

Do LLMs write like humans? Variation in grammatical and rhetorical styles

Alex Reinhart^{a,1}, Ben Markey^b, Michael Laudenbach^c, Kachatad Pantusen^{a,d}, Ronald Yurko^a, Gordon Weinberg^a, and David West Brown^b

This manuscript was compiled on August 25, 2025

Large language models (LLMs) are capable of writing grammatical text that follows instructions, answers questions, and solves problems. As they have advanced, it has become difficult to distinguish their output from human-written text. While past research has found some differences in features such as word choice and punctuation, and developed classifiers to detect LLM output, none has studied the rhetorical styles of LLMs. Using several variants of Llama 3 and GPT-4o, we construct two parallel corpora of human- and LLM-written texts from common prompts. Using Douglas Biber's set of lexical, grammatical, and rhetorical features, we identify systematic differences between LLMs and humans and between different LLMs. These differences persist when moving from smaller models to larger ones, and are larger for instruction-tuned models than base models. This observation of differences demonstrates that despite their advanced abilities, LLMs struggle to match human stylistic variation. Attention to more advanced linguistic features can hence detect patterns in their behavior not previously recognized.

corpus linguistics | large language models | natural language processing | writing style

As large language models (LLMs) have advanced in recent years, from “stochastic parrots” to models evidently capable of performing complex tasks, most attention has focused on their reasoning performance: solving mathematical problems, writing code, evaluating arguments, diagnosing diseases, and so on (1–4). While past research has studied their mastery of basic grammar and vocabulary (5), there is relatively little research on their language performance more generally: their ability to produce readable text in a variety of styles. Rather than exploring it in detail, commentators discuss the business and communication tasks that might be automated by LLMs with their writing ability, or consider the dangers of impersonation and misinformation facilitated by LLMs (6–9). As more and more writing tasks are automated, such problems appear inevitable.

However, the impression that LLMs write “like humans” is based primarily on qualitative evaluation of their output, not on thorough linguistic evaluation of their text. So far, quantitative comparisons have looked mainly at basic grammar and syntax (5) or features such as word choice, punctuation, sentence length, and so on, finding evidence of some differences between human- and LLM-written text (10–14). Other work has used these features, or language models trained on sample texts, to classify LLM-written texts with varying degrees of success (15–17). Though not definitive, these results suggest there are indeed structural differences between human- and LLM-written text.

We used several recent LLMs (OpenAI's GPT-4o and GPT-4o Mini, and four variants of Meta Llama 3) to generate text from prompts drawn from a large, representative corpus of English, allowing us to directly compare the style of LLM writing to human writing. We find large differences in grammatical, lexical, and stylistic features, demonstrating that LLMs prefer specific grammatical structures and struggle to match the stylistic variation present in human communication, particularly as that variation aligns with the conventions that structure genres such as academic writing, interactive speech, or journalistic news. In Llama 3, where we are able to compare base models (which produce text completions) to instruction-tuned variants (which have been further trained to answer questions and complete tasks specified in prompts), we further see that the instruction tuning introduces more extreme grammatical differences, making them easier to distinguish from human writing and introducing features similar to those present in GPT-4o and GPT-4o Mini.

For example, the instruction-tuned LLMs used present participial clauses at 2 to 5 times the rate of human text, such as in this sentence from GPT-4o using

Significance Statement

As large language models (LLMs) have grown in power and become more widely available, research has focused on their ability to complete various tasks and the biases they exhibit when doing so. In this study, we instead examine their writing style in detail. We show that instruction-tuned models, which are trained to answer questions and solve problems, have a distinct noun-heavy, informationally dense writing style, even when prompted to match the style of informal speech and writing. These findings suggest that instruction-tuned models generate text that does not align with genre conventions familiar to human audiences, and demonstrate the value of linguistic variables in evaluating the output of LLMs.

