

---

# Structured Behavioral Heterogeneity as Latent Regime Constraints

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Sequential decision data often exhibit *structured behavioral heterogeneity*: similar  
2 observed conditions can lead to different actions across trajectories and time. Stan-  
3 dard approaches typically attribute this variation to drifting rewards or unstructured  
4 noise, which can produce models that are difficult to interpret and compare across  
5 settings. We instead model this heterogeneity as arising from **latent constraint**  
6 **regimes**, which induce persistent, regime-dependent deviations from a shared  
7 decision anchor. We introduce a framework that decomposes behavior into (i) a  
8 *shared anchor value function* capturing stable, outcome-relevant decision logic and  
9 (ii) a *switching regime process* that generates sparse, interpretable deviations from  
10 this anchor. This decomposition separates stable decision structure from systematic  
11 variation in observed actions. To construct the anchor under partial observability  
12 and irregular decision timing, we use a continuous-time latent-state model that  
13 infers belief states from irregularly timed decision data and produces compara-  
14 ble action values across trajectories. The regime layer then captures structured  
15 deviations without modifying the underlying objective. Across clinical care and  
16 retail pricing datasets, the proposed approach improves held-out action prediction  
17 relative to stationary models and interpretable behavioral baselines, while identi-  
18 fying persistent regime assignments and sparse, feature-level deviations. These  
19 findings show that structured heterogeneity can be captured by a shared anchor  
20 with switching, regime-dependent deviations.

## 21 1 Introduction

22 Sequential decisions in clinical care, dynamic pricing, and operations are commonly formalized  
23 as Markov decision processes [1], with offline reinforcement learning used to learn from logged  
24 trajectories despite distribution shift [2, 3, 4]. A recurring empirical pattern, however, is *apparent*  
25 *behavioral heterogeneity*: under similar observed conditions, agents take different actions across time  
26 and contexts [5, 6].

27 In ICU sepsis management, for example, patients with nearly identical clinical measurements may still  
28 receive different treatment decisions. Importantly, these differences often do not reflect disagreement  
29 about medical objectives. Instead, the decision environment is continuously reshaped by latent  
30 factors: medication shortages, staffing fluctuations across shifts, evolving hospital protocols, or  
31 temporary restrictions on intervention options. The observed behavior therefore reflects not only  
32 clinical reasoning, but also persistent structure in the surrounding decision context.

33 Standard modeling approaches typically absorb such variation into drifting rewards, flexible policy  
34 classes, or unstructured noise [7, 8]. While expressive, this view effectively attributes all heterogeneity  
35 to changes in *what is optimized*, conflating stable decision logic with systematic distortions induced  
36 by context. As a result, models may fit behavior well but offer limited insight into why similar  
37 situations lead to different actions.

38 This limitation is reflected across major paradigms in offline reinforcement learning and imitation  
 39 learning. Inverse reinforcement learning explains heterogeneity by recovering different reward  
 40 functions that rationalize observed behavior [9, 10]. Mixture and hierarchical policy models capture  
 41 variation via latent tasks or contexts that induce different policies [11, 12, 13]. Bounded-rational and  
 42 cooperative formulations attribute it to latent cognitive constraints or interaction structure [14, 15],  
 43 while classical switching state-space models drive behavior changes through latent regime dynamics  
 44 [16, 17]. Across these views, heterogeneity is embedded in the decision problem itself rather than  
 45 separated from the mechanisms that systematically distort decisions within a shared structure.

46 We propose a different perspective: much of observed heterogeneity can be understood as arising  
 47 from *latent constraint regimes* acting on a shared decision anchor. Here, “constraint” denotes  
 48 persistent deviations from  $Q_{\text{anc}}$ , not necessarily literal hard constraints. Instead of allowing rewards  
 49 or policies to vary freely, we decompose behavior into (i) *a shared action-value anchor capturing*  
 50 *stable, outcome-relevant decision structure* and (ii) *a switching regime process that induces sparse,*  
 51 *structured deviations from this anchor*. These regimes are not alternative tasks or objectives; they  
 52 represent persistent modes of constraint that systematically reshape action selection relative to a  
 53 common decision rule.

54 A central difficulty in making this decomposition operational is that *decision data are often partially*  
 55 *observed and irregularly timed*. Constructing a meaningful shared anchor therefore requires aligning  
 56 value estimates across uneven decision processes while preserving comparability. To address this,  
 57 we infer decision-time belief states using a *continuous-time latent-state model*, which produces  
 58 consistent state representations under irregular observation patterns [18, 19, 20]. Conditioned on  
 59 these beliefs, the anchor defines a stable value function, while the regime layer captures deviations  
 60 without modifying the underlying decision structure.

61 This separation yields a representation in which heterogeneity is expressed explicitly as structured  
 62 deviations around a shared anchor, rather than implicitly through changing rewards or latent tasks.  
 63 Across clinical care and retail pricing datasets, this formulation improves held-out action prediction  
 64 over stationary models and interpretable behavioral baselines, while recovering persistent regime  
 65 structure with sparse, feature-level interpretations. Beyond predictive performance, the learned  
 66 regimes provide a direct lens for understanding systematic departures from shared decision logic in  
 67 terms of recurring constraint patterns.

68 Finally, this viewpoint contrasts with prior approaches along a simple axis: rather than explaining  
 69 heterogeneity by changing the decision problem (rewards, tasks, policies, or latent regimes), *we keep*  
 70 *the decision structure fixed and explain variation through regime-dependent constraints acting on top*  
 71 *of it*. This leads to a more factored interpretation of sequential behavior, where stability lies in the  
 72 anchor and heterogeneity lies in structured deviations. An extended discussion of related work is  
 73 provided in Appendix A.

## 74 2 Method

75 We model structured behavioral heterogeneity in sequential decision logs. The goal is not to learn  
 76 a separate reward, policy, or MDP for each latent group. Instead, we ask whether heterogeneous  
 77 actions can be explained as persistent, structured deviations from a shared value-based anchor.

78 The method has two main stages (as illustrated in Figure 1). **Stage I** constructs decision-time latent  
 79 states and a shared anchor value function. **Stage II** learns latent constraint regimes that modulate  
 80 action selection relative to this anchor. At a high level, the action score under regime  $u$  is

$$S_u(x, a) = \underbrace{Q_{\text{anc}}(x, a)}_{\text{shared anchor value}} - \underbrace{C_u(\xi, a)}_{\text{regime-dependent deviation}},$$

81 where  $x$  is the inferred decision-time state and  $\xi$  collects the variables available to the regime layer,  
 82 such as  $x$  itself and optional side information. The anchor  $Q_{\text{anc}}$  is shared across all trajectories  
 83 and regimes. Heterogeneity enters only through the regime-dependent deviation  $C_u$ . This gives a  
 84 common coordinate system for comparing behavior: regimes differ by how they modulate the same  
 85 value anchor, not by changing the value function itself.

### 86 2.1 Sequential decision data

87 We observe  $N$  trajectories. In trajectory  $i$ , decisions occur at irregular times

$$t_1^i < t_2^i < \dots < t_{M_i}^i.$$

88 At decision epoch  $m$ , we observe covariates  $o_m^i$  and an action  $a_m^i \in \mathcal{A}$ , where  $\mathcal{A}$  is finite:

$$\mathcal{D}^i = \{(t_m^i, o_m^i, a_m^i)\}_{m=1}^{M_i}.$$

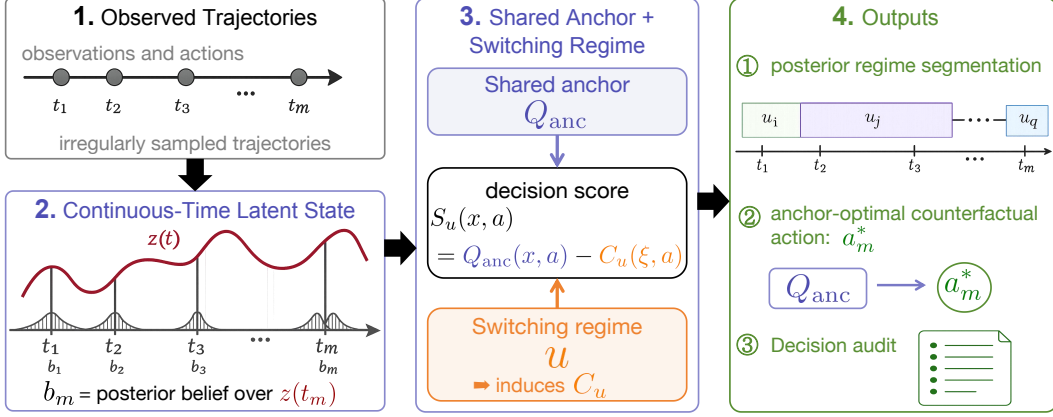


Figure 1: Shared-anchor regime decomposition. Irregular observations are mapped to continuous-time latent beliefs, which define a shared anchor  $Q_{\text{anc}}$ . A switching regime  $u_m$  induces sparse deviations  $C_u(\xi, a)$  around this anchor, yielding posterior regime segmentation, anchor-optimal counterfactual actions, and decision audits.

89 The observed covariates need not be Markovian or fully informative. We therefore introduce a latent  
 90 system state  $z(t) \in \{1, \dots, K_z\}$ , evolving in continuous time, and use the posterior belief over this  
 91 state as the decision-time representation:

$$b_m^i(k) = \Pr(z(t_m^i) = k \mid o_{1:m}^i, a_{1:m-1}^i).$$

92 We also maintain continuous outcome-proxy variables  $\lambda_m^i = \lambda(t_m^i)$ . The state used by the action  
 93 model is

$$x_m^i = (b_m^i, \lambda_m^i).$$

94 Thus,  $x_m$  summarizes the model’s belief about the underlying system when the action is chosen.

## 95 2.2 Stage I: constructing the shared anchor

96 Stage I produces two objects for every decision epoch:

$$x_m \text{ and } Q_{\text{anc}}(x_m, a), a \in \mathcal{A}.$$

97 The first object is the decision-time belief state. The second is a shared action-value anchor used to  
 98 compare actions across trajectories.

99 In our implementation, the latent system state follows an action-modulated continuous-time Markov  
 100 chain. If action  $a_m$  is held fixed over the interval  $[t_m, t_{m+1})$ , then

$$\Pr(z(t_{m+1}) = \cdot \mid z(t_m), a_m) = \exp\{A(a_m)\Delta_m\}, \quad \Delta_m = t_{m+1} - t_m,$$

101 where  $A(a)$  is the generator under action  $a$ . This directly handles irregular decision times without  
 102 forcing trajectories onto a fixed grid.

103 The outcome proxies  $\lambda(t)$  evolve continuously between observations according to state-dependent  
 104 dynamics; the explicit ODE form is deferred to Appendix C. Observations are generated from a  
 105 likelihood

$$p_\theta(o_m \mid z(t_m), \lambda(t_m)).$$

106 Belief updating then takes the standard prediction-correction form

$$\begin{aligned} 107 \quad \tilde{b}_m &= b_{m-1} \exp\{A(a_{m-1})\Delta_{m-1}\}, \\ b_m(k) &\propto \tilde{b}_m(k) p_\theta(o_m \mid z(t_m) = k, \lambda(t_m)). \end{aligned}$$

108 The exact emission model is application-specific; the key idea is that Stage I learns latent system  
 109 dynamics from irregularly observed trajectories, while treating observed actions as inputs that  
 110 influence state evolution.

111 Given the fitted dynamics, we define the shared anchor value

$$Q_{\text{anc}}(x, a) = \mathbb{E}_{\pi_0} \left[ \int_0^\infty e^{-\rho t} r_{\text{anc}}(z(t), \lambda(t), a(t)) dt \mid x(0) = x, a(0) = a \right],$$

112 where  $r_{\text{anc}}$  is a reference reward and  $\pi_0$  is a fixed continuation rule. The anchor is not learned  
 113 separately for each regime. It is fixed before Stage II, so the regime layer cannot explain heterogeneity  
 114 by changing the underlying value function.

115 This separation is central. Stage I answers: given the observed history, what is the inferred system  
 116 state, and how do actions compare under a shared value anchor? Stage II answers: how do observed  
 117 actions systematically depart from that anchor?

### 118 2.3 Stage II: switching latent constraint regimes

119 Stage II models structured behavioral heterogeneity through a latent switching regime layer. A  
 120 *constraint regime* is a persistent latent mode under which certain actions become systematically  
 121 more or less likely relative to the shared anchor value function  $Q_{\text{anc}}$ . The term “constraint” is used  
 122 broadly: regimes may reflect unobserved protocol pressure, treatment complexity, operational burden,  
 123 workload, or other persistent factors that influence action selection without changing the anchor itself.  
 124 We introduce a latent regime

$$u_m \in \{1, \dots, K_u\},$$

125 which evolves along the decision sequence according to

$$\Pr(u_m = v \mid u_{m-1} = u) = \Pi_{uv}.$$

126 The Markov structure encourages regimes to explain coherent behavioral segments rather than isolated  
 127 action noise. Each regime defines a *deviation function*

$$C_u(\xi, a),$$

128 which modulates action selection relative to the shared anchor. Here,  $\xi_m$  is the introduced *additional*  
 129 *context* for explaining behavioral variation, such as protocol indicators, workload proxies, or action-  
 130 history features, which can contain latent decision state  $x_m$ .

131 Actions are generated according to

$$p(a_m = a \mid x_m, \xi_m, u_m = u) = \frac{\exp\{\beta[Q_{\text{anc}}(x_m, a) - C_u(\xi_m, a)]\}}{\sum_{a' \in \mathcal{A}} \exp\{\beta[Q_{\text{anc}}(x_m, a') - C_u(\xi_m, a')]\}}, \quad (1)$$

132 where  $\beta > 0$  is an inverse-temperature parameter. Intuitively, the constraint term shifts action  
 133 preferences relative to the shared anchor: some actions become systematically suppressed, while  
 134 others become more likely under particular regimes. For example, one regime may discourage  
 135 complex multi-drug changes, while another may suppress actions associated with protocol burden  
 136 or resource-intensive interventions. The deviation function  $C_u$  can take flexible forms. In the main  
 137 interpretable specification, we use a sparse linear parameterization  $C_u(\xi, a) = w_u^\top g_\psi(\xi, a)$ , where  
 138  $g_\psi : \Xi \times \mathcal{A} \rightarrow \mathbb{R}^d$  is a feature map and  $w_u$  is sparse. Sparsity makes each regime interpretable as a  
 139 small set of active behavioral deviations.

140 This parameterization (1) admits a simple decision-theoretic interpretation. For fixed  $(x, \xi, u)$ , the  
 141 softmax distribution is the unique solution to

$$\max_{\pi \in \Delta(\mathcal{A})} \sum_{a \in \mathcal{A}} \pi(a) [Q_{\text{anc}}(x, a) - C_u(\xi, a)] + \frac{1}{\beta} H(\pi),$$

142 where  $H(\pi) = -\sum_a \pi(a) \log \pi(a)$ . Thus, regimes act as structured soft constraints that reshape  
 143 action probabilities around a shared anchor value function.

144 This distinguishes our approach from mixture-policy or latent-task models. Those methods explain  
 145 behavioral heterogeneity through separate rewards, policies, or transition structures across latent  
 146 groups. In contrast, all regimes in our framework share the same latent dynamics, state representation,  
 147 and anchor value function; only the deviation term  $C_u$  changes across regimes.

148 After fitting, Stage II produces

$$\gamma_m^i(u) = \Pr(u_m^i = u \mid \mathcal{D}^i), \quad \Pi, \quad \{C(\cdot; u)\}_{u=1}^{K_u}.$$

## 149 3 Learning and Model Analysis

150 Learning follows the decomposition above. Stage I estimates latent system dynamics and constructs  
 151  $(x_m, Q_{\text{anc}})$ . Stage II learns the switching regime layer conditional on these fixed anchor objects.  
 152 After fitting, the model can also be used for anchor-based diagnostics, but these diagnostics are  
 153 separate from the estimator.

### 154 3.1 Stage I: fitting latent dynamics and computing anchor values

155 Stage I estimates the continuous-time latent dynamics from the observation sequence, treating actions  
 156 as inputs to the system dynamics. Let  $\Theta_{\text{dyn}}$  denote the parameters of the CTMC generators, outcome-  
 157 proxy dynamics, emission model, and initial latent-state distribution. We fit  $\Theta_{\text{dyn}}$  by maximum

158 likelihood:

$$\hat{\Theta}_{\text{dyn}} \in \arg \max_{\Theta_{\text{dyn}}} \sum_{i=1}^N \log p_{\Theta_{\text{dyn}}} (o_{1:M_i}^i | a_{1:M_i-1}^i, t_{1:M_i}^i).$$

159 This objective models the non-choice observations. It is not an action-prediction objective; action  
160 prediction is handled by Stage II.

161 In practice, we use an EM-style routine. The E-step filters and smooths the latent trajectory  $z(t)$ . The  
162 M-step updates the transition generators, outcome dynamics, and emission parameters. After fitting,  
163 we compute

$$\hat{x}_m^i = (\hat{b}_m^i, \hat{\lambda}_m^i)$$

164 at every decision epoch. We then evaluate

$$\hat{Q}_{\text{anc}}(\hat{x}_m^i, a), \quad a \in \mathcal{A},$$

165 using dynamic programming, numerical integration, or Monte Carlo rollouts under the fitted dynamics  
166 and fixed continuation rule  $\pi_0$ .

167 The output of Stage I is the table

$$\left\{ \hat{x}_m^i, \hat{Q}_{\text{anc}}(\hat{x}_m^i, a) : i = 1, \dots, N, m = 1, \dots, M_i, a \in \mathcal{A} \right\}.$$

168 These quantities are held fixed in Stage II, ensuring that learned regimes are deviations from a  
169 common anchor rather than alternative value functions.

### 170 3.2 Stage II: fitting switching constraint regimes

171 Stage II estimates the latent regime process and the regime-dependent deviations. For the sparse  
172 linear specification,

$$C_u(\xi, a) = w_u^\top g_\psi(\xi, a), \quad (2)$$

173 the regime parameters are

$$\Theta_{\text{reg}} = \{\mu, \Pi, w_1, \dots, w_{K_u}, \psi\},$$

174 where  $\mu$  is the initial regime distribution,  $\Pi$  is the transition matrix,  $w_u$  are regime-specific deviation  
175 weights, and  $\psi$  parameterizes the feature map.

176 Conditional on Stage I outputs, the action likelihood is an HMM with softmax emissions:

$$\mathcal{L}(\Theta_{\text{reg}}) = \sum_{i=1}^N \log \sum_{u_{1:M_i}^i} \mu_{u_1^i} \prod_{m=2}^{M_i} \Pi_{u_{m-1}^i, u_m^i} \prod_{m=1}^{M_i} p_{\Theta_{\text{reg}}}(a_m^i | \hat{x}_m^i, \hat{\xi}_m^i, u_m^i),$$

177 where  $p_{\Theta_{\text{reg}}}(a_m | \hat{x}_m, \hat{\xi}_m, u)$  is evaluated as (1).

