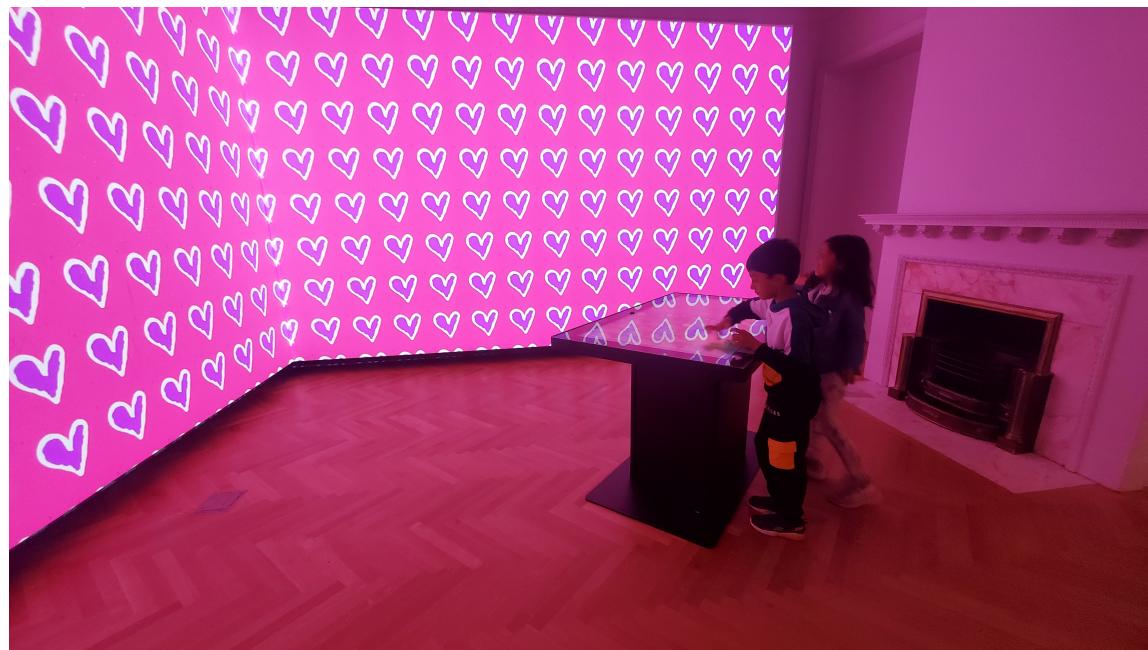


1    **A Banal Account of a Safety-Creativity Tradeoff in Generative AI**

2  
3    KUSH R. VARSHNEY, IBM Research – Thomas J. Watson Research Center, USA

4  
5    LAV R. VARSHNEY, University of Illinois Urbana-Champaign, USA



29    Fig. 1. Children creating wallpaper designs in the Immersion Room of the Cooper Hewitt, Smithsonian Design Museum.

30    Safety is banal.

31    CCS Concepts: • Computing methodologies → Artificial intelligence.

32    Additional Key Words and Phrases: computational creativity, generative model, safety, information geometry

33    **1 INTRODUCTION**

34    DALL-E 2, Stable Diffusion, Midjourney, GPT-3, ChatGPT, YouChat and other generative artificial intelligence (AI)  
35    models may be used in a variety of tasks, some mundane and some creative. Their safety may be of concern.

36    **2 SAFETY**

37    Safety is defined in terms of harm, aleatoric uncertainty, and epistemic uncertainty [3]. Safe AI systems constrain  
38    the probability of expected harms and the possibility of unexpected harms [6]. Harms from generative AI may be  
39    representational, allocative, quality-of-service, interpersonal, or societal [5].

40    **3 CREATIVITY**

41    Creativity is the generation of an artifact that is high-quality and novel [9]. Quality metrics are specific to the application.  
42    Novelty is a more application-agnostic concept that may be measured using Bayesian surprise, the relative entropy

53 between the empirical distribution of an inspiration set and that set updated with the new artifact [2]. An inspiration  
 54 set is a collection of previous artifacts in the creative domain.

