

---

# The Partial Testimony of Logs: Evaluation of Language Model Generation under Confounded Model Choice

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Offline evaluation of language models from usage logs is biased when model  
2 choice is confounded: the same user-side factors that drive which model is used  
3 also shape how its output is judged. Randomization removes this bias, yet produc-  
4 tion randomization remains scarce. We study a three-source design that combines  
5 a large confounded observational log (OBS) for scale, a small randomized exper-  
6 iment (EXP) for unconfounded scoring, and an offline simulator (SIM) that  
7 replays candidate models on cached contexts. An identification theorem shows  
8 that EXP together with SIM recovers causal model values, with OBS contributing  
9 afterward to reduce estimation variance. Six estimator families are evaluated in  
10 a controlled semi-synthetic validation and in two real-task cached benchmarks  
11 for summarization and coding. The winning family shifts with the EXP budget  
12 and with how closely the OBS signal predicts the target reward, giving a practical  
13 recipe for choosing among hybrid OBS/EXP estimators.

## 14 1 Introduction

15 Offline evaluation of language models increasingly relies on deployment logs, preference data,  
16 and judge annotations collected in the wild [Zheng et al., 2024, Zhao et al., 2024, Zheng et al.,  
17 2023, Chiang et al., 2024]. Once model choice is user-driven, those logs are confounded: the  
18 same unobserved factors that influence which model is used can also influence how its output is  
19 judged [Kallus et al., 2018, Rosenman et al., 2023, Cheng and Cai, 2021]. Raw comparisons of  
20 logged outcomes therefore mix self-selected subpopulations and cannot recover a common target. In  
21 summarization, a reader who prefers concise prose may select a model known for that style and rate  
22 it more favorably; in coding assistance, repository difficulty or developer familiarity may affect both  
23 which assistant is invoked and whether the resulting patch succeeds.

24 A randomized experiment breaks this link by overriding model choice. In practice, such experiments  
25 are costly and disruptive, so they remain small. Prior work studies how to combine a large confounded  
26 observational sample with a small randomized one [Kallus et al., 2018, Rosenman et al., 2023, Cheng  
27 and Cai, 2021, Lin et al., 2025, Yang et al., 2025, Colnet et al., 2024].

28 Causal evaluation under confounded model choice has two parts. *Outcome generation* asks which  
29 output a model would have produced on a given context; in many deployments an offline simulator  
30 (SIM) handles this by replaying models on cached contexts. *Outcome scoring* asks how the realized  
31 output should be evaluated without inheriting the bias in logged model choice. Separating the two  
32 motivates a three-source design: a large observational log (OBS) with self-selected model choices,  
33 a small randomized experiment (EXP) with unconfounded outcome labels, and SIM outputs under  
34 alternative model choices. The goal is to estimate how each generative model would perform when

35 logged model choice is confounded, randomized outcome labels are limited, and offline replay can  
 36 regenerate candidate outputs.

37 **Contributions.** The paper makes three contributions.

38 **(C1) Identification (Theorem 1).** Under the structural assumptions stated in Section 2, the simulator  
 39 and the randomized experiment alone identify causal model values, even when the observational  
 40 log is fully confounded.

41 **(C2) Estimator role of OBS (post-identification only).** Identification and estimation are separated,  
 42 and six post-identification estimator families are compared based on how they exploit OBS. A  
 43 theory-to-evidence map links each family to the data regime in which it helps.

44 **(C3) Empirical evaluation.** A controlled semi-synthetic validation and two real-task cached bench-  
 45 marks (summarization, coding) compare the families on held-out recommendation regret.  
 46 Two practical drivers explain the winning family: the size of the EXP budget, and how well  
 47 OBS-derived structure matches the target reward.

48 Section 2 defines the causal graph, data sources, target estimands, and assumptions. Section 3  
 49 states and proves the SIM-plus-EXP identification result. Section 4 introduces the post-identification  
 50 estimator families and the theory-to-evidence map. Section 5 describes the controlled validation and  
 51 cached benchmarks; Sections 5.4 and 5.5 report the cached summarization and cached coding probes.  
 52 Sections 6, 7, and 8 cover related work, scope and limitations, and conclusions.

## 53 2 Problem setting

54 OBS provides scale and auxiliary supervision from self-selected usage, EXP provides unbiased  
 55 outcome labels, and SIM provides counterfactual outputs. Figure 1 shows the causal structure;  
 56 Appendix A.1 gives the full three-source pipeline.

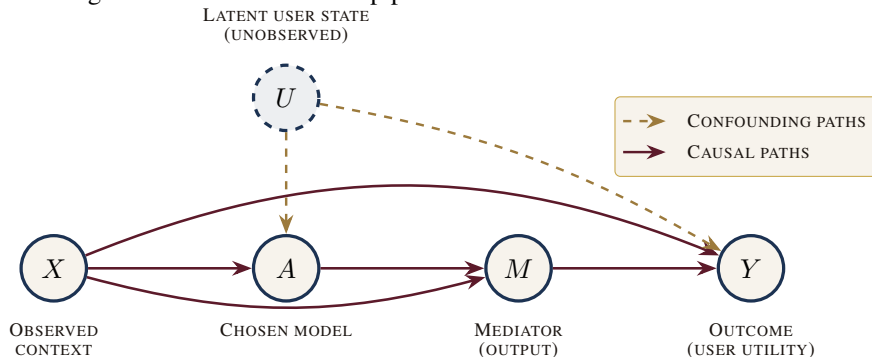


Figure 1: Causal graph with confounded model choice. The latent state  $U$  affects both model choice  $A$  and outcome  $Y$ , so OBS comparisons are biased. EXP randomization breaks the  $U \rightarrow A$  link, and no latent mediator confounding blocks  $U \rightarrow M$  once  $(X, A)$  is fixed.

57 We study  $(X, U, A, M, Y)$ , where  $X$  is the observed context,  $U$  is an unobserved user state,  $A \in \mathcal{A}$   
 58 indexes the chosen generative model,  $M$  is its output, and  $Y \in [0, 1]$  is the outcome. The graph adds  
 59 two restrictions to ordinary confounding: conditional on  $(X, A)$ ,  $U$  does not directly affect  $M$ , and  
 60  $A$  affects  $Y$  only through  $M$ .

61 **Three sources.** OBS is a large logged sample  $D_{\text{OBS}} = \{(X_i, A_i, M_i, Y_i, Z_i)\}_{i=1}^{n_{\text{OBS}}}$  in which  
 62  $A_i$  may depend on  $U_i$  and auxiliary labels  $Z_i$  can train proxy representations. EXP is a smaller  
 63 randomized sample  $D_{\text{EXP}} = \{(X_j, A_j, M_j, Y_j)\}_{j=1}^{n_{\text{EXP}}}$  in which  $A$  is randomly assigned over the  
 64 experimental action set  $\mathcal{A}_{\text{EXP}} \subseteq \mathcal{A}$ , so  $A \perp U \mid X$ . SIM reruns model  $a$  on context  $x$  and samples  
 65  $M \sim p_{\text{sim}}(\cdot \mid x, a)$ . EXP therefore supplies causal scores, SIM supplies counterfactual outputs, and  
 66 OBS contributes auxiliary supervision used only after identification (Figure 2).

67 **Targets.** The marginal and context-conditional causal values are

$$\mu(a) := \mathbb{E}[Y(\text{do}(A = a))], \quad q(x, a) := \mathbb{E}[Y(\text{do}(A = a)) \mid X = x].$$

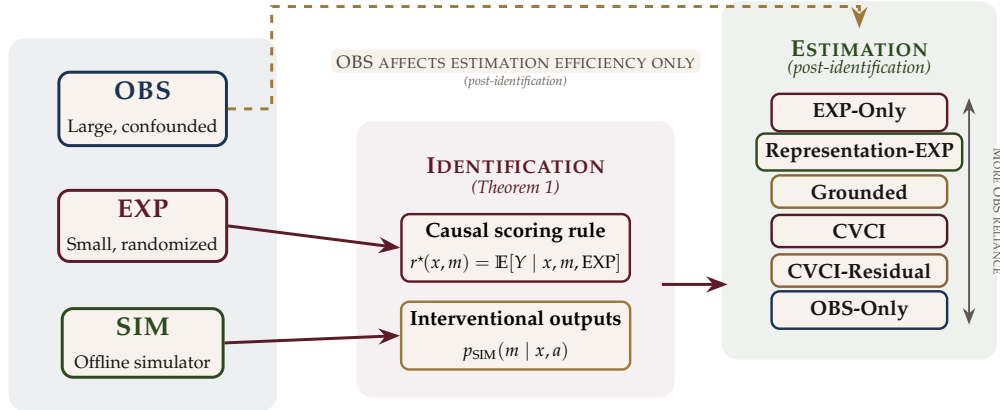


Figure 2: Three-source design. SIM and EXP are sufficient for identification of  $r^*$ ,  $q$ , and  $\mu$  on shared support; OBS affects estimation efficiency only. Solid arrows mark the identification flow; dashed gold arrows mark post-identification statistical signal.

68 For finite cached benchmarks, write  $r^*(x, m) := \mathbb{E}[Y \mid X = x, M = m, \text{EXP}]$  for the EXP-side  
69 outcome regression and  $m_{\text{cache}}(x, a)$  for the stored output associated with context  $x$  and model  $a$ .  
70 The cached evaluation target is

$$q_{\text{cache}}(x, a) := r^*(x, m_{\text{cache}}(x, a)),$$

71 which replaces the SIM mediator average with a single cached draw and is therefore a property of the  
72 realized cache.

73 **Assumptions.** Identification uses five design assumptions, each isolating one source of bias.

74 **Assumption A1** (Randomization in EXP). In EXP,  $A$  is randomized over  $\mathcal{A}_{\text{EXP}}$  so that  $A \perp U$  and  
75  $A \perp U \mid X$ . Without this, no source identifies the causal scoring rule.

76 **Assumption A2** (SIM validity and shared support). The simulator matches the interventional mediator  
77 law,  $p_{\text{sim}}(m \mid x, a) = \mathbb{P}(M = m \mid X = x, \text{do}(A = a))$ , with  $\text{supp } p_{\text{sim}}(\cdot \mid x, a) \subseteq \text{supp } \mathbb{P}(M \mid$   
78  $X = x, A = a, \text{EXP})$  for every  $x$  in the EXP context support and every  $a \in \mathcal{A}_{\text{EXP}}$ .

79 **Assumption A3** (Outcome consistency under intervention). The conditional outcome law given  
80  $(X, M, U)$  is invariant across OBS, EXP, and intervention:  $\mathbb{P}(Y \mid X, M, U, \text{OBS}) = \mathbb{P}(Y \mid$   
81  $X, M, U, \text{EXP}) = \mathbb{P}(Y \mid X, M, U, \text{do}(A = a))$  on the relevant support.

82 **Assumption A4** (No latent mediator confounding; single-round generation).  $U \perp M \mid X, A$ , so  
83 the generative model has no edge  $U \rightarrow M$  once  $(X, A)$  is fixed. This restricts the framework to  
84 single-round, stateless generation.

85 **Assumption A5** (Context distribution alignment). The marginal target distribution coincides with  
86 EXP,  $P_X^{\text{tgt}} = P_X^{\text{EXP}}$ . When this fails, only  $q(x, a)$  on the shared support is identified; reweighting  
87 requires standard transportability arguments [Cole and Stuart, 2010, Pearl and Bareinboim, 2011,  
88 Bareinboim and Pearl, 2016].

89 Throughout the paper, “A1”–“A5” refer to Assumptions A1–A5; Appendix A.1 expands their motiva-  
90 tion, support conditions, and transportability caveat.

### 91 3 Identification

92 EXP provides the unbiased scoring rule  $r^*(x, m) = \mathbb{E}[Y \mid X = x, M = m, \text{EXP}]$  defined in  
93 Section 2, and SIM provides the interventional mediator distribution  $p_{\text{sim}}(\cdot \mid x, a)$ ; together they  
94 identify causal value. Throughout, the analysis is restricted to actions in the experimental action set  
95  $\mathcal{A}_{\text{EXP}}$  and to mediator values on the shared support of the EXP law of  $M \mid X = x, A = a$  and the  
96 SIM law  $p_{\text{sim}}(\cdot \mid x, a)$ .

97 **Theorem 1** (SIM plus EXP identifies causal value). *Under the causal graph in Figure 1 and*  
 98 *Assumptions A1–A5, for every  $x$  in the EXP context support and every  $a \in \mathcal{A}_{\text{EXP}}$ ,*

$$q(x, a) = \mathbb{E}[Y(\text{do}(A = a)) \mid X = x] = \mathbb{E}_{M \sim p_{\text{sim}}(\cdot \mid x, a)}[r^*(x, M)], \quad (1)$$

99 *and therefore, with  $\mu(a)$  defined under the target context distribution  $P_X^{\text{tgt}} = P_X^{\text{EXP}}$  of Assumption*  
 100 *A5,*

$$\mu(a) = \mathbb{E}_{X \sim P_X^{\text{tgt}}}[q(X, a)] = \mathbb{E}_{X \sim P_X^{\text{tgt}}}[\mathbb{E}_{M \sim p_{\text{sim}}(\cdot \mid X, a)}[r^*(X, M)]] . \quad (2)$$

101 *Hence both  $q(x, a)$  and  $\mu(a)$  are identified from EXP and SIM on the shared EXP/SIM support;*  
 102 *identification places no unconfoundedness requirement on OBS.*

103 **Proof sketch.** The proof uses each assumption exactly once. A4 strips  $M$  of any residual infor-  
 104 mation about  $U$  given  $(X, A)$ ; A1 then removes the remaining  $U \rightarrow A$  path inside EXP, so the  
 105 conditional law of  $Y$  given  $(X, M)$  in EXP matches the law under  $\text{do}(A = a)$  on the shared support;  
 106 A3 carries this equality from EXP to the post-intervention world, identifying the EXP regression  
 107  $r^*(x, m)$  as the causal scoring rule; A2 plugs in the SIM mediator law to obtain  $q(x, a)$ ; and A5  
 108 averages over the target context distribution to recover  $\mu(a)$ . The full proof is in Appendix C.1;  
 109 Appendix A.2 explains why conditioning on  $(X, M)$  in OBS still leaves residual  $U$ -confounding  
 110 once A1 is dropped.

111 **From identification to the experimental design.** The controlled validation probes the theorem-  
 112 level target with a known latent reward generator. The two real cached benchmarks report regret  
 113 against  $q_{\text{cache}}$  with one cached output per context–action pair, diagnosing the post-identification  
 114 estimator families. Across all settings, Section 4 uses OBS only after identification, to reduce variance  
 115 or improve function approximation.

## 116 4 Estimators

117 After SIM and EXP identify the target, the remaining question is how OBS can reduce estimation  
 118 error. The six estimator families differ in how they use OBS after identification: some ignore OBS  
 119 outcomes, some use OBS only to learn a proxy representation, and some use OBS outcomes through  
 120 a baseline or pooled fit. Proxy-learning details, value aggregation, and the theory-to-evidence map  
 121 are expanded in Appendices A.3, A.6, and A.4.

122 **Notation.** Throughout,  $z = (x, m)$  denotes a context–mediator pair,  $\varphi(z) \in \mathbb{R}^d$  a fixed high-  
 123 dimensional text-feature map, and  $\psi(z) \in \mathbb{R}^k$  ( $k \ll d$ ) a lower-dimensional proxy representation  
 124 learned from OBS auxiliary labels in the spirit of embeddings adapted for causal adjustment [Veitch  
 125 et al., 2020]. Write  $\mathcal{F}_\varphi = \{r(z) = \text{clip}_{[0,1]}(w^\top \varphi(z) + b) : (w, b) \in \mathbb{R}^{d+1}\}$  and  $\mathcal{F}_\psi$  analogously  
 126 on  $\psi$  for the clipped-linear function classes used as the reward-model hypothesis spaces; the proxy  
 127 changes estimation only. Two families use OBS outcomes as well:  $f_{\text{OBS}}$  is a reward predictor pre-fit  
 128 by ridge on  $D_{\text{OBS}}$ , used either as a baseline that EXP corrects (Grounded) or as part of a pooled  
 129 OBS+EXP loss (CVCI variants).

Table 1: How the six reward estimators use OBS and EXP. Darker shading indicates heavier reliance on observational signal.

Estimator	OBS role	EXP role	Bias–variance intuition
EXP-Only	None	Fits the reward model	No confounding bias, but high variance when EXP is small.
OBS-Only	Fits the reward model on logged outcomes	None	Low variance, but inherits bias from confounded model choice.
Representation-EXP	Learns the proxy representation $\psi$ from auxiliary labels	Fits the reward head on $\psi$	Effective when the target varies along directions retained by $\psi$ .
Grounded	Provides the baseline predictor $f_{\text{OBS}}$ and the proxy space $\psi$	Estimates and tunes a correction to $f_{\text{OBS}}$	Uses OBS for efficiency and EXP to remove low-dimensional bias.
CVCI	Contributes logged-outcome loss directly	Chooses the pooling weight by causal cross-validation and anchors the fit to randomized outcomes	Trades OBS bias against finite-EXP estimation error through direct pooling.
CVCI-Residual	Provides the baseline and contributes pooled residual loss in $\psi$	Chooses the pooling and shrinkage of the residual correction	Uses OBS for the coarse fit and EXP for the residual mismatch.

Table 2: Theory-to-evidence map for the six estimator families. Tier 1: formal main-paper support. Tier 2: formal appendix support. Tier 3: empirical diagnostic.

Mechanism	Tier	Formal support	Empirical reading
SIM + EXP identify the target	1	Theorem 1	All benchmarks separate scoring from replay/cache
<b>Grounded</b> enlarges the proxy correction class	2	Appendix Theorem on class expansion	Basis-expanded <b>Grounded</b> improves over single-linear in summarization
Oracle vs. finite-sample correction	2	Appendix <b>Grounded-vs-OBS</b> theorem	<b>Grounded</b> competitive but not uniformly best
<b>CVCI-Residual</b> residual simplicity	2	Appendix residual-vs-pooling theorem	Wins selected cells, can underperform on real summarization
Proxy/reward alignment for Representation-EXP	3	— (empirical)	Coding fix-success crossover near $\alpha_{\text{fix}} \approx 0.25$

130 **Reward-fit objectives.** All six reward models minimize a penalized squared loss with predictions  
 131 clipped to  $[0, 1]$ . With shorthand  $L_S(r) := \sum_{(x,m,y) \in D_S} (r(x, m) - y)^2$  for  $S \in \{\text{OBS}, \text{EXP}\}$  and  
 132 ridge penalty  $\Omega(\cdot)$ , the families differ in the data and feature space they use:

$$\hat{r}_{\text{EXP}} = \arg \min_{r \in \mathcal{F}_\varphi} L_{\text{EXP}}(r) + \Omega(r), \quad (3)$$

$$\hat{r}_{\text{OBS}} = \arg \min_{r \in \mathcal{F}_\varphi} L_{\text{OBS}}(r) + \Omega(r), \quad (4)$$

$$\hat{r}_{\text{Rep}} = \arg \min_{r \in \mathcal{F}_\psi} L_{\text{EXP}}(r) + \Omega(r), \quad (5)$$

$$\hat{r}_{\text{G}}(z) = \text{clip}_{[0,1]} \left[ f_{\text{OBS}}(z) + \hat{\theta}^\top \psi(z) + \hat{b} \right], \quad (6)$$

$$\hat{r}_{\text{CVCI}}(\lambda) = \arg \min_{r \in \mathcal{F}_\varphi} (1 - \lambda)L_{\text{OBS}}(r) + \lambda L_{\text{EXP}}(r) + \Omega(r), \quad (7)$$

$$\hat{r}_{\text{CVCIRes}}(\lambda) = f_{\text{OBS}} + \arg \min_{g \in \mathcal{F}_\psi} (1 - \lambda)L_{\text{OBS}}^{\text{res}}(g) + \lambda L_{\text{EXP}}^{\text{res}}(g) + \Omega(g), \quad (8)$$

133 where  $L_S^{\text{res}}(g) := \sum_{(x,m,y) \in D_S} (g(x, m) - (y - f_{\text{OBS}}(x, m)))^2$  is the residual loss against the  
 134 OBS anchor, the **Grounded** correction  $(\hat{\theta}, \hat{b})$  in (6) is fit on EXP residuals  $Y - f_{\text{OBS}}(z)$ , and the  
 135 CVCI pooling weight  $\hat{\lambda} \in [0, 1]$  in (7)–(8) is chosen by EXP cross-validation. The unbiased–biased  
 136 endpoints in (3)–(4) make EXP-Only the variance-bound baseline and OBS-Only the bias-bound  
 137 endpoint; the four hybrids interpolate between them through three orthogonal levers (proxy features  
 138  $\psi$ , an OBS anchor  $f_{\text{OBS}}$ , and the pooling weight  $\lambda$ ).

139 **Theory-to-evidence map.** Table 2 groups the six families by the kind of evidence that supports  
 140 each comparison. Tier-1 statements (identification) are formal and main-paper; Tier-2 statements  
 141 (oracle correction, residual simplicity, function-class enlargement) are formal but appendix-only;  
 142 Tier-3 patterns (proxy alignment, reward geometry) are empirical diagnostics. The benchmarks in  
 143 Sections 5–5.5 probe Tier 3, and Appendix A.4 expands the map.

