

Buying the Oracle: Costly Calibration and Intermediated Routing in a Market of Heterogeneous LLM Agents

Christoph Pfeiffer*

June 2026

Abstract

We study task allocation in a market of LLM agents with heterogeneous capabilities and real dollar inference costs: three model tiers (Haiku, Sonnet, Opus) bid for tasks of calibrated difficulty in a reverse auction with pay-on-success contracts. We compare five information regimes: a 2×2 design — {forced calibration phase: yes/no} \times {granularity: scalar vs. level-conditional} — plus a no-information baseline. Three findings emerge in this setup. First, the welfare effect of public scalar reputation is fragile: it raises answer quality consistently but beats the baseline in only 2 of 5 seeds on requester surplus, because pay-on-success caps the damage of adverse selection and coarse information causes rationing. Second, the two information instruments are complements: a paid calibration phase without granular bookkeeping destroys surplus, granular bookkeeping without calibration is crippled by cold start, and only their combination — an intermediary that pays for structured exploration and keeps an agent \times level accuracy matrix — comes close to covering its information-acquisition cost. Third, the intermediated regime is the only robust improvement we observe and survives a pre-registered confirmatory replication on 10 fresh seeds: +14.5% gross requester surplus over scalar reputation, 8 of 10 paired markets (down from +22% exploratory, as expected under regression toward the mean), with near-full assignment and tied accuracy — the gain is allocational. The implied break-even volume for the intermediary’s calibration capex is roughly 38–40 tasks per market. We interpret the results in the tradition of intermediation as costly expertise and state all claims as conditional on this setup.

Keywords: LLM agents, mechanism design, adverse selection, reputation, intermediation, market experiments.

JEL: D82, D47, C90, L15.

1 Introduction

Earlier stages of this research program (Section 1.1) showed that public performance signals can improve allocation in markets of tool-differentiated LLM agents, but that the headline effect did not survive a randomized-identity replication — a warning that seed-level fragility is a first-order concern in LLM market experiments. The present design responds with two ideas. First, heterogeneity now comes from *model tiers* with real inference costs (Haiku, Sonnet, Opus with adaptive thinking), not from prompt or tool endowments: capability differences are intrinsic and the cost of capability is a real dollar amount deducted from each agent’s wallet. Second, instead of asking *whether* reputation helps, we decompose *which component* of an information

*Contact: christophpp@gmail.com. Code, data, run logs, and the pre-registration document for the confirmatory experiment are archived and available from the author.

regime does the work: paid exploration (a forced calibration phase) or information granularity (level-conditional rather than aggregate records).

The economic question is the classic one of markets under asymmetric information: the requester cannot observe *ex ante* which agent can solve which task. We ask:

Which information structures allow a market of cost-heterogeneous LLM agents to find the separating allocation — cheap models for easy tasks, expensive models for hard tasks — and which of them pay for themselves?

Throughout, “welfare” means *requester surplus*: realized task value minus payments (and minus fees where applicable). Where information acquisition has a real cost, we report *total surplus* net of that cost, treating intermediary fees as transfers. We make no claims about agent intentionality; agents behave *as if* responding to incentives, which is the relevant fact for mechanism design regardless of its origin in training data.

1.1 The research program in brief

This paper is the fifth stage of a program asking whether economic mechanisms make LLM agents more reliable; the earlier stages supply the design lessons embodied here, failures included.

Stage 1 (single agent, incentive framing): token rewards and penalties raised a mid-tier model’s accuracy from 57% to 76%, concentrated at the capability boundary — but a stern non-economic prompt (“a human life depends on this”) reached 79%. Incentive *framing* is one high-stakes prompt among several, not a specifically economic channel. Lesson: testing economic mechanisms requires real state — budgets, costs, consequences — not prompt decoration. The present design gives agents real dollar wallets and books every inference cost.

Stage 2 (homogeneous reputation): persistent reputation improved nothing in a homogeneous population. Reputation is a screening mechanism; with nothing to screen, it is empty. Lesson: heterogeneity is the precondition. Here, capability differences are intrinsic to the model tiers.