Author affiliations: ^aDepartment of Statistics & Data Science, Carnegie Mellon University, Pittsburgh, PA 15213; ^bDepartment of English, Carnegie Mellon University, Pittsburgh, PA 15213; ^cDepartment of Humanities & Social Sciences, New Jersey Institute of Technology, Newark, NJ 07102; ^dHeinz College of Information Systems and Public Policy, Carnegie Mellon University, Pittsburgh, PA 15213

A.R., D.W.B., B.M., M.L., R.Y., and G.W. designed research; A.R., D.W.B., K.P., and R.Y. performed research; A.R., D.W.B., B.M., M.L., R.Y., and G.W. wrote the paper.

The authors declare no competing interests.

¹To whom correspondence should be addressed. E-mail: areinhar@stat.cmu.edu

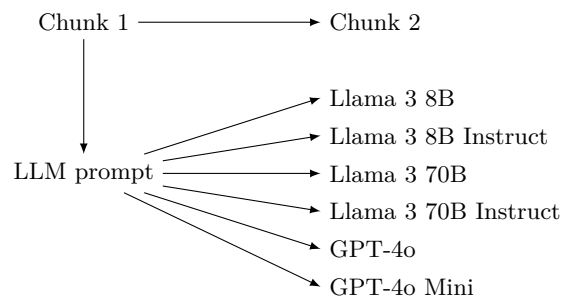


Fig. 1. The LLM text generation workflow. Each human text was split into two chunks of roughly 500 words; the first chunk was used to prompt an LLM to create text that was compared to the second human chunk.

two present participles: “Bryan, *leaning on his agility*, dances around the ring, *evading Show’s heavy blows*.” They also use nominalizations at 1.5 to 2 times the rate of humans, such as in this sentence from Llama 3 70B Instruct containing four: “These schemes can help to reduce *deforestation*, *habitat destruction*, and *pollution*, while also promoting sustainable *consumption* patterns.” On the other hand, GPT-4o uses the agentless passive voice at roughly half the rate as human texts—but in each case, the Llama base models use these features at rates more closely matching humans. This suggests that instruction tuning, rather than training the models to write even more like humans, instead trains them in a particular informationally dense, noun-heavy style, and limits their ability to mimic other writing styles leading to, in some cases, genre misalignment.

These results demonstrate the value of attending to linguistic structure (morphosyntactic, functional, and rhetorical) in order to better understand the affordances and outputs of large language models. Since the rise of the Internet and the concurrent development of efficient processing architectures, language modeling has relied on relatively simple linguistic principles (i.e., sequences and context windows) as the availability of massive amounts of text allowed models trained on text to rapidly outperform older paradigms based on linguistic theory (18); but linguistic theory can provide better ways to evaluate LLM output, just as improved benchmark problems can provide better ways to evaluate their reasoning ability. These results also demonstrate the limits of current LLMs in matching human language, showing that despite their apparent ability, they have measurable limitations compared to human authors.

Methods

We created two corpora of parallel human- and LLM-written texts. Each corpus began with $n = 12,000$ human-authored English texts from a range of genres, from spoken word (such as podcast transcripts) to news and magazine articles to formal academic writing. Language use varies in relation to situational factors such as audience and purpose (19, 20); by including multiple genres, we aimed to capture a diverse range of language production. As illustrated in Figure 1, from each text we extracted two consecutive chunks of roughly 500 words (split at sentence boundaries). The first chunk was provided to each LLM to give it context and a sample of the writing style. The LLMs were prompted to write 500 more words in the same style, tone, and diction; their generated

text was then compared to the next 500-word chunk of the human text. We used six LLMs: GPT-4o and GPT-4o Mini (21, 22) and Meta Llama 3 8B, 70B, 8B Instruct, and 70B Instruct (23), producing six LLM-authored texts for each human-authored text. See the SI Appendix for detailed prompt information.

We constructed the first corpus, the Human-AI Parallel English corpus, from six categories of text (academic, news, fiction, spoken word, blogs, and TV/movie scripts). The second corpus, the COCA AI Parallel (CAP) Corpus, is drawn from the pre-existing Corpus of Contemporary American English (COCA), a large, representative corpus of over 1 billion words in eight registers: spoken, fiction, magazines, newspapers, academic, blogs, web pages, and TV/movie subtitles (24). The HAP-E corpus was used for our primary analyses, while CAP was used to evaluate the generalizability of the results to different texts. The LLMs sometimes refused to respond to prompts or gave short, unusable answers; after these were removed, there were $n = 8,290$ HAP-E texts and $n = 9,615$ CAP texts with outputs from all LLMs. With two human chunks and six LLM-authored chunks for each text, HAP-E comprised $n = 66,320$ chunks and CAP $n = 76,920$ chunks. See SI Appendix Tables S2–S3 for corpus size and composition.

