# Model-Based Reinforcement Learning Under Confounding

Nishanth Venkatesh[1], *Student Member, IEEE*, and Andreas A. Malikopoulos[1,2], *Senior Member, IEEE*

*Abstract*—We investigate model-based reinforcement learning in contextual Markov decision processes (C-MDPs) in which the context is unobserved and induces confounding in the offline dataset. In such settings, conventional model-learning methods are fundamentally inconsistent, as the transition and reward mechanisms generated under a behavioral policy do not correspond to the interventional quantities required for evaluating a state-based policy. To address this issue, we adapt a proximal off-policy evaluation approach that identifies the confounded reward expectation using only observable state-action-reward trajectories under mild invertibility conditions on proxy variables. When combined with a behavior-averaged transition model, this construction yields a surrogate MDP whose Bellman operator is well defined and consistent for state-based policies, and which integrates seamlessly with the maximum causal entropy (MaxCausalEnt) model-learning framework. The proposed formulation enables principled model learning and planning in confounded environments where contextual information is unobserved, unavailable, or impractical to collect.

*Index Terms*—Offline reinforcement learning, causality, and confounding.

## I. INTRODUCTION

Model-based reinforcement learning (MBRL) has long been recognized as a data-efficient approach to sequential decision-making, as it explicitly estimates the system dynamics and reward mechanism to enable planning through simulated rollouts [1]. This paradigm typically achieves greater sample efficiency than model-free methods. Recent developments include the use of Gaussian process models for dynamics learning [2] and neural dynamics models with model-free fine-tuning [3], both of which demonstrate improved data efficiency in policy optimization. Additional advances based on probabilistic ensembles with trajectory sampling [4] and latent-dynamics planning [5] further illustrate that learned models can support competitive performance on complex control tasks. A common requirement across these methodologies is the ability to learn a model that accurately captures the underlying transition and reward structures, as model fidelity is essential for reliable planning.

Most MBRL methods implicitly assume that the data-generating process is fully observed and free of hidden factors that jointly influence decisions and outcomes. In many real-world domains, this assumption is violated. Examples include human–robot interaction, routing and maneuvering connected and automated vehicles in mixed traffic [6], medical records, and personalized decision systems [7], [8]. In such settings, the offline dataset available for learning may omit latent variables that affect both the environment and the data-generating policy. When actions depend on these hidden variables, the collected data encode arbitrary correlations among state, action, next state, and reward. In causal inference, this phenomenon is referred to as *confounding*. Consequently, transition and reward models learned directly from such data could misrepresent the effect of candidate policies and introduce systematic bias into planning.

For example, consider historical medical records in which each entry contains a doctor's chosen treatment and the patient's eventual outcome. Although the dataset includes observable patient characteristics, the doctor's decisions were based on additional information—such as severity of symptoms, early warning signs, or clinical intuition—that was not recorded. This latent factor affects both the treatment selected and the patient's likelihood of recovery, producing correlations that reflect the doctor's unobserved reasoning rather than the causal effect of the treatment itself. Learning treatment effects directly from such observational data leads to biased conclusions, since the behavioral policy exploited information unavailable to the learner. This phenomenon mirrors the confounding encountered in contextual decision processes, where unobserved context simultaneously influences the behavioral policy and the system response.

Reinforcement learning under unobserved confounding has been studied primarily from a model-free perspective. Identification and proximal causal inference frameworks [9]–[11] provide tools to recover causal effects using proxy variables when hidden confounders are present. These include robust off-policy evaluation (OPE) under bounded confounding [12], non-identifiability and hardness results for confounded environments [13], and proximal reinforcement learning methods that exploit proxy variables [14]. Recent developments provide tensor-based identification policies for offline data [15], [16] and proximal off-policy evaluation in partially observed Markov decision processes [17]. These contributions provide powerful tools for off-policy evaluation and policy optimization. However, they largely operate in a model-free setting and do not directly address learning a surrogate transition-reward model that can be embedded into standard MBRL pipelines.

Contextual Markov decision processes (C-MDPs) provide a natural setting for confounded MBRL, as an unobserved context variable can influence both the transition law and the reward function. When the behavioral policy that generated

the dataset depends on this hidden context, the observed transitions and rewards reflect correlations induced by the unobserved variable rather than the causal effect of the actions. As a result, models learned directly from such behavioral data fail to represent the system from the perspective of a state-based evaluation policy, and the Bellman equations associated with that policy are no longer consistent with the available observations. A central challenge, therefore, is to construct a surrogate MDP whose transition dynamics and reward expectations are identifiable from the restricted dataset and faithfully capture the environment as experienced by a state-based agent.

In this paper, we develop an MBRL framework for contextual MDPs in which the context is unobserved in the available dataset and acts as a confounder. Our approach combines a behavior-averaged transition model with a proximal reconstruction of the reward component of the Bellman equation. We adapt a proximal OPE construction technique [17] to the C-MDP setting to express the confounded reward expectation as a functional of observable proxy variables under standard invertibility assumptions from proximal causal inference. This yields a deconfounded Bellman operator that is compatible with state-based policies and can be evaluated from offline data. We adapt a well-known data-based model-learning framework- maximum causal entropy (MaxCausalEnt), introduced in [18]. Inverse reinforcement learning and belief-dynamics extensions to the MaxCausalEnt framework in [19], [20] demonstrate how such formulations can capture purposeful behavior in stochastic environments. We integrate the proximal reward correction with the MaxCausalEnt model-learning formulation. In our setting, the MaxCausalEnt framework is used to jointly learn a Q-function and a behavior-averaged transition model that are Bellman-consistent with respect to the deconfounded reward expectation. The resulting surrogate MDP supports model-based planning in confounded environments where contextual information is unavailable or impractical to record.

The remainder of the paper is organized as follows. In Section II-A, we formalize the learning problem that arises when data are generated by an unknown, context-dependent behavioral policy, and we introduce the resulting surrogate model-learning formulation. In Section III, we develop our solution approach: we reinterpret the C–MDP through a causal inference perspective, reformulate the system as a POMDP to reveal the role of the hidden context, and derive a proximal off-policy evaluation method that identifies the confounded reward expectation using observable proxy variables. In Section IV, we present a numerical example that illustrates the necessity and impact of the proximal correction. Finally, in Section V, we draw concluding remarks and outline several directions for future research.

