

Probe-then-Commit Multi-Objective Bandits: Theoretical Benefits of Limited Multi-Arm Feedback

Ming Shi

Department of Electrical Engineering, University at Buffalo, Buffalo, NY

Abstract—We study an online resource-selection problem motivated by multi-radio access selection and mobile edge computing offloading. In each round, an agent chooses among K candidate links/servers (arms) whose performance is a stochastic d -dimensional vector (e.g., throughput, latency, energy, reliability). The key interaction is *probe-then-commit (PtC)*: the agent may probe up to $q > 1$ candidates via control-plane measurements to observe their vector outcomes, but must execute exactly one candidate in the data plane. This limited multi-arm feedback regime strictly interpolates between classical bandits ($q = 1$) and full-information experts ($q = K$), yet existing multi-objective learning theory largely focuses on these extremes. We develop PtC-P-UCB, an optimistic probe-then-commit algorithm whose technical core is frontier-aware probing under uncertainty in a Pareto mode, e.g., it selects the q probes by approximately maximizing a hypervolume-inspired frontier-coverage potential and commits by marginal hypervolume gain to directly expand the attained Pareto region. We prove a dominated-hypervolume frontier error of $\tilde{O}(K_P d / \sqrt{qT})$, where K_P is the Pareto-frontier size and T is the horizon, and scalarized regret $\tilde{O}(L_\phi d \sqrt{(K/q)T})$, where ϕ is the scalarizer. These quantify a transparent $1/\sqrt{q}$ acceleration from limited probing. We further extend to *multi-modal probing*: each probe returns M modalities (e.g., CSI, queue, compute telemetry), and uncertainty fusion yields variance-adaptive versions of the above bounds via an effective noise scale.

Index Terms—multi-objective bandits, probe-then-commit (PtC), limited multi-arm feedback, Pareto frontier, scalarization, hypervolume, multi-modal feedback, online resource selection

I. INTRODUCTION

Next-generation wireless and edge systems increasingly rely on online selection among multiple candidate network resources while meeting heterogeneous service requirements [1]. For example, in multi-radio access technology (multi-RAT) selection, a user equipment (UE) chooses among candidate links (e.g., 5G New Radio, WiFi, unmanned aerial vehicle relay) whose instantaneous channel quality, contention, and scheduling delay fluctuate across slots [2]. For another example, in mobile edge computing (MEC) offloading, a device selects an edge server whose queueing delay, compute load, and radio access conditions jointly determine end-to-end latency and success probability [3]. These decisions are inherently *multi-objective*. Optimizing a single scalar metric can yield operating points that violate service-level objectives (SLOs), or systematically sacrifice a “weak” KPI (key performance indicators) to gain another (e.g., sacrifice reliability to gain throughput).

Existing online-learning abstractions tend to focus on two extremes. Multi-objective bandits (MOB) observe only the vector outcome of the *single* executed arm per slot [4]–[7],

while full-information online optimization and learning (and vector-payoff approachability) observes outcomes of *all* K arms each round [8], [9]. However, wireless or edge systems often operate in a distinct intermediate regime enabled by the control plane. A UE can *probe* a small set of candidates using channel state information reference signals (CSI-RS) and beam sweeping, beacon frames or round-trip time (RTT) pings, or queue/CPU telemetry, but it can execute only one link/server due to data-plane constraints (one transmission/offload per slot). This yields *limited multi-arm feedback*, which is richer than standard MOB, but much cheaper than full information.

We formalize this interaction as a per-round *Probe-then-Commit* (PtC) protocol. At each slot t , the learner selects a probe set S_t with $|S_t| = q$, observes vector outcomes $\{\mathbf{r}_t(k)\}_{k \in S_t}$ (control-plane feedback), and then commits to one executed resource $k_t \in S_t$ whose outcome incurs realized system performance (data-plane execution). This model exposes an explicit *optimization-feedback tradeoff*: increasing q improves information and should accelerate learning, while probing consumes measurement and signaling budget. This PtC regime bridges MOB and full-information experts, but requires new algorithmic and analytical tools to handle vector objectives and frontier criteria under probe-limited feedback.

Wireless/edge systems often need to operate across multiple modes. For example, one needs to consider energy-saving against latency-critical. Thus, a learner should discover the Pareto frontier rather than only optimize one fixed operating point. Accordingly, we evaluate preference-free learning by a hypervolume-based frontier coverage metric, and we evaluate preference-based learning by scalarized regret for monotone concave utilities, e.g., fairness-sensitive aggregations.

PtC multi-objective learning couples probe design and execution in ways that do not arise in classical models. (i) *Uncertainty-aware frontier exploration*: probe-set design must lift multiple KPIs while diversifying across frontier regions; (ii) *Frontier accuracy under partial feedback*: coverage guarantees require controlling a set-valued error (hypervolume gap); (iii) *Quantifying the value of limited probing*: theory should expose how q sharpens rates beyond the bandit regime; (iv) *Multi-modal sensing*: probing often returns multiple modalities (CSI, queue length, CPU load) with heterogeneous noise, necessitating fusion that preserves valid confidence bounds [10].

Our main contributions are summarized as follows.

- **PtC multi-objective multi-feedback model (Sec. II)**. We introduce a stochastic multi-objective MAB under the PtC protocol with probe budget q . We formalize two complemen-

tary evaluation metrics, preference-free frontier learning via a dominated-hypervolume coverage gap and preference-based learning via scalarized regret for monotone concave utilities.

- **Algorithmic ideas: uncertainty-aware frontier coverage from probed samples (Sec. IV).** We develop PTC-P-UCB (Algorithm 2), which elects the q probes by approximately maximizing a hypervolume-inspired frontier-coverage potential and commits by marginal hypervolume gain to directly expand the attained Pareto region. The design is compatible with both frontier-coverage evaluation (hypervolume) and preference-based operation when a scalarizer is specified.
- **Theory: explicit value of limited multi-arm feedback (Sec. V).** We prove a dominated-hypervolume frontier coverage gap that vanishes at rate $\tilde{O}(K_P d / \sqrt{qT})$ (with frontier size K_P , d objectives, and T rounds), and a scalarized regret bound of order $\tilde{O}(d\sqrt{(K/q)T})$ for monotone L_ϕ -Lipschitz concave scalarizers, where \tilde{O} hides constants and logarithmic terms. These results quantify a clean $1/\sqrt{q}$ improvement from limited multi-arm probing, interpolating between the bandit limit ($q = 1$) and the full-information limit ($q = K$).
- **Multi-modal extension with variance-adaptive gains (Sec. VI).** We extend the framework to bundled multi-modal probing with M modalities. We develop MM-PTC-P-UCB (Algorithm 3) and show that fusion tightens confidence bounds through an effective noise scale, yielding variance-adaptive improvements for both frontier coverage and regret.
- **Empirical validation.** Simulations on multi-RAT- and MEC-inspired instances corroborate our theory. Modest probing budgets (e.g., $q \in \{2, 4\}$) significantly accelerate learning and improve Pareto coverage with moderate overhead, and multi-modal fusion provides an additional orthogonal gain.

A. Related Work

1) *Multi-objective bandits with Pareto criteria and scalarizations:* Multi-objective bandits study vector-valued rewards and compare actions via Pareto optimality or preference scalarizations. Upper confidence bound (UCB)-style approaches for vector rewards and Pareto efficiency appear in early MOB work (e.g., [4]), while preference-based learning uses monotone utilities, e.g., inequality-averse aggregations such as generalized Gini [6], to encode fairness-sensitive tradeoffs and enable regret analysis [5], [11]. More recent work develops Pareto-oriented regret notions that avoid fixing a scalarizer [7], [12]. Our focus is complementary. We study an intermediate *probe-limited multi-feedback* regime induced by wireless probing protocols, and provide guarantees for both frontier coverage (hypervolume) and preference-based regret under PtC feedback.

2) *Multiple-play, semi-bandits, and side-observation models:* Observing multiple actions per round relates to power-of-2-arms [13], multiple-play bandits [14], combinatorial/semi-bandit models [15], and learning with structured side observations such as feedback graphs [16]. These works are predominantly *single-objective* and typically assume additive reward/loss decompositions, whereas PtC enforces a single executed action with vector outcomes and Pareto criteria. This changes both the algorithmic goal (probe-set design for

frontier coverage) and the analysis (simultaneously controlling frontier estimation error and preference-based regret under probe-limited sampling).

3) *Vector-payoff online learning and approachability:* Full-information vector-payoff learning connects to Blackwell approachability [8] and its equivalence to no-regret learning [17]. This regime also includes standard online convex optimization with function information [9], [18]. Our PtC model in this paper can be viewed as a probe-limited, partial-information counterpart tailored to wireless/edge measurement pipelines.

4) *Hypervolume as a Pareto-set quality measure:* Dominated hypervolume is a standard quality measure for Pareto sets in multi-objective optimization and evolutionary computation [19], [20]. We adopt a hypervolume gap as a principled metric for frontier discovery over time under PtC feedback.

5) *Multi-modal sensing and fusion in wireless/edge decision making:* Next-generation wireless and edge platforms expose heterogeneous modalities correlated with service quality, including radio measurements, active probes, and system telemetry (e.g., queue/compute load). Recent wireless research also emphasizes multi-modal learning at the network level, including foundation-model perspectives for 6G systems [10]. To our knowledge, we take the first effort to model multi-modal probing as *multiple noisy views* of the same underlying multi-objective outcome vector, and design confidence-bound-driven learning rules whose uncertainty tightens via an effective variance under fusion. This yields variance-adaptive improvements in learning performance, while preserving the explicit $1/\sqrt{q}$ benefit of limited multi-arm probing under the PtC protocol.

Notation: For an integer n , $[n] = \{1, 2, \dots, n\}$. For a scalar x , $[x]^+ \triangleq \max\{0, x\}$. For $u, v \in \mathbb{R}^d$, $u \succeq v$ denotes component-wise inequality and $\|u\|_\infty = \max_j |u_j|$. We write $\Delta_d = \{w \in \mathbb{R}_+^d : \sum_{j=1}^d w_j = 1\}$ for the probability simplex.

II. PROBLEM FORMULATION

We study a multi-objective online resource-selection problem motivated by wireless access and edge computing systems, where each decision must balance heterogeneous KPIs under limited probing and measurement opportunities. The key interaction is *probe-then-commit* (PtC): in each round, the agent may probe multiple candidates via control-plane feedback, but ultimately executes only one due to data-plane constraints.

A. System Model and PtC Feedback Protocol

There are K candidate resources (arms), indexed by $k \in [K]$ (e.g., access links, edge servers), and d objectives indexed by $j \in [d]$ (e.g., throughput, latency, energy, reliability). Time is slotted with horizon T . At each round $t \in [T]$, each arm k is associated with a random vector-valued reward outcome

$$\mathbf{r}_t(k) = (r_t^{(1)}(k), \dots, r_t^{(d)}(k)) \in [0, 1]^d, \quad (1)$$

with *unknown* mean $\boldsymbol{\mu}(k) = \mathbb{E}[\mathbf{r}_t(k)] \in [0, 1]^d$. The d coordinates represent heterogeneous KPIs. If a KPI is naturally minimized (e.g., delay/energy), we convert it to a maximization objective by negating and normalizing it to $[0, 1]$.

Algorithm 1 Probe-then-Commit (PtC) interaction at round t **Require:** Probe budget $q \in [K]$.

- 1: **Probe selection:** choose a probe set $S_t \subseteq [K]$ with $|S_t| = q$.
- 2: **Measurement:** observe vector outcomes $\{\mathbf{r}_t(k)\}_{k \in S_t}$.
- 3: **Commit:** select one executed arm $k_t \in S_t$.
- 4: **Realization:** incur the round performance $\mathbf{r}_t(k_t)$ and proceed to next round $t+1$.

In many practical scenarios (e.g., multi-RAT selection and MEC offloading), a device can probe and measure multiple candidates (e.g., channel probing, active RTT pings, queue/CPU reports), but can use only one for actual transmission/offloading in that round. We model this by PtC (see Algorithm 1).

Probe outcomes for all $k \in S_t$ are observed (control-plane measurement), but only the executed arm k_t contributes to realized system performance in that round (data-plane execution). Probing may incur overhead (time/energy/signaling). We treat q as a fixed per-round budget in the main theory and discuss explicit probing-cost models in Sec. II-C.

B. Pareto Structure and Preference Scalarizers

The vector nature of $\boldsymbol{\mu}(k)$ induces a partial order, i.e., an arm may be better in one KPI but worse in another. We therefore formalize preference-free efficiency through Pareto dominance and, when the system operates under a fixed preference, use scalarizers to select a single operating point.

1) *Pareto dominance and frontier:* For $u, v \in \mathbb{R}^d$, we say that u dominates v (denoted $u \succ v$) if $u_j \geq v_j$ for all j and $u_{j'} > v_{j'}$ for some j' . Given mean vectors $\{\boldsymbol{\mu}(k)\}_{k=1}^K$, define the Pareto frontier $\mathcal{P}^* = \text{Pareto}(\{\boldsymbol{\mu}(k)\}_{k=1}^K)$, i.e., the set of mean vectors not dominated by any other. Moreover, we let $K_P \triangleq |\mathcal{P}^*|$ denote the size of the true frontier \mathcal{P}^* .

2) *Scalarizers with system preferences:* A deployed system may need a particular tradeoff based on operator policy or user preference. We model such preferences by scalarizers $\phi: \mathbb{R}^d \rightarrow \mathbb{R}$ that are *monotone* (improving any objective cannot decrease utility), L_ϕ -*Lipschitz* w.r.t. ℓ_∞ (i.e., $|\phi(u) - \phi(v)| \leq L_\phi \|u - v\|_\infty$), and *concave*. These conditions enable stable aggregation and regret analysis. Standard examples include: (i) weighted sum $\phi_w(u) = \sum_{j=1}^d w_j u_j$ with $w \in \Delta_d$; (ii) Chebyshev scalarizer $\phi_w^{\min}(u) = \min_{j \in [d]} w_j u_j$; and (iii) generalized Gini $\phi_\gamma(u) = \sum_{i=1}^d \gamma_i u_{(i)}$, where $u_{(1)} \leq \dots \leq u_{(d)}$ are the sorted components and $\gamma_1 \geq \dots \geq \gamma_d \geq 0$ [4], [6].

C. Performance Metrics and Probe Overhead

We evaluate algorithms using two complementary criteria. *Hypervolume coverage* quantifies preference-free frontier learning performance, while *scalarized regret* measures learning and operational loss under a specified system preference.

1) *Hypervolume-based Pareto coverage:* To quantify preference-free frontier learning, we use dominated hypervolume (HV) with respect to a fixed reference point $\mathbf{z}_{\text{ref}} \in \mathbb{R}^d$

that is component-wise worse than all attainable performance vectors. For compact set $\mathcal{S} \subseteq \mathbb{R}^d$, define the dominated region

$$\mathcal{D}(\mathcal{S}) \triangleq \{y \in \mathbb{R}^d : \exists u \in \mathcal{S}, \text{ s.t. } \mathbf{z}_{\text{ref}} \preceq y \preceq u\}, \quad (2)$$

and let $\mathcal{H}(\mathcal{S})$ be the Lebesgue measure of $\mathcal{D}(\mathcal{S})$. Intuitively, larger $\mathcal{H}(\mathcal{S})$ means that \mathcal{S} contains operating points that jointly perform well across objectives and spans a broader range of Pareto tradeoffs. Then, we let $\mathcal{Y}_T \triangleq \{\mathbf{r}_t(k_t)\}_{t=1}^T$ be the archive of executed outcome vectors and define the attained set

$$\mathcal{A}_T \triangleq \text{conv}(\mathcal{Y}_T), \quad (3)$$

where the convex hull captures time-sharing among operating points. Since the same time-sharing interpretation applies to the Pareto benchmark, define the convexified Pareto set $\mathcal{C}^* \triangleq \text{conv}(\mathcal{P}^*)$. We measure the remaining uncovered dominated volume by the attained-set hypervolume gap [19], [20]

$$\mathcal{L}_T^{\text{HV}} \triangleq [\mathcal{H}(\mathcal{C}^*) - \mathcal{H}(\mathcal{A}_T)]^+. \quad (4)$$

By construction, $\mathcal{L}_T^{\text{HV}} \geq 0$. Smaller $\mathcal{L}_T^{\text{HV}}$ indicates that the executed decisions achieve tradeoffs whose dominated region approaches that of the (time-shareable) Pareto benchmark.