55 Creativity by modern generative AI is implicitly or explicitly combinatorial. It generates unfamiliar combinations of  
 56 familiar ideas [1]. Combinatorial creativity has precise information-theoretic limits on the tradeoff between quality and  
 57 novelty [7]. On average, higher quality implies lower novelty and vice versa.

58 The more immature a creative domain is, the smaller the size of the inspiration set is. Creativity is easier because  
 59 many concepts are unexplored. The feasible region bounded by the quality-novelty tradeoff curve is larger.  
 60

61 When creative artifacts are constrained, for example by requiring intentionality, the region becomes smaller and cre-  
 62 ativity becomes more difficult [8]. (This statistical phenomenon of optimal creativity systems contrasts the computational  
 63 phenomenon of humans often being more creative with more constraints [4].)  
 64

## 66 4 SAFETY AND CREATIVITY

67 Safety is a constraint on artifacts. Like other constraints, safety makes the feasible region under the quality-novelty  
 68 tradeoff curve smaller and creativity more difficult. Thus, banality, the lack of creativity, follows from safety. There is a  
 69 tradeoff between safety and creativity.  
 70

## 72 5 IMPLICATIONS

73 Some applications of generative AI, like autonomously writing boilerplate, require safety whereas others, like inspiring  
 74 a human poet, do not. Some applications of generative AI, like writing poetry, require creativity and others, like writing  
 75 boilerplate do not. Applications requiring safety tend to also be ones not requiring creativity. Applications not requiring  
 76 safety tend to also be ones requiring creativity.  
 77

## 78 6 CONCLUSION

79 Information theory tells us that most natural applications of combinatorial creativity with modern generative AI are  
 80 feasible in terms of the safety-creativity tradeoff. Future work requires constructive algorithms for placing safety  
 81 constraints on generative AI. The end.  
 82

## 83 REFERENCES

- 84 [1] Payel Das and Lav R. Varshney. 2022. Explaining Artificial Intelligence Generation and Creativity. *IEEE Signal Processing Magazine* 39, 4 (2022),  
 85–95.
- 85 [2] Laurent Itti and Pierre Baldi. 2005. Bayesian Surprise Attracts Human Attention. In *Advances in Neural Information Processing Systems*.
- 86 [3] Niklas Möller. 2012. The Concepts of Risk and Safety. In *Handbook of Risk Theory*. Springer, Dordrecht, Netherlands, 55–85.
- 87 [4] R. Keith Sawyer. 2012. *Explaining Creativity: The Science of Human Innovation*. Oxford University Press, New York, NY, USA.
- 88 [5] Renee Shelby, Shalaleh Rismani, Kathryn Henne, AJung Moon, Negar Rostamzadeh, Paul Nicholas, N'Mah Yilla, Jess Gallegos, Andrew Smart, Emilio  
 89 Garcia, and Gurleen Virk. 2022. Sociotechnical Harms: Scoping a Taxonomy for Harm Reduction. arXiv:2210.05791.
- 90 [6] Kush R. Varshney and Homa Alemzadeh. 2017. On the Safety of Machine Learning: Cyber-Physical Systems, Decision Sciences, and Data Products.  
 91 *Big Data* 5, 3 (2017), 246–255.
- 92 [7] Lav R. Varshney. 2019. Mathematical Limit Theorems for Computational Creativity. *IBM Journal of Research and Development* 63, 1 (2019), 2.
- 93 [8] Lav R. Varshney. 2020. Limits Theorems for Creativity with Intentionality. In *Proceedings of the International Conference on Computational Creativity*.  
 94 390–393.
- 95 [9] Lav R. Varshney, Florian Pinel, Kush R. Varshney, Debarun Bhattacharjya, Angela Schörgendorfer, and Yi-Min Chee. 2019. A Big Data Approach to  
 96 Computational Creativity: The Curious Case of Chef Watson. *IBM Journal of Research and Development* 63, 1 (2019), 7.

97 Received 10 January 2023