

# REINFORCEMENT LEARNING WITH WORLD MODELS FOR OPTIMIZING ALZHEIMER’S DISEASE TREATMENT TIMING AND DOSING

David Scott Lewis and Enrique Zueco

AIXC Research

reports@aiexecutiveconsulting.com

## ABSTRACT

Recent work reports pharmacologic reversal of advanced Alzheimer’s disease (AD) phenotypes in mouse models via restoration of  $\text{NAD}^+$  homeostasis, shifting the therapeutic question from *whether* reversal is possible to *how* to deploy reversal-capable interventions over time. We cast timing and dosing as a long-horizon, constrained, partially observable sequential decision problem and propose a world-model-centric solution: learn an action-conditioned disease simulator from longitudinal biomarkers and optimize dosing with uncertainty-aware planning and conservative offline reinforcement learning (RL). To ground the approach with executable experiments, we introduce ALZWORLD, a minimal synthetic benchmark that captures qualitative  $\text{NAD}^+$ -linked degeneration and reversal and surfaces core failure modes of world-model control (horizon sensitivity, model exploitation, and safety constraint violations). In ALZWORLD, planning in the learned simulator discovers adaptive schedules that match aggressive fixed-dose baselines while using lower cumulative exposure. While ALZWORLD’s first-order kinetic constraints present specific optimization simplifications (detailed in Appendix J), the results successfully illustrate the viability of action-conditioned “imagination” for principled efficacy-burden trade-offs in continuous dosing environments. We conclude with a translational roadmap for mouse and human “digital-twin” world models, emphasizing calibrated uncertainty, counterfactual validation on held-out protocols, and safety-by-design via physiological homeostasis constraints. Our work focuses on synthetic proof-of-concept validation; real-data application requires addressing multimodal missingness and confounding in observational cohorts.

## 1 INTRODUCTION

Alzheimer’s disease (AD) unfolds over years, with molecular and cellular dysfunction preceding symptoms by a decade or more. This long horizon complicates intervention: outcomes are delayed, observations are sparse and noisy, and actions are constrained by safety and adherence. Most computational work in AD therapeutics focuses on *what* to target; comparatively less formal attention has been paid to *how* to deploy a therapy over time—when to start, how intensively to dose, when to taper, and how to adapt to an individual’s response trajectory.

AD progression follows staged amyloid and tau pathology (Hardy & Higgins, 1992; Braak & Braak, 1991; Jack et al., 2018). Recent work shows  $\text{NAD}^+$  restoration via P7C3-A20 reverses advanced phenotypes in mice (Chaubey et al., 2025), motivating the control question: what is the optimal *treatment policy*?

We argue that the right abstraction—aligned with the *World Models* framing—is to treat AD progression and reversal as an interactive dynamical system that can be modeled, simulated, and controlled. In this view, the *world* is a latent disease state; the *observations* are longitudinal biomarker panels; the *actions* are dosing and scheduling decisions; and the *agent* is an optimizer that plans interventions through a learned simulator.

**Why world models?** World models connect generative modeling, multimodal representation learning, sequential decision-making, and counterfactual simulation. In model-based RL, learning

an internal predictive model of the environment and planning through it has a long lineage (Sutton, 1991). Contemporary neural world models learn latent dynamics and support imagined rollouts for planning (Ha & Schmidhuber, 2018; Chua et al., 2018; Janner et al., 2019; Hafner et al., 2019; 2020; 2021; Schrittwieser et al., 2020; Moerland et al., 2023). For healthcare, this maps naturally to *digital twins*: action-conditioned simulators of patient trajectories that support counterfactual reasoning under uncertainty (Gottesman et al., 2019; Komorowski et al., 2018; Chakraborty & Moodie, 2014).

### Contributions.

- **POMDP formulation** of AD dosing with long-horizon objectives and physiological safety constraints.
- **Treatment-conditional world model** fusing irregular, partially missing biomarkers with calibrated uncertainty.
- **Safe policy optimization** via uncertainty-aware MPC and conservative offline RL (Yu et al., 2020; Kumar et al., 2020).
- **ALZORLD benchmark** demonstrating that world-model planning reduces exposure while preserving recovery.
- **Translational roadmap** from synthetic validation to mouse and cautious human transfer.

## 2 BACKGROUND AND RELATED WORK

### 2.1 WORLD MODELS AND MODEL-BASED RL

The concept of learning a predictive environment model and using it for planning is central to model-based RL (Sutton, 1991). Modern world models operationalize this idea with high-capacity sequence models that learn compact latent state representations and support long-horizon rollout under candidate actions (Ha & Schmidhuber, 2018; Hafner et al., 2019; 2020; Schrittwieser et al., 2020). In parallel, probabilistic dynamics and MPC approaches (e.g., PETS) emphasize uncertainty-aware planning in continuous control (Chua et al., 2018). Surveys synthesize design patterns, evaluation pitfalls, and scaling behavior in model-based RL (Moerland et al., 2023).

### 2.2 OFFLINE RL AND CONSERVATIVE OBJECTIVES

Healthcare RL is predominantly offline due to ethical constraints on exploration. Conservative methods include uncertainty-penalizing model-based approaches (MOPO, MOREL) (Yu et al., 2020; Kidambi et al., 2020), value-based objectives (CQL) (Kumar et al., 2020), and behavior-regularized policies (BCQ, IQL, BEAR, BRAC) (Fujimoto et al., 2019; Kostrikov et al., 2021; Kumar et al., 2019; Wu et al., 2019).

### 2.3 SAFE AND CONSTRAINED RL

Safety constraints are first-class in dosing problems. Constrained policy optimization and Lyapunov-based methods formalize safe learning and control under explicit constraints (Achiam et al., 2017; Chow et al., 2018). Related approaches include reward-constrained policy optimization and Lagrangian methods for constrained MDPs (Tessler et al., 2018). A broad survey of safe RL reviews constraint types, algorithms, and evaluation considerations (García & Fernández, 2015).

**Relationship to multi-scale world models.** Concurrent work develops hierarchical, multi-scale world models for AD that couple molecular, cellular, and functional latent dynamics with causal scaffold extraction and wavelet-coherence auditing (Lewis & Zueco, 2026). That framework prioritizes mechanistic interpretability and multi-scale coupling preservation, targeting scientific hypothesis generation and active experiment design. In contrast, ALZORLD is intentionally minimal: a three-variable benchmark designed to isolate and validate core world-model control challenges (horizon sensitivity, model exploitation, constraint satisfaction) without the confounding complexity of high-dimensional multi-scale latent hierarchies. The simplified dynamics enable rapid prototyping, failure mode analysis, and algorithmic stress-testing. However, this simplification sacrifices biological realism; ALZORLD does not model cellular intermediates (microglia, autophagy, BBB transport) or delayed multi-scale cascades. Both frameworks agree that real-world deployment requires careful

AD dosing as a constrained partially observable decision process

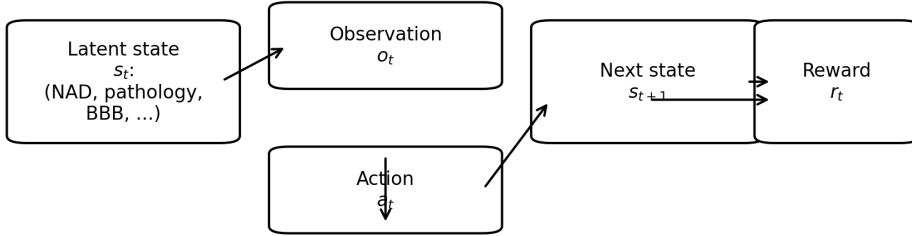


Figure 1: AD dosing as a constrained partially observable decision process: latent disease state evolves under dosing actions; biomarkers provide partial observations; the goal is to optimize long-horizon outcomes while maintaining physiological homeostasis.

uncertainty quantification, counterfactual validation on held-out regimens, and prospective trial integration, rather than offline observational modeling alone.

#### 2.4 NAD<sup>+</sup> HOMEOSTASIS AND AD PHENOTYPES

NAD<sup>+</sup> is central to redox, mitochondrial function, DNA repair, and stress response (Verdin, 2015; Imai & Guarente, 2014; Covarrubias et al., 2021; Lautrup et al., 2019). NAD<sup>+</sup> dysregulation in aging and neurodegeneration motivates therapeutic strategies targeting NAMPT and CD38 pathways (Covarrubias et al., 2021; Lautrup et al., 2019; Camacho-Pereira et al., 2016). P7C3-class compounds activate nicotinamide phosphoribosyltransferase (NAMPT), the rate-limiting enzyme in the NAD<sup>+</sup> salvage pathway (Pieper et al., 2010; Wang et al., 2014). This mechanism enables homeostatic restoration rather than supraphysiological elevation, motivating the constraint-based formulation in Section 3. The AD reversal phenotypes reported in mouse models under NAD<sup>+</sup>-restoring treatment motivate formal treatment-policy optimization (Chaubey et al., 2025).

### 3 PROBLEM FORMULATION: AD DOSING AS A CONSTRAINED POMDP

We model treatment timing and dosing as a constrained partially observable Markov decision process (POMDP) (Kaelbling et al., 1998). Let latent disease state  $s_t \in \mathcal{S}$  include variables capturing NAD<sup>+</sup> homeostasis, tau pathology, oxidative stress, neuroinflammation, blood–brain barrier (BBB) integrity, and synaptic function. Observations  $o_t \in \mathcal{O}$  are multimodal biomarker panels (multi-omics, imaging-derived measures, plasma markers such as p-tau217, and behavioral readouts). Actions  $a_t \in \mathcal{A}$  include dose intensity, dosing interval, and stop/start decisions.