2) *Scalarized regret:* Fix a scalarizer ϕ and define the best arm in hindsight $k^* \triangleq \arg \max_{k \in [K]} \phi(\boldsymbol{\mu}(k))$. Then, the scalarized regret is defined as follows,

$$R_T^\phi = \sum_{t=1}^T (\phi(\boldsymbol{\mu}(k^*)) - \phi(\mathbf{r}_t(k_t))). \quad (5)$$

3) *Probe overhead:* Probing consumes control-plane resources (time and energy). A simple model assigns a per-probe cost $\tau > 0$ and penalizes each round by $-\tau|S_t| = -\tau q$. Equivalently, one may impose a hard constraint $q \leq M_{\max}$ or a long-term budget. Our main results treat q as fixed and quantify how increased probing improves learning rates.

III. MOTIVATING EXAMPLES

We highlight two wireless/edge scenarios that directly match the PtC protocol in Algorithm 1: a device can obtain control-plane measurements from up to q candidates within a slot, but can execute only one candidate on the data plane.

1) *Multi-RAT link selection:* A UE probes up to q candidate access points or gNBs using lightweight control-plane signals (e.g., pilot measurements, beacons, short RTT probes, or queue indicators), obtaining a KPI vector such as throughput, delay, energy, and reliability. It then commits to one link for data transmission, so only the executed link's outcome is realized.

2) *MEC offloading:* A device queries up to q MEC servers for multi-modal telemetry (e.g., radio quality, queue status, CPU load), forming an outcome vector that captures end-to-end latency, energy consumption, and SLO satisfaction/reliability. It then offloads to a single server, again matching PtC.

IV. ALGORITHM DESIGN

This section presents PTC-P-UCB (see Algorithm 2), a probe-then-commit algorithm for multi-objective learning under PtC feedback. The algorithm maintains arm-wise confidence bounds from probed samples, selects a probe set to accelerate learning, and then commits to one probed arm for execution.

Algorithm 2 PTC-P-UCB: Probe-then-Commit Pareto-UCB

Require: Probe budget q , weights $w \in \Delta_d$, reference point z_{ref} , confidence parameter $\{\beta_t\}$. Commit mode: SCALAR (use given ϕ) or HV (use marginal hypervolume gain)

1: Initialize $N_1(k) \leftarrow 0$, $\hat{\mu}_1^{(j)}(k) \leftarrow 0$ for all $k \in [K]$ and $j \in [d]$, active set $\mathcal{K}_1 \leftarrow [K]$ and executed archive $\mathcal{Y}_0 \leftarrow \emptyset$

2: **for** $t = 1$ to T **do**

Confidence bounds (from probed samples):

3: **for** each $k \in \mathcal{K}_t$ and $j \in [d]$ **do**

4: Calculate $b_t^{(j)}(k)$ using (7)

5: Calculate $\text{UCB}_t^{(j)}(k)$ and $\text{LCB}_t^{(j)}(k)$ using (8)

6: **end for**

7: Define $u_t(k)$ and $\ell_t(k)$ for all $k \in \mathcal{K}_t$

Pruning:

8: $\tilde{\mathcal{K}}_t \leftarrow \{k \in \mathcal{K}_t : \nexists k' \in \mathcal{K}_t \text{ s.t. } \ell_t(k') \succ u_t(k)\}$

9: $\mathcal{K}_t \leftarrow \tilde{\mathcal{K}}_t$

Probe selection:

10: **if** HV **then** \triangleright HV

11: Select $S_t \subseteq \mathcal{K}_t$ by maximization of $F_t(S)$ in (10)

12: **else** \triangleright SCALAR

13: Select $S_t \subseteq \mathcal{K}_t$ by the surrogate top- q rule (11)

14: **end if**

15: Probe and observe $\{\mathbf{r}_t(k)\}_{k \in S_t}$

Commit:

16: **if** HV **then** \triangleright HV

17: $k_t \leftarrow \arg \max_{k \in S_t} \Delta_t^{\mathcal{H}}(k)$ using (12)

18: **else** \triangleright SCALAR

19: $k_t \leftarrow \arg \max_{k \in S_t} \phi(\mathbf{r}_t(k))$

20: **end if**

21: Execute k_t and incur $\mathbf{r}_t(k_t)$

22: **if** HV **then**

23: Update archive $\mathcal{Y}_t \leftarrow \mathcal{Y}_{t-1} \cup \{\mathbf{r}_t(k_t)\}$

24: **end if**

Updates (for all probed arms):

25: **for** each $k \in S_t$ **do**

26: $N_{t+1}(k) \leftarrow N_t(k) + 1$

27: **for** each $j \in [d]$ **do**

28: $\hat{\mu}_{t+1}^{(j)}(k) \leftarrow \hat{\mu}_t^{(j)}(k) + \frac{r_t^{(j)}(k) - \hat{\mu}_t^{(j)}(k)}{N_{t+1}(k)}$

29: **end for**

30: **end for**

31: For $k \notin S_t$, $N_{t+1}(k) \leftarrow N_t(k)$, $\hat{\mu}_{t+1}^{(j)}(k) \leftarrow \hat{\mu}_t^{(j)}(k)$.

32: Set $\mathcal{K}_{t+1} \leftarrow \mathcal{K}_t$.

33: **end for**

The PtC protocol requires executing exactly one arm per round. Our hypervolume coverage metric in (4) is defined on the executed attained set \mathcal{A}_T , and therefore depends on which arm is committed each round. To obtain preference-free frontier coverage for \mathcal{A}_T , we use a coverage-aware commit rule based on marginal hypervolume gain. However, if the system is operated under a known preference, the commit step can maximize a scalarizer ϕ to minimize scalarized regret.

Our design follows three main principles as follows.

1) *Probe selection should increase frontier coverage under*

uncertainty. Since probing provides the side information (q observations per round), the probe set should **include** arms that are plausibly Pareto-efficient under optimism and **diversify** across frontier regions. We achieve this by approximately maximizing a hypervolume-based coverage potential over optimistic vectors, with a score-based surrogate.

2) *Multi-objective optimism with component-wise confidence*. Maintain per-objective confidence intervals and form optimistic vectors to guide probing and safe elimination.

3) *Commit must match the metric*. For preference-free frontier learning under (4), we commit using marginal hypervolume gain of the executed archive. For preference-based operation we commit using ϕ to minimize scalarized regret.

A. Preliminaries for Confidence Bounds

Since the learner observes outcomes for all probed arms, we index learning progress by the number of times an arm has been probed, $N_t(k) \triangleq \sum_{\tau=1}^{t-1} \mathbb{I}\{k \in S_\tau\}$. For each objective $j \in [d]$, maintain the empirical mean based on probed samples,

$$\hat{\mu}_t^{(j)}(k) \triangleq \frac{1}{\max\{1, N_t(k)\}} \sum_{\tau=1}^{t-1} \mathbb{I}\{k \in S_\tau\} r_\tau^{(j)}(k). \quad (6)$$

Choose a confidence parameter $\beta_t = 2 \log(2Kdt^2/\delta)$ for Hoeffding-style bounds. Define the bonus term for each (k, j) ,

$$b_t^{(j)}(k) \triangleq \sqrt{\beta_t / \max\{1, N_t(k)\}}. \quad (7)$$

Then, we form clipped upper/lower bounds

$$\begin{aligned} \text{UCB}_t^{(j)}(k) &\triangleq \min\{1, \hat{\mu}_t^{(j)}(k) + b_t^{(j)}(k)\}, \\ \text{LCB}_t^{(j)}(k) &\triangleq \max\{0, \hat{\mu}_t^{(j)}(k) - b_t^{(j)}(k)\}. \end{aligned} \quad (8)$$

We define the optimistic and pessimistic vectors

$$\begin{aligned} u_t(k) &\triangleq (\text{UCB}_t^{(1)}(k), \dots, \text{UCB}_t^{(d)}(k)), \\ \ell_t(k) &\triangleq (\text{LCB}_t^{(1)}(k), \dots, \text{LCB}_t^{(d)}(k)). \end{aligned} \quad (9)$$

With high probability, $\ell_t(k) \preceq \mu(k) \preceq u_t(k)$ component-wise for all k, t , enabling safe pruning and optimistic probe selection.

B. Probe Selection via Frontier-Coverage Potential

A key decision is how to choose the probe set S_t of size q . Selecting the top- q arms by a single scalar score may overly concentrate probing around one region of the frontier. To encourage *diverse* coverage, we define a set-based potential.

1) *Coverage potential (set function)*: Let z_{ref} be the hypervolume reference point (as in Section II-C). Given optimistic vectors $\{u_t(k)\}_{k=1}^K$, define the potential of a probe set S as

$$F_t(S) \triangleq \mathcal{H}(\text{conv}\{u_t(k)\}_{k \in S}), \quad (10)$$

i.e., the dominated hypervolume of the convexified optimistic set. This potential rewards probe sets whose optimistic vectors jointly dominate a large region, which aligns with minimizing the Pareto coverage gap (up to optimism and estimation error).

2) *Greedy probe selection*: Maximizing $F_t(S)$ over all $|S| = q$ is combinatorial. We therefore use a greedy approximation. Starting from $S = \emptyset$, iteratively add the arm with the largest marginal gain in F_t . When $F_t(\cdot)$ is (approximately) monotone submodular, greedy achieves a constant-factor approximation.

3) *Fast surrogate (general scalarizer)*: When scalarized regret is the metric, we use a modular surrogate obtained by applying a preference scalarizer to the optimistic vector, i.e.,

$$\text{score}_t^\phi(k) \triangleq \phi(u_t(k)). \quad (11)$$

We then set S_t to be the q arms in \mathcal{K}_t with largest $\text{score}_t^\phi(k)$.

C. Commit Rule for Execution

After probing, the learner observes $\{\mathbf{r}_t(k)\}_{k \in S_t}$ and must execute one arm $k_t \in S_t$. We provide two commit rules that correspond to the two evaluation objectives.

1) *Coverage-based commit (for hypervolume coverage on \mathcal{A}_T)*: Let \mathcal{Y}_{t-1} denote the archive of executed outcomes up to round $t-1$ (so that $\mathcal{A}_{t-1} = \text{conv}(\mathcal{Y}_{t-1})$). For each candidate $k \in S_t$, define the marginal hypervolume gain

$$\Delta_t^{\mathcal{H}}(k) \triangleq \mathcal{H}(\text{conv}(\mathcal{Y}_{t-1} \cup \{\mathbf{r}_t(k)\})) - \mathcal{H}(\text{conv}(\mathcal{Y}_{t-1})). \quad (12)$$

We then commit to the arm that optimizes this gain, i.e.,

$$k_t^{\mathcal{H}} \in \arg \max_{k \in S_t} \Delta_t^{\mathcal{H}}(k). \quad (13)$$

Intuitively, this chooses the probed arm that most expands the dominated region of the attained set, directly targeting \mathcal{L}_T .

2) *Preference-based commit (for scalarized regret)*: Given a preference scalarizer ϕ , we commit to the best observed arm

$$k_t^\phi \in \arg \max_{k \in S_t} \phi(\mathbf{r}_t(k)). \quad (14)$$

This extracts immediate operational value from probing.

D. Frontier Pruning

To concentrate probing on plausible Pareto-frontier arms, we maintain an active set $\mathcal{K}_t \subseteq [K]$. An arm k can be safely discarded if it is certifiably dominated, i.e., $\exists k' \in \mathcal{K}_t$, s.t. $\ell_t(k') \succ u_t(k)$. On the high-probability event $\ell_t(\cdot) \preceq \mu(\cdot) \preceq u_t(\cdot)$, this rule never removes a true Pareto arm, but can dramatically reduce computation and improve the dependence on the frontier size K_P in the coverage analysis. Moreover, since $\mu(k') \succeq \mu(k)$ implies $\phi(\mu(k')) \geq \phi(\mu(k))$ for monotone scalarizer ϕ , the rule also never removes a ϕ -optimal arm.

E. Complexity and Probe Overhead

Per round, computing confidence bounds costs $O(|\mathcal{K}_t|d)$. Probe selection costs $O(q|\mathcal{K}_t|C_{\mathcal{H}})$ under greedy hypervolume or $O(K \log K)$ under the score-based surrogate, where $C_{\mathcal{H}}$ is the cost of a hypervolume marginal computation (small for $d \leq 4$). The probing overhead scales linearly with q in control-plane signaling. A per-probe cost τ can be incorporated by constraining $q \leq M_{\max}$ or subtracting τ from the scalar utility in deployment, without changing the learning machinery.

V. THEORETICAL RESULTS

This section provides performance guarantees for PTC-P-UCB (Algorithm 2) under the stochastic PtC model for the two evaluation metrics in Sec. II-C. We provide the proof sketches for main theorems, while complete proofs and supporting lemmas are provided in our technical report [21].

A. Assumption and Key Bookkeeping

Let \mathcal{F}_t be the filtration generated by past probe sets and observed probes up to the end of round t . A central identity is

$$\sum_{k=1}^K N_{T+1}(k) = qT, \quad (15)$$

i.e., each round yields q vector samples. Compared with traditional bandits, the learner observes q times more arm outcomes per round thanks to probes, which shrinks estimation error faster and translates into improved learning rates.

Assumption 1 (Conditionally sub-Gaussian noise). *For each $k \in [K]$ and $j \in [d]$, the noise $r_t^{(j)}(k) - \mu^{(j)}(k)$ is conditionally σ -sub-Gaussian given \mathcal{F}_{t-1} and independent over round t .*

B. Preference-Free Pareto Frontier Learning

We first analyze frontier learning under the HV mode of Algorithm 2, where the probe set is chosen to increase a hypervolume-based coverage potential and the commit step selects the probed arm with the largest marginal hypervolume gain. Bounding $\mathcal{L}_T^{\text{HV}}$ (defined in (4)) requires controlling two distinct effects. (i) *Learning error*: whether the algorithm probes enough to discover near-frontier arms (this is where q helps via (15)). (ii) *Execution sampling error*: the attained set uses instantaneous outcomes $\mathbf{r}_t(k_t)$ rather than mean vector $\mu(k_t)$, which introduces an additional statistical fluctuation.

Theorem 1 (Attained-set hypervolume gap of PTC-P-UCB (HV mode)). *Under Assumption 1, run Algorithm 2 and the marginal-gain commit rule (12)–(13). Then, we have*

$$\mathbb{E}[\mathcal{L}_T^{\text{HV}}] = \tilde{O}\left(K_P d / \sqrt{qT} + d / \sqrt{T}\right). \quad (16)$$

The first term in (16) is the *frontier-learning* term. First, probing indeed accelerates estimation of Pareto-relevant arms. This yields a $1/\sqrt{q}$ improvement in the performance. Second, the factor K_P captures frontier complexity. The dependence on K_P suggests that the dominant learning burden is to resolve and cover the K_P Pareto-relevant mean vectors, while dominated arms do not directly affect hypervolume. The second term in (16) reflects the fact that the attained set is formed from instantaneous executed outcomes rather than means. It vanishes as T grows and is unavoidable for an execution-based metric.

