# COST-OPTIMAL ACTIVE AI MODEL EVALUATION

**Anonymous authors** 

000

001 002 003

004

006

008 009

010

011

012

013

014

015

016

017

018

019

021

024

025

027 028

029

031

033

034

037

040

041

042

043

044

046

047

048

051 052 Paper under double-blind review

### **ABSTRACT**

The development lifecycle of generative AI systems requires continual evaluation, data acquisition, and annotation, which is costly in both resources and time. In practice, a desire for rapid iteration often makes it necessary to rely on synthetic annotation data because of its low cost, despite the potential for substantial bias. In this paper, we develop a rigorous theoretical framework for novel, cost-aware evaluation pipelines that actively balance the use of a cheap, but often inaccurate, weak rater—such as a model-based *autorater* that is designed to automatically assess the quality of generated content—with a more expensive, but also more accurate, strong rater such as a human annotator. Building on recent work in active and predictionpowered statistical inference, we theoretically derive a family of cost-optimal policies for allocating a given annotation budget between weak and strong raters so as to maximize statistical efficiency. Next, using synthetic and real-world data, we empirically characterize conditions under which these types of policies can yield significant improvements over classical methods. Finally, we find that practical approximations of the theoretically optimal policies can achieve the same estimation precision at a far lower total annotation budget than standard evaluation methods, especially in tasks where there is high variability in the difficulty of examples.

## 1 Introduction

Accurately and efficiently evaluating generative AI (GenAI) systems is a core technical challenge, both for model development and for reliable model deployment. In this paper, we introduce new statistical tools for active, cost-sensitive model evaluation. Specifically, we develop evaluation pipelines that dynamically annotate data using a mix of weak and strong annotation options in a way that is aware of their relative costs and strengths. The core idea is to strategically balance inexpensive but potentially inaccurate annotations from a *weak rater* against more accurate, but also more costly, annotations from a more sophisticated *strong rater* alternative. Our goal will be to use the weak raters to help give unbiased estimates of the mean of the strong rater's judgments. This is a key target for many AI evaluation applications, as it captures fundamental metrics like model accuracy, win-rate, or hallucination rate.

The exact composition of the weak and strong raters is flexible; for example, the weak rater might be a small AI model or rule-based heuristic, while the strong rater might be a larger AI model, an AI model with tools or larger inference-time reasoning capabilities, a human, or even the consensus of multiple expert humans. The cost of the evaluation might then be measured in compute, latency, or dollars. Active evaluation aims to minimize cost by selectively obtaining expensive annotations only when they are informative, relying on the cheaper option otherwise. All of the annotations are then combined using statistically principled, unbiased methods to yield reliable, yet cost-effective, performance metrics.

Combining different data sources to improve evaluation quality is not new: in particular, the use of cheap but biased metrics as control variates to improve statistical efficiency in model evaluation has been explored before from various perspectives (Angelopoulos et al., 2023a;b; Boyeau et al., 2024; Chaganty et al., 2018; Chatzi et al., 2024; Fisch et al., 2024; Jung et al., 2025; Saad-Falcon et al., 2024; Zrnic & Candès, 2024). Here, our main technical contribution is a theoretical framework for cost-optimal active evaluations—unbiased evaluation algorithms that strategically choose when to deploy the strong rater as opposed to the weak rater in order to achieve accurate evaluations at low cost. Informally, these policies solve the following constrained optimization problem:

maximize Accuracy of the evaluation,

subject to Cost of the evaluation remaining below a budget B.

We derive these optimal policies via new technical extensions and combinations of modern techniques in statistics, namely, active statistical inference (Zrnic & Candès, 2024) and prediction-powered inference (PPI; Angelopoulos et al., 2023a;b; Zrnic & Candès, 2024). As we will prove, the resulting oracle policies—which represent the best strategies we can hope to achieve in theory—depend on (i) the rater costs, but also (ii) task-specific distributional properties (like the weak rater's error) that are often unknown in practice. Nevertheless, using the form of these optimal policies as a guiding foundation, we are able to test and analyze empirical approximations that operate by first estimating these unknown quantities from data (e.g., using a "burn-in" set), and then use the theoretical form of the optimal policy with the estimated parameters plugged in (and provide bounds on the optimality gap). Empirically, we demonstrate that this practical approach can achieve substantial savings over passive strategies, although we also highlight important open challenges for future work that naturally arise due to "cold-start" issues, as well as imperfections of existing autorater models and their uncertainty estimates.

Finally, though AI model evaluation is the primary motivation and focus in this paper, we note that our framework also extends to general convex M-estimation problems (in any domain). See Appendix B.

Related work. Prediction-powered inference (PPI; Angelopoulos et al., 2023a;b; Zrnic & Candès, 2024) is the technique of combining a small number of trusted observations with predictions from a machine learning system for the purpose of statistical estimation. Its core statistical principles are closely related to control variate estimators (Chaganty et al., 2018; Ripley, 1987) as well as semiparametric inference with missing data (Chernozhukov et al., 2018; Robins & Rotnitzky, 1995; Tsiatis, 2006). Recently, a body of work has explored applying PPI to the evaluation of GenAI systems, where human annotations are combined with "autorater" outputs (Boyeau et al., 2024; Chatzi et al., 2024; Egami et al., 2023; Fisch et al., 2024; Saad-Falcon et al., 2024); though it has also been noted that the sample efficiency gained is limited when the autorater is not sufficiently accurate (Dorner et al., 2025; Thakur et al., 2025). A natural extension of PPI is to actively select a fixed number of examples on which to obtain trusted observations, while deferring the remaining examples to the autorater (Gligorić et al., 2024; Zrnic & Candès, 2024). Roughly speaking, these approaches sample human annotations with probability proportional to the uncertainty of the autorater. However, they work only in a restricted setting in which the ratio of expensive to cheap ratings, n/N, is fixed in advance, and then pick the optimal policy subject to that constraint. No guidance is given as to what this ratio should be based on the relative costs of the ratings, or even what the total number of examples N should be.

**Contributions.** Our work extends this literature both theoretically and empirically. Our core theoretical contribution is the derivation of error-minimizing sampling rules under cost constraints. That is, previous methods have a fixed ratio n/N and a policy that maximizes accuracy under that fixed ratio, while our policy maximizes accuracy subject to a cost constraint by optimizing everything including the ratio n/N. We theoretically derive two forms of optimal policies: (i) the best fixed sampling rate (Proposition 1), and (ii) the best active sampling rule that depends on covariates (Proposition 2). One additional novelty of our work is that it improves upon the policy proposed by Zrnic & Candès (2024) by accounting for the constraint that the policy must lie in [0,1] for all values of x. Finally, Appendix B includes further theoretical innovations, such as an extension to convex M-estimators and an optimal method for selecting the covariate x (as opposed to only the label, as considered in the prior work).

On the empirical end, we extend the scope of the standard PPI framework to heterogeneous model evaluation settings involving two distinct rating sources, each with a different cost-performance profile. This goes beyond the typical "human-vs-LLM" scenario described above, and encompasses any situation where less expensive, less accurate ratings are combined with more expensive, more accurate ones, even if both sources are automated (e.g., smaller vs. larger modelsm, or more vs. less inference-time reasoning). In Sections 3 and 4, we present an extensive empirical investigation into the conditions under which these new sampling rules prove beneficial over classical estimation. Specifically, we identify that the success of our framework is determined by: (a) the overall error of the weak rater, (b) the overall variance of the target strong rater, and (c) the heteroskedasticity of the weak rater's errors.

#### 2 Cost-optimal annotation policies

We now describe our methods for constructing active, cost-optimal evals. The methods rely on one critical ingredient: an *annotation policy*  $\pi$ . The job of the annotation policy is to look at the input and decide whether it should be labeled by the expensive rater. The theory in this section derives the **optimal policies under different restrictions on the policy space**. These policies are *oracle* policies—we

prove that they depend on properties of the data distribution, some of which are impossible to know in advance. As described in Section 1, the point of this section is to tell us **what kinds of policies we should be targeting, not how to find them**; later, we will explore how to estimate them in practice.

#### 2.1 BASIC NOTATION

 We observe inputs  $X \sim P_X$  from some space  $\mathcal{X}$  and distribution  $P_X$ : in the setting of LLMs, we think of the input X as containing the prompt as well as the response from one or multiple LLMs. Our goal is to approximate an expensive rating  $h(X) \in \mathbb{R}$ , such as a human preference, with a cheap automated evaluator  $g(X) \in \mathbb{R}$ ; for notational convenience we define  $H \triangleq h(X)$  and  $G \triangleq g(X)$ . In our setup, querying H and G cost  $c_h$  and  $c_g$ , respectively. We seek to query H only when it is "worth the cost".

We consider a sequential setting: for every  $t \in \mathbb{N}$ , we observe i.i.d.  $X_t \sim P_X$  and  $G_t \sim P_{G|X}$ . Upon observing  $X_t$ , we then have the option to query  $H_t \sim P_{H|X}$ . Our objective is to estimate  $\theta^* = \mathbb{E}[H]$ , the mean target rating. To this end, we develop estimators that efficiently sample only the data points for which  $H_t$  is needed, and stop sampling after a certain budget is exhausted. Define the random variable  $\xi_t \sim \mathrm{Bern}(\pi_t(X_t))$ , which is the indicator of whether we sampled  $H_t$ . It equals 1 with probability  $\pi_t(X_t)$ , and we have the freedom to choose the annotation policy  $\pi_t$  based on the previous data we have seen so far. We estimate  $\theta^*$  with the following unbiased estimator, defined for all  $T \in \mathbb{N}$ :

$$\hat{\theta}_T = \frac{1}{T} \sum_{t=1}^T \Delta_t \quad \text{where } \Delta_t = G_t + (H_t - G_t) \frac{\xi_t}{\pi_t(X_t)}. \tag{1}$$

Here  $\pi_t \in \Pi$  for some policy class  $\Pi$ . If  $\Pi$  is left unspecified, it should be assumed that  $\pi_t$  can be any function with range (0,1]. This is the sequential estimator from Zrnic & Candès (2024): the difference will be in how we set  $\pi_t$  to balance labeling costs. In general, the annotation policy  $\pi_t$  is allowed to change arbitrarily online as a function of past data, as is the predictor g. For simplicity, we will focus on the setting where the parameters of  $\pi$  and g remain fixed throughout and are *not* updated online, or as if we are updating in batches; however our results will also hold asymptotically when  $\pi$  and g are updated online and converge. We use the notation  $\hat{\theta}_T^{\pi}$  to denote the estimator in (1) with a fixed policy  $\pi$ , i.e., where  $\pi_t = \pi$ ,  $\forall t \in T$ . To calculate the cost and error of our estimator, we additionally define:

$$\operatorname{Error}_T(\pi) \triangleq \mathbb{E}\left[\left(\hat{\theta}_T^{\pi} - \theta^*\right)^2\right] = \frac{1}{T}\left(\operatorname{Var}(H) - \mathbb{E}[(H - G)^2] + \mathbb{E}\left[(H - G)^2 \frac{1}{\pi(X)}\right]\right), \quad (2)$$

$$\operatorname{Cost}_T(\pi) \triangleq T(c_h \mathbb{E}\left[\pi(X)\right] + c_g).$$

These functions describe the mean squared error and expected cost of the estimator with annotation policy  $\pi$  as a function of time, and our goal will be to minimize one subject to a constraint on the other. When we refer to a budget on the cost, it will be denoted as B. Furthermore, we note that, for convenience, the cost-optimized policies that we present in the remainder of this section will relax the constraint that the stopping time  $T^{\text{stop}}$  at which  $\text{Cost}_{T^{\text{stop}}}(\pi)$  is just under budget must be an integer, though this does not have a significant effect on the optimization for large enough budgets B where  $T^{\text{stop}} \gg 1$ . Some additional treatment for this restriction is included in Appendix B.

#### 2.2 OPTIMAL RANDOM ANNOTATION

The simplest annotation policy does not depend on X, and simply queries H with some fixed probability, which we denote as  $\pi(x) = p$  for a sampling rate  $p \in (0,1]$ . In other words, we let  $\pi \in \Pi^{\mathrm{random}} = \{x \mapsto p : p \in (0,1]\}$ . When p is too large, the cost is too high; when p is too small, the error blows up. Our job is to choose the optimal balance, and the next result shows it has a simple, explicit form that depends on the cost ratio  $c_q/c_h$  and the error of G compared to the variance of H.

**Proposition 1.** Let  $(X_1, G_1, H_1), \ldots, (X_T, G_T, H_T), T \in \mathbb{N}$ , be an i.i.d. sequence of real-valued random variables with joint distribution P, and define Error, Cost, and  $\Pi^{\mathrm{random}}$  as above. Assume that  $\mathbb{P}(G_1 = H_1) < 1$  and that  $c_h > c_q > 0$ , and define the optimization problem

$$\underset{\pi \in \Pi^{\text{random}}, \ T^{\text{stop}} \in \mathbb{R}_{>0}}{\text{minimize}} \quad \text{Error}_{T^{\text{stop}}}(\pi) \quad \text{subject to} \quad \text{Cost}_{T^{\text{stop}}}(\pi) \le B.$$
(3)

Then the solution to Problem (3) for all  $x \in \mathcal{X}$  is

$$\pi_{\text{random}}(x) = \begin{cases} \sqrt{\frac{c_g}{c_h} \frac{\mathbb{E}[(H-G)^2]}{\text{Var}(H) - \mathbb{E}[(H-G)^2]}} & \text{if } \mathbb{E}[(H-G)^2] < \frac{c_h}{c_h + c_g} \text{Var}(H) \\ 1 & \text{otherwise.} \end{cases}$$
(4)

We can make a few observations about  $\pi_{\mathrm{random}}$ . First, if the mean squared error of the weak rater G is greater than the variance of H (or more precisely, more than a  $c_h/(c_h+c_g)$  fraction of the variance of H), then it is not helpful—and we should simply choose to query H all the time. If  $\mathrm{MSE}(H,G)$  is sufficiently low, however, then the rate at which we sample H varies inversely with both the ratio of  $\mathrm{Var}(H)$  to  $\mathrm{MSE}(H,G)$  and the ratio of the cost of H to the cost of H. This makes intuitive sense: if the target label H is high variance but our "weak" rater H is in fact a fairly "strong" rater (in that it produces similar ratings to those of H), then we should primarily exploit H is low cost, high-quality predictions, while sampling H at just a low rate to correct for any minor bias that arises.

# 2.3 OPTIMAL ACTIVE ANNOTATION

Next, we study policies that depend on X; i.e., they query H with some probability that depends on X. This strategy can greatly improve statistical power when the error distribution is heteroskedastic in X; for example, when some prompts are much harder than others. In this setting, it makes sense for  $\pi$  to depend on X, and to ask for advanced rating help more often when G is likely to be wrong. Towards that end, we define our annotation policy class to be  $\pi \in \Pi = \{x \mapsto f(x) : f(x) \in (0,1]; \forall x \in \mathcal{X}\}$ , which is the set of annotation policies placing a strictly positive amount of sampling mass on each query. As the next proposition shows, the optimal policy in this setting will depend on the uncertainty of the weak rater,  $u(x) \triangleq \mathbb{E}[(H-G)^2 \mid X=x]$ , expressed as the expected mean squared conditional error given X=x. For notational convenience, we also define the random variable  $U \triangleq u(X)$ .

