# Causally Correct Partial Models for Reinforcement Learning

**Danilo J. Rezende** [* 1]  **Ivo Danihelka** [* 1 2]  **George Papamakarios** [1]  **Nan Rosemary Ke** [3]  **Ray Jiang** [1]
**Theophane Weber** [1]  **Karol Gregor** [1]  **Hamza Merzic** [1]  **Fabio Viola** [1]  **Jane Wang** [1]  **Jovana Mitrovic** [1]
**Frederic Besse** [1]  **Ioannis Antonoglou** [1 2]  **Lars Buesing** [1]

## Abstract

In reinforcement learning, we can learn a model of future observations and rewards, and use it to plan the agent's next actions. However, jointly modeling future observations can be computationally expensive or even intractable if the observations are high-dimensional (e.g. images). For this reason, previous works have considered partial models, which model only part of the observation. In this paper, we show that partial models can be causally incorrect: they are confounded by the observations they don't model, and can therefore lead to incorrect planning. To address this, we introduce a general family of partial models that are provably causally correct, yet remain fast because they do not need to fully model future observations.

## 1. Introduction

The ability to predict future outcomes of hypothetical decisions is a key aspect of intelligence. One approach to capture this ability is via *model-based reinforcement learning* (MBRL) (Munro, 1987; Werbos, 1987; Nguyen & Widrow, 1990; Schmidhuber, 1991). In this framework, an agent builds an internal representation $s_t$ by sensing an environment through observational data $y_t$ (such as rewards, visual inputs, proprioceptive information) and interacts with the environment by taking actions $a_t$ according to a policy $\pi(a_t|s_t)$. The sensory data collected is used to build a model that typically predicts future observations $y_{>t}$ from past actions $a_{\leq t}$ and past observations $y_{\leq t}$. The resulting model may be used in various ways, e.g. for planning (Oh et al., 2017; Schrittwieser et al., 2019), generation of synthetic training data (Weber et al., 2017), better credit assignment (Heess et al., 2015), learning useful internal representations and belief states (Gregor et al., 2019; Guo et al., 2018), or

*Equal contribution  [1]DeepMind  [2]University College London  [3]Mila, Université de Montréal.  Correspondence to: <{danilor,danihelka}@google.com>.

exploration via quantification of uncertainty or information gain (Pathak et al., 2017).

Within MBRL, commonly explored methods include action-conditional, next-step models (Oh et al., 2015; Ha & Schmidhuber, 2018; Chiappa et al., 2017; Schmidhuber, 2010; Xie et al., 2016; Deisenroth & Rasmussen, 2011; Lin & Mitchell, 1992; Li et al., 2015; Diuk et al., 2008; Igl et al., 2018; Ebert et al., 2018; Kaiser et al., 2019; Janner et al., 2019). However, it is often not tractable to accurately model all the available information. This is both due to the fact that conditioning on high-dimensional data such as images would require modeling and generating images in order to plan over several timesteps (Finn & Levine, 2017), and to the fact that modeling images is challenging and may unnecessarily focus on visual details which are not relevant for acting. These challenges have motivated researchers to consider simpler models, henceforth referred to as *partial models*, i.e. models which are neither conditioned on, nor generate the full set of observed data (Oh et al., 2017; Amos et al., 2018; Guo et al., 2018; Gregor et al., 2019). A notable example of a partial model is the MuZero model (Schrittwieser et al., 2019).

In this paper, we demonstrate that partial models will often fail to make correct predictions under a new policy, and link this failure to a problem in causal reasoning. Prior to this work, there has been a growing interest in combining causal inference with RL research in the directions of non-model-based bandit algorithms (Bareinboim et al., 2015; Forney et al., 2017; Zhang & Bareinboim, 2017; Lee & Bareinboim, 2018; Bradtke & Barto, 1996; Lu et al., 2018) and causal discovery with RL (Zhu & Chen, 2019). Contrary to previous works, in this paper we focus on model-based approaches and propose a novel framework for learning better partial models. A key insight of our methodology is the fact that any piece of information about the state of the environment that is used by the policy to make a decision, but is not available to the model, acts as a confounding variable for that model. As a result, the learned model is causally incorrect. Using such a model to reason may lead to the wrong conclusions about the optimal course of action as we demonstrate in this paper.
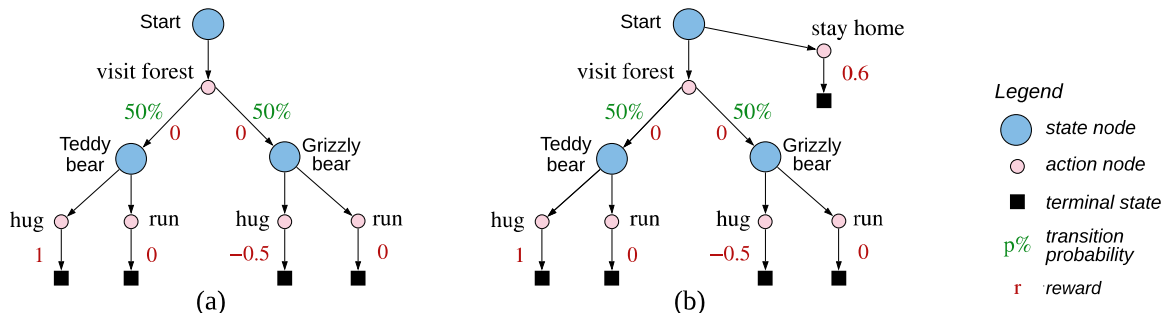
*Figure 1.* Examples of stochastic MDPs. (a) **FuzzyBear:** after visiting a forest, the agent meets either a teddy bear or a grizzly bear with 50% chance and can either hug the bear or run away. (b) **AvoidFuzzyBear:** here, the agent has the extra option to stay home.

We address these issues of partial models by combining general principles of causal reasoning, probabilistic modeling and deep learning. Our contributions are as follows.

- We identify and clarify a fundamental problem of partial models from a causal-reasoning perspective and illustrate it using simple, intuitive Markov Decision Processes (MDPs) (Section 2).

- In order to tackle these shortcomings we examine the following question: What is the minimal information that we have to condition a partial model on such that it will be causally correct with respect to changes in the policy? (Section 4)

- We answer this question by proposing a family of viable solutions and empirically investigate their effects on models learned in illustrative environments (simple MDPs, MiniPacman and 3D environments). Our method is described in Section 4, with experiments in Section 5.

## 2. A simple example: FuzzyBear

We illustrate the issues with partial models using a simple example. Consider the *FuzzyBear* MDP shown in Figure 1(a): an agent at initial state $s_0$ transitions into an encounter with either a teddy bear or a grizzly bear with 50% random chance, and can then take an action to either hug the bear or run away. In order to plan, the agent may learn a partial model $q_\theta(r_2|s_0, a_0, a_1)$ that predicts the reward $r_2$ after performing actions $\{a_0, a_1\}$ starting from state $s_0$. This model is *partial* because it conditions on a sequence of actions without conditioning on the intermediate state $s_1$. The model is suitable for deterministic environments, but it will have problems on stochastic environments, as we shall see. Such a reward model is usually trained on the agent's experience which consists of sequences of past actions and associated rewards.

Now, suppose the agent wishes to evaluate the sequence of actions $\{a_0 = visit forest, a_1 = hug\}$ using the average reward under the model $q_\theta(r_2|s_0, a_0, a_1)$. From Figure 1(a),

we see that the correct average reward is $0.5 \times 1 + 0.5 \times (-0.5) = 0.25$. However, if the model has been trained on past experience in which the agent has mostly hugged the teddy bear and ran away from the grizzly bear, it will learn that the sequence $\{visit forest, hug\}$ is associated with a reward close to 1, and that the sequence $\{visit forest, run\}$ is associated with a reward close to 0. Mathematically, the model will learn the following conditional probability:

$$p(r_2|s_0, a_0, a_1) = \sum_{s_1} p(s_1|s_0, a_0, a_1)p(r_2|s_1, a_1)$$

$$= \sum_{s_1} \frac{p(s_1|s_0, a_0)\pi(a_1|s_1)}{\sum_{s_1'} p(s_1'|s_0, a_0)\pi(a_1|s_1')}p(r_2|s_1, a_1),$$

where $s_1$ is the state corresponding to either *teddy bear* or *grizzly bear*. In the above expression, $p(s_1|s_0, a_0)$ and $p(r_2|s_1, a_1)$ are the transition and reward dynamics of the MDP, and $\pi(a_1|s_1)$ is the agent's behavior policy that generated its past experience. As we can see, the behavior policy affects what the model learns.

The fact that the reward model $q_\theta(r_2|s_0, a_0, a_1)$ is not robust to changes in the behavior policy has serious implications for planning. For example, suppose that instead of visiting the forest, the agent could have chosen to stay at home as shown in Figure 1(b). In this situation, the optimal action is to stay home as it gives a reward of 0.6, whereas visiting the forest gives at most a reward of $0.5 \times 1 + 0.5 \times 0 = 0.5$. However, an agent that uses the above reward model to plan will overestimate the reward of going into the forest as being close to 1 and choose the suboptimal action.[1]

One way to avoid this bias is to use a behavior policy that doesn't depend on the state $s_1$, i.e. $\pi(a_1|s_1) = \pi(a_1)$. Unfortunately, this approach does not scale well to complex environments as it requires an enormous amount of training data for the behavior policy to explore interesting states. A better approach is to make the model robust to changes in

---

[1] This problem is not restricted to toy examples. In a medical domain, a model could learn that leaving the hospital increases the probability of being healthy.
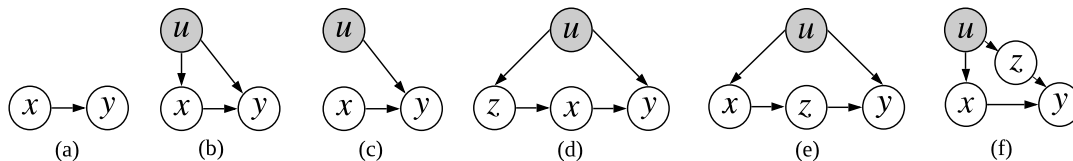
*Figure 2.* Illustration of various causal graphs. (a) Simple dependence without confounding. This is the prevailing assumption in many machine-learning applications. (b) Graph with confounding. (c) Intervention on graph (b) equivalent to setting the value of $x$ and observing $y$. (d) Graph with a backdoor $z$ blocking all paths from $u$ to $x$. (e) Graph with a frontdoor $z$ blocking all paths from $x$ to $y$. (f) Graph with a variable $z$ blocking the direct path from $u$ to $y$.

the behavior policy. Fundamentally, the problem is due to *causally incorrect reasoning*: the model learns the *observational conditional* $p(r_2|s_0, a_0, a_1)$ instead of the *interventional conditional* given by:

$$p(r_2|s_0, \mathrm{do}(a_0), \mathrm{do}(a_1)) = \sum_{s_1} p(s_1|s_0, a_0)p(r_2|s_1, a_1),$$

where the *do-operator* $\mathrm{do}(\cdot)$ means that the actions are performed *independently* of the unspecified context (i.e. independently of $s_1$). The interventional conditional is robust to changes in the policy and is a more appropriate quantity for planning.

In contrast, the observational conditional quantifies the statistical association between the actions $a_0, a_1$ and the reward $r_2$ regardless of whether the actions caused the reward or the reward caused the actions. In Section 3, we review relevant concepts from causal reasoning, and based on them we propose solutions that address the problem.

Finally, although using $p(r_2|s_0, \mathrm{do}(a_0), \mathrm{do}(a_1))$ leads to causally correct planning, it is not optimal either: it predicts a reward of $0.25$ for the sequence $\{visit\,forest, hug\}$ and $0$ for the sequence $\{visit\,forest, run\}$, whereas the optimal policy obtains a reward of $0.5$. The optimal policy makes the decision after observing $s_1$ (teddy bear vs grizzly bear); it is *closed-loop* as opposed to *open-loop*. The solution is to make the intervention at the *policy* level instead of the *action* level, as we discuss in the following sections.

## 3. Background on causal reasoning

Many applications of machine learning involve predicting a variable $y$ (target) from a variable $x$ (covariate). A standard way to make such a prediction is by fitting a model $q_\theta(y|x)$ to a dataset of $(x, y)$-pairs. Then, if we are given a new $x$ and the data-generation process hasn't changed, we can expect that a well trained $q_\theta(y|x)$ will make an accurate prediction of $y$.

**Confounding:** In many situations however, we would like to use the data to make different kinds of predictions. For example, what prediction of $y$ should we make, if something in the environment has changed, or if we set $x$ our-

selves? In the latter case $x$ didn't come from the original data-generation process. This may cause problems in our prediction, because there may be unobserved variables $u$, known as *confounders*, that affected both $x$ and $y$ during the data-generation process. That is, the actual process was of the form $p(u)p(x|u)p(y|x, u)$ where we only observed $x$ and $y$ as shown in Figure 2(b). Under this assumption, a model $q_\theta(y|x)$ fitted on $(x, y)$-pairs will converge to the target $p(y|x) \propto \int p(u)p(x|u)p(y|x, u)du$. However, if at prediction time we set $x$ ourselves, the actual distribution of $y$ will be $p(y|\mathrm{do}(x)) = \int p(u)p(y|x, u)du$. This is because setting $x$ ourselves changes the original graph from Figure 2(b) to the one in Figure 2(c).

**Interventions:** The operation of setting $x$ to a fixed value $x_0$ independently of its parents, known as the *do-operator* (Pearl et al., 2016), changes the data-generation process to $p(u)\delta(x - x_0)p(y|x, u)$, where $\delta(x - x_0)$ is the delta-function. As explained above, this results in a different target distribution $\int p(u)p(y|x_0, u)du$, which we refer to as $p(y|\mathrm{do}(x = x_0))$, or simply $p(y|\mathrm{do}(x))$ when $x_0$ is implied. Let $\mathrm{par}_j$ be the parents of $x_j$. The do-operator is a particular case of the more general concept of an *intervention*: given a generative process $p(x) = \prod_j p_j(x_j|\mathrm{par}_j)$, an intervention is defined as a change that replaces one or more factors by new factors. For example, the intervention $p_k(x_k|\mathrm{par}_k) \to \psi_k(x_k|\mathrm{par}'_k)$ changes $p(x)$ to $p(x)\frac{\psi_k(x_k|\mathrm{par}'_k)}{p_k(x_k|\mathrm{par}_k)}$. The do-operator is a "hard" intervention whereby we replace a node by a delta function; that is, $p(x_{/k}, \mathrm{do}(x_k = v)) = p(x)\frac{\delta(x_k - v)}{p_k(x_k|\mathrm{par}_k)}$, where $x_{/k}$ denotes the collection of all variables except $x_k$.