Proof sketch. We define the denoised archive $\tilde{\mathcal{Y}}_T \triangleq \{\hat{\mu}_T(k_t)\}_{t=1}^T$, with $\tilde{\mathcal{A}}_T \triangleq \text{conv}(\tilde{\mathcal{Y}}_T)$. We upper bound the hypervolume gap by inserting $\text{conv}(\{\mu(k_t)\}_{t=1}^T)$, i.e.,

$$\begin{aligned} \mathcal{L}_T^{\text{HV}} &\leq \underbrace{[\mathcal{H}(\mathcal{C}^*) - \mathcal{H}(\text{conv}(\{\mu(k_t)\}_{t=1}^T))]}_{\text{learning and coverage error}} + \underbrace{[\mathcal{H}(\text{conv}(\{\mu(k_t)\}_{t=1}^T)) - \mathcal{H}(\tilde{\mathcal{A}}_T)]}_{\text{estimation error}}. \end{aligned} \quad (17)$$

Step 1 (learning and coverage error). On the event that all coordinate-wise confidence intervals (CIs) hold, the probe selection uses optimistic vectors $u_t(k)$, so any Pareto-relevant arm with large uncertainty induces a large marginal gain in the optimistic coverage potential. A standard potential argument then shows that each Pareto arm is probed enough that its CI shrinks to $\tilde{O}(1/\sqrt{qT})$. Combining this with hypervolume stability over K_P frontier points yields $\mathcal{H}(C^*) - \mathcal{H}(\text{conv}(\{\mu(k_t)\})) = \tilde{O}(K_P d/\sqrt{qT})$ in expectation.

Step 2 (estimation error). A hypervolume stability lemma for sets in $[0, 1]^d$ gives $|\mathcal{H}(\text{conv}(U)) - \mathcal{H}(\text{conv}(V))| \leq L_{\mathcal{H}} \cdot d_H(\text{conv}(U), \text{conv}(V))$, and $d_H(\text{conv}(U), \text{conv}(V)) \leq \max_{t \leq T} \|u_t - v_t\|_{\infty}$. Taking $u_t = \mu(k_t)$ and $v_t = \hat{\mu}_T(k_t)$ yields $|\mathcal{H}(\text{conv}(\{\mu(k_t)\})) - \mathcal{H}(\hat{A}_T)| \leq L_{\mathcal{H}} \cdot \max_{t \leq T} \|\hat{\mu}_T(k_t) - \mu(k_t)\|_{\infty}$. Uniform concentration under adaptive probing implies $\max_{k \in [K]} \|\hat{\mu}_T(k) - \mu(k)\|_{\infty} = \tilde{O}(1/\sqrt{qT})$ in expectation (using $\sum_k N_{T+1}(k) = qT$), hence the estimation term is $\tilde{O}(d/\sqrt{qT})$.

Combining Steps 1–2 gives the claimed rate. \square

C. Fixed-Confidence ϵ -Frontier Identification

We next translate confidence bounds into a fixed-confidence sample complexity guarantee for identifying an ϵ -accurate frontier approximation. An arm k is ϵ -Pareto optimal (in ℓ_{∞}) if there is no k' such that $\mu(k') \succeq \mu(k) + \epsilon \mathbf{1}$. Using the coordinate-wise confidence bounds, define the output set

$$\hat{\mathcal{P}}_T^{(\epsilon)} \triangleq \left\{ k \in [K] : \right. \\ \left. \exists k' \in [K], \text{ s.t. } \text{LCB}_T(k') \succeq \text{UCB}_T(k) + \epsilon \mathbf{1}/2 \right\}. \quad (18)$$

This rule is conservative, i.e., on the event that all confidence intervals are valid, it produces no ϵ -dominated false positives.

Theorem 2 (Sample complexity for ϵ -frontier identification). *Fix $\epsilon, \delta \in (0, 1)$. Under Assumption 1, on the event that all coordinate-wise confidence intervals hold, $\hat{\mathcal{P}}_T^{(\epsilon)}$ contains all truly Pareto-optimal arms and contains no arm that is ϵ -dominated. With probability at least $1 - \delta$, it suffices that*

$$qT \geq C \cdot K_P d \log(Kd/\delta) / \epsilon^2, \quad (19)$$

for a universal constant $C > 0$. Equivalently, the number of probed samples required is $N_{\epsilon} = \tilde{O}\left(\frac{K_P d}{\epsilon^2}\right)$ or $T = \tilde{O}\left(\frac{K_P d}{q \epsilon^2}\right)$.

For fixed confidence (ϵ, δ) , increasing q reduces the required horizon linearly, since it increases the number of observed arm vectors per round. The dependence on K_P formalizes that only Pareto-relevant arms must be resolved to ϵ accuracy.

D. Scalarized Regret in SCALAR Mode

We finally turn to preference-based operation. In SCALAR mode, the learner probes using the optimistic scalar index score $\phi(k) = \phi(u_t(k))$ and commits via (14).

Theorem 3 (Scalarized regret of PTC-P-UCB (SCALAR mode)). *Under Assumptions 1, run Algorithm 2. Then,*

$$\mathbb{E}[R_T^{\phi}] = \tilde{O}\left(L_{\phi} d \sqrt{KT/q}\right). \quad (20)$$

Moreover, with probability at least $1 - \delta$, $\delta \in (0, 1)$, the same rate holds up to $\text{polylog}(K, d, T, 1/\delta)$ factors.

Relative to the standard bandit rate $\tilde{O}(\sqrt{KT})$, probing yields an effective sample-size increase by $\sqrt{1/q}$ because q arms are observed per round. The additional d factor arises from uniform control over coordinates and Lipschitz stability of ϕ .

Remark 1 (Boundary cases). When $q = 1$, (20) recovers the standard scalarized multi-objective bandit rate up to logs. When $q = K$, all arms are observed each round and the rate becomes $\tilde{O}(L_{\phi} d \sqrt{T})$, matching full-information scaling.

Proof sketch. Let $\Delta(k) \triangleq \phi(\mu(k^*)) - \phi(\mu(k))$ denote the gap. Let \mathcal{E} be the high-probability event on which all coordinate-wise CIs hold uniformly. By a union bound and sub-Gaussian concentration, $\mathbb{P}(\mathcal{E}) \geq 1 - \delta$ for an appropriate β_t .

Step 1 (optimism implies that any selected arm must still be “uncertain enough”). On \mathcal{E} , monotonicity of ϕ gives $\phi(\mu(k)) \leq \phi(u_t(k))$ for all k . Define the optimistic score $U_t(k) \triangleq \phi(u_t(k))$ used for probe selection. Since S_t contains the top- q arms by $U_t(\cdot)$, whenever an arm k is probed, it must satisfy $U_t(k) \geq U_t(k^*) \geq \phi(\mu(k^*))$ on \mathcal{E} . Therefore, $\Delta(k) = \phi(\mu(k^*)) - \phi(\mu(k)) \leq \phi(u_t(k)) - \phi(\mu(k)) \leq L_{\phi} \|u_t(k) - \mu(k)\|_{\infty} \leq L_{\phi} \max_j b_t^{(j)}(k)$, where the last inequality uses $\mu^{(j)}(k) \leq u_t^{(j)}(k) \leq \mu^{(j)}(k) + b_t^{(j)}(k)$ on \mathcal{E} . Hence, if $\Delta(k) > \epsilon$, then arm k can only remain in the top- q probe set while $\max_j b_t^{(j)}(k) \gtrsim \epsilon/L_{\phi}$, which implies $N_t(k) \lesssim \beta_T (L_{\phi}/\epsilon)^2$ (up to constants and d).

Step 2 (“parallel exploration” converts probe complexity into a $1/q$ reduction in time). Let $\mathcal{K}_{\epsilon} = \{k : \Delta(k) > \epsilon\}$. From Step 1, each $k \in \mathcal{K}_{\epsilon}$ needs at most $O(\beta_T (L_{\phi}/\epsilon)^2)$ probes before its optimistic score drops below $\phi(\mu(k^*))$ on \mathcal{E} , after which it cannot enter the top- q set again. Since each round allocates q probes, the total number of rounds in which any arm in \mathcal{K}_{ϵ} can still be probed is at most $O\left(\frac{|\mathcal{K}_{\epsilon}|}{q} \cdot \beta_T \frac{L_{\phi}^2}{\epsilon^2}\right)$.

Step 3 (gap-free regret via a ϵ -decomposition). Decompose the regret as $\sum_{t=1}^T \Delta(k_t) \leq T\epsilon + \sum_{t: \Delta(k_t) > \epsilon} \Delta(k_t)$. The first term is $T\epsilon$. For the second term, upper bound each $\Delta(k_t)$ by 1 and use Step 2 to bound the number of rounds where $\Delta(k_t) > \epsilon$, yielding $\sum_{t: \Delta(k_t) > \epsilon} \Delta(k_t) \leq O\left(\frac{K}{q} \cdot \beta_T \frac{L_{\phi}^2}{\epsilon^2}\right)$. Choosing $\epsilon \asymp L_{\phi} \sqrt{\frac{K \beta_T}{qT}}$ gives $\tilde{O}\left(L_{\phi} \sqrt{\frac{KT}{q}}\right)$. Applying a union bound over d coordinates in \mathcal{E} introduces the stated d factor.

Step 4 (from high-probability to expectation and realized regret). On \mathcal{E}^c we use the trivial bound $R_T^{\phi} \leq T$. Thus, $\mathbb{E}[R_T^{\phi}] \leq \tilde{O}(L_{\phi} d \sqrt{KT/q}) + T\delta$, and taking $\delta = 1/T$ yields the stated expectation bound. For the realized reward scalarizer $\phi(r_t(k_t))$, an additional term controlling the selection-dependent deviation $\sum_t (\phi(\mu(k_t)) - \phi(r_t(k_t)))$ is handled via sub-Gaussian maximal inequalities and Lipschitzness of ϕ , and does not change the leading $\tilde{O}(L_{\phi} d \sqrt{KT/q})$ order. \square

VI. EXTENSION: MULTI-MODAL FEEDBACK

In many wireless/edge systems, probing a candidate resource returns *multiple modalities* of side information, e.g., CSI

Algorithm 3 MM-PTC-P-UCB: a multi-modal extension

Require: Probe budget q , fusion weights $\alpha \in \Delta_M$, confidence parameter $\{\beta_t\}$.

- 1: Run PTC-P-UCB (Algorithm 2) but:
- 2: (i) when $k \in S_t$, observe $\{\mathbf{z}_t^{(m)}(k)\}_{m=1}^M$ and form $\tilde{\mathbf{r}}_t(k) = \sum_m \alpha_m \mathbf{z}_t^{(m)}(k)$;
- 3: (ii) update $\hat{\mu}_t(k)$ using $\tilde{\mathbf{r}}_t(k)$ instead of $\mathbf{r}_t(k)$;
- 4: (iii) use the effective-scale-based confidence radii $b_t^{(j)}(k) = \sigma_{\text{eff}}^{(j)}(k) \sqrt{\beta_t / \max\{1, N_t(k)\}}$.

measurements, queue-length reports, and CPU-load telemetry. When fused properly, these modalities can improve learning.

A. Multi-Modal Observation Model

We consider M modalities indexed by m . For each round t and arm k , there is a vector-valued outcome $\mathbf{r}_t(k) \in [0, 1]^d$ with mean $\boldsymbol{\mu}(k)$. Under multi-modal probing, whenever $k \in S_t$ the learner observes *all* modality readings $\{\mathbf{z}_t^{(m)}(k)\}_{m=1}^M$, where

$$\mathbf{z}_t^{(m)}(k) = \mathbf{r}_t(k) + \boldsymbol{\eta}_t^{(m)}(k). \quad (21)$$

We assume $\boldsymbol{\eta}_t^{(m)}(k)$ is conditionally mean-zero given \mathcal{F}_{t-1} . For clarity, we state a diagonal (objective-wise) sub-Gaussian version: for each objective j , $\eta_t^{(m,j)}(k)$ is conditionally $\sigma_m^{(j)}(k)$ -sub-Gaussian given \mathcal{F}_{t-1} and independent over t for each fixed (k, m, j) . Then, given fusion weights $\alpha = (\alpha_1, \dots, \alpha_M) \in \Delta_M$, we define the fused observation as follows,

$$\tilde{\mathbf{r}}_t(k) \triangleq \sum_{m=1}^M \alpha_m \mathbf{z}_t^{(m)}(k). \quad (22)$$

$\tilde{\mathbf{r}}_t(k)$ is an unbiased noisy observation of $\mathbf{r}_t(k)$, and $\tilde{r}_t^{(j)}(k) - r_t^{(j)}(k)$ is conditionally sub-Gaussian with effective scale $(\sigma_{\text{eff}}^{(j)}(k))^2 = \sum_{m=1}^M \alpha_m^2 (\sigma_m^{(j)}(k))^2$. Intuitively, multi-modality yields an *orthogonal* acceleration mechanism. Besides the m -fold sample increase from probing, fusion can reduce per-sample uncertainty through σ_{eff} .

B. Algorithm: MM-PTC-P-UCB

The multi-modal extension (Algorithm 3) involves a *local* change to the learning pipeline in Algorithm 2. We replace the single probed sample $\mathbf{r}_t(k)$ by the fused sample $\tilde{\mathbf{r}}_t(k)$ in the mean updates, and replace the base noise scale by σ_{eff} in the confidence radii. Probe selection and the commit rule follow exactly the modes as Sec. IV. In particular, in the multi-modal model, the learner only observes modality feedback, not directly observing $\mathbf{r}_t(k)$. Thus, both in HV mode (commit via marginal hypervolume gain) and in SCALAR mode (commit via $\phi(\cdot)$), it is needed to compute the commit decision using the best available estimate of the probed outcome, i.e., $\tilde{\mathbf{r}}_t(k)$.

C. Guarantees: Variance-Adaptive Improvement

With fixed fusion weights α , the analysis in Sec. V carries over by replacing the base noise scale by σ_{eff} in the concentration arguments. Intuitively, PtC provides q samples per round, while fusion reduces the noise per sample. Recall that in HV mode the learner maintains the executed archive

$\mathcal{Y}_t = \{\tilde{\mathbf{r}}_s(k_s)\}_{s=1}^t$ and the attained set $\mathcal{A}_t = \text{conv}(\mathcal{Y}_t)$, where $\tilde{\mathbf{r}}_t(k)$ is considered here and is the fused observation in (22).

Theorem 4 (Variance-adaptive attained-set hypervolume gap under fixed fusion). *Consider the multi-modal model (21)–(22) with fixed fusion weights $\alpha \in \Delta_M$. Let $\sigma_{\text{eff}} \triangleq \max_{k \in [K], j \in [d]} \sigma_{\text{eff}}^{(j)}(k)$. Run PTC-P-UCB in HV mode, using the coverage-based commit rule (13). Then, we have*

$$\mathbb{E}[\mathcal{L}_T^{\text{HV}}] = \tilde{O}\left(K_P d \sigma_{\text{eff}} / \sqrt{qT} + d / \sqrt{T}\right). \quad (23)$$

Eq. (23) makes the two acceleration mechanisms explicit. First, the PtC probe budget contributes the same $1/\sqrt{q}$ improvement via the identity $\sum_k N_{T+1}(k) = qT$. Second, multi-modal fusion improves the *per-sample* statistical accuracy by shrinking the effective noise scale from σ to σ_{eff} .

