

---

# Learning to Orchestrate Heterogeneous Agents under Uncertainty

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Adaptive orchestration of heterogeneous agents requires making sequential dele-  
2 gation decisions under uncertain and evolving agent behaviour, e.g., coordinating  
3 specialised AI models with varying reliability, cost, and response quality. While  
4 prior work on agent orchestration focuses on performance or cost, uncertainty in  
5 agent reliability and output distributions is typically not modelled explicitly at the  
6 orchestration level. In this work, we study the problem of adaptive orchestration  
7 of heterogeneous agents under uncertainty, where a meta-controller must decide  
8 when to delegate to a single agent, accounting for reliability, cost, and uncer-  
9 tainty. We propose BOT-Orch, a lightweight framework that casts orchestration  
10 as a bandit problem over agents, regularized by OT distances between agent out-  
11 put distributions and task-specific reference distributions. We show that the regu-  
12 larised orchestration enjoys  $\mathcal{O}(\sqrt{T})$  regret under standard assumptions, and prov-  
13 ably induces preference ordering among agents with identical mean rewards but  
14 differing distributional alignment. Empirically, we demonstrate that BOT-Orch  
15 outperforms standard bandit and heuristic baselines in synthetic but adversarial  
16 task allocation settings with heterogeneous, non-i.i.d. agent behaviour.

## 17 1 Introduction

18 Incorporating uncertainty when coordinating multiple decision-making entities is a fundamental  
19 challenge across machine learning and autonomous systems [Wurman et al., 2007, Olfati-Saber  
20 et al., 2007, Shalev-Shwartz, 2012]. In real-world environments, where stochasticity is dominant  
21 (e.g., partial and noisy observations), variability in agent capabilities makes centralized control im-  
22 practical and can break classical solutions. Classical multi-agent frameworks (e.g., Dec-POMDPs,  
23 stochastic games, DCOPs) explicitly model uncertainty in multi-agent decision-making but often as-  
24 sume known system models, have limited scalability, or simplify agent interactions, restricting appli-  
25 cability in complex coordination settings [Bernstein et al., 2002, Oliehoek and Amato, 2016, Hansen  
26 et al., 2004, Heifetz et al., 2006, Shoham and Leyton-Brown, 2008, Fioretto et al., 2018]. In practice,  
27 heterogeneity in agent reliability, operational cost, and resource constraints further complicates coordi-  
28 nation, especially at scale where performance, risk, and cost must be jointly managed [Rizk et al.,  
29 2019, Arjun et al., 2025]. Surveys in distributed decision-making and heterogeneous multi-agent  
30 coordination consistently identify these issues as key barriers to scalable deployment [Olfati-Saber  
31 et al., 2007, Rizk et al., 2019].

32 Multi-agent reinforcement learning (MARL) provides tools for learning cooperative behaviors under  
33 partial observability and decentralised execution [Buşoniu et al., 2008, Hernandez-Leal et al., 2019].  
34 Value decomposition and centralized training with decentralized policies can maximize joint rewards  
35 even with distinct roles and observations [Sunehag et al., 2018, Rashid et al., 2018]. However, many

36 MARL methods implicitly assume similar reliability and reaction costs, and rarely address adaptive  
37 selection among heterogeneous agents of uncertain outputs.

38 Modern AI systems and multi-agent orchestration frameworks sharpen these challenges. In plan-  
39 ning, reasoning, and AI agent teams, systems increasingly rely on specialized agents with differ-  
40 ent expertise, reliability, and computational cost. Effective orchestration requires deciding which  
41 agent(s) to invoke and when to combine multiple predictions for robustness. Recent work shows  
42 that dynamic selection and composition conditioned on context can outperform static pipelines [Park  
43 et al., 2023, Liang et al., 2023, Cheng et al., 2023].

44 A natural abstraction for sequential decision-making under uncertainty is the multi-armed bandit  
45 (MAB) framework, capturing the exploration–exploitation trade-off [Lattimore and Szepesvári,  
46 2020, Bubeck and Cesa-Bianchi, 2012]. Multi-agent extensions study cooperation, regret minimiza-  
47 tion, and communication among learners [Landgren et al., 2016, Gupta et al., 2021]. Yet most bandit  
48 models treat agents as interchangeable up to mean reward and ignore query costs, limiting their  
49 suitability for orchestration with heterogeneous reliability, non-stationarity, and explicit invocation  
50 costs.

51 Optimal Transport (OT) offers a complementary way to compare distributions and quantify discrep-  
52 ancies between uncertain outcomes [Villani, 2003, Peyré and Cuturi, 2019]. OT is widely used  
53 for uncertainty-aware distributional comparison in domain adaptation [Courty et al., 2016], gener-  
54 ative modeling [Arjovsky et al., 2017, Bousquet et al., 2017], and statistical inference [Panaretos  
55 and Zemel, 2020]. However, OT remains underused in multi-agent coordination as a mechanism  
56 to compare agent output distributions and guide adaptive orchestration. Recent work suggests OT  
57 can be fruitfully combined with multi-agent reinforcement learning for scalability and alignment in  
58 complex environments [Baheri and Kochenderfer, 2024].

## 59 **Our contribution.**

- 60 • **Bandit-based orchestration with OT alignment:** We cast delegation over heterogeneous agents  
61 as a stochastic bandit regularized by OT distances between agent output distributions and task-  
62 specific references, enabling uncertainty-aware alignment and adaptive decision-making.
- 63 • **Theoretical guarantees:** We establish sublinear OT-regularized regret ( $\epsilon$  of Theorem 4.1), struc-  
64 tural optimality and robustness to noisy alignment ( $\delta$ ), and convergence and consistency proper-  
65 ties (5–4).
- 66 • **Empirical validation under heterogeneity and shift:** We evaluate BOT-Orch across synthetic  
67 and semi-synthetic settings, including a human–AI triage scenario under deployment shift (Sec-  
68 tion 6-7), demonstrating consistent improvements over standard bandit and heuristic baselines in  
69 both i.i.d. and non-i.i.d. environments.

## 70 **2 Related Work**

71 **Agent orchestration.** Prior work on agent orchestration largely focuses on selecting agents based  
72 on expected utility or offline accuracy [Keswani et al., 2021, Lai et al., 2022, Rasal and Hauer,  
73 2024]. However, these approaches often overlook practical constraints such as availability, cost, and  
74 capability. Recent work on human–AI orchestration highlights how inter-agent interactions shape  
75 system-level decision networks [Collins et al., 2024], yet uncertainty-aware adaptive orchestration  
76 under realistic constraints remains underexplored. We address this gap via uncertainty-aware OT-  
77 based orchestration. In a related direction, DiscoPOP [Lu et al., 2024] learns loss functions without  
78 human input, optimizing over objectives. This suggests that automatic objective discovery could  
79 extend to orchestration.

80 **Bandit approaches.** Multi-armed bandits provide a natural framework for sequential decision  
81 making under uncertainty, balancing exploration and exploitation [Lattimore and Szepesvári, 2020,  
82 Chen, 2024, Tong, 2024, Anonymous, 2024]. Recent surveys summarize advances in classical and  
83 contextual bandits and their applications [Chen, 2024, Anonymous, 2024], while empirical studies  
84 highlight design choices affecting performance [Bietti et al., 2021]. Emerging work connects bandits  
85 with large language models in complex environments [Xie et al., 2026]. Here, we use bandits  
86 to learn orchestrations that maximize expected utility across heterogeneous agents with uncertain  
87 performance.

88 **OT for uncertainty-aware orchestration.** Optimal transport provides a principled framework  
 89 for comparing probability distributions via geometrically meaningful discrepancies [Villani et al.,  
 90 2008]. We use OT to quantify uncertainty in heterogeneous agent outputs through distributional  
 91 disagreement and variability, building on recent work in OT-based uncertainty quantification [Oliver  
 92 et al., 2025a,c,b]. This yields uncertainty-aware weights capturing both confidence and disagree-  
 93 ment. The framework connects to distributional and uncertainty-aware reinforcement learning [Os-  
 94 band et al., 2013, Bellemare et al., 2017], providing a unified approach to adaptive multi-agent  
 95 coordination with theoretical guarantees and strong empirical performance.

## 96 3 Preliminaries and Notation

### 97 3.1 Agents, Tasks, and Task Space

98 Let  $\mathcal{A} := \{a_1, \dots, a_M\}$  denote a finite set of  $M$  agents, and let  $\mathcal{X} \subset \mathbb{R}^d$  be a measurable task  
 99 space endowed with its Borel  $\sigma$ -algebra. Time evolves over a finite horizon of  $T$  rounds indexed  
 100 by  $t = 1, \dots, T$ . At each round  $t$ , a task  $x_t \in \mathcal{X}$  arrives, forming a stochastic process  $(x_t)_{t=1}^T$ .  
 101 For each round  $s$  and each agent  $a_i \in \mathcal{A}$ , let  $R_s^i \in \mathbb{R}$  denote the reward that agent  $a_i$  would  
 102 obtain from handling task  $x_s$ , and let  $W_s^i \in \mathbb{R}$  denote the task-agent alignment cost (or mismatch  
 103 cost) associated with assigning task  $x_s$  to agent  $a_i$ . We collect these quantities into vectors  $\mathbf{R}_s :=$   
 104  $(R_s^1, \dots, R_s^M)^\top \in \mathbb{R}^M$  and  $\mathbf{W}_s := (W_s^1, \dots, W_s^M)^\top \in \mathbb{R}^M$ . Task arrivals may depend on past  
 105 observations, environmental conditions, or previous assignments. Define the history up to (but not  
 106 including) time  $t$  as  $\mathcal{H}_t := \{(x_s, \mathbf{R}_s, \mathbf{W}_s)\}_{s=1}^{t-1}$ . Then each task is drawn from a history-dependent  
 107 conditional distribution  $x_t \sim \mathcal{P}(\cdot \mid \mathcal{H}_t)$ ,  $t = 1, \dots, T$ , where  $\mathcal{P}(\cdot \mid \mathcal{H}_t)$  is a probability measure  
 108 over  $\mathcal{X}$  conditioned on the past history.

### 109 3.2 Correlated Multi-Agent Rewards

110 Recall that  $\mathbf{R}_t \in \mathbb{R}^M$  denotes the vector of rewards for all agents at round  $t$  (defined in Section 2.1).  
 111 In general, rewards may depend on both the current task and past system outcomes. We model this  
 112 via the conditional first and second moments:

$$\mathbb{E}[\mathbf{R}_t \mid x_t, \mathcal{H}_t] = \mathbf{r}(x_t) + \mathbf{f}(\mathbf{R}_{1:t-1}), \quad \text{Cov}[\mathbf{R}_t \mid x_t, \mathcal{H}_t] = \Sigma_t.$$

113 Here  $\mathbf{r} : \mathcal{X} \rightarrow \mathbb{R}^M$  is the baseline task-dependent mean reward vector,  $\mathbf{f} : \mathbb{R}^{M \times (t-1)} \rightarrow \mathbb{R}^M$   
 114 captures temporal dependence on past realized rewards (e.g., learning effects, fatigue, or system  
 115 congestion), and  $\Sigma_t \in \mathbb{R}^{M \times M}$  is a time-varying covariance matrix that models uncertainty and  
 116 cross-agent correlations at time  $t$ .

117 **I.I.D. tasks as a special case.** A common simplifying assumption is that tasks are independent and  
 118 identically distributed:

$$x_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{P}_X, \quad t = 1, \dots, T,$$

119 for some fixed distribution  $\mathcal{P}_X$  over  $\mathcal{X}$ . If, in addition, rewards do not depend on past outcomes,  
 120 then the history  $\mathcal{H}_t$  provides no additional information once  $x_t$  is known. In this case, the temporal  
 121 dependence term vanishes and the covariance becomes time-invariant:

$$\mathbb{E}[\mathbf{R}_t \mid x_t] = \mathbf{r}(x_t), \quad \text{Cov}[\mathbf{R}_t \mid x_t] = \Sigma,$$

122 where  $\Sigma$  is a constant matrix. This corresponds exactly to the special case of the correlated reward  
 123 model above with  $\mathbf{f} \equiv 0$  and  $\Sigma_t \equiv \Sigma$  for all  $t$ . Thus, the i.i.d. setting removes both temporal  
 124 dependence in the mean rewards and time variation in the covariance structure.

### 125 3.3 Survival-Based Rewards with Latent Frailty

126 In many real-world systems, agent performance is naturally measured by the *time to complete a task*  
 127 or *time until failure*, as in service response times, human or robotic workflows, or component life-  
 128 times. Such outcomes are commonly modelled using *survival analysis*, where the primary variable  
 129 is a *time-to-event*. Formally, let  $T_i(i) > 0$  denote the time-to-event for agent  $a_i$  on task  $x_t$ . Since  
 130 event times may be only partially observed, we introduce a *censoring indicator*  $\delta_t(i) \in \{0, 1\}$ ,  
 131 where  $\delta_t(i) = 1$  indicates a fully observed event and  $\delta_t(i) = 0$  indicates right-censoring.

132 Let  $S_i(\tau | x_t) = \mathbb{P}(T_t(i) > \tau | x_t)$  denote the *survival function* of agent  $a_i$ , capturing the proba-  
 133 bility that the task is not completed by time  $\tau$ , and allowing for heterogeneous performance across  
 134 agents and tasks. To model unobserved task-level factors affecting all agents, such as difficulty or  
 135 system load, we introduce a *latent frailty variable*  $\theta_t > 0$ , a shared random effect that multiplicat-  
 136 ively scales task difficulty.

137 The survival model is given by  $\mathbb{P}(T_t(i) > \tau | x_t, \theta_t) = S_i(\tau | x_t)^{\theta_t}$ ,  $R_t(i) = \delta_t(i) S_i(T_t(i) | x_t)^{\theta_t}$ .  
 138 Under this formulation, agents on the same task share  $\theta_t$ , inducing positive correlation in outcomes:  
 139 more difficult tasks lead to longer completion times for all agents. This survival-based reward model  
 140 is well-suited to settings involving reliability, speed, or risk, where hidden task-level factors jointly  
 141 influence performance.

