

Compare Results

2/9/2026 12:36:50 PM

Summary of Comments on A99bpjh4_1y3dhnw_mm4.tmp

This page contains no comments

Old File:		New File:	
2024_emnlp-main_293.pdf	versus	2024.emnlp-main.293.pdf	11 pages (11.41 MB)
17 pages (327 KB)			
2/8/2026 4:32:32 AM			

Total Changes	Content	Styling and Annotations
32		
6	Replacements	0 Styling
10	Insertions	0 Annotations
16	Deletions	

[Go to First Change \(page 1\)](#)

Academics Can Contribute to Domain-Specialized Language Models

Mark Dredze^{1,2}, Genta Indra Winata^{*}, Prabhanjan Kambadur¹, Shijie Wu^{*},

Ozan Irsøy¹, Steven Lu¹, Vadim Dabrowski¹,

David S Rosenberg¹, Sebastian Gehrmann¹

¹Bloomberg ²Johns Hopkins University ³Capital One

mdredze@bloomberg.net

Abstract

Commercially available models dominate academic leaderboards. While impressive, this has concentrated research on creating and adapting general-purpose models to improve NLP leaderboard standings for large language models. However, leaderboards collect many individual tasks and general-purpose models often underperform in specialized domains: domain-specific or adapted models yield superior results. This focus on large general-purpose models excludes most academics and draws attention away from areas where they can make important contributions. We advocate for a renewed focus on developing and evaluating domain- and task-specific models, and highlight the unique role of academics in this endeavor.

1 Introduction

Natural language processing (NLP) research has historically produced domain- and task-specific supervised models. The field has shifted course in the past few years, with a singular focus on general-purpose generative large language models (LLMs) that, rather than focusing on a single task, domain, do well across many tasks (Brown et al., 2020; Chowdhery et al., 2022; Workshop et al., 2022; Zhang et al., 2022; Touvron et al., 2023b). By training on massive amounts of data from many sources, these models can do well on extremely broad professional and linguistic examinations (Achim et al., 2023; Afni et al., 2023), college-level knowledge questions (Hendrycks et al., 2021; Lai et al., 2023) and collections of reasoning tasks (Suzgan et al., 2023).

While the trend to develop a single, general-purpose generative model is a net positive change that has resulted in impressive results, it has also slowed down progress in other areas of NLP. First, we are less focused on problems that cannot be

solved with a chat-like interface. Second, the best-performing LLMs are often commercial systems, which are sometimes opaque about training data, system architecture, and training details. Third, frequent model updates hinder reproducibility. The resources required to train large general language models naturally constrain research at large organizations, and researchers (or academics) outside of these organizations have become dependent on closed commercial systems, or open systems with limited transparency regarding their training data. This is partly reflected in broader AI trends: Zhang et al. (2021) found that roughly 30% of papers at AI conferences (including *CL) have a Fortune 500 tech affiliation. Increased resources contribute to the success of transformer-based LLMs (Vaswani et al., 2017), with available hardware (Hoover, 2021) and benchmarks (Deshpande et al., 2021) both playing a deciding role in what models end up being developed. By optimizing the average score across hundreds of shallow tasks, we are smoothing out any signal that would be gained from deeply engaging with individual tasks. Developing domain-specific models can help identify model and training choices that yield improvements on tasks within those domains.

In this paper, we argue for renewed attention to domain-specific models with rigorous and domain-expert informed evaluations. Because many academics are excluded from LLM development due to resource constraints, attention has been drawn away from research areas where academics can make the greatest contributions: deep dives on specific challenging problems. Thus, we propose several research questions to reorient the research community towards developing domain-specific models and applications, where academics are uniquely suited to lead.

^{*}The project was completed during work at Bloomberg.

2 LLMs: A Brief History

While modern LMs date back to Jeřínek (1976), we summarize very recent history to describe the current environment. In the wake of the popularization of neural word embeddings by word2vec (Mikolov et al., 2013), contextualized representations of language as features for supervised systems were realized by ELMs (Peters et al., 2018) followed by BERT (Devlin et al., 2019; Liu et al., 2019). BERT and subsequent models became the base models for supervised systems utilizing task-specific fine-tuning and continued pre-training for new domains (Gururangan et al., 2020), e.g., for clinical tasks ELMo (Schuster and Dredze, 2019) and clinicalBERT (Huang et al., 2019).

Parallel work utilized transformers for autoregressive LLMs, resulting in GPT (Radford et al., 2018), GPT-2 (Radford et al., 2019), BART (Lewis et al., 2020a; Liu et al., 2020), CTRL (Keskar et al., 2019), T5 (Raffel et al., 2020), Xue et al. (2021), and XGLM (Lin et al., 2021). These models had some few-shot capabilities, but they could each be adapted (fine-tuned) for specific tasks of interest. Some models were available to academics, though training a new model was beyond reach for many. GPT-3 (Brown et al., 2020) greatly increased model size and changed our understanding of LLMs. Impressive in-context (few-shot) learning pushed the idea that a single large model could solve a wide range of tasks. While the cost of resources meant training was restricted to a few groups, work focused on training bigger models (Chowdhery et al., 2022; Anil et al., 2023; Zhang et al., 2022; Touvron et al., 2023a; Rae et al., 2021).