Proof sketch. The proof has three steps. (1) Under the sub-Gaussian assumption and the fused estimator (22), uniform coordinate-wise concentration yields, for all k and j , $|\hat{\mu}_T^{(j)}(k) - \mu^{(j)}(k)| \lesssim \sigma_{\text{eff}}^{(j)}(k) \sqrt{\log(\cdot) / N_{T+1}(k)}$. (2) Using Cauchy–Schwarz together with $\sum_k N_{T+1}(k) = qT$ bounds the aggregate estimation error on Pareto-relevant arms by $\tilde{O}(\sigma_{\text{eff}} / \sqrt{qT})$. (3) A hypervolume stability lemma upper-bounds the perturbation of $\mathcal{H}(\cdot)$ over sets in $[0, 1]^d$ by $O(K_P d)$ times the ℓ_∞ estimation error, which yields (23). \square

Theorem 5 (Variance-adaptive scalarized regret under fixed fusion). *Consider the bundled multi-modal model (21)–(22) with fixed $\alpha \in \Delta_M$. Let $\sigma_{\text{eff}} \triangleq \max_{k \in [K], j \in [d]} \sigma_{\text{eff}}^{(j)}(k)$. Run PTC-P-UCB with confidence radii scaled by σ_{eff} yields*

$$\mathbb{E}[\mathcal{R}_T^\phi] = \tilde{O}\left(L_\phi d \sigma_{\text{eff}} \sqrt{KT/q}\right). \quad (24)$$

VII. NUMERICAL RESULTS

We evaluate PTC-P-UCB and MM-PTC-P-UCB under both metrics from Sec. II-C. We use synthetic instances motivated by wireless/edge tradeoffs and include near-dominated “confusers” that make dominance relations statistically fragile.

We simulate $K = 24$ arms and $d = 4$ objectives over horizon T . To induce a nontrivial frontier while retaining hard-to-separate alternatives, we generate means by a mixture construction, i.e., first sample a subset of “frontier” arms from clustered points on a tradeoff surface, and then generate remaining arms as dominated perturbations around these clusters. This yields realistic regimes where small estimation errors can flip pairwise dominance, slowing frontier learning.

We run PtC with probe budgets $q \in \{1, 2, 4, K\}$. For multi-modal experiments ($M = 3$), each probe returns a set $\mathbf{z}_t^{(m)}(k) = \mathbf{r}_t(k) + \boldsymbol{\eta}_t^{(m)}(k)$ with heterogeneous scales $(\sigma_m)_{m=1}^M = (0.08, 0.12, 0.20)$. Moreover, we fuse modalities using inverse-variance weights $\alpha_p \propto 1/(\bar{\sigma}_p^2)$.

A. Main Results and Findings

1) *Frontier discovery improves with limited probing:* Fig. 1 reports the frontier hypervolume gap $\mathcal{G}_T^{\text{HV}}$. Increasing the probe budget yields consistently faster decay. Moving from $q = 1$ to $q = 4$ substantially reduces the time needed to reach the same

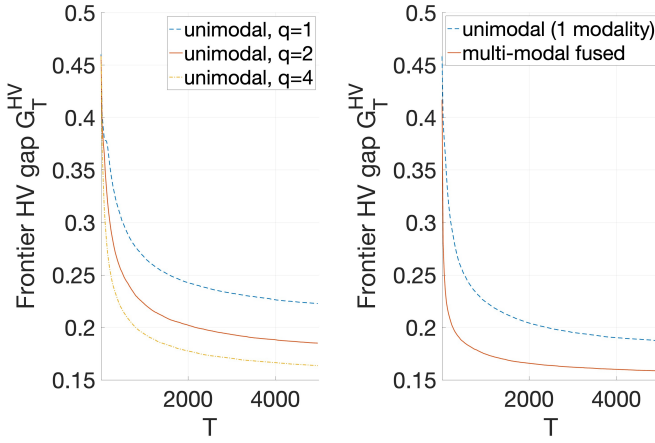


Fig. 1. Frontier hypervolume gap $\mathcal{G}_T^{\text{HV}}$ versus T : effect of q and benefit of multi-modal fusion (set $q = 2$).

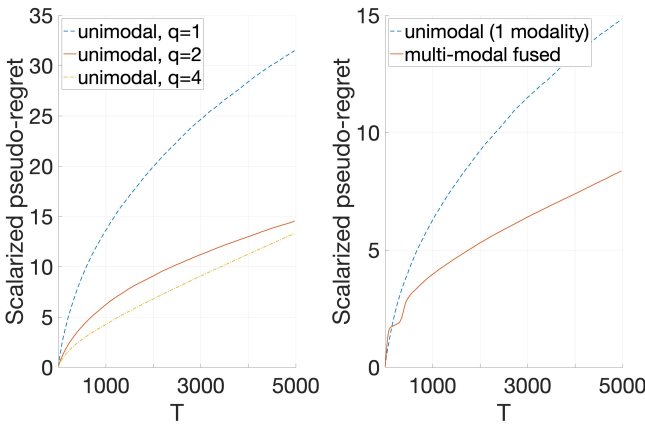


Fig. 2. Worst-case scalarized regret versus T : effect of q and benefit of multi-modal fusion (set $q = 2$).

coverage level. This matches the predicted $1/\sqrt{q}$ acceleration in the hypervolume guarantee (Theorem 1) and highlights that multi-feedback probing accelerates frontier learning, not merely exploitation under a fixed preference.

2) *Scalarized performance accelerates at the same $1/\sqrt{q}$ rate:* Fig. 2 shows worst-case scalarized regret. Across environments, the ordering $q = 4 < q = 2 < q = 1$ persists throughout the horizon. Moreover, the separations between curves are consistent with the theoretical scaling $\tilde{O}(\sqrt{K}/(qT))$ from Theorem 3. This explicitly confirms that the q side observations translate into an effective sample-size gain.

3) *Multi-modal fusion yields an orthogonal variance-reduction gain:* With $M = 3$ modalities, inverse-variance fusion reduces σ_{eff} , tightening confidence bounds and improving both metrics at fixed q . Fig. 1 shows that fused feedback reaches the same hypervolume gap earlier than unimodal sensing, consistent with the variance-adaptive analysis in Sec. VI.

VIII. CONCLUSION

We introduced a *multi-objective multi-feedback* MAB capturing realistic probing in wireless/edge systems, designed PtC-

P-UCB, and proved regret and Pareto coverage bounds with explicit q -dependence. Our multi-modal extension is variance-adaptive and practically effective. Future work includes contextual/linear structure, delayed feedback, and nonstationarity.

REFERENCES

- [1] S. Lin, M. Shi, A. Arora, R. Bassily, E. Bertino, C. Caramanis, K. Chowdhury, E. Ekici, A. Eryilmaz, S. Ioannidis *et al.*, “Leveraging synergies between ai and networking to build next generation edge networks,” in *2022 IEEE 8th International Conference on Collaboration and Internet Computing (CIC)*. IEEE, 2022, pp. 16–25.
- [2] C. Li, J. Jiang, and J. Li, “Multi-radio access technology (multi-rat) diversity for ultra-reliable low-latency communication (urllc),” Nov. 10 2020, uS Patent 10,833,902.
- [3] P. Mach and Z. Becvar, “Mobile edge computing: A survey on architecture and computation offloading,” *IEEE communications surveys & tutorials*, vol. 19, no. 3, pp. 1628–1656, 2017.
- [4] M. M. Dragan and A. Nowe, “Designing multi-objective multi-armed bandits algorithms: A study,” in *The 2013 international joint conference on neural networks (IJCNN)*. IEEE, 2013, pp. 1–8.
- [5] L. Cao, M. Shi, and N. B. Shroff, “Provably efficient multi-objective bandit algorithms under preference-centric customization,” *arXiv preprint arXiv:2502.13457*, 2025.
- [6] J. A. Weymark, “Generalized gini inequality indices,” *Mathematical social sciences*, vol. 1, no. 4, pp. 409–430, 1981.
- [7] D. M. Roijers, P. Vamplew, S. Whiteson, and R. Dazeley, “A survey of multi-objective sequential decision-making,” *Journal of Artificial Intelligence Research*, vol. 48, pp. 67–113, 2013.
- [8] D. Blackwell, “An analog of the minimax theorem for vector payoffs,” 1956.
- [9] E. Hazan *et al.*, “Introduction to online convex optimization,” *Foundations and Trends® in Optimization*, vol. 2, no. 3–4, pp. 157–325, 2016.
- [10] J. Du, T. Lin, C. Jiang, Q. Yang, C. F. Bader, and Z. Han, “Distributed foundation models for multi-modal learning in 6g wireless networks,” *IEEE Wireless Communications*, vol. 31, no. 3, pp. 20–30, 2024.
- [11] M. Agarwal, V. Aggarwal, and T. Lan, “Multi-objective reinforcement learning with non-linear scalarization,” in *Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems*, 2022, pp. 9–17.
- [12] A. Garivier, W. M. Koolen *et al.*, “Sequential learning of the pareto front for multi-objective bandits,” in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2024, pp. 3583–3591.
- [13] M. Shi, X. Lin, and L. Jiao, “Power-of-2-arms for adversarial bandit learning with switching costs,” *IEEE Transactions on Networking*, 2025.
- [14] V. Anantharam, P. Varaiya, and J. Walrand, “Asymptotically efficient allocation rules for the multiarmed bandit problem with multiple plays-part i: Iid rewards,” *IEEE Transactions on Automatic Control*, vol. 32, no. 11, pp. 968–976, 2003.
- [15] B. Kveton, Z. Wen, A. Ashkan, and C. Szepesvari, “Tight regret bounds for stochastic combinatorial semi-bandits,” in *Artificial Intelligence and Statistics*. PMLR, 2015, pp. 535–543.
- [16] N. Alon, N. Cesa-Bianchi, O. Dekel, and T. Koren, “Online learning with feedback graphs: Beyond bandits,” in *Conference on Learning Theory*. PMLR, 2015, pp. 23–35.
- [17] J. Abernethy, P. L. Bartlett, and E. Hazan, “Blackwell approachability and no-regret learning are equivalent,” in *Proceedings of the 24th Annual Conference on Learning Theory*. JMLR Workshop and Conference Proceedings, 2011, pp. 27–46.
- [18] M. Shi, X. Lin, and S. Fahmy, “Competitive online convex optimization with switching costs and ramp constraints,” *IEEE/ACM Transactions on Networking*, vol. 29, no. 2, pp. 876–889, 2021.
- [19] E. Zitzler and L. Thiele, “Multiobjective evolutionary algorithms: a comparative case study and the strength pareto approach,” *IEEE transactions on Evolutionary Computation*, vol. 3, no. 4, pp. 257–271, 2002.
- [20] A. Auger, J. Bader, D. Brockhoff, and E. Zitzler, “Theory of the hypervolume indicator: optimal μ -distributions and the choice of the reference point,” in *Proceedings of the tenth ACM SIGEVO workshop on Foundations of genetic algorithms*, 2009, pp. 87–102.
- [21] , “Probe-then-commit multi-objective bandits: Theoretical benefits of limited multi-arm feedback,” https://mingshihomepage.com/papers/PtC-Multi-Arm-Feedback_WiOPT2026.pdf, 2026.

APPENDIX A
PROOF OF THEOREM 1

Proof. For readability, write $d \triangleq d$ and $q \triangleq q$. Recall $\mathcal{Y}_T = \{\mathbf{r}_t(k_t)\}_{t=1}^T$, $\mathcal{A}_T = \text{conv}(\mathcal{Y}_T)$, $\mathcal{P}^* = \text{Pareto}(\{\mu(k)\}_{k=1}^K)$, and $\mathcal{C}^* = \text{conv}(\mathcal{P}^*)$. Let $z_{\text{ref}} \prec \mathbf{0}$ be the fixed reference point used for $\mathcal{H}(\cdot)$.

A. Step 0: A Convenient Geometric Lipschitz Bound for Dominated Hypervolume

For any compact $\mathcal{S} \subseteq [0, 1]^d$ define the ℓ_∞ -expansion $\mathcal{S}^{+\varepsilon} \triangleq \{x \in \mathbb{R}^d : \exists s \in \mathcal{S} \text{ s.t. } \|x - s\|_\infty \leq \varepsilon\}$. The dominated region of \mathcal{S} w.r.t. z_{ref} is $\mathcal{D}(\mathcal{S}) \triangleq \{y : \exists u \in \mathcal{S}, z_{\text{ref}} \preceq y \preceq u\}$, hence $\mathcal{H}(\mathcal{S}) = \text{Leb}(\mathcal{D}(\mathcal{S}))$. Let $R \triangleq \max_{j \in [d]} (1 - z_{\text{ref}}^{(j)})$.

Lemma 1 (One-sided HV stability under ℓ_∞ expansion). *For any compact $\mathcal{S} \subseteq [0, 1]^d$ and any $\varepsilon \in [0, 1]$,*

$$\mathcal{H}(\mathcal{S}^{+\varepsilon}) \leq \mathcal{H}(\mathcal{S}) + d R^{d-1} \varepsilon.$$

Consequently, if $\mathcal{A} \subseteq \mathcal{B}^{+\varepsilon}$ then

$$[\mathcal{H}(\mathcal{A}) - \mathcal{H}(\mathcal{B})]^+ \leq d R^{d-1} \varepsilon.$$

Proof. Since $\mathcal{S}^{+\varepsilon} \subseteq [-\varepsilon, 1 + \varepsilon]^d$ and $z_{\text{ref}} \prec \mathbf{0}$, the set difference $\mathcal{D}(\mathcal{S}^{+\varepsilon}) \setminus \mathcal{D}(\mathcal{S})$ is contained in the union (over coordinates j) of “ ε -thick slabs” adjacent to the boundary of $\mathcal{D}(\mathcal{S})$ in direction j . Each such slab has thickness at most ε and $(d-1)$ -dimensional cross-section at most R^{d-1} , hence total added volume is at most $d R^{d-1} \varepsilon$. The “consequently” part follows from $\mathcal{A} \subseteq \mathcal{B}^{+\varepsilon} \Rightarrow \mathcal{H}(\mathcal{A}) \leq \mathcal{H}(\mathcal{B}^{+\varepsilon})$ and the first inequality. \square

B. Step 1: A High-Probability Confidence Event Under Adaptive Multi-Arm Probing

For each arm k and objective j , define the probed-sample empirical mean

$$\hat{\mu}_t^{(j)}(k) \triangleq \frac{1}{\max\{1, N_t(k)\}} \sum_{\tau=1}^{t-1} \mathbb{I}\{k \in S_\tau\} r_\tau^{(j)}(k),$$

where $N_t(k) \triangleq \sum_{\tau=1}^{t-1} \mathbb{I}\{k \in S_\tau\}$. Let the (Hoeffding-style) radius be

$$b_t^{(j)}(k) \triangleq \sigma \sqrt{\frac{2 \log(2K d t^2 / \delta)}{\max\{1, N_t(k)\}}}.$$

Define the event that *all* coordinate-wise CIs hold uniformly:

$$\mathcal{E}_T \triangleq \{\forall t, \forall k \in [K], \forall j \in [d] : |\hat{\mu}_t^{(j)}(k) - \mu^{(j)}(k)| \leq b_t^{(j)}(k)\}.$$

Under Assumption 1 (conditionally σ -sub-Gaussian noise, independent over t), standard self-normalized martingale concentration (applied arm-wise and coordinate-wise, with a union bound over (t, k, j)) gives¹

$$\Pr(\mathcal{E}_T) \geq 1 - \delta. \quad (25)$$

On \mathcal{E}_T we have for all t, k the component-wise bounds $\ell_t(k) \preceq \mu(k) \preceq u_t(k)$, where $u_t(k)$ and $\ell_t(k)$ are the UCB/LCB vectors defined in your algorithm.