### 142 3.4 OT-Based Alignment Costs

143 Beyond rewards, we also model how well an agent’s capabilities align with the requirements of a  
 144 task. In many applications, both tasks and agents are naturally described by *distributions* rather than  
 145 single feature vectors. For example, an agent may produce a distribution of outcomes (quality levels,  
 146 response times, error types), while a task may specify a desired target distribution over outcomes.

147 Let  $\mu_i \in \mathcal{P}(\mathcal{Y})$  denote the outcome distribution induced by agent  $a_i$  over a measurable space  $\mathcal{Y}$ ,  
 148 and let  $\nu_t \in \mathcal{P}(\mathcal{Y})$  represent the reference or desired outcome distribution associated with task  $x_t$ .  
 149 We quantify the mismatch between an agent and a task using the Wasserstein distance  $W_c(\nu_t, \mu_i)$ ,  
 150 which measures the minimal cost of transporting mass from one distribution to the other under  
 151 ground cost  $c$ . This provides a geometrically meaningful notion of alignment that accounts for the  
 152 full distribution of outcomes rather than just summary statistics.

153 To incorporate randomness and modeling noise, we define the stochastic alignment cost  $\mathcal{W}_t(i) :=$   
 154  $W_c(\nu_t, \mu_i) + \epsilon_t(i)$ , with  $\epsilon_t(i) \sim \mathcal{N}(0, \sigma_t^2)$ . This formulation captures both systematic mismatch  
 155 (via the Wasserstein term) and unpredictable variability (via  $\epsilon_t(i)$ ). Introducing alignment costs at  
 156 this stage allows us to jointly model (i) how well an agent is suited to a task and (ii) the stochastic  
 157 rewards that result from performing it, providing a unified framework for assignment decisions under  
 158 uncertainty.

### 159 3.5 Orchestration Policy and Objective

160 We now formalize the decision-making problem faced by the orchestrator. At each round  $t$ , after  
 161 observing the task  $x_t$  and past history  $\mathcal{H}_t$ , the orchestrator selects a *randomized assignment policy*  
 162  $\pi_t \in \Delta(\mathcal{A})$ , where  $\Delta(\mathcal{A})$  denotes the probability simplex over agents. An agent  $i_t \sim \pi_t$  is then  
 163 sampled and assigned to handle task  $x_t$ . Randomized policies allow exploration and robustness to  
 164 uncertainty in agent performance.

165 Given a policy  $\pi_t$ , the expected net reward at round  $t$  is defined as

$$R_t(\pi_t) := \pi_t^\top (\mathbf{r}(x_t) - \lambda \mathbb{E}[\mathcal{W}_t]), \quad \lambda > 0,$$

166 where  $\mathbf{r}(x_t) \in \mathbb{R}^M$  is the vector of expected task-dependent rewards and  $\mathcal{W}_t \in \mathbb{R}^M$  is the vector  
 167 of alignment costs. The scalar parameter  $\lambda > 0$  controls the trade-off between maximizing perfor-  
 168 mance and minimizing mismatch costs. Thus,  $R_t(\pi_t)$  represents the expected utility of assigning  
 169 the task according to  $\pi_t$ , balancing reward quality against alignment penalties.

170 The overall objective is to maximize the cumulative expected net reward over the time horizon:

$$J(\pi_{1:T}) := \sum_{t=1}^T \mathbb{E}[R_t(\pi_t)],$$

171 where the expectation is taken over task arrivals, reward randomness, alignment noise, and the  
 172 orchestrator’s own randomization. This objective formalizes the goal of learning an assignment  
 173 strategy that performs well on average while accounting for uncertainty and task-agent compatibility.

174 We may also consider a finite set of candidate orchestration strategies  $\mathcal{S} = \{s_1, \dots, s_K\}$ , where  
 175 each strategy  $s_k$  specifies a rule for selecting policies  $\pi_t^{(k)}$  based on the observed history.

176 **3.6 Regret and Optimal Policy**

177 The optimal policy for known rewards and alignment costs is

$$\pi_t^* := \arg \max_{\pi \in \Delta(\mathcal{A})} \pi^\top \left( \mathbf{r}(x_t) - \lambda \mathbb{E}[\mathcal{W}_t] \right), \quad (1)$$

178 with cumulative regret

$$\mathcal{R}_T := \sum_{t=1}^T (\pi_t^*)^\top \left( \mathbf{r}(x_t) - \lambda \mathbb{E}[\mathcal{W}_t] \right) - \sum_{t=1}^T \mathbb{E}[R_t(\pi_t)]. \quad (2)$$

179 For i.i.d. tasks, the regret simplifies due to stationarity. For non-i.i.d. tasks,  $\mathcal{R}_T$  captures efficiency  
 180 loss arising from temporal correlations, history-dependent distributions, and uncertainty in strategy  
 181 selection.

182 **4 Theoretical Properties**

183 We begin by stating the standing assumptions that remain in force throughout.

- 184 **(A1)** The ground cost  $c : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$  is  $L$ -Lipschitz and bounded.
- 185 **(A2)** Rewards are uniformly bounded:  $0 \leq R_t(i) \leq R_{\max}$  for all  $t$  and all actions  $i \in \mathcal{A}$ .
- 186 **(A3)** The frailty variables  $(\theta_t)_{t \geq 1}$  admit finite exponential moments.
- 187 **(A4)** Conditional on  $(x_t, \theta_t)$ , censoring is independent of survival time.
- 188 **(A5)** The learning rate satisfies  $\eta_t = O(t^{-1/2})$ .

189 Assumptions (A1)–(A5) ensure: (i) stability of Wasserstein distances under empirical perturbations;  
 190 (ii) sub-exponential behaviour of frailty-adjusted rewards; (iii) well-posedness of the induced soft-  
 191 max stochastic approximation dynamics.

192 We remark that these conditions are standard in the literature on stochastic approximation,  
 193 Wasserstein-based learning, and frailty-adjusted reward models [Ambrosio et al., 2008, Shalev-  
 194 Shwartz and Ben-David, 2014, Trillos and Slepčev, 2016]. In practice, they are not restrictive:  
 195 Lipschitz and bounded costs are typical in OT and multi-agent learning applications, bounded re-  
 196 wards naturally arise in reinforcement learning, and learning rates of order  $t^{-1/2}$  are widely used to  
 197 guarantee convergence. Moreover, our framework generalizes several prior works by allowing frailty  
 198 variables with arbitrary distributions admitting finite exponential moments, rather than restricting to  
 199 specific parametric forms [Del Barrio and Loubes, 2019, Wang et al., 2020].

200 **4.1 Main Results**

201 Following the assumptions in Section 4, we now present the main theoretical guarantees for the  
 202 OT-regularised bandit model. The regret analysis is conducted for a general exponential-weights  
 203 (softmax) procedure applied to bounded, OT-regularised reward signals. Importantly, the resulting  
 204 guarantees depend only on the boundedness of these rewards and are independent of the specific  
 205 generative model.

206 The additional modelling components introduced in the setup, namely, correlated rewards, survival-  
 207 based frailty, and non-i.i.d. task generation—serve as a motivating probabilistic framework in which  
 208 such bounded reward processes naturally arise. They are not directly used in the regret derivation,  
 209 but instead provide one possible instantiation of the abstract reward model. Proofs are provided in  
 210 Appendices B.1–B.3.

211 **Theorem 4.1** (Performance and Convergence of OT-Regularized Bandits). *Let  $\mathcal{A}$  be a finite set of*  
 212 *agents, and let  $(\tilde{r}_t(i))_{i \in \mathcal{A}}$  denote the BOT–Orch alignment-adjusted rewards defined by*

$$\tilde{r}_t(i) = \hat{r}_t(i) - \lambda W_c(\mu_i, \nu_t),$$

213 *where  $\hat{r}_t(i)$  is the empirical reward,  $\mu_i$  is the agent-specific outcome distribution,  $\nu_t$  is the target*  
 214 *distribution at time  $t$ , and  $W_c$  is the ground-cost Wasserstein distance. Assume (A1)–(A5) hold.*

Environment	BOT-Orch	No-OT ( $\lambda=0$ )	Random	UCB1 (MAB)
<b>IID</b>				
IID-G	<b>537.43±16.57</b>	656.31±22.64	673.99±19.08	664.89±29.85
IID-M	<b>459.84±9.27</b>	545.85±17.74	553.07±19.16	549.48±7.40
<b>Non-IID</b>				
NonIID-BB	<b>410.04±71.06</b>	503.14±91.24	520.28±101.61	496.60±79.57
NonIID-PS	<b>571.43±18.71</b>	713.81±25.79	734.65±30.08	729.51±28.34
NonIID-SD	<b>564.13±18.11</b>	717.24±22.51	707.69±18.33	708.60±16.06

Table 1: Cumulative Alignment Cost (mean  $\pm$  95% CI across 5 seeds;  $T = 200$ ). Best in bold.

215 Let  $i_t^* = \arg \max_{i \in \mathcal{A}} \tilde{r}_t(i)$  denote the optimal agent at round  $t$ , and define the OT-regularized  
216 cumulative regret by

$$\mathcal{R}_T := \sum_{t=1}^T (\tilde{r}_t(i_t^*) - \tilde{r}_t(i_t)),$$

217 where  $i_t$  is the agent chosen by the BOT-Orch policy at round  $t$ . Let  $\phi_t = \pi_t(\cdot)$  denote the strategy-  
218 mixing vector. Suppose that the initial conditions are well-defined, i.e.,  $\hat{r}_0(i)$  is finite for all  $i \in \mathcal{A}$ ,  
219 and  $\phi_0 \in \Delta^{|\mathcal{A}|-1}$ . Then, the following statements hold:

- 220 1. **Sublinear OT-Regret.** The cumulative OT-regularized regret satisfies  $\mathcal{R}_T = O(\sqrt{T})$ .
- 221 2. **Structural OT-Optimality.** For any two agents  $i, j \in \mathcal{A}$  with  $\hat{r}_t(i) = \hat{r}_t(j)$ ,  
222 if  $W_c(\mu_i, \nu_t) < W_c(\mu_j, \nu_t)$ , then  $\tilde{r}_t(i) > \tilde{r}_t(j)$ .
- 223 3. **Margin Robustness under Noisy Alignment.** Suppose the observed Wasserstein distances are  
224 perturbed by independent Gaussian noise  $\epsilon_t(i) \sim \mathcal{N}(0, \sigma^2)$ ,  
225 defining  $\tilde{W}_c(\mu_i, \nu_t) = W_c(\mu_i, \nu_t) + \epsilon_t(i)$ , and let  $\Delta_{ij} = W_c(\mu_j, \nu_t) - W_c(\mu_i, \nu_t)$ . If it holds  
226  $|\Delta_{ij}| > \sigma\sqrt{2 \log 2}$ , then  $\mathbb{P}(\tilde{r}_t(i) < \tilde{r}_t(j)) < 1/4$ .
- 227 4. **Convergence of Orchestration Weights.** The sequence of policy vectors  $(\phi_t)_{t \geq 1}$  converges al-  
228 most surely to a stable equilibrium  $\phi_\infty$  of the limiting ODE  $\dot{\phi} = \text{Softmax}(\mathbb{E}[\tilde{r}]) - \phi$ .
- 229 5. **Uniform Consistency of Empirical Rewards.** Let  $\hat{r}_t(i)$  denote the empirical average of frailty-  
230 adjusted rewards for agent  $i$ . Then,

$$\sup_{i \in \mathcal{A}} |\hat{r}_t(i) - \mathbb{E}[R_t(i) | x_t]| \xrightarrow[t \rightarrow \infty]{\text{a.s.}} 0.$$

231 We note that Theorem 4.1 provides a foundational characterization of OT-regularized bandit learn-  
232 ing. It integrates distributional alignment into the reward structure, yielding principled agent dif-  
233 ferentiation 2, sublinear cumulative regret 1, robustness to noisy observations 3, convergence of  
234 orchestration weights 4, and uniform consistency of empirical rewards 5. The framework enables  
235 heterogeneous agents to learn coordinated policies under distributionally-aware uncertainty, sup-  
236 porting robust, adaptive orchestration in stochastic, partially observable environments. When task-  
237 specific reference distributions are unavailable, they can be constructed using Wasserstein barycenters  
238 [Chewi et al., 2025] of observed empirical task distributions. We refer the reader to Appendix  
239 B for a proof of Theorem 4.1.

## 240 5 The proposed algorithm

241 BOT-Orch combines bandit-based selection with OT alignment and survival-based rewards to han-  
242 dle heterogeneous, non-stationary agents. In the *i.i.d. task setting* (Algorithm 1 in the Appendix),  
243 a Boltzmann policy selects agents using exponentially smoothed rewards penalized by OT mis-  
244 alignment, balancing exploitation and alignment, while survival rewards capture latent difficulty,  
245 censoring, and reliability. The *non-i.i.d. extension* (Algorithm 2 in the Appendix) allows history-  
246 dependent task distributions and reward updates, handling temporal correlations, non-stationarity,  
247 and regime shifts via a correction term that encodes memory effects. Overall, BOT-Orch is a *risk-*  
248 *aware, alignment-regularized bandit algorithm* in distributional space, where OT enforces task-  
249 agent compatibility and survival rewards provide robustness to censoring and heterogeneity. Its  
250 modular design accommodates alternative OT solvers, survival models, and exploration schemes,  
251 and motivates regret analysis under composite reward-cost objectives as well as questions of con-  
252 vergence and adaptivity in non-stationary environments.

