

Using Large Language Models as Beneficial Tools in Education

Anonymous ACL submission

Abstract

Public fame and easy open access to the ChatGPT, and the following wide use, or what could be considered misuse and abuse, of the model by some in the education and research communities, caused initially sharp negative reaction in the education and academic institutions and publishing services, aimed at detection and ban of the LLM (Large Language Models) generated texts, under efforts to combat plagiarism and chatting. Later, upon realising that such a blanket prohibition is technically problematic with the desired degree of reliability and confidence, as well as that LLMs can be legitimately used as tools for increasing productivity by taking on mundane writing tasks, the communities' attitude relaxed. The most remarkable changes in the public discourse are related to rethinking the very aims of the education system: "If some of the areas of the intellectual labour could be automated and become obsolete by LLM, maybe it is time for education to concentrate on teaching students to think and behave not like LLMs"? Such a Constructivist view on education, considered unrealistic a century ago, now may become the only sound way forward.

1 Introduction

Large Language Models (LLM) are posed to replace a significant part of so-called intellectual labour. Students, being taught by the current education system primarily to memorise, or at least to obtain pre-packaged "knowledge", will risk being outcompeted by the more efficient LLMs on routine and trivial tasks, which require extensive information search and mundane text generation. Therefore, new education adapted to the LLMs' presence needs to find intellectual labour niches in which humans are superior to LLMs, and needs to teach students to be not like LLMs to maintain competitiveness in the new market. Hence, significant changes are needed in education, preliminaries

to which, and changes themselves, we discuss in this position paper.

The contribution is organized in the following manner: Section 2 gives a brief overview of the attitude development on the LLM emergence; Section 3 discusses LLM flaws ; Section 4 outlines potential education changes to incorporate into the teaching process; and Section 5 concludes the discussion.

2 Large Language Models - a Friend or a Foe?

An explosive debut in public of the ChatGPT (Bib, 2023a) and the following similar Large Language Models (LLM) (Bib, 2023c; Chowdhery et al., 2022; Bib, 2023b; Touvron et al., 2023) also initiated a debate on LLMs' effects on education. An obvious first reaction was concern about abusing the LLMs' ability to generate human-like texts for cheating and plagiarism (Orenstrakh et al., 2023) in such examinations and tests that evaluate students in such faculties as memorisation, summarisation, reviewing, and basic analysis. Various methods of detection and prevention of using LLMs in education and academia were proposed (Tang et al., 2023; Khalil and Er, 2023; Rodriguez et al., 2022; Savelka et al., 2023).

However, the next wave of publications on the place of LLMs in education started to contemplate the thought that even if education shut the doors before LLMs, the industry would not, such as putting graduates who are not accustomed to the use of LLMs at a disadvantage. The publications started coming to the conclusion that education itself should change, not pursuing obsolete goals and not executing obsolete practices (Anders, 2023; Rudolph et al., 2023), but instead concentrating more on the areas where human-lead education (even armed with LLMs as tools) has advantages over mere LLMs in themselves (Fuchs, 2023; Cope and Kalantzis, 2019).

082 From the literary text analysis perspective, the
083 generated by LLMs, though usually syntactically
084 correct, are effete, emotionless washed-up texts,
085 lacking linguistic variability and distinctness, and
086 pragmatic intercity and originality (Gao et al.,
087 2022; Chaves and Gerosa, 2021; Wilkenfeld et al.,
088 2022; Mitrović et al., 2023). On the dynamic de-
089 bating or deliberation text generation, LLMs also
090 perform far from ideal. For example, on detecting
091 discourse move, ChatGPT performed even worse
092 than simple BERT models (Wang et al., 2023). De-
093 bates with ChatGPT, as everybody can see using the
094 OpenAI interface, suffer from circular arguments,
095 self-contradiction, and evasiveness - tendencies to
096 please human preferences in Reinforcement Learn-
097 ing (RL) (Ramamurthy et al., 2022; Carta et al.,
098 2023) - exactly those practices that nobody wants
099 to foster in students. When used to detect manipu-
100 lative discussion tactics of cyberattacks, ChatGPT
101 also scored significantly worse than simple BERT
102 models (Fayyazi and Yang, 2023).