## II. Problem Formulation

We consider a system that can be modeled as a C-MDP given by the tuple $\mathcal{M} = (\mathcal{X}, \mathcal{Z}, \mathcal{U}, P, r, \eta, \nu, T)$. The system evolves over discrete time steps $t = 0, \ldots, T$, where $T \in \mathbb{N}$ is the finite horizon. The state space $\mathcal{X}$ and action space $\mathcal{U}$ are

finite. At each time $t$, the random variables $X_t \in \mathcal{X}$ and $U_t \in \mathcal{U}$ denote the system state and control action, respectively, and the initial state satisfies $X_0 \sim \eta$. The system is influenced by an unobserved context $Z_t \in \mathcal{Z}$, with $Z_0 \sim \nu$. The joint process $(X_t, Z_t)$ is Markovian with respect to $(X_t, Z_t, U_t)$. The state transition is governed by

$$P_X : \mathcal{X} \times \mathcal{U} \times \mathcal{Z} \to \Delta(\mathcal{X}), \qquad X_{t+1} \sim P_X(\cdot \mid x_t, z_t, u_t),$$

and the context evolves according to

$$P_Z : \mathcal{Z} \times \mathcal{X} \times \mathcal{U} \to \Delta(\mathcal{Z}), \qquad Z_{t+1} \sim P_Z(\cdot \mid x_t, z_t, u_t).$$

The context $Z_t$ is unobserved by the planner, and it may influence both dynamics and rewards. Furthermore, the reward function is given by the mapping $r : \mathcal{X} \times \mathcal{U} \times \mathcal{Z} \to \mathbb{R}$. At each $t = 0, \ldots, T-1$, we let $\tau_t$ be the realization of a trajectory of the system up to time $t$, which we define as

$$\tau_t = (x_0, z_0, u_0, r_0, \ldots, x_t, z_t, u_t, r_t), \qquad (1)$$

where $\tau = (x_0, z_0, u_0, r_0, \ldots, x_T, z_T, u_T, r_T)$ denotes the entire trajectory up to the time of horizon. Let $\mathcal{T}$ be the space of all feasible full trajectories $\tau$ of the system. Having specified the C-MDP model, we now formalize the decision-making problem of interest within this framework.

### A. State-based Control policy

A planner has access to a dataset generated by an expert who interacts with the C-MDP described above. We consider that the expert follows a behavioral policy $\boldsymbol{g^b} = \{g_t^{\boldsymbol{b}}\}_{t=0}^T$, where $g_t^{\boldsymbol{b}}$ is the expert's control law at time $t$. Each control law $g_t^{\boldsymbol{b}}$ may depend on the unobserved context $Z_t$ and can therefore be of the form $g_t^{\boldsymbol{b}} : \mathcal{X} \times \mathcal{Z} \to \Delta(\mathcal{U})$. Since contextual information is not recorded, the planner observes only $(x_t, u_t, r_t)$ tuples. Let $\tau^o = (x_0, u_0, r_0, \ldots, x_T, u_T, r_T)$ denote the observable component of a trajectory $\tau$, and let $\mathcal{D}^{\boldsymbol{b}} = \{\tau^{o,i}\}_{i=1}^N$ be the dataset.

Since the expert's behavior may depend on the latent context, two key components of the C-MDP become unidentifiable from the observable data: (i) the context-dependent transition mapping $p_x(x_{t+1} \mid x_t, z_t, u_t)$, and (ii) the expert's full state-context-based control law $g_t^{\boldsymbol{b}}(u_t \mid x_t, z_t)$ at each $t$. However, although the context is unobserved, the behavioral policy induces a well-defined *state-action next-state* distribution, which is identifiable from the dataset. Thus, based on the triplets $(x_t, u_t, x_{t+1})$ from the dataset, the context-marginalized transition model

$$p^\phi(x_{t+1} \mid x_t, u_t) = \sum_{z_t \in \mathcal{Z}} p(x_{t+1} \mid x_t, z_t, u_t)\, p(z_t \mid x_t, u_t).$$

$$(2)$$

is identifiable and parameterized by a stochastic matrix $\phi$ of appropriate dimension, which serves as a surrogate that approximates the context-marginalized dynamics. Each element of the matrix $\phi$ is denoted by $p^\phi(x_{t+1} \mid x_t, u_t)$. As the context is unobserved, the planner cannot recover the true context-dependent model $p_X(x_{t+1} \mid x_t, z_t, u_t)$. Instead, the learned model represents the *context-marginalized* evolution of the state, induced by the expert.

In parallel, the behavioral policy induces a distribution over actions conditioned on the observable state. The planner, therefore, seeks to learn a state-based control policy that best explains the observed actions under a suitable modeling framework. In particular, the planner aims to learn (i) a transition model $p^\phi(x_{t+1} \mid x_t, u_t)$ from observable transitions, and (ii) a state-based policy $\boldsymbol{g}^{\boldsymbol{\theta}} = \{g_t^{\boldsymbol{\theta}}\}_{t=0}^{T-1}$ indexed by $\theta$, where $\theta$ identifies a particular member of the policy class under consideration. For each $t$, the mapping $g_t^\theta : \mathcal{X} \to \Delta(\mathcal{U})$ specifies a distribution over actions conditioned only on the observable state.

### B. State-Based Policy and Value/Q Parameterization

We restrict attention to state-based policies $\boldsymbol{g}^{\boldsymbol{\theta}} = \{g_t^{\boldsymbol{\theta}}\}_{t=0}^{T}$, where the control law at each $t$ is a mapping $g_t^{\boldsymbol{\theta}} : \mathcal{X} \to \Delta(\mathcal{U})$. For any state $x_t$ and action $u_t$, we denote the action probability by $p^{\boldsymbol{\theta}}(u_t \mid x_t)$. Given a surrogate transition model $p^\phi$, the value of policy $\boldsymbol{g}^{\boldsymbol{\theta}}$ at time $t$ is

$$V_t^{\boldsymbol{\theta},\phi}(x_t) = \mathbb{E}^{\boldsymbol{\theta},\phi}\left[\sum_{k=t}^{T-1} r(X_k, Z_k, U_k) \,\Big|\, X_t = x_t\right]. \quad (3)$$

The corresponding Q-function is defined as

$$Q_t^{\boldsymbol{\theta},\phi}(x_t, u_t) = \mathbb{E}^{\boldsymbol{\theta}}\left[r(X_t, Z_t, U_t) \,\big|\, X_t = x_t, U_t = u_t\right]$$
$$+ \mathbb{E}^{\phi}\left[V_{t+1}^{\boldsymbol{\theta},\phi}(X_{t+1}) \,\big|\, X_t = x_t, U_t = u_t\right], \quad (4)$$

where $\mathbb{E}^{\boldsymbol{\theta}}[\cdot]$ denotes expectation with respect to the distribution of the latent context $Z_t$ induced by the policy $\boldsymbol{g}^{\boldsymbol{\theta}}$, and $\mathbb{E}^{\phi}[\cdot]$ denotes expectation with respect to the stochastic evolution of $X_{t+1}$ under the surrogate transition model $\phi$. The value and Q-functions satisfy

$$V_t^{\boldsymbol{\theta},\phi}(x_t) = \sum_{u \in \mathcal{U}} p^{\boldsymbol{\theta}}(u \mid x_t)\, Q_t^{\boldsymbol{\theta},\phi}(x_t, u). \quad (5)$$

To jointly learn the surrogate transition model and the state-based policy, we adopt the maximum causal entropy framework [18]. This framework provides a data-driven methodology to infer both the dynamics and the behavioral policy from a dataset of trajectories. Among all models that explain the observed actions, this framework selects one that (i) matches the expert's action distribution as closely as possible, while (ii) avoiding additional structure not justified by the data. In a finite-horizon setting, this principle leads to a softmax (Boltzmann) representation of the policy in terms of the Q-function under the surrogate model.

