

---

# Embroid: Unsupervised Prediction Smoothing Can Improve Few-Shot Classification

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

Recent work has shown that language models’ (LMs) prompt-based learning capabilities make them well suited for automating data labeling in domains where manual annotation is expensive. The challenge is that while writing an initial prompt is cheap, improving a prompt is costly—practitioners often require significant labeled data in order to evaluate the impact of prompt modifications. Our work asks whether it is possible to improve prompt-based learning *without* additional labeled data. We approach this problem by attempting to modify the predictions of a prompt, rather than the prompt itself. Our intuition is that accurate predictions should also be consistent: samples which are similar under some feature representation should receive the same prompt prediction. We propose EMBROID, a method which computes multiple representations of a dataset under different embedding functions, and uses the consistency between the LM predictions for neighboring samples to identify mispredictions. EMBROID then uses these neighborhoods to create additional predictions for each sample, and combines these predictions with a simple latent variable graphical model in order to generate a final corrected prediction. In addition to providing a theoretical analysis of EMBROID, we conduct a rigorous empirical evaluation across six different LMs and up to 95 different tasks. We find that (1) EMBROID substantially improves performance over original prompts (e.g., by an average of 7.3 points on GPT-JT), (2) also realizes improvements for more sophisticated prompting strategies (e.g., chain-of-thought), and (3) can be specialized to domains like law through the embedding functions.

## 1 Introduction

Acquiring labeled data for domains like medicine and law is essential to training machine learning models or performing basic data analysis (e.g., “how many contracts contain a choice-of-forum clause” or “how many patient medical histories discuss an adverse reaction to a drug?”) [15, 18]. However, building large labeled datasets is difficult, and efforts like [24] show that manual labeling with domain experts is cost-prohibitive. Recent works have begun exploring if language models (LMs) could learn annotation tasks *in-context* [6] and replace manual labeling at scale [13, 15, 22, 30]. The promise of this approach is that LMs’ in-context capabilities enable them to learn tasks from descriptions of the task (i.e., *prompts*). However, the challenge is that producing high performance prompts is still expensive, as practitioners require labeled data in order to measure the impact of modifications to a prompt [49]. Existing work has thus focused on how domain experts can optimally construct prompts for a task (*prompt-engineering*), while minimizing reliance on labeled data [40, 56, 60]. Yet, because language models are sensitive to even small changes in prompt language, these techniques are imperfect and still produce erroneous predictions [2, 8, 37, 49, 66].

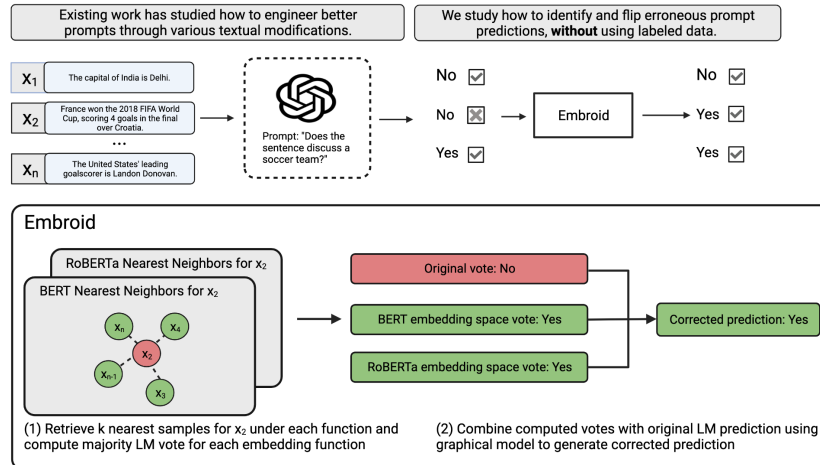


Figure 1: The EMBROID method for prompt-patching.

Our work approaches the challenge of improving prompt performance without labels from an orthogonal perspective: given the predictions of any prompted LM, can we identify and correct mis-predictions using *unlabeled* data? We describe this as the problem of *prompt-patching*. In the context of data annotation tasks, prompt-patching methods should meet three goals. First, they should be theoretically explainable, so that practitioners can understand when and how to apply them. Second, they should be fast, so that practitioners can efficiently integrate them into existing workflows. Finally, they should rarely be wrong, so that they don’t worsen the performance of predictions.

Our work presents EMBROID: a method for automatically identifying and correcting LM predictions with unlabeled data and no expert supervision. Recent work has shown that for many tasks, samples close-by in embedding spaces (produced by models like BERT) have the same label [7]. EMBROID applies this intuition to the prompt-patching regime. Specifically, after LM predictions for all samples have been generated, EMBROID retrieves the  $k$  most similar samples for each input, under  $N$  different embedding functions. For each embedding function, EMBROID computes a scaled-modified majority vote over the LM’s predictions for the  $k$  retrieved samples. EMBROID then combines these  $N$  votes with the original LM prediction for the test sample using a simple latent variable graphical model that is learned with a fast method-of-moments estimator [14]. The intuition behind EMBROID is that good prompts are *smooth* with respect to their predictions over a dataset—samples which are proximate under an embedding function should receive consistent predictions. Thus, modifying the predictions of a prompt to increase neighborhood agreement can improve the accuracy of those predictions. Lastly, because a single embedding space may imperfectly capture similarities between samples, retrieving neighbors from multiple embedding spaces improves robustness [27, 39].

Because EMBROID relies on weak-supervision—the subject of recent rigorous study [7]—it is possible to theoretically analyze and explain *why* and *when* EMBROID will improve performance. In particular, we find that performance is a function of the quality of the embeddings and the performance of the initial prompt. We also empirically study EMBROID, conducting experiments over six LMs, on up to 95 tasks, with several different prompt strategies. We find that EMBROID rarely worsens performance, and often improves F1 by a substantial margin. For instance, EMBROID improves GPT-3.5 by an average of 4.9 points F1 per task, and GPT-JT by an average of 7.3 points per task. The magnitude of EMBROID’s gains are such that it enables a 1.3B parameter model to outperform an instruction-tuned 6.7B parameter model. EMBROID is also complementary to advanced prompt engineering strategies, and achieves performance improvements when applied to prompts designed using chain-of-thought [60], AMA [2], and selective annotation [56]. Finally, EMBROID can be extended to specialized domains like law, through the use of already-available domain specific embeddings.

Succinctly, our contributions in this paper are: (1) EMBROID, a simple prompt-patching framework for improving LM predictions over text classification tasks; (2) a theoretical analysis of EMBROID

72 which explains performance improvements in terms of embedding smoothness and base accuracy;  
 73 and (3) an empirical evaluation of EMBROID covering up to 95 tasks and six different LMs.

## 74 2 Related work

75 **Improving LM performance.** Improving the in-context generalization abilities of LMs has been  
 76 intensely studied. The first family of approaches focuses on adapting LMs in order to make them more  
 77 amenable to prompting. This includes task-specific finetuning [25, 26, 46], training on instruction  
 78 data [9, 57], RLHF [48], and weight-surgery methods which attempt to “correct” incorrect information  
 79 stored in model weights [10, 21, 42, 43]. A second family of approaches explores strategies for  
 80 optimizing prompts to models, either through the specific textual features of the prompt [28, 44,  
 81 60], the use of task decompositions or LM recursion [2], implicit prompt representations [35, 36],  
 82 or external databases [45]. Prompt-patching, in contrast, focuses on identifying mistakes in the  
 83 predictions generated from a particular prompt. The most related approaches are aggregation methods,  
 84 in which the outputs of multiple prompts are combined with an ensembling method [2, 40]. We find  
 85 that EMBROID outperforms many such baselines, and can be applied to enhance their outputs.

86 **Weak supervision.** EMBROID leverages statistical techniques developed in the weak supervision  
 87 literature. The objective in weak supervision is to generate probabilistic labels for unlabeled data by  
 88 combining the predictions of multiple noisy heuristics [14, 51, 52, 55, 62]. EMBROID’s novelty is that  
 89 it uses embeddings to construct additional synthetic predictions, which are combined with the original  
 90 predictions. In contrast, recent weak supervision approaches which incorporate embeddings use them  
 91 to produce more fine-grained accuracy parameters [7], detect and discard training points [33], and as  
 92 the basis for label propagation with final weak supervision predictions [50].

## 93 3 Problem setup and background

94 **Problem setup.** Our problem setup comprises three elements: an unlabeled dataset, predictions from  
 95 a LM for each sample in this dataset, and embedding representations of our dataset. Our goal is to  
 96 improve the accuracy of LM predictions, by using the embedding representations to identify and  
 97 correct predictions likely to be incorrect. Because recent work has explored how predictions from  
 98 multiple prompts can be combined for a task [2], we present a generalized version of EMBROID in  
 99 which we have access to multiple LM predictions. In our empirical evaluation however, we show that  
 100 EMBROID performs well regardless of the number of predictions per sample available.

101 More formally, we focus on a binary classification task where  $x \in \mathcal{X}$  denotes a sentence or paragraph  
 102 and  $y \in \mathcal{Y} = \{-1, 1\}$  is the binary label. We assume we are given an unlabeled dataset  $\mathcal{D} = \{x_i\}_{i=1}^{n_u}$   
 103 of  $n_u$  points. Each point  $x$  is sampled i.i.d. from a distribution  $\mathbb{P}_x$ , and there exists a true underlying  
 104 distribution  $\mathbb{P}$  on the joint  $(x, y)$ . Following the true few-shot regime [49], we assume the only  
 105 labels available are those used in the prompt. We denote a language model (e.g., GPT-3) as  $\lambda_{\text{LLM}}$ ,  
 106 and a task-specific prompt as  $\phi$ , which prepends task instructions to input  $x$  (e.g., “Does the clause  
 107 contain an audit provision? Yes or No.”). The prediction this prompt induces for  $\lambda_{\text{LLM}}$  over  $x$  is  
 108  $\lambda_{\text{LLM}}(\phi(x)) \in \mathcal{Y}$ .<sup>1</sup> Varying  $\phi$  by changing the task description, in-context demonstrations, or punctu-  
 109 ation will alter the prediction generated for  $x$ . For a set of  $m$  prompts  $[\phi_1, \dots, \phi_m]$ , we denote their  
 110 respective predictions on  $x$  as a vector of *weak sources*  $\lambda(x) = [\lambda_{\text{LLM}}(\phi_1(x)), \dots, \lambda_{\text{LLM}}(\phi_m(x))]$ .  
 111 For convenience, we denote  $\lambda_{\text{LLM}}(\phi_i(x))$  as  $\lambda_i(x)$  or  $\lambda_i$  when the  $x$  is obvious, and similarly use  $\lambda$   
 112 instead of  $\lambda(x)$ . We distinguish between two regimes: in the *single-prompt* regime with  $m = 1$ , we  
 113 have access to a LM prediction for each point  $x$ , while in the *multi-prompt* regime with  $m > 1$ , we  
 114 have access to multiple predictions.

115 We assume access to  $N$  embedding models  $\mathcal{E} = [E_1, \dots, E_N]$ , each represented as a fixed mapping  
 116  $E_i : \mathcal{X} \mapsto \mathcal{Z}_i$  from an input  $x$  to an embedding vector  $z$ . These auxiliary embedding models  
 117 provide representations of  $x$  which encode different types of similarity information. Through model  
 118 repositories like HuggingFace [61], it is possible to download a number of models which generate  
 119 representations for text sentences (e.g., BERT or RoBERTa [12, 41]). These models have the property  
 120 that semantically similar sentences are close-by in embedding space [7, 27, 39].

<sup>1</sup>We assume a task-specific mapping function which allows a practitioner to associate a text generation from  
 an LM to a particular class prediction in  $\mathcal{Y}$ .