103 General LLMs' problems with functional do-
104 mains such as mathematics, reasoning, and logic
105 (Frieder et al., 2023), emotional expressivity, wit,
106 humour and ethics (Borji, 2023; Arkoudas, 2023),
107 factual data, privacy, and false, bias and discrim-
108 ination (Basta et al., 2019; Kurita et al., 2019;
109 Sheng et al., 2019; Gehman et al., 2020; Bib,
110 2022; Bianchi et al., 2022; Weidinger et al., 2021;
111 Tang et al.; Goldstein et al., 2023) are well doc-
112 umented. Machine Learning (ML) specific prob-
113 lems of LLMs add such issues as lack of inter-
114 pretableity and understanding, (Bender and Koller,
115 2020; Lake and Murphy, 2020; Marcus et al., 2022;
116 Ouyang et al., 2022; Leivada et al., 2022; Ruis
117 et al., 2022), and catastrophic ageing and forgetting
118 by LLMs (Lazaridou et al., 2021; Hombaiyah et al.,
119 2021; Dhingra et al., 2022; McCloskey and Cohen,
120 1989; Parisi et al., 2019; Ratcliff, 1990; Kirkpatrick
121 et al., 2017). When using LLMs in education, their
122 shortcomings may not only be accounted for in the
123 real-life application but also can be used as a foun-
124 dation of fresh approaches to education to foster
125 those qualities and skills of students that will not
126 be made obsolete by the use of LLMs, and on the
127 opposite, give students a competitive edge.

128 3 Fundamental Foundations of the LLMs' 129 Flaws

130 Although implementation details of the latest mod-
131 els are kept proprietary, previously published re-

132 search shows that LLM models are built and trained
133 using three main principles. Traditional Natural
134 Language Processing (NLP) tokenizing techniques
135 include the preprocessing stage, on which "stop-
136 words" are removed, remaining words are stemmed
137 and lemmatized (converted to canonical dictionary
138 form), and the Bag of Words (BoW) algorithm is
139 used to map lemmatized words into a linear vector
140 space, spanned on the most frequent and impor-
141 tant words dictionary basis. The whole sentence
142 or a bigger text is represented as a linear sum of
143 all token vectors (or also so-called "embeddings")
144 (Zhang et al., 2010). Such an approach is very re-
145 source usage effective but does not count in the
146 sentence or larger text structure. For example, such
147 sentences as: "A dog bites a man", "A man bites a
148 dog", and "Dogs bite men" would be represented
149 by the same embedding.

150 To introduce implicit elements of the linguis-
151 tic structures, modern NLP models frequently use
152 context tokenizers (Taylor, 1953) of the BERT-like
153 family (Devlin et al., 2018). A simple illustration of
154 the BoW and BERT embedding differences would
155 be the former creating "DOG", "BITE", "MAN",
156 and the latter - "nullDOGbite", "dogBITEman",
157 "biteMANnull", "nullMANbite", "manBITEdog",
158 "biteDOGnull". That solves the BoW's structure
159 blindness problem but greatly increases the dimen-
160 sionality of the embedding space, which is the start-
161 ing point of LLMs' high computational demands
162 and size.

163 The second foundation technology the LLMs
164 use is based on the statistical n-gram approach
165 (Brown et al., 1992). The supervised training of the
166 Machine Learning (ML) models has a bottleneck
167 in the manual labelling of the training data sets.
168 To process high amounts of text and other media,
169 LLM uses a self-supervised approach based on the
170 Masked Language Model (MLM) (Salazar et al.,
171 2019; Besag, 1975). In such a paradigm, part of
172 the words are kept hidden from the ML model in
173 training, and the purpose of the training is to find
174 words with the highest probability of being in the
175 hidden positions. Again, such an approach does not
176 directly model linguistic structures but implicitly
177 stochastically takes them into account.

178 To keep with the human reader's attention span
179 and produce a coherent flow of text, LLMs have
180 to use long context windows for MLM training
181 of thousands of words. The brute force use of
182 the whole continuous windows is computationally

183 problematic; therefore, another technique of ex- 234
184 tracting the most valuable and influential context 235
185 words on the predicted word gave birth to compu- 236
186 tationally tractable but still huge LLMs - Attention 237
187 mechanism (Bahdanau et al., 2014; Luong et al., 238
188 2015; Gehring et al., 2016) and its Transformer im- 239
189 plementation (Vaswani et al., 2017). In such an ap- 240
190 proach of “attention”, learnable matrices are used 241
191 to compute cosine or Euclidean distances between 242
192 the word relevance to the projected prediction over 243
193 the context window sliding, and the most consistent 244
194 contributor over time is kept and used, in such a 245
195 way, reducing computational demand. 246

196 The stochastic nature of the LLMs in modelling 247
197 structured natural languages has been a point of 248
198 fierce debate since the LLMs introduction (Ben- 249
199 der et al., 2021; Schick and Schütze, 2020; Mar- 250
200 cus, 2018; Blodgett and Madaio, 2021; Bommasani 251
201 et al., 2021). 252

202 Another obvious problem of LLMs is the naivety 253
203 of their language representation from the theoret- 254
204 ical linguistics perspective that operates with cat- 255
205 egories of syntactic and semantic structures. The 256
206 former are various kinds or relations in the mathe- 257
207 matical sense (Combe et al., 2022; Marcolli et al., 258
208 2023), specific to particular languages, which en- 259
209 dow non-ordered multi-sets of the morphing lex- 260
210 emes and are continuously mapped to the univer- 261
211 sal semantic structures (of meaning or of thought) 262
212 (Chomsky, 2023) (or, possibly, to universal gram- 263
213 mar) (Watumull and Chomsky, 2020). 264

