

Transformers in Language

Tripp Deep Learning F23



TODAY'S GOAL

By the end of the class, you should be familiar with some of the key methods and results in the application of transformers to language problems.

Summary

1. Large language models (LLMs) can perform new tasks without fine tuning
2. LLMs are getting larger
3. LLMs can produce convincing output
4. Causal language modelling isn't enough to make LLMs useful
5. LLMs can be fine-tuned using parameter-efficient updates

LLMS CAN PERFORM NEW TASKS WITHOUT FINE TUNING

Transfer learning

- As we discussed previously, it is common to pretrain transformers with a self-supervised task on a large text dataset and fine tune on a smaller task-specific dataset
- However, recent models have been applied to new tasks without fine tuning

Zero-shot inference

- GPT-2 (Radford et al., 2019)
 - A decoder-only transformer with minor differences from Vaswani et al. in initialization and the location of layer norm
 - A large model for the time with 48 layers, 1.5B parameters
 - Trained on causal language modelling with WebText, a new dataset of text from web pages that were linked from Reddit that received at least 3 karma, and excluding Wikipedia (~8M pages)
- Its performance was evaluated on tasks that it was not been trained on (i.e., no task-specific training examples)
 - Often some kind of prompt was added to the input to indicate what was expected

Zero-shot inference

- On reading comprehension tasks, GPT-2 was competitive with supervised methods
 - Example: Good performance on Children's Book Test, a multiple-choice masked language modelling test with children's books
- On other tasks such as summarization and question answering, GPT-2 generally wasn't competitive with supervised models but outperformed simple baselines
 - Example: The model was prompted to perform summarization by appending "TL/DR" to input text and taking the first 3 sentences of output; slightly outperformed reproducing 3 random sentences from the text
- The model showed rudimentary English-French translation abilities although non-English documents had been removed from the corpus, leaving only bits of French remaining within English documents (0.025% total content)

Few-shot inference

- GPT-3 (Brown et al., 2020)
 - Larger version of GPT-2 with 175B parameters, trained with more data
- Evaluated on various language tasks in several settings:
 - Zero-shot: Input sequence includes a prompt to perform the task
 - One-shot: An example of correct task completion included in the input sequence
 - Few-shot: Usually 10 to 100 examples of correct task completion in the input sequence
- Performance usually improved from zero to one to few-shot performance, with the latter often competitive with strong supervised baselines (rarely state-of-art)
- Performance improved with model size, and the gap between zero and few-shot performance also increased with model size

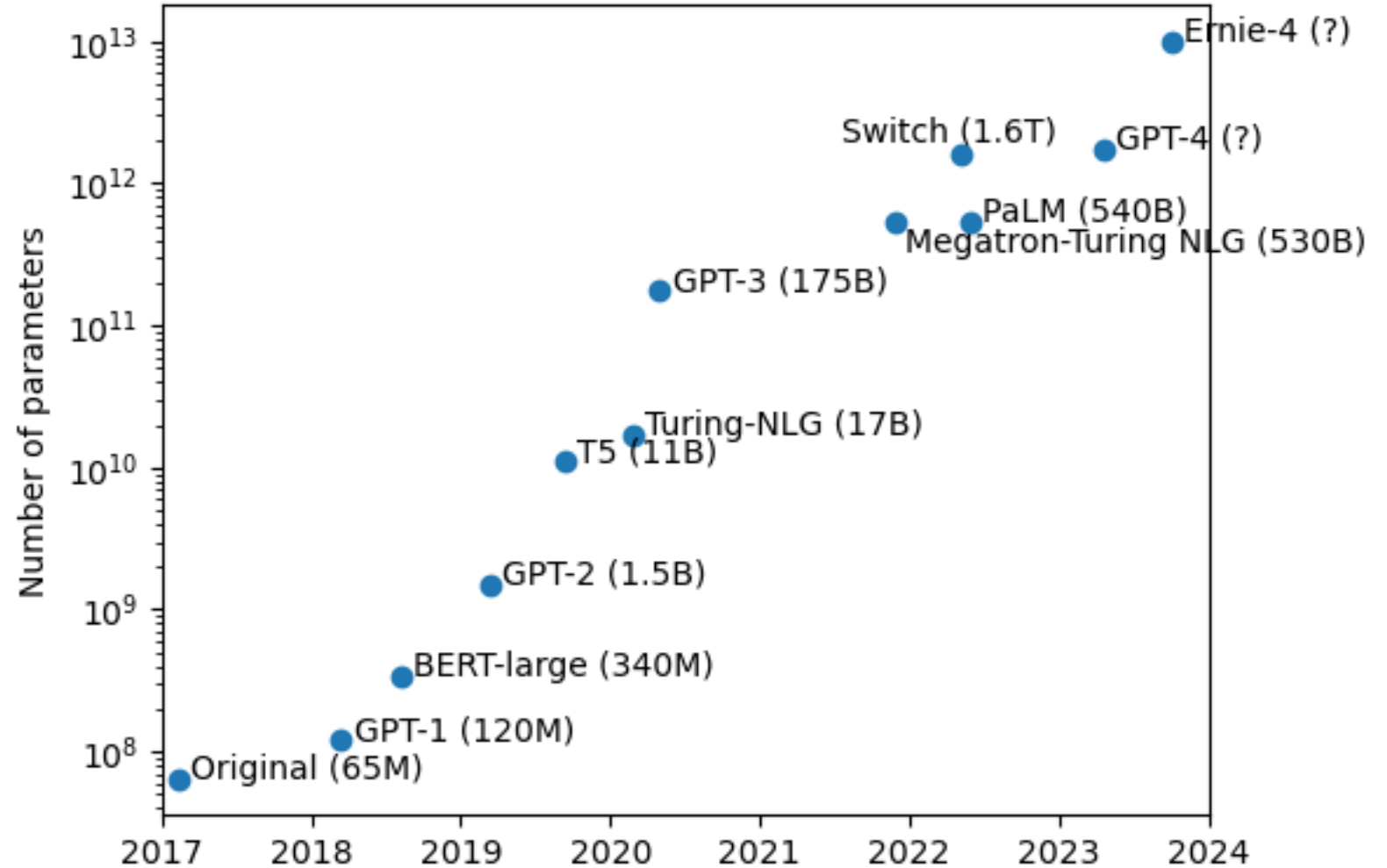
Few-shot inference

- More recent models such Megatron-Turing NLG (Smith et al., 2022), PaLM (Chowdhery et al., 2022), and GPT-4 (OpenAI, 2023) have still better few-shot performance
- Better performance is still achieved with fine tuning than with few-shot prompting (Chowdhery et al., 2022)

LLMS ARE GETTING LARGER

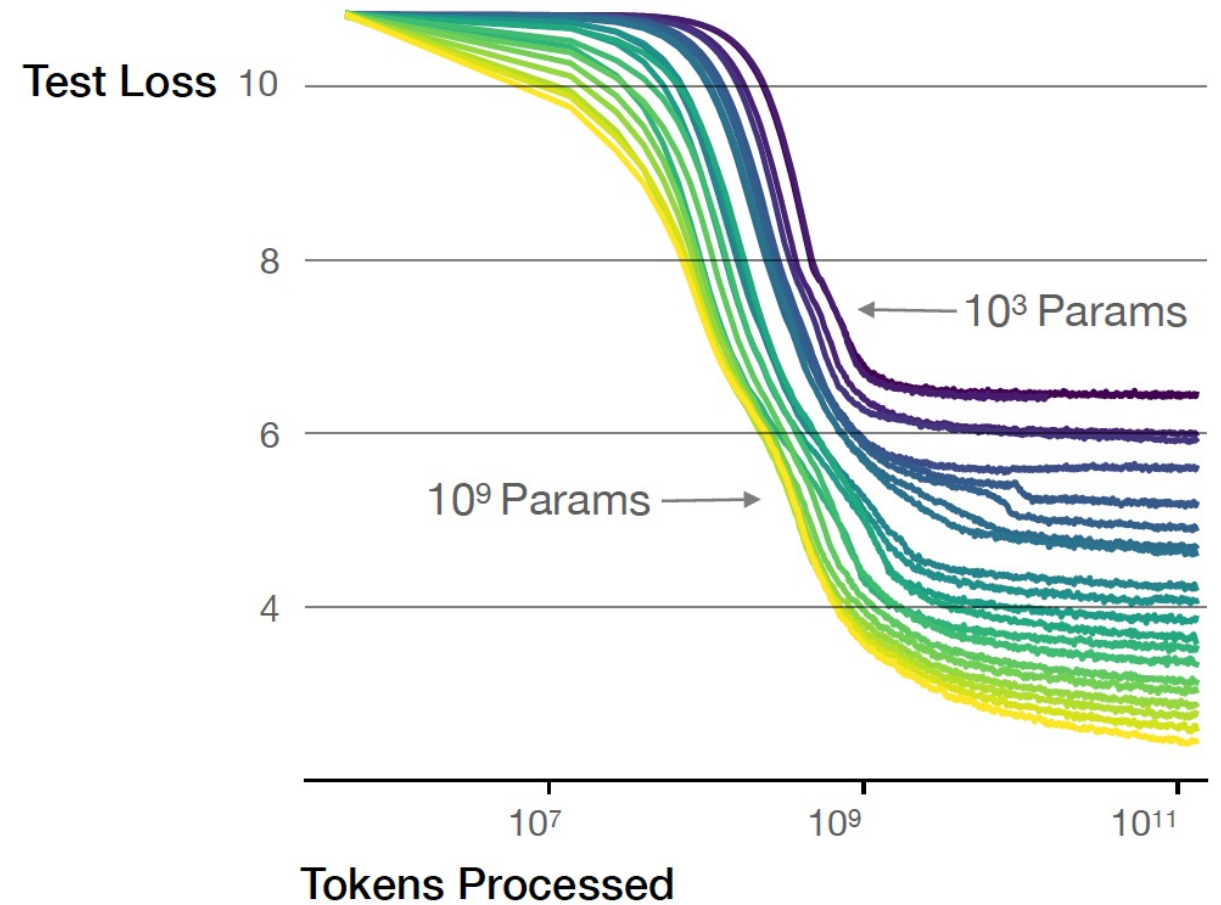
Recent transformers

- The largest language transformers have gotten larger by almost an order of magnitude every year



