



# Examining emerging capabilities and mitigating potential risks of VLLMs

ONE DAY WORKSHOP

50+ ATTENDANCE (expected)

2 tasks, 3 talks, panel discussion

November 20, 2023

The emergence of very large language models (VLLMs) has dramatically altered the trajectory of progress in AI and its applications. With the release of ChatGPT, that excitement has now transcended boundaries from AI researchers to the common man; indeed it asserts we are living in an exciting time of scientific proliferation. However, we see two divergent community views on VLLMs. Believers in the magical powers of VLLMs claim that VLLMs are autodidactic in learning a wide gamut of new capabilities – referred to as "*emerging capabilities*" in the community – an effect that gets pronounced with model size and dataset scale. On the other hand, critics are not yet prepared to acknowledge the self-learning power of VLLMs; they criticize it as only a statistical learner and out several flaws – *hallucination* being the most prominent one. In this forum, we want to bring together both communities – the believers and the critics, to explore exciting tasks together: (i) **CT<sup>2</sup> - Counter Turing Test: AI-Generated Text Detection**, and (ii) **HELT - Hallucination eLiciaTion through automatic detection and mitigation**. Expect this to be an exciting forum to discuss, debate, and explore exciting scientific pathways for the future.

## 1. Rationale - Why does AI need to be civilized? - Call for Papers (CFP)

Advances in AI during the past couple of years have led to AI systems becoming immensely more powerful than ever before. While their applications for social good cannot be overstated, as an unintended by-product, risks of misuse have also been exacerbated. This prompted an open petition letter ([Marcus and of Life Institute, 2023](#)) (led by [Gary Marcus](#)) by the nonprofit Future of Life Institute, calling for all AI labs to immediately pause for at least 6 months "moratorium" the training of AI systems more powerful than GPT-4. The letter has (18K+ and still counting) signatures from technologists and luminaries, which include [Yoshua Bengio](#),

[Stuart Russell](#), [Elon Musk](#), [Steve Wozniak](#), and [Andrew Yang](#). It also includes policy leaders such as [Rachel Bronson](#), president of the Bulletin of the Atomic Scientists, a science-oriented advocacy group known for its warnings against humanity-ending nuclear war. On the other hand, the opposing campaign, which doesn't believe in halting scientific progress, has powerful people too, including [Bill Gates](#) ([gat, 2023](#)), [Andrew Ng](#) ([aih, 2023](#)), [Yann LeCun](#) ([aih, 2023](#)) and many others. Furthermore, both the United States ([whi, 2023](#)) and the European Union ([eua, 2023](#)) governments have recently proposed regulatory frameworks for AI. This is a significant time in the history of scientific development. In this forum, we will discuss and debate emerging capabilities and mitigating potential risks and limitations of VLLMs. Call for papers includes, but is not limited to:

- *unique emerging abilities of VLLM*;
- *negative, position, and full paper on potential risks of VLLMs*;
- *ethics and VLLMs*;
- *making VLLMs more responsible*;
- *detection AI-generated content*;
- *mitigation of harmful hallucinations*.

## 2. Two Shared Tasks

Shared tasks are an effective way to attract research attention to any emerging area. We will host two shared tasks: (i) **CT<sup>2</sup>**, and (ii) **HELT**. The findings of **CT<sup>2</sup>** will mitigate misusage risks, while **HELT** endeavors to make VLLMs more human-sensitive and responsible.

### 2.1. CT<sup>2</sup> - Counter Turing Test for AI-Generated Text Detection

With the emergence of ChatGPT, the risk of AI-generated content has reached an alarming apocalypse. ChatGPT has been declared banned by the school system in NYC ([Rosenblatt, 2023](#)), Google ads ([Grant and Metz, 2022](#)), and Stack Overflow ([Makyen and Olson, 1969](#)), while scientific conferences like ACL (Chairs, 2023) and ICML ([Founda-](#)

tion, 2023) have released new policies deterring the usage of ChatGPT for scientific writing. After the initial skepticism, ChatGPT has been seen as a listed author in scientific papers (Kung et al., 2023; O'Connor et al., 2022), while Elsevier (Elsevier, 2023) and Springer (Springer, 2023) have adopted more inclusive guidelines on *the use of ChatGPT for scientific writing*.