**Backdoors and frontdoors:** In general, for graphs of the form of Figure 2(b), $p(y|x)$ does not equal $p(y|\mathrm{do}(x))$. As a consequence, it is not generally possible to recover $p(y|\mathrm{do}(x))$ using observational data, i.e. $(x, y)$-pairs sampled from $p(x, y)$, regardless of the amount of data available or the expressivity of the model. However, recovering $p(y|\mathrm{do}(x))$ from observational data alone becomes possible if we assume additional structure in the data-generation process. Suppose there exists another observed variable $z$ that blocks all paths from the confounder $u$ to the covariate $x$ as shown in Figure 2(d). This variable is a particular case of

the concept of a *backdoor* (Pearl et al., 2016, Chapter 3.3) and is said to be a backdoor for the pair $x - y$. In this case, we can express $p(y|\text{do}(x))$ entirely in terms of distributions that can be obtained from the observational data as:

$$p(y|\text{do}(x)) = \mathbb{E}_{p(z)}[p(y|z, x)]. \qquad (1)$$

This formula holds as long as $p(x|z) > 0$ and is referred to as *backdoor adjustment*. The same formula applies when $z$ blocks the effect of the confounder $u$ on $y$ as in Figure 2(f). More generally, we can use $p(z)$ and $p(y|z, x)$ from Equation (1) to compute the marginal distribution $p(y)$ under an arbitrary intervention of the form $p(x|z) \to \psi(x|z)$ on the graph in Figure 2(b). We refer to the new marginal as $p_{\text{do}(\psi)}(y)$ and obtain it by:

$$p_{\text{do}(\psi)}(y) = \mathbb{E}_{\psi(x|z)p(z)}[p(y|z, x)]. \qquad (2)$$

A similar formula can be derived when there is a variable $z$ blocking the effect of $x$ on $y$, which is known as a *frontdoor*, shown in Figure 2(e). Derivations for the backdoor and frontdoor adjustment formulas are provided in Appendix A.

**Causally correct models:** Given data generated by an underlying generative process $p(x)$, we say that a learned model $q_\theta(x)$ is causally correct with respect to a set of interventions $\mathcal{I}$ if the model remains accurate after any intervention in $\mathcal{I}$. That is, if $q_\theta(x) \approx p(x)$ and $q_\theta(x)$ is causally correct with respect to $\mathcal{I}$, then $q_{\theta,\text{do}(\psi)}(x) \approx p_{\text{do}(\psi)}(x)$ for all $\text{do}(\psi)$ in $\mathcal{I}$.

**Connection to MBRL:** As we will see in much greater detail in the next section, there is a direct connection between partial models in MBRL and the causal concepts discussed above. In MBRL we are interested in making predictions about some aspect of the future (observed frames, rewards, etc.); these would be the dependent variables $y$. Such predictions are conditioned on actions which play the role of the covariates $x$. When using partial models, the models will not have access to the full state of the policy and so the policy's state will be a confounding variable $u$. Any variable in the computational graph of the policy that mediates the effect of the state in the actions will be a backdoor with respect to the action-prediction pair.

**Backdoor-adjustment and importance sampling:** Importance-sampling has been extensively explored for off-policy policy evaluation (Sutton & Barto, 2018; Precup, 2000; Precup et al., 2001; Munos et al., 2016; Espeholt et al., 2018; Nachum et al., 2018) as well as for counter-factual policy evaluation (Buesing et al., 2019) and for counter-factual model evaluation (Bottou et al., 2013). Given a dataset of $N$ tuples $(z_n, x_n, y_n)$ generated from the joint distribution $p(u)p(z|u)p(x|z)p(y|x, u)$, we could alternatively approximate the marginal distribution $p_{\text{do}(\psi)}(y)$ after an intervention $p(x|z) \to \psi(x|z)$ by fitting

a distribution $q_\theta(y)$ to maximize the re-weighted likelihood:

$$L(\theta) = \mathbb{E}_{p(u)p(z|u)p(x|z)p(y|x,u)}[w(x, z) \log q_\theta(y)]$$
$$\approx \frac{1}{N} \sum_n w(x_n, z_n) \log q_\theta(y_n), \quad (3)$$

where $w(x, z) = \psi(x|z)/p(x|z)$ are the importance weights. While this solution is a mathematically sound way of obtaining $p_{\text{do}(\psi)}(y)$, it requires re-fitting of the model for any new $\psi(x|z)$. Moreover, if $\psi(x|z)$ is very different from $p(x|z)$ the importance weights $w(x, z)$ will have high variance. By fitting the conditional distribution $p(y|z, x)$ and using Equation (2) we can avoid these limitations.

## 4. Learning causally correct models

We consider environments with a hidden state $e_t$ and dynamics specified by an unknown transition probability of the form $p(e_t|e_{t-1}, a_{t-1})$. At each step $t$, the environment receives an action $a_{t-1}$, updates its state to $e_t$ and produces observable data $y_t \sim p(y_t|e_t)$ which includes a reward $r_t$ and potentially other forms of data such as images. An agent with internal state $s_t$ interacts with the environment via actions $a_t$ produced by a policy $\pi(a_t|s_t)$ and updates its state using the observations $y_{t+1}$ by $s_{t+1} = f_s(s_t, a_t, y_{t+1})$, where $f_s$ can for instance be implemented with an RNN. Figure 3(a) illustrates the interaction between the agent and the environment.

Consider an agent at an arbitrary point in time and whose current state[2] is $s_0$, and assume we are interested in generative models that can predict the outcome[3] $y_T$ of a sequence of actions $\{a_0, \ldots, a_{T-1}\}$ on the environment, for an arbitrary time $T$. A first approach, shown in Figure 3(c), would be to use an action-conditional autoregressive model of observations; initializing the model state $h_1$ to a function of $(s_0, a_0)$, sample $y_1$ from $p(.|h_1)$, update the state $h_2 = f_s(h_1, a_1, y_1)$, sample $y_2$ from $p(.|h_2)$, and so on until $y_T$ is sampled. In other words, the prediction of observation $y_T$ is conditioned on all available observations ($s_0$, $y_{<T}$) and actions $a_{<T}$. This approach is for instance found in (Oh et al., 2015).

In contrast, another approach is to predict observation $y_T$ given actions but using no observation data beyond $s_0$. This family of models, sometimes called *models with overshoot*, can for instance be found in (Oh et al., 2017; Luo et al., 2019; Guo et al., 2018; Hafner et al., 2018; Gregor et al.,

---

[2]We reindex time for notational simplicity, but recall that $s_0$ is indeed a function of past observations $y_{\leq 0}$.

[3]For full generality, it may be that the predicted observation is only a subset or simple function of the full observation $y_T$; for instance one could predict only future rewards. For notational simplicity we make no difference between the full observation and the prediction.
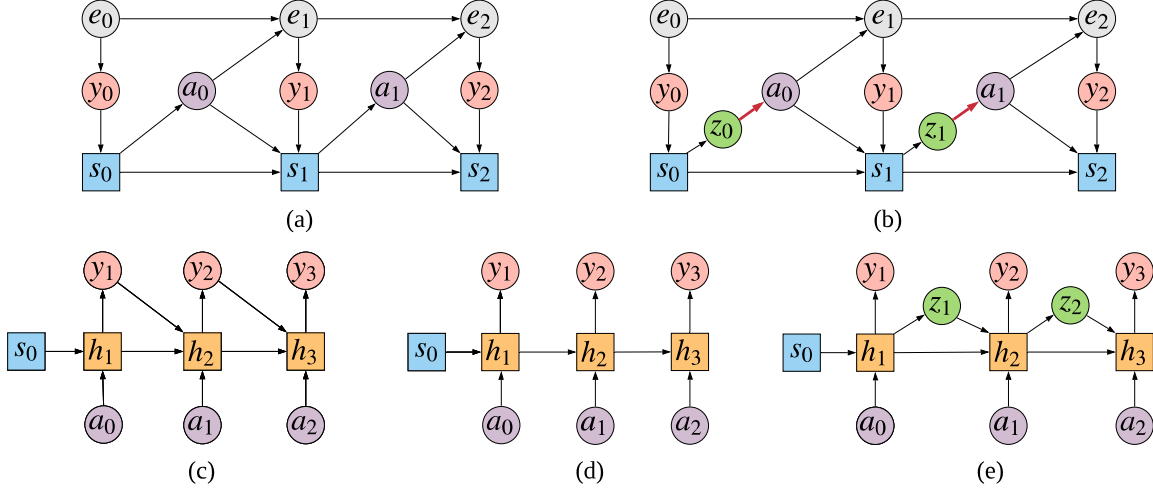
*Figure 3.* Graphical representations of the environment, the agent, and the various models. Circles are stochastic nodes, rectangles are deterministic nodes. (a) Agent interacting with the environment, generating a trajectory $\{y_t, a_t\}_{t=0}^{T}$. These trajectories are the training data for the models. (b) Same as (a) but also including the backdoor $z_t$ in the generated trajectory. The red arrows indicate the locations of the interventions. (c) Standard autoregressive generative model of observations. The model predicts the observation $y_t$ which it then feeds into $h_{t+1}$. (d) Example of a Non-Causal Partial Model (NCPM) that predicts the observation $y_t$ without feeding it into $h_{t+1}$. (e) Proposed Causal Partial Model (CPM), with a backdoor $z_t$ for the actions.

2019; Asadi et al., 2019) and MuZero (Schrittwieser et al., 2019) and is illustrated in Figure 3(d). The model deterministically updates its state $h_{t+1} = f_h(h_t, a_t)$, and generates $y_T$ from $p(.|h_T)$. An advantage of those models is that they can generate $y_T$ directly without generating intermediate observations.

More generally, we define a *partial view* $v_t$ as any function of past observations $y_{\leq t}$ and actions $a_{\leq t}$. We define a *partial model* as a generative model whose predictions are only conditioned on $s_0$, the partial views $v_{<t}$ and the actions $a_{<t}$: to generate $y_T$, the agent generates $v_1$ from $p(.|h_1)$, updates the state to $h_2 = f_h(h_1, v_1, a_1)$, and so on, until it has computed $h_T$ and sampled $y_T$ from $p(.|h_T)$. Both previous examples can be seen as special cases of a partial model, with $v_t = y_t$ and $v_t = \varnothing$ respectively.

A subtle consequence of conditioning the model only on a partial view $v_t$ is that the variables $y_{<T}$ become confounders for predicting $y_T$. The $y_{<T}$ are confounders because they are used by the policy to produce the actions $a_{<T}$, while the partial model is missing the $y_{<T}$ in its input. In Section 3 we showed that the presence of confounders makes it impossible to correctly predict the target distribution after changes in the covariate distribution. In the context of partial models, the covariates are the actions $a_{<T}$ executed by the agent and the agent's initial state $s_0$, whereas the targets are the predictions $y_T$ we want to make at time $T$. A corollary of this is that the learned partial model may not be robust against changes in the policy and thus cannot be used to make predictions under different policies $\pi$, and therefore should not be used for planning.

In Section 3 we saw that if there was a variable blocking the influence of the confounders on the covariates (a backdoor) or a variable blocking the influence of the covariates on the targets (a frontdoor), it may be possible to make predictions under a broad range of interventions if we learn the correct components from data, e.g. using the backdoor-adgustment formula in Equation (2). In general it may not be straightforward to apply the backdoor-adjustment formula because we may not have enough access to the graph details to know which variable is a backdoor. In reinforcement learning however, we can fully control the agent's graph. This means that *we can choose any node in the agent's computational graph that is between its internal state $s_t$ and the produced action $a_t$ as a backdoor variable for the actions*. Given the backdoor $z_t$, the action $a_t$ is conditionally independent of the agent state $s_t$.

To make partial models causally correct, **we propose to choose the partial view $v_t$ to be equal to the backdoor** $z_t$. This allows us to learn all components we need to *make predictions under an arbitrary new policy $\psi(a_t|h_t, z_t)$*. In the rest of this paper we will refer to such models as *Causal Partial Models* (CPM), and all other partial models will be henceforth referred to as *Non-Causal Partial Models* (NCPM). We assume the backdoor $z_t$ is sampled from a distribution $m(z_t|s_t)$ and the policy is a distribution conditioned on $z_t$, $\pi(a_t|z_t)$. This is illustrated in Figure 3(b) and described in more details in Table 1(right). We can perform a simulation under a new policy $\psi(a_t|h_t, z_t)$ by directly applying the backdoor-adjustment formula, Equation (1), to

*Table 1.* Comparison between non-causal partial model and the proposed architecture. The shaded cells indicate the key differences in architectures.

| | | NCPM architecture (overshoot) | | CPM architecture | |
|---|---|---|---|---|---|
| **Agent** | | | | Backdoor | $z_t \sim m(z_t\|s_t)$ |
| | Action | $a_t \sim \pi(a_t\|s_t)$ | | Action | $a_t \sim \pi(a_t\|z_t)$ |
| | State Update | $s_{t+1} = \mathrm{RNN}_s(s_t, a_t, y_{t+1})$ | | State Update | $s_{t+1} = \mathrm{RNN}_s(s_t, a_t, y_{t+1})$ |
| **Model** | State Init. | $h_1 = g(s_0, a_0)$ | | State Init. | $h_1 = g(s_0, a_0)$ |
| | | | | Backdoor | $z_t \sim p(z_t\|h_t)$ |
| | State Update | $h_{t+1} = \mathrm{RNN}_h(h_t, a_t)$ | | State Update | $h_{t+1} = \mathrm{RNN}_h(h_t, z_t, a_t)$ |
| | Prediction | $y_t \sim p(y_t\|h_t)$ | | Prediction | $y_t \sim p(y_t\|h_t)$ |

the RL graph as follows:

$$p_{\mathrm{do}(\psi(a_t|h_t,z_t))}(y_{t+1}|h_t) =$$
$$\mathbb{E}_{\psi(a_t|h_t,z_t)p(z_t|h_t)}[p(y_{t+1}|h_{t+1})], \qquad (4)$$

where the components $p(z_t|h_t)$ and $p(y_{t+1}|h_{t+1})$ with $h_{t+1} = f_h(h_t, z_t, a_t)$ can be learned from observational data produced by the agent.

