



Examining emerging capabilities and mitigating potential risks of VLLMs

ONE DAY WORKSHOP | 50+ ATTENDANCE (expected)

2 tasks, 3 talks, panel discussion

Abstract

1 The emergence of very large language models (VLLMs) has dramatically altered
2 the trajectory of progress in AI and its applications. With the release of ChatGPT,
3 that excitement has now transcended boundaries from AI researchers to the common
4 man; indeed it asserts we are living in an exciting time of scientific proliferation.
5 However, we see two divergent community views on VLLMs. Believers in the
6 magical powers of VLLMs claim that VLLMs are autodidactic in learning a wide
7 gamut of new capabilities – referred to as "*emerging capabilities*" in the community
8 – an effect that gets pronounced with model size and dataset scale. On the other hand,
9 critics are not yet prepared to acknowledge the self-learning power of VLLMs;
10 they criticize it as only a statistical learner and out several flaws – *hallucination*
11 being the most prominent one. In this forum, we want to bring together both
12 communities – the believers and the critics, to explore exciting tasks together:
13 (i) **CT²** - *Counter Turing Test: AI-Generated Text Detection*, and (ii) **HECT** -
14 *Hallucination eLiciTation through automatic detection and mitigation*. Expect this
15 to be an exciting forum to discuss, debate, and explore exciting scientific pathways
16 for the future.

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

18 Advances in AI during the past couple of years have led to AI systems becoming immensely more
19 powerful than ever before. While their applications for social good cannot be overstated, as an
20 unintended by-product, risks of misuse have also exacerbated. This prompted an open petition
21 letter [19] (led by Gary Marcus) by the nonprofit Future of Life Institute, calling for all AI labs to
22 immediately pause for at least 6 months "moratorium" the training of AI systems more powerful
23 than GPT-4. The letter has (18K+ and still counting) signatures from technologists and luminaries,
24 which include Yoshua Bengio, Stuart Russell, Elon Musk, Steve Wozniak, and Andrew Yang. It also
25 includes policy leaders such as Rachel Bronson, president of the Bulletin of the Atomic Scientists,
26 a science-oriented advocacy group known for its warnings against humanity-ending nuclear war.
27 On the other hand, the opposing campaign, which doesn't believe in halting scientific progress, has
28 powerful people too, including Bill Gates [3], Andrew Ng [33], Yann LeCun [33] and many others.

29 This is a significant time in the history of scientific development. In this forum, we will discuss
30 and debate emerging capabilities and mitigating potential risks and limitations of VLLMs. Call for
31 papers includes, but is not limited to: • *unique emerging abilities of VLLM*; • *negative, position, and*
32 *full paper on potential risks of VLLMs*; • *ethics and VLLMs*; • *making VLLMs more responsible*; •
33 *detection AI-generated content*; • *mitigation of harmful hallucinations*.

34 2 Two Shared Tasks

35 Shared tasks are an effective way to attract research attention to any emerging area. We will host
36 two shared tasks: (i) **CT²**, and (ii) **HECT**. The findings of **CT²** will mitigate misusage risks, while
37 **HECT** endeavours to make VLLMs more human sensitive and responsible.

38 **2.1 CT² - Counter Turing Test for AI-Generated Text Detection**

39 With the emergence of ChatGPT, the risk of AI-generated content has reached an alarming apocalypse.
40 ChatGPT has been declared banned by the school system in NYC [25], Google ads [10], and Stack
41 Overflow [17], while scientific conferences like ACL [5] and ICML [8] have released new policies
42 deterring the usage of ChatGPT for scientific writing. After the initial scepticism, ChatGPT has been
43 seen as a listed author in scientific papers [14, 21], while Elsevier [7] and Springer [26] have adopted
44 more inclusive guidelines on *the use of ChatGPT for scientific writing*.