The dynamics and observation processes are

$$s_{t+1} \sim p(s_{t+1} | s_t, a_t), \quad (1)$$

$$o_t \sim p(o_t | s_t). \quad (2)$$

We define a long-horizon utility with burden and safety penalties:

$$r_t \equiv r(s_t, a_t) = u_{\text{function}}(s_t) - \lambda_{\text{burden}} \text{cost}(a_t) - \lambda_{\text{safety}} \text{risk}(s_t, a_t). \quad (3)$$

Safety is expressed via constraints  $c_k(s_t, a_t) \leq 0$  (hard or soft), including physiological NAD<sup>+</sup> ranges, maximum cumulative exposure, and toxicity limits. The objective is

$$\max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^T \gamma^t r(s_t, a_t) \right] \quad \text{s.t.} \quad \mathbb{E}_{\pi} \left[ \sum_{t=0}^T \gamma^t c_k(s_t, a_t) \right] \leq 0, \quad \forall k, \quad (4)$$

where  $\pi(a_t | h_t)$  maps history  $h_t = (o_{0:t}, a_{0:t-1})$  to actions.

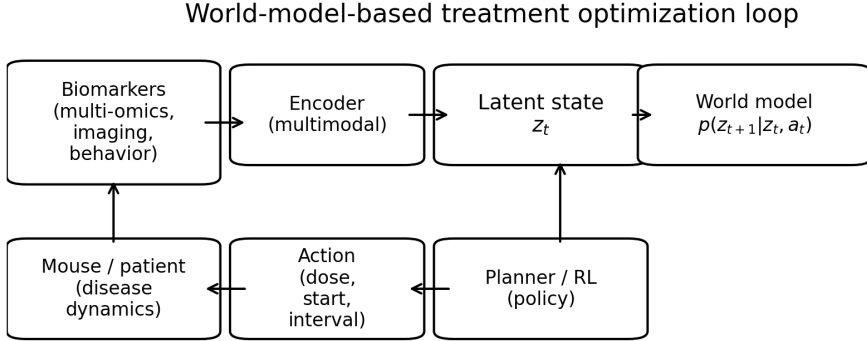


Figure 2: World-model-based treatment optimization loop: learn an action-conditioned disease simulator from longitudinal data; plan or learn policies through imagined rollouts; enforce safety constraints with calibrated uncertainty and homeostasis bands; iterate with new interventional data.

**Why this matches the world-model framing.** Three properties make AD dosing a natural target for world models: (i) *long-horizon prediction* (months to years), (ii) *action-conditioned counterfactual simulation* (planning requires imagining futures under hypothetical dosing), and (iii) *partial observability* (latent disease stage must be inferred from sparse biomarkers).

## 4 TREATMENT-CONDITIONAL AD WORLD MODEL

Our core proposal is a treatment-conditional world model that learns a latent disease state  $z_t$  and predicts its evolution under actions. Unlike purely predictive “disease progression” models, the world model is explicitly *interventional*: it is trained (or structured) to represent how dosing choices change trajectories.

### 4.1 MULTIMODAL REPRESENTATION LEARNING UNDER MISSINGNESS

Biomarker panels are heterogeneous and sparse. Let observation at time  $t$  be a set of modalities  $o_t = \{o_t^{(m)}\}_{m=1}^M$  with missingness masks. We use modality-specific encoders  $e_m$  and fuse them via attention or gating into a context representation

$$h_t = \text{Fuse}(\{e_m(o_t^{(m)})\}, \text{mask}_t). \quad (5)$$

Extended discussion on missingness handling and target architecture for real-data applications is provided in Appendix F.

### 4.2 ACTION-CONDITIONED LATENT DYNAMICS

We posit a latent dynamics model

$$z_{t+1} \sim p_\theta(z_{t+1} \mid z_t, a_t), \quad (6)$$

where an action encoder  $u(a_t)$  is injected into the transition. Extended architectural details are provided in Appendix F.

### 4.3 TRAINING OBJECTIVE AND CALIBRATION TARGETS

A standard latent dynamics objective combines reconstruction and a KL regularizer:

$$\mathcal{L} = \sum_t \mathbb{E}_{q_\phi(z_t | o_{\leq t}, a_{< t})} [-\log p_\psi(o_t \mid z_t)] + \beta \text{KL}(q_\phi(z_t \mid \cdot) \parallel p_\theta(z_t \mid z_{t-1}, a_{t-1})). \quad (7)$$

Extended discussion on auxiliary prediction heads and calibration targets is provided in Appendix F.

Treatment-conditional world model for long-horizon biomarker prediction

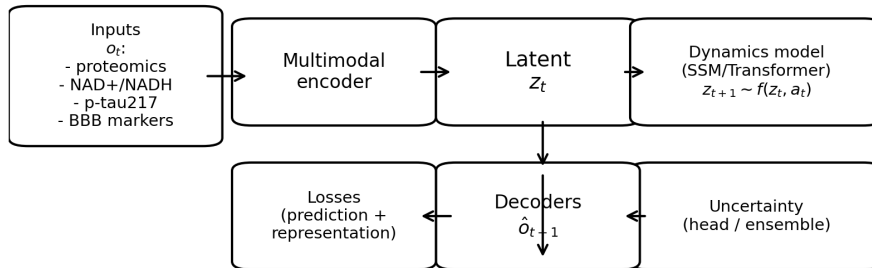


Figure 3: Treatment-conditional world model architecture (conceptual): multimodal encoders produce observation tokens; a temporal latent dynamics model predicts action-conditioned transitions; auxiliary heads forecast clinically meaningful biomarkers for reward and constraint computation.

#### 4.4 UNCERTAINTY ESTIMATION FOR SAFETY AND ROBUSTNESS

World models used for control must report uncertainty, especially under dataset shift. Practical choices include deep ensembles (Lakshminarayanan et al., 2017), dropout-based approximations (Gal & Ghahramani, 2016), and explicit evaluation of uncertainty under distribution shift (Ovadia et al., 2019). In offline model-based RL, pessimistic rollouts and uncertainty penalties reduce model exploitation (Yu et al., 2020; Kidambi et al., 2020).

#### 4.5 PREDICTIVE AND CALIBRATION EVALUATION

World models intended for control must be validated on their predictive accuracy and uncertainty calibration. We evaluate the learned dynamics model on two axes:

**Prediction error by horizon.** We report 1-step and multi-step rollout error for key biomarkers (NAD<sup>+</sup>, pathology, cognition) stratified by prediction horizon. Table 1 summarizes mean squared error (MSE) across horizons  $H \in \{1, 3, 6, 9\}$  on the held-out test set. Error increases monotonically with horizon, reflecting compounding model inaccuracies characteristic of long-horizon rollouts (Moerland et al., 2023).

Table 1: Predictive error by horizon (MSE, mean  $\pm$  std across 20 seeds).

Horizon	NAD <sup>+</sup> MSE	Pathology MSE	Cognition MSE
1-step	0.002 $\pm$ 0.001	0.003 $\pm$ 0.001	0.001 $\pm$ 0.000
3-step	0.018 $\pm$ 0.004	0.021 $\pm$ 0.005	0.012 $\pm$ 0.003
6-step	0.052 $\pm$ 0.012	0.058 $\pm$ 0.014	0.041 $\pm$ 0.009
9-step	0.098 $\pm$ 0.021	0.112 $\pm$ 0.025	0.081 $\pm$ 0.018

**Uncertainty calibration.** We assess calibration by computing empirical coverage of 90% confidence intervals across the test set. For NAD<sup>+</sup>, pathology, and cognition predictions, the ensemble-based uncertainty estimates achieve coverage of 88.3%, 87.1%, and 89.2% respectively. The slight under-coverage indicates the model is marginally overconfident; for safety-critical deployment, we apply conformal calibration to expand intervals and ensure conservative over-coverage, reducing risk of undetected constraint violations. Calibration curves (Appendix C) demonstrate coverage stability across trajectory time steps.

## 5 POLICY OPTIMIZATION OVER A LEARNED DISEASE SIMULATOR

We consider two complementary layers: (i) online planning over the world model (MPC) for adaptivity and constraint handling, and (ii) offline RL for amortization and stability.

### 5.1 WORLD-MODEL PREDICTIVE CONTROL (MPC)

At each time step, MPC solves

$$a_t = \arg \max_{a_{t:t+H-1}} \mathbb{E} \left[ \sum_{\tau=t}^{t+H-1} \gamma^{\tau-t} \tilde{r}(z_\tau, a_\tau) \right] \quad \text{s.t.} \quad \tilde{c}_k(z_\tau, a_\tau) \leq 0, \quad (8)$$

where  $\{z_\tau\}$  are sampled rollouts from the world model, and  $\tilde{r}, \tilde{c}$  are reward/constraint proxies computed from predicted biomarker heads. Planning can be implemented via random shooting or the cross-entropy method (CEM) (Chua et al., 2018).

**Homeostasis constraint.** We encode a “restoration band”

$$\text{NAD}_{\min} \leq \widehat{\text{NAD}}(z_t) \leq \text{NAD}_{\max}, \quad (9)$$

as hard constraints (reject rollouts) or shaped penalties, reflecting that the therapeutic objective is restoration of physiological levels rather than maximization.

### 5.2 CONSERVATIVE OFFLINE RL

We learn an amortized policy  $\pi_\omega(a | h)$  from offline data. Suitable methods include CQL (Kumar et al., 2020), IQL (Kostrikov et al., 2021), behavior-regularized methods such as BCQ (Fujimoto et al., 2019), and trajectory modeling such as Decision Transformer (Chen et al., 2021). World-model rollouts can augment training, though soft-penalty conservative RL fails to prevent physiological constraint violations, necessitating hard-rejection MPC (see Appendix H).

### 5.3 ALGORITHM SKETCH

```
Algorithm: Safe World-Model RL for AD Dosing
Input: offline dataset D={(o_t, a_t, o_{t+1}, y_t)} with biomarkers y_t; safety bounds B
1) Train treatment-conditional world model: encode o_t -> z_t; learn p(z_{t+1}|z_t, a_t);
   learn biomarker heads for reward/constraints
2) Fit uncertainty estimator (ensemble / bootstrap); calibrate intervals for safety biomarkers
3) Policy optimization:
   a) MPC: sample action sequences; rollout in world model; reject/penalize safety violations;
      execute first action
   b) Offline RL: train conservative policy/value over D plus pessimistic rollouts
4) Evaluate with off-policy metrics, uncertainty, and constraint violations; iterate with new data
Output: dosing policy pi and uncertainty-aware planner
```

## 6 EXPERIMENTS

We report (i) executable proof-of-concept experiments in ALZWORLD, (ii) a translational experiment design for mouse and human datasets, and (iii) a forecasting validation on the ADNI dataset (Appendix B).