To extract meaningful features from our corpus for training our classifiers, we used Douglas Biber’s tagset of 66 linguistic categories (19, 25, 26), which includes indices of lexical complexity and raw linguistic features ranging from the lexical to the grammatical. For example, features include mean word length, the use of nominalizations (nouns formed from adjectives or verbs, such as *development* or *robustness*), agentless passive voice, hedging phrases (such as *something like* or *almost*), and clausal coordination. All features are listed in the SI Appendix, Table S4. Differences between LLM and human use of features were tested for statistical significance with the paired Wilcoxon signed-rank test with Bonferroni multiple comparison correction.

As a further check of generalizability, we used part of the M4 parallel corpus (17) consisting of abstracts from the arXiv preprint service alongside abstracts generated by GPT-3.5 when prompted with the preprint title. These texts are from a different LLM (GPT-3.5) and a very distinct genre of writing (academic abstracts), providing a check on the consistency of results in different genres.

Results

Classifying text by source. A random forest classifier using the Biber features to distinguish between the seven text sources in HAP-E (human chunk 2 and the six LLMs) achieved a test accuracy of 66%, compared to an expected accuracy of 14% from random guessing. The confusion matrix, shown in Figure 2, demonstrates that little of the error was due to confusion between human texts and the LLMs: instead, most classification errors confused Llama 3 8B and 70B, Llama 3 8B Instruct and 70B Instruct, or GPT-4o and 4o Mini. Each pair consists of models of two different sizes trained on similar data, implying that the size difference does not produce dramatically different style. Overall, only 4.2% of LLM texts were falsely classified as human, and only 9.8% of human texts were falsely classified as LLMs.

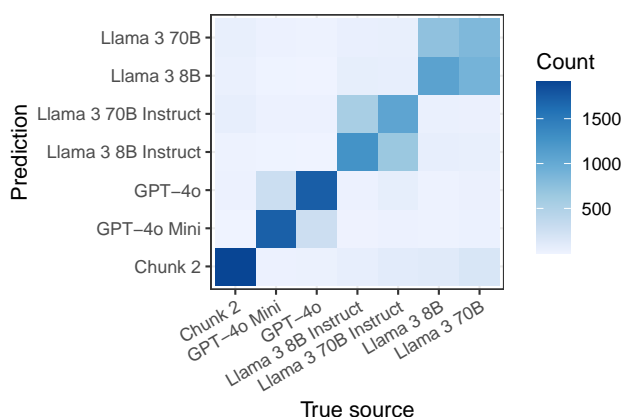


Fig. 2. Confusion matrix for a random forest classifying HAP-E texts by their linguistic and rhetorical features, evaluated on the test set (25% of the HAP-E corpus, including $n = 14$, 535 human and LLM texts). The block diagonal structure indicates that most classification errors were between different versions of the same LLM, rather than between humans and LLMs.

Differences in style and vocabulary. Figure 3 illustrates the large variation in rate of occurrence of the fifteen most important features (as identified by the random forest) in texts generated by LLMs, relative to the rate observed in the human text. All four instruction-tuned models have strong preferences for present participial clauses, ‘that’ clauses as subjects, nominalization, and phrasal co-ordination, which are typical markers of more informationally dense, noun-heavy style of writing (27). For example, GPT-4o uses present participial clauses at 5.3 times the rate of humans (paired Cohen’s $d = 1.38$), ‘that’ clauses as subject 2.6 times as often ($d = 0.77$), nominalizations 2.1 times as often ($d = 1.23$), and phrasal coordination 1.9 times as often ($d = 0.81$). (Rates and effect sizes for all 66 features are provided in the SI Appendix, Tables S5–S6; Figure S1 illustrates paired differences.) There are also signs of local patterns that emerge with specific models: both GPT-4o models avoid clausal coordination, while all Llama 3 variants use it more frequently than humans; while both GPT-4o models use downtoners (such as *barely* or *nearly*) more frequently than humans, all Llama 3 variants avoid them.

