



Reply to Comment

## Ecological causal cognition needs disentanglement: Reply to comments on “Disentangled representations for causal cognition”

Filippo Torresan <sup>a</sup>, Manuel Baltieri <sup>a,b,\*</sup><sup>a</sup> Araya Inc., Tokyo, Japan<sup>b</sup> University of Sussex, Falmer, Brighton, BN1 9RH, United Kingdom

### ARTICLE INFO

Editor: J. Fontanari

#### Keywords:

disentanglement  
causal cognition  
ecological cognition  
embodiment  
causality

One of the main goals of our recent article [1] was to highlight mutual points of interest between fields that conduct seemingly different investigations on causal learning and reasoning, and to suggest relevant directions for future research. To achieve this goal, we offered a unifying computational perspective on animal causal cognition by elaborating upon the three dimensions of this faculty, i.e., (1) explicitness, (2) sources, and (3) integration, proposed by [2]. A key move in our treatment was to identify explicitness with the machine learning notion of disentanglement, i.e., the factorisation of a representation, and view it as the central building block of our account, next to which the other two dimensions would naturally fit.

With a different emphasis, both commentaries suggest that several of the approaches and theoretical tools described in the target paper seem to not take embodied and ecological perspectives on cognition enough into consideration. Here we want to clarify our position and stress some of the opportunities afforded by our proposal, as one of the first of its kind, with respect to some of the points raised by the commentaries.

More specifically, one of the commentaries (**Goddu**) observes that, in causal machine learning, perception tends to be modelled purely as information processing of observations viewed as raw perceptual data. This approach, **Goddu** suggests, seemingly ignores that (1) perceptual learning is always intimately linked to what is valuable for a certain agent (value-laden perception), (2) perception and action co-organise and co-transform in relation to what serves the needs of the agents (integration of perception and action), and (3) perception is more about disclosing a world of affordances than reconstructing a veridical picture of reality characterised by objective, agent-independent, properties.

These have long been recognised as indispensable for a proper study of natural cognition, and form the bedrock of many lines of research in contemporary cognitive science [3,4]. It is also true that causal machine learning has sometimes ignored those points by developing techniques with a different aim, related to the discovery of causal relationships in various scientific domains (e.g., healthcare) and therefore many of those techniques and their properties might not directly speak to methods more in line with ecological perception and action in agents interacting with an environment.

\* Corresponding author.

E-mail addresses: [filippot.research@protonmail.com](mailto:filippot.research@protonmail.com) (F. Torresan), [manuel\\_baltieri@araya.org](mailto:manuel_baltieri@araya.org) (M. Baltieri).

To address precisely these concerns, the second half of the target article underlines the relevance of reinforcement learning as a mathematical and computational framework that is focussed on modelling agents that learn from experience, i.e., by interacting with an environment [5]. More concretely, contemporary deep reinforcement learning methods with learned value functions, such as actor-critic algorithms, already model a sense in which perception is value-laden (1). This is because the critic is a deep neural network that converts observations into value predictions for the agent [6–9]. Furthermore, since value predictions are used to train the actor (or policy) network, whose actions in turn determine the trajectory of observations on which the critic is trained, a certain kind of action-perception integration is arguably *the* starting point of several approaches in deep reinforcement learning (2). Even more so if we consider critics that make q-value predictions, i.e., predictions of state-action pairs' values, whose training depends even more closely on the actor network. One common critique of this framing involves the fact that, in general, the approach sketched so far presupposes the existence of a *reward function*, from which the value functions themselves can be learned. To tackle this, several proposals have introduced some form of intrinsic motivation that can help agents to explore and solve tasks in complex environment with sparse or, even, no reward, see for instance [10–14]. Similarly, different proposals have adapted parts of a theory of affordances to reinforcement learning (3), showing how more technical formulations of 'affordance' can be beneficial for planning and learning in reinforcement learning, e.g., in robotics applications [15–17].

The second commentary (Finke & Raja) objects, on the other hand, that disentangled representations might not even be necessary for explaining certain causal abilities. This is because, according to the authors and again stressing key aspects of the embodied and enactive traditions, natural agents are more inclined to adopt solutions that reduce their cognitive load. For instance, one point made in the commentary is that forms of *associative* learning aiming to track causal relationships (see, e.g., [18]) could be a sufficient cognitive strategy for resource-bound adaptive agents most of the time, freeing them from the costs of building explicit causal models when there is a need to act promptly.

