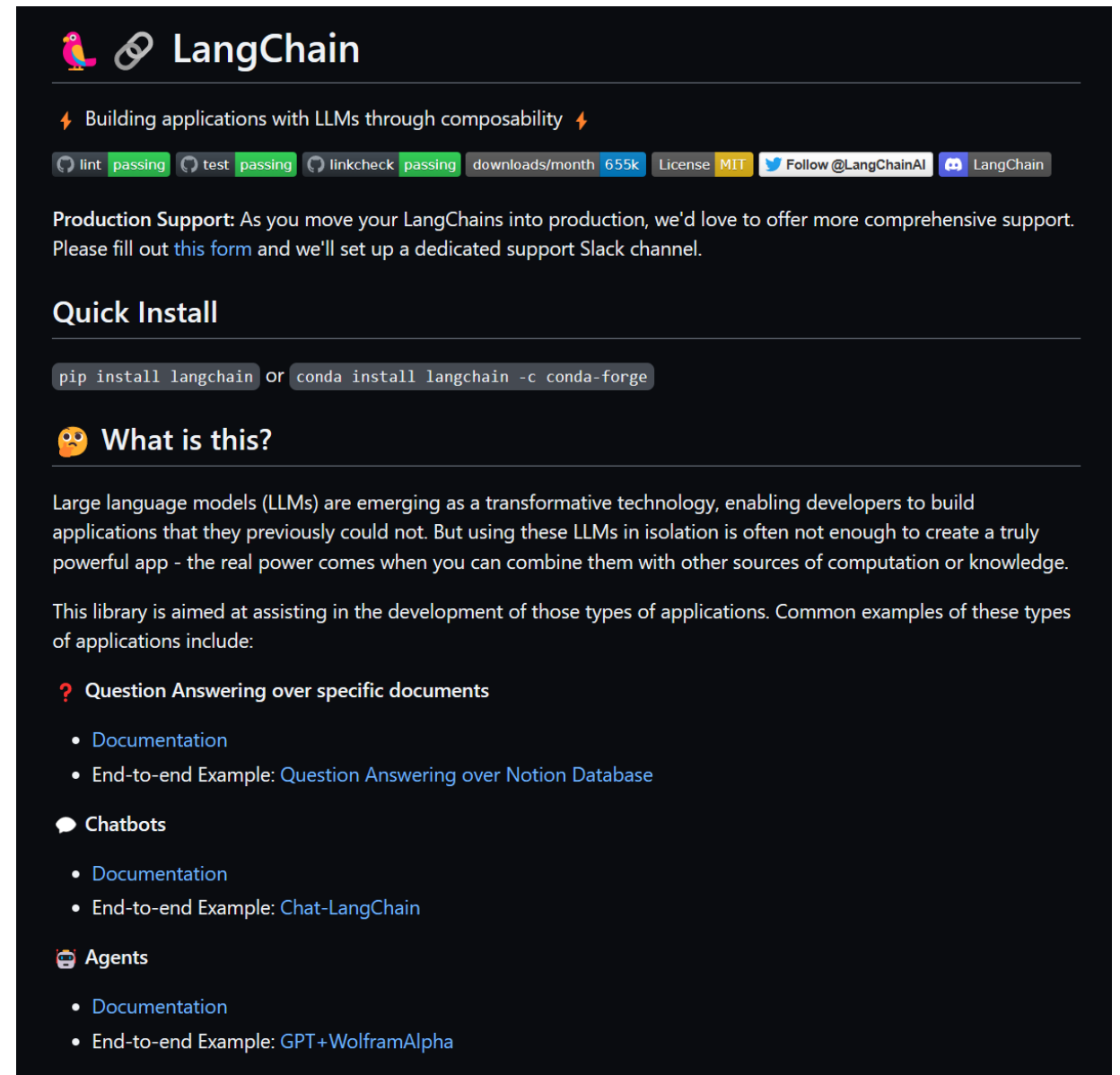
Frameworks, Tools and Agents


Extending Large Language Models

- Frameworks
 - Tools: LangChain
 - Models: JARVIS, HuggingGPT
- Automation: AutoGPT
- Embodiment: PaLM-E
- Integration: Tools, Agents and World

LangChain

- Prompt Template
- Memory
- Tools
- Agent
- Agent Executor (LLM)



 **LangChain**

⚡ Building applications with LLMs through composability ⚡

`lint` `passing` `test` `passing` `linkcheck` `passing` `downloads/month` `655k` `License` `MIT` `Follow @LangChainAI` `LangChain`

Production Support: As you move your LangChains into production, we'd love to offer more comprehensive support. Please fill out [this form](#) and we'll set up a dedicated support Slack channel.