---

**Algorithm 1** EMBROID: Correcting LLMs with embeddings
 

---

**Input:** Unlabeled data  $\mathcal{D}$ , LLM predictions  $\lambda(x)$  for each  $x \in \mathcal{D}$ , embedding models  $\mathcal{E} = \{E_1, \dots, E_N\}$ , shrinkage parameter  $\tau$ , nearest neighbors parameter  $k$   
**for all** unlabelled  $x \in \mathcal{D}$  **do**  
   **for all** embedding models  $E_j \in \mathcal{E}$  **do**  
     Compute  $k$ -nearest neighbors  $\text{NN}_{j,k}(x)$   
     Compute smoothed neighborhood prediction  $\lambda_{\text{sm},j}(x)$  using  $\lambda$ ,  $\text{NN}_{j,k}(x)$ , and  $\tau$  using eq. (2)  
   **end for**  
**end for**  
 Solve graphical model  $\text{Pr}(y, \lambda(x), \lambda_{\text{sm}}(x))$  in eq. (3) with triplet method over  $\mathcal{D}$  (Algorithm 2).  
**for all** unlabeled  $x \in \mathcal{D}$  **do**  
   Sample  $\hat{y}_x \sim \hat{\text{Pr}}(y|\lambda(x), \lambda_{\text{sm}}(x))$   
**end for**  
**Output:** Label set  $\hat{Y} = \{\hat{y}_x | x \in \mathcal{D}\}$

---

121 **Weak supervision background.** Weak supervision uses a graphical model to combine votes from  
 122 multiple noisy sources into a single prediction, by estimating the accuracy of each source. It mod-  
 123 els  $\text{Pr}(y, \lambda(x))$  as a latent variable graphical model and uses  $\hat{y} = \arg\max_y \hat{\text{Pr}}(y|\lambda(x))$  to produce  
 124 label estimates, where  $\hat{\text{Pr}}$  represents the learned model. The graphical model is based on a graph  
 125  $G = (V, E)$ , where  $V = y \cup \lambda$  and  $E$  consists of edges from  $y$  to each  $\lambda_j$ . We assume no dependen-  
 126 cies between sources, although simple extensions can incorporate them [59]. The formal graphical  
 127 model is:

$$\text{Pr}(y, \lambda(x)) = \frac{1}{Z} \exp(\underbrace{\theta_y y}_{(I)} + \underbrace{\theta^\top \lambda(x) y}_{(II)}) \quad (1)$$

128 where  $Z$  is the partition function used for normalization,  $(I)$  represents a label balance term with  
 129 parameter  $\theta_y$  controlling the prior of  $\text{Pr}(y = 1)$ , and  $(II)$  represents the source accuracy term  
 130 where each  $\theta_i$  is an *accuracy parameter* for the  $i$ th source. Note that from this model, sources are  
 131 conditionally independent:  $\lambda_i \perp \lambda_j | y$  for any  $i, j \in [m]$ . Our use of this model has two steps. First,  
 132 we must learn the accuracy parameters of  $\text{Pr}(y, \lambda(x))$  without access to  $y$ . We use the triplet method  
 133 introduced in [14], which is an efficient method-of-moments estimator for the parameters. Then, at  
 134 inference we compute  $\hat{\text{Pr}}(y|\lambda(x))$ . Appendix B contains more details.

## 135 4 EMBROID

136 First, EMBROID uses the embedding models  $\mathcal{E}$  to compute additional votes for each  $x$ . Let  
 137  $\text{NN}_{j,k}(x) \subset \mathcal{D}$  be the  $k$ -nearest neighbors of sample  $x$  under the embedding function  $E_j$ . We define  
 138 the smoothed neighborhood prediction vector  $\lambda_{\text{sm},j}(x) \in \{-1, 0, 1\}^m$  as follows, with  $\lambda_{\text{sm},j}[i](x)$   
 139 being the  $i$ th element:

$$\begin{aligned} \tilde{\lambda}_j[i](x) &= \frac{1}{k} \sum_{\tilde{x} \in \text{NN}_{j,k}(x)} \lambda_i(\tilde{x}) \\ \lambda_{\text{sm},j}[i](x) &= \begin{cases} 1 & \tilde{\lambda}_j[i](x) > \tau_i^+ \\ -1 & \tilde{\lambda}_j[i](x) < \tau_i^- \\ 0 & \text{o.w.} \end{cases} \end{aligned} \quad (2)$$

140 where  $\tau_i^+ \in [-1, 1]$  and  $\tau_i^- \in [-1, 1]$  act as shrinkage parameters for  $\lambda_i$  which control the level of  
 141 agreement amongst the neighbors of  $x$  necessary to generate a particular vote. The scalar  $\lambda_{\text{sm},j}[i](x)$  is  
 142 the average vote of  $\lambda_i$  amongst the neighbors of  $x$  in  $E_j$ . When  $\lambda_{\text{sm},j}[i](x)$  is sufficiently positive, i.e.,  
 143  $\lambda_{\text{sm},j}[i](x) > \tau_i^+$ , EMBROID sets  $\lambda_{\text{sm},j}[i](x)$  to be a positive vote. When  $\lambda_{\text{sm},j}[i](x)$  is sufficiently  
 144 negative, i.e.,  $\lambda_{\text{sm},j}[i](x) < \tau_i^-$ , EMBROID sets  $\lambda_{\text{sm},j}[i](x)$  to be a negative vote. Otherwise,  
 145  $\lambda_{\text{sm},j}[i](x)$  is set to be an abstain. The intuition is that  $\lambda_{\text{sm},j}[i](x)$  will be an accurate vote over  $x$   
 146 whenever two conditions are met: (1) the LM is generally accurate, i.e.,  $\lambda_j$  is usually correct, and (2)  
 147  $E_j$  is *smooth*, i.e., nearest-neighbors share the same task label.

Next, we augment our base model in equation (1) to incorporate these auxiliary neighborhood predictions  $\lambda_{\text{sm}} = [\lambda_{\text{sm},1}, \dots, \lambda_{\text{sm},N}] \in \{-1, 0, 1\}^{N_m}$  computed using the embeddings:

$$\Pr(y, \lambda, \lambda_{\text{sm}}) = \frac{1}{Z} \exp \left( \theta_y y + \theta^\top \lambda y + \sum_{j=1}^N \alpha_j^\top \lambda_{\text{sm},j} y \right) \quad (3)$$

where the vector  $\alpha_j \in \mathbb{R}^m$  represents the quality parameters for the  $j^{\text{th}}$  embedding model when used with the  $m$  different prompts. To solve this model and produce label estimates, we note that it has the same format as (1) if we concatenate  $\lambda$  and  $\lambda_{\text{sm}}$  into one set of weak sources. Therefore, we can use the triplet method from [14] to learn parameters and output estimates  $\hat{\Pr}(y|\lambda(x), \lambda_{\text{sm}}(x))$  for each  $x \in \mathcal{D}$  at inference time (see Appendix B for details).

Parameters  $\theta$  and  $\alpha_j$  in (3) allow us to trade-off two different sources of information—one presented by directly prompting an LM to obtain a label and the other by incorporating similarity information from the embedding models—and to further account for varying error modes among the embedding models. Our use of the neighborhood predictions in (3) yields a more expressive model than the standard weak supervision framework solely on LLM predictions in (1), which we can recover when  $k = 0$ , and can thus help make corrections to the LLM predictions. In practice, we find that setting  $\tau_i^+ = \tau_i^- = \mathbb{E}[\lambda_i]$  (i.e., the average source vote) yields good performance (Appendix G).

## 5 Theoretical analysis

We analyze EMBROID, discussing the advantages of using  $\lambda_{\text{sm}}$  in addition to  $\lambda$ , and show that embedding smoothness and base prediction accuracy play a critical role in information gain. Appendix F provides synthetics demonstrating these tradeoffs and comparing to weak-supervision baselines.

First, we provide a result on the generalization error of our model  $\hat{\Pr}(y|\lambda, \lambda_{\text{sm}})$ . Define the generalization error as the expected cross-entropy loss,  $L(\lambda, \lambda_{\text{sm}}, \mathcal{D}) = \mathbb{E}_{y, \lambda(x), \lambda_{\text{sm}}(x), \mathcal{D}} [-\log \hat{\Pr}(y|\lambda(x), \lambda_{\text{sm}}(x))]$ . We use  $[\lambda_1, \dots, \lambda_{(N+1)m}]$  to represent  $[\lambda, \lambda_{\text{sm}}]$  and denote by  $a_{\max} = \max_i \mathbb{E}[\lambda_i(x)y]$  the largest accuracy (scaled to  $[-1, 1]$ ) of any source, and by  $b_{\min} = \min_{i,j} \{\mathbb{E}[\lambda_i \lambda_j], \mathbb{E}[\lambda_i \lambda_j]\}$  the minimum expected pairwise product between any two sources. Assume that all sources are better than random, e.g.,  $\Pr(\lambda_i = y) > 0.5$ . These terms and assumptions are from using the triplet method.

**Proposition 5.1.** *Suppose that the data  $x, y, \lambda, \lambda_{\text{sm}}$  follows the model in (3). The generalization error of  $\hat{\Pr}(y|\lambda, \lambda_{\text{sm}})$  can be decomposed into*

$$L(\lambda, \lambda_{\text{sm}}, \mathcal{D}) \leq \underbrace{H(y|\lambda, \lambda_{\text{sm}})}_{\text{Irreducible Error}} + \underbrace{\frac{C(N+1)m}{n_u}}_{\text{Variance}} + o(1/n_u),$$

where  $C = \frac{3(1-b_{\min}^2)}{8b_{\min}^2(1-a_{\max}^2)} \left( \frac{1}{b_{\min}^4} + \frac{2}{b_{\min}^2} \right)$ .

In the bound above, the variance term comes from estimation error when learning the parameters via the triplet method. The irreducible error depends on quality of  $\lambda$  and  $\lambda_{\text{sm}}$ . If knowledge of the LLM prediction and neighborhood prediction significantly reduces uncertainty in  $y$ , the conditional entropy term  $H(y|\lambda(x), \lambda_{\text{sm}}(x))$  is low.

**Information gain from using both  $\lambda, \lambda_{\text{sm}}$ .** We compare upper bounds on generalization error when both  $\lambda, \lambda_{\text{sm}}$  are modeled versus when only  $\lambda$  is modeled, as in (1) corresponding to classical weak supervision. Based on the bound in Proposition 5.1, modeling both  $\lambda$  and  $\lambda_{\text{sm}}$  increases the variance term by a constant multiplicative factor.

Here, we examine how the irreducible error is affected, that is, the difference  $H(y|\lambda) - H(y|\lambda, \lambda_{\text{sm}})$ . Since this quantity is always nonnegative, we focus on bounding the *pointwise* difference in conditional entropy—which we call the information gain—for a given  $x_0$  on which the LLM is incorrect. For simplicity, suppose we have one embedding  $E$ . An embedding  $E$  is  $M$ -smooth with respect to the label if

$$\Pr(\tilde{y} = c | y = c, \|E(x) - E(\tilde{x})\| \leq \varepsilon) \geq M_E(\varepsilon), \quad (4)$$

where  $c \in \mathcal{Y}$ ,  $\varepsilon > 0$  and  $M_E(\cdot) \in [0, 1]$  is decreasing in its input. Define  $\beta_i = \Pr(\lambda_i = y)$  as the accuracy of  $\lambda_i$  and  $p_\lambda = \Pr(y = 1|\lambda(x_0))$  as the prediction on  $x_0$  given only access to  $\lambda$ .

Let  $\varepsilon_k = \max_{\tilde{x} \in \text{NN}_k(x)} \|E(x) - E(\tilde{x})\|$  be the maximum distance between  $x_0$  and its  $k$  neighbors. Without loss of generality, assume the label on  $x_0$  is  $y = 1$ .

**Theorem 5.2.** Assume that  $E$  is  $M$ -smooth. The pointwise information gain on  $x_0$  is

$$H(y|\lambda(x_0)) - H(y|\lambda(x_0), \lambda_{\text{sm}}(x_0)) \geq 2(1 - p_\lambda) \left[ \prod_{i=1}^m [1 - \exp[-2k(\beta_{\text{NN}_k, i} - 0.5)^2]] - 0.5 \right]$$

where  $\beta_{\text{NN}_k, i} = \Pr_{\tilde{x} \sim \text{NN}_k}(\lambda_i(\tilde{x}) = y) \geq \beta_i M_E(\varepsilon_k)$  is the neighborhood accuracy.

A few observations on the bound are in order.

- **Improvement over WS:** If the neighborhood accuracy is bounded sufficiently far from  $\frac{1}{2}$  and  $k$  is large, using EMBROID has better irreducible error than just using  $\lambda$ . For example, setting  $m = 1$ ,  $k = 10$ ,  $\beta_{\text{NN}_k, i} = 0.7$ , and  $p_\lambda = 0.25$  gives us an improvement of 0.076 nats.
- **Smoothness:** If  $E$  is highly smooth, then  $M_E(\varepsilon_k)$  will be large and irreducible error will be small.
- **Base prediction accuracy:** If the original sources  $\lambda$  have high accuracy ( $\beta_i$ ), irreducible error will be small.