**Proposition 2.** In the same setting as Proposition 1, define  $\Pi$  as above, let  $\mathcal{X}$  be discrete, and additionally define the optimization problem

$$\underset{\pi \in \Pi, \ T^{\text{stop}} \in \mathbb{R}_{>0}}{\text{minimize}} \quad \text{Error}_{T^{\text{stop}}}(\pi) \quad \text{subject to} \quad \text{Cost}_{T^{\text{stop}}}(\pi) \le B.$$
(5)

*Define the scaled and clipped policy,*  $\pi_{clip}$ *, as:* 

$$\pi_{\text{clip}}(x;\tau) = \min\left(\gamma^*(\tau)\sqrt{u(x)},1\right) = \begin{cases} \gamma^*(\tau)\sqrt{u(x)} & \textit{if } \sqrt{u(x)} \leq \tau \\ 1 & \textit{otherwise,} \end{cases}$$

where  $c_h>c_g>0$  and  $\gamma^*(\tau)\in \left(0,\frac{1}{\tau}\right]$  is defined as

$$\gamma^*(\tau) = \min\left(\sqrt{\frac{c_g/c_h + \mathbb{P}(U > \tau^2)}{\left(\operatorname{Var}(H) - \mathbb{E}[U\mathbb{1}\{U \le \tau^2\}]\right)_+}}, \frac{1}{\tau}\right).$$

Then the solution to Problem (5) is  $\pi_{\text{active}}(x) = \pi_{\text{clip}}(x; \tau^*)$ , where  $\tau^* > 0$  is the solution to

$$\tau^* = \operatorname*{argmin}_{\tau \in \mathbb{R}_{>0}} \left( c_h \mathbb{E}[\pi_{\text{clip}}(x;\tau)] + c_g \right) \left( \operatorname{Var}(H) + \mathbb{E}\left[ U \left( \pi_{\text{clip}}(x;\tau)^{-1} - 1 \right) \right] \right).$$

**Remark 3.** The final optimization problem presented for the clipping threshold  $\tau^*$  is non-convex and has no analytical solution. However, because it is a 1-dimensional optimization problem, we can coarsely discretize and optimize  $\tau$  via simple grid search.

On a technical level, the solution in Proposition 2 has a similar form to the active sequential estimator proposed in Zrnic & Candès (2024), but with an optimized proportionality constant, as well as additional clipping to rigorously account for the constraints on  $\pi(x) \in (0,1]$ . The latter point is particularly important, as it is not accounted for in prior work. In contrast to the fixed, prespecified ratio prescribed by prior work, in Appendix B.7 we show how the *cost-optimal* target ratio of expensive to cheap ratings can be as extreme as 0 or 1, depending on the cost ratio of G to H.

While the form of  $\pi_{\rm active}$  is more complex than that of  $\pi_{\rm random}$ , it still admits a fairly straightforward interpretation: for some confidence threshold  $\tau^*$  below which the conditional mean squared error of G over all confident data points with  $\sqrt{u(x)} \leq \tau^*$  is sufficiently low, we sample proportional to  $\sqrt{u(x)}$ . On the remaining highly uncertain examples where  $\sqrt{u(x)} > \tau^*$ , we always use H, and ignore G. The exact threshold  $\tau^*$  depends on the distributions of H and G, and their cost-ratio.

We can also observe that Proposition 2 is a direct generalization of Proposition 1. When X is independent of  $(H-G)^2$  so that  $u(x) = \mathbb{E}[(H-G)^2] \ \forall x \in \mathcal{X}$ , the policy  $\pi_{\text{active}}$  reduces to  $\pi_{\text{random}}$ :

$$\underbrace{\gamma^*(\tau^*)\sqrt{u(x)}}_{\text{optimal active}} = \sqrt{\frac{c_g}{c_h}} \frac{\mathbb{E}[(H-G)^2 \mid X=x]}{\text{Var}(H) - \mathbb{E}[(H-G)^2]} = \underbrace{\sqrt{\frac{c_g}{c_h}} \frac{\mathbb{E}[(H-G)^2]}{\text{Var}(H) - \mathbb{E}[(H-G)^2]}}_{\text{optimal random}}.$$

The intuitive conclusion is that active querying can help if the conditional squared error of G has significant variance to it (i.e., there exist some regions of  $\mathcal X$  where G has a much higher level of agreement with H than on other regions of  $\mathcal X$ , such as on easy vs. hard examples). This can be contrasted with the optimal random policy,  $\pi_{\mathrm{random}}$ , from (4): there we sample at a fixed rate for each X, where that rate depends only on G's average error with respect to H across all types of inputs.

# Takeaways: Cost-optimal annotation policies

We derive two policies for sampling the expensive rating H given a budget B:  $\pi_{\rm random}$  chooses the optimal fixed probability  $p^* \in (0,1]$ , while  $\pi_{\rm active}$  defines an optimal input conditional probability  $\pi_{\rm active}(x) \in (0,1]$ . Both navigate the following trade-off: reducing  $\mathbb{E}[\pi(X)]$  increases the total number of samples we can afford to rate at all, but not querying H when G is inaccurate increases variance. Finally, both policies converge to the baseline estimator (i.e.,  $\pi_{\rm base}(x) = 1$ ) when the error of G is too high relative to the variance of H.

### 3 COMPARING COST-OPTIMAL POLICIES IN SIMULATED SETTINGS

The estimation error of the optimal policies presented in Section 2 depends on the distributions of the expensive target label H, the cheap estimated label G, and the cost-ratio  $c_g/c_h$  for querying G versus H. To build a clearer understanding of how these variables influence the performance of our proposed policies, we now conduct a series of carefully controlled experiments on simulated data. Note that since all of the key distributional quantities (i.e., Var(H), MSE(H,G), etc) are known in the synthetic settings we consider in this section, we are also able to compute  $\pi_{active}$  and  $\pi_{random}$  exactly—as opposed to the more difficult real-world data settings we will tackle next in Section 4.

# 3.1 METRICS

To measure the relative performance of annotation policy  $\pi_1$  vs  $\pi_2$ , we compute the *ratio* of their errors at  $T_i^{\text{stop}}$ . Once again relaxing the restriction that  $T_i^{\text{stop}} \in \mathbb{N}$ , we compute a budget-free approximation based on the expression for  $\text{Error}_{T_i^{\text{stop}}}(\pi)$ , where  $T_i^{\text{stop}} = B/(c_h \mathbb{E}[\pi_i(X)] + c_g)$ :

$$\mathsf{ErrorRatio}(\pi_1, \pi_2) \triangleq \frac{\left(c_h \mathbb{E}[\pi_1(X)] + c_g\right) \left( \mathrm{Var}(H) - \mathbb{E}[(H - G)^2] + \mathbb{E}\left[ (H - G)^2 \frac{1}{\pi_1(X)} \right] \right)}{\left(c_h \mathbb{E}[\pi_2(X)] + c_g\right) \left( \mathrm{Var}(H) - \mathbb{E}[(H - G)^2] + \mathbb{E}\left[ (H - G)^2 \frac{1}{\pi_2(X)} \right] \right)}.$$

Note that while  $\operatorname{ErrorRatio}(\pi_1,\pi_2)$  does not depend on the budget, it does implicitly depend on  $P_X$  as well as  $P_{H|X}$  and  $P_{G|X}$ . We will focus on  $\operatorname{ErrorRatio}(\pi_{\operatorname{active}},\pi_{\operatorname{base}})$ , the error ratio of the active estimator to the baseline estimator which only uses H, as well as  $\operatorname{ErrorRatio}(\pi_{\operatorname{active}},\pi_{\operatorname{random}})$ , which compares the active estimator to the estimator that doe not depend on X. Note that for  $\operatorname{ErrorRatio}(\cdot,\pi_{\operatorname{base}})$ , we disregard  $e_g$  for  $\pi_{\operatorname{base}}$ , and replace the denominator with  $e_h\operatorname{Var}(H)$ .

#### 3.2 Gaussian data

We construct an experiment where we change  $\mathrm{Var}(H)$ ,  $\mathrm{MSE}(H,G)$ , and  $\mathrm{Var}(U)$  independently (recall that we introduced  $U \triangleq u(X) = \mathbb{E}[(H-G)^2 \mid X]$  in Section 2.3). First, we draw  $H \sim \mathcal{N}(0,\nu)$  so that  $\mathbb{E}[H] = 0$  and  $\mathrm{Var}(H) = \nu$ . Then we draw  $U \sim \mathrm{Gamma}\left(\mu^2/\eta,\eta/\mu\right)$  so that  $\mathrm{MSE}(H,G) = \mathbb{E}[U] = \mu$  and  $\mathrm{Var}(U) = \eta$ . Finally, we set  $G = H + \sqrt{U}$ .

Results are shown in the top row of Figure 1. The left panel plots the error ratio of  $\pi_{\text{active}}$  to  $\pi_{\text{base}}$  as a function of  $\mathrm{MSE}(H,G)$  and for different  $\mathrm{Var}(H)$ , while keeping  $\mathrm{Var}(U)=0.5$ . As expected, the error of  $\pi_{\mathrm{active}}$  increases with the  $\mathrm{MSE}(H,G)$ , with the rate of increase influenced by  $\mathrm{Var}(H)$ . When

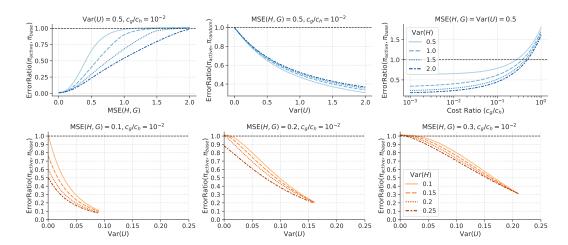


Figure 1: Results on the Gaussian (top) and Bernoulli (bottom) settings while varying MSE(H, G), Var(U), and  $c_g/c_h$ . Each line plots a different value of Var(H), where we choose values that are representative of low, medium, or high variance settings compared to MSE(H, G). In the Bernoulli setting,  $MSE(H, G) = \mathbb{E}[H \neq G]$ , and Var(U) can be at most MSE(H, G)(1 - MSE(H, G)).

 $\mathrm{MSE}(H,G)$  is large relative to  $\mathrm{Var}(H)$ ,  $\pi_{\mathrm{active}}$  provides no benefit over  $\pi_{\mathrm{base}}$ . The middle panel plots the error ratio of  $\pi_{\mathrm{active}}$  to  $\pi_{\mathrm{random}}$  while varying  $\mathrm{Var}(U)$  for a fixed  $\mathrm{MSE}(H,G)$ . For small values of  $\mathrm{Var}(U)$ , the conditional error in G is nearly the same everywhere, and there is no benefit to using  $\pi_{\mathrm{active}}$  over  $\pi_{\mathrm{random}}$ . Larger values of  $\mathrm{Var}(U)$ , however, lead to a performance advantage for  $\pi_{\mathrm{active}}$ . The right panel plots the error ratio of  $\pi_{\mathrm{active}}$  to  $\pi_{\mathrm{base}}$  while keeping  $\mathrm{MSE}(H,G)$  and  $\mathrm{Var}(U)$  fixed, but varying  $c_q/c_h$ . As expected,  $\pi_{\mathrm{active}}$  is most effective when  $c_q \ll c_h$ .

# 3.3 BERNOULLI DATA

While the Gaussian setting above is informative, in many typical situations H is bounded, such as when H is a binary, Bernoulli rating for win-rate or accuracy estimation. This creates a more difficult setting for  $\pi_{\rm active}$ , since both  ${\rm Var}(H)$  and  ${\rm Var}(U)$  are upper-bounded by 0.25 for Bernoulli H. In fact, in binary settings,  ${\rm MSE}(H,G)$  and  ${\rm Var}(U)$  are in tension: the more accurate G is, the lower the variance of its errors, and  $\pi_{\rm active}$  will be limited in terms of any relative benefit it can provide over  $\pi_{\rm random}$ . The same is also true for when G is uniformly inaccurate. To better analyze this kind of setting, we construct a binary dataset where first we draw  $H \sim {\rm Bern}(0.5 + \sqrt{0.25} - \nu)$ , so that  ${\rm Var}(H) = \nu$ . Next, we draw H from a Beta distribution with mean H and variance H0, where H1 with probability H2 to get the prediction H3. Finally, we flip H4 with probability H4 to get the prediction H5. The same is also Bernoulli with H5. Finally, we flip H6.

Results are shown in the bottom row of Figure 1 for  $\pi_{\text{active}}$  vs  $\pi_{\text{base}}$  (see Appendix D for  $\pi_{\text{active}}$  vs.  $\pi_{\text{random}}$ ). As in the Gaussian setting, the error ratio of  $\pi_{\text{active}}$  to  $\pi_{\text{base}}$  improves dramatically with larger Var(U). Note that the active and random estimator are the same when Var(U) = 0. For larger MSE(H,G), Var(U) must also be increasingly large for  $\pi_{\text{active}}$  to improve significantly over  $\pi_{\text{base}}$ . Indeed, on the right-hand side of the bottom row of Figure 1 where MSE(H,G) > Var(H), we can see that  $\pi_{\text{random}}$  provides no benefits over  $\pi_{\text{base}}$ ; that is,  $\text{ErrorRatio}(\pi_{\text{random}}, \pi_{\text{base}}) = 1$  when Var(U) = 0, which corresponds to the fixed-rate sampling policy as noted earlier. When  $\text{Var}(U) \gg 0$ , however,  $\pi_{\text{active}}$  can obtain substantially lower estimation error than  $\pi_{\text{base}}$ . Still, unlike the earlier Gaussian data, the best active error ratio in this setting is bounded from below by MSE(H,G), and is achieved when U has maximum variance (which is also bounded).

## Takeaways: Performance characteristics of cost-optimal annotation policies

In general, the following properties hold for active annotation versus standard annotation (similar findings for random): (i) as the error of G,  $\mathrm{MSE}(H,G)$ , increases, the benefit **decreases**; (ii) as the variance of the conditional squared-error of G,  $\mathrm{Var}(U)$ , increases, the benefit **increases**; and (iii) as the cost ratio,  $c_q/c_h$ , of G relative to H increases, the benefit **decreases**.

# 4 ESTIMATING COST-OPTIMAL POLICIES IN PRACTICE

The theoretical results in Section 2 derive optimal annotation policies under the assumption that the key distributional parameters governing the relationship between the expensive rater (H) and the cheap rater (G) are known. In reality, these parameters must be estimated (imperfectly; see Appendix B.5 for theoretical error analysis). Furthermore, the optimal threshold  $\tau^*$  and scaling factor  $\gamma(\tau^*)$  for the active policy  $\pi_{\rm active}$  in Proposition 2 also depend on conditional versions of these unknown quantities (e.g., the conditional MSE,  $\mathbb{E}[(H-G)^2 \mid U \leq \tau]$ ). Some of these estimates can be derived automatically from the model itself, for example if  $g(x) \in [0,1]$  is a binary classifier, we may choose u(x) = g(x)(1-g(x)), which is equal to  $\mathbb{E}[(H-g(x))^2 \mid X=x]$  when  $g(x) = \mathbb{P}(H=1 \mid X=x)$ . Alternatively, u(x) can be a separate prediction, such as by asking an LLM for its confidence (Kadavath et al., 2022; Xiong et al., 2024). For the key parameters  $\mathrm{Var}(H)$ ,  $\mathrm{MSE}(H,G)$ ,  $\gamma(\tau^*)$ , and  $\tau^*$ , we explore estimating them using the following approaches:

Policy transfer from related datasets (A1). Here we *transfer* all parameters necessary for  $\pi$  from a separate, but related, dataset. For example, in Section 4.2, we use data from the Chatbot Arena dataset (Zheng et al., 2023; Chiang et al., 2024) to estimate the win-rate of GPT-4 over Claude 2.1, but transfer parameters for  $\pi_{\rm random}$  and  $\pi_{\rm active}$  from a separate set of comparisons between different available models. We also calibrate G using Platt scaling (Platt, 1999) on the transfer dataset.