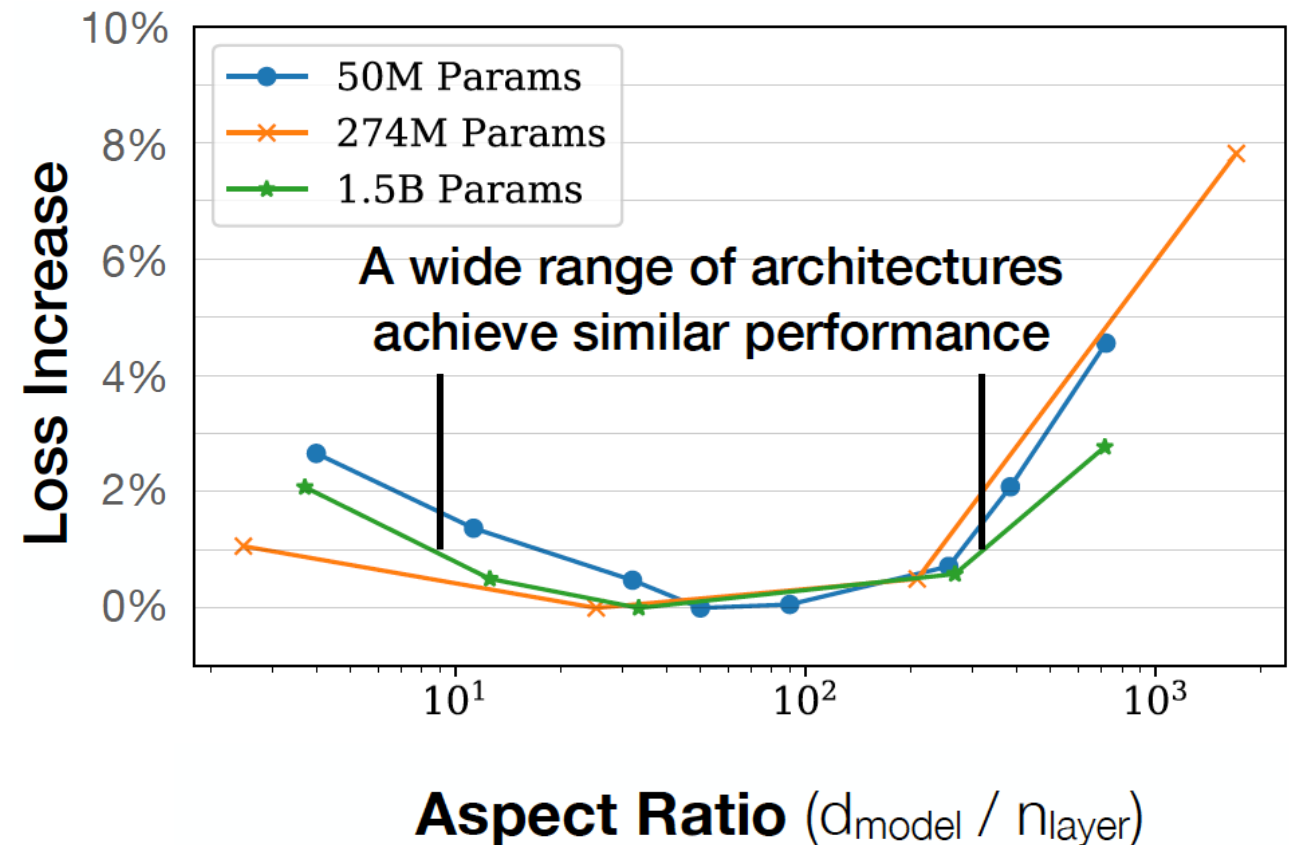
Performance improves with scale

- Kaplan et al. (2020) showed that transformer performance improves reliably with the number of parameters and dataset size
- They found that larger models required fewer training samples to reach the same performance



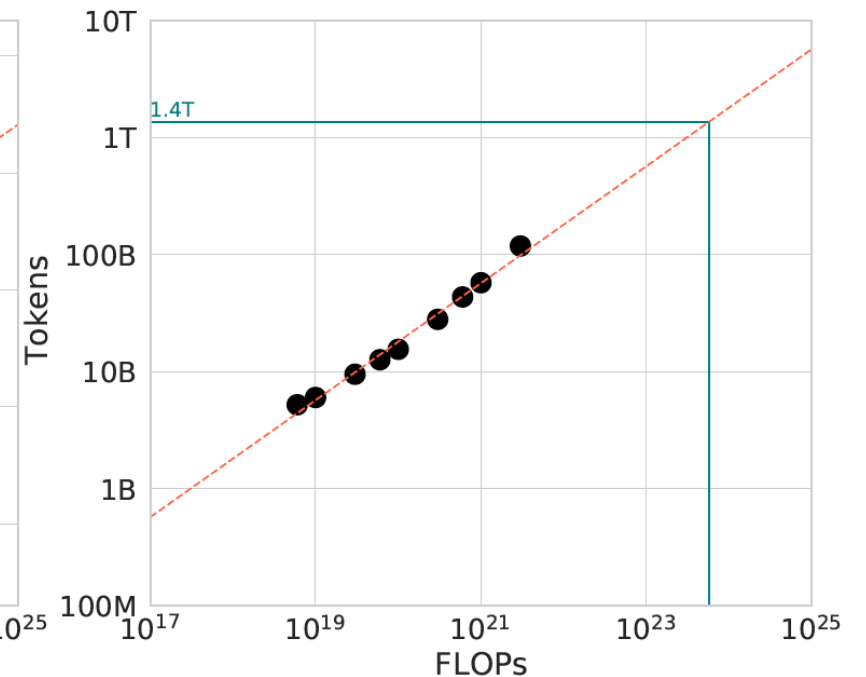
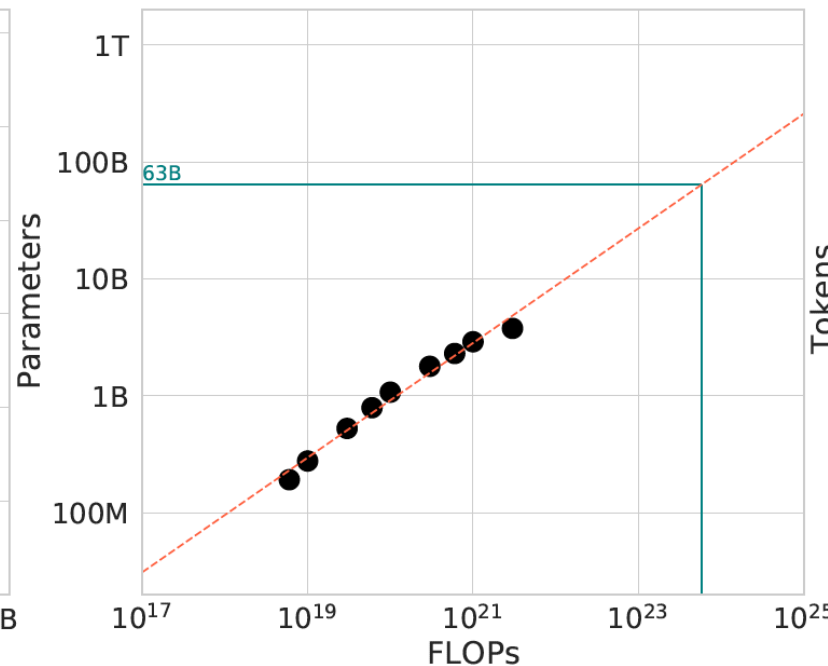
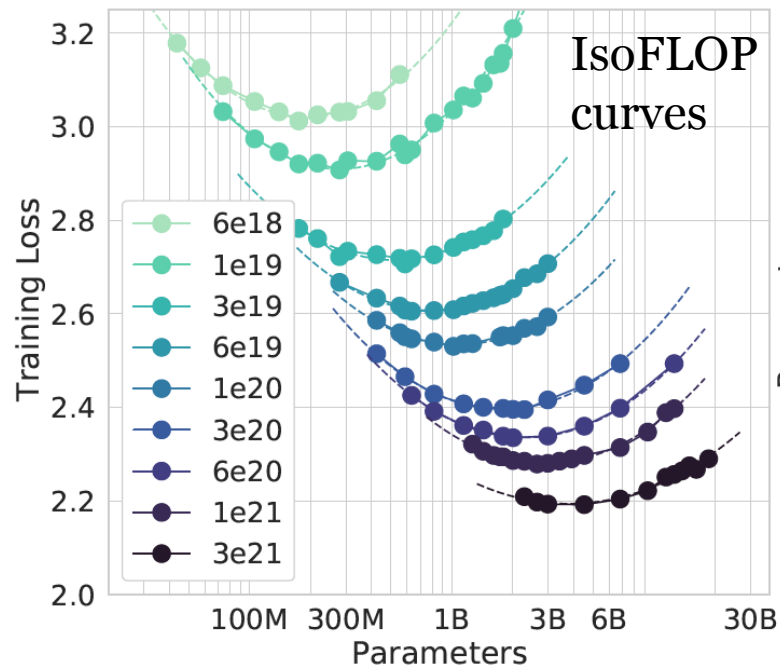
Performance improves with scale