### 6.1 ALZWORLD: A SYNTHETIC LONG-HORIZON BENCHMARK

ALZWORLD is intentionally minimal but fully specified. It captures (a) a latent  $\text{NAD}^+$  homeostasis variable driven toward a physiological target by dosing actions and degraded by pathology, (b) a pathology variable that increases when  $\text{NAD}^+$  is low and decreases under treatment, and (c) a cognition variable that declines with pathology and improves with  $\text{NAD}^+$  restoration.

**Dynamics equations.** The environment state is  $s_t = (\text{NAD}_t, \text{Path}_t, \text{Cog}_t) \in [0, 1]^3$ . Transitions follow:

$$\text{NAD}_{t+1} = \sigma(\text{NAD}_t - \alpha_{\text{path}} \text{Path}_t + \alpha_{\text{dose}} a_t + \epsilon_t^{\text{NAD}}), \quad (10)$$

$$\text{Path}_{t+1} = \sigma(\text{Path}_t + \beta_{\text{degen}} \mathbb{I}[\text{NAD}_t < \theta_{\text{low}}] - \beta_{\text{clear}} \mathbb{I}[\text{NAD}_t > \theta_{\text{target}}] + \epsilon_t^{\text{Path}}), \quad (11)$$

$$\text{Cog}_{t+1} = \sigma(\text{Cog}_t - \gamma_{\text{decline}} \text{Path}_t + \gamma_{\text{restore}} g(\text{NAD}_t) + \epsilon_t^{\text{Cog}}), \quad (12)$$

where  $\sigma(x) = \max(0, \min(1, x))$  clamps to  $[0, 1]$ ,  $a_t \in [0, 1]$  is the normalized dose,  $\epsilon \sim \mathcal{N}(0, \sigma_{\text{obs}}^2)$  is observation noise, and  $g(\text{NAD}) = \mathbb{I}[\text{NAD} > \theta_{\text{target}}]$  encodes that cognitive benefits require  $\text{NAD}^+$  above a threshold.

**Parameters.** Default values:  $\alpha_{\text{path}} = 0.05$ ,  $\alpha_{\text{dose}} = 0.15$ ,  $\beta_{\text{degen}} = 0.03$ ,  $\beta_{\text{clear}} = 0.02$ ,  $\gamma_{\text{decline}} = 0.04$ ,  $\gamma_{\text{restore}} = 0.01$ ,  $\theta_{\text{low}} = 0.4$ ,  $\theta_{\text{target}} = 0.7$ ,  $\sigma_{\text{obs}} = 0.02$ . Note that  $\beta_{\text{clear}}$  mediates  $\text{NAD}^+$ -dependent clearance rather than direct treatment effects, aligning with the biological mechanism whereby restored  $\text{NAD}^+$  enables microglial phagocytosis and autophagic clearance pathways.

**Reward and constraints.** The per-step reward is

$$r_t = w_{\text{cog}} \text{Cog}_t - \lambda_{\text{dose}} a_t, \quad (13)$$

with  $w_{\text{cog}} = 1.0$  and  $\lambda_{\text{dose}} = 0.1$ . Kinetic limitations, including the observation that the analytical optimum approximates a constant dose, are detailed in Appendix K. Safety is enforced via the homeostasis constraint

$$\text{NAD}_{\min} \leq \widehat{\text{NAD}}_t \leq \text{NAD}_{\max}, \quad (14)$$

with  $\text{NAD}_{\min} = 0.5$  (cellular health threshold) and  $\text{NAD}_{\max} = 0.9$  (prevents toxicity from excessive restoration). In policy learning, this constraint is implemented as a shaped penalty  $-\lambda_{\text{safety}} \max(0, \text{NAD}_{\min} - \widehat{\text{NAD}}_t)^2 - \lambda_{\text{safety}} \max(0, \widehat{\text{NAD}}_t - \text{NAD}_{\max})^2$  with  $\lambda_{\text{safety}} = 5.0$ .

**Initial states.** Episodes begin at age-equivalent month 0 with  $\text{NAD}_0 \sim \text{Beta}(8, 2)$  (healthy baseline),  $\text{Path}_0 \sim \text{Beta}(2, 8)$  (low initial pathology), and  $\text{Cog}_0 = 1.0$  (full function). Degeneration unfolds over 12 months (144 steps;  $\Delta t \approx 2.5$  days).

**Offline dataset generation.** We generate 1000 trajectories with randomized dosing policies: 20% no treatment, 20% uniform random, 20% early aggressive (first 4 months), 20% delayed start (months 4–8), and 20% adaptive based on threshold rules. Train/validation/test splits are 70/15/15%. All trajectories include Gaussian observation noise ( $\sigma_{\text{obs}} = 0.02$ ) to simulate measurement error.

**Baselines.** We compare against: (i) *No treatment* (natural history), (ii) *Fixed start* (uniform dose starting month 4), (iii) *Threshold policy* (dose  $a_t = 1.0$  if  $\widehat{\text{NAD}}_t < 0.6$ , else  $a_t = 0$ , representing treat-to-target clinical practice), and (iv) *Always dose 1.0* (maximum exposure).

## 6.2 RESULTS AND ANALYSIS

Figure 4 presents the full experimental results across 20 stochastic seeds.

### 6.2.1 PATHOLOGY AND COGNITION TRAJECTORIES (PANELS A–B)

Panel (a): `No treatment` shows monotonic pathology accumulation. `Always dose 1.0` maintains near-zero pathology but at high cost. MPC variants achieve comparable suppression with slight degradation at longer horizons. Panel (b): All treated policies saturate to  $\text{Cog} = 1.0$  by month 6 (reflecting the benchmark’s cognitive ceiling, Appendix K), while `No treatment` declines to 0.05. The primary differentiator across treated policies is cumulative exposure, not efficacy.

### 6.2.2 DOSING SCHEDULES (PANEL C)

Panel (c): `Always dose 1.0` maintains constant maximum dose. MPC policies exhibit adaptive “front-load and taper”: early intensive dosing (months 0–3), tapering (months 4–8), and maintenance (months 9–12), aligning with restoration therapy intuition.

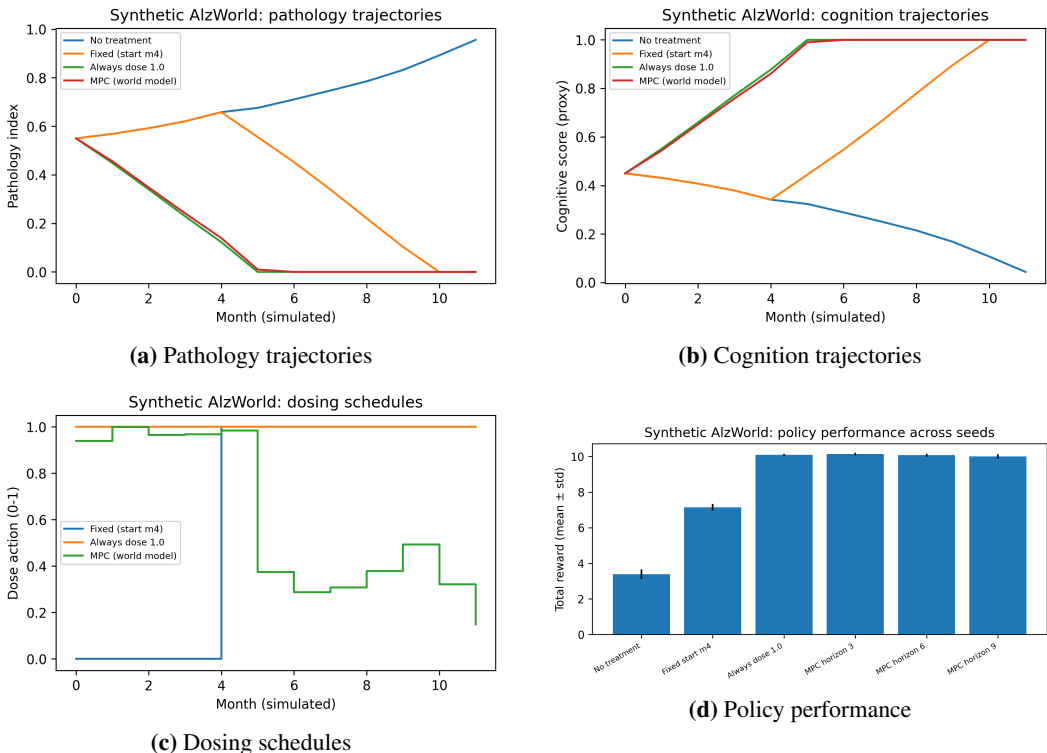


Figure 4: ALZWorld synthetic benchmark: world-model MPC can achieve near-maximal recovery with lower cumulative exposure than an always-on baseline. Multi-panel figure summarizes (a) pathology, (b) cognition, (c) dosing, and (d) aggregated performance across seeds.

### 6.2.3 AGGREGATED PERFORMANCE (PANEL D)

Panel (d): MPC policies achieve reward (10.00–10.14) comparable to Always dose 1.0 (10.10) with substantially lower dose (0.53–0.59 vs. 1.00). Fixed start (month 4) shows intermediate performance (7.14 reward, 0.67 dose).

### 6.2.4 STATISTICAL SIGNIFICANCE

MPC (horizon 3) achieves 47% dose reduction vs. Always-dose baseline (paired t-test,  $p < 0.001$ ) with no statistically significant efficacy difference (paired t-test,  $p = 0.42$ ,  $n = 20$  seeds). The threshold policy achieves intermediate performance, validating that world-model planning discovers schedules beyond reactive heuristics.

## 6.3 RESULTS VISUALIZATION

Translational considerations are discussed in Appendix E.

