Commentary: Computational models of intrinsic motivation for curiosity and creativity^{*}

Sophia Becker^{1,2}, Alireza Modirshanechi^{1,2}, and Wulfram Gerstner^{1,2}

¹School of Computer and Communication Sciences, EPFL ²School of Life Sciences, EPFL

Abstract

We link Ivancovsky et al.'s novelty-seeking model (NSM) to computational models of intrinsically motivated behavior and learning. We argue that dissociating different forms of curiosity, creativity, and memory based on the involvement of distinct intrinsic motivations (e.g., surprise and novelty) is essential to empirically test the conceptual claims of the NSM.

Human and animal behavior is driven not only by extrinsically available rewards like food and money but also by various intrinsic motivations, such as the desire to experience novelty or surprise.^{1,2} Curiosity and creativity are two modes of cognitive processing where such intrinsic motivations have a significant influence. Ivancovsky et al.'s novelty-seeking model (NSM) creates a valuable conceptual link between these intuitively related modes, and divides the shared cognitive processes underlying curiosity and creativity into four phases.[?] However, the model's high-level conceptual nature makes it challenging to give quantitative explanations and derive experimentally testable hypotheses. To address this problem, we relate each of the four phases of the NSM to computational models of intrinsically motivated behavior and learning. We discuss (i) in which ways computational models support or contradict the NSM's core claims, and illustrate (ii) how computational models make the conceptual explanations and predictions of the NSM empirically testable.

First, the NSM posits that curiosity and creativity share brain networks and mechanisms to detect 'novelty', either in the external space of sensory stimuli (curiosity) or in the internal space of associations (creativity). Second, these shared mechanisms initiate downstream processing of the 'novel' stimulus or association.³ However, while Ivancovsky et al. use 'novelty' as a general notion, distinct intrinsic motivations contributing to curiosity (e.g., novelty, surprise, information-gain) are mathematically well-defined,^{4,5} have different neural signatures,^{6–9} and are triggered by different statistical regularities of the task or environment¹⁰ (see¹¹ for a review). For example, novelty signals are triggered by unfamiliar stimuli and situations, both when the unfamiliarity is expected and when it is unexpected.¹² Surprise signals, on the other hand, arise in the face of unexpected stimuli, both familiar and unfamiliar ones.⁹ In line with that, different neuromodulatory signals are thought to communicate expected vs. unexpected novelty or uncertainty;^{13,14} and computational models suggest different network mechanisms for the detection of novelty and surprise.^{15,16} Despite the partial overlap in the processing of novelty and surprise,⁹ we can thus not simply speak of

^{*}See Sophia Becker, Alireza Modirshanechi, Wulfram Gerstner. Computational models of intrinsic motivation for curiosity and creativity. Behavioral and Brain Sciences. 2024;47:e94. doi:10.1017/S0140525X23003424

'novelty' detection as a homogeneous process as assumed in the NSM. When empirically testing shared neural mechanisms of curiosity- and creativity-related signal detection and downstream processing, we should therefore consider how the neural correlates of curiosity and creativity may vary across environments and experimental tasks.

Third, the NSM proposes that both curiosity and creativity require a balance of exploratory and exploitatory states of mind (SoM), and that this balance is mediated by cognitive control processes. This NSM prediction agrees with reinforcement learning-based (RL) models that arbitrate between intrinsic motivations (curiosity/exploratory SoM) and extrinsic motivations (reward/exploitatory SoM).^{5,17} Importantly, these RL models quantify the respective contributions of exploration and exploitation to behavior, and allow to test which mechanisms regulate the trade-off between the exploratory and exploitatory states. For example, a recent model that arbitrates exploration and exploitation based on the agent's reward optimism¹⁸ provides a concrete computational implementation of Ivancovsky et al.'s conceptual links between curiosity, creativity and the SoM dimension of openness to experience. We propose that this modeling approach is a useful tool to experimentally validate links between curiosity/creativity and different SoM dimensions as suggested by the NSM.