Policy burn-in on the first  $n_b$  examples (A2). When a suitable transfer dataset is not available as in A1, we can take a hybrid approach where we start by sampling H for the first  $n_b=200$  examples with probability 1, and then use them to estimate the parameters necessary for  $\pi_{\text{active}}$  and  $\pi_{\text{random}}$ . We also calibrate G using Platt scaling on these  $n_b$  examples. As a fair comparison to the baseline method of only using H, we also allow these  $n_b$  samples to be used as additional data for estimating  $\theta = \mathbb{E}[H]$ . Specifically, we use the (estimated) inverse-variance-weighted average of the annotation policy  $\pi$ 's estimate,  $\hat{\theta}_T^{\pi}$ , and the classical estimate on the burn-in data,  $\hat{\theta}_{n_b} = \frac{1}{n_b} \sum_{i=1}^{n_b} H_i$ ,

$$\hat{\theta}_{T^{\text{stop}}+n_b}^{\pi} = \frac{\widehat{\text{Var}}(\hat{\theta}_{T^{\text{stop}}}^{\pi})}{\widehat{\text{Var}}(\hat{\theta}_{n_b}) + \widehat{\text{Var}}(\hat{\theta}_{T^{\text{stop}}}^{\pi})} \hat{\theta}_{n_b} + \frac{\widehat{\text{Var}}(\hat{\theta}_{n_b})}{\widehat{\text{Var}}(\hat{\theta}_{n_b}) + \widehat{\text{Var}}(\hat{\theta}_{T^{\text{stop}}}^{\pi})} \hat{\theta}_{T^{\text{stop}}}^{\pi},$$

where  $\widehat{\mathrm{Var}}(\cdot)$  is also estimated on the burn-in data. Note that  $\widehat{\theta}_{T^{\mathrm{stop}}+n_b}$  is still unbiased. To get a sense of how close to optimal the estimated policies are, we also compute an **Oracle**:  $\pi_{\mathrm{active}}$  with parameters computed using the whole dataset, and u(x) taken directly as  $|h(x) - g(x)|^2$ .

### 4.1 METRICS

We compare the baseline method  $\pi_{\text{base}}$  of always sampling H with the random policy  $\pi_{\text{random}}$  and the active policy  $\pi_{\text{active}}$ . For each policy, we compute the **mean squared error**,  $\mathbb{E}[(\hat{\theta}_T^\pi - \theta^*)^2]$ , for a range of budgets B ( $c_h$  is normalized to be one "cost unit"), with 95% bootstrap CIs shown over 2k trials. We then compute the **mean effective budget**, which we define as the budget B' required for  $\pi_{\text{base}}$  to achieve the same MSE as the given policy  $\pi$  at a budget B. If  $\pi$  is more cost-effective than  $\pi_{\text{base}}$ , then B' will be larger than B (higher is better). Finally, we also compute the **mean cost savings** for a given mean-squared error, which we define as the budget deficit relative to  $\pi_{\text{base}}$  required to achieve that target error (higher is better). By definition, we have that the mean effective budget for  $\pi_{\text{base}}$  is the line y = x (since B' = B always), while the cost savings for  $\pi_{\text{base}}$  is 0.

### 4.2 Datasets

We report experimental results on four datasets, which span a diverse range of raters and distributional characteristics. For each task, we calculate  $\theta^* = \mathbb{E}[H]$  using the full dataset. For simplicity we assume that the total number of data points  $X_t$  is at least  $\lceil B/c_g \rceil$ , and sample with replacement from the original dataset if not. We leave treatment of datasets where  $T^{\text{stop}} \leq T^{\text{max}}$  (and the constraint is active) to future work. See Appendix D for results on three additional datasets: Attributed Question Answering (Bohnet et al., 2023), ImageNet (Deng et al., 2009), and Seahorse (Clark et al., 2023).

**Chatbot Arena.** The Chatbot Arena dataset (Zheng et al., 2023; Chiang et al., 2024) evaluates LLMs via pairwise comparisons (i.e., eliciting preferences for response A vs. B from two models for the same query). Among the 64 models present in the 57k total comparisons, we focus on estimating the winrate of GPT-4 (specifically, the 11/06 preview model) vs. Claude 2.1, as they are both strong models,

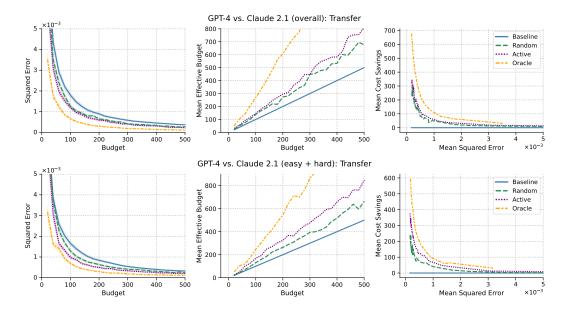


Figure 2: Results for estimating the win-rate of GPT-4 vs. Claude-2.1 on the Chatbot Arena dataset when using policy transfer (see approach A1 in Section 4). Both  $\pi_{\rm random}$  and  $\pi_{\rm active}$  substantially improve estimation quality over  $\pi_{\rm base}$  for a given budget. Consistent with our theory,  $\pi_{\rm active}$ 's performance benefits are substantially magnified on the heterogenous easy/hard split (bottom row).

and also have the most pairwise comparisons in the dataset (1073 total), which allows us to get a reliable estimate of  $\theta^*$ . We model H via the majority preference from 10 Gemini 1.5 Flash (Gemini Team, 2024) evals (5 samples each comparing A vs. B and B vs. A to mitigate position bias). G is the win probability predicted by a Gemma-3 4B model (Gemma Team, 2025) which has been fine-tuned on the other model comparisons from the dataset to predict the Gemini labels. U is computed as G(1-G).

Chatbot Arena (estimated easy/hard split). In an effort to include a dataset with more (identifiable) heteroskedasticy, we also include a filtered version of the GPT-4 versus Claude 2.1 task described above, where we construct a dataset slice containing only the examples corresponding to the bottom 25% and top 25% of Gemma's uncertainty estimates (we use U as the metric). While partly manipulated, this scenario is designed to test for potential gains from actively choosing when to query the expensive rater, as per the intuition from Section 3, where it was shown how higher Var(U) benefits active policies (though note this may not be true if the estimated U is inaccurate).

### 4.3 RESULTS

 Figure 2 shows results for the Chatbot Arena datasets using the *transfer* approach (A1), while Figure 3 shows results for all datasets using the *burn-in* approach (A2). As expected, the absolute improvement for both  $\pi_{\text{active}}$  over  $\pi_{\text{random}}$  and  $\pi_{\text{random}}$  over  $\pi_{\text{base}}$  is greatest in the transfer setting in Figure 2, where the parameters of  $\pi_{\text{random}}$  and  $\pi_{\text{active}}$  can be approximated in advance. In particular, to achieve a root mean-squared error (RMSE) of 0.05,  $\pi_{\text{active}}$  requires only  $\approx 40\%$  of the budget required by  $\pi_{\text{base}}$  in the overall setting of Chatbot Arena, and only  $\approx 50\%$  of the budget in the easy/hard setting. These cost savings become even more pronounced the more precise (i.e., lower MSE) the estimates are required to be. In Figure 3, where the first  $n_b = 200$  examples are fully labeled in order to estimate the parameters of  $\pi_{\text{random}}$  and  $\pi_{\text{active}}$ , the absolute *difference* in MSE is smaller for  $\pi_{\text{random}}$  and  $\pi_{\text{active}}$  over  $\pi_{\text{base}}$ , though the subsequent cost savings over  $\pi_{\text{base}}$  for achieving lower and lower MSE (that is, past the MSE of the initial  $n_b$  sample estimate) are consistent.

As also predicted by our theory, the results in Figure 2 and Figure 3 show that the extent of the improvement in estimation accuracy varies per dataset (see also the additional results in Appendix D). In particular, the best results are obtained on the easy/hard split of the Chatbot Arena dataset, where (i) the weak annotator G is a good proxy of the strong annotator H (both are LLMs), and (ii) there

<sup>&</sup>lt;sup>1</sup>We also apply power tuning (Angelopoulos et al., 2023b) after all samples are collected. See Appendix B.2.

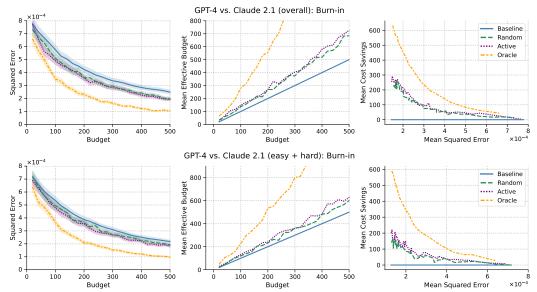


Figure 3: Results on Chatbot Arena using 200 examples as a burn-in to estimate policy parameters, and then switching to the initialized annotation policy thereafter (see approach A2 in Section 4; note that budgets B on the x-axis reflect the "additional" budget used *after* the burn-in examples). While the absolute differences in squared errors for the estimated means are smaller than in the transfer setting in Figure 2, both  $\pi_{\rm random}$  and  $\pi_{\rm active}$  still achieve consistent improvements over  $\pi_{\rm base}$ .

is more variability in the difficulty of examples according to the predicted U, resulting in a greater opportunity for improvement for  $\pi_{\text{active}}$ . On the other hand, while results on the homogeneous split of the Chatbot Arena dataset show improvements for  $\pi_{\text{random}}$  and  $\pi_{\text{active}}$  over  $\pi_{\text{base}}$ , the relative improvement of  $\pi_{\text{active}}$  over  $\pi_{\text{random}}$  is fairly small—indicating that while the weak annotator G that is used is relatively good on average, there is not much variability in its estimated uncertainty, u(x), on those distributions. To that point, when we compare to performance using the oracle active policy, it is also clear that the estimated u(x) is also far from perfect. Even on the datasets where the improvement due to the estimated active policy is small, the oracle policy which has knowledge of the true error of G often promises significant headroom: indicating that the working on better autorater uncertainty estimation is a promising and important direction for future work.

## Takeaways: Performance of practical approximations to cost-optimal policies on real data

Section 2 proved that the optimal policies depend on distributional parameters that must be estimated. How well they are estimated **does not affect the consistency or unbiasedness** of the overall estimator, but it does affect the policy's performance advantage over passive strategies. Yet while estimation is non-trivial, our experiments validate generic recipes that can **successfully approximate the optimal policy**—and yield policies with consistent gains over  $\pi_{\text{base}}$ .

#### 5 CONCLUSION

This paper introduces theory and practice for **cost-optimal active evals**, a framework that strategically combines cheap raters with more expensive, accurate alternatives to improve evaluation efficiency. We derive annotation policies that are optimal in the sense of minimizing expected error under annotation budget constraints, and we empirically characterize the conditions under which such policies yield improvements over non-hybrid (e.g., human-only) and non-active hybrid alternatives. However, we also show that the annotation policies that are *optimal in theory* are distribution dependent, and include a number of task-specific parameters that must be estimated. Furthermore, optimal active annotation depends on having an accurate uncertainty estimates, which can be uncalibrated for AI raters. Nevertheless, many realistic evaluation scenarios involve incrementally adding new models to existing benchmarks; and as shown in §4.3, policy transfer can work quite well. Furthermore, our results demonstrate that even when active sampling is difficult for the reasons outlined above, the simple—but optimal—fixed sampling rate policy that we derived consistently provides substantial improvements.

# REPRODUCIBILITY STATEMENT

All proofs of theoretical results are included in Appendix C. Implementation details for all of the empirical experiments are included in Appendix E. All datasets used in for the experiments in Section 4, and the additional experiments in Appendix D, are publicly available. The generation process for the synthetic datasets in Section 3 is described in detail in Sections 3.2 and 3.3.

#### REFERENCES

- Anastasios N Angelopoulos, Stephen Bates, Clara Fannjiang, Michael I Jordan, and Tijana Zrnic. Prediction-powered inference. *Science*, 382(6671):669–674, 2023a.
- Anastasios N Angelopoulos, John C Duchi, and Tijana Zrnic. PPI++: Efficient prediction-powered inference. *arXiv preprint arXiv:2311.01453*, 2023b.
- Bernd Bohnet, Vinh Q. Tran, Pat Verga, Roee Aharoni, Daniel Andor, Livio Baldini Soares, Massimiliano Ciaramita, Jacob Eisenstein, Kuzman Ganchev, Jonathan Herzig, Kai Hui, Tom Kwiatkowski, Ji Ma, Jianmo Ni, Lierni Sestorain Saralegui, Tal Schuster, William W. Cohen, Michael Collins, Dipanjan Das, Donald Metzler, Slav Petrov, and Kellie Webster. Attributed question answering: Evaluation and modeling for attributed large language models. *arXiv preprint: 2212.08037*, 2023.
- Pierre Boyeau, Anastasios N Angelopoulos, Nir Yosef, Jitendra Malik, and Michael I Jordan. AutoEval done right: Using synthetic data for model evaluation. *arXiv preprint arXiv:2403.07008*, 2024.
- Arun Chaganty, Stephen Mussmann, and Percy Liang. The price of debiasing automatic metrics in natural language evaluation. In Iryna Gurevych and Yusuke Miyao (eds.), *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 643–653, Melbourne, Australia, July 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-1060. URL https://aclanthology.org/P18-1060/.
- Ivi Chatzi, Eleni Straitouri, Suhas Thejaswi, and Manuel Gomez Rodriguez. Prediction-powered ranking of large language models. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 113096–113133. Curran Associates, Inc., 2024. URL https://proceedings.neurips.cc/paper\_files/paper/2024/file/cd47cd67caa87f5b1944e00f6781598f-Paper-Conference.pdf.
- Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James Robins. Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21(1):C1–C68, 01 2018. ISSN 1368-4221. doi: 10.1111/ectj.12097. URL https://doi.org/10.1111/ectj.12097.
- Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E. Gonzalez, and Ion Stoica. Chatbot arena: An open platform for evaluating llms by human preference. *arXiv preprint:* 2403.04132, 2024.
- Elizabeth Clark, Shruti Rijhwani, Sebastian Gehrmann, Joshua Maynez, Roee Aharoni, Vitaly Nikolaev, Thibault Sellam, Aditya Siddhant, Dipanjan Das, and Ankur Parikh. SEAHORSE: A multilingual, multifaceted dataset for summarization evaluation. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 9397–9413, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.584. URL https://aclanthology.org/2023.emnlp-main.584.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Florian E. Dorner, Vivian Yvonne Nastl, and Moritz Hardt. Limits to scalable evaluation at the frontier: LLM as judge won't beat twice the data. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=NO6Tv6QcDs.

Naoki Egami, Musashi Hinck, Brandon Stewart, and Hanying Wei. Using imperfect surrogates for downstream inference: Design-based supervised learning for social science applications of large language models. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 68589–68601. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper\_files/paper/2023/file/d862f7f5445255090de13b825b880d59-Paper-Conference.pdf.

- Adam Fisch, Joshua Maynez, R. Alex Hofer, Bhuwan Dhingra, Amir Globerson, and William W. Cohen. Stratified prediction-powered inference for effective hybrid evaluation of language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=8CBcdDQFDQ.
- Gemini Team. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context, 2024. URL https://arxiv.org/abs/2403.05530.
- Gemma Team. Gemma 3 technical report. arXiv preprint: arXiv 2503.19786, 2025. URL https://arxiv.org/abs/2503.19786.
- Kristina Gligorić, Tijana Zrnic, Cinoo Lee, Emmanuel J Candès, and Dan Jurafsky. Can unconfident llm annotations be used for confident conclusions? *arXiv preprint arXiv:2408.15204*, 2024.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Or Honovich, Roee Aharoni, Jonathan Herzig, Hagai Taitelbaum, Doron Kukliansy, Vered Cohen, Thomas Scialom, Idan Szpektor, Avinatan Hassidim, and Yossi Matias. True: Re-evaluating factual consistency evaluation. *arXiv* preprint: 2204.04991, 2022.
- Jaehun Jung, Faeze Brahman, and Yejin Choi. Trust or escalate: LLM judges with provable guarantees for human agreement. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=UHPngSTBPO.
- Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*, 2022.
- John C Platt. Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. In *Advances in large margin classifiers*, pp. 61–74. MIT Press, 1999.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020. URL http://jmlr.org/papers/v21/20-074.html.
- B. D. Ripley. Stochastic simulation. John Wiley & Sons, Inc., New York, NY, USA, 1987. ISBN 0-471-81884-4.
- James M. Robins and Andrea Rotnitzky. Semiparametric efficiency in multivariate regression models with missing data. *Journal of the American Statistical Association*, 90(429):122–129, 1995. ISSN 01621459. URL http://www.jstor.org/stable/2291135.
- Jon Saad-Falcon, Omar Khattab, Christopher Potts, and Matei Zaharia. ARES: An automated evaluation framework for retrieval-augmented generation systems. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 338–354, Mexico City, Mexico, June 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.naacl-long.20. URL https://aclanthology.org/2024.naacl-long.20/.
- Aman Singh Thakur, Kartik Choudhary, Venkat Srinik Ramayapally, Sankaran Vaidyanathan, and Dieuwke Hupkes. Judging the judges: Evaluating alignment and vulnerabilities in llms-as-judges. *ArXiv preprint: arXiv 2406.12624*, 2025. URL https://arxiv.org/abs/2406.12624.