This objection, however, neglects the fact that biological constraints, such as *minimising* some form of activity energy, have been linked to the emergence of disentangled representations [19]. Furthermore, we note that studies in deep reinforcement learning agents also provide evidence that some proxies for disentanglement, such as orthogonality and sparsity, together with other properties like complexity reduction, facilitate transfer learning [20]. This kind of empirical work suggests that disentanglement promotes an efficient use of learning resources by producing "good" representations, i.e. representations that can be reused in other contexts. In other words, in contrast to Finke & Raja's point that causal models are costly (e.g., from a metabolic point of view), evidence suggests that disentanglement is a consequence of agents having limited resources, providing an advantage by enabling an efficient reuse of their cognitive abilities.

Furthermore, theoretical studies have shown that agents robust to distributional shifts, i.e., capable of adapting to changes in the data-generation process (*cf.*, *strong generalisation*, see [21]), must have learned an approximate causal model of the environment [22,23]. This is evidence that, while in some sense more costly and difficult to achieve compared to its weak counterpart, strong disentanglement (see [1, 360]) is required for the kind of causal learning and reasoning that enables an agent, especially in ecological situations, to continually adapt to environmental perturbations. In other words, while these studies focus only on simplified agent-environment setups, they suggest that causal models are important for generalisation even in more complex, ecological settings, and that disentanglement has a crucial role to play in an account of causal learning and decision-making (see also [24–26]).

Overall, these considerations only highlight one instance of the separation between conceptual accounts of causal cognition and the mechanistic and empirical accounts that our work aims to bring closer together. This separation is due partly to an underspecification of how such conceptual understandings can be operationalised, and partly to a focus on accounts of simple associative learning rules that have so far largely failed to explain any form of robust generalisation across contexts and higher order reasoning. For this reason, our aim is to provide a clear mathematical formulation, together with precise computational tools, to *ground* questions about the forms of causal cognition exhibited by several living organisms. Our belief is that this approach could also benefit embodied and enactive accounts of cognition, by formally characterising what (causal) models can be attributed to agents [27–30], and by allowing us to explore a computational space of causal learning strategies of different complexity, and empirically access what circumstances require an agent to prefer simpler ones over more complex ones, or *vice versa*, as we suggest in the target article.

According to both commentaries, our proposal appears to underestimate the centrality of an agent's adaptive actions for causal learning, e.g., as seemingly indicated by the importance given to one-shot problems. However, this would be true only if *observational* causal learning was the core of our proposal, which would certainly make the discovery of causal regularities exceedingly difficult and costly for natural agents.

By focusing on deep reinforcement learning instead, the proposed computational framework regards adaptive actions as crucial, and has the resources to accommodate their role in the exploration of "causally embedded context" (Finke & Raja) and in "making things happens" by means of "time-locked attention" to specific features and/or outcomes (Goddu). On this view, observational causal learning, e.g., one-shot problems, should only be seen as a sophisticated skill built on top of experiences acquired and organised via sustained complex interactions with an environment.

Overall, our proposal aims to inspire computational investigations on agents' exploration and learning in open-ended (multi-agent) environments [31–33]. Our goal is to answer more specific questions about when and how, in those settings, agents would be compelled to learn causal models, progressively refine them through continual interactions, and develop more sophisticated causal abilities, e.g., enabling the solution of one-shot problems (for a recent study along these lines, see [34]). To do so, we believe that an operationalisation like the one we proposed, or one that reaches a comparable amount of detail, is paramount to ground questions like the ones we [1] and the authors of the two commentaries put forward.

## Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Acknowledgements

This work was supported by JST, Moonshot R&D, Grant Number JPMJMS2012.

