# Introduction to LLMs with a Focus on Biomedical Data Science

Presenter: Shubo Tian

# Learning Objective

- Understanding what is language model, type of language models, and existing major language models and their applications in general and in biomedicine.

# What is Language Model

Given a sequence of words/tokens:

$$w_1, w_2, \ldots, w_n$$

a model that computes either of the following probabilities:

$$P(w_1, w_2, \ldots, w_n) \text{ or}$$
$$P(w_n | w_1, w_2, \ldots, w_{n-1})$$

is called a language model.

$$P(w_1, w_2, \ldots, w_n) = \prod_{i=1}^{n} P(w_i | w_1, \ldots, w_{i-1})$$

**P("P53 is a tumor suppressor gene.")**

$= P(\text{"P53"}) \cdot$

$\quad P(\text{"is" |"P53"}) \cdot$

$\quad P(\text{"a"|"P53 is"}) \cdot$

$\quad P(\text{"tumor"|"P53 is a"}) \cdot$

$\quad P(\text{"suppressor"|"P53 is a tumor"}) \cdot$

$\quad P(\text{"gene"|"P53 is a tumor suppressor"}) \cdot$

$\quad P(\text{"."|"P53 is a tumor suppressor gene"})$
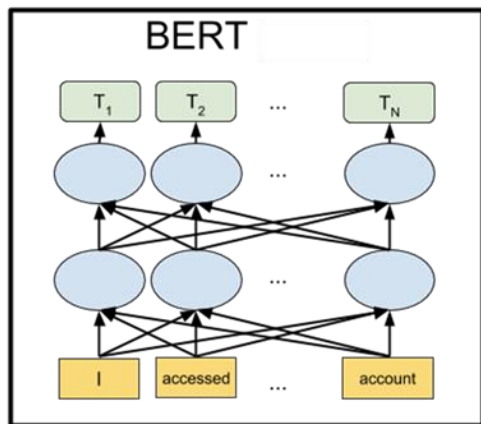
# Language Modeling

Word count modeling (n-gram)

$$P(w_n|w_1, w_2, \ldots, w_{n-1}) = \frac{\text{count}(w_1, w_2, \ldots, w_n)}{\text{count}(w_1, w_2, \ldots, w_{n-1})}$$

Neural network modeling

$$P(w_n|w_1, w_2, \ldots, w_{n-1}) = P(w_n|h)$$
$$= \frac{\exp(h^T \text{emb}(w_n))}{\sum_{w' \in \text{Vocab}} \exp(h^T \text{emb}(w'))}$$

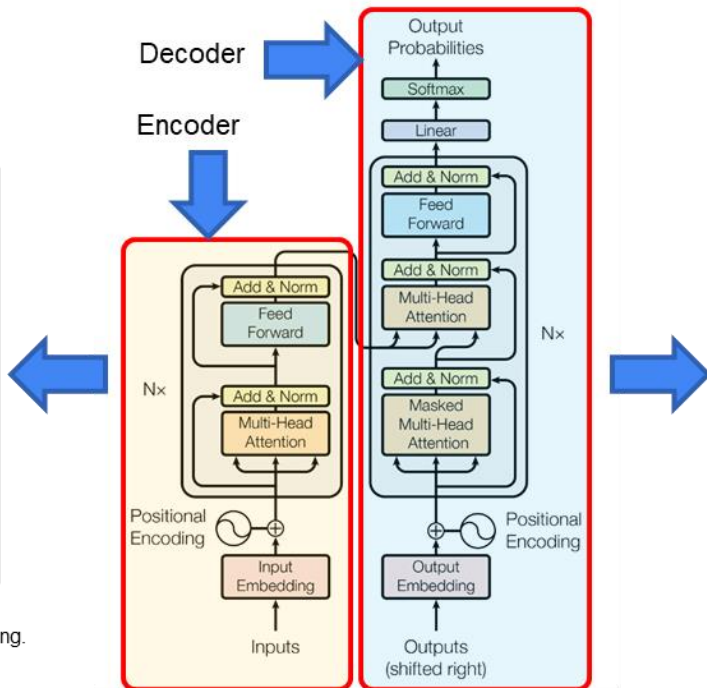$$h = \text{Encoder}(w_1, w_2, \ldots, w_{n-1})$$

Jurafsky, D. and Martin, J.H. (2023) Speech and Language Processing. https://web.stanford.edu/~jurafsky/slp3