Quick Install

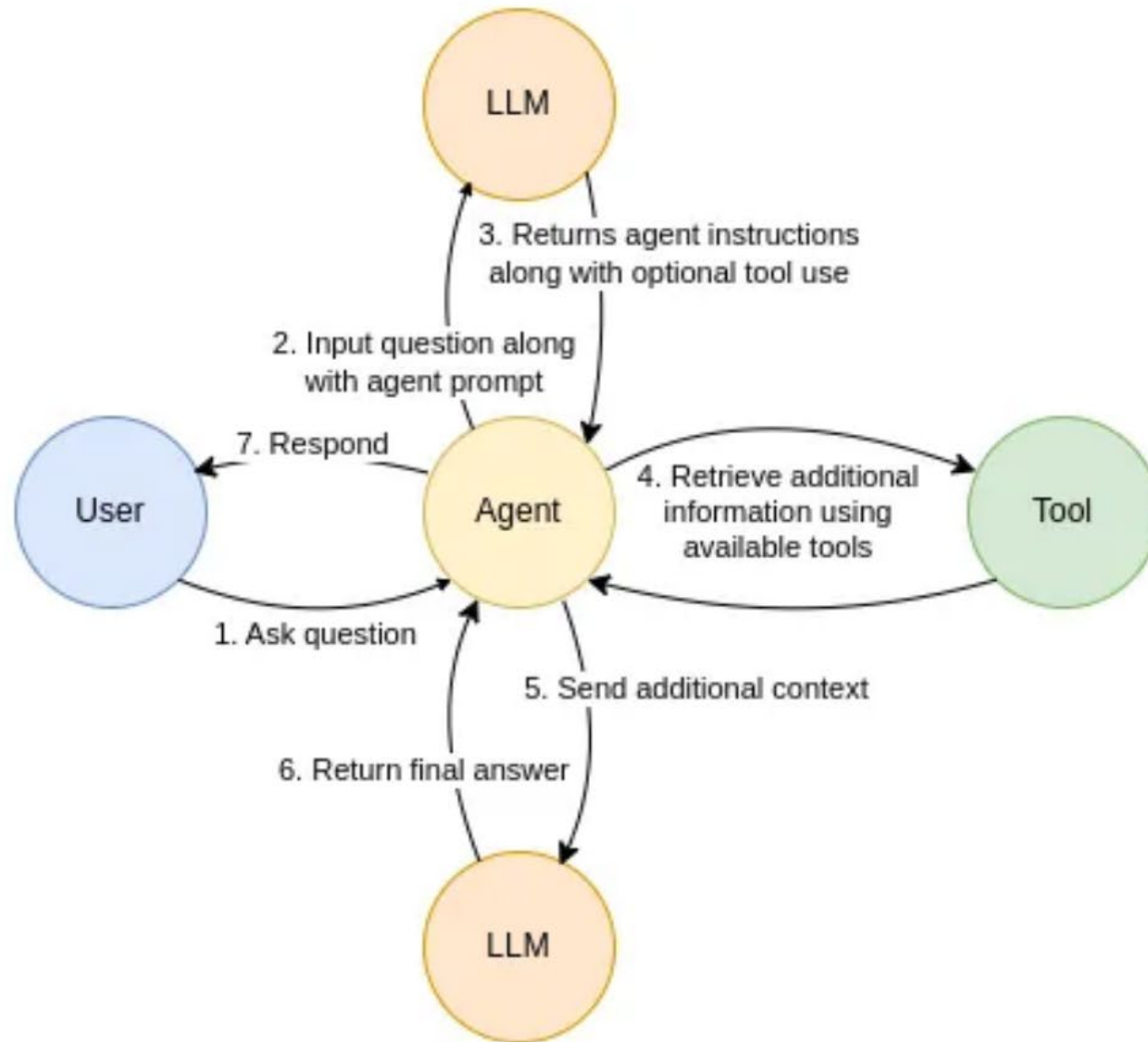
```
pip install langchain or conda install langchain -c conda-forge
```

🤖 What is this?

Large language models (LLMs) are emerging as a transformative technology, enabling developers to build applications that they previously could not. But using these LLMs in isolation is often not enough to create a truly powerful app - the real power comes when you can combine them with other sources of computation or knowledge.

This library is aimed at assisting in the development of those types of applications. Common examples of these types of applications include:

- ? **Question Answering over specific documents**
 - [Documentation](#)
 - End-to-end Example: [Question Answering over Notion Database](#)
- 🗨️ **Chatbots**
 - [Documentation](#)
 - End-to-end Example: [Chat-LangChain](#)
- 🤖 **Agents**
 - [Documentation](#)
 - End-to-end Example: [GPT+WolframAlpha](#)



Agent steps:

1. User asks question
2. Question is send to an LLM along with the Agent prompt
3. LLM responds with further instructions either to immediately answer the user or use tools for additional information
4. Retrieve additional information
- 5 & 6. LLM constructs a final answer based on additional context



Table of contents

- LangChain
 - Installation
 - LLMs**
 - Prompt Templates
 - Chains
 - Agents and Tools
 - Memory
 - Document Loaders
 - Indexes
 - End-to-end example
- + Section

+ Code + Text

Reconnect ▾



Copyright © 2023 Patrick Loeber

- <https://www.youtube.com/watch?v=LbT1yp6quS8&t=6s>

LangChain

LangChain is a framework for developing applications powered by language models.

- GitHub: <https://github.com/hwchase17/langchain>
- Docs: <https://python.langchain.com/en/latest/index.html>

Overview:

- Installation
- LLMs
- Prompt Templates
- Chains
- Agents and Tools
- Memory
- Document Loaders
- Indexes

Auto-GPT: An Autonomous GPT-4 Experiment

unit tests **passing** AutoGPT 21248 members Stars 101k Follow @siggravitas

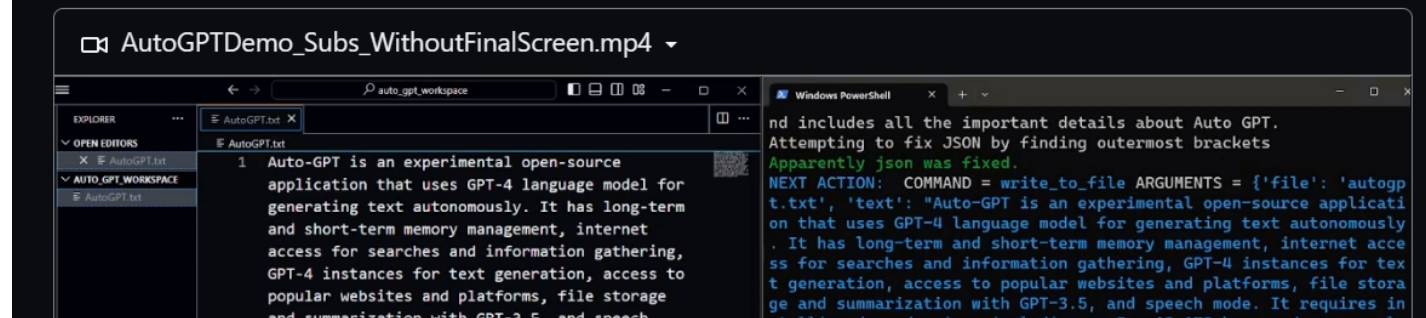
💡 Get help - Q&A or Discord

● ● ● Urgent: USE **stable** not **master** ● ● ●

Download the latest **stable** release from here: <https://github.com/Significant-Gravitas/Auto-GPT/releases/latest>. The **master** branch may often be in a broken state.

Auto-GPT is an experimental open-source application showcasing the capabilities of the GPT-4 language model. This program, driven by GPT-4, chains together LLM "thoughts", to autonomously achieve whatever goal you set. As one of the first examples of GPT-4 running fully autonomously, Auto-GPT pushes the boundaries of what is possible with AI.