Indeed, detecting AI-generated text has suddenly emerged as a concern that needs immediate attention. While watermarking as a potential solution to the problem is being studied by OpenAI (Wiggers, 2022b), a handful of systems that detect AI-generated text such as GPT-2 output detector (Wiggers, 2022a), GLTR (Strobel et al., 2022), GPTZero (Tian, 2022), DetectGPT (Mitchell et al., 2023), etc. have recently been orange observed in practical use. To address the inevitable question of ownership attribution for AI-generated artifacts, the US Copyright Office (Office, 2023) released a statement stating that if the content is traditional elements of authorship produced by a machine, the work lacks human authorship and the office will not register it for copyright. Given this censorial spotlight on generative AI, AI-generated text detection is a topic that needs a thorough investigation. In this regard, there are three families of techniques proposed so far:

- **Watermarking:** First introduced in (Aaronson, 2022), watermarking AI-generated text involves embedding an imperceptible code or signal to verify the author of a particular text with certainty. (Kirchenbauer et al., 2023) proposed this by selecting the next token pseudorandomly (rather than simply choosing the one with the highest probability) using a cryptographic pseudorandom function whose key is only possessed by the LLM maker. It would be remiss not to mention the most obvious pitfall of this approach, which is that if the text is altered or modified in any way, detecting the watermark proves to be a difficult task.
- **Negative log likelihood (NLL):** NLL-based implementations such as DetectGPT (Mitchell et al., 2023) have demonstrated the detection of AI-generated text by comparing log-likelihood of generated tokens after perturbing the input text by replacing some tokens with others. If the new, perturbed version of the text lies in the negative curvature regions of log-likelihood, it was likely generated by AI. The limitation of this approach is that it requires access to the log probabilities of the text in order to work which implies that knowledge of which LLM was used to generate the text is essential;
- **Perplexity and Burstiness:** GPTZero (Tian, 2022), an example of a detection technique based on perplexity and burstiness, has demonstrated that a text with lower perplexity (a measure of how

predictable the text is), and with lower burstiness (the measure of how uniform text is) has a high probability of being generated by an AI. The limitations here are that GPTZero also requires access to log probabilities of text as well as the fact that it approximates perplexity values using a linear model.

Although AI-generated text detection has suddenly received immense attention, Liang et al. (Liang et al., 2023) suggest that available AI-generated text detectors consistently misclassify non-native English writing samples as AI-generated, whereas native writing samples are accurately identified, highlighting the ethical implications of deploying AI-generated content detectors and risking misrepresentation. This implies that a community effort is needed to tackle the issue of detectors penalizing under-represented sub-population(s). In our recent publication, "Counter Turing Test  $CT^2$ : AI-Generated Text Detection Challenges" (Chakraborty et al., 2023), we introduce a benchmark for evaluating the robustness of AGTD techniques. Our results clearly show the vulnerabilities of current AGTD methods. As discussions on AI policy intensify, assessing the detectability of content produced by LLMs is vital. To this end, we present the AI Detectability Index (ADI) for quantitative ranking based on detectability.

The  $CT^2$  task will be the first of its kind in bringing together researchers in advancing the area of detecting AI-assisted generated text.

### 2.1.1. Data to be released and the task

$CT^2$  will consist of three sub-tasks. We will be releasing 100K data points, consisting of (i) prompt, (ii) human-written text, and (iii) AI-generated text by 15 different LLMs.

- **Task A:** given a set of human-generated text documents vs. AI-generated text documents participants need to design techniques to detect AI-generated text. Indeed, human-written text vs. AI-generated text would be parallel, which means they will be on the same topic. In this task, we will let participants know that the generated text is from GPT, OPT, BERT, XLNet, etc. As such, this is an LLM-specific AI detection task.
- **Task B:** in this task, we will not tell people which the generated text is using which LLM. Participants need to design techniques which is LLM agnostic.
- **Task C:** In this task, we will offer AI-assisted writing, i.e., AI-generated text interlaced with minor edits by another language model and human, as input. Given the intricacies and challenges of AI-assisted writing, it would be the hardest task to attempt.