Modern deep-learning agents (as in e.g. Espeholt et al., 2018; Gregor et al., 2019; Ha & Schmidhuber, 2018) have complex graphs, which means that there are many possible choices for the backdoor $z_t$. So an important question is: what are the simplest choices of $z_t$? Below we list a few of the simplest choices we can use and discuss their advantages and trade-offs; more choices for $z_t$ are listed in Appendix C.

**Agent state:** Identifying $z_t$ with the agent's state $s_t$ can be very informative about the future, but this comes at a cost. As part of the generative model, we have to learn the component $p(z_t|h_t)$. This may be difficult in practice when $z_t = s_t$ due to the high-dimensionality of $s_t$, hence performing simulations would be computationally expensive. Another issue is the presence of irrelevant distractors, which would just increase the variance of the simulations.

**Policy probabilities:** The $z_t$ can be the vector of probabilities produced by a policy when we have discrete actions. The vector of probabilities is informative about the underlying state, if different states produce different probabilities.

**Intended action:** The $z_t$ can be the *intended* action before using some form of exploration, e.g. $\varepsilon$-greedy exploration. This is an interesting choice when the actions are discrete, as it is simple to model and, when doing planning, results in a low branching factor which is independent of the complexity of the environment (e.g. in visually rich 3D environments).

The causal correction methods presented in this section can be applied to any partial model. In our experiments, we focus on environment models of the form proposed by Gregor et al. (2019). These models consist of a deterministic "backbone" RNN that integrates actions and other contextual

information. The states of this RNN are then used to condition a generative model of the observed data $y_t$, but the observations are not fed back to the model autoregressively, as shown in Table 1(left). This corresponds to learning a model of the form $p(y_t|s_0, a_0, a_1, \ldots, a_{t-1})$.

We will compare this against our proposed model of the form $p(y_t|s_0, a_0, z_1, a_1, \ldots, z_{t-1}, a_{t-1}) = p(y_t|h_t)$. In this setup, a policy network produces $z_t$ before an action $a_t$. For example, if the $z_t$ is the *intended* action before $\varepsilon$-exploration, $z_t$ will be sampled from a policy $m(z_t|s_t)$ and the *executed* action $a_t$ will then be sampled from an $\varepsilon$-exploration policy $\pi(a_t|z_t) = (1-\varepsilon)\delta_{z_t,a_t} + \varepsilon\frac{1}{n_a}$, where $n_a$ is the number of actions and $\varepsilon$ is in $(0,1)$. It is imperative that we use some form of exploration to ensure that $\pi(a_t|z_t) > 0$ for all $a_t$ and $z_t$ as this is a necessary to allow the model to learn the effects of the actions. Acting with the sampled actions is diagrammed in Figure 3(b) and the mathematical description is provided in Table 1.

The model components $p(z_t|h_t)$ and $p(y_t|h_t)$ are trained via maximum likelihood on observational data collected by the agent. The partial model does not need to model all parts of the $y_t$ observation. For example, a model to be used for planning can model just the reward and the expected return. The model usage is summarized in Algorithms 1 and 2 in Appendix D and we discuss the model properties and limitations in Appendix E.

## 5. Experiments

We analyse the effect of the proposed corrections on a variety of models and environments. When the enviroment is an MDP, such as the FuzzyBear MDP from Section 2, we can compute exactly both the non-causal and the causal model directly from the MDP transition matrix and the behavior policy. In Section 5.1, we compare the optimal policies computed from the non-causal and the causal model via value iteration. For this analysis, we used the intended action as the partial view, since it's compatible with a tabular representation. In Section 5.2, we use learned models instead.

Finally in Section 5.3, we provide an analysis of the model rollouts in a visually rich 3D environment.

## 5.1. Value-iteration analysis on MDPs

Given an MDP and a behavior policy $\pi$, the optimal values $V^*_{M(\pi)}$ of planning based on a NCPM and CPM are derived in Appendix I. The theoretical analysis of the MDP does not use empirically trained models from the policy data, but rather assumes that the transition probabilities of the MDP and the policy from which training data are collected and accurately learned by the model. This allows us to isolate the quality of planning using the model from how accurate the model is.
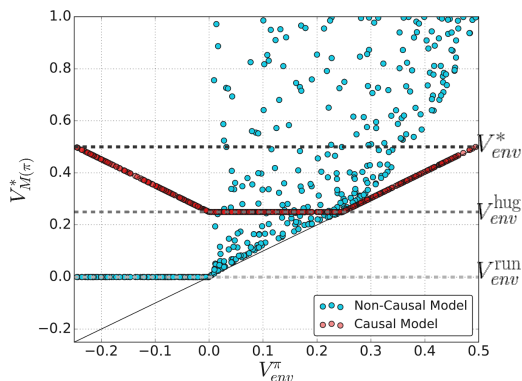


*Figure 4.* MDP Analysis: In the FuzzyBear environment, we randomly generate 500 policies and scatter-plot them with x-axis showing the quality of the behavior policy $V^\pi_{env}$ and y-axis showing corresponding model optimal evaluations $V^*_{M(\pi)}$. For each policy, we derive the corresponding converged model $M(\pi)$ equivalent to training on data generated by the policy. We then compute the optimal evaluation $V^*_{M(\pi)}$ using this model. We contrast the unrealistic optimism of the non-causal model evaluations $V^*_{NCPM(\pi)}$ with the more realistic causal model evaluations $V^*_{CPM(\pi)}$ for good policies $\pi$, as well as the over-pessimism of the non-causal model compared to the causal model for bad policies.

**Sub-optimal behavior policies:** The optimal policy of the FuzzyBear MDP (Figure 1(a)) is to always hug the teddy bear and run away from the grizzly bear. We empirically show the difference between the causal and non-causal models when learning from randomly generated policies. For each policy, we derive the corresponding converged model $M(\pi)$ using training data generated by the policy. We then compute the optimal value of $V^*_{M(\pi)}$ using this model. On FuzzyBear (Figure 4), we see that the causal model always produces a value greater than or equal to the value of the behavior policy. The value estimated by the causal model can always be achieved in the real environment. If the behavior policy was already good, the simulation policy used inside the model can reproduce the behavior policy by respecting the intended action. If the behavior policy is random, the intended action is uninformative about the underlying state,

so the simulation policy has to choose the most rewarding action, independently of the state. And if the behavior policy is bad, the simulation policy can choose the opposite of the intended action. This allows to find a very good simulation policy, when the behavior policy is very bad. To further improve the policy, the search for better policies should be done also in state $s_1$. And the model can then be retrained on data from the improved policies.

If we look at the non-causal model, we see that it displays the unfortunate property of becoming unrealistically optimistic as the behavior policy becomes better.

Similarly, the worse the policy is, i.e. the lower $V^\pi_{env}$ is, the non-causal model becomes less able to improve the policy.

## 5.2. Experiments with Expectimax and MCTS

We will now describe experiments with learned models trained by gradient descent. The models will be compared based on their ability to support a tree search. The tree search is used to find an improved policy, given the model. We will use the classical expectimax search (Michie, 1966; Russell & Norvig, 2009) or a variant of MuZero Monte-Carlo Tree Search (MCTS) (Schrittwieser et al., 2019). When using a causal partial model, we will use the intended action as the partial view $z_t$. The intended action is a good match for the discrete tree search because the intended action has a small number of categorical values, is easy to model and has a low variance and a low branching factor.

**On AvoidFuzzyBear:** On the simple AvoidFuzzyBear MDP (Figure 1(b)), it is enough to use expectimax with a search depth of 3: a decision node, a chance node and a decision node. The policy found by the search was used to produce the next action for the real environment.

Only the non-causal model was not able to solve the task. Expectimax with the non-causal model consistently preferred the sub-optimal stochastic path with the fuzzy bear, as predicted by our theoretical analysis with value iteration. Results with MuZero-style MCTS are in Figure 5(a). MuZero has a number of properties that mitigate the negative effects of the non-causal partial model. We discuss MuZero properties in Appendix F and the experimental setup is described in Appendix G.

**On MiniPacman:** We used the 2D MiniPacman environment (Guez et al., 2019) to test whether the non-causal models have problems also in other stochastic environments. To reduce variance and to remove differences in exploration, we trained all models on data from the same pretrained policy. In Figure 5(b) and Figure 5(c) we indeed see that the non-causal partial model (NCPM) achieved visibly smaller reward when used with expectimax or MCTS. Combining MCTS with chance nodes mitigated some of the negative effects of the NCPM, as discussed in Appendix F.
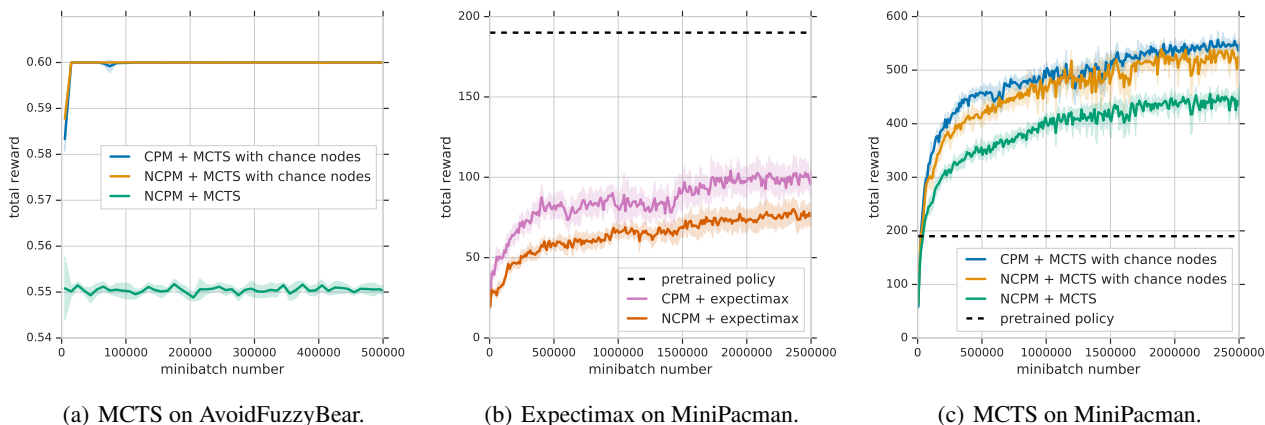
(a) MCTS on AvoidFuzzyBear.  (b) Expectimax on MiniPacman.  (c) MCTS on MiniPacman.

*Figure 5.* (a) MCTS on AvoidFuzzyBear with $p(\textit{Teddy bear}) = 0.55$. The optimal policy should achieve reward 0.6. (b) The non-causal partial model (NCPM) produced visibly worse expectimax search. (c) The causal partial model (CPM) was compatible with MuZero-style MCTS. The search was able to find a much better policy than the pretrained behavior policy.

## 5.3. Visually rich 3D Environment

The setup for these experiments is similar to (Gregor et al., 2019), where the agent is trained using the IMPALA algorithm (Espeholt et al., 2018), and the model is trained alongside the agent via ELBO optimization on the data collected by the agent. For these experiments, the partial view $z_t$ was chosen to be the policy probabilities, and $p(z_t|h_t)$ was parametrized as a mixture of Dirichlet distributions. See Appendix J for more details. We demonstrate the effect of the causal correction on the 3D T-Maze environment where an agent walks around in a 3D world with the goal of collecting the reward blocks (food). The layout of this environment is shown in Figure 6(a). From our previous results, we expect NCPMs to be unrealistically optimistic. This is indeed what we see in Figure 6(b). Compared to NCPM, CPM with generated $z$ generates food at the end of a rollout with around 50% chance, as expected given that the environment randomly places the food on either side. In Figure 6(c) and Figure 6(d, left) we show subsets of rollouts from NCPM and CPM respectively (see Figures 12–14 in Appendix K for full rollouts).

## 6. Conclusion

We have characterized and explained some of the issues of partial models in terms of causal reasoning. We proposed a simple, yet effective, modification to partial models so that they can still make correct predictions under changes in the behavior policy, which we validated theoretically and experimentally. The proposed modifications address the correctness of the model against policy changes, but don't address the correctness/robustness against other types of intervention in the environment. We will explore these aspects in future work.
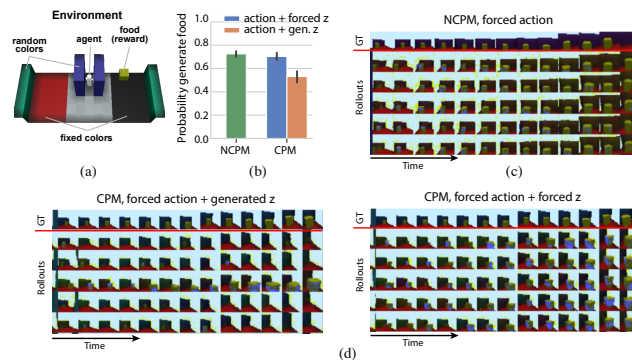


*Figure 6.* Causal partial model (CPM) and non-causal partial model (NCPM) rollouts in a 3D T-Maze environment. (a) The agent always begins in a walled-off corridor, and can obtain food reward that spawns randomly on either the left (red) or the right (black) side. The colors of the corridor and side walls are randomized every episode. (b) Probability of the model generating food in rollouts for NCPM and CPM. (c)–(d) Subset of frames from example rollouts using (c) NCPM and (d) CPM. In all rollouts depicted, the top row shows the real frames observed by an agent following a fixed policy (Ground Truth, GT). Bottom 5 rows indicate model rollouts, conditioned on 3 previous frames without revealing the location of the food. CPM and NCPM differ in their state-update formula and action generation (see Table 1), but frame generation $y_t \sim p(y_t|h_t)$ is the same for both, as introduced in (Gregor et al., 2019). For CPM, we compare rollouts with forced actions and generated $z$ to rollouts with forced actions and forced $z$ from the ground truth. We can observe that rollouts with the generated $z$ (left) respect the randomness in the food placement (with and without food), while the rollouts with forced $z$ (right) consistently generate food blocks, if following actions consistent with the partial view $z$ from the well-trained ground-truth policy.