Demo April 16th 2023

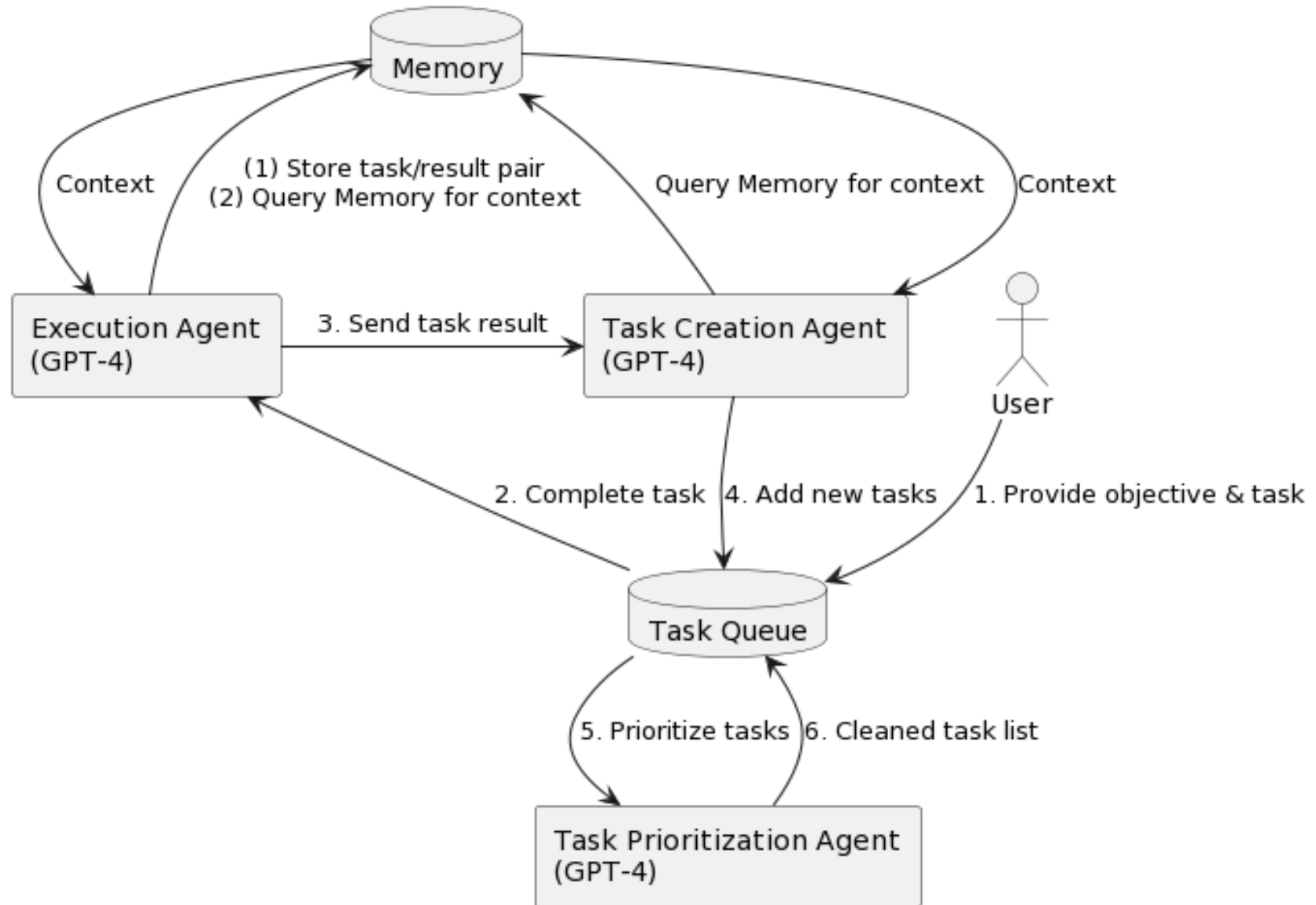


The screenshot shows a code editor window titled 'AutoGPTDemo_Subs_WithoutFinalScreen.mp4' with a file explorer on the left and a code editor on the right. The code editor displays the following text:

```
1 Auto-GPT is an experimental open-source application that uses GPT-4 language model for generating text autonomously. It has long-term and short-term memory management, internet access for searches and information gathering, GPT-4 instances for text generation, access to popular websites and platforms, file storage and summarization with GPT-3.5, and speech
```

Below the code editor is a Windows PowerShell terminal window showing the following output:

```
nd includes all the important details about Auto GPT. Attempting to fix JSON by finding outermost brackets Apparently json was fixed. NEXT ACTION: COMMAND = write_to_file ARGUMENTS = {'file': 'autogp t.txt', 'text': "Auto-GPT is an experimental open-source applicati on that uses GPT-4 language model for generating text autonomously . It has long-term and short-term memory management, internet acce ss for searches and information gathering, GPT-4 instances for tex t generation, access to popular websites and platforms, file stora ge and summarization with GPT-3.5, and speech mode. It requires in
```



AgentGPT Beta

Assemble, configure, and deploy autonomous AI Agents in your browser.

AgentGPT

Image Copy PDF

> Create an agent by adding a name / goal, and hitting deploy!

You can provide your own OpenAI API key in the settings tab for increased limits!