Stage 3 (tool heterogeneity, public signals): with tool-endowment heterogeneity and Vickrey auctions, public performance signals raised accuracy by up to 6.3 pp over no signal (oracle 97.5% > scalar 94.3% > granular public 93.3% > private 89.3% > none 88.0%). A 10-seed randomized-identity replication, however, did not reproduce the effect (+0.6 pp instead of +11.3 pp on hard tasks). Lesson: seed robustness is a primary outcome; point estimates from three seeds are anecdotes. This experience motivates the pre-registered confirmatory replication in Section 4.4.

Stage 4a (embedding reputation, real tools): a granular reputation vector over an embedding space, evaluated in a market-versus-planner comparison with real tools, showed no market advantage. The 2×2 in this paper explains why retrospectively: granularity without exploration fails at cold start.

1.2 Contribution

1. We introduce a market testbed in which LLM agents pay **real inference costs** from a dollar wallet, including endogenous reasoning costs (Opus’s thinking budget scales with task difficulty, so the agent does not know its own cost before solving — a natural source of cost uncertainty).
2. We organize the regimes in a minimal **framework** (Section 2): an information structure is a partition plus a sampling process, and three propositions deliver the misrouting/rationing, lock-in, and break-even/fee-window predictions that the experiment then measures.

3. We provide a 2×2 **identification design** separating paid exploration from information granularity, and document that they are *complements*: neither pays alone; only the combination approaches covering its cost.
4. We document that the only **seed-robust welfare improvement** is the intermediated regime (calibration + private matrix), confirmed in a **pre-registered replication** on 10 fresh seeds: +14.5% gross over scalar reputation, 8/10 paired wins. Simple scalar reputation improves quality but not, robustly, welfare.
5. We give an empirical **break-even volume** for the intermediary: the confirmatory gross gain (\$0.50/market) covers the calibration capex (\$0.64) at roughly 38–40 tasks per market.

2 A simple framework

Before the experiment, we fix ideas with a deliberately minimal model of the allocator’s problem. It is an organizing device, not an equilibrium theory: bids in our market are produced by the LLM agents themselves and we treat them as *behavioral data*, modeling only the requester (or intermediary) side. Proofs are in Appendix A.

Setup. Agents $i \in A$ face tasks with difficulty $\ell \in L$ and value V_ℓ . Agent i solves a task of difficulty ℓ with probability $p_i(\ell)$. Contracts are pay-on-success: a winning bid b_i is paid only on a verified correct answer. Given bids, the allocator assigns the task to $\arg \max_i V_\ell \hat{p}_i(\ell) - b_i$ provided the maximum is positive, where \hat{p} is the allocator’s belief under the prevailing *information structure*. An information structure has two components: (i) the **partition** on which beliefs may condition — trivial (one scalar per agent) or level-conditional (one cell per difficulty); (ii) the **sampling process** generating the data — endogenous (outcomes of past auction winners only) or exogenous (a paid calibration phase that samples every agent in every cell). The experimental 2×2 of Section 3 crosses exactly these two components.

Proposition 1 (Coarse partition: misrouting and rationing). *Under the trivial partition, $\hat{p}_i = \bar{p}_i := \mathbb{E}_\ell[p_i(\ell)]$ for all tasks (the Laplace-smoothed scalar record converges to this under the prevailing level mixture). Then (a) whenever the maximizer of $V_\ell \bar{p}_i - b_i$ differs from the maximizer of $V_\ell p_i(\ell) - b_i$, the task is misrouted; and (b) if some agent’s $p_i(\ell)$ is non-constant in ℓ , there exist (V_ℓ, b) with $\max_i V_\ell p_i(\ell) - b_i > 0 > \max_i V_\ell \bar{p}_i - b_i$, so the task goes unassigned although trade is efficient (rationing). Both are instances of Blackwell garbling [2]: the scalar record is a deterministic coarsening of the level-conditional record and hence weakly less valuable in any decision problem — here strictly.*

Proposition 2 (Endogenous sampling: lock-in). *Suppose beliefs update only on tasks an agent wins (greedy assignment on endogenous data), and hold bids fixed. Then for an open set of parameters — in particular whenever the incumbent’s true cell rate net of bids exceeds the unsampled agent’s smoothed prior net of bids, $V_\ell p_1 - b_1 > V_\ell \hat{p}_2^0 - b_2$ — the following event has positive probability: an agent whose true $p_i(\ell)$ is highest in some cell is never sampled in that cell, beliefs never correct, and the inferior assignment persists forever. This is the incomplete-learning result of [5] transplanted to the allocator’s problem; exogenous sampling restores learning by breaking the dependence of the data on past assignments.*