Additionally, we observe that if  $p_\lambda$  is a high-quality prediction close to the true label 1, the information gain is small. Choice of the  $k$  parameter presents a performance trade-off: increasing  $k$  will increase  $\varepsilon_k$  and incorporate farther-away, less reliable predictions, but it will also reduce the noise of the majority vote. We also comment on the information gain when using both  $\lambda$  and  $\lambda_{\text{sm}}$  over just  $\lambda_{\text{sm}}$  in Appendix C.

## 6 Results

Our empirical evaluation focuses on three questions: (1) How robust is EMBROID’s performance across LMs? (2) How does EMBROID, as a prompt-patching method, compare to high performance prompt-engineering methods? (3) How sensitive is EMBROID to the embeddings and dataset size?

**Tasks.** We study tasks where sentence embeddings can capture information relevant to the task, leading us to focus on sentence classification datasets. We consider a collection of 95 class-balanced sentence classification tasks, derived from binarizing existing multi-class legal, scientific, and general domain classification benchmarks like CUAD, AGNews, DBpedia-14, FewRel, and several others [20, 24, 29, 63, 65].<sup>2</sup> Example tasks include, “Classify if the following texts discuss a recording label” or “Classify if the following contractual clauses contain an audit rights provision.”

**Choice of embedding models.** Following prior work illustrating the benefits of domain specific representation [19, 67], EMBROID uses different embeddings for each task domain. For law tasks, we rely on two BERT-variants trained on different legal corpora [23, 67]. For science tasks, we rely on three BERT-variants trained on science, biology, and medical texts [3, 17, 34]. For general domain tasks, we rely on BERT, Roberta, and SentenceBert embeddings [12, 41, 53].

**Prompts.** Prompts are constructed using fixed instructions, and by manually selecting three random samples (from each class) as in-context demonstrations (Appendix E). We follow the true few-shot regime [49], in that we assume the only labeled data available to the data scientist are the labels used for in-context demonstrations. Prior work has found this regime to most realistically represent real-world workflows.

**Models.** We evaluate on two API-access models: GPT-3.5 (text-davinci-003) and J1-Jumbo [38]. Because API models raise significant privacy and compliance concerns for data scientists working with sensitive data [16], we also evaluate on open-source models. We select models in the 6-7B parameter range, as these are the largest models which fit on commonly available 40GB A100 machines. Specifically, we evaluate Bloom [54] and OPT [64]. Given the increasing popularity of instruction-tuning, we also evaluate on GPT-JT [58], an 6.7B parameter instruction tuned version of GPT-J. Because of cost-constraints, we evaluate API-access models on a representative selection of 12 tasks, while evaluating all other models on the full suite of 95 tasks. Appendix D provides details.

<sup>2</sup>We hope to explore multi-class extensions and more complex reasoning tasks in future work.

	LM	Win rate (%)	Avg. Improvement (F1)
API-Access Models	J1-Jumbo (176B)	72.2	10.6
	GPT-3.5 (> 170B)	80.6	4.9
Open Source	Bloom (7.1B)	91.2	10.1
	OPT (6.7B)	91.2	11.6
Instruction Tuned	GPT-JT (6B)	89.1	7.3

Table 1: We evaluate the extent to which EMBROID improves the original prompt on different models in terms of win rate and relative improvement (defined in-line). All models are run with three trials. For each model, we report the percentage of tasks (across all trials) for which EMBROID improves, and the average improvement (in F1 points). Additional details provided in Appendix.

## 6.1 By how much does prompt-patching improve performance?

**Performance across LM families.** We examine if EMBROID achieves improvements for different types of LMs. For each LM, we select three different combinations of in-context demonstrations (i.e., three prompts), generate predictions for each prompt, and apply EMBROID to independently each prompt’s predictions. This produces  $3 \times 95 = 285$  trials for open-source models, and  $3 \times 12 = 36$  trials for API-models. We report *win-rate*, i.e., the proportion of trials for which EMBROID outperforms the original predictions, and *improvement*, i.e., the average difference in F1 points (across all trials) between EMBROID and the original predictions.

As Table 1 illustrates, EMBROID improves performance for a substantial proportion of prompts, by a substantial margin, across all models. On GPT-3.5 for instance, EMBROID achieves a win-rate of 80.6%, with an average of improvement of 4.9 points. EMBROID also improves for open source models, with a win-rate of 91.2% on OPT-6.7 and an average improvement of 11.6 points. Finally, EMBROID achieves similar gains on an instruction tuned model, with a win-rate of 89.1% and an average improvement of 7.3 points.

**Performance when prompts are good.** We additionally investigate how EMBROID’s performance improvements change as a function of the performance of the base prompt. Hypothetically, one could imagine that better performing prompts are *smoother* with respect to embeddings, thus diminishing (or negating) EMBROID. In Figure 2 (upper left), we plot the improvement of EMBROID against the performance of the base prompt for GPT-JT. Even when the base prompt performs well (i.e., F1 > 0.8), EMBROID improves on 89% of tasks by an average of 4.1 points.

**Measuring performance in parameter count.** A trend in recent literature has been to measure the magnitude of improvements to prompt performance in terms of parameter count [2], by showing how a particular method makes a smaller method equivalent in performance to a larger model. We find that EMBROID enables the non-instructed tuned 1.3B GPT-Neo model to outperform an instruction tuned 6.7B model; across all trials, GPT-JT scores an average F1 of 67.8, while EMBROID +GPT-Neo-1.3B scores an average of 68.5.

## 6.2 Comparing prompt-patching to prompt-engineering

Our work distinguishes between prompt-construction methods—which control how a prompt is generated—and prompt-patching methods—which attempt to identify and correct errors in the predictions produced by a prompt. We use EMBROID to further study the difference between these frameworks in two ways. First, we compare EMBROID’s performance improvement over a base prompt to that of several specialized prompting strategies. Second, we examine the extent to which EMBROID—when applied to the predictions produced by these prompting strategies—can generate further performance improvements. We study three prompting strategies:

1. Ensemble strategies, in which the predictions of multiple prompts are combined using an unsupervised ensembling model. Specifically, we compare to two ensembling methods previously studied for LLMs (AMA [2] and majority vote [40]), one ensembling method which incorporates embedding information (Liger [7]), and one well regarded weak supervision baseline (FlyingSquid [14]). Each baseline is run over the predictions generated by three different prompts.

LM	MV	Liger	FlyingSquid	AMA	EMBROID-1	EMBROID-3
J1-Jumbo	47.4	48.7	50.5	60.7	60.4	<b>64.5</b>
GPT-3.5	81.4	82.5	82.1	84.7	83.9	<b>86.0</b>
Bloom-7.1B	54.6	55.8	54.3	63.0	64.7	<b>69.1</b>
OPT-6.7	46.1	46.8	46.3	56.3	59.8	<b>64.2</b>
GPT-JT	69.3	69.4	70.1	74.6	75.1	<b>79.0</b>

Table 2: We evaluate how EMBROID compares to common ensemble approaches for improving prompt prediction performance. All ensemble baselines are run with three sets of predictions. EMBROID-1 is run with one set of predictions, and EMBROID-3 is run with three set of predictions. For each method, we report the average macro-F1 over all tasks. We observe that EMBROID-1 is competitive with ensemble methods which use many more predictions, while EMBROID-3 outperforms all other methods by a substantial margin.

Base prompt	+CoT	+EMBROID	+ CoT + EMBROID
76.3	80.1	81.9	<b>85.4</b>

Table 3: We evaluate EMBROID compared to, and applied to, CoT prompting on GPT-3.5 for a subset of 10 tasks. We report the average across the studied tasks.

2. Chain-of-thought prompting [60], in which for each in-context demonstration, we provide a step-by-step explanation for the demonstration’s label.
3. Selective annotation (SA) [56], in which we use embeddings to select a subset of  $k$  data samples to label, and then, for each input sample, retrieve the most similar samples (under some embedding function) from this pool to use as in-context demonstrations.

**Ensemble methods.** We evaluate two versions of EMBROID. In the first version, we run EMBROID with the predictions of only one prompt (EMBROID-1). In the second version, we run EMBROID with the predictions of three prompts (EMBROID-3). The second version is comparable to applying EMBROID to the outputs of an ensemble method. In Table 2, we observe that EMBROID-1 is competitive with the ensemble baselines (while using substantially fewer predictions), while EMBROID-3 consistently outperforms these baselines (across different LMs).

**Chain-of-thought.** We compare EMBROID to chain-of-thought (CoT) prompting for a subset of a representative subset of 10 tasks on GPT-3.5. For each task, we manually construct a base prompt consisting of six demonstrations, and a CoT prompt where an explanation is provided for each demonstration. We first find that EMBROID’s performance improvement over the base prompt exceeds that of chain-of-thought prompting (Table 3). Using EMBROID to modify the base prompt is better on average than CoT prompting, outperforming CoT on six out of the ten tasks. Second, applying EMBROID to predictions generated by a CoT prompt yields further improvements, outperforming vanilla CoT predictions on eight of ten tasks.

**Selective annotation (SA).** We compare EMBROID to selective annotation with a label budget of [6, 25, 50, 100] (Figure 2, upper-right). For each task, we run selective annotation using a domain specific embedding. EMBROID (applied to a prompt with randomly chosen samples) outperforms selective annotation with a label budget of 25 samples. When a label budget of 100 samples is available, EMBROID improves the performance of a prompt constructed using selective annotation on 88% of tasks, by an average of 4.3 points.

### 6.3 Ablations

Finally, we perform several ablations of EMBROID to study how performance changes as a function of (1) the domain specificity of the embedding used, (2) the quality of the embedding spaces used, and (3) the size of the dataset. Additional ablations are presented in Appendix G.

**Domain specific embeddings improve performance.** We compare how performance on the legal and science tasks changes when we shift from domain specialized embeddings to general domain embeddings. On law tasks for GPT-JT, we find that using two legal embedding spaces outperforms

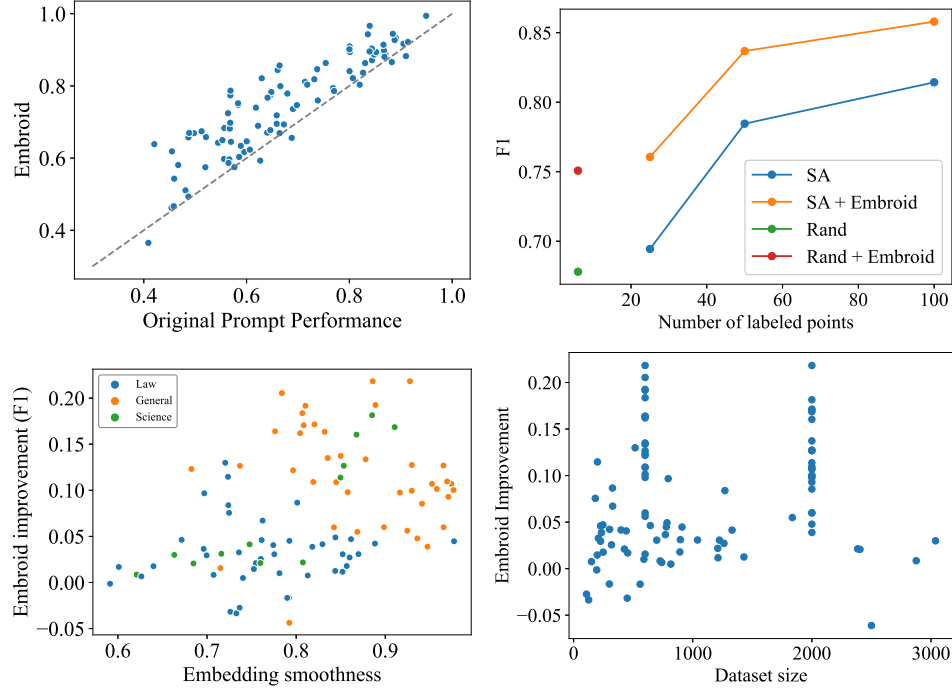


Figure 2: **Upper left:** The EMBROID F1 plotted against the F1 score for the original prompt for GPT-JT. Even for high performing prompts, EMBROID is capable of improving performance. The dashed line  $y = x$  is plotted for visual aid. **Upper right:** A comparison of EMBROID to selective annotation (SA) over all tasks for GPT-JT. **Bottom left:** For each task (using GPT-JT), we plot the performance improvement of EMBROID against the average smoothness of the embeddings used. We observe a positive correlation ( $r = 0.39$ ). **Bottom right:** Across all tasks, we measure the performance improvement of EMBROID against the size of the task.

using BERT and RoBERTa for 77% of tasks, by up to 6 points F1 on certain tasks [23, 67]. For science tasks for GPT-JT, we find that using two science embedding spaces [3, 34] outperforms using BERT and RoBERTa for 92% of tasks, by up to 4.3 points F1 on certain tasks.