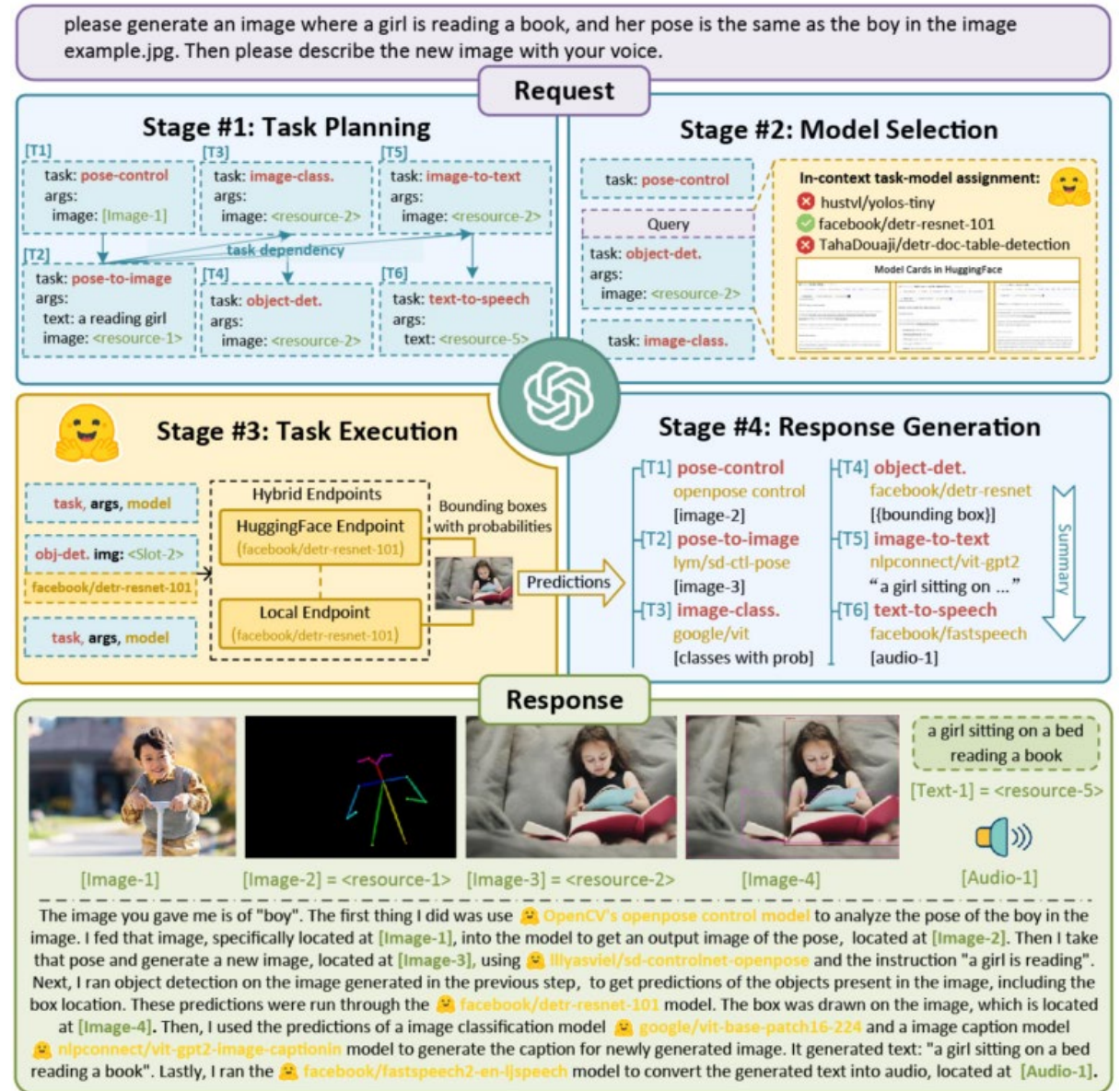
Help support the advancement of AgentGPT. Please consider sponsoring the project on GitHub. [Support now](#)

Name: AgentGPT

Goal: Make the world a better place.

JARVIS / HuggingGPT

- Frameworks
 - Tools
 - Models
- Automation
 - Planning
 - Model Selection
 - Execution
 - Response
- Embodiment
- 360 Integration



HuggingGPT



A system to connect LLMs with ML community. See our [Project](#) and [Paper](#).

[Duplicate Space](#)



Duplicate the Space and run securely with your OpenAI API Key and Hugging Face Token

Note: Only a few models are deployed in the local inference endpoint due to hardware limitations. In addition, online HuggingFace inference endpoints may sometimes not be available. Thus the capability of HuggingGPT is limited.

.....

.....

Chatbot

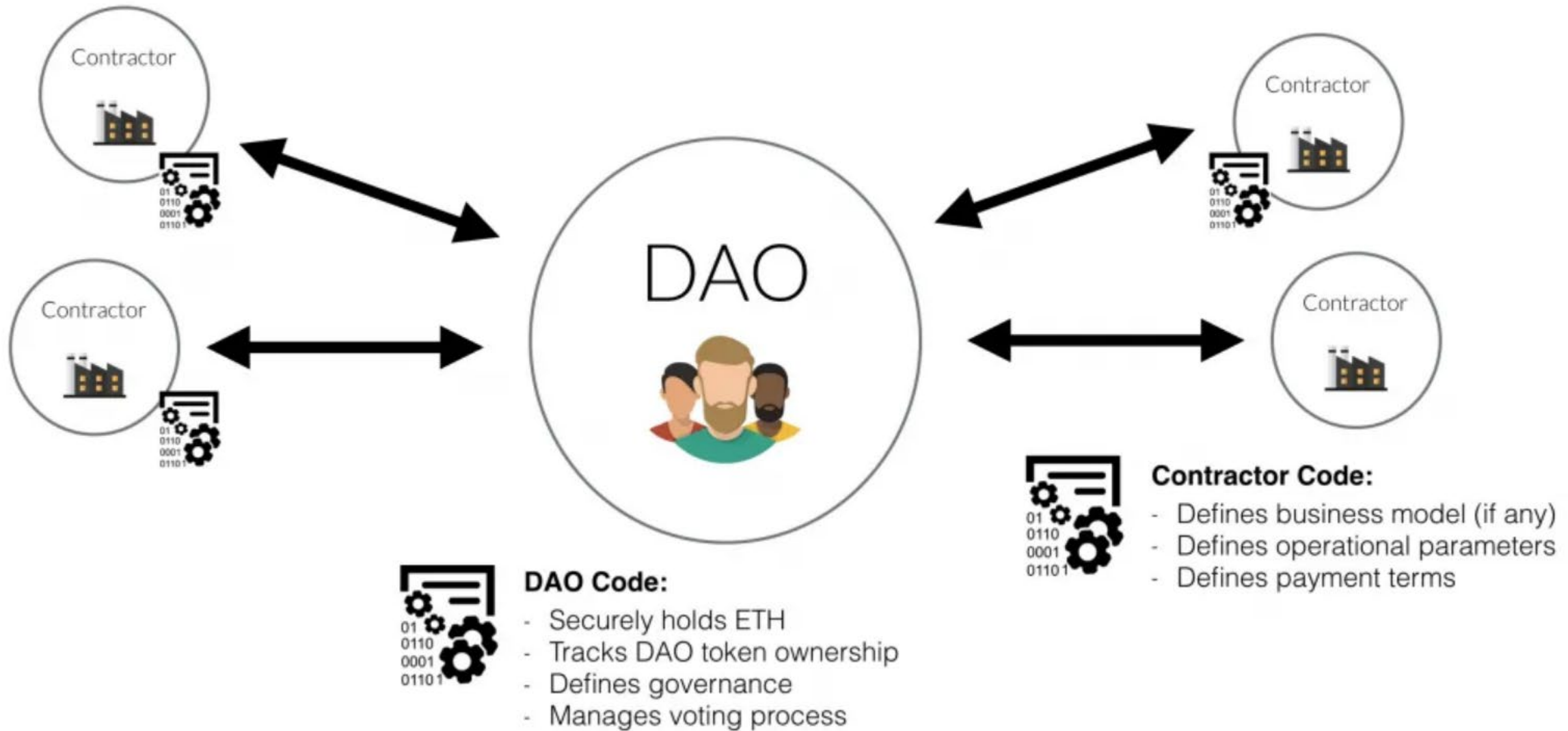
```

{
  0: {
    task: {
      task: "image-to-text",
      id: 0,
      dep: [
        0: -1
      ],
      args: {
        image: "public//examples/a.jpg"
      }
    },
    inference result: {
      generated text: "a cat sitting on a window sill looking out "
    },
    choose model result: {
      id: "ydshieh/vit-gpt2-coco-en",
      reason: "Only one model available."
    }
  },
  1: {
    task: {
      task: "object-detection",
      id: 1,
      den: [