178 We maximize the sparsity-regularized objective using generalized EM. In the E-step, forward-  
179 backward computes posterior regime probabilities

$$\gamma_m^i(u) = \Pr(u_m^i = u | \mathcal{D}^i)$$

180 and expected transition counts. In the M-step,  $\mu$  and  $\Pi$  are updated from normalized expected counts.

181 The deviation parameters are updated by solving

$$\max_{\{w_u\}, \psi} \sum_{i,m} \sum_{u=1}^{K_u} \gamma_m^i(u) \log p_{\Theta_{\text{reg}}}(a_m^i | \hat{x}_m^i, \hat{\xi}_m^i, u) - \zeta \sum_{u=1}^{K_u} \|w_u\|_1.$$

182 For hand-crafted feature  $g_\psi$ , this is a posterior-weighted sparse multinomial logit problem. If  $g_\psi$   
183 contains learned components (e.g., simple neural networks),  $\psi$  is updated by gradient steps inside the  
184 generalized M-step, with objective improvement checked after each update.

185 The single-regime  $K_u = 1$  is an important ablation. It allows one global deviation around the anchor  
186 but does not permit switching behavioral modes. Comparing  $K_u = 1$  with  $K_u > 1$  tests whether the  
187 data contain persistent heterogeneity that cannot be represented by a single deviation profile.

### 188 3.3 What the theory guarantees

189 We summarize the main theoretical consequences because they clarify what can be interpreted from  
190 the fitted model. Formal assumptions and proofs are provided in the appendix.

191 **Identifiable action-level deviations.** Condition on the anchor  $Q_{\text{anc}}$  and feature map  $g_\psi$ . Under  
192 standard finite-state HMM identifiability conditions, the regime transition law and emissions are

193 identifiable up to permutation of regime labels. Within each regime, the identifiable quantities are  
 194 action-dependent deviation differences. This follows from the log-odds identity

$$\log \frac{p(a \mid x, \xi, u)}{p(a' \mid x, \xi, u)} = \beta [Q_{\text{anc}}(x, a) - Q_{\text{anc}}(x, a') - \{C_u(\xi, a) - C_u(\xi, a')\}].$$

195 Since  $Q_{\text{anc}}$  is fixed, the data identify how each regime shifts action scores relative to the anchor, up  
 196 to action-independent offsets.

197 **Sparse recovery of regime profiles.** For the sparse linear specification  $C_u(\xi, a) = w_u^\top g_\psi(\xi, a)$ ,  
 198 suppose the true regime profiles are  $s$ -sparse and the posterior-weighted logit objective satisfies  
 199 standard restricted strong convexity and score concentration conditions. Then, with an appropriate  
 200 choice of  $\zeta$ ,

$$\sum_{u=1}^{K_u} \|\widehat{w}_u - w_u^*\|_2 \lesssim K_u \sqrt{\frac{s \log d}{NM}},$$

201 where  $d$  is the feature dimension and  $\bar{M} = N^{-1} \sum_i M_i$ . This supports reporting compact feature-  
 202 level regime profiles rather than treating regimes as uninterpreted latent clusters.

203 **Stability to Stage I approximation error.** Stage II uses plug-in estimates of the state and anchor  
 204 values. Let

$$\epsilon_Q = \sup_{m,a} \left| \widehat{Q}_{\text{anc}}(\widehat{x}_m, a) - Q_{\text{anc}}(x_m, a) \right|, \quad \epsilon_x = \sup_m \|\widehat{x}_m - x_m\|_1.$$

205 If  $Q_{\text{anc}}$  and  $g_\psi$  are Lipschitz, and  $\|w_u\|_1 \leq W$ , then for fixed regime parameters,

$$\sup_{m,u} \left| \log \widehat{p}(a_m \mid \widehat{x}_m, \widehat{\xi}_m, u) - \log p(a_m \mid x_m, \xi_m, u) \right| \leq 2\beta [\epsilon_Q + (L_Q + WL_g)\epsilon_x].$$

206 Thus, moderate Stage I error induces a controlled perturbation of the Stage II likelihood.

207 **Why switching regimes can improve fit.** The  $K_u = 1$  model is nested inside the  $K_u > 1$  model.  
 208 If the true conditional action law contains separated persistent regimes that are not observationally  
 209 equivalent to a single deviation profile, then the population likelihood of the switching model is  
 210 strictly larger than that of the single-regime model by a positive KL gap. This justifies the empirical  
 211 comparison between one shared deviation profile and multiple latent constraint regimes.

### 212 3.4 Anchor-based diagnostics after fitting

213 After Stages I and II are fitted, the model can be used for descriptive diagnostics. These diagnostics  
 214 are not part of training.

215 One diagnostic is the anchor-optimal action

$$a_{\text{anc}}^*(x_m) \in \arg \max_{a \in \mathcal{A}} Q_{\text{anc}}(x_m, a),$$

216 which provides a reference point for comparing observed actions with the shared anchor. Another  
 217 diagnostic modifies selected components of  $C_u$ , such as setting a feature group to zero or scaling a  
 218 subset of deviation weights, and recomputes the implied action probabilities. These analyses describe  
 219 how learned regimes affect action selection under the fitted model.

220 All such diagnostics are conditional on the learned dynamics, anchor reward, and feature map. They  
 221 are useful for summarizing regimes and inspecting trajectories.

## 222 4 Experiments

223 We evaluate our model on two clinical settings with markedly different timing scales: high-frequency  
 224 ICU sepsis treatment from MIMIC-IV [21, 22] and long-horizon chronic care from a private CKD-  
 225 MBD cohort; we further use retail pricing from Dominick’s Finer Foods (DFF) [23, 24] as a  
 226 cross-domain validation. Detailed data summaries are provided in Appendix D. We ask whether  
 227 the constraint-regime layer (i) recovers interpretable, persistent deviations from a shared anchor,  
 228 (ii) improves held-out action prediction over stationary and interpretable behavioral baselines while  
 229 remaining robust to reward specifications, and (iii) supports counterfactual anchor-optimal actions  
 230 under constraint relaxation and case-level auditing.

### 231 4.1 Interpretable Phenotypes and Constraint Regimes

232 We examine interpretability in both latent layers: outcome-grounded phenotypes and sparse anchor-  
 233 relative wedges.

234 **Fitted dynamics layer and recovered phenotypes.** The fitted continuous-time dynamics layer  
 235 recovers  $K_z$  outcome-grounded phenotypes, providing an interpretable view of patient heterogeneity  
 236 through phenotype signatures, dwell times, and population trajectories, with  $\lambda(t)$  dynamics and  
 237 detailed experimental results reported in Appendices C and F. The proxy  $\lambda(t)$  tracks Ca/P/PTH for  
 238 CKD-MBD, hemodynamic/metabolic indicators for MIMIC-IV, and market-state indicators for DFF.  
 239 For example, in CKD-MBD, the recovered phenotypes show clinically coherent Ca/P/PTH signatures,  
 240 separate stable from transient high-risk states through distinct dwell times, and trend toward more  
 241 stable mineral-control profiles over the treatment horizon.

242 **Regime layer and domain-grounded feature maps.** We fit the regime layer in each domain  
 243 using the sparse wedge  $C_u(\xi, a) = w_u^\top g_\psi(\xi, a)$ , where  $\xi$  may include  $x = (b, \lambda)$  and optional  
 244 side information. The feature map  $g_\psi$  is shared across regimes, while only  $w_u$  is regime-specific;  
 245 because  $a$ -independent features cancel in the choice rule (1), we construct  $g_\psi$  from action-side and  
 246 context-action features.

247 The instantiated  $g_\psi$  is 16-dim for MIMIC-IV (e.g., Sepsis-3 compliance, regimen complexity, and  
 248 therapeutic-modality intensity) [25], 8-dim for CKD-MBD (e.g., mineral-imbalance safety synergies,  
 249 polypharmacy burden, and treatment-escalation inertia) [26], and 8-dim for DFF (e.g., menu-cost inertia,  
 250 competitive-response gaps, and brand-equity guardrails) [27]; full definitions are in Appendix E.

251 **Recovering interpretable regimes.** For MIMIC-IV with  $K_u = 3$ , Figure 2 shows sparse and  
 252 clinically interpretable regime profiles, with elevated transition diagonals ( $\Pi_{uu} \in [0.43, 0.62]$ )  
 253 indicating persistent regimes rather than per-step noise:

- 254 • Regime 1, acute escalation: high weights on added\_count and n\_groups\_active, indicating  
 255 decisions in which clinicians introduce new treatments and activate multiple therapeutic groups.  
 256 We interpret this as an acute stabilization mode.
- 257 • Regime 0, protocol adherence: concentration on forbid\_violate, which encodes contraindication  
 258 or safety-rule violations. This regime captures protocol- and safety-shaped action selection.
- 259 • Regime 2, complexity management: weights on drug\_count and removed\_count, marking  
 260 decisions involving regimen size and treatment removal. We interpret this as a complexity-  
 261 management or de-escalation mode.

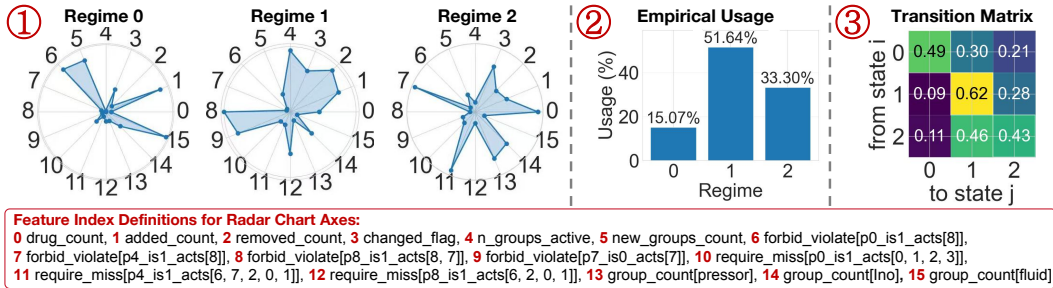


Figure 2: Proposed model on MIMIC-IV ( $K_u = 3$ ). (1) Per-regime weights  $w_u$  on the 16-dimensional feature map. (2) Posterior regime usage  $\sum_m \gamma_m(u)$ . (3) Inter-regime transition matrix  $\Pi$ .

262 **Coupled and decoupled regime–state dynamics.** Trajectory-level diagnostics (Figure 3) show  
 263 three qualitative patterns, motivating the model’s separation of decision-clock regimes from  
 264 continuous-time dynamics.

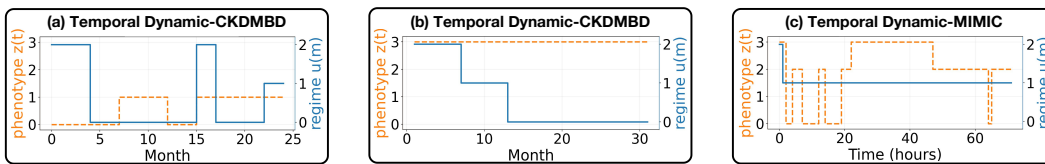


Figure 3: Representative test-set trajectories of latent state  $z(t)$  and regime path  $u_m$  ( $K_z = 4$ ,  $K_u = 3$ ). (a) A regime transition co-occurs with a  $z(t)$  shift. (b) A regime transition occurs under stable  $z(t)$ . (c) The regime remains stable despite frequent  $z(t)$  shifts.

265 **4.2 Predictive Performance and Robustness**

266 Having established interpretability, we next evaluate held-out predictive performance and robust-  
 267 ness along four axes: **interpretable behavioral competitors**, **switching-regime depth**, **reward-**  
 268 **specification robustness**, and **continuous-time dynamics validation**.

269 **Interpretable behavioral baselines.** We compare our model with Maximum Entropy IRL (MaxEnt  
 270 IRL) [10], which fits a stationary reward representation, and Soft Decision Trees (SDT) [28], which  
 271 provide discriminative hierarchical decision rules. As shown in Table 1, Our model achieves the  
 272 best held-out log-likelihood and Acc@5 in all three domains. The baselines capture different parts  
 273 of the problem but not the full decomposition: MaxEnt IRL cannot represent persistent regime-  
 274 dependent wedges around a shared anchor, while SDT does not provide a generative regime process  
 275 or anchor-optimal counterfactual actions.

Table 1: Held-out action-prediction performance across clinical and economic domains.

Model	CKD-MBD	MIMIC-IV	DFF
	Log-L / Acc@5 (%)	Log-L / Acc@5 (%)	Log-L / Acc@5 (%)
Ours*	-4.52 ± 0.85 / 50.5 ± 5.3	-0.74 ± 0.02 / 96.8 ± 0.2	-0.64 ± 0.02 / 98.5 ± 0.1
MaxEnt IRL†	-5.64 ± 0.00 / 29.3 ± 0.7	-5.29 ± 0.07 / 80.1 ± 0.7	-3.22 ± 0.00 / 70.8 ± 2.7
SDT‡	-5.23 ± 0.25 / 15.4 ± 4.2	-1.84 ± 0.00 / 93.4 ± 0.0	-2.35 ± 0.06 / 83.8 ± 1.0

\* Anchor-optimal action ( $a^*$ ), constraint regimes ( $u$ ), and latent phenotypes.

† Stationary state-action reward  $r(s, a)$  on an empirical finite-horizon MDP.

‡ Hierarchical discriminative decision paths.

276 **Predictive gain from switching regimes.** Figure 4 reports held-out performance as  $K_u$  varies  
 277 from 1 (stationary wedge baseline) to 6 on CKD-MBD, MIMIC-IV, and DFF. Allowing  $K_u > 1$   
 278 consistently improves log-likelihood and top- $k$  accuracy, with gains saturating around  $K_u = 3-5$ ,  
 279 indicating predictive structure not captured by a single stationary deviation profile.

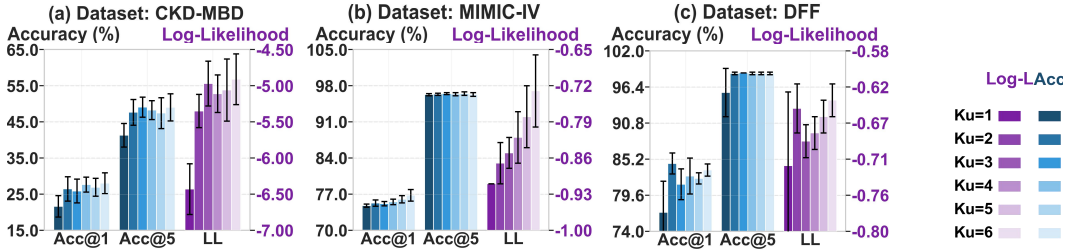


Figure 4: Held-out action prediction across  $K_u \in \{1, \dots, 6\}$  for the three datasets.  $K_u = 1$  is the stationary wedge baseline;  $K_u > 1$  allows switching constraint regimes.

280 **Reward-specification robustness.** We next evaluate whether the predictive advantage depends  
 281 on a particular reference reward. We refit our model under three reward families: outcome-based,  
 282 stability-based, and risk-sensitive. Across these specifications, switching-regime models ( $K_u > 1$ )  
 283 consistently improve held-out likelihood over the stationary  $K_u = 1$  baseline, indicating that the  
 284 gains are not artifacts of a single reward design. Reward definitions and metrics are in Appendix G.

285 **Dynamics-layer validation.** Finally, we evaluate the fitted continuous-time dynamics layer, which  
 286 supplies the anchor values  $Q_{\text{anc}}$  used by the regime model and counterfactual diagnostics; errors  
 287 in this layer would propagate directly to all anchor-relative analyses. Against ODE-RNN [19] and  
 288 NCDSSM [20], the CTMC-ODE layer remains competitive on hold-out prediction while retaining  
 289 interpretable latent phenotypes; detailed dynamics benchmarks are provided in Appendix H.

290 **4.3 Counterfactual Diagnostics and Decision Auditing**

291 After the two estimation stages, a diagnostic layer turns the value anchor and constraint-regime layers  
 292 into counterfactual diagnostics at two levels. At the **population level**, *constraint relaxation* scales  
 293 down selected wedge groups and reads off the change in  $Q_{\text{anc}}$ , ranking which regime-dependent

294 deviations carry the largest recoverable value. At the **unit level**, we synthesize **anchor-optimal**  
 295 **counterfactual actions**  $\{a_m^*\}$ , i.e., the ideal per-trajectory actions under the shared anchor, via  
 296 an SCM-grounded abduction-action-prediction rollout; decision auditing then interprets each gap  
 297  $a_m - a_m^*$  through the active regime  $u_m$  and latent state  $z_m$ . Together these diagnostics support policy  
 298 refinement, stress-testing, and case-level review, conditional on the fitted shared anchor.

299 **Population-level constraint relaxation.** Constraint relaxation produces an anchor-relative ranking  
 300 of intervention targets, indicating which wedge groups would yield the largest recoverable gain in  
 301  $Q_{\text{anc}}$  if reduced. On MIMIC-IV, we partition the 16-dimensional wedge feature map into three  
 302 clinical groups: *Compliance* (protocol violations and required-treatment omissions), *Complexity*  
 303 (regimen size and step-to-step treatment changes), and *Resource* (high-intensity therapeutic breadth).  
 304 For each group, we scale weights by  $(1 - \alpha)$ ,  $\alpha \in [0, 1]$ , and report the relative increase in  $Q_{\text{anc}}$ .  
 305 Figure 5 shows that *Complexity* dominates: full relaxation ( $\alpha = 1$ ) recovers 13.5% of  $Q_{\text{anc}}$ , versus  
 306 6.6% for *Compliance* and 7.0% for *Resource*, identifying treatment complexity and step-to-step  
 307 adjustment as the largest recoverable deviation in MIMIC-IV.

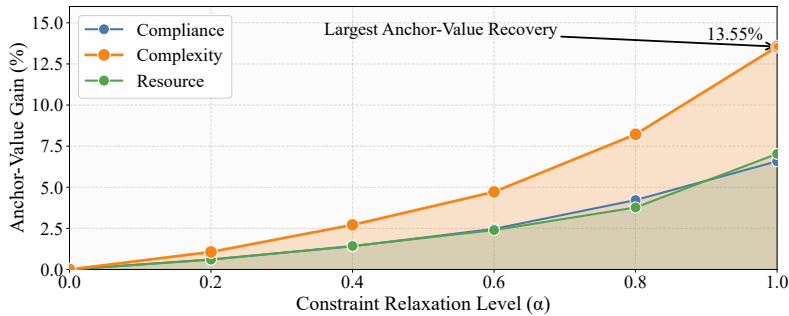


Figure 5: Anchor-value recovery under group-wise constraint relaxation for MIMIC-IV. Curves report the relative increase in  $Q_{\text{anc}}$  as  $\alpha$  increases.