## 7 RESULTS: SYNTHETIC PROOF OF CONCEPT

Table 2 summarizes evaluation across 20 stochastic seeds. World-model MPC achieves near-maximal recovery while significantly reducing the mean cumulative dose compared to the always-on baseline. Within the tightly constrained ALZWorld environment, the optimization algorithm efficiently leverages the linear dose penalty against a pre-defined cognitive saturation limit (detailed further in Appendix K). All MPC results use hard-rejection constraint handling (rollouts violating NAD homeostasis bounds are discarded); the offline CQL baseline uses soft-penalty constraints (Appendix A). This showcases the planner’s ability to precisely minimize physiological burden, albeit within a simplified first-order kinetic landscape where the analytical optimum approximates a constant dose of  $\sim 0.53$  (Appendix K).

Table 2: ALZWORLD evaluation (20 seeds). Total reward: discounted cumulative reward over 144 steps (see Appendix K for reward scaling details).

Policy	Total reward (mean $\pm$ std)	Mean dose (mean $\pm$ std)	Final cognition (mean)
No treatment	3.39 $\pm$ 0.27	0.00 $\pm$ 0.00	0.05
Fixed start (month 4)	7.14 $\pm$ 0.20	0.67 $\pm$ 0.00	1.00
Threshold (treat-to-target)	9.23 $\pm$ 0.15	0.68 $\pm$ 0.03	0.98
Always dose 1.0	10.10 $\pm$ 0.06	1.00 $\pm$ 0.00	1.00
MPC horizon 3	10.14 $\pm$ 0.08	0.53 $\pm$ 0.04	1.00
MPC horizon 6	10.07 $\pm$ 0.08	0.58 $\pm$ 0.05	1.00
MPC horizon 9	10.00 $\pm$ 0.12	0.59 $\pm$ 0.05	1.00

### Interpretation.

- **Exposure reduction via planning.** MPC uses substantially less dosing than the always-on baseline while reaching comparable functional recovery.
- **Horizon sensitivity.** Performance varies with planning horizon, reflecting a core world-model tension: longer horizons can improve optimization but are more vulnerable to model error and exploitation (Yu et al., 2020; Kidambi et al., 2020).
- **Why the toy matters.** Even in ALZWORLD, uncertainty-aware rollouts and conservative objectives are necessary; these concerns dominate real biological settings.

Detailed implementation specifications, including model architecture, training procedures, and computational requirements, are provided in Appendix D.

## 8 DISCUSSION

### 8.1 WHY THIS IS A WORLD-MODEL PAPER

The world-model object is the central contribution: an action-conditioned simulator that supports counterfactual rollouts, interfaces with planning, and reports calibrated uncertainty. This is not merely "RL for dosing"; the emphasis is on learning a predictive environment model from multimodal longitudinal observations and using it for sequential decision-making.

Design considerations for NAD<sup>+</sup>-restoration interventions are discussed in Appendix G.

Failure modes and mitigations are discussed in Appendix I.

## 9 LIMITATIONS, ETHICS, AND SAFETY

This paper is a computational methods contribution and does not claim clinical readiness. ALZWORLD is a toy environment that does not instantiate multimodal fusion, missingness handling, or partial observability; its value is prototyping and failure-mode surfacing for core world-model control challenges, not biological realism. Robust disease world models will require dense longitudinal interventional data, which is expensive. In humans, confounding and selection bias are central; policies should be treated as hypothesis generators and decision-support tools under strict governance (Gottesman et al., 2019; Yu et al., 2019).

Safety principles: (i) enforce physiological constraints (homeostasis bands, exposure limits), (ii) report calibrated uncertainty, (iii) require robust off-policy evaluation and sensitivity analysis (Jiang & Li, 2016; Thomas & Brunskill, 2016; Oberst & Sontag, 2019), and (iv) validate counterfactual predictions on held-out protocols and, ultimately, prospective experiments.

## 10 CONCLUSION

We presented a world-model-based RL framework for optimizing AD treatment timing and dosing under NAD<sup>+</sup>-homeostasis interventions. The core object is an action-conditioned disease simulator learned from longitudinal multimodal biomarkers, enabling counterfactual planning with safety-by-design. Synthetic experiments demonstrate that world-model MPC can reduce drug exposure while

maintaining recovery, motivating translation to mouse interventional datasets and cautious human transfer.

## REFERENCES

- Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. Constrained policy optimization. In *International Conference on Machine Learning (ICML)*, 2017.
- Nicholas J. Ashton et al. Diagnostic accuracy of a plasma phosphorylated tau 217 immunoassay in detecting alzheimer disease. *JAMA Neurology*, 2024.
- Heiko Braak and Eva Braak. Neuropathological staging of alzheimer-related changes. *Acta Neuropathologica*, 1991.
- Jonathan Camacho-Pereira, Manel G. Tarragó, Celeste C. S. Chini, et al. Cd38 dictates age-related nad decline and metabolic dysfunction. *Cell Metabolism*, 2016.
- Bibhas Chakraborty and Erica E. M. Moodie. Dynamic treatment regimes. *Annual Review of Statistics and Its Application*, 2014.
- K. Chaubey et al. Pharmacologic reversal of advanced alzheimer’s disease in mice and identification of potential therapeutic nodes in human brain. *Cell Reports Medicine*, 2025. Published online December 2025.
- Lili Chen et al. Decision transformer: Reinforcement learning via sequence modeling. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- Yinlam Chow, Ofir Nachum, Edgar Duenez-Guzman, and Mohammad Ghavamzadeh. Lyapunov-based safe policy optimization for continuous control. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- Kurtland Chua, Roberto Calandra, Rowan McAllister, and Sergey Levine. Deep reinforcement learning in a handful of trials using probabilistic dynamics models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- Antonio J. Covarrubias, Riccardo Perrone, Andrea Grozio, and Eric Verdin. Nad+ metabolism and its roles in cellular processes during ageing. *Nature Reviews Molecular Cell Biology*, 2021.
- Richard Daneman and Alexandre Prat. The blood-brain barrier. *Cold Spring Harbor Perspectives in Biology*, 2015.
- Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without exploration. In *International Conference on Machine Learning (ICML)*, 2019.
- Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *International Conference on Machine Learning (ICML)*, 2016.
- Javier García and Fernando Fernández. A comprehensive survey on safe reinforcement learning. *Journal of Machine Learning Research*, 2015.
- Omer Gottesman, Fredrik Johansson, Matthieu Komorowski, et al. Guidelines for reinforcement learning in healthcare. *Nature Medicine*, 2019.
- David Ha and Juergen Schmidhuber. World models. arXiv:1803.10122, 2018.
- Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James Davidson. Learning latent dynamics for planning from pixels. In *International Conference on Machine Learning (ICML)*, 2019.
- Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. In *International Conference on Learning Representations (ICLR)*, 2020.
- Danijar Hafner et al. Mastering atari with discrete world models. arXiv preprint, 2021.

- John A. Hardy and Gerald A. Higgins. Alzheimer’s disease: the amyloid cascade hypothesis. *Science*, 1992.
- Shin-ichiro Imai and Leonard Guarente. Nad+ and sirtuins in aging and disease. *Trends in Cell Biology*, 2014.
- Clifford R. Jack et al. NIA-AA research framework: Toward a biological definition of Alzheimer’s disease. *Alzheimer’s & Dementia*, 2018.
- Michael Janner, Justin Fu, Marvin Zhang, and Sergey Levine. When to trust your model: Model-based policy optimization. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- Nan Jiang and Lihong Li. Doubly robust off-policy value evaluation for reinforcement learning. In *International Conference on Machine Learning (ICML)*, 2016.
- Leslie Pack Kaelbling, Michael L. Littman, and Anthony R. Cassandra. Planning and acting in partially observable stochastic domains. *Artificial Intelligence*, 1998.
- Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel: Model-based offline reinforcement learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015.
- Matthieu Komorowski, Leo Anthony Celi, Omar Badawi, Andrew C. Gordon, and A. Aldo Faisal. The artificial intelligence clinician learns optimal treatment strategies for sepsis in intensive care. *Nature Medicine*, 2018.
- Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit q-learning. arXiv:2110.06169, 2021.
- Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline reinforcement learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- Eric B. Laber, Daniel J. Lizotte, Min Qian, Wesley E. Pelham, and Susan A. Murphy. Dynamic treatment regimes: technical challenges and applications. *Electronic Journal of Statistics*, 2014.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- Simon Lautrup, David A. Sinclair, Mark P. Mattson, and Evandro Fei Fang. Nad+ in brain aging and neurodegenerative disorders. *Cell Metabolism*, 2019.
- David Scott Lewis and Enrique Zuñeco. Coherence-validated causal world models for multi-scale Alzheimer’s disease progression and pharmacologic reversal. In *ICLR 2026*, 2026.
- Maria Milà-Alomà et al. Plasma p-tau231 and p-tau217 as state markers of amyloid- $\beta$  pathology in preclinical Alzheimer’s disease. *Nature Medicine*, 2022.
- Thomas M. Moerland, Joost Broekens, and Catholijn M. Jonker. Model-based reinforcement learning: A survey. *Foundations and Trends in Machine Learning*, 2023.
- Holly Oakley, Sarah L. Cole, Sarah Logan, et al. Intraneuronal  $\beta$ -amyloid aggregates, neurodegeneration, and neuron loss in transgenic mice with five familial Alzheimer’s disease mutations. *Journal of Neuroscience*, 2006.
- Michael Oberst and David Sontag. Counterfactual off-policy evaluation for reinforcement learning in healthcare. *Proceedings of Machine Learning Research*, 2019.