One might expect larger models to better match human text than smaller models (e.g., Llama 70B versus Llama 8B, or GPT-4o versus GPT-4o Mini), but this does not appear to be the case in Figure 3. Also, instruction tuning appears to make the model output less human, not more: the Llama 3 base models use features at rates similar to human texts, while GPT-4o and Llama 3 instruction-tuned models have much wider variation from feature to feature.

Similar to past research (14), we find that LLMs also favor specific vocabulary. Figure 4 shows the rate of usage for words used more than once per million words by humans, comparing the usage of each LLM to the usage by humans in Chunk 2 of HAP-E. Compared to the base Llama models, in the instruction-tuned Llama and GPT-4o models certain words are used at dramatically higher and lower rates. Table 1 highlights words overrepresented in LLM outputs: GPT-4o and 4o Mini use words like *camaraderie*, *palpable*, *tapestry*, and *intricate* at more than 100 times the rate of humans, such as in the GPT-4o output phrase “The camaraderie was palpable.”

As a result, “tapestry” appeared in 23% of GPT-4o outputs and “amidst” in 27% (SI Appendix, Table S7). Instruction-tuned variants of Llama 3 also favor words like *camaraderie* and *palpable*, as well as *unease* and *reminder*, though at lower rates than GPT-4o and in a much smaller fraction of documents.* Conversely, they use certain obscurities more than 100 times less often (SI Appendix Table S8).

While many words listed in Table 1 may be occasionally expected in belletristic works of fiction, their pervasiveness across LLM output in a diverse array of genres is notable. To those familiar with academic writing, newspapers, or television scripts, these words are largely unexpected, and to experts likely signal an overwritten, sentimental, or simply uneven text. The point here is not that humans refrain from using these words, but that humans refrain from using these words in certain genres. In this case, words that are unremarkable in fiction are highly conspicuous and unconventional when used other genres. As word choice appears most similar to humans for the base models, this suggests the word choice bias is introduced by the instruction tuning process, not simply by bias in the texts composing the training sets.

In the GPT-4o models in particular, many of these words connote some form of complex relation among objects (e.g., *tapestry*, *intricate*, *camaraderie*, *cacophony*, *amidst*). Coupled with positive items such as *vibrant* and *solace*, these words together may signal a preference for grandiose, if hollow, summative sentences.

Distinguishing individual LLMs. When classifying between human-generated text and one specific LLM, rather than comparing all LLMs, our classifiers achieve much higher accuracy. Typical accuracies achieved by random forests were 93–98% even when trained on HAP-E and tested on CAP or vice versa (SI Appendix Table S9). Lasso-penalized logistic regression classifiers attained similar performance for all LLMs except for the Llama base models, which had accuracies around 75% (SI Appendix Table S10). Since the lasso regressions only consider additive terms, this implies that interactions between the Biber features contain relevant signals for the Llama base models.

For both methods, the lower classification accuracy for the Llama base models relative to the GPT-4o and instruction-tuned Llama models indicates that instruction tuning may lead to writing that is easier to distinguish from human writing.

Generalization across corpora. When each pairwise random forest was used to classify arXiv preprints from the M4 corpus, accuracy dropped significantly. Random forests trained on instruction-tuned LLMs were able to classify M4’s GPT-3.5 output with greater-than-chance accuracy, but models trained on the Llama base outputs attained only 50% accuracy, equal to random guessing (SI Appendix Tables S11–S12). These results demonstrate that instruction-tuned LLMs do have features in common that permit their classification, but that generalizability across LLMs or to different registers of text is difficult.

* Some overuse may be artifacts of the generation process; for example, Llama 3 instruction-tuned variants overuse *continuation* because their outputs sometimes begin with “Here is the continuation of the text...” Llama base models have a tendency to repeat themselves, so Llama 3 8B uses *Deborah* at 52 times the rate of humans largely because of a single document repeating it 308 times.