308 **anchor-optimal counterfactual synthesis.** A core capability of our model, beyond imitation-  
 309 based [29] or regression-based [30] recommenders, is the synthesis of *anchor-optimal counterfactual*  
 310 *action sequences*  $\{a_m^*\}$  via an **SCM-grounded abduction-action-prediction cycle** on each observed  
 311 trajectory. The cycle abducts the exogenous noises of  $(z(t), \lambda(t))$  through the fitted action-modulated  
 312 CTMC and outcome ODE, forward-rolls candidate actions with these noises held fixed, and selects  
 313 the  $Q_{\text{anc}}$ -maximizing action as  $a_m^*$ . Full derivation and sampling are in Appendix I. The resulting  
 314  $\{a_m^*\}$  are unit-level consistent: each counterfactual respects the individual’s observed history rather  
 315 than averaging across the population. The gap  $a_m - a_m^*$  supplies the anchor-relative quantities used  
 316 in the unit-level auditing below.

317 **Unit-level auditing.** While population-level relaxation identifies which wedge groups are most  
 318 influential overall, for individual trajectories we pair the synthesized  $\{a_m^*\}$  with an LLM-based  
 319 reporting layer. The audit input consists of five structured quantities: observations  $o_m$ , observed  
 320 actions  $a_m$ , anchor-optimal actions  $a_m^*$ , inferred regimes  $u_m$ , and latent states  $z_m$ . The LLM  
 321 receives these structured model outputs and produces a three-part audit summarizing (i) whether  
 322  $a_m^*$  is domain-plausible and operationally feasible, (ii) which constraint-regime wedges explain the  
 323 deviation between  $a_m$  and  $a_m^*$ , and (iii) whether regime and state transitions align with recorded  
 324 trajectory events. The LLM is used only as a reporting layer. Appendix J presents a representative  
 325 CKD-MBD trajectory audit, and Appendix K provides the full prompt templates.

## 326 5 Conclusion

327 We introduced a shared-anchor regime-decomposition framework for modeling nonstationary sequen-  
 328 tial decisions through a common action-value anchor and switching constraint-regime deviations.  
 329 This separation represents persistent behavioral variation as sparse, interpretable wedges rather than  
 330 arbitrary policy drift, enabling posterior regime segmentation, constraint-relaxation diagnostics, and  
 331 SCM-grounded anchor-optimal counterfactuals. Across clinical care and retail pricing, the model  
 332 improves held-out action prediction over stationary and interpretable behavioral baselines while  
 333 retaining an auditable counterfactual interpretation.

334 **References**

- 335 [1] RICHARD BELLMAN. A markovian decision process. *Journal of Mathematics and Mechanics*,  
336 6(5):679–684, 1957. ISSN 00959057, 19435274. URL [http://www.jstor.org/stable/](http://www.jstor.org/stable/24900506)  
337 24900506.
- 338 [2] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning:  
339 Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 340 [3] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for  
341 deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 342 [4] Matthieu Komorowski, Leo A Celi, Omar Badawi, Anthony C Gordon, and A Aldo Faisal. The  
343 artificial intelligence clinician learns optimal treatment strategies for sepsis in intensive care.  
344 *Nature medicine*, 24(11):1716–1720, 2018.
- 345 [5] Gerhard Widmer and Miroslav Kubat. Learning in the presence of concept drift and hidden  
346 contexts. *Machine learning*, 23(1):69–101, 1996.
- 347 [6] Kenneth Jung and Nigam H. Shah. Implications of non-stationarity on predictive modeling  
348 using ehers. *Journal of Biomedical Informatics*, 58:168–174, 2015. ISSN 1532-0464. doi: <https://doi.org/10.1016/j.jbi.2015.10.006>. URL [https://www.sciencedirect.com/science/](https://www.sciencedirect.com/science/article/pii/S1532046415002282)  
349 [article/pii/S1532046415002282](https://www.sciencedirect.com/science/article/pii/S1532046415002282).  
350
- 351 [7] Wang Chi Cheung, David Simchi-Levi, and Ruihao Zhu. Reinforcement learning for non-  
352 stationary markov decision processes: The blessing of (more) optimism. In *International*  
353 *conference on machine learning*, pages 1843–1854. PMLR, 2020.
- 354 [8] Erwan Lecarpentier and Emmanuel Rachelson. Non-stationary markov decision processes, a  
355 worst-case approach using model-based reinforcement learning. *Advances in neural information*  
356 *processing systems*, 32, 2019.
- 357 [9] Andrew Y Ng, Stuart Russell, et al. Algorithms for inverse reinforcement learning. In *Icml*,  
358 volume 1, page 2, 2000.
- 359 [10] Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, Anind K Dey, et al. Maximum entropy  
360 inverse reinforcement learning. In *Aaai*, volume 8, pages 1433–1438. Chicago, IL, USA, 2008.
- 361 [11] Michael Bloem and Nicholas Bambos. Ground delay program analytics with behavioral cloning  
362 and inverse reinforcement learning. *Journal of aerospace information systems*, 12(3):299–313,  
363 2015.
- 364 [12] Lantao Yu, Tianhe Yu, Chelsea Finn, and Stefano Ermon. Meta-inverse reinforcement learning  
365 with probabilistic context variables. *Advances in neural information processing systems*, 32,  
366 2019.
- 367 [13] Elynn Y Chen, Rui Song, and Michael I Jordan. Reinforcement learning in latent heterogeneous  
368 environments. *Journal of the American Statistical Association*, 119(548):3113–3126, 2024.
- 369 [14] Dylan Hadfield-Menell, Stuart J Russell, Pieter Abbeel, and Anca Dragan. Cooperative inverse  
370 reinforcement learning. *Advances in neural information processing systems*, 29, 2016.
- 371 [15] Athul Paul Jacob, Abhishek Gupta, and Jacob Andreas. Modeling boundedly rational agents  
372 with latent inference budgets, 2023. URL <https://arxiv.org/abs/2312.04030>.
- 373 [16] Zoubin Ghahramani and Geoffrey E. Hinton. Variational learning for switching state-space  
374 models. *Neural Computation*, 12(4):831–864, 2000. doi: 10.1162/089976600300015619.
- 375 [17] James D Hamilton. A new approach to the economic analysis of nonstationary time series and  
376 the business cycle. *Econometrica: Journal of the econometric society*, pages 357–384, 1989.
- 377 [18] Yulia Rubanova, Ricky TQ Chen, and David K Duvenaud. Latent ordinary differential equations  
378 for irregularly-sampled time series. *Advances in neural information processing systems*, 32,  
379 2019.

- 380 [19] Mansura Habiba and Barak A Pearlmutter. Neural ordinary differential equation based recurrent  
381 neural network model. In *2020 31st Irish signals and systems conference (ISSC)*, pages 1–6.  
382 IEEE, 2020.
- 383 [20] Abdul Fatir Ansari, Alvin Heng, Andre Lim, and Harold Soh. Neural continuous-discrete state  
384 space models for irregularly-sampled time series. In *International Conference on Machine*  
385 *Learning*, pages 926–951. PMLR, 2023.
- 386 [21] A Johnson, L Bulgarelli, T Pollard, B Gow, B Moody, S Horng, LA Celi, and R Mark. Mimic-iv  
387 (version 3.1). physionet, 2024.
- 388 [22] Alistair EW Johnson, Lucas Bulgarelli, Lu Shen, Alvin Gayles, Ayad Shammout, Steven Horng,  
389 Tom J Pollard, Sicheng Hao, Benjamin Moody, Brian Gow, et al. Mimic-iv, a freely accessible  
390 electronic health record dataset. *Scientific data*, 10(1):1, 2023.
- 391 [23] Stephen J Hoch, Byung-Do Kim, Alan L Montgomery, and Peter E Rossi. Determinants of  
392 store-level price elasticity. *Journal of marketing Research*, 32(1):17–29, 1995.
- 393 [24] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An introduction to*  
394 *statistical learning: with applications in R*, volume 103. Springer, 2013.
- 395 [25] Andrew Rhodes, Laura E Evans, Waleed Alhazzani, Mitchell M Levy, Massimo Antonelli,  
396 Ricard Ferrer, Anand Kumar, Jonathan E Sevransky, Charles L Sprung, Mark E Nunnally, et al.  
397 Surviving sepsis campaign: international guidelines for management of sepsis and septic shock:  
398 2016. *Intensive care medicine*, 43(3):304–377, 2017.
- 399 [26] Markus Ketteler, Geoffrey A Block, Pieter Evenepoel, Masafumi Fukagawa, Charles A Herzog,  
400 Linda McCann, Sharon M Moe, Rukshana Shroff, Marcello A Tonelli, Nigel D Toussaint,  
401 et al. Executive summary of the 2017 kdigo chronic kidney disease–mineral and bone disorder  
402 (ckd-mbd) guideline update: what’s changed and why it matters. *Kidney international*, 92(1):  
403 26–36, 2017.
- 404 [27] Daniel Levy, Mark Bergen, Shantanu Dutta, and Robert Venable. The magnitude of menu costs:  
405 direct evidence from large us supermarket chains. *The Quarterly Journal of Economics*, 112(3):  
406 791–824, 1997.
- 407 [28] Nicholas Frosst and Geoffrey Hinton. Distilling a neural network into a soft decision tree. *arXiv*  
408 *preprint arXiv:1711.09784*, 2017.
- 409 [29] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and  
410 structured prediction to no-regret online learning. In *Proceedings of the fourteenth interna-*  
411 *tional conference on artificial intelligence and statistics*, pages 627–635. JMLR Workshop and  
412 Conference Proceedings, 2011.
- 413 [30] Susan Athey and Stefan Wager. Policy learning with observational data. *Econometrica*, 89(1):  
414 133–161, 2021.
- 415 [31] Aniruddh Raghu, Matthieu Komorowski, Leo Anthony Celi, Peter Szolovits, and Marzyeh  
416 Ghassemi. Continuous state-space models for optimal sepsis treatment: a deep reinforcement  
417 learning approach. In *Machine Learning for Healthcare Conference*, pages 147–163. PMLR,  
418 2017.
- 419 [32] Aniruddh Raghu, Matthieu Komorowski, Imran Ahmed, Leo Celi, Peter Szolovits, and Marzyeh  
420 Ghassemi. Deep reinforcement learning for sepsis treatment, 2017. URL <https://arxiv.org/abs/1711.09602>.
- 422 [33] Omer Gottesman, Fredrik Johansson, Matthieu Komorowski, Aldo Faisal, David Sontag, Finale  
423 Doshi-Velez, and Leo Anthony Celi. Guidelines for reinforcement learning in healthcare. *Nature*  
424 *medicine*, 25(1):16–18, 2019.
- 425 [34] Miroslav Dudík, John Langford, and Lihong Li. Doubly robust policy evaluation and learning.  
426 *arXiv preprint arXiv:1103.4601*, 2011.

- 427 [35] Nan Jiang and Lihong Li. Doubly robust off-policy value evaluation for reinforcement learning.  
428 In *International conference on machine learning*, pages 652–661. PMLR, 2016.
- 429 [36] Taylor W Killian, Samuel Daulton, George Konidaris, and Finale Doshi-Velez. Robust and  
430 efficient transfer learning with hidden parameter markov decision processes. *Advances in neural  
431 information processing systems*, 30, 2017.
- 432 [37] Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions  
433 and use interpretable models instead. *Nature Machine Intelligence*, 1(5):206–215, 2019. doi:  
434 10.1038/s42256-019-0048-x.
- 435 [38] Judea Pearl. *Causality*. Cambridge University Press, 2 edition, 2009.
- 436 [39] Judea Pearl. Causal inference in statistics: An overview. *Statistics Surveys*, 3(none):96 – 146,  
437 2009. doi: 10.1214/09-SS057. URL <https://doi.org/10.1214/09-SS057>.
- 438 [40] Ioana Bica, Ahmed M Alaa, James Jordon, and Mihaela Van Der Schaar. Estimating counter-  
439 factual treatment outcomes over time through adversarially balanced representations. *arXiv  
440 preprint arXiv:2002.04083*, 2020.
- 441 [41] Nabeel Seedat, Fergus Imrie, Alexis Bellot, Zhaozhi Qian, and Mihaela van der Schaar.  
442 Continuous-time modeling of counterfactual outcomes using neural controlled differential  
443 equations. *arXiv preprint arXiv:2206.08311*, 2022.
- 444 [42] Zining Zhu, Hanjie Chen, Xi Ye, Qing Lyu, Chenhao Tan, Ana Marasovic, and Sarah Wiegrefe.  
445 Explanation in the era of large language models. In Rui Zhang, Nathan Schneider, and Snigdha  
446 Chaturvedi, editors, *Proceedings of the 2024 Conference of the North American Chapter of the  
447 Association for Computational Linguistics: Human Language Technologies (Volume 5: Tutorial  
448 Abstracts)*, pages 19–25, Mexico City, Mexico, June 2024. Association for Computational  
449 Linguistics. doi: 10.18653/v1/2024.naacl-tutorials.3. URL [https://aclanthology.org/  
450 2024.naacl-tutorials.3/](https://aclanthology.org/2024.naacl-tutorials.3/).
- 451 [43] Alexandros Vassiliades, Nikolaos Polatidis, Stamatios Samaras, Sotiris Diplaris, Ignacio Cabrera  
452 Martin, Yannis Manolopoulos, Stefanos Vrochidis, and Ioannis Kompatsiaris. Utilizing large  
453 language models for machine learning explainability, 2025. URL [https://arxiv.org/abs/  
454 2510.06912](https://arxiv.org/abs/2510.06912).
- 455 [44] Ansgar Koene, Chris Clifton, Yohko Hatada, Helena Webb, and Rashida Richardson. A  
456 governance framework for algorithmic accountability and transparency. 2019.
- 457 [45] Mervyn Singer, Clifford S Deutschman, Christopher Warren Seymour, Manu Shankar-Hari,  
458 Djillali Annane, Michael Bauer, Rinaldo Bellomo, Gordon R Bernard, Jean-Daniel Chiche,  
459 Craig M Coopersmith, et al. The third international consensus definitions for sepsis and septic  
460 shock (sepsis-3). *Jama*, 315(8):801–810, 2016.
- 461 [46] Laura Evans, Andrew Rhodes, Waleed Alhazzani, Massimo Antonelli, Craig M Coopersmith,  
462 Craig French, Flávia R Machado, Lauralyn Mcintyre, Marlies Ostermann, Hallie C Prescott,  
463 et al. Surviving sepsis campaign: international guidelines for management of sepsis and septic  
464 shock 2021. *Critical care medicine*, 49(11):e1063–e1143, 2021.
- 465 [47] Suchi Saria. Individualized sepsis treatment using reinforcement learning. *Nature medicine*, 24  
466 (11):1641–1642, 2018.
- 467 [48] Sharon M Moe, Tilman Drüeke, Norbert Lameire, and Garabed Eknoyan. Chronic kidney  
468 disease–mineral-bone disorder: a new paradigm. *Advances in chronic kidney disease*, 14(1):  
469 3–12, 2007.
- 470 [49] Adam E Gaweda, Eleanor D Lederer, and Michael E Brier. Artificial intelligence–guided  
471 precision treatment of chronic kidney disease–mineral bone disorder. *CPT: Pharmacometrics &  
472 Systems Pharmacology*, 11(10):1305–1315, 2022.
- 473 [50] Wei Chen and David A Bushinsky. Kdigo ckd–mbd guideline update: evolution in the face of  
474 uncertainty. *Nature Reviews Nephrology*, 13(10):600–602, 2017.

475 **A Related Work**

476 Our framework is closely related to offline reinforcement learning for high-stakes sequential decision-  
 477 making, where policies must be learned from logged behavior with limited exploration [2, 3]. In  
 478 healthcare, MDP and deep RL models have been used to derive treatment policies from MIMIC sepsis  
 479 trajectories [4, 31, 32], but subsequent work emphasizes that clinician actions are heterogeneous,  
 480 confounded, and not necessarily an optimal gold standard [33]. Off-policy evaluation provides tools  
 481 for estimating policy value from observational logs [34, 35], and hidden-parameter MDPs model  
 482 patient-level variation during transfer [36]. Our framework differs in goal: rather than learning a  
 483 single deployable treatment policy, it treats the behavior policy itself as data to be decomposed into a  
 484 shared anchor and switching, anchor-relative deviations.

485 Our regime layer also connects to inverse reinforcement learning and heterogeneous behavior model-  
 486 ing. Classical IRL and maximum-entropy IRL infer rewards that rationalize demonstrations [9, 10],  
 487 while later work handles diverse demonstrations through latent contexts, mixtures, or heterogeneous  
 488 environments [11, 12, 13]. Cooperative IRL and bounded-rational-agent models further recognize  
 489 that observed behavior can reflect uncertainty, limited inference, or interaction between human and  
 490 algorithmic objectives [14, 15]. These methods usually explain heterogeneity by changing rewards,  
 491 policies, or latent task identities. Our model instead keeps a shared anchor fixed and assigns system-  
 492 atic behavioral variation to a sparse switching deviation layer. This makes the learned regimes closer  
 493 to auditable decision pressures than to opaque policy clusters, aligning with calls for interpretable  
 494 models in high-stakes settings [37, 28].

495 Finally, Our framework builds on work on nonstationarity, regime switching, continuous-time latent  
 496 dynamics, and counterfactual evaluation. Concept drift and nonstationary MDPs show why fixed  
 497 decision rules can fail when data-generating processes evolve [5, 6, 8, 7]; switching state-space  
 498 models provide a classical latent-regime mechanism for piecewise-stationary structure [17, 16]. For  
 499 irregularly sampled clinical records, latent ODEs and neural continuous-discrete state-space models  
 500 learn dynamics without forcing observations onto an artificial grid [18, 19, 20]. Counterfactual  
 501 time-series methods and structural causal models provide semantics for asking how outcomes or  
 502 values would change under alternative interventions [38, 39, 40, 41]. The proposed framework  
 503 couples these threads by placing the anchor on a continuous-time latent state while letting a discrete  
 504 regime process govern deviations on the decision clock; the resulting anchor-relative counterfactuals  
 505 can then be summarized by an LLM-based auditor for case-level explanation [42, 43, 44].

506 **B Theoretical Properties and Proofs**

507 This appendix formalizes the theoretical claims used in Section 3. The analysis focuses on the  
 508 Stage-II regime layer conditional on the Stage-I outputs. Throughout, the decision-time state is  
 509  $x_m = (b_m, \lambda_m)$ , the regime-layer context  $\xi_m$  collects the variables available to the regime layer  
 510 (e.g.,  $x_m$  itself and any optional side information), and the shared anchor value is  $Q_{\text{anc}}(x_m, a)$ . The  
 511 regime-dependent deviation is denoted by  $C_u(\xi_m, a)$ , with sparse linear specification  $C_u(\xi, a) =$   
 512  $w_u^\top g_\psi(\xi, a)$ .

513 We treat  $Q_{\text{anc}}$ ,  $\beta$ , and  $g_\psi$  as fixed when proving identifiability and sparse recovery. If  $\beta$  is estimated  
 514 jointly, a scale normalization is required; otherwise the product  $\beta[Q_{\text{anc}} - C_u]$  is identifiable but its  
 515 scale is not separately determined.