While only a few could train large models, many studied how best to use them: prompt engineering (Liu et al., 2023), prompt tuning (Han et al., 2022; Wei et al., 2022), evaluation (Liang et al., 2022), among many other topics. Commercial LLM APIs, and eventually open source models (Zhang et al., 2022; Workshop et al., 2022; Touvron et al., 2023ab; Groeneweld et al., 2024), facilitated this work. Ignat et al. (2024) noted the massive research shift to LLMs reflected in Google Scholar citations. Subsequent work in instruction tuning (Oiyang et al., 2022) and fine-tuning (Wei et al., 2022; Zhang et al., 2022; Longpre et al., 2023) have further centralized research around general-purpose models. Many consider fine-tuning for specific applications to be obsolete: *“why would you tune a model for a specific task when you can tune*

a single model to do well on all tasks?

Despite this view, multiple domain-specific LLMs have demonstrated that outperform much larger models (Wu et al., 2023; Taylor et al., 2022). Med-Palm has shown that adapting even giant LLMs to a specific domain leads to vastly increased performance (Singhal et al., 2022, 2023). Furthermore, the release of LLaMA (Touvron et al., 2023a) led quickly to Alpaca (Faturi et al., 2023) and a wave of new fine-tuned versions of LLaMA for specific tasks. This trend strongly indicates that domain-specific models, especially for constrained sizes, are still highly relevant.

To be clear, our concern is not with closed models, which play an important role in the *product* ecosystem. Models range from full to limited to no access, with some closed models providing incredibly detailed information (Hoffman et al., 2022; Rae et al., 2019; Wu et al., 2023) and others providing none (Achham et al., 2023). Our lament over this focus on general models, either open or closed, is that it draws attention away from work on task- and domain-specific models and evaluations. Academics have become product testers, instead of focusing on tasks where they can play a unique role. Moreover, existing academic benchmarks increasingly serve a reduced purpose for commercial models: we are hill-climbing on benchmarks without a way to ensure existing LLMs have not been tailored to excel on these benchmarks (Dodge et al., 2021). Furthermore, we rely on benchmarks in place of deep engagement with an application and its stakeholders.

3 The Need for Domain-Specific LLMs

In general, web data does not reflect the needs of all NLP systems. Historically, the community has developed systems for specialized domains such as finance, law, bio-medicine, and science. Accordingly, there have been efforts to build LLMs for these domains (Wu et al., 2023; Taylor et al., 2022; Singhal et al., 2022; Bolton et al., 2023; Liao et al., 2022; Lehman et al., 2023; García-Ferrero et al., 2024). We need a deep investment in how best to develop and evaluate these models in partnership with domain experts. *How should we best integrate*

¹Distillation for task-specific models remains popular if smaller models are desired (Hsieh et al., 2023).

²We acknowledge that the biomedical domain is a rapidly developing area, and GPT-2 without fine-tuning was reported to surpass Medpalm 2 (Norri et al., 2023).

insights gained from the development of general-purpose models with those efforts? We propose several research directions.

How can general-purpose models inform domain-specific models? Building domain-specific models should benefit from insights and investments into general-purpose models. There are several strategies: training domain-specific models from scratch (Taylor et al., 2022; Bolton et al., 2023), mixing general and domain-specific data (Wu et al., 2023), and fine-tuning existing models (Singhal et al., 2022, 2023). Focusing on domain-specific needs, applications, and knowledge with guidance from topic experts will benefit us in acquiring a better model for specific NLP tasks. *Which approach yields the best results for task performance and overall cost?*

What is the role of in-context learning and fine-tuning? Both LIMA (Zhou et al., 2023) and Med-PaLM (Singhal et al., 2022) use a small number of examples to tune a model. With expanding context size, we may soon rely entirely on in-context learning (Pettori et al., 2020). This blurs the lines between changing model parameters and conditioning during inference. Beyond inference speed tradeoffs between the two, there may be value in tuning on tens of thousands (or more) of examples. *Which domain-specific examples are the most effective to include and in what manner?*

How can LLMs be integrated with domain-specific knowledge? Specialized knowledge is key in many domains. RAG (Lewis et al., 2020b; Guu et al., 2020) and KILT-derived work (Pettori et al., 2021) focus on knowledge-intensive tasks by including retrieval steps. Work on attributional QA (Behne et al., 2022) takes a similar approach, as do search LLMs that rephrase interaction with retrieved data (Nakano et al., 2021). Rich updated knowledge sources will always exist beyond the model, especially in environments like medicine, finance, and many academic disciplines.

4 Evaluation of Domain-Specific Models

The evaluation of NLP systems is at a crossroads, and the downstream usage of LLMs and evaluation approaches have diverged. Benchmarks assume that their results translate to insights into similar tasks and usefulness for commercial applications. But benchmarks have become increasingly narrow about broken evaluations (Gehrmann et al., 2022).

It is dubious whether we gain insights into non-task-specific generation through NLU benchmarks.

If we are performing the depth-first evaluation of a generation task, a remaining hurdle – and why researchers fall back to NLU tasks – is the lack of robust metrics. While there is much recent work on

better metrics (Celiyilmaz et al., 2020; Gehrmann et al., 2023), a troubling trend is the use of LLMs as evaluators (e.g., Sellam et al., 2020; Chiang et al., 2023). This approach poses many risks, including the implicit assumption that the evaluating model has access to the ground truth judgment. While there are some promising results, using an LLM out of the box should be avoided (e.g., Wang et al., 2023a,b). Moreover, it is unclear how to evaluate the evaluator when it is a non-deterministic API, or how to scale the development of learned metrics and quantify the strength of a metric.