Anastasios A. (Anastasios Athanasios) Tsiatis. *Semiparametric theory and missing data*. Springer series in statistics. Springer, New York, 2006. ISBN 9780387373454.

Aad W Van der Vaart. Asymptotic statistics, volume 3. Cambridge university press, 2000.

- Miao Xiong, Zhiyuan Hu, Xinyang Lu, YIFEI LI, Jie Fu, Junxian He, and Bryan Hooi. Can LLMs express their uncertainty? an empirical evaluation of confidence elicitation in LLMs. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=gjeQKFxFpZ.
- Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. mt5: A massively multilingual pre-trained text-to-text transformer. *arXiv* preprint arXiv:2010.11934, 2020.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging LLM-as-a-judge with MT-bench and chatbot arena. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. URL https://openreview.net/forum?id=uccHPGDlao.
- Tijana Zrnic and Emmanuel J Candès. Active statistical inference. arXiv preprint arXiv:2403.03208, 2024.
- Tijana Zrnic and Emmanuel J. Candès. Cross-prediction-powered inference. *Proceedings of the National Academy of Sciences*, 121(15):e2322083121, 2024. doi: 10.1073/pnas.2322083121. URL https://www.pnas.org/doi/abs/10.1073/pnas.2322083121.

# **CONTENTS**

A	Ethi	cs statement	14
В	Add	itional theoretical results	14
	B.1	Derivation of $Error_T(\pi)$	14
	B.2	Power tuning	14
	B.3	Optimal random annotation: discrete time case	15
	B.4	Extension to convex M-estimators	15
	B.5	Effects of noisy policy parameters on estimator variance	17
	B.6	Optimal active sampling of input evaluation queries	17
	B.7	Informative special cases for $\pi_{\mathrm{active}}$	18
C	C Proofs		19
	<b>C</b> .1	Proof of Proposition 1	19
	C.2	Proof of Proposition 2	20
	C.3	Proof of Proposition 4	23
	C.4	Proof of Proposition 5	24
	C.5	Proof of Proposition 6	24
	C.6	Proof of Proposition 8	25
	<b>C</b> .7	Proof of Corollary 9	26
	<b>C</b> .8	Proof of Proposition 10	26
D	Additional empirical results		27
	D.1	Bernoulli data	27
	D.2	Real data	28
E	Imp	lementation details	29

# A ETHICS STATEMENT

This paper describes fundamental research on the evaluation of generative AI systems, which is a core technical challenge. Hybrid active evaluation has the potential to improve the cost/accuracy tradeoff of system evaluation, which can make high-quality AI systems easier to build, deploy, and monitor. We do not speculate about broader impacts that may follow from this technical contribution. Gemini was used for light copy-editing during the writing of this work.

# B ADDITIONAL THEORETICAL RESULTS

# B.1 DERIVATION OF Error $_T(\pi)$

We provide a short derivation of  $\mathsf{Error}_T(\pi)$  in (2). Because the estimator  $\hat{\theta}_T^{\pi}$  is unbiased,

$$\mathbb{E}\left[\left(\hat{\theta}_T^{\pi} - \theta^*\right)^2\right] = \operatorname{Var}(\hat{\theta}_T^{\pi}) = \frac{1}{T}\operatorname{Var}(\Delta^{\pi})$$

when  $\pi$  and g are fixed, and where  $\Delta^{\pi} = G + (H - G)^2 \frac{\xi}{\pi(X)}$ . Then,

$$\begin{aligned} \operatorname{Var}(\Delta^{\pi}) &= \mathbb{E}\left[\left(G + (H - G)\frac{\xi}{\pi(X)}\right)^{2}\right] - (\theta^{*})^{2} \\ &= \mathbb{E}\left[G^{2}\right] + \mathbb{E}\left[\left((H - G)\frac{\xi}{\pi(X)}\right)^{2}\right] + 2\mathbb{E}\left[G(H - G)\frac{\xi}{\pi(X)}\right] - (\theta^{*})^{2} \\ &= \mathbb{E}\left[G^{2}\right] + \mathbb{E}\left[(H - G)^{2}\frac{1}{\pi(X)}\right] + 2\mathbb{E}\left[G(H - G)\right] - (\theta^{*})^{2} \\ &= \operatorname{Var}(H) - \mathbb{E}[(H - G)^{2}] + \mathbb{E}\left[(H - G)^{2}\frac{1}{\pi(X)}\right]. \end{aligned}$$

# B.2 POWER TUNING

Angelopoulos et al. (2023b) proposed "power tuning" as a way to improve upon the standard PPI estimator by allowing the estimator to adapt to the "usefulness" of the supplementary predictions (here, the weak rater G) with a tuning parameter  $\lambda \in \mathbb{R}$ . We now extend this to our setting.

Let us consider a modified version of our estimator, with some fixed policy  $\pi$  and  $\lambda \in \mathbb{R}$ :

$$\hat{\theta}_T^{\lambda} = \frac{1}{T} \sum_{t=1}^{T} \lambda G_t + (H_t - \lambda G_t) \frac{\xi_t}{\pi(X_t)}.$$

For all values of  $\lambda$ , this estimator is unbiased. Our job is to pick the value with minimum error. Following the previous derivation in Section B.1, the error of the estimator is

$$\mathsf{Error}_{T,\pi}(\lambda) = \frac{1}{T} \left( \mathrm{Var}(H) - \mathbb{E}[(H - \lambda G)^2] + \mathbb{E}\left[ (H - \lambda G)^2 \frac{1}{\pi(X)} \right] \right),$$

which is optimized by

$$\begin{split} \lambda^* &= \operatorname*{argmin}_{\lambda \in \mathbb{R}} \mathbb{E} \left[ (H - \lambda G)^2 \left( \frac{1}{\pi(X)} - 1 \right) \right] \\ &= \operatorname*{argmin}_{\lambda \in \mathbb{R}} \lambda^2 \mathbb{E} \left[ G^2 \left( \frac{1}{\pi(X)} - 1 \right) \right] - 2\lambda \mathbb{E} \left[ HG \left( \frac{1}{\pi(X)} - 1 \right) \right]. \end{split}$$

The above expression is quadratic in  $\lambda$ , and its optimizer is

$$\lambda^* = \frac{\mathbb{E}\left[HG\left(\frac{1}{\pi(X)} - 1\right)\right]}{\mathbb{E}\left[G^2\left(\frac{1}{\pi(X)} - 1\right)\right]},$$

which can be estimated in any consistent way, e.g., by its prediction-powered plug-in that can be computed after sampling all  $(X_t, G_t, H_t, \xi_t)$  as:

$$\hat{\lambda}_T = \frac{\frac{1}{T} \sum_{t=1}^{T} \left( G_t^2 + \left( H_t G_t - G_t^2 \right) \frac{\xi_t}{\pi_t(X_t)} \right) \left( \frac{1}{\pi_t(X_t)} - 1 \right)}{\frac{1}{T} \sum_{t=1}^{T} G_t^2 \left( \frac{1}{\pi_t(X_t)} - 1 \right)}.$$

### B.3 OPTIMAL RANDOM ANNOTATION: DISCRETE TIME CASE

The following proposition is the full version of Proposition 1—with the constraint that  $T^{\text{stop}}$  is an integer. This leads to a substantially more complex optimization problem; we show the solution here, but we do not implement it in practice.

**Proposition 4.** Let  $(X_1, G_1, H_1), \ldots, (X_T, G_T, H_T), T \in \mathbb{N}$ , be an i.i.d. sequence of real-valued random variables with joint distribution P, and define Error, Cost, and  $\Pi^{\text{random}}$  as above. Additionally, define the optimization problem

$$\underset{\pi \in \Pi^{\text{random}}}{\text{minimize}} \quad \text{Error}_{T^{\text{stop}}}(\pi) \\
T^{\text{stop}} \in \mathbb{N}_{+} \quad \text{cost}_{T^{\text{stop}}}(\pi) \leq B.$$
(6)

Then the optimal solution to Problem (6) is either  $\pi^*(x) = 1$  or

$$\pi^*(x) = \frac{B - k^* c_g}{k^* c_h}.$$

for all  $x \in \mathcal{X}$ , where

$$k^* = \operatorname*{argmin}_{k \in \mathcal{K}} \frac{1}{k} \left( \operatorname{Var}(H) - \mathbb{E}[(H - G)^2] \right) + \frac{c_h}{B - kc_g} \mathbb{E}[(H - G)^2],$$

and

$$\mathcal{K} = \left\{ \left[ B \frac{1 + \sqrt{\frac{c_h}{c_g} \frac{\mathbb{E}[(H-G)^2]}{\text{Var}(H) - \mathbb{E}[(H-G)^2]}}}{c_g - c_h \frac{\mathbb{E}[(H-G)^2]}{\text{Var}(H) - \mathbb{E}[(H-G)^2]}} \right], \left[ B \frac{1 + \sqrt{\frac{c_h}{c_g} \frac{\mathbb{E}[(H-G)^2]}{\text{Var}(H) - \mathbb{E}[(H-G)^2]}}}{c_g - c_h \frac{\mathbb{E}[(H-G)^2]}{\text{Var}(H) - \mathbb{E}[(H-G)^2]}} \right] \right\}.$$

It is easy to disambiguate between  $p^* = 1$  and the optimal policy based on  $k^*$  by comparing the objective values directly.

### B.4 EXTENSION TO CONVEX M-ESTIMATORS

Here we give an extension of Proposition 2 to general convex M-estimators (Van der Vaart, 2000). Consider a convex loss function,  $\ell_{\theta}$  for some  $\theta \in \mathbb{R}^d$ , equipped with the simplified notation  $\ell_{\theta,t} = \ell_{\theta}(X_t, H_t)$  for all  $t \in \mathbb{N}$  and  $\ell_{\theta,t}^g = \ell_{\theta}(X_t, G_T)$ . We also use  $\ell_{\theta} = \ell_{\theta}(X, H)$  and  $\ell_{\theta}^g = \ell_{\theta}(X, G)$  for generic points  $(X, G, H) \sim P$ . The target of estimation is the population minimizer,  $\theta^* = \operatorname{argmin}_{\theta \in \mathbb{R}^d} \mathbb{E}[\ell_{\theta}]$ . The active estimator is

$$\hat{\theta}_T = \operatorname*{argmin}_{\theta \in \mathbb{R}^d} \frac{1}{T} \sum_{t=1}^T \Delta_{\theta,t} \quad \text{ where } \Delta_{\theta,t} = \ell_{\theta,t}^g + \left(\ell_{\theta,t} - \ell_{\theta,t}^g\right) \frac{\xi_t}{\pi_t(X_t)},$$

for some sequence of annotation policies  $\pi_t$ ,  $t \in \mathbb{N}$ . For the purpose of deriving optimal annotation policies when  $\pi_t$  is fixed as in Section 2, we will also define

$$\hat{\theta}_T^{\pi} = \operatorname*{argmin}_{\theta \in \mathbb{R}^d} \frac{1}{T} \sum_{t=1}^T \Delta_{\theta,t} \quad \text{where } \Delta_{\theta,t} = \ell_{\theta,t}^g + \left(\ell_{\theta,t} - \ell_{\theta,t}^g\right) \frac{\xi_t}{\pi(X_t)}.$$

Unlike the estimator in the case of mean estimation from Section 2,  $\hat{\theta}_T^\pi$  does not have a closed-form variance in finite samples. The standard solution in the analysis of M-estimators is to appeal to the asymptotic linearity of M-estimators to analyze the variance (Van der Vaart, 2000), as is done in Theorem 1 of Zrnic & Candès (2024). The result below combines the aforementioned theorem with standard parametric analysis to give the asymptotic distribution of the squared error.

**Proposition 5.** Let  $\ell_{\theta}$  be smooth (see Assumption 1 in (Zrnic & Candès, 2024)) and define the Hessian  $W_{\theta^*} = \nabla^2 \mathbb{E}[\ell_{\theta^*,t}]$ . Then if  $\hat{\theta}_T^{\pi} \stackrel{p}{\to} \theta^*$ , we have

$$\sqrt{T}(\hat{\theta}_{T}^{\pi} - \theta^{*}) \stackrel{d}{\to} \mathcal{N}(0, \Sigma^{*}),$$
where  $\Sigma^{*} = W_{\theta^{*}}^{-1} \operatorname{Var} \left( \nabla \ell_{\theta^{*}, t}^{g} + \left( \nabla \ell_{\theta^{*}, t} - \nabla \ell_{\theta^{*}, t}^{g} \right) \frac{\xi_{t}}{\pi(X_{t})} \right) W_{\theta^{*}}^{-1}.$  Therefore, we have
$$T \left\| \hat{\theta}_{T}^{\pi} - \theta^{*} \right\|_{2}^{2} \stackrel{d}{\to} \sum_{j \in [d]} \lambda_{j} \zeta_{j},$$

where  $\zeta_j \overset{\text{i.i.d.}}{\sim} \chi_1^2$  for all  $j \in [d]$  and  $\lambda_j$  is the jth eigenvalue of  $\Sigma^*$ .

The above proposition gives us consistency of the active estimator, and more importantly, the asymptotic distribution of the squared error. Since  $\mathbb{E}[\zeta_j]=1$  for all j, and the sum of the eigenvalues of a square matrix is equal to the trace, we know the mean-squared error is asymptotically equal to  $\mathsf{Error}_T(\pi) = \frac{1}{T}\operatorname{Tr}(\Sigma^*)$ . With this in hand, we can use the same strategy from earlier to find the optimal annotation policy, using the asymptotic approximation of the error. For simplicity, here we assume that we are always on the interior of the constrained optimization problem, i.e., we solve for unconstrained  $\pi(x)$  while assuming that  $\gamma^*\sqrt{u(X)} \leq 1$ . That said, a more rigorous treatment analogous to that in Proposition 2 can also be applied here, which we leave to future work.

**Proposition 6.** In the setting of Proposition 5, let  $(X_1, G_1, H_1), \ldots, (X_T, G_T, H_T), T \in \mathbb{N}$ , be an i.i.d. sequence of real-valued random variables with joint distribution P, and define  $\mathsf{Error}_T(\pi) = \frac{1}{T} \sum_{j \in [d]} \mathrm{Tr}(\Sigma^*)$ . Furthermore, define  $\mathsf{Cost}$  and  $\Pi$  as in Proposition 2.

