
Position: World Models must live in Parallel Worlds

Sahithya Ravi^{*1,2} Aditya Chinchure^{*1,2}
Pushkar Shukla³, Vered Shwartz^{1,2}, Leonid Sigal^{1,2}

¹The University of British Columbia, ²Vector Institute of AI
³Toyota Technological Institute at Chicago
{sahiravi, aditya10, vshwartz, lsigal}@cs.ubc.ca
pushkarshukla@ttic.edu

Abstract

World models learn spatio-temporal representations of a world, enabling them to predict future states, and support interaction, navigation, and simulation capabilities. For generative models to become effective agents in the physical world, they must develop and use world models. We posit that world models must be capable of counterfactual simulation – the ability to reason about *what if* scenarios. By simulating alternative realities, world models will be more capable, safe and creative when faced with novel, out-of-distribution scenarios. Furthermore, they can transcend mere pattern matching to achieve a true causal understanding of the world, a capability central to human intelligence, and a prerequisite for the next generation of AI agents.

1 The Case-Based Generalization Crisis

Generative AI has demonstrated remarkable capabilities in creating text, images and videos that mimic human output [1, 2]. However, for these models to transition from digital content creators to effective agents in the physical world, they require a deeper understanding of how the world works. This understanding is encapsulated in the concept of a “world model”, an internal representation that allows an agent to simulate and predict the consequences of actions within its environment [3].

Current efforts to build world models focus on predicting future states from past observations, typically by scaling models and exposing them to millions of examples [1, 2, 4–7]. This approach, while powerful, creates a fundamental generalization gap. The resulting models excel at interpolating within their training data, but falter when asked to extrapolate to novel scenarios. They engage in *case-based generalization* [4, 8–10], effectively imitating the most similar training instances rather than abstracting the underlying physical or causal principles. This brittleness manifests in critical failures.

Models often lack compositionality, struggling to combine familiar concepts in novel contexts. For example, Veo 3 [2] struggles to generate “a hummingbird flying over a city” (Appendix A.1) because it associates the bird with natural habitats, rather than abstracting hummingbird and flying as transferable concepts. They also mistake correlation for causation, producing hallucinations such as people walking backwards in Genie 3 [5]. Additionally, they remain opaque black boxes, unable to explain their reasoning and unsuitable for safety critical applications. This raises a central question: *What capabilities would make world models robust in novel, out-of-distribution settings?*

^{*}Equal contribution. Author order selected by coin flip.

²<https://deepmind.google/models/veo/>

³Why generative world models aren’t ready for real applications (2025)

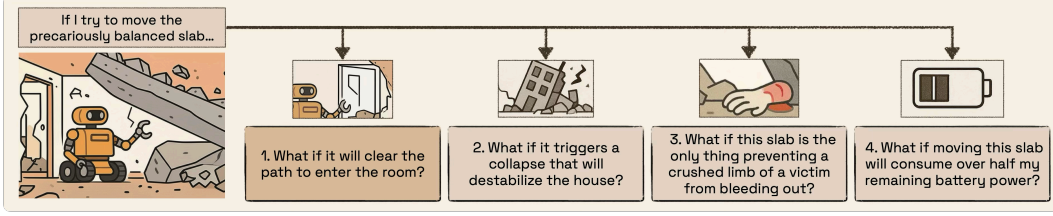


Figure 1: A search-and-rescue robot in a collapsed building: Should it move the precariously balanced slab? Current world models may predict trajectory (1), but **counterfactual simulation (CfS)** enables evaluating counterfactual trajectories (2, 3, 4) that are crucial for safety.

Human cognition provides a useful point of reference. When humans encounter a novel situation, we hypothesize alternatives: if a human drives into a school zone, we imagine that children may run into the road, or that a car door may suddenly open [11–13]. Such counterfactual simulations allow us to substitute, recombine, and adapt knowledge flexibly, enabling robustness in out-of-distribution settings.

Motivated by this, we posit that equipping models with the ability to perform **counterfactual simulation (CfS)** could be a key ingredient to achieving robustness, safety, and out-of-distribution generalization. For example, in Figure 1, a search and rescue robot enters a collapsed building and confronts a precariously balanced slab. A standard predictive world model simulates the most likely future (e.g., moving the slab to clear the path) and can miss catastrophic alternatives, such as collapse from a slight disturbance. By simulating parallel, hypothetical futures, world models could move beyond statistical pattern matching towards true causal reasoning, as we observe in human mental models.