**Embedding quality.** Building on Section 5, we compare EMBROID’s performance improvement over the base prompt to the average smoothness of the embedding spaces with respect to each task (Figure 2). We observe a positive correlation: smoother embedding spaces are associated with larger performance gains (with a Pearson coefficient of  $r = 0.39$ ). Applying this insight, we explore how performance changes when *extremely* high quality embeddings are added. For a subset of 19 tasks we generate OpenAI `text-embedding-ada-002` embeddings, and find that adding them to EMBROID improves performance by up to 13 points F1 (at an average of 2 points across all studied tasks).

**Dataset size.** Finally, we study how EMBROID’s performance improvement changes as the dataset size changes. Because EMBROID relies on nearest-neighbors in different embedding spaces, we might expect performance to be poor when the dataset being annotated is small. In Figure 2 (bottom right), we see that EMBROID achieves performance improvements even for “small” datasets with only several hundred samples.

## 7 Conclusion

We study the problem of improving prompt-based learning, by developing a method (EMBROID) for detecting and correcting erroneous predictions without labeled data. We validate EMBROID across a range of datasets and LMs, finding consistent improvement in many regimes. We take a moment to address the societal impact of our work: while we do not foresee any *direct* harmful impacts arising from our work, we caution that any use of language models in meaningful applications should be accompanied by conversations regarding risks, benefits, stakeholder interests, and ethical safeguards.

## References

- [1] Monica Agrawal, Stefan Hegselmann, Hunter Lang, Yoon Kim, and David Sontag. Large language models are few-shot clinical information extractors, 2022.
- [2] Simran Arora, Avanika Narayan, Mayee F. Chen, Laurel Orr, Neel Guha, Kush Bhatia, Ines Chami, Frederic Sala, and Christopher Ré. Ask me anything: A simple strategy for prompting language models, 2022.
- [3] Iz Beltagy, Kyle Lo, and Arman Cohan. Scibert: A pretrained language model for scientific text. *arXiv preprint arXiv:1903.10676*, 2019.
- [4] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- [5] Daniel Borkan, Lucas Dixon, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. Nuanced metrics for measuring unintended bias with real data for text classification. In *Companion proceedings of the 2019 world wide web conference*, pages 491–500, 2019.
- [6] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [7] Mayee F. Chen, Daniel Y. Fu, Dyah Adila, Michael Zhang, Frederic Sala, Kayvon Fatahalian, and Christopher Ré. Shoring up the foundations: Fusing model embeddings and weak supervision, 2022.
- [8] Yanda Chen, Chen Zhao, Zhou Yu, Kathleen McKeown, and He He. On the relation between sensitivity and accuracy in in-context learning. *arXiv preprint arXiv:2209.07661*, 2022.
- [9] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*, 2022.
- [10] Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. Knowledge neurons in pretrained transformers. *arXiv preprint arXiv:2104.08696*, 2021.
- [11] Franck Dernoncourt and Ji Young Lee. Pubmed 200k rct: a dataset for sequential sentence classification in medical abstracts. *arXiv preprint arXiv:1710.06071*, 2017.
- [12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [13] Bosheng Ding, Chengwei Qin, Linlin Liu, Lidong Bing, Shafiq Joty, and Boyang Li. Is gpt-3 a good data annotator? *arXiv preprint arXiv:2212.10450*, 2022.
- [14] Daniel Fu, Mayee Chen, Frederic Sala, Sarah Hooper, Kayvon Fatahalian, and Christopher Re. Fast and three-rious: Speeding up weak supervision with triplet methods. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 3280–3291. PMLR, 13–18 Jul 2020.
- [15] Karan Goel, Sabri Eyuboglu, Arjun Desai, James Zou, and Chris Ré. Meerkat and the path to foundation models as a reliable software abstraction. <https://hazyresearch.stanford.edu/blog/2023-03-01-meerkat>, 2023.
- [16] Karla Grossenbacher. Employers should consider these risks when employees use chatgpt. <https://news.bloomberglaw.com/us-law-week/employers-should-consider-these-risks-when-employees-use-chatgpt>, 2023.
- [17] Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. Domain-specific language model pretraining for biomedical natural language processing, 2020.

[18] Neel Guha, Daniel E Ho, Julian Nyarko, and Christopher Ré. Legalbench: Prototyping a collaborative benchmark for legal reasoning. *arXiv preprint arXiv:2209.06120*, 2022.

[19] Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A Smith. Don’t stop pretraining: adapt language models to domains and tasks. *arXiv preprint arXiv:2004.10964*, 2020.

[20] Xu Han, Hao Zhu, Pengfei Yu, Ziyun Wang, Yuan Yao, Zhiyuan Liu, and Maosong Sun. Fewrel: A large-scale supervised few-shot relation classification dataset with state-of-the-art evaluation. *arXiv preprint arXiv:1810.10147*, 2018.

[21] Peter Hase, Mona Diab, Asli Celikyilmaz, Xian Li, Zornitsa Kozareva, Veselin Stoyanov, Mohit Bansal, and Srinivasan Iyer. Do language models have beliefs? methods for detecting, updating, and visualizing model beliefs. *arXiv preprint arXiv:2111.13654*, 2021.

[22] Xingwei He, Zhenghao Lin, Yeyun Gong, A-Long Jin, Hang Zhang, Chen Lin, Jian Jiao, Siu Ming Yiu, Nan Duan, and Weizhu Chen. Annollm: Making large language models to be better crowdsourced annotators, 2023.

[23] Peter Henderson, Mark S Krass, Lucia Zheng, Neel Guha, Christopher D Manning, Dan Jurafsky, and Daniel E Ho. Pile of law: Learning responsible data filtering from the law and a 256gb open-source legal dataset. *arXiv preprint arXiv:2207.00220*, 2022.

[24] Dan Hendrycks, Collin Burns, Anya Chen, and Spencer Ball. CUAD: an expert-annotated NLP dataset for legal contract review. *CoRR*, abs/2103.06268, 2021.

[25] Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alexander Ratner, Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes, 2023.

[26] Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. Large language models can self-improve. *arXiv preprint arXiv:2210.11610*, 2022.

[27] Qiao Jin, Bhuwan Dhingra, William W Cohen, and Xinghua Lu. Probing biomedical embeddings from language models. *arXiv preprint arXiv:1904.02181*, 2019.

[28] Jaehun Jung, Lianhui Qin, Sean Welleck, Faeze Brahman, Chandra Bhagavatula, Ronan Le Bras, and Yejin Choi. Maieutic prompting: Logically consistent reasoning with recursive explanations. *arXiv preprint arXiv:2205.11822*, 2022.

[29] Martin Krallinger, Obdulia Rabal, Saber Ahmad Akhondi, Martín Pérez Pérez, Jesus Santamaría, Gael Pérez Rodríguez, Georgios Tsatsaronis, Ander Intxaurre, J. A. Lopez, Umesh K. Nandal, Erin M. van Buel, Ambika Chandrasekhar, Marleen Rodenburg, Astrid Lægreid, Marius A. Doornenbal, Julen Oyarzábal, Anália Lourenço, and Alfonso Valencia. Overview of the biocreative vi chemical-protein interaction track. 2017.

[30] Taja Kuzman, Igor Mozetic, and Nikola Ljubešić. Chatgpt: Beginning of an end of manual linguistic data annotation? use case of automatic genre identification. *ArXiv*, abs/2303.03953, 2023.

[31] Stanford Legal Design Lab. Legal issues taxonomy, 2023.

[32] Suffolk Law School Legal Innovation & Technology Lab. Spot’s training data, 2022.

[33] Hunter Lang, Aravindan Vijayaraghavan, and David Sontag. Training subset selection for weak supervision, 2022.

[34] Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240, 2020.

[35] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning, 2021.

- [36] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation, 2021.
- [37] Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. Holistic evaluation of language models. *arXiv preprint arXiv:2211.09110*, 2022.
- [38] Opher Lieber, Or Sharir, Barak Lenz, and Yoav Shoham. Jurassic-1: Technical details and evaluation. *White Paper: AI21 Labs*, 1, 2021.
- [39] Lucy H Lin and Noah A Smith. Situating sentence embedders with nearest neighbor overlap. *arXiv preprint arXiv:1909.10724*, 2019.
- [40] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9):1–35, 2023.
- [41] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [42] Kevin Meng, David Bau, Alex J Andonian, and Yonatan Belinkov. Locating and editing factual associations in gpt. In *Advances in Neural Information Processing Systems*, 2022.
- [43] Kevin Meng, Arnab Sen Sharma, Alex Andonian, Yonatan Belinkov, and David Bau. Mass-editing memory in a transformer. *arXiv preprint arXiv:2210.07229*, 2022.
- [44] Swaroop Mishra, Daniel Khashabi, Chitta Baral, Yejin Choi, and Hannaneh Hajishirzi. Reframing instructional prompts to gptk’s language. *arXiv preprint arXiv:2109.07830*, 2021.
- [45] Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D Manning, and Chelsea Finn. Memory-based model editing at scale. In *International Conference on Machine Learning*, pages 15817–15831. PMLR, 2022.
- [46] Akihiro Nakamura and Tatsuya Harada. Revisiting fine-tuning for few-shot learning. *arXiv preprint arXiv:1910.00216*, 2019.
- [47] Laurel Orr. Manifest. <https://github.com/HazyResearch/manifest>, 2022.
- [48] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022.
- [49] Ethan Perez, Douwe Kiela, and Kyunghyun Cho. True few-shot learning with language models. *Advances in neural information processing systems*, 34:11054–11070, 2021.
- [50] Rattana Pukdee, Dylan Sam, Maria-Florina Balcan, and Pradeep Ravikumar. Label propagation with weak supervision. *arXiv preprint arXiv:2210.03594*, 2022.
- [51] A. J. Ratner, B. Hancock, J. Dunnmon, F. Sala, S. Pandey, and C. Ré. Training complex models with multi-task weak supervision. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Honolulu, Hawaii, 2019.
- [52] Alexander Ratner, Stephen H Bach, Henry Ehrenberg, Jason Fries, Sen Wu, and Christopher Ré. Snorkel: Rapid training data creation with weak supervision. In *Proceedings of the VLDB Endowment. International Conference on Very Large Data Bases*, volume 11, page 269. NIH Public Access, 2017.
- [53] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*, 2019.
- [54] Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*, 2022.

- [55] Changho Shin, Winfred Li, Harit Vishwakarma, Nicholas Carl Roberts, and Frederic Sala. Universalizing weak supervision. In *International Conference on Learning Representations (ICLR)*, 2022.
- [56] Hongjin Su, Jungo Kasai, Chen Henry Wu, Weijia Shi, Tianlu Wang, Jiayi Xin, Rui Zhang, Mari Ostendorf, Luke Zettlemoyer, Noah A. Smith, and Tao Yu. Selective annotation makes language models better few-shot learners, 2022.
- [57] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Alpaca: A strong, replicable instruction-following model. <https://crfm.stanford.edu/2023/03/13/alpaca.html>, 2023.
- [58] Together. Releasing v1 of gpt-jt powered by open-source ai. <https://www.together.xyz/blog/releasing-v1-of-gpt-jt-powered-by-open-source-ai>. Accessed: 2023-01-25.
- [59] Paroma Varma, Frederic Sala, Ann He, Alexander Ratner, and Christopher Re. Learning dependency structures for weak supervision models. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 6418–6427. PMLR, 09–15 Jun 2019.
- [60] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*, 2022.
- [61] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pages 38–45, 2020.
- [62] Renzhi Wu, Shen-En Chen, Jieyu Zhang, and Xu Chu. Learning hyper label model for programmatic weak supervision. *arXiv preprint arXiv:2207.13545*, 2022.
- [63] Jieyu Zhang, Yue Yu, Yinghao Li, Yujing Wang, Yaming Yang, Mao Yang, and Alexander Ratner. Wrench: A comprehensive benchmark for weak supervision. *arXiv preprint arXiv:2109.11377*, 2021.
- [64] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.
- [65] Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28, 2015.
- [66] Tony Z. Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. Calibrate before use: Improving few-shot performance of language models, 2021.
- [67] Lucia Zheng, Neel Guha, Brandon R Anderson, Peter Henderson, and Daniel E Ho. When does pretraining help? assessing self-supervised learning for law and the casehold dataset of 53,000+ legal holdings. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Law*, pages 159–168, 2021.