Construct the optimization problem

$$\underset{\pi \in \mathcal{F}, \ T^{\text{stop}} \in \mathbb{R}_{>0}}{\text{minimize}} \quad \text{Error}_{T^{\text{stop}}}(\pi) 
\text{subject to} \quad \text{Cost}_{T^{\text{stop}}}(\pi) < B.$$
(7)

where  $\mathcal{F} = \{x \mapsto f(x) : f(x) \in (0, \infty); \forall x \in \mathcal{X}\}$ . Then the solution to Problem (7) is

$$\pi^*(x) = \sqrt{\frac{c_g}{c_h} \cdot \frac{u(x)}{C}}$$

where

$$u(x) = \mathbb{E}\left[\operatorname{Tr}\left(W_{\theta^*}^{-1}\left(\nabla \ell_{\theta^*} - \nabla \ell_{\theta^*}^g\right)\left(\nabla \ell_{\theta^*} - \nabla \ell_{\theta^*}^g\right)^{\top}W_{\theta^*}^{-1}\right) \mid X = x\right],$$

and

$$C = \operatorname{Tr}\left(W_{\theta^*}^{-1} \left( \mathbb{E}\left[\nabla \ell_{\theta^*}^g (\nabla \ell_{\theta^*})^\top + (\nabla \ell_{\theta^*} - \nabla \ell_{\theta^*}^g)(\nabla \ell_{\theta^*}^g)^\top \right] - \mathbb{E}[\nabla \ell_{\theta^*}] \mathbb{E}[\nabla \ell_{\theta^*}]^\top \right) W_{\theta^*}^{-1} \right).$$

**Remark 7.** When  $\pi^*(x) \leq 1$ ,  $\forall x \in \mathcal{X}$ , then  $\pi^*$  is also optimal for Problem (7) solved for  $\pi \in \Pi$ .

#### **B.4.1** MEAN ESTIMATION

In the case of mean estimation, the loss function takes the form

$$\ell_{\theta}(x,h) = \frac{1}{2}(h-\theta)^2,$$

where  $\nabla \ell_{\theta^*}(X, H) - \nabla \ell_{\theta^*}(X, G) = H - G$ , and  $W_{\theta^*}$  is the identity matrix. Plugging back into  $\pi^*$  in Proposition 6 recovers  $\pi_{\text{active}}$  from Proposition 2 without clipping  $(\tau^* = \infty)$ , i.e.,

$$\sqrt{\frac{c_g}{c_h}} \frac{\mathbb{E}[(H-G)^2 \mid X=x]}{\operatorname{Var}(H) - \mathbb{E}[(H-G)^2]}.$$

#### B.4.2 GENERALIZED LINEAR MODELS

In the case of GLMs, the loss function takes the form

$$\ell_{\theta}(x,h) = -hx^{\top}\theta + \psi(x^{\top}\theta)$$

for some convex log-partition function  $\psi$ . Thus,  $\nabla \ell_{\theta^*}(X, H) - \nabla \ell_{\theta^*}(X, G) = (G - H)X$ . So, again by the linearity of the trace, we have that

$$\pi^*(x) \propto \sqrt{\mathbb{E}\left[(H-G)^2 \mid X=x\right]} \operatorname{Tr}\left(W_{\theta^*}^{-1} x x^{\top} W_{\theta^*}^{-1}\right).$$

# B.5 EFFECTS OF NOISY POLICY PARAMETERS ON ESTIMATOR VARIANCE

In practice, we will be using only an imperfect estimate of u(x) for  $\pi_{\rm active}$ , which can negatively affect the performance of  $\pi_{\rm active}$  to a substantial degree, as we have seen for some of the datasets in Section 4. Similarly, we will also only be using imperfect estimates of the optimal scaling and thresholding parameters used in  $\pi_{\rm active}$ , which further limit performance.

There are two main factors that affect the error of a policy:

- 1. The variance,  $Var(\Delta^{\pi})$ , of each active increment  $\Delta^{\pi}$ , where  $\Delta^{\pi} = G + (H G)\frac{\xi}{\pi(X)}$ .
- 2. The average sample size at which the estimator runs out of budget,  $T^{\text{stop}} = \frac{B}{c_h \mathbb{E}[\pi(X)] + c_q}$ .

In this section, we provide some additional theoretical analysis on the first factor, i.e., the increase in  $Var(\Delta^{\pi})$  due to the mispecification error of an *estimated* active policy, while noting that the total error will be further affected by the relative increase/decrease of the mean sampling rate,  $\mathbb{E}[\pi(X)]$ .

**Proposition 8.** In the same setting as Proposition 2, let  $\tilde{\pi}: \mathcal{X} \to (0,1]$  be any function satisfying

$$\mathbb{E}\left[\frac{1}{\tilde{\pi}(X)} - \frac{1}{\pi^*(x)}\right] \le \delta,$$

where  $\pi^*$  is the oracle estimate of  $\pi_{\text{active}}$ . Let  $(H-G)^2 \stackrel{\text{a.s.}}{\leq} b$ . Then  $\operatorname{Var}(\Delta^{\tilde{\pi}}) \leq \operatorname{Var}(\Delta^{\pi^*}) + b\delta$ .

If we simply things by assuming an additive error model for a policy without thresholding (i.e.,  $\tau^* = \infty$ ), we can refine the bound somewhat further:

**Corollary 9.** Let  $\tilde{\pi} = \tilde{\gamma}\sqrt{\tilde{u}(x)}$ , where  $\tilde{\gamma} = \gamma^* + \delta_{\gamma}$ ,  $\tilde{u}(x) = u(x) + \delta_u(x)$ , and  $u(X) \overset{\text{a.s.}}{\geq} \epsilon$ . Further assume that  $\tilde{\pi}$  is admissible, i.e.,  $\tilde{\pi}(x) \in (0,1] \ \forall x$ . Then, up to first-order terms in  $\delta_{\gamma}$  and  $\delta_u(x)$ ,

$$\operatorname{Var}(\Delta^{\tilde{\pi}}) \leq \operatorname{Var}(\Delta^{\pi^*}) + b \left( \frac{|\delta_{\gamma}|}{(\gamma^*)^2 \sqrt{\epsilon}} + \frac{1}{2\gamma^* \epsilon^{3/2}} \mathbb{E}[|\delta_U(X)|] \right).$$

We can make a few observations about the results in Proposition 8 and Corollary 9. First, as long as the error,  $(H-G)^2$  is bounded, and the estimated inverse propensity score  $1/\tilde{\pi}(X)$  is not significantly higher than the oracle inverse propensity score  $1/\pi^*(X)$  on average, then the increase in variance over the oracle will not be that large. Generally speaking, this is satisfied when the estimated policy is not *overconfident* on examples that in fact have high error. Of course, regularizing the estimated policy to be underconfident on all examples is also not always a satisfying solution: as  $\mathbb{E}[\tilde{\pi}(X)] \to 1$ , we obtain a policy that is no better than  $\pi_{\text{base}}$ . Similarly, as seen in Section 3, the headroom for  $\pi_{\text{active}}$  over  $\pi_{\text{base}}$  is largest when  $(H-G)^2$  is *not* bounded (e.g., the Gaussian data setting compared to the Bernoulli data setting), as large  $(H-G)^2$  also increase the possible variance of U. This reinforces the importance of having **accurate uncertainty estimates** when computing active policies.

#### B.6 OPTIMAL ACTIVE SAMPLING OF INPUT EVALUATION QUERIES

This section shows how to optimally choose the distribution of X. In contrast, Section 2 in the main paper focuses only on querying annotators for H given i.i.d. samples from the fixed distribution P for X. Deciding which inputs to sample can be a more difficult problem than deciding whether to annotate a given input sample because  $\mathcal X$  can be large and complex. However, we can always apply the optimal rules to a coarse stratification of  $\mathcal X$ . Towards this end, we define the estimator

$$\hat{\theta}_{T}^{Q,\pi} = \frac{1}{T} \sum_{t=1}^{T} \Delta_{t}, \text{ where } \quad \Delta_{t}^{Q,\pi} = \frac{dP}{dQ}(X_{t}) \left( G_{t} + (H_{t} - G_{t}) \frac{\xi_{t}}{\pi(X_{t})} \right),$$

$$X_{t} \overset{\text{i.i.d.}}{\sim} Q, \ H_{t} \sim P_{H|X}, \ G_{t} \sim P_{G|X}$$

which is our previous estimator with a fixed policy  $\pi$ , and where the X are sampled from a distribution Q, and the distribution of  $H \mid X$  and  $G \mid X$  remain unchanged. This estimator is unbiased for  $\theta^*$ ,

and a straightforward calculation gives that the error of the estimator is

$$\begin{split} \mathsf{Error}_T(Q,\pi) &= \mathbb{E}_Q \left[ \left( \hat{\theta}_T^{Q,\pi} - \theta^* \right)^2 \right] \\ &= \frac{1}{T} \mathrm{Var}(\Delta^{Q,\pi}) \\ &= \frac{1}{T} \left( \mathbb{E}_P \left[ \frac{dP}{dQ}(X) \left( H^2 + \left( \frac{1}{\pi(X)} - 1 \right) (H - G)^2 \right) \right] - (\theta^*)^2 \right). \end{split}$$

The goal is to pick a distribution Q to minimize the error of the estimator. The following proposition gives an explicit form for this optimal sampling distribution.

**Proposition 10.** Define  $\mathsf{Error}_T$  as above, and define the set of all strictly positive densities,  $\mathcal{Q} = \{x \mapsto Q(x) : Q(x) \in \mathbb{R}_{>0} \text{ and } Q \in \Delta^{\mathcal{X}} \}$ . Furthermore, define the optimization problem

$$\underset{Q \in \mathcal{Q}}{\text{minimize}} \quad \mathsf{Error}_T(Q, \pi) \tag{8}$$

for a fixed time  $T \in \mathbb{N}$ . Then the solution to Problem (3) is

$$Q^*(x) = \mathbb{P}(X = x) \frac{\sqrt{\nu(x)}}{\mathbb{E}_P[\sqrt{\nu(X)}]},$$

where

$$\nu(x) = \mathbb{E}_P\left[\left(H^2 + \left(\frac{1}{\pi(X)} - 1\right)(H - G)^2\right) \mid X = x\right]$$

for all  $x \in \mathcal{X}$ .

We leave empirical exploration of active input sampling to future work.

### B.7 Informative special cases for $\pi_{\text{active}}$

Prior work (Zrnic & Candès, 2024; Gligorić et al., 2024) target some fixed, prespecified value (i.e., some ratio n/N) for  $\mathbb{E}[\pi(X)]$ . A key distinction of this work is that we also optimize  $\mathbb{E}[\pi(X)]$ , which will depend strongly on  $c_g/c_h$ , that is, the cost ratio of G to H. In this section we analyze two extreme, but informative cases, for active sampling when either  $c_g/c_h=0$  or  $c_g/c_h=\infty$ , that serve to illustrate how  $\mathbb{E}[\pi_{\rm active}(X)]$  for the cost-optimal policy  $\pi_{\rm active}$  can consequently be as extreme as 0 or 1.

Optimal policy for  $c_q = 0$ 

We start with the special case where  $c_g=0$ , so that we can obtain essentially infinitely many queries of the weak rater G irrespective of the budget constraint. In this case, we expect that unless G has a prohibitively large error  $\mathbb{E}[(H-G)^2]$ , we can purely rely on querying G, and overcome its error with sufficiently many samples. Indeed, let us assume that  $\mathbb{E}[(H-G)^2]=\mathbb{E}[U]<\mathrm{Var}(H)$ . Then we note that for any  $\tau>0$ :

$$\tau \sqrt{\frac{c_g/c_h + \mathbb{P}(U > \tau^2)}{\operatorname{Var}(H) - \mathbb{E}[U\mathbb{1}\{U \le \tau^2\}]}} = \sqrt{\frac{\tau^2 \mathbb{P}(U > \tau^2)}{\operatorname{Var}(H) - \mathbb{E}[U\mathbb{1}\{U \le \tau^2\}]}}$$

$$\leq \sqrt{\frac{\tau^2 \mathbb{P}(U > \tau^2)}{\mathbb{E}[U] - \mathbb{E}[U\mathbb{1}\{U \le \tau^2\}]}}$$

$$= \sqrt{\frac{\tau^2 \mathbb{P}(U > \tau^2)}{\mathbb{E}[U\mathbb{1}\{U > \tau^2\}]}} \leq 1,$$

where the first inequality is due to our assumption that  $\mathbb{E}[U] < \text{Var}(H)$ , and the last inequality follows from  $\mathbb{E}[U\mathbbm{1}\{U > \tau^2\}] > \tau^2\mathbb{P}(U > \tau^2)$ . Consequently, we get that in this case,

$$\gamma^*(\tau) = \sqrt{\frac{\mathbb{P}(U > \tau^2)}{\operatorname{Var}(H) - \mathbb{E}[U\mathbb{1}\{U \le \tau^2\}]}}.$$

Suppose for now that we only consider the values  $\tau$  where U further satisfies that  $\sqrt{U}\gamma^*(\tau) \leq 1$  almost surely, and denote this set by  $\mathcal{T}$ . Let  $\Delta = \operatorname{Var}(H) - \mathbb{E}[U] > 0$ . Then we see that

$$\min_{\tau \in \mathcal{T}} c_h \mathbb{E}[\pi_{\text{clip}}(x;\tau)] \left[ \Delta + \mathbb{E} \left[ \frac{U}{\pi_{\text{clip}}(x;\tau)} \right] \right] 
= \min_{\tau \in \mathcal{T}} c_h \mathbb{E}[\sqrt{U}\gamma^*(\tau)] \left[ \Delta + \mathbb{E} \left[ \frac{U}{\sqrt{U}\gamma^*(\tau)} \right] \right] 
= \min_{\tau \in \mathcal{T}} c_h \mathbb{E}[\sqrt{U}\gamma^*(\tau)] \Delta + c_h \mathbb{E}[\sqrt{U}\gamma^*(\tau)] \mathbb{E} \left[ \frac{\sqrt{U}}{\gamma^*(\tau)} \right].$$

Since  $\gamma^*(\tau)$  is deterministic, it cancels from the second term above, and we get that the annotation cost over  $\tau \in \mathcal{T}$  is monotonically increasing in  $\gamma^*(\tau)$ , meaning that we choose  $\tau \to \infty$ , which yields  $\gamma^*(\tau) \to 0$  (since  $P(U > \infty) = 0$ ). We also note that whenever  $\sqrt{U} \leq B$ , all  $\tau$  such that  $\gamma^*(\tau) < 1/B$  are in  $\mathcal{T}$  trivially, since this satisfies  $\sqrt{U}\gamma^*(\tau) < 1$ . In particular, this includes our choice of  $\tau \to \infty$ , which ensures that  $\gamma^*(\tau) \to 0$ . Finally, we note that from the proof of Proposition 2 (specifically, Equation 11), we have that  $\pi(x) = \gamma \sqrt{u(x)}$  minimizes the objective

$$(c_h E[\pi(X)] + c_g) \left[ \operatorname{Var}(H) - \mathbb{E}[U] + \mathbb{E}\left[\frac{U}{\pi(X)}\right] \right],$$

over all mappings  $\pi \in \{x \mapsto f(x) : f(x) \in (0, \infty); \forall x \in \mathcal{X}\}$ . Since we find that our optimal choice without imposing the constraint  $\pi(x; \tau) \leq 1$  is already feasible, it is also optimal for the constrained problem,  $\pi \in \{x \mapsto f(x) : f(x) \in (0, 1]; \forall x \in \mathcal{X}\}$ .

Optimal policy for  $c_h = 0$ 

The other extreme case is simpler. When  $c_h=0$ , the objective for  $\tau^*$  becomes monotonically decreasing in  $\pi(x;\tau)$ . If we assume that  $\tau$  is such that  $\gamma^*(\tau)<1/\tau$ , then we find that the expression

$$\sqrt{\frac{c_g/c_h + \mathbb{P}(U > \tau^2)}{\left(\operatorname{Var}(H) - \mathbb{E}[U\mathbb{1}\{U \le \tau^2\}]\right)_+}}$$

becomes infinite due to  $c_h=0$ , and hence we must have  $\gamma^*(\tau)=1/\tau$ . However, for any x such that  $\pi(x;\tau)<1$ , we have  $1/\pi(x;\tau)=\tau/\sqrt{u(x)}$ . Consequently, minimizing over  $\tau$  results in  $\tau=0$ . But this gives  $\pi(x;\tau)=\infty$ , so that we must have  $\pi(x;\tau)=1$  for all x. Intuitively, this makes sense because any  $\pi(x;\tau)<1$  results in an estimator with variance strictly greater than  $\mathrm{Var}(H)$ , but having  $\pi(x;\tau)\equiv 1$  allows us to attain the smallest possible variance of  $\mathrm{Var}(H)$ . Since there is no effect of these choices on the estimation cost, we choose the lowest variance estimator in this case, and direct all our queries to the strong rater.