Environment	Cumulative Net Utility				Oracle Regret			
	BOT-Orch	No-OT ( $\lambda=0$ )	Random	UCB1 (MAB)	BOT-Orch	No-OT ( $\lambda=0$ )	Random	UCB1 (MAB)
<b>IID</b>								
IID-G	<b>-467.528±17.60</b>	-588.34±23.19	-605.71±18.91	-603.31±29.13	<b>122.65±5.68</b>	243.47±14.26	260.83±11.59	258.43±18.11
IID-M	<b>-386.10±10.04</b>	-478.54±20.49	-482.81±18.65	-484.57±9.35	<b>74.70±5.45</b>	166.25±11.99	170.52±10.99	172.28±8.40
<b>Non-IID</b>								
NonIID-BB	<b>-335.25±70.42</b>	-434.01±86.69	-449.34±100.73	-426.09±79.48	<b>76.85±7.42</b>	175.61±16.44	190.94±31.44	167.69±19.20
NonIID-PS	<b>-498.28±22.49</b>	-642.77±25.89	-666.10±29.31	-663.22±27.77	<b>126.77±11.68</b>	271.26±11.84	294.58±20.25	291.71±14.45
NonIID-SD	<b>-492.92±18.74</b>	-644.68±21.62	-637.82±18.54	-643.69±17.20	<b>128.77±13.30</b>	280.53±10.81	273.67±10.74	279.54±11.52

Table 2: Cumulative Net Utility and Oracle Regret (mean  $\pm$  95% CI across 5 seeds;  $T = 200$ ). Best in bold.

## 253 6 Synthetic Experiments

### 254 6.1 Dataset and Task Description

255 We consider both i.i.d and non-i.i.d. regimes (see Fig. 2-3 in the Appendix):

- 256 • **IID-G:** Tasks  $x_t$  sampled i.i.d. over  $\mathcal{X}$ . Agent rewards drawn i.i.d. from fixed Gaussian distribu-  
257 tions with matched mean ( $\approx 0.5$ ) but heterogeneous higher-order shape.
- 258 • **IID-M:** Tasks  $x_t$  sampled i.i.d. from a half-moons distribution. Agent rewards drawn i.i.d. from  
259 fixed distributions with matched mean ( $\approx 0.5$ ) but differing variance/skewness/bimodality.
- 260 • **NonIID-PS:** Piecewise-stationary rewards: mean fixed, variance shifts at unknown changepoints,  
261 inducing distributional shifts without mean change.
- 262 • **NonIID-SD:** Smooth non-stationarity via sinusoidal drift in reward means.
- 263 • **NonIID-BB:** Latent reward means follow a temporally correlated Brownian-bridge path with fixed  
264 endpoints.

265 All experiments were conducted on a standard x86\_64 CPU platform using 2 CPU cores, 13.6 GB  
266 RAM, and 107 GB disk storage.

### 267 6.2 Baselines

- 268 • **BOT-Orch (ours).** Uses alignment-adjusted rewards  $\tilde{r}_t(i) = \hat{r}_t(i) - \lambda W_t(i)$  and samples the  
269 selected agent  $i_t$  from a softmax policy over  $\tilde{r}_t(i)$  (Algorithm 1–2). This couples OT alignment  
270 with sequential exploration/exploitation.
- 271 • **No-OT ( $\lambda = 0$ ).** Ablation removing the OT alignment term. The policy is computed from  $\hat{r}_t(i)$   
272 only. This isolates the contribution of distributional alignment.
- 273 • **Random.** Uniformly samples an agent each round. This serves as a naive lower bound that does  
274 not learn from observations.
- 275 • **UCB1 (MAB).** A classical multi-armed bandit baseline that selects the agent with the highest up-  
276 per confidence bound based on empirical reward estimates, balancing exploration and exploitation  
277 without OT alignment.

278 We omit Greedy-EMA, which picks at each round the agent with the highest exponentially-smoothed  
279 reward estimate, because typical implementations assume *full-information* rewards for all agents  
280 each round, unlike our *bandit feedback* setting (only the chosen agent is observed). A bandit-  
281 compatible “greedy” would require exploration/confidence bounds and is essentially covered by  
282 standard bandit baselines (e.g., UCB1).

### 283 6.3 Evaluation Metrics

- 284 • **Cumulative net utility.** We evaluate the OT-regularized net utility  $U_t(i_t) = R_t(i_t) - \lambda W_t(i_t)$ ,  
285 and report the cumulative net utility  $\sum_{t=1}^T U_t(i_t)$ . This captures the overall performance when  
286 explicitly trading off reward quality and alignment cost via  $\lambda$ .
- 287 • **Cumulative alignment cost.** We report  $\sum_{t=1}^T W_t(i_t)$ , the cumulative OT alignment cost incurred  
288 by the selected agent assignments. Lower values indicate more distributionally aligned agent-task  
289 matching.

	Cumulative Net Utility				Oracle Regret			
	BOT-Orch	No-OT ( $\lambda=0$ )	Random	UCB1	BOT-Orch	No-OT ( $\lambda=0$ )	Random	UCB1
<b>IID</b>	<b>108.84±2.22</b>	103.37±2.51	83.98±5.92	80.55±4.88	<b>2.29±2.11</b>	9.80±1.61	28.72±6.03	31.31±3.98
<b>Non-IID</b>	<b>110.61±1.03</b>	103.17±1.80	79.78±5.52	85.82±4.91	<b>0.59±0.93</b>	10.14±1.17	32.97±4.87	26.26±4.28

	Team Accuracy				Cumulative Alignment Cost			
	BOT-Orch	No-OT ( $\lambda=0$ )	Random	UCB1	BOT-Orch	No-OT ( $\lambda=0$ )	Random	UCB1
<b>IID</b>	<b>0.980±0.010</b>	0.907±0.022	0.917±0.020	0.902±0.026	<b>1.28±0.53</b>	10.51±1.61	10.28±2.01	11.14±1.33
<b>Non-IID</b>	<b>0.993±0.007</b>	0.905±0.016	0.905±0.024	0.919±0.025	<b>0.85±0.23</b>	10.84±1.18	11.69±1.62	9.46±1.43

Table 3: Deployment metrics across all four methods and both experimental conditions (Mean  $\pm$  95% CI across 30 seeds,  $T=114$ ,  $\lambda=3.0$ ). IID condition uses Algorithm 1; Non-IID uses Algorithm 2 (ID patients rounds 1–57, shifted rounds 58–114). Best in bold.

- **Oracle regret.** We report OT-regularized cumulative regret relative to the best alignment-adjusted agent at each round:  $\sum_{t=1}^T (\max_i U_t(i) - U_t(i_t))$  measuring efficiency loss due to suboptimal selection under uncertainty and non-stationarity.

## 6.4 Results

BOT-Orch achieves the highest cumulative net utility and lowest oracle regret across all environments (Table 1-2), consistently outperforming No-OT, Random, and UCB1. The improvement is particularly pronounced in non-i.i.d. settings, where BOT-Orch maintains low regret under distributional shifts such as piecewise variance changes and smooth drift. These results indicate that incorporating OT-based alignment enables more effective adaptation to heterogeneous and non-stationary agent behaviour.

In Appendix D, we evaluate BOT-Orch on additional synthetic benchmarks using survival-based metrics, showing that BOT-Orch achieves best performance. We further conduct an ablation study on the parameter  $\lambda$  (Appendix D.5). As  $\lambda$  increases, BOT-Orch exhibits performance improvements, with large gains relative to  $\lambda = 0$  and diminishing returns at higher values.

## 7 Semi-Synthetic Experiments: Human-AI Triage Under Deployment Shift

### 7.1 Dataset and Task Description

We use the Breast Cancer Wisconsin (Diagnostic) dataset, consisting of 569 patient cases with 30 numerical features and a binary target (malignant vs. benign). The data is split into 60% train, 20% calibration, 10% Test-ID, and 10% Test-Shift. To simulate deployment shift, we perturb a subset of features in Test-Shift with additive Gaussian noise  $\mathcal{N}(0, 0.64)$  and a +0.5 standard-unit bias, modelling a change in patient population at deployment. At each round  $t$ , a patient  $x_t$  arrives and the selected agent receives a reward of 1 if it classifies the patient correctly and 0 otherwise, under standard bandit feedback.

We consider two agents: an AI classifier and a human proxy with complementary accuracy, where the human performs better on shifted cases while the AI performs better in-distribution. Full model details and accuracy statistics are provided in Appendix E. All experiments were conducted on a standard x86\_64 CPU platform using 2 CPU cores, 13.6 GB RAM, and 107 GB disk storage.

### 7.2 Baselines

We compare the same four methods used in the synthetic experiments, now applied to the clinical triage task. We use  $\alpha = 0.90$ ,  $\eta = 5.0$ ,  $\lambda = 3.0$ , and, for the non-i.i.d. variant, the same history correction form as Algorithm 2 with  $\beta = 0.05$ . Each episode has  $T = 114$  rounds, and we average over  $n = 30$  random seeds (patient orderings). We report mean  $\pm$  95% CI across seeds.

### 7.3 Evaluation Metrics

From Section 6.3, we use cumulative net utility, cumulative alignment cost, and oracle regret, defined as above. We additionally report metrics specific to the human–AI deferral setting:

- **Team accuracy:** fraction of patients correctly classified by the selected agent (AI or human).

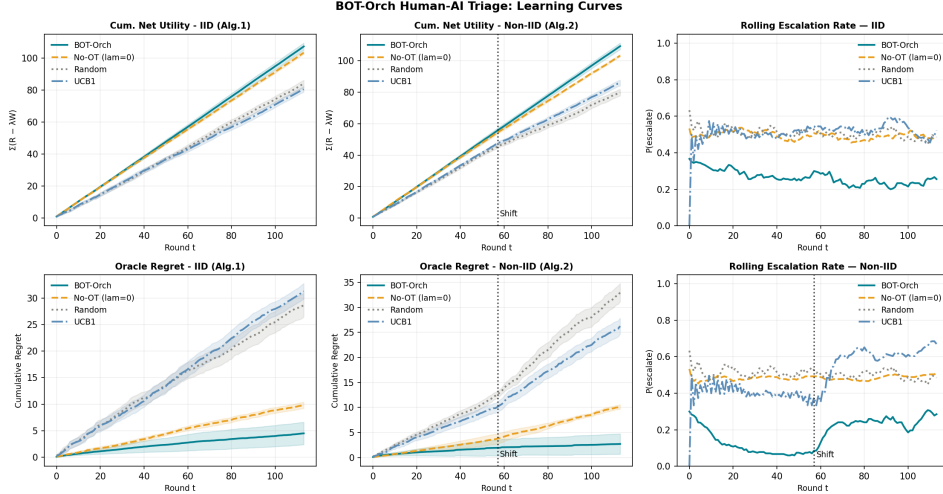


Figure 1: **Cumulative learning curves.** *Top row:* cumulative net utility. *Bottom row:* oracle regret. *Right column:* rolling escalation rate (window  $w=8$  rounds). Left panels: IID condition (Algorithm 1); middle panels: Non-IID condition (Algorithm 2), with the dotted vertical line marking the shift onset at round 57; right panels: escalation rate evolution.

326 • **Escalation rate:** fraction of patients routed to the human expert.

## 327 7.4 Results

328 Table 3 summarizes the results. BOT-Orch achieves the highest cumulative net utility and lowest  
 329 oracle regret in both regimes, outperforming all baselines. In the IID condition, it improves net  
 330 utility by +5.5 over No-OT while reducing alignment cost by an order of magnitude. In the Non-  
 331 IID setting, the gain is larger, reaching  $110.61 \pm 1.03$  net utility and  $0.59 \pm 0.93$  regret compared to  
 332  $103.17 \pm 1.80$  and  $10.14 \pm 1.18$  for No-OT.

333 The No-OT ablation confirms that the OT term drives the improvement: removing it increases regret  
 334 and alignment cost. BOT-Orch also exhibits targeted escalation under shift, routing fewer patients  
 335 overall while increasing escalation on shifted cases, indicating effective adaptation to distributional  
 336 mismatch. Figure 1 further shows that BOT-Orch separates from all baselines early and maintains  
 337 the performance gap over time. In the Non-IID setting, the escalation rate increases after the shift,  
 338 demonstrating online adaptation to the changing population. Furthermore, the rolling escalation rate  
 339 in Figure 1 shows how the policy adapts over time, increasing reliance on the human expert when  
 340 distributional shift or uncertainty rises.

341 Additional results are presented in Appendix F, including more results on the escalation rate, diag-  
 342 nostic analysis, and an ablation study on the parameter  $\lambda$ . We observe that, as  $\lambda$  increases, perfor-  
 343 mance improves markedly peaking around  $\lambda \approx 3.0$ , after which gains saturate or slightly diminish  
 344 at very high values.

## 345 8 Discussion and Conclusion.

346 This paper formulates heterogeneous agent orchestration under uncertain as a sequential decision  
 347 problem and introduces BOT-Orch, a bandit-based framework that jointly learns delegation while  
 348 incorporating optimal transport alignment costs. While BOT-Orch provides a principled approach  
 349 with strong theoretical and empirical performance, it has various practical limitations. First, com-  
 350 puting Wasserstein distances can be expensive, especially in high-dimensional settings, limiting  
 351 real-time use. Second, the reward–alignment trade-off is controlled by a scalar parameter that may  
 352 require tuning across environments. Third, the framework assumes access to task-specific refer-  
 353 ence distributions, which may be unavailable in open or non-stationary settings. Addressing these  
 354 challenges remains an important direction for future work.

