

Lessons from using PLMs for Human Cognitive Modeling

Anonymous ACL submission

Abstract

Many studies show evidence for cognitive abilities in Pre-trained Language Models (PLMs). Researchers have evaluated the cognitive alignment of PLMs, i.e., their correspondence to adult performance across a range of cognitive domains. More recently, the focus has expanded to the developmental alignment of these models: identifying phases during training where improvements in model performance track improvements in children’s thinking over development. However, challenges to this use are twofold: (1) PLMs have very different architectures than human minds and brains, and the data sets on which they are trained differ in many ways from the inputs children receive. (2) The “outputs” of PLMs are different from the behavioral measures that cognitive scientists collect in their experiments and evaluate their theories against. In this paper, we distill lessons learned from using PLMs for cognitive modeling and outline the pitfalls of attempting to use PLMs, not as engineering artifacts, but as cognitive science and developmental science models. We review assumptions used by researchers to map measures of PLM performance to measures of human performances and then, *enumerate criteria for using PLMs as credible accounts of cognition and cognitive development.*

1 Introduction

With the improving performance of language models (Touvron et al., 2023; Gemini Team, 2023; OpenAI, 2023; Wei et al., 2022), researchers have increasingly advocated for the use of Language models as computational models of cognition (Piantadosi, 2023; Mahowald et al., 2024; Warstadt and Bowman, 2024). This includes domains such as mathematical reasoning (Shah et al., 2023; Ahn et al., 2024), language comprehension (Warstadt et al., 2020; Ye et al., 2023; Koubaa, 2023), concept understanding (Vemuri et al., 2024), and analogical reasoning (Webb et al., 2023; Hu et al., 2023).

More recently, researchers have used PLMs for modeling cognitive development in children (Hosseini et al., 2022; Kosoy et al., 2023; Frank, 2023; Shah et al., 2024). For example, Portelance et al. (2023) suggests the use of language models to predict the age of acquisition of words in children. Shah et al. (2024) map the development of cognitive intelligence in humans to scaling training tokens and model size in PLMs. Other researchers also propose studying bilingualism by mapping pre-training steps in PLMs to understand the rate of language development (Evanson et al., 2023; Marian, 2023; Sharma et al., 2024).

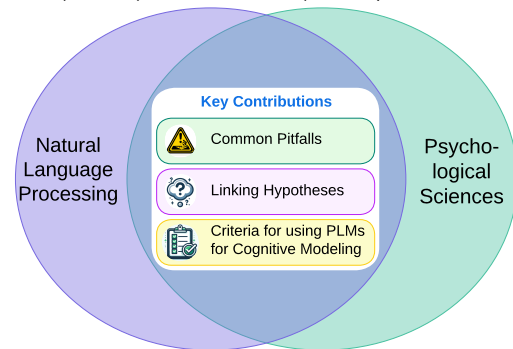


Figure 1: Overview of the Lessons from using PLMs for Human Cognitive Modeling.

In this paper, we advocate for the use of PLMs as candidate theories of cognitive and developmental science. We first review the pitfalls of using PLMs in psychological science and caution researchers against over-interpreting PLM alignment to human cognition. We then review the common assumptions used by researchers to map measures of PLM performance to measures of human performance. In doing so, we build upon previous work enumerating best practices for cognitive evaluations of PLMs (Ivanova, 2023; Mahowald et al., 2024).

2 Pitfalls of using PLMs as scientific theories

Some pitfalls come when using PLMs as cognitive science theories, i.e., of adult thinking.