# References

Amos, B., Dinh, L., Cabi, S., Rothörl, T., Muldal, A., Erez, T., Tassa, Y., de Freitas, N., and Denil, M. Learning awareness models. In *International Conference on Learning Representations*, 2018.

Asadi, K., Misra, D., Kim, S., and Littman, M. L. Combating the compounding-error problem with a multi-step model. *arXiv preprint arXiv:1905.13320*, 2019.

Bareinboim, E., Forney, A., and Pearl, J. Bandits with unobserved confounders: A causal approach. In *Advances in Neural Information Processing Systems*, pp. 1342–1350, 2015.

Bottou, L., Peters, J., Quiñonero-Candela, J., Charles, D. X., Chickering, D. M., Portugaly, E., Ray, D., Simard, P., and Snelson, E. Counterfactual reasoning and learning systems: The example of computational advertising. *The Journal of Machine Learning Research*, 14(1):3207–3260, 2013.

Bradtke, S. J. and Barto, A. G. Linear least-squares algorithms for temporal difference learning. *Machine learning*, 22(1-3):33–57, 1996.

Browne, C. B., Powley, E., Whitehouse, D., Lucas, S. M., Cowling, P. I., Rohlfshagen, P., Tavener, S., Perez, D., Samothrakis, S., and Colton, S. A survey of Monte Carlo tree search methods. *IEEE Transactions on Computational Intelligence and AI in Games*, 4(1):1–43, 2012.

Buesing, L., Weber, T., Zwols, Y., Heess, N., Racanière, S., Guez, A., and Lespiau, J.-B. Woulda, coulda, shoulda: Counterfactually-guided policy search. In *International Conference on Learning Representations*, 2019.

Chiappa, S., Racanière, S., Wierstra, D., and Mohamed, S. Recurrent environment simulators. *arXiv preprint arXiv:1704.02254*, 2017.

Deisenroth, M. P. and Rasmussen, C. E. PILCO: A model-based and data-efficient approach to policy search. In *International Conference on Machine Learning*, 2011.

Diuk, C., Cohen, A., and Littman, M. L. An object-oriented representation for efficient reinforcement learning. In *Proceedings of the 25th International Conference on Machine learning*, pp. 240–247, 2008.

Ebert, F., Finn, C., Dasari, S., Xie, A., Lee, A., and Levine, S. Visual foresight: Model-based deep reinforcement learning for vision-based robotic control. *arXiv preprint arXiv:1812.00568*, 2018.

Espeholt, L., Soyer, H., Munos, R., Simonyan, K., Mnih, V., Ward, T., Doron, Y., Firoiu, V., Harley, T., Dunning, I., et al. IMPALA: Scalable distributed deep-RL with importance weighted actor-learner architectures. *arXiv preprint arXiv:1802.01561*, 2018.

Finn, C. and Levine, S. Deep visual foresight for planning robot motion. In *IEEE International Conference on Robotics and Automation*, pp. 2786–2793, 2017.

Forney, A., Pearl, J., and Bareinboim, E. Counterfactual data-fusion for online reinforcement learners. In *Proceedings of the 34th International Conference on Machine Learning*, pp. 1156–1164, 2017.

Gregor, K., Rezende, D. J., Besse, F., Wu, Y., Merzic, H., and van den Oord, A. Shaping belief states with generative environment models for RL. *arXiv preprint arXiv:1906.09237*, 2019.

Guez, A., Mirza, M., Gregor, K., Kabra, R., Racaniere, S., Weber, T., Raposo, D., Santoro, A., Orseau, L., Eccles, T., et al. An investigation of model-free planning. In *International Conference on Machine Learning*, pp. 2464–2473, 2019.

Guo, Z. D., Azar, M. G., Piot, B., Pires, B. A., Pohlen, T., and Munos, R. Neural predictive belief representations. *arXiv preprint arXiv:1811.06407*, 2018.

Ha, D. and Schmidhuber, J. World models. *arXiv preprint arXiv:1803.10122*, 2018.

Hafner, D., Lillicrap, T., Fischer, I., Villegas, R., Ha, D., Lee, H., and Davidson, J. Learning latent dynamics for planning from pixels. *arXiv preprint arXiv:1811.04551*, 2018.

Heess, N., Wayne, G., Silver, D., Lillicrap, T., Erez, T., and Tassa, Y. Learning continuous control policies by stochastic value gradients. In *Advances in Neural Information Processing Systems*, pp. 2944–2952, 2015.

Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.

Igl, M., Zintgraf, L., Le, T. A., Wood, F., and Whiteson, S. Deep variational reinforcement learning for POMDPs. *arXiv preprint arXiv:1806.02426*, 2018.

Janner, M., Fu, J., Zhang, M., and Levine, S. When to trust your model: Model-based policy optimization. In *Advances in Neural Information Processing Systems*, pp. 12498–12509, 2019.

Kaiser, L., Babaeizadeh, M., Milos, P., Osinski, B., Campbell, R. H., Czechowski, K., Erhan, D., Finn, C., Kozakowski, P., Levine, S., et al. Model-based reinforcement learning for Atari. *arXiv preprint arXiv:1903.00374*, 2019.

Lee, S. and Bareinboim, E. Structural causal bandits: Where to intervene? In *Advances in Neural Information Processing Systems*, pp. 2568–2578, 2018.

Li, X., Li, L., Gao, J., He, X., Chen, J., Deng, L., and He, J. Recurrent reinforcement learning: a hybrid approach. *arXiv preprint arXiv:1509.03044*, 2015.

Lin, L. and Mitchell, T. Memory approaches to reinforcement learning in non-Markovian domains. Technical report, Carnegie Mellon University, 1992.

Lu, C., Schölkopf, B., and Hernández-Lobato, J. M. Deconfounding reinforcement learning in observational settings. *arXiv preprint arXiv:1812.10576*, 2018.

Luo, Y., Xu, H., Li, Y., Tian, Y., Darrell, T., and Ma, T. Algorithmic framework for model-based deep reinforcement learning with theoretical guarantees. In *International Conference on Learning Representations*, 2019.

Michie, D. Game-playing and game-learning automata. In *Advances in programming and non-numerical computation*, pp. 183–200. Elsevier, 1966.

Munos, R., Stepleton, T., Harutyunyan, A., and Bellemare, M. Safe and efficient off-policy reinforcement learning. In *Advances in Neural Information Processing Systems*, pp. 1054–1062, 2016.

Munro, P. A dual back-propagation scheme for scalar reward learning. In *9th Annual Conference of the Cognitive Science Society*, pp. 165–176, 1987.

Nachum, O., Gu, S. S., Lee, H., and Levine, S. Data-efficient hierarchical reinforcement learning. In *Advances in Neural Information Processing Systems*, pp. 3303–3313, 2018.

Nguyen, D. and Widrow, B. The truck backer-upper: An example of self-learning in neural networks. In *Advanced neural computers*, pp. 11–19. Elsevier, 1990.

Oh, J., Guo, X., Lee, H., Lewis, R. L., and Singh, S. Action-conditional video prediction using deep networks in Atari games. In *Advances in Neural Information Processing Systems*, pp. 2863–2871, 2015.

Oh, J., Singh, S., and Lee, H. Value prediction network. In *Advances in Neural Information Processing Systems*, pp. 6118–6128, 2017.

Pathak, D., Agrawal, P., Efros, A. A., and Darrell, T. Curiosity-driven exploration by self-supervised prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 16–17, 2017.

Pearl, J., Glymour, M., and Jewell, N. P. *Causal inference in statistics: A primer*. John Wiley & Sons, 2016.

Precup, D. Eligibility traces for off-policy policy evaluation. *Computer Science Department Faculty Publication Series*, pp. 80, 2000.

Precup, D., Sutton, R. S., and Dasgupta, S. Off-policy temporal difference learning with function approximation. In *Proceedings of the 18th International Conference on Machine Learning*, pp. 417–424, 2001.

Russell, S. and Norvig, P. *Artificial Intelligence: A Modern Approach*. Prentice Hall Press, 3rd edition, 2009. ISBN 0136042597, 9780136042594.

Schmidhuber, J. Curious model-building control systems. In *Proceedings of the IEEE International Joint Conference on Neural Networks*, pp. 1458–1463, 1991.

Schmidhuber, J. Formal theory of creativity, fun, and intrinsic motivation (1990–2010). *IEEE Transactions on Autonomous Mental Development*, 2(3):230–247, 2010.

Schrittwieser, J., Antonoglou, I., Hubert, T., Simonyan, K., Sifre, L., Schmitt, S., Guez, A., Lockhart, E., Hassabis, D., Graepel, T., Lillicrap, T., and Silver, D. Mastering Atari, go, chess and shogi by planning with a learned model. *arXiv preprint arXiv:1911.08265*, 2019.

Silver, D., Sutton, R. S., and Müller, M. Sample-based learning and search with permanent and transient memories. In *Proceedings of the 25th International Conference on Machine learning*, pp. 968–975, 2008.

Sutton, R. S. Integrated architectures for learning, planning, and reacting based on approximating dynamic programming. In *Machine Learning Proceedings*, pp. 216–224. Elsevier, 1990.

Sutton, R. S. and Barto, A. G. *Reinforcement Learning: An Introduction*. The MIT Press, 2nd edition, 2018. ISBN 0262039249.

Weber, T., Racanière, S., Reichert, D. P., Buesing, L., Guez, A., Rezende, D. J., Badia, A. P., Vinyals, O., Heess, N., Li, Y., et al. Imagination-augmented agents for deep reinforcement learning. *arXiv preprint arXiv:1707.06203*, 2017.

Werbos, P. J. Learning how the world works: Specifications for predictive networks in robots and brains. In *Proceedings of the IEEE International Conference on Systems, Man and Cybernetics*, 1987.

Xie, C., Patil, S., Moldovan, T., Levine, S., and Abbeel, P. Model-based reinforcement learning with parametrized physical models and optimism-driven exploration. In

*IEEE International Conference on Robotics and Automation*, pp. 504–511, 2016.

Zhang, J. and Bareinboim, E. Transfer learning in multi-armed bandits: A causal approach. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pp. 1340–1346, 2017.

Zhu, S. and Chen, Z. Causal discovery with reinforcement learning. *arXiv preprint arXiv:1906.04477*, 2019.

# Supplementary to Causally Correct Partial Models for Reinforcement Learning

## A. Backdoor and frontdoor adjustment formulas

Starting from a data-generation process of the form illustrated in Figure 2(b), $p(x, y, u) = p(u)p(x|u)p(y|x, u)$, we can use the do-operator to compute $p(y|\text{do}(x)) = \int p(u)p(y|x, u)du$.

Without assuming any extra structure in $p(x|u)$ or in $p(y|x, u)$ it is not possible to compute $p(y|\text{do}(x))$ from the knowledge of the joint $p(x, y)$ alone.

If there was a variable $z$ blocking all the effects of $u$ on $x$, as illustrated in Figure 2(d), then $p(y|\text{do}(x))$ can be derived as follows:

$$\text{Joint density} \qquad p(u)p(z|u)p(x|z)p(y|x, u) \tag{5}$$

$$\text{Intervention} \qquad p(x|z) \to \psi(x) \tag{6}$$

$$\text{Joint after intervention} \qquad p(u)p(z|u)\psi(x)p(y|x, u) \tag{7}$$

$$\text{Conditioning the new joint} \qquad p(y|\text{do}(x)) = \frac{\int p(u)p(z|u)\psi(x)p(y|x, u)dudz}{\psi(x)} \tag{8}$$

$$= \int p(z)p(y|x, z)dz \tag{9}$$

$$= \mathbb{E}_{p(z)}[p(y|x, z)], \tag{10}$$

where we used the formula

$$\int p(u)p(z|u)p(y|x, u)du = p(z)\int p(u|z)p(y|x, u)du \tag{11}$$

$$= p(z)p(y|x, z). \tag{12}$$

If instead of just fixing the value of $x$, we perform a more general intervention $p(x|z) \to \psi(x|z)$, then $p_{\text{do}(\psi(x|z))}(y)$ can be derived as follows:

$$\text{Joint density} \qquad p(u)p(z|u)p(x|z)p(y|x, u) \tag{13}$$

$$\text{Intervention} \qquad p(x|z) \to \psi(x|z) \tag{14}$$

$$\text{Joint after intervention} \qquad p(u)p(z|u)\psi(x|z)p(y|x, u) \tag{15}$$

$$\text{New marginal} \qquad p_{\text{do}(\psi(x|z))}(y) = \int p(u)p(z|u)\psi(x|z)p(y|x, u)dudzdx \tag{16}$$

$$= \int p(z)\psi(x|z)p(y|x, z)dzdx \tag{17}$$

$$= \mathbb{E}_{p(z)\psi(x|z)}[p(y|x, z)]. \tag{18}$$

Applying the same reasoning to the graph shown in Figure 2(e), we obtain the formula

$$p(y|\text{do}(x)) = \mathbb{E}_{p(z|x)}[p(y|\text{do}(z))] = \mathbb{E}_{p(z|x)p(x')}[p(y|x', z)], \tag{19}$$

where $p(z|x)$, $p(x')$ and $p(y|x', z)$ can be directly measured from the available $(x, y, z)$ data. This formula holds as long as $p(z|x) > 0, \forall x, z$ and it is a simple instance of *frontdoor adjustment* (Pearl et al., 2016).

## B. Derivation of causal correctness for the models in Figure 3

Here, we will show in more detail that the models (c) and (e) in Figure 3 are causally correct, whereas model (d) is causally incorrect. Specifically, we will show that given an initial state $s_0$ and after setting the actions $a_0$ to $a_T$ to specific values, models (c) and (e) make the same prediction about the future observation $y_{T+1}$ as performing the intervention in the real world, whereas model (d) does not.