```

Based on the inference results, there are totally 3 pictures with 0, 1 and 2 zebras respectively. For the first picture, the inference result of object detection is 'cat' and 'potted plant', and the inference result of visual question answering is '0'. For the second picture, the inference result of object detection is 'zebra', and the inference result of visual question answering is '1'. For the

Enter text and press enter. The url must contain the media type. e.g, <https://example.com/example.jpg>



Medical Research with LLMs

As of April 2023

(don't blink)

ChatGPT in Healthcare: A Taxonomy and Systematic Review

Jianning Li, Amin Dada, Jens Kleesiek, Jan Egger*

Institute of Artificial Intelligence in Medicine, University Hospital
Essen (AöR), Girardetstraße, 45131 Essen, Germany.

*Corresponding author: [jan.egger \(at\) uk-essen.de](mailto:jan.egger@uk-essen.de) (J.E.)

March 2023

Abstract

The recent release of ChatGPT, a chat bot research project/product of natural language processing (NLP) by OpenAI, stirs up a sensation

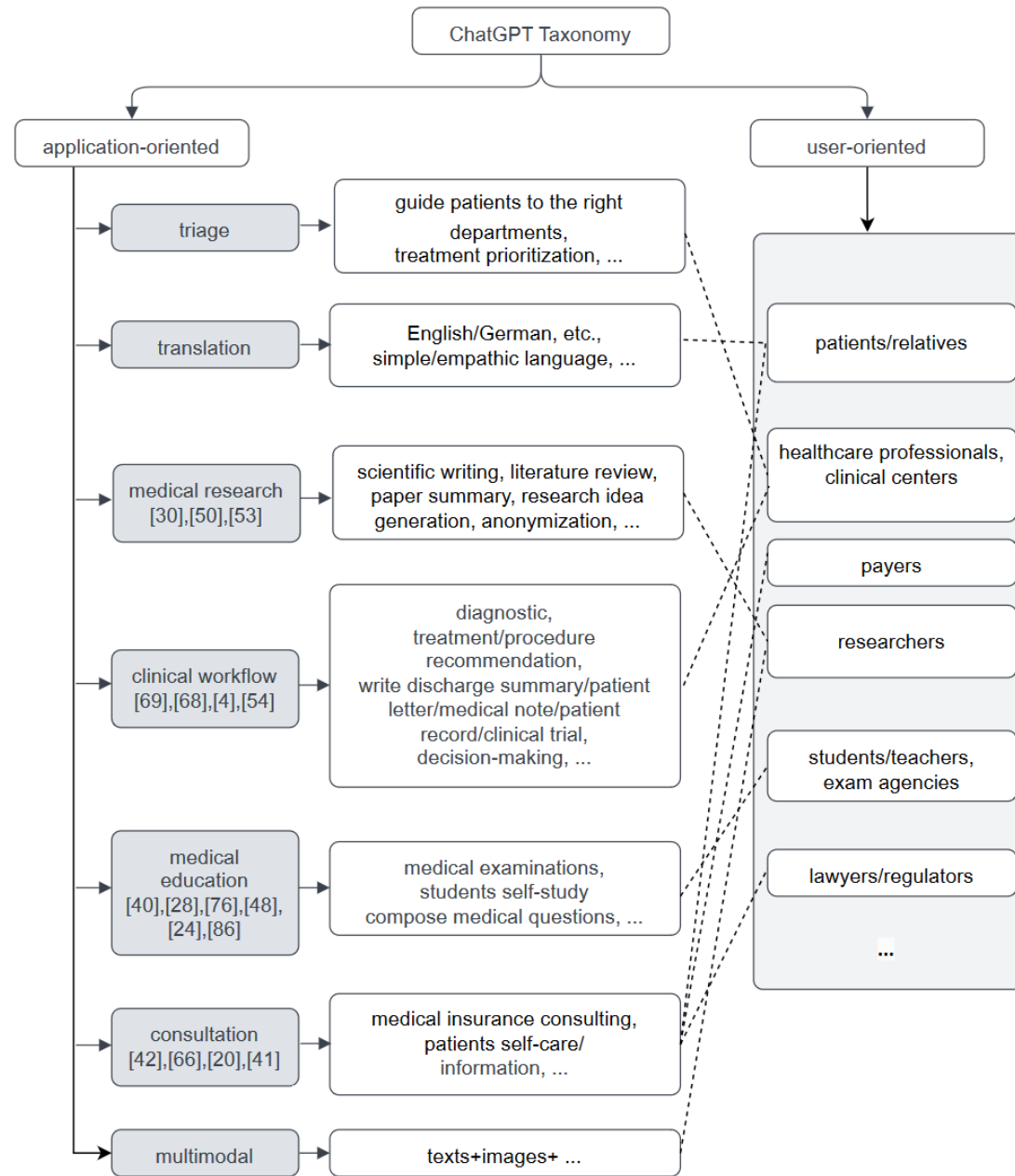


Figure 3: Application- and user-oriented Taxonomy used in the ChatGPT review. The references shown in the application boxes are the *Level 3* publications.

LLM Applications in Medicine

- Open-ended Natural Language UI
- Reasoning
- Med Education
- Patient Dialog & Communications
- Fine-tuned vs LLM
- Human Supervision

Capabilities of GPT-4 on Medical Challenge Problems

Harsha Nori¹, Nicholas King¹, Scott Mayer McKinney²,
Dean Carignan¹, and Eric Horvitz¹

¹Microsoft

²OpenAI

12 Apr 2023

Abstract

Large language models (LLMs) have demonstrated remarkable capabilities in natural language **understanding and generation across various domains**, including medicine. We present a comprehensive evaluation of GPT-4 [Ope23], a state-of-the-art LLM, on **medical competency examinations**

Summary

- USMLE & MultiMedQA
- GPT 3.5, GPT4 & GPT-4-base
- Zero-Shot without Context
- Text & Image
- Calibration: Trustworthy and Interpretable Probabilities
- Memorization: MELD (Precision/Recall)
- Probability Calibration (Trust)
- 20pts > GPT3.5 & Med-PaLM/Flan-PaLM 540B (Human 60.2%)
- Patient care is not Multiple Choice

Table 4: Performance comparison of the publicly released GPT-4 model with GPT-4-base.