45 Indeed, detecting AI-generated text has suddenly emerged as a concern that needs immediate attention.
46 While watermarking as a potential solution to the problem is being studied by OpenAI [32], a handful
47 of systems that detect AI-generated text such as GPT-2 output detector [31], GLTR [27], GPTZero
48 [29], DetectGPT [20], etc. have recently been orange observed in practical use. To address the
49 inevitable question of ownership attribution for AI-generated artifacts, the US Copyright Office [22]
50 released a statement stating that if the content is traditional elements of authorship produced by a
51 machine, the work lacks human authorship and the office will not register it for copyright. Given
52 this cynosural spotlight on generative AI, AI-generated text detection is a topic that needs a thorough
53 investigation. In this regard, there are three families of techniques proposed so far:

- 54 • **Watermarking:** First introduced in [1], watermarking AI-generated text involves embedding an
55 imperceptible code or signal to verify the author of a particular text with certainty. [13] proposed this
56 by selecting the next token pseudorandomly (rather than simply choosing the one with the highest
57 probability) using a cryptographic pseudorandom function whose key is only possessed by the LLM
58 maker. It would be remiss not to mention the most obvious pitfall of this approach, which is that if
59 the text is altered or modified in any way, detecting the watermark proves to be a difficult task.
- 60 • **Negative log likelihood (NLL):** NLL-based implementations such as DetectGPT [20] have demon-
61 strated the detection of AI-generated text by comparing log-likelihood of generated tokens after
62 perturbing the input text by replacing some tokens with others. If the new, perturbed version of
63 the text lies in the negative curvature regions of log-likelihood, it was likely generated by AI. The
64 limitation of this approach is that it requires access to the log probabilities of the text in order to
65 work which implies that knowledge of which LLM was used to generate the text is essential;
- 66 • **Perplexity and Burstiness:** GPTZero [29], an example of a detection technique based on perplexity
67 and burstiness, has demonstrated that a text with lower perplexity (a measure of how predictable
68 the text is), and with lower burstiness (the measure of how uniform text is) has a high probability
69 of being generated by an AI. The limitations here are that GPTZero also requires access to log
70 probabilities of text as well as the fact that it approximates perplexity values using a linear model.

71 Although AI-generated text detection has suddenly received immense attention, Liang et al. [16]
72 suggest that available AI-generated text detectors consistently misclassify non-native English writing
73 samples as AI-generated, whereas native writing samples are accurately identified, highlighting the
74 ethical implications of deploying AI-generated content detectors and risking misrepresentation. This
75 implies that a community effort is needed to tackle the issue of detectors penalizing under-represented
76 sub-population(s). The CT² task will be the first of its kind in bringing together researchers in
77 advancing the area of detecting AI-assisted generated text.

78 **2.1.1 Data to be released and the task**

79 CT² will consist of three sub-tasks. We will be releasing 100K data points, consisting of (i) prompt,
80 (ii) human-written text, and (iii) AI-generated text by 5 different LLMs.

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

MELT

2.2 MELT: 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.

2.2.1 Detection Methods

There are mainly two broad ways: (i) **Black-box factuality assessment**: Black-box hallucination detection approaches, such as SelfCheckGPT [18], endeavour 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 [24], 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 [28, 30, 9, 15, 12, 11, 23, 2], 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 Debate [tentative]

Talks: We wish to have 3 speakers from industry/academia, people who have experience in building LLMs. We (tentatively) plan to invite Prof. Christopher Manning, Professor of Linguistics and Computer Science Director, Stanford Artificial Intelligence Laboratory (SAIL); Dr. Vinodkumar Prabhakaran, Senior Research Scientist, Google LLC, co-author of PaLM [6]; and Nick Ryder who is a principal researcher at is the co-author of the GPT3 [4], works at OpenAI.



145 Workshop Organizers

Name	Web, Email, GScholar	Research Interest & Organizing activities
Dr. Amitava Das	 Web Email Google Scholar	<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>Research interests: Code-Mixing and Social Computing.</p> <p>Organizing Activities [selective]</p> <ul style="list-style-type: none"> - DEFACTIFY 1.0 @AAAI2022 - DEFACTIFY 2.0 @AAAI2023 - Memotion @SemEval2020 - SentiMix @SemEval2020 - Computational Approaches to Linguistic Code-Switching @ LREC 2020 - CONSTRAINT @AAAI2021
Dr. Amit Sheth	 Web Email Google Scholar	<p>Dr. Amit Sheth is the founding Director of the Artificial Intelligence Institute, NSC Chair and a Professor of CSE at the University of South Carolina. He was awarded 2023 IEEE Wallace McDowell award, and is a fellow of IEEE, ACM, AAAS, and AAAI.</p> <p>Research interests: Organized over 40 workshops, given over 40 tutorials, Neurosymbolic Knowledge Graph, NLP/NLU, translational research.</p> <p>Organizing Activities [selective]</p> <ul style="list-style-type: none"> - Cysoc2021 @ ICWSM2021 - Emoji2021 @ICWSM2021 - KiLKG 2021 @KGC21
Aman Chadha	 Web Email Google Scholar	<p>Aman Chadha is an Applied Science Manager at Amazon Alexa AI and a Researcher at Stanford AI.</p> <p>Research interests: Multimodal AI, On-device AI, and Human-Centered AI.</p>
Vinija Jain	 Web Email Google Scholar	<p>Vinija Jain is a Machine Learning Lead at Amazon Music.</p> <p>Research interests: Recommender Systems and NLP.</p>