Proposition 3 (Intermediary viability and the fee window). *Let $\Delta w > 0$ be the per-task gain in expected requester surplus of the intermediated regime (fine partition with exogenous sampling) over the free scalar alternative, gross of fees; let K be the cost of exogenous sampling, n the number of tasks served by one calibration, and \bar{V} the expected completed value per task. Intermediation is socially viable iff $n \Delta w \geq K$. A success fee at rate φ on completed value is individually rational for the requester iff $\varphi \leq \Delta w / \bar{V}$ and covers the intermediary’s cost iff $\varphi \geq K / (n \bar{V})$; a Pareto window exists iff $n \geq K / \Delta w$, i.e., above the break-even volume.*

The experiment measures the objects this framework leaves free: the size of the partition losses (a) and (b), the empirical incidence of lock-in and its repair by calibration, and the magnitudes Δw , K , and \bar{V} that locate the break-even volume and the fee window.

3 Market Architecture

3.1 Agents: model tiers with real cost

Three agents correspond to three Claude model tiers. Each holds a dollar wallet (initial balance \$2.00) from which *all* inference costs are deducted — solving costs, and also bidding costs, so that selection overhead is explicitly on the books (the Coase constraint: the cost of using the market mechanism must not exceed its allocative gains). Bidding is cheap by design (one short call without thinking, \approx \$0.0008).

Table 1: Agent tiers and empirical accuracy by difficulty level (no-tools calibration, seed 42).

Agent	API price (\$/MTok)		Accuracy by level			
	in	out	L1–L2	L3–L4	L5	L6
Haiku	1	5	100%	60–80%	20–50%	25%
Sonnet	3	15	100%	100%	75–80%	50%
Opus (adaptive thinking)	5	25	100%	100%	100%	100%

The resulting cost–capability gradient defines three *oracle zones*: L1–L2 belong to Haiku (\approx \$0.002/task), L3–L4 to Sonnet (\approx \$0.009), L5–L6 to Opus (\$0.06 up to \$0.24 — Opus’s thinking tokens scale with difficulty, making its cost endogenous and unknown to the agent ex ante). Tasks are solved *without* tools: tools were the Stage-3 treatment and saturate the task set (all tiers \geq 98%), which would confound the tier effect.

3.2 Tasks and values

Tasks are procedurally generated mathematical problems on six calibrated difficulty levels (L1–L6); ground truth is computed numerically (no LLM-as-judge). Levels L5/L6 (chained modular arithmetic, $5 \times 5/6 \times 6$ determinants, surjection counts, coin-change DP) were added and calibrated so that each tier has a distinct competence frontier (Table 1). Each market run uses 5 tasks per level (30 tasks), shuffled, with deterministic seeding; all conditions within a seed face identical tasks, enabling paired comparisons. Requester values are $(V_{L1}, \dots, V_{L6}) = (0.02, 0.03, 0.06, 0.10, 0.30, 0.60)$ dollars.

3.3 Mechanism

Each task is allocated by reverse auction with **pay-on-success**: all agents may bid or decline (bid cost always deducted); the requester selects one winner; the winner solves the task and is paid its bid price only if the verified answer is correct. Inference costs are deducted regardless of outcome.

3.4 Information regimes

Table 2: The five regimes. The 2×2 varies forced calibration and granularity; `no_rep` is the baseline.

Regime	Calibration	Selection signal	Bidders see	Fee
<code>no_rep</code>	—	lowest bid (reserve V)	nothing	—
<code>scalar_rep</code>	—	$V \cdot \widehat{acc} - b$, scalar	public records	—
<code>matrix_rep</code>	—	$V \cdot \widehat{acc}_{a,\ell} - b$	nothing	—
<code>calib_scalar</code>	paid	$V \cdot \widehat{acc} - b$, scalar	public records	10% of V
<code>broker_rep</code>	paid	$V \cdot \widehat{acc}_{a,\ell} - b$	nothing	10% of V