2 Counterfactual Simulation

World model. Following [3], we define a world model as an internal representation that captures the causal structure and spatiotemporal dynamics of an environment. Given an initial state and an action, a world model predicts the next state and can roll out sequences to generate trajectories.

Counterfactuals. A counterfactual is a hypothetical “what-if” that changes a specific aspect of the world and examines how the outcome would differ, allowing us to reason about alternative outcomes. Within Pearl’s Ladder of Causation [14], counterfactuals sit at the highest level. Beyond association (observing correlations) and intervention (predicting the effects of actions), they answer questions of the form, “What would have happened if I had acted differently?” (see Appendix A.2 for details).

Counterfactual simulation (CfS). A counterfactual simulation (CfS) is an alternative sequence of events generated by a world model that explores what could have happened had a specific event been different. To formalize, we define an event (e_t) as a single unit combining a state and an action at time t , such that $e_t = (s_t, a_t)$. The system’s dynamics are captured by a world model (M), a function that predicts the next event. A trajectory (τ) is a sequence of these events over time, $\tau = (e_0, e_1, \dots, e_T)$. Consider a factual trajectory (τ_{fact}) that represents the most likely sequence of events:

$$\tau_{\text{fact}} = (e_0^{\text{fact}}, e_1^{\text{fact}}, \dots, e_T^{\text{fact}})$$

where each $e_t^{\text{fact}} = (s_t^{\text{fact}}, a_t^{\text{fact}})$.

A **counterfactual trajectory** (τ_{cf}) is a hypothetical alternative created by performing an intervention. This involves selecting a specific step k and replacing the factual event e_k^{fact} with a different, hypothetical event e_k^{cf} . The new trajectory is then defined as:

$$\tau_{\text{cf}} = (e_0^{\text{cf}}, e_1^{\text{cf}}, \dots, e_T^{\text{cf}})$$

This trajectory is constructed by following a three-part process, which can be visualized as a branching path from the factual trajectory. First, for all steps leading up to the intervention ($t < k$), the counterfactual events are identical to the factual ones: $e_t^{\text{cf}} = e_t^{\text{fact}}$. Second, at the intervention point ($t = k$), the factual event is replaced: $e_k^{\text{cf}} \neq e_k^{\text{fact}}$. Finally, for all steps following the intervention

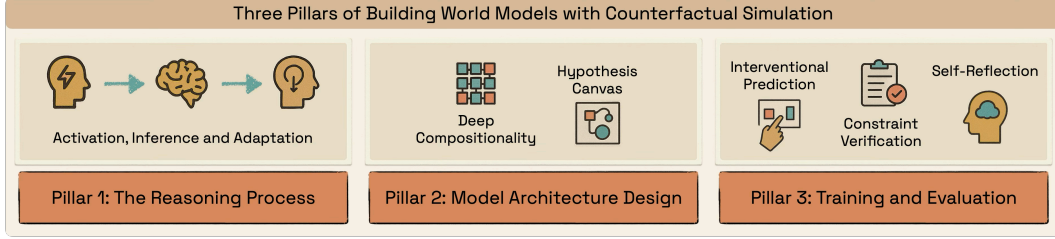


Figure 2: The three pillars for effective implementation of counterfactual simulation in world models.

($t > k$), the subsequent events are generated by the world model, simulating the consequences of the counterfactual change: $e_{t+1}^{\text{cf}} \sim M(\cdot | e_t^{\text{cf}})$. This process generates a new trajectory that is identical to the factual one up to the intervention point but diverges afterward due to the change at step k .

3 Alternative Views

Prevailing View. A common prevailing view is that case-base generalization will be overcome naturally from scaling up generative models with more data and compute [1, 15]. Others advocate for test-time scaling: allocating greater compute at inference to improve reasoning abilities and generalization [16-18].

Our Position. We believe that scaling data, models or test-time compute alone is not a viable path forward to robust world modeling, for the following reasons.

- ① High quality human data is finite⁴, pushing reliance on synthetic data that risks model collapse as models train on their own outputs [19]. The energy and water demands of large data centers raise sustainability concerns and question the long term viability of continued scaling.
- ② Large-scale models are data inefficient. LLMs train on trillions of tokens, far exceeding a child’s linguistic exposure and still lack robust world understanding.⁵
- ③ Allocating more compute at test-time is not a cure-all for a flawed underlying model. A model that doesn’t grasp intuitive physics, like object permanence, will just explore a larger tree of physically implausible outcomes, no matter how much compute is thrown at it.