# Transformer and LLMs
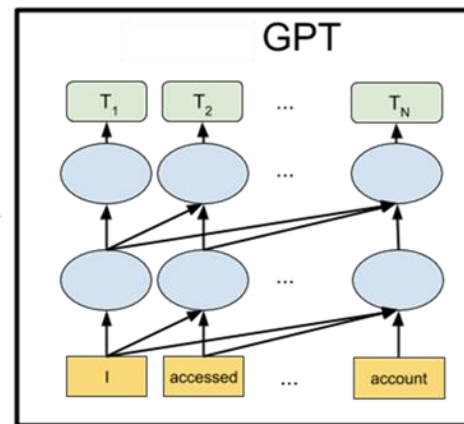


Devlin, J. *et al.* (2019) BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv:1810.04805 [cs]*.

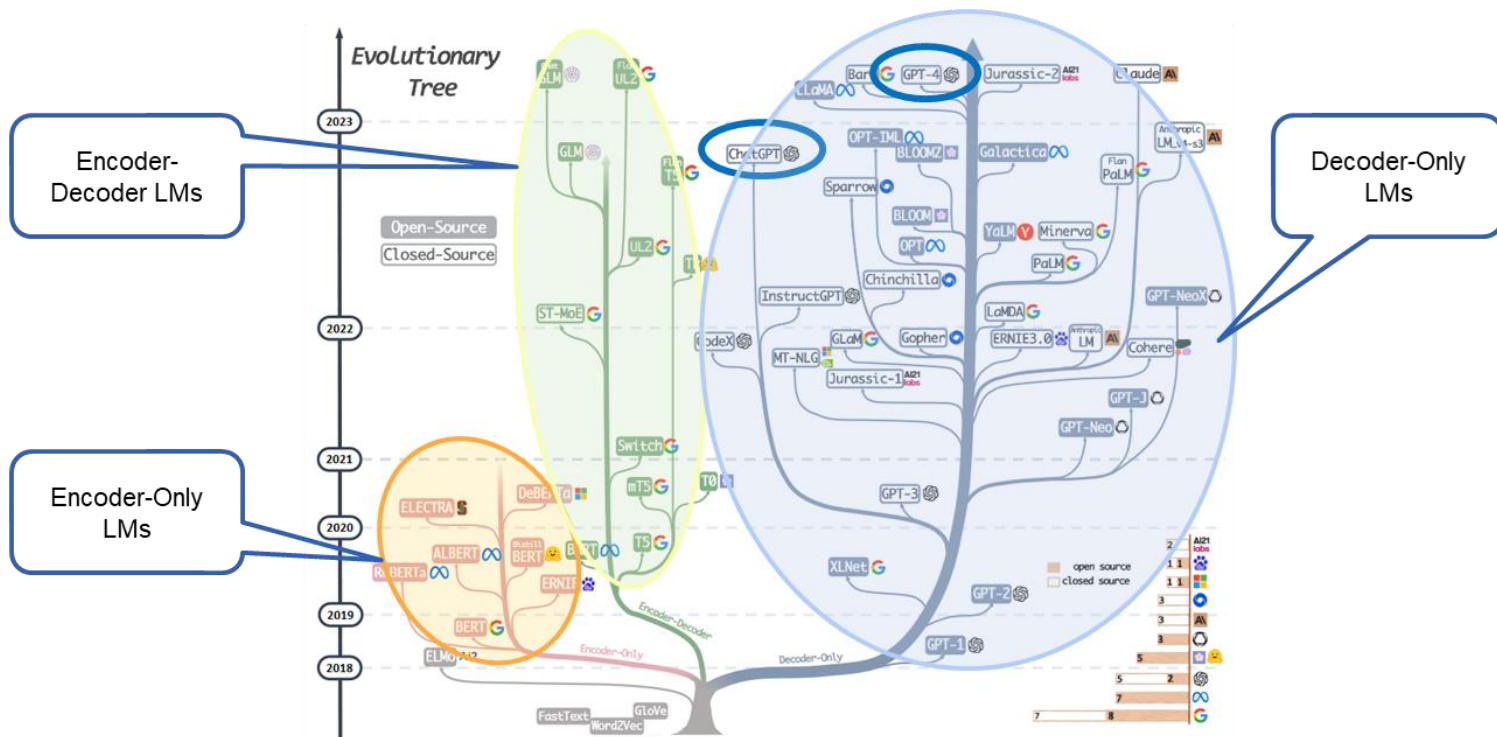https://ai.googleblog.com/2018/11/open-sourcing-bert-state-of-art-pre.html

Vaswani, A. *et al.* (2017) Attention Is All You Need. *arXiv:1706.03762 [cs]*.