355 **References**

- 356 Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré. *Gradient Flows: In Metric Spaces and in*  
357 *the Space of Probability Measures*. Lectures in Mathematics ETH Zürich. Birkhäuser Basel, 2nd  
358 edition, 2008.
- 359 Anonymous. Bandit algorithms: A comprehensive review and their dynamic selection from a portfo-  
360 lio for multicriteria top-k recommendation. *Expert Systems with Applications*, 246:123151, 2024.  
361 doi: 10.1016/j.eswa.2024.123151.
- 362 Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein gan. In *ICML*, 2017.
- 363 Krishna Arjun, David Parlevliet, Hai Wang, and Amirmehdi Yazdani. Optimizing coalition forma-  
364 tion strategies for scalable multi-robot task allocation: A comprehensive survey. *Robotics*, 14(7):  
365 93, 2025. doi: 10.3390/robotics14070093. URL [https://www.mdpi.com/2218-6581/14/7/](https://www.mdpi.com/2218-6581/14/7/93)  
366 [93](https://www.mdpi.com/2218-6581/14/7/93).
- 367 Ali Baheri and Mykel J. Kochenderfer. The synergy between optimal transport theory and multi-  
368 agent reinforcement learning. *arXiv preprint arXiv:2401.10949*, 2024. Explores integration  
369 of optimal transport theory with MARL for policy alignment, resource distribution, and non-  
370 stationarity mitigation.
- 371 Marc G Bellemare, Will Dabney, and Remi Munos. A distributional perspective on reinforcement  
372 learning. In *International Conference on Machine Learning*, pages 449–458, 2017.
- 373 Michel Benaim. Dynamics of stochastic approximation algorithms. *Annals of Probability*, 27(1):  
374 361–414, 1999. doi: 10.1214/aop/1019160366.
- 375 Daniel S. Bernstein, Robert Givan, Neil Immerman, and Shlomo Zilberstein. The complexity of  
376 decentralized control of markov decision processes. *Mathematics of Operations Research*, 27(4):  
377 819–840, 2002.
- 378 Alberto Bietti, Alekh Agarwal, and John Langford. A contextual bandit bake-off. In *Journal of*  
379 *Machine Learning Research*, volume 22, pages 1–49, 2021.
- 380 Olivier Bousquet, Sylvain Gelly, Ilya Tolstikhin, Carl-Johann Simon-Gabriel, and Bernhard  
381 Schoelkopf. From optimal transport to generative modeling: the vegan cookbook. *arXiv preprint*  
382 *arXiv:1705.07642*, 2017.
- 383 Sébastien Bubeck and Nicolò Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-  
384 armed bandit problems. *Foundations and Trends in Machine Learning*, 5(1):1–122, 2012.
- 385 Lucian Buşoniu, Robert Babuška, and Bart De Schutter. A comprehensive survey of multiagent  
386 reinforcement learning. *IEEE Transactions on Systems, Man, and Cybernetics, Part C*, 38(2):  
387 156–172, 2008.
- 388 Nicolò Cesa-Bianchi and Gábor Lugosi. *Prediction, Learning, and Games*. Cambridge University  
389 Press, 2006.
- 390 Qiufan Chen. A survey on contextual multi-armed bandits. *Applied and Computational Engineering*,  
391 53:287–295, 2024. doi: 10.54254/2755-2721/53/20241593.
- 392 Yao Cheng et al. Autogen: Enabling next-gen llm applications via multi-agent conversation. *arXiv*  
393 *preprint arXiv:2308.08155*, 2023.
- 394 Sinho Chewi, Jonathan Niles-Weed, and Philippe Rigollet. *Wasserstein Barycenters*, pages 211–  
395 227. Springer Nature Switzerland, Cham, 2025. ISBN 978-3-031-85160-5. doi: 10.1007/  
396 978-3-031-85160-5\_8. URL [https://doi.org/10.1007/978-3-031-85160-5\\_8](https://doi.org/10.1007/978-3-031-85160-5_8).
- 397 Katherine M Collins, Valerie Chen, Ilia Sucholutsky, Hannah Rose Kirk, Malak Sadek, Holli  
398 Sargeant, Ameet Talwalkar, Adrian Weller, and Umang Bhatt. Modulating language model expe-  
399 riences through frictions. *arXiv preprint arXiv:2407.12804*, 2024.

- 400 Nicolas Courty, Rémi Flamary, Devis Tuia, and Alain Rakotomamonjy. Optimal transport for domain  
401 adaptation. *IEEE transactions on pattern analysis and machine intelligence*, 39(9):1853–  
402 1865, 2016.
- 403 Eustasio Del Barrio and Jean-Michel Loubes. Frailty models and empirical process theory. *Annals  
404 of Statistics*, 47(5):2519–2543, 2019.
- 405 Ferdinando Fioretto, Enrico Pontelli, and William Yeoh. Distributed constraint optimization prob-  
406 lems and applications: A survey. *Journal of Artificial Intelligence Research*, 61:623–698, 2018.
- 407 Samarth Gupta, Shreyas Chaudhari, Gauri Joshi, and Osman Yağan. Multi-armed bandits with  
408 correlated arms. *IEEE Transactions on Information Theory*, 67(10):6711–6732, 2021.
- 409 Eric A. Hansen, Daniel S. Bernstein, and Shlomo Zilberstein. Dynamic programming for partially  
410 observable stochastic games. In *Proceedings of the National Conference on Artificial Intelligence  
411 (AAAI)*, 2004.
- 412 Aviad Heifetz, Martin Meier, and Burkhard Schipper. Interactive unawareness. *Journal of Economic  
413 Theory*, 130(1):78–94, 2006.
- 414 Pablo Hernandez-Leal, Bilal Kartal, and Matthew E. Taylor. A survey and critique of multiagent  
415 deep reinforcement learning. *Autonomous Agents and Multi-Agent Systems*, 33:750–797, 2019.
- 416 Vijay Keswani, Matthew Lease, and Krishnaram Kenthapadi. Towards unbiased and accurate def-  
417erral to multiple experts. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and  
418 Society*, pages 154–165, 2021.
- 419 Vivian Lai, Samuel Carton, Rajat Bhatnagar, Q Vera Liao, Yunfeng Zhang, and Chenhao Tan.  
420 Human-ai collaboration via conditional delegation: A case study of content moderation. In *Pro-  
421 ceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, pages 1–18,  
422 2022.
- 423 Peter Landgren, Vaibhav Srivastava, and Naomi E. Leonard. Distributed cooperative decision-  
424 making in multiarmed bandits: Frequentist and bayesian algorithms. *CDC*, 2016.
- 425 Tor Lattimore and Csaba Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2020.
- 426 Jacky Liang, Wenjia Wang, et al. Taskmatrix.ai: Completing tasks by connecting foundation models  
427 with millions of apis. *arXiv preprint arXiv:2303.16434*, 2023.
- 428 Chris Lu, Samuel Holt, Claudio Fanconi, Alex Chan, Jakob Foerster, Mihaela van der Schaar, and  
429 Robert Lange. Discovering preference optimization algorithms with and for large language mod-  
430 els. *Advances in Neural Information Processing Systems*, 37:86528–86573, 2024.
- 431 Reza Olfati-Saber, J. Alex Fax, and Richard M. Murray. Consensus and cooperation in networked  
432 multi-agent systems. *Proceedings of the IEEE*, 95(1):215–233, 2007. URL [https://doi.org/  
433 10.1109/JPR0C.2006.887293](https://doi.org/10.1109/JPR0C.2006.887293).
- 434 Frans A. Oliehoek and Christopher Amato. *A Concise Introduction to Decentralized POMDPs*.  
435 Springer, 2016.
- 436 Mary Chriselda Antony Oliver, Matthew Graham, Ioanna Manolopoulou, Graham F Medley,  
437 Lorenzo Pellis, Koen B Pouwels, Matthew Thorpe, and T Deirdre Hollingsworth. Uncertainty  
438 quantification in cost-effectiveness analysis for stochastic-based infectious disease models: In-  
439 sights from surveillance on lymphatic filariasis. *Journal of Theoretical Biology*, page 112197,  
440 2025a.
- 441 Mary Chriselda Antony Oliver, Emmanuel Hartman, and Tom Needham. Conic formulations  
442 of transport metrics for unbalanced measure networks and hypernetworks. *arXiv preprint  
443 arXiv:2508.10888*, 2025b.
- 444 Mary Chriselda Antony Oliver, Michael Roberts, Carola-Bibiane Schönlieb, and Matthew Thorpe.  
445 Laplace learning in wasserstein space. *arXiv preprint arXiv:2511.13229*, 2025c.

- 446 Ian Osband, Daniel Russo, and Benjamin Van Roy. (more) efficient reinforcement learning via  
447 posterior sampling. In *Advances in Neural Information Processing Systems*, pages 3003–3011,  
448 2013.
- 449 Victor M. Panaretos and Yoav Zemel. *Statistical aspects of Wasserstein distances*. Annual Review  
450 of Statistics and Its Application, 2020.
- 451 Joon Sung Park, Joseph O’Brien, Carrie J. Cai, et al. Generative agents: Interactive simulacra of  
452 human behavior. *UIST*, 2023.
- 453 Gabriel Peyré and Marco Cuturi. *Computational Optimal Transport*. Foundations and Trends in  
454 Machine Learning, 2019.
- 455 Sumedh Rasal and EJ Hauer. Navigating complexity: Orchestrated problem solving with multi-  
456 agent llms. *arXiv preprint arXiv:2402.16713*, 2024.
- 457 Tabish Rashid, Mikayel Samvelyan, Christian Schroeder, et al. Qmix: Monotonic value function  
458 factorisation for deep multi-agent reinforcement learning. In *ICML*, 2018.
- 459 Yara Rizk, Mariette Awad, and Edward W. Tunstel. Cooperative heterogeneous multi-robot systems:  
460 A survey. *ACM Computing Surveys*, 52(2):1–31, 2019. doi: 10.1145/3303848. URL <https://doi.org/10.1145/3303848>.
- 462 Shai Shalev-Shwartz. Online learning and online convex optimization. *Foundations and Trends in*  
463 *Machine Learning*, 4:107–194, 2012.
- 464 Shai Shalev-Shwartz and Shai Ben-David. *Understanding Machine Learning: From Theory to*  
465 *Algorithms*. Cambridge University Press, 2014.
- 466 Yoav Shoham and Kevin Leyton-Brown. *Multiagent Systems: Algorithmic, Game-Theoretic, and*  
467 *Logical Foundations*. Cambridge University Press, 2008. ISBN 9780521899437.
- 468 Peter Sunehag, Guy Lever, Audrunas Gruslys, et al. Value-decomposition networks for cooperative  
469 multi-agent learning. In *AAMAS*, 2018.
- 470 Ruoyi Tong. A survey of the application and technical improvement of the multi-armed bandit.  
471 *Applied and Computational Engineering*, 77:25–31, 2024.
- 472 Nicolas Garcia Trillos and Dejan Slepčev. A variational approach to the consistency of spectral  
473 clustering. *Applied and Computational Harmonic Analysis*, 40(2):274–319, 2016.
- 474 Cédric Villani. *Topics in Optimal Transportation*. American Mathematical Society, 2003.
- 475 Cédric Villani et al. *Optimal transport: old and new*, volume 338. Springer, 2008.
- 476 Peter Walters. *An Introduction to Ergodic Theory*, volume 79 of *Graduate Texts in Mathematics*.  
477 Springer, 1982. ISBN 978-0-387-90626-6.
- 478 Qiang Wang et al. Stochastic approximation for survival models with frailty terms. *Statistics in*  
479 *Medicine*, 39(10):1405–1420, 2020.
- 480 Peter R. Wurman, Raffaello D’Andrea, and Mick Mountz. Coordinating hundreds of cooperative,  
481 autonomous vehicles in warehouses. In *AAAI Conference on Artificial Intelligence*, 2007. URL  
482 <https://www.aaai.org/Library/AAAI/2007/aaai07-282.php>.
- 483 Miao Xie, Siguang Chen, and Chunli Lv. A component-based survey of interactions between large  
484 language models and multi-armed bandits. *arXiv preprint*, 2026.

# SUPPLEMENTARY MATERIALS

485

## 486 Contents

487	<b>A Algorithm: Bandit–OT Orchestration with Survival-Based Rewards</b>	<b>13</b>
488	A.1 i.i.d. Task Arrivals . . . . .	13
489	A.2 Non-i.i.d. Task Arrivals . . . . .	14
490	<b>B Appendix: Missing Proofs</b>	<b>16</b>
491	B.1 Proof of 5 in Theorem 4.1 (Lipschitz stability of alignment-adjusted rewards) . . . . .	16
492	B.2 Proof of 1 in Theorem 4.1 (Sublinear OT-Regret) . . . . .	17
493	B.3 Proof of 4 in Theorem 4.1 (Convergence of orchestration weights) . . . . .	19
494	B.4 Proof of 3 in Theorem 4.1 (Margin robustness under Gaussian noise) . . . . .	19
495	B.5 Proof of 5 in Theorem 4.1 (Consistency of reward estimates) . . . . .	20
496	<b>C Visualization of the Synthetic Datasets used for the Experiments</b>	<b>20</b>
497	<b>D Additional Synthetic Experiments</b>	<b>20</b>
498	D.1 Dataset and Task Description . . . . .	20
499	D.2 Baselines . . . . .	20
500	D.3 Evaluation Metrics . . . . .	22
501	D.4 Results. . . . .	22
502	D.5 Ablation Study: Sensitivity to the Alignment Penalty $\lambda$ . . . . .	22
503	<b>E Semi-Synthetic Experiment Settings</b>	<b>22</b>
504	<b>F Additional Semi-Synthetic Experiments and Figures</b>	<b>24</b>
505	F.1 Diagnostic Analysis . . . . .	24
506	F.2 Escalation Rate and Escalation Rate on Shifted Patients . . . . .	24
507	F.3 Ablation Study: Sensitivity to the Alignment Penalty $\lambda$ . . . . .	25
508	<b>A Algorithm: Bandit–OT Orchestration with Survival-Based Rewards</b>	
509	<b>A.1 i.i.d. Task Arrivals</b>	