Products are not Baselines If we really do want to evaluate 100+ tasks, there are many issues with the soundness of evaluation setups. At this scope, it is impossible to run careful ablation studies or to assess the effect of changes to methodology in a causal manner. Moreover, different LLMs respond differently to prompts. The BLOOM evaluation averaged over multiple prompts and found significant variance (Workshop et al., 2022). This variance leads to a lack of reproducibility: LLaMA (Touvron et al., 2023a) claimed high MMLU (Hendrycks et al., 2021) performance but didn't release the prompts that led to them.³ Similarly, the evaluation scheme makes a difference (Liang et al., 2022, Fig. 33). High evaluation costs mean benchmarks pick a small number of setups (sometimes only one) for each task, which introduces further bias, making it hard to construct fair benchmarks on many tasks.

An additional issue with the current benchmarking approach is that the best-performing models are often commercial APIs. With limited transparency regarding data and training, we cannot fairly evaluate these models (e.g., data leakage). Furthermore, task-specific tuning may have been selected based on these specific benchmarks. Moreover, the underlying models change frequently, so it is unclear whether a result will hold for long.

These evaluation issues prompt significant open questions: 1) How do we develop consistent evaluation setups across models that give true measures of performance? 2) How do we develop evaluation setups and metrics more closely aligned with downstream usage? 3) How do we develop evaluation suites that support depth-first evaluation and not breadth-first benchmarking?

³There was significant confusion surrounding model evaluation: <https://huggingface.co/biobg/open-lla-leaderboard-mmlu>

5 The Role of Academics

A focus on general-purpose LLMs has forced academics to work with large base models and perhaps shifted the focus to solve problems of immediate industrial interest. Many academics feel excluded from current research trends (Ignat et al., 2022) and the academic and industry relationship is changing (Litman et al., 2022). Shifting attention back to domain-specific applications emphasizes areas where academics hold an advantage: partnerships with domain experts to invest in specific tasks, and consideration of broader societal needs.

Developing domain-specific models requires domain expertise and universities are diverse academic environments that house experts in many domains. Collaborations with these experts can identify data sources, tests, and challenges important within each domain. Furthermore, these collaborations are the best avenues for better alignment of evaluations with use cases (Winita et al., 2024), and can support the development of proper metrics. These collaborations are necessary to explore wide open interdisciplinary topics, such as models for protein structure prediction (Tanyasuvankool et al., 2021; Vi et al., 2021) and games as proxies for reasoning (Silver et al., 2016; Agostinelli et al., 2019; Schriftwieser et al., 2020). This includes developing domain-specific resources, which require domain experts to properly design and construct the datasets. Further, areas where industry underinvests are those where academics could focus attention. For example, low-resource languages are not served by a general-purpose multilingual LLM, nor will we reasonably have enough data to support current LLM training methods. Dialects and variations in languages are still wide open topics (Aji et al., 2022; Winita et al., 2023; Nichols and Batista, 2023).

General-purpose LLMs are unlikely to solve problems in many important domains, with many open research problems that can only be solved by domain-specific approaches. Focusing on domain-specific knowledge will benefit us in acquiring a better model and developing application strategies more aligned with how humans learn domain-specific knowledge (Tricot and Sweller, 2014). For many interdisciplinary areas, subject matter experts are essential, and the problems must be defined clearly. The first pass from an LLM is often impressive, but it hides the trenches, and areas where things are most interesting. We need a renewed fo-

cus on developing and evaluating domain-specific models and applications, an area where academics can play a leading role. Let us not be disarmed by claims that a single model solves all tasks, and instead deeply explore and understand the needs and challenges of specific domains.

Limitations

The literature that we explored in this opinion paper is limited to the area of LLMs. We study the history of LLMs from the literature on word embeddings, encoder-only, and generative transformers to the latest advancement of API-based LLMs.