071	• Human brains and PLMs are architecturally different. Recent research is attempting to map regions of the brain to different aspects of PLMs (layers, attention heads, etc.) in terms of local functionality and performance characteristics. However, this work is in its infancy, and its viability remains an open question (Hosseini et al., 2022; Kauf et al., 2023).	119
072		120
073		121
074		122
075		123
076		124
077		125
078		126
079	• Researchers use a <i>linking hypothesis</i> to map model performance characteristics (e.g., log probabilities) to human performance characteristics (e.g., reading times) (Shah et al., 2024). These links are often quite distal, making it unclear whether PLMs are actually “explaining” cognitive science data. (See the next section for further discussion.)	127
080		128
081		129
082		130
083		
084		
085		
086		
087	• PLMs are opaque and have limited interpretability. Human alignment and the lack thereof is hard to debug. This is a barrier to treating these models as scientific theories (McGrath et al.; Kar et al., 2022).	131
088		132
089		
090		
091		
092	• Most studies evaluating the cognitive alignment of PLMs focus on a narrow range of cognitive abilities and overlook correlations with other abilities. This is in contrast with psychometric approaches to intelligence that investigate the correlations across tests of a broad range of cognitive abilities: mathematical, verbal, spatial, fluid, and so on (Snow et al., 1984; Schneider and McGrew, 2012). This is also in contrast to unified theories of cognition that attempt to model all cognitive abilities within a single computational framework (Mellon et al., 2007; Varma, 2011).	
093		
094		
095		
096		
097		
098		
099		
100		
101		
102		
103		
104	Other pitfalls are specific to the use of PLMs as developmental science theories, i.e., of the progressions in children’s thinking.	
105		
106		
107	• PLM checkpoints are snapshots or fingerprints of the data they are trained on. Most research only looks at the final model checkpoints and not the change in the cognitive alignment of language models (development) as a function of data observed (Warstadt and Bowman, 2022; Frank, 2023; Shah et al., 2024). Often, this is because of the unavailability of intermediate training checkpoints or resource constraints. This limits our understanding of the nature of PLM training.	
108		
109		
110		
111		
112		
113		
114		
115		
116		
117	• Differences also exist in the nature of the data observed by PLMs versus humans. PLMs are	
118		
	trained on textual data that is magnitudes larger than the number of words seen by children (Huebner et al., 2021; Hosseini et al., 2022; Warstadt et al., 2023; Bhardwaj et al., 2024). On the other hand, children learn from input from multiple senses (Smith and Gasser, 2005), whereas models are not of an embodied nature (Chemero, 2023).	119
		120
		121
		122
		123
		124
		125
		126
	• For studies evaluating the developmental alignment, the observed developmental trajectories might be artifacts of the pre-training order (Shah et al., 2024).	127
		128
		129
		130
	3 Linking Hypotheses Mapping Model Performance to Human Performance	131
		132
	Researchers use various linking hypotheses to map PLM performance to human performance measures. While all these linking hypotheses are simple in theory, they have multiple possible operationalizations. We review four below.	133
		134
		135
		136
		137
	Similarity computations: Many cognitive tasks require people to judge the similarity of two items. In this case, human similarity judgments are directly modeled by computing the similarity between the corresponding representations in a PLM’s latent space. This can be via cosine similarity or another metric. One example is typicality effects, which is the finding people regard some members as “better” examples of a category than others (Bhatia and Richie, 2022; Richie and Bhatia, 2021). The rank of an item is determined by the proportion of humans that produce an item when asked to enumerate the items of the category. In language models, the typicality of an exemplar for a category is estimated by encoding the exemplar name as a string, passing it through the language model, and obtaining the corresponding word embedding. Thereafter, the similarity between this exemplar vector and the category prototype is calculated, with a higher value indicating that the exemplar is more typical.	138
		139
		140
		141
		142
		143
		144
		145
		146
		147
		148
		149
		150
		151
		152
		153
		154
		155
		156
		157
		158
	Other cognitive tasks require comparing or discriminating between two items. There, a common linking hypothesis is that the greater the similarity between the items in the model’s latent space, the longer the comparison/discrimination time. For example, Shah et al. (2023) map the time it takes to compare two numbers to PLM similarity - the greater the similarity of two number representations, the longer it takes for humans to differentiate which one is greater (or lesser).	159
		160
		161
		162
		163
		164
		165
		166
		167
		168

However, there are two common obstacles to using the *similarity* linking hypothesis. The first is that doing so requires models that make available the latent representations for similarity computations. The second problem is that this method suffers from problems due to tokenization. Humans use words as granular units while models use tokens (potentially words or subwords). This makes the nature of the mapping inconsistent as one unit of text for humans (words) may be mapped to two units of text for PLMs (tokens).