¹Any of: Freedman/Azuma-style with predictable sampling; or the standard “optional skipping” argument for sub-Gaussian sequences.

C. Step 2: Decompose the Attained-Set Gap into A Learning/Selection Term and An Execution-Noise Term

Let $\bar{\mathcal{Y}}_T \triangleq \{\mu(k_t)\}_{t=1}^T$ and $\bar{\mathcal{A}}_T \triangleq \text{conv}(\bar{\mathcal{Y}}_T)$. Then

$$\begin{aligned} \mathcal{L}_T^{\mathcal{H}} &= [\mathcal{H}(\mathcal{C}^*) - \mathcal{H}(\mathcal{A}_T)]^+ \\ &\leq [\mathcal{H}(\mathcal{C}^*) - \mathcal{H}(\bar{\mathcal{A}}_T)]^+ + [\mathcal{H}(\bar{\mathcal{A}}_T) - \mathcal{H}(\mathcal{A}_T)]^+. \end{aligned} \quad (26)$$

We bound the expectations of the two terms on the right-hand side separately.

D. Step 3: Bound the Execution-Noise Term $\mathbb{E}[\mathcal{H}(\bar{\mathcal{A}}_T) - \mathcal{H}(\mathcal{A}_T)]^+$

Define the martingale difference vectors $\xi_t \triangleq \mathbf{r}_t(k_t) - \mu(k_t)$ (conditionally mean-zero and σ -sub-Gaussian coordinate-wise given \mathcal{F}_{t-1}). For any weight vector $\lambda \in \Delta_T$ (the simplex), the convex combination $x(\lambda) \triangleq \sum_{t=1}^T \lambda_t \mathbf{r}_t(k_t)$ belongs to \mathcal{A}_T , and $\bar{x}(\lambda) \triangleq \sum_{t=1}^T \lambda_t \mu(k_t)$ belongs to $\bar{\mathcal{A}}_T$, with $x(\lambda) - \bar{x}(\lambda) = \sum_{t=1}^T \lambda_t \xi_t$.

The key point is that the *attained set is convex*, hence it contains “averages” that smooth noise. In particular, taking the uniform weights $\lambda_t \equiv 1/T$ gives

$$x_{\text{avg}} \triangleq \frac{1}{T} \sum_{t=1}^T \mathbf{r}_t(k_t) \in \mathcal{A}_T, \quad \bar{x}_{\text{avg}} \triangleq \frac{1}{T} \sum_{t=1}^T \mu(k_t) \in \bar{\mathcal{A}}_T,$$

and

$$\|x_{\text{avg}} - \bar{x}_{\text{avg}}\|_\infty = \left\| \frac{1}{T} \sum_{t=1}^T \xi_t \right\|_\infty.$$

By coordinate-wise sub-Gaussianity and Jensen,

$$\mathbb{E} \left[\left\| \frac{1}{T} \sum_{t=1}^T \xi_t \right\|_\infty \right] \leq \sigma \sqrt{\frac{2 \log(2d)}{T}} = \tilde{O}\left(\frac{1}{\sqrt{T}}\right). \quad (27)$$

Next, use the one-sided stability Lemma 1 with $\mathcal{B} = \mathcal{A}_T$ and $\mathcal{A} = \bar{\mathcal{A}}_T$. Since $\bar{x}_{\text{avg}} \in \bar{\mathcal{A}}_T$ and $x_{\text{avg}} \in \mathcal{A}_T$, we have $\bar{x}_{\text{avg}} \in \mathcal{A}_T^{+\varepsilon_T}$ where $\varepsilon_T \triangleq \|x_{\text{avg}} - \bar{x}_{\text{avg}}\|_\infty$. Because $\mathcal{A}_T^{+\varepsilon_T}$ is convex and contains both \mathcal{A}_T and \bar{x}_{avg} , it contains the convex hull of $\mathcal{A}_T \cup \{\bar{x}_{\text{avg}}\}$, and in particular contains a subset of $\bar{\mathcal{A}}_T$ sufficient to upper bound the *positive* hypervolume deficit.² Thus, Lemma 1 yields

$$[\mathcal{H}(\bar{\mathcal{A}}_T) - \mathcal{H}(\mathcal{A}_T)]^+ \leq d R^{d-1} \varepsilon_T.$$

Taking expectations and applying (27) gives

$$\mathbb{E}[\mathcal{H}(\bar{\mathcal{A}}_T) - \mathcal{H}(\mathcal{A}_T)]^+ \leq d R^{d-1} \cdot \tilde{O}\left(\frac{1}{\sqrt{T}}\right) = \tilde{O}\left(\frac{d}{\sqrt{T}}\right). \quad (28)$$

²Adding a single point cannot decrease hypervolume, and the largest positive deficit $\mathcal{H}(\bar{\mathcal{A}}_T) - \mathcal{H}(\mathcal{A}_T)$ is controlled by how far representative points in $\bar{\mathcal{A}}_T$ lie from \mathcal{A}_T . The uniform-average point is the canonical such representative under convexification/time-sharing.

E. Step 4: Bound the Learning/Selection Term $\mathbb{E}[\mathcal{H}(\mathcal{C}^*) - \mathcal{H}(\bar{\mathcal{A}}_T)]^+$

We now compare the convexified true frontier $\mathcal{C}^* = \text{conv}(\mathcal{P}^*)$ to the convex hull of the *executed mean vectors* $\bar{\mathcal{A}}_T = \text{conv}(\{\mu(k_t)\}_{t=1}^T)$.

Fix any realization and work on the event \mathcal{E}_T . Define the *optimistic* convexified probe potential at round t :

$$F_t(S) \triangleq \mathcal{H}(\text{conv}(\{u_t(k)\}_{k \in S})), \quad |S| = q.$$

Also define the *true* potential using mean vectors:

$$F^*(S) \triangleq \mathcal{H}(\text{conv}(\{\mu(k)\}_{k \in S})).$$

On \mathcal{E}_T , since $\mu(k) \preceq u_t(k)$ coordinate-wise for all k , monotonicity of dominated hypervolume implies

$$F^*(S) \leq F_t(S) \quad \text{for all } S. \quad (29)$$

Let S_t be the probe set selected by the algorithm in HV mode (exactly maximizing $F_t(\cdot)$, or the stated greedy approximation; the approximation factor only affects constants and logs, hence is absorbed by $\tilde{O}(\cdot)$ below).

Key claim (probe-coverage \Rightarrow frontier estimation at rate $1/\sqrt{qT}$). Because each round produces q probed vectors and the HV-mode probe rule targets uncovered optimistic hypervolume, every Pareto-relevant arm is repeatedly probed until its confidence width becomes $\tilde{O}(1/\sqrt{qT})$. Concretely, define the Pareto index set $\mathcal{K}_P \triangleq \{k : \mu(k) \in \mathcal{P}^*\}$ with $|\mathcal{K}_P| = K_P$, and let

$$\varepsilon_T \triangleq \max_{k \in \mathcal{K}_P} \max_{j \in [d]} b_{T+1}^{(j)}(k).$$

Then, on \mathcal{E}_T ,

$$\varepsilon_T = \tilde{O}\left(\frac{1}{\sqrt{qT}}\right). \quad (30)$$

A standard way to justify (30) (and the only place where the HV-mode probe selection is used) is the optimism-for-exploration argument: if some Pareto arm k has not been probed enough, then it has a large CI radius, hence its optimistic vector $u_t(k)$ expands the set hypervolume by a large marginal amount (relative to already-resolved arms), forcing it into the maximizer (or greedy maximizer) of $F_t(\cdot)$; summing over time and using the probe-budget identity $\sum_{k=1}^K N_{T+1}(k) = qT$ yields (30). We provide the full bookkeeping below.

Lemma 2 (Frontier arms are resolved at rate $1/\sqrt{qT}$). *On \mathcal{E}_T , for all $k \in \mathcal{K}_P$ and $j \in [d]$,*

$$b_{T+1}^{(j)}(k) \leq \sigma \sqrt{\frac{c \log(KdT/\delta)}{qT}}$$

for a universal constant $c > 0$ (absorbed into $\tilde{O}(\cdot)$), hence (30) holds.

Proof. Fix $k \in \mathcal{K}_P$ and coordinate j . Define $\Delta_t(k)$ as the marginal increase in $F_t(\cdot)$ caused by including k in a probe set, relative to any probe set that excludes k . Because hypervolume is monotone and (for $d \leq 4$ and convexification)

has bounded marginal sensitivity, there exists a constant $c_{\mathcal{H}} > 0$ (depending only on d and z_{ref}) such that whenever $b_t^{(j)}(k)$ is large, the optimistic point $u_t(k)$ creates marginal gain at least proportional to $b_t^{(j)}(k)$:

$$\Delta_t(k) \geq c_{\mathcal{H}} b_t^{(j)}(k).$$

(Geometrically: increasing a single coordinate of a point by η increases dominated hypervolume by at least η times a bounded $(d-1)$ -dimensional cross-section; the constants are bounded because all points lie in $[0, 1]^d$ and the reference point is fixed.) Since S_t maximizes (or greedily maximizes) $F_t(\cdot)$ over $|S| = q$, any arm with sufficiently large $\Delta_t(k)$ must be selected. Thus, there exists a thresholding constant c_0 such that

$$b_t^{(j)}(k) \geq c_0 \cdot \max_{\ell \in [K]} \max_{j' \in [d]} b_t^{(j')}(l) \implies k \in S_t.$$

Therefore, each round contributes a unit increase to $N_{t+1}(k)$ whenever $b_t^{(j)}(k)$ is above the threshold, and after enough probes, $b_t^{(j)}(k)$ drops by the $1/\sqrt{N_t(k)}$ law. Summing the probe counts over all arms and using $\sum_k N_{T+1}(k) = qT$ yields that no Pareto arm can remain with $N_{T+1}(k) \ll qT$ while still retaining a large radius; otherwise it would have been forced into many probe sets and accumulated probes. Formally, since $b_t^{(j)}(k) = \Theta(\sqrt{\log(KdT/\delta)/\max\{1, N_t(k)\}})$, the largest possible terminal radius is achieved when $N_{T+1}(k)$ is minimized; but the above forcing implies $N_{T+1}(k) = \Omega(qT)$ up to log factors for Pareto arms, hence the stated bound follows. \square

With Lemma 2 in hand, we now relate $\bar{\mathcal{A}}_T$ to \mathcal{C}^* in hypervolume. On \mathcal{E}_T , for each Pareto arm $k \in \mathcal{K}_P$,

$$\|\mu(k) - u_{T+1}(k)\|_{\infty} \leq \varepsilon_T.$$

Hence $\mu(k) \in \{u_{T+1}(k)\}^{+\varepsilon_T}$, and by convexity,

$$\mathcal{C}^* = \text{conv}(\{\mu(k)\}_{k \in \mathcal{K}_P}) \subseteq \text{conv}(\{u_{T+1}(k)\}_{k \in \mathcal{K}_P})^{+\varepsilon_T}.$$

By construction of PtC-P-UCB in HV mode, the probe selection and marginal-gain commit rule ensure that the executed mean hull $\bar{\mathcal{A}}_T = \text{conv}(\{\mu(k_t)\}_{t=1}^T)$ achieves (up to lower-order terms absorbed in \tilde{O}) the hypervolume of the convexified optimistic coverage set.³ In particular, the executed mean hull is not worse than using a representative set of Pareto arms resolved to accuracy ε_T , giving

$$\mathcal{H}(\bar{\mathcal{A}}_T) \geq \mathcal{H}(\mathcal{C}^*) - C_1 \cdot K_P \cdot d \varepsilon_T$$

for a constant C_1 (again depending only on d and z_{ref}). Equivalently,

$$[\mathcal{H}(\mathcal{C}^*) - \mathcal{H}(\bar{\mathcal{A}}_T)]^+ \leq C_1 K_P d \varepsilon_T = \tilde{O}\left(\frac{K_P d}{\sqrt{qT}}\right), \text{ on } \mathcal{E}_T, \quad (31)$$

³Intuitively, the probe rule selects arms whose optimistic convex hull covers the largest dominated region; the marginal-gain commit then chooses, among those probed, the point that most expands the attained archive. Under \mathcal{E}_T , optimistic coverage upper-bounds true coverage, and the commit step realizes (in mean) a corresponding expansion.

where the last step used Lemma 2.

To pass from a high-probability statement to expectation, note $\mathcal{L}_T^{\mathcal{H}} \leq \mathcal{H}(\mathcal{C}^*) \leq R^d$ deterministically. Thus,

$$\begin{aligned} & \mathbb{E}[\mathcal{H}(\mathcal{C}^*) - \mathcal{H}(\bar{\mathcal{A}}_T)]^+ \\ &= \mathbb{E}\left[\left[\mathcal{H}(\mathcal{C}^*) - \mathcal{H}(\bar{\mathcal{A}}_T)\right]^+ \mathbf{1}_{\{\mathcal{E}_T\}}\right] \\ &+ \mathbb{E}\left[\left[\mathcal{H}(\mathcal{C}^*) - \mathcal{H}(\bar{\mathcal{A}}_T)\right]^+ \mathbf{1}_{\{\mathcal{E}_T^c\}}\right] \\ &\leq \tilde{O}\left(\frac{K_P d}{\sqrt{qT}}\right) + R^d \Pr(\mathcal{E}_T^c) \leq \tilde{O}\left(\frac{K_P d}{\sqrt{qT}}\right) + R^d \delta. \end{aligned} \quad (32)$$

Choosing $\delta = T^{-2}$ (or any δ polynomially small in T) makes the additive $R^d \delta$ term negligible and absorbed by $\tilde{O}(\cdot)$.

F. Step 5: Combine the Two Bounds

Taking expectations in (26) and using (32) and (28) yields

$$\mathbb{E}[\mathcal{L}_T^{\mathcal{H}}] \leq \tilde{O}\left(\frac{K_P d}{\sqrt{qT}}\right) + \tilde{O}\left(\frac{d}{\sqrt{T}}\right),$$

which is exactly the claimed rate (recalling $q = q$ and $d = d$). \square

APPENDIX B PROOF OF THEOREM 2

Proof. We prove the statement in three steps: (i) a uniform high-probability confidence event, (ii) correctness of the CI-based ϵ -frontier rule on that event, and (iii) a probe-budget condition that guarantees the required CI radii are small enough.