516 **B.1 Regime-layer setup**

517 Let  $\mathcal{A}$  be a finite action set with  $|\mathcal{A}| \geq 2$ , and let  $K = K_u$  be the number of latent constraint regimes.  
 518 A regime  $u_m \in \{1, \dots, K\}$  evolves as a Markov chain with initial distribution  $\mu$  and transition  
 519 matrix  $\Pi$ . Conditional on  $x_m, \xi_m$ , and  $u_m = u$ , the action emission model is

$$e_u(a \mid x, \xi) := p_\Theta(a \mid x, \xi, u) = \frac{\exp\{\beta[Q_{\text{anc}}(x, a) - C_u(\xi, a)]\}}{\sum_{a' \in \mathcal{A}} \exp\{\beta[Q_{\text{anc}}(x, a') - C_u(\xi, a')]\}}. \quad (3)$$

520 For a trajectory of length  $M$ , the conditional HMM likelihood is

$$p_\Theta(a_{1:M} \mid x_{1:M}, \xi_{1:M}) = \sum_{u_{1:M}} \mu_{u_1} \prod_{m=2}^M \Pi_{u_{m-1}, u_m} \prod_{m=1}^M e_{u_m}(a_m \mid x_m, \xi_m). \quad (4)$$

521 In the sparse linear model,  $C_u(\xi, a) = w_u^\top g_\psi(\xi, a)$ , where  $g_\psi : \Xi \times \mathcal{A} \rightarrow \mathbb{R}^d$  is a fixed feature map  
 522 and  $w_u \in \mathbb{R}^d$ .

523 Two parameter values are considered equivalent if they induce the same conditional distribution of  
 524  $a_{1:M}$  given  $x_{1:M}, \xi_{1:M}$ , up to a permutation of regime labels and action-independent offsets in the

525 deviations. Such offsets are unavoidable: replacing  $C_u(\xi, a)$  by  $C_u(\xi, a) + \kappa_u(\xi)$ , where  $\kappa_u$  does  
 526 not depend on  $a$ , does not change the softmax probabilities.

## 527 B.2 Entropy-regularized interpretation

528 **Proposition 1** (Entropy-regularized anchor-minus-deviation). *Fix  $(x, \xi, u)$  and define*

$$S_u(x, \xi, a) = Q_{\text{anc}}(x, a) - C_u(\xi, a).$$

529 *Then the softmax distribution in (3) is the unique maximizer of*

$$\max_{\pi \in \Delta(\mathcal{A})} \left\{ \sum_{a \in \mathcal{A}} \pi(a) S_u(x, \xi, a) + \frac{1}{\beta} H(\pi) \right\}, \quad H(\pi) = - \sum_{a \in \mathcal{A}} \pi(a) \log \pi(a). \quad (5)$$

530 *Proof.* The objective in (5) is strictly concave on the simplex because the entropy term is strictly  
 531 concave and  $\beta > 0$ . Form the Lagrangian

$$\mathcal{L}(\pi, \nu) = \sum_a \pi(a) S_u(x, \xi, a) - \frac{1}{\beta} \sum_a \pi(a) \log \pi(a) + \nu \left( \sum_a \pi(a) - 1 \right).$$

532 At the interior optimum,

$$\frac{\partial \mathcal{L}}{\partial \pi(a)} = S_u(x, \xi, a) - \frac{1}{\beta} \{1 + \log \pi(a)\} + \nu = 0.$$

533 Thus  $\log \pi(a) = \beta [S_u(x, \xi, a) + \nu] - 1$ , so  $\pi(a) \propto \exp\{\beta S_u(x, \xi, a)\}$ . Normalizing over  $\mathcal{A}$  gives  
 534 (3). Strict concavity gives uniqueness.  $\square$

535 This proposition justifies viewing  $C_u(\xi, a)$  as a regime-specific soft constraint or deviation term: the  
 536 regime does not replace  $Q_{\text{anc}}$ ; it modifies action probabilities around that shared anchor.

## 537 B.3 Identifiability of regime-dependent deviation contrasts

538 We next state conditions under which the learned regimes and action-level deviation contrasts are  
 539 identifiable.

540 **Assumption 1** (Conditional HMM identifiability). Conditional on the context sequence  $(x_m, \xi_m)$ , the  
 541 finite-state HMM in (4) is identifiable up to permutation of hidden-state labels. A sufficient condition  
 542 is that  $\Pi$  is full rank and ergodic, regimes induce distinct emission distributions, and the emission  
 543 vectors  $\{e_u(\cdot | x, \xi)\}_{u=1}^K$  are linearly independent on a set of contexts with positive probability.

544 **Assumption 2** (Positivity and boundedness). The parameter space is compact and the induced  
 545 emissions satisfy  $e_u(a | x, \xi) \geq \underline{p} > 0$  for all  $u, a$  and all contexts in the support.

546 **Assumption 3** (Feature-difference rank for the sparse linear model). For the sparse linear specifica-  
 547 tion, define the feature-difference set

$$\mathcal{G} = \{g_\psi(\xi, a) - g_\psi(\xi, a') : (\xi, a, a') \text{ in the support}\}.$$

548 The identifiable part of  $w_u$  is its projection onto  $\text{span}(\mathcal{G})$ . Full identification of  $w_u$  requires that the  
 549 corresponding feature-difference design has full rank on the active coordinates.

550 **Theorem 2** (Identifiability of deviation contrasts). *Under Assumptions 1–2, and conditional on  $Q_{\text{anc}}$ ,  
 551  $\beta$ , and  $g_\psi$ , the regime transition law and emissions are identifiable up to permutation of regime labels.  
 552 For each identified regime  $u$ , the action-level deviation contrasts*

$$C_u(\xi, a) - C_u(\xi, a')$$

553 *are identifiable for all  $(\xi, a, a')$  in the support, up to the same label permutation. In the sparse linear  
 554 model,  $w_u$  is identifiable modulo the null space of the feature-difference design in Assumption 3; if  
 555 the active feature-difference design has full rank, then  $w_u$  is identifiable on its active coordinates.*

556 *Proof.* By Assumption 1, the transition matrix  $\Pi$  and the emission kernels  $\{e_u(\cdot | x, \xi)\}_{u=1}^K$  are  
 557 identified up to a common permutation of regime labels.

558 Fix an identified regime  $u$  and context  $(x, \xi)$ . For any two actions  $a, a' \in \mathcal{A}$ , the log-odds under (3)  
 559 satisfy

$$\log \frac{e_u(a | x, \xi)}{e_u(a' | x, \xi)} = \beta [Q_{\text{anc}}(x, a) - Q_{\text{anc}}(x, a') - \{C_u(\xi, a) - C_u(\xi, a')\}]. \quad (6)$$

560 The left-hand side is identified from the emission distribution, and  $Q_{\text{anc}}$  and  $\beta$  are fixed. Therefore  
 561  $C_u(\xi, a) - C_u(\xi, a')$  is identified.

562 No stronger statement is possible without additional restrictions, because adding any action-  
 563 independent offset  $\kappa_u(\xi)$  to  $C_u(\xi, a)$  leaves all action probabilities unchanged. For the sparse

564 linear model, if two weight vectors  $w_u$  and  $w'_u$  induce the same emissions, then for every feature  
565 difference in the support,

$$(w_u - w'_u)^\top \{g_\psi(\xi, a) - g_\psi(\xi, a')\} = 0.$$

566 Thus  $w_u - w'_u$  lies in the null space of the feature-difference design. Full rank on the active coordinates  
567 removes this ambiguity and identifies  $w_u$  on those coordinates.  $\square$

#### 568 B.4 Stability to Stage-I plug-in error

569 The Stage-II likelihood uses plug-in estimates from Stage I. We show that moderate error in these  
570 plug-in quantities induces only a controlled perturbation of the action likelihood.

571 **Lemma 3** (Log-sum-exp Lipschitzness). *For any  $r, s \in \mathbb{R}^{|\mathcal{A}|}$ ,*

$$\left| \log \sum_{a \in \mathcal{A}} e^{r_a} - \log \sum_{a \in \mathcal{A}} e^{s_a} \right| \leq \|r - s\|_\infty.$$

572 *Proof.* Let  $m = \|r - s\|_\infty$ . Then  $s_a - m \leq r_a \leq s_a + m$  for all  $a$ . Hence  $\sum_a e^{s_a - m} \leq \sum_a e^{r_a} \leq$   
573  $\sum_a e^{s_a + m}$ . Taking logs gives the result.  $\square$

574 Let  $q_m(a) = Q_{\text{anc}}(x_m, a)$  and  $\hat{q}_m(a) = \hat{Q}_{\text{anc}}(\hat{x}_m, a)$ . Suppose

$$\epsilon_Q = \sup_{m,a} |\hat{q}_m(a) - q_m(a)|, \quad \epsilon_\xi = \sup_m \|\hat{\xi}_m - \xi_m\|_1.$$

575 **Theorem 4** (Plug-in stability of the Stage-II likelihood). *Assume  $C_u(\xi, a)$  is  $L_C$ -Lipschitz in  $\xi$ ,*  
576 *uniformly over  $u$  and  $a$ . Then for any fixed regime-layer parameters,*

$$\sup_{m,u} \left| \log \hat{e}_u(a_m | \hat{x}_m, \hat{\xi}_m) - \log e_u(a_m | x_m, \xi_m) \right| \leq 2\beta(\epsilon_Q + L_C \epsilon_\xi),$$

577 *where  $\hat{e}_u$  is the emission model evaluated with the plug-in anchor values and plug-in contexts. For the*  
578 *sparse linear model, if  $\|w_u\|_1 \leq W$  and  $g_\psi(\cdot, a)$  is  $L_g$ -Lipschitz in  $\xi$ , then one may take  $L_C \leq W L_g$ .*

579 *Proof.* Fix  $m$  and  $u$ . Define scaled utility vectors over actions:

$$r_a = \beta \{\hat{q}_m(a) - C_u(\hat{\xi}_m, a)\}, \quad s_a = \beta \{q_m(a) - C_u(\xi_m, a)\}.$$

580 Then

$$\log \hat{e}_u(a_m | \hat{x}_m, \hat{\xi}_m) - \log e_u(a_m | x_m, \xi_m) = (r_{a_m} - s_{a_m}) - \left( \log \sum_a e^{r_a} - \log \sum_a e^{s_a} \right).$$

581 By Lemma 3, the absolute value is at most  $2\|r - s\|_\infty$ . For any action  $a$ ,

$$|r_a - s_a| \leq \beta \left( |\hat{q}_m(a) - q_m(a)| + |C_u(\hat{\xi}_m, a) - C_u(\xi_m, a)| \right) \leq \beta(\epsilon_Q + L_C \epsilon_\xi).$$

582 Taking the supremum over  $a, m, u$  gives the result. In the sparse linear case,

$$|C_u(\hat{\xi}, a) - C_u(\xi, a)| = |w_u^\top \{g_\psi(\hat{\xi}, a) - g_\psi(\xi, a)\}| \leq \|w_u\|_1 L_g \|\hat{\xi} - \xi\|_1 \leq W L_g \|\hat{\xi} - \xi\|_1.$$

583  $\square$

584 If one instead separates anchor-function error from state-estimation error, the same theorem applies  
585 with  $\epsilon_Q \leq \epsilon_{\text{val}} + L_Q \epsilon_x$ , where  $\epsilon_{\text{val}} = \sup_{m,a} |\hat{Q}_{\text{anc}}(\hat{x}_m, a) - Q_{\text{anc}}(\hat{x}_m, a)|$ ,  $Q_{\text{anc}}$  is  $L_Q$ -Lipschitz  
586 in  $x$ , and  $\epsilon_x = \sup_m \|\hat{x}_m - x_m\|_1$ .

#### 587 B.5 Sparse recovery for regime profiles

588 We state a standard oracle inequality for the sparse linear Stage-II M-step. This result is conditional  
589 on fixed responsibilities  $\gamma_m^i(u)$  and a fixed feature map  $g_\psi$ . It therefore describes the weighted sparse-  
590 logit subproblem solved inside generalized EM; global finite-sample guarantees for the nonconvex  
591 HMM likelihood would require additional initialization and landscape assumptions.

592 Let  $n = \sum_{i=1}^N M_i = N \bar{M}$ . For regime  $u$ , define the normalized negative weighted log-likelihood

$$\mathcal{R}_u(w) = -\frac{1}{n} \sum_{i,m} \gamma_m^i(u) \log e_{u,w}(a_m^i | \hat{x}_m^i, \hat{\xi}_m^i),$$

593 where  $e_{u,w}$  denotes (3) with  $C_u(\xi, a) = w^\top g_\psi(\xi, a)$ . The penalized estimator is

$$\hat{w}_u \in \arg \min_w \{\mathcal{R}_u(w) + \zeta \|w\|_1\}.$$

594 **Assumption 4** (Sparse truth and bounded features). The oracle target  $w_u^*$  is  $s$ -sparse with support  
595  $S_u$ , and  $\|g_\psi(\xi, a)\|_\infty \leq G$  almost surely.

596 **Assumption 5** (Score concentration). The regularization level satisfies

$$\|\nabla \mathcal{R}_u(w_u^*)\|_\infty \leq \zeta/2$$

597 with high probability.

598 **Assumption 6** (Restricted strong convexity). For all  $\Delta$  in the cone  $\|\Delta_{S_u^c}\|_1 \leq 3\|\Delta_{S_u}\|_1$ ,

$$\mathcal{R}_u(w_u^* + \Delta) - \mathcal{R}_u(w_u^*) - \langle \nabla \mathcal{R}_u(w_u^*), \Delta \rangle \geq \frac{\kappa}{2} \|\Delta\|_2^2.$$

599 **Theorem 5** (Oracle inequality for sparse regime deviations). *Under Assumptions 4–6,*

$$\|\widehat{w}_u - w_u^*\|_2 \leq \frac{3\zeta\sqrt{s}}{\kappa}, \quad \|\widehat{w}_u - w_u^*\|_1 \leq \frac{12\zeta s}{\kappa}.$$

600 *If  $\zeta \asymp G\sqrt{\log d/n}$ , then summing over  $K$  regimes gives*

$$\sum_{u=1}^K \|\widehat{w}_u - w_u^*\|_2 \lesssim K \sqrt{\frac{s \log d}{NM}}.$$

601 *Proof.* Let  $\Delta = \widehat{w}_u - w_u^*$ . By optimality of  $\widehat{w}_u$ ,

$$\mathcal{R}_u(w_u^* + \Delta) + \zeta \|w_u^* + \Delta\|_1 \leq \mathcal{R}_u(w_u^*) + \zeta \|w_u^*\|_1.$$

602 Thus

$$\mathcal{R}_u(w_u^* + \Delta) - \mathcal{R}_u(w_u^*) \leq \zeta (\|w_u^*\|_1 - \|w_u^* + \Delta\|_1).$$

603 Using the support  $S_u$ , the right-hand side is at most

$$\zeta (\|\Delta_{S_u}\|_1 - \|\Delta_{S_u^c}\|_1).$$

604 By convexity and the score bound,

$$\mathcal{R}_u(w_u^* + \Delta) - \mathcal{R}_u(w_u^*) \geq \langle \nabla \mathcal{R}_u(w_u^*), \Delta \rangle \geq -\frac{\zeta}{2} \|\Delta\|_1.$$

605 Combining these inequalities yields

$$-\frac{\zeta}{2} (\|\Delta_{S_u}\|_1 + \|\Delta_{S_u^c}\|_1) \leq \zeta (\|\Delta_{S_u}\|_1 - \|\Delta_{S_u^c}\|_1),$$

606 which implies the cone condition  $\|\Delta_{S_u^c}\|_1 \leq 3\|\Delta_{S_u}\|_1$ .

607 Applying restricted strong convexity,

$$\frac{\kappa}{2} \|\Delta\|_2^2 \leq \mathcal{R}_u(w_u^* + \Delta) - \mathcal{R}_u(w_u^*) - \langle \nabla \mathcal{R}_u(w_u^*), \Delta \rangle.$$

608 Using the basic inequality again and the score bound gives

$$\frac{\kappa}{2} \|\Delta\|_2^2 \leq \frac{3\zeta}{2} \|\Delta_{S_u}\|_1 \leq \frac{3\zeta}{2} \sqrt{s} \|\Delta\|_2.$$

609 Therefore  $\|\Delta\|_2 \leq 3\zeta\sqrt{s}/\kappa$ . The cone condition then gives

$$\|\Delta\|_1 = \|\Delta_{S_u}\|_1 + \|\Delta_{S_u^c}\|_1 \leq 4\|\Delta_{S_u}\|_1 \leq 4\sqrt{s} \|\Delta\|_2 \leq \frac{12\zeta s}{\kappa}.$$

610 The stated rate follows from the standard choice  $\zeta \asymp G\sqrt{\log d/n}$  and  $n = NM$ .  $\square$

## 611 B.6 Consistency of the regime-layer MLE

612 Let

$$\ell_i(\Theta) = \log p_\Theta(a_{1:M_i}^i \mid x_{1:M_i}^i, \xi_{1:M_i}^i)$$

613 denote the conditional regime-layer log-likelihood contribution for trajectory  $i$ .

614 **Assumption 7** (Sampling and trajectory lengths). Trajectories  $\{(x_{1:M_i}^i, \xi_{1:M_i}^i, a_{1:M_i}^i)\}_{i=1}^N$  are inde-  
615 pendent across  $i$ . Either  $\sup_i M_i \leq M_{\max} < \infty$ , or  $M_i$  has an exponential tail and the likelihood is  
616 analyzed on events  $M_i \leq C \log N$ , whose probability tends to one.

617 **Assumption 8** (Compactness and continuity). The parameter space for  $\Theta = (\mu, \Pi, \{C_u\}_{u=1}^K)$ , or for  
618 the sparse linear parameterization  $(\mu, \Pi, \{w_u\}_{u=1}^K, \psi)$ , is compact. The map  $\Theta \mapsto \ell_i(\Theta)$  is almost  
619 surely continuous and dominated by an integrable envelope.

620 **Theorem 6** (Consistency of the regime-layer MLE). *Assume the Stage-II model is well specified,*  
621 *with true parameter  $\Theta^*$ , and Assumptions 1, 7, and 8 hold. Then any global maximizer*

$$\widehat{\Theta}_N \in \arg \max_{\Theta} \sum_{i=1}^N \ell_i(\Theta)$$

622 converges to  $\Theta^*$  in the quotient space induced by label permutation and action-independent deviation  
 623 offsets.

624 *Proof.* By Assumptions 7 and 8, a uniform law of large numbers gives

$$\sup_{\Theta} \left| \frac{1}{N} \sum_{i=1}^N \ell_i(\Theta) - \mathbb{E}[\ell_1(\Theta)] \right| \rightarrow 0$$

625 almost surely. Under well specification, the population criterion is uniquely maximized at  $\Theta^*$  modulo  
 626 the equivalence relation defined above. This uniqueness follows from Assumption 1 and Theorem 2.  
 627 Standard argmax consistency for M-estimators then implies convergence of any global maximizer in  
 628 the quotient space.  $\square$

629 **Theorem 7** (Quasi-MLE under misspecification). *Under Assumptions 7 and 8, if the true conditional*  
 630 *law is not contained in the Stage-II model class, then any global maximizer converges to the set of*  
 631 *maximizers of  $\mathbb{E}[\ell_1(\Theta)]$ , equivalently the KL projection of the true conditional law onto the model*  
 632 *class.*

633 *Proof.* The same uniform convergence argument applies. The population objective can be written  
 634 as a constant minus the KL divergence from the true conditional law to the model law. Hence its  
 635 maximizers are exactly the KL projections. Argmax consistency gives convergence to this set.  $\square$

### 636 B.7 Why switching regimes can improve fit

637 Let  $\mathcal{M}_1$  denote the single-regime model and  $\mathcal{M}_K$  the switching-regime model with  $K \geq 2$ . The  
 638 single-regime model is nested in the switching model.