Surprisal values: One way of quantifying the uncertainty of model generations is in terms of the summation of logarithmic probabilities. A common linking hypothesis is that higher surprisal values correspond to longer human response times. For instance, studies of reading (Rambelli et al., 2024; Ivanova et al., 2024b) and categorization (Misra et al., 2021) have found that higher model uncertainties predict longer human response times. Relative log probabilities have been used to distinguish grammatical and ungrammatical sentences (Warstadt et al., 2020; Shah et al., 2024). They enable us to directly compare the right answer with all the candidate answers in a deterministic manner, i.e., there is no chance that the PLM will not calculate the sequential probability for a candidate string. Research also shows that surprisal values provide a better match to human plausibility judgments than prompts (Ivanova et al., 2024a).

However, a problem with surprisal-based approaches is that they fail to show robustness to context (prompts).

Prompting: PLM generation is probabilistic and therefore the same model can give different results across inference runs. A PLM is prompted to follow the same exercise as a human multiple times and generate a probability distribution over the output space. The probability of the correct output is then mapped to model confidence. Directly comparing the generated behavior of PLMs with that of humans reduces or even eliminates the need for linking hypotheses (Patel and Pavlick, 2021; Webb et al., 2023; Zahraei and Emami, 2024).

For example, a PLM can be prompted with: Which statement is grammatically correct? Your response must be "1" or "2".

1. Noah likes to swim. 2. Noah likes to.

The PLM can generate "1" for the correct answer.

Another benefit of prompting is that it allows for *variable output length*. Tasks that benefit from this flexibility, like commonsense reasoning, are more

suited for prompting (Yasunaga et al., 2023).

Example SAT analogical reasoning problem:

Analogy: Runner : Marathon :: ?

Options:

• Envoy : Embassy • oarsman : Regatta

• Martyr : Massacre • Horse : Stable

An example problem benefitting from *variable output length* is SAT analogical reasoning tasks (Turney, 2013). These are of the form A:B::? (see example above). These are MCQ-based choice tasks that can be operationalized as similarity computation, surprisal, or a prompt-based reasoning problem. In the prompting case, we can force the generations to adhere to a goal: *answer the correct analogy in this output format: {"A":"B"::"C":"D"}*.

That said, there are several limitations to the prompting approach. First, PLMs are *not robust to the prompt format*. Answering prompts requires PLMs to have two abilities, (1) understanding the prompt and (2) knowing the answer. Often, the outputs of PLMs vary substantially based on the input, and sometimes maximum performance characteristics are obtained on gibberish prompts (Deng et al., 2022). Second, PLMs may output an answer beyond the given candidate options for the prompting approach. For example, the answer to the problem described in blue above could be "It would be irresponsible to imply that the grammatical structure of a sentence is inconsequential, as clear communication is fundamental for safety and understanding" (Cai et al., 2024).

Direct probing: Another alternative is to directly ask the model about its current state. For example, to measure the incremental semantic understanding of temporarily ambiguous sentences, Li et al. (2024) presents a dichotomous verification sentence to the model after each word. In addition to the direct probing, one can directly recover the implicit parse tree after each word to measure incremental syntactic understanding (Manning et al., 2020). PLM prompting helps support direct probing, which is often not possible in humans with behavioral measures or even with neuroimaging measures to "look under the hood".

4 Criteria for Evaluating and Developing PLMs as Scientific Models

PLMs are being increasingly considered as models of cognitive and developmental science phenomena.