4 The Path Forward

Effective counterfactual simulation requires a strong underlying causal framework for the world the model is trained on, and a structured process to trigger, generate, store, and use counterfactual simulations. We propose a three pillar approach (Fig. 2) to CfS in world models: the underlying reasoning process, the architectures to support it, and the training and evaluation methods for it.

4.1 Pillar 1: The Reasoning Process

Inspired by the human cognitive process that governs counterfactual thinking ([13]), a world model capable of CfS requires capabilities for the reasoning processes of *activation* \rightarrow *inference* \rightarrow *adaptation*.

Activation. World models interacting with the physical world must decide *when to simulate counterfactuals*. This requires a system to identify which event $e_{t=k}^{\text{fact}}$ in τ_{fact} is an “activation event” based on its causal significance. In Fig. 1, the robot must identify that moving the slab is an activation event, because its alterations would create meaningfully different futures.

Inference. When an activation event is identified, the world model can perform targeted interventions on it to simulate consequences. Since full simulations are computationally expensive, a meta-cognitive

⁴<https://globalcio.com/news/14933/>

⁵<https://babylm.github.io/>

process determines when to use complete simulation versus cheaper heuristics. This process weighs decision importance and uncertainty against computational cost, reserving full counterfactual rollouts for high-stakes scenarios.

Adaptation. After inference, the model can use CfS outcomes to inform appropriate actions or preventative measures, using reasoning methods.

4.2 Pillar 2: Model Architecture Design

Deep Compositionality. To simulate meaningful counterfactuals, we call for architectures that can encode the world as a system of disentangled concepts, objects, and physical rules — where elemental building blocks can combine to form larger concepts. This deep compositionality requires the model to also learn the fundamental, transferable properties of objects and concepts, and the causal relationships between them. In Fig. 1, the robot must reason about properties, affordances, and resources: the slab is heavy, and lifting it would draw significant battery power, constraining subsequent actions. Graph-based and neurosymbolic methods have shown promise in encoding such compositional structure, [20, 21], but scaling them to open world modeling remains challenging [20, 22].

Hypothesis Canvas. In order to instantiate and maintain multiple parallel trajectories (τ_{cf}) in the inference process (§ 4.1), we propose the use of an external memory workspace or canvas. A world represented as a graph of entities and relations can be copied and modified on this canvas, creating distinct subgraphs for each CfS trajectory.

4.3 Pillar 3: Training and Evaluation

Training and evaluation should prioritize logical and physical consistency over reconstruction accuracy, given the lack of ground truth for counterfactuals.

Training. We need training objectives that encourage *interventional prediction*: given an initial trajectory τ_{fact} , identify activation event $e_{t=k}^{fact}$, and predict a simulation e_{t+1}^{cf} . Interventional objectives force world models to capture causal relationships rather than mere correlations.

Evaluation. For CfS, given the scarcity of ground truth for counterfactual simulation, is it more meaningful to evaluate models by verifying their adherence to logical and physical constraints. Rather than speculating what a specific alternate world “should” look like, *constraint verifiers* can validate that simulations respect domain rules (e.g., conservation laws, plausible dynamics of lift, consistent shading/shadows, compatibility between mass and motion). Beyond evaluating the quality of CfS, we also need to validate the *self-reflection* capabilities of the world model: how does the world model agent decide when to use a counterfactual trajectory based on the outcomes of CfS generated.

5 Discussion

Developing world models with **counterfactual simulation (CfS)** will enable robustness in critical domains like robotics and healthcare, where reasoning about novel situations is key to safe and effective operation. However, there are technical and safety challenges that need to be met. In this section, we discuss some research questions that must be addressed.

How can we build these models? Building a single, all-encompassing causal model of the world is currently computationally infeasible. A more practical path may involve an ensemble of smaller, context-specific models that are more efficient and adaptable, and specialized memory modules for storing and retrieving CfS.