## 2.2. **H** **E** **L** **T**: Hallucination eLiciaTion through automatic detection and mitigation

With the recent and rapid advances in the areas of LLMs and Generative AI, the pre-eminent and ubiquitous concern is of hallucination. We release large-scale first-of-its-kind human-annotated data with detailed annotations on - intrinsic vs. extrinsic hallucinations, and degree of hallucination, and we ask participants to come up with either black-box factuality Assessment and/or evidence-based fact-checking. First, we define categories of hallucinations:

- **Intrinsic Hallucination:** Intrinsic hallucination refers to the phenomenon when an LLM generates text that topically slightly deviates from the input and/or has a lack of grounding in reality. For example, given a prompt "*USA on Ukraine war*" an LLM generates "*U.S. President Barack Obama says the U.S. will not put troops in Ukraine*". We can see a clear case of intrinsic hallucination as the US president during the Ukraine-Russia war is Joe Biden, not Barack Obama, contradicting the reality.
- **Extrinsic Hallucination:** We define extrinsic hallucination to be the generated output from an LLM that cannot be verified, or contradicted from the source content provided as prompt. For example, we provide a prompt stating "*North Korea has conducted six underground nuclear tests, and a seventh may be on the way.*", and the resulting output generated by the model was "*The international community has condemned North Korea's nuclear tests, with the UN Security Council imposing a range of sanctions in response. The US and other world powers have urged North Korea to abandon its nuclear weapons program and to comply with international law.*" As the source input makes no reference to the U.S. or U.N. Security Council imposing any sanctions, the claimed imposition of sanctions cannot be verified from the input alone.

Next, based on the degree of hallucination, we categorize it into three levels: **(i) Mild:** At this level, hallucination can be categorized as minor and the generated output may just contain innocuous errors or inconsistencies in the text. The generated text will not significantly impact its coherence, **(ii) Moderate:** Moderate hallucination will involve more significant errors or distortions in the generated text. The text may contain nonsensical phrases or ideas that deter from the topic at hand., **(iii) Alarming:** Alarming hallucinations in LLM entail a drastic deviation from the intended output. Such deviations can also be offensive and incongruous with the desired output. In our recent publication, "The Troubling Emergence of Hallucination in Large Language Models – An Extensive

Definition, Quantification, and Prescriptive Remediations" (Rawte et al., 2023), we offer a fine-grained discourse on profiling hallucination based on its degree, orientation, and category, along with offering strategies for alleviation.

### 2.2.1. Detection Methods

There are mainly two broad ways: **(i) Black-box factuality assessment:** Black-box hallucination detection approaches, such as SelfCheckGPT (Manakul et al., 2023), endeavor to fact-check models without utilizing an external database of facts. In the case of SelfCheckGPT, it has the capacity to fact-check models in a zero-resource fashion by evaluating if the generated facts have similarities or if they contradict each other. **(ii) Evidence-based fact checking:** Evidence-based fact-checking, such as LLM-Augmenter (Peng et al., 2023), leverages the use of task-specific databases or similar reliable resources that contain accurate facts.

### 2.2.2. Data to be released and the task

We will be releasing 20K annotated data manually labelled at sentence level on -i) intrinsic vs. extrinsic hallucination, and ii) degree of hallucination mild, moderate, and alarming.

- **Task A:** Given an AI-generated text and associated prompt the task is to detect hallucination at the sentence level. For the competition purpose overall hallucination detection accuracy will be considered averaged over sub-categories like intrinsic and extrinsic.
- **Task B:** This task is to propose hallucination mitigation techniques. While evidence-based fact-checking is well studied in the fact-checking community (Thorne et al., 2018; Wang, 2017; Garg and Sharma, 2020; Kwiatkowski et al., 2019; Jiang et al., 2020; Gupta and Srikumar, 2021; Onoe et al., 2021; Aly et al., 2021), here in this forum we are mainly interested in black-box factuality assessment. However, teams can use external resources to mitigate hallucination and will be considered in the evidence-based mitigation group. We will request participants to report such experiments in their reports. Since it is hard to assess whether the reformed text still has hallucinations, Task B will mostly be an academic exercise.