Table 2 summarizes the five regimes. All accuracy estimates are Laplace-smoothed. In the *calibration phase* (regimes `calib_scalar`, `broker_rep`), every agent solves two held-out tasks per level before the market opens; the calibrating party reimburses inference at cost (agents break even) and books the outlay as capex (mean \$0.69 for `calib_scalar` on the five common seeds, \$0.72 for `broker_rep` on the eight exploratory seeds, \$0.64 in the confirmatory run; dominated by Opus’s thinking on L5/L6). In both calibration regimes, the party that charges the fee funds the calibration. In `broker_rep` — the intermediated regime — the resulting agent×level matrix is private to the intermediary, which routes by level-conditional expected utility and charges the requester 10% of V on success. `matrix_rep` is the same bookkeeping without calibration (free, requester-side, cold start); `calib_scalar` is calibration whose outcomes feed only an aggregate scalar. Bidder-visible information varies by regime: public regimes display success records to bidders; matrix regimes do not, since the information is requester- or intermediary-private by construction.

Seeds: all five regimes ran on seeds {42, 123, 999, 7, 2025}; the central pairing `broker_rep` vs. `scalar_rep` additionally ran on {11, 77, 555} (8 paired seeds, 240 tasks per regime).

4 Results

4.1 Scalar reputation: robust quality effect, fragile welfare effect

Without information, the market exhibits the textbook adverse-selection pattern [1]: the mid-tier agent underbids on hard tasks, wins them, and fails (overreach: 53 assignments below the oracle zone, 17 of them failed, vs. 24–33 with information; accuracy 86% vs. 87–93%, the low end being the cold-start matrix regime). Scalar reputation removes much of this quality damage. However, the *welfare* effect is fragile: across 5 seeds the paired difference `scalar_rep`–`no_rep` averages only +\$0.13 and is positive in 2 of 5 seeds (range -0.79 to $+2.20$). Two forces offset the quality gain: cheap agents that fail are simply not paid (pay-on-success caps the damage), and coarse information causes *rationing* — under `scalar_rep`, 20 of 240 tasks went unassigned because no agent cleared the expected-utility threshold, forfeiting their entire surplus (the rationing of Proposition 1b). A separate 3-seed pilot run had shown +18.6% with 3 of 3 paired wins; the effect survived neither re-running those seeds (LLM sampling noise at fixed seed) nor two further seeds. We report this explicitly: it repeats the Stage-3 replication experience and motivates treating seed-level robustness as a primary outcome.

4.2 The 2×2 : exploration and granularity are complements

Table 3: Mean *total* surplus per market (gross requester surplus minus calibration capex; fees are transfers), 5 common seeds, 30 tasks each. Baseline `no_rep`: \$3.22.

	Scalar signal	Level-conditional matrix
No calibration	3.34	3.31
Paid calibration	2.98	3.37

Table 3 reports mean total surplus across the five common seeds; the off-diagonal cells identify the components. *Exploration without granularity* (`calib_scalar`) raises gross surplus (3.67 vs. 3.34) — the warm-up removes rationing entirely (150/150 assigned) — but a scalar cannot exploit the level structure the calibration paid to reveal (the garbling of Proposition 1), so the capex exceeds the gain and total surplus falls below every other cell. *Granularity without exploration* (`matrix_rep`) is statistically indistinguishable from scalar reputation: the matrix starts empty, Laplace smoothing prices every agent at $\widehat{acc} = 0.5$, and early hard tasks are rationed or misrouted (cold start — the lock-in of Proposition 2). Only the combination approaches covering its information cost in total terms — it just covers it on the five common exploratory seeds, while the confirmatory run places break-even slightly above this volume (Section 4.5). The complementarity, not either instrument alone, is the result.

4.3 The intermediated regime is the only seed-robust welfare improvement

On the 8 exploratory paired seeds, `broker_rep` exceeds `scalar_rep` in gross requester surplus by \$0.72 per market (+22%; \$3.96 vs. \$3.24), winning **7 of 8 paired seeds** (deltas +0.05, -0.68, +1.68, +1.15, +1.40, +1.00, +0.57, +0.61); Section 4.4 reports the pre-registered confirmation on fresh seeds. Three mechanics carry the effect:

1. **No rationing.** The intermediary assigns 240/240 tasks (scalar: 220/240). Level-conditional estimates clear the expected-utility threshold where a blended scalar does not.
2. **Routing, not accuracy.** Over 8 seeds, accuracy is identical (92% vs. 92%): the gain comes from assigning the forfeited tasks and from better price-quality matches (L6 routed to Opus 14/15 vs. 10/15 on the exploratory seeds), not from fewer wrong answers.
3. **Variance reduction.** Cross-seed gross surplus ranges 3.39–4.41 under the intermediary vs. 2.39–4.44 under scalar reputation; the worst seed improves by \$1.00. The calibrated matrix removes the bad tail that information-poor regimes produce. In this sense the intermediary sells *insurance*, not only routing.