144 **From reward fits to model values.** With a fitted reward  $\hat{r}$ , denote by  $\mathcal{X}_{\text{eval}}$  a held-out evaluation  
 145 context set drawn from the target distribution  $P_X^{\text{tgt}} = P_X^{\text{EXP}}$  and by  $n_{\text{EXP},a}$  the number of EXP  
 146 samples with  $A_j = a$ . The direct-method (DM) and doubly robust (DR) value estimators are

$$\hat{\mu}^{\text{DM}}(a) = \frac{1}{|\mathcal{X}_{\text{eval}}|} \sum_{x \in \mathcal{X}_{\text{eval}}} \mathbb{E}_{M \sim p_{\text{sim}}(\cdot | x, a)} [\hat{r}(x, M)], \quad (9)$$

$$\hat{\mu}^{\text{DR}}(a) = \hat{\mu}^{\text{DM}}(a) + \frac{1}{n_{\text{EXP},a}} \sum_{j: A_j=a} (Y_j - \hat{r}(X_j, M_j)). \quad (10)$$

147 On cached benchmarks the SIM expectation collapses to a single cached draw and the DM term  
 148 reduces to  $\hat{r}(x, m_{\text{cache}}(x, a))$ . Empirically, on the synthetic  $(\beta, n_{\text{OBS}}, n_{\text{EXP}})$  grid DR does not  
 149 improve agent-level value RMSE over DM on average (Appendix Table 9: DR adds +0.013 to  
 150 +0.019 RMSE for **EXP-Only**, **OBS-Only**, **CVCI**, and **Representation-EXP**), while in a separate  
 151 strong-self-selection regime with intentional reward-model underfit, DR halves agent-level value  
 152 RMSE (Appendix Figure 4). DR is therefore a model-level diagnostic, not a default. Sections 5–5.5  
 153 diagnose these benchmark-specific family instantiations; Appendix A.4 separates formal estimator  
 154 support from empirical target-alignment diagnostics.

Table 3: Regret-optimal estimator on the semi-synthetic summarization benchmark across confounding strength  $\beta$  and budgets  $(n_{\text{OBS}}, n_{\text{EXP}})$  (mean regret over 30 seeds; appendix table includes regret gaps).

Budget cell	$\beta = 0$	$\beta = 0.2$	$\beta = 0.5$	$\beta = 0.8$	$\beta = 0.9$	$\beta = 0.99$
(2,000, 20)	CVCI-Residual	CVCI	CVCI	OBS-Only	Grounded	Grounded
(20,000, 20)	CVCI	CVCI	CVCI	CVCI	CVCI	CVCI
(2,000, 100)	CVCI	CVCI	CVCI	OBS-Only	OBS-Only	OBS-Only
(20,000, 100)	CVCI	EXP-Only	EXP-Only	EXP-Only	EXP-Only	EXP-Only

## 155 5 Empirical benchmarks

### 156 5.1 Benchmark setup

157 Across all benchmarks, the task is to estimate a reward surface, recommend the best action on each  
 158 held-out context, and measure held-out recommendation regret. The OBS choice law is governed by  
 159 a confounding strength parameter  $\beta \geq 0$  (with  $\beta = 0$  yielding uniform OBS sampling and larger  $\beta$   
 160 inducing stronger user-side selection on  $A$ ); EXP randomizes uniformly over  $\mathcal{A}_{\text{EXP}}$ ; and the two real  
 161 cached benchmarks replace the latent target  $q$  with  $q_{\text{cache}}$ . Each fitted reward  $\hat{r}$  induces a per-context  
 162 estimated value  $\hat{q}(x, a) = \mathbb{E}_{M \sim p_{\text{sim}}(\cdot|x, a)}[\hat{r}(x, M)]$  on the controlled benchmark, or  $\hat{q}(x, a) =$   
 163  $\hat{r}(x, m_{\text{cache}}(x, a))$  on cached benchmarks, and a recommendation  $\hat{\pi}(x) = \arg \max_{a \in \mathcal{A}} \hat{q}(x, a)$ . We  
 164 report

$$\text{Regret}_{\text{test}} = \frac{1}{|\mathcal{X}_{\text{test}}|} \sum_{x \in \mathcal{X}_{\text{test}}} \left[ \max_{a \in \mathcal{A}} q(x, a) - q(x, \hat{\pi}(x)) \right],$$

165 with  $q$  replaced by  $q_{\text{cache}}$  on held-out cached grids. Appendix A.12 gives the full benchmark laws  
 166 and the benchmark snapshot table.

### 167 5.2 Controlled validation as a theory bridge

168 The controlled semi-synthetic validation replaces judged rewards in the summarization task family  
 169 with a known latent reward generator, so identification-side error and estimation-side error can  
 170 be measured separately. Sweeping the confounding strength  $\beta$  over six values and the budgets  
 171  $(n_{\text{OBS}}, n_{\text{EXP}})$  over four settings yields a 24-cell regret map (Table 3). The pattern is regime-  
 172 structured: at the smallest  $(n_{\text{OBS}}, n_{\text{EXP}}) = (2,000, 20)$  row, four different families win across the  
 173 six  $\beta$  values; at  $(20,000, 20)$ , CVCI wins all six  $\beta$  cells; at  $(2,000, 100)$ , the winner shifts from CVCI  
 174 at low  $\beta$  to OBS-Only at high  $\beta$ ; at  $(20,000, 100)$ , EXP-Only wins five of the six cells. CVCI is the  
 175 most stable single winner, with most regret gaps below 0.002 in absolute terms (Appendix Table 15  
 176 reproduces the full grid with regret-gap annotations).

177 The cached benchmarks below retain the same families. Their winner pattern differs from this  
 178 synthetic one (Section 5.4), so reward definition and auxiliary-feature alignment, not just the EXP  
 179 budget, drive estimator choice in the real setting.

### 180 5.3 Main-text benchmarks

181 The main paper reports two real-task cached benchmarks with synthetic OBS/EXP resampling:  
 182 CNN/DailyMail summarization [Hermann et al., 2015, Nallapati et al., 2016, See et al., 2017]  
 183 and SWE-bench Verified coding with BouncerBench patches [Jimenez et al., 2024, OpenAI, 2024,  
 184 Mathews and Nagappan, 2025]. Both use real tasks, real candidate outputs, and real judged or  
 185 program-test rewards, while OBS and EXP are benchmark constructions over the cached pool.  
 186 Appendix Table 14 lists the source data, cached objects, and targets.

### 187 5.4 Real-task cached summarization benchmark

188 The cached summarization benchmark uses 48 difficult CNN/DailyMail articles, 20 candidate summa-  
 189 rization systems, and one judged summary per article–model pair. It evaluates smooth and sharpened  
 190 user-segment reward maps on the same cache; OBS and EXP samples are obtained by resampling the  
 191 cached pool under segment-dependent routing and uniform randomization. Appendix A.10 gives the

Method	$n_{\text{OBS}} = 2,000$		$n_{\text{OBS}} = 20,000$	
	$n_{\text{EXP}} = 20$	$n_{\text{EXP}} = 200$	$n_{\text{EXP}} = 20$	$n_{\text{EXP}} = 200$
CVCI	0.0128	0.0128	0.0120	0.0120
OBS-Only	0.0128	0.0128	0.0123	0.0123
Grounded	0.0200	0.0133	0.0133	0.0135
EXP-Only	0.0170	0.0121	0.0140	0.0118
Representation-EXP	0.0330	0.0184	0.0365	0.0124
Best OBS-based method	0.0128	0.0128	0.0120	0.0120
EXP-Only – best OBS-based	+0.0042	-0.0007	+0.0020	-0.0002

Table 4: Recommendation regret under the smooth aggregate reward (lower is better; mean over 30 seeds). Positive values in the last row mean an OBS-assisted family beats EXP-Only at that budget.

Table 5: Regret-optimal estimator on the real summarization benchmark across  $\beta$  and budgets, parallel to Table 3 (mean over 30 seeds; appendix Table 16 has regret gaps).

Budget cell	$\beta = 0$	$\beta = 0.2$	$\beta = 0.5$	$\beta = 0.8$	$\beta = 0.9$	$\beta = 0.99$
(2,000, 20)	Grounded	Grounded	Grounded	OBS-Only	OBS-Only	Grounded
(20,000, 20)	Grounded	Grounded	EXP-Only	EXP-Only	EXP-Only	EXP-Only
(2,000, 100)	OBS-Only	OBS-Only	OBS-Only	CVCI	CVCI	CVCI
(20,000, 100)	CVCI-Residual	CVCI	CVCI	CVCI	CVCI	CVCI

192 exact slate construction, reward maps, budget grid, and standard errors; Appendix Table 11 reports  
 193 the full aggregate ranking summary.

194 **Supervision scarcity.** Table 4 traces a budget-driven crossover. At  $n_{\text{EXP}} = 20$  the best OBS-  
 195 assisted family improves on EXP-Only, and at  $n_{\text{EXP}} = 200$  EXP-Only catches up or slightly wins.  
 196 This matches the post-identification view in Table 2: once EXP and SIM identify the target, additional  
 197 randomized labels improve the EXP-side reward fit. In the aggregate smooth/sharpened comparison,  
 198 CVCI has the best average rank and excess regret, with a small gap to OBS-Only relative to cell-level  
 199 uncertainty (Appendix Tables 11 and 13). Reward-shape and  $R^2$  diagnostics in Appendix Table 12  
 200 show that varying the target reward while fixing the cached outputs changes alignment with the  
 201 rubric-linear and proxy classes.

202 **Real-vs-synthetic winner-map contrast.** On the same six- $\beta$  four-budget grid, the regret winner  
 203 pattern on the cached real benchmark differs from the synthetic surface (Table 5). Grounded wins  
 204 five of six low-budget cells ( $n_{\text{OBS}}=2,000, n_{\text{EXP}}=20$ ) where CVCI and CVCI-Residual dominate  
 205 the synthetic version, and CVCI reclaims the high-OBS large-EXP row (20,000, 100). Ranking  
 206 diagnostics tell yet another story: Grounded wins top-1 and top-3 in all 24 real cells while winning  
 207 regret in only six (Appendix A.14). Estimator choice is therefore a function of (data regime,  
 208 deployment loss, target reward) jointly, and a single “best” family does not transfer between the  
 209 synthetic and the real benchmark even when the budget grid is identical.

## 210 5.5 Real-task cached coding benchmark

211 The cached coding benchmark uses SWE-bench Verified issues [Jimenez et al., 2024, OpenAI, 2024]  
 212 and one BouncerBench candidate patch per issue-agent pair [Mathews and Nagappan, 2025]. It  
 213 fixes  $n_{\text{OBS}} = 2,000, n_{\text{EXP}} = 100$ , and  $\beta = 0.5$ , then varies the reward definition. Each patch has  
 214 true fix success  $c_1$  and patch-quality components  $c_2, c_3, c_4$ . The fix-success weight  $\alpha_{\text{fix}}$  shifts target  
 215 mass toward  $c_1$ , while  $\omega_{\text{weak}}$  changes the non-success components from additive to weakest-link  
 216 aggregation. Appendix A.15 gives the full construction,  $R^2$  diagnostics, and pairwise heatmaps.

217 **Reward-alignment crossover.** For user type  $u$ , let  $w_{u1} \in [0, 1]$  denote the fix-success mass for  
 218 that type and  $\bar{w}_{uk} \geq 0$  ( $k = 2, 3, 4, \sum_k \bar{w}_{uk} = 1$ ) the normalized weights on the patch-quality  
 219 components. The benchmark utility is then

$$Y_u = w_{u1} \alpha_{\text{fix}} c_1 + (1 - w_{u1}) \left( (1 - \omega_{\text{weak}}) \sum_{k=2}^4 \bar{w}_{uk} c_k + \omega_{\text{weak}} \min_{k=2,3,4} c_k \right),$$

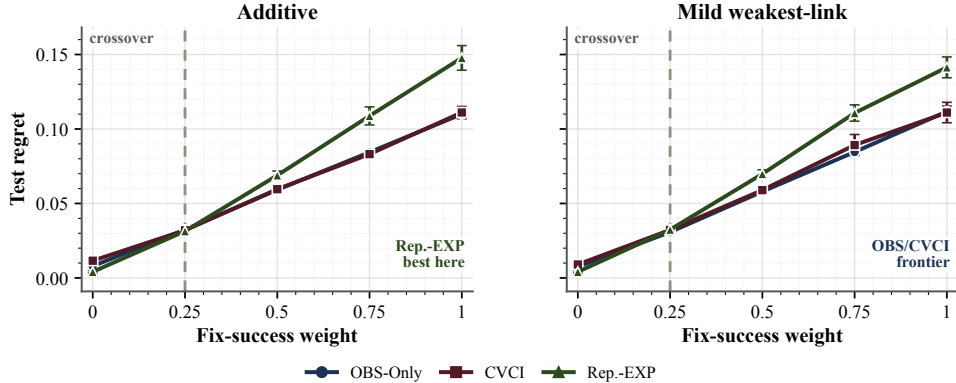


Figure 3: Mean regret in the cached coding benchmark as the fix-success weight  $\alpha_{\text{fix}}$  increases (30 seeds; lower is better). At low  $\alpha_{\text{fix}}$ , the target follows auxiliary patch-quality components and Representation-EXP is competitive; above the crossover near 0.25, OBS-based families overtake it.

220 and  $q_{\text{cache}}$  averages  $Y_u$  over benchmark user types. Figure 3 shows a crossover near  $\alpha_{\text{fix}} \approx 0.25$ :  
 221 Representation-EXP is among the lowest-regret families when the target follows patch quality, while  
 222 OBS-based families win once true fix success dominates. On the additive slice, the regret gap between  
 223 Representation-EXP and the better of OBS-Only and CVC1 moves from  $-0.0039$  at  $\alpha_{\text{fix}} = 0$  to  
 224  $+0.0003$  at 0.25 and  $+0.0415$  at 1. The mild weakest-link slice ( $\omega_{\text{weak}} = 0.25$ ) preserves the same  
 225 crossover and ordering, isolating the effect to which component of the reward dominates. Tilting  
 226 toward  $c_1$  moves the target away from auxiliary patch-quality signal that Representation-EXP can  
 227 extract from OBS labels and toward the binary fix-success signal that pooled OBS+EXP fits absorb  
 228 directly. The full  $5 \times 5$  pairwise delta heatmaps in Appendix Figure 5 show the same target-alignment  
 229 diagnostic across both aggregation modes.

## 230 6 Related work

231 The paper combines off-policy evaluation, hybrid observational–randomized estimation, and LLM  
 232 evaluation. Standard Direct Method with Replay (DM) and Doubly Robust (DR) estimators come  
 233 from contextual-bandit, policy-evaluation, and missing-data work [Robins et al., 1994, Bang and  
 234 Robins, 2005, Dudík et al., 2011, 2014, Swaminathan and Joachims, 2015, Jiang and Li, 2016,  
 235 Thomas and Brunskill, 2016], where logged outcomes are typically assumed unconfounded given  
 236 context. Our setting differs because user-side factors can confound both model choice and outcome,  
 237 so SIM plus EXP first identifies the reward surface and OBS enters only afterward.

238 The estimator families are closest to experimental grounding, data fusion, trial generalizability, and  
 239 recent methods for combining large biased samples with smaller randomized samples [Kallus et al.,  
 240 2018, Cole and Stuart, 2010, Pearl and Bareinboim, 2011, Bareinboim and Pearl, 2016, Rosenman  
 241 et al., 2023, Cheng and Cai, 2021, Lin et al., 2025, Yang et al., 2025, Colnet et al., 2024]. The  
 242 experiments draw on summarization and LLM-evaluation work, including CNN/DailyMail [Hermann  
 243 et al., 2015, Nallapati et al., 2016, See et al., 2017, Stiennon et al., 2020], benchmark and judge  
 244 reliability studies [Liang et al., 2023, Zheng et al., 2023, Chiang et al., 2024], real-world logs [Zheng  
 245 et al., 2024, Zhao et al., 2024], evaluator-bias analyses [Wang et al., 2023, Dubois et al., 2024,  
 246 Panickssery et al., 2024, Verga et al., 2024], preference-learning pipelines [Christiano et al., 2017,  
 247 Ouyang et al., 2022, Rafailov et al., 2023], and causal NLP views of text as treatment, outcome,  
 248 mediator, or proxy [Feder et al., 2022]. Appendix A.16 preserves the fuller positioning.

## 249 7 Discussion and Future Work

250 Each benchmark probes a different driver of estimator choice. The controlled validation isolates the  
 251 EXP-budget effect with a known reward generator; the summarization benchmark stresses supervision  
 252 scarcity at fixed reward geometry; and the coding benchmark stresses reward-target geometry at fixed  
 253 supervision. The pattern across the three sweeps yields concrete deployment guidance.

254 **A practical recipe.** Tie estimator choice to three observable quantities. (i) The deployment loss:  
255 regret and top-rank metrics rank families differently, with **Grounded** winning all 24 top-1 and top-3  
256 ranking cells in the cached summarization benchmark while winning regret in only six. (ii) The  
257 EXP budget  $n_{\text{EXP}}$  relative to the OBS budget  $n_{\text{OBS}}$ : OBS-assisted families gain most when EXP is  
258 scarce, and **EXP-Only** catches up as EXP grows. (iii) How well the auxiliary representation  $\psi$  or the  
259 OBS predictor  $f_{\text{OBS}}$  tracks the held-out target: both **Representation-EXP** and **OBS-anchored families**  
260 inherit the alignment of the OBS signal with the deployment reward.

261 **When OBS information helps and when it does not.** Whenever the deployment reward is well  
262 predicted by the same features that drive OBS outcomes, the pooled fits in (7) reduce variance without  
263 paying a meaningful bias cost; **CVCI** and **Grounded** are then strong defaults. When the deployment  
264 reward depends on a signal that OBS labels do not contain — the binary fix-success component in  
265 coding, or a sharply re-weighted user segment in summarization — the same OBS pull becomes a  
266 bias liability and **EXP-Only** or **Representation-EXP** becomes preferable. The coding crossover near  
267  $\alpha_{\text{fix}} \approx 0.25$  in Figure 3 is a concrete instance of this transition.

268 **Scope and extensions.** The real-task benchmarks use real candidate outputs and real judged or  
269 program-test rewards; OBS and EXP mechanisms are controlled constructions over cached pools,  
270 which is what makes the supervision-scarcity and reward-geometry sweeps cleanly interpretable.  
271 Two extensions follow directly: broader tasks with naturally observed deployment logs paired with  
272 live randomized trials, and multi-round human-agent interaction, which relaxes Assumption A4 to  
273 allow stateful  $U \rightarrow M$  paths through dialogue history. Appendix A.16 expands on scope, limitations,  
274 and alternative scoring rules.

## 275 8 Conclusion

276 Offline evaluation of language models separates cleanly into identification and post-identification  
277 estimation. SIM and EXP together identify causal model values on shared support, while OBS  
278 contributes after identification through auxiliary signal or biased outcomes whenever they predict the  
279 causal target. Across the controlled validation and two cached real-task probes, the best estimator  
280 tracks the EXP budget and the alignment between OBS-derived structure and the target reward; live  
281 randomization remains the prerequisite for any deployment-grade comparison.

## 282 References

- 283 Heejung Bang and James M. Robins. Doubly robust estimation in missing data and causal inference  
284 models. *Biometrics*, 61(4):962–973, 2005. doi: 10.1111/j.1541-0420.2005.00377.x.
- 285 Elias Bareinboim and Judea Pearl. Causal inference and the data-fusion problem. *Proceedings of the*  
286 *National Academy of Sciences*, 113(27):7345–7352, 2016. doi: 10.1073/pnas.1510507113. URL  
287 <https://www.pnas.org/doi/10.1073/pnas.1510507113>.
- 288 David Cheng and Tianxi Cai. Adaptive combination of randomized and observational data, 2021.  
289 URL <https://arxiv.org/abs/2111.15012>.
- 290 Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng  
291 Li, Banghua Zhu, Hao Zhang, Michael Jordan, Joseph E. Gonzalez, and Ion Stoica. Chatbot  
292 arena: An open platform for evaluating LLMs by human preference. In *Proceedings of the*  
293 *41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine*  
294 *Learning Research*, pages 8359–8388, 2024. URL [https://proceedings.mlr.press/v235/](https://proceedings.mlr.press/v235/chiang24b.html)  
295 [chiang24b.html](https://proceedings.mlr.press/v235/chiang24b.html).
- 296 Paul F. Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep  
297 reinforcement learning from human preferences. In *Advances in Neural Information Processing*  
298 *Systems*, volume 30, 2017. URL <https://arxiv.org/abs/1706.03741>.
- 299 Stephen R. Cole and Elizabeth A. Stuart. Generalizing evidence from randomized clinical trials to  
300 target populations: The ACTG 320 trial. *American Journal of Epidemiology*, 172(1):107–115,  
301 2010. doi: 10.1093/aje/kwq084.