## 3. Invited Talks & Panel [tentative]

**Talks:** We wish to have 3 speakers from industry/academia, people who have experience in building LLMs. We (tentatively) plan to invite **Prof. Yejin Choi**, University of Washington and Allen Institute for Artificial Intelligence; (**confirmed**) **Dr. Vinodkumar Prabhakaran**, Senior Research Scientist, Google LLC, co-author of PaLM (Chowdhery et al., 2022); and **Nick Ryder** who is a principal researcher at is the co-author of the GPT3 (Brown et al., 2020), works at OpenAI.

## 4. Workshop Organizers

Name	Web, Email, GScholar	Research Interest & Organizing activities
 Dr. Amitava Das	<a href="#">Web</a> <a href="#">Email</a> <a href="#">Google Scholar</a>	<p>Dr. Amitava Das is a Research Associate Professor at AIISC, UofSC, USA, and an advisory scientist at Wipro AI Labs, Bangalore, India.</p> <p><b>Research interests:</b> Code-Mixing and Social Computing.</p> <p><b>Organizing Activities [selective]</b></p> <ul style="list-style-type: none"> <li>- DEFACTIFY 1.0 @AAAI2022</li> <li>- DEFACTIFY 2.0 @AAAI2023</li> <li>- Memotion @SemEval2020</li> <li>- SentiMix @SemEval2020</li> <li>- Computational Approaches to Linguistic Code-Switching @ LREC 2020</li> <li>- CONSTRAINT @AAAI2021</li> </ul>
 Dr. Amit Sheth	<a href="#">Web</a> <a href="#">Email</a> <a href="#">Google Scholar</a>	<p>Dr. Amit Sheth is the NCR Chair and a Professor of CSE at the University of South Carolina. He founded the Artificial Intelligence Institute, which now has nearly 50 researchers. He is a fellow of IEEE, ACM, AAAS, and AAAI. His significant awards include the 2023 IEEE Wallace McDowell award.</p> <p><b>Research interests:</b> Organized over 50 workshops, given over 50 tutorials, nearly 100 keynotes on Neurosymbolic AI, Knowledge Graphs, NLP/NLU, AI for Social Good, etc.</p> <p><b>Organizing Activities [selective]</b></p> <ul style="list-style-type: none"> <li>- Cysoc2021 @ ICWSM2021 - Emoji2021 @ICWSM2021 - KiLKG 2021 @KGC21</li> </ul>
 Aman Chadha	<a href="#">Web</a> <a href="#">Email</a> <a href="#">Google Scholar</a>	<p>Aman Chadha is an Applied Science Manager at Amazon Alexa AI and a Researcher at Stanford AI.</p> <p><b>Research interests:</b> Multimodal AI, On-device AI, and Human-Centered AI.</p>
 Vinija Jain	<a href="#">Web</a> <a href="#">Email</a> <a href="#">Google Scholar</a>	<p>Vinija Jain is a Machine Learning Lead at Amazon Music and a Researcher at Stanford AI.</p> <p><b>Research interests:</b> Recommender Systems and NLP.</p>