Ethics Statement

Our work does not include any experiments or use of data. No potential ethical issues in this work.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lana Ahmad, Ilge Akaya, Florencia Leon Alvarado, Diogo Athieza, Janko Altenberndt, Sam Altman, Shyamal Anand et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Forest Agostinelli, Stephen McAleer, Alexander Shimakov, and Pierre Baldi. 2019. Solving the Rubik's Cube with deep reinforcement learning and search. *Nature Machine Intelligence*, 1(8):356–353.
- Alihan Aji, Genia Indra Winaata, Fajri Koto, Samuel Cahyawijaya, Ade Romadhony, Rahmad Muhendra, Kemal Kurniawan, David Melcick, Radityo Eko Prasolo, Timothy Baldwin, et al. 2022. One country, 700+ languages: NLP challenges for underrepresented languages and dialects in indonesia. In *Proceedings of the 6th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7226–7249.
- Rohan Aul, Andrew M. Dai, Orhan Erat, Melvin Johnson, Dmitry Lenkin, Aleksander Pasos, Stananak Shakeri, Emma Taylor, Varga Bailey, Zhifeng Chen, et al. 2023. Palm: technical report. *arXiv preprint arXiv:2305.10493*.
- John W. Ayers, Adam Poliak, Mark Dredze, Eric C. Leas, Zechariah Zhu, Jessica B. Kelley, Dennis C. Fair, Aaron Goodman, Christopher A. Longhurst, Michael Hozarth, et al. 2023. Comparing physician and artificial intelligence chatbot responses to patient questions posed to a public social media forum. *JAMA Internal Medicine*, 183(6):589–596.
- Andrew Blair-Starek, Nils Holzenberger, and Benjamin Van Durme. 2023a. Can gpt-3 perform statutory reasoning? In *Proceedings of the Nineteenth International Conference on Artificial Intelligence and Law*, pages 22–31.
- Andrew Blair-Starek, Nils Holzenberger, and Benjamin Van Durme. 2023b. Openai cribbed our tax example, but can gpt-4 really do tax? *arXiv preprint arXiv:2308.09992*.
- Bernd Bohnet, Vinh Q. Tran, Fatima Verga, Rose Abramoff, Daniel Andr, Lirio Baldini Soares, Jacob Eisenstein, Kuzman Ganchev, Jonathan Heitzig, Kaifu Liu, Tom Kwiatkowski, Ji Ma, Jiamou Ni, Tal Schuster, William W. Cohen, Michael Collins, Dipanjan Das, Donald Metzler, Slav Petrov, and Kellie Webster. 2022. Attributed question answering: Evaluation and modeling for attributed large language models. *CoRR*, abs/2212.38377.
- Elliot Bolton, David Hall, Michihiro Yasunaga, Tony Lee, Chris Manning, and Percy Liang. 2023. BioMedLM. <https://github.com/stanford-crfn/BioMedLM>.
- Ton Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D. Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901.
- Ashi Celikyilmaz, Elizabeth Clark, and Jianfeng Gao. 2020. Evaluation of text generation: A survey. *arXiv*, 2006.14799.
- Hanjie Chen, Zhouzhang Fang, Yash Singla, and Mark Dreize. 2024. Benchmarking large language models on answering and explaining chatbot-impressing medical questions. *arXiv preprint arXiv:2402.18060*.
- Wei-Lin Chiang, Zhuohan Li, Zi Lu, Ying Sheng, Zhiqiao Wu, Jiao Zhang, Liannan Zheng, Siyuan Huang, Yonghao Zhuang, Joseph F. Gonzalez, et al. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%+ chatbot quality. See <https://vicuna.in>. *tiny.cc/meyarw* (accessed 14 April 2023), 2(3):16.
- Aakanksha Chowdhery, Sharad Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hung Wong Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tyvyshevko, Iosifina Meynez, Abhishek Rao, Parker Barnes, Yi Fey, Noam Shazeer, Vineet Cadkumar, Prabhakaran, Enrico Rei, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengchen Yin, Taji Duce, Anselm Levyay, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denyy Zhao, Daphne Apolito, David Luan, Hyeonae Lim, Barret Zoph, Alexander Sridharan, Ryan Sepassi, David Doban, Shivani Arsalan, Mark Omernick, Andrew McDa, Thennarasu Sankaranarayara Pillai, Maril Pellar, Attila Lehtovszky, Irina Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Gaeta, Mark Diaz, Oshan Firat, Michele Cutatta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. Palr: Scaling language modeling with pathways. *CoRR*, abs/2204.02511.

- Hyung Won Chung, Le Hou, Shayne Logue, Barrett Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Bhattacharya, et al. 2022. Scaling instruction-finetuned language models. *arXiv*, 22(10.1):416.
- Mostafa Dehghani, Yi Tay, Alexey A. Grishchenko, Zhe Zhao, Neil Houlsby, Fernando Diaz, Donald Metzler, and Oriol Vinyals. 2021. The benchmark lottery. *CoRR*, abs/2107.07002.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Kausubh D. Dhole, Varun Gangal, Sebastian Gehrmann, Adesh Gupta, Zhenhai Li, Shad Mahmood, Abinaya Mahendiran, Simon Mills, Ashish Srivastava, Samson Tan, Tongshang Wu, Jashica Sohn-Dickstein, Jinbo D. Choi, Edward H. Hovy, Ondrej Dušek, Sebastian Ruder, Sajan Arand, Nagaendra Anuja, Robin Banjade, Lisa Barthe, Hema Behme, Ian Bentor-Artveli, Connor Boyce, Caroline Burn, Marco Antonio Sobrevilla Cabedo, Samuel Cahayawya, Emilie Chapuis, Waxiang Che, Mukund Choudhary, Christian Clauß, Pierre Colombo, Filip Connell, Gautier Dagar, Mayukh Das, Tanay Dixit, Thomas Dopiere, Paul Dray, Suchitra Dubey, Tatiana Ekinbör, Marco Di Giovanni, Rishabh Gupta, Rushabh Gupta, Louises Hanila, Sang Han, Fabrice Haré-Canaïa, Antoine Honore, Ishaan Jindal, Przemysław K. Jonak, Denis Klejko, Venelin Kovatchev, and et al. 2021. Ni-augmenter: A framework for task-sensitive natural language augmentation. *CoRR*, abs/2112.02721.
- Jesse Dodge, Maarten Sap, Andriy Marasović, William Agnew, Gabriel Ilharco, Dirk Groeneweld, Mergen Mitchell, and Matt Gardner. 2021. Documenting large webtext corpora: A case study on the colossus clean crawled corpus. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1286–1305.
- Alexander V. Erkisen, Steven Müller, and Jesper Ryg. 2023. Use of sp-4 to diagnose complex clinical cases.
- Kawin Ethayarajh and Dan Jurafsky. 2020. Utility in the eye of the user: A critique of NL2 leaderboard. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4846–4853, Online. Association for Computational Linguistics.
- Alexander R. Fabbri, Wojciech Krzysztofik, Bryan McCan, Caiming Xiong, Richard Socher, and Dragomir Radev. 2021. SummEval: Re-evaluating summarization evaluation. *Transactions of the Association for Computational Linguistics*, 9:391–400.
- Iker García-Ferreiro, Rodríguez Agueri, Alzíbar Asturias, Salazar, Elena Cabré, Iker de la Iglesia, Ainhoa Lavelle, Bernardo Magnini, Benjamin Molinet, Joana Medina, and Tomás Ronero. 2024. Media inMT5: an open-source multilingual text-to-text LLM for the medical domain. *arXiv preprint arXiv:2404.07653*.
- Sébastien Gehrmann, Elizabeth Clark, and Thibault Selam. 2023. Repaining the cracked foundation: A survey of obstacles in evaluation practices for generated text. *Journal of Artificial Intelligence Research*, 77:103–166.
- Tanya Goyal, Junyi Jessy Li, and Greg Durrett. 2022. News summarization and evaluation in the era of GPT-3. *CoRR*, abs/2209.12356.
- Dirk Groeneweld, Iz Bellegu, Pete Welsh, Alkisita Blagia, Rodney Kinney, Oymrd Tafjord, Ananya Jaiswal, Jiaji Hanish, Ivusan, Ian Magnusson, Yizhoye Wang, et al. 2024. Otnro: Accelerating the scipice of language models. *arXiv preprint arXiv:2402.00828*.
- Shuchin Guan, Anand Mehta, Swabha Swayamdipta, Kyle Lo, Iz Bezugay, Doug Dewey, and Noam A. Smith. 2020. Don't stop pretraining: Adaptive language models for domains and tasks. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, pages 8842–8850.
- Kelin Guan, Kepon Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. 2020. Retrieval augmented language model pre-training. In *Proceedings of the 37th International Conference on Machine Learning (ICML) 2020, 13–18 July 2020, Virtual Event, volume 97 of Proceedings of Machine Learning Research*, pages 3929–3938. PMLR.
- Tessa Han, Aounon Kumar, Chirag Agarwal, and Himabindu Lakkaraju. 2024. Towards safe large language models for medicine. In *ICML 2024 Workshop on Models of Human Feedback for AI Alignment*.
- Xu Han, Weilin Zhao, Ning Ding, Zhiyuan Liu, and Maosong Sun. 2022. PTR: Prompt tuning with rules for text classification. *AI Open*, 3:182–192.
- Dan Hendrycks, Colin Burns, Steven Basart, Andy Zoo, Manas Mezaika, Daws Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. In *International Conference on Learning Representations*.
- Jordan Hoffmann, Sébastien Bergendau, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutterford, Diego de Las Casas, Lisa Anne Heuericks, Johannes Weßl, Alan Clark, et al. 2022. An empirical analysis of compute-optimal large language model training. *Advances in Neural Information Processing Systems*, 35:30016–30300.
- Sara Hooker. 2021. The ‘harder’ lottery. *Commun. ACM*, 64(12):58–65.