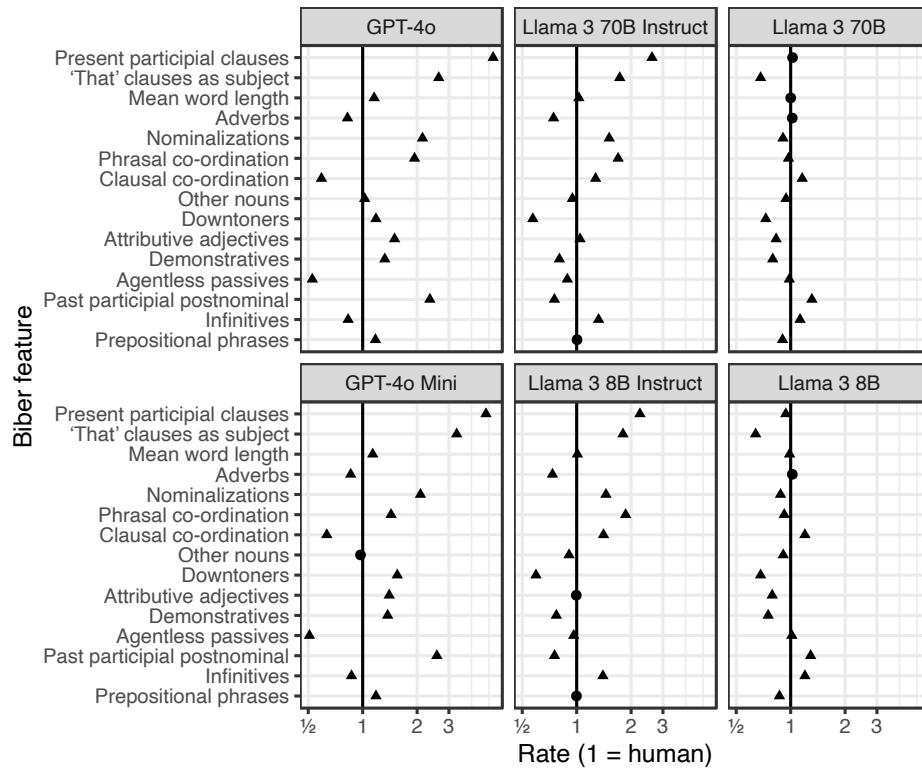


Fig. 3. Rate of Biber feature use by different LLMs, relative to the human usage of each feature, for the top 15 most important features in the HAP-E corpus. Note the log scale. GPT-4o and GPT-4o Mini show the largest variation from human texts, while the base variants of Llama 3 most closely resemble human grammar and style. Larger models (top row) generally show the same stylistic differences as their smaller counterparts (bottom row), despite performing better on other benchmark tasks. Triangles indicate statistically significant differences from human usage. Figure S1 of the SI Appendix gives the distributions of paired differences.

Table 1. Most overrepresented words in LLM-generated texts, relative to human usage rates

GPT-4o		GPT-4o Mini		Llama 3 70B Instruct		Llama 3 8B Instruct		Llama 3 70B		Llama 3 8B	
Word	Rate	Word	Rate	Word	Rate	Word	Rate	Word	Rate	Word	Rate
camaraderie	162	camaraderie	171	unease	63	unease	101	bananas	31	deborah	52
tapestry	155	tapestry	147	palpable	47	continuation	52	paperback	30	rambo	22
intricate	119	palpable	145	continuation	29	palpable	48	bam	26	matty	20
underscore	107	grapple	131	shoutout	28	reminder	33	verona	25	goodnight	18
unspoken	102	intricate	129	intricate	27	pang	29	filth	19	ml	15
amidst	100	fleeting	124	pang	25	rut	29	rekall	17	merlin	13
palpable	95	ignite	122	camaraderie	24	waft	28	denis	14	worcester	11
solace	95	vibrant	92	polycymaker	24	prioritize	27	darry	12	fay	10
fleeting	84	amidst	90	prioritize	24	grapple	24	ebook	12	missy	10
unravel	83	cacophony	89	reminder	24	camaraderie	23	janice	12	elisa	10