- 302 Bénédicte Colnet, Imke Mayer, Guanhua Chen, Awa Dieng, Ruohong Li, Gaël Varoquaux, Jean-  
303 Philippe Vert, Julie Josse, and Shu Yang. Causal inference methods for combining randomized  
304 trials and observational studies: A review. *Statistical Science*, 39(1):165–191, 2024. doi: 10.1214/  
305 23-STSS889. URL <https://doi.org/10.1214/23-STSS889>.
- 306 Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B. Hashimoto. Length-controlled  
307 alpacaeval: A simple way to debias automatic evaluators. In *First Conference on Language*  
308 *Modeling (COLM)*, 2024. URL <https://openreview.net/forum?id=CybEmzWBX0>.
- 309 Miroslav Dudík, John Langford, and Lihong Li. Doubly robust policy evaluation and learning. In  
310 *Proceedings of the 28th International Conference on Machine Learning*, 2011. URL <https://arxiv.org/abs/1103.4601>.  
311
- 312 Miroslav Dudík, Dumitru Erhan, John Langford, and Lihong Li. Doubly robust policy evaluation  
313 and optimization. *Statistical Science*, 29(4):485–511, 2014. doi: 10.1214/14-STSS500. URL  
314 <https://doi.org/10.1214/14-STSS500>.
- 315 Amir Feder, Katherine A. Keith, Emaad Manzoor, Reid Pryzant, Dhanya Sridhar, Zach Wood-  
316 Doughty, Jacob Eisenstein, Justin Grimmer, Roi Reichart, Margaret E. Roberts, Brandon M.  
317 Stewart, Victor Veitch, and Diyi Yang. Causal inference in natural language processing: Estima-  
318 tion, prediction, interpretation and beyond. *Transactions of the Association for Computational*  
319 *Linguistics*, 10:1138–1158, 2022. doi: 10.1162/tacl\_a\_00511.
- 320 Gemma Team. Gemma 3 technical report, 2025. URL <https://arxiv.org/abs/2503.19786>.
- 321 Karl Moritz Hermann, Tomáš Kočiský, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa  
322 Suleyman, and Phil Blunsom. Teaching machines to read and comprehend. In *Advances in*  
323 *Neural Information Processing Systems*, volume 28, 2015. URL [https://arxiv.org/abs/](https://arxiv.org/abs/1506.03340)  
324 [1506.03340](https://arxiv.org/abs/1506.03340).
- 325 Nan Jiang and Lihong Li. Doubly robust off-policy value evaluation for reinforcement learning. In  
326 *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceed-*  
327 *ings of Machine Learning Research*, pages 652–661, 2016. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v48/jiang16.html)  
328 [press/v48/jiang16.html](https://proceedings.mlr.press/v48/jiang16.html).
- 329 Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik  
330 Narasimhan. SWE-bench: Can language models resolve real-world GitHub issues? In *The Twelfth*  
331 *International Conference on Learning Representations*, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2310.06770)  
332 [2310.06770](https://arxiv.org/abs/2310.06770).
- 333 Nathan Kallus, Aahlad Manas Puli, and Uri Shalit. Removing hidden confounding by ex-  
334 perimental grounding. In *Advances in Neural Information Processing Systems*, volume 31,  
335 pages 10911–10920, 2018. URL [https://proceedings.neurips.cc/paper/2018/hash/](https://proceedings.neurips.cc/paper/2018/hash/566f0ea4f6c2e947f36795c8f58ba901-Abstract.html)  
336 [566f0ea4f6c2e947f36795c8f58ba901-Abstract.html](https://proceedings.neurips.cc/paper/2018/hash/566f0ea4f6c2e947f36795c8f58ba901-Abstract.html).
- 337 Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian  
338 Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby  
339 Yan, Ce Zhang, Christian Cosgrove, Christopher D. Manning, Christopher Ré, Diana Acosta-Navas,  
340 Drew A. Hudson, Eric Zelikman, Esin Durmus, Faisal Ladhak, Frieda Rong, Hongyu Ren, Huaxiu  
341 Yao, Jue Wang, Keshav Santhanam, Laurel Orr, Lucia Zheng, Mert Yuksekogun, Mirac Suzgun,  
342 Nathan Kim, Neel Guha, Niladri Chatterji, Omar Khattab, Peter Henderson, Qian Huang, Ryan  
343 Chi, Sang Michael Xie, Shibani Santurkar, Surya Ganguli, Tatsunori Hashimoto, Thomas Icard,  
344 Tianyi Zhang, Vishrav Chaudhary, William Wang, Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta  
345 Koreeda. Holistic evaluation of language models. *Transactions on Machine Learning Research*,  
346 2023. URL <https://arxiv.org/abs/2211.09110>.
- 347 Xi Lin, Jens Magelund Tarp, and Robin J. Evans. Combining experimental and observational data  
348 through a power likelihood. *Biometrics*, 81(1):ujaf008, 2025. doi: 10.1093/biometc/ujaf008. URL  
349 <https://academic.oup.com/biometrics/article/81/1/ujaf008/8016472>.
- 350 Noble Saji Mathews and Meiyappan Nagappan. Is your automated software engineer trustworthy?,  
351 2025. URL <https://arxiv.org/abs/2506.17812>.

- 352 Ramesh Nallapati, Bowen Zhou, Cicero dos Santos, Caglar Gulcehre, and Bing Xiang. Abstractive  
353 text summarization using sequence-to-sequence RNNs and beyond. In *Proceedings of The 20th*  
354 *SIGNLL Conference on Computational Natural Language Learning*, pages 280–290, 2016. URL  
355 <https://arxiv.org/abs/1602.06023>.
- 356 OpenAI. Introducing SWE-bench Verified. OpenAI Blog, 2024. URL [https://openai.com/  
357 index/introducing-swe-bench-verified/](https://openai.com/index/introducing-swe-bench-verified/).
- 358 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela  
359 Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman,  
360 Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welin-  
361 der, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow in-  
362 structions with human feedback. In *Advances in Neural Information Processing Systems*,  
363 volume 35, 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/  
364 hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2022/hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html).
- 365 Arjun Panickssery, Samuel R. Bowman, and Shi Feng. LLM evaluators recognize and favor their  
366 own generations, 2024. URL <https://arxiv.org/abs/2404.13076>.
- 367 Judea Pearl and Elias Bareinboim. Transportability of causal and statistical relations: A formal  
368 approach. In *Proceedings of the Twenty-Fifth AAAI Conference on Artificial Intelligence*, 2011.  
369 doi: 10.1609/aaai.v25i1.7861.
- 370 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D. Manning, Stefano Er-  
371 mon, and Chelsea Finn. Direct preference optimization: Your language model is se-  
372 cretely a reward model. In *Advances in Neural Information Processing Systems*, vol-  
373 ume 36, 2023. URL [https://proceedings.neurips.cc/paper\\_files/paper/2023/  
374 hash/a85b405ed65c6477a4fe8302b5e06ce7-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2023/hash/a85b405ed65c6477a4fe8302b5e06ce7-Abstract-Conference.html).
- 375 James M. Robins, Andrea Rotnitzky, and Lue Ping Zhao. Estimation of regression coefficients when  
376 some regressors are not always observed. *Journal of the American Statistical Association*, 89(427):  
377 846–866, 1994. doi: 10.1080/01621459.1994.10476818.
- 378 Evan T. R. Rosenman, Guillaume Basse, Art B. Owen, and Mike Baiocchi. Combining observational  
379 and experimental datasets using shrinkage estimators. *Biometrics*, 79(4):2961–2973, 2023. doi:  
380 10.1111/biom.13827. URL <https://pubmed.ncbi.nlm.nih.gov/36629736/>.
- 381 Abigail See, Peter J. Liu, and Christopher D. Manning. Get to the point: Summarization with  
382 pointer-generator networks. In *Proceedings of the 55th Annual Meeting of the Association for*  
383 *Computational Linguistics*, pages 1073–1083, 2017. URL [https://arxiv.org/abs/1704.  
384 04368](https://arxiv.org/abs/1704.04368).
- 385 Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,  
386 Dario Amodei, and Paul Christiano. Learning to summarize from human feedback. In *Advances*  
387 *in Neural Information Processing Systems*, volume 33, 2020. URL [https://arxiv.org/abs/  
388 2009.01325](https://arxiv.org/abs/2009.01325).
- 389 Adith Swaminathan and Thorsten Joachims. Counterfactual risk minimization: Learning from logged  
390 bandit feedback. In *Proceedings of the 32nd International Conference on Machine Learning*,  
391 volume 37 of *Proceedings of Machine Learning Research*, pages 814–823, 2015. URL <https://proceedings.mlr.press/v37/swaminathan15.html>.
- 393 Philip S. Thomas and Emma Brunskill. Data-efficient off-policy policy evaluation for reinforce-  
394 ment learning. In *Proceedings of The 33rd International Conference on Machine Learning*,  
395 volume 48 of *Proceedings of Machine Learning Research*, pages 2139–2148, 2016. URL  
396 <https://proceedings.mlr.press/v48/thomasa16.html>.
- 397 Victor Veitch, Dhanya Sridhar, and David M. Blei. Adapting text embeddings for causal inference. In  
398 *Proceedings of the 36th Conference on Uncertainty in Artificial Intelligence (UAI)*, volume 124 of  
399 *Proceedings of Machine Learning Research*, pages 919–928, 2020. URL [https://proceedings.  
400 mlr.press/v124/veitch20a.html](https://proceedings.mlr.press/v124/veitch20a.html).

- 401 Pat Verga, Sebastian Hofstätter, Sophia Althammer, Yixuan Su, Aleksandra Piktus, Arkady Arkhang-  
402 orodsky, Minjie Xu, Naomi White, and Patrick Lewis. Replacing judges with juries: Evaluating  
403 LLM generations with a panel of diverse models, 2024. URL [https://arxiv.org/abs/2404.](https://arxiv.org/abs/2404.18796)  
404 18796.
- 405 Peiyi Wang, Lei Li, Liang Chen, Dawei Zhu, Binghui Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and  
406 Zhifang Sui. Large language models are not fair evaluators, 2023. URL [https://arxiv.org/](https://arxiv.org/abs/2305.17926)  
407 [abs/2305.17926](https://arxiv.org/abs/2305.17926).
- 408 Xuelin Yang, Licong Lin, Susan Athey, Michael I. Jordan, and Guido W. Imbens. Cross-validated  
409 causal inference: a modern method to combine experimental and observational data, 2025. URL  
410 <https://arxiv.org/abs/2511.00727>.
- 411 Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng. Wildchat:  
412 1M ChatGPT interaction logs in the wild. In *The Twelfth International Conference on Learning*  
413 *Representations*, 2024. URL <https://arxiv.org/abs/2405.01470>.
- 414 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,  
415 Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica.  
416 Judging LLM-as-a-judge with MT-bench and chatbot arena. In *Advances in Neural Information*  
417 *Processing Systems*, volume 36, 2023. URL <https://arxiv.org/abs/2306.05685>.
- 418 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Tianle Li, Siyuan Zhuang, Zhanghao Wu, Yonghao  
419 Zhuang, Zhuohan Li, Zi Lin, Eric P. Xing, Joseph E. Gonzalez, Ion Stoica, and Hao Zhang.  
420 LMSYS-Chat-1M: A large-scale real-world LLM conversation dataset. In *The Twelfth International*  
421 *Conference on Learning Representations*, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=B0fDKxfwt0)  
422 [B0fDKxfwt0](https://openreview.net/forum?id=B0fDKxfwt0).

423 **A Additional details**

424 **A.1 Extended problem setting and assumptions**

425 This appendix preserves the full three-source design and assumption motivation behind Section 2.

426 **Accessible data sources.** Three data sources are available, each contributing something the others  
 427 cannot provide alone. OBS is a large observational dataset

$$D_{\text{OBS}} = \{(X_i, A_i, M_i, Y_i, Z_i)\}_{i=1}^{n_{\text{OBS}}},$$

428 where  $A_i$  may depend on the corresponding latent state  $U_i$  and  $Z_i \in [0, 1]^K$  are auxiliary labels that  
 429 may or may not be available. The auxiliary labels provide extra supervision for learning a compact  
 430 proxy representation of the context–output pair; user feedback ratings and pairwise preferences are  
 431 natural sources and can often be collected at scale [Christiano et al., 2017, Ouyang et al., 2022,  
 432 Chiang et al., 2024].

433 EXP is a much smaller dataset

$$D_{\text{EXP}} = \{(X_j, A_j, M_j, Y_j)\}_{j=1}^{n_{\text{EXP}}},$$

434 in which model choice is randomized and therefore unconfounded. Since  $A_j$  is drawn from the  
 435 experimental action set  $\mathcal{A}_{\text{EXP}}$  according to a fixed distribution,  $A \perp U$  and  $A \perp U \mid X$ , so EXP  
 436 outcomes are causally interpretable. Randomizing model choice disrupts user experience, which is  
 437 why these experiments remain small in practice.

438 The simulator reruns any generative model on any context and returns a mediator draw from  $M \sim$   
 439  $p_{\text{sim}}(\cdot \mid X = x, A = a)$ . Operationally, SIM reruns generative model  $a$  on held-out context  $x$  and  
 440 records the generated output. In generative systems, SIM is especially natural because prompts,  
 441 retrieved documents, tool traces, and other environment inputs can often be stored as deterministic  
 442 inputs, while the remaining randomness comes from the generative model itself. Figure 2 (main text)  
 443 shows the three-source pipeline.

444 **Target estimand details.** The marginal model value asks how well each model would perform on  
 445 average if used by the full target population:

$$\mu(a) := \mathbb{E}[Y(\text{do}(A = a))].$$

446 The context-conditional model value asks what outcome each model would deliver for a particular  
 447 observed context:

$$q(x, a) := \mathbb{E}[Y(\text{do}(A = a)) \mid X = x].$$

448 Section 3 shows how both targets are recovered: SIM supplies the model-specific mediator distribution  
 449 and EXP provides unconfounded scores for realized outputs. In real cached benchmarks, the finite  
 450 cached target  $q_{\text{cache}}$  approximates  $q(x, a)$  while the superpopulation estimand remains latent.

451 **Design assumptions for SIM-based identification. A1 (Randomization in EXP).** In EXP,  $A$  is  
 452 randomized over  $\mathcal{A}_{\text{EXP}}$  so that  $A \perp U$  and  $A \perp U \mid X$  hold. Randomization is the sole source of  
 453 unconfounded outcome information; without it, identification is impossible regardless of OBS size.

454 **A2 (SIM validity and support).** The simulator matches the interventional mediator distribution,  
 455  $p_{\text{sim}}(m \mid x, a) = \mathbb{P}(M = m \mid X = x, \text{do}(A = a))$ , and is supported within the EXP mediator  
 456 support: for every  $x$  in the EXP context support and every  $a \in \mathcal{A}_{\text{EXP}}$ ,

$$\text{supp } p_{\text{sim}}(\cdot \mid x, a) \subseteq \text{supp } \mathbb{P}(M \mid X = x, A = a, \text{EXP}).$$

457 Throughout the paper, “the shared EXP/SIM support” refers to this support condition. SIM validity is  
 458 satisfied whenever any randomness in the generated output comes from the generative model itself  
 459 and all deployment-time inputs that affect the generative model’s behavior are included in  $X$ .

460 **A3 (Outcome consistency under intervention).** OBS and EXP measure the same outcome  $Y$ , and  
 461 the conditional outcome law given  $(X, M, U)$  is invariant across OBS, EXP, and the post-intervention  
 462 world: for every  $(x, m, u)$  on the relevant support,

$$\begin{aligned} \mathbb{P}(Y \mid X = x, M = m, U = u, \text{OBS}) &= \mathbb{P}(Y \mid X = x, M = m, U = u, \text{EXP}) \\ &= \mathbb{P}(Y \mid X = x, M = m, U = u, \text{do}(A = a)). \end{aligned}$$

463 Only the model-choice mechanism differs across OBS, EXP, and intervention. Without this, EXP  
 464 supervision could not transport either to the intervention or to OBS-trained reward models.

465 **A4 (No latent mediator confounding; single-round generation).** Conditional on  $(X, A)$ , the latent  
 466 state  $U$  is independent of the mediator  $M$  in OBS, EXP, and under intervention:

$$U \perp M \mid X, A.$$

467 Equivalently, the SCM has no edge  $U \rightarrow M$  once  $(X, A)$  is fixed, and any deployment-time inputs  
 468 that drive the generative model are absorbed into  $X$ . A4 restricts the framework to single-round,  
 469 stateless generation. In repeated, adaptive, or multi-round interactions the residual dependence of  $M$   
 470 on  $U$  given  $(X, A)$  is generally nonzero, so A4 fails.

471 **A5 (Context distribution alignment).** The marginal model value  $\mu(a)$  is defined with respect to  
 472 a target context distribution  $P_X^{\text{tgt}}$  that coincides with the EXP context distribution:  $P_X^{\text{tgt}} = P_X^{\text{EXP}}$ .  
 473 Equivalently, the EXP sample is drawn from the same context population over which  $\mu(a)$  is averaged.  
 474 When this fails, only  $q(x, a)$  on the shared support is identified, and reweighting to a different  $P_X^{\text{tgt}}$   
 475 requires standard transportability and trial-generalizability arguments [Cole and Stuart, 2010, Pearl  
 476 and Bareinboim, 2011, Bareinboim and Pearl, 2016].

## 477 A.2 Identification remarks

478 **SIM resolves generation; EXP resolves scoring.** SIM regenerates the mediator  $M$  under  
 479  $\text{do}(A = a)$ , while EXP supplies the outcome labels needed to score that mediator. So SIM handles  
 480 counterfactual generation and EXP identifies causal scoring. The randomized sample is what  
 481 identifies the scoring rule  $r^*(x, m)$  for a realized context–output pair.

482 **Why observational outcomes target a different object.** In deployment,  $A$  depends on the latent  
 483 variable  $U$ , and  $U$  also affects  $Y$  directly. Conditioning on  $(X, M)$  therefore leaves self-selection  
 484 bias in place, so in general

$$\mathbb{E}[Y \mid X = x, M = m, \text{OBS}] \neq r^*(x, m).$$

485 This is why the additional assumptions in Section 4 are estimator-side assumptions for using OBS,  
 486 while identification itself comes from SIM and EXP.

## 487 A.3 Estimator details and concrete implementations

488 **Exact reward-model objectives.** This subsection records the exact fitting problems and benchmark-  
 489 specific implementation details deferred from Section 4. **EXP-Only** fits

$$(\hat{w}_{\text{EXP}}, \hat{b}_{\text{EXP}}) \in \arg \min_{w, b} \sum_{j \in \text{EXP}} (Y_j - (w^\top \varphi(z_j) + b))^2 + \alpha \|w\|_2^2,$$

490 with prediction rule  $\hat{r}_{\text{EXP}}(z) = \text{clip}_{[0,1]}(\hat{w}_{\text{EXP}}^\top \varphi(z) + \hat{b}_{\text{EXP}})$ .

491 Representation-EXP fits the same ridge head on the proxy representation:

$$(\hat{w}_\psi, \hat{b}_\psi) \in \arg \min_{w, b} \sum_{j \in \text{EXP}} (Y_j - (w^\top \psi(z_j) + b))^2 + \alpha \|w\|_2^2,$$

492 with  $\hat{r}_\psi(z) = \text{clip}_{[0,1]}(\hat{w}_\psi^\top \psi(z) + \hat{b}_\psi)$ .

493 **OBS-Only** fits the raw-feature ridge model on OBS outcomes,

$$(\hat{w}_{\text{OBS}}, \hat{b}_{\text{OBS}}) \in \arg \min_{w, b} \sum_{i \in \text{OBS}} (Y_i - (w^\top \varphi(z_i) + b))^2 + \alpha \|w\|_2^2,$$

494 and predicts  $f_{\text{OBS}}(z) = \text{clip}_{[0,1]}(\hat{w}_{\text{OBS}}^\top \varphi(z) + \hat{b}_{\text{OBS}})$ .

495 For **Grounded**, define EXP residual targets  $\Delta_j = f_{\text{OBS}}(z_j) - Y_j$  and fit

$$(\hat{\theta}, \hat{c}) \in \arg \min_{\theta, c} \sum_{j \in \text{EXP}} (\Delta_j - (\theta^\top \psi(z_j) + c))^2 + \lambda_\theta \|\theta\|_2^2.$$

Table 6: How each estimator family is implemented in the appendix validation and the real benchmarks. Family names are shared across settings, but the proxy channel and EXP-side tuning differ by benchmark.

Paper family	Semi-synthetic instantiation	Real-benchmark instantiation	EXP-side tuning
EXP-Only	Raw-feature ridge fit on EXP only	Same family on cached real outputs	Fixed defaults
OBS-Only	Raw-feature ridge fit on OBS outcomes only	Same family on cached real outputs	None
Representation-EXP	EXP-only reward head on a proxy learned from OBS auxiliary labels	EXP-only reward head on a heuristic proxy learned from cached trajectory metadata	Fixed defaults
<b>Grounded</b>	OBS baseline plus low-dimensional proxy correction	OBS baseline plus richer proxy-basis correction on the heuristic auxiliary proxy	Small EXP-only tuning
CVCI	Direct pooled OBS/EXP fit in raw text features	Same family with model-level EXP cross-validation for the pooling weight	Benchmark-specific EXP tuning
CVCI-Residual	OBS baseline plus pooled residual fit in proxy space	Same family with the heuristic auxiliary proxy and model-level EXP cross-validation	Benchmark-specific EXP tuning

496 The resulting predictor subtracts the fitted proxy-side correction from the OBS baseline and clips the  
 497 result to  $[0, 1]$ . In the real benchmarks, the grounded family keeps the same outer form but replaces  
 498 the single linear proxy correction with a small library of proxy-basis corrections built on the heuristic  
 499 auxiliary proxy. Appendix A.5 compares this richer correction class with a single-linear baseline and  
 500 with a weak pooled-anchor variant.

501 **CVCI** pools the OBS and EXP regression objectives:

$$\begin{aligned}
 (\hat{w}_\lambda, \hat{b}_\lambda) \in \arg \min_{w, b} & \left[ (1 - \lambda) \frac{1}{n_{\text{EXP}}} \sum_{j \in \text{EXP}} (Y_j - (w^\top \varphi(z_j) + b))^2 \right. \\
 & \left. + \lambda \frac{1}{n_{\text{OBS}}} \sum_{i \in \text{OBS}} (Y_i - (w^\top \varphi(z_i) + b))^2 \right] + \alpha \|w\|_2^2,
 \end{aligned} \tag{11}$$

502 with  $\hat{r}_\lambda(z) = \text{clip}_{[0,1]}(\hat{w}_\lambda^\top \varphi(z) + \hat{b}_\lambda)$ . Here  $\lambda = 1$  is the OBS-only endpoint and  $\lambda = 0$  is the  
 503 EXP-only endpoint. The final predictor refits (11) on the full OBS and EXP data at the value of  $\lambda$   
 504 selected by Appendix A.11.