---

**Algorithm 1** BOT-Orch: i.i.d. Task Version

---

- 1: **Input:** Agents  $\mathcal{A} = \{a_1, \dots, a_M\}$ , OT cost  $c$ , alignment weight  $\lambda \geq 0$ , horizon  $T$ , survival model  $S_i(\cdot)$ , learning rate  $\alpha$ , inverse temperature  $\eta_t$
- 2: Initialize estimated rewards  $\hat{r}_0(i) = 0$  and frailty  $\theta_0 = 1$  for all  $i \in \mathcal{A}$
- 3: **for**  $t = 1$  **to**  $T$  **do**
- 4:     Sample task  $x_t \sim \mathcal{P}_X$  ▷ i.i.d. task sampling
- 5:     Sample latent frailty  $\theta_t \sim p(\theta)$
- 6:     **for**  $i = 1$  **to**  $M$  **do**
- 7:         Compute agent output distribution  $\mu_i$
- 8:         Compute OT alignment cost  $\mathcal{W}_t(i) = W_c(\nu_t, \mu_i) + \epsilon_t(i)$
- 9:         Sample survival time  $T_t(i)$  and censoring  $\delta_t(i)$
- 10:         Compute reward:  $R_t(i) = \delta_t(i) S_i(T_t(i) | x_t)^{\theta_t}$
- 11:         Update estimated reward:  $\hat{r}_t(i) = \alpha \hat{r}_{t-1}(i) + (1 - \alpha) R_t(i)$
- 12:     **end for**
- 13:     Compute orchestration policy:

$$\pi_t(i) = \frac{\exp(\eta_t[\hat{r}_t(i) - \lambda \mathcal{W}_t(i)])}{\sum_{j=1}^M \exp(\eta_t[\hat{r}_t(j) - \lambda \mathcal{W}_t(j)])}$$

- 14:     Sample and execute agent  $i_t \sim \pi_t$
  - 15: **end for**
  - 16: **Output:** Policies  $\{\pi_t\}_{t=1}^T$  and cumulative reward
- 

510 **A.2 Non-i.i.d. Task Arrivals**

---

**Algorithm 2** BOT-Orch: Non-i.i.d. Task Version

---

- 1: **Input:** Same as Algorithm 1
- 2: Initialize  $\hat{r}_0(i) = 0$  and  $\theta_0 = 1$  for all  $i \in \mathcal{A}$
- 3: **for**  $t = 1$  **to**  $T$  **do**
- 4:     Sample task  $x_t \sim \mathcal{P}(x_t | \mathcal{H}_t)$  ▷ history-dependent, non-i.i.d.
- 5:     Sample latent frailty  $\theta_t \sim p(\theta_t | \mathcal{H}_t)$
- 6:     **for**  $i = 1$  **to**  $M$  **do**
- 7:         Compute agent output distribution  $\mu_i$
- 8:         Compute OT alignment cost  $\mathcal{W}_t(i) = W_c(\nu_t, \mu_i) + \epsilon_t(i)$
- 9:         Sample survival time  $T_t(i)$  and censoring  $\delta_t(i)$
- 10:         Compute reward:  $R_t(i) = \delta_t(i) S_i(T_t(i) | x_t)^{\theta_t}$
- 11:         Update estimated reward using history: ▷ temporal dependence from past rewards

$$\hat{r}_t(i) = \alpha \hat{r}_{t-1}(i) + (1 - \alpha) R_t(i) + f_i(\mathbf{R}_{1:t-1})$$

- 12:     **end for**
- 13:     Compute orchestration policy:

$$\pi_t(i) = \frac{\exp(\eta_t[\hat{r}_t(i) - \lambda \mathcal{W}_t(i)])}{\sum_{j=1}^M \exp(\eta_t[\hat{r}_t(j) - \lambda \mathcal{W}_t(j)])}$$

- 14:     Sample and execute agent  $i_t \sim \pi_t$
  - 15:     Update history  $\mathcal{H}_{t+1} = \mathcal{H}_t \cup \{(x_t, \mathbf{R}_t, \mathcal{W}_t)\}$
  - 16: **end for**
  - 17: **Output:** Policies  $\{\pi_t\}_{t=1}^T$  and cumulative reward
- 

511 **Derivation of Algorithms from the Theoretical Framework (Special Case Realisation).** The  
512 BOT-Orch algorithms in Algorithms 1 and 2 can be formally interpreted as special cases of the  
513 abstract framework introduced in Theorem 4.1. At the theoretical level, the model is defined in  
514 terms of an unobserved alignment-adjusted reward process

$$\tilde{r}_t(i) = \mathbb{E}[R_t(i) | x_t] - \lambda W_c(\mu_i, \nu_t),$$

515 over which the learning dynamics are characterised via exponential-weights updates on the simplex.  
516 The algorithms instantiate this framework by specifying an explicit stochastic realisation of the re-  
517 ward process together with a consistent estimator of its conditional expectation. In particular, the  
518 survival–frailty construction generates bounded random variables  $R_t(i)$  whose conditional expect-  
519 ation coincides with the abstract reward functional assumed in the theory, thereby embedding the  
520 model within a well-defined stochastic process satisfying assumptions (A2)–(A4).

521 The empirical quantity  $\hat{r}_t(i)$ , defined via exponential smoothing in the i.i.d. case and augmented with  
522 a history-dependent correction term in the non-i.i.d. case, constitutes a Robbins–Monro stochastic  
523 approximation of  $\mathbb{E}[R_t(i) \mid x_t]$ , ensuring asymptotic consistency under the respective dependence  
524 structures. Substituting this estimator into the theoretical objective yields a computable approxima-  
525 tion of  $\tilde{r}_t(i)$ , while the entropy-regularised optimisation over  $\Delta^{|\mathcal{A}|-1}$  induces the softmax policy  
526 used in the algorithm.

527 Consequently, the i.i.d. algorithm corresponds to the stationary special case in which  $(x_t, \theta_t)$  are in-  
528 dependent draws and the induced reward process is temporally homogeneous, whereas the non-i.i.d.  
529 algorithm generalises this construction to an adapted filtration  $\mathcal{H}_t$ , allowing for history-dependent  
530 task and frailty evolution while preserving boundedness and measurability of the reward sequence.  
531 In both cases, the algorithm complements by providing an explicit implementation of the main theo-  
532 rem, showing that BOT-Orch is a realised instance of the general OT-regularised exponential-weights  
533 framework under different assumptions on the data-generating process.

534 **B Appendix: Missing Proofs**

535 This appendix contains full proofs of the statements in Section 4.

536 **B.1 Proof of 5 in Theorem 4.1 (Lipschitz stability of alignment-adjusted rewards)**

537 **Lemma B.1** (Stability of alignment-adjusted rewards). *Assume (A1)–(A2). Let  $\nu, \nu' \in \mathcal{P}(\mathcal{Y})$  be two*  
 538 *probability measures on  $\mathcal{Y}$ , and let  $c : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$  be  $L$ -Lipschitz. For each agent  $i \in \mathcal{A}$ , define*  
 539 *the alignment-adjusted reward*

$$\tilde{r}(i; \nu) := \hat{r}(i) - \lambda W_c(\mu_i, \nu),$$

540 where  $W_c$  is the Wasserstein distance with ground cost  $c$ . Then,

$$|\tilde{r}(i; \nu) - \tilde{r}(i; \nu')| \leq \lambda L W_1(\nu, \nu'),$$

541 where  $W_1$  is the 1-Wasserstein distance induced by  $c$ . In particular, if  $W_c(\nu, \nu') \leq \varepsilon$ , then

$$|\tilde{r}(i; \nu) - \tilde{r}(i; \nu')| \leq \lambda L \varepsilon.$$

542 *Proof.* By definition,

$$\tilde{r}(i; \nu) - \tilde{r}(i; \nu') = -\lambda(W_c(\mu_i, \nu) - W_c(\mu_i, \nu')).$$

543 To bound the difference, we use the standard Lipschitz property of the Wasserstein distance: if the  
 544 ground cost  $c$  is  $L$ -Lipschitz, then for any measures  $\nu, \nu'$ ,

$$|W_c(\mu_i, \nu) - W_c(\mu_i, \nu')| \leq L W_1(\nu, \nu'),$$

545 where  $W_1$  is the 1-Wasserstein distance (see, e.g., [Villani et al., 2008]).

546 Combining the two inequalities yields

$$|\tilde{r}(i; \nu) - \tilde{r}(i; \nu')| \leq \lambda L W_1(\nu, \nu') \leq \lambda L \varepsilon,$$

547 as claimed. □

548 **Lemma B.2** (Concentration under multiplicative frailty). *Assume (A2)–(A4), and let  $S_i(t | x)$  de-*  
 549 *note the baseline survival function, uniformly bounded:  $0 \leq S_i(t | x) \leq 1$  for all  $t$  and  $x \in \mathcal{X}$ . Let*  
 550  *$R_t(i)$  denote the frailty-adjusted reward for agent  $i$  at round  $t$ . Then, for every  $\varepsilon > 0$ ,*

$$\mathbb{P}\left(|R_t(i) - \mathbb{E}[R_t(i) | x_t]|\right) > \varepsilon \leq 2 \exp\left(-\frac{C \varepsilon^2}{1 + \text{Var}(\theta_t)}\right),$$

551 for some constant  $C > 0$  independent of  $t$ .

552 *Proof.* Fix  $x_t = x$  and an agent  $i \in \mathcal{A}$ . Let  $(\theta, T, \delta)$  denote a realization of the generative process,  
 553 where  $\theta$  is the frailty,  $T$  is the survival time with conditional survival function  $S_i(\cdot | x)$ , and  $\delta$  is the  
 554 censoring indicator, independent of  $T$  conditional on  $(x, \theta)$ . Define

$$R := \delta S_i(T | x)^\theta.$$

555 Since  $0 \leq S_i(T | x) \leq 1$  and  $\delta \in \{0, 1\}$ , it follows that

$$0 \leq R \leq 1.$$

556 Define the conditional expectation given  $\theta$ :

$$g(\theta) := \mathbb{E}[R | x, \theta] = \mathbb{E}[\delta S_i(T | x)^\theta | x, \theta],$$

557 so that  $0 \leq g(\theta) \leq 1$ . Then

$$R - \mathbb{E}[R | x] = (R - g(\theta)) + (g(\theta) - \mathbb{E}[R | x]). \tag{1}$$

558 Conditional on  $\theta$ ,  $R$  is bounded in  $[0, 1]$ , and hence by Hoeffding's lemma,

$$\mathbb{E}\left[e^{\lambda(R - g(\theta))} \mid x, \theta\right] \leq \exp\left(\frac{\lambda^2}{8}\right), \quad \forall \lambda \in \mathbb{R}.$$

559 Observe that

$$g(\theta) - \mathbb{E}[R | x] = \mathbb{E}[R | x, \theta] - \mathbb{E}[R | x].$$

560 Since  $0 \leq S_i(T | x) \leq 1$ , the map  $\theta \mapsto g(\theta)$  is Lipschitz with constant at most 1, i.e.,

$$|g(\theta_1) - g(\theta_2)| \leq |\theta_1 - \theta_2|, \quad \forall \theta_1, \theta_2 \geq 0.$$

561 By assumption (A3),  $\theta$  has finite variance and finite exponential moments. Therefore,  $g(\theta) - \mathbb{E}[R | x]$   
562 is sub-Gaussian with variance proxy  $\nu_2^2 \leq \text{Var}(\theta_t)$ .

563 By (1),  $R - \mathbb{E}[R | x]$  is the sum of two independent sub-Gaussian variables:  $(R - g(\theta))$  and  
564  $(g(\theta) - \mathbb{E}[R | x])$ . Therefore,  $R - \mathbb{E}[R | x]$  is sub-Gaussian with variance proxy

$$\nu^2 := \nu_1^2 + \nu_2^2 \leq \frac{1}{4} + \text{Var}(\theta_t) \leq C(1 + \text{Var}(\theta_t)),$$

565 for some universal constant  $C > 0$ .

566 Hence, by the standard sub-Gaussian tail bound, for all  $\varepsilon > 0$ ,

$$\mathbb{P}\left(|R - \mathbb{E}[R | x]| > \varepsilon\right) \leq 2 \exp\left(-\frac{\varepsilon^2}{2\nu^2}\right) \leq 2 \exp\left(-\frac{C\varepsilon^2}{1 + \text{Var}(\theta_t)}\right),$$

567 which proves the claim for a single observation.

568 The extension to empirical averages over independent draws follows from standard concentration  
569 inequalities for sums of independent sub-Gaussian variables.  $\square$

## 570 **B.2 Proof of 1 in Theorem 4.1 (Sublinear OT-Regret)**

571 **Theorem B.3** (Sublinear OT-Regret). *Under assumptions (A1)–(A5), the BOT-Orch softmax policy*  
572 *with an appropriate temperature schedule satisfies*

$$\mathcal{R}_T = O(\sqrt{T}).$$

573 *Proof.* We give a potential-based analysis of the exponential-weights (softmax) policy applied to  
574 bounded adjusted rewards. Let  $\tilde{r}_t(i)$  denote the (possibly noisy) adjusted reward received at time  $t$   
575 by agent  $i$ . By assumption (A2) and the boundedness of the cost  $W_c$  from (A1), there exist constants  
576  $a < b$  such that

$$\tilde{r}_t(i) \in [a, b] \quad \text{for all } t \text{ and } i.$$

577 Define  $\Delta := b - a$ .

578 The algorithm maintains weights

$$w_{t+1}(i) = w_t(i) \exp(\eta_t \tilde{r}_t(i)), \quad \pi_t(i) = \frac{w_t(i)}{\sum_j w_t(j)},$$

579 and we introduce the log-partition potential

$$\Phi_t := \log \sum_i w_t(i).$$

580 **Step 1: Constant inverse temperature.** We first consider the case  $\eta_t \equiv \eta > 0$ . Standard calcula-  
581 tions for the exponential-weights forecaster [Cesa-Bianchi and Lugosi, 2006, Section 2.2], based on  
582 Hoeffding's inequality, yield the one-step bound

$$\Phi_{t+1} - \Phi_t = \log \sum_i \pi_t(i) \exp(\eta \tilde{r}_t(i)) \leq \eta \sum_i \pi_t(i) \tilde{r}_t(i) + \frac{\eta^2 \Delta^2}{8}.$$

583 Summing over  $t = 1, \dots, T$  gives

$$\Phi_{T+1} - \Phi_1 \leq \eta \sum_{t=1}^T \mathbb{E}_{i \sim \pi_t} [\tilde{r}_t(i)] + \frac{\eta^2 T \Delta^2}{8}.$$

584 On the other hand, for any fixed agent  $i^*$ ,

$$\Phi_{T+1} \geq \log w_{T+1}(i^*) = \log w_1(i^*) + \eta \sum_{t=1}^T \tilde{r}_t(i^*).$$

585 Assuming uniform initialization  $w_1(i) = 1$  for all  $i$ , we have  $\log w_1(i^*) = 0$ . Combining the two  
586 inequalities and rearranging yields

$$\sum_{t=1}^T \tilde{r}_t(i^*) - \sum_{t=1}^T \mathbb{E}_{i \sim \pi_t}[\tilde{r}_t(i)] \leq \frac{\log M}{\eta} + \frac{\eta T \Delta^2}{8},$$

587 where  $M$  denotes the number of agents.