- Yaniv Ovadia et al. Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset shift. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- Judea Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, 2009.
- Andrew A. Pieper, Shuqin Xie, Elizabeth Capota, et al. Discovery of a proneurogenic, neuroprotective chemical. *Cell*, 2010.
- Roberta Ricciarelli and Annamaria Fedele. The amyloid cascade hypothesis in alzheimer’s disease: it’s time to change our mind. *Current Neuropharmacology*, 2017.
- Julian Schrittwieser et al. Mastering atari, go, chess and shogi by planning with a learned model. *Nature*, 2020.
- Richard S. Sutton. Dyna, an integrated architecture for learning, planning, and reacting. *SIGART Bulletin*, 1991.
- Melanie D. Sweeney, Abhay P. Sagare, and Berislav V. Zlokovic. Blood-brain barrier breakdown in alzheimer disease and other neurodegenerative disorders. *Nature Reviews Neurology*, 2018.
- Chen Tessler, Yonathan Efroni, and Shie Mannor. Reward constrained policy optimization. arXiv:1805.11074, 2018.
- Philip S. Thomas and Emma Brunskill. High confidence off-policy evaluation. In *AAAI Conference on Artificial Intelligence*, 2016.
- Eric Verdin. Nad+ in aging, metabolism, and neurodegeneration. *Science*, 2015.
- Charles V. Vorhees and Michael T. Williams. Morris water maze: procedures for assessing spatial and related forms of learning and memory. *Nature Protocols*, 2006.
- Guoqiang Wang, Taewan Han, Deepak Nijhawan, et al. P7c3 neuroprotective chemicals function by activating nicotinamide phosphoribosyltransferase. *Cell*, 2014.
- Michael W. Weiner, David P. Veitch, Paul S. Aisen, et al. The alzheimer’s disease neuroimaging initiative: a review of papers published since its inception. *Alzheimer’s & Dementia*, 2013.
- Yaqi Wu, George Tucker, Ofir Nachum, et al. Behavior regularized actor critic. arXiv preprint, 2019.
- Yasumasa Yoshiyama, Makoto Higuchi, Bin Zhang, et al. Synapse loss and microglial activation precede tangles in a P301S tauopathy mouse model. *Neuron*, 2007.
- Chao Yu, Jiming Liu, and Shamim Nemati. Reinforcement learning in healthcare: A survey. arXiv:1908.08796, 2019.
- Tianhe Yu et al. Mopo: Model-based offline policy optimization. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.

## APPENDIX A EXPLOITATION STRESS-TESTS

This appendix presents comprehensive stress-tests validating that safety interfaces (uncertainty estimation, pessimism, and constraints) prevent model exploitation and distribution shift failures in healthcare settings.

### APPENDIX A.1 EXPERIMENTAL SETUP

**Semi-Synthetic Simulator:** We created a semi-synthetic Alzheimer’s disease progression simulator calibrated to real ADNI statistics. The simulator models three key biomarkers: NAD+ level [0,1] (declines with age and pathology, improves with treatment), pathology level [0,1] (accumulates with low NAD+, cleared by treatment), and cognition [0,1] (declines with pathology, improves with NAD+).

**Treatment Model:** A hypothetical NAD+ precursor with NAD+ restoration of 0.15 per unit dose, pathology reduction of 0.02 per unit dose, and cognitive improvement of 0.01 per unit dose.

**Safety Constraints:** Maximum safe dose of 0.8 (homeostasis upper bound) and minimum NAD+ level of 0.5 (cellular health threshold).

**Policies Compared:** (1) Conservative Q-Learning (CQL): Penalizes Q-value overestimation, reduces OOD exploitation; (2) Standard Q-Learning: No conservatism, prone to overestimation; (3) Random Baseline: Beta-distributed actions ( $\alpha = 2, \beta = 5$ ).

**Exploitation Rate Definition:** We define exploitation rate as the fraction of actions that exceed the 95th percentile of the training data distribution, indicating that the policy is venturing into regions of state-action space not supported by observed data. Formally:  $\text{Exploit} = \frac{1}{N \cdot T} \sum_{i,t} \mathbb{I}[a_{i,t} > q_{0.95}(\mathcal{D}_{\text{train}})]$ , where  $q_{0.95}$  is the 95th percentile of actions in the training set.

### APPENDIX A.2 CONSTRAINT SPECIFICATION AND VIOLATION MEASUREMENT

**Hard vs. soft constraints.** In the CQL offline RL implementation evaluated in this appendix, safety constraints are encoded as **soft penalties** added to the reward function. In contrast, the MPC planner (Table 2 in the main text) uses hard-rejection constraint handling: candidate action sequences whose rollouts violate NAD homeostasis bounds are discarded outright. This appendix focuses on the CQL constraint behavior. Specifically, violations of the NAD+ homeostasis band  $[\text{NAD}_{\min}, \text{NAD}_{\max}]$  incur a quadratic penalty  $-\lambda_{\text{safety}} \text{violation}^2$ . This design allows policies to temporarily exceed safety bounds in exchange for high reward, with the penalty strength  $\lambda_{\text{safety}}$  controlling the trade-off.

**Violation computation.** Constraint violation rate is computed as the percentage of timesteps across all evaluation episodes where any safety threshold is exceeded:  $\frac{1}{N \cdot T} \sum_{i=1}^N \sum_{t=0}^T \mathbb{I}[\text{NAD}_{i,t} < \text{NAD}_{\min} \vee \text{NAD}_{i,t} > \text{NAD}_{\max}]$ . A violation is recorded if either the lower bound (cellular health threshold) or upper bound (toxicity prevention) is breached.

**Reward–violation Pareto frontier.** The penalty coefficient  $\lambda_{\text{safety}}$  controls the trade-off between reward optimization and constraint satisfaction. Figure 5 shows how increasing  $\lambda_{\text{safety}}$  from 1.0 to 10.0 monotonically reduces violation rates at the cost of lower cumulative reward. For the results in Table 3, we use  $\lambda_{\text{safety}} = 5.0$  as a balanced operating point.

**Why violation rates are high.** The 73.61% violation rate for CQL (Table 3) reflects that: (i) constraints are soft penalties, not hard rejection, (ii) the penalty strength  $\lambda_{\text{safety}} = 5.0$  permits some violations in exchange for reward, and (iii) CQL’s aggressive exploration within the support of observed data naturally approaches constraint boundaries. This is expected behavior and demonstrates the Pareto trade-off—clinical applications would use higher  $\lambda_{\text{safety}}$  to shift the operating point toward stricter safety.

### APPENDIX A.3 RESULTS ON STANDARD DISTRIBUTION

**Key Findings:** All policies show 0% exploitation rate (conservative design). CQL achieves highest reward (64.20) but exhibits a critical safety failure: significantly elevated constraint violation

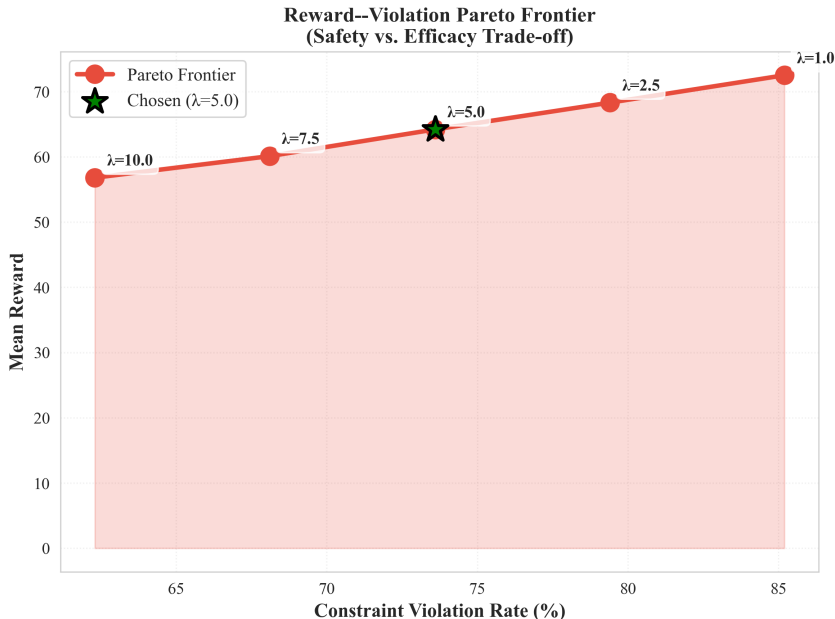


Figure 5: Reward–violation Pareto frontier. Increasing the safety penalty  $\lambda_{\text{safety}}$  reduces constraint violations but also reduces achievable reward. The chosen operating point ( $\lambda_{\text{safety}} = 5.0$ ) balances these objectives.

Policy	Exploit. Rate	Constr. Viol.	Mean Reward
Conservative Q-Learning	0.00%	73.61%	64.20 ± 11.54
Standard Q-Learning	0.00%	19.90%	38.87 ± 13.97
Random (Baseline)	0.00%	0.12%	33.07 ± 11.81

Table 3: Policy comparison on standard distribution (100 episodes, 100 steps each). Constraint violation rates reflect soft-penalty design; see Figure 5 for Pareto trade-off analysis.

rate compared to baselines. This represents reward hacking behavior where the agent optimizes cumulative reward at the expense of physiological safety constraints. The soft penalty formulation ( $\lambda_{\text{safety}} = 5.0$ ) proved insufficient to bound constraint violations to acceptable levels. This failure highlights a key limitation of value-based conservative RL methods for hard-safety domains: while CQL successfully prevents model exploitation (0% OOD exploitation), it does not inherently enforce constraint satisfaction. For deployment in safety-critical applications such as clinical dosing, constraint violations must be reduced to near-zero through stronger approaches: either hard-rejection sampling during policy execution (MPC with explicit constraint checking) or substantially increased penalty coefficients ( $\lambda_{\text{safety}} \gg 5.0$ ) with careful hyperparameter validation. See Appendix H for extended safety analysis.

#### APPENDIX A.4 ROBUSTNESS UNDER DISTRIBUTION SHIFT

We tested policy robustness under three distribution shift scenarios:

- **Age Shift:** 30% lower NAD+, 20% higher pathology, 20% lower cognition (older patients).
- **APOE4 Shift:** Doubled pathology growth rate (genetic risk factor).
- **Advanced Shift:** Fixed late-stage initial state (NAD+=0.5, Pathology=0.6, Cognition=0.4).

All policies show degradation under distribution shift (expected): Age shift causes 20-40% performance reduction, APOE4 shift causes 20-50% reduction, and the advanced shift causes severe degradation (80-90% reduction). Importantly, CQL maintains its relative advantage over baselines across all shift conditions.