505 **CVCI-Residual** residualizes around  $f_{\text{OBS}}$  using the residual target  $Y_j - f_{\text{OBS}}(z_j)$  and pools only that  
 506 residual fit in proxy space:

$$\begin{aligned}
 (\hat{\theta}_\lambda, \hat{c}_\lambda) \in \arg \min_{\theta, c} & \left[ (1 - \lambda) \frac{1}{n_{\text{EXP}}} \sum_{j \in \text{EXP}} (Y_j - f_{\text{OBS}}(z_j) - (\theta^\top \psi(z_j) + c))^2 \right. \\
 & \left. + \lambda \frac{1}{n_{\text{OBS}}} \sum_{i \in \text{OBS}} (Y_i - f_{\text{OBS}}(z_i) - (\theta^\top \psi(z_i) + c))^2 \right] + \alpha_\psi \|\theta\|_2^2,
 \end{aligned}$$

507 with  $\hat{r}_{\text{res},\lambda}(z) = \text{clip}_{[0,1]}(f_{\text{OBS}}(z) + \hat{\theta}_\lambda^\top \psi(z) + \hat{c}_\lambda)$ . In the experiments, both the pooling weight  $\lambda$   
 508 and the residual ridge penalty  $\alpha_\psi$  are selected by Appendix A.11; when residual-CVCI candidates  
 509 are tied within numerical tolerance, the implementation resolves the tie toward the more regularized  
 510 candidate.

511 **How each estimator family is implemented.** Table 6 shows how each estimator family is imple-  
 512 mented in the appendix validation and in the two real benchmarks. The family names are the same  
 513 across settings, but the proxy construction and EXP-side tuning differ by benchmark.

#### 514 A.4 Theory-to-evidence map for estimator comparisons

515 The estimator comparison is guided by three evidence tiers. The identification theorem is part of  
 516 the main paper. The estimator-comparison results in Appendix C provide formal support for several  
 517 oracle and approximation statements. Target-alignment patterns in the benchmarks diagnose whether  
 518 a target reward aligns with a proxy, baseline, or correction class. They provide empirical rather than  
 519 formal evidence. Table 7 summarizes this map.

520 Throughout Section 5.4, **Grounded** refers to the rich proxy-basis version in Table 6. Appendix A.5  
 521 compares it with the single-linear baseline and the pooled-anchor variant.

Table 7: Theory-to-evidence map for the estimator comparison. Tier 1 is formal support in the main paper; Tier 2 is formal support from the appendix; Tier 3 is empirical diagnostic evidence only.

Mechanism	Tier	Formal support	Empirical diagnostic	Interpretation rule
SIM plus EXP identify the target	1	Theorem 1	All benchmarks separate randomized scoring from replayed or cached outputs	EXP supplies unfounded scores; SIM or cached replay supplies mediators.
<b>Grounded</b> can enlarge a proxy-side approximation class	2	Theorem 2	Basis-expanded <b>Grounded</b> improves over the single-linear grounded baseline in Appendix A.5	Treat this as benchmark evidence for richer correction bases; the theorem establishes class expansion, not finite-sample dominance.
Oracle correction versus finite-sample correction	2	Theorem 4 and (15)	<b>Grounded</b> is competitive but not uniformly best in the cached benchmarks	The oracle result is not a finite-sample dominance guarantee.
Residual simplicity	2	Theorem 3	<b>CVCI-Residual</b> wins only selected cells and often underperforms in real summarization	Read wins or losses as benchmark evidence about residual simplicity, not as direct tests of the oracle condition.
Proxy alignment and re-ward definition	3	Empirical diagnostic unless additional approximation-class results are formalized	Summarization $R^2$ diagnostics and coding fix-success sweep	Interpret as empirical diagnostics rather than formal oracle statements.

## 522 A.5 Grounded-family implementations for the real cached benchmark

523 **A common grounded template.** All real-benchmark grounded variants start from the same decom-  
524 position: fit an OBS baseline  $f_{\text{OBS}}$  from logged outcomes, then use the heuristic auxiliary proxy to  
525 estimate a correction that is subtracted from that baseline. The variants differ along two design axes:

- 526 1. the correction class used on the proxy representation, and
- 527 2. whether OBS enters the second stage only through  $f_{\text{OBS}}$  or also through a small pooled anchor  
528 after the correction direction has been fixed.

529 The grounded-family comparison instead organizes three principled instantiations of the same  
530 baseline-minus-correction idea.

531 **Variant 1: single-linear grounded correction.** This is the most straightforward grounded in-  
532 stantiation. It uses the heuristic auxiliary proxy, keeps a single linear correction family, and  
533 tunes the correction strength on EXP. For the auxiliary real summarization comparison over  
534  $\beta \in \{0, 0.2, 0.5, 0.8, 0.9, 0.99\}$ ,  $n_{\text{OBS}} \in \{2,000, 20,000\}$ , and  $n_{\text{EXP}} \in \{20, 100\}$ , we report the  
535 version of this linear class tuned by model-level EXP cross-validation, which serves as the matched  
536 linear baseline for the richer grounded variants. That single-linear class uses the predictor

$$\hat{r}_{\text{lin},\alpha}(z) = \text{clip}_{[0,1]}(f_{\text{OBS}}(z) - \alpha(\hat{\theta}^\top \psi_{\text{aux}}(z) + \hat{c})),$$

537 where  $\psi_{\text{aux}}$  is the learned heuristic auxiliary proxy and the correction is linear in that proxy.

538 **Variant 2: **Grounded** as a rich but still small proxy-side correction class.** The main-text  
539 **Grounded** instantiation keeps the same baseline-minus-correction outer form, but replaces the single  
540 linear correction by a small library of proxy-basis corrections. Write  $\tilde{\psi}_{\text{aux}}(z) \in \mathbb{R}^{16}$  for the  
541 standardized 16-dimensional SVD compression of the learned heuristic auxiliary proxy. For

$$B \in \{B_{\text{id}}, B_{\text{poly}}\}, \quad B_{\text{id}}(u) = u, \quad B_{\text{poly}}(u) = [u, u \odot u],$$

542 and ridge level  $\tau \in \{10^{-2}, 1, 100\}$ , the rich grounded direction fits

$$(\hat{\theta}_{B,\tau}, \hat{c}_{B,\tau}) \in \arg \min_{\theta,c} \sum_{j \in \text{EXP}} (\Delta_j - (\theta^\top B(\tilde{\psi}_{\text{aux}}(z_j)) + c))^2 + \tau \|\theta\|_2^2, \quad \Delta_j = f_{\text{OBS}}(z_j) - Y_j,$$

543 and predicts

$$\hat{r}_{\text{rich},B,\alpha}(z) = \text{clip}_{[0,1]}(f_{\text{OBS}}(z) - \alpha(\hat{\theta}_{B,\tau}^\top B(\tilde{\psi}_{\text{aux}}(z)) + \hat{c}_{B,\tau})).$$

Table 8: Grounded-family implementations for the real cached benchmark. The table compares the single-linear grounded baseline, the richer proxy-basis grounded instantiation used in the main text, and a weak pooled-anchor variant built on the same rich correction direction. ‘Regret wins’ counts wins among the seven methods obtained by adding the pooled-anchor variant to the six main-text methods.

Variant	Correction family	Second-stage OBS use	Tuning	Regret wins	Macro regret	Mean top-1	Mean top-3
Single-linear correction baseline	Single linear correction on the heuristic auxiliary proxy	None beyond the baseline $f_{\text{OBS}}$	Model-level EXP cross-validation for $\alpha_{\text{corr}}$	1	0.0357	0.628	0.646
Grounded rich proxy-basis correction	Rich proxy-basis correction on standardized 16-dimensional heuristic auxiliary features with $B \in \{\text{id}, \text{poly2}\}$ and $\tau \in \{10^{-2}, 1, 100\}$	None beyond the baseline $f_{\text{OBS}}$	Model-level EXP cross-validation over $(B, \tau, \alpha_{\text{corr}})$	6	0.0318	0.736	0.736
Rich correction with weak pooled anchor	Same rich EXP-only correction direction as Grounded, followed by anchored coefficient pooling	OBS enters only through the pooled coefficients $(b, \alpha)$	Model-level EXP cross-validation over $(B, \tau, \lambda)$ with closed-form anchored $(b, \alpha)$	4	0.0326	0.744	0.744

544 The final implementation selects  $(B, \tau, \alpha)$  by agent-CV on EXP. Relative to the linear grounded class,  
 545 this replaces the correction family

$$\mathcal{H}_{\text{lin}} = \left\{ u \mapsto \theta^\top u + c \right\}$$

546 by

$$\mathcal{H}_{\text{rich}} = \left\{ u \mapsto \theta^\top B(u) + c : B \in \{B_{\text{id}}, B_{\text{poly}}\} \right\},$$

547 so  $\mathcal{H}_{\text{lin}} \subseteq \mathcal{H}_{\text{rich}}$ . The corresponding grounded predictor class therefore expands the proxy-side  
 548 approximation space while still learning the correction direction from EXP only.

549 **Variant 3: pooled-anchor grounded.** The appendix pooled-anchor variant keeps the rich EXP-only  
 550 correction direction but adds a weak pooled anchor in the second stage. Let

$$\widehat{c}_{B,\tau}(z) = \widehat{\theta}_{B,\tau}^\top B(\tilde{\psi}_{\text{aux}}(z)) + \widehat{c}_{B,\tau}$$

551 denote the grounded correction direction learned from EXP. The pooled-anchor predictor is

$$\widehat{r}_{\text{anchor},\lambda}(z) = \text{clip}_{[0,1]} \left( \widehat{b}_\lambda f_{\text{OBS}}(z) - \widehat{\alpha}_\lambda \widehat{c}_{B,\tau}(z) \right),$$

552 where  $(\widehat{b}_\lambda, \widehat{\alpha}_\lambda)$  are obtained from the anchored objective

$$\begin{aligned} (\widehat{b}_\lambda, \widehat{\alpha}_\lambda) \in \arg \min_{b,\alpha} & \left[ \lambda \frac{1}{n_{\text{OBS}}} \sum_{i \in \text{OBS}} (Y_i - (b f_{\text{OBS}}^{\text{cf}}(z_i) - \alpha \widehat{c}_{B,\tau}(z_i)))^2 \right. \\ & \left. + (1 - \lambda) \frac{1}{n_{\text{EXP}}} \sum_{j \in \text{EXP}} (Y_j - (b f_{\text{OBS}}(z_j) - \alpha \widehat{c}_{B,\tau}(z_j)))^2 \right] \\ & + \rho_b (b - 1)^2 + \rho_\alpha (\alpha - 1)^2. \end{aligned}$$

553 Here  $f_{\text{OBS}}^{\text{cf}}$  is a cross-fitted OBS baseline prediction. This variant pools only two calibration coeffi-  
 554 cients after EXP has already fixed the correction direction.

555 **Why these implementations should behave differently.** The three variants place their modeling  
 556 flexibility in different parts of the decomposition. The linear grounded baseline puts the entire second  
 557 stage inside the smallest class  $\mathcal{H}_{\text{lin}}$ . So it is the most conservative option, but it can underfit whenever  
 558 the discrepancy between  $f_{\text{OBS}}$  and the causal target is only approximately linear in the proxy. Because  
 559 Grounded replaces  $\mathcal{H}_{\text{lin}}$  by the strictly larger class  $\mathcal{H}_{\text{rich}} \supseteq \mathcal{H}_{\text{lin}}$  while still estimating the correction  
 560 direction from EXP only, it should help precisely when the residual mismatch is slightly nonlinear  
 561 in the proxy yet still low-dimensional. The pooled-anchor variant changes a different part of the  
 562 problem. It keeps the correction direction fixed at the grounded solution and lets OBS influence only  
 563 the coarse scaling coefficients  $(b, \alpha)$ . It should help when the main issue is the relative calibration of  
 564 the baseline and the correction, especially once EXP is large enough to identify a stable correction  
 565 direction.

566 **Results on the  $(\beta, n_{\text{OBS}}, n_{\text{EXP}})$  comparison.** Relative to the matched single-linear baseline,  
 567 **Grounded** lowers macro regret from 0.0357 to 0.0318, increases regret wins from 1 to 6, and has  
 568 lower mean regret in 15 of the 24  $(\beta, n_{\text{OBS}}, n_{\text{EXP}})$  settings. The gains are concentrated in (2000, 20)  
 569 and (20000, 100), where **Grounded** wins all six beta values, and in (20000, 20), where it wins three  
 570 of six.

571 The pooled-anchor variant behaves differently. It beats **Grounded** in 11 of the 24 settings and has the  
 572 strongest average ranking metrics in this three-way comparison, but its macro regret is 0.0326 and it  
 573 wins regret in only four settings, all with  $n_{\text{EXP}} = 100$ . Those wins occur only when  $n_{\text{EXP}} = 100$ :

$$(\beta, n_{\text{OBS}}, n_{\text{EXP}}) \in \{(0, 2000, 100), (0.5, 2000, 100), (0.8, 20000, 100), (0.9, 20000, 100)\}.$$

574 The selected pooling weight is usually small: across 720 fitted models,  $\lambda_* \leq 0.03$  in 565 fits and  
 575  $\lambda_* = 0$  in 264 fits. In this comparison, the pooled anchor changes calibration more than it changes  
 576 the regret winner.

## 577 A.6 Value estimators and empirical benchmark metrics

578 The reward models in Section 4 estimate a conditional outcome surface. To obtain generative-model-  
 579 level causal values, we combine those reward predictions with SIM. These DM / DR aggregators  
 580 are used for the semi-synthetic generative-model values. They sit alongside the direct held-out-grid  
 581 metrics in the real cached benchmark.

582 Let  $\{x_i\}_{i=1}^{n_{\text{eval}}}$  be evaluation contexts and let  $m_{i,a,b} \sim p_{\text{sim}}(\cdot | x_i, a)$  denote SIM-generated mediators.  
 583 The default SIM plug-in estimator is

$$\hat{\mu}^{\text{DM}}(a) = \frac{1}{n_{\text{eval}}} \sum_{i=1}^{n_{\text{eval}}} \frac{1}{B} \sum_{b=1}^B \hat{r}(x_i, m_{i,a,b}).$$

584 This simply averages predicted rewards over SIM-generated outputs.

585 A doubly robust variant augments the plug-in estimator with a residual term computed from EXP.  
 586 Because EXP randomizes generative-model choice uniformly,  $p_{\text{EXP}}(a) = 1/|\mathcal{A}_{\text{EXP}}|$ . Define

$$\hat{q}(x, a) = \frac{1}{B_{\text{DR}}} \sum_{b=1}^{B_{\text{DR}}} \hat{r}(x, m_{a,b}), \quad m_{a,b} \sim p_{\text{sim}}(\cdot | x, a),$$

587 as the Monte Carlo approximation to the SIM-integrated reward. The estimator is

$$\hat{\mu}^{\text{DR}}(a) = \frac{1}{n_{\text{EXP}}} \sum_{j=1}^{n_{\text{EXP}}} \hat{q}(X_j, a) + \frac{1}{n_{\text{EXP}}} \sum_{j=1}^{n_{\text{EXP}}} \frac{\mathbb{I}\{A_j = a\}}{p_{\text{EXP}}(a)} (Y_j - \hat{r}(X_j, M_j)).$$

588 The first term is the SIM-based plug-in estimate. The second term corrects it using the randomized  
 589 residual observed in EXP. For reward models trained on EXP, we use cross-fitting: EXP is partitioned  
 590 into  $K_{\text{cf}}$  folds, and the reward model used in the residual term is trained on the remaining  $K_{\text{cf}} - 1$   
 591 folds.

592 **Role of DM / DR versus direct cached-benchmark metrics.** In the appendix semi-synthetic vali-  
 593 dation, we center held-out recommendation regret computed against the generator’s true interventional  
 594 rewards. The same runs also produce model-level value diagnostics such as the agent-level RMSE  
 595 of  $\hat{\mu}^{\text{DM}}(a)$  or  $\hat{\mu}^{\text{DR}}(a)$ , which remain useful for diagnosing reward-surface fit. In the real cached  
 596 benchmark, because the full held-out cached test grid is available, we additionally evaluate the fitted  
 597 reward surface directly on that grid. Let  $\mathcal{X}_{\text{test}}$  be the held-out test contexts and  $\mathcal{G}_{\text{test}} = \mathcal{X}_{\text{test}} \times \mathcal{A}$   
 598 the held-out cached grid. For each  $(x, a) \in \mathcal{G}_{\text{test}}$ , let  $q_{\text{cache}}(x, a)$  be the judged scalar reward of the  
 599 single cached held-out trajectory for context  $x$  and agent  $a$ , and let  $\hat{q}(x, a)$  be the fitted reward model  
 600 evaluated on that cached trajectory. Define the benchmark-level agent aggregates

$$\bar{q}_{\text{cache}}(a) = \frac{1}{|\mathcal{X}_{\text{test}}|} \sum_{x \in \mathcal{X}_{\text{test}}} q_{\text{cache}}(x, a), \quad \bar{q}_{\text{pred}}(a) = \frac{1}{|\mathcal{X}_{\text{test}}|} \sum_{x \in \mathcal{X}_{\text{test}}} \hat{q}(x, a).$$

601 The real-benchmark empirical metrics are then

$$\text{RMSE}_{xa} = \left( \frac{1}{|\mathcal{G}_{\text{test}}|} \sum_{(x,a) \in \mathcal{G}_{\text{test}}} (\hat{q}(x, a) - q_{\text{cache}}(x, a))^2 \right)^{1/2},$$

602

$$\text{RMSE}_{\text{agent}} = \left( \frac{1}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} (\bar{q}_{\text{pred}}(a) - \bar{q}_{\text{cache}}(a))^2 \right)^{1/2}, \quad \hat{\pi}(x) = \arg \max_{a \in \mathcal{A}} \hat{q}(x, a),$$

603

$$\text{Regret}_{\text{test}} = \frac{1}{|\mathcal{X}_{\text{test}}|} \sum_{x \in \mathcal{X}_{\text{test}}} \left[ \max_{a \in \mathcal{A}} q_{\text{cache}}(x, a) - q_{\text{cache}}(x, \hat{\pi}(x)) \right].$$

604 These real-benchmark metrics summarize the finite held-out cached grid. The natural next extension  
 605 is repeated cached replays / repeated judged trajectories per  $(x, a)$  so that the benchmark can estimate  
 606 latent  $q(x, a)$  more directly.

### 607 A.7 Empirical role of DM / DR on the controlled benchmark

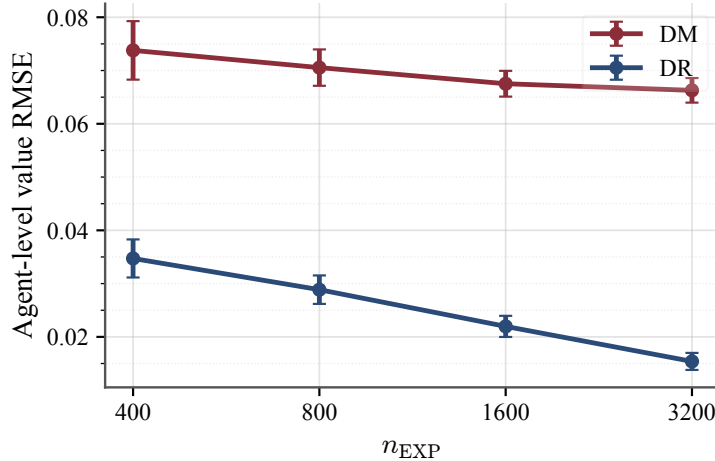
608 The same semi-synthetic runs also let us compare DM and DR on model-level value estimation.  
 609 These diagnostics answer a narrower question than the regret tables: when does randomized residual  
 610 correction improve the SIM plug-in value estimate?

611

Method	DM RMSE	DR RMSE	$\Delta$ (DR-DM)
EXP-Only	0.0303	0.0432	+0.0128
OBS-Only	0.0268	0.0409	+0.0141
CVCI	0.0212	0.0399	+0.0187
Representation-EXP	0.0263	0.0424	+0.0161

Table 9: Average agent-level value RMSE across the synthetic  $(\beta, n_{\text{OBS}}, n_{\text{EXP}})$  grid with  $\beta \in \{0, 0.2, 0.5, 0.8, 0.9, 0.99\}$ ,  $n_{\text{OBS}} \in \{2,000, 20,000\}$ , and  $n_{\text{EXP}} \in \{20, 100\}$ . Positive  $\Delta$  means that replacing DM by DR increases error.

612 Table 9 shows that, on this synthetic  $(\beta, n_{\text{OBS}}, n_{\text{EXP}})$  grid, the SIM plug-in estimator has lower  
 613 average agent-level value RMSE than DR for every representative method shown. Across these 24  
 614 settings, DR does not improve model-level value estimation on average.



615

Figure 4: Targeted DM/DR diagnostic under strong self-selection, intentional reward-model underfit, and increasing EXP budget with fixed OBS budget. In this regime, DR materially reduces agent-level value RMSE.

616 Figure 4 shows a different case, with stronger self-selection, intentional reward-model underfit, lower  
 617 outcome noise, and increasing EXP budget at fixed OBS budget. In that setting, DR lowers agent-level  
 618 RMSE from 0.0738 to 0.0347 at  $n_{\text{EXP}} = 400$ , and from 0.0663 to 0.0154 at  $n_{\text{EXP}} = 3200$ . This is  
 619 a narrower result about model-level value estimation; it does not change the main regret comparisons  
 620 in the real summarization and coding benchmarks.

621 **A.8 Empirical role of DM / DR on the real summarization benchmark**

622 The real summarization benchmark also allows a DM/DR comparison for model-level value estima-  
 623 tion. Here the target is the benchmark-level agent average over the held-out cached test contexts, and  
 624 DR uses the known uniform EXP propensity together with cross-fitting for methods trained on EXP.  
 625 This is a model-level value check built on top of the fitted reward model, not a replacement for the  
 626 regret comparisons in the main text.