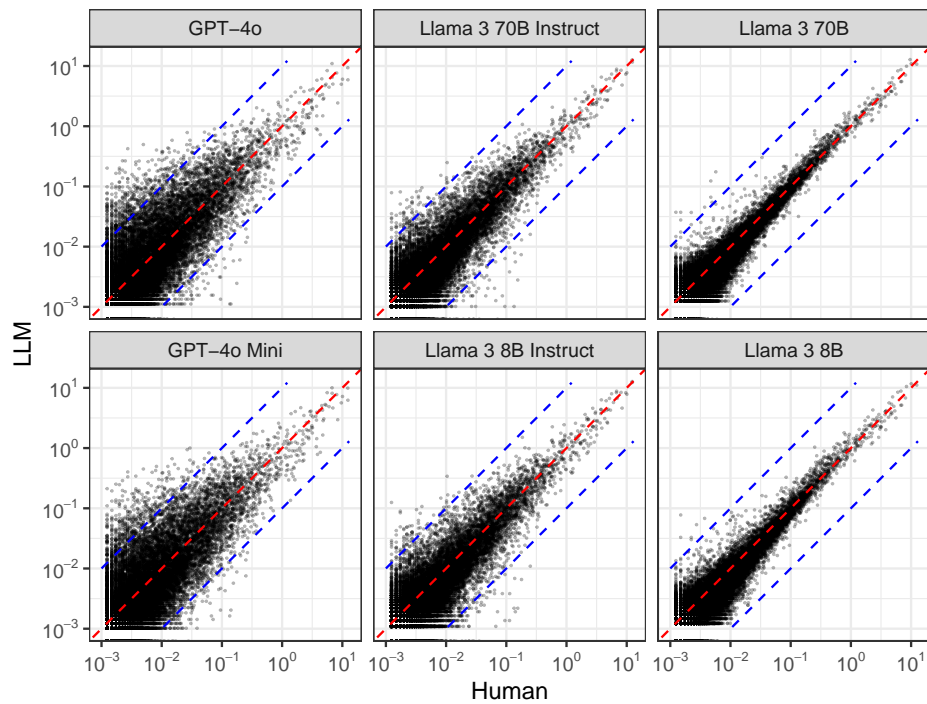


Fig. 4. Rates of word use by different LLMs (per 1,000 words) compared to the human use of each word in chunk 2, in the HAP-E corpus (log scale). Includes all words used more than once per million words in chunk 2. Words are lemmatized to group together inflected forms. Words on the diagonal are used equally often in human and LLM texts. Dashed blue lines indicate the range between $10\times$ more and $10\times$ less than human use. Note that the instruction-tuned models show more variation from the diagonal, indicating more deviation in vocabulary use relative to humans.

Discussion

This study identifies salient differences both between human and LLM-generated texts and among various models. The features that distinguish between humans and different LLMs include present participial clauses, ‘that’ clauses as sentence subject, passive voice, and nominalizations, to name a few. These findings corroborate other research that points out the ways these features produce informationally dense prose (28). Other research links these features to increased lexical diversity in generated text, as well as human judgments of linguistic mastery (13). Lastly, prior work found that ChatGPT-4 text evidences more nominalizations, and fewer human subjects and epistemic stance markers (29), findings we see reproduced in our list of distinguishing features.

A second major finding of this research is the apparently central role of instruction tuning in creating these discrepancies between human and model general texts. While we do not have access to untuned versions of GPT, comparisons between Llama’s base and tuned models emphasize the degree to which instruction tuning pushes models to produce text that reads *unlike* a human. This suggests that differences in style are not simply due to the selection of texts for training the base models, but due to the instruction-tuning process. Similarly, differences between GPT-4o and the instruction-tuned Llama variants may be due to differences in instruction tuning, either through different human preferences in rating responses or differences in the tasks (such as summarization) used to tune the models. (As the instruction tuning processes are not publicly documented, it is not possible to determine the cause more precisely.) While instruction tuning has previously been shown to introduce cognitive biases (30), to our knowledge, these changes in writing style are not discussed elsewhere in similar research.