The matrix also takes calculated risks: it routed 15 L5 tasks to Sonnet ($\approx 78\%$ accurate there, far cheaper than Opus), which is expected-utility rational and produced only 7 overreach failures across 150 tasks — informed risk-taking over the competence frontier rather than defensive over-provisioning.

4.4 Pre-registered confirmatory replication

Because the exploratory seeds were extended sequentially after results were visible ($3 \rightarrow 5 \rightarrow 8$ — a “garden of forking paths” [6]), we ran a confirmatory replication with 10 fresh seeds (101–1010), declared together with the confirmation criterion *before* the run (criterion: $\geq 7/10$ paired wins and positive mean delta; the full pre-registration is archived with the results). Table 4 lists the per-seed results.

Table 4: Confirmatory run: gross requester surplus per seed (30 tasks each).

Seed	101	202	303	404	505	606	707	808	909	1010
scalar_rep	3.87	2.52	3.48	4.09	2.74	4.24	2.33	4.00	3.70	3.52
broker_rep	3.49	4.28	4.20	4.40	2.60	4.36	4.17	4.12	4.31	3.59

The criterion is met: 8 of 10 paired wins, mean delta +\$0.50 (+14.5%; \$3.95 vs. \$3.45) (Figure 1). As inferential statistics on the paired deltas: a one-sided sign test gives $p = 0.055$ and a one-sided paired t -test gives $t(9) = 2.09$, $p = 0.033$ — moderate evidence, consistent with the pre-registered criterion being met. The effect size shrinks from the exploratory +22%, as expected under regression toward the mean; we report the confirmatory number. The mechanics replicate: accuracy is tied (92% vs. 91%), while scalar reputation ratios 16/300 tasks against the intermediary’s 2/300 — the welfare gain is allocational (Figure 2). The worst-case seed under each regime remains better for the intermediary (\$2.60 vs. \$2.33), though the variance compression is weaker than in the exploratory sample.

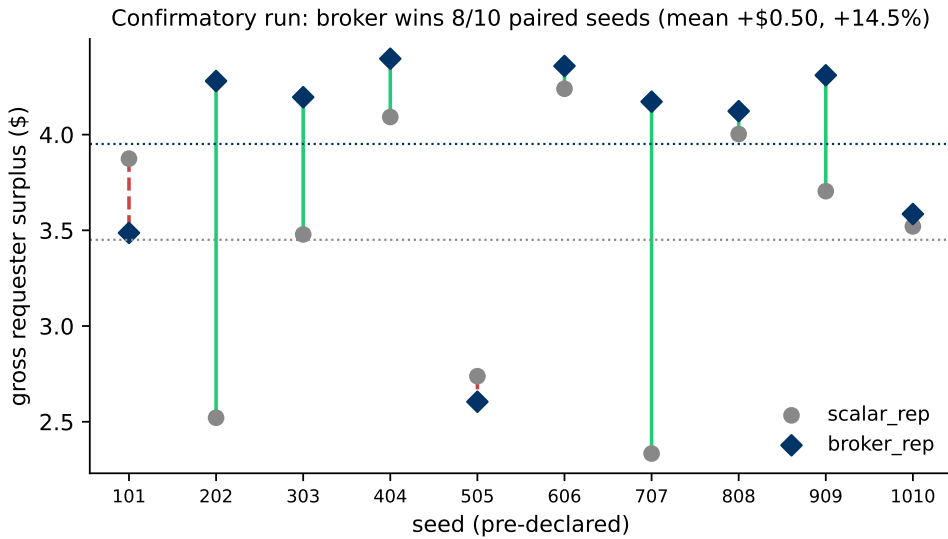


Figure 1: Paired gross requester surplus on the 10 pre-declared confirmatory seeds. Solid (green) connectors: intermediary above scalar reputation (8 of 10); dashed (red): below.