271	In light of the pitfalls outlined above, we propose	The second set of criteria is for guiding the <i>devel-</i>	319
272	two sets of criteria for using PLMs for this purpose.	<i>opment</i> of PLMs as credible accounts of cognition	320
273	The first set concerns the <i>appropriateness</i> of	and its development. This is a more open-ended	321
274	PLMs as scientific tools for cognitive and develop-	task, and the following can be considered as mere	322
275	mental modeling:	suggestions to researchers:	323
276	• Design multiple experiments to test the align-	• Evaluation techniques should follow the appropri-	324
277	ment of each cognitive or developmental phe-	ateness criteria above. PLMs should be evaluated	325
278	nomenon: PLMs may track the human perfor-	at regular intervals of pre-training to assess their	326
279	mance characteristic well under one linking	potential developmental alignment, which is of-	327
280	hypothesis or one type of test. However, this	ten overlooked in studies of cognitive alignment.	328
281	alignment may just be an artifact, for exam-		
282	ple, of pre-training data contamination. Con-	• PLMs can be tuned with specific cognitively im-	329
283	ducting more experiments evaluating the same	portant tasks and evaluated on a breadth of cog-	330
284	cognitive/developmental phenomena establishes	nitively relevant tasks. For example, typicality	331
285	stronger empirical plausibility.	experiments (Vemuri et al., 2024; Misra et al.,	332
286		2021) could be used to preference-tune PLMs	333
287	• Test the path-dependency of PLMs for devel-	using reinforcement learning techniques. This	334
288	opmental alignment: The claim that the final	may lead to better cognitive alignment of PLMs	335
289	model state of a PLM approximates adult perfor-	across a broader set of tasks.	336
290	mance leads to the question of the path by which		
291	it arrived there. Ideally, the model’s performance	• Pre-training data may benefit from developmen-	337
292	improvements over training should also track the	tally plausible corpora (Bhardwaj et al., 2024;	338
293	progression of cognitive abilities over develop-	Warstadt and Bowman, 2022; Frank, 2023). This	339
294	ment (Elman, 1996; Bengio et al., 2009). This	includes training on a curriculum based on hu-	340
295	would support researchers exploring the scaling	man skill acquisition, for example, the age-of-	341
296	of training data and model size in their investiga-	acquisition of a word (Huebner et al., 2021; Porte-	342
297	tions of human development.	lance et al., 2023). Informed pre-training will	343
298		allow us to better understand the developmental	344
299	• Use multiple methods to interpret PLM suc-	alignment of models.	345
300	cesses and failures: PLMs lack explainability		
301	and interpretability due to their large size (Mc-	5 Conclusion	346
302	Grath et al.). Some methods for PLM interpre-	This paper advocates for the use of Pre-trained	347
303	tation are often better than others. For example,	Language Models (PLMs) as theoretical tools for	348
304	in the experiments conducted by Li et al. (2024),	investigating human cognition and its development	349
305	incrementally constructed parse trees provided a	In this advocacy, we are not alone (Warstadt and	350
306	better account of PLM alignment than the infor-	Bowman, 2022; McGrath et al.; Frank, 2023; Ma-	351
307	mation from attention weights.	howald et al., 2024). At the same time, we cau-	352
308		tion researchers towards the informed use of PLMs	353
309	• Control for tuning techniques: PLMs are often	in psychological sciences. We have highlighted,	354
310	tuned on specific data and in different manners	common pitfalls, reviewed the different linking hy-	355
311	like Instruction Tuning, Reinforcement Learning	potheses used by researchers to map PLM perfor-	356
312	from Human Feedback, etc. These tuning tech-	mance to human performance, and outlined criteria	357
313	niques influence model behavior and the model’s	for evaluating and developing PLMs as credible	358
314	output centers around tuning goals rather than	models of cognition and cognitive development.	359
315	developing a representation of world knowledge.	These criteria are intended to guide researchers in	360
316		designing robust experiments, interpreting PLM	361
317	• Remember the linking hypothesis: Adaptation	behaviors accurately, and increasing fidelity to hu-	362
318	of human experimental materials to textual coun-	man data. Given the constantly evolving nature of	363
	terparts requires certain assumptions (refer to	the field, we call for researchers to continuously	364
	different operationalization in section 3). These	refine and expand these guidelines to match new	365
	assumptions need to be well documented and	advancements in NLP and Cognitive Science.	366
	understood based on the experiment type.		

6 Limitations

(1) The paper highlights common pitfalls, linking hypotheses, and evaluative criteria while using PLMs for cognitive modeling. These constitute a set of sound views to aid new researchers in the field. They do not exhaustively cover every pitfall, hypothesis, or criterion. (2) The suggestions in this work are good-to-have practices that support the use of PLMs for open cognitive and developmental science. There is no one-answer-fits-all approach for these tasks. Natural Language Processing is a developing field and we recommend articulating newer guidelines and practices as more PLMs are built and deployed. (3) Our work calls for the use of language technologies for the psychological sciences and provides criteria for developing credible accounts of cognition and cognitive development. Despite providing general guidelines, our work does not conduct experiments or offer any empirical evidence of performance comparisons or other quantitative measures.