214 Noam Chomsky especially emphasises the non- 265
215 locality of such synthetic units. For example, in 266
216 inflectional languages such as Balto-Slavic, or ag- 267
217 glutinating such as Japanese, the non-locality is 268
218 obvious because of their free word order, but even 269
219 for the significantly sequential analytic English, 270
220 Chomsky referees at the semantic attachment of an 271
221 adverb to a correct verb regardless of their position 272
222 and order, for example in “Intuitively, birds that fly 273
223 swim” (Berwick and Chomsky, 2016). 274

224 Building models of such complex relations in 275
225 LLMs, capable of discovering and retrieving such 276
226 linguistic structures and, in such a way, achieving 277
227 explainability and interoperability of LLMs, is a 278
228 drastically undeveloped area of research (Delétang 279
229 et al., 2022), frequently limited to naive methods 280
230 of asking LLMs about their internals (Jiang et al., 281
231 2020). 282

232 These mechanisms introduce implicit naive syn- 283
233 tax emulation elements by projecting hierarchical

tree structures on flat sequences but with the loss of 234
complexity. For example, in Chomsky’s example, 235
“Intuitively” can become the sequential neighbour 236
of “swim” by dropping “fly”. 237

238 Even more complicated question of whether 239
LLMs can model thought and intelligence, al- 240
though receiving some optimistic answers (Kosin- 241
ski, 2023; Bubeck et al., 2023), predominately an- 242
swered negatively (Ullman, 2023; Sap et al., 2022). 243

244 From the linguistics view on natural human lan- 245
guages, universal semantic roles and relations be- 246
tween parts of a sentence, for example “Elmer 247
threw a porcupine to Hortense”, such as Actor 248
(Elmer), Patient (porcupine), and Beneficiary (Hort- 249
ense) could be mapped to syntactic roles and re- 250
lations, specific to particular languages (Marantz, 251
1981). In English, syntactic relations between Sub- 252
ject, Direct and Indirect Objects are marked by the 253
order and prepositions (to); in languages such as 254
Balto-Slavic - by the case (nominative, accusative, 255
dative) suffixes; in Japanese - by particles (を, に). 256

257 However, the question of what is the language of 258
semantics/meaning, or the “language of thought”, 259
and how it is externalised into syntactic structures, 260
is difficult even for linguistics and neuroscience of 261
the natural human languages (Gallistel, 2011). 262

263 Surprisingly, in the last years, the voices of the 264
critics of the limitations of the traditional narrow 265
ML (and LLMs as part of it), such as Noam Chom- 266
sky and Garry Marcus, were joined by such big 267
names of the narrow ML as Joshua Bengio (Lex 268
Clips, 2023), Yann LeCun (Bib, 2023d), and even 269
Geoffrey Hinton whose students built ChatGPT 270
(Metz, 2023). 271

272 4 Education Ameliorating Horizons In 273 274 the Context of LLMs 275

276 Although LLMs lack agency, structural represen- 277
tation of the language, and real-world picture 278
(Browning, 2022; Floridi, 2023), they, under hu- 279
man teacher supervision, could still be used to 280
help foster those abilities in students. Such non- 281
commodified abilities to behave not like LLM 282
(LLMs behaviour is described by Ben Goertzel 283
as “competent mediocrity” (Charrington, 2023)), 284
will remain in high value and demand. 285

286 Educational methodologies founded on initia- 287
tive, curiosity, and active actionable students’ con- 288
struction of knowledge (Vygotsky, 2012; Beilin, 289
1992; Shchedrovitsky, 1995), and therefore de- 290
manding high educator involvement, hence pro- 291

hibitively costly, with the routine and trivial tasks delegated to LLMs may become practically sound. We want students to be “competent”, for which goal LLMs may be useful tools and examples, but also not “mediocre”, for which LLMs may be used as counter-example tools. It’s been observed that LLM-generated scientific paper abstracts are easily identified by humans based on Goertzel’s “competent mediocrity” style, though such estimates have a noticeable false positive error - people also write papers in such a style (Gao et al., 2022).

Sporadic research in applying LLMs to education change in the active direction is visible in publications. For example, one of the routine tasks a competent educator may be released from, but a general eye on, is the trace of the students’ discourse flow (Wang et al., 2023), or teamwork feedback (Katz et al., 2023). Constant feedback, personalized and adaptive learning (Annuš, 2023), student initiative and psychometrics (Katz et al., 2023), collaborative, transparent and diverse intelligence (Cope et al., 2021). LLMs and other AI models are inherently student-driven, and it’s up to the education system, particularly up to its change, to view and experience that drive as a threat or benefit (Dai et al., 2023; Haensch et al., 2023).