Hence, for a fixed surrogate transition model $\phi$, the state-based control law at time $t$ is parameterized as

$$p^{\boldsymbol{\theta}}(u_t \mid x_t) = \frac{\exp\big(Q_t^{\boldsymbol{\theta},\phi}(x_t, u_t)\big)}{\sum_{u' \in \mathcal{U}} \exp\big(Q_t^{\boldsymbol{\theta},\phi}(x_t, u')\big)}. \quad (6)$$

For each fixed $\phi$, the Q-function is determined by the parameters $(\boldsymbol{\theta}, \phi)$, and the softmax relation (6) induces a unique state-based control law $g_t^{\boldsymbol{\theta}}$. The policy is indexed by $\boldsymbol{\theta}$ to emphasize that $\phi$ is treated as part of the learned surrogate environment, while $\boldsymbol{\theta}$ parameterizes the choice of policy. Thus, learning $(\boldsymbol{\theta}, \phi)$ amounts to learning a Q-function that is Bellman-consistent with respect to $p^\phi$ and then using (6) to obtain the corresponding state-based policy.

### C. MaxCausalEnt Model Learning

The MaxCausalEnt principle fits $(\boldsymbol{\theta}, \phi)$ by maximizing the likelihood of expert actions under the softmax policy while enforcing Q-based Bellman consistency.

**Problem 1** (MaxCausalEnt Model Learning)**.**

$$\max_{\boldsymbol{\theta}, \phi} \quad \mathbb{E}_{\mathcal{D}^b}\left[-\log\left(\prod_{t=0}^{T-1} p^{\boldsymbol{\theta}}(u_t \mid x_t)\right)\right] \quad (7)$$

$$\text{subject to} \quad Q_t^{\boldsymbol{\theta},\phi}(x_t, u_t) = \mathbb{E}^{\boldsymbol{\theta}}[r(X_t, Z_t, U_t) \mid x_t, u_t]$$
$$+ \mathbb{E}^{\phi}\left[V_{t+1}^{\boldsymbol{\theta},\phi}(X_{t+1}) \mid x_t, u_t\right], \quad (8)$$

$$V_t^{\boldsymbol{\theta}}(x_t) = \sum_{u_t \in \mathcal{U}} p^{\boldsymbol{\theta}}(u_t \mid x_t)\, Q_t^{\boldsymbol{\theta},\phi}(x_t, u_t),$$
$$\forall x_t \in \mathcal{X}, \; \forall u_t \in \mathcal{U},$$

where (8) expresses the Bellman consistency of the parametrized policy $\boldsymbol{g}^{\boldsymbol{\theta}}$ under the surrogate model $\phi$.

**Remark 1.** If contextual information were available, the reward and transition expectations in (8) could be computed by marginalizing over $(x_t, z_t, u_t, x_{t+1})$ tuples.

## III. SOLUTION APPROACH

In this section, we first reinterpret the C-MDP from a causal inference perspective. This highlights how the latent context acts as an unobserved confounder. We cast the system as a structural causal model (SCM) to clarify how the behavioral policy and any parameterized policy induce different causal graphs and, consequently, different trajectory distributions. This motivates the use of *proximal learning*—a causal inference tool designed to recover causal effects in the presence of hidden confounding—as a means to compute the reward term in the Bellman equation using only observable variables. We then introduce the notation needed for the analysis and reinterpret the C-MDP as a POMDP to make the role of the hidden context explicit. Later, we develop an off-policy evaluation-based modification of the MaxCausalEnt framework that uses proximal identification to obtain a deconfounded Bellman expectation specific to our problem formulation.

### A. Causal Inference Perspective

Causal inference provides tools for reasoning about how interventions change the behavior of a system. This allows analysis of the system beyond what is revealed by observational or correlational structure alone. This distinction is especially important in offline settings—such as the C-MDP that we consider. It is essential to note that the offline data collected under a behavioral policy may rely on additional information unavailable to the planner, as given in Problem 1. In such cases, functional relationships between variables fail to reveal how the system would evolve under a different (intervened) policy.

A SCM represents each variable in the system as the output of a structural equation, with a directed acyclic graph
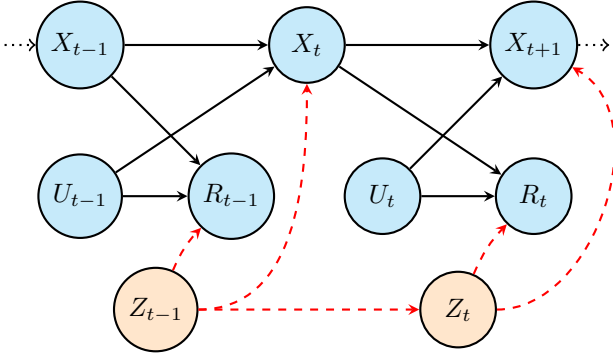
Fig. 1: Causal graph of the C-MDP.

(DAG) specifying the underlying causal dependencies [21]. For the C–MDP introduced in Section II-A, these structural equations coincide with the stochastic dynamics governing the state, context, and reward processes. Figure 1 illustrates a representative case in which the context is *persistent*, i.e., $Z_{t+1} = Z_t$. The DAG shows how past states and actions influence future evolution while highlighting that the latent context $Z_t$ is a causal parent of both the next state $X_{t+1}$ and the reward $R_t$. Consequently, when $Z_t$ is unobserved, the behavioral data contain *confounded* samples of the transition and reward pairs, as their observed variation reflects both the effects of the action and the hidden context.

Under the behavioral policy $g^b$, the action $U_t$ may depend on both $X_t$ and $Z_t$. This alters the system's SCM to introduce at each $t$, causal arrows from the context $Z_t$ and state $X_t$ to the action $U_t$ into the DAG. By contrast, under a state-based parameterized policy $g^\theta$, the action depends only on $X_t$ and the arrow $Z_t \to U_t$ is absent. Thus, the behavioral and parameterized policies induce *different* SCMs, therefore different distributions over trajectories. These issues motivate the need for an OPE procedure capable of expressing the reward expectations. In the next subsection, we reinterpret the C-MDP as a partially observed Markov decision process (POMDP), which reveals precisely how the trajectory distribution depends on hidden variables and enables the proximal OPE construction that follows.

### B. POMDP Representation

We recast the C-MDP as a POMDP whose full state at time $t$ is defined by $S_t = (X_t, Z_t)$, where $X_t$ is the observed component and $Z_t$ is the unobserved context. The set of full states is defined as $\mathcal{S} := \mathcal{X} \times \mathcal{Z}$. The agent observes only $Y_t = X_t$, with the set of observations given by $\mathcal{Y} := \mathcal{X}$. At each $t$, the agent receives a reward

$$R_t = r(X_t, Z_t, U_t) = r(S_t, U_t). \tag{9}$$

A complete trajectory and the corresponding observable component up to time $t$ is therefore written as

$$\tau_t = (s_{0:t}, y_{0:t}, u_{0:t}), \qquad \tau_t^o = (y_{0:t}, u_{0:t}) \tag{10}$$

Any parameterized policy $g^\theta$ induces a probability measure over the space of such trajectories. This measure determines the expected reward term in the Bellman consistency constraint

(8). Thus, to enforce Bellman consistency for a state-based policy, we express the reward distribution under the trajectory distribution generated by $g^\theta$ in terms of the observational data collected under $g^b$.

In the following subsections, we show the steps to achieve this via a proximal OPE formulation. Essentially, we rewrite the confounded reward expectation using only observable quantities. Before presenting the proximal analysis, we introduce the probabilistic notation used throughout.