146 **References**

- 147 [1] Scott Aaronson. [My Projects at OpenAI](https://scottaaronson.blog/?p=6823). 2022. URL: <https://scottaaronson.blog/?p=6823>.
- 148 [2] Rami Aly et al. “FEVEROUS: Fact Extraction and VERification Over Unstructured and Structured information”. In: 2021. URL: <https://arxiv.org/abs/2106.05707>.
- 149 [3] Bill Gates says calls to pause AI won't 'solve challenges'. Apr. 2023. URL: <https://www.reuters.com/technology/bill-gates-says-calls-pause-ai-wont-solve-challenges-2023-04-04/>.
- 150 [4] Tom Brown et al. “Language Models are Few-Shot Learners”. In: [Advances in Neural Information Processing Systems](#). Ed. by H. Larochelle et al. Vol. 33. Curran Associates, Inc., 2020, pp. 1877–1901. URL: https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf.
- 151 [5] Program Chairs. [ACL 2023 policy on Ai Writing Assistance](#). Jan. 2023. URL: <https://2023.aclweb.org/blog/ACL-2023-policy/>.
- 152 [6] Aakanksha Chowdhery et al. [PaLM: Scaling Language Modeling with Pathways](#). 2022. arXiv: 2204.02311 [cs.CL].
- 153 [7] Elsevier. [The Use of AI and AI-assisted Technologies in Scientific Writing](#). 2023. URL: <https://www.elsevier.com/about/policies/publishing-ethics?fbclid=IwAR2DBCQShp05yS7y7BT0LUxZBTVLego78m4j2tOKshCiWlCcXQpwHADkals#Authors>.
- 154 [8] Neural Information Processing Systems Foundation. 2023. URL: <https://icml.cc-Conferences/2023/11m-policy>.
- 155 [9] Sonal Garg and Dilip Kumar Sharma. “New Politifact: A Dataset for Counterfeit News”. In: 2020 9th International Conference System Modeling and Advancement in Research Trends (SMART), 2020, pp. 17–22. DOI: 10.1109/SMART50582.2020.9337152.
- 156 [10] Nico Grant and Cade Metz. [A new chat bot is a 'code red' for google's search business](#). Dec. 2022. URL: <https://www.nytimes.com/2022/12/21/technology/ai-chatgpt-google-search.html>.
- 157 [11] Ashim Gupta and Vivek Srikumar. “X-FACT: A New Benchmark Dataset for Multilingual Fact Checking”. In: arxiv, 2021. URL: <https://arxiv.org/abs/2106.09248>.
- 158 [12] Yichen Jiang et al. “HoVer: A Dataset for Many-Hop Fact Extraction And Claim Verification”. In: Findings of the Conference on Empirical Methods in Natural Language Processing (EMNLP), 2020. URL: <https://aclanthology.org/2020.findings-emnlp.309.pdf>.
- 159 [13] John Kirchenbauer et al. [A Watermark for Large Language Models](#). 2023. arXiv: 2301.10226 [cs.LG].
- 160 [14] Tiffany H Kung et al. “Performance of ChatGPT on USMLE: Potential for AI-assisted medical education using large language models”. In: [PLoS digital health](#) 2.2 (2023), e0000198.
- 161 [15] Tom Kwiatkowski et al. “Natural Questions: A Benchmark for Question Answering Research”. In: [Transactions of the Association for Computational Linguistics](#) 7 (2019), pp. 452–466. DOI: 10.1162/tacl_a_00276. URL: <https://aclanthology.org/Q19-1026>.
- 162 [16] Weixin Liang et al. [GPT detectors are biased against non-native English writers](#). 2023.
- 163 [17] Mod Makyen and Peter Olson. [Temporary policy: Chatgpt is banned](#). Dec. 1969. URL: <https://meta.stackoverflow.com/questions/421831/temporary-policy-chatgpt-is-banned>.
- 164 [18] Potsawee Manakul, Adian Liusie, and Mark JF Gales. “SelfCheckGPT: Zero-Resource Black-Box Hallucination Detection for Generative Large Language Models”. In: (2023). URL: <https://arxiv.org/abs/2303.08896>.
- 165 [19] Gary Marcus and Future of Life Institute. [Pause Giant AI Experiments: An Open Letter](#). Mar. 2023. URL: <https://futureoflife.org/open-letter/pause-giant-ai-experiments/>.
- 166 [20] Eric Mitchell et al. “Detectgpt: Zero-shot machine-generated text detection using probability curvature”. In: [arXiv preprint arXiv:2301.11305](#) (2023).