We propose systematic research on the use of LLMs and other AI methods in practical implementation methods of education of constructing knowledge and understanding, such as (but not limited to):

- Fostering a big picture view, understanding, and based on them, first-hand actionable application, experimentation and implementation of the knowledge.
- Continuous, recursive (i.e. changing assignments) feedback (aizuchi - a rare Japanese loan into English linguistic jargon (Kita and Ide, 2007)).
- Pursuit of student questions and interests. Interactive (i.e. self-assigning) and co-acting (together with pedagogue) learning.
- Non-disciplinary or non-didactic learning, self-involved assessment.
- Dynamic knowledge acquisition, with each step in it being a challenge for the student, seemingly impossible, but with guidance and work achievable, building confidence in own abilities.

- Collaborative, social learning - learning through teaching other students.

- Emotion and sentiment expression aware and competent learning and teaching.

5 Discussion and Conclusions

5.1 Limitations

The presented review is in no way comprehensive and exhaustive - a number of publications on various aspects of LLM creation and use are published at an astonishing rate, and the very LLM landscape is changing quickly, outpacing academic publishing cycles. The research results are frequently contradicting, not merely because some of them are not rigorous - the research field is so vast that available results are fragmented and patchy, depending on the initial conditions that hardly can cover comprehensively all possible aspects of the LLM use. Inevitably, this opinion piece is incomplete in its foundations and subjective in proposals.

5.2 Risks

A significant change in the education system, especially if it is related to a significant cost increase, and hence, applied to limited society strata, can lead to further societal disparity. However, the risks of keeping the outdated education system that produces an incompetent and unneeded workforce can be even greater.

5.3 Conclusions

Under the likely perspective of LLMs taking on a significant share of the previously thought of “intellectual” labour, education needs to shift its goals and methods to fostering students’ abilities and habits that differentiate them from LLMs. That requires gaining a better understanding of what LLMs can not successfully do, not only from the empirical perspective but also from the first principles laying in the foundations of LLM. Building the education system from human strengths, such as agency, individual initiative and interest, social collaboration, emotional involvement, and structural view of the language and world picture, would likely require significant and expensive education system change, the core of which would likely align with the Constructivist view on it Of Vygotsky and Piaget.