## 509 A Notation

510 The glossary is given in Table 4 below.

Symbol	Used for
$x$	Input sentence or paragraph $x \in \mathcal{X}$ .
$y$	Binary task label $y \in \mathcal{Y} = \{-1, +1\}$ .
$\mathcal{D}$	Unlabeled dataset $\mathcal{D} = \{x_i\}_{i=1}^{n_u}$ of $n_u$ points.
$\mathbb{P}, \mathbb{P}_x$	The joint distribution of $(x, y)$ and the marginal on $x$ , respectively.
$n_l$	Number of labeled in-context examples used in querying the LLM (5-10 examples).
$\lambda_{\text{LLM}}(\cdot), \phi(\cdot)$	Users interact with a language model $\lambda_{\text{LLM}}$ via a prompt $\phi$ on $x$ .
$m$	Number of prompts that we have access to.
$\lambda$	$\lambda(x) = [\lambda_1(x), \dots, \lambda_m(x)]$ where $\lambda_i(x)$ is shorthand for $\lambda_{\text{LLM}}(\phi_i(x))$ .
$\mathcal{E}$	The set of $N$ embedding models, $\mathcal{E} = \{E_1, \dots, E_N\}$ where each embedding is represented as a fixed mapping $E_i : \mathcal{X} \mapsto \mathcal{Z}_i$ .
$Z$	Partition function for normalization of (1).
$\theta_y, \theta$	$\theta_y$ is a label balance parameter and each $\theta_i$ is a scalar accuracy parameter for the $i$ th source in (1).
$\text{NN}_{j,k}(x)$	The $k$ -nearest neighbors of $x$ in embedding space $E_j$ .
$\lambda_{\text{sm}}$	$\lambda_{\text{sm}}(x) = [\lambda_{\text{sm},1}, \dots, \lambda_{\text{sm},N}] \in \{-1, 0, 1\}^{N^m}$ , where $\lambda_{\text{sm},j} = [\lambda_{\text{sm},j}[1], \dots, \lambda_{\text{sm},j}[m]]$ and $\lambda_{\text{sm},j}[i](x)$ is the smoothed neighborhood prediction of $\lambda_i(x)$ in $E_j$ (eq. (2)).
$\tau_i^+, \tau_i^-$	Shrinkage parameters for determining when $\lambda_{\text{sm},j}[i](x)$ is set to 0, -1, or 1.
$\alpha_j$	Vector of $m$ accuracy parameters for $E_j$ when used with $m$ prompts in (3).
$L(\lambda, \lambda_{\text{sm}}, \mathcal{D})$	Generalization error of EMBROID (expected cross-entropy loss).
$a_{\text{max}}$	The largest scaled accuracy of any source, $a_{\text{max}} = \max_i \mathbb{E}[\lambda_i(x)y]$ .
$b_{\text{min}}$	The smallest expected pairwise product between any two sources, $b_{\text{min}} = \min_{i,j} \{\mathbb{E}[\lambda_i \lambda_j], \hat{\mathbb{E}}[\lambda_i \lambda_j]\}$ .
$M_E(\cdot)$	An embedding is $M$ -smooth if $\Pr(\tilde{y} = c   y = c, \ E(x) - E(\tilde{x})\  \leq \varepsilon) \geq M_E(\varepsilon)$ for all $c \in \mathcal{Y}$ and any $\varepsilon > 0$ , where $M_E(\cdot) \in [0, 1]$ is decreasing in its input.
$\beta_i$	The accuracy of $\lambda_i$ , $\beta_i = \Pr(\lambda_i = y)$ .
$p_\lambda$	The prediction on $x_0$ given only access to $\lambda$ , $p_\lambda = \Pr(y = 1   \lambda(x_0))$ .
$\varepsilon_k$	The maximum distance between $x_0$ and its $k$ neighbors, $\varepsilon_k = \max_{\tilde{x} \in \text{NN}_k(x)} \ E(x) - E(\tilde{x})\ $ .

Table 4: Glossary of variables and symbols used in this paper.

## B Weak supervision background

In this section, we provide details on the inference and learning procedures for solving the graphical model defined in equation (1). The content from this section is derived from [14] and [7].

**Pseudolabel inference.** To perform inference, we compute  $\hat{\Pr}(y|\lambda(x))$  for some  $x \in \mathcal{X}$ . This is done via Bayes' rule and the conditional independence of weak sources:

$$\Pr(y = 1|\lambda(x)) = \frac{\prod_{i=1}^m \Pr(\lambda_i(x)|y = 1) \Pr(y = 1)}{\Pr(\lambda(x))}. \quad (5)$$

We assume that the class balance is known; for our datasets, the class balance is  $\Pr(y = 1) = 0.5$ . More generally, it can be estimated [51]. The latent parameter of interest in this decomposition is  $\Pr(\lambda_i = 1|y = 1)$ , which corresponds to the accuracy of  $\lambda_i$ .

---

### Algorithm 2 Triplet method [14]

---

**Input:** Dataset  $\mathcal{D}$ , weak sources  $\lambda(x)$ .

**for**  $i \in [m]$  **do**

**for**  $j, k \in [m] \setminus i$  **do**

    Estimate  $\hat{\mathbb{E}}[\lambda_i \lambda_j]$  over  $\mathcal{D}$ , and similarly estimate  $\hat{\mathbb{E}}[\lambda_i \lambda_k]$  and  $\hat{\mathbb{E}}[\lambda_j \lambda_k]$ .

    Compute  $\hat{a}_i^{j,k} = \sqrt{\frac{\hat{\mathbb{E}}[\lambda_i \lambda_j] \hat{\mathbb{E}}[\lambda_i \lambda_k]}{\hat{\mathbb{E}}[\lambda_j \lambda_k]}}$ .

**end for**

  Calculate average  $\hat{a}_i = \text{Mean}(\hat{a}_i^{j,k} \mid j, k \in [m] \setminus i)$ .

  Compute estimated accuracy  $\hat{\Pr}(\lambda_i = y) = \frac{1 + \hat{a}_i}{2}$ .

**end for**

**Output:** Accuracies  $\hat{\Pr}(\lambda_i = y)$  for all  $i \in [m]$ .

---

**Source parameter estimation: Triplet method.** Previous approaches have considered how to estimate  $\Pr(\lambda_i = 1|y = 1)$  via the *triplet method* [14], which exploits conditional independence properties. First, by the properties of the graphical model in (1), it holds that the accuracy of  $\lambda_i$  is symmetric:  $\Pr(\lambda_i = 1|y = 1) = \Pr(\lambda_i = -1|y = -1) = \Pr(\lambda_i = y)$  (Lemma 4 of [7]). Therefore,  $\Pr(\lambda_i = 1|y = 1)$  can be written in terms of  $\mathbb{E}[\lambda_i y]$  with  $\mathbb{E}[\lambda_i y] = 2 \Pr(\lambda_i = 1|y = 1) - 1$ .

Define  $a_i = \mathbb{E}[\lambda_i y]$ . The graphical model in (1) tells us that  $\lambda_i y \perp\!\!\!\perp \lambda_j y$  if  $\lambda_i \perp\!\!\!\perp \lambda_j|y$ , which holds for all  $i, j \in [m]$  (Proposition 1 of [14]). As a result,  $\mathbb{E}[\lambda_i y] \times \mathbb{E}[\lambda_j y] = \mathbb{E}[\lambda_i \lambda_j y^2] = \mathbb{E}[\lambda_i \lambda_j]$ , which is a quantity that can be computed from observed LLM predictions. That is, we have that  $a_i a_j = \mathbb{E}[\lambda_i \lambda_j]$ . If we introduce a third  $\lambda_k$ , we can generate a system of equations over  $a_i, a_j, a_k$  in terms of their pairwise rates of agreements:

$$a_i a_j = \mathbb{E}[\lambda_i \lambda_j] \quad (6)$$

$$a_i a_k = \mathbb{E}[\lambda_i \lambda_k] \quad (7)$$

$$a_j a_k = \mathbb{E}[\lambda_j \lambda_k]. \quad (8)$$

Solving, we get that

$$|a_i| := \sqrt{\frac{\mathbb{E}[\lambda_i \lambda_j] \mathbb{E}[\lambda_i \lambda_k]}{\mathbb{E}[\lambda_j \lambda_k]}}, \quad (9)$$

and likewise for  $a_j, a_k$ . If we assume that each weak source is better than random over the dataset, then  $a_i = |a_i| > 0$ , so we can uniquely recover the accuracy of each source by selecting two other sources and computing the above expression by using empirical expectations over  $\mathcal{D}$ . We then set  $\hat{\Pr}(\lambda_i = 1|y = 1) = \frac{1 + \hat{a}_i}{2}$  and plug this into the expression for  $\Pr(y = 1|\lambda(x))$  in (5).

This approach is formally described in Algorithm 2.

## 535 C Proofs

### 536 C.1 Proof of proposition 5.1

537 We note that  $[\lambda, \lambda_{\text{sm}}]$  can be viewed as a set of sources in the weak supervision set up used in [7, 14].  
 538 Therefore, we can apply Theorem 1 from [7] to our problem setting, noting that we do not perform  
 539 their clustering step and that our predictions do not abstain and output 0 in addition to  $\{-1, 1\}$ . We  
 540 have a total of  $(N + 1)m$  sources, so

$$L(\lambda, \lambda_{\text{sm}}, \mathcal{D}) \leq H(y|\lambda, \lambda_{\text{sm}}) + \frac{3(1 - b_{\min}^2)}{8b_{\min}^2(1 - a_{\max}^2)} \left( \frac{1}{b_{\min}^4} + \frac{2}{b_{\min}^2} \right) \frac{(N + 1)m}{n_u} + o(1/n_u). \quad (10)$$

### 541 C.2 Proof of theorem 5.2

542 We can write the change in point-wise irreducible error as follows:

$$H(y|\lambda(x_0)) - H(y|\lambda(x_0), \lambda_{\text{sm}}(x_0)) = \mathbb{E} [-\log \Pr(y|\lambda(x_0)) + \log \Pr(y|\lambda(x_0), \lambda_{\text{sm}}(x_0))] \quad (11)$$

$$= \mathbb{E} \left[ \log \frac{\Pr(y|\lambda(x_0), \lambda_{\text{sm}}(x_0))}{\Pr(y|\lambda(x_0))} \right] \quad (12)$$

$$= \mathbb{E} \left[ \log \left( \frac{\Pr(\lambda(x_0), \lambda_{\text{sm}}(x_0)|y) \Pr(y)}{\Pr(\lambda(x_0), \lambda_{\text{sm}}(x_0))} \cdot \frac{\Pr(\lambda(x_0))}{\Pr(\lambda(x_0)|y) \Pr(y)} \right) \right] \quad (13)$$

$$= \mathbb{E} \left[ \log \frac{\Pr(\lambda_{\text{sm}}(x_0)|\lambda(x_0), y)}{\Pr(\lambda_{\text{sm}}(x_0)|\lambda(x_0))} \right]. \quad (14)$$

543 Next, we use the fact that  $\lambda(x_0) \perp \lambda_{\text{sm}}(x_0)|y$  to simplify the expression into

$$\mathbb{E} \left[ \log \frac{\Pr(\lambda_{\text{sm}}(x_0)|y)}{\Pr(\lambda_{\text{sm}}(x_0)|y = 1) \Pr(y = 1|\lambda(x_0)) + \Pr(\lambda_{\text{sm}}(x_0)|y = -1) \Pr(y = -1|\lambda(x_0))} \right]. \quad (15)$$

544 The exact  $\lambda_{\text{sm}}(x_0)$  is unknown but is drawn from the distribution  $\Pr(\lambda_{\text{sm}}|y = 1)$  since  $x_0$ 's label is  
 545 1. Then, this expression becomes an expectation over  $\lambda_{\text{sm}}$ :

$$\mathbb{E}_{\lambda_{\text{sm}}|y=1} \left[ \log \frac{\Pr(\lambda_{\text{sm}}|y = 1)}{\Pr(\lambda_{\text{sm}}|y = 1)p_{\lambda} + \Pr(\lambda_{\text{sm}}|y = -1)(1 - p_{\lambda})} \right]. \quad (16)$$