### C. Notations

We represent probabilities over latent states, observations, and rewards using vector and matrix notation. For any realizations $s_t, s_{t+1} \in \mathcal{S}$ and any action $u_t$, the transition probability $p(s_{t+1} \mid s_t, u_t)$ is collected into the column vector

$$P(S_{t+1} \mid s_t, u_t) = \left(p(s_{t+1}^1 \mid s_t, u_t), \ldots, p(s_{t+1}^{|\mathcal{S}|} \mid s_t, u_t)\right)^\top,$$

and into the row vector

$$P(s_{t+1} \mid S_t, u_t) = \left(p(s_{t+1} \mid s_t^1, u_t), \ldots, p(s_{t+1} \mid s_t^{|\mathcal{S}|}, u_t)\right).$$

In our analysis, the multiplication of a pair of vectors of appropriate dimension is their *scalar product*. The full transition kernel is expressed as the matrix

$$P(S_{t+1} \mid S_t, U_t) \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|},$$

with analogous conventions for matrices involving joint variables or rectangular conditional distributions. In our analysis, the multiplication of two matrices refers to their *algebraic matrix multiplication*. Our analysis also makes use of the conditional independencies among variables in the C-MDP, based on their causal relationship illustrated in Fig. 1. As an example, the causal relationship corresponding to the Markovian evolution of the state given by

$$p(s_{t+1} \mid s_{0:t}, u_{0:t}) = p(s_{t+1} \mid s_t, u_t)$$

is denoted compactly by

$$s_{t+1} \perp\!\!\!\perp (s_{0:t-1}, u_{0:t-1}) \mid (s_t, u_t).$$

For any policy $g^i, i \in \{b, \theta\}$, the distribution induced by $g^i$ on any tuple of random variables is denoted by a $p^i(\cdot)$.

These conventions allow us to express the reward distribution as matrix products, which is crucial for the proximal OPE construction.

### D. OPE to Compute the Bellman Equation

We now show how the Bellman equation in Problem 1 can be evaluated from the restricted dataset $\mathcal{D}^b$ by using the POMDP formulation introduced earlier. Recall that the full state is $S_t = (X_t, Z_t)$, the observation is $Y_t = X_t$, and the reward is $R_t = r(S_t, U_t)$. Each policy $g^\theta$ induces a probability measure on full trajectories $\tau_t = (s_{0:t}, y_{0:t}, u_{0:t})$, while the dataset contains only $\tau_t^o = (y_{0:t}, u_{0:t})$. The expectations in the Bellman equation arise from this trajectory distribution, but the hidden context $Z_t$ makes the reward component confounded. Recall the Bellman consistency constraint

$$Q_t^{\theta, \phi}(x_t, u_t) = \mathbb{E}^\theta \left[ r(X_t, Z_t, U_t) \mid X_t = x_t, U_t = u_t \right]$$

$$+ \mathbb{E}^{\boldsymbol{\phi}}\Big[V_{t+1}^{\boldsymbol{\theta},\boldsymbol{\phi}}(X_{t+1}) \,\big|\, X_t = x_t, U_t = u_t\Big], \tag{11}$$

For a parameterized policy $\boldsymbol{g}^{\boldsymbol{\theta}}$, the second term depends only on the surrogate transition model $p^{\boldsymbol{\phi}}(x_{t+1} \mid x_t, u_t)$ and can be computed from observable transitions $(x_t, u_t, x_{t+1})$ in the dataset.

The first term, $\mathbb{E}^{\boldsymbol{\theta}}[r(X_t, Z_t, U_t) \mid x_t, u_t]$, depends on the unobserved context $Z_t$ and is therefore not identifiable by data-based averaging. Let $\mathcal{R}$ denote the finite set of possible reward realizations. Thus, we can express

$$\mathbb{E}^{\boldsymbol{\theta}}\left[r(X_t, Z_t, U_t)\right] = \mathbb{E}^{\boldsymbol{\theta}}\left[R_t\right] = \sum_{r_t \in \mathcal{R}} r_t \cdot p^{\boldsymbol{\theta}}(r_t). \tag{12}$$

We consider $p^{\boldsymbol{\theta}}(r_t)$ and use the law of total probability and Bayes' theorem to get

$$p^{\boldsymbol{\theta}}(r_t) = \sum_{\tau_t \in \mathcal{T}_t} p^{\boldsymbol{\theta}}(r_t|\tau_t)\, p^{\boldsymbol{\theta}}(\tau_t), \tag{13}$$

$$= \sum_{\tau_t \in \mathcal{T}_t} p^{\boldsymbol{\theta}}(r_t|s_t, u_t)\, p^{\boldsymbol{\theta}}(\tau_t), \tag{14}$$

where, in the last equality, we use the conditional independence $r_t \perp\!\!\!\perp \tau_{t-1}|(s_t, u_t)$ between the reward $r_t$ and past trajectory $\tau_{t-1}$.

Thus, in the POMDP representation, the reward component of the Bellman equation is confounded, while the future-value term remains identifiable from observable transitions. To proceed with model learning, we therefore require a method to express $\mathbb{E}^{\boldsymbol{\theta}}[r(X_t, Z_t, U_t) \mid x_t, u_t]$ using only the observable trajectory prefixes $\tau_t^o$ available in $D^{\boldsymbol{b}}$. This is precisely where proximal learning becomes essential. Proximal OPE provides a method to replace the confounded reward expectation with observable quantities from data and a sequence of weight matrices that are proxy-based. The past and current observations are used as proxies to recover the average effect of the hidden context on rewards without ever observing $Z_t$.

### E. Proximal Identification of the Confounded Reward Term

In this subsection, we extend the proximal learning framework of [17] to the POMDP formulation developed above. Our objective is to express the confounded reward distribution in (14), originally defined with respect to the full latent state, in an equivalent form that depends only on observable variables. This reformulation is essential to ensure that the Bellman operator used in Problem 1 can be evaluated without access to the hidden context $Z_t$. To this end, we establish a sequence of lemmas that progressively eliminate the dependence on the full state $S_t$, thereby yielding a reward model identifiable from observational data. Throughout this subsection, we assume the following assumptions hold.

**Assumption 1.** Access is available to an auxiliary "null" observation $Y_N$, distributed according to the prior over the initial observation $Y_0$. This variable serves as a proxy for the latent initial state $S_0$ and enables identification of the confounded reward distribution at $t = 0$.

**Assumption 2.** For each $t = 1, \dots, T$, the matrices $P^{\boldsymbol{b}}(Y_t \mid Y_{t-1}, u_t)$ and $P^{\boldsymbol{b}}(S_t \mid Y_{t-1}, u_t)$ are invertible.

This invertibility ensures that past and present observations carry sufficient information about the latent state $S_t$ so that hidden quantities can be inferred using proxy variables.

**Remark 2.** We emphasize that Assumption 2 does *not* hold in settings where the context $Z_t$ is resampled independently at each time step. In such cases, there is no temporal carryover of contextual information. Hence, the past observations cannot contain the proxy signal required to recover $Z_t$.

Each lemma isolates a structural step in the transformation from the original reward expansion to a fully observable representation.

**Step 1: Factorizing the trajectory distribution.** We first characterize the trajectory distribution to distinguish policy-dependent factors from those governed exclusively by the system dynamics. This separation will subsequently allow the application of proximal methods to the terms affected by unobserved variables.