378	References	430	
379	2022. Doctor GPT-3: hype or reality? - Nabla,	431	
380	https://www.nabla.com/blog/gpt-3 . [Online; ac-		
381	cessed 5. Sep. 2022].		
382	2023a. GPT-4, https://openai.com/research/gpt-4 . [On-		
383	line; accessed 16. Mar. 2023].		
384	2023b. Introducing LLaMA: A founda-		
385	tional, 65-billion-parameter language model,		
386	https://ai.facebook.com/blog/large-language-model-		
387	llama-meta-ai . [Online; accessed 17. Mar. 2023].		
388	2023c. Pathways Language Model (PaLM): Scaling		
389	to 540 Billion Parameters for Breakthrough Perform-		
390	ance, https://ai.googleblog.com/2022/04/pathways-		
391	language-model-palm-scaling-to.html . [Online; ac-		
392	cessed 17. Mar. 2023].		
393	2023d. Post LinkedIn. [Online; accessed 11. Aug.		
394	2023].		
395	Brent A Anders. 2023. Is using chatgpt cheating, pla-		
396	giarism, both, neither, or forward thinking? <i>Patterns</i> ,		
397	4(3).		
398	Norbert Annuš. 2023. Chatbots in education: The		
399	impact of artificial intelligence based chatgpt on		
400	teachers and students. <i>International Journal of</i>		
401	<i>Advanced Natural Sciences and Engineering Re-</i>		
402	<i>searches</i> , 7(4):366–370.		
403	Konstantine Arkoudas. 2023. ChatGPT is no stochastic		
404	parrot. But it also claims $1 > 1$. Medium. <i>Medium</i> .		
405	Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Ben-		
406	gio. 2014. Neural machine translation by jointly		
407	learning to align and translate. <i>arXiv preprint</i>		
408	<i>arXiv:1409.0473</i> .		
409	Christine Basta, Marta R. Costa-jussà, and Noe Casas.		
410	2019. Evaluating the underlying gender bias in con-		
411	textualized word embeddings.		
412	Harry Beilin. 1992. Piaget’s enduring contribution to de-		
413	velopmental psychology. <i>Developmental psychology</i> ,		
414	28(2):191.		
415	Emily M. Bender, Timnit Gebru, Angelina McMillan-		
416	Major, and Shmargaret Shmitchell. 2021. On the		
417	dangers of stochastic parrots: Can language mod-		
418	els be too big? . In <i>Proceedings of the 2021 ACM</i>		
419	<i>Conference on Fairness, Accountability, and Trans-</i>		
420	<i>parency</i> , FAccT ’21, page 610–623, New York, NY,		
421	USA. Association for Computing Machinery.		
422	Emily M. Bender and Alexander Koller. 2020. Climbing		
423	towards NLU: On meaning, form, and understanding		
424	in the age of data. In <i>Proceedings of the 58th Annual</i>		
425	<i>Meeting of the Association for Computational Lin-</i>		
426	<i>guistics</i> , pages 5185–5198, Online. Association for		
427	Computational Linguistics.		
428	Robert C Berwick and Noam Chomsky. 2016. <i>Why only</i>		
429	<i>us: Language and evolution</i> . MIT press.		
	Julian Besag. 1975. Statistical analysis of non-lattice		
	data. <i>Journal of the Royal Statistical Society: Series</i>		
	<i>D (The Statistician)</i> , 24(3):179–195.	432	
	Federico Bianchi, Pratyusha Kalluri, Esin Durmus,	433	
	Faisal Ladhak, Myra Cheng, Debora Nozza, Tat-	434	
	sunori Hashimoto, Dan Jurafsky, James Zou, and	435	
	Aylin Caliskan. 2022. Easily accessible text-to-	436	
	image generation amplifies demographic stereotypes	437	
	at large scale.	438	
	Su Lin Blodgett and Michael Madaio. 2021. Risks	439	
	of AI foundation models in education. <i>CoRR</i> ,	440	
	abs/2110.10024.	441	
	Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ	442	
	Altman, Simran Arora, Sydney von Arx, Michael S.	443	
	Bernstein, Jeannette Bohg, Antoine Bosselut, Emma	444	
	Brunskill, Erik Brynjolfsson, Shyamal Buch, Dal-	445	
	las Card, Rodrigo Castellon, Niladri S. Chatterji,	446	
	Annie S. Chen, Kathleen Creel, Jared Quincy	447	
	Davis, Dorottya Demszky, Chris Donahue, Moussa	448	
	Doumbouya, Esin Durmus, Stefano Ermon, John	449	
	Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea	450	
	Finn, Trevor Gale, Lauren Gillespie, Karan Goel,	451	
	Noah D. Goodman, Shelby Grossman, Neel Guha,	452	
	Tatsunori Hashimoto, Peter Henderson, John He-	453	
	witt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing	454	
	Huang, Thomas Icard, Saahil Jain, Dan Jurafsky,	455	
	Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keel-	456	
	ing, Fereshte Khani, Omar Khattab, Pang Wei Koh,	457	
	Mark S. Krass, Ranjay Krishna, Rohith Kuditipudi,	458	
	and et al. 2021. On the opportunities and risks of	459	
	foundation models. <i>CoRR</i> , abs/2108.07258.	460	
	Ali Borji. 2023. A categorical archive of chatgpt fail-	461	
	ures. <i>arXiv preprint arXiv:2302.03494</i> .	462	
	Peter F Brown, Vincent J Della Pietra, Peter V Desouza,	463	
	Jennifer C Lai, and Robert L Mercer. 1992. Class-	464	
	based n-gram models of natural language. <i>Computa-</i>	465	
	<i>tional linguistics</i> , 18(4):467–480.	466	
	Jacob Browning. 2022. AI And The Limits Of Lan-	467	
	guage. <i>NOEMA</i> .	468	
	Sébastien Bubeck, Varun Chandrasekaran, Ronen El-	469	
	dan, Johannes Gehrke, Eric Horvitz, Ece Kamar,	470	
	Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lund-	471	
	berg, et al. 2023. Sparks of artificial general intelli-	472	
	gence: Early experiments with gpt-4. <i>arXiv preprint</i>	473	
	<i>arXiv:2303.12712</i> .	474	
	Thomas Carta, Clément Romac, Thomas Wolf, Sylvain	475	
	Lamprier, Olivier Sigaud, and Pierre-Yves Oudeyer.	476	
	2023. Grounding large language models in interac-	477	
	tive environments with online reinforcement learning.	478	
	<i>arXiv preprint arXiv:2302.02662</i> .	479	
	The Twiml Ai Podcast with Sam Charrington. 2023.	480	
	Are Large Language Models a Path to AGI? with	481	
	Ben Goertzel - 625. [Online; accessed 9. Aug. 2023].	482	