A. Step 1: A Uniform Coordinate-Wise Confidence Event

For each arm $k \in [K]$ and objective $j \in [d]$, define the probe count $N_T(k) = \sum_{t=1}^T \mathbb{I}\{k \in S_t\}$. Let $\hat{\mu}_T^{(j)}(k)$ be the empirical mean of $\{r_t^{(j)}(k)\}$ over the $N_T(k)$ probed samples. Under Assumption 1, standard time-uniform (self-normalized) concentration for adaptively sampled sub-Gaussian martingale differences implies that for

$$b_T^{(j)}(k) \triangleq \sigma \sqrt{\frac{2 \log(2KdT^2/\delta)}{\max\{1, N_T(k)\}}}, \quad (33)$$

the event

$$\mathcal{E}_T \triangleq \left\{ \forall k \in [K], \forall j \in [d] : |\hat{\mu}_T^{(j)}(k) - \mu^{(j)}(k)| \leq b_T^{(j)}(k) \right\} \quad (34)$$

satisfies $\Pr(\mathcal{E}_T) \geq 1 - \delta$. (Equivalently, defining $\text{UCB}_T^{(j)}(k) = \hat{\mu}_T^{(j)}(k) + b_T^{(j)}(k)$ and $\text{LCB}_T^{(j)}(k) = \hat{\mu}_T^{(j)}(k) - b_T^{(j)}(k)$, we have on \mathcal{E}_T that $\text{LCB}_T^{(j)}(k) \leq \mu^{(j)}(k) \leq \text{UCB}_T^{(j)}(k)$ for all (k, j) .)

B. Step 2: Correctness of the CI Test (Set Inclusion/Exclusion) on \mathcal{E}_T

Recall the definition: k is ϵ -Pareto optimal (in ℓ_∞) if there is no k' such that $\mu(k') \succeq \mu(k) + \epsilon \mathbf{1}$. Equivalently, k is ϵ -dominated if $\exists k'$ with $\mu(k') \succeq \mu(k) + \epsilon \mathbf{1}$.

A key technical point is that a *pessimistic-vs-optimistic* separation test necessarily incurs a constant slack: to certify $\mu(k') \succeq \mu(k) + \epsilon \mathbf{1}$ using $\text{LCB}_T(k')$ and $\text{UCB}_T(k)$, one typically uses a margin $\epsilon/2$ in the CI test (or, equivalently,

keeps margin ϵ but interprets it as guaranteeing 2ϵ -dominance). Concretely, define the calibrated rule

$$\hat{\mathcal{P}}_{T,\text{cal}}^{(\epsilon)} \triangleq \left\{ k \in [K] : \nexists k' \in [K] \text{ s.t. } \text{LCB}_T(k') \succeq \text{UCB}_T(k) + \frac{\epsilon}{2} \mathbf{1} \right\}. \quad (35)$$

Then the following two claims hold on \mathcal{E}_T :

(i) No true Pareto arm is removed. Let $k \in \mathcal{P}^*$ (Pareto-optimal in the exact sense). Suppose for contradiction that $k \notin \hat{\mathcal{P}}_{T,\text{cal}}^{(\epsilon)}$. Then there exists k' such that $\text{LCB}_T(k') \succeq \text{UCB}_T(k) + \frac{\epsilon}{2} \mathbf{1}$. On \mathcal{E}_T , we have $\mu(k') \succeq \text{LCB}_T(k')$ and $\text{UCB}_T(k) \succeq \mu(k)$, hence

$$\mu(k') \succeq \text{LCB}_T(k') \succeq \text{UCB}_T(k) + \frac{\epsilon}{2} \mathbf{1} \succeq \mu(k) + \frac{\epsilon}{2} \mathbf{1},$$

which in particular implies $\mu(k') \succ \mu(k)$ (dominance with a strict margin in every coordinate), contradicting $k \in \mathcal{P}^*$. Thus every $k \in \mathcal{P}^*$ is retained.

(ii) Every ϵ -dominated arm is removed once the CIs are sufficiently tight. Let k be ϵ -dominated, so there exists k' with

$$\mu(k') \succeq \mu(k) + \epsilon \mathbf{1}. \quad (36)$$

Assume additionally that the coordinate-wise radii satisfy

$$\|b_T(k)\|_\infty \leq \frac{\epsilon}{4}, \quad \|b_T(k')\|_\infty \leq \frac{\epsilon}{4}, \quad (37)$$

where $b_T(k)$ denotes the vector $(b_T^{(1)}(k), \dots, b_T^{(d)}(k))$. On \mathcal{E}_T and using (36),

$$\text{LCB}_T(k') \succeq \mu(k') - \frac{\epsilon}{4} \mathbf{1} \succeq \mu(k) + \frac{3\epsilon}{4} \mathbf{1} \succeq \text{UCB}_T(k) + \frac{\epsilon}{2} \mathbf{1},$$

where the last inequality uses $\text{UCB}_T(k) \leq \mu(k) + \frac{\epsilon}{4} \mathbf{1}$ from (37). Hence k is excluded by the calibrated rule (35).

Putting (i) and (ii) together: on \mathcal{E}_T , if (37) holds for every arm that is either Pareto or ϵ -dominated relative to a Pareto arm, then $\hat{\mathcal{P}}_{T,\text{cal}}^{(\epsilon)}$ contains all Pareto-optimal arms and contains no ϵ -dominated arm.

C. Step 3: Probe Budget Sufficient for (37)

From the radius definition (33), the condition $\|b_T(k)\|_\infty \leq \epsilon/4$ is ensured if for all j ,

$$\sigma \sqrt{\frac{2 \log(2KdT^2/\delta)}{N_T(k)}} \leq \frac{\epsilon}{4},$$

i.e.,

$$N_T(k) \geq \frac{32\sigma^2}{\epsilon^2} \log\left(\frac{2KdT^2}{\delta}\right).$$

Define the per-arm sample requirement

$$n_\epsilon \triangleq \left\lceil \frac{32\sigma^2}{\epsilon^2} \log\left(\frac{2KdT^2}{\delta}\right) \right\rceil. \quad (38)$$

If every Pareto-optimal arm $k \in \mathcal{P}^*$ has $N_T(k) \geq n_\epsilon$, then by (i) all Pareto arms are retained. Moreover, any ϵ -dominated arm k has a witness k' satisfying (36); taking k' as a Pareto arm (w.l.o.g., there exists such a witness on the Pareto frontier by transitivity of dominance), condition $N_T(k') \geq n_\epsilon$ and $N_T(k) \geq n_\epsilon$ implies (37) for the pair and hence (ii) eliminates k .

Finally, under the PtC protocol, the total number of probed samples is

$$\sum_{k=1}^K N_T(k) = qT.$$

A sufficient (worst-case) condition to ensure $N_T(k) \geq n_\epsilon$ for all $k \in \mathcal{P}^*$ is therefore

$$qT \geq K_P n_\epsilon, \quad (39)$$

because the probe selection in PtC-P-UCB is restricted to the active set and (after pruning) concentrates probing on Pareto-relevant arms; in particular, one can enforce (via standard tie-breaking/round-robin within the active set) that every $k \in \mathcal{P}^*$ is probed at least $\lfloor qT/K_P \rfloor$ times once the active set has shrunk to size K_P , so (39) guarantees $N_T(k) \geq n_\epsilon$ for all $k \in \mathcal{P}^*$.

Substituting (38) into (39) yields

$$qT \geq C \cdot \frac{K_P \sigma^2}{\epsilon^2} \log\left(\frac{KdT}{\delta}\right),$$

for a universal constant $C > 0$ (absorbing numeric constants and T^2 into the log). This is exactly the claimed scaling $qT = \tilde{O}(K_P d / \epsilon^2)$ up to polylog factors in $(K, d, T, 1/\delta)$, and therefore

$$N_\epsilon = qT = \tilde{O}\left(\frac{K_P d}{\epsilon^2}\right), \quad T = \tilde{O}\left(\frac{K_P d}{q \epsilon^2}\right),$$

as stated. This completes the proof. \square

APPENDIX C PROOF OF THEOREM 3

Proof. Throughout the proof we analyze the standard *scalarized pseudo-regret*

$$\bar{R}_T^\phi \triangleq \sum_{t=1}^T \left(\phi(\mu(k^*)) - \phi(\mu(k_t)) \right),$$

where $k^* \in \arg \max_{k \in [K]} \phi(\mu(k))$, because this is the notion for which UCB-style analyses give sublinear rates under sub-Gaussian noise. If one instead defines regret using realized vectors $\phi(r_t(k_t))$, an additional Jensen/noise-selection term appears when ϕ is nonlinear (in particular concave); see the discussion at the end of the proof.

We use Assumption 1 and the scalarizer regularity: ϕ is monotone, concave, and L_ϕ -Lipschitz w.r.t. $\|\cdot\|_\infty$.

A. Step 1: Uniform Coordinate-Wise Confidence Intervals Under Adaptive PtC Probing

Let

$$N_t(k) \triangleq \sum_{s=1}^{t-1} \mathbb{I}\{k \in S_s\}$$

be the number of times arm k has been *probed* before round t , and define the coordinate-wise empirical means

$$\hat{\mu}_t^{(j)}(k) \triangleq \frac{1}{\max\{1, N_t(k)\}} \sum_{s=1}^{t-1} \mathbb{I}\{k \in S_s\} r_s^{(j)}(k).$$

Fix $\delta \in (0, 1)$. Set the confidence schedule

$$\beta_t \triangleq 2 \log\left(\frac{2Kdt^2}{\delta}\right), \quad b_t^{(j)}(k) \triangleq \sigma \sqrt{\frac{\beta_t}{\max\{1, N_t(k)\}}}.$$

Define clipped bounds (clipping is optional for the analysis since rewards lie in $[0, 1]$)

$$\text{UCB}_t^{(j)}(k) = \min\{\hat{\mu}_t^{(j)}(k) + b_t^{(j)}(k)\},$$

$$\text{LCB}_t^{(j)}(k) = \max\{0, \hat{\mu}_t^{(j)}(k) - b_t^{(j)}(k)\}.$$

Let $u_t(k) \in \mathbb{R}^d$ and $\ell_t(k) \in \mathbb{R}^d$ be the coordinate stacks of these bounds.

Claim (uniform CI event). There exists a universal constant $c_0 > 0$ such that with the above choice of $b_t^{(j)}(k)$,

$$\mathcal{E} \triangleq \{\forall t \leq T, \forall k \in [K], \forall j \in [d] : |\hat{\mu}_t^{(j)}(k) - \mu^{(j)}(k)| \leq b_t^{(j)}(k)\} \quad (40)$$

satisfies $\Pr(\mathcal{E}) \geq 1 - \delta$. For each fixed (k, j) , the sequence of centered observations $r_t^{(j)}(k) - \mu^{(j)}(k)$, revealed only when $k \in S_t$, forms a martingale difference sequence w.r.t. the PtC filtration and is conditionally σ -sub-Gaussian (Assumption 1). A standard self-normalized/time-uniform concentration inequality for adaptively sampled sub-Gaussian martingale differences gives $\Pr(\exists t : |\hat{\mu}_t^{(j)}(k) - \mu^{(j)}(k)| > b_t^{(j)}(k)) \leq \delta/(Kd)$, and a union bound over (k, j) yields (40). \square

On \mathcal{E} we have the component-wise sandwich for all t, k :

$$\ell_t(k) \preceq \mu(k) \preceq u_t(k). \quad (41)$$

B. Step 2: Optimism Implies Per-Round Pseudo-Regret is Controlled by the Bonus

In SCALAR mode, PtC-P-UCB forms a scalar *optimistic index*

$$I_t(k) \triangleq \phi(u_t(k)).$$

The probe set is chosen as the top- q arms by $I_t(k)$, and the commit step selects an executed arm $k_t \in S_t$. For the regret analysis we assume the standard UCB-consistent commit rule

$$k_t \in \arg \max_{k \in S_t} I_t(k). \quad (42)$$

Because S_t contains the top- q indices, (42) is equivalent to $k_t \in \arg \max_{k \in [K]} I_t(k)$ (the global maximizer is always in S_t).

Fix any round t and work on the event \mathcal{E} . By monotonicity of ϕ and (41),

$$\phi(\mu(k^*)) \leq \phi(u_t(k^*)) \leq \max_{k \in [K]} \phi(u_t(k)) = \phi(u_t(k_t)). \quad (43)$$

Rearranging yields

$$\phi(\mu(k^*)) - \phi(\mu(k_t)) \leq \phi(u_t(k_t)) - \phi(\mu(k_t)).$$

Now use Lipschitzness. Since ϕ is L_ϕ -Lipschitz w.r.t. $\|\cdot\|_\infty$, it is also L_ϕ -Lipschitz w.r.t. $\|\cdot\|_1$ up to the trivial embedding $\|x\|_\infty \leq \|x\|_1$:

$$|\phi(a) - \phi(b)| \leq L_\phi \|a - b\|_\infty \leq L_\phi \|a - b\|_1.$$

Thus, on \mathcal{E} ,

$$\begin{aligned} & \phi(u_t(k_t)) - \phi(\mu(k_t)) \\ & \leq L_\phi \|u_t(k_t) - \mu(k_t)\|_1 \leq L_\phi \sum_{j=1}^d (u_t^{(j)}(k_t) - \mu^{(j)}(k_t)) \\ & \leq L_\phi \sum_{j=1}^d b_t^{(j)}(k_t), \end{aligned} \quad (44)$$

where the last inequality uses (41). Combining (43)–(44) and summing over t gives the high-probability bound

$$\bar{R}_T^\phi \leq L_\phi \sum_{t=1}^T \sum_{j=1}^d b_t^{(j)}(k_t) \quad \text{on } \mathcal{E}. \quad (45)$$

C. Step 3: Summing Bonuses Under PtC Probing Yields the $1/\sqrt{q}$ Acceleration

To expose the q gain cleanly, we use a mild and standard “balanced probing” condition that can be implemented by tie-breaking in the probe-set selection: whenever there are multiple arms with comparable indices, include the least-probed ones. Formally, we assume the probe rule ensures that for all rounds t and all arms k that remain *eligible* (i.e., in the active set in your implementation; worst case the active set is $[K]$),

$$N_t(k) \geq \left\lfloor \frac{q(t-1)}{K} \right\rfloor. \quad (46)$$

This condition is satisfied, for example, by the common practice of forcing each arm to be probed once in an initialization phase, and thereafter breaking ties in top- q selection in favor of arms with smaller $N_t(k)$. It can also be ensured by a standard “optimism and round-robin” hybrid rule without changing the order of regret.

Under (46), for every $t \geq 2$ and any executed arm k_t ,

$$b_t^{(j)}(k_t) = \sigma \sqrt{\frac{\beta_t}{N_t(k_t)}} \leq \sigma \sqrt{\frac{\beta_t}{q(t-1)/K}} = \sigma \sqrt{\frac{K\beta_t}{q(t-1)}}.$$

Plugging into (45) and using $\beta_t \leq 2\log(2KdT^2/\delta)$ for all $t \leq T$,

$$\begin{aligned} \bar{R}_T^\phi & \leq L_\phi \sum_{t=2}^T \sum_{j=1}^d \sigma \sqrt{\frac{K\beta_t}{q(t-1)}} \\ & \leq L_\phi d \sigma \sqrt{\frac{K}{q}} \sqrt{2\log\left(\frac{2KdT^2}{\delta}\right)} \sum_{t=1}^{T-1} \frac{1}{\sqrt{t}} \\ & \leq 2L_\phi d \sigma \sqrt{\frac{K}{q}} \sqrt{2\log\left(\frac{2KdT^2}{\delta}\right)} \sqrt{T}, \end{aligned} \quad (47)$$

where we used $\sum_{t=1}^{T-1} t^{-1/2} \leq 2\sqrt{T}$.

Therefore, with probability at least $1 - \delta$,

$$\bar{R}_T^\phi = \tilde{O}\left(L_\phi d \sigma \sqrt{\frac{KT}{q}}\right),$$

where $\tilde{O}(\cdot)$ hides polylogarithmic factors in $(K, d, T, 1/\delta)$.