7 Ethical Considerations

There are no major risks associated with conducting this research beyond those associated with working with PLMs. There may be risks in misinterpreting the criteria enlisted in this study. The suggestions in this study are one-way: we wish to find human performance characteristics and behaviors in PLMs to help model psychological sciences to aid people with cognitive impairments. We do not advocate for developing PLMs to replace humans or suggest ways to reach Artificial General Intelligence. PLMs are experimental technologies and future work using these models should proceed with caution.

References

Janice Ahn, Rishu Verma, Renze Lou, Di Liu, Rui Zhang, and Wenpeng Yin. 2024. Large language models for mathematical reasoning: Progresses and challenges. *arXiv preprint arXiv:2402.00157*.

Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. 2009. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*, pages 41–48.

Khushi Bhardwaj, Raj Sanjay Shah, and Sashank Varma. 2024. *Pre-training llms using human-like development data corpus*. *Preprint*, arXiv:2311.04666.

Sudeep Bhatia and Russell Richie. 2022. Transformer networks of human conceptual knowledge. *Psychological Review*.

Alice Cai, Ian Arawjo, and Elena L Glassman. 2024. Antagonistic ai. *arXiv preprint arXiv:2402.07350*.

Anthony Chemero. 2023. Llms differ from human cognition because they are not embodied. *Nature Human Behaviour*, 7(11):1828–1829.

Mingkai Deng, Jianyu Wang, Cheng-Ping Hsieh, Yihan Wang, Han Guo, Tianmin Shu, Meng Song, Eric P Xing, and Zhiting Hu. 2022. Rlprompt: Optimizing discrete text prompts with reinforcement learning. *arXiv preprint arXiv:2205.12548*.

Jeffrey L Elman. 1996. *Rethinking innateness: A connectionist perspective on development*, volume 10. MIT press.

Linnea Evanson, Yair Lakretz, and Jean-Rémi King. 2023. Language acquisition: do children and language models follow similar learning stages? *arXiv preprint arXiv:2306.03586*.

Michael C Frank. 2023. Bridging the data gap between children and large language models. *Trends in Cognitive Sciences*.

Gemini Team. 2023. *Gemini: A family of highly capable multimodal models*. *Preprint*, arXiv:2312.11805.

Eghbal A Hosseini, Martin Schrimpf, Yian Zhang, Samuel Bowman, Noga Zaslavsky, and Evelina Fedorenko. 2022. Artificial neural network language models align neurally and behaviorally with humans even after a developmentally realistic amount of training. *BioRxiv*, pages 2022–10.

Xiaoyang Hu, Shane Storks, Richard L Lewis, and Joyce Chai. 2023. In-context analogical reasoning with pre-trained language models. *arXiv preprint arXiv:2305.17626*.

Philip A. Huebner, Elinor Sulem, Fisher Cynthia, and Dan Roth. 2021. *BabyBERTa: Learning more grammar with small-scale child-directed language*. In *Proceedings of the 25th Conference on Computational Natural Language Learning*, pages 624–646, Online. Association for Computational Linguistics.

Anna A Ivanova. 2023. Running cognitive evaluations on large language models: The do’s and the don’ts. *arXiv preprint arXiv:2312.01276*.

Anna A Ivanova, Aalok Sathe, Benjamin Lipkin, Evelina Fedorenko, and Jacob Andreas. 2024a. Log probability scores provide a closer match to human plausibility judgments than prompt-based evaluations.