488	Ana Paula Chaves and Marco Aurelio Gerosa. 2021. The impact of chatbot linguistic register on user perceptions: a replication study. In <i>International Workshop on Chatbot Research and Design</i> , pages 143–159. Springer.	
489	Noam Chomsky. 2023. Genuine explanation and the strong minimalist thesis. <i>Cognitive Semantics</i> , 8(3):347–365.	
490		
491	Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. Palm: Scaling language modeling with pathways. <i>arXiv preprint arXiv:2204.02311</i> .	
492		
493	Noemie Combe, Yuri I Manin, and Matilde Marcolli. 2022. Geometry of information: Classical and quantum aspects. <i>Theoretical Computer Science</i> , 908:2–27.	
494		
495	Bill Cope and Mary Kalantzis. 2019. Education 2.0: Artificial intelligence and the end of the test. <i>Beijing International Review of Education</i> , 1(2-3):528–543.	
496		
497	Bill Cope, Mary Kalantzis, and Duane Sears-Smith. 2021. Artificial intelligence for education: Knowledge and its assessment in ai-enabled learning ecologies. <i>Educational Philosophy and Theory</i> , 53(12):1229–1245.	
498		
499	Yun Dai, Ang Liu, and Cher Ping Lim. 2023. Reconceptualizing chatgpt and generative ai as a student-driven innovation in higher education.	
500		
501	Grégoire Delétang, Anian Ruoss, Jordi Grau-Moya, Tim Genewein, Li Kevin Wenliang, Elliot Catt, Chris Cundy, Marcus Hutter, Shane Legg, Joel Veness, et al. 2022. Neural networks and the chomsky hierarchy. <i>arXiv preprint arXiv:2207.02098</i> .	
502		
503	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. <i>arXiv preprint arXiv:1810.04805</i> .	
504		
505	Bhuwan Dhingra, Jeremy R. Cole, Julian Martin Eisenschlos, Daniel Gillick, Jacob Eisenstein, and William W. Cohen. 2022. Time-aware language models as temporal knowledge bases. <i>Transactions of the Association for Computational Linguistics</i> , 10:257–273.	
506		
507	Reza Fayyazi and Shanchieh Jay Yang. 2023. On the uses of large language models to interpret ambiguous cyberattack descriptions. <i>arXiv preprint arXiv:2306.14062</i> .	
508		
509	Luciano Floridi. 2023. Ai as agency without intelligence: on chatgpt, large language models, and other generative models. <i>Philosophy & Technology</i> , 36(1):15.	
510		
511	Simon Frieder, Luca Pinchetti, Ryan-Rhys Griffiths, Tommaso Salvatori, Thomas Lukasiewicz, Philipp Christian Petersen, Alexis Chevalier, and Julius Berner. 2023. Mathematical capabilities of chatgpt. <i>arXiv preprint arXiv:2301.13867</i> .	
512		
513	Kevin Fuchs. 2023. Exploring the opportunities and challenges of nlp models in higher education: is chatgpt a blessing or a curse? In <i>Frontiers in Education</i> , volume 8, page 1166682. Frontiers.	540
514		541
515	Charles Randy Gallistel. 2011. Prelinguistic thought. <i>Language learning and development</i> , 7(4):253–262.	542
516		543
517	Catherine A Gao, Frederick M Howard, Nikolay S Markov, Emma C Dyer, Siddhi Ramesh, Yuan Luo, and Alexander T Pearson. 2022. Comparing scientific abstracts generated by chatgpt to original abstracts using an artificial intelligence output detector, plagiarism detector, and blinded human reviewers. <i>BioRxiv</i> , pages 2022–12.	544
518		545
519	Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. 2020. Realltoxicityprompts: Evaluating neural toxic degeneration in language models.	546
520		547
521	Jonas Gehring, Michael Auli, David Grangier, and Yann N Dauphin. 2016. A convolutional encoder model for neural machine translation. <i>arXiv preprint arXiv:1611.02344</i> .	548
522		549
523	Josh A. Goldstein, Girish Sastry, Micah Musser, Renee DiResta, Matthew Gentzel, and Katerina Sedova. 2023. Generative language models and automated influence operations: Emerging threats and potential mitigations.	550
524		551
525	Anna-Carolina Haensch, Sarah Ball, Markus Herklotz, and Frauke Kreuter. 2023. Seeing chatgpt through students’ eyes: An analysis of tiktok data. <i>arXiv preprint arXiv:2303.05349</i> .	552
526		553
527	Spurthi Amba Hombaiyah, Tao Chen, Mingyang Zhang, Michael Bendersky, and Marc Najork. 2021. Dynamic language models for continuously evolving content. In <i>Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining</i> . ACM.	554
528		555
529	Zhengbao Jiang, Frank F Xu, Jun Araki, and Graham Neubig. 2020. How can we know what language models know? <i>Transactions of the Association for Computational Linguistics</i> , 8:423–438.	556
530		557
531	Andrew Katz, Siqing Wei, Gaurav Nanda, Christopher Brinton, and Matthew Ohland. 2023. Exploring the efficacy of chatgpt in analyzing student teamwork feedback with an existing taxonomy. <i>arXiv preprint arXiv:2305.11882</i> .	558
532		559
533	Mohammad Khalil and Erkan Er. 2023. Will chatgpt get you caught? rethinking of plagiarism detection. <i>arXiv preprint arXiv:2302.04335</i> .	560
534		561
535	James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. 2017. Overcoming catastrophic forgetting in neural networks.	562
536		563

Proceedings of the national academy of sciences,
114(13):3521–3526. 645

594 Sotaro Kita and Sachiko Ide. 2007. Nodding, aizuchi,
595 and final particles in japanese conversation: How
596 conversation reflects the ideology of communica-
597 tion and social relationships. *Journal of Pragmatics*,
598 39(7):1242–1254.