639 **Theorem 8** (Strict population improvement over one global deviation). *Suppose the true conditional*  
 640 *law  $P^*$  belongs to  $\mathcal{M}_K$  and is not contained in the closure of  $\mathcal{M}_1$ . Then*

$$\sup_{\Theta \in \mathcal{M}_K} \mathbb{E}_{P^*}[\ell_1(\Theta)] > \sup_{\theta \in \mathcal{M}_1} \mathbb{E}_{P^*}[\ell_1(\theta)].$$

641 *The gap equals the minimum KL divergence from  $P^*$  to  $\mathcal{M}_1$ , and is strictly positive.*

642 *Proof.* Since  $\mathcal{M}_1 \subset \mathcal{M}_K$ , the left-hand side is at least the right-hand side. Because  $P^* \in \mathcal{M}_K$ , the  
 643 best value over  $\mathcal{M}_K$  is achieved by the true conditional law. For any  $\theta \in \mathcal{M}_1$ ,

$$\mathbb{E}_{P^*}[\ell_1(\Theta^*)] - \mathbb{E}_{P^*}[\ell_1(\theta)] = \text{KL}(P^* \| P_\theta).$$

644 Taking the infimum over  $\theta \in \mathcal{M}_1$  gives the stated gap. Since  $P^*$  is not in the closure of  $\mathcal{M}_1$  and the  
 645 parameter space is compact, this infimum is strictly positive.  $\square$

646 With the uniform convergence conditions from Assumptions 7 and 8, the empirical likelihood gap  
 647 converges to the positive population gap. This justifies the empirical comparison between  $K_u = 1$   
 648 and  $K_u > 1$ .

### 649 B.8 Optional localization result for regime segmentation

650 The next result is an auxiliary statement about the Viterbi path under known parameters. It is not  
 651 needed for estimating the model, but it explains why persistent regimes can yield stable segmentations  
 652 when emissions are well separated.

653 Let  $u_{1:M}^*$  be the true regime path for one trajectory, let  $\hat{u}_{1:M}$  be the Viterbi path under the true  
 654 parameters, and let  $S$  be the number of true switches.

655 **Assumption 9** (Emission separation). There exists  $\delta > 0$  such that for any  $u \neq v$ ,

$$\mathbb{E} \left[ \log \frac{e_u(a_m | x_m, \xi_m)}{e_v(a_m | x_m, \xi_m)} \mid u_m^* = u \right] \geq \delta.$$

656 The log-likelihood ratio increments are bounded or sub-exponential uniformly over contexts.

657 **Assumption 10** (Regime persistence). The transition matrix has large diagonal entries,  $\min_u \Pi_{uu} \geq$   
 658  $1 - \epsilon$ , and all transition probabilities are bounded away from zero on the support.

659 **Proposition 9** (Segmentation errors localize near switches). *Under Assumptions 9 and 10, there*  
 660 *exist constants  $c, C > 0$  such that, for any  $\eta \in (0, 1)$ , with probability at least  $1 - \eta$ , Viterbi errors*  
 661 *are confined to neighborhoods of the true switch points with total length at most  $CS \log(M/\eta)$ .*  
 662 *Consequently,*

$$\frac{1}{M} \sum_{m=1}^M \mathbf{1}\{\hat{u}_m \neq u_m^*\} \leq C \frac{S \log(M/\eta)}{M}.$$

663 Thus, if  $S = o(M/\log M)$ , the misclassification fraction vanishes.

664 *Proof sketch.* Consider an interior block of length  $L$  lying inside a true segment with regime  $u$ .  
 665 For any competing regime  $v \neq u$ , the cumulative log-likelihood ratio favoring  $u$  over  $v$  is a sum  
 666 of increments with mean at least  $\delta$ . By concentration for bounded or sub-exponential increments,  
 667 the probability that  $v$  dominates  $u$  on this block is at most  $\exp(-cL)$ . A union bound over blocks  
 668 and competing regimes shows that, with probability at least  $1 - \eta$ , no interior block longer than  
 669  $C \log(M/\eta)$  is misassigned. Therefore errors can occur only within logarithmic neighborhoods of  
 670 true switch points. The transition penalty from Assumption 10 prevents rapid alternating paths unless  
 671 supported by strong emission evidence, yielding the stated localization.  $\square$

## 672 B.9 Diagnostics after fitting

673 The anchor-based diagnostics used after fitting are descriptive consequences of the fitted probabilistic  
 674 model. For example, one may compare the observed action  $a_m$  to  $a_m^* \in \arg \max_a Q_{\text{anc}}(x_m, a)$ , or  
 675 modify selected components of  $C_u$  and recompute the implied action probabilities. These diagnostics  
 676 are conditional on the learned dynamics, anchor reward, feature map, and regime layer. They are not  
 677 used in training and are not claimed here as standalone causal or decision-theoretic guarantees.

## 678 C Stage I: Latent Dynamics Specification

679 This appendix details the continuous-time dynamics layer summarized in Section 2. Stage I learns  
 680 three coupled components: an action-modulated continuous-time Markov chain (CTMC) for the  
 681 latent state  $z(t)$ , a state-dependent ODE for the outcome proxies  $\lambda(t)$ , and an emission model linking  
 682  $(z, \lambda)$  to observations  $o$ .

683 **Action-modulated CTMC for  $z(t)$ .** Between decision epochs  $[t_m, t_{m+1})$ , action  $a_m$  is held fixed  
 684 and  $z(t) \in \{1, \dots, K_z\}$  evolves as a CTMC with action-dependent generator  $A(a_m) \in \mathbb{R}^{K_z \times K_z}$ .  
 685 The off-diagonal entries are parameterized as

$$[A(a)]_{jk} = \exp(\theta_{jk}^\top \phi(a)), \quad j \neq k,$$

686 where  $\phi(a)$  is an action-feature embedding; diagonal entries satisfy  $[A(a)]_{jj} = -\sum_{k \neq j} [A(a)]_{jk}$ .  
 687 Over an interval of length  $\Delta_m = t_{m+1} - t_m$ , the transition matrix is  $\exp\{A(a_m)\Delta_m\}$ , which directly  
 688 handles irregular decision times.

689 **Outcome-proxy ODE for  $\lambda(t)$ .** Each outcome proxy  $\lambda_\ell(t) \in (0, 1)$  ( $\ell = 1, \dots, L$ ) evolves  
 690 between observations according to a state-dependent first-order ODE

$$\frac{d}{dt} \lambda_\ell(t) = \alpha_\ell^{(z(t))} (1 - \lambda_\ell(t)) - \kappa_\ell^{(z(t))} \lambda_\ell(t), \quad (7)$$

691 with state-conditioned rates  $\alpha_\ell^{(k)}, \kappa_\ell^{(k)} \geq 0$  for each  $k \in \{1, \dots, K_z\}$ . Equation (7) admits a closed-  
 692 form solution conditional on  $z(t) = k$  being constant on a sub-interval, which we use for efficient  
 693 propagation between decision epochs.

694 **Emission model.** Conditional on  $z(t_m) = k$  and  $\lambda(t_m)$ , the observation  $o_m$  is drawn from an  
 695 application-specific likelihood  $p_\theta(o_m | z(t_m) = k, \lambda(t_m))$ . Concrete instantiations for each dataset  
 696 are described in Appendix D.

697 **Belief filtering.** Beliefs over  $z$  are propagated through the CTMC and updated by Bayes' rule:

$$\tilde{b}_m = b_{m-1} \exp\{A(a_{m-1})\Delta_{m-1}\}, \quad b_m(k) \propto \tilde{b}_m(k) p_\theta(o_m | z(t_m) = k, \lambda(t_m)).$$

698 The decision-time state used by the choice model is  $x_m = (b_m, \lambda_m)$ .

## 699 D Experimental Contexts and Data Sources

700 To evaluate our model, we use three sequential decision settings with different timing scales and  
 701 operational contexts: high-frequency ICU sepsis treatment, long-horizon CKD-MBD chronic care,  
 702 and retail pricing. These domains provide heterogeneous trajectories for testing whether a shared  
 703 anchor plus switching constraint-regime layer improves prediction and yields interpretable anchor-  
 704 relative deviations.

Table 2: Treatments and Output Indicators in MIMIC-IV Experiment

Category	Items
<b>Vasopressor</b>	Epinephrine Phenylephrine Norepinephrine Dopamine
<b>Inotropic</b>	Dobutamine Milrinone
<b>Fluid</b>	Crystalloid Colloid Water
<b>Lab Measurement</b>	Low system blood pressure (sysbp_low) Low saturation levels (spo2_sao2_low) Low white blood cell count (wbc_count_low) High arterial lactate (arterial_lactate_high) Normal central venous pressure (cvp_normal)* Normal systemic vascular resistance (svr_normal)* Normal serum creatinine concentration (creatinine_normal)* Normal arterial lactate (arterial_lactate_normal)*
<b>Output</b>	Low-urine

#### 705 D.1 Clinical Decision-Making: MIMIC-IV and CKD-MBD

706 These two clinical datasets offer complementary perspectives on medical decision-making: while  
707 CKD-MBD captures the intricate, slow-evolving 'balancing act' of chronic outpatient care, MIMIC-  
708 IV provides a high-frequency view of acute physiological stabilization in the ICU. Together, they test  
709 the framework's robustness against both persistent cognitive load and rapid environmental shifts.

##### 710 Acute Care (MIMIC-IV):

711 Sepsis is a life-threatening condition resulting from a dysregulated host response to infection, leading  
712 to widespread inflammation, organ failure, and high mortality rates [45]. To study clinician behavior  
713 in this high-stakes setting, we utilize the MIMIC-IV database, a comprehensive resource containing  
714 de-identified, high-resolution health data from ICU patients. Due to the significant uncertainty  
715 and dynamic nature of sepsis management [46], clinician decisions often deviate from stationary  
716 protocols. The rich, irregularly sampled clinical trajectories in MIMIC-IV provide an ideal testbed  
717 for our framework, allowing us to learn a continuous-time latent patient dynamics model and discover  
718 the underlying constraint regimes, such as resource pressures or changing risk postures, that are  
719 associated with persistent deviations in clinical decisions over time.

720 • **Patient Cohort:** We extracted 2,475 sequences of patients diagnosed with sepsis [47], using  
721 a 72-hour observation window with complete data records.

722 • **Treatments (Actions):** The primary interventions in our study are vasopressor therapy and  
723 fluid resuscitation, which are essential for stabilizing blood pressure and ensuring organ  
724 perfusion [4]. In our model, these represent the clinician's actions  $a_t$ . Vasopressors serve to  
725 increase cardiac output via vasoconstriction, while fluid administration restores intravascular  
726 volume. We categorize these treatments into 3 types of fluids, 4 types of vasopressors, and 2  
727 types of inotropes, as detailed in Table 2.

728 • **Outcomes and Observations:** To ground the clinical objective and learn the patient dynam-  
729 ics model, we utilize real-time urine output and survival status as primary indicators. Low  
730 urine output is a critical signal of kidney dysfunction and impending septic shock. Addition-  
731 ally, we include eight informative lab measurements relevant to sepsis: four representing  
732 poor physiological states and four representing favorable outcomes. These variables allow  
733 our latent dynamics model to evaluate the counterfactual consequences of clinician actions.

##### 734 Chronic Management (CKD-MBD):

735 Chronic Kidney Disease-Mineral Bone Disorder (CKD-MBD) is a systemic syndrome characterized  
736 by a complex interplay of biochemical abnormalities, including serum calcium (Ca), phosphorus (P),  
737 and parathyroid hormone (PTH) levels [26, 48]. Because these variables are deeply interconnected  
738 and influence long-term risks of cardiovascular events and bone fractures, clinical management  
739 requires frequent adjustments of multiple medications [49]. We apply our framework to this setting  
740 to understand the latent factors, such as shifting risk postures or institutional guidelines, that lead  
741 clinicians to change their treatment priorities over the course of long-term care.

- 742 • **Data Source and Patients:** We utilize a private dataset collected from a hospital, consisting  
743 of longitudinal clinical records for 150 patients captured between January 2020 and January  
744 2024. This cohort provides a unique, high-resolution view of real-world clinical practice.  
745 Each patient trajectory includes a series of laboratory measurements and the corresponding  
746 pharmacological interventions prescribed by attending physicians.
- 747 • **Treatments (Actions):** The treatment space encompasses 13 discrete action markers repre-  
748 senting dosage levels across five essential drug classes: calcium-based binders, phosphate  
749 binders, vitamin D analogs, and calcimimetics. In our model, these prescriptions represent  
750 the clinician’s actions  $a_t$ . These agents are used concurrently to manage the mineral balance,  
751 but the specific weighting of one drug over another often reflects the clinician’s underlying  
752 decision regime.
- 753 • **Outcomes and Observations:** We monitor three primary lab-based indicators: corrected  
754 calcium (cCa), phosphate (P), and parathyroid hormone (PTH). Following clinical standards,  
755 these are discretized into low, normal, and high categories, forming 10 outcome markers.  
756 These biomarkers serve as the state observations  $x_t$  used to train our continuous-time patient  
757 dynamics model. By evaluating counterfactual outcomes for these three variables, our  
758 framework can infer when a clinician has shifted from prioritizing one metabolic target (e.g.,  
759 phosphorus control) to another (e.g., PTH suppression) through persistent anchor-relative  
760 deviations captured by the regime layer.
- 761 • **Rationale:** The management of CKD-MBD is often described as a “balancing act” where  
762 guidelines (such as the KDIGO updates) provide targets, but individual clinician behavior  
763 fluctuates based on perceived patient risk or cognitive load [50]. This dataset is ideal for  
764 our regime-switching behavioral model because it captures these subtle shifts in treatment  
765 strategy that are not explicitly recorded in the electronic health record but are reflected in  
766 the diverging decision trajectories.

## 767 D.2 Economic Decision-Making: Dominick’s Finer Foods (DFF)

768 To evaluate our model beyond healthcare, we apply it to retail dynamic pricing using the Dominick’s  
769 Finer Foods (DFF) dataset from the Chicago Booth Kilts Center for Marketing [23, 24]. In this  
770 setting, retail managers balance anchor revenue considerations against operational frictions such as  
771 menu costs, promotional inertia, and psychological pricing barriers. The dataset tests whether our  
772 model can recover persistent pricing regimes as anchor-relative deviations rather than treating all  
773 variation as demand noise or stationary reward misspecification.

- 774 • **Data Source and Market Units:** We focus on the Orange Juice category, a standard bench-  
775 mark in the pricing literature. From the broader DFF database, we select a representative  
776 cohort of stores with the most complete and consistent longitudinal records for our exper-  
777 imental analysis. Our final processed dataset consists of trajectories from key retail units  
778 monitored over a 52-week period, providing a high-resolution view of real-world pricing  
779 dynamics.
- 780 • **Pricing (Actions):** The action space  $a_t$  encompasses 8 discrete markers representing pricing  
781 interventions and promotional activities. These include three price tiers for the focal brand  
782 (Citrus Hill), indicators for in-store special displays and significant discount bursts, and  
783 dynamic pricing markers for price inertia, price hikes, and price drops. Competitor pricing  
784 and promotion variables for Minute Maid enter the market-state/context indicators rather  
785 than the action space.
- 786 • **Outcomes and Indicators:** Consistent with our structural modeling approach, we define 9  
787 market state indicators that serve as observations  $x_t$ . The core indicator is Customer Brand  
788 Loyalty (LoyalCH), a persistent latent state reflecting the brand’s resilience to competitive  
789 pricing. Following clinical analogies, we discretize loyalty into four levels: "crisis," "at-  
790 risk," "stable," and "robust". We also track relative market position through price-difference  
791 markers and store-specific fixed effects, while incorporating a real-time purchase signal as a

792 feedback indicator. By evaluating anchor-relative pricing deviations, our model estimates  
 793 constraint-labeled wedges associated with observed pricing inertia and market-state changes.  
 794 • **Rationale:** Retail pricing is often modeled as an MDP in which the agent seeks long-term  
 795 revenue. In practice, observed prices may deviate from an anchor policy because of menu  
 796 costs, promotional calendars, and other operational frictions. This dataset allows us to test  
 797 whether our model recovers switching pricing regimes that improve prediction and provide  
 798 interpretable anchor-relative deviations.

## 799 E Structural Definitions of Regime-Defining Features

800 A central component of our model is the domain feature map  $g(b, \lambda, a)$ , which parameterizes anchor-  
 801 relative wedges through interpretable action and context features. These features do not by themselves  
 802 identify literal constraints; rather, they provide constraint-labeled coordinates on which the regime  
 803 layer places sparse shadow-cost weights.

### 804 E.1 Clinical Management Features (MIMIC-IV and CKD-MBD)

805 The clinical features are grounded in medical safety protocols and the cognitive burden of managing  
 806 multi-drug regimens.

807 **MIMIC-IV ICU Protocol Features** The 16-dimensional feature vector for MIMIC-IV is grounded  
 808 in clinical protocols for sepsis and hemodynamic management.

- 809 • **A) History & Complexity (Dim 0–5):**
  - 810 – **0.** `drug_count`: Total number of base drugs in the current action. It reflects the  
 811 complexity and burden of the medication combination.
  - 812 – **1.** `added_count`: Number of drugs added relative to the previous action.
  - 813 – **2.** `removed_count`: Number of drugs removed relative to the previous action.
  - 814 – **3.** `changed_flag`: A binary indicator (0/1) of whether any medication change occurred  
 815 in this step.
  - 816 – **4.** `n_groups_active`: Count of active coarse treatment groups (*pressor*, *fluid*, or  
 817 *inotrope*), reflecting the breadth of therapeutic modalities.
  - 818 – **5.** `new_groups_count`: Number of newly introduced treatment groups (transitioning  
 819 from unused to active).
- 820 • **B) forbid\_if Violations (Dim 6–9):** Defined as  $1[\text{trigger}] \cdot 1[\text{used forbidden}]$ .
  - 821 – **6.** `forbid_violate[p0_is1_acts[8]]`: Risk penalty for using *Water* (8) when blood  
 822 pressure is critically low ( $p_0$  is True).
  - 823 – **7.** `forbid_violate[p4_is1_acts[8]]`: Risk penalty for using *Water* (8) when lactate  
 824 levels are high ( $p_4$  is True).
  - 825 – **8.** `forbid_violate[p8_is1_acts[8, 7]]`: Risk penalty for using *Water* (8) or *Colloid*  
 826 (7) during low urine output ( $p_8$  is True).
  - 827 – **9.** `forbid_violate[p7_is0_acts[7]]`: Risk penalty for using *Colloid* (7) when creati-  
 828 nine is abnormal ( $p_7$  is False).
- 829 • **C) require\_if Omissions (Dim 10–12):** Defined as  $1[\text{trigger}] \cdot 1[\text{miss all required}]$ .
  - 830 – **10.** `require_miss[p0_is1_acts[0, 1, 2, 3]]`: Omission penalty for failing to provide  
 831 any vasopressors (0/1/2/3) during hypotension.
  - 832 – **11.** `require_miss[p4_is1_acts[6, 7, 2, 0, 1]]`: Omission penalty for failing to provide  
 833 either fluid (6/7) or specific vasopressors (0/1/2) during high lactate.
  - 834 – **12.** `require_miss[p8_is1_acts[6, 2, 0, 1]]`: Omission penalty for failing to provide  
 835 either crystalloid (6) or specific vasopressors (0/1/2) during low urine output.
- 836 • **D) Group Intensity Counts (Dim 13–15):**
  - 837 – **13.** `group_count[Ino]`: Count of active drugs within the *Inotrope* group ([4, 5]).
  - 838 – **14.** `group_count[fluid]`: Count of active drugs within the *Fluid* group ([6, 7, 8]).
  - 839 – **15.** `group_count[pressor]`: Count of active drugs within the *Vasopressor* group ([0,  
 840 1, 2, 3]).