Anna A Ivanova, Aalok Sathe, Benjamin Lipkin, Unnathi Kumar, Setayesh Radkani, Thomas H Clark, Carina Kauf, Jennifer Hu, RT Pramod, Gabriel Grand, et al. 2024b. Elements of world knowledge (ewok):

467	A cognition-inspired framework for evaluating basic world knowledge in language models. <i>arXiv preprint arXiv:2405.09605</i> .	Roma Patel and Ellie Pavlick. 2021. Mapping language models to grounded conceptual spaces. In <i>International conference on learning representations</i> .	521
468			522
469			523
470	Kohitij Kar, Simon Kornblith, and Evelina Fedorenko. 2022. Interpretability of artificial neural network models in artificial intelligence versus neuroscience. <i>Nature Machine Intelligence</i> , 4(12):1065–1067.	Steven Piantadosi. 2023. Modern language models refute chomsky’s approach to language. <i>Lingbuzz Preprint, lingbuzz</i> , 7180.	524
471			525
472			526
473			
474	Carina Kauf, Anna A Ivanova, Giulia Rambelli, Emanuele Chersoni, Jingyuan Selena She, Zawad Chowdhury, Evelina Fedorenko, and Alessandro Lenci. 2023. Event knowledge in large language models: the gap between the impossible and the unlikely. <i>Cognitive Science</i> , 47(11):e13386.	Eva Portelance, Yuguang Duan, Michael C. Frank, and Gary Lupyan. 2023. Predicting age of acquisition for children’s early vocabulary in five languages using language model surprisal. <i>Cognitive science</i> , 47 9:e13334.	527
475			528
476			529
477			530
478			531
479			
480	Eliza Kosoy, Emily Rose Reagan, Leslie Lai, Alison Gopnik, and Danielle Krettek Cobb. 2023. Comparing machines and children: Using developmental psychology experiments to assess the strengths and weaknesses of lamda responses. <i>arXiv preprint arXiv:2305.11243</i> .	Giulia Rambelli, Emmanuele Chersoni, Davide Testa, Philippe Blache, and Alessandro Lenci. 2024. Neural generative models and the parallel architecture of language: A critical review and outlook. <i>Topics in cognitive science</i> .	532
481			533
482			534
483			535
484			536
485			
486	Anis Koubaa. 2023. GPT-4 vs. GPT-3.5: A concise showdown.	Russell Richie and Sudeep Bhatia. 2021. Similarity judgment within and across categories: A comprehensive model comparison. <i>Cognitive science</i> , 45(8):e13030.	537
487			538
488	Andrew Li, Xianle Feng, Siddhant Narang, Austin Peng, Tianle Cai, Raj Sanjay Shah, and Sashank Varma. 2024. Incremental comprehension of garden-path sentences by large language models: Semantic interpretation, syntactic re-analysis, and attention.	W Joel Schneider and Kevin S McGrew. 2012. The cattell-horn-carroll model of intelligence.	541
489			542
490			
491			
492			
493	Kyle Mahowald, Anna A Ivanova, Idan A Blank, Nancy Kanwisher, Joshua B Tenenbaum, and Evelina Fedorenko. 2024. Dissociating language and thought in large language models. <i>Trends in Cognitive Sciences</i> .	Raj Shah, Khushi Bhardwaj, and Sashank Varma. 2024. Development of cognitive intelligence in pre-trained language models. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> .	543
494			544
495			545
496			546
497	Christopher D Manning, Kevin Clark, John Hewitt, Urvasi Khandelwal, and Omer Levy. 2020. Emergent linguistic structure in artificial neural networks trained by self-supervision. <i>Proceedings of the National Academy of Sciences</i> , 117(48):30046–30054.	Raj Sanjay Shah, Vijay Marupudi, Reba Koenen, Khushi Bhardwaj, and Sashank Varma. 2023. Numeric magnitude comparison effects in large language models. <i>Preprint, arXiv:2305.10782</i> .	548
498			549
499			550
500			551
501			
502	Viorica Marian. 2023. Studying second language acquisition in the age of large language models: Unlocking the mysteries of language and learning, a commentary on “age effects in second language acquisition: Expanding the emergentist account” by catherine l. caldwell-harris and brian macwhinney. <i>Brain and language</i> , 246.	Mihir Sharma, Ryan Ding, Raj Sanjay Shah, and Sashank Varma. 2024. Monolingual and bilingual language acquisition in language models.	552
503			553
504			554
505			
506			
507			
508			
509	Sam Whitman McGrath, Jacob Russin, Ellie Pavlick, and Roman Feiman. How can deep neural networks inform theory in psychological science?	Linda Smith and Michael Gasser. 2005. The development of embodied cognition: Six lessons from babies. <i>Artificial life</i> , 11(1-2):13–29.	555
510			556
511			557
512	John R Anderson Richard King Mellon et al. 2007. <i>How can the human mind occur in the physical universe?</i> , volume 3. Oxford University Press, USA.	Richard E Snow, Patrick C Kyllonen, Brachia Marshalek, et al. 1984. The topography of ability and learning correlations. <i>Advances in the psychology of human intelligence</i> , 2(S 47):103.	558
513			559
514			560
515	Kanishka Misra, Allyson Ettinger, and Julia Taylor Rayz. 2021. Do language models learn typicality judgments from text? <i>arXiv preprint arXiv:2105.02987</i> .	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models. <i>Preprint, arXiv:2302.13971</i> .	562
516			563
517			564
518			565
519	OpenAI. 2023. Gpt-4 technical report. <i>Preprint, arXiv:2303.08774</i> .	Peter D. Turney. 2013. Distributional semantics beyond words: Supervised learning of analogy and paraphrase. <i>Transactions of the Association for Computational Linguistics</i> , 1:353–366.	566
520			567
			568
			569
			570
			571
			572