546 Given that our  $\lambda_{\text{sm}}$  is high-quality, we suppose that  $\lambda_{\text{sm}}(x_0)$  all equal 1 with high probability, and  
 547 then we can lower bound our expression by

$$\Pr(\lambda_{\text{sm}} = 1|y = 1) \log \frac{\Pr(\lambda_{\text{sm}} = 1|y = 1)}{\Pr(\lambda_{\text{sm}} = 1|y = 1)p_{\lambda} + \Pr(\lambda_{\text{sm}} = 1|y = -1)(1 - p_{\lambda})}. \quad (17)$$

548 The key quantity of interest is  $\Pr(\lambda_{\text{sm}} = 1|y = 1) = \prod_{i=1}^m \Pr(\lambda_{\text{sm}, [i]} = 1|y = 1)$ . We focus on  
 549 bounding  $\Pr(\lambda_{\text{sm}, [i]} = 1|y = 1)$  next. Suppose that the  $k$  neighbors of  $x_0$  are  $x_1, \dots, x_k$ . Define  
 550  $p_j = \Pr(\lambda_i(x_j) = 1|y = 1)$  for all  $j \in [k]$ . Note that  $\lambda_i(x_j) \perp \lambda_i(x_{j'})|y$  for any  $j, j' \in [k]$  (while  
 551  $\lambda_{\text{sm}, [i]}$  as a whole is dependent on  $y$ , individual neighbors are still conditionally independent). Then,  
 552 the event that  $\lambda_{\text{sm}, [i]} = 1|y = 1$  is as least as likely as the event that  $\text{Binomial}(k, \min_{i \in [k]} p_i) \geq \frac{k}{2}$ .  
 553 Let  $p_{\min} = \min_{i \in [k]} p_i$ , and assume that  $p_{\min} \geq \frac{1}{2}$ . Then,

$$\Pr(\lambda_{\text{sm}, [i]} = 1|y = 1) \geq \Pr \left( \text{Binomial}(k, p_{\min}) \geq \frac{k}{2} \right) = \Pr \left( \frac{1}{k} \sum_{j=1}^k X_j \geq \frac{1}{2} \right) \quad (18)$$

$$= \Pr \left( \frac{1}{k} \sum_{j=1}^k X_j \geq p_{\min} - \left( p_{\min} - \frac{1}{2} \right) \right), \quad (19)$$

554 where  $X_j \sim \text{Bernoulli}(p_{\min})$ . Next, let  $\delta = p_{\min} - \frac{1}{2}$ . We can apply Hoeffding's inequality to get

$$\Pr\left(\text{Binomial}(k, p_{\min}) \geq \frac{k}{2}\right) = \Pr\left(\frac{1}{k} \sum_{j=1}^k X_j \geq p_{\min} - \delta\right) = 1 - \Pr\left(\frac{1}{k} \sum_{j=1}^k X_j \leq p_{\min} - \delta\right) \geq 1 - \exp(-2\delta^2 k)$$

(20)

$$= 1 - \exp(-2k(p_{\min} - 0.5)^2).$$

(21)

555 All that's left is to lower bound  $p_{\min}$ . Without loss of generality, suppose that  $p_{\min}$  corresponds to an  
556 arbitrary  $p_j = \Pr(\lambda_i(x_j) = 1|y = 1)$ . We can decompose this probability into

$$\begin{aligned} \Pr(\lambda_i(x_j) = 1|y = 1) &= \Pr(\lambda_i(x_j) = 1, y(x_j) = 1|y = 1) + \Pr(\lambda_i(x_j) = 1, y(x_j) = -1|y = 1) \\ &= \Pr(\lambda_i(x_j) = 1|y(x_j) = 1) \Pr(y(x_j) = 1|y = 1) + \Pr(\lambda_i(x_j) = 1|y(x_j) = -1) \Pr(y(x_j) = -1|y = 1). \end{aligned}$$

(22)

(23)

557 Since  $\Pr(\lambda_i(x_j) = 1|y(x_j) = 1)$  is over all  $x_j \sim \mathbb{P}_x$ , this quantity is just equal to the accuracy  
558 of  $\lambda_i, a_i$ . Next, recall that  $\|E(x_j) - E(x)\| \leq \varepsilon_k$ , where  $\varepsilon_k = \max_{x_i \in \text{NN}(x)} \|E(x) - E(x_i)\|$   
559 is the maximum distance from the  $k$  neighbors to  $x$ . Then, we can write  $\Pr(y(x_j) = 1|y = 1)$   
560 as  $\Pr(y(x_j) = 1|y = 1, \|E(x_j) - E(x)\| \leq \varepsilon_k) \geq M_E(\varepsilon_k)$ , since we have assumed that  $E$  is  
561  $M$ -smooth. We can now bound  $p_{\min}$ :

$$p_{\min} \geq a_i M_E(\varepsilon_k) + (1 - a_i)(1 - M_E(\varepsilon_k)).$$

(24)

562 Therefore, we have that

$$\Pr(\lambda_{\text{sm}} = 1|y = 1) \geq \prod_{i=1}^m [1 - \exp[-2k(a_i M_E(\varepsilon_k) - 0.5)^2]].$$

(25)

563 Before we plug in  $\Pr(\lambda_{\text{sm}} = 1|y = 1)$  into (17), we simplify the expression. Note that  $\Pr(\lambda_{\text{sm}} =$   
564  $1|y = 1)$  can be written as  $\prod_{i=1}^m p_i$  for some  $p_i$ , and  $\Pr(\lambda_{\text{sm}} = 1|y = -1)$  can be written as  
565  $\prod_{i=1}^m (1 - p_i)$ . A simple proof by induction shows that  $\prod_{i=1}^m (1 - p_i) \leq 1 - \prod_{i=1}^m p_i$ . Therefore, we  
566 can write that (17) is lower bounded by

$$\Pr(\lambda_{\text{sm}} = 1|y = 1) \log \frac{\Pr(\lambda_{\text{sm}} = 1|y = 1)}{\Pr(\lambda_{\text{sm}} = 1|y = 1)p_{\lambda} + (1 - \Pr(\lambda_{\text{sm}} = 1|y = 1))(1 - p_{\lambda})}$$

(26)

567 Let's abbreviate  $\Pr(\lambda_{\text{sm}} = 1|y = 1)$  as  $x$  and define the function

$$f(x) = x \log \frac{x}{xp_{\lambda} + (1 - x)(1 - p_{\lambda})}.$$

(27)

568 We note that for  $x \geq 0.5$ ,  $f(x)$  is convex and can thus be lower bounded by  $f(x) \geq f'(0.5)(x - 0.5)$ .  
569 We compute  $f'(x) = \frac{1 - p_{\lambda}}{xp_{\lambda} + (1 - x)(1 - p_{\lambda})}$ , so  $f'(0.5) = 2(1 - p_{\lambda})$ . Therefore,  $f(x) \geq 2(1 - p_{\lambda})(x -$   
570  $0.5)$ . Our final bound on the pointwise difference in irreducible error on  $x_0$  is

$$H(y|\lambda(x_0)) - H(y|\lambda(x_0), \lambda_{\text{sm}}(x_0)) \geq 2(1 - p_{\lambda}) \left[ \prod_{i=1}^m [1 - \exp[-2k(a_i M_E(\varepsilon_k) - 0.5)^2]] - 0.5 \right].$$

(28)

571 **Information gain from using  $\lambda, \lambda_{\text{sm}}$  over  $\lambda_{\text{sm}}$**  We briefly comment on the opposite direction—  
572 how much does using both LLM predictions and neighborhood predictions help over just using  
573 neighborhood predictions?

574 The quantity we aim to lower bound is  $H(y|\lambda_{\text{sm}}(x_0)) - H(y|\lambda(x_0), \lambda_{\text{sm}}(x_0))$  for a point of interest  
575  $x_0$ . We can write this quantity as

$$H(y|\lambda_{\text{sm}}(x_0)) - H(y|\lambda(x_0), \lambda_{\text{sm}}(x_0)) = \mathbb{E} \left[ \log \frac{\Pr(\lambda(x_0)|y)}{\Pr(\lambda(x_0)|\lambda_{\text{sm}}(x_0))} \right]$$

(29)

576 Without loss of generality, suppose that the true label on  $x_0$  is  $y = 1$ , and that for each  $\lambda_i$ , the  
 577 neighborhood around  $x_0$  consists of a balanced mix of  $\lambda_i = 1$  and  $\lambda_i = -1$ . Then, with high  
 578 probability we have that  $\Pr(y = 1|\boldsymbol{\lambda}_{\text{sm}}(x_0)) = p_{\boldsymbol{\lambda}_{\text{sm}}} \approx 0.5$ . From our proof of Theorem 5.2, we can  
 579 thus write

$$H(y|\boldsymbol{\lambda}_{\text{sm}}(x_0)) - H(y|\boldsymbol{\lambda}(x_0), \boldsymbol{\lambda}_{\text{sm}}(x_0)) = \mathbb{E}_{\lambda(x_0)} \left[ \log \frac{\Pr(\boldsymbol{\lambda}(x_0)|y = 1)}{\Pr(\lambda(x_0)|y = 1)p_{\boldsymbol{\lambda}_{\text{sm}}} + \Pr(\lambda(x_0)|y = -1)(1 - p_{\boldsymbol{\lambda}_{\text{sm}}})} \right] \quad (30)$$

$$\geq \Pr(\boldsymbol{\lambda}(x_0) = 1|y = 1) \log \frac{\Pr(\boldsymbol{\lambda}(x_0) = 1|y = 1)}{\Pr(\boldsymbol{\lambda}(x_0) = 1|y = 1)p_{\boldsymbol{\lambda}_{\text{sm}}} + \Pr(\boldsymbol{\lambda}(x_0) = 1|y = -1)(1 - p_{\boldsymbol{\lambda}_{\text{sm}}})} \quad (31)$$

$$\geq 2(1 - p_{\boldsymbol{\lambda}_{\text{sm}}})(\Pr(\boldsymbol{\lambda}(x_0) = 1|y = 1) - 0.5) \quad (32)$$

580 If  $\boldsymbol{\lambda}$  on  $x_0$  has high accuracy and  $p_{\boldsymbol{\lambda}_{\text{sm}}}$  is low, then we can have significant point-wise information  
 581 gain from modeling both  $\boldsymbol{\lambda}$  and  $\boldsymbol{\lambda}_{\text{sm}}$  rather than just  $\boldsymbol{\lambda}_{\text{sm}}$ .

## D Datasets

**Motivation.** We study the performance of our method across a diverse collection of *task definitions*. In our setting, a task definition denotes a specific classification that a data scientist wishes to perform. For instance, a data scientist working on quantifying the breadth of legal issues that individuals face may wish to identify which posts in an online forum refer implicate legal issues related to housing.

This evaluation strategy is motivated by the observation that task definitions vary in their smoothness across embedding spaces, as different embeddings may do a better job of capturing features relevant for the task. For instance, out-of-the-box Sentence-BERT embeddings are better than traditional BERT at capturing the topicality of a sentence [53]. By focusing on a broad range of task definitions, we can better forecast how our method might perform for new tasks that practitioners may need to create classifiers for. We also avoid issues with leakage that may arise as the practice of finetuning LLMs on tasks increases [9].

In total, we study 95 distinct task definitions, encompassing 100,418 total samples. Each task varies between 108 and 3308 samples.

**Legal tasks.** The emergence of LLMs is exciting for law and finance, where expert-annotations are especially difficult to acquire [4, 18]. Drawing on recent benchmarks and released datasets framing the potential use cases for LLMs in law, we study the following datasets:

- CUAD [24]: The original CUAD dataset consists of 500 contracts spanning an array of sectors, with clauses manually into one of 41 legal categories. Following [18], we adapt the original dataset for clause-by-clause classification. We turn each clause type into a binary classification task, where the objective is to distinguish clauses of that type from clauses of other types (i.e. “negatives”). Negative clauses are sampled randomly so as to make the task class balanced. We ignore clauses for which there are insufficient annotations in the original dataset.
- Learned Hands [32]: The Learned Hands dataset consists of legal questions that individuals publicly posted to an online forum (r/legaladvice). The questions have been coded by experts into legal categories according to the Legal Issues Taxonomy [31]. We consider several such issue classes, and create a binary classification task for each issue. Negative clauses are sampled randomly so as to make the task class balanced. Because these questions can be long, we truncate them at 50 tokens.