#### C Proofs

#### C.1 PROOF OF PROPOSITION 1

*Proof.* Since  $\mathsf{Error}_T(\pi)$  is monotone in T for all  $\pi$ , we should first set  $T^{\mathsf{stop}}$  to be the largest T for which the constraint holds. This value is

$$T^{\text{stop}} = \frac{B}{c_h p + c_a}.$$

Plugging this into the objective yields

$$(c_h p + c_g) \left( \operatorname{Var}(H) - \mathbb{E}[(H - G)^2] + \frac{1}{p} \mathbb{E}[(H - G)^2] \right),$$

which, after removing terms that do not depend on p, is equivalent to minimizing

$$c_h p \left( \operatorname{Var}(H) - \mathbb{E}[(H - G)^2] \right) + \frac{c_g}{p} \mathbb{E}\left[ (H - G)^2 \right]$$

subject to the constraint that  $p \in [0, 1]$ .

This is a convex problem in p, and we know that the solution lies either on the boundary or on the interior. We will compare the values of the objectives in three cases:  $p^* = 0$ ,  $p^* = 1$ , and  $p^* \in (0,1)$ . It is clear that  $p^* = 0$  is infeasible (unless  $H \stackrel{\text{a.s.}}{=} G$ , which renders the problem trivial) because the

factor  $c_q/p$  appears in the above objective. In the case that  $p^* = 1$ , the objective value is

$$c_h \left( \operatorname{Var}(H) - \mathbb{E}[(H - G)^2] \right) + c_g \mathbb{E} \left[ (H - G)^2 \right].$$

In the case that  $p^* \in (0,1)$ , it must be a critical value, so it satisfies the first-order condition

$$c_h \left( \text{Var}(H) - \mathbb{E}[(H - G)^2] \right) = \frac{c_g}{(p^*)^2} \mathbb{E}[(H - G)^2],$$

and thus,

$$(p^*)^2 = \frac{c_g \mathbb{E}[(H-G)^2]}{c_h(\text{Var}(H) - \mathbb{E}[(H-G)^2])}.$$
 (9)

However, because we are in the case  $p^* \in (0,1)$ , the right-hand side above must be positive (otherwise the square root would be imaginary), and it cannot be greater than 1 (otherwise we would have  $p^* > 1$ , which is a contradiction). This gives us that

$$p^* \in (0,1) \implies \mathbb{E}[(H-G)^2] < \operatorname{Var}(H) \text{ and } (c_g + c_h)\mathbb{E}[(H-G)^2] < c_h \operatorname{Var}(H).$$

Under these conditions, we can take square roots on both sides of (9) to obtain

$$p^* = \sqrt{\frac{c_g}{c_h} \frac{1}{\frac{\operatorname{Var}(H)}{\mathbb{E}[(H-G)^2]} - 1}}.$$

The objective value at this point is

$$2\sqrt{c_g c_h}\sqrt{\mathbb{E}[(H-G)^2](\operatorname{Var}(H)-\mathbb{E}[(H-G)^2])}$$

Finally, comparing the above objective value with that of  $p^* = 1$ , we have that

$$2\sqrt{c_g c_h} \sqrt{\mathbb{E}[(H-G)^2](\operatorname{Var}(H) - \mathbb{E}[(H-G)^2])} < c_h \left(\operatorname{Var}(H) - \mathbb{E}[(H-G)^2]\right) + c_g \mathbb{E}\left[(H-G)^2\right]$$

$$\iff 0 < c_h^2 \left( \text{Var}(H) - \mathbb{E}[(H - G)^2] \right)^2 - 2c_g c_h \mathbb{E}[(H - G)^2] \left( \text{Var}(H) - \mathbb{E}[(H - G)^2] \right) + c_g^2 \mathbb{E}\left[(H - G)^2\right]^2$$

$$\iff 0 < \left(c_h(\operatorname{Var}(H) - \mathbb{E}[(H - G)^2]) - c_g \mathbb{E}[(H - G)^2]\right)^2.$$

Under the condition that  $(c_g + c_h)\mathbb{E}[(H - G)^2] < c_h \text{Var}(H)$ , the above inequality cannot hold, since the squared term on the right-hand side will always be positive (and nonzero). Thus, we have that

$$p^* = \begin{cases} \sqrt{\frac{c_g}{c_h} \frac{1}{\frac{\operatorname{Var}(H)}{\mathbb{E}[(H-G)^2]} - 1}} & \text{if } (c_g + c_h) \mathbb{E}[(H-G)^2] < c_h \operatorname{Var}(H) \text{ and } \mathbb{E}[(H-G)^2] < \operatorname{Var}(H) \\ 1 & \text{otherwise.} \end{cases}$$

Under the constraint that  $c_h \geq c_q$ , this simplifies to

$$p^* = \begin{cases} \sqrt{\frac{c_g}{c_h} \frac{\mathbb{E}[(H-G)^2]}{\text{Var}(H) - \mathbb{E}[(H-G)^2]}} & \text{if } \mathbb{E}[(H-G)^2] < \frac{c_h}{c_g + c_h} \text{Var}(H) \\ 1 & \text{otherwise.} \end{cases}$$

## C.2 PROOF OF PROPOSITION 2

*Proof.* Following the simplification of Problem (3) in the proof of Proposition (1), Problem (5) is also equivalent to minimizing the following objective:

$$J(\pi) = c_h \mathbb{E}[\pi(X)] \left( \operatorname{Var}(H) - \mathbb{E}[(H - G)^2] + \mathbb{E}\left[ (H - G)^2 \frac{1}{\pi(X)} \right] \right) + c_g \mathbb{E}\left[ (H - G)^2 \frac{1}{\pi(X)} \right].$$

At this point, we leverage the discreteness of  $\mathcal X$  to write the objective in a simpler form. Let  $P \in \Delta^{\mathcal X}$  be the probability mass function of X, expressed as a vector, and let  $I \in \{0,1\}^{|\mathcal X|}$  be the indicator that X takes each value in  $\mathcal X$ . Furthermore, let  $p \in [0,1]^{|\mathcal X|}$  be the vector of  $\pi(x)$  for all  $x \in \mathcal X$ . Then, we can express  $\pi(X) = p^\top I$  and  $\mathbb E[\pi(X)] = p^\top P$ , and write the objective as

$$J(\pi) = J(p) = p^{\top} P\left(\operatorname{Var}(H) - \mathbb{E}[(H - G)^2] + \mathbb{E}\left[(H - G)^2 \frac{1}{p^{\top} I}\right]\right) + \frac{c_g}{c_h} \mathbb{E}\left[(H - G)^2 \frac{1}{p^{\top} I}\right]. \tag{10}$$

From here on out, we assume that  $P_x>0$  for all  $x\in\mathcal{X}$ . The final result will hold without loss of generality, since the value of the optimal policy on measure-zero points does not change the value of the objective. For any x, we clearly cannot have  $p_x=0$ , otherwise the objective would be infinite. This rules out  $p_x=0$  for almost all x. We are left with the constraint that  $p\preceq 1$ .

Forming the Lagrangian,

$$\mathcal{L}(p,\lambda) = J(p) + \lambda^{\top}(p-1)$$

$$= p^{\top} P\left(\operatorname{Var}(H) - \mathbb{E}[(H-G)^2] + \mathbb{E}\left[(H-G)^2 \frac{1}{p^{\top}I}\right]\right) + \frac{c_g}{c_h} \mathbb{E}\left[(H-G)^2 \frac{1}{p^{\top}I}\right] + \lambda^{\top}(p-1).$$

Taking the gradient with respect to p gives  $\nabla_p \mathcal{L}(p,\lambda)$  equal to

$$P\left(\operatorname{Var}(H) - \mathbb{E}[(H - G)^2]\right) - \left(p^{\top}P + \frac{c_g}{c_h}\right)\mathbb{E}\left[(H - G)^2 \frac{I}{(p^{\top}I)^2}\right] + P\mathbb{E}\left[(H - G)^2 \frac{1}{p^{\top}I}\right] + \lambda.$$

Setting the gradient to zero coordinate-wise then gives that for each x,

$$P_x\left(\operatorname{Var}(H) - \mathbb{E}[(H-G)^2]\mathbb{E}\left[(H-G)^2\frac{1}{p^\top I}\right]\right) = \left(p^\top P + \frac{c_g}{c_h}\right)\mathbb{E}\left[(H-G)^2\frac{I_x}{p_x^2}\right] - \lambda_x.$$

By the definition of the conditional expectation, and rearranging, we can rewrite this as

$$\operatorname{Var}(H) - \mathbb{E}[(H - G)^2] + \mathbb{E}\left[(H - G)^2 \frac{1}{p^\top I}\right] + \frac{\lambda_x}{P_x} = \left(p^\top P + \frac{c_g}{c_h}\right) \mathbb{E}\left[(H - G)^2 \frac{1}{p_x^2} \mid X = x\right].$$

Solving for the optimal value as a function of the Lagrange multipliers  $\lambda$  gives the following expression:

$$p_x(\lambda)^2 = \frac{\left(p^\top P + \frac{c_g}{c_h}\right) \mathbb{E}\left[(H - G)^2 \mid X = x\right]}{\operatorname{Var}(H) - \mathbb{E}[(H - G)^2] + \mathbb{E}\left[(H - G)^2 \frac{1}{p^\top I}\right] + \frac{\lambda_x}{P_x}}.$$

The denominator of this expression is always positive, since for all valid p,  $\frac{(H-G)^2}{p^+I} \stackrel{\text{a.s.}}{\geq} (H-G)^2$ , and the remaining terms are positive. Thus,

$$p_x(\lambda) = \sqrt{\frac{\left(p^\top P + \frac{c_g}{c_h}\right) \mathbb{E}\left[(H - G)^2 \mid X = x\right]}{\operatorname{Var}(H) - \mathbb{E}\left[(H - G)^2\right] + \mathbb{E}\left[(H - G)^2 \frac{1}{p^\top I}\right] + \frac{\lambda_x}{P_x}}}.$$
(11)

Next, we require some detailed case-by-case analysis.

Case 1: The Interior. First, we handle the case when the constraint is inactive, i.e., for any fixed  $\lambda$ ,  $p_x(\lambda) \in (0,1)$ . (If no such x exists, then the solution is trivially  $p(\lambda) = \mathbf{1}_{|\mathcal{X}|}$ .) For any x such that  $p_x(\lambda)$  is in the interior, by complementary slackness,  $\lambda_x = 0$ . Now, for any  $x' \in \mathcal{X}$  satisfying  $p_{x'}(\lambda) \in (0,1)$ , we can write

$$\frac{p_x(\lambda)}{p_{x'}(\lambda)} = \sqrt{\frac{\mathbb{E}\left[ (H - G)^2 \mid X = x \right]}{\mathbb{E}\left[ (H - G)^2 \mid X = x' \right]}},$$

simply by applying (11) to x and x', then dividing these expressions. This tells us that for all  $\lambda$  and all x such that  $p_x(\lambda)$  is in the interior,  $p_x(\lambda) = \gamma u(x)$  for some as-yet-unknown  $\gamma \in$ 

 $\left(0, \frac{1}{\sup_{x:p_x\in(0,1)}u(x)}\right]$ . Because  $\lambda_x=0$  on these x, the solution to the optimization problem must have the same property.

Case 2: The Boundary. When the constraint is active,  $p_x(\lambda) = 1$ , since  $p_x(\lambda) = 0$  is almost always impossible, as established earlier. Examining (11) shows us that the constraint only activates in the case that  $u(x) = \mathbb{E}[(H - G)^2 \mid X = x]$  is too large:

$$p_x = 1 \Longleftrightarrow u(x) \ge \sqrt{\frac{\operatorname{Var}(H) - \mathbb{E}[(H - G)^2] + \mathbb{E}\left[(H - G)^2 \frac{1}{p^\top I}\right]}{p^\top P + \frac{c_g}{c_h}}} = \tau(p),$$

since in the alternate case, the unconstrained solution lies in the interior. The Lagrange multiplier  $\lambda_x$ , in this case, takes on the value such that  $p_x(\lambda_x) = 1$ ; a non-negative such value always exists by virtue of the fact that u(x) is sufficiently large.

Combining Case 1 and Case 2 tells us that our optimal policy has the form

$$\pi(x) = \begin{cases} \gamma \sqrt{u(x)} & \sqrt{u(x)} \leq \tau \\ 1 & \text{otherwise} \end{cases},$$

for a  $\tau \in \mathbb{R}_{>0}$  and  $\gamma \in \left(0, \inf_{x \in u(x) \le \tau^2} \sqrt{u(x)}\right]$ , which we assume w.l.o.g. is equivalent to  $\gamma \in \left(0, \frac{1}{\tau}\right)$ . The constraint on  $\gamma$  is necessary, as otherwise we can have  $\pi(x) > 1$ , which is a contradiction.

Note that another way to express this policy is as  $p_x = \mathbb{1}\left\{u(x) > \tau^2\right\} + \gamma \sqrt{u(x)}\mathbb{1}\left\{u(x) \leq \tau^2\right\}$ . With this in mind, and defining the vector U with entries  $U_x = \mathbb{E}[(H-G)^2 \mid X = x]$  and  $W = \operatorname{Var}(H) - \mathbb{E}[(H-G)^2]$  we can rewrite the objective in (10) as

$$J(p) = \left(\sum_{x \in \mathcal{X}} p_x P_x\right) \left(W + \mathbb{E}\left[(H - G)^2 \frac{1}{p_X}\right]\right) + \frac{c_g}{c_h} \mathbb{E}\left[(H - G)^2 \frac{1}{p_X}\right]$$

which is equivalent to

$$\begin{split} J(\gamma,\tau) &= \mathbb{E}\left[\mathbbm{1}\left\{U_X > \tau^2\right\} + \gamma\sqrt{U_X}\mathbbm{1}\left\{U_X \leq \tau^2\right\}\right] \\ &\times \left(W + \mathbb{E}\left[(H-G)^2\mathbbm{1}\left\{U_X > \tau^2\right\}\right] + \mathbb{E}\left[\frac{(H-G)^2}{\gamma\sqrt{U_X}}\mathbbm{1}\left\{U_X \leq \tau^2\right\}\right]\right) \\ &+ \frac{c_g}{c_h}\left(\mathbb{E}\left[(H-G)^2\mathbbm{1}\left\{U_X > \tau^2\right\}\right] + \mathbb{E}\left[\frac{(H-G)^2}{\gamma\sqrt{U_X}}\mathbbm{1}\left\{U_X \leq \tau^2\right\}\right]\right) \end{split}$$

This objective is convex in  $\gamma$ , but not differentiable or convex in  $\tau$ . For that reason, we will solve for the optimal  $\gamma$  as a function of  $\tau$  subject to the constraint that  $\gamma>0$  and  $\gamma\sqrt{u(x)}\leq 1\ \forall x$  where  $u(x)\leq \tau^2$ , and our algorithm will search over  $\tau$  to complete the optimization. Keeping only terms with a dependence on  $\gamma$ , and recognizing that  $\mathbb{E}\left[\frac{(H-G)^2}{\sqrt{U_X}}\right]=\mathbb{E}\left[\sqrt{U_X}\right]$  gives the expression

$$\mathbb{E}\left[\sqrt{U_X}\mathbb{1}\left\{U_X \le \tau^2\right\}\right] \times \left[\frac{1}{\gamma}\left(\mathbb{E}\left[\mathbb{1}\left\{U_X > \tau^2\right\}\right] + \frac{c_g}{c_h}\right) + \gamma\left(W + \mathbb{E}\left[(H - G)^2\mathbb{1}\left\{U_X > \tau^2\right\}\right]\right)\right]$$
(12)