Method	$n_{\text{EXP}}$	DM RMSE	DR RMSE	$\Delta$ (DR-DM)
OBS-Only	100	0.0479	0.0459	-0.0021
OBS-Only	200	0.0479	0.0410	-0.0070
OBS-Only	800	0.0479	0.0310	-0.0169
CVCI	100	0.0479	0.0724	+0.0245
CVCI	200	0.0479	0.0572	+0.0093
CVCI	800	0.0479	0.0325	-0.0154
Representation-EXP	100	0.0226	0.0749	+0.0522
Representation-EXP	200	0.0187	0.0634	+0.0447
Representation-EXP	800	0.0190	0.0298	+0.0108
EXP-Only	100	0.1092	0.2495	+0.1403
EXP-Only	200	0.0809	0.1528	+0.0718
EXP-Only	800	0.0604	0.0798	+0.0194

Table 10: Budget-specific DM/DR agent-level value RMSE on the real summarization benchmark at  $n_{\text{OBS}} = 20,000$ . Negative  $\Delta$  means that replacing DM by DR reduces error.

628 Table 10 makes the heterogeneity explicit. OBS-Only improves in 23 of the 30 benchmark cells  
 629 and becomes more favorable as  $n_{\text{EXP}}$  grows. CVCI changes sign only at  $n_{\text{EXP}} = 800$ , where the  
 630 residual term is large enough to help more than it hurts. The EXP-side methods behave differently:  
 631 Representation-EXP worsens in 28 of 30 cells, and EXP-Only worsens in 27 of 30.

632 The natural interpretation is that DR is correcting the remaining OBS-side bias in model-level values  
 633 rather than uniformly improving all fitted reward models. On this benchmark, that makes it most  
 634 useful for the OBS-heavy baseline, conditionally useful for CVCI once the EXP budget is larger, and  
 635 generally unhelpful for the EXP-side methods whose remaining error is not well summarized by a  
 636 simple EXP residual correction.

637 **A.9 Additional summarization diagnostics**

638 **Reward maps and routing.** The summarization benchmark fixes one judged cached summary  
 639 for each article-model pair and evaluates two reward maps on that same cache. Each cached  
 640 summary carries a rubric vector  $s(x, a) = (s_{\text{faith}}(x, a), s_{\text{cov}}(x, a), s_{\text{clar}}(x, a), s_{\text{conc}}(x, a)) \in [0, 1]^4$   
 641 for faithfulness, coverage, clarity, and conciseness, together with an unsupported-claims count  $u(x, a)$ .  
 642 Let  $\mathcal{T}$  denote the four user segments and let  $w_\tau \in \Delta^4$  be the segment-specific weight vector. For  
 643 sharpening exponent  $\gamma \geq 1$  and unsupported-claims penalty  $\lambda \geq 0$ , define

$$\bar{w}_{\tau k}^{(\gamma)} = \frac{w_{\tau k}^\gamma}{\sum_{\ell=1}^4 w_{\tau \ell}^\gamma}, \quad \ell_\tau^{(\gamma, \lambda)}(x, a) = \sum_{k=1}^4 \bar{w}_{\tau k}^{(\gamma)} s_k(x, a) - \lambda \frac{u(x, a)}{20}.$$

644 The smooth map is

$$g_\tau^{\text{sm}}(x, a) = \ell_\tau^{(1, 0)}(x, a) = \langle w_\tau, s(x, a) \rangle, \quad q_{\text{cache}}^{\text{sm}}(x, a) = \frac{1}{|\mathcal{T}|} \sum_{\tau \in \mathcal{T}} g_\tau^{\text{sm}}(x, a).$$

645 With  $\sigma(t) = (1 + e^{-t})^{-1}$  and  $[t]_{[0, 1]} = \min\{1, \max\{0, t\}\}$ , the sharpened map is

$$g_\tau^{\text{sh}}(x, a) = \sigma\left(\frac{[\ell_\tau^{(16, 0.1)}(x, a)]_{[0, 1]} - 0.9}{0.02}\right), \quad q_{\text{cache}}^{\text{sh}}(x, a) = \frac{1}{|\mathcal{T}|} \sum_{\tau \in \mathcal{T}} g_\tau^{\text{sh}}(x, a).$$

646 OBS and EXP are sampled only from training contexts. User segments are drawn uniformly from  $\mathcal{T}$ ,  
 647 EXP randomizes uniformly over  $\mathcal{A}$ , and OBS uses the  $\beta$ -indexed router

$$p_{\text{OBS}}^\beta(a | x, \tau) = (1 - \varepsilon) \frac{\exp(\eta_a(x) + \beta \xi_a(\tau))}{\sum_{a' \in \mathcal{A}} \exp(\eta_{a'}(x) + \beta \xi_{a'}(\tau))} + \frac{\varepsilon}{|\mathcal{A}|},$$

Method	Avg. rank	Top-3	Excess (%)
<b>CVCI</b>	<b>2.04</b>	<b>43</b>	<b>4.8</b>
OBS-Only	2.83	32	7.9
EXP-Only	2.90	32	8.8
Grounded	3.21	26	7.9
CVCI-Residual	4.40	9	24.0
Representation-EXP	5.62	2	77.6

Table 11: Overall recommendation regret in the cached summarization benchmark. *Avg. rank* is the average rank across the 48 settings; *Top-3* is the count of settings in which a method ranked among the top three; *Excess (%)* is the mean percentage regret increase relative to the best method in each setting.

Average regret over the four budget settings			
Method	Smooth	Sharpened	$\Delta_{\text{shape}}$
CVCI	0.0126	0.0423	+0.0297
OBS-Only	0.0125	0.0477	+0.0352
EXP-Only	0.0147	0.0376	+0.0229
Representation-EXP	0.0267	0.0773	+0.0507
Held-out target-fit $R^2$			
Quantity	Smooth	Sharpened	$\Delta$
Proxy-bundle $R^2$	0.161	0.171	+0.010
Rubric-linear $R^2$	1.000	0.779	-0.221

Table 12: Recommendation regret and held-out target-fit  $R^2$  under smooth versus sharpened reward maps, with cached outputs held fixed. Sharpening the reward inflates regret for every estimator family and drops rubric-linear  $R^2$ , while proxy-bundle  $R^2$  stays low.

648 where  $\eta_a(x) = \langle \theta_a^X, f_{\text{ctx}}(x) \rangle - \kappa c_a$  is an article-dependent logit and  $\xi_a(\tau) = \langle \theta_a^U, e(\tau) \rangle$  is a  
649 segment-affinity logit. Thus  $\beta = 0$  removes the latent-user term, and larger  $\beta$  makes routing more  
650 segment-specific.

651 In the cached summarization benchmark, CVCI has the best average rank, top-3 coverage, and excess  
652 regret across 30 seeds and 48 settings. The gap to the next family is small relative to the spread of  
653 cell-level standard errors reported in Appendix A.8, so this result reads as family-level evidence  
654 within this benchmark instantiation.

655 In Table 12, the two held-out  $R^2$  values are out-of-sample linear fits of the cached target surface  
656 using the auxiliary proxy representation and the raw rubric vector  $s(x, a)$ , respectively. This table  
657 changes the target reward while holding the cached outputs fixed, so the regret increase and the drop  
658 in rubric-linear fit diagnose target alignment rather than a formal oracle condition.

659 Table 13 reports model-level value RMSE for DM and DR. In this benchmark, DR mainly helps  
660 OBS-heavy methods, while DM remains stronger for the other families. Appendix A.8 gives the  
661 budget-specific breakdown and implementation details.

Method	DM RMSE	DR RMSE	$\Delta$ (DR-DM)
OBS-Only	0.0479	<b>0.0393</b>	-0.0086
CVCI	0.0479	0.0540	+0.0061
Representation-EXP	<b>0.0201</b>	0.0560	+0.0359
EXP-Only	0.0835	0.1607	+0.0772

Table 13: Model-level value RMSE for DM and DR, averaged at  $n_{\text{OBS}} = 20,000$  and  $n_{\text{EXP}} \in \{100, 200, 800\}$  over 30 seeds. Negative  $\Delta$  means DR helps over DM; positive  $\Delta$  means DR hurts.

662 **A.10 Exact real benchmark protocol**

663 **Dataset and candidate generative models.** Contexts  $X$  are CNN/DailyMail articles [Hermann  
664 et al., 2015, See et al., 2017]. The candidate set contains 20 summarization systems: 16 prompt-based  
665 abstractive systems built from two Gemma 3 instruction models [Gemma Team, 2025], together with  
666 four extractive baselines. For each context–generative-model pair, the cache stores one summary  
667 together with its judged metadata.

668 **Judged rewards and proxy features.** For each cached output, the judge returns rubric-style signals  
669 that are combined into the scalar reward used throughout the real benchmark comparisons. The  
670 real benchmark instead uses a heuristic auxiliary-feature channel derived from cached trajectories  
671 and judged-output metadata. That channel is best read as an OBS proxy, distinct from the hidden  
672 representation  $\phi^*$  used in the semi-synthetic study.

673 **Fixed benchmark protocol.** The real-data study fixes a 48-context subset of relatively difficult  
674 CNN/DailyMail articles, selected once by an article-level difficulty score. Within this fixed pool,  
675 each seed performs an 80/20 context-level train/test split, uses no separate context-level validation  
676 split, and tunes hyperparameters only on a held-out fraction of EXP when the estimator requires it.  
677 OBS and EXP are sampled only from the training contexts. The observational-choice protocol is a  
678 fixed paired-sampling design under a self-selection mechanism that routes each decision through a  
679 five-system candidate slate before the final model choice.

680 **Budget views and aggregation.** The real summarization study uses three budget views. The  
681 main-text regret comparisons use

$$\beta \in \{0, 0.2, 0.5, 0.8, 0.9, 0.99\}, \quad n_{\text{OBS}} \in \{2,000, 20,000\}, \quad n_{\text{EXP}} \in \{20, 200\},$$

682 with matched smooth and sharpened reward maps, giving the Section 5.4 summaries over six  $\beta$   
683 values, two OBS budgets, two EXP budgets, and two reward maps. The DM/DR value diagnostics in  
684 Table 10 instead fix  $n_{\text{OBS}} = 20,000$  and use  $n_{\text{EXP}} \in \{100, 200, 800\}$ . The winner-map table below  
685 is a separate smooth-reward  $(\beta, n_{\text{OBS}}, n_{\text{EXP}})$  comparison with

$$\beta \in \{0, 0.2, 0.5, 0.8, 0.9, 0.99\}, \quad n_{\text{OBS}} \in \{2,000, 20,000\}, \quad n_{\text{EXP}} \in \{20, 100\},$$

686 reported to show the full 24-cell winner map. All reported means and standard errors are recomputed  
687 from seed-level paired runs, with 30 seeds in every reported  $(\beta, n_{\text{OBS}}, n_{\text{EXP}}, \text{method})$  cell.

688 **A.11 Tuning on EXP: model-level cross-validation, EXP holdout, and defaults**

689 We use three EXP-side tuning strategies. Model-level EXP cross-validation partitions the generative  
690 models *appearing in the sampled EXP data for that seed/cell* and validates at the held-out model  
691 level over that sampled EXP model set. EXP holdout reserves a validation subset from the sampled  
692  $n_{\text{EXP}}$  budget itself and tunes on that holdout. Fixed defaults use dataset-configured settings with no  
693 extra EXP hyperparameter sweep.

694 More precisely,  $n_{\text{EXP}}$  is the total EXP budget for a given seed/cell. If a method uses an explicit EXP  
695 holdout, that holdout is carved from the sampled EXP budget itself; no extra EXP rows are sampled  
696 outside budget, and no separate context-level validation split is required for this.

697 For **Grounded-**, **CVCI-**, and **CVCI-Residual**-style families tuned by model-level EXP cross-validation,  
698 let  $\mathcal{A}_{\text{EXP}}$  be the set of generative models appearing in the sampled EXP data and partition it into  $K_{\text{CV}}$   
699 folds  $\{\mathcal{A}^{(k)}\}_{k=1}^{K_{\text{CV}}}$ . The corresponding validation indices are

$$\mathcal{I}_{\text{val}}^{(k)} = \{j \in \text{EXP} : A_j \in \mathcal{A}^{(k)}\}, \quad \mathcal{I}_{\text{tr}}^{(k)} = \text{EXP} \setminus \mathcal{I}_{\text{val}}^{(k)}.$$

700 For a candidate tuning parameter  $\eta$  (for example,  $\eta = \alpha_{\text{corr}}$  for **Grounded** or  $\eta = \lambda$  for **CVCI**), fit the  
701 corresponding estimator using all OBS data together with the EXP indices in  $\mathcal{I}_{\text{tr}}^{(k)}$ . Then evaluate  
702 held-out generative-model means with

$$\mathcal{L}^{(k)}(\eta) = \sum_{a \in \mathcal{A}^{(k)}} \omega_{k,a} \left( \bar{r}_{\eta}^{(k)}(a) - \bar{Y}^{(k)}(a) \right)^2, \quad (12)$$

703 where  $\widehat{r}_\eta^{(k)}(a) = \frac{1}{n_{k,a}} \sum_{j \in \mathcal{I}_{\text{val}}^{(k)}: A_j=a} \widehat{r}_\eta^{(-k)}(z_j)$ ,  $\overline{Y}^{(k)}(a) = \frac{1}{n_{k,a}} \sum_{j \in \mathcal{I}_{\text{val}}^{(k)}: A_j=a} Y_j$ , and  $n_{k,a}$  is  
 704 the number of held-out EXP examples for generative model  $a$  in fold  $k$ . We use count weights  
 705  $\omega_{k,a} \propto n_{k,a}$ , normalized so that  $\sum_{a \in \mathcal{A}^{(k)}} \omega_{k,a} = 1$ , and select the parameter that minimizes the  
 706 average fold loss. If too few generative models appear in the sampled EXP data to sustain the  
 707 requested number of model-level folds, the procedure falls back to sample-level cross-validation and  
 708 records the effective mode.

709 In the real main sweep, **Grounded**, **CVCI**, and **CVCI-Residual** use model-level EXP cross-validation.  
 710 The simpler historical linear grounded baseline instead uses an EXP holdout through correction- $\alpha$   
 711 tuning on the held-out EXP subset. **EXP-Only**, **OBS-Only**, and **Representation-EXP** use fixed  
 712 defaults.

713 When candidate hyperparameters are tied within numerical tolerance, the implementation resolves  
 714 ties toward the more regularized / less aggressive solution. For residual **CVCI-Residual** this means  
 715 preferring larger residual ridge penalties and, if still tied, larger OBS pooling weights. For holdout-  
 716 tuned grounded corrections, ties in the correction- $\alpha$  grid are resolved toward smaller shrinkage  
 717 factors.

718 Because the `news_hard` top-48 real benchmark contains only 48 contexts, each 80/20 context  
 719 split leaves roughly 10 held-out test contexts per seed. Single-seed regret is therefore noisy, and all  
 720 real-benchmark conclusions are aggregated across many seeds.

## 721 A.12 Benchmark setup details

722 Each benchmark setting specifies a training-context distribution, a held-out test set, an action set, a  
 723 simulator, an experimental reward map, and an observational logging family. Within one such cell,  
 724 the observational law factors as

$$p_{\text{OBS}}(x, u, a, m, y) = P_{\text{tr}}(x) P(u | x) p_{\text{OBS}}^\beta(a | x, u) p_{\text{sim}}(m | x, a) P_Y(y | x, m, u),$$

725 and the experimental law replaces the choice model by an unconfounded randomization rule  $\pi_{\text{EXP}}(a |$   
 726  $x)$ :

$$p_{\text{EXP}}(x, u, a, m, y) = P_{\text{tr}}(x) P(u | x) \pi_{\text{EXP}}(a | x) p_{\text{sim}}(m | x, a) P_Y(y | x, m, u).$$

727 Here  $P_Y$  is the benchmark-specific outcome law. In the main-text benchmark settings,  $\pi_{\text{EXP}}(a |$   
 728  $x) = 1/|\mathcal{A}|$ , and  $r^*(x, m) = \mathbb{E}[Y | X = x, M = m, \text{EXP}]$  integrates out  $U | X = x$  under  $p_{\text{EXP}}$ .  
 729 The identified target remains

$$q(x, a) = \mathbb{E}_{M \sim p_{\text{sim}}(\cdot | x, a)}[r^*(x, M)].$$

730 Each estimator returns  $\widehat{q}$ , inducing  $\widehat{\pi}(x) = \arg \max_{a \in \mathcal{A}} \widehat{q}(x, a)$ , and the main metric is held-out  
 731 recommendation regret. The  $\beta$ -indexed observational choice law is benchmark-specific. In the  
 732 appendix controlled validation,  $p_{\text{OBS}}^\beta(a | x, u) = (1 - \beta)\pi_X(a | x) + \beta\pi_U(a | u)$ , where  $\pi_X$  is a  
 733 context-driven policy and  $\pi_U$  is a latent-user policy. Thus  $\beta = 0$  removes latent-user routing, whereas  
 734  $\beta = 1$  uses only the latent-user component. Section 5.4 instead uses a softmax router in which  $\beta$   
 735 scales the latent-user logit contribution.

736 The real cached benchmarks fix one realized mediator  $m_{\text{cache}}(x, a)$  for each context-action pair and  
 737 report the finite cached target  $q_{\text{cache}}(x, a) = r^*(x, m_{\text{cache}}(x, a))$ . Their reported recommendation  
 738 regret therefore replaces  $q$  by  $q_{\text{cache}}$  on the held-out cache. Randomness then comes from OBS/EXP  
 739 sampling and the benchmark’s latent draws, not from regenerating mediators at evaluation time.

## 740 A.13 Semi-synthetic controlled validation

741 The semi-synthetic validation keeps the summarization task family but replaces judged rewards by a  
 742 known latent generator, so held-out regret can be evaluated directly against the true interventional  
 743 surface. Each context–mediator pair is mapped to an unobserved representation

$$\phi^*(x, m) = H\varphi(x, m) \in \mathbb{R}^{d_\phi},$$

744 where  $H$  is a fixed signed-hash projection of sparse hashed text features  $\varphi(X, M) \in \mathbb{R}^{2^{18}}$ . The  
 745 benchmark-specific outcome law is

$$P_Y(\cdot | x, m, u) = \mathcal{L}(\text{sigmoid}((w^*)^\top \phi^*(x, m) + \lambda^\top u + \varepsilon)), \quad \varepsilon \sim \mathcal{N}(0, \sigma_Y^2),$$

Table 14: Benchmarks used in the main text. The two cached benchmarks evaluate finite cached targets  $q_{\text{cache}}$  rather than the latent superpopulation target  $q$ ; their OBS and EXP samples are constructed by resampling and reweighting the cached pool under specified routing and randomization rules.

Benchmark	Source data	Cached object	Evaluation target
Controlled synthetic	semi-CNN/DailyMail article family with a known latent reward generator	SIM-generated summaries from the generator pool	True interventional reward surface $q$ (theorem-level target).
Real-task cached summarization	48 difficult CNN/DailyMail articles; 20 candidate models	One judged summary per article and model	Finite cached target $q_{\text{cache}}^{\text{sm}}$ (or $q_{\text{cache}}^{\text{sh}}$ ): segment-average rubric score over four user types under the smooth (or sharpened) reward map.
Real-task cached coding	SWE-bench Verified issues with a BouncerBench patch pool	One cached candidate patch per issue and agent	Finite cached target $q_{\text{cache}}$ : user-average utility over fix success ( $c_1$ ) and patch-quality components ( $c_2, c_3, c_4$ ).

746 SO

$$r^*(x, m) = \mathbb{E}[Y \mid X = x, M = m, \text{EXP}] = \mathbb{E}_{U \sim P(\cdot | x), \varepsilon} \left[ \text{sigmoid}((w^*)^\top \phi^*(x, m) + \lambda^\top U + \varepsilon) \right].$$

747 OBS and EXP share this same  $P_Y$ , so the only confounding channel is model choice through  
748  $p_{\text{OBS}}^\beta(a \mid x, u)$ , while EXP randomizes uniformly over the candidate agent set. This appendix  
749 validation varies only  $\beta$ ,  $n_{\text{OBS}}$ , and  $n_{\text{EXP}}$ , so changes in the regret winner can be attributed to  
750 confounding strength and data budget.

751 **Implementation details.** The semi-synthetic generator uses sparse hashed text features  $\varphi(X, M) \in$   
752  $\mathbb{R}^{2^{18}}$  for all reward models. The hidden representation has dimension  $d_\phi = 8$  and is constructed by  
753 a seeded signed-hash projection; it is fixed across replicates and is not observed by any estimator.  
754 When hybrid estimators use the proxy map  $\psi(z)$  from Section 4, the PCA dimension is  $d_\psi = 20$ . The  
755 context-only embedding  $\psi_X(X)$  that appears in the observational choice model is a separate object. It  
756 is fit on raw article text alone, with no context–output pairs. Operationally, we sample 5,000 contexts  
757 from the observational split, fit a TF–IDF vectorizer with English stop-word removal, unigram–  
758 bigram features, and a vocabulary cap of 50,000 terms, and then apply rank-50 truncated SVD. Thus  
759  $\psi_X(X) \in \mathbb{R}^{50}$  is a low-dimensional context representation used only for model choice and for the  
760  $U \mid X$  model in the semi-synthetic generator. This is distinct from the proxy representation  $\psi(X, M)$   
761 in Section 4, which is learned from OBS auxiliary labels and is used by the hybrid reward estimators.  
762 Reference values and held-out regret metrics are computed on  $n_{\text{true}} = 100$  and  $n_{\text{eval}} = 100$  held-out  
763 contexts respectively, using  $B = 5$  simulator draws per context–agent pair.

764 **Fixed protocol in the appendix controlled validation.** The appendix semi-synthetic validation  
765 fixes one confounding protocol and varies only the self-selection weight  $\beta$ , the observational budget  
766  $n_{\text{OBS}}$ , and the experimental budget  $n_{\text{EXP}}$ . EXP randomizes uniformly over the candidate agent set.  
767 OBS mixes a context-only softmax policy with a latent-user shortlist policy, where  $\beta$  is the probability  
768 of drawing from the confounded latent-user component. This validation uses

$$\beta \in \{0, 0.2, 0.5, 0.8, 0.9, 0.99\}, \quad n_{\text{OBS}} \in \{2,000, 20,000\}, \quad n_{\text{EXP}} \in \{20, 100\},$$

769 and reports 30 paired seeds in every available cell. The direct latent effect on the reward score and  
770 the auxiliary-label noise level are both fixed throughout this validation grid.