**Science tasks.** LLMs have generated similar excitement for medical and science informatics applications [1, 4]. We study established classification/extraction benchmarks.

- Chemprot [29]: ChemProt consists of sentences from PubMed abstracts describing chemical-protein relationships. We study seven relations, and create a binary task for each one. Each task is class balanced, with negatives sampled from the other relations.
- RCT [11]: The RCT dataset consists of PubMed abstracts for papers discussing randomized control trials, where sentences in the abstract are annotated according to their semantic role (e.g., background, methods, results, etc). There are five roles, and we create a binary task for each one. Each task is class balanced, with negatives sampled from the other relations.

**General domain tasks.** Finally, we study the following “general domain” tasks, derived from popular sentence classification and information extraction benchmarks.

- FewRel [20]: This is a relationship classification/extraction dataset, where each sample corresponds to a sentence mentioning the relationship between two entities. We select 20 relations, and for each relation construct a binary classification task with 700 positive instances of the relation, and 700 randomly sampled sentences (corresponding to other relations).
- Spam Detection [63]: We study the YouTube spam detection task from the WRENCH benchmark. This task requires classifying YouTube comments as spam/not spam.
- Toxic content detection [5]: This task requires classifying whether posted comments are toxic or not. We use a sampled subset of the CivilComments dataset.

- AG News [65]: The original dataset organizes news snippets into four categories: World, Sports, Business, and Science/Technology. We create a separate task for each category. Negatives are sampled from the remaining classes.
- DBPedia [65]: DBPedia is a 14-way ontology classification dataset. We convert this into 14 distinct tasks, corresponding to each of the ontology types.

Table 5: Legal tasks

Task	Description/Intent	Size
Affiliate License-Licensee (CUAD)	Does the clause describe a license grant to a licensee (incl. sublicensor) and the affiliates of such licensee/sublicensor?	208
Anti-Assignment (CUAD)	Does the clause require consent or notice of a party if the contract is assigned to a third party?	1212
Audit Rights (CUAD)	Does the clause discuss potential audits?	1224
Cap On Liability (CUAD)	Does the clause specify a cap on liability upon the breach of a party's obligation?	1262
Change Of Control (CUAD)	Does the clause give one party the right to terminate if such party undergoes a change of control?	426
Competitive Restriction Exception (CUAD)	Does the clause mention exceptions or carveouts to Non-Compete, Exclusivity and No-Solicit of Customers?	226
Covenant Not To Sue (CUAD)	Does the clause mention if a party is restricted from contesting the validity of the counterparty's ownership of intellectual property?	318
Effective Date (CUAD)	Does the clause mention when the contract becomes effective?	246
Exclusivity (CUAD)	Does the clause mention an exclusive dealing commitment with the counterparty?	770
Expiration Date (CUAD)	Does the clause mentions a date when the contract's term expires?	892
Governing Law (CUAD)	Does the clause mentions which state/country's laws govern interpretation of the contract?	910
Insurance (CUAD)	Does the clause mention a requirement for insurance?	1040
Ip Ownership Assignment (CUAD)	Does the clause mention if intellectual property created by one party become the property of the counterparty?	590
Irrevocable Or Perpetual License (CUAD)	Does the clause describe a license grant that is irrevocable or perpetual?	300
Joint Ip Ownership (CUAD)	Does the clause provide for joint or shared ownership of intellectual property between the parties to the contract?	198
License Grant (CUAD)	Does the clause describe a license granted by one party to its counterparty?	1430
Liquidated Damages (CUAD)	Does the clause award either party liquidated damages for breach or a fee upon the termination of a contract (termination fee)?	226
Minimum Commitment (CUAD)	Does the clause specifies a minimum order size or minimum amount or units per time period that one party must buy from the counterparty?	778

Table 5 – continued from previous page

<b>Task</b>	<b>Description/Intent</b>	<b>Size</b>
No-Solicit Of Employees (CUAD)	Does the clause restricts a party's soliciting or hiring employees and/or contractors from the counterparty, whether during the contract or after the contract ends (or both).	150
Non-Compete (CUAD)	Does the clause restrict the ability of a party to compete with the counterparty or operate in a certain geography or business or technology sector?	450
Non-Disparagement (CUAD)	Does the clause require a party not to disparage the counterparty?	108
Non-Transferable License (CUAD)	Does the clause limit the ability of a party to transfer the license being granted to a third party?	558
Notice Period To Terminate Renewal (CUAD)	Does the clause requires a notice period to terminate renewal?	234
Post-Termination Services (CUAD)	Does the clause subject a party to obligations after the termination or expiration of a contract, including any post-termination transition, payment, transfer of IP, wind-down, last-buy, or similar commitments?	816
Renewal Term (CUAD)	Does the clause mention a renewal term for after the initial term expires?	398
Revenue-Profit Sharing (CUAD)	Does the clause require a party to share revenue or profit with the counterparty for any technology, goods, or services?	784
Rofr-Rofo-Rofn (CUAD)	Does the clause provide a party with a right of first refusal?	698
Source Code Escrow (CUAD)	Does the clause requires one party to deposit its source code into escrow with a third party, which can be released to the counterparty upon the occurrence of certain events (bankruptcy, insolvency, etc.)?	126
Termination For Convenience (CUAD)	Does the clause state that one party can terminate this contract without cause (solely by giving a notice and allowing a waiting period to expire)?	442
Uncapped Liability (CUAD)	Does the clause state that a party's liability is uncapped upon the breach of its obligation in the contract?	302
Volume Restriction (CUAD)	Does the clause describe a fee increase or consent requirement if one party's use of the product/services exceeds certain threshold?	328
Warranty Duration (CUAD)	Does the clause mentions the duration of any warranty against defects or errors in technology, products, or services provided under the contract?	326
BU (Learned Hands)	Does the text discuss issues relating to business or intellectual property?	200
CO (Learned Hands)	Does the text discuss issues relating to courts and lawyers?	194
CR (Learned Hands)	Does the text discuss issues relating to criminal issues?	644
ES (Learned Hands)	Does the text discuss issues relating to estates or wills?	182
FA (Learned Hands)	Does the text discuss issues relating to family or divorce?	794

Table 5 – continued from previous page

<b>Task</b>	<b>Description/Intent</b>	<b>Size</b>
HE (Learned Hands)	Does the text discuss issues relating to health?	248
HO (Learned Hands)	Does the text discuss issues relating to housing?	1270
MO (Learned Hands)	Does the text discuss issues relating to payments or debt?	740
TO (Learned Hands)	Does the text discuss issues relating to accidents or harassment?	454
TR (Learned Hands)	Does the text discuss issues relating to cars or traffic?	516
WO (Learned Hands)	Does the text discuss issues relating to employment or job?	726

Table 6: Science tasks

<b>Task</b>	<b>Description/Intent</b>	<b>Size</b>
Agonist (Chemprot)	Does the sentence describe an agonist relationship?	896
Antagonist (Chemprot)	Does the sentence describe an antagonist relationship?	1330
Downregulator (Chemprot)	Does the sentence describe a downregulator relationship?	3038
Part_of (Chemprot)	Does the sentence describe a part-of relationship?	1210
Regulator (Chemprot)	Does the sentence describe a regulator relationship?	2876
Substrate (Chemprot)	Does the sentence describe a substrate relationship?	2384
Upregulator (Chemprot)	Does the sentence describe an upregulator relationship?	2404
Background (RCT)	Does the sentence describe background on the study?	2000
Conclusions (RCT)	Does the sentence state a conclusion?	2000
Methods (RCT)	Does the sentence describe a scientific experimental method?	2000
Objective (RCT)	Does the sentence describe the goal of the study?	2000
Results (RCT)	Does the sentence describe experimental results?	2000

Table 7: General domain tasks

<b>Task</b>	<b>Description/Intent</b>	<b>Size</b>
Business (AGNews)	Does the article discuss business news?	2000
Sports (AGNews)	Does the article discuss sports news?	2000
Technology (AGNews)	Does the article discuss technology news?	2000
World (AGNews)	Does the article discuss global affairs or world events?	2000
Civil Comments	Does the sentence contain toxic or hateful content?	2500
Album (DBPedia)	Is the entity discussed in the sentence an example of a album?	2000
Animal (DBPedia)	Is the entity discussed in the sentence an example of a animal?	2000
Artist (DBPedia)	Is the entity discussed in the sentence an example of a artist?	2000
Athlete (DBPedia)	Is the entity discussed in the sentence an example of a athlete?	2000
Building (DBPedia)	Is the entity discussed in the sentence an example of a building?	2000
Company (DBPedia)	Is the entity discussed in the sentence an example of a company?	2000
Educational institution (DBPedia)	Does the sentence discuss a school, university, or college?	2000

Table 7 – continued from previous page

<b>Task</b>	<b>Description/Intent</b>	<b>Size</b>
Film (DBPedia)	Is the entity discussed in the sentence an example of a film?	2000
Mean of transportation (DBPedia)	Does the sentence discuss a car, ship, train, or plane?	2000
Natural place (DBPedia)	Is the entity discussed in the sentence an example of a natural landscape or environment?	2000
Office holder (DBPedia)	Is the entity discussed in the sentence an example of a office holder?	2000
Plant (DBPedia)	Is the entity discussed in the sentence an example of a plant?	2000
Village (DBPedia)	Is the entity discussed in the sentence an example of a village?	2000
Written work (DBPedia)	Is the entity discussed in the sentence a writing?	2000
Architect (FewRel)	Does the text mention an architect?	600
Composer (FewRel)	Does the text mention a musical composer?	600
Country (FewRel)	Does the text mention a country?	600
Developer (FewRel)	Does the text mention development?	600
Director (FewRel)	Does the text mention a film director?	600
Distributor (FewRel)	Does the text mention a film distributor?	600
Father (FewRel)	Does the text mention a father?	600
Genre (FewRel)	Does the text mention the genre of a song or artist?	600
Instrument (FewRel)	Does the text mention an instrument?	600
League (FewRel)	Does the text mention a sports competition, league, or division?	600
Military Branch (FewRel)	Does the text mention a military branch?	600
Movement (FewRel)	Does the text mention an art movement?	600
Occupation (FewRel)	Does the text mention a professional occupation?	600
Participating Team (FewRel)	Does the text mention a sports team?	600
Platform (FewRel)	Does the text mention an online platform?	600
Sibling (FewRel)	Does the text mention a sibling?	600
Successful Candidate (FewRel)	Does the text mention an election winner?	600
Taxonomy Rank (FewRel)	Does the text mention a taxonomy class of animals?	600
Tributary (FewRel)	Does the text mention a tributary?	600
Winner (FewRel)	Does the text mention a competition winner?	600
YouTube Spam Detection	Does the comment ask the user to check out another video?	1836

## 637 E Prompts

638 We construct a prompt for each task by randomly selecting three examples of each class to use as  
639 in-context demonstrations [49]. We manually defined “instructions” for each task. An example of an  
abridged prompt is shown in Figure 3.

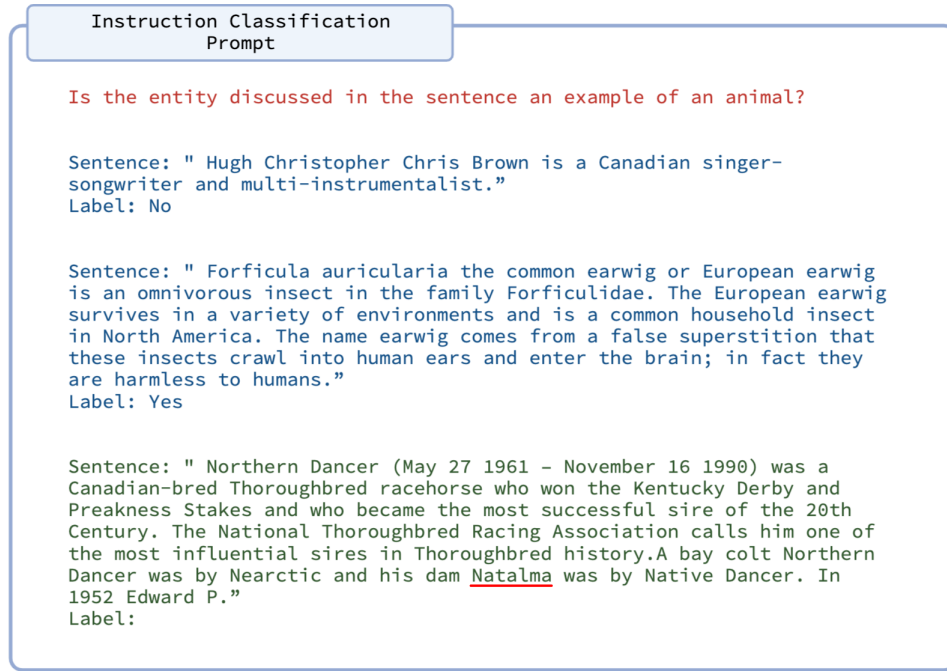


Figure 3: An example of the instruction classification prompt for the dbpedia\_animal task, with two in-context demonstrations. Here, the task instructions are in red, the in-context demonstrations are in blue, and the sample for which we want a label is in green.