- 201 [21] Siobhan O'Connor et al. "Open artificial intelligence platforms in nursing education: Tools for
202 academic progress or abuse?" In: *Nurse Education in Practice* 66 (2022), pp. 103537–103537.
- 203 [22] Copyright Office. "Copyright Registration Guidance: Works Containing Material Generated
204 by Artificial Intelligence". In: Library of Congress, Mar. 2023. URL: <https://public-inspection.federalregister.gov/2023-05321.pdf>.
- 205 [23] Yasumasa Onoe et al. "CREAK: A Dataset for Commonsense Reasoning over Entity Knowledge". In: arXiv, 2021. doi: 10.48550/ARXIV.2109.01653. URL: <https://arxiv.org/abs/2109.01653>.
- 206 [24] Baolin Peng et al. "Check Your Facts and Try Again: Improving Large Language Models with
207 External Knowledge and Automated Feedback". In: (2023).
- 208 [25] Kalhan Rosenblatt. "CHATGPT banned from New York City public schools' devices and
209 Networks". In: NBCUniversal News Group, Jan. 2023. URL: <https://www.nbcnews.com/tech/tech-news/new-york-city-public-schools-ban-chatgpt-devices-networks-rcna64446>.
- 210 [26] Springer. Guidance on the use of Large Language Models (LLM) e.g. ChatGPT. 2023. URL:
211 <https://www.springer.com/journal/10584/updates/24013930>.
- 212 [27] Hendrik Strobelt, Sebastian Gehrmann, and Alexander Rush.
Giant Language model Test Room. 2022. URL: <http://gltr.io/dist/index.html>.
- 213 [28] James Thorne et al. "FEVER: a Large-scale Dataset for Fact Extraction and VERification". In:
214 New Orleans, Louisiana: Association for Computational Linguistics, June 2018, pp. 809–819.
215 DOI: 10.18653/v1/N18-1074. URL: <https://aclanthology.org/N18-1074>.
- 216 [29] Edward Tian. GPTZero. 2022. URL: <https://gptzero.me/>.
- 217 [30] William Yang Wang. "Liar, Liar Pants on Fire: A New Benchmark Dataset for Fake News
218 Detection". In: Vancouver, Canada: Association for Computational Linguistics, July 2017,
219 pp. 422–426. DOI: 10.18653/v1/P17-2067. URL: <https://aclanthology.org/P17-2067>.
- 220 [31] Kyle Wiggers. GPT-2 Output Detector Demo. 2022. URL: <https://openai-openai-detector.hf.space/>.
- 221 [32] Kyle Wiggers. OpenAI's attempts to watermark AI text hit limits. 2022. URL: <https://techcrunch.com/2022/12/10/openais-attempts-to-watermark-ai-text-hit-limits>.
- 222 [33] Yann LeCun and Andrew Ng: Why the 6-month AI Pause is a Bad Idea. Apr. 2023. URL:
223 https://www.youtube.com/watch?v=BY9KV8uCtj4&ab_channel=DeepLearningAI.
- 224
- 225
- 226
- 227
- 228
- 229
- 230
- 231
- 232
- 233
- 234