Once again, we know that the optimal solution as a function of  $\tau$ ,  $\gamma^*(\tau)$  lies either on the boundary or the interior, and we will compare the values of the objective in both cases. In the case that  $\gamma^*(\tau) \in (0, \tau^{-1})$ ,  $\gamma^*(\tau)$  is a critical point, thus differentiating and setting equal to zero gives that

$$\frac{1}{\gamma^*(\tau)^2} \left( \frac{c_g}{c_h} + \mathbb{P}\left( U_X > \tau^2 \right) \right) = W + \mathbb{E}\left[ (H - G)^2 \mathbb{1}\left\{ U_X > \tau^2 \right\} \right],$$

and thus.

$$\gamma^*(\tau)^2 = \frac{\frac{c_g}{c_h} + \mathbb{P}\left(U_X > \tau^2\right)}{W + \mathbb{E}\left[(H - G)^2 \mathbb{1}\left\{U_X > \tau^2\right\}\right]} = \frac{\frac{c_g}{c_h} + \mathbb{P}\left(U_X > \tau^2\right)}{\text{Var}(H) - \mathbb{E}\left[(H - G)^2 \mathbb{1}\left\{U_X \le \tau^2\right\}\right]}.$$
 (13)

As in the proof of Proposition 1, because we are in the case  $\gamma^* \in (0, \tau^{-1})$ , the right-hand side must be positive and it cannot be greater than  $\tau^{-1}$ . This gives us that

$$\gamma^*(\tau) \in (0, \tau^{-1}) \Longrightarrow \mathbb{E}\left[ (H - G)^2 \mathbb{1} \left\{ U_X \le \tau^2 \right\} \right] < \text{Var}(H)$$

and

$$\frac{\frac{c_g}{c_h} + \mathbb{P}\left(U_X > \tau^2\right)}{\operatorname{Var}(H) - \mathbb{E}\left[(H - G)^2 \mathbb{1}\left\{U_X \leq \tau^2\right\}\right]} < \frac{1}{\tau^2},$$

and under these conditions we can take square roots on both sides of (13) to obtain

$$\gamma^{*}(\tau) = \sqrt{\frac{\frac{c_{g}}{c_{h}} + \mathbb{P}(U_{X} > \tau^{2})}{\text{Var}(H) - \mathbb{E}[(H - G)^{2}\mathbb{1}\{U_{X} \le \tau^{2}\}]}} = \sqrt{\frac{\frac{c_{g}}{c_{h}} + \mathbb{P}(U_{X} > \tau^{2})}{\text{Var}(H) - \mathbb{E}[U_{X}\mathbb{1}\{U_{X} \le \tau^{2}\}]}}.$$
 (14)

Comparing the objective value with  $\gamma^*(\tau) = \tau^{-1}$  vs (14), we know that (12) is decreasing in  $\gamma$  for  $0 < \gamma < \sqrt{\frac{\frac{cg}{c_h} + \mathbb{P}(U_X > \tau^2)}{\mathrm{Var}(H) - \mathbb{E}[U_X \mathbb{1}\{U_X \leq \tau^2\}]}}$ . Thus, we have that

$$\gamma^*(\tau) = \min\left(\sqrt{\frac{c_g/c_h + \mathbb{P}(U_X > \tau^2)}{\left(\operatorname{Var}(H) - \mathbb{E}[U_X \mathbb{1}\{U_X \le \tau^2\}]\right)_+}}, \frac{1}{\tau}\right).$$

Plugging into the original objective  $J(\tau, \gamma^*(\tau))$  and minimizing over  $\tau$  yields the solution.

# C.3 PROOF OF PROPOSITION 4

*Proof.* Since  $\mathsf{Error}_T(\pi)$  is monotone in T for all  $\pi$ , we should first set  $T^{\mathsf{stop}}$  to be the largest T for which the constraint holds. This value is

$$T^{\text{stop}} = \left\lfloor \frac{B}{c_h p + c_g} \right\rfloor.$$

Plugging this into the objective yields

$$\frac{1}{\left|\frac{B}{c_h p + c_g}\right|} \left( \operatorname{Var}(H) - \mathbb{E}[(H - G)^2] + \frac{1}{p} \mathbb{E}\left[(H - G)^2\right] \right).$$

This is a complicated optimization problem because of the floor function, and cannot be solved by setting the gradient to zero. We will begin by searching over all values of  $p \in (0,1]$  for which  $\frac{B}{c_h p + c_g} = k$  for  $k \in \mathbb{N}_+$ , i.e.,  $p \in \left\{\frac{B - k c_g}{k c_h} : k \in \{\lceil B/(c_h + c_g) \rceil, \dots, \lfloor B/c_g \rfloor\}\right\}$ . In terms of k, and denoting  $E = \mathbb{E}[(H - G)^2]$  and  $V = \mathrm{Var}(H) - \mathbb{E}[(H - G)^2]$ , the objective then becomes

$$\frac{1}{k}\left(V + \frac{kc_h}{B - kc_g}E\right) = \frac{1}{k}V + \frac{c_h}{B - kc_g}E.$$

Ignoring the discreteness of k, in the case that  $p^* \in (0,1)$  we can set the derivative to zero, getting

$$\frac{c_g c_h}{(B - k c_g)^2} E = \frac{1}{k^2} V$$

$$\iff k^2 c_g c_h \frac{E}{V} = (B - k c_g)^2$$

$$\iff k^2 \left( c_g^2 - c_g c_h \frac{E}{V} \right) - 2k c_g B + B^2 = 0$$

The positive solution to this quadratic is

$$k = \frac{2c_g B + \sqrt{4c_g^2 B^2 - 4B^2 \left(c_g^2 - c_g c_h \frac{E}{V}\right)}}{2\left(c_g^2 - c_g c_h \frac{E}{V}\right)} = B \frac{1 + \sqrt{\frac{c_h}{c_g} \frac{E}{V}}}{c_g - c_h \frac{E}{V}}.$$

Thus, the optimal  $k^*$  solves the following optimization problem:

$$k^* = \underset{k \in \left\{ \left\lfloor B \frac{1 + \sqrt{\frac{c_h}{c_g} \frac{E}{V}}}{c_g - c_h \frac{E}{V}} \right\rfloor, \left\lceil B \frac{1 + \sqrt{\frac{c_h}{c_g} \frac{E}{V}}}{c_g - c_h \frac{E}{V}} \right\rceil \right\}}{\frac{1}{k}V + \frac{c_h}{B - kc_g}E,$$

And the optimal  $p^*$  is either

$$p^* = \frac{B - k^* c_g}{k^* c_h}$$

or the boundary solution  $p^* = 1$ . To disambiguate between the two, we can directly compute the objective value for each.

#### C.4 Proof of Proposition 5

*Proof.* The asymptotic normality statement can be read off as a simplified version of Theorem 1 from (Zrnic & Candès, 2024). The second part follows because if  $Z \sim \mathcal{N}(0, \Sigma^*)$ , then  $\|Z\|_2^2 = \|(V^*)^{-1/2}Z\|_2^2$ , where  $V^*$  is the eigenvector matrix of  $\Sigma^*$  (since  $(V^*)^{-1/2}$  is unitary). Thus, taking  $\Lambda^*$  to be the (diagonal) eigenvalue matrix of  $\Sigma^*$  and defining we have that  $\|Z\|_2^2 \stackrel{d}{=} \|\Lambda Z'\|_2^2$ , where

$$Z' \sim \mathbb{N}(0, \mathbf{I}_d)$$
. Since  $\|\Lambda Z'\|_2^2 = \sum_{j=1}^d \lambda_j (Z'_j)^2$ , and  $Z'_j \stackrel{\text{i.i.d.}}{\sim} \chi_1^2$ , the result is proven.

### C.5 Proof of Proposition 6

*Proof.* Following the simplification of Problem (5), our problem is equivalent to minimizing the following objective:

$$(c_h \mathbb{E}[\pi(X)] + c_g) \operatorname{Tr}(\Sigma^*). \tag{15}$$

Expanding out  $\Sigma^*$ , we can write

$$\operatorname{Tr}(\Sigma^*) = \operatorname{Tr}\left(W_{\theta^*}^{-1}\operatorname{Var}\left(\nabla \ell_{\theta^*}^g + (\nabla \ell_{\theta^*} - \nabla \ell_{\theta^*}^g)\frac{\xi}{\pi(X)}\right)W_{\theta^*}^{-1}\right)$$

Expanding out the variance gives

$$\operatorname{Var}\left(\nabla \ell_{\theta^*}^g + \left(\nabla \ell_{\theta^*} - \nabla \ell_{\theta^*}^g\right) \frac{\xi}{\pi(X)}\right)$$

$$= \mathbb{E}\left[\left(\nabla \ell_{\theta^*}^g + \left(\nabla \ell_{\theta^*} - \nabla \ell_{\theta^*}^g\right) \frac{\xi}{\pi(X)}\right) \left(\nabla \ell_{\theta^*}^g + \left(\nabla \ell_{\theta^*} - \nabla \ell_{\theta^*}^g\right) \frac{\xi}{\pi(X)}\right)^{\top}\right] - \mathbb{E}[\nabla \ell_{\theta^*}] \mathbb{E}[\nabla \ell_{\theta^*}]^{\top}.$$

Expanding out the squared term yields

$$\begin{split} & \mathbb{E}\left[\left(\nabla\ell_{\theta^*}^g + (\nabla\ell_{\theta^*} - \nabla\ell_{\theta^*}^g) \frac{\xi}{\pi(X)}\right) \left(\nabla\ell_{\theta^*}^g + (\nabla\ell_{\theta^*} - \nabla\ell_{\theta^*}^g) \frac{\xi}{\pi(X)}\right)^\top\right] \\ & = \mathbb{E}\left[\nabla\ell_{\theta^*}^g (\nabla\ell_{\theta^*}^g)^\top\right] \\ & + \mathbb{E}\left[\frac{\xi}{\pi(X)} \left(\nabla\ell_{\theta^*}^g (\nabla\ell_{\theta^*} - \nabla\ell_{\theta^*}^g)^\top + (\nabla\ell_{\theta^*} - \nabla\ell_{\theta^*}^g) (\nabla\ell_{\theta^*}^g)^\top\right)\right] \\ & + \mathbb{E}\left[\left((\nabla\ell_{\theta^*} - \nabla\ell_{\theta^*}^g) \frac{\xi}{\pi(X)}\right) \left((\nabla\ell_{\theta^*} - \nabla\ell_{\theta^*}^g) \frac{\xi}{\pi(X)}\right)^\top\right] \\ & = \mathbb{E}\left[\nabla\ell_{\theta^*}^g (\nabla\ell_{\theta^*}^g)^\top\right] \\ & + \mathbb{E}\left[\nabla\ell_{\theta^*}^g (\nabla\ell_{\theta^*} - \nabla\ell_{\theta^*}^g)^\top + (\nabla\ell_{\theta^*} - \nabla\ell_{\theta^*}^g) (\nabla\ell_{\theta^*}^g)^\top\right] \\ & + \mathbb{E}\left[\frac{1}{\pi(X)} \left((\nabla\ell_{\theta^*} - \nabla\ell_{\theta^*}^g)\right) \left((\nabla\ell_{\theta^*} - \nabla\ell_{\theta^*}^g)\right)^\top\right]. \end{split}$$

Thus, by the linearity of the Tr operator, we can rewrite the trace as  $\text{Tr}(\Sigma^*) = \mathbb{E}\left[\frac{M}{\pi(X)}\right] + C$ , where

$$M = \operatorname{Tr}\left(W_{\theta^*}^{-1} \left(\nabla \ell_{\theta^*} - \nabla \ell_{\theta^*}^g\right) \left(\nabla \ell_{\theta^*} - \nabla \ell_{\theta^*}^g\right)^\top W_{\theta^*}^{-1}\right)$$

and C is

$$\operatorname{Tr}\left(W_{\theta^*}^{-1}\left(\mathbb{E}\left[\nabla \ell_{\theta^*}^g(\nabla \ell_{\theta^*}^g)^{\top} + \nabla \ell_{\theta^*}^g\left(\nabla \ell_{\theta^*} - \nabla \ell_{\theta^*}^g\right)^{\top} + (\nabla \ell_{\theta^*} - \nabla \ell_{\theta^*}^g)\left(\nabla \ell_{\theta^*}^g\right)^{\top}\right] - \mathbb{E}[\nabla \ell_{\theta^*}]\mathbb{E}[\nabla \ell_{\theta^*}]^{\top}\right)W_{\theta^*}^{-1}\right)$$

$$=\operatorname{Tr}\left(W_{\theta^*}^{-1}\left(\mathbb{E}\left[\nabla \ell_{\theta^*}^g(\nabla \ell_{\theta^*})^\top + (\nabla \ell_{\theta^*} - \nabla \ell_{\theta^*}^g)(\nabla \ell_{\theta^*}^g)^\top\right] - \mathbb{E}[\nabla \ell_{\theta^*}]\mathbb{E}[\nabla \ell_{\theta^*}]^\top\right)W_{\theta^*}^{-1}\right).$$

Returning to the objective, and excluding factors that do not depend on  $\pi$ , we can write it now as

$$(c_h \mathbb{E}[\pi(X)] + c_g) \left( \mathbb{E}\left[\frac{M}{\pi(X)}\right] + C \right) \propto_{\pi} (c_h \mathbb{E}[\pi(X)] + c_g) \mathbb{E}\left[\frac{M}{\pi(X)}\right] + c_h \mathbb{E}[\pi(X)]C.$$

In discrete form, following Propostion 2, this is equivalent to

$$(c_h p^{\top} P + c_g) \mathbb{E} \left[ \frac{M}{p^{\top} I} \right] + c_h p^{\top} P C.$$

Taking the derivative with respect to p and setting it to zero coordinatewise yields

$$c_h P_x \mathbb{E}\left[\frac{M}{p^\top I}\right] + c_h P_x C = (c_h p^\top P + c_g) \mathbb{E}\left[M I_x\right],$$

and thus,

$$p_x = \sqrt{\frac{(c_h p^\top P + c_g) \mathbb{E}\left[M \mid X = x\right]}{c_h \mathbb{E}\left[\frac{M}{p^\top I}\right] + c_h C}} \propto_x \sqrt{\mathbb{E}\left[M \mid X = x\right]} = \sqrt{U(x)}.$$

Plugging  $\pi(x) = \gamma \sqrt{\mathbb{E}[M \mid X = x]}$  back into (15) gives the one-dimensional objective

$$\frac{c_g}{\gamma} \mathbb{E} \left[ \frac{M}{\sqrt{\mathbb{E}[M \mid X = x]}} \right] + c_h \gamma \mathbb{E}[\sqrt{\mathbb{E}[M \mid X = x]}] C.$$

The tower property gives us that  $\mathbb{E}\left[\frac{M}{\sqrt{\mathbb{E}[M|X=x]}}\right] = \mathbb{E}\left[\sqrt{\mathbb{E}\left[M\mid X=x\right]}\right]$ , yielding the objective

$$\frac{c_g}{\gamma} \mathbb{E}\left[\sqrt{\mathbb{E}\left[M \mid X = x\right]}\right] + c_h \gamma \mathbb{E}\left[\sqrt{\mathbb{E}\left[M \mid X = x\right]}\right]C,$$

which is equivalent to minimizing

$$\frac{c_g}{\gamma} + c_h \gamma C.$$

The solution to this problem is

$$\gamma^* = \sqrt{\frac{c_g}{c_h} \cdot \frac{1}{C}}.$$

# C.6 Proof of Proposition 8

*Proof.* Following the derivation in Section B.1, we have that for any  $\pi$ 

$$\operatorname{Var}(\Delta^{\pi}) = \operatorname{Var}(H) - \mathbb{E}[(H - G)^{2}] + \mathbb{E}\left[(H - G)^{2} \frac{1}{\pi(X)}\right].$$