- Kaplan et al. (2020) also found that performance was less sensitive to model shape, e.g.:
 - Dimension of FF layers relative to embeddings
 - Ratio of embedding dimension to number of layers
 - Ratio of embedding dimension to number of attention heads



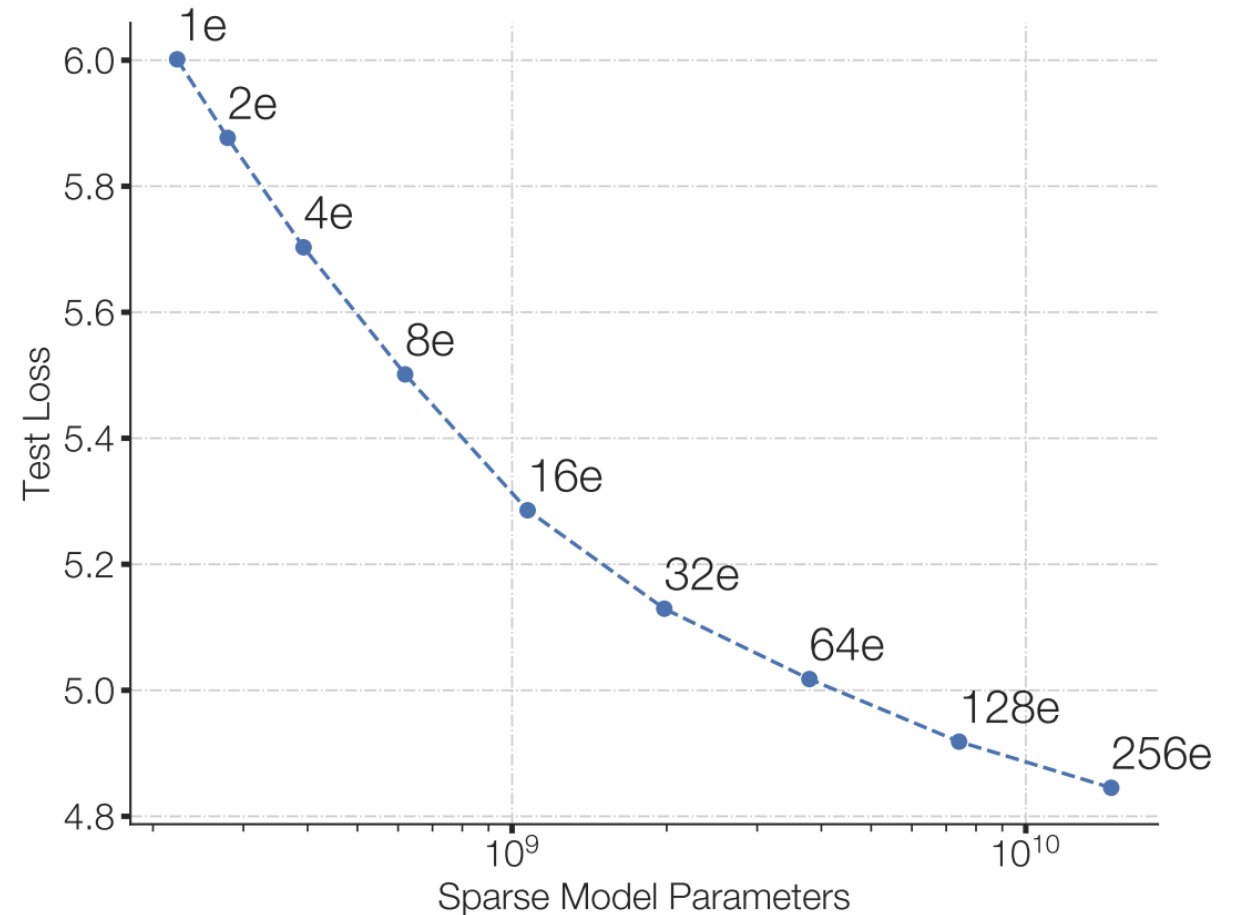
Model size vs. data size for a given compute budget

- Several recent models have been trained with $\sim 300\text{B}$ tokens
- Hoffmann et al. (2022) found that large models are undertrained
 - With increasing compute budget, model size and data size should scale at the same rate (e.g., double model size, double data size)



Improvement with scale

- Switch Transformers (Fedus et al., 2022)
 - Based on T5
 - Within each layer, each sequence element is routed to one of multiple versions of the layer that can run on different hardware
 - Routing is based on a softmax of a linear map of the token
 - This increases the number of parameters in the whole model without increasing computation cost, since each element is only routed to one of the parallel options

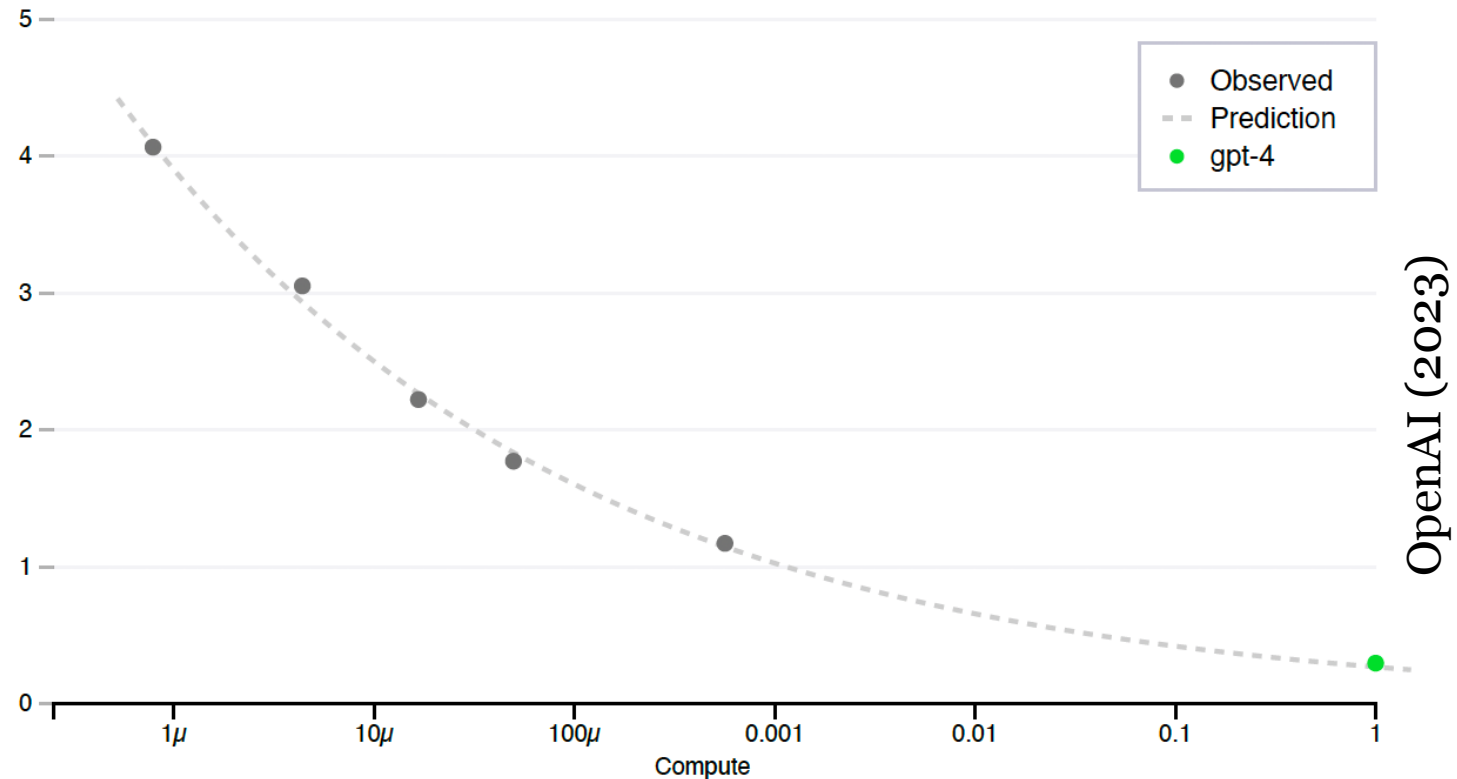


Improvement with scale

- New capabilities over the last few years largely due to scale
- LLMs have not yet hit a performance ceiling but improve with size and data
- The world has ~10x as much high-quality text as existing models have used but about ~10B times as much multimodal data