**Model (c)**   Using the do-operator, a hard intervention in the model is given by:

$$q_\theta(y_{T+1}|s_0, \mathrm{do}(a_{0:T})) = q_\theta(y_{T+1}|s_0, a_{0:T}) = \int \prod_{t=1}^{T+1} q_\theta(y_t|h_t) \, dy_{1:T}, \tag{20}$$

where $h_t$ is a deterministic function of $s_0$, $a_{0:t-1}$ and $y_{1:t-1}$. The same hard intervention in the real world is given by:

$$p(y_{T+1}|s_0, \mathrm{do}(a_{0:T})) = \int p(y_{1:T+1}|s_0, \mathrm{do}(a_{0:T})) \, dy_{1:T} \tag{21}$$

$$= \int \prod_{t=1}^{T+1} p(y_t|s_0, \mathrm{do}(a_{0:T}), y_{1:t-1}) \, dy_{1:T} \tag{22}$$

$$= \int \prod_{t=1}^{T+1} p(y_t|s_0, a_{0:t-1}, y_{1:t-1}) \, dy_{1:T}. \tag{23}$$

If the model is trained perfectly, the factors $q_\theta(y_t|h_t)$ will become equal to the conditionals $p(y_t|s_0, a_{0:t-1}, y_{1:t-1})$. Hence, an intervention in a perfectly trained model makes the same prediction as in the real world, which means that the model is causally correct.

**Model (d)**   The interventional conditional in the model is simply:

$$q_\theta(y_{T+1}|s_0, \mathrm{do}(a_{0:T})) = q_\theta(y_{T+1}|s_0, a_{0:T}) = q_\theta(y_{T+1}|h_{T+1}), \tag{24}$$

where $h_{T+1}$ is a deterministic function of $s_0$ and $a_{0:T}$. In a perfectly trained model, we have that $q_\theta(y_{T+1}|h_{T+1}) = p(y_{T+1}|s_0, a_{0:T})$. However, the observational conditional $p(y_{T+1}|s_0, a_{0:T})$ is not generally equal to the inverventional conditional $p(y_{T+1}|s_0, \mathrm{do}(a_{0:T}))$, which means that the model is causally incorrect.

**Model (e)**   Finally, the interventional conditional in this model is:

$$q_\theta(y_{T+1}|s_0, \mathrm{do}(a_{0:T})) = q_\theta(y_{T+1}|s_0, a_{0:T}) = \int q_\theta(y_{T+1}|h_{T+1}) \prod_{t=1}^{T} q_\theta(z_t|h_t) \, dz_{1:T}, \tag{25}$$

where $h_t$ is a deterministic function of $s_0$, $a_{0:t-1}$ and $z_{1:t-1}$. The same intervention in the real world can be written as follows:

$$p(y_{T+1}|s_0, \mathrm{do}(a_{0:T})) = \int p(y_{T+1}|s_0, \mathrm{do}(a_{0:T}), z_{1:T}) p(z_{1:T}|s_0, \mathrm{do}(a_{0:T})) \, dz_{1:T} \tag{26}$$

$$= \int p(y_{T+1}|s_0, \mathrm{do}(a_{0:T}), z_{1:T}) \prod_{t=1}^{T} p(z_t|s_0, \mathrm{do}(a_{0:T}), z_{1:t-1}) \, dz_{1:T} \tag{27}$$

$$= \int p(y_{T+1}|s_0, a_{0:T}, z_{1:T}) \prod_{t=1}^{T} p(z_t|s_0, a_{0:t-1}, z_{1:t-1}) \, dz_{1:T}. \tag{28}$$

In a perfectly trained model, we have that $q_\theta(y_{T+1}|h_{T+1}) = p(y_{T+1}|s_0, a_{0:T}, z_{1:T})$ and $q_\theta(z_t|h_t) = p(z_t|s_0, a_{0:t-1}, z_{1:t-1})$. That means that the intervention in a perfectly trained model makes the same prediction as the same intervention in the real world, hence the model is causally correct.

## C. Additional choices of the backdoor $z_t$

The first alternative backdoor we consider is the empty backdoor:

**Empty backdoor $z_t = \varnothing$:** This backdoor is in general not appropriate; it is however appropriate when the behavior policy does in fact depend on no information, i.e. is not a function of the state $s_t$. For example, the policy can be uniformly random (or any non-state dependent distribution over actions). This severely limits the behavior policy. Because the backdoor contains no information about the observations, the simulations are open-loop, i.e. we can only consider plans which consist of a sequence of fixed actions, not policies.

**An intermediate layer:** In principle, the $z_t$ can be any layer from the policy. To model the layer with a $p(z_t|h_t)$ distribution, we would need to know the needed numerical precision of the considered layer. For example, a quantized layer can be modeled by a discrete distribution. Alternatively, if the layer is produced by a variational encoder or variational information bottleneck, we can train $p(z_t|h_t)$ to minimize the $\mathrm{KL}(p_{encoder}(z_t|s_t) \parallel p(z_t|h_t))$.

Finally, if a backdoor is appropriate, we can combine it with additional information:

**Combinations:** It is possible to combine a layer with information from other layers. For example, the intended action can be combined with extra bits from the input layer. Such $z_t$ can be more informative. For example, the extra bits can hold a downsampled and quantized version of the input layer.

## D. Algorithms for training and simulating from the model

Algorithms 1 and 2 describe how the model is trained and used to simulate trajectories. The algorithm for training assumes a distributed actor-learner setup (Espeholt et al., 2018).

---

**Algorithm 1** Model training

---

    **Data collection on an actor:**
    For each step:
        $z_t \sim m(z_t|s_t)$ ... sample the backdoor (e.g., the partial view with the intended action)
        $a_t \sim \pi(a_t|z_t)$ ... sample the executed action (e.g., add $\varepsilon$-exploration)
    Collect:
        $s_t$ ... agent state
        $z_t$ ... partial view
        $a_t$ ... executed action
        $y_{t+1}$ ... targets (rewards, returns, ... )

    **Model training on a learner:**
    Require a trajectory: $s_0, a_{<T}, z_{<T}, y_{\leq T}$
    $h_1 = g(s_0, a_0)$ ... initialize the model state
    For each trajectory step:
        Train $p(y_t|h_t)$ to model $y_t$.
        Train $p(z_t|h_t)$ to model $z_t$.
        $h_{t+1} = \mathrm{RNN}_h(h_t, z_t, a_t)$ ... update the model state

---

---

**Algorithm 2** Using the model to generate a simulation under a new policy $\psi$

---

    Require an agent state: $s_0$
    $a_0 = \psi(a_0|s_0)$ ... choose the first action
    $h_1 = g(s_0, a_0)$ ... initialize the model state
    For each trajectory step:
        Predict the wanted targets $p(y_t|h_t)$ (e.g., rewards, returns, ... ).
        $z_t \sim p(z_t|h_t)$ ... sample the partial view
        $a_t \sim \psi(a_t|h_t, z_t)$ ... choose the next action
        $h_{t+1} = \mathrm{RNN}_h(h_t, z_t, a_t)$ ... update the model state

---

## E. Discussion of model properties

Table 2 provides an overview of properties of autoregressive models, deterministic non-causal models and the causal partial models. The causal partial models have to generate only a partial view. The partial view can be small and easy to model. For example, a partial view with the discrete intended action can be flexibly modeled by a categorical distribution.

The causal partial models are fast and causally correct in stochastic environments. The causal partial models have a low simulation variance, because they do not need to model and generate unimportant background distractors. If the environment has deterministic regions, the model can quickly learn to ignore the small partial view and collect information only from the executed action.

It is interesting that the causal partial models are invariant of the $\pi(a_t|z_t)$ distribution. For example, if the partial view $z_t$ is the intended action, the optimally learned model would be invariant of the used $\varepsilon$-exploration: $\pi(a_t|z_t)$. Analogously, the autoregressive models are invariant of the whole policy $\pi(a_t|s_t)$. This allows the autoregressive models to evaluate any other policy inside of the model. The causal partial model can run inside the simulation only policies conditioned on the starting state $s_0$, the actions $a_{<t}$ and the partial views $z_{\leq t}$. If we want to evaluate a policy conditioned on different features, we can collect trajectories from the policy and retrain the model. The model can always evaluate the policy used to produce the training data. We can also improve the policy, because the model allows to estimate the return for an initial $(s_0, a_0)$ pair, so the model can be used as a critic for a policy improvement.

*Table 2.* Models and their properties.

| | **Autoregressive** | **Deterministic** | **Causal Partial Model** |
|---|---|---|---|
| **Generates** | observation | nothing | partial view: $z_t$ |
| **Speed** | slow | fast | fast |
| **Causally correct** | always | in deterministic environments or with on-policy simulations | always |
| **Simulation variance** | high (distracted) | lowest | low |
| **Extra branching** | huge | 0 | controlled by $z_t$ size |
| **Invariant of** | $\pi(a_t|s_t)$ | - | $\pi(a_t|z_t)$ |
| **Evaluable policies** | any | $\psi(a_t|s_0, a_{<t})$ | $\psi(a_t|s_0, a_{<t}, z_{\leq t})$ |
| **Training** | once | iterative with policy | iterative with policy |

## F. MuZero properties on stochastic environments

MuZero performed surprisingly well on the stochastic environments, even when using a deterministic action-conditioned non-causal model. We will provide an explanation here. First, let us denote by $\hat{V}(s_0, a_0, a_1, \ldots, a_t)$ the output of the learned value network, given $s_0, a_0, a_1, \ldots, a_t$. The $\hat{V}(s_0, a_0)$ will correctly model the expected return, given $s_0, a_0$. But the next $\hat{V}(s_0, a_0, a_1)$ can lead to causally incorrect planning, because $s_1$ is a confounder here. When planning with a new policy $\psi(a_1|s_0, a_0)$, the $\sum_{a_1} \psi(a_1|s_0, a_0)\hat{V}(s_0, a_0, a_1)$ can be biased on stochastic environments. The planning with a causally correct model would instead compute the expected return by:

$$\sum_{z_1} p(z_1|s_0, a_0) \sum_{a_1} \psi(a_1|s_0, a_0, z_1)\hat{V}(s_0, a_0, z_1, a_1), \tag{29}$$

where $z_1$ is a backdoor to make $a_1$ independent of $s_1$, given $z_1$.

By analyzing the search trees on the AvoidFuzyBear MDP, we were able to find a number of MuZero properties that mitigate the negative effects of the non-causal model:

1. The correctly estimated value $\hat{V}(s_0, a_0)$ discourages opening a tree branch that is suboptimal in the real environment.

2. The learned prior policy $p(a_t|h_t)$ discourages opening the suboptimal branch, if the action leading to the suboptimal branch is not the most probable action. E.g., MuZero has less problems on AvoidFuzzyBear, if $p(teddy) < 0.5$.

3. The averaging of the value-network values from all nodes of the search tree assigns a small weight to the correct $\hat{V}(s_0, a_0)$ only after many simulations.

4. If using chance nodes inside a search tree, the bigger branching factor leads to a more shallow search. The shallow search bootstraps more from the causally correct $\hat{V}(s_0, a_0)$.

## G. Tree search experiments

We bootstrap the MCTS search from a learned value function and use pUCT to select the simulation actions, as done in MuZero (Schrittwieser et al., 2019). Additionally, we enhance the MCTS with chance nodes (Browne et al., 2012). The value of a chance node is the expected value of its children, based on the learned $p(z_t|h_t)$ probabilities. The chance nodes do not increase the number of network evaluations per a simulation, because if a child of the chance node is not expanded yet, we bootstrap from the child value.

When using expectimax or MCTS with chance nodes, we use the intended actions $z_t$ as the chance outcomes. During a simulation, we sample the chance outcomes from the learned $p(z_t|h_t)$ model. Conveniently, we are able to reuse the MuZero policy network as the $p(z_t|h_t)$ model, because the policy network is trained to model the intended actions. The deterministic non-causal partial model (NCPM) simply ignores the outcomes of the chance nodes. We do not recommend to use the non-causal partial model together with MCTS with chance nodes. We still see that the MCTS with chance nodes helps the non-causal model, because it reduces the search depth and the agent bootstraps more from the correct $\hat{V}(s_0, a_0)$.

The best performing hyper-parameters were found by trying multiples of 3 for the learning rate, the probability of exploration and the policy cost weight. The final used hyper-parameters are listed in Table 3. The pUCT hyper-parameters from MuZero (Schrittwieser et al., 2019) worked well. Each reported experiment was run with at least 8 independent random seeds. The shaded area in the plots indicates 95% confidence intervals.

*Table 3.* Hyper-parameters for the MDP and MiniPacman experiments.

| Hyper-parameter | Description | Value |
|:---:|:---|:---|
| $\mu$ | Learning rate | 0.0003 |
| $\beta_1$ | Adam $\beta_1$ | 0.9 |
| $\beta_2$ | Adam $\beta_2$ | 0.999 |
| `batch_size` | Mini-batch size | 512 |
| `max_depth` | Expectimax search depth | 3 |
| `num_simulations` | Number of MCTS simulations | 50 |
| `buffer_size` | Replay buffer size | 500000 |
| $\varepsilon$ | Probability of exploration | 0.01 |
| $c_{reward}$ | Reward cost weight | 1.0 |
| $c_{value}$ | Value cost weight | 1.0 |
| $c_{policy}$ | Policy cost weight | 300.0 |
| $L_o$ | Overshoot Length | 5 |
| $L_u$ | Length of n-step returns | 10 |
| $\gamma$ | Discount factor | 0.995 |

### G.1. Details of experiments on MDPs

In Figure 7(b), we see the results for the different models, when using expectimax. The models with clustered probabilities and clustered observations approximate modeling of the probabilities or observations. These models are described in Appendix G.4.

### G.2. Details of experiments on MiniPacman

We use the same stochastic version of MiniPacman as released by Guez et al. (2019). To make the task harder, we included 2 ghosts from the start of the game. And we increase the number of ghosts by one after solving two 2 levels or 4 levels or other multiplies of 2. The game options are in Table 4. To avoid having to stack frames, we draw the ghost movement direction to a pixel before the ghost. This makes the environment more Markovian.

The pretrained policy was trained by expectimax with search depth 1. This is similar to Q-learning, except that the network consists of a reward model and a value network. This forms a strong baseline agent. Understandably, when training a model

Table 4. The used options for the MiniPacman environments.

| Argument | Description | Value |
|---|---|---|
| frame_cap | The maximum number of steps in an episode. | 3000 |
| mode | The game mode. | 'regular' |
| npills | The number of power pills. | 2 |
| nghosts | The number of ghosts at the start of the game. | 2 |
| ghost_speed_init | The ghost probability of moving. | 0.5 |
| ghost_speed_increase | An increase to the ghost speed after a level. | 0 |

on data from the pretrained policy, expectimax performed worse than when training on on-policy data. The expectimax trained on the data from the pretrained policy was not able to test the proposed action in the real environment.