Lastly, a central component of the NSM is the bidirectional link between memory and curiosity/creativity.³ However, there are different forms of memory and distinct synaptic learning rules that are influenced by intrinsic motivational signals (three-factor learning rules^{19,20}). While we agree with the bidirectional link between curiosity/creativity and memory systems, we propose that the respective memory system with which curiosity and creativity engage could differ (e.g., recognition vs. episodic memory). More importantly, distinct forms of curiosity and creativity may link to different learning rules and roles of memory. For example, novelty is particularly important for initial memory formation,^{21,22} while surprise signals the violation of known rules and expectations^{?,?,23} and might therefore be more important for targeted memory updates. Another relevant distinction that the NSM is currently abstracting is between (i) memory systems that support the detection of intrinsic motivational signals and (ii) memory systems that are downstream targets of curiosity/creativity-related signals. These memory systems may – but do not have to – be identical. For example, novelty detection relies on state representations in sensory areas and recognition memory,^{12,24} but downstream novelty signals are also involved in updating semantic or episodic memories.^{21,22,25} To empirically determine how memory is shared by curiosity and creativity, it is necessary to experimentally test how different memory systems are involved at each stage and in each type of curiosity/creativity-related processing.

To conclude, we illustrated how the high-level cognitive NSM framework relates to concrete computational models of intrinsically motivated behavior and learning. While computational models and the NSM align on the general structure of curiosity and creativity-related processing, computational models suggest important distinctions within each phase of the NSM. In particular, different forms of curiosity and creativity arising from the contribution of distinct intrinsic motivational signals, like novelty and surprise, could differ in the specifics of how they are detected, signaled to downstream targets and interacting with memory systems. Linking the NSM to computational models is thus a necessary step to empirically test the NSM's conceptual predictions and gain insights into the neural correlates and network mechanisms underlying curiosity and creativity.

1 Competing Interests

The authors declare no competing interests in relation to this work.

2 Financial Support

This work was supported by the Swiss National Science Foundation No. 200020_207426.

References

- Jacqueline Gottlieb and Pierre-Yves Oudeyer. Towards a neuroscience of active sampling and curiosity. *Nature Reviews Neuroscience*, 19(12):758–770, 2018.
- [2] Alireza Modirshanechi, Kacper Kondrakiewicz, Wulfram Gerstner, and Sebastian Haesler. Curiosity-driven exploration: foundations in neuroscience and computational modeling. *Trends in Neurosciences*, 46:1054–1066, 2023.
- [3] Tal Ivancovsky, Shira Baror, and Moshe Bar. A shared novelty-seeking basis for creativity and curiosity. *Behavioral and Brain Sciences*, 47:e89, 2024.
- [4] Andrew Barto, Marco Mirolli, and Gianluca Baldassarre. Novelty or surprise? Frontiers in psychology, 4:907, 2013.
- [5] Alireza Modirshanechi, Johanni Brea, and Wulfram Gerstner. A taxonomy of surprise definitions. Journal of Mathematical Psychology, 110:102712, 2022.
- [6] He A. Xu, Alireza Modirshanechi, Marco P. Lehmann, Wulfram Gerstner, and Michal H. Herzog. Novelty is not Surprise: Human exploratory and adaptive behavior in sequential decision-making. *PLoS Computational Biology*, 17:e1009070, 2021.
- [7] Joachim Morrens, Çağatay Aydin, Aliza Janse van Rensburg, José Esquivelzeta Rabell, and Sebastian Haesler. Cue-evoked dopamine promotes conditioned responding during learning. *Neuron*, 106(1):142 – 153.e7, 2020.
- [8] Korleki Akiti, Iku Tsutsui-Kimura, Yudi Xie, Alexander Mathis, Jeffrey E. Markowitz, Rockwell Anyoha, Sandeep Robert Datta, Mackenzie Weygandt Mathis, Naoshige Uchida, and Mitsuko Watabe-Uchida. Striatal dopamine explains novelty-induced behavioral dynamics and individual variability in threat prediction. *Neuron*, 110:3789–3804, 2022.
- [9] Kaining Zhang, Ethan S. Bromberg-Martin, Fatih Sogukpinar, Kim Kocher, and Ilya E. Monosov. Surprise and recency in novelty detection in the primate brain. *Current Biology*, 32(10):2160–2173.e6, 2022.
- [10] Maxime Maheu, Stanislas Dehaene, and Florent Meyniel. Brain signatures of a multiscale process of sequence learning in humans. *eLife*, 8:e41541, 2019.
- [11] Alireza Modirshanechi, Sophia Becker, Johanni Brea, and Wulfram Gerstner. Surprise and novelty in the brain. *Current Opinion in Neurobiology*, 82:102758, 2023.
- [12] Jan Homann, Sue A. Koay, Kevin S. Chen, David W. Tank, and Michael J. Berry. Novel stimuli evoke excess activity in the mouse primary visual cortex. *Proceedings of the National Academy of Sciences*, 119:e2108882119, 2022.
- [13] Angela J. Yu and Peter Dayan. Uncertainty, neuromodulation, and attention. Neuron, 46:681–692, 2005.
- [14] M. Meeter J. Schomaker. Short- and long-lasting consequences of novelty, deviance and surprise on brain and cognition. *Neuroscience & Biobehavioral Reviews*, 55:268–279, 2015.