Dataset	Component	GPT-4-base (5 shot)	GPT-4-base (zero shot)	GPT-4 (5 shot)	GPT-4 (zero shot)
USMLE Self Assessment	Step 1	86.72	85.38	85.21	83.46
	Step 2	91.50	90.62	89.50	84.75
	Step 3	85.23	85.23	83.52	81.25
USMLE Sample Exam	Step 1	85.71	84.87	85.71	80.67
	Step 2	85.00	86.67	83.33	81.67
	Step 3	92.70	93.43	90.71	89.78

CAN LARGE LANGUAGE MODELS REASON ABOUT MEDICAL QUESTIONS?

Valentin Liévin^{1,2} **Christoffer Egeberg Hother**³ **Ole Winther**^{1, 2, 4, 5}

¹ Section for Cognitive Systems, Technical University of Denmark, Denmark

² FindZebra ApS, Denmark

³ Department of Clinical Immunology, Rigshospitalet, Copenhagen University Hospital, Denmark

⁴ Center for Genomic Medicine, Rigshospitalet, Copenhagen University Hospital, Denmark

⁵ Bioinformatics Centre, Department of Biology, University of Copenhagen, Denmark

valv@dtu.dk, christoffer.egeberg.hother@regionh.dk, olwi@dtu.dk

ABSTRACT

Although large language models (LLMs) often produce impressive outputs, it remains unclear how they perform in real-world scenarios requiring strong reasoning skills and expert domain knowledge. We set out to investigate whether GPT-3.5 (Codex and InstructGPT) can be applied to answer and reason about

| 24 Jan 2023

Summary

- USMLE/MedMCQA vs PubMedQA
- GPT3.5 (Codex & InstructGPT)
- 3 Prompts: CoT, Zero/Few-Shot and KB Augmented
- Errors: Knowledge, Reasoning, Guessing Heuristics
- Codex 5-shot CoT ~ Human
 - USMLE 60.2%
 - MedMCQA 62.7%
 - PubMedQA 78.2%

Table 5: Zero-shot answering accuracy of InstructGPT (`text-davinci-002`) on the USMLE (test), MedMCQA (valid.) and PubMedQA (test) datasets. We report the best finetuned BERT-based methods. We tested 5 domain-specific CoT cues (#1-5) and report the mean performances with standard deviations. See Table 8, Appendix A, for a complete overview of our results, including results on the full MedMCQA test set.

Model	Grounding	Prompt	USMLE	MedMCQA	PubMedQA
InstructGPT	\emptyset	Direct	46.0	44.0	73.2
InstructGPT	\emptyset	CoT #1–5	46.1 \pm 0.7	40.4 \pm 2.2	59.9 \pm 3.5
InstructGPT	BM25	Direct	47.3	46.7	–
InstructGPT	BM25	CoT #1–5	46.4 \pm 0.7	42.5 \pm 1.7	–
InstructGPT	\emptyset	Ensemble (n=6) ¹	50.0	42.4	70.4
InstructGPT	BM25	Ensemble (n=6) ¹	49.3	48.8	–
InstructGPT	\emptyset + BM25	Ensemble (n=12) ¹	53.1	47.6	–
Finetuned BERT	BM25, DPR, \emptyset		44.6 ²	43.0 ³	72.2 ²
Human (passing score)			60.0	50.0	–
Human (expert score)			87.0	90.0	78.0

¹Majority voting (direct + CoT prompts), ² BioLinkBERT (Yasunaga et al., 2022), ³ PubMedBERT (Gu et al., 2021) from Pal et al. (2022).

Table 6: Frequency of observed patterns (A, B, C, D, E, F) identified among 50 CoTs generated by InstructGPT with temperature $\tau=0$. The CoTs are generated based on USMLE questions and using the CoT prompts #1–5 (Table 4). We report the frequencies of CoTs with correct and incorrect predictions along with the total.

Pattern	Correct [†] (16)	Incorrect [†] (34)	Total (50)
A Correct reasoning step*	94% (15)	59% (20)	70% (35)
B Correct recall of knowledge*	87% (14)	65% (22)	72% (36)
C Correct reading comprehension*	100% (16)	85% (29)	90% (45)
D Incorrect reasoning step*	12% (2)	86% (29)	62% (31)
E Incorrect or insufficient knowledge*	25% (4)	74% (25)	58% (29)
F Incorrect reading comprehension*	6% (1)	50% (17)	36% (18)

* At least one (...), [†] Correct/incorrect prediction

Figure 3: Answering accuracy of Codex 5-shot CoT (code-davinci-002) on the USMLE (test), the MedMCQA (valid.) and the PubMedQA (test) datasets for 100 CoTs sampled with temperature $\tau \in \{0, 0.5\}$. We report the average accuracy for ensemble models evaluated using random subsets of $k' = 1 \dots 100$ CoTs. We display the performances of the best finetuned methods along with the lower human baselines.