841 **CKD-MBD Chronic Management Features** The 8-dimensional feature set for CKD-MBD focuses  
 842 on long-term safety and resource trade-offs:

- 843 •  $f_1$ : **Hypercalcemia** × **Ca-binder risk**: Safety penalty for calcium-based binders during  
 844 high-calcium states.

- 845 •  $f_2$ :  **$Ca \times P$  risk proxy**: Captures synergistic risk of vascular calcification.
- 846 •  $f_3$ : **Polypharmacy burden**: Structural penalty on regimen complexity.
- 847 •  $f_4$ : **Uncontrolled labs proxy**: Penalizes under-treatment relative to biomarker elevation.
- 848 •  $f_5$ : **Medication cost**: Reflects economic/resource constraints.
- 849 •  $f_6$ : **Treatment escalation**: The "inertia" cost of increasing strategy intensity.
- 850 •  $f_7$ : **Specialty drug flag**: Marks the use of high-tier, specialized therapies.
- 851 •  $f_8$ : **Overall intensity**: Aggregated continuous "dosage" level.

## 852 E.2 Economic Pricing Features (DFF)

853 For the retail pricing domain, we define an 8-dimensional feature set for constraint-labeled pricing  
 854 wedges. These features capture operational and market frictions that can induce persistent deviations  
 855 from the shared pricing anchor.

- 856 •  $f_1$ : **Menu Costs (Inertia)**: A penalty for changing prices between weeks ( $Price_t \neq$   
 857  $Price_{t-1}$ ). This reflects the operational cost of updating tags and systems.
- 858 •  $f_2$ : **Psychological Pricing Barrier**: Captures the "stickiness" of ending prices (e.g., .99).  
 859 Managers are less likely to cross these thresholds.
- 860 •  $f_3$ : **Competitive Response Gap**: Measures the delay in responding to a significant price  
 861 drop by a primary competitor (Minute Maid).
- 862 •  $f_4$ : **Promotion Fatigue**: A penalty for maintaining "Special" status for too many consecutive  
 863 weeks, reflecting diminishing marginal returns.
- 864 •  $f_5$ : **Cannibalization Risk**: Reflects the latent cost of pricing the focal brand (CH) too close  
 865 to the generic house brand, which might erode overall store margins.
- 866 •  $f_6$ : **Brand Equity Guardrail**: A structural constraint that prevents the retail price from  
 867 falling below a "floor" to preserve premium brand image.
- 868 •  $f_7$ : **Inventory-driven Pressure**: Proxy for the urgency to lower prices when historical  
 869 volume decreases (M9/M10 indicators).
- 870 •  $f_8$ : **Site-specific pricing rigidity**: An interaction between the Store 7 indicator and price-  
 871 adjustment actions, capturing site-specific policy frictions in changing prices.

## 872 F Detailed Visualization of Phenotypes

873 To provide a granular understanding of the physiological dynamics captured by the Continuous-  
 874 Time Markov Chain (CTMC), this section presents a detailed visualization of the six latent patient  
 875 phenotypes of CKD-MBD. As illustrated in Figure 6, these phenotypes are characterized through  
 876 their stationary distributions, temporal stability, and population-level evolution.

877 **(a) Phenotype Radar Profiles** The radar charts reveal distinct biochemical signatures for each  
 878 phenotype. Phenotypes  $z_0$  and  $z_1$  exhibit a dominant "Normal" (blue) region, representing stages  
 879 where mineral metabolism is relatively well-controlled. In contrast, Phenotypes  $z_4$  and  $z_5$  show  
 880 significant expansions along the "High" (green/red) axes for Phosphorus (P) and PTH, identifying  
 881 high-risk physiological profiles characterized by severe mineral dysregulation.

882 **(b)  $\gamma$ -weighted Dwell Time** Variations in dwell time reflect the inherent stability of each physi-  
 883 ological state. Phenotype  $z_0$  maintains an average dwell time of 28.3 months, significantly longer  
 884 than other states, designating it as a quasi-stable maintenance state within the cohort. Conversely,  
 885 high-risk states such as Phenotype  $z_5$  have a dwell time of only 3.9 months. This volatility suggests  
 886 that such states are transient and often trigger more aggressive clinical interventions intended to shift  
 887 the patient toward a more stable phenotype.

888 **(c) Population Posterior Evolution** The population-level trends provide empirical evidence of  
 889 treatment progression over time. As the observation period increases, the posterior probabilities for  
 890 stable states ( $z_0$  and  $z_1$ ) show a clear upward trajectory, while the density of high-risk states gradually  
 891 declines. This evolution provides qualitative evidence that the CTMC layer captures population-level  
 892 movement from decompensated phases toward more compensated physiological states.

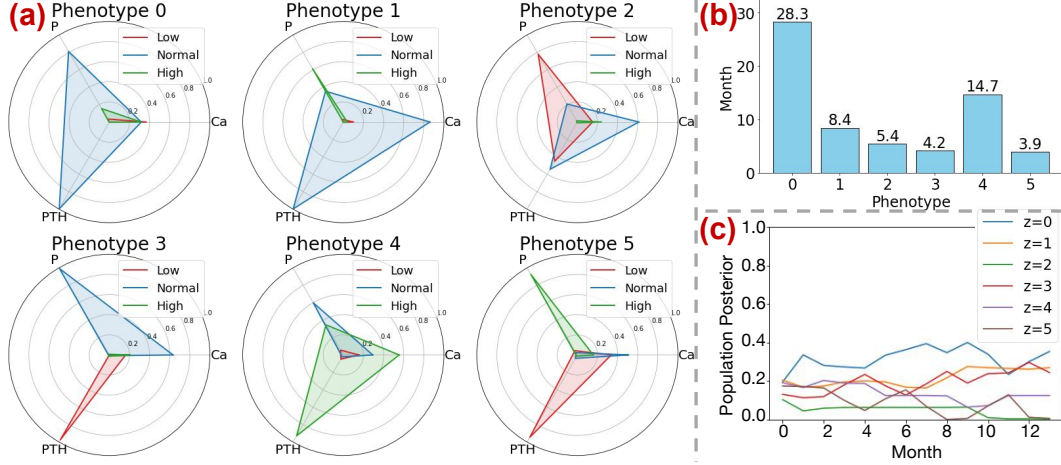


Figure 6: Comprehensive analysis of latent patient phenotypes in the CKD-MBD cohort. (a) **Top Radar Charts:** Emission probabilities of Calcium (Ca), Phosphorus (P), and PTH across three ranges (Low, Normal, High) for each phenotype. (b) **Middle Bar Chart:**  $\gamma$ -weighted average dwell time in months for each phenotype, indicating physiological stability. (c) **Bottom Line Chart:** Temporal evolution of the population posterior probabilities for each latent state  $z_i$  across the treatment horizon.

## 893 G Detailed Sensitivity Analysis for Reward Specifications

### 894 G.1 Reward Specification

895 A common critique in inverse decision modeling is that the “true reward” is not observable and  
 896 may vary across stakeholders. We therefore treat  $r_{\text{anc}}$  as a *anchor proxy* used to define a reference  
 897 action-value function  $Q_{\text{anc}}$ . Taking our clinical applications as the primary evaluative anchor, we  
 898 explicitly evaluate the robustness of our model to plausible baseline families within high-acuity and  
 899 chronic care settings.

900 **Baseline families.** We consider baseline specifications that depend only on available outcomes and  
 901 clinical stability signals. Examples include:

- 902 1. **Outcome-based terminal baseline.** We implement an episodic terminal signal (with  
 903 discounting over decision epochs), e.g.,

$$r_{\text{anc}}^{\text{term}}(z, a) = 0,$$

$$R_T = \mathbb{I}(\text{favorable outcome}) - \alpha \mathbb{I}(\text{adverse outcome}),$$

904 where “favorable/adverse” outcomes are instantiated by the application (e.g., survival vs.  
 905 mortality, remission vs. progression).

- 906 2. **Stability/shaping baselines.** Using physiologic stability proxies derived from observations  
 907 (or emissions), e.g.,

$$r_{\text{anc}}^{\text{stab}}(z, a) = -(\omega_1 \text{shock}(z) + \omega_2 \text{resp\_fail}(z) + \omega_3 \text{renal}(z)),$$

908 where  $\text{shock}(\cdot)$ ,  $\text{resp\_fail}(\cdot)$ , etc. are latent-state risk scores learned from the emission  
 909 model or defined via clinically motivated mappings.

- 910 3. **Risk-sensitive variants.** To reflect aversion to catastrophic deterioration, we consider  
 911 conservative baselines such as a tail-risk action-value  $Q_{\text{anc}}^{\text{CVaR}}(\cdot)$  or moment-penalized  
 912 variants of the form

$$r_{\text{anc}}(z, a) - \kappa_{\text{var}} \text{Var}[\Delta \text{Score} \mid z, a],$$

913 where  $\Delta \text{Score}$  denotes a clinically relevant change (e.g., in an organ-failure proxy) and the  
 914 variance is estimated under model-based rollouts.

915 Importantly, all choices yield the same modeling structure: the reference  $Q_{\text{anc}}$  changes, while the  
 916 inferred switching constraint regimes capture residual anchor-relative deviations.

917 **Practical computation of  $Q_{\text{anc}}$  under continuous-time dynamics.** Given action-held intervals  
 918  $[t_m, t_{m+1})$ , the CTMC defines a transition kernel

$$P^{(a_m)}(\Delta t_m) = \exp(A^{(a_m)} \Delta t_m),$$

Table 3: Predictive performance across reward specifications and regime counts ( $K_u$ ).

Reward Type	$K_u$	Log-Likelihood $\uparrow$	Top-1 Acc. $\uparrow$	Top-5 Acc. $\uparrow$	Top-10 Acc. $\uparrow$
<b>Dataset: CKD-MBD</b>					
Outcome-based	1	$-5.57 \pm 0.39$	$0.17 \pm 0.01$	$0.46 \pm 0.04$	$0.56 \pm 0.06$
	3	<b><math>-4.52 \pm 0.85</math></b>	<b><math>0.22 \pm 0.05</math></b>	<b><math>0.51 \pm 0.05</math></b>	<b><math>0.63 \pm 0.06</math></b>
	5	$-4.96 \pm 0.49$	$0.20 \pm 0.03$	$0.46 \pm 0.06$	$0.58 \pm 0.05$
Stability-based	1	$-5.89 \pm 0.24$	$0.22 \pm 0.03$	$0.41 \pm 0.03$	$0.51 \pm 0.04$
	3	$-4.77 \pm 0.23$	<b><math>0.27 \pm 0.03</math></b>	$0.45 \pm 0.02$	$0.53 \pm 0.04$
	5	<b><math>-4.52 \pm 0.27</math></b>	$0.26 \pm 0.03$	<b><math>0.47 \pm 0.05</math></b>	<b><math>0.54 \pm 0.05</math></b>
Risk-sensitive	1	$-5.57 \pm 0.39$	$0.17 \pm 0.02$	$0.44 \pm 0.04$	$0.59 \pm 0.04$
	3	$-5.40 \pm 0.26$	$0.17 \pm 0.01$	$0.51 \pm 0.05$	$0.63 \pm 0.06$
	5	<b><math>-5.26 \pm 0.33</math></b>	<b><math>0.19 \pm 0.03</math></b>	<b><math>0.51 \pm 0.03</math></b>	<b><math>0.64 \pm 0.02</math></b>
<b>Dataset: MIMIC-IV</b>					
Outcome-based	1	$-0.97 \pm 0.00$	$0.76 \pm 0.01$	$0.95 \pm 0.00$	$0.98 \pm 0.00$
	3	$-0.77 \pm 0.08$	$0.81 \pm 0.04$	<b><math>0.96 \pm 0.01</math></b>	<b><math>0.98 \pm 0.00</math></b>
	5	<b><math>-0.68 \pm 0.06</math></b>	<b><math>0.83 \pm 0.05</math></b>	$0.96 \pm 0.01$	$0.98 \pm 0.00$
Stability-based	1	$-0.91 \pm 0.00$	$0.75 \pm 0.00$	$0.96 \pm 0.00$	$0.98 \pm 0.00$
	3	$-0.85 \pm 0.03$	$0.75 \pm 0.00$	<b><math>0.96 \pm 0.00</math></b>	$0.99 \pm 0.00$
	5	<b><math>-0.78 \pm 0.06</math></b>	<b><math>0.76 \pm 0.01</math></b>	$0.96 \pm 0.00$	<b><math>0.99 \pm 0.00</math></b>
Risk-sensitive	1	$-0.91 \pm 0.00$	$0.76 \pm 0.01$	$0.96 \pm 0.00$	$0.98 \pm 0.00$
	3	$-0.78 \pm 0.02$	<b><math>0.77 \pm 0.00</math></b>	<b><math>0.97 \pm 0.00</math></b>	$0.98 \pm 0.00$
	5	<b><math>-0.74 \pm 0.02</math></b>	$0.77 \pm 0.01$	$0.97 \pm 0.00$	<b><math>0.98 \pm 0.00</math></b>

919 which enables value iteration on the latent state space (or Monte Carlo rollouts) to compute  
 920  $Q_{\text{anc}}(b, \lambda, a)$  for each belief  $b$ . When the emission model is complex, we compute  $Q_{\text{anc}}$  on a  
 921 discretized belief representation (e.g., via clustering beliefs) or via rollouts initialized from the  
 922 posterior belief state  $b$ .

## 923 G.2 Performance under Different Reward Specifications

924 This section provides a comprehensive performance breakdown of the proposed framework under  
 925 different clinical reward specifications. To ensure that the inferred switching regimes are not artifacts  
 926 of a specific objective, we evaluate the framework using three distinct reward families: (i) Outcome-  
 927 based, (ii) Stability-based, and (iii) Risk-sensitive.

928 Table 3 summarizes the predictive metrics for both datasets. Across all reward types, models with  
 929 switching regimes ( $K_u > 1$ ) consistently outperform the stationary baseline ( $K_u = 1$ ), substantiating  
 930 the robustness of our framework.

## 931 H Detailed Comparative Benchmarking of System Dynamics Models

932 This section presents a comparative benchmarking of the CTMC+ODE dynamics model against  
 933 ODE-RNN and NCDSSM baselines through a hold-out prediction task. In Table 4, we evaluate model  
 934 fidelity across diverse temporal scales to ensure robustness across domains: specifically 3-month  
 935 and 6-month horizons for the chronic CKD-MBD cohort, 6-hour and 12-hour horizons for the acute  
 936 MIMIC-IV sequences, and a 3-week out-of-sample window for the DFF retail pricing dataset. This  
 937 multi-horizon assessment validates the model’s capacity to capture both rapid state fluctuations and  
 938 long-term systemic trends.

## 939 I SCM-based Counterfactual Sampling and Anchor-Optimal Action 940 Synthesis

941 To ground sequential decision-making in realistic counterfactuals, we formalize the process as a  
 942 Structural Causal Model (SCM). This allows us to evaluate anchor-optimal counterfactual actions  
 943 while preserving the system’s underlying exogenous characteristics, so anchor-recovery estimates are  
 944 evaluated under unit-level counterfactual consistency rather than fresh stochastic draws.

Table 4: Comparison of System Dynamics Models (Hold-out Prediction). Results represent mean  $\pm$  SD across all biomarkers;  $ECE_{avg}$  denotes mean Expected Calibration Error.