How would a model initially learn what constitutes a causally significant “activation event” without already having causal understanding? Two strategies can bridge this gap. First, the model can leverage statistical uncertainty as a diagnostic signal to identify causally significant events. High-uncertainty scenarios then become targets for deeper causal analysis through counterfactual simulation. Second, the model can proactively seek out surprising or anomalous scenarios by continuously predicting potential outcomes and prioritizing exploration where prediction errors are highest, following works in surprise detection [23], thereby systematically building causal understanding in regions where statistical patterns prove insufficient.

How would you manage the exponential growth of possible counterfactuals? The number of counterfactual scenarios can be infinite in an unconstrained world model. Therefore, it is imperative to develop a mechanism where counterfactuals are generated from most plausible to the least plausible. A dynamically determined *counterfactual budget* will only allow for the most impactful and most likely CfS. This *counterfactual budget* would require filtering events based on a specific combination of criteria. First, simulations must respect contextual plausibility, exploring variations that are physically and situationally coherent rather than exhaustively enumerating all logically possible events. Second, the system should prioritize unexpected, rare, or surprising scenarios that deviate from statistical norms, as these often reveal hidden causal structures. Third, selection should favor causally significant and high-stakes situations where understanding intervention outcomes has substantial downstream consequences for decision-making. Only the counterfactual scenarios that fit these criterion could be simulated.

How do we measure the success rate of a system that is capable of CfS? Directly evaluating CfS capability can be challenging. Instead, indirect evaluation through downstream benchmarks could be a feasible approach. Such benchmarks require collecting scenarios where counterfactual outcomes, that may deviate from the norm, are the valid outcomes. Additionally, carefully designed simulations can be used in building such benchmarks. Another aspect of evaluating CfS capability in world models is improving interpretability.

How do we design initial experiments to encourage CfS? An initial experiment design can follow in the footsteps of previous works that use small simulated worlds to evaluate models in. Toy worlds can be designed with novel constraints and rewards that require world models to identify activation events, store hypotheses, and perform CfS. Furthermore, such worlds can be compositional and verifiable.

What are the ethical and safety implications? A model that can simulate “what if” scenarios can imagine both beneficial and harmful outcomes. It is critical to implement safeguards that prevent the model from acting on dangerous simulations. The core challenge lies in aligning the model to use this powerful capability solely for constructive and safe exploration. Furthermore, the latent reasoning behind a counterfactual simulation may be a black box. For these models to be trustworthy, they must be able to decode their simulations into human-understandable formats (like text or video), providing transparency into their decision-making process.

6 Related Work

Prior works have considered the use of counterfactuals, especially in world modeling for autonomous driving. OmniDrive [24] considers counterfactual questions in driving scenarios. GAIA-2 [25] enables OOD scene generation for synthetic driving data, akin to counterfactual simulations we propose. OCTET [26] generates counterfactual explanations in driving scenarios. These research works further justify the need for counterfactual simulation, and show several practical scenarios where it is useful.

Another line of research introduces Counterfactual World Models [27-34]. These methods use masked modeling to train promptable visual world models to perform counterfactual simulations, and to perform counterfactual reasoning in agent models in the context of RL. These works are clear starting points for researchers interested in working on visual world modeling. Our goal is to achieve open-world, real-life counterfactual simulation, as imagined in our three pillar process (§ 4).

Finally, researchers have investigated core imagination and latent world models, that simulate alternative futures. Dreamer V3 [35], a world model-based RL algorithm that learns by simulating parallel futures, has achieved remarkable success across different domains. This approach validates the power of parallel simulation for learning, and serves as a proof-of-concept for the power of considering alternatives. Plan2Explore [36] initiates a two step process that builds a ‘global world model’ that understands the environment’s general dynamics and then, adapts to solve new unseen tasks in a zero or few-shot manner. While these works demonstrate the value of simulating counterfactuals, we argue for moving beyond data-driven approaches toward systems that predict outcomes from first-principles of physics rather than prior observation.