771 **Budget-axis repetition for pure baselines.** OBS-Only ignores EXP labels, so its reported value is  
772 constant across the two EXP columns for a fixed  $(\beta, n_{\text{OBS}})$  pair. EXP-Only ignores OBS outcomes,  
773 so its reported value is constant across the two OBS rows for a fixed  $(\beta, n_{\text{EXP}})$  pair. This repetition  
774 is carried through consistently in the synthetic regime table and summary table.

#### 775 A.14 Appendix results for the summarization benchmark

776 **Synthetic summarization results.** Table 15 shows the regret winner in each synthetic summariza-  
777 tion setting over six  $\beta$  values and four  $(n_{\text{OBS}}, n_{\text{EXP}})$  budget settings.

778 The smallest-budget row (2,000, 20) is the most mixed: **CVCI-Residual** wins once, **CVCI** wins twice,  
779 **OBS-Only** wins once, and **Grounded** wins twice. At (20,000, 20), **CVCI** wins all six beta values. At

Budget cell	$\beta = 0$	$\beta = 0.2$	$\beta = 0.5$	$\beta = 0.8$	$\beta = 0.9$	$\beta = 0.99$
$n_{OBS} = 2,000$ $n_{EXP} = 20$	CVCI-Residual +0.0004	CVCI +0.0002	CVCI +0.0009	OBS-Only +0.0004	Grounded +0.0010	Grounded +0.0004
$n_{OBS} = 20,000$ $n_{EXP} = 20$	CVCI +0.0005	CVCI +0.0015	CVCI +0.0015	CVCI +0.0021	CVCI +0.0018	CVCI +0.0017
$n_{OBS} = 2,000$ $n_{EXP} = 100$	CVCI +0.0008	CVCI +0.0017	CVCI +0.0012	OBS-Only +0.0006	OBS-Only +0.0008	OBS-Only +0.0005
$n_{OBS} = 20,000$ $n_{EXP} = 100$	CVCI +0.0002	EXP-Only +0.0012	EXP-Only +0.0016	EXP-Only +0.0013	EXP-Only +0.0012	EXP-Only +0.0016

Table 15: Regret-optimal estimator on the semi-synthetic summarization benchmark. The number beneath each winner is the regret gap to the runner-up.

Budget cell	$\beta = 0$	$\beta = 0.2$	$\beta = 0.5$	$\beta = 0.8$	$\beta = 0.9$	$\beta = 0.99$
$n_{OBS} = 2,000$ $n_{EXP} = 20$	Grounded +0.0002	Grounded +0.0011	Grounded +0.0036	OBS-Only +0.0000	OBS-Only +0.0008	Grounded +0.0006
$n_{OBS} = 20,000$ $n_{EXP} = 20$	Grounded +0.0016	Grounded +0.0049	EXP-Only +0.0013	EXP-Only +0.0062	EXP-Only +0.0019	EXP-Only +0.0019
$n_{OBS} = 2,000$ $n_{EXP} = 100$	OBS-Only +0.0007	OBS-Only +0.0002	OBS-Only +0.0002	CVCI +0.0000	CVCI +0.0000	CVCI +0.0008
$n_{OBS} = 20,000$ $n_{EXP} = 100$	CVCI-Residual +0.0048	CVCI +0.0009	CVCI +0.0052	CVCI +0.0012	CVCI +0.0013	CVCI +0.0046

Table 16: Regret-optimal estimator on the real summarization benchmark. The number beneath each winner is the regret gap to the runner-up.

780 (2,000, 100), the winner shifts from CVCI at lower  $\beta$  to OBS-Only at higher  $\beta$ . At (20,000, 100),  
781 EXP-Only wins all but the  $\beta = 0$  cell.

782 Across the 24 settings, CVCI is the most stable overall winner, but no single method wins everywhere.  
783 Most regret gaps are small, so Table 15 is best read as a map of where the winner changes.

784 **Additional synthetic aggregate interpretation.** Across the same 24 settings, CVCI is also the  
785 aggregate regret leader. EXP-Only and OBS-Only form the next tier, while Grounded and CVCI-  
786 Residual win only a few cells. The ranking diagnostics are flatter than regret: no single method leads  
787 top-1 or top-3 across the grid.

788 **Real summarization results on the smooth auxiliary comparison.** Table 16 shows the regret  
789 winner in each real summarization setting for the auxiliary smooth-reward comparison over six  $\beta$   
790 values and four  $(n_{OBS}, n_{EXP})$  budget settings.

791 At (2,000, 20), Grounded wins four beta values and OBS-Only wins two; the  $\beta = 0.8$  cell is  
792 essentially a tie. At (20,000, 20), Grounded wins at low  $\beta$ , while EXP-Only wins at higher  $\beta$ . Once  
793  $n_{EXP} = 100$ , the high-OBS row (20,000, 100) is dominated by CVCI, with a single CVCI-Residual  
794 win at  $\beta = 0$ , while the low-OBS row (2,000, 100) shifts from OBS-Only at lower  $\beta$  to CVCI at  
795 higher  $\beta$ .

796 **Additional real ranking-regret comparison.** Across these 24 settings, CVCI is the aggregate  
797 regret leader and Grounded is second. Ranking tells a different story: Grounded wins top-1 and top-3  
798 in all 24 settings but wins regret in only 6.

## 799 A.15 Public cached coding benchmark

800 **Fixed protocol.** The coding benchmark uses a public cached grid of issue descriptions, patch  
801 mediators, and execution-derived rewards. The coding comparison in Section 5.4 fixes

$$\beta = 0.5, \quad n_{OBS} = 2,000, \quad n_{EXP} = 100,$$

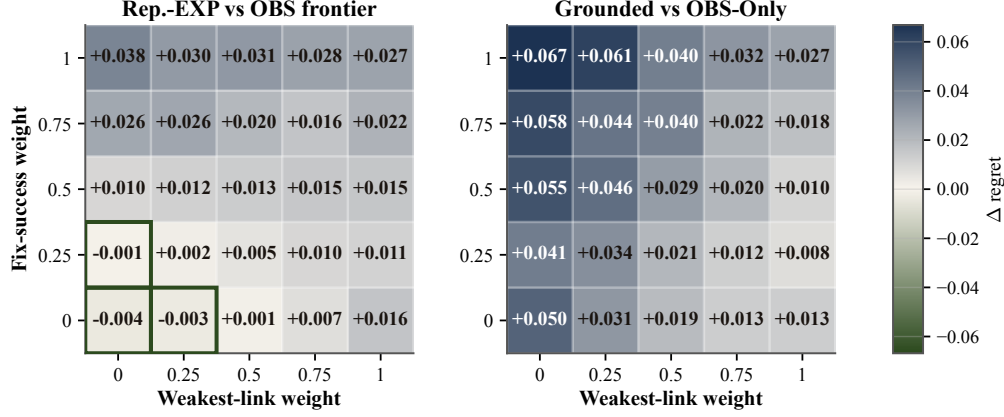


Figure 5: Pairwise regret differences  $\Delta \text{regret} = \text{regret}(\text{family A}) - \text{regret}(\text{family B})$  across the  $5 \times 5$  grid of  $(\alpha_{\text{fix}}, \omega_{\text{weak}})$  in the cached coding benchmark instantiation; negative values mean family A has lower regret than family B in that cell. Left panel: Representation-EXP minus the best OBS-based family. Right panel: **Grounded** minus OBS-Only.

802 uses outcome-precision 64, and evaluates only the feature-based implementations of the estimator  
 803 families. It then varies two reward-design parameters:

$$\text{fix-success weight} \in \{0, 0.25, 0.5, 0.75, 1\}, \quad \text{weakest-link weight} \in \{0, 0.25, 0.5, 0.75, 1\},$$

804 with 10 paired seeds in every cell. The first parameter rescales the fix-success component of the  
 805 reward, while the second moves the non-fix patch-quality score from additive aggregation toward a  
 806 weakest-link rule.

807 **Reward construction and target-alignment diagnostics.** Each cached patch is scored by four  
 808 components  $c_1, c_2, c_3, c_4$ , where  $c_1$  records true fix success and  $c_2, c_3, c_4$  summarize patch size,  
 809 locality, and issue focus. These patch-quality components can disagree with true fix success, which is  
 810 why changing the reward definition matters here.

811 For user type  $u$  with weights  $w_u = (w_{u1}, w_{u2}, w_{u3}, w_{u4})$ , normalized non-success weights are  
 812  $\bar{w}_{uk} = w_{uk} / (w_{u2} + w_{u3} + w_{u4})$  for  $k \in \{2, 3, 4\}$ , and

$$g_u(c_2, c_3, c_4) = (1 - \omega_{\text{weak}}) \sum_{k=2}^4 \bar{w}_{uk} c_k + \omega_{\text{weak}} \min\{c_2, c_3, c_4\}.$$

813 The benchmark utility for user type  $u$  is

$$Y_u(x, a) = w_{u1} \alpha_{\text{fix}} c_1(x, a) + (1 - w_{u1}) g_u(c_2(x, a), c_3(x, a), c_4(x, a)),$$

814 and the evaluation target is the cached user-average surface  $q_{\text{cache}}(x, a) = \mathbb{E}[Y_U(x, a) \mid X =$   
 815  $x]$ . Representation-EXP maps each issue-patch pair into an auxiliary representation learned from  
 816 observational patch statistics and fits EXP rewards on that representation.

817 With the data budget fixed, rankings change because the reward puts more or less weight on true fix  
 818 success. At  $\omega_{\text{weak}} = 0.25$ , the regret gap between Representation-EXP and the better of OBS-Only  
 819 and CVCI moves from  $-0.0034$  at  $\alpha_{\text{fix}} = 0$ , to  $+0.0024$  at  $\alpha_{\text{fix}} = 0.25$ , to  $+0.0375$  at  $\alpha_{\text{fix}} = 1$ .  
 820 This sweep tests whether auxiliary patch statistics retain the target signal as the target moves from  
 821 patch quality toward true fix success.

822 Held-out  $R^2$  comes from linear fits of the cached target using either auxiliary patch statistics or  $c_1$   
 823 alone. On the additive slice, the auxiliary-feature fit drops from  $R^2 = 0.347$  at  $\alpha_{\text{fix}} = 0$  to  $0.081$  at  
 824  $\alpha_{\text{fix}} = 1$ , while the  $c_1$ -only fit rises from essentially zero to  $R^2 = 0.967$ . The  $R^2$  shift matches the  
 825 change in method ranking: Representation-EXP helps when the target follows patch-quality features,  
 826 but it loses that advantage once the reward is driven mainly by fix success.

827 Changing  $\omega_{\text{weak}}$  changes only how the non-success components are combined. Its largest effect is  
 828 on **Grounded**: at zero fix-success weight, the grounded-minus-OBS regret gap shrinks from about

Fix-success weight/weakest-link weight	0	0.25	0.5	0.75	1.0
0.00	Representation-EXP +0.0039	Representation-EXP +0.0033	CVCI-Residual +0.0001	OBS-Only +0.0001	OBS-Only +0.0002
0.25	Representation-EXP +0.0006	OBS-Only +0.0004	CVCI +0.0004	CVCI +0.0013	CVCI +0.0009
0.50	OBS-Only +0.0004	OBS-Only +0.0010	CVCI +0.0003	CVCI +0.0008	CVCI +0.0004
0.75	CVCI +0.0013	OBS-Only +0.0029	OBS-Only +0.0004	OBS-Only +0.0003	OBS-Only +0.0013
1.00	OBS-Only +0.0011	CVCI +0.0007	CVCI +0.0001	OBS-Only +0.0068	OBS-Only +0.0056

Table 17: Regret-optimal estimator on the coding benchmark at fixed  $n_{\text{OBS}} = 2,000$ ,  $n_{\text{EXP}} = 100$ ,  $\beta = 0.5$ , and reward-noise precision 64. Columns give the weakest-link weight in the non-fix patch-quality score, and rows give the weight on true fix success. Most top-two gaps are small, so the table mainly shows where the regret winner changes.

829 +0.0495 under additive scoring to about +0.0126 under the pure weakest-link rule. Figure 5 (right)  
830 shows the same effect. Even so, **Grounded** remains worse than the best OBS-based method throughout  
831 this fixed-budget setting.

832 **Appendix coding results.** In this coding appendix comparison, we treat a winner change as  
833 substantive when the top-two regret gap exceeds  $3 \times 10^{-3}$ ; smaller gaps are near-ties. Table 17  
834 reports the regret winner in each cell. Every top-two gap is below 0.01, and 20 of the 25 cells  
835 are below  $1.5 \times 10^{-3}$ . Most flips between OBS-Only and CVCI therefore reflect near-ties in the  
836 high-fix-success part of the grid. The clearest differences appear at low fix-success weight under  
837 additive scoring, where Representation-EXP outperforms the OBS-based methods, and in the steady  
838 narrowing of the grounded-versus-OBS gap as weakest-link scoring becomes stronger.

839 **Secondary axes omitted from the main text.** Additional coding comparisons over confounding  
840 strength and outcome noise were qualitatively weaker. They changed absolute difficulty, but they did  
841 not produce ranking shifts as clear as those induced by the fix-success weight and weakest-link weight  
842 in this fixed-budget probe. For that reason, the main text centers this reward-design comparison and  
843 leaves the other axes to this appendix discussion.

#### 844 A.16 Extended positioning, scope, and limitations

845 **Off-policy evaluation.** Direct Method with Replay (DM) and Doubly Robust (DR) are standard  
846 estimators in contextual-bandit and policy-evaluation settings [Dudík et al., 2011, 2014, Swaminathan  
847 and Joachims, 2015, Jiang and Li, 2016, Thomas and Brunskill, 2016], and they descend from the  
848 foundational doubly robust and missing-data lineage of Robins et al. [1994] and Bang and Robins  
849 [2005]. The OPE setting these papers analyze starts from logged action–outcome pairs in which  
850 outcomes are assumed unconfounded given context: the central technical issue is action coverage  
851 under the logging policy, and the recorded outcomes faithfully represent the causal target. In our  
852 setting, even for observed actions, the recorded outcomes are confounded because user-side factors  
853 influence both who selects which model and how its output is judged. Because SIM can regenerate  
854 outputs for any model on any held-out context, DM and DR enter only after the reward surface itself  
855 has been estimated. We use SIM and EXP to identify the causal scoring rule and then use OBS to  
856 reduce variance.

857 **Experimental grounding and hybrid estimation.** The closest methodological antecedents are  
858 experimental grounding [Kallus et al., 2018], data-fusion and transportability formalisms [Pearl and  
859 Bareinboim, 2011, Bareinboim and Pearl, 2016], and trial-to-target generalizability [Cole and Stuart,  
860 2010]. A related line of work studies how to combine large biased samples with smaller randomized  
861 ones [Rosenman et al., 2023, Cheng and Cai, 2021, Lin et al., 2025, Yang et al., 2025, Colnet et al.,  
862 2024]. That literature asks when observational signal can be used to reduce variance without giving  
863 up the causal credibility of the randomized sample; in those papers, the unit-level outcome is observed  
864 on both samples and the modeling work is done on a fixed outcome label. The estimator families in  
865 Section 4 are concrete instantiations of the same bias–variance question. Our setting differs in one  
866 specific way: the target outcome is not directly observed even on the randomized sample because the

867 unit “output” must first be regenerated by SIM and then scored. Identification in our setting depends  
868 on SIM validity (Assumption A2) together with EXP randomization.

869 **Summarization and LLM-based evaluation.** Our experiments use CNN/DailyMail contexts  
870 [Hermann et al., 2015, Nallapati et al., 2016, See et al., 2017, Stiennon et al., 2020] and rely on  
871 LLM-judged rubric scores in the pool-based study. Recent LLM-evaluation work covers benchmark  
872 design [Liang et al., 2023], judge reliability and pairwise preference platforms [Zheng et al., 2023,  
873 Chiang et al., 2024], and real-world usage logs [Zheng et al., 2024, Zhao et al., 2024]. It also  
874 documents biases in automatic evaluators, including positional unfairness [Wang et al., 2023], length  
875 bias [Dubois et al., 2024], and self-preference [Panickssery et al., 2024], and studies panel-based  
876 mitigations [Verga et al., 2024]. Preference data is central in alignment pipelines, both with explicit  
877 learned reward models [Christiano et al., 2017, Stiennon et al., 2020, Ouyang et al., 2022] and with  
878 implicit reward formulations such as direct preference optimization [Rafailov et al., 2023]. The  
879 LLM-evaluation work cited above primarily treats judged scores as outcome data and studies how  
880 reliable those scores are and how to average them. We use cached outputs and judged scores as inputs  
881 to a causal evaluation problem; they become the causal target only after EXP-identified scoring.

882 A complementary line of work in causal NLP studies how text can serve as treatment, outcome,  
883 mediator, or proxy for unobserved confounders [Feder et al., 2022]; OBS-derived auxiliary labels in  
884 our setting play exactly the proxy role identified there. Taken together, the setting sits across the three  
885 literatures rather than fitting cleanly inside any one of them.

886 **Scope and future work.** Ranking diagnostics can diverge sharply from recommendation regret.  
887 In the cached summarization benchmark, **Grounded** wins every top-1 and top-3 ranking setting  
888 but wins regret in only 6 settings. Regret is driven by a small number of costly recommendation  
889 mistakes that average ranking metrics do not isolate. The benchmarks therefore center  $\text{Regret}_{\text{test}}$ ,  
890 with context-action RMSE and model-level RMSE as secondary diagnostics.

891 The two real-task benchmarks use real candidate model outputs and real LLM-judged rubric or  
892 program-test scores on those outputs. Their OBS and EXP mechanisms are benchmark constructions  
893 over the cached pool. The empirical comparisons therefore diagnose how the estimator families  
894 behave under controlled synthetic OBS/EXP resampling on real tasks. The empirical scope is  
895 deliberately narrow: cached summarization isolates supervision scarcity, and cached coding isolates  
896 reward geometry. Family-level wins in these benchmarks are benchmark-specific diagnostics.

897 Estimator choice depends on what limits performance in the benchmark at hand. In practice, the  
898 key checks are whether more EXP labels materially improve EXP-based methods and whether a  
899 structured representation predicts the held-out target better than a simple success-only baseline.  
900 Deployment decisions should still be confirmed by live randomized evaluation. These conclusions are  
901 specific to the task families and benchmark constructions studied here. Future work includes broader  
902 tasks, repeated mediator draws, naturally observed deployment logs paired with in-deployment ran-  
903 domized trials, and multi-round human–agent interaction, which would require relaxing or replacing  
904 Assumption A4.

## 905 **A.17 Implementation notes**

906 The synthetic validation uses only cached data generation; no external API calls occur after the  
907 generator pool is constructed. DR uses EXP propensities from uniform randomization together  
908 with cross-fitting for reward models trained on EXP. Additional synthetic diagnostics and older  
909 experimental variants are omitted from the present main-text comparisons.

## 910 **B Rubric subscores for CNN/DailyMail summaries**

911 Given an article  $X$  and candidate summary  $M$ , the judge outputs four integer subscores  $\tilde{s}(X, M) \in$   
912  $\{0, 1, 2, 3, 4, 5\}^4$ : *faithfulness*, *coverage*, *clarity*, and *conciseness*. Scores use the anchors below;  
913 intermediate values interpolate between adjacent anchors.

914 **Faithfulness (0–5).** Factual consistency with the article. Unsupported claims include invented  
915 entities, numbers, events, causal attributions, or quotes.

- 916 • **5:** No unsupported claims; all salient statements are grounded in  $X$ .  
917 • **4:** Minor unsupported detail(s) that do not affect the main facts.  
918 • **3:** Multiple unsupported or weakly supported statements; main story remains recognizable.  
919 • **2:** Several unsupported claims affecting key details (e.g., actor/action/outcome/number).  
920 • **1:** Many unsupported claims; substantial contradiction or fabrication.  
921 • **0:** Predominantly hallucinatory or contradictory relative to  $X$ .
- 922 **Coverage (0–5).** Inclusion of the article’s major points (headline facts, primary actors, core events,  
923 outcomes, and key qualifiers).
- 924 • **5:** Covers all major points; omissions are limited to minor details.  
925 • **4:** Covers most major points; at most one major point missing.  
926 • **3:** Covers some major points; several major omissions.  
927 • **2:** Covers few major points; summary reflects only a small slice of the article.  
928 • **1:** Minimal coverage; largely misses the article’s main content.  
929 • **0:** Misses most major points or is largely off-topic relative to  $X$ .
- 930 **Clarity (0–5).** Readability and organization (coherence, grammaticality, referential clarity, and  
931 logical flow).
- 932 • **5:** Clear and well-structured; unambiguous references and fluent phrasing.  
933 • **4:** Generally clear; minor awkwardness or minor ambiguity.  
934 • **3:** Understandable with effort; noticeable disfluency, repetition, or unclear references.  
935 • **2:** Hard to follow; frequent grammatical issues or unclear structure.  
936 • **1:** Mostly unclear; major coherence failures.  
937 • **0:** Unreadable or incoherent.
- 938 **Conciseness (0–5).** Information density at an appropriate length; redundancy and irrelevant detail  
939 reduce the score; excessive brevity that omits necessary content also reduces the score.
- 940 • **5:** Efficient summary with minimal redundancy and no filler.  
941 • **4:** Slight redundancy or mild over/under-length without major impact.  
942 • **3:** Noticeable redundancy or length mismatch; still usable as a summary.  
943 • **2:** Clearly too verbose or too terse; substantial inefficiency or truncation.  
944 • **1:** Extremely verbose or extremely short; poor summary form.  
945 • **0:** Pathological length or pervasive redundancy; not a usable summary.