588 Optimizing the right-hand side by choosing

$$\eta = \sqrt{\frac{8 \log M}{T \Delta^2}}$$

589 gives

$$\sum_{t=1}^T \tilde{r}_t(i^*) - \sum_{t=1}^T \mathbb{E}_{i \sim \pi_t}[\tilde{r}_t(i)] \leq \Delta \sqrt{\frac{T \log M}{2}}.$$

590 This bound holds deterministically for any reward sequence bounded in  $[a, b]$ . Taking expectations  
591 in the presence of stochastic rewards yields

$$\mathbb{E}[\mathcal{R}_T] \leq \Delta \sqrt{\frac{T \log M}{2}},$$

592 which implies  $\mathcal{R}_T = O(\sqrt{T})$  since  $M$  is fixed.

593 **Step 2: Time-varying inverse temperature.** We now consider a time-varying schedule  $\eta_t \asymp t^{-1/2}$ ,  
594 as used in Algorithms 1–2. In this case, the potential increment satisfies

$$\Phi_{t+1} - \Phi_t \leq \eta_t \sum_i \pi_t(i) \tilde{r}_t(i) + \frac{\eta_t^2 \Delta^2}{8}.$$

595 Summing over  $t$  and comparing with  $\log w_{T+1}(i^*)$  yields

$$\sum_{t=1}^T \tilde{r}_t(i^*) - \sum_{t=1}^T \mathbb{E}_{i \sim \pi_t}[\tilde{r}_t(i)] \leq \frac{\log M}{\eta_T} + \sum_{t=1}^T \frac{\eta_t \Delta^2}{8}.$$

596 Choosing  $\eta_t = c/\sqrt{t}$  for a suitable constant  $c > 0$  gives

$$\frac{\log M}{\eta_T} = O(\sqrt{T}), \quad \sum_{t=1}^T \eta_t = O(\sqrt{T}),$$

597 and therefore

$$\mathcal{R}_T = \tilde{O}(\sqrt{T}).$$

598 Finally, suppose that the adjusted rewards  $\tilde{r}_t(i)$  are stochastic and admit the decomposition

$$\tilde{r}_t(i) = \mathbb{E}[\tilde{r}_t(i) \mid \mathcal{F}_{t-1}] + \xi_t(i),$$

599 where  $\{\mathcal{F}_t\}_{t \geq 0}$  is the natural filtration generated by the history of the algorithm and  $\{\xi_t(i)\}$  is a  
600 martingale difference sequence satisfying

$$\mathbb{E}[\xi_t(i) \mid \mathcal{F}_{t-1}] = 0, \quad |\xi_t(i)| \leq \Delta \quad \text{a.s.}$$

601 uniformly over  $t$  and  $i$ . Such a decomposition covers randomness arising from empirical reward  
602 estimation, survival-time sampling, and censoring.

603 Applying the Azuma–Hoeffding inequality yields, for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ ,

$$\left| \sum_{t=1}^T \xi_t(i) \right| \leq \Delta \sqrt{2T \log(1/\delta)}.$$

604 Consequently, the deviation between cumulative realized rewards and their conditional expectations  
605 is of order  $O(\sqrt{T})$  almost surely up to logarithmic factors. Incorporating these fluctuations into the  
606 regret analysis contributes only lower-order terms and does not alter the overall  $\sqrt{T}$  regret rate.  $\square$

607 **Proposition B.4** (OT dominance). *If  $\hat{r}_t(i) = \hat{r}_t(j)$  and  $W_c(\mu_i, \nu_t) < W_c(\mu_j, \nu_t)$  then for any  $\lambda > 0$*   
 608 *we have  $\tilde{r}_t(i) > \tilde{r}_t(j)$ .*

609 *Proof.* By direct algebra:

$$\tilde{r}_t(i) - \tilde{r}_t(j) = (\hat{r}_t(i) - \hat{r}_t(j)) - \lambda(W_c(\mu_i, \nu_t) - W_c(\mu_j, \nu_t)).$$

610 The hypothesis  $\hat{r}_t(i) = \hat{r}_t(j)$  makes the first term zero; since  $W_c(\mu_i, \nu_t) - W_c(\mu_j, \nu_t) < 0$ , multi-  
 611 plying by  $-\lambda$  gives a strictly positive number. Thus  $\tilde{r}_t(i) - \tilde{r}_t(j) > 0$ , proving the claim.  $\square$

### 612 **B.3 Proof of 4 in Theorem 4.1 (Convergence of orchestration weights)**

613 **Theorem B.5** (Almost-sure convergence of orchestration weights). *Assume (A1)–(A5) hold. Let*  
 614  *$(\gamma_t)_{t \geq 0}$  be a deterministic step-size sequence satisfying*

$$\gamma_t > 0, \quad \sum_{t=0}^{\infty} \gamma_t = \infty, \quad \sum_{t=0}^{\infty} \gamma_t^2 < \infty.$$

615 *Then the stochastic process  $(\phi_t)$  converges almost surely to an asymptotically stable equilibrium of*  
 616 *the ordinary differential equation*

$$\dot{\phi} = \text{Softmax}(\mathbb{E}[\tilde{r}]) - \phi.$$

617 *Proof.* Let  $(\mathcal{F}_t)_{t \geq 0}$  denote the natural filtration of  $(\phi_t)$ . The stochastic recursion can be written in  
 618 Robbins–Monro form:

$$\phi_{t+1} = \phi_t + \gamma_t(H(\phi_t) + \xi_{t+1}),$$

619 where

$$H(\phi) := \text{Softmax}(\mathbb{E}[\tilde{r} \mid \phi]) - \phi, \quad \xi_{t+1} := \text{Softmax}(\tilde{r}_t) - \text{Softmax}(\mathbb{E}[\tilde{r}_t \mid \mathcal{F}_t]).$$

620 Then  $(\xi_t)$  is a martingale-difference sequence with respect to  $(\mathcal{F}_t)$  and satisfies  $\sup_t \mathbb{E}[\|\xi_{t+1}\|^2 \mid$   
 621  $\mathcal{F}_t] < \infty$  by boundedness of  $\tilde{r}$ . The iterates remain in the compact simplex  $\Delta(\mathcal{A})$ , and  $H$  is  
 622 Lipschitz continuous on  $\Delta(\mathcal{A})$ .

623 Let  $\bar{\phi}$  denote the continuous-time piecewise-linear interpolation of  $(\phi_t)$ . By Theorem 3.2 of  
 624 [Benaïm, 1999],  $\bar{\phi}$  is an asymptotic pseudo-trajectory of the ODE

$$\dot{\phi} = H(\phi),$$

625 and by Theorem 3.9 of [Benaïm, 1999], the almost-sure limit set of  $(\phi_t)$  is contained in the set of  
 626 internally chain-transitive invariant sets of this ODE.

627 If the ODE admits only isolated asymptotically stable equilibria (e.g., when  $\mathbb{E}[\tilde{r}]$  induces a strictly  
 628 concave potential), each internally chain-transitive set reduces to a single equilibrium. Hence,  $(\phi_t)$   
 629 converges almost surely to an asymptotically stable equilibrium of

$$\dot{\phi} = \text{Softmax}(\mathbb{E}[\tilde{r}]) - \phi.$$

630  $\square$

### 631 **B.4 Proof of 3 in Theorem 4.1 (Margin robustness under Gaussian noise)**

632 **Lemma B.6** (Margin robustness, ). *Let  $\widetilde{W}_c(\mu_i, \nu) = W_c(\mu_i, \nu) + \epsilon_i$ , with independent  $\epsilon_i \sim$*   
 633  *$\mathcal{N}(0, \sigma^2)$ . Let  $\Delta_{ij} = W_c(\mu_j, \nu) - W_c(\mu_i, \nu)$ . If  $|\Delta_{ij}| > \sigma\sqrt{2 \log 2}$  then*

$$\mathbb{P}(\tilde{r}(i) < \tilde{r}(j)) < 1/4.$$

634 *Proof.* Assume without loss of generality that  $\Delta_{ij} > 0$  (the other sign is symmetric). The noisy  
 635 observed difference is

$$\widetilde{\Delta}_{ij} = \widetilde{W}_c(\mu_j, \nu) - \widetilde{W}_c(\mu_i, \nu) = \Delta_{ij} + Z,$$

636 where  $Z = \epsilon_j - \epsilon_i \sim \mathcal{N}(0, 2\sigma^2)$  by independence. The event that  $\tilde{r}(i) < \tilde{r}(j)$  (ignoring small  
 637 differences from empirical reward estimation) is equivalent to  $\tilde{\Delta}_{ij} > 0$ . Thus

$$\mathbb{P}(\tilde{r}(i) < \tilde{r}(j)) = \mathbb{P}(Z > -\Delta_{ij}) = \Phi\left(\frac{\Delta_{ij}}{\sqrt{2}\sigma}\right),$$

638 where  $\Phi$  is the standard normal CDF. The condition  $\Delta_{ij} > \sigma\sqrt{2\log 2}$  gives

$$\frac{\Delta_{ij}}{\sqrt{2}\sigma} > \sqrt{\log 2},$$

639 hence

$$\Phi\left(\frac{\Delta_{ij}}{\sqrt{2}\sigma}\right) < \Phi(\sqrt{\log 2}) < \frac{1}{4},$$

640 where the last inequality follows from numerical evaluation of the Gaussian CDF (or standard Gaus-  
 641 sian tail bounds such as  $\Phi(x) \leq \frac{1}{2}e^{-x^2/2}$  for  $x > 0$ ). This concludes the proof.  $\square$

## 642 B.5 Proof of 5 in Theorem 4.1 (Consistency of reward estimates)

643 **Theorem B.7** (Consistency of reward estimates). *Under (A2)–(A4) and continuity of  $S_i(\cdot | x_t)$ , the*  
 644 *empirical estimator  $\hat{r}_t(i)$  satisfies*

$$\sup_{i \in \mathcal{A}} |\hat{r}_t(i) - \mathbb{E}[R_t(i) | x_t]| \xrightarrow[t \rightarrow \infty]{\text{a.s.}} 0.$$

645 *Proof. Case 1: i.i.d. tasks.* Let  $x_t \equiv x$ ,  $\{R_t(i)\}_{t \geq 1}$  i.i.d., bounded by  $R_{\max}$ . By the strong law of  
 646 large numbers:

$$\hat{r}_t(i) = \frac{1}{t} \sum_{s=1}^t R_s(i) \xrightarrow[t \rightarrow \infty]{\text{a.s.}} \mathbb{E}[R(i) | x], \quad \forall i.$$

647 Uniformity over finite  $\mathcal{A}$  follows from the union bound:

$$\sup_{i \in \mathcal{A}} |\hat{r}_t(i) - \mathbb{E}[R(i) | x]| \xrightarrow[t \rightarrow \infty]{\text{a.s.}} 0.$$

648 **Case 2: stationary ergodic (non i.i.d) tasks.** Let  $\{(x_t, R_t(i))\}_{t \geq 1}$  stationary ergodic. By  
 649 Birkhoff's ergodic theorem [Walters, 1982, Theorem 1.14]

$$\frac{1}{t} \sum_{s=1}^t R_s(i) \xrightarrow[t \rightarrow \infty]{\text{a.s.}} \mathbb{E}[R(i) | x_t] \quad \forall i.$$

650 Uniformity over finite  $\mathcal{A}$  via the union bound.

651 Let  $\hat{r}_t(i) = \sum_{s=1}^t \gamma_{s-1} R_s(i)$  with  $\sum \gamma_s = \infty$ ,  $\sum \gamma_s^2 < \infty$ . Ergodicity and continuity of  $S_i$  imply

$$\hat{r}_t(i) \xrightarrow[t \rightarrow \infty]{\text{a.s.}} \mathbb{E}[R(i) | x_t], \quad \sup_{i \in \mathcal{A}} |\hat{r}_t(i) - \mathbb{E}[R(i) | x_t]| \rightarrow 0.$$

652  $\square$

## 653 C Visualization of the Synthetic Datasets used for the Experiments

### 654 D Additional Synthetic Experiments

#### 655 D.1 Dataset and Task Description

656 Datasets and tasks are as in 6.

#### 657 D.2 Baselines

658 Baselines are as in Section 6.

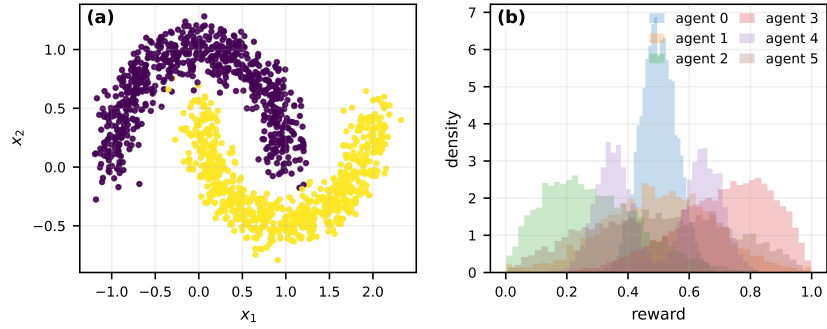


Figure 2: **IID data generation** (a) IID task contexts sampled from a half-moons distribution, illustrating a stationary but structured task space. (b) IID reward distributions for  $K$  agents, where rewards are independently drawn over time from fixed distributions with a matched mean (approximately 0.5) but heterogeneous higher-order properties (e.g., variance, skewness, and bimodality).