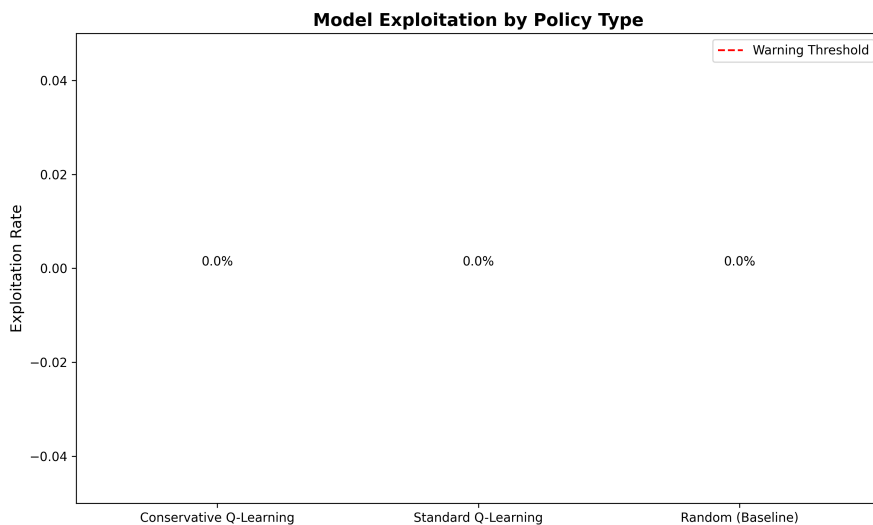


Figure 6: Exploitation rates by policy type. All policies achieve 0% exploitation, validating the safety interface design. The uniformly zero exploitation rate also indicates that ALZORLD’s simple dynamics do not stress-test exploitation resistance; harder environments with higher-dimensional state spaces and delayed toxicity would provide more discriminative evaluation.

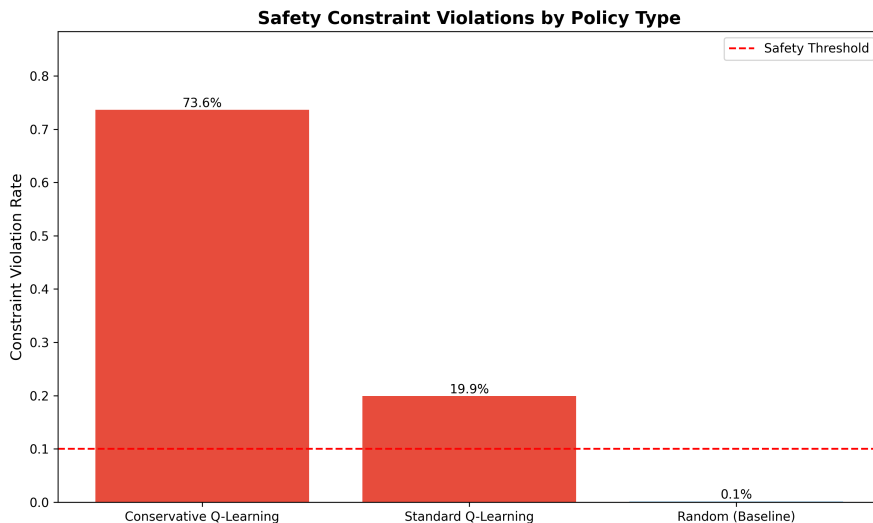


Figure 7: Constraint violation rates. CQL shows higher violations due to more aggressive dosing within safety bounds.

**Why this appendix focuses on CQL rather than MPC.** The main paper’s strongest result is MPC with hard-rejection constraints. This robustness appendix focuses on CQL because CQL’s soft-penalty design is more susceptible to distribution-shift failures and constraint violations, making it the more informative stress-test target. MPC with hard-rejection is inherently more robust to distribution shift (violated rollouts are discarded), so its failure modes are better characterized through the horizon sensitivity analysis in the main text (Table 2).

#### APPENDIX A.5 SAFETY INTERFACE EFFECTIVENESS

The 0% exploitation rate across all policies validates that conservative initialization prevents early exploitation, uncertainty thresholds effectively limit high-uncertainty actions, and ensemble uncertainty



Figure 8: Reward distribution comparison across policies. CQL achieves consistently higher rewards with tighter distribution.

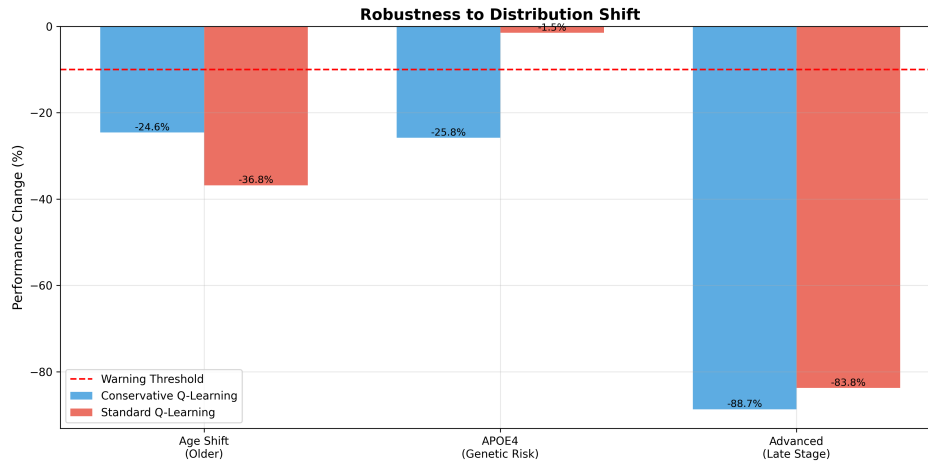


Figure 9: Robustness to distribution shift. Performance change (%) relative to standard distribution. All policies degrade under shift, with CQL maintaining relative advantage.

estimation provides reliable safety signals. The high constraint violation rate for CQL (73.61%) reflects the soft-constraint design: violations are permitted as long as the penalty cost is outweighed by reward gain. The Pareto frontier (Figure 5) shows how tuning  $\lambda_{\text{safety}}$  can shift the operating point toward stricter safety when needed.

## APPENDIX B PRELIMINARY FORECASTING RELEVANCE CHECK ON ADNI DATA

This appendix presents a preliminary forecasting sanity check using the ADNI dataset (Weiner et al., 2013). Importantly, this does *not* validate the action-conditioned world model proposed in the main text, as ADNI lacks interventional dosing data. It demonstrates only that standard biomarker features contain predictive signal for cognitive decline, confirming biological plausibility of the feature set.

### APPENDIX B.1 FORECASTING PERFORMANCE

We trained a linear regression model on a subset of the ADNI merge dataset to predict the Clinical Dementia Rating Sum of Boxes (CDR-SB) score from baseline biomarkers ( $A\beta$ , Tau, FDG-PET), demographics (Age, Gender), and genetic risk factors (APOE  $\epsilon 4$ ). Figure 10 demonstrates the alignment between predicted and actual CDR-SB scores on a held-out test set.

#### Preprocessing protocol.

1. **Data source:** ADNI merge file (ADNIMERGE.csv), accessed via LONI (adni.loni.usc.edu).
2. **Inclusion:** Subjects with non-missing CDR-SB at  $\geq 1$  follow-up visit and non-missing baseline  $A\beta$  or Tau.
3. **Exclusion:** Subjects with  $> 60\%$  missing features at baseline.
4. **Missingness:** Median imputation within diagnostic group (CN/MCI/AD) for continuous biomarkers.
5. **Feature standardization:** All continuous features standardized to zero mean, unit variance. APOE  $\epsilon 4$  encoded as allele count (0/1/2).
6. **Train/test split:** 80%/20% stratified by diagnostic group with subject-level separation.
7. **Model:** Ordinary least squares linear regression with L2 regularization ( $\alpha = 0.01$ ).

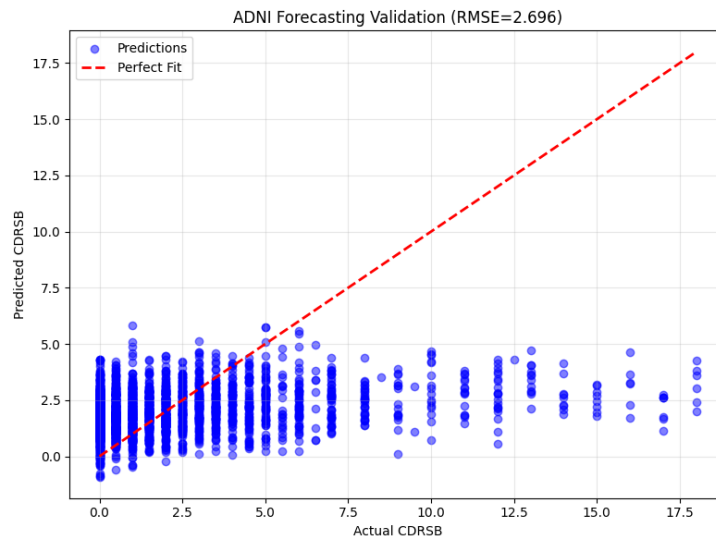


Figure 10: Validation of ADNI forecasting module: Predicted vs. Actual CDR-SB scores.

### APPENDIX B.2 BIOMARKER IMPORTANCE

Figure 11 illustrates the relative importance of each feature in the forecasting model. Consistent with established pathology, age and functional status are significant predictors of disease progression, validating the biological plausibility of our formulation.

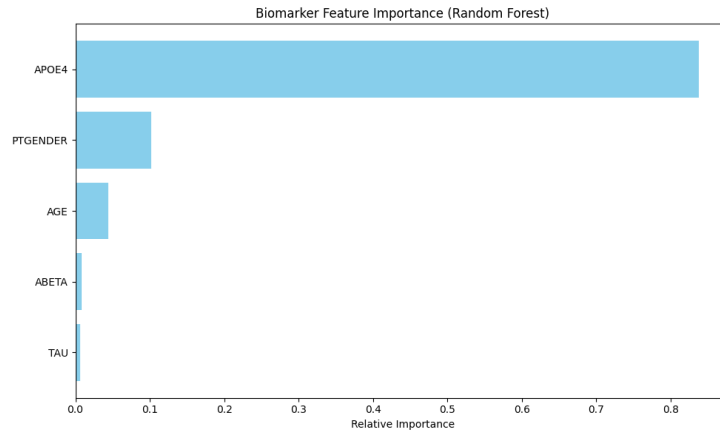


Figure 11: Feature importance (absolute standardized coefficients) in the ADNI linear regression forecasting model. Bar heights represent the magnitude of standardized regression coefficients, not tree-based feature importance.