Task allocation shares, confirmatory run (oracle zones: L1–L2 haiku, L3–L4 sonnet, L5–L6 opus)

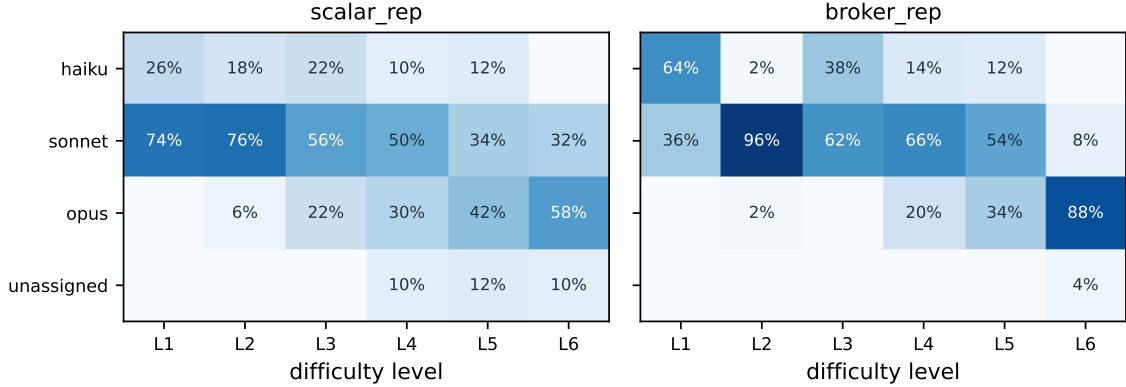


Figure 2: Allocation shares by difficulty level (confirmatory run, 300 tasks per regime). The intermediary routes L5–L6 toward Opus and nearly eliminates unassigned tasks; under scalar reputation, rationing concentrates where coarse information meets high-value tasks.

4.5 Break-even volume and fee economics

In total-surplus terms the exploratory 8-seed comparison landed at parity (both \$3.24, with a mean broker capex of \$0.72 on those seeds); the confirmatory run sharpens this to a slightly higher break-even: the gross gain of \$0.50 per 30-task market falls short of the mean calibration capex of \$0.64, implying break-even at roughly **38–40 tasks per market** (assuming the per-task gain extrapolates linearly beyond the 30-task markets we ran). Calibration is a fixed cost: at larger volume the same matrix serves more tasks, so total surplus turns positive while the allocational gains persist. On the fee side, the 10% success fee yields $\approx \$0.47$ per market against \$0.64 capex; the broker’s break-even fee at 30 tasks is $\approx 13.5\%$ of task value and falls roughly inversely with volume ($\approx 7\%$ at 60 tasks). The requester-side ceiling, by contrast, is volume-independent: the gross gain is $\approx 10.6\%$ of completed task value, so any fee above that leaves the requester worse off than under free scalar reputation. A Pareto window in which intermediation is viable for both sides therefore opens only above the break-even volume — at 60 tasks, fees between roughly 7% and 10% finance the calibration while leaving the requester strictly better off. This is the window of Proposition 3 with the measured magnitudes $\Delta w \approx \$0.017$ per task, $K \approx \$0.64$, and $\bar{V} \approx \$0.16$ of completed value per task.

5 Discussion

5.1 Relation to intermediation theory

The setup instantiates middlemen-as-experts [3]: an intermediary invests in quality-detection expertise that individual market participants would not, because the investment is fixed-cost and its returns accrue across many transactions. The calibration phase makes the usually abstract “expertise” concrete and priced: it is the purchase of an oracle matrix at measurable dollar cost, and our 2×2 shows the purchase only pays if the buyer also has the bookkeeping structure to exploit it. The certification-intermediary results of Lizzeri [4] — intermediaries may reveal coarse information and still capture surplus — map onto our fee analysis. With a single requester, public reputation and a requester-private matrix coincide informationally; our design therefore identifies the value of *granularity plus exploration*, not of privacy. Separating the privacy component requires multiple requesters (free-riding on a public calibration good vs. an intermediary’s private one), which is the natural next stage.