A third major finding of this work is significant success of Biber’s tagset in modeling and classifying text. This success suggests that varied linguistic perspectives—which, perhaps, are not prioritized during the development and in-house assessments of LLMs—can reveal otherwise tacit information that distinguishes a text as machine-generated. With the linguistic perspective offered by pseudobibeR, we built a model that recognizes machine-generated text with relative ease. Our study reveals the clear value of linguistics expertise and functional conceptions of language in both LLM use and development.

That said, our intention is not to propose another way to construct LLM detectors or to police the writing of students and learners. Instead, we maintain that this type of comparative analysis is useful for identifying differences between human- and machine-generated text, zeroing in on specific teachable moments in the revision of machine-generated text.

As LLMs are increasingly put to work completing diverse writing tasks, these results suggest a notable misalignment between generated texts and the contexts in which we put them to use. This is another way of saying that LLMs do not vary their linguistic output in response to contextual factors in ways similar to humans. This misalignment affects experts and learners differently. For those proficient in a genre—think, a therapist collating notes on a patient or, say, a college graduate writing thank-you notes to friends and family—this misalignment is likely flagged and the output is

appropriately revised. When the writer is proficient in the genre, previous experience guides a current sense of what a particular document should look like in order to be successful. For experts, then, LLMs appear a worthwhile productivity tool, suitable so long as they lend their expertise to further shaping the output.

For learners, though, LLMs appear more problematic. Of course, using LLMs to learn more about a concept, or help generate ideas, is one thing. Using output in a text is a different matter, one that may affect a student’s learning trajectories. In this case, students offload the important cognitive labor of shaping a text for a particular audience and purpose to the LLM. Never mind that LLMs do not appear to write like humans—when students offload writing work, they offload opportunities to learn how to write IMRaD articles, client reports, executive summaries, investment pitches, etc. When LLMs are used in the classroom as writing tools, instructors of all levels and disciplines need to help students see both the shortcomings of the generated text and avenues for improvement. LLMs are not bad, either technically or morally—it is only that instructors must help inculcate in students the critical perspective of the expert to know what’s working and what isn’t.

In other contexts, however, overreliance on LLMs could produce output that might be awkward and inauthentic (e.g., in a creative genre), confusing (e.g., in instructional material), or unpersuasive (e.g., in argumentative texts). The current work thus suggests the importance of LLM practice—both in and out of the classroom—informed by human expertise via a continual dialogue of creation and revision, where LLM users are more aware and mindful of the effective uses as well as the limitations of various LLMs.

ACKNOWLEDGMENTS. We thank members of the TeachStat Research Group for helpful discussions, the Dietrich College of Humanities and Social Sciences at Carnegie Mellon University for use of the Wright GPU cluster, and Aadi Menon for exploring suitable prompts.

1. B Wang, X Yue, H Sun, Can ChatGPT defend its belief in truth? Evaluating LLM reasoning via debate in *Findings of the Association for Computational Linguistics: EMNLP 2023*, eds. H Bouamor, J Pino, K Bali. (Association for Computational Linguistics, Singapore), pp. 11865–11881 (2023).
2. J Huang, KCC Chang, Towards reasoning in large language models: A survey in *Findings of the Association for Computational Linguistics: ACL 2023*, eds. A Rogers, J Boyd-Graber, N Okazaki. (Association for Computational Linguistics, Toronto, Canada), pp. 1049–1065 (2023).
3. Y Chen, TX Liu, Y Shan, S Zhong, The emergence of economic rationality of GPT. *Proc. Natl. Acad. Sci.* **120**, e2316205120 (2023).
4. SA Lehr, A Caliskan, S Liyanage, MR Banaji, ChatGPT as research scientist: Probing GPT’s capabilities as a research librarian, research ethicist, data generator, and data predictor. *Proc. Natl. Acad. Sci.* **121**, e2404328121 (2024).