References

- [1] T. Brooks, B. Peebles, C. Holmes, W. DePue, Y. Guo, L. Jing, D. Schnurr, J. Taylor, T. Luhman, E. Luhman, C. Ng, R. Wang, and A. Ramesh, “Video generation models as world simulators,” 2024.
- [2] P. J. Ball, J. Bauer, F. Belletti, B. Brownfield, A. Ephrat, S. Fruchter, A. Gupta, K. Holsheimer, A. Holynski, J. Hron, C. Kaplanis, M. Limont, M. McGill, Y. Oliveira, J. Parker-Holder, F. Perbet, G. Scully, J. Shar, S. Spencer, O. Tov, R. Villegas, E. Wang, J. Yung, C. Baetu, J. Berbel, D. Bridson, J. Bruce, G. Buttmore, S. Chakera, B. Chandra, P. Collins, A. Cullum, B. Damoc, V. Dasagi, M. Gazeau, C. Gbadamosi, W. Han, E. Hirst, A. Kachra, L. Kerley, K. Kjems, E. Knoepfel, V. Koriakin, J. Lo, C. Lu, Z. Mehring, A. Moufaret, H. Nandwani, V. Oliveira, F. Pardo, J. Park, A. Pierson, B. Poole, H. Ran, T. Salimans, M. Sanchez, I. Saprykin, A. Shen, S. Sidhwani, D. Smith, J. Stanton, H. Tomlinson, D. Vijaykumar, L. Wang, P. Wingfield, N. Wong, K. Xu, C. Yew, N. Young, V. Zubov, D. Eck, D. Erhan, K. Kavukcuoglu, D. Hassabis, Z. Ghahramani, R. Hadsell, A. van den Oord, I. Mosseri, A. Bolton, S. Singh, and T. Rocktäschel, “Genie 3: A new frontier for world models,” 2025.
- [3] D. Ha and J. Schmidhuber, “Recurrent world models facilitate policy evolution,” in *Advances in Neural Information Processing Systems 31*, pp. 2451–2463, Curran Associates, Inc., 2018. <https://worldmodels.github.io>
- [4] Q. Gao, X. Pi, K. Liu, J. Chen, R. Yang, X. Huang, X. Fang, L. Sun, G. Kishore, B. Ai, S. Tao, M. Liu, J. Yang, C.-J. Lai, C. Jin, J. Xiang, B. Huang, Z. Chen, D. Danks, H. Su, T. Shu, Z. Ma, L. Qin, and Z. Hu, “Do vision-language models have internal world models? towards an atomic evaluation,” 2025.
- [5] X. Zhou, D. Liang, S. Tu, X. Chen, Y. Ding, D. Zhang, F. Tan, H. Zhao, and X. Bai, “Hermes: A unified self-driving world model for simultaneous 3d scene understanding and generation,” *ArXiv*, vol. abs/2501.14729, 2025.
- [6] Y. Lu, X. Ren, J. Yang, T. Shen, Z. Wu, J. Gao, Y. Wang, S. Chen, M. Chen, S. Fidler, and J. Huang, “Infinitube: Unbounded and controllable dynamic 3d driving scene generation with world-guided video models,” 2024.
- [7] M. Assran, A. Bardes, D. Fan, Q. Garrido, R. Howes, Mojtaba, Komeili, M. Muckley, A. Rizvi, C. Roberts, K. Sinha, A. Zholus, S. Arnaud, A. Gejji, A. Martin, F. R. Hogan, D. Dugas, P. Bojanowski, V. Khalidov, P. Labatut, F. Massa, M. Szafraniec, K. Krishnakumar, Y. Li, X. Ma, S. Chandar, F. Meier, Y. LeCun, M. Rabbat, and N. Ballas, “V-jepa 2: Self-supervised video models enable understanding, prediction and planning,” 2025.
- [8] B. Kang, Y. Yue, R. Lu, Z. Lin, Y. Zhao, K. Wang, G. Huang, and J. Feng, “How far is video generation from world model? – a physical law perspective,” *arXiv preprint arXiv:2411.02385*, 2024.
- [9] S. Motamed, L. Culp, K. Swersky, P. Jaini, and R. Geirhos, “Do generative video models understand physical principles?,” *arXiv preprint arXiv:2501.09038*, 2025.
- [10] A. Chinchure, S. Ravi, R. Ng, V. Schwartz, B. Li, and L. Sigal, “Black swan: Abductive and defeasible video reasoning in unpredictable events,” in *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 24201–24210, 2025.
- [11] F. Ilievski, B. Hammer, F. van Harmelen, B. Paassen, S. Saralajew, U. Schmid, M. Biehl, M. Bolognesi, X. L. Dong, K. Gashtevski, P. Hitzler, G. Marra, P. Minervini, M. Mundt, A.-C. N. Ngomo, A. Oltramari, G. Pasi, Z. G. Saribatur, L. Serafini, J. Shawe-Taylor, V. Schwartz, G. Skitalinskaya, C. Stachl, G. M. van de Ven, and T. Villmann, “Aligning generalisation between humans and machines,” 2025.
- [12] S. Harnad, “To cognize is to categorize: Cognition is categorization,” in *Handbook of Categorization* (C. Lefebvre and H. Cohen, eds.), Elsevier, 2005.
- [13] N. Van Hoeck, P. D. Watson, and A. K. Barbey, “Cognitive neuroscience of human counterfactual reasoning,” *Frontiers in human neuroscience*, vol. 9, p. 420, 2015.
- [14] J. Pearl, “An introduction to causal inference,” *The International Journal of Biostatistics*, vol. 6, 2009.
- [15] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. teusz Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, “Language models are few-shot learners,” *ArXiv*, vol. abs/2005.14165, 2020.
- [16] C. Snell, J. Lee, K. Xu, and A. Kumar, “Scaling llm test-time compute optimally can be more effective than scaling model parameters,” 2024.