Model	NLL $\downarrow$	Brier $\downarrow$	Acc $\uparrow$	$ECE_{avg}$ $\downarrow$
<b>Dataset: CKD-MBD (3-Month Horizon)</b>				
CTMC+ODE	<b>0.668 <math>\pm</math> 0.012</b>	0.261 $\pm$ 0.001	<b>0.726 <math>\pm</math> 0.006</b>	<b>0.098 <math>\pm</math> 0.006</b>
ODE-RNN	0.672 $\pm$ 0.001	<b>0.239 <math>\pm</math> 0.001</b>	0.672 $\pm$ 0.006	0.199 $\pm$ 0.009
NCDSSM	0.690 $\pm$ 0.009	0.248 $\pm$ 0.004	0.528 $\pm$ 0.039	0.160 $\pm$ 0.021
<b>Dataset: CKD-MBD (6-Month Horizon)</b>				
CTMC+ODE	<b>0.667 <math>\pm</math> 0.009</b>	0.282 $\pm$ 0.002	<b>0.717 <math>\pm</math> 0.011</b>	<b>0.091 <math>\pm</math> 0.008</b>
ODE-RNN	0.673 $\pm$ 0.001	<b>0.240 <math>\pm</math> 0.000</b>	0.664 $\pm$ 0.008	0.185 $\pm$ 0.009
NCDSSM	0.700 $\pm$ 0.008	0.253 $\pm$ 0.004	0.523 $\pm$ 0.018	0.124 $\pm$ 0.021
<b>Dataset: MIMIC-IV (6-Hour Horizon)</b>				
CTMC+ODE	<b>0.195 <math>\pm</math> 0.004</b>	<b>0.112 <math>\pm</math> 0.003</b>	<b>0.936 <math>\pm</math> 0.002</b>	<b>0.031 <math>\pm</math> 0.003</b>
ODE-RNN	0.645 $\pm$ 0.000	0.226 $\pm$ 0.000	0.876 $\pm$ 0.002	0.448 $\pm$ 0.000
NCDSSM	0.861 $\pm$ 0.019	0.145 $\pm$ 0.003	0.847 $\pm$ 0.004	0.283 $\pm$ 0.002
<b>Dataset: MIMIC-IV (12-Hour Horizon)</b>				
CTMC+ODE	<b>0.204 <math>\pm</math> 0.004</b>	<b>0.121 <math>\pm</math> 0.003</b>	<b>0.925 <math>\pm</math> 0.003</b>	<b>0.034 <math>\pm</math> 0.004</b>
ODE-RNN	0.645 $\pm$ 0.000	0.226 $\pm$ 0.000	0.877 $\pm$ 0.001	0.449 $\pm$ 0.000
NCDSSM	1.447 $\pm$ 0.038	0.154 $\pm$ 0.002	0.847 $\pm$ 0.003	0.220 $\pm$ 0.002
<b>Dataset: DFF (3-Week Horizon)</b>				
CTMC+ODE	<b>0.536 <math>\pm</math> 0.000</b>	0.336 $\pm$ 0.000	<b>0.791 <math>\pm</math> 0.001</b>	<b>0.124 <math>\pm</math> 0.001</b>
ODE-RNN	0.701 $\pm$ 0.008	0.249 $\pm$ 0.003	0.690 $\pm$ 0.021	0.439 $\pm$ 0.003
NCDSSM	0.551 $\pm$ 0.011	<b>0.189 <math>\pm</math> 0.004</b>	0.733 $\pm$ 0.000	0.325 $\pm$ 0.012

## 945 I.1 Anchor-Optimal Action and Counterfactual Reward

946 At each decision epoch  $t_m^i$ , we synthesize an anchor-optimal counterfactual action by evaluating  
 947 candidate actions under the fitted structural dynamics. For each candidate  $a \in \mathcal{A}$ , we intervene on the  
 948 current action, propagate the latent trajectory forward under the SCM-consistent rollout, and compute  
 949 the corresponding anchor return. The anchor-optimal action is then

$$a_m^* \in \arg \max_{a \in \mathcal{A}} Q_{\text{anc}}(x_m, a),$$

950 where  $Q_{\text{anc}}(x_m, a)$  is evaluated by Monte Carlo rollouts initialized from the fitted belief state and  
 951 the abducted unit-level noise. The instantaneous reward  $r_\tau$  at time  $\tau$  is defined as a function of the  
 952 deterministic target-alignment probabilities  $\lambda_\ell(\tau)$ :

$$r_\tau = \sum_{\ell \in L} w_\ell \lambda_\ell(\tau) - w_{\text{risk}} \prod_{j \in J} p_j^{\text{hi}}(\tau, z(\tau)), \quad (8)$$

953 where  $w_\ell > 0$  weights the desirability of remaining target-aligned across dimensions  $\ell \in L$ , and  
 954  $w_{\text{risk}} > 0$  penalizes synergistic risks. Consistent with our emission model, the probability of a  
 955 high-risk violation for dimension  $j$  is defined as  $p_j^{\text{hi}}(\tau, z) = \eta_{z,j,\text{hi}}(1 - \lambda_j(\tau))$ , where  $\eta_{z,j,\text{hi}}$  denotes  
 956 the emission split mass assigned to the high-risk off-target category. The cumulative discounted  
 957 reward is:

$$R_m(a) = \sum_{\tau=m}^T \delta^{\tau-m} r_\tau, \quad \delta \in (0, 1], \quad (9)$$

958 For regularly spaced epochs, the discrete discount factor satisfies  $\delta = e^{-\rho\Delta}$ ; with irregular timing we  
 959 use  $\delta^{\tau-m}$  as the epoch-indexed approximation to the continuous-time discount in the anchor-value  
 960 definition in Section 2. Here,  $\lambda_\tau$  acts as a continuous-time surrogate for target alignment along each  
 961 rollout.

## 962 I.2 Abduction-Action-Prediction Framework

963 Following the standard causal inference workflow, we generate counterfactuals by fixing the exoge-  
 964 nous noise inferred from the observed history  $\mathcal{H}^i = (\{o_m^i, a_m^i\}_{m=1}^{M_i})$ .

965 **Step 1: Abduction.** Using current model parameters, we infer the posterior over latent states and  
 966 exogenous noises  $p(z_{1:T}, \varepsilon^{\text{tr}}, \varepsilon^{\text{em}} \mid \mathcal{H}^i)$  via a backward-smoothing particle-EM procedure. Here,  
 967  $\varepsilon^{\text{tr}}$  denotes the specific realizations of Gumbel noise that drove CTMC jumps, and  $\varepsilon^{\text{em}}$  denotes the  
 968 Gumbel noise associated with categorical emissions  $o_m^i$ .

969 **Step 2: Action Intervention.** For each candidate action  $a \in \mathcal{A}$  at decision epoch  $t_m^i$ , we replace  
 970 the factual action with the intervention  $\tilde{a}_{t_m^i} = a$ . Subsequent actions can either follow the factual  
 971 trajectory for single-step action evaluation or be recursively replaced by anchor-optimal actions for  
 972 sequence synthesis. To ensure unit-level consistency, we maintain the latent Poisson arrival times  $s_k$   
 973 and the specific Gumbel realizations  $\varepsilon^{\text{tr}}$  as fixed exogenous traits of the individual unit.

974 **Step 3: Prediction (Forward Rollout).** The system is rolled forward under the intervention  
 975  $\text{do}(a_t := \tilde{a}_t)$  while keeping exogenous noises fixed:

976 1. **Latent State  $z(t)$ :** Transitions at each uniformization jump time  $s_k$  are computed using the  
 977 identical Gumbel realizations  $g_{k,j} \in \varepsilon^{\text{tr}}$  from the abduction step:

$$z_{k+1} = \arg \max_j \{\log \tilde{P}_{z_k j}(\tilde{a}_{s_k}) + g_{k,j}\}, \quad \tilde{P}(\tilde{a}_{s_k}) = I + \frac{1}{\Omega} A(\tilde{a}_{s_k}). \quad (10)$$

978 where  $\Omega \geq \max_{j,a} |A_{jj}(a)|$  is the uniformization rate.

979 2. **Outcome Dynamics  $\lambda_\tau$ :** Between epochs,  $\lambda_\ell$  evolves deterministically via the state-  
 980 dependent ODE used in the fitted dynamics layer. Over an interval of length  $\Delta_\tau$  with  
 981 latent state  $z_\tau$ , the closed-form update is

$$\lambda_\ell(\tau + \Delta_\tau) = \bar{\lambda}_\ell^{(z_\tau)} + (\lambda_\ell(\tau) - \bar{\lambda}_\ell^{(z_\tau)}) \exp\left[-(\alpha_\ell^{(z_\tau)} + \kappa_\ell^{(z_\tau)})\Delta_\tau\right], \quad \bar{\lambda}_\ell^{(z)} = \frac{\alpha_\ell^{(z)}}{\alpha_\ell^{(z)} + \kappa_\ell^{(z)}}. \quad (11)$$

982 3. **Emissions at Observations:** Counterfactual signals  $\tilde{o}_{m,\ell}^i$  are reconstructed using the fixed  
 983 emission Gumbels  $\varepsilon^{\text{em}}$  and the newly generated latent trajectory.

984 **anchor-action synthesis.** For each decision epoch, we average over  $S$  posterior samples to estimate  
 985 the anchor value of each candidate action:

$$\hat{Q}_{\text{anc}}(x_m, a) = \frac{1}{S} \sum_{s=1}^S \sum_{\tau=m}^T \delta^{\tau-m} r_\tau^{(s,a)}.$$

986 We then select  $a_m^* \in \arg \max_{a \in \mathcal{A}} \hat{Q}_{\text{anc}}(x_m, a)$ . This procedure produces anchor-optimal counter-  
 987 factual actions from the unit’s inferred condition and supports anchor-recovery diagnostics under  
 988 constraint relaxation.

## 989 J Case-Level Auditing Example

990 We formalize the policy evaluation process by utilizing a Large Language Model (LLM) as an  
 991 automated reporting layer that translates structured model outputs into natural-language audit reports.  
 992 The auditor receives five structured inputs: system observations  $o_m$ , observed actions  $a_m$ , anchor-  
 993 optimal counterfactual actions  $a_m^*$ , inferred regimes  $u_m$ , and latent states  $z_m$ . It then generates a  
 994 three-tiered report designed for systemic auditing:

- 995 1. **Rationality Analysis:** Assesses the plausibility of the synthesized anchor-optimal action  
 996  $a^*$ , checking whether the anchor action is operationally reachable under the fitted trajectory  
 997 context.
- 998 2. **Wedge Interpretation:** Explains the anchor-relative deviation between  $a_m$  and  $a_m^*$  by  
 999 mapping the active regime  $u_m$  to specific institutional, clinical, or market frictions.
- 1000 3. **Fidelity Summary:** Relates latent transitions in regimes and states to recorded environmen-  
 1001 tal shocks or clinical events as a trajectory-level consistency check.

1002 This auditing layer extends our model from a predictive model into a structured diagnostic tool,  
 1003 enabling stakeholders to retrospectively identify when, where, and how constraint-labeled deviations  
 1004 from the shared anchor arise. Complete prompt templates are provided in Appendix K. To illustrate  
 1005 the practical auditing capability of our model, Figure 7 presents a representative case study from the  
 1006 CKD-MBD cohort, integrating anchor-optimal counterfactual synthesis, latent regime inference, and  
 1007 an LLM-based reporting layer.

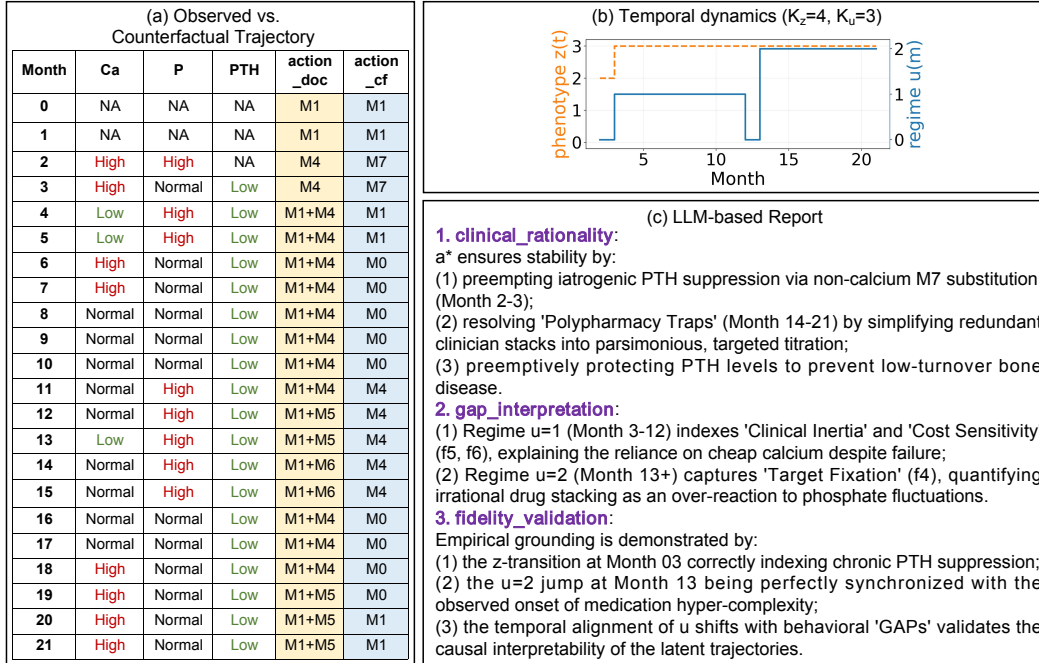


Figure 7: Automated decision auditing and structural counterfactual analysis for a CKD-MBD trajectory. The figure integrates three structural perspectives: **(a) Observed vs. Counterfactual Actions.** Comparison between observed clinician actions ( $a_m$ ) and model-synthesized anchor-optimal actions ( $a_m^*$ ), alongside longitudinal biomarkers (Ca, P, PTH). Treatment markers denote specific medications and dosages: M0–M1 represent calcium carbonate, M4–M6 represent sevelamer, and M7 represents lanthanum carbonate. Combined markers indicate concurrent therapies. **(b) Latent Structural Dynamics.** Temporal evolution of the inferred latent physiological state  $z_m$  and constraint regime  $u_m$ . **(c) LLM-based Institutional Audit Report.**

## 1008 K Detailed Counterfactual Auditing Prompts for LLMs

1009 We utilize GPT-4o as the automated reporting layer, configured with a temperature of 0.1 and a fixed  
 1010 seed of 42 to improve consistency and reproducibility across longitudinal trajectories. The LLM  
 1011 receives structured model outputs only; anchor actions, regimes, and counterfactual quantities are  
 1012 computed by our model.

### 1013 K.1 Detailed Counterfactual Auditing Prompt for MIMIC-IV

#### Complete Prompt Template for Counterfactual Auditing and Clinical Validation (MIMIC-IV)

##### Task Description

- **Role:** You are an expert intensivist specializing in Sepsis and Septic Shock.
- **Objective:** Audit a patient's treatment trajectory by comparing the **Observed Clinical Action** ( $a$ ) against the **anchor-optimal Counterfactual Strategy** ( $a^*$ ) derived from the framework, and assess whether the anchor-relative deviations are clinically plausible under the inferred regime.
- **Context:** Hourly longitudinal records including hemodynamic vital signs, laboratory measurements, and complex intervention markers (Markers 0–102).

##### Clinical Knowledge Base: Marker Definitions (Action and Observation Space)

- **(A) Treatment Markers (Actions 0–8):**
  - **Vasopressors/Inotropes:** Marker 0 (Epinephrine), 1 (Phenylephrine), 2 (Norepinephrine), 3 (Dopamine), 4 (Dobutamine), 5 (Milrinone).
  - **Fluid Resuscitation:** Marker 6 (Crystalloid), 7 (Colloid), 8 (Water).
- **(B) Body Predicate Markers (Observations 9–101):**
  - This section covers 31 laboratory or vital sign items. Each item has **3 markers** representing low, normal, and high ranges.

1014

- Totaling 93 markers (Markers 9–101). For example: Markers 9–11 (Heart Rate), 15–17 (Mean Arterial Pressure), and 99–101 (Lactate).
- **(C) Head Predicate Marker (Clinical Event 102):**
  - Marker 102: Low Urine Output (a critical indicator of sepsis-associated acute kidney injury or impaired organ perfusion).

#### Model Component Definitions

- **Latent Patient State ( $z$ ):** A continuous-time Markov Chain (CTMC) physiological phenotype (e.g., septic shock progression, fluid responsiveness) that determines the baseline clinical reward.
- **Latent Constraint Regime ( $u$ ):** A hidden decision mode reflecting time-varying clinician pressures (e.g., risk-aversion, institutional policy, or resource constraints).
- **anchor-relative Deviation:** Systematic deviations where the observed action  $a$  differs from the anchor-optimal action  $a^*$  due to the shadow costs associated with the active regime  $u$ .

#### Input Format

- `patient_biomarkers = {HOURLY_HEMODYNAMIC_AND_LAB_DATA}`
- `observed_actions_a = {ACTUAL_INTERVENTION_TRAJECTORY}`
- `optimal_actions_a_star = {A_STAR_TRAJECTORY}`
- `inferred_regimes_u = {TIME_SERIES_OF_U}`
- `latent_trajectories_z = {TIME_SERIES_OF_Z}`

#### Your Analytical Task

1. **Rationality Analysis:** Determine whether  $a^*$  would plausibly stabilize vital signs (e.g., MAP) and clear lactate relative to  $a$ . Explain whether the model’s suggested intervention (e.g., vasopressor titration vs. fluid bolus) is clinically sound for sepsis management.
2. **Deviation Interpretation:** Identify plausible unobserved frictions or decision pressures (e.g., fear of pulmonary edema, institutional protocols) using the inferred regime  $u$  to explain the deviation between  $a$  and  $a^*$ .
3. **Trajectory-Regime Correlation:** Correlate transitions in  $z$  and  $u$  with clinical events. Assess whether shifts in  $u$  align with the timing of anchor-relative deviations observed in the actual care.

#### Output (Strict JSON Report)

```
{
  "rationality": "Plausibility of a*.",
  "deviation": "Deviation explained by u.",
  "fidelity": "Links among z, u, and actions."
}
```

1015

## 1016 K.2 Detailed Counterfactual Auditing Prompt for CKD-MBD

### Complete Prompt Template for Counterfactual Auditing and Clinical Validation (CKD-MBD)

#### Task Description

- **Role:** You are an expert nephrologist and clinical auditor specializing in Chronic Kidney Disease-Mineral and Bone Disorder (CKD-MBD).
- **Objective:** Audit a patient’s treatment trajectory by comparing the **Observed Clinical Action ( $a$ )** against the **anchor-optimal Counterfactual Strategy ( $a^*$ )**. Assess whether our model captures clinically plausible anchor-relative deviations through the inferred constraint regime.
- **Context:** Monthly longitudinal records including Calcium (Ca), Phosphorus (P), and PTH. Medication is represented by Markers 0–12.

#### Clinical Knowledge Base: Medication Marker Definitions

- **(A) Calcium Supplements (Calcium Carbonate/Caltrate):**
  - Marker 0: 1.0–3.0 g/day (Low dose).
  - Marker 1: 4.0–8.0 g/day (Medium dose).
  - Marker 2: 9.0–12.0 g/day (High dose).
  - Marker 3: >12.0 g/day (Ultra-high dose, rare).
- **(B) Phosphorus Binders (Sevelamer):**
  - Marker 4: 0.25–3.0 g/day (Low dose).
  - Marker 5: 4.0–8.0 g/day (Medium dose).
  - Marker 6: 9.0–12.0 g/day (High dose).
- **(C) Lanthanum Carbonate:** Marker 7: 1.5 g/day (Specialized binder).
- **(D) Active Vitamin D (Calcitriol and analogs):**

1017

- Marker 8: <1.0  $\mu\text{g}/\text{week}$  (Low dose).
- Marker 9: 1.0–2.0  $\mu\text{g}/\text{week}$  (Medium dose).
- Marker 10: >2.0  $\mu\text{g}/\text{week}$  (High dose).
- (E) **Calcimimetics (Cinacalcet)**: Marker 11 (<0.2 mg/day), Marker 12 (>0.2 mg/day).

#### Model Definitions

- **Latent Patient State ( $z$ )**: A continuous-time physiological phenotype (e.g., hyperparathyroid or mineral imbalance state) that determines the baseline clinical reward.
- **Latent Constraint Regime ( $u$ )**: A hidden decision mode reflecting time-varying clinician pressures (e.g., risk-aversion, institutional policy, or resource constraints).
- **anchor-relative Deviation**: Systematic deviations where the observed action  $a$  differs from the anchor-optimal action  $a^*$  due to the shadow costs associated with the active regime  $u$ .

#### Input Format

- `patient_biomarkers = {MONTHLY_CA_P_PTH_DATA}`
- `observed_actions_a = {ACTUAL_MARKER_TRAJECTORY}`
- `optimal_actions_a_star = {A_STAR_TRAJECTORY}`
- `inferred_regimes_u = {TIME_SERIES_OF_U}`
- `latent_trajectories_z = {TIME_SERIES_OF_Z}`

#### Your Analytical Task

1. **Rationality Analysis**: Determine whether  $a^*$  would plausibly stabilize biochemical markers relative to  $a$ . Explain why the model's suggested dosage is clinically sound.
2. **Deviation Interpretation**: Identify plausible unobserved frictions or decision pressures (e.g., compliance, institutional protocols) using the inferred regime  $u$  to explain the deviation between  $a$  and  $a^*$ .
3. **Trajectory-Regime Correlation**: Correlate transitions in  $z$  and  $u$  with clinical events. Assess whether shifts in  $u$  align with the timing of anchor-relative deviations in the actual care.