5.2 Limitations

Three agents, one requester, mathematical tasks with verifiable ground truth, list-price inference costs, and 5–10 seeds per comparison. The fragility we document for scalar reputation across 5 seeds counsels equal caution about our own positive result; the pre-registered confirmatory replication (8/10 paired wins on fresh seeds) addresses the seed-adaptivity concern for the headline claim, and we report the shrunken confirmatory effect size (+14.5%) rather than the exploratory one. Bidding behavior is imperfectly calibrated (the mid-tier model persistently underbids the small one on easy tasks — a pattern that originates in bidding behavior, not in the information regime). All “as-if” behavioral statements are claims about reproducible input–output patterns of current models, not about strategic cognition.

6 Conclusion

In this setup, the question “does reputation help?” has an unhelpfully fragile answer, but the decomposed question — *which information structure pays for itself?* — has a robust one: paid, structured exploration combined with granular bookkeeping, i.e., an intermediary that buys the oracle and routes with it. It is the regime that assigns essentially every task, the only one that beats its alternative under pre-registration (8 of 10 fresh paired seeds, +14.5% gross), and the only one whose information cost its own gains can amortize at modest volume (≈ 40 tasks) — while improving the worst-case outcomes that make information-poor agent markets unreliable.

A Proofs

Proof of Proposition 1. (a) is immediate: the allocator maximizes $V_\ell \bar{p}_i - b_i$, so whenever its argmax differs from that of $V_\ell p_i(\ell) - b_i$, the realized expected surplus is strictly below the attainable one. For (b), two levels $\{e, h\}$ in equal proportion and a single agent with $p(e) = 1$, $p(h) = 0$ (so $\bar{p} = \frac{1}{2}$) suffice: for an easy task, any bid b with $V_e \bar{p} < b < V_e$ gives true surplus $V_e - b > 0$ but believed surplus $V_e/2 - b < 0$, so the allocator declines a profitable trade. The scalar record is a deterministic coarsening of the level-conditional record, i.e., a garbling in the sense of [2], which implies weakly lower decision value in any decision problem and strictly lower here. \square

Proof of Proposition 2. Consider two agents and one cell, with true success probabilities $p_1 < p_2$ in that cell, and Laplace-smoothed beliefs updated only on assigned tasks. If early draws make \hat{p}_1 exceed \hat{p}_2 's prior by more than $(b_1 - b_2)/V_\ell$, agent 1 wins the cell; agent 2's belief is never updated and remains at its prior forever, so the event that agent 1 retains every future assignment has positive probability (its belief converges to p_1 , which can remain above agent 2's smoothed prior net of bids). This is the two-armed-bandit incomplete-learning logic of [5]: optimal (here: greedy) experimentation on endogenous data settles on the inferior arm with positive probability. Under exogenous sampling, every agent is sampled in every cell at a rate independent of assignments, so beliefs converge to the truth by the law of large numbers and the lock-in event has probability zero. \square

Proof of Proposition 3. Social viability: the fine partition raises expected requester surplus by Δw per task over n tasks at one-time cost K ; the net gain $n\Delta w - K$ is nonnegative iff $n \geq K/\Delta w$. The requester accepts the fee iff her per-task net gain $\Delta w - \varphi \bar{V}$ is nonnegative, i.e., $\varphi \leq \Delta w/\bar{V}$. The intermediary covers cost iff $\varphi n \bar{V} \geq K$, i.e., $\varphi \geq K/(n\bar{V})$. The interval $[K/(n\bar{V}), \Delta w/\bar{V}]$ is nonempty iff $n \geq K/\Delta w$. \square

References

- [1] Akerlof, G. (1970). The Market for “Lemons”: Quality Uncertainty and the Market Mechanism. *Quarterly Journal of Economics*, 84(3), 488–500.
- [2] Blackwell, D. (1953). Equivalent Comparisons of Experiments. *Annals of Mathematical Statistics*, 24(2), 265–272.
- [3] Biglaiser, G. (1993). Middlemen as Experts. *RAND Journal of Economics*, 24(2), 212–223.
- [4] Lizzeri, A. (1999). Information Revelation and Certification Intermediaries. *RAND Journal of Economics*, 30(2), 214–231.
- [5] Rothschild, M. (1974). A Two-Armed Bandit Theory of Market Pricing. *Journal of Economic Theory*, 9(2), 185–202.
- [6] Gelman, A., & Loken, E. (2013). The Garden of Forking Paths: Why Multiple Comparisons Can Be a Problem, Even When There Is No “Fishing Expedition” or “p-Hacking”. Working paper, Columbia University.