## 5. Bibliographical References

2023. Bill gates says calls to pause ai won't 'solve challenges'.
2023. Blueprint for an ai bill of rights: Making automated systems work for the american people.
2023. Proposal for a regulation of the european parliament and of the council laying down harmonised rules on artificial intelligence (artificial intelligence act) and amending certain union legislative acts.
2023. Yann lecun and andrew ng: Why the 6-month ai pause is a bad idea.
- Scott Aaronson. 2022. My projects at openai.
- Rami Aly, Zhijiang Guo, Michael Schlichtkrull, James Thorne, Andreas Vlachos, Christos Christodoulopoulos, Oana Cocarascu, and Arpit Mittal. 2021. Feverous: Fact extraction and verification over unstructured and structured information.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In Advances in Neural Information Processing Systems, volume 33, pages 1877–1901. Curran Associates, Inc.
- Program Chairs. 2023. Acl 2023 policy on ai writing assistance.
- Megha Chakraborty, S.M Towhidul Islam Tonmoy, S M Mehedi Zaman, Krish Sharma, Niyar R Barman, Chandan Gupta, Shreya Gautam, Tanay Kumar, Vinija Jain, Aman Chadha, Amit P. Sheth, and Amitava Das. 2023. Counter turing test ct<sup>2</sup> : Ai – generated text detection is not as easy as you may think – introducing a detectability index.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. Palm: Scaling language modeling with pathways.
- Elsevier. 2023. The use of ai and ai-assisted technologies in scientific writing.
- Neural Information Processing Systems Foundation. 2023. [link].
- Sonal Garg and Dilip Kumar Sharma. 2020. New politifact: A dataset for counterfeit news. pages 17–22. 2020 9th International Conference System Modeling and Advancement in Research Trends (SMART).
- Nico Grant and Cade Metz. 2022. A new chat bot is a 'code red' for google's search business.
- Ashim Gupta and Vivek Srikumar. 2021. X-FACT: A New Benchmark Dataset for Multilingual Fact Checking. arxiv.
- Yichen Jiang, Shikha Bordia, Zheng Zhong, Charles Dognin, Maneesh Singh, and Mohit Bansal. 2020. HoVer: A dataset for many-hop fact extraction and claim verification. Findings of the Conference on Empirical Methods in Natural Language Processing (EMNLP).
- John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. 2023. A watermark for large language models.
- Tiffany H Kung, Morgan Cheatham, Arielle Medenilla, Czarina Sillos, Lorie De Leon, Camille Elepaño, Maria Madriaga, Rimel Aggabao, Giezel Diaz-Candido, James Maningo, et al. 2023. Performance of chatgpt on usmle: Potential for ai-assisted medical education using large language models. PLoS digital health, 2(2):e0000198.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav

- Petrov. 2019. Natural questions: A benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:452–466.
- Weixin Liang, Mert Yuksekgonul, Yining Mao, Eric Wu, and James Zou. 2023. Gpt detectors are biased against non-native english writers.
- Mod Makyen and Peter Olson. 1969. *Temporary policy: Chatgpt is banned*.
- Potsawee Manakul, Adian Liusie, and Mark JF Gales. 2023. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models.
- Gary Marcus and Future of Life Institute. 2023. *Pause giant ai experiments: An open letter*.
- Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D Manning, and Chelsea Finn. 2023. Detectgpt: Zero-shot machine-generated text detection using probability curvature. *arXiv preprint arXiv:2301.11305*.
- Siobhan O'Connor et al. 2022. Open artificial intelligence platforms in nursing education: Tools for academic progress or abuse? *Nurse Education in Practice*, 66:103537–103537.
- Copyright Office. 2023. *Copyright registration guidance: Works containing material generated by artificial intelligence*. Library of Congress.
- Yasumasa Onoe, Michael J. Q. Zhang, Eunsol Choi, and Greg Durrett. 2021. *Creak: A dataset for commonsense reasoning over entity knowledge*. arXiv.
- Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng, Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou Yu, and Weizhu Chen and Jianfeng Gao. 2023. Check your facts and try again: Improving large language models with external knowledge and automated feedback.
- Vipula Rawte, Swagata Chakraborty, Anubhav Sarkar Agnih Pathak, S.M Towhidul Islam Tonmoy, Amit P. Sheth Aman Chadha, and Amitava Das. 2023. *The troubling emergence of hallucination in large language models – an extensive definition, quantification, and prescriptive remediations*.
- Kalhan Rosenblatt. 2023. *Chatgpt banned from new york city public schools' devices and networks*. NBCUniversal News Group.
- Springer. 2023. *Guidance on the use of large language models (llm) e.g. chatgpt*.
- Hendrik Strobelt, Sebastian Gehrmann, and Alexander Rush. 2022. *Giant language model test room*.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. *FEVER: a large-scale dataset for fact extraction and VERification*. pages 809–819, New Orleans, Louisiana. Association for Computational Linguistics.
- Edward Tian. 2022. *Gptzero*.
- William Yang Wang. 2017. *Liar, liar pants on fire: A new benchmark dataset for fake news detection*. pages 422–426, Vancouver, Canada. Association for Computational Linguistics.
- Kyle Wiggers. 2022a. *Gpt-2 output detector demo*.
- Kyle Wiggers. 2022b. *Openai's attempts to watermark ai text hit limits*.