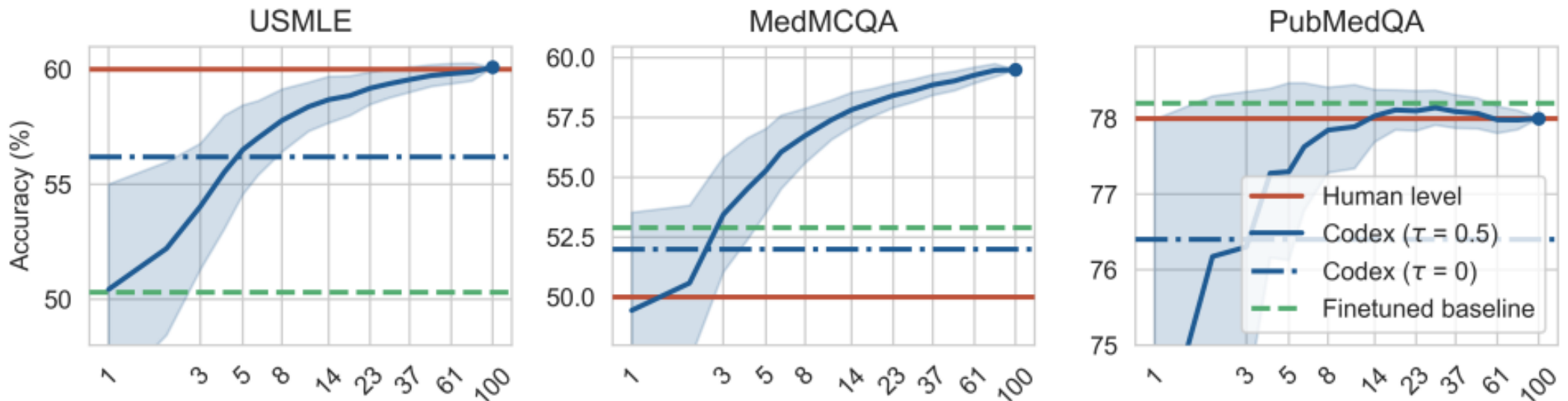
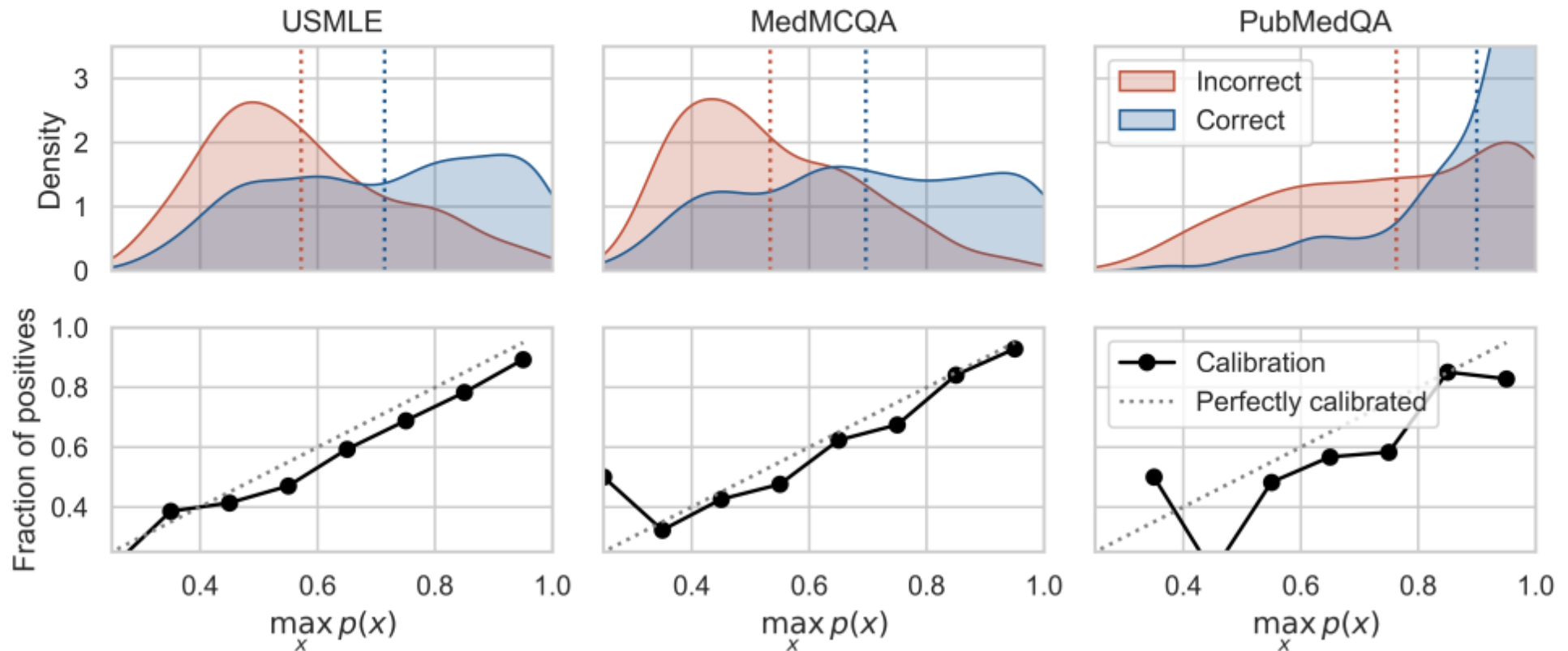


Figure 4: First row: distribution of the probability assigned to the correct label for correct predictions and incorrect predictions (see Equation 1). Second row: calibration plot. The probabilities are obtained using Codex 5-shot CoT and an ensemble of $k = 100$ predictions sampled with temperature $\tau = 0.5$.



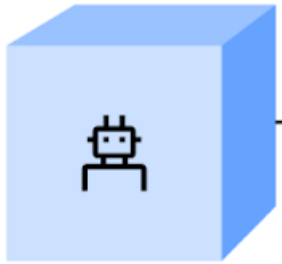
MedTech Startups and LLMs

As of April 2023

(Cambrian explosion, many under the radar)

Monitoring Generative AI Models

AI APPLICATION

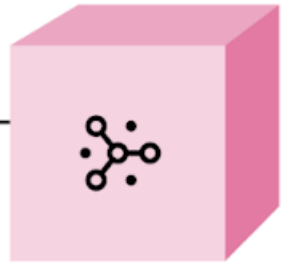


Prompt →

- ✓ **Correctness** – How accurate is the response?
- 📈 **Performance** – Is the model decaying?
- 💰 **Cost** – Where is the best ROI?
- 💬 **Prompt Monitoring** – How are my prompts changing?
- 🕒 **Latency** – How long is the model taking?
- 🔍 **Transparency** – Why did the model say that?
- 📈 **Bias** – Is the model's response biased?
- 📊 **AB Test** – Is the model changing across versions?
- 🔒 **Safety Monitoring** – Is the model's response safe?