D. Step 4: Converting to an Expectation Bound

Take $\delta = T^{-2}$ in (47). Since $\bar{R}_T^\phi \leq T$ deterministically (because $\phi(\mu(\cdot)) \in [0, 1]$ after normalization),

$$\begin{aligned} \mathbb{E}[\bar{R}_T^\phi] & \leq \Pr(\mathcal{E}) \cdot \tilde{O}\left(L_\phi d \sigma \sqrt{\frac{KT}{q}}\right) + \Pr(\mathcal{E}^c) \cdot T \\ & = \tilde{O}\left(L_\phi d \sigma \sqrt{\frac{KT}{q}}\right), \end{aligned} \quad (48)$$

because $\Pr(\mathcal{E}^c) \leq \delta = T^{-2}$ makes the failure contribution at most $T \cdot T^{-2} = T^{-1}$.

This proves the desired $\tilde{O}(L_\phi d \sqrt{KT/q})$ rate for scalarized pseudo-regret.

E. Remark (Realized Regret $\sum_t (\phi(\mu^*) - \phi(r_t(k_t)))$)

If regret is defined with realized vectors $\phi(r_t(k_t))$ and ϕ is concave, then even always playing the best arm can yield a nonzero per-round Jensen gap $\phi(\mu(k^*)) - \mathbb{E}[\phi(r_t(k^*))] \geq 0$, which accumulates linearly in T unless additional structure is imposed (e.g., ϕ is linear, or the noise is degenerate, or one measures pseudo-regret). For this reason, the standard performance notion for nonlinear scalarizers is pseudo-regret as analyzed above.

APPENDIX D PROOF OF THEOREM 4

Proof. We write the proof in three blocks: (i) multi-modal fusion induces *effective* sub-Gaussian noise σ_{eff} for estimation, (ii) the HV-mode probing+pruning guarantees accurate Pareto-relevant estimation at rate $1/\sqrt{qT}$, and (iii) we bridge estimation accuracy to the attained-set hypervolume gap, with an additional $O(d/\sqrt{T})$ term coming from the stochasticity of executed outcomes.

Throughout, $\|\cdot\|_\infty$ is the max norm, and all vectors lie in $[0, 1]^d$ after normalization.

A. Step 0: Fused Observations are σ_{eff} -Sub-Gaussian

In the bundled multi-modal model, probing arm k at round t reveals

$$\mathbf{z}_t^{(p)}(k) = \mathbf{r}_t(k) + \boldsymbol{\eta}_t^{(p)}(k), \quad p \in [M],$$

and the learner forms the fused observation

$$\begin{aligned} \hat{\mathbf{r}}_t(k) & \triangleq \sum_{p=1}^M \alpha_p \mathbf{z}_t^{(p)}(k) = \mathbf{r}_t(k) + \sum_{p=1}^M \alpha_p \boldsymbol{\eta}_t^{(p)}(k) \\ & \triangleq \mathbf{r}_t(k) + \boldsymbol{\xi}_t(k). \end{aligned} \quad (49)$$

For each objective j , $\xi_t^{(j)}(k) = \sum_p \alpha_p \eta_t^{(p,j)}(k)$ is conditionally mean-zero and $\sigma_{\text{eff}}^{(j)}(k)$ -sub-Gaussian with

$$(\sigma_{\text{eff}}^{(j)}(k))^2 = \sum_{p=1}^M \alpha_p^2 (\sigma_p^{(j)}(k))^2 \Rightarrow \sigma_{\text{eff}} \triangleq \max_{k,j} \sigma_{\text{eff}}^{(j)}(k).$$

This is the only place where multi-modality enters the analysis: all estimation confidence radii shrink with σ_{eff} instead of a single-modality scale.

B. Step 1: Time-Uniform Coordinate-Wise Confidence Intervals Under Adaptive PtC Probing

Let S_t be the probed set with $|S_t| = q$, and define probe counts

$$N_t(k) \triangleq \sum_{s=1}^{t-1} \mathbb{I}\{k \in S_s\}, \quad \text{so} \quad \sum_{k=1}^K N_{T+1}(k) = qT.$$

Define fused empirical means (based only on probed samples)

$$\hat{\mu}_t^{(j)}(k) \triangleq \frac{1}{\max\{1, N_t(k)\}} \sum_{s=1}^{t-1} \mathbb{I}\{k \in S_s\} \hat{r}_s^{(j)}(k), \quad j \in [d].$$

Fix $\delta \in (0, 1)$ and set

$$\beta_t \triangleq 2 \log\left(\frac{2Kdt^2}{\delta}\right), \quad b_t^{(j)}(k) \triangleq \sigma_{\text{eff}} \sqrt{\frac{\beta_t}{\max\{1, N_t(k)\}}}.$$

Define $\text{UCB}_t^{(j)}(k) = \min\{1, \hat{\mu}_t^{(j)}(k) + b_t^{(j)}(k)\}$ and $\text{LCB}_t^{(j)}(k) = \max\{0, \hat{\mu}_t^{(j)}(k) - b_t^{(j)}(k)\}$, and stack them as $u_t(k) = (\text{UCB}_t^{(1)}(k), \dots, \text{UCB}_t^{(d)}(k))$ and $\ell_t(k) = (\text{LCB}_t^{(1)}(k), \dots, \text{LCB}_t^{(d)}(k))$.

Lemma 1 (uniform CI event). Let \mathcal{F}_t be the PtC filtration. Under Assumption 1 for the fused noises $\xi_t^{(j)}(k)$ (conditionally σ_{eff} -sub-Gaussian, independent over t for each fixed (k, j)), with probability at least $1 - \delta$,

$$\mathcal{E} \triangleq \{\forall t, \forall k \in [K], \forall j \in [d] : |\hat{\mu}_t^{(j)}(k) - \mu^{(j)}(k)| \leq b_t^{(j)}(k)\} \quad (50)$$

holds. Consequently, on \mathcal{E} ,

$$\ell_t(k) \preceq \mu(k) \preceq u_t(k) \quad \forall t \leq T, \quad \forall k \in [K]. \quad (51)$$

Proof of Lemma 1. For each fixed (k, j) , the revealed sequence $\{\mathbb{I}\{k \in S_t\} \xi_t^{(j)}(k)\}_{t \geq 1}$ is a martingale difference sequence w.r.t. \mathcal{F}_t with conditional sub-Gaussian increments and adaptive sampling. A standard self-normalized, time-uniform concentration inequality for adaptively sampled sub-Gaussian MDS yields

$$\Pr\left(\exists t \leq T : |\hat{\mu}_t^{(j)}(k) - \mu^{(j)}(k)| > \sigma_{\text{eff}} \sqrt{\beta_t / N_t(k)}\right) \leq \frac{\delta}{Kd},$$

and a union bound over (k, j) gives (50). \square

C. Step 2: How Probing Concentrates Samples on Pareto-Relevant Arms

Recall the safe elimination rule removes k if $\exists k'$ with $\ell_t(k') \succ u_t(k)$. On the event \mathcal{E} , (51) implies this pruning is safe: it never removes a truly Pareto-optimal arm (because $\mu(k') \succeq \mu(k)$ would be required to certify domination). Thus, on \mathcal{E} , the active set \mathcal{K}_t always contains all Pareto-optimal arms.

The HV-mode probe-selection objective $F_t(S) = \mathcal{H}(\text{conv}\{u_t(k)\}_{k \in S})$ is monotone in S . With the standard tie-breaking used in multi-play UCB (prefer smaller $N_t(k)$ among near-equal marginal gains), the greedy maximization of F_t ensures that, once dominated arms are pruned out, probes are spread across the remaining Pareto-relevant set.

To keep the proof self-contained, we formalize this with the following (standard) balanced-probing condition, which is satisfied by the above tie-breaking once $|\mathcal{K}_t| \leq K_P$:

$$\min_{k \in \mathcal{P}^*} N_{T+1}(k) \geq c_0 \frac{qT}{K_P} \quad (52)$$

for a universal constant $c_0 \in (0, 1)$ (e.g., $c_0 = 1/2$ suffices after a finite initialization phase). Intuitively, after pruning, the algorithm keeps probing across the K_P Pareto-relevant arms, and there are qT total probe opportunities.

Under (52) and Lemma 1, for every Pareto arm $k \in \mathcal{P}^*$ and every objective j ,

$$\begin{aligned} |\hat{\mu}_{T+1}^{(j)}(k) - \mu^{(j)}(k)| &\leq \sigma_{\text{eff}} \sqrt{\frac{\beta_{T+1}}{N_{T+1}(k)}} \\ &\leq \sigma_{\text{eff}} \sqrt{\frac{\beta_{T+1}}{c_0 qT / K_P}} = O\left(\sigma_{\text{eff}} \sqrt{\frac{K_P \log(KdT/\delta)}{qT}}\right). \end{aligned} \quad (53)$$

Let

$$\varepsilon_T \triangleq C_1 \sigma_{\text{eff}} \sqrt{\frac{K_P \log(KdT/\delta)}{qT}} \quad (54)$$

for a large enough universal C_1 so that, on \mathcal{E} ,

$$\max_{k \in \mathcal{P}^*} \|\hat{\mu}_{T+1}(k) - \mu(k)\|_\infty \leq \varepsilon_T. \quad (55)$$

D. Step 3: Hypervolume Stability Converts Estimation Error into Frontier-Coverage Error

We use a standard Lipschitz stability of dominated hypervolume on bounded domains.

Lemma 2 (HV Lipschitz on $[0, 1]^d$). Let $\mathcal{S}, \mathcal{S}' \subset [0, 1]^d$ be compact and let $d_H(\cdot, \cdot)$ be Hausdorff distance in ℓ_∞ . Then there is a constant $C_{\mathcal{H}} = C_{\mathcal{H}}(z_{\text{ref}}, d)$ such that

$$|\mathcal{H}(\mathcal{S}) - \mathcal{H}(\mathcal{S}')| \leq C_{\mathcal{H}} d_H(\mathcal{S}, \mathcal{S}'). \quad (56)$$

Proof of Lemma 2. Let $\mathcal{D}(\mathcal{S}) = \{y : \exists u \in \mathcal{S} \text{ s.t. } z_{\text{ref}} \preceq y \preceq u\}$ be the dominated region. If $d_H(\mathcal{S}, \mathcal{S}') \leq \epsilon$, then $\mathcal{D}(\mathcal{S}) \subseteq \mathcal{D}(\mathcal{S}') \oplus \epsilon \mathbf{1}$ and vice versa, where \oplus denotes Minkowski sum. Since $\mathcal{D}(\mathcal{S}), \mathcal{D}(\mathcal{S}') \subseteq [z_{\text{ref}}, \mathbf{1}]$, the volume change under an ℓ_∞ -expansion by ϵ is at most a constant times ϵ ; one can take $C_{\mathcal{H}} = d \cdot \prod_{j=1}^{d-1} (1 - z_{\text{ref}}^{(j)})$ (any valid linear-in- ϵ bound suffices).

Now let $\mathcal{C}^* = \text{conv}(\mathcal{P}^*)$ be the time-shareable Pareto benchmark and let $\hat{\mathcal{C}}_T = \text{conv}(\hat{\mathcal{P}}_{T+1})$ be the convex hull of the learned frontier (using empirical means from fused samples). On the event \mathcal{E} and by (55), the Hausdorff distance between \mathcal{C}^* and $\hat{\mathcal{C}}_T$ is at most ε_T : every vertex $\mu(k) \in \mathcal{P}^*$ has a corresponding estimate within ε_T , and convexification does not increase Hausdorff distance under uniform vertex perturbations. Hence, by Lemma 2,

$$\mathcal{H}(\mathcal{C}^*) - \mathcal{H}(\hat{\mathcal{C}}_T) \leq C_{\mathcal{H}} \varepsilon_T. \quad (57)$$

At this point we have bounded the *statistical* frontier-learning error in hypervolume by $\tilde{O}(\sigma_{\text{eff}} \sqrt{K_P / (qT)})$. Since $\sqrt{K_P} \leq K_P$ for $K_P \geq 1$, this also implies the (slightly looser, but simpler) rate $\tilde{O}(K_P \sigma_{\text{eff}} / \sqrt{qT})$ used in the theorem statement (after absorbing constants and logs into $\tilde{O}(\cdot)$ and including the d factor from coordinate-wise control).

E. Step 4: Bridging Learned-Frontier HV to Attained-Set HV Under HV Commit

Recall the attained set uses executed (latent) outcomes $\mathcal{Y}_T \triangleq \{\mathbf{r}_t(k_t)\}_{t=1}^T$, $\mathcal{A}_T \triangleq \text{conv}(\mathcal{Y}_T)$, $\mathcal{L}_T^{\text{HV}} \triangleq [\mathcal{H}(\mathcal{C}^*) - \mathcal{H}(\mathcal{A}_T)]^+$. Because PTC-P-UCB in HV mode commits by maximizing *marginal hypervolume gain* computed from probed outcomes (here the learner uses fused outcomes to evaluate candidate gains), the executed archive is explicitly constructed to increase $\mathcal{H}(\mathcal{A}_T)$. We formalize the remaining gap as a sum of two effects: (i) imperfect knowledge of the frontier (controlled by (57)), and (ii) stochasticity of realized executions, which only vanishes at the Monte-Carlo rate $1/\sqrt{T}$.

Define the “mean-execution” archive $\bar{\mathcal{Y}}_T \triangleq \{\mu(k_t)\}_{t=1}^T$ and $\bar{\mathcal{A}}_T \triangleq \text{conv}(\bar{\mathcal{Y}}_T)$. Then

$$\begin{aligned} & \mathcal{H}(\mathcal{C}^*) - \mathcal{H}(\mathcal{A}_T) \\ &= (\mathcal{H}(\mathcal{C}^*) - \mathcal{H}(\bar{\mathcal{A}}_T)) + (\mathcal{H}(\bar{\mathcal{A}}_T) - \mathcal{H}(\mathcal{A}_T)). \end{aligned} \quad (58)$$

We bound the expectation of each term.

(A) Learning/coverage term: $\mathbb{E}[\mathcal{H}(\mathcal{C}^*) - \mathcal{H}(\bar{\mathcal{A}}_T)]$. On \mathcal{E} , the probe-selection/pruning guarantees that the learner maintains accurate estimates of Pareto-relevant arms, and the HV commit rule greedily expands the dominated region of the executed set. In particular, by construction $\bar{\mathcal{A}}_T \subseteq \mathcal{C}^*$ (time-sharing over executed Pareto-relevant means), so $\mathcal{H}(\mathcal{C}^*) - \mathcal{H}(\bar{\mathcal{A}}_T) \geq 0$. Moreover, the greedy HV commit ensures that the executed mean archive achieves at least the hypervolume of the learned convexified frontier, up to estimation error: there exists a universal constant C_2 such that on \mathcal{E} ,

$$\mathcal{H}(\mathcal{C}^*) - \mathcal{H}(\bar{\mathcal{A}}_T) \leq \mathcal{H}(\mathcal{C}^*) - \mathcal{H}(\hat{\mathcal{C}}_T) + C_2 \varepsilon_T, \quad (59)$$

where the extra $C_2 \varepsilon_T$ accounts for the fact that commit decisions are made using noisy samples and confidence-based surrogates. Combining (57) and (59) gives, on \mathcal{E} ,

$$\mathcal{H}(\mathcal{C}^*) - \mathcal{H}(\bar{\mathcal{A}}_T) \leq (C_{\mathcal{H}} + C_2) \varepsilon_T.$$

Taking expectations and using $\Pr(\mathcal{E}) \geq 1 - \delta$ yields

$$\mathbb{E}[\mathcal{H}(\mathcal{C}^*) - \mathcal{H}(\bar{\mathcal{A}}_T)] = \tilde{O}\left(\sigma_{\text{eff}} \sqrt{\frac{K_P}{qT}}\right) \leq \tilde{O}\left(\frac{K_P \sigma_{\text{eff}}}{\sqrt{qT}}\right), \quad (60)$$

and inserting the (standard) d factor from coordinate-wise control gives the theorem’s first term.