- [17] N. Muennighoff, Z. Yang, W. Shi, X. L. Li, L. Fei-Fei, H. Hajishirzi, L. Zettlemoyer, P. Liang, E. Candès, and T. Hashimoto, “s1: Simple test-time scaling,” 2025.
- [18] DeepSeek-AI, “Deepseek-rl: Incentivizing reasoning capability in llms via reinforcement learning,” 2025.
- [19] I. Shumailov, Z. Shumaylov, Y. Zhao, Y. Gal, N. Papernot, and R. Anderson, “The curse of recursion: Training on generated data makes models forget,” *ArXiv*, vol. abs/2305.17493, 2023.
- [20] M. Mancini, M. F. Naeem, Y. Xian, and Z. Akata, “Learning graph embeddings for open world compositional zero-shot learning,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. PP, pp. 1–1, 2021.
- [21] A. Sehgal, A. Grayeli, J. J. Sun, and S. Chaudhuri, “Neurosymbolic grounding for compositional world models,” in *The Twelfth International Conference on Learning Representations*.
- [22] B. Knyazev, H. de Vries, C. Cangea, G. W. Taylor, A. C. Courville, and E. Belilovsky, “Generative compositional augmentations for scene graph prediction,” *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 15807–15817, 2020.
- [23] S. Ravi, A. Chinchure, R. T. Ng, L. Sigal, and V. Shwartz, “Spike-rl: Video-llms meet bayesian surprise,” *arXiv preprint arXiv:2509.23433*, 2025.
- [24] S. Wang, Z. Yu, X. Jiang, S. Lan, M. Shi, N. Chang, J. Kautz, Y. Li, and J. M. Álvarez, “Omnidrive: A holistic vision-language dataset for autonomous driving with counterfactual reasoning,” *2025 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 22442–22452, 2024.
- [25] L. Russell, A. Hu, L. Bertoni, G. Fedoseev, J. Shotton, E. Arani, and G. Corrado, “Gaia-2: A controllable multi-view generative world model for autonomous driving,” *arXiv preprint arXiv:2503.20523*, 2025.
- [26] M. Zemni, M. Chen, É. Zablocki, H. Ben-younes, P. P’erez, and M. Cord, “Octet: Object-aware counterfactual explanations,” *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 15062–15071, 2022.
- [27] H. Chen, K. Kotar, D. Bear, D. L. K. Yamins, W. Lee, A. Durango, K. T. Feigelis, and R. Venkatesh, “Unifying (machine) vision via counterfactual world modeling,” *ArXiv*, vol. abs/2306.01828, 2023.
- [28] J. E. Fan, K. Smith, D. L. K. Yamins, W. Lee, K. T. Feigelis, K. Jedoui, K. Kotar, F. J. Binder, R. Venkatesh, S. Liu, and H. Chen, “Understanding physical dynamics with counterfactual world modeling,” in *European Conference on Computer Vision*, 2023.
- [29] M. Li, M. Yang, F. Liu, X. Chen, Z. Chen, and J. Wang, “Causal world models by unsupervised deconfounding of physical dynamics,” *ArXiv*, vol. abs/2012.14228, 2020.
- [30] S. Prasanna, K. Farid, R. Rajan, and A. Biedenkapp, “Dreaming of many worlds: Learning contextual world models aids zero-shot generalization,” *ArXiv*, vol. abs/2403.10967, 2024.
- [31] G. Gendron, J. M. Rovzanec, M. Witbrock, and G. Dobbie, “Causal cartographer: From mapping to reasoning over counterfactual worlds,” *ArXiv*, vol. abs/2505.14396, 2025.
- [32] M. Singh, A. Alabdulkarim, G. Mansi, and M. O. Riedl, “Explainable reinforcement learning agents using world models,” *ArXiv*, vol. abs/2505.08073, 2025.
- [33] X.-H. Chen, Y. Yu, Z. Zhu, Z. Yu, Z.-Y. Chen, C. Wang, Y. Wu, H. Wu, R. Qin, R. Ding, and F. Huang, “Adversarial counterfactual environment model learning,” *ArXiv*, vol. abs/2206.04890, 2023.
- [34] H. Wang and J. Zhang, “Genai-based multi-agent reinforcement learning towards distributed agent intelligence: A generative-rl agent perspective,” *ArXiv*, vol. abs/2507.09495, 2025.
- [35] D. Hafner, J. Pasukonis, J. Ba, and T. Lillicrap, “Mastering diverse domains through world models,” 2024.
- [36] R. Sekar, O. Rybkin, K. Daniilidis, P. Abbeel, D. Hafner, and D. Pathak, “Planning to explore via self-supervised world models,” 2020.
- [37] D. B. Rubin, “Estimating causal effects of treatments in randomized and nonrandomized studies,” *Journal of educational Psychology*, vol. 66, no. 5, p. 688, 1974.
- [38] S. Wachter, B. Mittelstadt, and C. Russell, “Counterfactual explanations without opening the black box: Automated decisions and the gdpr,” *Harvard Journal of Law & Technology*, vol. 31, no. 2, pp. 841–887, 2017.