### G.3. Simple extensions

**Exact KL.** Usually, we know the distribution of the used partial view: $z_t \sim m(z_t|s_t)$. When training the $p(z_t|h_t)$ model, we then minimize the exact $\mathrm{KL}(m(Z_t|s_t) \| p(Z_t|h_t))$.

**Replay with resampling.** To improve data efficiency, we use a shuffling replay buffer and replay each trajectory 4-times. The trajectory does not have to store the used $z_t$. We store the used $m(z_t|s_t)$ distribution and resample the $z_t$ from a posterior when replaying the trajectory. The posterior is:

$$p(z_t|s_t, a_t) = \frac{m(z_t|s_t)\pi(a_t|z_t)}{\sum_{z_t'} m(z_t'|s_t)\pi(a_t|z_t')} \tag{30}$$

where $\pi(a_t|z_t)$ is the $\varepsilon$-exploration distribution.

### G.4. Models trained by clustering

When using a tree-search, we want to have a small branching factor at the chance nodes. A good $z_t$ variable would be discrete with a small number of categories. This is satisfied, if the $z_t$ is the intended action and the number of the possible actions is small. We do not have such compact discrete $z_t$, if using as $z_t$ the observation, the policy probabilities or some other modeled layer. Here, we will present a model that approximates such causal partial models. The idea is to cluster the modeled layers and use just the cluster index as $z_t$. The cluster index is discrete and we can control the branching factor by choosing the the number of clusters.

Concretely, let's call the modeled layer $x_t$. We will model the layer with a mixture of components. The mixture gives us a discrete latent variable $z_t$ to represent the component index. To train the mixture, we use a clustering loss to train only the best component to model the $x_t$, given $h_t$ and $z_t$:

$$L_{\text{clustering}} = \min_{z_t}(-\beta_{\text{clustering}} \log p(z_t|h_t) - \log p(x_t|h_t, z_t)) \tag{31}$$

where $p(z_t|h_t)$ is a model of the categorical component index and $\beta_{\text{clustering}} \in (0, 1)$ is a hyper-parameter to encourage moving the information bits to the latent $z_t$. During training, we use the index of the best component as the inferred $z_t$. In theory, a better inference can be obtained by smoothing.

In contrast to training by maximum likelihood, the clustering loss uses just the needed number of the mixture components. This helps to reduce the branching factor in a search.

In general, the cluster index is not guaranteed to be sufficient as a backdoor, if the reconstruction loss $-\log p(x_t|h_t, z_t)$ is not zero. For example, if $x_t$ is the next observation, the number of mixture components may need to be unrealistically large, if the observation can contain many distractors.

# H. Dyna experiments on AvoidFuzzyBear

In this experiment we demonstrate that we can do a policy improvement step based on an off-policy experience. The algorithm is described in detail in Appendix H.1. In short, we simulate experiences from the partial model, and use policy gradient to learn the optimal policy on these experiences as if they were real experiences (this is possible since the policy gradient only needs action probabilities, values, predicted rewards and ends of episodes). We compare a non-causal model and a causal model where the backdoor $z_t$ is the intended action. For the environment we use AvoidFuzzyBear (Figure 1(b)). We collect experiences that are sub-optimal: half the time the agent visits the forest and half the time it stays home, but once in the forest it acts optimally with probability $0.9$. This is meant to simulate situations either where the agent has not yet learned the optimal policy but is acting reasonably, or where it is acting with a different objective (such as exploration or intrinsic reward), but would like to derive an improved policy. We expect the non-causal model to choose the sub-optimal policy of visiting the forest, since the sequence of actions of visiting the forest and hugging typically yields high reward.

This is what we indeed find, as shown in Figure 7(a). We see that the non-causal model indeed achieves a sub-optimal reward (less than $0.6$), but believes that it will achieve a high reward (more than $0.6$). On the other hand, the causal model achieves the optimal reward and correctly predicts that it will achieve the corresponding value.



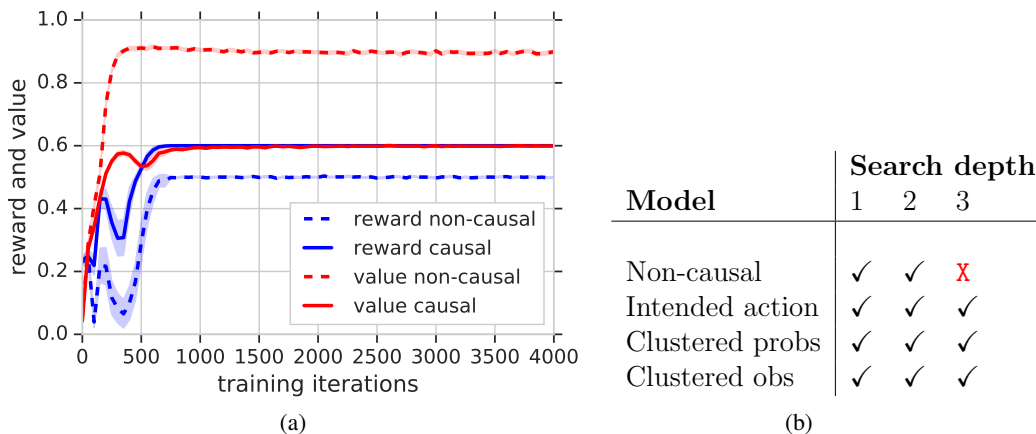| Model | Search depth | | |
|---|---|---|---|
| | 1 | 2 | 3 |
| Non-causal | ✓ | ✓ | X |
| Intended action | ✓ | ✓ | ✓ |
| Clustered probs | ✓ | ✓ | ✓ |
| Clustered obs | ✓ | ✓ | ✓ |

(a)  (b)

*Figure 7.* (a) Dyna on AvoidFuzzyBear. Models were trained on a sub-optimal behavior policy that explores both parts of the environment, and the evaluation policy was trained purely inside the model. The non-causal model (dotted lines) achieved sub-optimal reward, while expecting large reward. The causal model (solid lines) achieved optimal reward and had a correct expectation of the reward. The shaded area indicates 95% confidence intervals from 50 runs. (b) Models solving the AvoidFuzzyBear MDP with expectimax. The non-causal model misled the agent when using search depth 3 or higher.

## H.1. Dyna-style policy-gradient algorithm

In this section we derive an algorithm for learning an improved policy given a (non-optimal) experience that utilizes n-step returns from partial models presented in this paper. In general, a model of the environment can be used in a number of ways for reinforcement learning. In Dyna (Sutton, 1990), we sample experiences from the model, and apply a model-free algorithm (Q-learning in the original implementation, but more generally we could consider SARSA or policy gradient) as if these were real experiences. In Dyna-2 (Silver et al., 2008), the same process is applied but in the context the agent is currently in—starting the simulations from the current state—and adapting the policy locally (for example through separate fast weights). In MCTS, the model is used to build a tree of possibilities. Can we apply our model directly in these scenarios? While we don't have a full model of the environment, we can produce a causally correct simulation of rewards and values; one that should generalize to policies different from those the agent was trained on. Policy probabilities, values, rewards and ends of episodes are the only variables that the above RL algorithms need.

Here we propose a specific implementation of Dyna-style policy-gradient algorithm based on the models discussed in the paper. This is meant as a proof of principle, and more exploration is left for future work.

As the agent sees an observation $y_{t+1}$, it forms an internal agent state $s_t$ from this observation and the previous agent state: $s_{t+1} = \text{RNN}_s(s_t, a_t, y_{t+1})$. The agent state in our implementation is the state of the recurrent network, typically LSTM (Hochreiter & Schmidhuber, 1997). Next, let us assume that at some point in time with state $s_0$ the agent would

like to learn to do a simulation from the model. Let $h_t$ be the state of the simulation at time $t$. The agent first sets $h_1 = g(s_0, a_0)$ and proceeds with n-steps of the simulation recurrent network update $h_{t+1} = \text{RNN}(h_t, z_t, a_t)$. The agent learns the model $p(z_t|h_t)$ which it can use to simulate forward. We assume that the model was trained on some (non-optimal) policy/experience. We would like to derive an optimal policy and value function. Since these need to be used during acting (if the agent were to then act optimally in the real environment), they are functions of the agent state $s_t$: $\pi(a_t|s_t), V(s_t)$. Now in general, $h_t \neq s_t$ but we would like to use the simulation to train an optimal policy and value function. Thus we define a second pair of functions $\pi_h(a_t|h_t, z_t), V_h(h_t, z_t)$. Here the extra $z_t$'s are needed, since the $h_t$ has seen $z$'s only up to point $z_{t-1}$.

Next we are going to train these functions using policy gradients on simulated experiences. We start with some state $s_t$ and produce a simulation $h_{t+1}, \ldots, h_T$ by sampling $z_t$ from the model at each step and action $a_t \sim \pi_h(a_t|h_t, z_t)$. However at the initial point $t$, we sample from $\pi$, not $\pi_h$, and compute the value $V$, not $V_h$. The sequences of actions, values and policy parameters are the quantities needed to compute a policy gradient update. We use this update to train all these quantities.

There is one last element that the algorithm needs. The values and policy parameters are trained at the start state and along the simulation by n-step returns, computed from simulated rewards and the bootstrap value at the end of the simulation. However this last value is not trained in any way because it depends on the simulated state $V_h(h_T)$ not the agent state $s_T$. We would like this value to equal to what the agent state would produce: $V(s_T)$. Thus, during training of the model, we also train $V_h(h_T)$ to be close to $V(s_T)$ by imposing an $L_2$ penalty. In our implementation, we actually impose a penalty at every point $t$ during simulation but we haven't experimented with which choice is better.

## I. Value-iteration analysis on MDPs

### I.1. Optimal Value Derivations

We derive the following two model-based evaluation metrics for the MDP environments.

- $V^*_{\text{NCPM}(\pi)}(s_0)$: optimal value computed with the non-causal model, which is trained with training data from policy $\pi$, starting from state $s_0$.

- $V^*_{\text{CPM}(\pi)}(s_0)$: optimal value computed with the causal model, which is trained with training data from policy $\pi$, starting from state $s_0$.

The theoretical analysis of the MDP does not use empirically trained models from the policy data but rather assumes that the transition probabilities $p(s_{i+1} \mid s_i, a_i)$ of the MDP, and the policy, $\pi(a_i \mid s_i)$ or $\pi(z_i \mid s_i)$, from which training data are collected are accurately learned by the model.

**Computation of $V^*_{\text{NCPM}(\pi)}$ :**  For the non-causal model,

$$V^*_{\text{NCPM}(\pi)}(s_0) = \max_{a_0, \ldots, a_k} \sum_{i=0}^{k} \mathbb{E}_{s_i} \left[ r_{i+1}(s_i, a_i) \mid s_0, a_0, a_1, \ldots, a_i \right]$$

$$= \max_{a_0, \ldots, a_k} \sum_{i=0}^{k} \sum_{s_i} p(s_i \mid s_0, a_0, \ldots, a_i) r_{i+1}(s_i, a_i).$$

Notice that the probability of $s_i$ is affected by $a_i$ here, because the network gets $a_i$ as an input, when predicting the $r_{i+1}$. This will introduce the non-causal bias. The network implements the expectation implicitly by learning the mean of the reward seen in the training data. We can compute the expectation exactly, if we know the MDP. The $p(s_i \mid s_0, a_0, \ldots, a_i)$ can be computed recursively in two-steps as:

$$p(s_i \mid s_0, a_0, \ldots, a_i) = \frac{p(s_i \mid s_0, a_0, \ldots, a_{i-1})\pi(a_i \mid s_i)}{\sum_{s'_i} p(s'_i \mid s_0, a_0, \ldots, a_{i-1})\pi(a_i \mid s'_i)}. \tag{32}$$

Here, we see the dependency of the learned model on the policy $\pi$. The remaining terms can be expressed as:

$$p(s_i \mid s_0, a_0, \ldots, a_{i-1}) = \sum_{s_{i-1}} p(s_i, s_{i-1} \mid s_0, a_0, \ldots, a_{i-1}) \tag{33}$$

$$= \sum_{s_{i-1}} p(s_{i-1} \mid s_0, a_0, \ldots, a_{i-1}) p(s_i \mid s_{i-1}, a_{i-1}). \tag{34}$$

Denoting $p(s_i \mid s_0, a_0, \ldots, a_j)$ by $S_{i,j}$, we have the two-step recursion

$$S_{i,i} = \frac{S_{i,i-1}\, \pi(a_i \mid s_i)}{\sum_{s_i'} S_{i,i-1}'\, \pi(a_i \mid s_i')}, \tag{35}$$

$$S_{i,i-1} = \sum_{s_{i-1}} S_{i-1,i-1}\, p(s_i \mid s_{i-1}, a_{i-1}) \tag{36}$$

with $S_{1,0} = p(s_1 \mid s_0, a_0)$. We then compute $V_{\text{ncm}}^*(s_0)$ as $\max_{a_0,\ldots,a_k} \sum_{i=0}^{k} \sum_{s_i} S_{i,i} r_{i+1}(s_i, a_i)$.