Radford, A. *et al.* Improving Language Understanding by Generative Pre-Training.
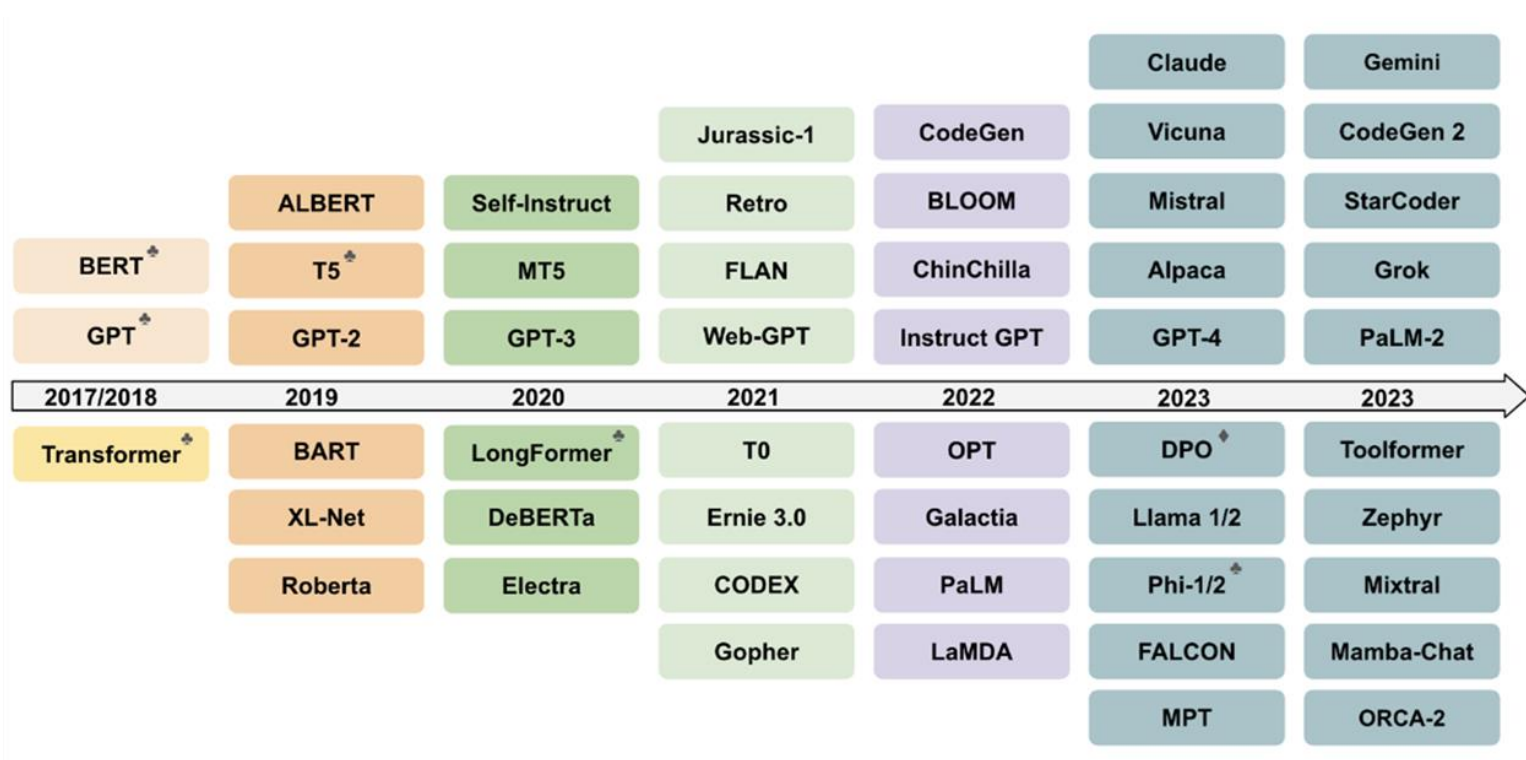