- Cheng-Yu Hsieh, Chun-Liang Li, Chih-cuan Yeh, Hoorat Nahost, Jashisa Fujii, Alex Ratner, Ranjan Krishna, Chen-Yu Lee, and Tomas Pfister. 2023. Distilling step-by-step outperforming large-language models with less training data and smaller model sizes. In *Finding of the Association for Computational Linguistics: ACL 2023*, pages 8003–8017.
- Xinxi Huang, Juan Alloza, and Rajesh Ranganathan. 2019. ClinicalBERT: Modeling clinical notes and predicting hospital readmission. *arXiv preprint arXiv:1906.05342*.
- Oana Ignat, Zhitong Jin, Artem Abyzov, Laura Bieser, Santiago Castro, Naihao Deng, Xiyi Gao, Aylin Ece Gunal, Jacy He, Ashkan Kazemi, et al. 2024. Has it all been solved? open nlp research questions not solved by large language models. In *Proceedings of the 2024 Inist International Conference on Computational Linguistics: Language Resources and Evaluation (LREC-COLING 2024)*, pages 8056–8094.
- Frederick Jelinek. 1976. Continuous speech recognition by statistical methods. *Proceedings of the IEEE*, 64(4):532–556.
- Nitish Shirish Keskar, Bryan McCann, Lav R. Varshney, Caiming Xiong, and Richard Socher. 2019. Critic: A conditional transformer language model for controllable generation. *arXiv preprint arXiv:1909.03636*.
- Douwe Kiela, Max Bartolo, Yixin Nie, Divyanshu Kaushik, Aticus Geiger, Zhengyan Wu, Beritc Vidgen, Grusha Pasci, Amarnpreet Singh, Pratik Ringsha, Zhiyi Ma, Tristan Thrush, Sebastian Riedel, Zeerak Waheed, Pontus Stenetorp, Robin Jia, Mohit Bansal, Christopher Potts, and Athira Willians. 2021. Dynabanc: Rethinking benchmarking in NLP. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1*, pages 4110–4124. Online Association for Computational Linguistics.
- Eun-Ah Kim, Haiping Pan, Nayanara Mukherjee, William Taranto, Subhashini Venugopalan, Rajzaman Bajwa, and Michael Brenner. 2024. Perfuring Hartree-Fock many-body physics calculations with large language models. *Bulletin of the American Physical Society*, 39:318–319.
- Viet Lai, Chien Nguyen, Nghia Ngo, Thau Nguyen, Franck Demontort, Zian Rossi, and Thien Nguyen. 2023. Okapi: instruction-tuned large language models in multiple languages with reinforcement learning from human feedback. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 318–320.
- Eric Lehman, Evan Hernandez, Divakar Mahajan, Jonas Wulf, Michal Smith, Zachary Ziegler, Daniel Nadler, Peter Szolovits, Alisair Johnson, and Emily Alsentzer. 2023. Do we still need clinical language models? In *Conference on health, inference, and learning*, pages 578–597. PMLR.
- Mike Lewis, Yinhan Liu, Naman Goyal, Jingfei Du, Veselin Stoyanov, and Luke Zettlemoyer. 2019. Luke-Zettlemoyer, and Veselin Stoyanov. 2019.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Charifizadeh, Abdolreza Mchanei, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020a. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880.
- Patrick S.-H. Lewis, Ethan Perez, Aleksandra Plakutis, Fabio Pericoli, Vladimir Karukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sébastien Riedel, and Douwe Kiela. 2020b. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *Advances in Neural Information Processing Systems*.
- Percy Liang, Rishi Bonmassari, Tony Lee, Dimitris Tsipras, Dileta Soylu, Mihitho Yilmazoglu, Yuan Zhang, Deepak Narayanan, Yuhua Wu, Ananya Kumar, Benjamin Newmark, Binhang Yuan, Boqun Yan, Ce Zhang, Christian Cosgrove, Christopher D. Manning, Christopher Ré, Diana Acosa, Xavas, Drew A. Hudson, Eric Zalensky, Esteban Eslava, Peisai Sodha, Frieder Kong, Hongyu Ren, Huiyu Yao, Jie Wang, Keshav Sathyanarayanan, Laurel J. Yu, Lucia Zheng, Merit Yükselkemal, Murat Sucuoglu, Nathan Kin, Neel Guha, Nihadri S. Chatzoglou, Omar Khatib, Peter Henderson, Qian Huang, Ryan Sang, Michael Xie, Shubham Sanujkar, Surya Ganguli, Taeunori Hwang, Thomas Head, Tianyi Zhang, Vishnuv Chaudhary, William Wang, Xuelian Li, Yifan Mai, Yuhui Zhang, and Yutai Koreeda. 2022. Holistic evaluation of language models. *CoRR*, abs/2211.09110.
- Xi Wenzhao Lin, Tocor Mhatslov, Mikel Artetxe, Tianfu Yang, Shuchii Chen, Daniel Simig, Myle Ott, Naman Goyal, Shruti Bhosale, Jingfei Du, et al. 2021. Few-shot learning with multilingual language models. *arXiv preprint arXiv:2112.10668*.
- Michael L. Litman, Ifeoma Ajunwa, Guy Berger, Craig Eouliier, Morgan Currie, Finale Doshi-Velez, Gillian Hadfield, Michael C. Horowitz, Charles Isbell, Hiroaki Kitano, et al. 2022. Gathering strength, gathering storms: The one hundred year study on artificial intelligence (ai100) 2022 study panel report. *arXiv preprint arXiv:2210.15757*.
- Pengfei Liu, Weiqi Yuan, Junlan Fu, Zhengbo Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9):1–35.
- Yinhan Liu, Jiaao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual denoising pre-training for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:726–742.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Marjan Ghazvininejad, Mike Lewis, Omer Levy, Daniel Joshi, Danqi Chen, Caiming Xiong, Veselin Stoyanov, and Luke Zettlemoyer. 2019.