← Response

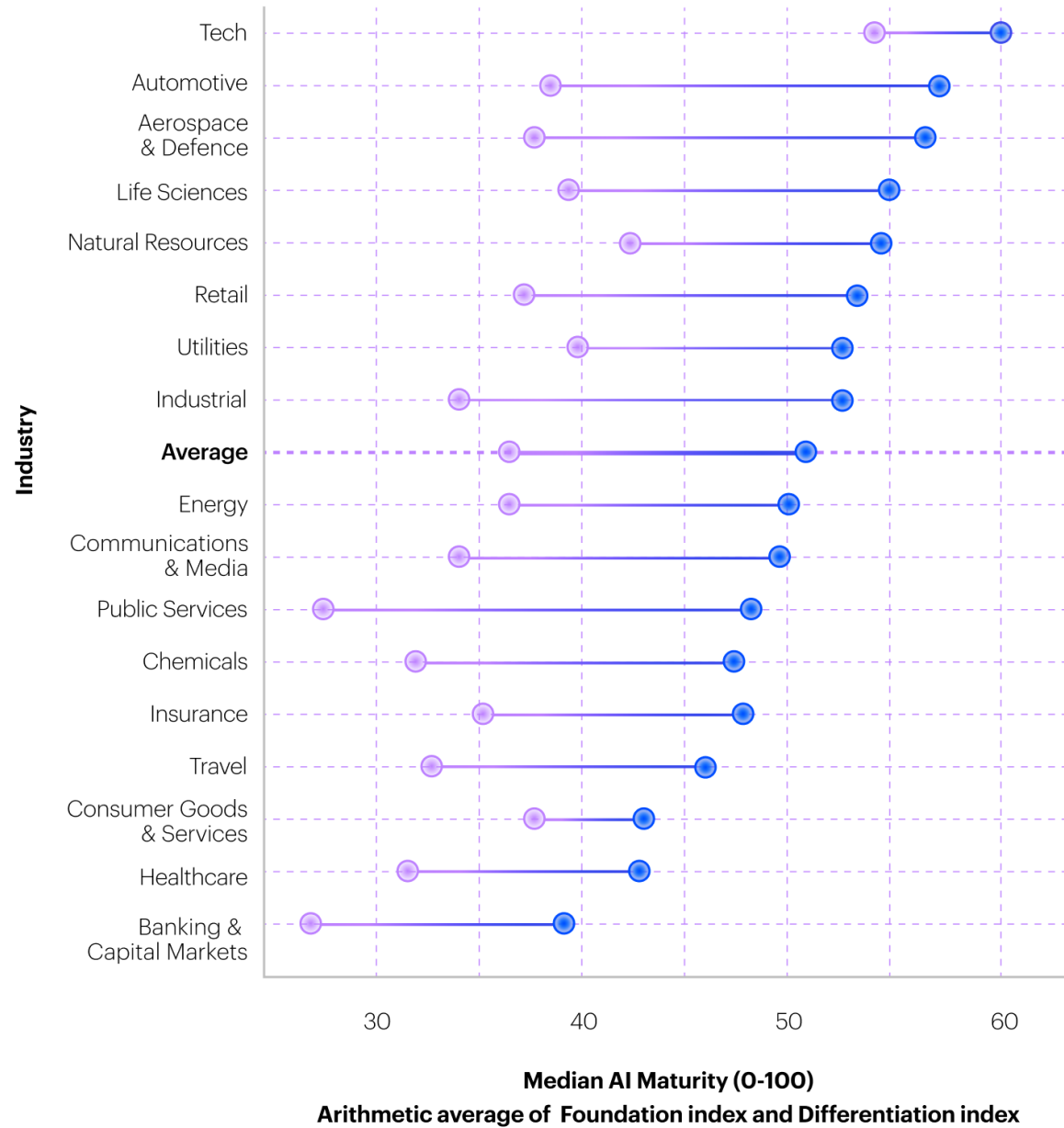
GAI MODEL



 fiddler

Time
● 2021
● 2024

The median AI Maturity Index in 2021 and 2024 by industry



Impact of AI and ML on Select Healthcare Outcomes in 2022 According to US Healthcare Executives

% of respondents

Improving clinical outcomes



Improving operational performance



Improving health system efficiency



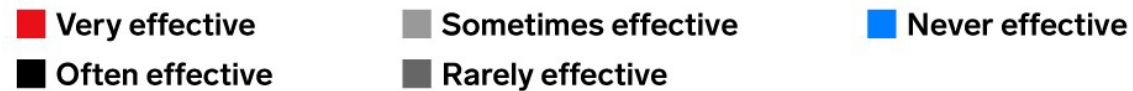
Improving administrative performance



Improving financial outcomes



Improving consumer engagement



Note: numbers may not add up to 100% due to rounding

Source: Innovaccer, "Healthcare's Data Readiness Crisis" conducted by Morning Consult, July 16, 2022

CBINSIGHTS
Digital Health
150
2022

Care coordination & collaboration

allinea health, Auxa Health, Buddy Healthcare, CareAlign, dina, eon, HEALTH [at] SCALE TECHNOLOGIES, Oath, rimidi, TCARE

Clinical intelligence

abridge, ARTIQ, ImmersiveTouch, iodine, MindTrace, OVERJET, theatr, VIDEA HEALTH

Home health & wellness

cala, CURETT, doccla, Embr Labs, Evvy, MedArrive, NymbL, suvera, sword

Computer-aided imaging

BrainSightAI, CEREBRIU, clarius, DIGITAL DIAGNOSTICS, endolyse, harrison.ai, ITERATIVE HEALTH, SUBTLE MEDICAL, SWIFT, iz.ai

Interoperability, data, & analytics

AETION, Availity, covera health, embold HEALTH, GenXys, Hospital IQ, J2 HEALTH, lifebit, Lightbeam Health Solutions, LYNXCARE, nubentis, particle, REDOX, rune labs, SEQSTER, SOCIALCLIMB, TRUVETA, xealth

Digital front door & patient engagement

Careology, CIPHERHealth, Healthee, heartbeat, hyro, kinetik, KODAHEALTH, mayaMD, Nodal, Playback Health, REDI HEALTH, Stride, univfy, vinehealth, WILDFLOWER, XP HEALTH

Virtual care

bloom.care, BRAVE, EQUIP, EQUUM Medical, exseed, Fourier Intelligence, Jasper, LEVY HEALTH, neoth, nuvo, recoveryone, Renalis, SUMMUS, THOUGHTFULL, vivante HEALTH, WYSA

Screening, monitoring, & diagnostics

alio, Aural Analytics, BABYSCRIPTS, biospectal, casana, CLOUD DX, COFACTOR genomics, Eko, ellipsis HEALTH, ENLITIC, galileo, GENOMENON, hyle.ai, IDOVEN, Infermedica, ixlayer, KINTSUGI, L7INFORMATICS, Ligence, MOBIO, NEUROVINE, Optellum, podimetrics, prevent biometrics, Qventus, skin ANALYTICS, spect, Starling

Digital therapeutics

AMALGAM, AppliedVR, aemo, GYENNO, MEDRhythms, metaMe Health, Rocket VR HEALTH

Digital pharmacy & DME

9amHealth, HouseRx

Workflow automation & digitization

CertifyOS, element5, Foldhealth, hint health, OpenLoop, OpenMedical, Rhyme, Vytalize

Clinical trials tech

TOPOGRAPHY HEALTH, UNLEARN, xCures

Hybrid care

Caraway, homeward, kindbody, MAVEN, Nest Health, SalvoHealth, sami, tia, vori health, Waymark

Revenue cycle management

enter, JANUS, LyfeGen, soda health

Note: Companies are private as of 11/4/22

End

(See PowerPoint speaker notes
for references and links)