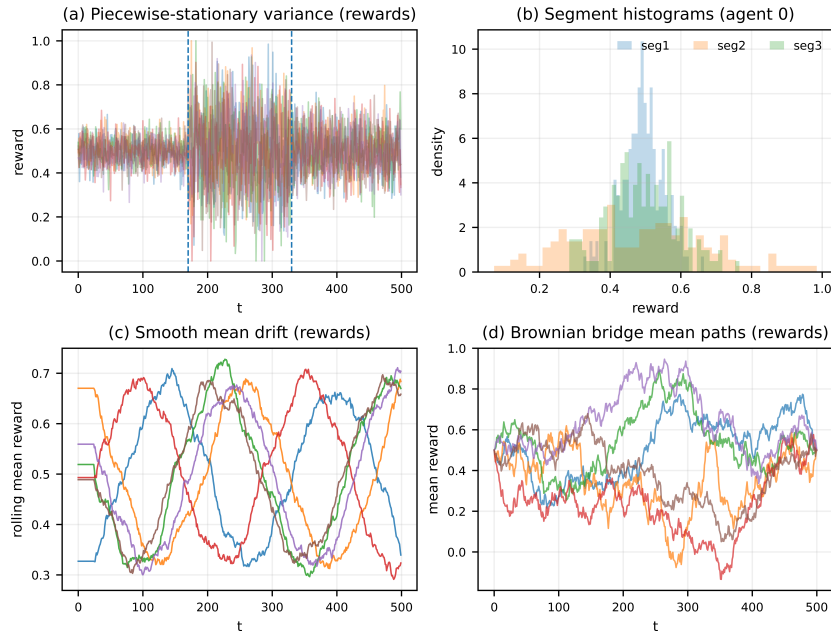


Figure 3: **Non-IID data generation** (a) Piecewise-stationary rewards where the variance changes at unknown changepoints while the mean remains fixed, shown as reward trajectories over time across agents. (b) Segment-wise reward histograms for a representative agent, revealing distributional shifts induced by variance changepoints. (c) Smooth-drift setting where agent reward means evolve gradually according to a sinusoidal drift (shown via rolling mean trajectories), inducing continuous non-stationarity. (d) Brownian-bridge setting illustrating temporally correlated latent mean paths constrained to fixed endpoints, producing structured stochastic dependence across time.

Environment	Event Rate				Mean Observed Time			
	BOT-Orch	No-OT ( $\lambda=0$ )	Random	UCB1 (MAB)	BOT-Orch	No-OT ( $\lambda=0$ )	Random	UCB1 (MAB)
<b>IID</b>								
IID-G	<b>0.63±0.02</b>	0.60±0.02	0.58±0.04	0.54±0.03	<b>0.59±0.04</b>	0.66±0.04	0.72±0.04	0.71±0.03
IID-M	<b>0.65±0.04</b>	0.60±0.04	0.60±0.03	0.59±0.03	<b>0.58±0.02</b>	0.66±0.04	0.69±0.03	0.70±0.05
<b>Non-IID</b>								
NonIID-BB	<b>0.67±0.04</b>	0.61±0.03	0.61±0.02	0.62±0.05	<b>0.55±0.09</b>	0.62±0.01	0.65±0.04	0.63±0.03
NonIID-PS	<b>0.66±0.04</b>	0.64±0.02	0.59±0.04	0.59±0.02	<b>0.56±0.04</b>	0.62±0.04	0.67±0.03	0.69±0.04
NonIID-SD	<b>0.65±0.03</b>	0.64±0.03	0.60±0.04	0.57±0.02	<b>0.59±0.03</b>	0.64±0.05	0.67±0.03	0.67±0.04

Table 4: Event Rate and Mean Observed Time (mean  $\pm$  95% CI across 5 seeds;  $T = 200$ ). Best in bold.

### 659 D.3 Evaluation Metrics

- 660 • **Event rate.** Under survival-style feedback with censoring indicator  $\delta_t(i_t) \in \{0, 1\}$ , we report  
661  $\frac{1}{T} \sum_{t=1}^T \delta_t(i_t)$ , i.e., the fraction of rounds with uncensored (fully observed) outcomes. Higher  
662 event rates indicate less censoring and more informative feedback.
- 663 • **Mean observed time.** We report the mean observed time  $\frac{1}{T} \sum_{t=1}^T T_t^{\text{obs}}(i_t)$ , where  
664  $T_t^{\text{obs}}(i_t) = \min\{T_t(i_t), C_t(i_t)\}$  under right censoring. This metric summarizes the typical ob-  
665 served completion time under the censoring mechanism.

### 666 D.4 Results.

667 Table 4 reports event rate and mean observed time. BOT-Orch achieves the highest event rate and  
668 lowest mean observed time across all environments, in both IID and non-IID settings, consistently  
669 outperforming No-OT, Random, and UCB1.

### 670 D.5 Ablation Study: Sensitivity to the Alignment Penalty $\lambda$

671 The alignment penalty weight  $\lambda$  controls the trade-off between exploitation of historical reward es-  
672 timates  $\hat{r}_t(i)$  and adherence to OT-based distributional alignment  $W_t(i)$ . We conduct a grid search  
673 over  $\lambda \in \{0.0, 0.5, 1.0, 1.25, 1.5, 1.75, 2.0, 3.0, 5.0, 10.0\}$ , running BOT-Orch for 30 seeds un-  
674 der both the IID (Algorithm 1) and Non-IID (Algorithm 2) conditions. The No-OT, Random, and  
675 UCB1 baselines serve as fixed reference lines since they are independent of  $\lambda$ . Table 5 summarise  
676 the results.

## 677 E Semi-Synthetic Experiment Settings

678 **Agents.** We set  $M = 2$  agents:

- 679 • **Agent 0 (AI):** a logistic regression classifier (L2,  $C=1.0$ ) trained on the training split and cal-  
680 ibrated using Platt scaling (isotonic regression) on the calibration split. Accuracy: 98.2% in-  
681 distribution, 80.7% under shift.
- 682 • **Agent 1 (Human):** a simulated clinical expert with complementary accuracy 88.0% on in-  
683 distribution patients, 94.7% on shifted patients.

684 The human expert is more accurate on the patients the AI handles worst, confirming a positive  
685 complementarity gap  $\Delta_{\text{comp}} > 0$ .

686 **Tasks and reward.** At each round  $t$ , a patient biopsy  $x_t \in \mathcal{X}$  arrives. The reward is binary  
687 correctness:  $R_t(i) = \mathbf{1}[\text{agent } i \text{ classifies patient } t \text{ correctly}]$ . This is bounded in  $[0, 1]$ , satisfying  
688 Assumption (A2). We use bandit feedback throughout and only the chosen agent’s reward is ob-  
689 served.

**OT alignment costs.** We use the output-space OT alignment cost on the binary label simplex  $\{0, 1\}$  with ground cost

$$C = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$\lambda$	IID (Algorithm 1)			Non-IID (Algorithm 2)		
	Cum. net	Regret	Team acc.	Cum. net	Regret	Team acc.
0	-311.34±3.43	122.61±2.49	0.177±0.013	-722.67±47.12	317.61±21.59	0.177±0.013
0.5	-271.42±2.79	82.68±1.92	0.336±0.011	-549.11±33.07	144.05±5.85	0.418±0.017
1	-262.29±2.69	73.56±1.78	0.380±0.011	-531.17±31.69	126.11±4.86	0.453±0.017
1.25	-259.93±2.57	71.19±1.78	0.386±0.011	-528.55±31.48	123.49±5.23	0.459±0.019
1.5	-258.52±2.62	69.78±1.72	0.389±0.011	-526.24±31.88	121.18±5.18	0.462±0.018
1.75	-257.98±2.27	69.24±1.47	0.398±0.009	-524.10±31.77	119.04±5.49	0.470±0.019
2	-256.20±2.55	67.47±1.64	0.403±0.010	-524.22±31.76	119.16±5.05	0.469±0.018
3	-253.76±2.49	65.02±1.60	0.412±0.011	-519.79±31.65	114.72±4.79	0.479±0.019
5	-253.28±2.04	64.54±1.67	0.413±0.010	-515.97±32.03	110.91±5.11	0.484±0.020
10	-251.99±2.16	63.25±1.46	0.420±0.009	-514.23±31.32	109.16±4.67	0.493±0.018
11	-252.06±2.22	63.32±1.47	0.420±0.009	-513.61±31.26	108.54±4.71	0.495±0.018
12	-252.05±2.28	63.32±1.51	0.420±0.009	-512.93±31.33	107.86±4.73	0.495±0.018
13	<b>-251.60±2.33</b>	<b>62.86±1.60</b>	<b>0.422±0.009</b>	<b>-513.01±31.01</b>	<b>107.95±5.00</b>	<b>0.494±0.018</b>
14	-251.71±2.42	62.97±1.50	0.423±0.009	-513.16±30.99	108.10±5.07	0.494±0.018
15	-251.80±2.45	63.06±1.56	0.423±0.008	-513.45±31.02	108.39±5.01	0.493±0.018
<i>Reference baselines (<math>\lambda</math>-independent)</i>						
No-OT	-311.34±3.43	122.61±2.49	0.177±0.013	-722.67±47.12	317.61±21.59	0.177±0.013
Random	-315.75±3.74	127.02±3.11	0.177±0.010	-731.00±50.00	325.93±22.38	0.169±0.010
UCB1	-313.22±3.92	124.48±2.93	0.177±0.012	-722.86±48.17	317.79±21.51	0.178±0.009

Table 5:  $\lambda$  grid search results for BOT-Orch on synthetic tasks. Mean  $\pm$  95% CI across 30 seeds,  $T = 114$ , evaluated with  $\lambda_{eval} = 1.0$ .  $\lambda = 0$  reproduces No-OT by construction (sanity check). Selected  $\lambda^* = 13$  (bold) maximizes average Cum. net across IID/Non-IID panels. Reference baselines at the bottom are  $\lambda$ -independent.

690 (0-1 loss). Under this cost, the Wasserstein distance between the true label one-hot  $\nu_t = \delta_{y_t}$  and  
691 each agent’s predictive distribution admits the closed form:

$$W_t(\text{AI}) = 1 - R_t(\text{AI}) \quad (\text{probability AI is wrong on patient } t), \quad (2)$$

$$W_t(\text{human}) = 1 - p_h(x_t) \quad (\text{probability human is wrong on patient } t), \quad (3)$$

692 where  $p_h(x_t)$  is the human’s accuracy for that patient’s shift status. This is a direct instantiation  
693 of the general alignment cost  $W_c(\nu_t, \mu_i)$  defined in Section 3.4. The closed form follows from the  
694 exactness of Sinkhorn transport on the  $2 \times 2$  binary simplex.

695 **OT Dominance verification.** With these costs, under distribution shift:

$$W_t(\text{AI}) \approx 0.193 \quad \text{vs.} \quad W_t(\text{human}) \approx 0.053.$$

696 The gap  $W_t(\text{AI}) - W_t(\text{human}) = 0.140$  is large and positive on shifted patients, confirming the  
697 precondition of Theorem 4.2 Part 2 that the human has strictly lower alignment cost on shifted  
698 patients and should be preferred by the BOT-Orch policy. On in-distribution patients the relationship  
699 is reversed ( $W_t(\text{AI}) \approx 0.018 < 0.120 = W_t(\text{human})$ ), so the AI is correctly preferred there. This is  
700 exactly the complementarity structure that BOT-Orch is designed to exploit without being told about  
701 the shift.

## 702 F Additional Semi-Synthetic Experiments and Figures

### 703 F.1 Diagnostic Analysis

704 Figure 4 and 5 provide diagnostic analyses.

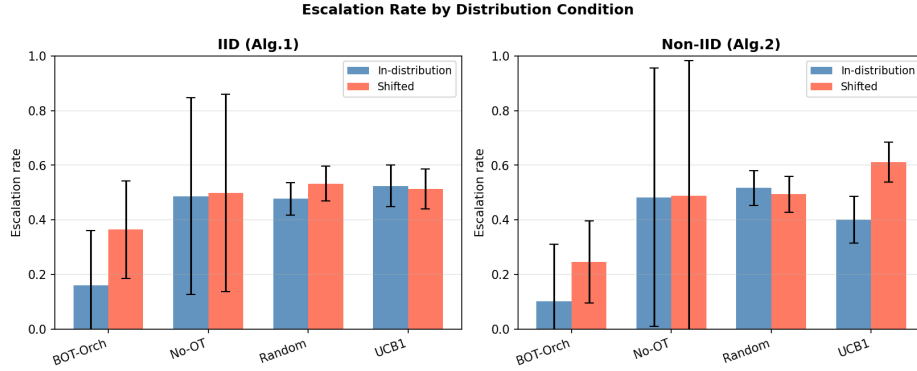


Figure 4: **Escalation rate by distribution condition.** Mean escalation rate for in-distribution patients (blue) and shifted patients (red) per method. Error bars show standard deviation across seeds. *Left*: IID condition; *right*: Non-IID condition. BOT-Orch achieves a higher escalation rate on shifted patients relative to in-distribution patients compared to all baselines, demonstrating targeted routing. No-OT’s large error bars reflect bimodal behaviour across seeds: some runs converge to always-AI and others to always-human, due to the cold-start problem in the binary bandit without OT regularisation.