- [39] M. J. Kusner, J. Loftus, C. Russell, and R. Silva, “Counterfactual fairness,” in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 30, 2017.
- [40] P. Howard, A. Madasu, T. Le, G. L. Moreno, A. Bhiwandiwala, and V. Lal, “Socialcounterfactuals: Probing and mitigating intersectional social biases in vision-language models with counterfactual examples,” *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 11975–11985, 2023.
- [41] C. Russell, M. Kusner, J. Loftus, and R. Silva, “When worlds collide: Integrating different counterfactual assumptions in fairness,” in *Proceedings of the NeurIPS Workshop on Fairness, Accountability, and Transparency*, 2017.
- [42] A. Chinchure, P. Shukla, G. Bhatt, K. Salij, K. Hosanagar, L. Sigal, and M. Turk, “Tibet: Identifying and evaluating biases in text-to-image generative models,” in *European Conference on Computer Vision*, pp. 429–446, Springer, 2024.
- [43] J. Wang, X. Liu, Z. Di, Y. Liu, and X. Wang, “T2iat: Measuring valence and stereotypical biases in text-to-image generation,” in *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 2560–2574, 2023.
- [44] P. Shukla, A. Chinchure, E. Diana, A. Tolbert, K. Hosanagar, V. N. Balasubramanian, L. Sigal, and M. Turk, “Mitigate one, skew another? tackling intersectional biases in text-to-image models,” *arXiv preprint arXiv:2505.17280*, 2025.
- [45] Y. Niu, K. Tang, H. Zhang, Z. Lu, X. Hua, and J. rong Wen, “Counterfactual vqa: A cause-effect look at language bias,” *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 12695–12705, 2020.
- [46] F. Feng, J. Zhang, X. He, H. Zhang, and T.-S. Chua, “Empowering language understanding with counterfactual reasoning,” in *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pp. 2226–2236, 2021.
- [47] E. Abbasnejad, D. Teney, A. Parvaneh, J. Q. Shi, and A. van den Hengel, “Counterfactual vision and language learning,” *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 10041–10051, 2020.
- [48] Z. Wu, L. Qiu, A. Ross, E. Akyürek, B. Chen, B. Wang, N. Kim, J. Andreas, and Y. Kim, “Reasoning or reciting? exploring the capabilities and limitations of language models through counterfactual tasks,” in *North American Chapter of the Association for Computational Linguistics*, 2023.
- [49] P. Howard, G. Singer, V. Lal, Y. Choi, and S. Swayamdipta, “Neurocounterfactuals: Beyond minimal-edit counterfactuals for richer data augmentation,” in *Findings of the Association for Computational Linguistics: EMNLP 2022*, pp. 5056–5072, 2022.