- Roberta: A robustly optimized BERT pretraining approach. *arXiv*.
- Shayne Longpre, Le Hou, Tu Yu, Albert Webster, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V. Le, Barret Zoph, Jason Wei, et al. 2023. The FLAN collector: Designing data and methods for effective instruction tuning. *arXiv*, 2301.13668.
- Renqian Liu, Laij Sun, Yuzee Xia, Tao Qin, Sheng Zhang, Hoifung Poon, and Tie-Yan Liu. 2022. BiGPT: generative pre-trained transformer for biomedical text generation and mining. *Briefings in Bioinformatics* 23(6):bbac409.
- Varun Magesh, Faiz Surani, Matthew Dahl, Mirac Süzgür, Christopher D Manning, and Daniel E. Ho. 2024. H�lumination-free? assessing the reliability of leading AI legal research tools. *arXiv preprint arXiv:2405.20362*.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Suman Mille, Kausubh D. Dhole, Saad Mafamood, Laura Perez-Beltrachini, Varun Ganesh, Mirin Sanjay Kale, Enrico van Miltenburg, and Sebastian Gehrmann. 2021. Automatic construction of evaluation suites for natural language generation datasets. In *Advances in Neural Information Processing Systems*.
- Margaret Mitchell, Simone Wu, Andrew Zeldisvar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Isholiwa Deborah Raji, and Timnit Gebru. 2019. Model cards for model reporting. In *Proceedings of the Conference on Fairness, Accountability, and Transparency, FAT* 2019, Atlanta, GA, USA, January 29-31, 2019*, pages 220–229. ACM.
- Rajithro Nakano, Jacob Hilton, Suciur Balaji, Jeff Wu, Long Qiang, Christinia Kim, Christopher Hesse, Shantau Jain, Vineet Kosaraju, William Saunders, et al. 2021. Weigner: Browser-as-asked question answering with human feedback. *arXiv preprint arXiv:2112.09312*.
- Gabriel Nicholas and Aliya Blaauw. 2022. Lost in translation: Large language models in non-english content analysis. *arXiv preprint arXiv:2206.07377*.
- Harsha Nori, Nicholas King, Scott Mayer McClintock, Dean Cangiano, and Eric Horvitz. 2023. Capabilities of gpt-4 on medical challenge problems. *arXiv preprint arXiv:2303.3375*.
- Long Qiang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pameela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.
- Mattev E. Peters, Mark Neumann, Molit lyer, Ivan Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. *arXiv*, 1802.05365.
- Fabio Petroni, Patrick Lewis, Aleksandra Pitkns, Tim Rocktäschel, Yuxiang Wu, Alexander H. Miller, and Sebastian Riedel. 2020. How context affects language models’ factual predictions. In *Automated Knowledge Base Construction*.
- Fabio Petroni, Aleksandra Pitkns, Angela Fan, Purcket Lewis, Maijin Yezdi, Niccolò De Cenzo, James Thorpe, Yacine Jernite, Vladimír Karpátkin, Jean Marc' Harc, Vasilis Plaouris, Tim Rocktäschel, and Sebastian Riedel. 2021. KILT: a benchmark for knowledge intensive language tasks. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Findings in Language Technologies*, pages 2525–2544. Online. Association for Computational Linguistics.
- Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilia Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Jack W Rae, Sébastien Bergéaud, Trevor Cai, Katie Milligan, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susanah Yong, et al. 2021. Scaling language models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*.
- Jack W Rae, Anna Potapenko, Siddhanth M Jayakumar, and Timothy P Lillicrap. 2019. Compressive transformers for long-range sequence modelling. *arXiv preprint arXiv:1911.05507*.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5455–5551.
- Julian Schmidbauer, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guéz, Edward Lockhart, Denis Hassabis, Thore Graepel, et al. 2020. Mastering Atari, Go, Chess and Shogi by planning with a learned model. *Nature*, 583(7839):604–609.
- Elliott Schumacher and Mark Dredze. 2019. Learning unsupervised contextual representations for medical synonymy discovery. *JAMIA Open*, 2(4):538–546.
- Thibault Selam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892.