- [15] Auguste Schulz, Christoph Miehl, Michael J. Berry II, and Julijana Gjorgjieva. The generation of cortical novelty responses through inhibitory plasticity. *eLife*, 10:e65309, 2021.
- [16] Martin Barry and Wulfram Gerstner. Fast adaptation to rule switching using neuronal surprise. bioRxiv, 2022.
- [17] Adrià Puigdomènech Badia, Bilal Piot, Steven Kapturowski, Pablo Sprechmann, Alex Vitvitskyi, Zhaohan Daniel Guo, and Charles Blundell. Agent57: Outperforming the Atari human benchmark. In Proceedings of the 37th International Conference on Machine Learning, volume 119 of Proceedings of Machine Learning Research, pages 507–517. PMLR, 2020.
- [18] Alireza Modirshanechi, He A. Xu, Wei-Hsiang Lin, Michael H. Herzog, and Wulfram Gerstner. The curse of optimism: a persistent distraction by novelty. *bioRxiv*, 2022.
- [19] F. Zenke, E. J. Agnes, and W. Gerstner. Diverse synaptic plasticity mechanisms orchestrated to form and retrieve memories in spiking neural networks. *Nature communications*, 6:1–13, 2015.
- [20] Wulfram Gerstner, Marco Lehmann, Vasiliki Liakoni, Dane Corneil, and Johanni Brea. Eligibility traces and plasticity on behavioral time scales: experimental support of neohebbian three-factor learning rules. *Frontiers in neural circuits*, 12, 2018.
- [21] Adrian J. Duszkiewicz, Colin G. McNamara, Tomonori Takeuchi, and Lisa Genzel. Novelty and dopaminergic modulation of memory persistence: A tale of two systems. *Trends in Neurosciences*, 42:102–114, 2019.
- [22] James B. Priestley, John C. Bowler, Sebi V. Rolotti, Stefano Fusi, and Attila Losonczy. Signatures of rapid plasticity in hippocampal CA1 representations during novel experiences. *Neuron*, 110:1978–1992, 2022.
- [23] He A Xu, Alireza Modirshanechi, and Marco P Lehmann. Novelty is not Surprise : Human exploratory and adaptive behavior in sequential decision-making. 2021.
- [24] Rafal Bogacz and Malcolm W. Brown. Comparison of computational models of familiarity discrimination in the perirhinal cortex. *Hippocampus*, 13:494–524, 2003.
- [25] B. C. Wittmann, N. Bunzeck, and R. J. Dolan. Anticipation of novelty recruits reward system and hippocampus while promoting recollection. *Neuroimage*, 38:194–202, 2007.