**Computation of $V_{\text{CPM}(\pi)}^*$ :** For the causal model,

$$V_{\text{CPM}(\pi)}^*(s_0) = \max_{a_0} \sum_{z_1} p(z_1 \mid s_0, a_0) \max_{a_1} \sum_{z_2} p(z_2 \mid s_0, a_0, z_1, a_1) \cdots \tag{37}$$

$$\max_{a_{k-1}} \sum_{z_k} p(z_k \mid s_0, a_0, z_1, a_1, \ldots, z_{k-1}, a_{k-1}) \tag{38}$$

$$\max_{a_k} \sum_{i=0}^{k} \mathbb{E}[r_{i+1}(s_i, a_i) \mid s_0, a_0, z_1, a_1, \ldots, a_i)], \tag{39}$$

where for any $i \in [1, k]$,

$$p(z_i \mid s_0, a_0, z_1, a_1, \ldots, z_{i-1}, a_{i-1}) = \sum_{s_i} p(s_i, z_i \mid s_0, a_0, z_1, a_1, \ldots, z_{i-1}, a_{i-1}) \tag{40}$$

$$= \sum_{s_i} p(s_i \mid s_0, a_0, z_1 \ldots, z_{i-1}, a_{i-1})\, \pi(z_i \mid s_i). \tag{41}$$

$$\tag{42}$$

Denoting $p(s_i \mid s_0, a_0, z_1 \ldots, z_{i-1}, a_{i-1})$ by $Z_i$, we have

$$Z_i = \sum_{s_{i-1}} p(s_{i-1}, s_i \mid s_0, a_0, z_1 \ldots, z_{i-1}, a_{i-1}) \tag{43}$$

$$= \sum_{s_{i-1}} p(s_i \mid s_{i-1}, a_{i-1}) p(s_{i-1} \mid s_0, a_0, z_1 \ldots, z_{i-1}, a_{i-1}) \tag{44}$$

$$= \sum_{s_{i-1}} p(s_i \mid s_{i-1}, a_{i-1}) p(s_{i-1} \mid s_0, a_0, z_1 \ldots, z_{i-1}), \tag{45}$$

where we used the fact that $s_{i-1}$ is independent of $a_{i-1}$, given $z_{i-1}$. Furthermore,

$$p(s_{i-1} \mid s_0, a_0, z_1 \ldots, z_{i-1}) = \frac{p(s_{i-1}, z_{i-1} \mid s_0, a_0, z_1 \ldots, a_{i-2})}{p(z_{i-1} \mid s_0, a_0, z_1 \ldots, a_{i-2})} \tag{46}$$

$$= \frac{\pi(z_{i-1} \mid s_{i-1}) p(s_{i-1} \mid s_0, a_0, z_1 \ldots, z_{i-2}, a_{i-2})}{\sum_{s_{i-1}'} \pi(z_{i-1} \mid s_{i-1}') p(s_{i-1}' \mid s_0, a_0, z_1 \ldots, z_{i-2}, a_{i-2})} \tag{47}$$

$$= \frac{\pi(z_{i-1} \mid s_{i-1}) Z_{i-1}}{\sum_{s_{i-1}'} \pi(z_{i-1} \mid s_{i-1}') Z_{i-1}'}. \tag{48}$$

Therefore we can compute $Z_i$ recursively,

$$Z_i = \sum_{s_{i-1}} p(s_i \mid s_{i-1}, a_{i-1}) \frac{\pi(z_{i-1} \mid s_{i-1}) Z_{i-1}}{\sum_{s'_{i-1}} \pi(z_{i-1} \mid s'_{i-1}) Z'_{i-1}} \tag{49}$$

with $Z_1 = p(s_1 \mid s_0, a_0)$. The last term to compute in the definition of $V^*_{\text{CPM}(\pi)}(s_0)$ is

$$\sum_{i=0}^{k} \mathbb{E}[r_{i+1}(s_i, a_i) \mid s_0, a_0, z_1, a_1, \ldots, a_i)] = \sum_{i=0}^{k} \mathbb{E}[r_{i+1}(s_i, a_i) \mid s_0, a_0, z_1, a_1, \ldots, z_i)] \tag{50}$$

$$= \sum_{i=0}^{k} \sum_{s_i} p(s_i \mid s_0, a_0, z_1, a_1, \ldots, z_i) r_{i+1}(s_i, a_i) \tag{51}$$

$$= \sum_{i=0}^{k} \sum_{s_i} \frac{\pi(z_i \mid s_i) Z_i}{\sum_{s'_i} \pi(z_i \mid s'_i) Z'_i} r_{i+1}(s_i, a_i). \tag{52}$$

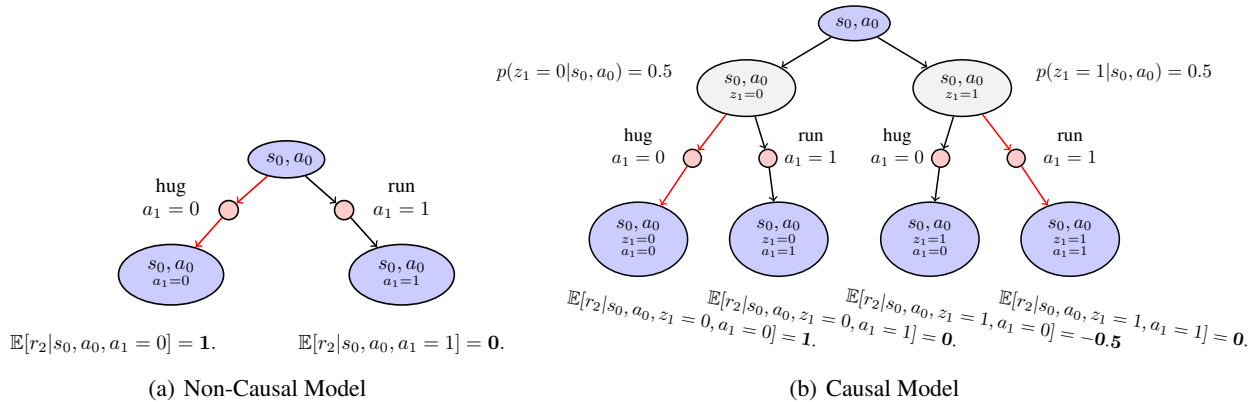## I.2. Planning with Non-Causal vs Causal Models on FuzzyBear



*Figure 8.* FuzzyBear decision trees evaluated with NCPM and CPM based on the optimal policy data with the intended action. Decision paths through red action nodes that give the maximum expected rewards are highlighted in red. The blue nodes are states and the gray ones are chance nodes.

In Figure 8(a), the non-causal agent always chooses **hug** at step $t = 1$, since it has learned from the optimal policy that a reward of $+1$ always follows after taking $a_1 = $ hug. Thus from the non-causal agent's point of view, the expected reward is always 1 after hugging. This is wrong since only hugging a teddy bear gives reward 1. Moreover it exceeds the maximum expected reward 0.5 of the FuzzyBear MDP. In Figure 8(b), the causal agent first samples the intention $z_1$ from the optimal policy, giving equal probability of landing in either of the two chance nodes. Then it chooses **hug** if $z_1 = 0$, indicating a teddy bear since the optimal policy intends to hug only if it observes a teddy bear. Likewise, it chooses **run** if $z_1 = 1$, indicating a grizzly bear. While the non-causal model expects unrealistically high reward, the causal model never over-estimates the expected reward.

## I.3. Learning with optimal policy and varying $\varepsilon$-exploration

The optimal policy of the FuzzyBear MDP (Figure 1(a)) is to always hug the teddy bear and run away from the grizzly bear. Using training data from this behavior policy, we show in Figure 8 the difference in the optimal planning based on the NCPM (Figure 3(d)) and CPM with the partial view $z_t$ being the *intended action* (Figure 3(e)). Learning from optimal policies with $\varepsilon$-exploration, the converged causal model is independent of the exploration parameter $\varepsilon$.

We analyze learning from optimal policy with varying amounts of $\varepsilon$-exploration for models on FuzzyBear (Figure 9(a)) and AvoidFuzzyBear (Figure 9(b)). As the parameter $\varepsilon$-exploration varies in range $(0, 1]$, the causal model has a constant
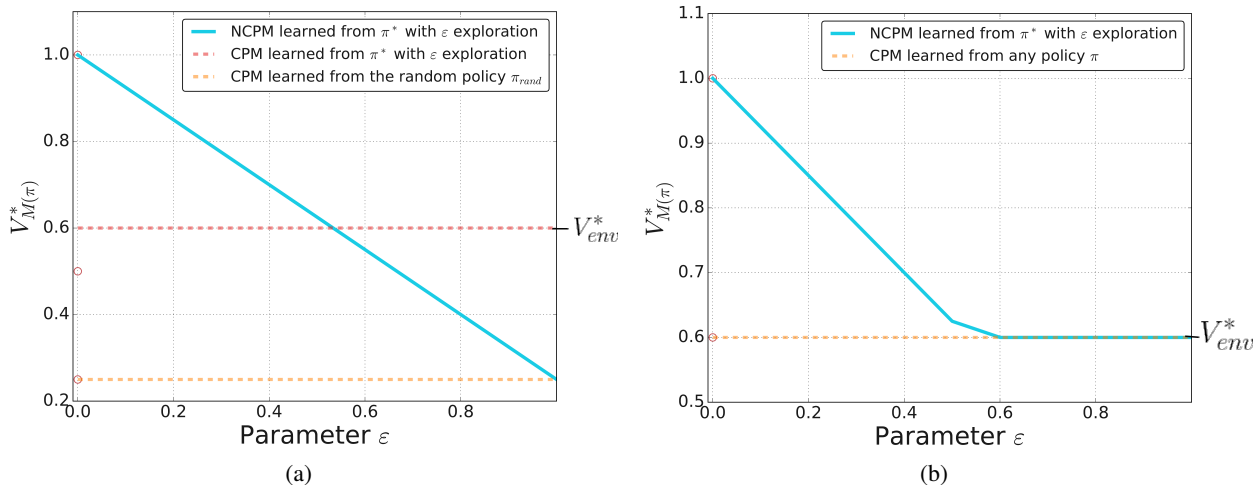
*Figure 9.* (a) In FuzzyBear, we use optimal policy with $\varepsilon$-exploration to generate training data for the Non-Causal Partial Model (NCPM) and Causal Partial Model (CPM). We vary the exploration parameter $\varepsilon \in (0, 1]$ and observe differences in found optimal values $V^*_{M(\pi)}$ under the model $M(\pi)$, where $M(\pi)$ denotes either CPM or NCPM trained on behavior policy $\pi$. The NCPM evaluation $V^*_{\text{NCPM}(\pi^*)}$ gives an unrealistically high value $1.0$ learned from the deterministic optimal policy ($\varepsilon = 0$). Expectantly, it decreases to the level of CPM optimal value $V^*_{\text{CPM}(\pi_{\text{rand}})}$ learned from the uniformly random policy as $\varepsilon \to 1$. The CPM optimal values $V^*_{\text{CPM}(\pi^*)}$ are constant for any value of $\varepsilon$ based on the theoretical analysis in Section I.1. (b) shows the same plots as (a) for the AvoidFuzzyBear environment. Learning from any policy $\pi$, the CPM optimal value always equals the maximum expected reward $0.6$, by correctly choosing to stay home.

evaluation since the intended action is not affected by the randomness in exploration. The non-causal model, on the other hand, evaluates based on the deterministic optimal policy data (i.e. at $\varepsilon = 0$) at an unrealistically high value of $1.0$ when the maximum expected reward is $0.5$. As $\varepsilon \to 1$, the training data becomes more random and its optimal evaluation expectantly goes down to match the causal evaluation based on a uniformly random policy. The causal evaluation based on the optimal policy $V^*_{\text{CM}(\pi^*)}$ converges to the ground truth environmental evaluation $V^*_{env}$ as $\varepsilon \to 0$.

### I.4. Sub-optimal behavior policies on AvoidFuzzyBear

On AvoidFuzzyBear (Figure 10), the optimal policy is to stay at home. Learning from data generated by random policies, the causal model indeed always prefers to stay home with any sampled intentions, resulting in a constant evaluation for all policies. On the other hand, the non-causal model gives varied, overly-optimistic evaluations, while choosing the wrong action (visit forest).

## J. Details for 3D experiments

### J.1. Conditional Dirichlet Mixture

When the backdoor variable $z_t$ was chosen to be the action probabilities, the distribution $p(z_t|h_t)$ was chosen as a mixture-network with $N_c$ Dirichlet components. The concentration parameters $\alpha_k(h_t)$ of each component were parametrized as $\alpha_k(h_t) = \alpha \, \text{softmax}(f_k(h_t))$, where $f_k$ is the output of a relu-MLP with layer sizes $[256, 64, N_c \times N_a]$, $\alpha$ is a total concentration parameter and $N_a$ is the number of actions.

### J.2. Hyper Parameters and Training

The hyper-parameter value ranges used in our 3D experiments are similar to (Gregor et al., 2019) and are shown in Table 5.

To speed up training, we interleaved training on the T-maze level with a simple "Food" level, in which the agent simply had to walk around and eat food blocks (described by Gregor et al. (2019)).

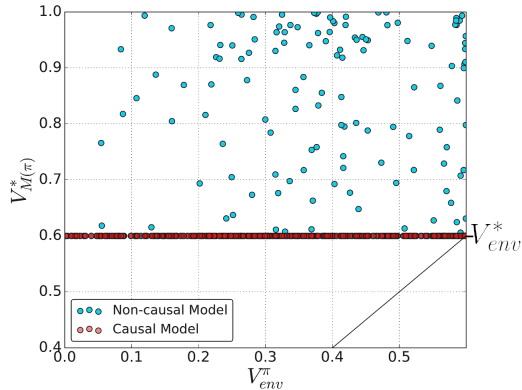In the bar plots, error bars represent standard error over 5 runs, each with 30 episodes.

*Figure 10.* MDP Analysis in the AvoidFuzzyBear environment. We randomly generate 500 policies and scatter plot them with x-axis showing the quality of the behavior policy $V^{\pi}_{env}$ and y-axis showing corresponding model optimal evaluations $V^{*}_{M(\pi)}$.