## References

- [1] Torresan F, Baltieri M. Disentangled representations for causal cognition. *Phys Life Rev* 2024. <https://doi.org/10.1016/j.plev.2024.10.003>
- [2] Starzak TB, Gray RD. Towards ending the animal cognition war: a three-dimensional model of causal cognition. *Biol Philos* 2021;36(2):1–24. <https://doi.org/10.1007/s10539-021-09779-1>
- [3] Clark A. *Surfing uncertainty: prediction, action, and the embodied mind*. Oxford, UK: Oxford University Press; 2016. ISBN 978-0-19-021703-7.
- [4] Baggs E, Chemero A. Radical embodiment in two directions. *Synthese* 2021;198(9):2175–90. <https://doi.org/10.1007/s11229-018-02020-9>
- [5] Sutton RS, Barto AG. *Reinforcement learning: an introduction*. Adaptive computation and machine learning; The MIT Press; 2018. ISBN 978-0-262-03924-6.
- [6] Barto AG, Sutton RS, Anderson CW. Neuronlike adaptive elements that can solve difficult learning control problems. *IEEE Trans Syst Man Cybern* 1983;SMC-13(5):834–46. <https://doi.org/10.1109/TSMC.1983.6313077>
- [7] Degris T, White M, Sutton RS, et al. Off-policy actor-critic. In: *Proceedings of the 29th international conference on machine learning*. ICML'12; Madison, WI, USA: Omnipress. ISBN 978-1-4503-1285-1; 2012, p. 179–86.
- [8] Haaroja T, Zhou A, Abbeel P, Levine S, et al. Soft actor-critic: off-policy maximum entropy deep reinforcement learning with a stochastic actor. In: *35th International conference on machine learning*, 2018;vol. 5: 2976–89.
- [9] Schulman J, Moritz P, Levine S, Jordan M, Abbeel P, High-dimensional continuous control using generalized advantage estimation. 2018. <https://doi.org/10.48550/arXiv.1506.02438>
- [10] Klyubin AS, Polani D, Nehaniv CL, et al. Empowerment: a universal agent-centric measure of control. In: *2005 IEEE Congress on evolutionary computation*; 2005;vol. 1: p. 128–35. . <https://doi.org/10.1109/CEC.2005.1554676>
- [11] Salge C, Glackin C, Polani D, et al. Changing the environment based on empowerment as intrinsic motivation. *Entropy* 2014;16(5):2789–819. . <https://doi.org/10.3390/e16052789>
- [12] Kim H, Kim J, Jeong Y, Levine S, Song HO, et al. EMI: exploration with mutual information. In: Chaudhuri K, Salakhutdinov R, editors. *Proceedings of the 36th international conference on machine learning*; vol. 97 of *Proceedings of Machine Learning Research*. Long Beach, California, USA: PMLR; 2019, p. 3360–9.
- [13] Rhinehart N, Wang J, Berseth G, Co-Reyes J, Hafner D, Finn C, et al., et al. Information is power: intrinsic control via information capture. In: Ranzato M, Beygelzimer A, Dauphin Y, Liang PS, Vaughan JW, editors. *Advances in neural information processing systems*; vol. 34. Curran Associates, Inc.; 2021, p. 10745–58.
- [14] Ramírez-Ruiz J, Grytskyy D, Mastrogiuseppe C, Habib Y, Moreno-Bote R, et al. Complex behavior from intrinsic motivation to occupy future action-state path space. *Nat Commun* 2024;15(1):6368. <https://doi.org/10.1038/s41467-024-49711-1>
- [15] Khetarpal K, Ahmed Z, Comanici G, Abel D, Precup D, et al. What can I do here? A theory of affordances in reinforcement learning. In: *Proceedings of the 37th international conference on machine learning*. PMLR; 2020, p. 5243–53.
- [16] Xu D, Mandlekar A, Martín-Martín R, Zhu Y, Savarese S, Fei-Fei L, et al. Deep affordance foresight: planning through what can be done in the future. In: *2021 IEEE International conference on robotics and automation (ICRA)*. 2021, pp. 6206–13. . <https://doi.org/10.1109/ICRA48506.2021.9560841>
- [17] Liao Y-C, Todi K, Acharya A, Keurulainen A, Howes A, Oulasvirta A, et al. Rediscovering affordance: a reinforcement learning perspective. In: *Proceedings of the 2022 CHI conference on human factors in computing systems*. CHI '22; New York, NY, USA: Association for Computing Machinery. ISBN 978-1-4503-9157-3; 2022, p. 1–15. <https://doi.org/10.1145/3491102.3501992>.