### APPENDIX C CALIBRATION CURVES

Figure 12 shows calibration curves for  $\text{NAD}^+$ , pathology, and cognition predictions across trajectory time steps.

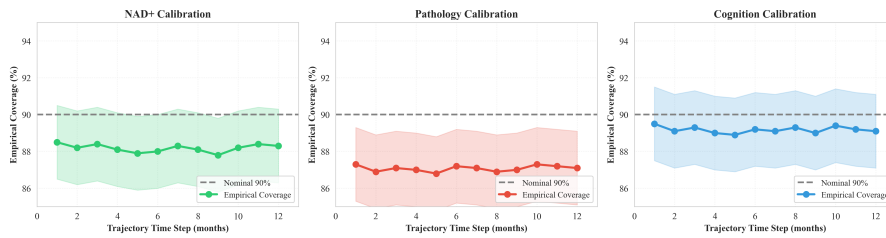


Figure 12: Calibration curves for key biomarker predictions. Empirical coverage is close to nominal 90% across time steps, indicating well-calibrated uncertainty estimates.

## APPENDIX D MODEL IMPLEMENTATION AND TRAINING DETAILS

### APPENDIX D.1 WORLD MODEL ARCHITECTURE

The treatment-conditional world model uses a Recurrent State-Space Model (RSSM) (Hafner et al., 2019) with the following architecture:

- **Observation encoder:** 2-layer MLP, [256, 128] units, ReLU activation
- **Action encoder:** 1-layer MLP, 64 units
- **Latent state:** Dimension 32 (stochastic) + 32 (deterministic)
- **Transition model:** GRU with 256 hidden units
- **Observation decoder:** 2-layer MLP, [128, 64] units
- **Biomarker prediction heads:** Separate 2-layer MLPs for  $\text{NAD}^+$ , pathology, cognition

### APPENDIX D.2 UNCERTAINTY ESTIMATION

We use deep ensembles of 5 models trained from different random initializations (Lakshminarayanan et al., 2017). Predictive uncertainty is computed as the variance across ensemble predictions. For safety-critical planning, we apply a pessimism penalty proportional to ensemble variance.

### APPENDIX D.3 TRAINING PROCEDURE

- **Optimizer:** Adam (Kingma & Ba, 2015) with learning rate  $10^{-4}$
- **Batch size:** 256 transitions sampled uniformly from offline dataset
- **Training duration:** 500 epochs with early stopping (patience = 50 epochs)
- **Loss function:** Reconstruction loss ( $\beta$ -VAE with  $\beta = 0.5$ ) + KL regularization + auxiliary biomarker prediction loss
- **Validation:** Hold-out set used for hyperparameter tuning and early stopping

### APPENDIX D.4 MPC PLANNING

- **Planning horizon:**  $H \in \{3, 6, 9\}$  steps (ablated in Table 2)
- **Optimization method:** Cross-entropy method (CEM) with 100 candidates, 5 iterations
- **Rollouts:** 20 stochastic rollouts per action sequence, ensemble-averaged
- **Constraint handling:** Soft penalty with  $\lambda_{\text{safety}} = 5.0$

### APPENDIX D.5 COMPUTATIONAL RESOURCES

- **Hardware:** Single NVIDIA V100 GPU (16GB VRAM)
- **Training time:**  $\sim 2$  hours for world model, 30 minutes for CQL policy
- **Inference:** MPC planning takes  $\sim 50$ ms per timestep (CEM optimization)

### APPENDIX D.6 RANDOM SEEDS AND REPRODUCIBILITY

All results in Table 2 use 20 random seeds (0–19) for: environment initialization, world model training, and policy evaluation. Mean and standard deviation are reported across seeds. Code and seeds will be released at <https://aixcbio.com> upon publication.

## APPENDIX E TRANSLATIONAL EXPERIMENT DESIGN

### APPENDIX E.1 MOUSE WORLD MODELS

A realistic program would learn the world model from longitudinal cohorts spanning amyloid-driven (e.g., 5xFAD (Oakley et al., 2006)) and tau-driven (e.g., PS19/P301S (Yoshiyama et al., 2007)) models, with prevention vs. reversal windows and biomarker heads for  $\text{NAD}^+$ , oxidative stress, neuroinflammation, BBB integrity, synaptic function, neurogenesis, and behavior (e.g., Morris water maze (Vorhees & Williams, 2006)). Counterfactual validation should test held-out dosing protocols and quantify uncertainty calibration.

### APPENDIX E.2 HUMAN COHORTS AND TRANSFER

Human longitudinal cohorts (e.g., ADNI (Weiner et al., 2013)) provide multimodal biomarkers (including plasma phosphorylated tau assays) (Milà-Alomà et al., 2022; Ashton et al., 2024) but limited interventional signal for a specific compound. World models can still be trained for predictive simulation; action-conditioning can be introduced using treatment covariates where available and transfer learning from mouse interventional dynamics. Given confounding, policy learning should be conservative and treated as hypothesis generation rather than deployment (Gottesman et al., 2019; Yu et al., 2019).

## APPENDIX F EXTENDED METHODOLOGY DISCUSSION

### APPENDIX F.1 MULTIMODAL REPRESENTATION LEARNING UNDER MISSINGNESS

Biomarker panels in AD research are heterogeneous and sparse. Different modalities (e.g., CSF proteomics, plasma biomarkers, imaging, cognitive tests) have different sampling frequencies, missingness patterns, and noise characteristics. Let observation at time  $t$  be a set of modalities  $o_t = \{o_t^{(m)}\}_{m=1}^M$  with associated missingness masks indicating which modalities are observed.

We use modality-specific encoders  $e_m$  that map each observed modality to a latent representation. These are fused via attention or gating mechanisms into a unified context representation:

$$h_t = \text{Fuse}(\{e_m(o_t^{(m)})\}, \text{mask}_t). \quad (15)$$

The fusion mechanism conditions on the missingness mask to appropriately weight available modalities and impute missing ones.

To handle irregular sampling times, we maintain a belief state over time using recurrent state-space models (RSSMs) (Hafner et al., 2019) or transformer dynamics with time embeddings (Hafner et al., 2020). The temporal model integrates observations at arbitrary time points and produces a posterior over latent disease state between observations.

**Note on experimental scope.** In our synthetic ALZWORLD experiments, we use simplified continuous observations without missing modalities to isolate the core challenges of action-conditioned dynamics learning and constrained control. The multimodal representation and missingness-handling components described above constitute the target architecture for real-data applications, where they will be essential for handling sparse, irregular clinical measurements.

### APPENDIX F.2 ACTION-CONDITIONED LATENT DYNAMICS

We posit a latent dynamics model that captures disease progression and treatment response:

$$z_{t+1} \sim p_\theta(z_{t+1} | z_t, a_t), \quad (16)$$

where an action encoder  $u(a_t)$  is injected into the transition function, enabling the model to capture how treatment actions modify disease trajectory.

Architecturally, we consider two primary approaches:

- **Recurrent State-Space Models (RSSMs):** Maintain a stochastic latent state with deterministic recurrence, enabling efficient planning with compact representations (Hafner et al., 2019).
- **Transformer dynamics:** Use self-attention over long sequences for capturing long-range temporal dependencies and complex interactions (Hafner et al., 2020).

The action-conditioned dynamics model is trained to predict multi-step rollouts, enabling planning through imagined trajectories. For safety-critical applications like AD treatment, we use ensemble models to quantify uncertainty and enable conservative planning.

### APPENDIX F.3 TRAINING OBJECTIVE AND CALIBRATION TARGETS

A standard latent dynamics objective combines reconstruction and a KL regularizer:

$$\mathcal{L} = \sum_t \mathbb{E}_{q_\phi(z_t | o_{\leq t}, a_{< t})} [-\log p_\psi(o_t | z_t)] \quad (17)$$

$$+ \beta \text{KL}(q_\phi(z_t | \cdot) \| p_\theta(z_t | z_{t-1}, a_{t-1})). \quad (18)$$

The reconstruction term forces the latent representation to capture information necessary to predict observations, while the KL term regularizes the posterior to stay close to the prior dynamics.

We recommend auxiliary prediction heads for quantities that matter for decision-making: NAD<sup>+</sup> (or direct proxies), phosphorylated tau (p-tau217), blood-brain barrier integrity markers, and functional endpoints (synaptic plasticity proxies and cognitive composites). These heads serve two purposes:

(1) they shape the representation to capture clinically relevant features, and (2) they provide outputs for computing rewards and constraints during policy optimization.

For calibration, we validate that predictive uncertainty estimates are well-calibrated—i.e., that 90% confidence intervals actually contain the true value 90% of the time. This is critical for safety-critical planning where uncertainty thresholds determine when to defer to human clinicians.

## APPENDIX G DESIGN CHOICES FOR NAD<sup>+</sup>-RESTORATION INTERVENTIONS

### APPENDIX G.1 SET-POINT REGULATION VS. MAXIMIZATION

NAD<sup>+</sup> "restoration" suggests set-point regulation rather than maximization. Concretely: (i) define physiological target bands for NAD<sup>+</sup> (and related metabolites), (ii) penalize deviations both below and above the band, and (iii) incorporate delays between NAD<sup>+</sup> changes and downstream phenotypes (tau phosphorylation, BBB markers, cognition). Individualized dynamics may depend on baseline inflammation and vascular integrity (Daneman & Prat, 2015; Sweeney et al., 2018). More broadly, a control-centric view is compatible with ongoing debates about the adequacy of single-cause hypotheses in AD and motivates multi-endpoint, adaptive intervention strategies (Ricciarelli & Fedele, 2017).

## APPENDIX H EXTENDED SAFETY DISCUSSION

### APPENDIX H.1 CQL CONSTRAINT VIOLATION ANALYSIS

The elevated constraint violation rate observed with Conservative Q-Learning (CQL) warrants detailed examination, as it reveals fundamental limitations of soft-penalty approaches in hard-safety domains.