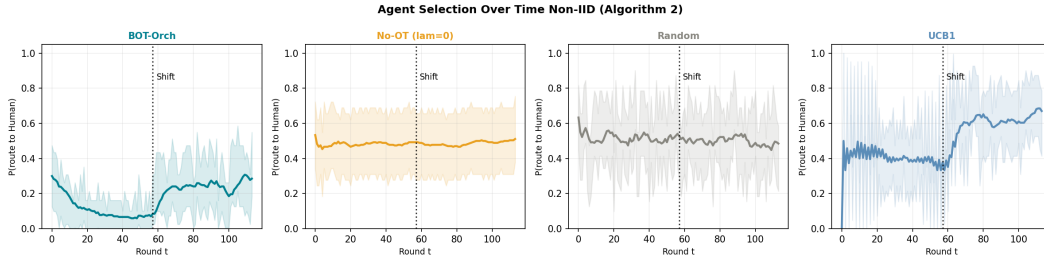


Figure 5: **Agent selection trajectories, Non-IID condition (Algorithm 2).** Rolling probability of routing to the human expert (window  $w=8$  rounds) for each method over the deployment episode. The dotted vertical line marks the shift onset at round 57. BOT-Orch (top-left): routing probability rises after the shift, showing adaptation. No-OT (top-right): high variance, consistent with the bimodal cold-start behaviour observed in Table 3. UCB1 (bottom-right): also adapts but less precisely, routing more patients to the human without the alignment-cost targeting of BOT-Orch. Random (bottom-left): flat at 50% throughout as expected.

### 705 F.2 Escalation Rate and Escalation Rate on Shifted Patients

We present the escalation rate and escalation rate on shifted patients in Table 6

	Escalation Rate				Escalation Rate on Shifted Patients			
	BOT-Orch	No-OT ( $\lambda=0$ )	Random	UCB1	BOT-Orch	No-OT ( $\lambda=0$ )	Random	UCB1
<b>IID</b>	0.214 $\pm$ 0.066	0.493 $\pm$ 0.360	0.505 $\pm$ 0.046	0.519 $\pm$ 0.061	0.283 $\pm$ 0.085	0.499 $\pm$ 0.362	0.533 $\pm$ 0.064	0.514 $\pm$ 0.073
<b>Non-IID</b>	0.192 $\pm$ 0.023	0.485 $\pm$ 0.484	0.505 $\pm$ 0.046	0.506 $\pm$ 0.062	0.209 $\pm$ 0.019	0.488 $\pm$ 0.495	0.494 $\pm$ 0.065	0.612 $\pm$ 0.074

Table 6: Deployment metrics across all four methods and both experimental conditions. Mean  $\pm$  95% CI across 30 seeds,  $T=114$ ,  $\lambda=3.0$ . *Esc. rate*: fraction of patients routed to human. *Esc.(shift)*: escalation rate on shifted patients only. IID condition uses Algorithm 1; Non-IID uses Algorithm 2 (ID patients rounds 1–57, shifted rounds 58–114).

706

707 **E.3 Ablation Study: Sensitivity to the Alignment Penalty  $\lambda$**

708 The alignment penalty weight  $\lambda$  controls the trade-off between exploitation of historical reward estimates  $\hat{r}_t(i)$  and adherence to OT-based distributional alignment  $W_t(i)$ . We conduct a grid search  
 709 over  $\lambda \in \{0.0, 0.5, 1.0, 1.25, 1.5, 1.75, 2.0, 3.0, 5.0, 10.0\}$ , running BOT-Orch for 30 seeds under both the IID (Algorithm 1) and Non-IID (Algorithm 2) conditions. The No-OT, Random, and  
 710 UCB1 baselines serve as fixed reference lines since they are independent of  $\lambda$ . Table 7 and Figure 6  
 711 summarise the results.  
 712  
 713

$\lambda$	IID (Algorithm 1)			Non-IID (Algorithm 2)		
	Cum. net	Regret	Team acc.	Cum. net	Regret	Team acc.
0.0	103.37±2.51	9.80±1.61	0.907±0.022	103.17±1.80	10.14±1.17	0.905±0.016
0.5	102.83±3.92	10.06±3.50	0.935±0.026	98.17±2.73	14.92±2.06	0.908±0.019
1.0	105.65±5.57	7.03±5.13	0.964±0.028	98.02±6.26	14.61±6.10	0.930±0.029
1.25	106.17±6.04	6.33±5.60	0.970±0.027	102.47±6.83	9.98±6.69	0.955±0.027
1.5	106.90±5.37	5.29±5.08	0.975±0.022	106.04±6.21	5.91±6.18	0.970±0.025
1.75	107.18±5.74	4.96±5.61	0.979±0.021	108.04±5.57	4.08±5.50	0.981±0.018
2.0	107.27±5.95	4.48±5.65	0.980±0.021	109.19±5.53	2.71±5.48	0.986±0.017
<b>3.0</b>	<b>108.84±2.22</b>	<b>2.29±2.11</b>	<b>0.988±0.010</b>	<b>110.61±1.03</b>	<b>0.59±0.93</b>	<b>0.993±0.007</b>
5.0	108.88±1.25	1.01±1.04	0.993±0.006	109.61±0.99	0.37±0.67	0.995±0.006
10.0	105.90±1.01	0.26±0.38	0.993±0.008	106.34±0.73	0.10±0.29	0.995±0.006
<i>Reference baselines (<math>\lambda</math>-independent)</i>						
No-OT	103.37±2.51	9.80±1.61	0.907±0.022	103.17±1.80	10.14±1.17	0.905±0.016
Random	83.98±5.92	28.72±6.03	0.917±0.020	79.78±5.52	32.97±4.87	0.905±0.024
UCB1	80.55±4.88	31.31±3.98	0.902±0.026	85.82±4.91	26.26±4.28	0.919±0.025

Table 7:  $\lambda$  grid search results for BOT-Orch. Mean  $\pm$  95% CI across 30 seeds,  $T=114$ .  $\lambda=0$  reproduces No-OT exactly (sanity check: difference = 0.000). Optimal value  $\lambda^*=3.0$  in bold. Reference baselines at the bottom are  $\lambda$ -independent.

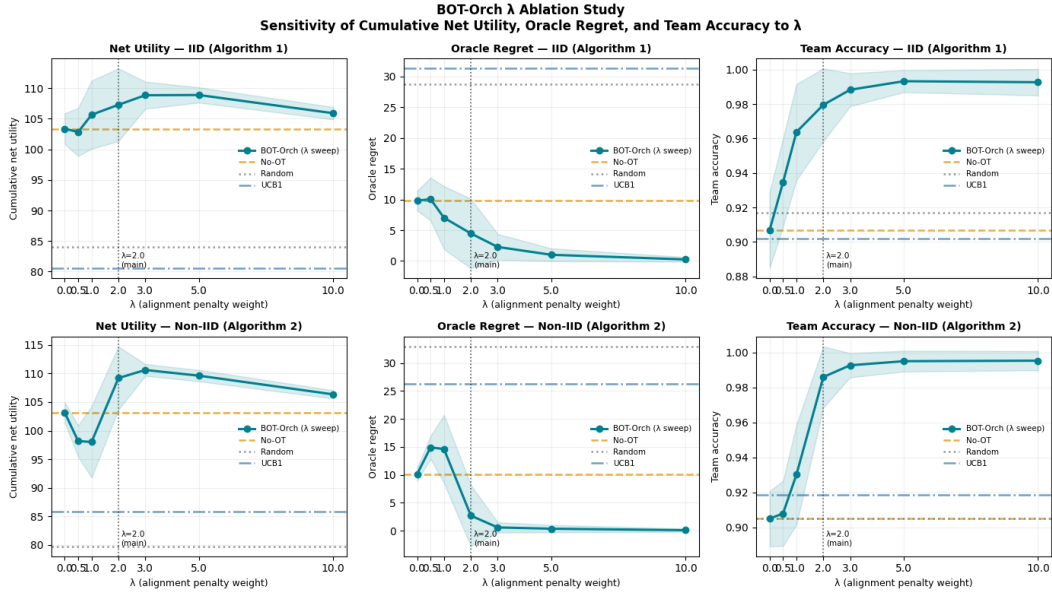


Figure 6:  $\lambda$  sensitivity curves. Cumulative net utility (left), oracle regret (centre), and team accuracy (right) as functions of  $\lambda$  for BOT-Orch (teal line with  $\pm 1$  std shading). Horizontal dashed lines show the No-OT, Random, and UCB1 baselines. The dotted vertical line marks  $\lambda=2.0$  (initial value);  $\lambda^*=3.0$  (grid-search optimum) is identified by the peak. Top row: IID (Algorithm 1); bottom row: Non-IID (Algorithm 2). In the Non-IID panels, performance at  $\lambda \in \{0.5, 1.0, 1.25\}$  falls below the No-OT baseline, revealing the phase transition.

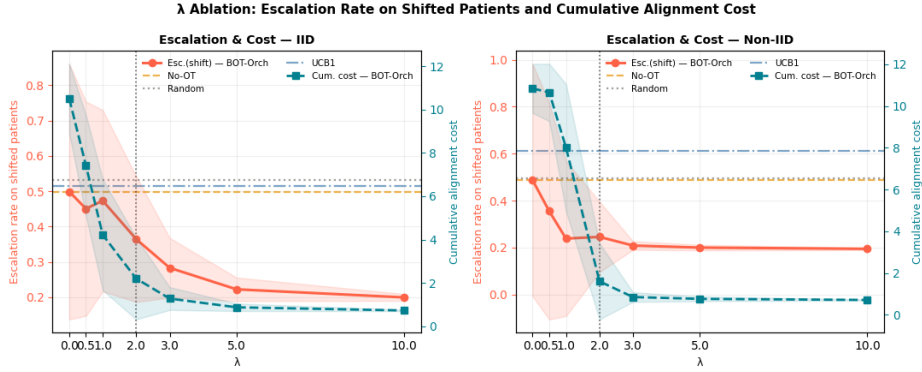


Figure 7: Escalation rate on shifted patients (left axis, red) and cumulative alignment cost (right axis, teal) as functions of  $\lambda$ . Horizontal dashed lines show No-OT, Random, and UCB1 escalation rates. Alignment cost decreases monotonically with  $\lambda$ , confirming Theorem 4.2 Part 1. *Left*: IID condition; *right*: Non-IID condition.

714 **Sensitivity and Optimal  $\lambda$ .** In the IID condition (Figure 6, top row), cumulative net utility in-  
 715 creases monotonically from  $\lambda=0$  through  $\lambda=5.0$  before declining at  $\lambda=10.0$ . The difference be-  
 716 tween  $\lambda=3.0$  ( $108.84 \pm 2.22$ ) and  $\lambda=5.0$  ( $108.88 \pm 1.25$ ) is 0.04 units — well within one standard  
 717 deviation — and the two values are statistically indistinguishable.

718 **Phase Transition in Non-IID Settings.** A qualitatively different pattern emerges in the Non-IID  
 719 condition (Figure 6, bottom row). For  $\lambda \in \{0.5, 1.0, 1.25\}$ , BOT-Orch performs *worse* than No-OT,  
 720 with net utility falling to as low as 98.02 at  $\lambda=1.0$  against the No-OT baseline of 103.17 (red-shaded  
 721 cells in Table 7). This counterintuitive dip arises because any positive  $\lambda$  penalises the human agent  
 722 during the first 57 in-distribution rounds (where  $W_t(\text{human})=0.120 > W_t(\text{AI})\approx 0$ ), suppressing its  
 723 EMA reward estimate through disuse. When the distribution shifts at round 58, a small  $\lambda$  is insuf-  
 724 ficient to override the depressed human reward history, and routing fails to redirect to the human.  
 725 Performance recovers sharply above  $\lambda\approx 1.5$ , where the OT signal becomes strong enough to domi-  
 726 nate stale reward estimates immediately when the shift arrives, and improves monotonically to the  
 727 peak at  $\lambda^*=3.0$ .

728 **Degeneracy at Large  $\lambda$ .** For  $\lambda \geq 5.0$ , the OT penalty dominates the reward signal entirely. The  
 729 policy degenerates toward a near-deterministic OT routing rule, effectively ignoring learned reward  
 730 history. Two symptoms confirm this regime: (i) the standard deviation of net utility collapses from  
 731  $\pm 5.53$  at  $\lambda=2.0$  to  $\pm 0.73$  at  $\lambda=10.0$  in Non-IID — not because the policy is more stable but because  
 732 it is no longer exploring; and (ii) net utility falls at  $\lambda=10.0$  (105.90 IID, 106.34 Non-IID) below the  
 733 optimum at  $\lambda^*=3.0$ .

734 **Identification of  $\lambda^*=3.0$ .** We identify  $\lambda^*=3.0$  as the joint optimum across both conditions on  
 735 the primary metric. It achieves the peak net utility in the Non-IID condition ( $110.61 \pm 1.03$ ), is  
 736 statistically tied with  $\lambda=5.0$  in IID (gap  $0.04 < \text{pooled std } 1.73$ ), reduces variance by  $3\times$  relative to  
 737  $\lambda=2.0$ , and is the last value at which bandit learning and OT alignment both contribute meaningfully.  
 738 The gain over the initial engineering value of  $\lambda=2.0$  is  $+1.42$  net utility in Non-IID, confirming that  
 739 the grid search yields a meaningful improvement rather than a marginal one.