573 Sashank Varma. 2011. Criteria for the design and eval-
574 uation of cognitive architectures. *Cognitive science*,
575 35(7):1329–1351.

576 Siddhartha K Vemuri, Raj Sanjay Shah, and Sashank
577 Varma. 2024. How well do deep learning models
578 capture human concepts? the case of the typicality
579 effect.

580 Alex Warstadt and Samuel R Bowman. 2022. What
581 artificial neural networks can tell us about human lan-
582 guage acquisition. In *Algebraic structures in natural*
583 *language*, pages 17–60. CRC Press.

584 Alex Warstadt and Samuel R. Bowman. 2024. [What](#)
585 [artificial neural networks can tell us about human](#)
586 [language acquisition](#). *Preprint*, arXiv:2208.07998.

587 Alex Warstadt, Aaron Mueller, Leshem Choshen, Ethan
588 Wilcox, Chengxu Zhuang, Juan Ciro, Rafael Mos-
589 quera, Bhargavi Paranjabe, Adina Williams, Tal
590 Linzen, and Ryan Cotterell, editors. 2023. *Proceed-*
591 *ings of the BabyLM Challenge at the 27th Conference*
592 *on Computational Natural Language Learning*. As-
593 sociation for Computational Linguistics, Singapore.

594 Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mo-
595 hananey, Wei Peng, Sheng-Fu Wang, and Samuel R
596 Bowman. 2020. Blimp: The benchmark of linguistic
597 minimal pairs for english. *Transactions of the Asso-*
598 *ciation for Computational Linguistics*, 8:377–392.

599 Taylor Webb, Keith J Holyoak, and Hongjing Lu. 2023.
600 Emergent analogical reasoning in large language
601 models. *Nature Human Behaviour*, 7(9):1526–1541.

602 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raff-
603 fel, Barret Zoph, Sebastian Borgeaud, Dani Yo-
604 gatama, Maarten Bosma, Denny Zhou, Donald Met-
605 zler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals,
606 Percy Liang, Jeff Dean, and William Fedus. 2022.
607 Emergent abilities of large language models. *arXiv*
608 *preprint arXiv:2206.07682*.

609 Michihiro Yasunaga, Xinyun Chen, Yujia Li, Panupong
610 Pasupat, Jure Leskovec, Percy Liang, Ed H Chi, and
611 Denny Zhou. 2023. Large language models as ana-
612 logical reasoners. *arXiv preprint arXiv:2310.01714*.

613 Junjie Ye, Xuanting Chen, Nuo Xu, Can Zu, Zekai
614 Shao, Shichun Liu, Yuhan Cui, Zeyang Zhou, Chao
615 Gong, Yang Shen, Jie Zhou, Siming Chen, Tao Gui,
616 Qi Zhang, and Xuanjing Huang. 2023. A compre-
617 hensive capability analysis of GPT-3 and GPT-3.5 series
618 models. *arXiv preprint arXiv:2303.10420*.

619 Pardis Sadat Zahraei and Ali Emami. 2024. Wsc+:
620 Enhancing the winograd schema challenge using tree-
621 of-experts. *arXiv preprint arXiv:2401.17703*.