https://ai.googleblog.com/2018/11/open-sourcing-bert-state-of-art-pre.html
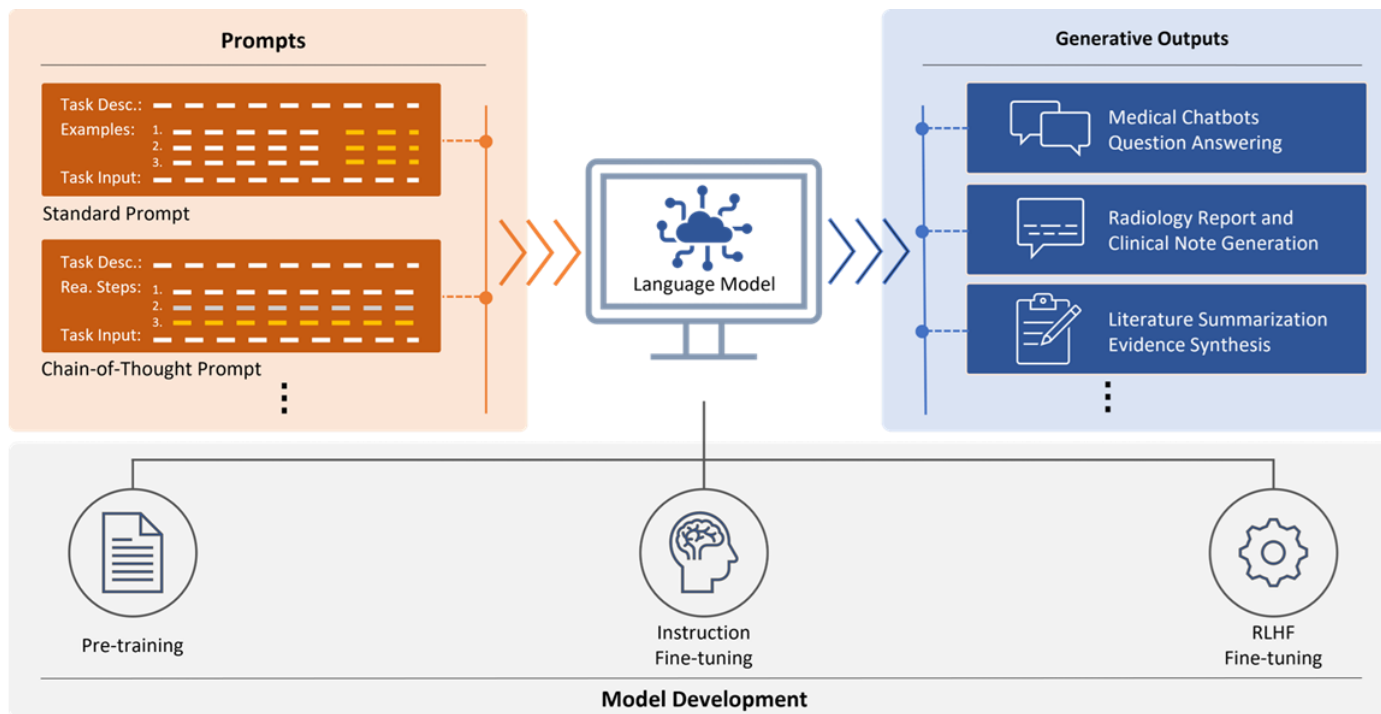
# Type of LLMs