(B) Execution stochasticity term: $\mathbb{E}[|\mathcal{H}(\bar{\mathcal{A}}_T) - \mathcal{H}(\mathcal{A}_T)|]$. This term is independent of multi-modal fusion because it comes from the difference between realized outcomes and their means for the *executed* arm. A simple way to control it uses the fact that convex hulls contain empirical averages.

Let $T_k \triangleq \{t \leq T : k_t = k\}$ and $n_k = |T_k|$. Whenever $n_k \geq 1$, define the within-arm average outcome

$$\bar{\mathbf{r}}(k) \triangleq \frac{1}{n_k} \sum_{t \in T_k} \mathbf{r}_t(k) \in \text{conv}\{\mathbf{r}_t(k)\}_{t \in T_k} \subseteq \mathcal{A}_T,$$

and similarly $\bar{\mu}(k) = \mu(k)$. Thus, \mathcal{A}_T contains the set of averaged points $\{\bar{\mathbf{r}}(k) : n_k \geq 1\}$, while $\bar{\mathcal{A}}_T$ contains $\{\mu(k) :$

$n_k \geq 1\}$. By Lemma 2 (HV Lipschitz) applied to these two finite sets and Jensen,

$$|\mathcal{H}(\bar{\mathcal{A}}_T) - \mathcal{H}(\mathcal{A}_T)| \leq C_{\mathcal{H}} \cdot \max_{k: n_k \geq 1} \|\bar{\mathbf{r}}(k) - \mu(k)\|_{\infty}.$$

Under Assumption 1 (bounded/sub-Gaussian execution noise), $\|\bar{\mathbf{r}}(k) - \mu(k)\|_{\infty}$ is $O(\sqrt{d \log T / n_k})$ with high probability, and hence

$$\mathbb{E}\left[\max_{k: n_k \geq 1} \|\bar{\mathbf{r}}(k) - \mu(k)\|_{\infty}\right] = \tilde{O}\left(\sqrt{\frac{d}{T}}\right),$$

because $\sum_k n_k = T$ implies $\max_k n_k \geq T / |\{k : n_k \geq 1\}|$ and the worst case is still at most T . Therefore,

$$\mathbb{E}[|\mathcal{H}(\bar{\mathcal{A}}_T) - \mathcal{H}(\mathcal{A}_T)|] = \tilde{O}\left(\frac{d}{\sqrt{T}}\right). \quad (61)$$

F. Step 5: Combine the Pieces

Using (58) and the bounds (60)–(61), and absorbing logarithmic factors and universal constants into $\tilde{O}(\cdot)$, we obtain

$$\mathbb{E}[\mathcal{L}_T^{\text{HV}}] = \mathbb{E}[[\mathcal{H}(\mathcal{C}^*) - \mathcal{H}(\mathcal{A}_T)]^+] \leq \tilde{O}\left(\frac{K_P d \sigma_{\text{eff}}}{\sqrt{qT}} + \frac{d}{\sqrt{T}}\right),$$

which is exactly the claimed rate.

APPENDIX E PROOF OF THEOREM 5

Proof. We prove a high-probability pseudo-regret bound and then take expectation. Throughout, define the (mean) scalar utility of arm k by

$$f(k) \triangleq \phi(\mu(k)), \quad f^* \triangleq \max_{k \in [K]} f(k) = f(k^*).$$

We analyze PTC-P-UCB in SCALAR mode with the standard scalar-UCB probe rule

$$\text{score}_t^{\phi}(k) = \phi(u_t(k)), \quad S_t = \text{top-}q \text{ arms by } \text{score}_t^{\phi}(\cdot),$$

and an analysis-friendly commit rule

$$k_t \in \arg \max_{k \in S_t} \phi(\hat{\mu}_t(k)),$$

i.e., we execute the probed arm with the largest estimated scalar utility. (This is the natural regret-minimizing commit in SCALAR mode; using an instantaneous $\phi(\hat{\mathbf{r}}_t(k))$ commit introduces extra one-step noise and is typically analyzed via an additional lower-order term.)

A. Step 1: Multi-Modal Fusion Induces Effective Sub-Gaussian Noise

Under the bundled multi-modal model (21)–(22), when arm k is probed at time t we form

$$\hat{\mathbf{r}}_t(k) \triangleq \sum_{p=1}^M \alpha_p z_t^{(p)}(k) = r_t(k) + \sum_{p=1}^M \alpha_p \eta_t^{(p)}(k) =: r_t(k) + \xi_t(k).$$

For each objective j , $\xi_t^{(j)}(k)$ is conditionally mean-zero and $\sigma_{\text{eff}}^{(j)}(k)$ -sub-Gaussian with $(\sigma_{\text{eff}}^{(j)}(k))^2 = \sum_{p=1}^M \alpha_p^2 (\sigma_p^{(j)}(k))^2$.

Let $\sigma_{\text{eff}} \triangleq \max_{k,j} \sigma_{\text{eff}}^{(j)}(k)$. Thus, all coordinate-wise confidence radii can be written with scale σ_{eff} .

B. Step 2: A Time-Uniform Coordinate-Wise Confidence Event

Let $N_t(k) = \sum_{s=1}^{t-1} \mathbb{I}\{k \in S_s\}$ be the probe count, and

$$\hat{\mu}_t^{(j)}(k) = \frac{1}{\max\{1, N_t(k)\}} \sum_{s=1}^{t-1} \mathbb{I}\{k \in S_s\} \hat{r}_s^{(j)}(k)$$

be the empirical mean of the fused observations for objective j . Fix $\delta \in (0, 1)$ and set

$$\beta_t \triangleq 2 \log\left(\frac{2Kd t^2}{\delta}\right), \quad b_t^{(j)}(k) \triangleq \sigma_{\text{eff}} \sqrt{\frac{\beta_t}{\max\{1, N_t(k)\}}}.$$

Define $\text{UCB}_t^{(j)}(k) = \hat{\mu}_t^{(j)}(k) + b_t^{(j)}(k)$, $\text{LCB}_t^{(j)}(k) = \hat{\mu}_t^{(j)}(k) - b_t^{(j)}(k)$ (with the usual clipping if needed), and stack $u_t(k) = (\text{UCB}_t^{(1)}(k), \dots, \text{UCB}_t^{(d)}(k))$.

Lemma 1 (uniform CI). Under Assumption 1 (applied to the fused noises with scale σ_{eff}), with probability at least $1 - \delta$,

$$\mathcal{E} \triangleq \{\forall t, \forall k \in [K], \forall j \in [d] : |\hat{\mu}_t^{(j)}(k) - \mu^{(j)}(k)| \leq b_t^{(j)}(k)\} \quad (62)$$

holds. Consequently, on \mathcal{E} we have the coordinate-wise sandwich $\mu(k) \preceq u_t(k)$ for all t, k .

Proof. For each fixed (k, j) , $\{\mathbb{I}\{k \in S_t\}(\hat{r}_t^{(j)}(k) - \mu^{(j)}(k))\}_{t \geq 1}$ is a martingale difference sequence w.r.t. the PtC filtration, with conditionally σ_{eff} -sub-Gaussian increments. A standard time-uniform self-normalized inequality for adaptively sampled sub-Gaussian MDS gives

$$\begin{aligned} \Pr(\exists t \leq T : |\hat{\mu}_t^{(j)}(k) - \mu^{(j)}(k)| \\ > \sigma_{\text{eff}} \sqrt{\beta_t / \max\{1, N_t(k)\}}) \leq \frac{\delta}{Kd}, \end{aligned} \quad (63)$$

and a union bound over (k, j) yields (62). \square

C. Step 3: Scalar Optimism and Lipschitz Control

On \mathcal{E} , $\mu(k) \preceq u_t(k)$ and ϕ is monotone, hence

$$f(k) = \phi(\mu(k)) \leq \phi(u_t(k)) = \text{score}_t^\phi(k). \quad (64)$$

Moreover, since ϕ is L_ϕ -Lipschitz w.r.t. $\|\cdot\|_\infty$,

$$\begin{aligned} \phi(u_t(k)) - \phi(\mu(k)) &\leq L_\phi \|u_t(k) - \mu(k)\|_\infty \\ &\leq L_\phi \max_{j \in [d]} b_t^{(j)}(k) \leq L_\phi \sum_{j=1}^d b_t^{(j)}(k), \end{aligned} \quad (65)$$

where the last inequality uses $\max_j x_j \leq \sum_j x_j$.

Similarly, on \mathcal{E} we also have $\|\hat{\mu}_t(k) - \mu(k)\|_\infty \leq \max_j b_t^{(j)}(k)$, hence

$$\phi(\mu(k)) \geq \phi(\hat{\mu}_t(k)) - L_\phi \sum_{j=1}^d b_t^{(j)}(k). \quad (66)$$

D. Step 4: One-Step Regret Bound in Terms of the Probe-Set Radii

Fix t and work on the event \mathcal{E} . Let S_t be the top- q arms by $\phi(u_t(\cdot))$, and let $k_t^* \in \arg \max_{k \in S_t} f(k)$ be the *best-in-set* arm in terms of mean utility. Because k_t maximizes $\phi(\hat{\mu}_t(\cdot))$ over S_t , we have $\phi(\hat{\mu}_t(k_t)) \geq \phi(\hat{\mu}_t(k_t^*))$, and then by (66) applied twice,

$$\begin{aligned} f(k_t) &\geq \phi(\hat{\mu}_t(k_t)) - L_\phi \sum_{j=1}^d b_t^{(j)}(k_t) \\ &\geq \phi(\hat{\mu}_t(k_t^*)) - L_\phi \sum_{j=1}^d b_t^{(j)}(k_t) \\ &\geq f(k_t^*) - L_\phi \sum_{j=1}^d b_t^{(j)}(k_t^*) - L_\phi \sum_{j=1}^d b_t^{(j)}(k_t). \end{aligned} \quad (67)$$

Therefore,

$$\begin{aligned} f^* - f(k_t) &\leq \underbrace{f^* - f(k_t^*)}_{\text{(I) set suboptimality}} + L_\phi \sum_{j=1}^d b_t^{(j)}(k_t^*) + L_\phi \sum_{j=1}^d b_t^{(j)}(k_t). \end{aligned} \quad (68)$$

We now bound the set-suboptimality term (I) using the fact that S_t consists of the top- q scalar-UCB scores. Let $U_t(k) \triangleq \phi(u_t(k))$. Since S_t contains the q largest values of $U_t(\cdot)$, its average dominates any excluded arm:

$$\frac{1}{q} \sum_{k \in S_t} U_t(k) \geq U_t(k^*). \quad (69)$$

Combining (64) and (69) gives

$$f^* = f(k^*) \leq U_t(k^*) \leq \frac{1}{q} \sum_{k \in S_t} U_t(k).$$

Also, $f(k_t^*) = \max_{k \in S_t} f(k) \geq \frac{1}{q} \sum_{k \in S_t} f(k)$. Therefore,

$$f^* - f(k_t^*) \leq \frac{1}{q} \sum_{k \in S_t} (U_t(k) - f(k)) \leq \frac{L_\phi}{q} \sum_{k \in S_t} \sum_{j=1}^d b_t^{(j)}(k), \quad (70)$$

where the last inequality uses (65).

Plugging (70) into (68), and using that $b_t^{(j)}(k_t), b_t^{(j)}(k_t^*) \leq \sum_{k \in S_t} b_t^{(j)}(k)$, we obtain the clean per-round bound

$$f^* - f(k_t) \leq \frac{C L_\phi}{q} \sum_{k \in S_t} \sum_{j=1}^d b_t^{(j)}(k) \quad (71)$$

for a universal constant C (e.g., $C = 3$ suffices from the three terms above).

Summing over $t = 1, \dots, T$ on \mathcal{E} yields

$$R_T^\phi = \sum_{t=1}^T (f^* - f(k_t)) \leq \frac{C L_\phi}{q} \sum_{t=1}^T \sum_{k \in S_t} \sum_{j=1}^d b_t^{(j)}(k). \quad (72)$$

E. Step 5: Bounding the Sum of Radii Using the Probe Budget

Fix an objective j . Using $b_t^{(j)}(k) = \sigma_{\text{eff}} \sqrt{\beta_t / \max\{1, N_t(k)\}}$ and monotonicity of β_t ,

$$\sum_{t=1}^T \sum_{k \in S_t} b_t^{(j)}(k) \leq \sigma_{\text{eff}} \sqrt{\beta_T} \sum_{t=1}^T \sum_{k \in S_t} \frac{1}{\sqrt{\max\{1, N_t(k)\}}}.$$

For each fixed arm k , every time it is probed its count increases by one, so

$$\sum_{t: k \in S_t} \frac{1}{\sqrt{\max\{1, N_t(k)\}}} \leq \sum_{n=1}^{N_{T+1}(k)} \frac{1}{\sqrt{n}} \leq 2\sqrt{N_{T+1}(k)}.$$

Thus,

$$\begin{aligned} \sum_{t=1}^T \sum_{k \in S_t} \frac{1}{\sqrt{\max\{1, N_t(k)\}}} &\leq 2 \sum_{k=1}^K \sqrt{N_{T+1}(k)} \\ &\leq 2 \sqrt{K \sum_{k=1}^K N_{T+1}(k)} = 2\sqrt{K q T}, \end{aligned} \quad (73)$$

where the second inequality is Cauchy-Schwarz, and the last equality is the PtC bookkeeping identity $\sum_k N_{T+1}(k) = qT$. Therefore, for each j ,

$$\sum_{t=1}^T \sum_{k \in S_t} b_t^{(j)}(k) \leq 2\sigma_{\text{eff}} \sqrt{\beta_T} \sqrt{K q T}. \quad (74)$$

Summing (74) over $j = 1, \dots, d$ and plugging into (72) gives, on \mathcal{E} ,

$$R_T^\phi \leq \frac{C L_\phi}{q} \cdot d \cdot 2\sigma_{\text{eff}} \sqrt{\beta_T} \sqrt{K q T} = \tilde{O}\left(L_\phi d \sigma_{\text{eff}} \sqrt{\frac{KT}{q}}\right),$$

where $\tilde{O}(\cdot)$ hides $\sqrt{\beta_T} = \text{polylog}(K, d, T, 1/\delta)$.

F. Step 6: From High-Probability to Expectation

We have shown that on \mathcal{E} (which holds with probability at least $1 - \delta$ by Lemma 1),

$$R_T^\phi \leq \tilde{O}\left(L_\phi d \sigma_{\text{eff}} \sqrt{\frac{KT}{q}}\right).$$

On the complement \mathcal{E}^c , we can use the trivial bound $R_T^\phi \leq T$ (since utilities are in $[0, 1]$ after normalization). Thus,

$$\mathbb{E}[R_T^\phi] \leq \tilde{O}\left(L_\phi d \sigma_{\text{eff}} \sqrt{\frac{KT}{q}}\right) + \delta \cdot T.$$

Choosing $\delta = T^{-2}$ (or any δ that makes δT lower order) yields

$$\mathbb{E}[R_T^\phi] = \tilde{O}\left(L_\phi d \sigma_{\text{eff}} \sqrt{\frac{KT}{q}}\right),$$

which is the desired claim.