**Lemma 1** (Expansion of the Policy–Dependent Trajectory Distribution). *At each $t$, for any trajectory $\tau_t = (u_{0:t}, y_{0:t}, s_{0:t})$, the policy-induced trajectory distribution satisfies*

$$p^{\boldsymbol{\theta}}(\tau_t) = \left(\prod_{k=0}^{t} p^{\boldsymbol{\theta}}(u_k \mid y_k)\right)\left(\prod_{k=0}^{t} p^{\boldsymbol{\theta}}(y_k \mid s_k)\right) \\ \cdot \left(\prod_{k=0}^{t-1} p^{\boldsymbol{\theta}}(s_{k+1} \mid s_k, u_k)\right) p^{\boldsymbol{\theta}}(s_0). \tag{15}$$

*Proof.*

$$p^{\boldsymbol{\theta}}(\tau_t)$$
$$= p^{\boldsymbol{\theta}}(u_t \mid y_t, s_t, \tau_{t-1})\, p^{\boldsymbol{\theta}}(y_t, s_t, \tau_{t-1}), \tag{16}$$
$$= p^{\boldsymbol{\theta}}(u_t \mid y_t)\, p^{\boldsymbol{\theta}}(y_t \mid s_t, \tau_{t-1})\, p^{\boldsymbol{\theta}}(s_t, \tau_{t-1}), \tag{17}$$
$$= p^{\boldsymbol{\theta}}(u_t \mid y_t)\, p^{\boldsymbol{\theta}}(y_t \mid s_t)\, p^{\boldsymbol{\theta}}(s_t \mid \tau_{t-1})\, p^{\boldsymbol{\theta}}(\tau_{t-1}), \tag{18}$$
$$= p^{\boldsymbol{\theta}}(u_t \mid y_t)\, p^{\boldsymbol{\theta}}(y_t \mid s_t)\, p^{\boldsymbol{\theta}}(s_t \mid s_{t-1}, u_{t-1})\, p^{\boldsymbol{\theta}}(\tau_{t-1}), \tag{19}$$

$$= \Pi_{k=0}^{t}\, p^{\boldsymbol{\theta}}(u_k \mid y_k)\, \Pi_{k=0}^{t}\, p^{\boldsymbol{\theta}}(y_k \mid s_k) \\ \cdot \Pi_{k=0}^{t-1}\, p^{\boldsymbol{\theta}}(s_{k+1} \mid s_k, u_k)\, p^{\boldsymbol{\theta}}(s_0), \tag{20}$$

where, in (17), we use the property that the parameterized policy depends only on the current observation. In (18) we use the conditional independence $y_t \perp\!\!\!\perp \tau_{t-1} \mid s_t$ as the observation is dependent purely on the state. Lastly, in (19) we use the Markovian nature of the evolution of the full state. $\square$

**Step 2: Decomposing the reward distribution.** Next, we substitute Lemma 1 into the reward expansion and isolate all policy-dependent terms.

**Lemma 2** (Reward Distribution Decomposition). *The reward distribution satisfies*

$$p^{\boldsymbol{\theta}}(r_t) = \sum_{\tau_t} \left(\prod_{k=0}^{t} p^{\boldsymbol{\theta}}(u_k \mid y_k)\right)\left(\prod_{k=0}^{t} p^{\boldsymbol{b}}(y_k \mid s_k)\right) \\ \cdot p^{\boldsymbol{b}}(r_t \mid s_t, u_t)\left(\prod_{k=0}^{t-1} p^{\boldsymbol{b}}(s_{k+1}, y_k \mid s_k, u_k)\right) p^{\boldsymbol{b}}(s_0), \tag{21}$$

*and equivalently,*

$$p^{\boldsymbol{\theta}}(r_t) = \sum_{\tau_t^o} \left( \prod_{k=0}^{t} p^{\boldsymbol{\theta}}(u_k \mid y_k) \right) p^{\boldsymbol{b}}(r_t, y_t \mid S_t, u_t)$$
$$\cdot \left( \prod_{k=0}^{t-1} p^{\boldsymbol{b}}(S_{k+1}, y_k \mid S_k, u_k) \right) p^{\boldsymbol{b}}(S_0). \quad (22)$$

*Proof.* We begin by noting a key implication of Lemma 1. Among the factors that appear in $p^{\boldsymbol{\theta}}(\tau_t)$, the only quantities that depend on the parameterized policy $\boldsymbol{g}^{\boldsymbol{\theta}}$ are the action probabilities $p^{\boldsymbol{\theta}}(u_k \mid y_k)$. In contrast, the observation model $p^{\boldsymbol{\theta}}(y_k \mid s_k)$ and the state transition model $p^{\boldsymbol{\theta}}(s_{k+1} \mid s_k, u_k)$ are structural properties of the underlying system and therefore do not change with the choice of policy. Hence, we can evaluate all the policy-independent terms based on any data-generating policy. This invariance is crucial, as it allows us to substitute these terms with their behavioral counterparts $p^{\boldsymbol{b}}(\cdot)$ as we expand the reward distribution.

Starting from

$$p^{\boldsymbol{\theta}}(r_t) = \sum_{\tau_t} p^{\boldsymbol{\theta}}(r_t \mid s_t, u_t) \, p^{\boldsymbol{\theta}}(\tau_t), \quad (23)$$

we substitute the factorization of Lemma 1. As a result of the invariance, at each $k = 0, \ldots, t-1$, we rewrite

$$p^{\boldsymbol{\theta}}(y_k \mid s_k) = p^{\boldsymbol{b}}(y_k \mid s_k), \quad (24)$$
$$p^{\boldsymbol{\theta}}(s_{k+1} \mid s_k, u_k) = p^{\boldsymbol{b}}(s_{k+1} \mid s_k, u_k). \quad (25)$$

Similarly, $p^{\boldsymbol{\theta}}(r_t \mid s_t, u_t) = p^{\boldsymbol{b}}(r_t \mid s_t, u_t)$ and $p^{\boldsymbol{\theta}}(s_0) = p^{\boldsymbol{b}}(s_0)$. We substitute these behavioral policy-based terms and collect the factors to produce (21). Finally, using our vector and matrix notation, we group each transition-observation pair into $p^{\boldsymbol{b}}(S_{k+1}, y_k \mid S_k, u_k)$ to produce (22). $\quad\square$

At this point, the decomposition isolates the core challenge, which is the terms that still depend on the latent state $S_k$. The next two lemmas show how these quantities can be rewritten using observable proxies. This removes all dependence on the unobserved context.

**Step 3: Introducing proxy variables.** Next, we show that past and present observations $(Y_{t-1}, Y_t)$ form valid proxies for $S_t$ and allow us to rewrite the terms based on the full-state.

**Lemma 3** (First Proxy: Incorporating $Y_{t-1}$). *The matrix $P^{\boldsymbol{b}}(S_{t+1}, y_t \mid S_t, u_t)$ satisfies*

$$P^{\boldsymbol{b}}(S_{t+1}, y_t \mid S_t, u_t)$$
$$= P^{\boldsymbol{b}}(S_{t+1}, y_t \mid Y_{t-1}, u_t) \, P^{\boldsymbol{b}}(S_t \mid Y_{t-1}, u_t)^{-1}. \quad (26)$$

*Proof.* Using the law of total probability,

$$P^{\boldsymbol{b}}(S_{t+1}, y_t \mid Y_{t-1}, u_t)$$
$$= P^{\boldsymbol{b}}(S_{t+1}, y_t \mid S_t, Y_{t-1}, u_t) \, P^{\boldsymbol{b}}(S_t \mid Y_{t-1}, u_t), \quad (27)$$
$$= P^{\boldsymbol{b}}(S_{t+1}, y_t \mid S_t, u_t) \, P^{\boldsymbol{b}}(S_t \mid Y_{t-1}, u_t), \quad (28)$$

where we use conditional independence $(S_{t+1}, y_t) \perp\!\!\!\perp Y_{t-1} \mid (S_t, u_t)$. Under Assumption 2, we right-multiply $P^{\boldsymbol{b}}(S_t \mid Y_{t-1}, u_t)^{-1}$ on both sides to obtain the desired expression. $\quad\square$

Thus, $Y_{t-1}$ can be incorporated as a proxy for $S_t$. To complete the construction, we also need to express $P^{\boldsymbol{b}}(S_t \mid Y_{t-1}, u_t)$ using $Y_t$.