*Table 5.* Hyper-parameters used. Each reported experiment was repeated at least 5 times with different random seeds.

| Hyper-parameter | Description | Value |
|---|---|---|
| $\mu_{policy}$ | Policy learning rate | 0.0001 |
| $\mu_{model}$ | Model learning rate | 0.0005 |
| $c$ | Policy entropy regularization | 0.0004 |
| $\beta_1$ | Adam $\beta_1$ | 0.9 |
| $\beta_2$ | Adam $\beta_2$ | 0.999 |
| $L_o$ | Overshoot Length | 8 |
| $L_u$ | Unroll Length | 100 |
| $N_t$ | Number of points used to evaluate the generative loss per trajectory | 6 |
| $N_g$ | Number of points used to evaluate the generative loss per overshoot | 2 |
| $N_s$ | Number of ConvDRAW Steps | 8 |
| $N_h$ | Number of units in LSTM | 512 |
| $\alpha$ | Total concentration of Dirichlet distributions | 100 |
| $N_c$ | Number of components of Dirichlet mixture | 10 |

## J.3. Analysis of rollouts

For each episode, 5 rollouts are generated after having observed the first 3 frames from the environment. For the 5 rollouts, we processed the first 25 frames to classify the presence of food blocks by performing color matching of RGB values, using K-means and assuming 7 clusters. Rollouts were generated shortly after the policy had achieved ceiling performance (15–20 million frames seen), but before the entropy of the policy reduces to the point that there is no longer sufficient exploration. See Figure 11 for these same results for later training.
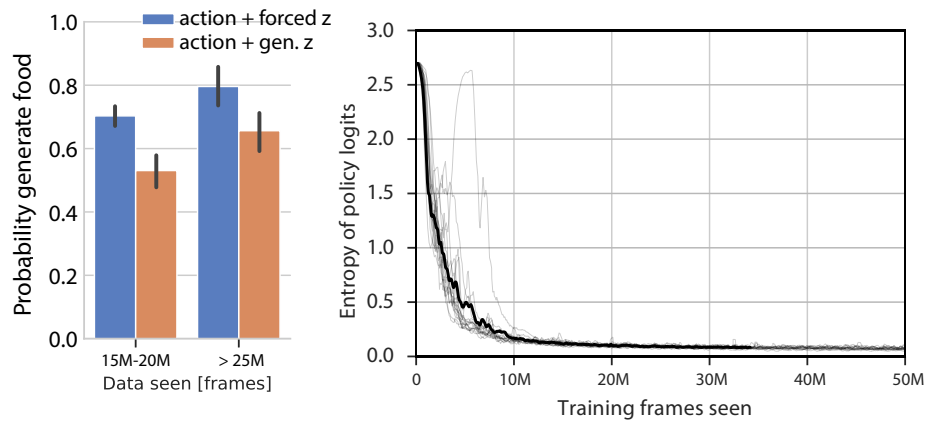
*Figure 11.* While earlier in training, CPM generates a diverse range of outcomes (food or no food), as the policy becomes more deterministic (as seen in the right plot of the policy entropy over training), CPM starts to generate more food and becomes overoptimistic, similar to NCPM. This can be avoided by training the model with non-zero $\varepsilon$-exploration.

## K. Model Rollouts

*Table 6.* Different types of model rollouts considered. Blue cells indicate variables that require interaction with the real environment, depending on the agent's state $s_t$. Orange cells indicate variables that can be computed from the model's state $h_t$.

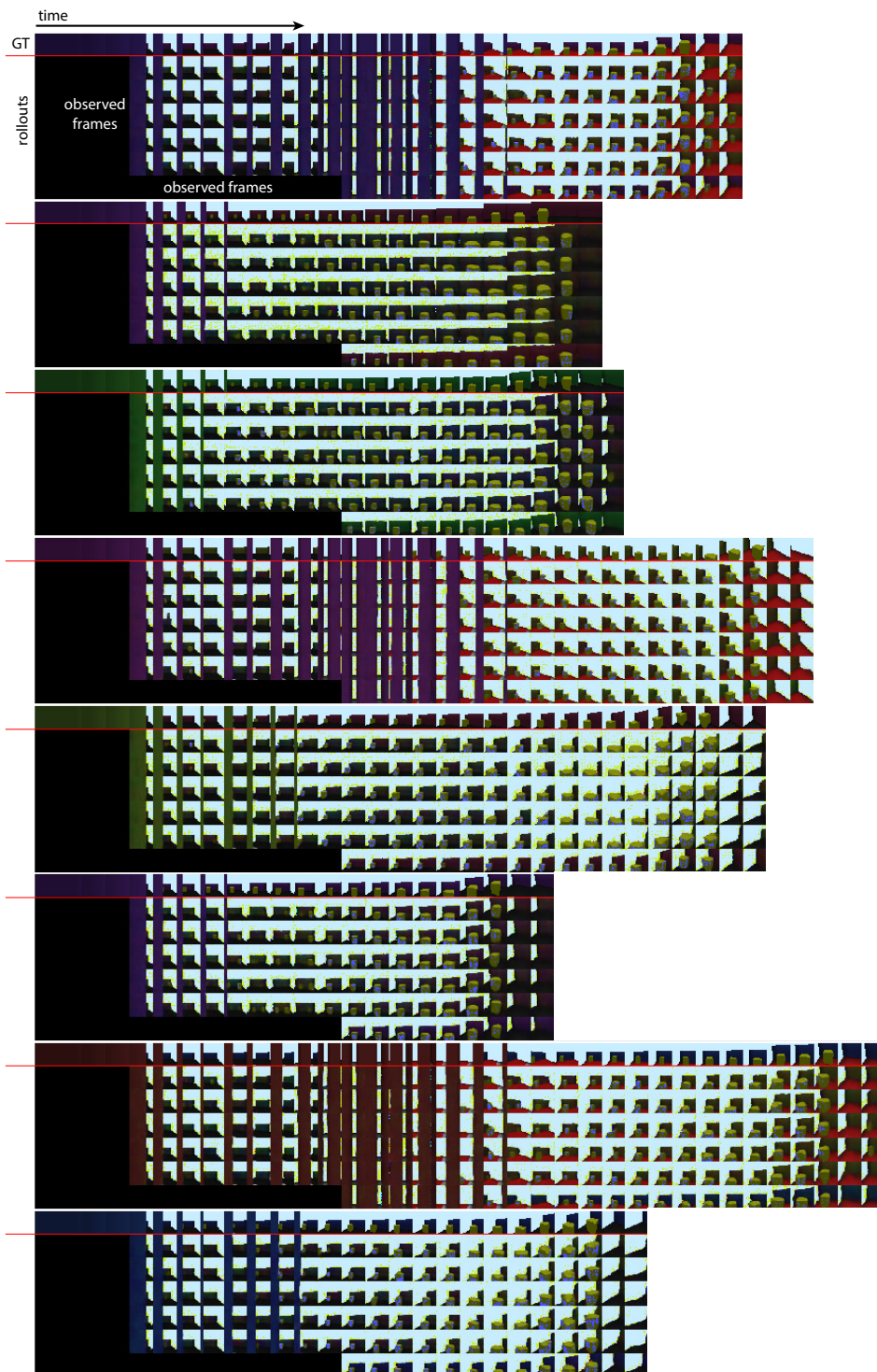| | Rollout | Rollout with forced actions | Rollout with forced actions and backdoor |
|---|---|---|---|
| **Backdoor** **Action** | $z_t \sim p(z_t\|h_t)$ $a_t \sim \pi(a_t\|z_t)$ | $z_t \sim p(z_t\|h_t)$ $a_t \sim \pi(a_t\|\hat{z}_t); \hat{z}_t \sim m(\hat{z}_t\|s_t)$ | $z_t \sim m(z_t\|s_t)$ $a_t \sim \pi(a_t\|z_t)$ |
| **State** | $h_t = \mathrm{RNN}_h(h_{t-1}, a_{t-1}, z_{t-1})$ | | |
| **Prediction** | $y_t \sim p(y_t\|h_t)$ | | |

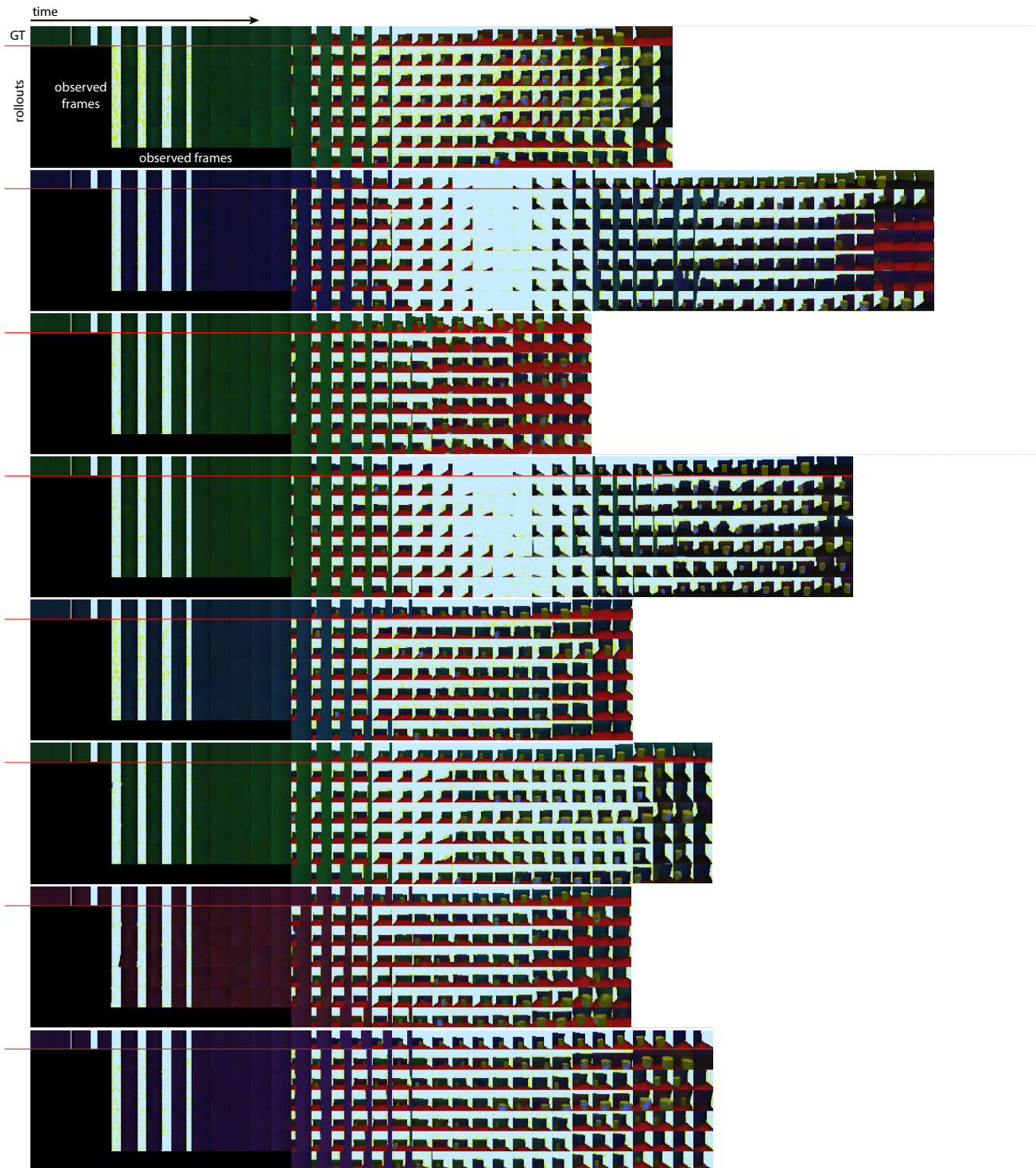*Figure 12.* Full rollouts for NCPM, conditioned on forced actions.

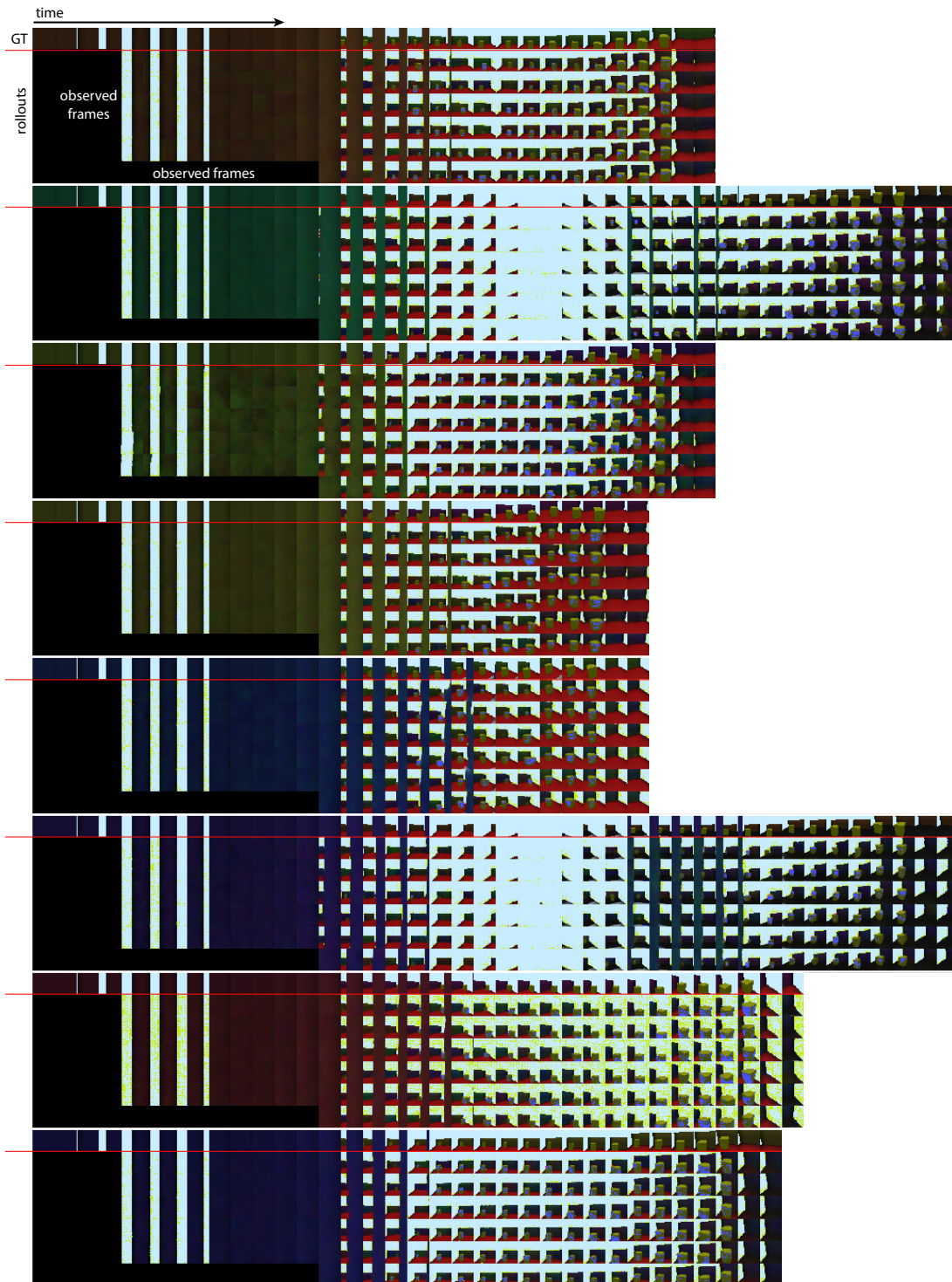*Figure 13.* Full rollouts for CPM, conditioned on forced actions and generated $z$.

*Figure 14.* Full rollouts for CPM, conditioned on forced actions and forced $z$.