599 Michal Kosinski. 2023. [Theory of mind may have spon-](#)
600 [taneously emerged in large language models.](#)

601 Keita Kurita, Nidhi Vyas, Ayush Pareek, Alan W Black,
602 and Yulia Tsvetkov. 2019. [Measuring bias in contex-](#)
603 [tualized word representations.](#)

604 Brenden M. Lake and Gregory L. Murphy. 2020. [Word](#)
605 [meaning in minds and machines.](#)

606 Angeliki Lazaridou, Adhiguna Kuncoro, Elena Gri-
607 bovskaya, Devang Agrawal, Adam Liska, Tayfun
608 Terzi, Mai Gimenez, C d M d’Autume, Sebas-
609 tian Ruder, Dani Yogatama, et al. 2021. Pit-
610 falls of static language modelling. *arXiv preprint*
611 *arXiv:2102.01951*.

612 Evelina Leivada, Elliot Murphy, and Gary Marcus. 2022.
613 [Dall-e 2 fails to reliably capture common syntactic](#)
614 [processes.](#)

615 Lex Clips. 2023. [Yoshua Bengio: From System 1 Deep](#)
616 [Learning to System 2 Deep Learning \(NeurIPS 2019\).](#)
617 [Online; accessed 11. Aug. 2023].

618 Minh-Thang Luong, Hieu Pham, and Christopher D
619 Manning. 2015. Effective approaches to attention-
620 based neural machine translation. *arXiv preprint*
621 *arXiv:1508.04025*.

622 Alec Marantz. 1981. *On the nature of grammatical*
623 *relations*. Ph.D. thesis, Massachusetts Institute of
624 Technology.

625 Matilde Marcolli, Noam Chomsky, and Robert Berwick.
626 2023. Mathematical structure of syntactic merge.
627 *arXiv preprint arXiv:2305.18278*.

628 Gary Marcus. 2018. [Deep learning: A critical appraisal.](#)
629 *CoRR*, abs/1801.00631.

630 Gary Marcus, Ernest Davis, and Scott Aaronson. 2022.
631 [A very preliminary analysis of dall-e 2.](#)

632 Michael McCloskey and Neal J. Cohen. 1989. [Catas-](#)
633 [trophic interference in connectionist networks: The](#)
634 [sequential learning problem.](#) volume 24 of *Psychol-*
635 *ogy of Learning and Motivation*, pages 109–165. Aca-
636 demic Press.

637 Cade Metz. 2023. [‘The Godfather of AI’ Quits Google](#)
638 [and Warns of Danger Ahead.](#) *N.Y. Times*.

639 Sandra Mitrović, Davide Andreoletti, and Omran Ay-
640 oub. 2023. Chatgpt or human? detect and explain.
641 explaining decisions of machine learning model for
642 detecting short chatgpt-generated text. *arXiv preprint*
643 *arXiv:2301.13852*.

Michael Sheinman Orenstrakh, Oscar Karnalim, Car-
los Anibal Suarez, and Michael Liut. 2023. Detect-
ing llm-generated text in computing education: A
arXiv preprint arXiv:2307.07411. 646
647
648

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Car-
roll L. Wainwright, Pamela Mishkin, Chong Zhang,
Sandhini Agarwal, Katarina Slama, Alex Ray, John
Schulman, Jacob Hilton, Fraser Kelton, Luke Miller,
Maddie Simens, Amanda Askell, Peter Welinder,
Paul Christiano, Jan Leike, and Ryan Lowe. 2022.
[Training language models to follow instructions with](#)
[human feedback.](#) 649
650
651
652
653
654
655
656

German I. Parisi, Ronald Kemker, Jose L. Part, Christo-
pher Kanan, and Stefan Wermter. 2019. [Continual](#)
[lifelong learning with neural networks: A review.](#)
Neural Networks, 113:54–71. 657
658
659
660

Rajkumar Ramamurthy, Prithviraj Ammanabrolu,
Kianté Brantley, Jack Hessel, Rafet Sifa, Christian
Bauckhage, Hannaneh Hajishirzi, and Yejin Choi.
2022. Is reinforcement learning (not) for natural
language processing?: Benchmarks, baselines, and
building blocks for natural language policy optimiza-
tion. *arXiv preprint arXiv:2210.01241*. 661
662
663
664
665
666
667

Roger Ratcliff. 1990. Connectionist models of recog-
nition memory: constraints imposed by learning
and forgetting functions. *Psychological review*,
97(2):285. 668
669
670
671

Juan Rodriguez, Todd Hay, David Gros, Zain Shamsi,
and Ravi Srinivasan. 2022. Cross-domain detection
of gpt-2-generated technical text. In *Proceedings of*
the 2022 Conference of the North American Chap-
ter of the Association for Computational Linguistics:
Human Language Technologies, pages 1213–1233. 672
673
674
675
676
677