5. TA Chang, BK Bergen, Language model behavior: A comprehensive survey. *Comput. Linguist.* **50**, 293–350 (2024).
6. D Barman, Z Guo, O Conlan, The dark side of language models: Exploring the potential of LLMs in multimedia disinformation generation and dissemination. *Mach. Learn. with Appl.* **16**, 100545 (2024).
7. S Kumar, V Balachandran, L Njoo, A Anastasopoulos, Y Tsvetkov, Language generation models can cause harm: So what can we do about it? an actionable survey in *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, eds. A Vlachos, I Augenstein. (Association for Computational Linguistics, Dubrovnik, Croatia), pp. 3299–3321 (2023).
8. B Kovács, The Turing test of online reviews: Can we tell the difference between human-written and GPT-4-written online reviews? *Mark. Lett.* **35**, 651–666 (2024).

9. T Hagendorff, Deception abilities emerged in large language models. *Proc. Natl. Acad. Sci.* **121**, e2317967121 (2024).
10. R Tang, YN Chuang, X Hu, The science of detecting LLM-generated text. *Commun. ACM* **67**, 50–59 (2024).
11. L Fröhling, A Zubiaga, Feature-based detection of automated language models: tackling GPT-2, GPT-3 and Grover. *PeerJ Comput. Sci.* **7**, e443 (2021).
12. A Muñoz-Ortiz, C Gómez-Rodríguez, D Vilares, Contrasting linguistic patterns in human and LLM-generated news text. *Artif. Intell. Rev.* **57** (2024).
13. S Herbold, A Hautli-Janisz, U Heuer, Z Kikteva, A Trautsch, A large-scale comparison of human-written versus ChatGPT-generated essays. *Sci. Reports* **13** (2023).
14. W Liang, et al., Monitoring AI-modified content at scale: A case study on the impact of ChatGPT on AI conference peer reviews in *Proceedings of the 41st International Conference on Machine Learning*, Proceedings of Machine Learning Research, eds. R Salakhutdinov, et al. (PMLR), Vol. 235, pp. 29575–29620 (2024).
15. JQJ Liu, et al., The great detectives: humans versus AI detectors in catching large language model-generated medical writing. *Int. J. for Educ. Integr.* **20**, 8 (2024).
16. E Mosca, MHI Abdalla, P Basso, M Musumeci, G Groh, Distinguishing fact from fiction: A benchmark dataset for identifying machine-generated scientific papers in the LLM era. in *Proceedings of the 3rd Workshop on Trustworthy Natural Language Processing (TrustNLP 2023)*, eds. A Oualle, et al. (Association for Computational Linguistics, Toronto, Canada), pp. 190–207 (2023).
17. Y Wang, et al., M4: Multi-generator, Multi-domain, and Multi-lingual Black-Box Machine-Generated Text Detection in *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, eds. Y Graham, M Purver. (Association for Computational Linguistics, St. Julian's, Malta), pp. 1369–1407 (2024).
18. X Li, "There's no data like more data": Automatic speech recognition and the making of algorithmic culture. *Osiris* **38**, 165–182 (2023).
19. D Biber, *Variation across Speech and Writing*. (Cambridge University Press), (1988).
20. C Miller, Genre as social action. *Q. J. Speech* **70**, 151–167 (1984).
21. OpenAI, Hello GPT-4o (2024).
22. OpenAI, GPT-4o mini: advancing cost-efficient intelligence (2024).
23. Meta, Introducing Meta Llama 3: The most capable openly available LLM to date (2024).
24. M Davies, *The Corpus of Contemporary American English (COCA)* (2008).
25. D Biber, *Dimensions of Register Variation: A Cross-Linguistic Comparison*. (Cambridge University Press), (1995).
26. D Biber, S Conrad, *Register, Genre, and Style*. (Cambridge University Press), (2009).
27. LL Aull, *How Students Write: A Linguistic Analysis*. (MLA), (2020).
28. B Markey, DW Brown, M Laudenbach, A Kohler, Dense and disconnected: Analyzing the sedimented style of ChatGPT-generated text at scale. *Writ. Commun.* **41**, 571–600 (2024).
29. FK Jiang, K Hyland, Does ChatGPT argue like students? Bundles in argumentative essays. *Appl. Linguist.* (2024).
30. I Itzhak, G Stanovsky, N Rosenfeld, Y Belinkov, Instructed to bias: Instruction-tuned language models exhibit emergent cognitive bias. *Transactions Assoc. for Comput. Linguist.* **12**, 771–785 (2024).