**Lemma 4** (Second Proxy: Incorporating the current observation $Y_t$). *The following expression holds:*

$$P^{\boldsymbol{b}}(S_t \mid Y_{t-1}, u_t) = P^{\boldsymbol{b}}(Y_t \mid S_t, u_t)^{-1} \, P^{\boldsymbol{b}}(Y_t \mid Y_{t-1}, u_t). \quad (29)$$

*Proof.* From the law of total probability,

$$P^{\boldsymbol{b}}(Y_t \mid Y_{t-1}, u_t) = P^{\boldsymbol{b}}(Y_t \mid S_t, u_t) \, P^{\boldsymbol{b}}(S_t \mid Y_{t-1}, u_t). \quad (30)$$

Under Assumption 2, we invert $P^{\boldsymbol{b}}(Y_t \mid S_t, u_t)$ to obtain the desired expression. $\quad\square$

**Step 5: Completing the proximal identification.** Next, we combine all previous steps to express the confounded reward distribution using only observable matrices and the evaluation policy.

**Theorem 1** (Proximal Identification of the Confounded Reward Term). *Under Assumptions 1 and 2,*

$$p^{\boldsymbol{\theta}}(r_t)$$
$$= \sum_{\tau_t^o} \left( \prod_{k=0}^{t} p^{\boldsymbol{\theta}}(u_k \mid y_k) \right)$$
$$\cdot p^{\boldsymbol{b}}(r_t, y_t \mid Y_{t-1}, u_t) \cdot \left( \prod_{k=0}^{t} W_k(\tau_k^o) \right). \quad (31)$$

*where, for each $k = 1, \ldots, t$, the weight matrix is given by*

$$W_k(\tau_k^o)$$
$$= P^{\boldsymbol{b}}(Y_k \mid Y_{k-1}, u_k)^{-1} \cdot P^{\boldsymbol{b}}(Y_k, y_{k-1} \mid Y_{k-2}, u_{k-1}), \quad (32)$$

*and for $k = 0$ we set*

$$W_0(\tau_0^o) = P^{\boldsymbol{b}}(Y_0 \mid u_0, Y_N)^{-1} \cdot P^{\boldsymbol{b}}(Y_0), \quad (33)$$

*with $Y_N$ denoting the null observation associated with the prior in Assumption 1. Hence, the confounded reward expectation $\mathbb{E}^{\boldsymbol{\theta}}[r_t]$ is identifiable from observable data.*

*Proof.* Lemma 2 expresses the reward distribution as

$$p^{\boldsymbol{\theta}}(r_t) = \sum_{\tau_t^o} \left( \prod_{k=0}^{t} p^{\boldsymbol{\theta}}(u_k \mid y_k) \right) p^{\boldsymbol{b}}(r_t, y_t \mid S_t, u_t)$$
$$\cdot \left( \prod_{k=0}^{t-1} p^{\boldsymbol{b}}(S_{k+1}, y_k \mid S_k, u_k) \right) p^{\boldsymbol{b}}(S_0). \quad (34)$$

Thus, the only remaining dependence on the latent state appears through the product

$$p^{\boldsymbol{b}}(r_t, y_t \mid S_t, u_t) \cdot \left( \prod_{k=0}^{t-1} p^{\boldsymbol{b}}(S_{k+1}, y_k \mid S_k, u_k) \right) p^{\boldsymbol{b}}(S_0). \quad (35)$$

Our goal is to show that this product can be rewritten in terms of observable quantities via the weight matrices $W_k(\tau_k^o)$.

At each $k$, we apply the proxy-based decompositions from Lemmas 3 and 4 to obtain

$$P^{\boldsymbol{b}}(S_{k+1}, y_k \mid S_k, u_k)$$
$$= P^{\boldsymbol{b}}(S_{k+1}, y_k \mid Y_{k-1}, u_k) \cdot P^{\boldsymbol{b}}(Y_k \mid Y_{k-1}, u_k)^{-1}$$
$$\cdot P^{\boldsymbol{b}}(Y_k \mid S_k, u_k). \tag{36}$$

Similarly, the reward-based term can also be decomposed to obtain

$$P^{\boldsymbol{b}}(r_t, y_t \mid S_t, u_t)$$
$$= P^{\boldsymbol{b}}(r_t, y_t \mid Y_{t-1}, u_t) \cdot P^{\boldsymbol{b}}(Y_t \mid Y_{t-1}, u_t)^{-1}$$
$$\cdot P^{\boldsymbol{b}}(Y_t \mid S_t, u_t). \tag{37}$$

Next, we consider the product of two consecutive full-state dependent terms:

$$P^{\boldsymbol{b}}(S_{k+1}, y_k \mid S_k, u_k) \cdot P^{\boldsymbol{b}}(S_k, y_{k-1} \mid S_{k-1}, u_{k-1})$$
$$= P^{\boldsymbol{b}}(S_{k+1}, y_k \mid Y_{k-1}, u_k) \cdot P^{\boldsymbol{b}}(Y_k \mid Y_{k-1}, u_k)^{-1}$$
$$\cdot \boxed{P^{\boldsymbol{b}}(Y_k \mid S_k, u_k) \cdot P^{\boldsymbol{b}}(S_k, y_{k-1} \mid Y_{k-2}, u_{k-1})}$$
$$\cdot P^{\boldsymbol{b}}(Y_{k-1} \mid Y_{k-2}, u_{k-1})^{-1} \cdot P^{\boldsymbol{b}}(Y_{k-1} \mid S_{k-1}, u_{k-1}), \tag{38}$$

where we have applied the same proxy expansions to the second factor. We now focus on the product of the highlighted inner full-state dependent terms in (38).

$$P^{\boldsymbol{b}}(Y_k \mid S_k, u_k) \cdot P^{\boldsymbol{b}}(S_k, y_{k-1} \mid Y_{k-2}, u_{k-1}).$$

Consider the observable matrix $P^{\boldsymbol{b}}(Y_k, y_{k-1} \mid Y_{k-2}, u_{k-1})$. Using the law of total probability, we expand the observable matrix

$$P^{\boldsymbol{b}}(Y_k, y_{k-1} \mid Y_{k-2}, u_{k-1})$$
$$= \sum_{s_k \in \mathcal{S}} P^{\boldsymbol{b}}(Y_k \mid s_k, y_{k-1}, Y_{k-2}, u_{k-1})$$
$$\cdot P^{\boldsymbol{b}}(s_k, y_{k-1} \mid Y_{k-2}, u_{k-1}), \tag{39}$$

where, by the observation-model-based conditional independence

$$Y_k \perp\!\!\!\perp (y_{k-1}, Y_{k-2}) \mid (S_k, u_k),$$

the conditional probability simplifies to

$$P^{\boldsymbol{b}}(Y_k \mid s_k, y_{k-1}, Y_{k-2}, u_{k-1}) = P^{\boldsymbol{b}}(Y_k \mid s_k, u_k).$$