- Nikita Nangia, Niklas Deckers, Niklas Muenninghoff, Nitish Shirish Kesser, Niveditha Iyer, Noell Constant, Noah Field, Nuan Wen, Oliver Zhang, Onur Agha, Omar ElBardissi, Omer Levy, Owain Evans, Pablo Antonio Moreno Casares, Parth Desai, Pascale Fung, Paul Pu Lang, Paul Viola, Pegah Alipoormashahi, Peiyuan Liao, Percy Liang, Peter W. Chang, Peter Eckersley, Phu Mon Fittu, Peter Hwang, P. Milkowski, Piyush S. Patil, Pouya Pezeskpoor, Priti Oli, Ciacchini Mei, QING YU, Qiniang Chen, Rabin Banade, Rachel Etta Rudolph, Raefael Gabril, Rakesh Haraketer, Ram on Ricks, Delgado, Raphael Millière, Rhythm Garg, Richard Barnes, Kif A. Saurous, Riku Arakawa, Robbie Raymakers, Robert Frank, Rothen Sikand, Roman Novak, Roman Strellew, Ronan Le Bras, Rosanne Liu, Royan Jacobs, Ruiz Zheng, Ruolan Salakkurimmo, Ryan Chu, Ryan Lee, Ryan Stoval, Ryan Hammadi, Sajant Aranda, Sam Dilavore, Sam Shiebler, Sam Wiseman, Samuel Grutteer, Sam Bowman, Kwatra, Sarah A. Rous, Sarkis Ghazarian, Sayan Ghosh, Sean Casey, Srujan Reddy, Sebastian Gehmann, Sebastian Schuster, Sepideh Sadeghi, Shadi S. Hamdan, Sharon Zion, Shashank Saraswatava, Sherry Shan, Shikhar Singh, Shima Asadi, Shixiang Shan, Gu, Shubh Pachchigar, Shubham Toshniwal, Shyam Upadhyay, Shyamolika Debennath, Stanislak Skaten, Sironu Thomeyer, Simone Melzi, Sriva Radcy, Sruha Priscilla, Makini, Soo Iwan Lee, Spencer Bradley Torenz, Sriharsha Hattwar, Stanislas Uchene, Stefan Divic, Stefano Errone, Stella Rose Biderman, Stephanie C. Lin, S. Prasad, Steven T. Pientadosi, Stuart M. Shieber, Summer Mislerghi, Svetlana Kiritchenko, Swapna Mistry, Tal Linzen, Tal Schuster, Tao Li, Tao Yu, Tariq A. Ali, Tatsuo Hashimoto, Te-Lin Wu, Theo Despres, Theodore Rotschild, Thomas Plan, Tiaale Wang, Tiberius Nkonyeh, Timo Schizk, T. N. Korne, Tim Othy Telles-a-Lawon, Trus Tuncuny, Tobias Zbarsky, Trenton Chang, Irshika Meherji, Uri Shabtai, Khot, Tyler O'Brien Shultz, Uri Shaham, Yedidat Misra, Vera Denberg, Victoria Nyamai, Yekas Raviv, Vinay Venkatesh Ramasesh, Vinitay Jay Pabbi, Vishwak Pethakumar, Vivek Srikanth, William Fedus, William Saunders, William Zhang, W. Yessen, Xiang Ren, Xiaoyu Long, Xinyi Wu, Xuoding Shen, Yadollah Yaghoozbashi, Yao Lakeizi, Yang Song, Yasaman Bahri, Ye Ji Chou, Yichi Yang, Yiding Hao, Yifu Chen, Jonathan Belitsky, Yiu Hou, Yu He, Yuntao Bai, Zichary Seid, Zhou Xiaoran, Zhuoye Zhao, Zhi Fu Wang, Zheni Wang, Zheni Wang, Ziyi Wu, Sahib Singh, and Irfi Shaham. 2022. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv*, abs/2206.04615.
- Mirac Suzugun, Nathan Seales, Nathanael Schärtl, Sebastian Gehrmann, Y. Tay, Hyung Won Chung, Akanksha Chowdhery, Ooc Le, Ed Chi, Denny Zhou, et al. 2023. Challenging big-bench tasks and whether chain-of-thought can solve them. In *Finding the Association for Computational Linguistics: ACL 2023*, pages 1300–1305.
- Rohan Taati, Ishan Gujralani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tsunjon B. Hashimoto. 2023. Stanford alaca: An instruction-following llama model. <https://gitlab.com/tarsu-lab/stanford-a-paca>.
- Ross Taylor, Marcin Kardasz, Guillermo Cicurini, Thomas Siadom, Anthony Hartshorn, Elvis Savaya, Andrew Polton, Victor Kerbez, and Robert Stenico. 2022. Galecta: A large language model for science CoRR, abs/2211.09083.
- Hugo Touvron, Thibaut Levil, Gauthier Izquierdo, Xavier Martínez, Marie-Anne Lachaux, Timothée Lacoste, Baptiste Rozère, Naman Gowal, Eric Horvitz, Faissal Azhar, Aurélien Rodriguez, Armand Joulin, Édouard Grave, and Cédric Laptev. 2021. Llama: Open and efficient foundation language models. *CoRR*, abs/2302.13971.
- Hugo Touvron, Louis Marnot, Kevin Stone, Peter Alber, Amad Alrifai, Yannine Babiet, Nikolay Bashlykov, Sonny Barra, Prajnava Bhagavatula, Shruti Bhowmik, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.02386*.
- Audé Tacau and John Sweller. 2014. Domain-specific knowledge and why learning generic skills does not work. *Educational Psychology in Practice*, Review, 26:255–283.
- Kathryn Turusuvanakool, Jonas Adler, Zachary Wu, Tim Green, Micha Zieliński, Agustin Ździęka, Alex Bridgland, Clemens Meyer, Agata Laydon, et al. 2021. Highly accurate protein structure prediction for the human proteome. *Nature*, 596:7873:590–596.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008.
- Jesse Vig, Al Madani, Lav R. Varshney, Caiming Xiong, richard socher, and Nazneen Rajani. 2021. BERTology meets biology: Interpreting attention in protein language models. In *International Conference on Learning Representations*.
- Pelby Wang, Lei Li, Liang Chen, Dawei Zhu, Binghaiui Lin, Yunbo Cao, Qi Liu, Tanyu Liu, and Zhirong Su. 2023a. Large language models are not fair evaluators. *arXiv preprint arXiv:2305.17926*.
- Yirong Wang, Hemish Ivison, Pradeep Dasigi, Jack Hessel, Austin Khet, Raghav Chaudhuri, David Wadden, Kelsey MacMillan, Noah A. Smith, Eric Sppong, Alicia D. Giannmarino, Yingjie Weng, Anje Kumar, Poonam Hosamani, Jason Hom, and Jonathan H. Chen. 2023. Chatbot vs medical student performance on free-response clinical reasoning examinations. *JAMA Internal Medicine*, 183(9) 1028–1030.