Jürgen Rudolph, Samson Tan, and Shannon Tan. 2023.
Chatgpt: Bullshit spewer or the end of traditional
assessments in higher education? *Journal of Applied*
Learning and Teaching, 6(1). 678
679
680
681

Laura Ruis, Akbir Khan, Stella Biderman, Sara Hooker,
Tim Rocktäschel, and Edward Grefenstette. 2022.
[Large language models are not zero-shot communi-](#)
[cators.](#) 682
683
684
685

Julian Salazar, Davis Liang, Toan Q Nguyen, and Katrin
Kirchhoff. 2019. Masked language model scoring.
arXiv preprint arXiv:1910.14659. 686
687
688

Maarten Sap, Ronan LeBras, Daniel Fried, and Yejin
Choi. 2022. Neural theory-of-mind? on the limits
of social intelligence in large lms. *arXiv preprint*
arXiv:2210.13312. 689
690
691
692

Jaromir Savelka, Arav Agarwal, Christopher Bogart,
Yifan Song, and Majd Sakr. 2023. Can generative pre-
trained transformers (gpt) pass assessments in higher
education programming courses? *arXiv preprint*
arXiv:2303.09325. 693
694
695
696
697

699 Timo Schick and Hinrich Schütze. 2020. [It's not just](#)
700 [size that matters: Small language models are also](#)
[few-shot learners](#). *CoRR*, abs/2009.07118:751

701 G P Shchedrovitsky. 1995. *Selected works*. Shola Kul-
702 turnoï Politiki.

703 Emily Sheng, Kai-Wei Chang, Premkumar Natarajan,
704 and Nanyun Peng. 2019. [The woman worked as a](#)
705 [babysitter: On biases in language generation](#).

706 Ruixiang Tang, Yu-Neng Chuang, and Xia Hu. 2023.
707 The science of detecting llm-generated texts. *arXiv*
708 *preprint arXiv:2303.07205*.

709 Ruixiang Tang, Xiaotian Han, Xiaoqian Jiang, and Xia
710 Hu. Does synthetic data generation of llms help
711 clinical text mining?

712 Wilson L Taylor. 1953. “cloze procedure”: A new
713 tool for measuring readability. *Journalism quarterly*,
714 30(4):415–433.

715 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier
716 Martinet, Marie-Anne Lachaux, Timothée Lacroix,
717 Baptiste Rozière, Naman Goyal, Eric Hambro,
718 Faisal Azhar, et al. 2023. Llama: Open and effi-
719 cient foundation language models. *arXiv preprint*
720 *arXiv:2302.13971*.

721 Tomer Ullman. 2023. Large language models fail on
722 trivial alterations to theory-of-mind tasks. *arXiv*
723 *preprint arXiv:2302.08399*.

724 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob
725 Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz
726 Kaiser, and Illia Polosukhin. 2017. [Attention is all](#)
727 [you need](#). *CoRR*, abs/1706.03762.

728 Lev S Vygotsky. 2012. *Thought and language*. MIT
729 press.

730 Deliang Wang, Dapeng Shan, Yaqian Zheng, Kai Guo,
731 Gaowei Chen, and Yu Lu. 2023. [Can chatgpt detect](#)
732 [student talk moves in classroom discourse? a pre-](#)
733 [liminary comparison with bert](#). In *Proceedings of*
734 *the 16th International Conference on Educational*
735 *Data Mining*, page 515–519. International Educa-
736 tional Data Mining Society.

737 Jeffrey Watumull and Noam Chomsky. 2020. Rethink-
738 ing universality. *Syntactic architecture and its conse-*
739 *quences II*, page 3.

740 Laura Weidinger, John Mellor, Maribeth Rauh, Conor
741 Griffin, Jonathan Uesato, Po-Sen Huang, Myra
742 Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh,
743 Zac Kenton, Sasha Brown, Will Hawkins, Tom
744 Stepleton, Courtney Biles, Abeba Birhane, Julia
745 Haas, Laura Rimell, Lisa Anne Hendricks, William
746 Isaac, Sean Legassick, Geoffrey Irving, and Iason
747 Gabriel. 2021. [Ethical and social risks of harm from](#)
748 [language models](#).

J Nan Wilkenfeld, Bei Yan, Jujun Huang, Guirong Luo,
and Kristina Algas. 2022. “ai love you”: Linguistic
convergence in human-chatbot relationship develop-
ment. In *Academy of Management Proceedings*, vol-
ume 2022, page 17063. Academy of Management
Briarcliff Manor, NY 10510. 752
753
754

Yin Zhang, Rong Jin, and Zhi-Hua Zhou. 2010. [Un-](#)
755 [derstanding bag-of-words model: a statistical frame-](#)
756 [work](#). *International Journal of Machine Learning*
757 *and Cybernetics*, 1(1):43–52. 758