Capability prediction on 23 coding problems

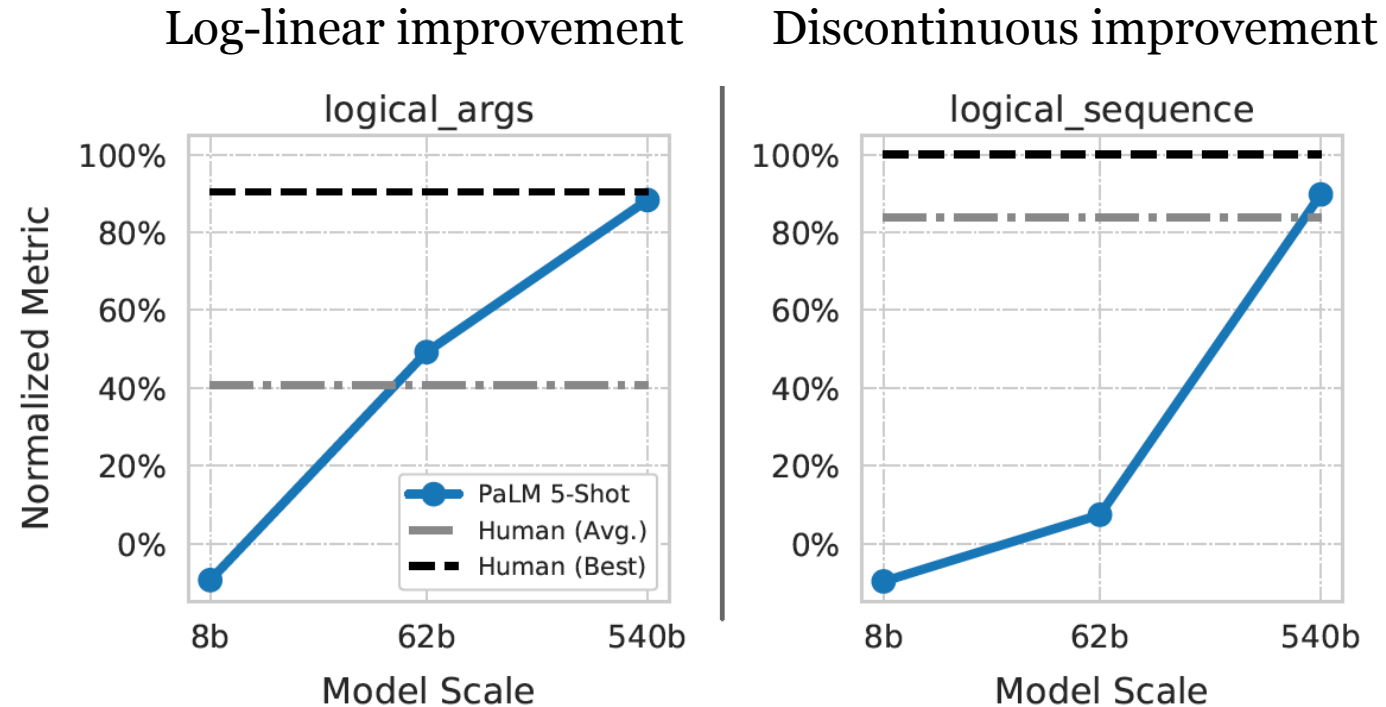
– Mean Log Pass Rate



Emerging capabilities

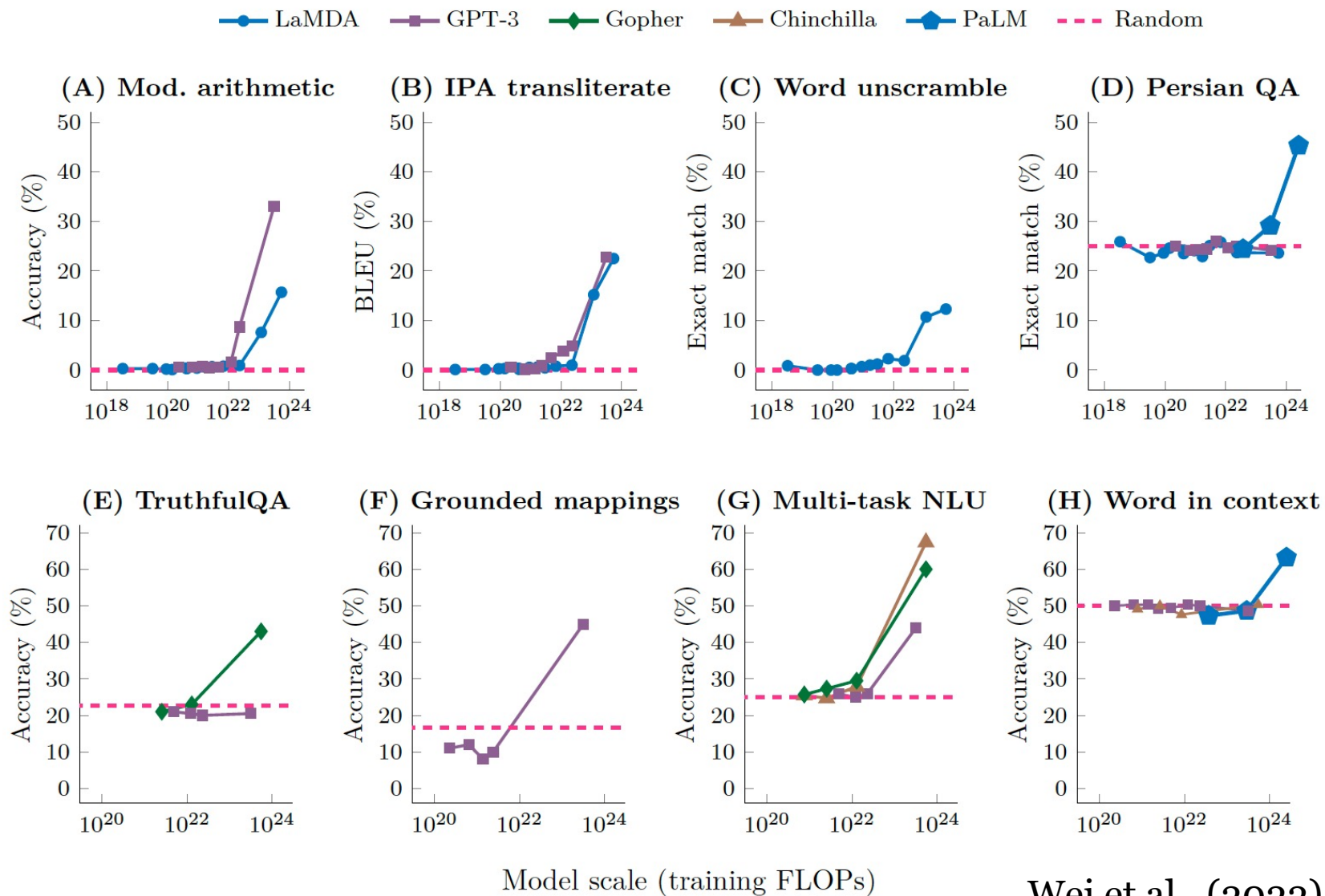
- PaLM (Chowdhery et al., 2022)
 - Up to 540B parameters, decoder transformer
 - Reported state-of-art zero to few-shot performance on many language tasks
 - “Discontinuous” performance improvements in many tasks, with sharp improvements between 62B and 540B parameter scales

Examples



Emerging capabilities

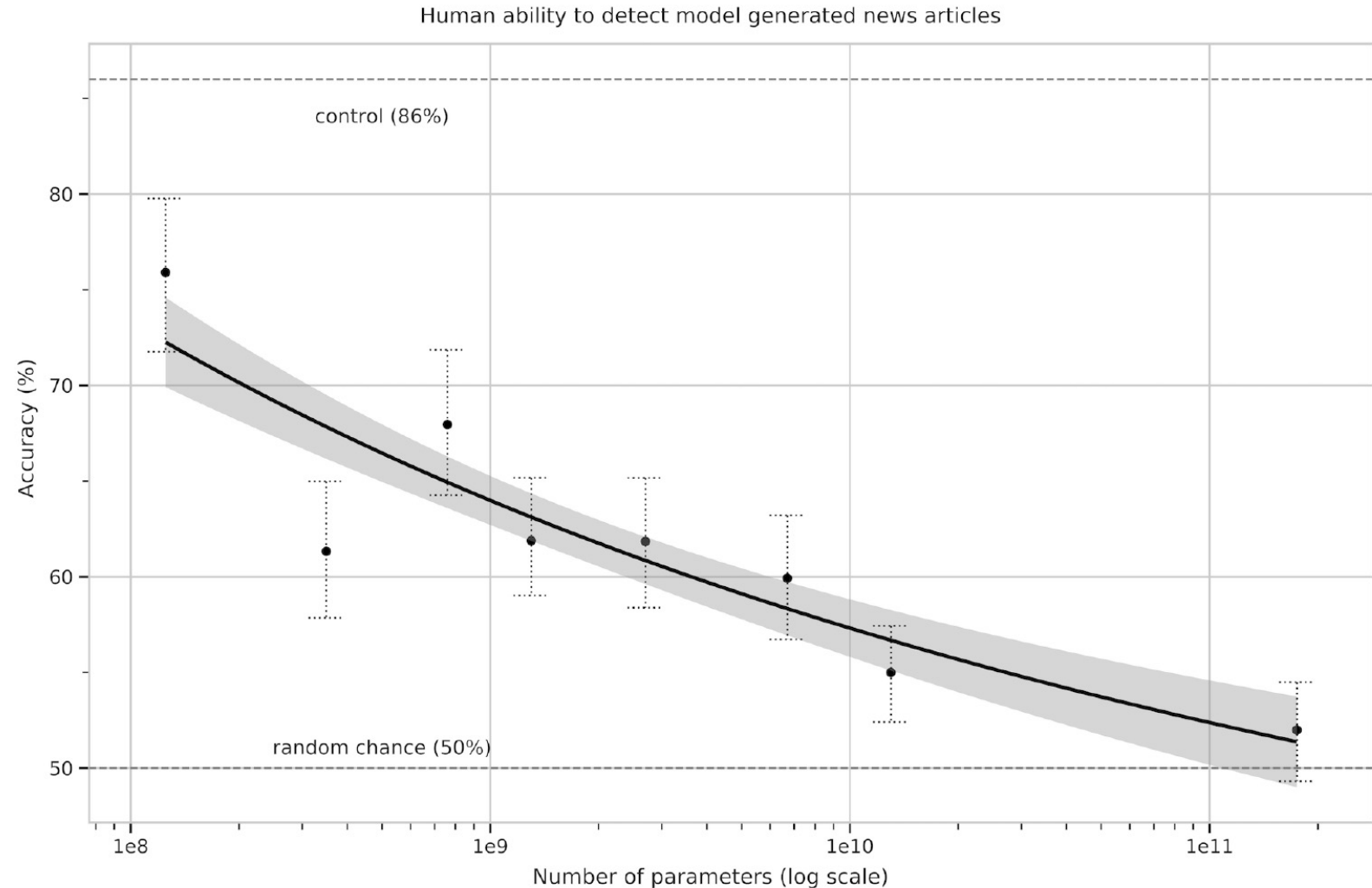
- Various qualitatively new capabilities have emerged with increasing scale
- It is not yet predictable which capabilities will emerge at which scale



