**Peer Response – Dalbir Singh**

Thanks, Dalbir, for your even-handed overview of what GPT-3 can do and where it worries you. Your points about its potential misuse in classrooms and creative fields really hit home, so I'd like to add some practical steps that could help curb those hazards.

To start, schools and companies could roll out short AI-literacy workshops showing people how to spot, question, and responsibly use machine-generated text. Trainees would learn, as Bender et al. 2021 remind us, that these models don't actually understand language; they simply guess the next word based on billions of patterns. Encouraging critical engagement, especially in academic settings, can reduce overreliance and ensure proper use.

Secondly, institutions should adopt AI usage policies that clearly define acceptable and unacceptable uses—especially in creative assessments, where originality and understanding are being tested. Automated detection tools can also help flag AI-generated submissions.

To address the risk of bias and misinformation, developers and researchers must prioritise transparent dataset documentation and auditing processes. As Weidinger et al. (2021) suggest, these practices can limit the propagation of harmful stereotypes or falsehoods.

Finally, a human-in-the-loop approach should be mandatory for any AI-driven creative or factual output—ensuring a person is accountable for verifying and contextualising the information. When handled carefully, these models can boost efficiency and spark new ideas; put them in the wild with no guardrails, though, and they can scramble learning, distort the truth, and shortchange fairness.

**References**

Bender, E.M. et al. (2021) 'On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?' ACM Conference on Fairness, Accountability, and Transparency, pp. 610–623. Available at: https://dl.acm.org/doi/pdf/10.1145/3442188.3445922.

Hutson, M. (2021) 'Robo-writers: the rise and risks of language-generating AI', Nature. Available at: https://www.nature.com/articles/d41586-021-00530-0.

Weidinger, L. et al. (2021) 'Ethical and social risks of harm from Language Models'. Available at: https://arxiv.org/pdf/2112.04359 .

Zhou, K., Prabhumoye, S. and Neubig, G. (2020) 'Exploring Ethical and Social Implications of Language Models'. Available at: https://arxiv.org/pdf/2010.12884 .

**Peer Response – Fahad Abdallah**

Thanks, Fahad, for the thorough and timely look at the ethical and social hurdles LLMs throw our way. Your observations about trust, misuse of authorship, and Big Techs heavy hand underline how quickly we need smart safeguards.

One obvious safeguard is a clear rule that students and researchers must say when AI helped their writing. Journals and schools could require authors to name the tool-was it GPT-3, Bard, or another-and explain what it did. Such a step fits recent publisher drives to keep research honest and to squeeze out ghost-writing or paper-mill scams (Kendall and Teixeira da Silva, 2024).

Institutions also need their own friendly, workable AI rules. Universities and labs should draft plain guidelines spelling out when, and how, LLMs can show up in homework, theses or papers. Policies might limit blind, unsupervised output, demand solid human review, and teach users to read AI text with a critical eye.

Pushing open-source LLM projects forward is another must-do if we want to fight bias and opacity. As Gibney (2022) points out, a clear model lets teams' peek under the hood, test, tweak, or toss out the faulty parts instead of just trusting a black box. With that kind of setup, outside experts and campus ethics boards could run regular audits to check that any AI-aided work is fair and truthful before it ever sees publication.

Promoting digital literacy and critical-thinking training will help readers, reviewers, and authors tell real research apart from text churned out by AI. An academic community that knows how tools work is harder for bad actors to fool.

When we pair smart rules, better training, and sound technology, we can enjoy what LLMs offer while keeping research trustworthy and serving the wider public good.

**References:**
Gibney, E. (2022) 'Open-source language AI challenges big tech's models', Nature. Available at: https://doi.org/10.1038/d41586-022-01705-z

Kendall, G. and Teixeira da Silva, J.A. (2024) 'Risks of abuse of large language models, like ChatGPT, in scientific publishing: Authorship, predatory publishing, and paper mills', Learned Publishing, 37(1). Available at: https://www.graham-kendall.com/papers/ktds2023a.pdf

**Peer Response – Ali Alshehhi**

Many thanks for your detailed and thought-provoking post, Ali. You capture well the twin sides of AI writing tools: they can turbocharge productivity, but they also slide into ethical and social trouble if we let them run free. Your worries about bias, random glitches, and careless use in sensitive areas definitely hit the mark.

One way to tame those dangers is to build strict validation and auditing steps into the workflow, especially in fields like medicine where a single wrong sentence can be deadly. As Omiye et al. (2024) remind us, large language models should sit in the room as helpers, not as decision-makers, and every high-stakes output needs a careful human check.

To fight algorithmic bias, teams should never stop diversifying training datasets and running bias-detection toolkits as part of the routine. Abburi et al. (2023) show that ensemble classifiers can catch harmful content before it leaves the lab, yet that guardrail only stays sharp if it is updated along with shifting social norms.

In schools, clear AI-usage policies—like always stating when AI was used and capping its role in creative tasks—can help protect the learning space (Anson and Straume, 2022). At the same time, training sessions show students and staff what these tools can and cannot do, while also pointing out the ethical bumps on the path.

Your remark about how hard it is to predict AI behaviour only strengthens the call for clear prompt-writing rules and public-facing documentation. Without consistent outputs, we cannot expect departments or publishers to accept language models in sensitive or high-stakes settings.

Taken together, responsible AI use hinges on ethical guardrails, open release practices, and users who understand both power and limit.

**References**

Abburi, H. et al. (2023) Generative AI Text Classification using Ensemble LLM Approaches , arXiv. Available at: https://arxiv.org/abs/2309.07755.

Anson, C.M. and Straume, I. (2022) Amazement and Trepidation: Implications of AI-Based Natural Language Production for Teaching Writing , Journal of Academic Writing , 12(1), pp. 1-9. Available at: https://doi.org/10.18552/joaw.v12i1.820.

Hutson, M. (2021) Robo-writers: The rise and risks of language-generating AI , Nature. Available at: https://doi.org/10.1038/d41586-021-00530-0.

Lingard, L. (2023) Writing with ChatGPT: An Illustration of Its Capacity, Limitations & Implications for Academic Writers , Perspectives on Medical Education , 12(1), p. 261. Available at: https://doi.org/10.5334/pme.1072.

Omiye, J.A. et al. (2024) Large Language Models in Medicine: The Potentials and Pitfalls: A Narrative Review , Annals of Internal Medicine , 177(2). Available at: https://doi.org/10.7326/M23-2772.