- [18] Le Pelley ME, Griffiths O, Beesley T. Associative accounts of causal cognition. In: Waldmann MR, editor. *The oxford handbook of causal reasoning*. Oxford, UK: Oxford University Press. ISBN 978-0-19-939955-0; 2017, p. 14–28.
- [19] Whittington J CR, Dorrell W, Ganguli S, Behrens T, et al. Disentanglement with biological constraints: a theory of functional cell types. In: *The eleventh international conference on learning representations*. 2023,.
- [20] Wang H, Miah E, White M, Machado MC, Abbas Z, Kumaraswamy R, et al. Investigating the properties of neural network representations in reinforcement learning. *Artif Intell* 2024;330:104100. . <https://doi.org/10.1016/j.artint.2024.104100>
- [21] Schölkopf B, Locatello F, Bauer S, Ke NR, Kalchbrenner N, Goyal A, et al. Toward causal representation learning. *Proc IEEE* 2021;109(5):612–34. <https://doi.org/10.1109/JPROC.2021.3058954>
- [22] Richens J, Everitt T. Robust agents learn causal world models. In: *The twelfth international conference on learning representations*. 2024.
- [23] Ceriscioli M, Mohan K. Agents robust to distribution shifts learn causal world models even under mediation. In: *The thirty-ninth annual conference on neural information processing systems*. 2025.
- [24] Mutti M, Santi RD, Rossi E, Calderon JF, Bronstein M, Restelli M, et al. Provably efficient causal model-based reinforcement learning for systematic generalization. *Proc AAAI Conf Artif Intell* 2023;37(8):9251–9. . <https://doi.org/10.1609/aaai.v37i8.26109>
- [25] Cao H, Feng F, Fang M, Dong S, Yang T, Huo J, et al. Towards empowerment gain through causal structure learning in model-based reinforcement learning. In: *The thirteenth international conference on learning representations*. 2025.
- [26] Sontakke SA, Mehrjou A, Itti L, Schölkopf B, et al. Causal curiosity: RL agents discovering self-supervised experiments for causal representation learning. In: Meila M, Zhang T, editors. *Proceedings of the 38th international conference on machine learning*; vol. 139. PMLR; 2021, p. 9848–58.
- [27] Baltieri M, Biehl M, Capucci M, Virgo N. A Bayesian interpretation of the internal model principle. 2025, arXiv preprint arXiv:250300511.
- [28] Virgo N, Biehl M, Baltieri M, Capucci M. A “good regulator theorem” for embodied agents. In: *Artificial life conference proceedings 37*; vol. 2025. MIT Press, Cambridge; 2025, p. 46.
- [29] Richens J, Everitt T, Abel D. General agents need world models. In: *Forty-second international conference on machine learning*. 2025.
- [30] Baltieri M, Suzuki K. Mathematical approaches to the study of agents. *Philos Trans R Soc B* 2025. <https://doi.org/10.1098/rstb.2025.0228>
- [31] Matthews M, Beukman M, Ellis B, Samvelyan M, Jackson MT, Coward S, et al. Craftax: a lightning-fast benchmark for open-ended reinforcement learning. In: *Proceedings of the 41st international conference on machine learning*. PMLR; 2024, p. 35104–37.
- [32] Voudouris K, Slater B, Cheke LG, Schellaert W, Hernández-Orallo J, Halina M, et al. The animal-AI environment: a virtual laboratory for comparative cognition and artificial intelligence research. *Behav Res Methods* 2025;57(4):107. <https://doi.org/10.3758/s13428-025-02616-3>
- [33] Ghugare R, Castanyer RC, Ji C, Wantlin K, Schofield J, Narasimhan K, et al. BuilderBench: the building blocks of intelligent agents. 2026. <https://doi.org/10.48550/arXiv.2510.06288>
- [34] Lidayan A, Du Y, Kosoy E, Rufova M, Abbeel P, Gopnik A, et al. Intrinsically-motivated humans and agents in open-world exploration. 2025. <https://doi.org/10.48550/arXiv.2503.23631>