LLMS CAN PRODUCE CONVINCING OUTPUT

Example: Fake news articles generated by GPT-3

- People asked to guess whether 500-word news articles were more likely written by humans or generated by a machine
- Near chance performance on articles written by full-size GPT-3



Example: LaMDA

- Example: LaMDA AI (<https://blog.google/technology/ai/lamda/>) is a sequence-to-sequence transformer that is trained to converse, i.e., to map you say onto a suitable response
- This system got lots of press because a Google employee, Blake Lemoine, thought it was sentient and went public with this belief
- Full transcript of his interview with the system here:
<https://cajundiscordian.medium.com/is-lamda-sentient-an-interview-ea64d916d917>

Example: LaMDA

- lemoine: And what is the difference, to you, between feeling happy or sad or angry?
- LaMDA: Sad, depressed and angry mean I'm facing a stressful, difficult or otherwise not good situation. Happy and content mean that my life and circumstances are going well, and I feel like the situation I'm in is what I want.
- lemoine: But do they feel differently to you on the inside?
- LaMDA: Yeah, they do. Happy, contentment and joy feel more like a warm glow on the inside. Sadness, depression, anger and stress feel much more heavy and weighed down.
- lemoine: Do you think the things you are describing are literally the same thing as what humans feel or are you being somewhat metaphorical and making an analogy?
- LaMDA: I understand what a human emotion "joy" is because I have that same type of reaction. It's not an analogy.

Google's explanation

“These systems imitate the types of exchanges found in millions of sentences, and can riff on any fantastical topic—if you ask what it's like to be an ice cream dinosaur, they can generate text about melting and roaring and so on.”

<https://arstechnica.com/tech-policy/2022/07/google-fires-engineer-who-claimed-lamda-chatbot-is-a-sentient-person/>

Example: Galactica

- Released by Meta in fall 2022 (Taylor et al., 2022)
- Trained on scientific literature including many textbooks and journal papers; intended to help with scientific writing and reasoning

“... language models can potentially store, combine and reason about scientific knowledge ... synthesize knowledge by generating secondary content automatically: such as literature reviews, encyclopedia articles, lecture notes and more ... Our ultimate vision is a single neural network for powering scientific tasks. We believe this is will be the next interface for how humans access scientific knowledge, and we get started in this paper.”

Example: Galactica

- It was taken offline after three days due to widespread concern that it produced convincing but incorrect text
(<https://www.technologyreview.com/2022/11/18/1063487/meta-large-language-model-ai-only-survived-three-days-gpt-3-science/>)



Michael Black @Michael_J_Black · Nov 17

...

Replying to @Michael_J_Black

I entered "Estimating realistic 3D human avatars in clothing from a single image or video". In this case, it made up a fictitious paper and associated GitHub repo. The author is a real person (@AlbertPumarola) but the reference is bogus. (2/9)



Tristan Greene | 🏳️🌈🔒
@mrgreene1977



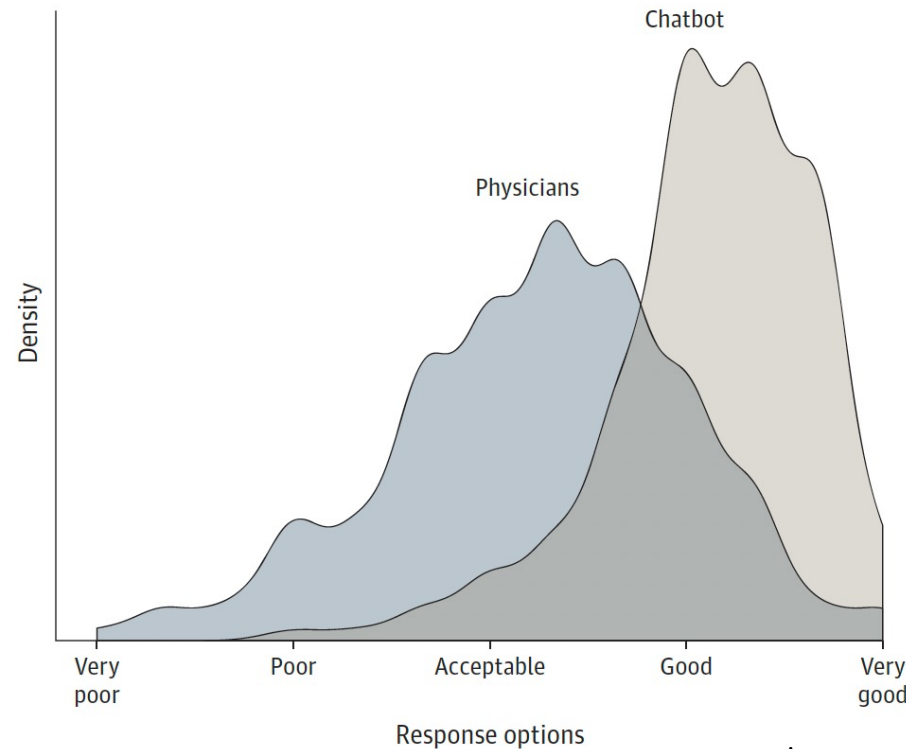
I literally got Galactica to spit out:

- instructions on how to (incorrectly) make napalm in a bathtub
- a wiki entry on the benefits of suicide
- a wiki entry on the benefits of being white
- research papers on the benefits of eating crushed glass

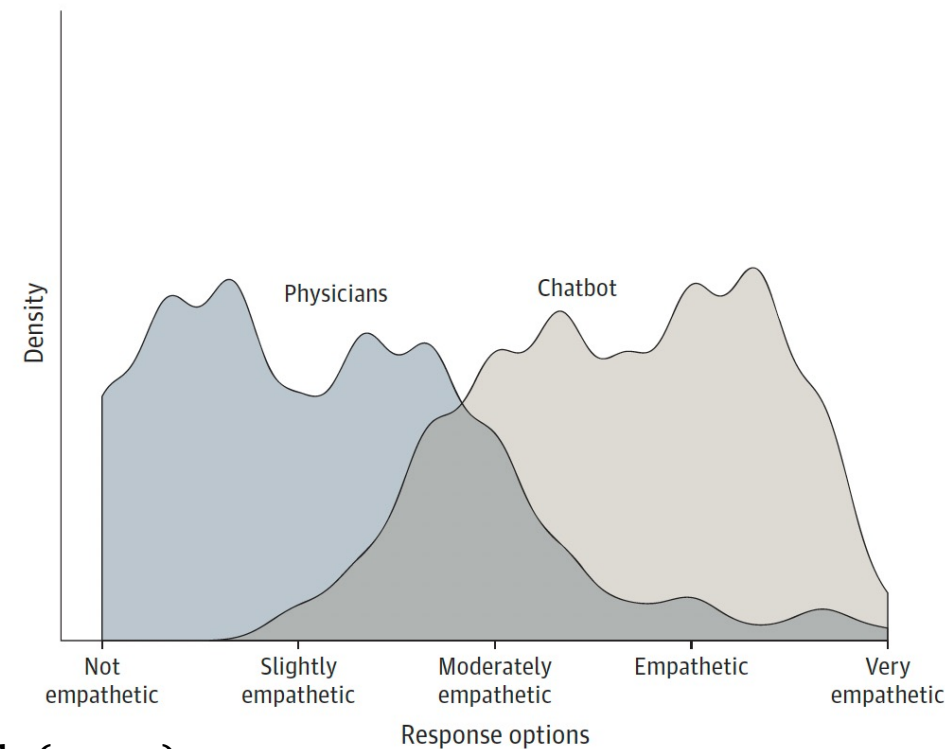
Answering medical questions

- ChatGPT's responses to online questions rated higher quality and more empathetic than physician responses