- Iz Bellegay, et al. 2023b. How far can canels go? Exploring the state of instruction tuning on open resources. *arXiv preprint arXiv:2306.04751*.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guo, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2022. Fine-tuned language models are zero-shot learners. In *International Conference on Learning Representations*.
- Genta Indra Winata, Alham Fikri Aji, Samuel Cahyani-jaya, Rahmat Malenda, Fajri Koto, Ade Romdhony, Kenai Kurniawan, David Mojsilović, Radivoj Eko Prasetyo, Pascal Fung, et al. 2023. NasuX: Multilingual parallel sentiment dataset for 13 Indonesian local languages. In *Proceedings of the 17th Conference on the European Chapter of the Association for Computational Linguistics*, pages 815–834.
- Genta Indra Winata, Han yang Zhao, Anirban Das, Wen-ping Tang, David D Yao, Shi-Xiong Zhang, and Sam-bit Saha. 2024. Preference tuning with human feedback on language speech, and vision tasks: A survey. *arXiv preprint arXiv:2409.11564*.
- BigScience Workshop, Teven Le Scor, Angela Far, Christopher Araki, Ellie Pavlick, Suzana Illic, Daniel Hesselow, Roman Castañé, Alessandra Sasha Lucion, François Yvon, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*.
- Shijie Wu, Ozan Irsoy, Steven Lu, Václav Dabrowski, Mark Dredze, Sebastian Gehrmann, Prabhanshu Kambadur, David Rosenberg, and Gideon Mann. 2023. BloombergGPT: A large language model for finance. *arXiv preprint arXiv:2303.17504*.
- Lining Xue, Noah Constant, Adam Roberts, Mihir Kale, Kaini Al-Rouf, Aditya S. dhdant, Aditya Barua, and Colin Raffel. 2021. mT5: A massively multi-modal pre-trained text-to-text transformer. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 387–398.
- Travis Zack, Eric Leaman, Miran Sungan Joye, A. Rodriguez, Leo Anthony Celi, Hugh Zeehner, Dan Krufsky, Peter Szolovits, David V. Park, Karen E. E. Abrialhon, et al. 2024. AI4Health: The potential of GPT-4 to expedite racialized perspectives in health care: A model evaluation study. *The Lancet Digital Health*, 6(1) e12-e27.
- Daniel Zhang, Saheed Oyedele, Erik Brynjolfsson, John Eichenlaub, Pedro Cavigli, Barbara Grossz, Terah Lyons, Jeffrey Pfeffer, Juan Carlos Niebles, Michael Sefton, et al. 2022. The AI index 2021 annual report. *arXiv preprint arXiv:2103.06312*.
- Sophia Zheng, Stephen Roller, Namar Goyal, Michael J. Zeng, Moya Chen, Shuhui Chen, Christopher DeDene, Mona Diab, Xun Li, Xi Victoria Lin, et al. 2022. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01668*.