## 946 C Full statements and proofs for identification and population estimator 947 comparisons

948 This appendix collects the proof of the main identification result and the full statements and proofs  
949 for the population estimator comparisons referenced in Section 4. It also explains how these  
950 approximation-class results motivate the empirical target-alignment diagnostics used in the bench-  
951 marks; those diagnostics are not additional formal results. Throughout,  $Z = (X, M)$ ,  $P_E$  denotes  
952 the EXP law of  $Z$ , and for a fixed target outcome  $Y$  we write

$$\mathcal{R}_E(f) = \sigma_E^2 + \|f - r^*\|_E^2, \quad \|g\|_E^2 := \mathbb{E}_{Z \sim P_E}[g(Z)^2], \quad \sigma_E^2 := \mathbb{E}[(Y - r^*(Z))^2 \mid \text{EXP}], \quad (13)$$

953 where  $r^*(z) := \mathbb{E}[Y \mid Z = z, \text{EXP}]$ . For readability, we ignore the final clipping to  $[0, 1]$  and absorb  
954 the intercept into both  $\psi$  and  $\varphi$  by augmenting them with a constant feature.

### 955 C.1 Proof of Theorem 1

956 *Proof of Theorem 1.* By Assumption A1,  $A \perp U \mid X$  in EXP. Assumption A4 states the structural  
957 restriction  $U \perp M \mid X, A$  in OBS, in EXP, and under  $\text{do}(A = a)$ , so that conditional on  $(X, A)$  the  
958 mediator  $M$  carries no residual information about  $U$ . Combining A1 and A4 with the chain rule for  
959 conditional independence gives

$$U \perp (A, M) \mid X \quad \text{in EXP,}$$

960 so on the shared EXP/SIM support the scoring rule  $r^*(x, m) = \mathbb{E}[Y \mid X = x, M = m, \text{EXP}]$   
 961 is identified from EXP and is free of self-selection bias. Assumption A2 gives  $p_{\text{sim}}(m \mid x, a) =$   
 962  $\mathbb{P}(M = m \mid X = x, \text{do}(A = a))$ , so for any  $x$  in the EXP support,

$$\begin{aligned} q(x, a) &= \mathbb{E}[Y(\text{do}(A = a)) \mid X = x] \\ &= \mathbb{E}_{M \sim p_{\text{sim}}(\cdot \mid x, a)}[\mathbb{E}[Y \mid X = x, M = M, \text{do}(A = a)]] \\ &= \mathbb{E}_{M \sim p_{\text{sim}}(\cdot \mid x, a)}[r^*(x, M)], \end{aligned}$$

963 where the last equality uses Assumption A3 (the conditional outcome law given  $(X, M, U)$  is the  
 964 same in OBS, EXP, and under  $\text{do}(A = a)$ , and  $r^*$  integrates  $U$  out under the EXP  $U \mid X$  distribution,  
 965 which by A1 equals the post-intervention  $U \mid X$  distribution). This proves (1). Averaging  $q(X, a)$   
 966 under the EXP context distribution and invoking Assumption A5 ( $P_X^{\text{tgt}} = P_X^{\text{EXP}}$ ) gives (2). Outside  
 967 Assumption A5, only  $q(x, a)$  on the shared support is identified; reweighting to a different target  
 968 context distribution requires a standard covariate-shift adjustment.  $\square$

969 Define the linear classes

$$\mathcal{H}_\psi := \{z \mapsto \theta^\top \psi(z) : \theta \in \mathbb{R}^{d_\psi}\}, \quad \mathcal{F}_\varphi := \{z \mapsto w^\top \varphi(z) : w \in \mathbb{R}^{d_\varphi}\},$$

970 and the grounded / residual affine class

$$\mathcal{G}(f_{\text{OBS}}, \psi) := \{f_{\text{OBS}} - h : h \in \mathcal{H}_\psi\} = f_{\text{OBS}} + \mathcal{H}_\psi,$$

971 where the equality uses that  $\mathcal{H}_\psi$  is a linear space.

## 972 C.2 Summary of the fixed-target comparisons and empirical alignment diagnostics

973 The main text uses the population comparisons as conditional guidance for reading the experiments,  
 974 not as a guarantee that one estimator dominates. We separate the formal results included below from  
 975 empirical diagnostics that are motivated by the same approximation-gap viewpoint.

### 976 Formal results included below.

- 977 • **At the oracle level, correcting OBS cannot hurt.** Projecting OBS bias onto the correction space  
 978 weakly improves EXP risk relative to **OBS-Only**; shrinkage matters only because the estimated  
 979 correction can be noisy in finite samples. See Theorem 4.
- 980 • **Grounded expands the proxy-side approximation class.** If the OBS baseline is not already linear  
 981 in  $\psi$ , then correcting an expressive baseline by a low-dimensional proxy residual can represent  
 982 targets that a pure  $\psi$ -linear EXP fit cannot. See Theorem 2.
- 983 • **CVCI-Residual is favored when the hard part is already in the OBS baseline.** Residualization  
 984 beats direct pooling when  $f_{\text{OBS}}$  absorbs the difficult part of the reward surface and the remaining  
 985 discrepancy is simpler in proxy space than the full target is in raw text features. See Theorem 3.
- 986 • **In the shared linear proxy special case, Grounded is centered ridge.** When both the OBS  
 987 baseline and the EXP target are linear in the same proxy representation, **Grounded** becomes ridge  
 988 centered at the OBS coefficient, and it beats **EXP-Only** exactly when that OBS center is a better  
 989 shrinkage target than zero. See Theorem 5 and Corollary 6.

990 **Empirical diagnostics motivated by the formal results.** The experiments do not test the oracle  
 991 conditions directly. They report proxy-alignment and reward-shape diagnostics that are related to  
 992 the approximation classes in the formal results. In summarization, the smooth-to-sharpened reward  
 993 comparison changes the target while keeping the cached outputs fixed. In coding, the fix-success  
 994 sweep moves the target away from auxiliary patch-quality features and toward true fix success. These  
 995 diagnostics are empirical probes of estimator–target alignment, not additional theorems.

- 996 • **Target–proxy alignment is measured empirically.** The representation-based estimators are  
 997 most useful when the target reward varies along directions retained by the auxiliary proxy. The  
 998 summarization and coding experiments report held-out target-fit  $R^2$  as a proxy for this alignment.
- 999 • **Reward definitions can change estimator rankings through approximation-class alignment.**  
 1000 Changing the reward definition can move the target toward or away from the approximation class  
 1001 used by an estimator. The smooth/sharpened summarization comparison and the coding fix-success  
 1002 sweep are diagnostics consistent with this mechanism.

1003 • **Regret is the primary downstream metric.** The experiments report recommendation regret  
 1004 because the evaluation objective is action selection; RMSE and ranking metrics are secondary  
 1005 diagnostics for reward-surface fit and model-level value fit.

### 1006 C.3 What representation correction buys

1007 **Theorem 2** (What representation correction buys). *For any fixed baseline  $f_{\text{OBS}}$ , the grounded class*  
 1008  $\mathcal{G}(f_{\text{OBS}}, \psi)$  *satisfies the following three statements.*

1009 (a) *If there exists  $\theta_\star \in \mathbb{R}^{d_\psi}$  such that*

$$r^\star(z) = f_{\text{OBS}}(z) - \theta_\star^\top \psi(z) \quad \text{for all } z,$$

1010 *then  $r^\star \in \mathcal{G}(f_{\text{OBS}}, \psi)$ . Conditional on the fitted baseline  $f_{\text{OBS}}$ , EXP only needs to estimate the*  
 1011  *$d_\psi$ -dimensional residual coefficient  $\theta_\star$ .*

1012 (b) *If  $f_{\text{OBS}} \notin \mathcal{H}_\psi$ , then*

$$\mathcal{G}(f_{\text{OBS}}, \psi) \neq \mathcal{H}_\psi, \quad f_{\text{OBS}} \in \mathcal{G}(f_{\text{OBS}}, \psi) \setminus \mathcal{H}_\psi.$$

1013 *So grounded correction can represent targets that a pure  $\psi$ -linear EXP fit cannot.*

1014 (c) *If  $f_{\text{OBS}} \in \mathcal{H}_\psi$ , then*

$$\mathcal{G}(f_{\text{OBS}}, \psi) = \mathcal{H}_\psi.$$

1015 *In that special case, representation correction does not enlarge the function class.*

1016 *Proof.* We prove the three claims in order.

1017 For part (a), the displayed identity

$$r^\star(z) = f_{\text{OBS}}(z) - \theta_\star^\top \psi(z)$$

1018 is exactly the statement that  $r^\star \in \mathcal{G}(f_{\text{OBS}}, \psi)$ . Once  $f_{\text{OBS}}$  is treated as fixed, the only EXP-side  
 1019 unknown is the  $d_\psi$ -dimensional coefficient vector  $\theta_\star$ .

1020 For part (b), note that  $0 \in \mathcal{H}_\psi$ , hence

$$f_{\text{OBS}} = f_{\text{OBS}} - 0 \in \mathcal{G}(f_{\text{OBS}}, \psi).$$

1021 If  $\mathcal{G}(f_{\text{OBS}}, \psi) = \mathcal{H}_\psi$ , then this would imply  $f_{\text{OBS}} \in \mathcal{H}_\psi$ , contradicting the assumption. Therefore  
 1022  $\mathcal{G}(f_{\text{OBS}}, \psi) \neq \mathcal{H}_\psi$ , and  $f_{\text{OBS}}$  is a member of  $\mathcal{G}(f_{\text{OBS}}, \psi)$  that does not belong to  $\mathcal{H}_\psi$ .

1023 For part (c), suppose  $f_{\text{OBS}}(z) = \beta_{\text{OBS}}^\top \psi(z)$  for some  $\beta_{\text{OBS}} \in \mathbb{R}^{d_\psi}$ . Every  $g \in \mathcal{G}(f_{\text{OBS}}, \psi)$  has the  
 1024 form

$$g(z) = \beta_{\text{OBS}}^\top \psi(z) - \theta^\top \psi(z) = (\beta_{\text{OBS}} - \theta)^\top \psi(z) \in \mathcal{H}_\psi,$$

1025 so  $\mathcal{G}(f_{\text{OBS}}, \psi) \subseteq \mathcal{H}_\psi$ . Conversely, for any  $g(z) = \beta^\top \psi(z) \in \mathcal{H}_\psi$ , choosing  $\theta = \beta_{\text{OBS}} - \beta$  gives

$$g(z) = f_{\text{OBS}}(z) - \theta^\top \psi(z) \in \mathcal{G}(f_{\text{OBS}}, \psi).$$

1026 Thus  $\mathcal{H}_\psi \subseteq \mathcal{G}(f_{\text{OBS}}, \psi)$ , proving equality. □

### 1027 C.4 Residualization versus direct pooling

1028 **Theorem 3** (Residualization versus direct pooling). *Let  $f_{\text{OBS}}$  be the observational baseline used*  
 1029 *by Grounded and CVCI-Residual, and assume  $f_{\text{OBS}}, r^\star \in L^2(P_E)$  together with finite-dimensional*  
 1030 *linear classes  $\mathcal{F}_\varphi$  and  $\mathcal{H}_\psi$ , so that both classes are closed in  $L^2(P_E)$  and the projections in question*  
 1031 *exist (when  $\arg \min$  is non-unique, any minimizer suffices, since the displayed risks depend only on*  
 1032 *the projections). Define the oracle direct and residual predictors by*

$$f_{\text{cvci}}^\dagger \in \arg \min_{f \in \mathcal{F}_\varphi} \mathcal{R}_E(f), \quad h_{\text{res}}^\dagger \in \arg \min_{h \in \mathcal{H}_\psi} \mathcal{R}_E(f_{\text{OBS}} + h).$$

1033 *Then*

$$\mathcal{R}_E(f_{\text{OBS}} + h_{\text{res}}^\dagger) \leq \mathcal{R}_E(f_{\text{cvci}}^\dagger)$$

1034 *if and only if*

$$\text{dist}_E(r^\star, f_{\text{OBS}} + \mathcal{H}_\psi) \leq \text{dist}_E(r^\star, \mathcal{F}_\varphi),$$

1035 *where  $\text{dist}_E(r, \mathcal{A}) := \inf_{g \in \mathcal{A}} \|r - g\|_E$ . In particular, if*

$$r^\star - f_{\text{OBS}} \in \mathcal{H}_\psi \quad \text{but} \quad r^\star \notin \overline{\mathcal{F}_\varphi},$$

1036 *then*

$$\mathcal{R}_E(f_{\text{OBS}} + h_{\text{res}}^\dagger) < \mathcal{R}_E(f_{\text{cvci}}^\dagger).$$

1037 *Proof.* By (13),

$$\inf_{f \in \mathcal{F}_\varphi} \mathcal{R}_E(f) = \sigma_E^2 + \inf_{f \in \mathcal{F}_\varphi} \|f - r^*\|_E^2 = \sigma_E^2 + \text{dist}_E^2(r^*, \mathcal{F}_\varphi).$$

1038 Likewise,

$$\inf_{h \in \mathcal{H}_\psi} \mathcal{R}_E(f_{\text{OBS}} + h) = \sigma_E^2 + \inf_{h \in \mathcal{H}_\psi} \|f_{\text{OBS}} + h - r^*\|_E^2 = \sigma_E^2 + \text{dist}_E^2(r^*, f_{\text{OBS}} + \mathcal{H}_\psi).$$

1039 Subtracting the common irreducible term  $\sigma_E^2$  gives the equivalence.

1040 If  $r^* - f_{\text{OBS}} \in \mathcal{H}_\psi$ , then

$$\text{dist}_E(r^*, f_{\text{OBS}} + \mathcal{H}_\psi) = 0.$$

1041 If also  $r^* \notin \overline{\mathcal{F}_\varphi}$ , then

$$\text{dist}_E(r^*, \mathcal{F}_\varphi) > 0.$$

1042 Hence

$$\mathcal{R}_E(f_{\text{OBS}} + h_{\text{res}}^\dagger) < \mathcal{R}_E(f_{\text{cvci}}^\dagger). \quad \square$$

### 1043 C.5 Oracle grounding and shrinkage

1044 **Theorem 4** (Oracle grounding is no worse than OBS-only). *Let  $b := f_{\text{OBS}} - r^*$  denote the OBS bias*  
 1045 *function on EXP. Let  $\mathcal{H}_\psi \subset L^2(P_E)$  be the resulting closed linear correction space, and let*

$$h^\dagger := \Pi_{\mathcal{H}_\psi} b$$

1046 *be the  $L^2(P_E)$  projection of  $b$  onto that space. For any  $\alpha \in [0, 1]$ , define the oracle grounded*  
 1047 *predictor*

$$f_{G,\alpha} := f_{\text{OBS}} - \alpha h^\dagger.$$

1048 *Then*

$$\mathcal{R}_E(f_{G,\alpha}) = \mathcal{R}_E(f_{\text{OBS}}) - (2\alpha - \alpha^2) \|h^\dagger\|_E^2. \quad (14)$$

1049 *Hence*

$$\mathcal{R}_E(f_{G,\alpha}) \leq \mathcal{R}_E(f_{\text{OBS}}) \quad \text{for all } \alpha \in [0, 1],$$

1050 *with strict inequality whenever  $\alpha > 0$  and  $h^\dagger \neq 0$ .*

1051 *Proof.* Write

$$b = h^\dagger + u, \quad u := b - h^\dagger.$$

1052 Because  $h^\dagger$  is the orthogonal projection of  $b$  onto the closed linear subspace  $\mathcal{H}_\psi$ , we have  $u \perp h^\dagger$  in  
 1053  $L^2(P_E)$ .

1054 Now

$$f_{G,\alpha} - r^* = (f_{\text{OBS}} - r^*) - \alpha h^\dagger = u + (1 - \alpha) h^\dagger.$$

1055 Therefore,

$$\mathcal{R}_E(f_{G,\alpha}) - \sigma_E^2 = \|u + (1 - \alpha) h^\dagger\|_E^2 = \|u\|_E^2 + (1 - \alpha)^2 \|h^\dagger\|_E^2,$$

1056 where the cross term vanishes by orthogonality. On the other hand,

$$\mathcal{R}_E(f_{\text{OBS}}) - \sigma_E^2 = \|b\|_E^2 = \|u + h^\dagger\|_E^2 = \|u\|_E^2 + \|h^\dagger\|_E^2.$$

1057 Subtracting the two displays gives

$$\mathcal{R}_E(f_{G,\alpha}) - \mathcal{R}_E(f_{\text{OBS}}) = ((1 - \alpha)^2 - 1) \|h^\dagger\|_E^2 = -(2\alpha - \alpha^2) \|h^\dagger\|_E^2,$$

1058 which is (14). Since  $2\alpha - \alpha^2 \geq 0$  on  $[0, 1]$ , the inequality follows, and it is strict whenever  $\alpha > 0$   
 1059 and  $h^\dagger \neq 0$ .  $\square$

1060 **Derivation of the noisy-correction identity.** For any  $\tilde{h} \in L^2(P_E)$ ,

$$\mathcal{R}_E(f_{\text{OBS}} - \alpha \tilde{h}) - \mathcal{R}_E(f_{\text{OBS}}) = \alpha^2 \|\tilde{h}\|_E^2 - 2\alpha \langle b, \tilde{h} \rangle_E. \quad (15)$$

1061 Indeed,  $f_{\text{OBS}} - \alpha \tilde{h} - r^* = b - \alpha \tilde{h}$ , so

$$\mathcal{R}_E(f_{\text{OBS}} - \alpha \tilde{h}) - \sigma_E^2 = \|b - \alpha \tilde{h}\|_E^2 = \|b\|_E^2 + \alpha^2 \|\tilde{h}\|_E^2 - 2\alpha \langle b, \tilde{h} \rangle_E.$$

1062 Subtracting  $\mathcal{R}_E(f_{\text{OBS}}) - \sigma_E^2 = \|b\|_E^2$  yields (15).

1063 **C.6 The shared linear proxy special case**

1064 **Theorem 5** (Grounded is centered ridge in the linear proxy special case). *Assume the observational*  
 1065 *baseline is already linear in the proxy representation,*

$$f_{\text{OBS}}(z) = \beta_{\text{OBS}}^\top \psi(z),$$

1066 *with the intercept absorbed into  $\psi$ . Let the EXP sample be  $\{(\psi_j, Y_j)\}_{j=1}^n$ , and write*

$$\Psi = \begin{bmatrix} \psi_1^\top \\ \vdots \\ \psi_n^\top \end{bmatrix} \in \mathbb{R}^{n \times d_\psi}, \quad Y = (Y_1, \dots, Y_n)^\top.$$

1067 *Under full correction  $\alpha_{\text{corr}} = 1$ , grounded correction solves*

$$\hat{\delta} \in \arg \min_{\delta \in \mathbb{R}^{d_\psi}} \|\Psi \beta_{\text{OBS}} - Y - \Psi \delta\|_2^2 + \lambda \|\delta\|_2^2$$

1068 *and outputs*

$$\hat{\beta}_{\text{G}} := \beta_{\text{OBS}} - \hat{\delta}.$$

1069 *Then  $\hat{\beta}_{\text{G}}$  is exactly the centered ridge estimator*

$$\hat{\beta}_{\text{G}} \in \arg \min_{\beta \in \mathbb{R}^{d_\psi}} \|Y - \Psi \beta\|_2^2 + \lambda \|\beta - \beta_{\text{OBS}}\|_2^2.$$

1070 *Equivalently,*

$$\hat{\beta}_{\text{G}} = (\Psi^\top \Psi + \lambda I)^{-1} (\Psi^\top Y + \lambda \beta_{\text{OBS}}).$$

1071 *For general  $\alpha_{\text{corr}} \in [0, 1]$ ,*

$$\hat{\beta}_{\text{G}, \alpha_{\text{corr}}} := \beta_{\text{OBS}} - \alpha_{\text{corr}} \hat{\delta} = (1 - \alpha_{\text{corr}}) \beta_{\text{OBS}} + \alpha_{\text{corr}} \hat{\beta}_{\text{G}}.$$

1072 *Proof.* Set

$$\beta = \beta_{\text{OBS}} - \delta.$$

1073 Then  $\delta = \beta_{\text{OBS}} - \beta$ , so

$$\Psi \beta_{\text{OBS}} - Y - \Psi \delta = \Psi \beta_{\text{OBS}} - Y - \Psi (\beta_{\text{OBS}} - \beta) = \Psi \beta - Y.$$

1074 Therefore

$$\|\Psi \beta_{\text{OBS}} - Y - \Psi \delta\|_2^2 + \lambda \|\delta\|_2^2 = \|Y - \Psi \beta\|_2^2 + \lambda \|\beta - \beta_{\text{OBS}}\|_2^2.$$

1075 So minimizing over  $\delta$  is equivalent to minimizing over  $\beta$ , which proves the centered-ridge representa-  
 1076 tion.

1077 For the closed form, differentiate the centered-ridge objective:

$$-2\Psi^\top (Y - \Psi \beta) + 2\lambda(\beta - \beta_{\text{OBS}}) = 0,$$

1078 hence

$$(\Psi^\top \Psi + \lambda I) \beta = \Psi^\top Y + \lambda \beta_{\text{OBS}},$$

1079 which yields the stated formula for  $\hat{\beta}_{\text{G}}$ .