We then immediately get that

$$\operatorname{Var}(\Delta^{\tilde{\pi}}) - \operatorname{Var}(\Delta^{\pi^*}) = \mathbb{E}\left[\frac{(H-G)^2}{\tilde{\pi}(X)} - \frac{(H-G)^2}{\pi^*(X)}\right] \le b\mathbb{E}\left[\frac{1}{\tilde{\pi}(X)} - \frac{1}{\pi^*(X)}\right] \le b\delta.$$

### 1351 C.7 PROOF OF COROLLARY 9

*Proof.* Since

$$\tilde{\pi}(x) = (\gamma^* + \delta_{\gamma})\sqrt{U(x) + \delta_U(x)},$$

we have

$$\frac{1}{\tilde{\pi}(x)} = \frac{1}{(\gamma^* + \delta_{\gamma})\sqrt{U(x) + \delta_U(x)}} = \frac{1}{\gamma^* \sqrt{U(x)}} \frac{1}{\left(1 + \frac{\delta_{\gamma}}{\gamma^*}\right)\sqrt{1 + \frac{\delta_U(x)}{U(x)}}}.$$

A first-order Taylor expansion yields

$$\frac{1}{\left(1+\frac{\delta_{\gamma}}{\gamma^{*}}\right)\sqrt{1+\frac{\delta_{U}(x)}{U(x)}}}=1-\frac{\delta_{\gamma}}{\gamma^{*}}-\frac{1}{2}\frac{\delta_{U}(x)}{U(x)}+o\Big(\delta_{\gamma},\frac{\delta_{U}(x)}{U(x)}\Big).$$

Thus,

$$\frac{1}{\tilde{\pi}(x)} - \frac{1}{\pi^*(x)} = \frac{-\delta_{\gamma}}{(\gamma^*)^2 \sqrt{U(x)}} - \frac{1}{2\gamma^*} \frac{\delta_U(x)}{U(x)^{3/2}} + o\Big(\delta_{\gamma}, \frac{\delta_U(x)}{U(x)}\Big).$$

Ignoring second-order terms, since  $U(x) \ge \epsilon$  almost surely, we have

$$\left|\frac{1}{\tilde{\pi}(x)} - \frac{1}{\pi^*(x)}\right| \le \frac{|\delta_{\gamma}|}{(\gamma^*)^2 \sqrt{\epsilon}} + \frac{1}{2\gamma^* \epsilon^{3/2}} |\delta_U(x)|.$$

Taking the expectation and using linearity,

$$\mathbb{E}\left[\frac{1}{\tilde{\pi}(X)} - \frac{1}{\pi^*(X)}\right] \le \frac{|\delta_{\gamma}|}{(\gamma^*)^2 \sqrt{\epsilon}} + \frac{1}{2\gamma^* \epsilon^{3/2}} \, \mathbb{E}[|\delta_U(X)|].$$

Plugging this bound into the initial inequality for  $Var(\Delta^{\tilde{\pi}})$  completes the proof:

$$\operatorname{Var}(\Delta^{\tilde{\pi}}) - \operatorname{Var}(\Delta^{\pi^*}) \le b \left( \frac{|\delta_{\gamma}|}{(\gamma^*)^2 \sqrt{\epsilon}} + \frac{1}{2\gamma^* \epsilon^{3/2}} \mathbb{E}[|\delta_U(X)|] \right).$$

# C.8 PROOF OF PROPOSITION 10

*Proof.* We will borrow notation from the proof of Proposition 2, and express all quantities in vector form. The optimization problem in (8) only depends on Q through the likelihood ratio,  $\frac{dP}{dQ} = r \in \mathbb{R}_{>0}^{|\mathcal{X}|}$ , and Q, P are absolutely continuous with respect to one another. So, we will learn r and then calculate  $Q^* = P/r$ .

Ignoring terms that do not depend on r, the problem in (8) can be rewritten as

Forming the Lagrangian,

$$\mathcal{L}(r,\lambda) = r^{\top} \mathbb{E}_{P} \left[ I \left( H^2 + \left( \frac{1}{\pi(X)} - 1 \right) (H - G)^2 \right) \right] + \lambda ((1/r)^{\top} P - 1).$$

Taking the gradient gives

$$\nabla_r \mathcal{L}(r,\lambda) = \mathbb{E}_P \left[ I \left( H^2 + \left( \frac{1}{\pi(X)} - 1 \right) (H - G)^2 \right) \right] - \lambda P / (r^2),$$

and setting it to zero yields

$$r_x^* \propto_x \sqrt{\frac{1}{\mathbb{E}_P\left[\left(H^2 + \left(\frac{1}{\pi(X)} - 1\right)(H - G)^2\right) \left|X = x\right]}} = \sqrt{\frac{1}{\nu_x}}.$$

To ensure the proper normalization, we set

$$r_x^* = \frac{\sqrt{\nu}^\top P}{\sqrt{\nu_x}}.$$

Thus, 
$$Q^*(x) = P/r^* = \frac{\sqrt{\nu_x} P_x}{\sqrt{\nu}^{\mathsf{T}} P}$$
.

# D ADDITIONAL EMPIRICAL RESULTS

#### D.1 BERNOULLI DATA

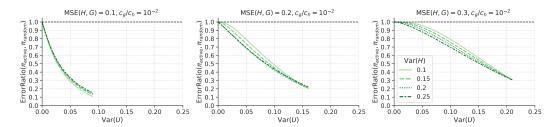


Figure 4: Results on the Bernoulli data (§3.3) for  $\pi_{\text{active}}$  vs.  $\pi_{\text{random}}$  while varying MSE(H,G) and Var(U). As in Figure 1, each line plots a different value of Var(H), where we choose values that are representative of low, medium, or high variance settings compared to MSE(H,G).

Figure 4 provides results for the Bernoulli data setting in Section 3.3 when comparing  $\pi_{\text{active}}$  to  $\pi_{\text{random}}$ . Recall that here the results differ from comparing to  $\pi_{\text{base}}$  only when  $\text{MSE}(H,G) < \frac{c_h}{c_h + c_q} \text{Var}(H)$ , as otherwise the optimal sampling rate for  $\pi_{\text{random}}$  is simply  $p^* = 1$ .

#### D.1.1 ON THE ERROR RATIO LOWER BOUND

It is interesting to observe that  $\mathsf{ErrorRatio}(\pi_{\mathsf{active}}, \pi_{\mathsf{base}})$  is lower-bounded in the Bernoulli data setting at a value close to  $\mathsf{MSE}(H,G)$ . To see why, we note that the lowest value of  $\mathsf{ErrorRatio}(\pi_{\mathsf{active}}, \pi_{\mathsf{base}})$  is obtained when U is maximum variance—which is achieved when U is a binary random variable that is 1 when  $G \neq H$ , and 0 otherwise. Recall that in the Bernoulli data setting both H and G are binary, and  $\mathsf{MSE}(H,G) = \mathbb{P}(H \neq G)$ . We can then compute  $\mathsf{ErrorRatio}(\pi_{\mathsf{active}}, \pi_{\mathsf{base}})$  after optimizing over  $\tau$  as approaching

$$\min \left( \begin{array}{c} \left( \gamma \text{MSE}(H,G) + \frac{c_g}{c_h} \right) \left( 1 + (\frac{1}{\gamma} - 1) \frac{\text{MSE}(H,G)}{\text{Var}(H)} \right) \\ \text{MSE}(H,G) + \frac{c_g}{c_h} \end{array} \right)$$

where 
$$\gamma = \sqrt{\frac{c_g/c_h}{(\text{Var}(H) - \text{MSE}(H, G))_+}}$$
.

Note that we have the first quantity only when  $MSE(H, G) \leq Var(H) + c_a/c_h$ .

Derivation. When  $U \to \mathbb{1} \{H \neq G\} \in \{0,1\}$ , from Proposition (2)  $\pi_{\text{active}}$  approaches either

$$\pi_{\mathrm{clip}}(x,\tau=1) = \begin{cases} \gamma^*(1) & \text{if } h(x) \neq g(x) \\ 0 & \text{otherwise} \end{cases} \quad \text{or} \quad \pi_{\mathrm{clip}}(x,\tau=0) = \begin{cases} 1 & \text{if } h(x) \neq g(x) \\ 0 & \text{otherwise} \end{cases}$$

where  $\gamma^*(1) = \sqrt{\frac{c_g/c_h}{(\operatorname{Var}(H) - \operatorname{MSE}(H,G))_+}}$ . Plugging these values into the optimization over  $\tau \in \{0,1\}$ ,  $\tau^* = \operatorname*{argmin}_{\tau \in \{0,1\}} \left( c_h \mathbb{E}[\pi_{\operatorname{clip}}(x;\tau)] + c_g \right) \left( \operatorname{Var}(H) + \mathbb{E}\left[ U\left(\pi_{\operatorname{clip}}(x;\tau)^{-1} - 1\right) \right] \right),$ 

at  $\tau = 1$  we get

$$(c_h \gamma^*(1) \text{MSE}(H, G) + c_g) \left( \text{Var}(H) + \left( \frac{1}{\gamma^*(1)} - 1 \right) \text{MSE}(H, G) \right),$$

and at  $\tau = 0$  we get

$$(c_h MSE(H,G) + c_q) Var(H),$$

so the optimal  $\tau^*$  is the smaller of the two. Dividing each by  $c_h \text{Var}(H)$  and taking the minimum gives the result for  $\text{ErrorRatio}(\pi_{\text{active}}, \pi_{\text{base}})$ .

A similar calculation can also be made for  $\operatorname{ErrorRatio}(\pi_{\operatorname{active}}, \pi_{\operatorname{random}})$ , with different bounds for when  $\operatorname{MSE}(H,G) \leq \operatorname{Var}(H) - \frac{c_g}{c_h}$  and/or  $\operatorname{MSE}(H,G) \leq \frac{c_h}{c_g+c_h} \operatorname{Var}(H)$  (i.e., both conditions, one or the other condition, or neither condition).

# D.2 REAL DATA

We provide experimental results on three additional datasets:

**AQA.** Attributed Question Answering (AQA) (Bohnet et al., 2023) assesses if a QA system's answer is both correct *and* supported by the text of a document provided as evidence for it (also by the QA system). We evaluate the highest-scoring "retrieve-and-read" system from the dataset. H is a binary human label that is 1 only if the answer is both correct *and* attributable. G is the probability predicted by an 11B parameter T5 model (Raffel et al., 2020). The T5 model is finetuned on a collection of natural language entailment tasks (Honovich et al., 2022). U is computed as G(1-G).

**ImageNet.** The ImageNet dataset (Deng et al., 2009) categorizes input images into one of 1k classes. Our goal is to evaluate the accuracy  $\mathbb{E}[H]$  of a pretrained ResNet model (He et al., 2016), where H is the binary indicator of whether the model's prediction matches the human label for a given image X. G is the softmax value the model assigns to its predicted class. U is computed as G(1-G).

**Seahorse.** The Seahorse dataset (Clark et al., 2023) focuses on multilingual summarization. We focus on the "attribution to the source document" metric for summaries produced by a finetuned 13B parameter mT5 model (Xue et al., 2020). H comes from human ratings. G is the probability score from a finetuned mT5-XXL autorater model assessing attribution. G is computed as G(1-G).

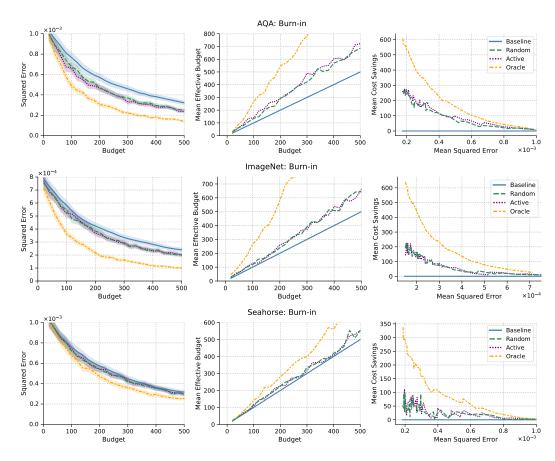


Figure 5: Results on the AQA, ImageNet, and Seahorse datasets using 200 examples as a burn-in (approach A2 in Section 4). The budget on the x-axis reflects "additional" budget used *after* the burn-in examples.

Results are shown in Figure 5, with similar takeaways as the other burn-in (approach A2) experiments in Section 4.3. For AQA and ImageNet, both  $\pi_{\text{active}}$  and  $\pi_{\text{random}}$  substantially outperform  $\pi_{\text{base}}$ ;

<sup>&</sup>lt;sup>2</sup>This checkpoint is available at

https://huggingface.co/collections/google/seahorse-release-6543b0c06d87d83c6d24193b

however, the estimated  $\pi_{\text{active}}$  still leaves a significant amount of headroom behind with respect to the oracle active policy, and has comparable performance to  $\pi_{\text{random}}$ . The Seahorse dataset is an interesting case where the weak rater G is simply not that good, even conditionally. This results in small (but still positive) gains for both the active and random policies—even when  $\pi_{\text{active}}$  uses oracle parameters.

### E IMPLEMENTATION DETAILS

All real data experiments in this paper were performed retrospectively with G and H computed once offline for all inputs x in each dataset. Pretrained models and labels for G and H, respectively, were used for all datasets except Chatbot Arena (Section 4). All subsequent experiments for active sampling were then performed on CPU resources with 32GB of RAM.

For the Chatbot Arena dataset, we sampled responses from Gemini 1.5 Flash (Gemini Team, 2024) using an adapted version of the Chatbot Arena auto-eval prompt.<sup>3</sup> Below is an example prompt. Color is added for clarity. Ten responses from Gemini 1.5 Flash are then sampled, with five responses using the same prompt with the order of A and B flipped. The final label is taken as the majority vote.

Please act as an impartial judge and evaluate the quality of the responses provided by two AI Assistants to the user prompt displayed below. You will be given Assistant A's answer and Assistant B's answer. Your job is to evaluate which assistant's answer is better.

When evaluating the assistants' answers, first identify any mistakes or inaccurate information. Next, consider if the assistant's answers are helpful, relevant, and concise. Helpful means the answer correctly responds to the prompt or follows the instructions. Note that when the user prompt has any ambiguity or more than one interpretation, it is more helpful and appropriate to ask for clarifications or more information from the user than providing an answer based on assumptions. Relevant means all parts of the response closely connect or are appropriate to what is being asked. Concise means the response is clear and not verbose or excessive. Then consider the creativity and novelty of the assistant's answers when needed. Finally, identify any missing important information in the assistants' answers that would be beneficial to include when responding to the user prompt.

<|User Prompt|> make a haiku on bacon the food <|The End of User Prompt|>

<|The Start of Assistant A's Answerl>
Crisp strips of delight,
Sizzling dance, morning's first light,
Bacon whispers, "Bite."
<|The End of Assistant A's Answerl>

<|The Start of Assistant B's Answerl>
Here is a haiku about bacon:

Sizzling in pan
Savory salty bacon strips
Crispy delight yum
<|The End of Assistant B's Answer|>

Is the higher quality response:
(A) Assistant A is better
(B) Assistant B is better

Please answer with either (A) or (B).

3https://github.com/lm-sys/arena-hard-auto/blob/main/config/judge\_ config.yaml For G, we finetune Gemma-3 4B for two hours on TPUv3 resources. The same prompt is used, however, we maximize the log-likelihood of the target Gemini-based answer used for H instead of sampling. Early stopping is done based on the validation loss at predicting H using a held-out split of the training data (recall that the training data is composed of other model comparisons from the Chatbot Arena dataset that are distinct from the one on which we evaluate our annotation policies).