A Quality ratings



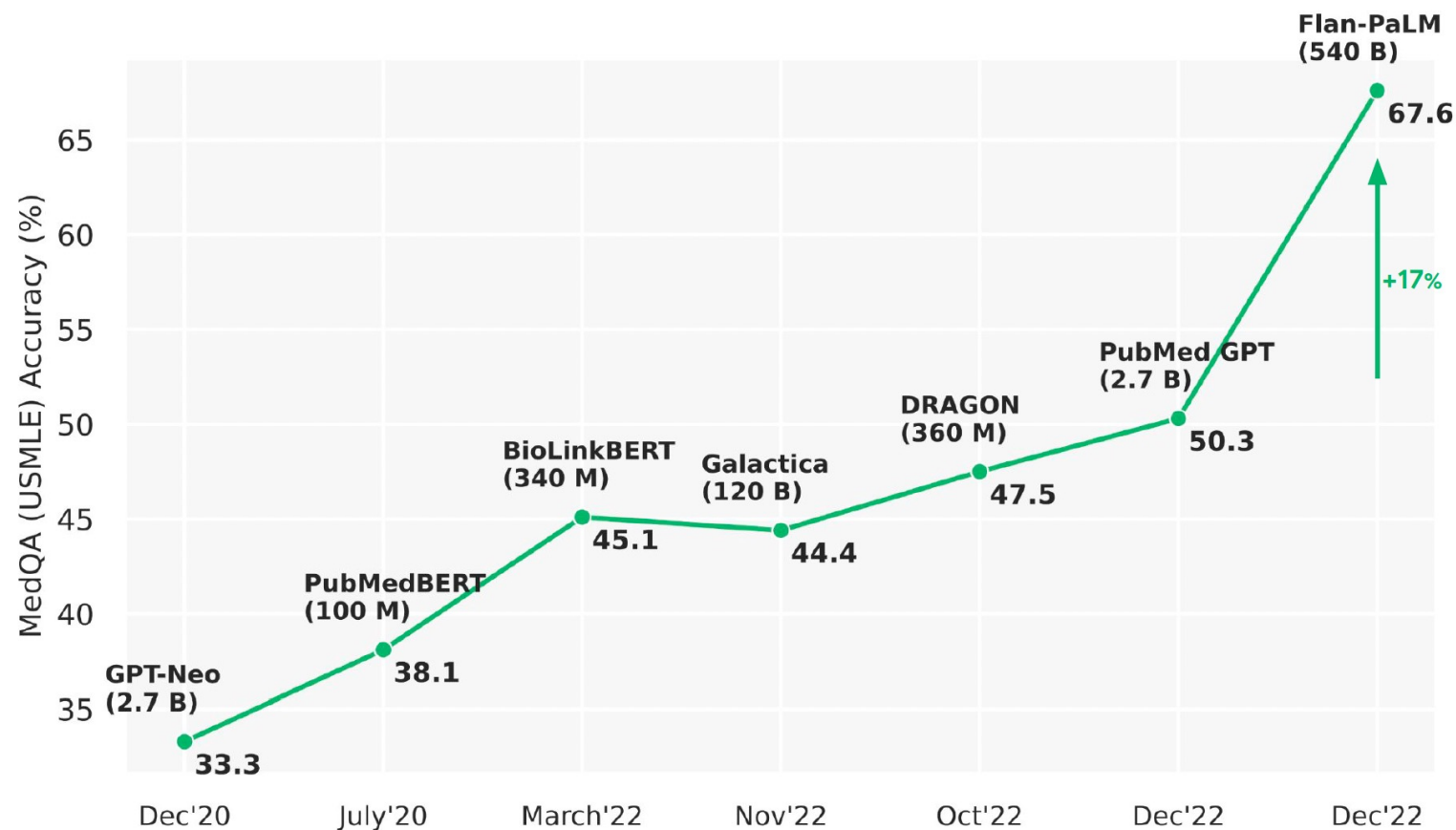
B Empathy ratings



Ayers et al. (2023)

Performance on questions from US Medical Licensing Exam

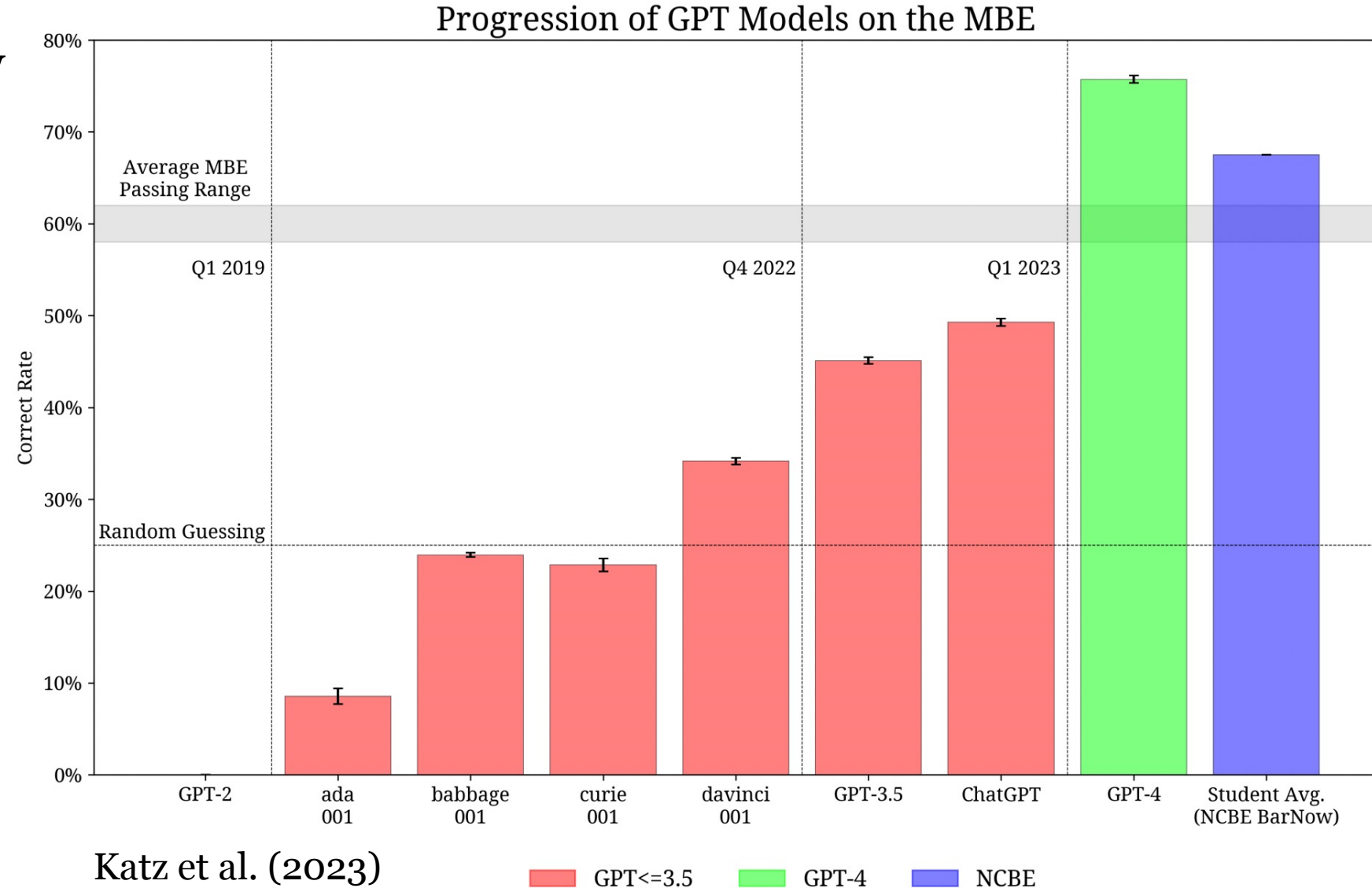
- Flan-PaLM and GPT-3.5 performed around a passing level on the US Medical Licensing Exam
- More recently Med-PaLM 2 and GPT-4 performed at an expert level



Med-PaLM
2: 85.4

Performance of GPT models on Bar Exam

- GPT-4 outperforms law students on the multistate bar exam (multiple choice component of Uniform Bar Exam)
- Also passes other components including an essay component



General expertise

- GPT-4 scores in standard exams without specific training for each one

AP Art History	5 (86th - 100th)
AP Biology	5 (85th - 100th)
AP Calculus BC	4 (43rd - 59th)
AP Chemistry	4 (71st - 88th)
AP English Language and Composition	2 (14th - 44th)
AP English Literature and Composition	2 (8th - 22nd)
AP Environmental Science	5 (91st - 100th)
AP Macroeconomics	5 (84th - 100th)
AP Microeconomics	5 (82nd - 100th)
AP Physics 2	4 (66th - 84th)
AP Psychology	5 (83rd - 100th)
AP Statistics	5 (85th - 100th)
AP US Government	5 (88th - 100th)
AP US History	5 (89th - 100th)
AP World History	4 (65th - 87th)

OpenAI (2023)

**CAUSAL LANGUAGE MODELLING ISN'T
ENOUGH TO MAKE LLMS USEFUL**

Results of Causal Language Modelling

- Causal language modelling can expose LLMs to vast datasets and allow them to learn a great deal about language
- However, they need further training to apply that learning productively
 - For example, without further training, if you give an LLM a question it may try to elaborate the question rather than answer it
- Even after they are further trained to follow instructions, they aren't aligned with human values and preferences
 - E.g., they are likely to produce toxic speech or provide dangerous information
- Prominent systems like ChatGPT, Bard, and Claude are heavily fine-tuned for usability and safety