Thus, the inner terms in (38) can now be rewritten to incorporate the proxies as

$$P^{\boldsymbol{b}}(Y_k \mid S_k, u_k) \cdot P^{\boldsymbol{b}}(S_k, y_{k-1} \mid Y_{k-2}, u_{k-1})$$
$$= P^{\boldsymbol{b}}(Y_k, y_{k-1} \mid Y_{k-2}, u_{k-1}). \tag{40}$$

*Defining the weight matrix:* Motivated by (36) and (40), we define, for $k \geq 1$, the observable weight matrix

$$W_k(\tau_k^o)$$

$$= P^{\boldsymbol{b}}(Y_k \mid Y_{k-1}, u_k)^{-1} \cdot P^{\boldsymbol{b}}(Y_k, y_{k-1} \mid Y_{k-2}, u_{k-1}). \tag{41}$$

Equation (40) shows that the factor involving the latent state $S_k$ in the inner product collapses to the observable term $P^{\boldsymbol{b}}(Y_k, y_{k-1} \mid Y_{k-2}, u_{k-1})$, while the remaining pre-multiplication by $P^{\boldsymbol{b}}(Y_k \mid Y_{k-1}, u_k)^{-1}$ is also observable. We collect these terms into $W_k(\tau_k^o)$. For $k = 0$, the same construction is applied at the initial time step, using Assumption 1, yields

$$W_0(\tau_0^o) = P^{\boldsymbol{b}}(Y_0 \mid u_0, Y_N)^{-1} \cdot P^{\boldsymbol{b}}(Y_0), \tag{42}$$

where $Y_N$ denotes the null observation associated with the prior.

Substituting (40) into the expression 38 and using the definition of $W_k(\tau_k^o)$, we obtain

$$P^{\boldsymbol{b}}(S_{k+1}, y_k \mid S_k, u_k) \cdot P^{\boldsymbol{b}}(S_k, y_{k-1} \mid S_{k-1}, u_{k-1})$$
$$= P^{\boldsymbol{b}}(S_{k+1}, y_k \mid Y_{k-1}, u_k) \cdot W_k(\tau_k^o)$$
$$\cdot P^{\boldsymbol{b}}(Y_{k-1} \mid Y_{k-2}, u_{k-1})^{-1} \cdot P^{\boldsymbol{b}}(Y_{k-1} \mid S_{k-1}, u_{k-1}),$$

so that all dependence on $S_k$ has been removed and its effect is encoded by the observable matrix $W_k(\tau_k^o)$ and boundary terms involving $Y_{k-1}$ and $S_{k-1}$.

We apply this reduction iteratively for $k = 0, \ldots, t$ to the full-state based product (35) to obtain a chain of observable weight matrices

$$p^{\boldsymbol{b}}(r_t, y_t \mid Y_{t-1}, u_t) \cdot \left( \prod_{k=0}^{t} W_k(\tau_k^o) \right). \tag{43}$$

multiplying boundary terms that only depend on $(Y_0, S_0)$ and the initial prior in Assumption 1. The remaining dependence on $S_0$ is captured by $p^{\boldsymbol{b}}(S_0)$ and is therefore compatible with the observable representation.

We Substitute the resulting expression for the product of the latent-state factors back into the decomposition of Lemma 2 to produce (31), completing the proof. $\square$

## IV. NUMERICAL SIMULATION

In this section, we illustrate the proposed framework on a synthetic clinical decision-making task. The example is motivated by intensive care unit (ICU) settings, where doctors prescribe treatments over time. In practice, treatments are repeatedly adjusted on the basis of recorded physiological measurements and unrecorded clinical judgment. We use this simulation to illustrate why it is necessary to learn a surrogate transition model that satisfies the proposed modified Bellman consistency equation to enable reliable model-based planning.

### A. Clinical Treatment Environment

We consider a finite-horizon C-MDP that represents a short treatment episode for a single patient over a horizon of $T = 9$ hours. The discrete state space $\mathcal{X}$ aggregates physiological measurements like heart rate, blood pressure, and oxygen saturation into coarse health states to represent a severity score:

$$\mathcal{X} = \{x^1, x^2, x^3, x^4\},$$

where $x^1$ denotes a *stable* state, $x^2$ a *borderline* state, $x^3$ a *deteriorating* state, and $x^4$ a *critical* state.

At each time $t$, the doctor chooses an action corresponding to the intensity of treatment, from the set

$$\mathcal{U} = \{u^1, u^2, u^3\},$$

where $u^1$ is *conservative treatment* (standard care), $u^2$ is *moderate escalation*, and $u^3$ is *aggressive escalation*.

The key feature is an unrecorded, patient-specific context $Z_t \in \mathcal{Z}$ that captures latent factors such as comorbidities and frailty, which the doctor can observe or access. We consider

$$\mathcal{Z} = \{z^L, z^H\},$$

where $z^L$ denotes a *low-risk* profile and $z^H$ a *high-risk* profile. We assume that the initial context $Z_0$ is drawn from a distribution $\nu$ and is persistent thereafter. Hence, $Z_{t+1} = Z_t$ for all $t$, which is consistent with the special-case causal graph illustrated in Fig. 1. For example, an aggressive treatment $U_t = U^2$ may move a low-risk critical patient from $x^4$ to $x^2$ with higher probability than a high-risk patient in the same observed state. We design the transition kernel to reflect this intuition. At each time $t$, the reward $R_t = r(X_t, Z_t, U_t)$ reflects clinical status and treatment burden. We use

$$r(x, z, u) = r_{\text{health}}(x, z) - \lambda\, c(u),$$

where $r_{\text{health}}(x, z)$ assigns higher reward to less severe states and may depend on the latent context. The term $c(u)$ is indicative of the cost or burden of treatment intensity and weight $\lambda \in \mathbb{R}_+$ controls the relative importance of this cost. In our setting, $r_{\text{health}}$ is the dominant component of the reward, while the penalty $\lambda c(u)$ discourages the unnecessary use of high-intensity treatments.

The clinician follows a behavioral policy $\boldsymbol{g}^{\boldsymbol{b}} = \{g_t^{\boldsymbol{b}}\}_{t=0}^{T-1}$ that depends on both $X_t$ and $Z_t$. Thus, high-risk patients (with $z^H$) may receive aggressive treatment $u^2$ more often than low-risk patients, even when $X_t$ is the same. The offline dataset $\mathcal{D}^b = \left\{ (x_0^i, u_0^i, r_0^i, \ldots, x_T^i, u_T^i, r_T^i) \right\}_{i=1}^N$ contains only observable state–action–reward tuples. The latent context $Z_t$ that influenced both $U_t$ and $R_t$ is never recorded.

In the experiments, we specify concrete transition and reward parameters for each $(x, z, u)$ and choose $g_b$ to mimic a risk-aware clinician. We then compare a naive model that uses data averaging and the proposed surrogate MDP with proximal reward correction, with an oracle model that observes $(X_t, Z_t)$. Although our primary focus is on model quality, we note that in this setting the policy learned along with the proximal surrogate attains an expected return approximately 1.4% higher than that of the naive learner (Monte Carlo estimate with $N = 1000$ episodes).