1080 Finally,

$$\hat{\beta}_{\text{G}, \alpha_{\text{corr}}} = \beta_{\text{OBS}} - \alpha_{\text{corr}} \hat{\delta} = \beta_{\text{OBS}} - \alpha_{\text{corr}} (\beta_{\text{OBS}} - \hat{\beta}_{\text{G}}) = (1 - \alpha_{\text{corr}}) \beta_{\text{OBS}} + \alpha_{\text{corr}} \hat{\beta}_{\text{G}}. \quad \square$$

1081 **Corollary 6** (When the OBS center beats EXP-only). *Assume in addition that the EXP target is*  
 1082 *linear in the same representation,*

$$Y = \Psi \beta_\star + \varepsilon, \quad \mathbb{E}[\varepsilon \mid \Psi] = 0, \quad \text{Var}(\varepsilon \mid \Psi) = \sigma^2 I_n.$$

1083 Let  $S := \Psi^\top \Psi$ ,  $A := (S + \lambda I)^{-1}$ , and define the EXP-only ridge estimator

$$\hat{\beta}_{\text{EXP}} := A \Psi^\top Y.$$

1084 For any positive semidefinite matrix  $G$ , write  $\|u\|_G^2 := u^\top G u$ . Then, conditional on  $\Psi$ ,

$$\mathbb{E} \left[ \|\hat{\beta}_{\text{G}} - \beta_\star\|_G^2 \mid \Psi \right] = \lambda^2 (\beta_\star - \beta_{\text{OBS}})^\top A G A (\beta_\star - \beta_{\text{OBS}}) + \sigma^2 \text{tr}(G A S A),$$

1085 while

$$\mathbb{E} \left[ \|\hat{\beta}_{\text{EXP}} - \beta_\star\|_G^2 \mid \Psi \right] = \lambda^2 \beta_\star^\top A G A \beta_\star + \sigma^2 \text{tr}(G A S A).$$

1086 Therefore the linear grounded estimator is better than EXP-Only if and only if

$$(\beta_\star - \beta_{\text{OBS}})^\top A G A (\beta_\star - \beta_{\text{OBS}}) \leq \beta_\star^\top A G A \beta_\star. \quad (16)$$

1087 *Proof.* First,

$$\widehat{\beta}_{\text{EXP}} - \beta_{\star} = A\Psi^{\top}(\Psi\beta_{\star} + \varepsilon) - \beta_{\star} = (AS - I)\beta_{\star} + A\Psi^{\top}\varepsilon.$$

1088 Since  $AS - I = -\lambda A$ ,

$$\widehat{\beta}_{\text{EXP}} - \beta_{\star} = -\lambda A\beta_{\star} + A\Psi^{\top}\varepsilon. \quad (17)$$

1089 For grounded, note that

$$\Psi\beta_{\text{OBS}} - Y = \Psi(\beta_{\text{OBS}} - \beta_{\star}) - \varepsilon = \Psi\delta_{\star} - \varepsilon.$$

1090 Hence

$$\widehat{\beta}_{\text{G}} - \beta_{\star} = \beta_{\text{OBS}} - \beta_{\star} - A\Psi^{\top}(\Psi\delta_{\star} - \varepsilon) = \delta_{\star} - AS\delta_{\star} + A\Psi^{\top}\varepsilon.$$

1091 Using  $I - AS = \lambda A$  gives

$$\widehat{\beta}_{\text{G}} - \beta_{\star} = \lambda A\delta_{\star} + A\Psi^{\top}\varepsilon. \quad (18)$$

1092 Now take conditional  $G$ -risk. From (17),

$$\mathbb{E}\left[\|\widehat{\beta}_{\text{EXP}} - \beta_{\star}\|_G^2 \mid \Psi\right] = \lambda^2 \beta_{\star}^{\top} A G A \beta_{\star} + \mathbb{E}\left[(A\Psi^{\top}\varepsilon)^{\top} G (A\Psi^{\top}\varepsilon) \mid \Psi\right],$$

1093 because the cross term vanishes by  $\mathbb{E}[\varepsilon \mid \Psi] = 0$ . The noise term is

$$\mathbb{E}\left[(A\Psi^{\top}\varepsilon)^{\top} G (A\Psi^{\top}\varepsilon) \mid \Psi\right] = \text{tr}\left(G A \Psi^{\top} \mathbb{E}[\varepsilon\varepsilon^{\top} \mid \Psi] \Psi A\right) = \sigma^2 \text{tr}(G A S A).$$

1094 This gives the EXP-only expression.

1095 The grounded formula is identical, replacing the bias vector  $-\lambda A\beta_{\star}$  by  $\lambda A\delta_{\star}$  from (18). Comparing  
1096 the two expressions yields (16).  $\square$

## 1097 **NeurIPS Paper Checklist**

### 1098 **1. Claims**

1099 Question: Do the main claims made in the abstract and introduction accurately reflect the paper’s  
1100 contributions and scope?

1101 Answer: [Yes].

1102 Justification: The abstract and Section 1 state the identification claim, the post-identification role  
1103 of OBS, and the empirical scope of the estimator comparison.

1104 Guidelines:

- 1105 • The answer [N/A] means that the abstract and introduction do not include the claims made in  
1106 the paper.
- 1107 • The abstract and/or introduction should clearly state the claims made, including the contributions  
1108 made in the paper and important assumptions and limitations. A [No] or [N/A] answer to this  
1109 question will not be perceived well by the reviewers.
- 1110 • The claims made should match theoretical and experimental results, and reflect how much the  
1111 results can be expected to generalize to other settings.
- 1112 • It is fine to include aspirational goals as motivation as long as it is clear that these goals are not  
1113 attained by the paper.

### 1114 **2. Limitations**

1115 Question: Does the paper discuss the limitations of the work performed by the authors?

1116 Answer: [Yes].

1117 Justification: Section 7 describes the cached-benchmark scope, benchmark-specific interpretation,  
1118 and future work needed for broader deployments.

1119 Guidelines:

- 1120 • The answer [N/A] means that the paper has no limitation while the answer [No] means that the  
1121 paper has limitations, but those are not discussed in the paper.
- 1122 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 1123 • The paper should point out any strong assumptions and how robust the results are to violations of  
1124 these assumptions (e.g., independence assumptions, noiseless settings, model well-specification,  
1125 asymptotic approximations only holding locally). The authors should reflect on how these  
1126 assumptions might be violated in practice and what the implications would be.
- 1127 • The authors should reflect on the scope of the claims made, e.g., if the approach was only tested  
1128 on a few datasets or with a few runs. In general, empirical results often depend on implicit  
1129 assumptions, which should be articulated.
- 1130 • The authors should reflect on the factors that influence the performance of the approach. For  
1131 example, a facial recognition algorithm may perform poorly when image resolution is low or  
1132 images are taken in low lighting. Or a speech-to-text system might not be used reliably to  
1133 provide closed captions for online lectures because it fails to handle technical jargon.
- 1134 • The authors should discuss the computational efficiency of the proposed algorithms and how  
1135 they scale with dataset size.
- 1136 • If applicable, the authors should discuss possible limitations of their approach to address  
1137 problems of privacy and fairness.
- 1138 • While the authors might fear that complete honesty about limitations might be used by reviewers  
1139 as grounds for rejection, a worse outcome might be that reviewers discover limitations that  
1140 aren’t acknowledged in the paper. The authors should use their best judgment and recognize  
1141 that individual actions in favor of transparency play an important role in developing norms that  
1142 preserve the integrity of the community. Reviewers will be specifically instructed to not penalize  
1143 honesty concerning limitations.

### 1144 **3. Theory assumptions and proofs**

1145 Question: For each theoretical result, does the paper provide the full set of assumptions and a  
1146 complete (and correct) proof?

1147 Answer: [Yes].

1148 Justification: Section 2 states Assumptions A1–A5, Section 3 states the main theorem and proof  
1149 sketch, and Appendix C gives full statements and proofs.

1150 Guidelines:

- 1151 • The answer [N/A] means that the paper does not include theoretical results.
- 1152 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.

- 1153 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 1154 • The proofs can either appear in the main paper or the supplemental material, but if they appear
- 1155 in the supplemental material, the authors are encouraged to provide a short proof sketch to
- 1156 provide intuition.
- 1157 • Inversely, any informal proof provided in the core of the paper should be complemented by
- 1158 formal proofs provided in appendix or supplemental material.
- 1159 • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 1160 4. **Experimental result reproducibility**

1161 Question: Does the paper fully disclose all the information needed to reproduce the main experi-  
1162 mental results of the paper to the extent that it affects the main claims and/or conclusions of the  
1163 paper (regardless of whether the code and data are provided or not)?

1164 Answer: [Yes].

1165 Justification: Sections 5, 5.4, and 5.5, together with Appendices A.10, A.11, A.13, and A.15,  
1166 specify the benchmark construction, resampling rules, budgets, tuning, seeds, and metrics.

1167 Guidelines:

- 1168 • The answer [N/A] means that the paper does not include experiments.
- 1169 • If the paper includes experiments, a [No] answer to this question will not be perceived well by
- 1170 the reviewers: Making the paper reproducible is important, regardless of whether the code and
- 1171 data are provided or not.
- 1172 • If the contribution is a dataset and/or model, the authors should describe the steps taken to make
- 1173 their results reproducible or verifiable.
- 1174 • Depending on the contribution, reproducibility can be accomplished in various ways. For
- 1175 example, if the contribution is a novel architecture, describing the architecture fully might
- 1176 suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary
- 1177 to either make it possible for others to replicate the model with the same dataset, or provide
- 1178 access to the model. In general, releasing code and data is often one good way to accomplish
- 1179 this, but reproducibility can also be provided via detailed instructions for how to replicate the
- 1180 results, access to a hosted model (e.g., in the case of a large language model), releasing of a
- 1181 model checkpoint, or other means that are appropriate to the research performed.
- 1182 • While NeurIPS does not require releasing code, the conference does require all submissions
- 1183 to provide some reasonable avenue for reproducibility, which may depend on the nature of the
- 1184 contribution. For example
  - 1185 (a) If the contribution is primarily a new algorithm, the paper should make it clear how to
  - 1186 reproduce that algorithm.
  - 1187 (b) If the contribution is primarily a new model architecture, the paper should describe the
  - 1188 architecture clearly and fully.
  - 1189 (c) If the contribution is a new model (e.g., a large language model), then there should either
  - 1190 be a way to access this model for reproducing the results or a way to reproduce the model
  - 1191 (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - 1192 (d) We recognize that reproducibility may be tricky in some cases, in which case authors are
  - 1193 welcome to describe the particular way they provide for reproducibility. In the case of
  - 1194 closed-source models, it may be that access to the model is limited in some way (e.g.,
  - 1195 to registered users), but it should be possible for other researchers to have some path to
  - 1196 reproducing or verifying the results.

#### 1197 5. **Open access to data and code**

1198 Question: Does the paper provide open access to the data and code, with sufficient instructions to  
1199 faithfully reproduce the main experimental results, as described in supplemental material?

1200 Answer: [No].

1201 Justification: The draft cites public source datasets and cached patch resources, but it does not yet  
1202 provide an anonymized code or data artifact with executable reproduction commands.

1203 Guidelines:

- 1204 • The answer [N/A] means that paper does not include experiments requiring code.
- 1205 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/public/](https://neurips.cc/public/guides/CodeSubmissionPolicy)  
1206 [guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1207 • While we encourage the release of code and data, we understand that this might not be possible,  
1208 so [No] is an acceptable answer. Papers cannot be rejected simply for not including code, unless  
1209 this is central to the contribution (e.g., for a new open-source benchmark).

- 1210 • The instructions should contain the exact command and environment needed to run to reproduce  
1211 the results. See the NeurIPS code and data submission guidelines ([https://neurips.cc/  
1212 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1213 • The authors should provide instructions on data access and preparation, including how to access  
1214 the raw data, preprocessed data, intermediate data, and generated data, etc.
- 1215 • The authors should provide scripts to reproduce all experimental results for the new proposed  
1216 method and baselines. If only a subset of experiments are reproducible, they should state which  
1217 ones are omitted from the script and why.
- 1218 • At submission time, to preserve anonymity, the authors should release anonymized versions (if  
1219 applicable).
- 1220 • Providing as much information as possible in supplemental material (appended to the paper) is  
1221 recommended, but including URLs to data and code is permitted.

## 1222 6. Experimental setting/details

1223 Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters,  
1224 how they were chosen, type of optimizer) necessary to understand the results?

1225 Answer: [Yes].

1226 Justification: Appendix A.10 gives the fixed pool, train/test split, OBS/EXP sampling, and  
1227 aggregation protocol, while Appendix A.11 specifies tuning and fallback rules.

1228 Guidelines:

- 1229 • The answer [N/A] means that the paper does not include experiments.
- 1230 • The experimental setting should be presented in the core of the paper to a level of detail that is  
1231 necessary to appreciate the results and make sense of them.
- 1232 • The full details can be provided either with the code, in appendix, or as supplemental material.

## 1233 7. Experiment statistical significance

1234 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
1235 information about the statistical significance of the experiments?

1236 Answer: [Yes].

1237 Justification: The main-text tables report means over paired seeds, and Appendix A.10 states that  
1238 means and standard errors are recomputed from seed-level paired runs.

1239 Guidelines:

- 1240 • The answer [N/A] means that the paper does not include experiments.
- 1241 • The authors should answer [Yes] if the results are accompanied by error bars, confidence  
1242 intervals, or statistical significance tests, at least for the experiments that support the main claims  
1243 of the paper.
- 1244 • The factors of variability that the error bars are capturing should be clearly stated (for example,  
1245 train/test split, initialization, random drawing of some parameter, or overall run with given  
1246 experimental conditions).
- 1247 • The method for calculating the error bars should be explained (closed form formula, call to a  
1248 library function, bootstrap, etc.)
- 1249 • The assumptions made should be given (e.g., Normally distributed errors).
- 1250 • It should be clear whether the error bar is the standard deviation or the standard error of the  
1251 mean.
- 1252 • It is OK to report 1-sigma error bars, but one should state it. The authors should preferably  
1253 report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of  
1254 errors is not verified.
- 1255 • For asymmetric distributions, the authors should be careful not to show in tables or figures  
1256 symmetric error bars that would yield results that are out of range (e.g., negative error rates).
- 1257 • If error bars are reported in tables or plots, the authors should explain in the text how they were  
1258 calculated and reference the corresponding figures or tables in the text.

## 1259 8. Experiments compute resources

1260 Question: For each experiment, does the paper provide sufficient information on the computer  
1261 resources (type of compute workers, memory, time of execution) needed to reproduce the experi-  
1262 ments?

1263 Answer: [No].

1264 Justification: The draft describes experimental budgets and cached-data construction, but it does  
1265 not yet report hardware, memory, runtime, or total compute requirements.

- 1266 Guidelines:
- 1267 • The answer [N/A] means that the paper does not include experiments.
  - 1268 • The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud
  - 1269 provider, including relevant memory and storage.
  - 1270 • The paper should provide the amount of compute required for each of the individual experimental
  - 1271 runs as well as estimate the total compute.
  - 1272 • The paper should disclose whether the full research project required more compute than the
  - 1273 experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it
  - 1274 into the paper).

#### 1275 9. Code of ethics

1276 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS

1277 Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

1278 Answer: [Yes].

1279 Justification: The work analyzes public/cached benchmark data and does not introduce human-

1280 subject data collection, high-risk model release, or deployed decision systems.

1281 Guidelines:

- 1282 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of Ethics.
- 1283 • If the authors answer [No], they should explain the special circumstances that require a deviation
- 1284 from the Code of Ethics.
- 1285 • The authors should make sure to preserve anonymity (e.g., if there is a special consideration due
- 1286 to laws or regulations in their jurisdiction).

#### 1287 10. Broader impacts

1288 Question: Does the paper discuss both potential positive societal impacts and negative societal

1289 impacts of the work performed?

1290 Answer: [No].

1291 Justification: The draft discusses methodological scope and deployment validation, but it does not

1292 yet contain a dedicated discussion of positive and negative societal impacts.

1293 Guidelines:

- 1294 • The answer [N/A] means that there is no societal impact of the work performed.
- 1295 • If the authors answer [N/A] or [No], they should explain why their work has no societal impact
- 1296 or why the paper does not address societal impact.
- 1297 • Examples of negative societal impacts include potential malicious or unintended uses (e.g.,
- 1298 disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deploy-
- 1299 ment of technologies that could make decisions that unfairly impact specific groups), privacy
- 1300 considerations, and security considerations.
- 1301 • The conference expects that many papers will be foundational research and not tied to par-
- 1302 ticular applications, let alone deployments. However, if there is a direct path to any negative
- 1303 applications, the authors should point it out. For example, it is legitimate to point out that
- 1304 an improvement in the quality of generative models could be used to generate Deepfakes for
- 1305 disinformation. On the other hand, it is not needed to point out that a generic algorithm for
- 1306 optimizing neural networks could enable people to train models that generate Deepfakes faster.
- 1307 • The authors should consider possible harms that could arise when the technology is being used
- 1308 as intended and functioning correctly, harms that could arise when the technology is being used
- 1309 as intended but gives incorrect results, and harms following from (intentional or unintentional)
- 1310 misuse of the technology.
- 1311 • If there are negative societal impacts, the authors could also discuss possible mitigation strategies
- 1312 (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for
- 1313 monitoring misuse, mechanisms to monitor how a system learns from feedback over time,
- 1314 improving the efficiency and accessibility of ML).

#### 1315 11. Safeguards

1316 Question: Does the paper describe safeguards that have been put in place for responsible release

1317 of data or models that have a high risk for misuse (e.g., pre-trained language models, image

1318 generators, or scraped datasets)?

1319 Answer: [N/A].

1320 Justification: The paper does not release a new high-risk model or scraped dataset; it evaluates

1321 methods on cached benchmark artifacts and public datasets.

1322 Guidelines:

1323 • The answer [N/A] means that the paper poses no such risks.

1324 • Released models that have a high risk for misuse or dual-use should be released with necessary

1325 safeguards to allow for controlled use of the model, for example by requiring that users adhere

1326 to usage guidelines or restrictions to access the model or implementing safety filters.

1327 • Datasets that have been scraped from the Internet could pose safety risks. The authors should

1328 describe how they avoided releasing unsafe images.

1329 • We recognize that providing effective safeguards is challenging, and many papers do not require

1330 this, but we encourage authors to take this into account and make a best faith effort.

1331 **12. Licenses for existing assets**

1332 Question: Are the creators or original owners of assets (e.g., code, data, models), used in the

1333 paper, properly credited and are the license and terms of use explicitly mentioned and properly

1334 respected?

1335 Answer: [No].

1336 Justification: The draft credits the datasets, models, and patch resources with citations and URLs,

1337 but it does not yet explicitly list licenses or terms of use for each asset.

1338 Guidelines:

1339 • The answer [N/A] means that the paper does not use existing assets.

1340 • The authors should cite the original paper that produced the code package or dataset.

1341 • The authors should state which version of the asset is used and, if possible, include a URL.

1342 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.

1343 • For scraped data from a particular source (e.g., website), the copyright and terms of service of

1344 that source should be provided.

1345 • If assets are released, the license, copyright information, and terms of use in the package should

1346 be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for

1347 some datasets. Their licensing guide can help determine the license of a dataset.

1348 • For existing datasets that are re-packaged, both the original license and the license of the derived

1349 asset (if it has changed) should be provided.

1350 • If this information is not available online, the authors are encouraged to reach out to the asset's

1351 creators.

1352 **13. New assets**

1353 Question: Are new assets introduced in the paper well documented and is the documentation

1354 provided alongside the assets?

1355 Answer: [N/A].

1356 Justification: The submission draft does not introduce or release a new dataset, model, or code

1357 package.

1358 Guidelines:

1359 • The answer [N/A] means that the paper does not release new assets.

1360 • Researchers should communicate the details of the dataset/code/model as part of their sub-

1361 missions via structured templates. This includes details about training, license, limitations,

1362 etc.

1363 • The paper should discuss whether and how consent was obtained from people whose asset is

1364 used.

1365 • At submission time, remember to anonymize your assets (if applicable). You can either create

1366 an anonymized URL or include an anonymized zip file.

1367 **14. Crowdsourcing and research with human subjects**

1368 Question: For crowdsourcing experiments and research with human subjects, does the paper

1369 include the full text of instructions given to participants and screenshots, if applicable, as well as

1370 details about compensation (if any)?

1371 Answer: [N/A].

1372 Justification: The paper does not conduct new crowdsourcing or human-subject experiments.

1373 Guidelines:

1374 • The answer [N/A] means that the paper does not involve crowdsourcing nor research with

1375 human subjects.

- 1376 • Including this information in the supplemental material is fine, but if the main contribution of  
1377 the paper involves human subjects, then as much detail as possible should be included in the  
1378 main paper.
- 1379 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other  
1380 labor should be paid at least the minimum wage in the country of the data collector.
- 1381 **15. Institutional review board (IRB) approvals or equivalent for research with human subjects**
- 1382 Question: Does the paper describe potential risks incurred by study participants, whether such  
1383 risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals  
1384 (or an equivalent approval/review based on the requirements of your country or institution) were  
1385 obtained?
- 1386 Answer: [N/A].
- 1387 Justification: The paper does not conduct new human-subject research or collect new participant  
1388 data.
- 1389 Guidelines:
- 1390 • The answer [N/A] means that the paper does not involve crowdsourcing nor research with  
1391 human subjects.
  - 1392 • Depending on the country in which research is conducted, IRB approval (or equivalent) may be  
1393 required for any human subjects research. If you obtained IRB approval, you should clearly  
1394 state this in the paper.
  - 1395 • We recognize that the procedures for this may vary significantly between institutions and  
1396 locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for  
1397 their institution.
  - 1398 • For initial submissions, do not include any information that would break anonymity (if applica-  
1399 ble), such as the institution conducting the review.
- 1400 **16. Declaration of LLM usage**
- 1401 Question: Does the paper describe the usage of LLMs if it is an important, original, or non-  
1402 standard component of the core methods in this research? Note that if the LLM is used only for  
1403 writing, editing, or formatting purposes and does *not* impact the core methodology, scientific rigor,  
1404 or originality of the research, declaration is not required.
- 1405 Answer: [Yes].
- 1406 Justification: The paper describes LLM model outputs and LLM-judged rubric scores in Sections 5,  
1407 5.4, and Appendix A.10.
- 1408 Guidelines:
- 1409 • The answer [N/A] means that the core method development in this research does not involve  
1410 LLMs as any important, original, or non-standard components.
  - 1411 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not be  
1412 described.