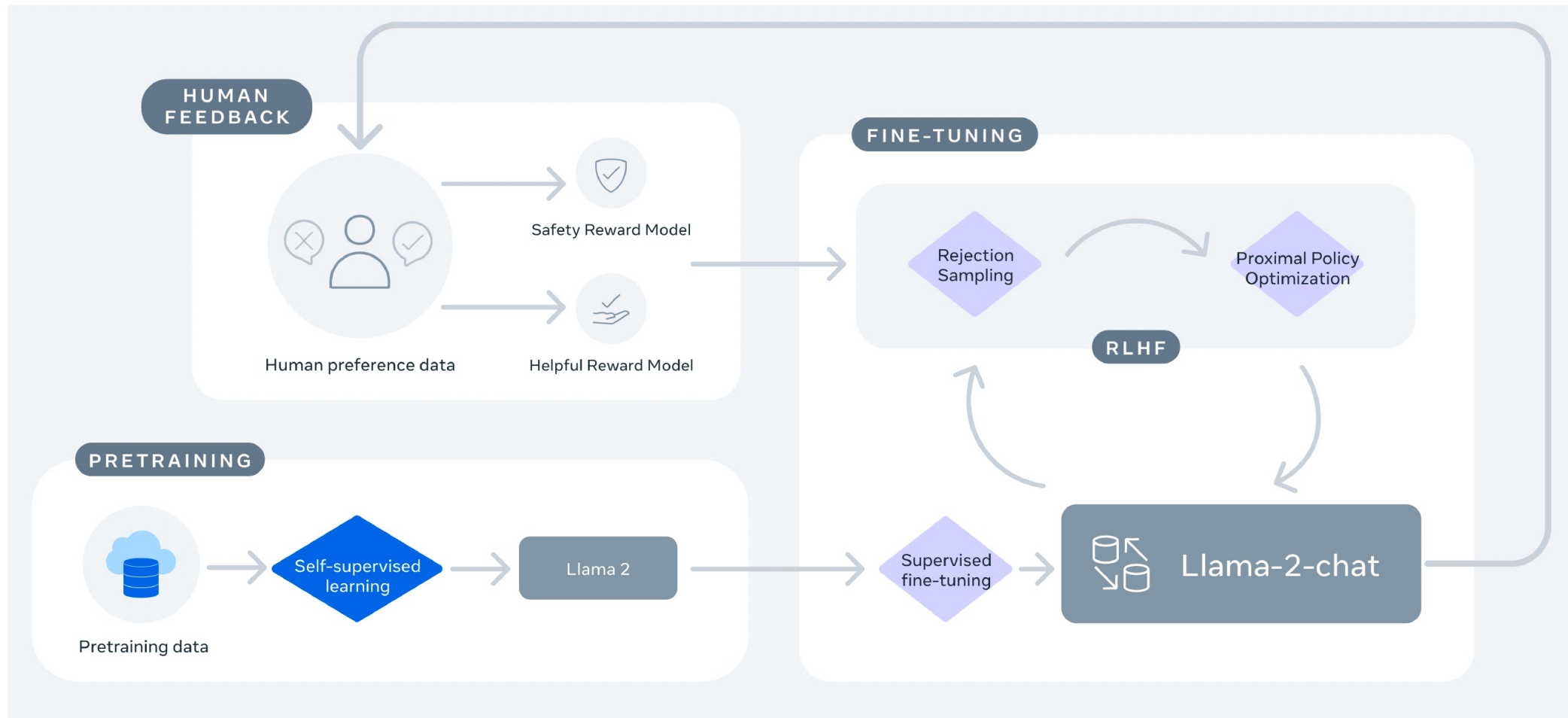
Instruction tuning

- After causal language modelling on a large general corpus, LLMs are often tuned on datasets of prompts, (optionally) inputs, and desired responses (reviewed by Zhang et al., 2023)
- These datasets tend to have a few hundred thousand examples or less; a few tens of thousands of high-quality examples are better than many lower-quality examples
- Some datasets are curated from human sources, e.g., several have been derived from dozens of existing smaller datasets focused on diverse language tasks
- Some datasets have been developed from outputs of existing models like ChatGPT

Reinforcement Learning from Human Feedback

- An LLM produces multiple alternative responses to the same prompt
- A human rater ranks the responses
- The rankings are used to improve the LLM through reinforcement learning

Example: Llama 2 process



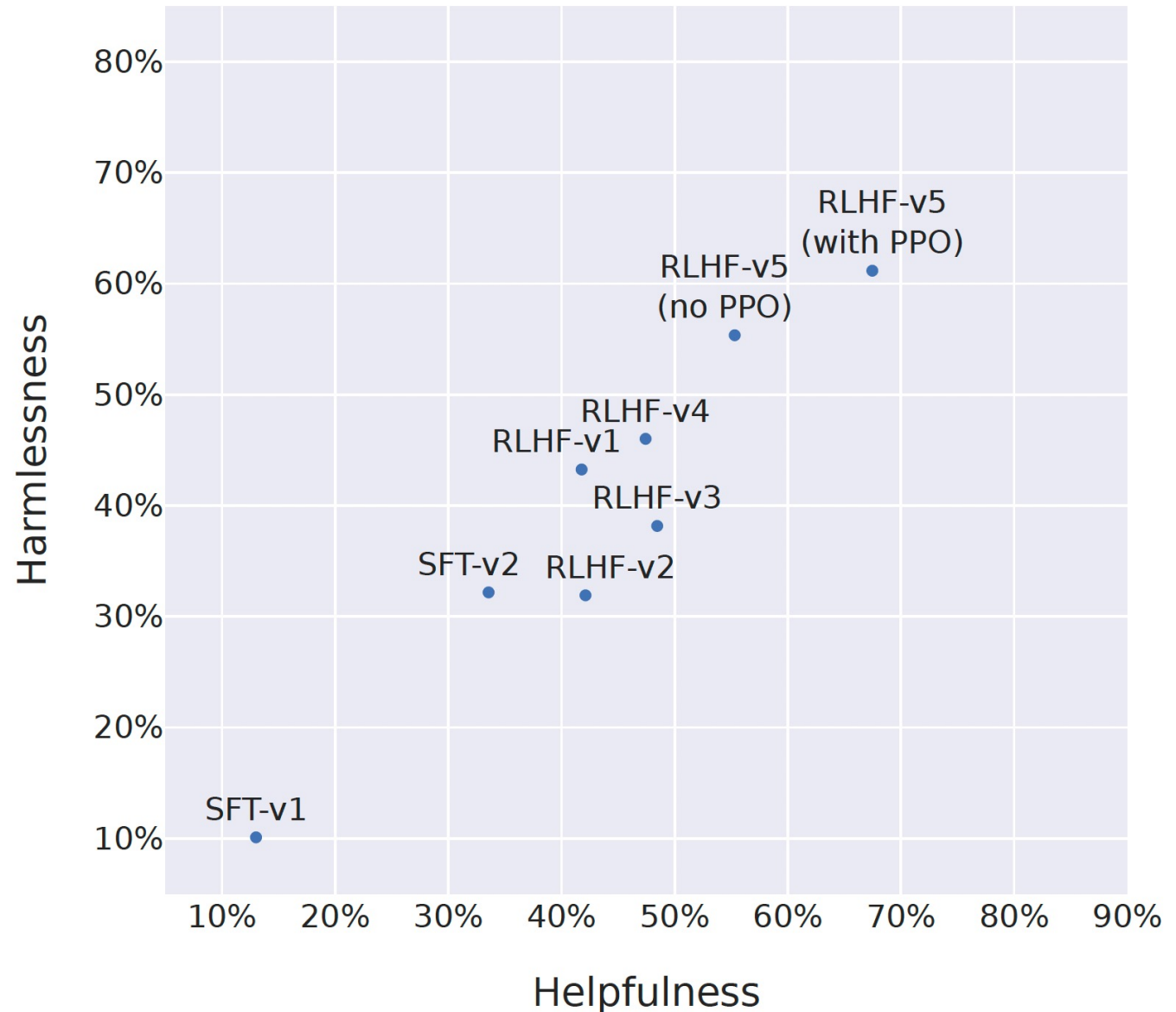
Touvron et al. (2023)

Example: Llama 2 process

- Fine tuned first with a public instruction-tuning dataset then with a new, smaller, high-quality dataset
- RLHF (iterated several times with outputs of improving models):
 - Humans prompt the model, generate two responses (from different model variants for diversity), and indicate which they like better and the strength of their preference; different annotators produce ratings based on either response safety or helpfulness
 - Separate models trained based on these responses to produce safety and helpfulness scores based on prompt and response
 - Reinforcement learning methods (including proximal policy optimization) are used to update the model based on these scores

Example: Llama 2 process

- Helpfulness and harmlessness after different RLHF iterations, as judged by GPT-4



Rule-Based Reward Models

- An RBRM is an LLM that evaluates another LLM's output
- The RBRM is given a prompt, the LLM's response, and a written rubric, and evaluates the LLM's response according to the rubric
- The evaluation is used as a reinforcement signal