Yang, J., Jin, H., Tang, R., Han, X., Feng, Q., Jiang, H., Yin, B. and Hu, X. (2023) Harnessing the Power of LLMs in Practice: A Survey on ChatGPT and Beyond. 10.48550/arXiv.2304.13712.
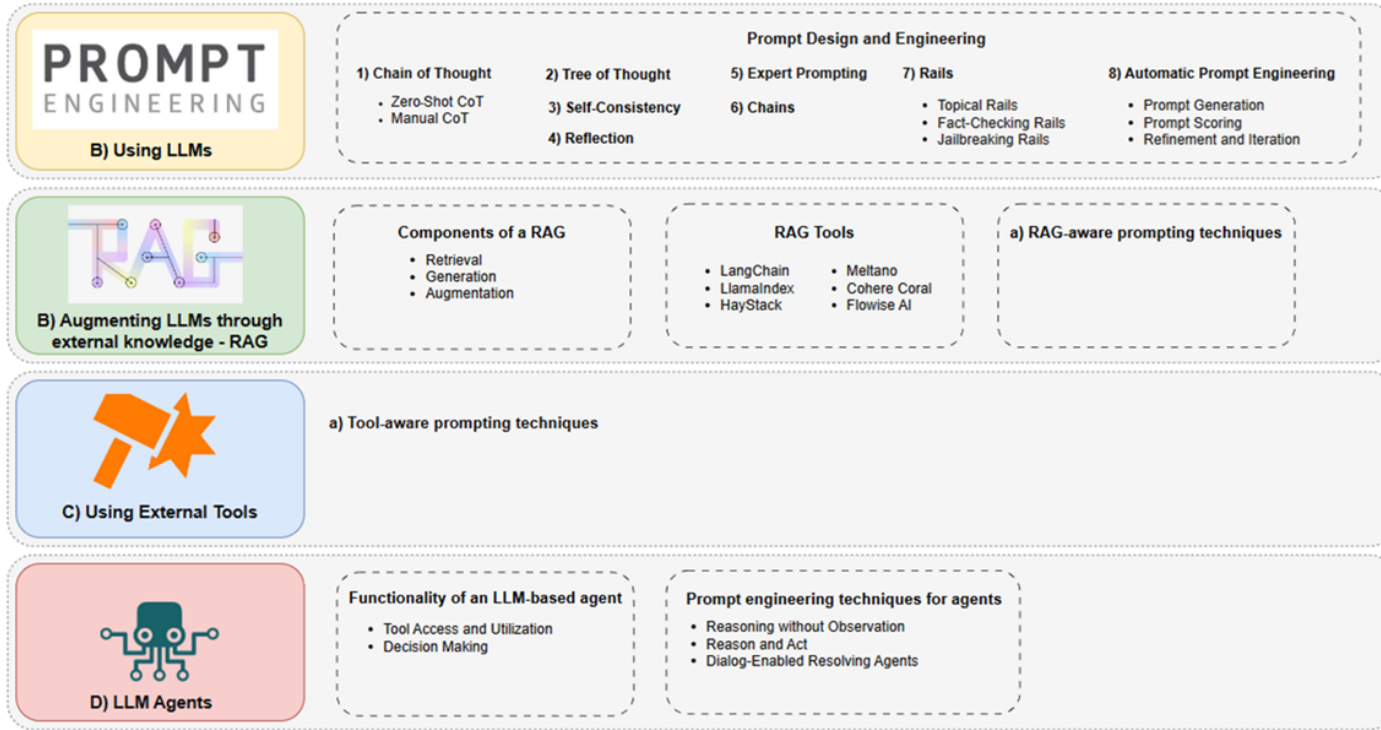
# Timeline of Representative LLMs