Figure 2 reports the multi-step rollout error in the C-MDP. We measure the $\ell_1$ distance between the state distribution induced by each learned model and that of the true system under a fixed policy. The empirical state distributions are estimated from $N = 1000$ Monte Carlo episodes. The naive model attains a good one-step fit to the confounded data, but the rollout error rapidly accumulates and remains large over the horizon. This behavior indicates that repeated application
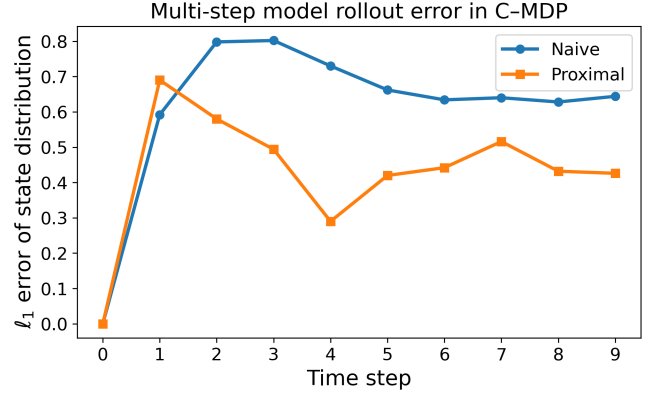


Fig. 2: Comparison of multi-step rollout error

of the biased transition kernel leads to distorted evolution of the state. In contrast, the surrogate model based on proximal learning exhibits consistently smaller $\ell_1$ error for $t \geq 2$. This suggests that the combination of the proximal reward construction and the surrogate dynamics yields a more accurate model of the system. In particular, the learned surrogate better reflects the long-horizon influence of the latent risk context on the observable state evolution, even though $Z_t$ is never observed in the offline dataset.

## V. CONCLUSION

In this paper, we presented a model-based reinforcement learning framework for contextual MDPs in which the unobserved context confounds both the transition and reward mechanisms. By recasting the problem as a POMDP and analyzing it from a causal perspective, we showed that the reward component of the Bellman equation is not identifiable from observational data alone. To overcome this issue, we adapt a proximal off-policy evaluation method that reconstructs the confounded reward expectation using observable proxy variables under standard invertibility conditions, yielding an identifiable surrogate reward model compatible with state-based policies.

Combining this proximal correction with a behavior-averaged transition model and a maximum causal entropy formulation produced a Bellman-consistent surrogate MDP suitable for planning in the presence of hidden confounding. The results extend the applicability of model-based reinforcement learning to settings in which contextual information is unavailable or costly to obtain. A potential direction for future research should explore extensions to continuous state spaces and relaxations of the invertibility requirements.

## REFERENCES

[1] R. S. Sutton, "Dyna, an integrated architecture for learning, planning, and reacting," *ACM SIGART Bulletin*, vol. 2, no. 4, pp. 160–163, 1991.
[2] M. Deisenroth and C. E. Rasmussen, "Pilco: A model-based and data-efficient approach to policy search," in *Proceedings of the 28th International Conference on Machine Learning (ICML)*, pp. 465–472, 2011.
[3] A. Nagabandi, G. Kahn, R. S. Fearing, and S. Levine, "Neural network dynamics for model-based deep reinforcement learning with model-free fine-tuning," in *IEEE International Conference on Robotics and Automation (ICRA)*, pp. 7559–7566, 2018.

[4] K. Chua, R. Calandra, R. McAllister, and S. Levine, "Deep reinforcement learning in a handful of trials using probabilistic dynamics models," in *Advances in Neural Information Processing Systems*, vol. 31, 2018.

[5] D. Hafner, T. Lillicrap, I. Fischer, R. Villegas, D. Ha, H. Lee, and J. Davidson, "Learning latent dynamics for planning from pixels," in *Proceedings of the 36th International Conference on Machine Learning (ICML)*, pp. 2555–2565, 2019.

[6] N. Venkatesh, V.-A. Le, A. Dave, and A. A. Malikopoulos, "Connected and automated vehicles in mixed-traffic: Learning human driver behavior for effective on-ramp merging," in *2023 62nd IEEE Conference on Decision and Control (CDC)*, pp. 92–97, IEEE, 2023.

[7] A. Dave, H. Bang, and A. A. Malikopoulos, "A framework for effective ai recommendations in cyber-physical-human systems," *IEEE Control Systems Letters*, vol. 8, pp. 1379–1384, 2024.

[8] W. Sun, H. Bang, and A. A. Malikopoulos, "Ai recommendation systems for lane-changing using adherence-aware reinforcement learning," in *28th IEEE International Conference on Intelligent Transportation Systems (ITSC)*, pp. 1060—1065, 2025.

[9] E. J. Tchetgen Tchetgen, Z. Ying, S. Yang, and Y. Cui, "Introduction to proximal causal inference," *arXiv preprint arXiv:2009.10982*, 2020.

[10] Y. Cui and E. J. Tchetgen Tchetgen, "Semiparametric proximal causal learning," *arXiv preprint arXiv:2008.06533*, 2020.

[11] W. Miao, Z. Geng, and E. J. Tchetgen Tchetgen, "Identifying causal effects with proxy variables of an unmeasured confounder," *Biometrika*, vol. 105, no. 4, pp. 987–993, 2018.

[12] N. Kallus and A. Zhou, "Confounding-robust policy evaluation in infinite-horizon reinforcement learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 22293–22304, 2020.

[13] M. Makar, H. Narasimhan, E. Brunskill, and S. Murphy, "Causal reinforcement learning under unobserved confounding: Theory and algorithms," *arXiv preprint arXiv:2206.15475*, 2022.

[14] A. Bennett and N. Kallus, "Proximal reinforcement learning: Efficient off-policy evaluation in partially observed markov decision processes," *Operations Research*, vol. 72, no. 3, pp. 1071–1086, 2024.

[15] C. Kausik, Y. Lu, K. Tan, M. Makar, Y. Wang, and A. Tewari, "Offline policy evaluation and optimization under confounding," in *International Conference on Artificial Intelligence and Statistics*, pp. 1459–1467, PMLR, 2024.

[16] C. Kausik, Y. Lu, K. Tan, M. Makar, Y. Wang, and A. Tewari, "Offline policy evaluation and optimization under confounding," *arXiv preprint arXiv:2311.14645*, 2024.

[17] G. Tennenholtz, U. Shalit, and S. Mannor, "Off-policy evaluation in partially observable environments," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, pp. 10276–10283, 2020.

[18] B. D. Ziebart and J. A. Bagnell, "Modeling interaction via the principle of maximum causal entropy," in *Proceedings of the 27th International Conference on Machine Learning (ICML)*, pp. 1147–1154, 2010.

[19] M. Herman, T. Gindele, J. Wagner, F. Schmitt, and W. Burgard, "Inverse reinforcement learning with simultaneous estimation of rewards and dynamics," in *Artificial intelligence and statistics*, pp. 102–110, PMLR, 2016.

[20] S. Reddy, A. D. Dragan, and S. Levine, "Where do you think you're going? inferring beliefs about dynamics from behavior," in *Advances in Neural Information Processing Systems*, vol. 31, 2018.

[21] J. Pearl, *Causality: Models, Reasoning, and Inference*. Cambridge University Press, 2009.