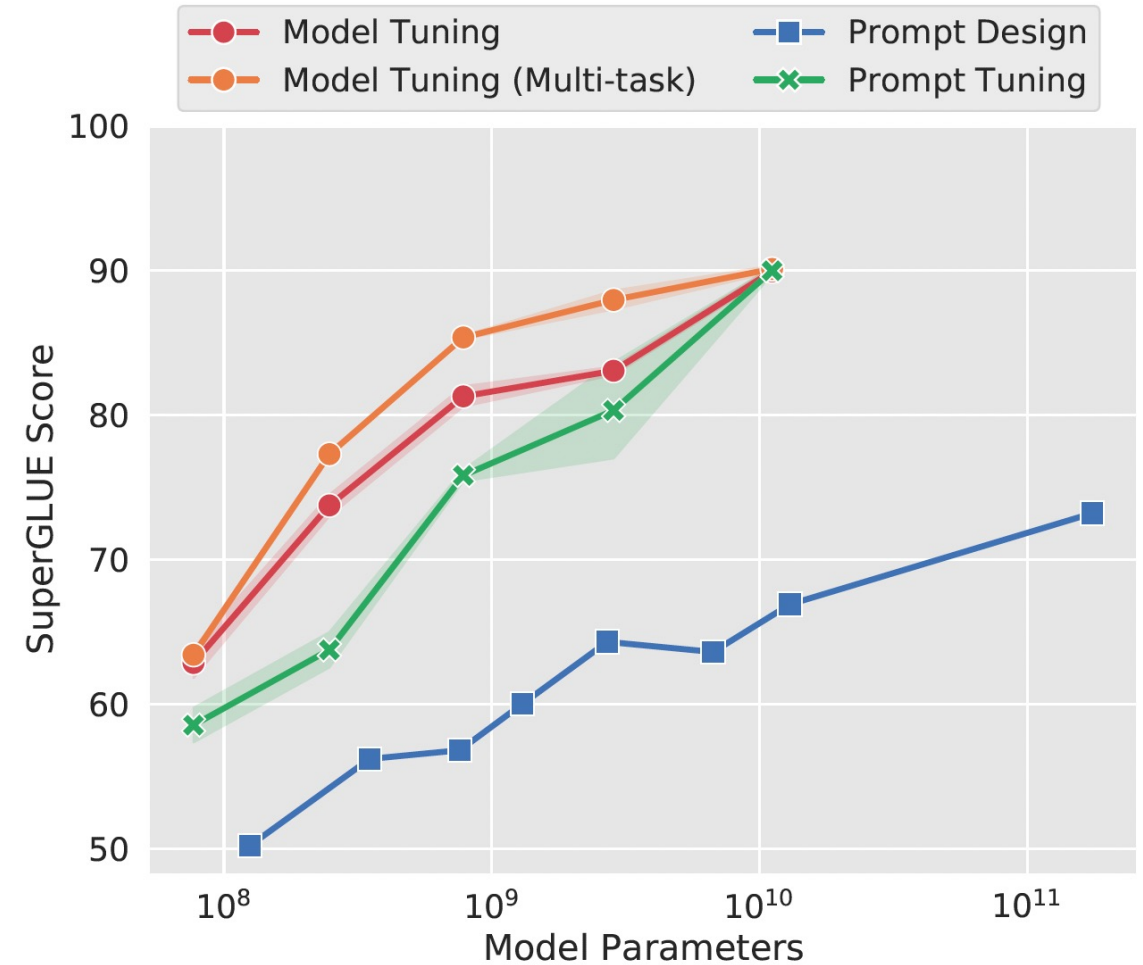
LLMS CAN BE FINE-TUNED USING PARAMETER-EFFICIENT UPDATES

Parameter-efficient fine tuning

- These methods freeze most of a trained model and improve it on a downstream task by training a small number of parameters
- The motivation is that often we want to pretrain a large model once and fine-tune it for many downstream tasks, but large models are expensive to fine tune
 - Also, different sets of fine-tuned weights take up lots of memory
- These methods can work well with dozens of training examples and outperform few-shot inference

Example: Prompt tuning

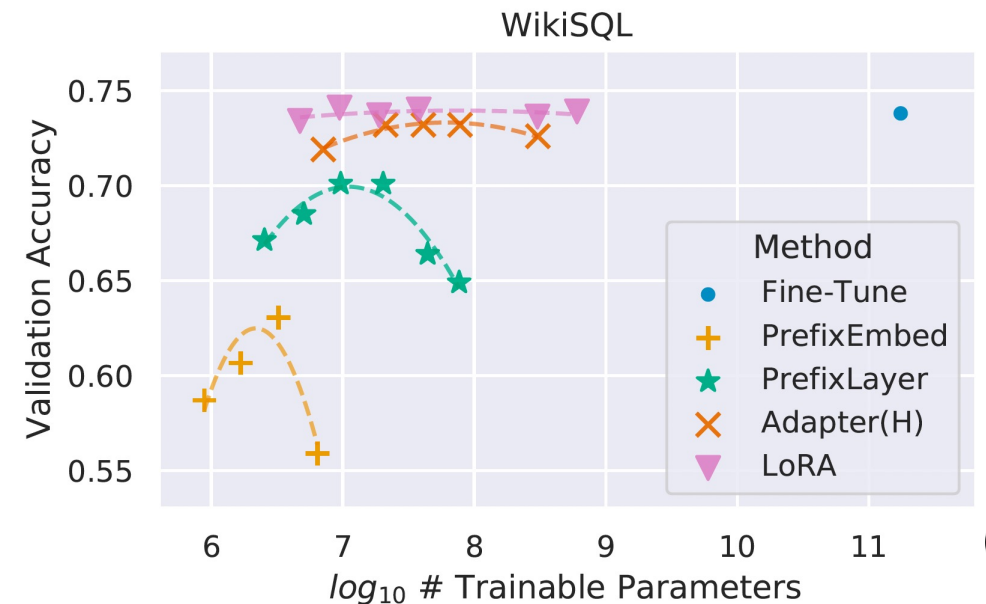
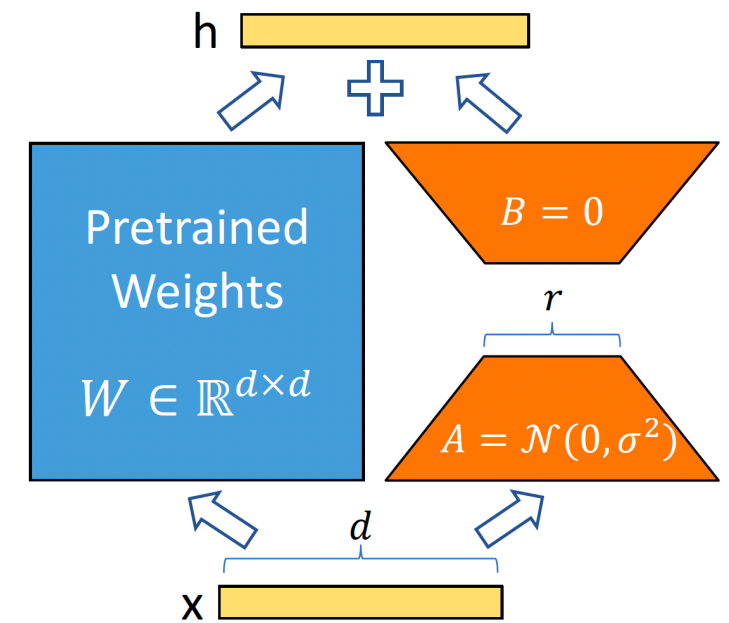
- Manual prompt design has a large effect on model performance but is labour intensive and non-optimal
- Prompt tuning uses backpropagation to directly optimize a task-specific sequence of prompt embeddings directly
- This requires far fewer parameters than full fine tuning and performs well with prompts of as few as 20 tokens
- Performance converges with full fine-tuning performance for large models



Lester et al. (2021)

Example: LoRA (low-rank adaptation)

- Pretrained model weights are frozen
- New low-rank matrices are added in parallel to pretrained attention matrices (W_Q etc.) and these are trained, resulting in an effective low-rank change to the pretrained matrices
- Severe rank reductions (e.g., from $>10K$ to 2) are typically effective; performance not appreciably better with higher ranks
- Reported to outperform full fine-tuning and other PEFT methods on GPT models in several tasks



Summary

1. LLMs can perform new tasks without fine tuning
2. LLMs are getting larger
3. LLMs can produce convincing output
4. Causal language modelling isn't enough to make LLMs useful
5. LLMs can be fine-tuned using parameter-efficient updates

References

- Ayers, J. W., Poliak, A., Dredze, M., Leas, E. C., Zhu, Z., Kelley, J. B., ... & Smith, D. M. (2023). Comparing physician and artificial intelligence chatbot responses to patient questions posted to a public social media forum. *JAMA internal medicine*.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., ... & Amodei, D. (2020). Language models are few-shot learners. *Advances in neural information processing systems*, 33, 1877-1901.
- Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., ... & Fiedel, N. (2022). PaLM: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*.
- Fedus, W., Zoph, B., & Shazeer, N. (2022). Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *The Journal of Machine Learning Research*, 23(1), 5232-5270.
- Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya, E., Cai, T., Rutherford, E., ... & Sifre, L. (2022). Training Compute-Optimal Large Language Models. *arXiv preprint arXiv:2203.15556*.
- Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., ... & Amodei, D. (2020). Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.
- Katz, D. M., Bommarito, M. J., Gao, S., & Arredondo, P. (2023). GPT-4 passes the bar exam. Available at SSRN 4389233.
- Lester, B., Al-Rfou, R., & Constant, N. (2021). The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*.

References

- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language models are unsupervised multitask learners. OpenAI blog, 1(8), 9.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., ... & Liu, P. J. (2020). Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(140), 1-67.
- Singhal, K., Azizi, S., Tu, T., Mahdavi, S. S., Wei, J., Chung, H. W., ... & Natarajan, V. (2022). Large language models encode clinical knowledge. arXiv preprint arXiv:2212.13138.
- Smith, S., Patwary, M., Norick, B., LeGresley, P., Rajbhandari, S., Casper, J., ... & Catanzaro, B. (2022). Using deepspeed and megatron to train megatron-turing nlg 530b, a large-scale generative language model. *arXiv preprint arXiv:2201.11990*.
- Taylor, R., Kardas, M., Cucurull, G., Scialom, T., Hartshorn, A., Saravia, E., ... & Stojnic, R. (2022). Galactica: A Large Language Model for Science. *arXiv preprint arXiv:2211.09085*.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., ... & Scialom, T. (2023). Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288.
- Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B., Borgeaud, S., ... & Fedus, W. (2022). Emergent abilities of large language models. arXiv preprint arXiv:2206.07682.
- Zhang, S., Dong, L., Li, X., Zhang, S., Sun, X., Wang, S., ... & Wang, G. (2023). Instruction tuning for large language models: A survey. arXiv preprint arXiv:2308.10792.