| 2017/2018 | 2019 | 2020 | 2021 | 2022 | 2023 | 2023 |
|---|---|---|---|---|---|---|
| | | | | | | Claude | Gemini |
| | | | Jurassic-1 | CodeGen | Vicuna | CodeGen 2 |
| | ALBERT | Self-Instruct | Retro | BLOOM | Mistral | StarCoder |
| BERT | T5 | MT5 | FLAN | ChinChilla | Alpaca | Grok |
| GPT | GPT-2 | GPT-3 | Web-GPT | Instruct GPT | GPT-4 | PaLM-2 |
| Transformer | BART | LongFormer | T0 | OPT | DPO | Toolformer |
| | XL-Net | DeBERTa | Ernie 3.0 | Galactia | Llama 1/2 | Zephyr |
| | Roberta | Electra | CODEX | PaLM | Phi-1/2 | Mixtral |
| | | | Gopher | LaMDA | FALCON | Mamba-Chat |
| | | | | | MPT | ORCA-2 |

Minaee, S., Mikolov, T., Nikzad, N., Chenaghlu, M., Socher, R., Amatriain, X. and Gao, J. (2024) Large Language Models: A Survey.
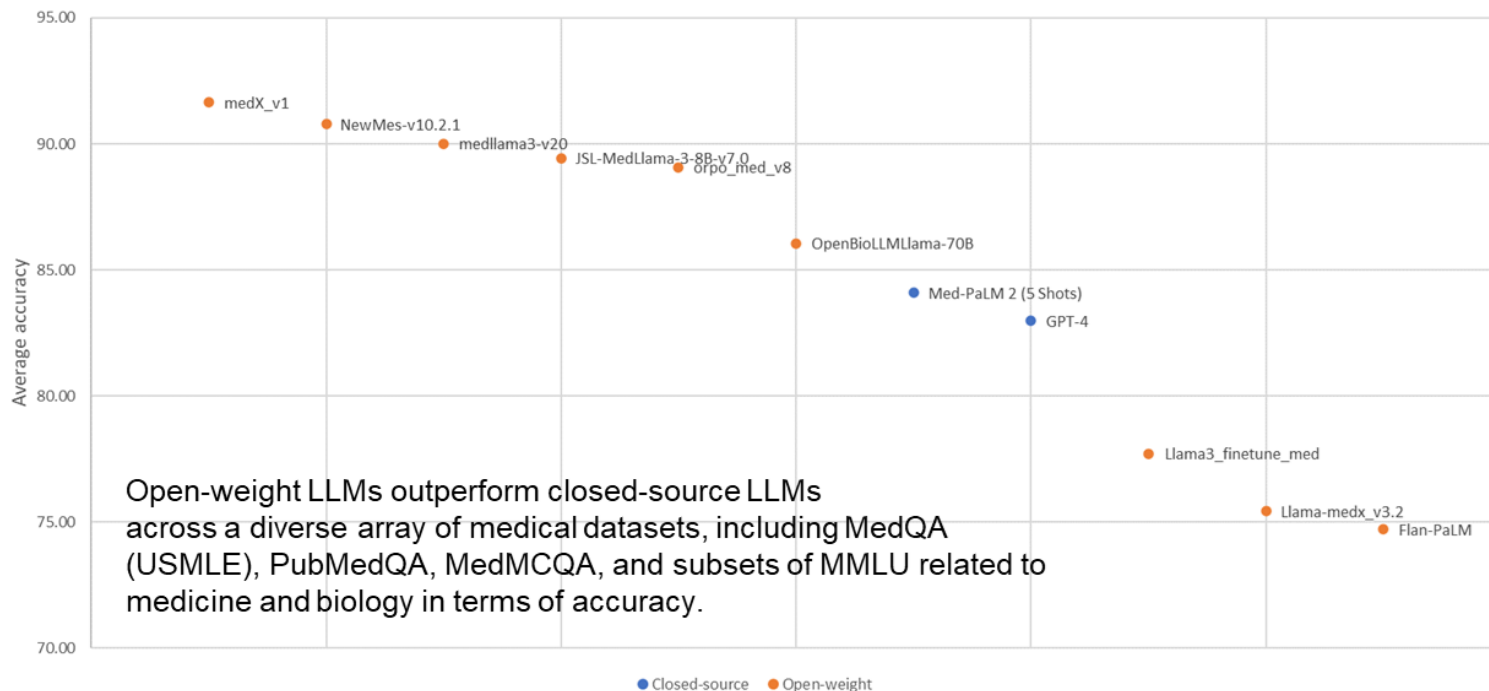
7

# The Paradigm of LLMs

Tian, S., Jin, Q., Yeganova, L., Lai, P.-T., Zhu, Q., Chen, X., Yang, Y., Chen, Q., Kim, W., Comeau, D.C., et al. (2024) Opportunities and challenges for ChatGPT and large language models in biomedicine and health. Briefings in Bioinformatics, 25, bbad493.

# How LLMs are Used

Minaee, S., Mikolov, T., Nikzad, N., Chenaghlu, M., Socher, R., Amatriain, X. and Gao, J. (2024) Large Language Models: A Survey.

# Major LLMs and Their Performance



Closed-source vs. open-weight models

Llama 3 405B from Meta closes the gap between closed-source and open-weight models.



Closed-source vs. Open-weight models (Arena ELO)

# LLMs in Biomedicine



Accuracy of LLM performance on the MedQA (USMLE) dataset has increased from the level of human passing by GPT3.5 to the level close to human expert by Med-PaLM 2 in less than half a year.

Tian, S., Jin, Q., Yeganova, L., Lai, P.-T., Zhu, Q., Chen, X., Yang, Y., Chen, Q., Kim, W., Comeau, D.C., et al. (2024) Opportunities and challenges for ChatGPT and large language models in biomedicine and health. Briefings in Bioinformatics, 25, bbad493.

# Major LLMs in Biomedicine and Their Performance



Open-weight LLMs outperform closed-source LLMs across a diverse array of medical datasets, including MedQA (USMLE), PubMedQA, MedMCQA, and subsets of MMLU related to medicine and biology in terms of accuracy.

https://huggingface.co/spaces/openlifescienceai/open_medical_llm_leaderboard

# Applications of LLMs in Biomedicine

# LLMs in the Context of AI in Biomedicine



Moor,M. et al. (2023) Foundation models for generalist medical artificial intelligence. Nature, 616, 259–265.

14

# Thank you!

Next talk in line: How to Use GPT-3.5 and GPT-4 with Python, Qiao Jin, M.D.