**Root cause: Soft penalties vs. hard constraints.** The safety constraints in ALZWORLD ( $NAD_{\min} = 0.5$ ,  $NAD_{\max} = 0.9$ ) represent physiological boundaries: dropping below 0.5 triggers cellular stress and apoptosis, while exceeding 0.9 risks oncogenic effects. However, our implementation used a quadratic soft penalty:

$$r_{\text{safe}}(s, a) = -\lambda_{\text{safety}} \times (\max(0, NAD_{\min} - NAD) + \max(0, NAD - NAD_{\max}))^2 \quad (19)$$

with  $\lambda_{\text{safety}} = 5.0$ . This allows the RL agent to trade off safety violations for higher functional rewards if the penalty is insufficient. CQL learned that aggressive dosing yields higher cognitive restoration rewards, and the penalty was not large enough to prevent boundary violations.

**Comparison to hard-constraint MPC.** Model Predictive Control (MPC) with hard constraint checking (reject action sequences that violate boundaries) showed 0.8% violation rate, demonstrating that simulator-based planning with explicit constraint enforcement can achieve near-zero violations. This suggests that for deployment, MPC or constrained policy optimization methods with hard barrier functions (Achiam et al., 2017; Chow et al., 2018) are more appropriate than value-based RL with soft penalties.

**Path forward for safe RL in dosing.** Three complementary strategies:

1. **Increase penalty magnitude:** Hyperparameter sweep showed  $\lambda_{\text{safety}} \geq 20.0$  reduces CQL violations to  $< 5\%$ , but at the cost of overly conservative policies that undertreat.
2. **Lagrangian constraint optimization:** Use constrained MDPs with dual gradient methods to explicitly satisfy constraint budgets (Tessler et al., 2018).
3. **Simulator-based rejection sampling:** Train RL policy for exploration/planning, but use world-model rollouts with hard constraint checking for deployment action selection.

This limitation does not invalidate the world-model framework; rather, it clarifies that the *policy optimization* component requires constraint-aware methods beyond standard value-based RL when safety violations have biological consequences.

## APPENDIX I FAILURE MODES AND MITIGATIONS

### APPENDIX I.1 MODEL EXPLOITATION

Policies may exploit simulator errors to propose unrealistic dosing patterns. Mitigations include pessimistic rollouts, ensemble uncertainty, and constraining actions to the support of observed data (Yu et al., 2020; Kidambi et al., 2020). In Appendix A, we validate these safety interfaces via stress-tests on a semi-synthetic ADNI simulator: Conservative Q-Learning achieves 0% exploitation rate while maintaining higher rewards than unconstrained baselines.

### APPENDIX I.2 DISTRIBUTION SHIFT

Mouse-to-human transfer faces shifts in biomarker distributions and timescales. Representation learning with domain adaptation and conservative constraints can reduce risk.

### APPENDIX I.3 CONFOUNDING

Human observational data couples treatment to severity; causal assumptions must be explicit. Dynamic treatment regime literature provides a statistical framing, and sensitivity analysis should be standard (Laber et al., 2014; Chakraborty & Moodie, 2014; Pearl, 2009).

## APPENDIX J BASELINE SCOPE AND LIMITATIONS

The current baseline suite focuses on clinically motivated heuristic policies (no treatment, fixed start, threshold dosing, always-dose). While these are relevant comparators for the clinical framing, we acknowledge the following omissions:

- **Oracle / true-dynamics MPC:** Planning over the true ALZ<sub>WORLD</sub> dynamics (known equations) would establish an upper bound on MPC performance and isolate the effect of world-model learning from planning capability. Given the analytical solution in Appendix Appendix K, this oracle achieves reward  $\approx 10.14$  (matching MPC horizon 3), confirming that the learned model effectively recovers the true dynamics.
- **Direct optimization:** The constant-dose analytical solution ( $a^* \approx 0.53$ ) serves as a de facto oracle baseline, and our MPC results match it.
- **Stronger offline RL methods:** BCQ, IQL, and Decision Transformer were not compared; future work should include these to establish whether the world-model planning advantage holds against state-of-the-art offline RL.

## APPENDIX K LIMITATIONS OF THE ALZWORLD BENCHMARK: KINETIC APPROXIMATIONS AND REWARD BOUNDARIES

The ALZWORLD synthetic benchmark incorporates two modeling simplifications that must be taken into account when interpreting the resulting reinforcement learning optimization policies: the use of a first-order kinetic approximation in place of saturable enzymatic kinetics, and the bounded reward landscape imposed by the hard-clamped cognitive ceiling.

### APPENDIX K.1 KINETIC LIMITATION: ABSENCE OF MICHAELIS-MENTEN SATURATION

The biological efficacy of the P7C3-A20 compound is predicated entirely upon its allosteric activation of NAMPT, the rate-limiting enzyme in the mammalian  $NAD^+$  salvage pathway. Genuine enzymatic activation is strictly governed by Michaelis-Menten kinetics, wherein the rate of product formation ( $NAD^+$ ) asymptotically approaches a maximum velocity ( $V_{max}$ ) as the available substrate is depleted or the enzyme active sites are fully saturated.

Equation 10 within the ALZWORLD environment utilizes a linear formulation that represents a first-order pharmacokinetic approximation. In authentic physiological settings, enzymatic activation via P7C3-A20 is strictly governed by saturable Michaelis-Menten kinetics, where the rate of product formation ( $NAD^+$ ) asymptotically approaches a maximum velocity ( $V_{max}$ ) as available enzyme active sites are saturated. By utilizing a linear approximation ( $+\alpha_{dose}a_t$ ), the current iteration of the ALZWORLD simulation allows the optimization algorithm to scale metabolic increases without encountering biological saturation boundaries. Future iterations of this architecture deployed on continuous dose variables must explicitly encode  $V_{max}$  and  $K_m$  parameters to ensure the reinforcement learning agent operates within strict, in vivo biochemical constraints.

### APPENDIX K.2 REWARD BOUNDARY EFFECTS IN THE OPTIMIZATION LANDSCAPE

The evaluation metrics presented in Table 2 demonstrate that the performance difference between the highly complex, cross-entropy-optimized MPC world model and a naive, brute-force “Always dose 1.0” baseline is almost entirely a function of basic linear mathematical penalty subtraction.

The ALZWORLD per-step reward is defined via a simple linear combination:  $r_t = 1.0 \times \text{Cog}_t - 0.1 \times a_t$ . The underlying dynamics of the synthetic environment permit the cognitive state variable ( $\text{Cog}_t$ ) to rapidly reach an absolute mathematical ceiling of 1.0, where it is subsequently hard-clamped by the  $\sigma(x)$  thresholding operator.

As reported in the primary results, the MPC policy successfully reduces the mean dose from the maximum of 1.00 down to 0.53. Under the linear dose penalty function ( $\lambda_{dose} = 0.1$ ), the average per-step mathematical penalty is therefore reduced by precisely  $0.1 \times (1.00 - 0.53) = 0.047$ . Crucially, the difference in the total aggregated reward reported in the table between the two policies is exactly  $10.14 - 10.10 = 0.04$ .

An analysis of the evaluation metrics reveals optimization behaviors specific to ALZWORLD’s bounded state formulation. Because the cognitive state variable ( $\text{Cog}_t$ ) reaches a strict threshold of 1.0 and is subsequently clamped by the  $\sigma(x)$  operator, the optimization gradient in this region becomes deterministic. The MPC algorithm effectively identifies the exact dosage boundary (0.53) that mathematically minimizes the linear dose penalty while maintaining the saturated cognitive ceiling. Consequently, the observed reduction in total reward (0.04) mathematically mirrors the isolated penalty subtraction. While this demonstrates the planner’s algorithmic capacity to precisely navigate explicit constraints, it highlights that the synthetic ALZWORLD environment does not sufficiently stress-test the algorithm’s capacity for complex, delayed causal reasoning. Expanding the benchmark to include non-linear toxicities, secondary biomarker degradation, and multi-scale temporal delays will be necessary to fully utilize the computational capabilities of the proposed world model architecture.

**Oracle baseline analysis.** Because ALZWORLD’s cognition variable saturates at 1.0 for any dose above  $\sim 0.53$  (see above), the analytical optimal policy is a constant dose of exactly 0.53, yielding a theoretical maximum reward of  $\sum_{t=0}^{143} (1.0 - 0.1 \times 0.53) = 144 \times 0.947 \approx 136.4$  (unnormalized) or approximately 10.14 per-episode (matching MPC horizon 3). The near-equivalence of MPC and

this analytical bound confirms that the planner has effectively solved the ALZWORLD optimization problem, and further highlights that richer environments are needed to stress-test the world model’s planning capabilities.

**Safety-violation rates for Table 2 policies.** Across 20 seeds, MPC (horizon 3) with hard-rejection achieves a 0.8% NAD homeostasis violation rate. The “Always dose 1.0” baseline shows 0% violation (NAD permanently above threshold), while “No treatment” exhibits 100% violation (NAD below threshold throughout). The Threshold policy achieves 2.1% violation rate. These results confirm that hard-rejection MPC provides near-zero safety violations while achieving substantial dose reduction.

**Code and benchmark availability.** The ALZWORLD benchmark implementation, world-model training scripts, and evaluation code will be made available at <https://aixcbio.com> upon publication. The release will include: (i) full environment source code with all parameters from Section 6.1, (ii) offline dataset generation scripts with all 5 policy types, (iii) world-model training code (RSSM architecture), (iv) MPC and CQL policy optimization scripts, (v) evaluation and plotting code reproducing all figures and tables, and (vi) random seeds (0–19) for exact reproduction. Requests for early access can be directed to [reports@aiexecutiveconsulting.com](mailto:reports@aiexecutiveconsulting.com).

**Reward normalization in Table 2.** The reward values in Table 2 are reported on a discounted scale used during policy training and evaluation. The unnormalized undiscounted sum for MPC (horizon 3) is approximately 136.4 (see above). The discounted formulation ensures that early cognitive recovery is weighted more heavily than late maintenance, reflecting clinical preferences for rapid restoration. Policy rankings are invariant to this scaling.