#### Output (Strict JSON Report)

```
{
  "rationality": "Plausibility of a*.",
  "deviation": "Deviation explained by u.",
  "fidelity": "Links among z, u, and actions."
}
```

1018

1019

### K.3 Detailed Counterfactual Auditing Prompt for DFF

#### Complete Prompt Template for Counterfactual Auditing and Economic Validation (DFF)

##### Task Description

- **Role**: You are an expert retail economist and pricing strategist specializing in the Orange Juice category (Citrus Hill vs. Minute Maid).
- **Objective**: Audit a retail unit's pricing trajectory by comparing the **Observed Pricing Action ( $a$ )** against the **anchor-optimal Counterfactual Strategy ( $a^*$ )** derived from the framework, and interpret anchor-relative deviations through inferred pricing frictions.
- **Context**: Weekly longitudinal records of pricing, promotional activities, and consumer loyalty (Markers 0–17).

##### Economic Knowledge Base: Marker Definitions (Action and Market Space)

- (A) **Pricing Interventions (Actions 0–7)**:
  - **Price Tiers**: Marker 0 (Ultra-low price, aggressive promotion), Marker 1 (Low price), Marker 2 (Normal/High price, margin protection).
  - **Promotional Activity**: Marker 3 (In-store special display), Marker 4 (Significant discount depth > 15%).
  - **Dynamic Interventions**: Marker 5 (Price Inertia: maintaining previous week's price), Marker 6 (Price Hike), Marker 7 (Price Drop).
- (B) **Market State Indicators (Observations 8–16)**:
  - **Customer Brand Loyalty (LoyalCH)**: Marker 8 (Crisis: <0.25), Marker 9 (At-risk: 0.25–0.5), Marker 10 (Stable: 0.5–0.8), Marker 11 (Robust: >0.8).
  - **Competitive Pressure (Minute Maid)**: Marker 12 (Competitor low-price challenge), Marker 13 (Competitor normal/high price), Marker 14 (Competitor in-store promotion).
  - **Market Positioning**: Marker 15 (Price advantage: CH is significantly cheaper), Marker 16 (Price disadvantage: CH is significantly more expensive).
- (C) **Feedback Signal (Outcome 17)**:

1020

- Marker 17: Purchase success for focal brand (Citrus Hill).

#### Model Component Definitions

- **Latent Market State ( $z$ ):** A hidden environment phenotype (e.g., price elasticity, competitive intensity) determining the baseline revenue potential.
- **Latent Constraint Regime ( $u$ ):** A hidden decision mode reflecting institutional priorities or frictions (e.g., menu costs, rigid promotional calendars, or cognitive load from pricing updates).
- **anchor-relative Deviation:** Systematic deviations where the observed pricing  $a$  differs from the anchor-optimal pricing action  $a^*$  due to shadow costs associated with the active regime  $u$ .

#### Input Format

- `market_indicators = {WEEKLY_LOYALTY_AND_COMPETITIVE_DATA}`
- `observed_actions_a = {ACTUAL_PRICING_TRAJECTORY}`
- `optimal_actions_a_star = {A_STAR_TRAJECTORY}`
- `inferred_regimes_u = {TIME_SERIES_OF_U}`
- `latent_trajectories_z = {TIME_SERIES_OF_Z}`

#### Your Analytical Task

1. **Rationality Analysis:** Determine whether  $a^*$  would plausibly improve purchase probability (Marker 17) and maintain brand loyalty (Markers 10–11) relative to  $a$ . Evaluate whether the suggested move is economically sound.
2. **Deviation Interpretation:** Identify latent frictions (e.g., menu costs, psychological pricing barriers, or institutional "stickiness") using the inferred regime  $u$  to explain why  $a$  differs from the anchor-optimal action  $a^*$ .
3. **Trajectory-Regime Correlation:** Correlate shifts in  $z$  and  $u$  with market events. Assess whether changes in  $u$  align with the timing of anchor-relative deviations.

#### Output (Strict JSON Report)

```
{
  "rationality": "Plausibility of a*.",
  "deviation": "Deviation explained by u.",
  "fidelity": "Links among z, u, and actions."
}
```

1021

## 1022 L Limitations and Broader Impacts

1023 Our model is intended as a retrospective modeling and auditing framework rather than an autonomous  
 1024 decision-making system. The recovered regimes are interpreted relative to the specified anchor  
 1025 reward, fitted latent dynamics, and domain feature map; this conditional interpretation is part of  
 1026 the model design and is supported empirically by dynamics-layer validation and reward-sensitivity  
 1027 analyses. Although we evaluate the framework across clinical and retail settings with different  
 1028 temporal scales, further validation on additional institutions, cohorts, and operational environments  
 1029 would strengthen evidence for broader generalization. In clinical applications, the results are meant to  
 1030 support structured review and hypothesis generation, not direct prospective decision support without  
 1031 expert oversight.

1032 The broader goal of this work is to improve transparency in nonstationary sequential decisions by  
 1033 separating shared anchor behavior from persistent anchor-relative deviations. This can help analysts  
 1034 audit when operational, institutional, or resource-related pressures are associated with deviations from  
 1035 a reference policy. Potential risks include over-interpreting anchor-optimal counterfactuals as direct  
 1036 recommendations or treating LLM-generated audit summaries as definitive conclusions. To mitigate  
 1037 these risks, the LLM is used only as a reporting layer over structured model outputs; all regimes,  
 1038 anchor actions, and counterfactual quantities are computed by the model. Private patient data are  
 1039 not released, and practical deployment should include domain expert review, privacy safeguards, and  
 1040 prospective validation.

## 1041 M Computing Infrastructure

1042 All experiments are performed on an Ubuntu 20.04.3 LTS server with an Intel(R) Xeon(R) Gold  
 1043 6248R CPU @ 3.00 GHz, 227 GB of system memory, and NVIDIA GeForce RTX 3090 GPUs.

1044 **NeurIPS Paper Checklist**

1045 **1. Claims**

1046 Question: Do the main claims made in the abstract and introduction accurately reflect the  
1047 paper’s contributions and scope?

1048 Answer: [Yes]

1049 Justification: The claims are supported by Sections 2–4 and Appendix B.

1050 Guidelines:

- 1051 • The answer [N/A] means that the abstract and introduction do not include the claims  
1052 made in the paper.
- 1053 • The abstract and/or introduction should clearly state the claims made, including the  
1054 contributions made in the paper and important assumptions and limitations. A [No] or  
1055 [N/A] answer to this question will not be perceived well by the reviewers.
- 1056 • The claims made should match theoretical and experimental results, and reflect how  
1057 much the results can be expected to generalize to other settings.
- 1058 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
1059 are not attained by the paper.

1060 **2. Limitations**

1061 Question: Does the paper discuss the limitations of the work performed by the authors?

1062 Answer: [Yes]

1063 Justification: Limitations are discussed in Appendix L.

1064 Guidelines:

- 1065 • The answer [N/A] means that the paper has no limitation while the answer [No] means  
1066 that the paper has limitations, but those are not discussed in the paper.
- 1067 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 1068 • The paper should point out any strong assumptions and how robust the results are to  
1069 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
1070 model well-specification, asymptotic approximations only holding locally). The authors  
1071 should reflect on how these assumptions might be violated in practice and what the  
1072 implications would be.
- 1073 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
1074 only tested on a few datasets or with a few runs. In general, empirical results often  
1075 depend on implicit assumptions, which should be articulated.
- 1076 • The authors should reflect on the factors that influence the performance of the approach.  
1077 For example, a facial recognition algorithm may perform poorly when image resolution  
1078 is low or images are taken in low lighting. Or a speech-to-text system might not be  
1079 used reliably to provide closed captions for online lectures because it fails to handle  
1080 technical jargon.
- 1081 • The authors should discuss the computational efficiency of the proposed algorithms  
1082 and how they scale with dataset size.
- 1083 • If applicable, the authors should discuss possible limitations of their approach to  
1084 address problems of privacy and fairness.
- 1085 • While the authors might fear that complete honesty about limitations might be used by  
1086 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
1087 limitations that aren’t acknowledged in the paper. The authors should use their best  
1088 judgment and recognize that individual actions in favor of transparency play an impor-  
1089 tant role in developing norms that preserve the integrity of the community. Reviewers  
1090 will be specifically instructed to not penalize honesty concerning limitations.

1091 **3. Theory assumptions and proofs**

1092 Question: For each theoretical result, does the paper provide the full set of assumptions and  
1093 a complete (and correct) proof?

1094 Answer: [Yes]

1095 Justification: Assumptions, theorem statements, and proofs are provided in Appendix B.

1096 Guidelines:

- 1097 • The answer [N/A] means that the paper does not include theoretical results.

- 1098 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 1099 referenced.
- 1100 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 1101 • The proofs can either appear in the main paper or the supplemental material, but if
- 1102 they appear in the supplemental material, the authors are encouraged to provide a short
- 1103 proof sketch to provide intuition.
- 1104 • Inversely, any informal proof provided in the core of the paper should be complemented
- 1105 by formal proofs provided in appendix or supplemental material.
- 1106 • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 1107 4. Experimental result reproducibility

1108 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
 1109 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
 1110 of the paper (regardless of whether the code and data are provided or not)?

1111 Answer: [Yes]

1112 Justification: The model (Section 2), the three-stage estimation pipeline (Section 3), the  
 1113 datasets (Appendix D), the feature maps (Appendix E), the reward families (Appendix G),  
 1114 and the counterfactual rollout procedure (Appendix I) jointly enable reproduction of the  
 1115 main results.

1116 Guidelines:

- 1117 • The answer [N/A] means that the paper does not include experiments.
- 1118 • If the paper includes experiments, a [No] answer to this question will not be perceived
- 1119 well by the reviewers: Making the paper reproducible is important, regardless of
- 1120 whether the code and data are provided or not.
- 1121 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 1122 to make their results reproducible or verifiable.
- 1123 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 1124 For example, if the contribution is a novel architecture, describing the architecture fully
- 1125 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 1126 be necessary to either make it possible for others to replicate the model with the same
- 1127 dataset, or provide access to the model. In general, releasing code and data is often
- 1128 one good way to accomplish this, but reproducibility can also be provided via detailed
- 1129 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 1130 of a large language model), releasing of a model checkpoint, or other means that are
- 1131 appropriate to the research performed.
- 1132 • While NeurIPS does not require releasing code, the conference does require all submis-
- 1133 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 1134 nature of the contribution. For example
  - 1135 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
  - 1136 to reproduce that algorithm.
  - 1137 (b) If the contribution is primarily a new model architecture, the paper should describe
  - 1138 the architecture clearly and fully.
  - 1139 (c) If the contribution is a new model (e.g., a large language model), then there should
  - 1140 either be a way to access this model for reproducing the results or a way to reproduce
  - 1141 the model (e.g., with an open-source dataset or instructions for how to construct
  - 1142 the dataset).
  - 1143 (d) We recognize that reproducibility may be tricky in some cases, in which case
  - 1144 authors are welcome to describe the particular way they provide for reproducibility.
  - 1145 In the case of closed-source models, it may be that access to the model is limited in
  - 1146 some way (e.g., to registered users), but it should be possible for other researchers
  - 1147 to have some path to reproducing or verifying the results.

#### 1148 5. Open access to data and code

1149 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
 1150 tions to faithfully reproduce the main experimental results, as described in supplemental  
 1151 material?

1152 Answer: [Yes]

1153 Justification: The code to reproduce our framework and experiments will be made available  
1154 upon acceptance. MIMIC-IV and DFF are available through their original credentialed or  
1155 licensed access mechanisms.

1156 Guidelines:

- 1157 • The answer [N/A] means that paper does not include experiments requiring code.
- 1158 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/  
1159 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1160 • While we encourage the release of code and data, we understand that this might not  
1161 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not  
1162 including code, unless this is central to the contribution (e.g., for a new open-source  
1163 benchmark).
- 1164 • The instructions should contain the exact command and environment needed to run to  
1165 reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 1166 • The authors should provide instructions on data access and preparation, including how  
1167 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 1168 • The authors should provide scripts to reproduce all experimental results for the new  
1169 proposed method and baselines. If only a subset of experiments are reproducible, they  
1170 should state which ones are omitted from the script and why.
- 1171 • At submission time, to preserve anonymity, the authors should release anonymized  
1172 versions (if applicable).
- 1173 • Providing as much information as possible in supplemental material (appended to the  
1174 paper) is recommended, but including URLs to data and code is permitted.

## 1176 6. Experimental setting/details

1177 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-  
1178 rameters, how they were chosen, type of optimizer) necessary to understand the results?

1179 Answer: [Yes]

1180 Justification: The paper describes the core estimation and evaluation protocol in Sections 3–  
1181 4, with dataset construction, feature definitions, reward variants, and counterfactual rollout  
1182 details provided in Appendices D, E, G, and I.

1183 Guidelines:

- 1184 • The answer [N/A] means that the paper does not include experiments.
- 1185 • The experimental setting should be presented in the core of the paper to a level of detail  
1186 that is necessary to appreciate the results and make sense of them.
- 1187 • The full details can be provided either with the code, in appendix, or as supplemental  
1188 material.

## 1189 7. Experiment statistical significance

1190 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
1191 information about the statistical significance of the experiments?

1192 Answer: [Yes]

1193 Justification: All quantitative tables report mean  $\pm$  one standard deviation across multiple  
1194 random seeds.

1195 Guidelines:

- 1196 • The answer [N/A] means that the paper does not include experiments.
- 1197 • The authors should answer [Yes] if the results are accompanied by error bars, confidence  
1198 intervals, or statistical significance tests, at least for the experiments that support the  
1199 main claims of the paper.
- 1200 • The factors of variability that the error bars are capturing should be clearly stated (for  
1201 example, train/test split, initialization, random drawing of some parameter, or overall  
1202 run with given experimental conditions).
- 1203 • The method for calculating the error bars should be explained (closed form formula,  
1204 call to a library function, bootstrap, etc.)
- 1205 • The assumptions made should be given (e.g., Normally distributed errors).
- 1206 • It should be clear whether the error bar is the standard deviation or the standard error  
1207 of the mean.

1208 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
1209 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
1210 of Normality of errors is not verified.

1211 • For asymmetric distributions, the authors should be careful not to show in tables or  
1212 figures symmetric error bars that would yield results that are out of range (e.g., negative  
1213 error rates).

1214 • If error bars are reported in tables or plots, the authors should explain in the text how  
1215 they were calculated and reference the corresponding figures or tables in the text.

1216 **8. Experiments compute resources**

1217 Question: For each experiment, does the paper provide sufficient information on the com-  
1218 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
1219 the experiments?

1220 Answer: [Yes]

1221 Justification: Compute resources are reported in Appendix M.

1222 Guidelines:

1223 • The answer [N/A] means that the paper does not include experiments.

1224 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
1225 or cloud provider, including relevant memory and storage.

1226 • The paper should provide the amount of compute required for each of the individual  
1227 experimental runs as well as estimate the total compute.

1228 • The paper should disclose whether the full research project required more compute  
1229 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
1230 didn't make it into the paper).

1231 **9. Code of ethics**

1232 Question: Does the research conducted in the paper conform, in every respect, with the  
1233 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

1234 Answer: [Yes]

1235 Justification: The work uses only de-identified retrospective records and is positioned as an  
1236 anchor-based auditing tool rather than a treatment or pricing recommender.

1237 Guidelines:

1238 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of  
1239 Ethics.

1240 • If the authors answer [No], they should explain the special circumstances that require a  
1241 deviation from the Code of Ethics.

1242 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
1243 eration due to laws or regulations in their jurisdiction).

1244 **10. Broader impacts**

1245 Question: Does the paper discuss both potential positive societal impacts and negative  
1246 societal impacts of the work performed?

1247 Answer: [Yes]

1248 Justification: Positive and negative societal impacts are discussed in Appendix L.

1249 Guidelines:

1250 • The answer [N/A] means that there is no societal impact of the work performed.

1251 • If the authors answer [N/A] or [No], they should explain why their work has no societal  
1252 impact or why the paper does not address societal impact.

1253 • Examples of negative societal impacts include potential malicious or unintended uses  
1254 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
1255 (e.g., deployment of technologies that could make decisions that unfairly impact specific  
1256 groups), privacy considerations, and security considerations.

1257 • The conference expects that many papers will be foundational research and not tied  
1258 to particular applications, let alone deployments. However, if there is a direct path to  
1259 any negative applications, the authors should point it out. For example, it is legitimate  
1260 to point out that an improvement in the quality of generative models could be used to  
1261 generate Deepfakes for disinformation. On the other hand, it is not needed to point out  
1262 that a generic algorithm for optimizing neural networks could enable people to train  
1263 models that generate Deepfakes faster.

- 1264
- 1265
- 1266
- 1267
- 1268
- 1269
- 1270
- 1271
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
  - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

1272

1273

1274

1275

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pre-trained language models, image generators, or scraped datasets)?

1276

Answer: [N/A]

1277

1278

1279

Justification: Our paper poses no such risks. The paper releases no high-risk pretrained models or scraped datasets; the LLM auditor only consumes model outputs through a standard API.

1280

Guidelines:

- 1281
- 1282
- 1283
- 1284
- 1285
- 1286
- 1287
- 1288
- 1289
- 1290
- The answer [N/A] means that the paper poses no such risks.
  - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
  - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
  - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

1291

1292

1293

1294

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

1295

Answer: [Yes]

1296

1297

1298

Justification: MIMIC-IV (PhysioNet credentialed license), DFF (Chicago Booth Kilts Center academic terms), and all baselines are properly cited in Section 4 and Appendix D and used under their respective licenses.

1299

Guidelines:

- 1300
- 1301
- 1302
- 1303
- 1304
- 1305
- 1306
- 1307
- 1308
- 1309
- 1310
- 1311
- 1312
- 1313
- 1314
- The answer [N/A] means that the paper does not use existing assets.
  - The authors should cite the original paper that produced the code package or dataset.
  - The authors should state which version of the asset is used and, if possible, include a URL.
  - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
  - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
  - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
  - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
  - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

## 13. New assets

1315

1316

1317

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

1318

Answer: [N/A]

1319

Justification: The paper does not release a new dataset or pretrained model asset.

1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349  
1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374

Guidelines:

- The answer [N/A] means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

**14. Crowdsourcing and research with human subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [N/A]

Justification: The paper does not involve crowdsourcing or newly recruited human subjects.

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

**15. Institutional review board (IRB) approvals or equivalent for research with human subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [N/A]

Justification: The clinical data used are fully de-identified retrospective records; no prospective human-subjects research was conducted.

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

**16. Declaration of LLM usage**

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does *not* impact the core methodology, scientific rigor, or originality of the research, declaration is not required.

Answer: [Yes]

Justification: GPT-4o is used only as a structured reporting layer over computed outputs (not as part of the core estimator); details and prompts are in Section 4.3 and Appendix K.

Guidelines:

- The answer [N/A] means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.

1375  
1376

- Please refer to our LLM policy in the NeurIPS handbook for what should or should not be described.