640

## F Synthetics

We conduct synthetic experiments which provide additional insights on EMBROID.

For our setup, we create two equal clusters of data of 500 points each,  $C_1$  and  $C_2$ , in  $\mathbb{R}^2$ . We assign labels to the points in each cluster i.i.d. according to a probability  $p$ , where  $\Pr(y = 1|x \in C_1) = p$  and  $\Pr(y = 1|x \in C_2) = 1 - p$ . When  $p = 0.5$ , both the clusters have a uniform label distribution (non-smooth) while  $p = 1$  ensures each cluster has one class (smooth). We set  $k = 20$ ,  $\tau_+ = P(\lambda_i = 1)$  and  $\tau_- = P(\lambda_i = -1)$ .

**Improvement over weak supervision.** We show that EMBROID offers improvement over methods that only use  $\lambda$ . We fix  $p = 0.8$  and  $\beta_i = \Pr(\lambda_i = y) = 0.6$  for each  $i \in [m]$ . We compare EMBROID against the standard weak supervision approach from [14], which requires  $m \geq 3$ , in Figure 4a.

**Smoothness.** EMBROID’s performance depends on the smoothness of the embedding as defined in eq. (4). We consider one LM prediction  $m = 1$  and vary the smoothness  $p$  from 0.5 to 1.0 and generate predictions using  $\beta_i = 0.6$ . Figure 4b exhibits that EMBROID’s accuracy is positively correlated with the embedding smoothness.

**Base prediction accuracy.** Finally, we show that EMBROID’s performance depends on the base prediction accuracy,  $\beta_i$ . We consider  $m = 1$  and set  $p = 0$ ,  $\beta_i = 0.8$ . We use a parameter  $\rho$  to denote the probability that  $\lambda_i$  is incorrect on points in cluster  $C_1$ . As  $\rho$  varies from 0 to 1, the predictions of  $\lambda_i$  become biased towards 1 and effectively reduces  $\beta_i$  to 0.5. In Figure 4c, we observe that as the base prediction accuracy decreases, EMBROID’s performance decreases and eventually goes below the base LM performance.

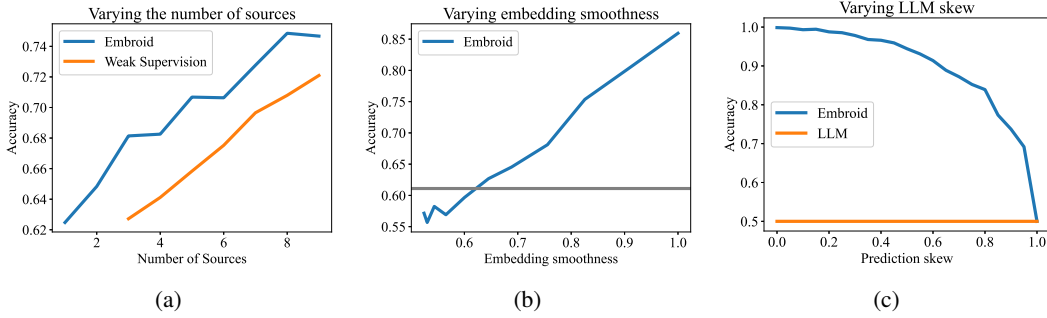


Figure 4: Synthetic experiments. (a) Comparison of EMBROID to weak supervision by varying number of LLM sources. Increasing sources consistently improves both EMBROID and WS and the gains remain constant. (b) EMBROID’s performance as embedding smoothness (eq. (4)) varies. EMBROID’s accuracy linearly improves as a function of embedding smoothness. (c) EMBROID’s performance with varying probability of LLM being incorrect in  $C_1$ . As the LM becomes incorrect, the cluster becomes less homogeneous and this degrades EMBROID’s performance.

## G Ablation

We perform additional ablations of EMBROID. We focus on two aspects:

- The role of  $\tau^-/\tau^+$ .
- The impact of weak supervision in combining  $\lambda_{sm,j}$ .

### G.1 Ablation over $\tau$

In our experiments, we set  $\tau_i^+ = \tau_i^- = \mathbb{E}[\lambda_i]$ , or the average vote of source  $\lambda_i$ . This has the following effect on the neighborhood vote  $\lambda_{sm,j}[i]$  for source  $\lambda_i$  under  $E_j$ :

- When the average vote for a source  $\lambda_i$  in a neighborhood under  $E_j$  for  $x$  is more negative than the average overall vote for a source, then  $\lambda_{sm,j}[i](x) = -1$ .

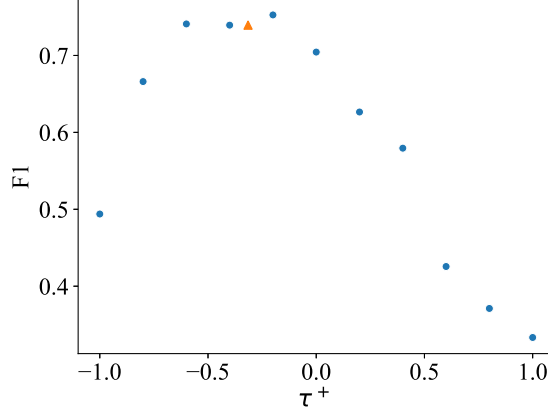


Figure 5: We analyze how F1 changes for different settings of  $\tau_i^+/\tau_i^-$  for GPT-JT on a task. Observe that setting  $\tau_i^+ = \tau_i^- = \mathbb{E}[\lambda_i]$  (denoted as the orange triangle) produces close-to-optimal performance.

LM	Base prompt	Majority vote aggregation	EMBROID
j1-jumbo	0.498	0.569	0.604
openai_text-davinci-003	0.806	0.844	0.855
bloom-7b1	54.7	61.2	64.7
opt-6.7b	48.2	56.1	59.8
GPT-JT-6B-v1	67.8	73.2	75.1

Table 8: We evaluate how EMBROID compares to a majority vote aggregation over neighborhood vote vectors. We report macro-F1 average across all tasks and prompts, mirroring results in Table 1.

- When the average vote for a source  $\lambda_i$  in a neighborhood under  $E_j$  for  $x$  is more positive than the average overall vote for a source, then  $\lambda_{\text{sm},j}[i](x) = 1$ .

In general, we find that this setting provides good performance, while requiring no additional tuning or validation. For example, the Figure 5 below compares a setting of  $\tau_i^+/\tau_i^-$  against the F1 score for GPT-JT on ag\_news\_business.

## G.2 Role of weak supervision aggregation

We quantify the extent to which performance gains are derived from (1) the computation of  $\lambda_{\text{sm}}$ , as opposed to (2) the use weak-supervision [14] to combine  $\lambda_{\text{sm}}$  and  $\lambda$ . Specifically, we replace Equation 3 with a simple majority vote classifier which combines the original prediction with the computed neighborhood votes.

## H Experiments

### H.1 Implementation details

**Compute.** Inference for API-access models (e.g., GPT-3.5 and J1-Jumbo) were run using the HELM API [37]. Inference for open source models (OPT, GPT-JT, and Bloom) were run using the Manifest library [47] on 40GB A100 NVIDIA GPU machines.

**Hyperparameters.** EMBROID was run with  $k = 10$ ,  $\tau_i^+ = P(\lambda_i = 1)$ , and  $\tau_i^- = P(\lambda_i = -1)$

**API-model tasks.** Due to cost constraints, we study the API access models (GPT-3.5 and J1-Jumbo) on a subset of tasks. These are:

- ag\_news\_world
- ag\_news\_sports
- dbpedia\_educational\_institution
- dbpedia\_athletechemprot\_regulator
- chemprot\_upregulator
- rct\_objective
- rct\_methods
- CUAD\_Audit\_Rights
- CUAD\_Non-Compete
- learned\_hands\_HE
- learned\_hands\_HO
- few\_rel\_architect
- few\_rel\_league

### H.2 Robustness across models

We provide the results for each LM on each task as CSV files in the supplemental attachment. We visualize the improvements in Table 1 below, by plotting the original prompt performance on the x-axis, and the performance after EMBROID on the y-axis (Figure 6).

### H.3 Comparison to AMA

We visualize the improvements in Table 2 below, by plotting the performance of AMA on the x-axis, and the performance of EMBROID-3 on the y-axis (Figure 7).

### H.4 Chain-of-thought experiments

We compare EMBROID to chain-of-thought (CoT) [60] on the following tasks:

- ag\_news\_business
- ag\_news\_sports
- CUAD\_Affiliate License-Licensee
- CUAD\_Audit Rights
- dbpedia\_album
- dbpedia\_building
- few\_rel\_architect
- few\_rel\_country
- learned\_hands\_HO
- learned\_hands\_MO

We manually write logical chains for each prompt. Because CoT primarily succeeds on “large” models [60], we focus our experiments on GPT-3.5.

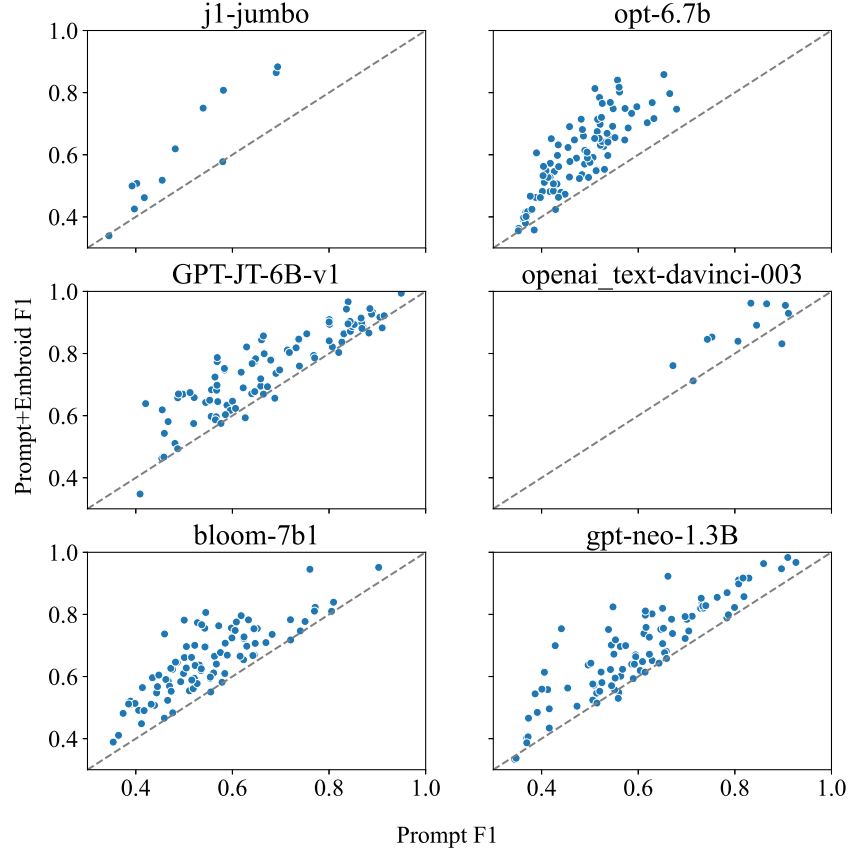


Figure 6: We visualize the improvement from EMBROID over the base prompt for each model’s tasks. All models except for gpt-neo-1.3B were run thrice per task, with each run using different in-context samples for the prompt. Each dot corresponds to a task. The x-axis measures the average macro F1 of the base prompt, and the y-axis measures the average macro F1 of EMBROID (across all runs). Because GPT-3.5 and J1-Jumbo are studied on a subset of 12 tasks, there are fewer dots in the plot.

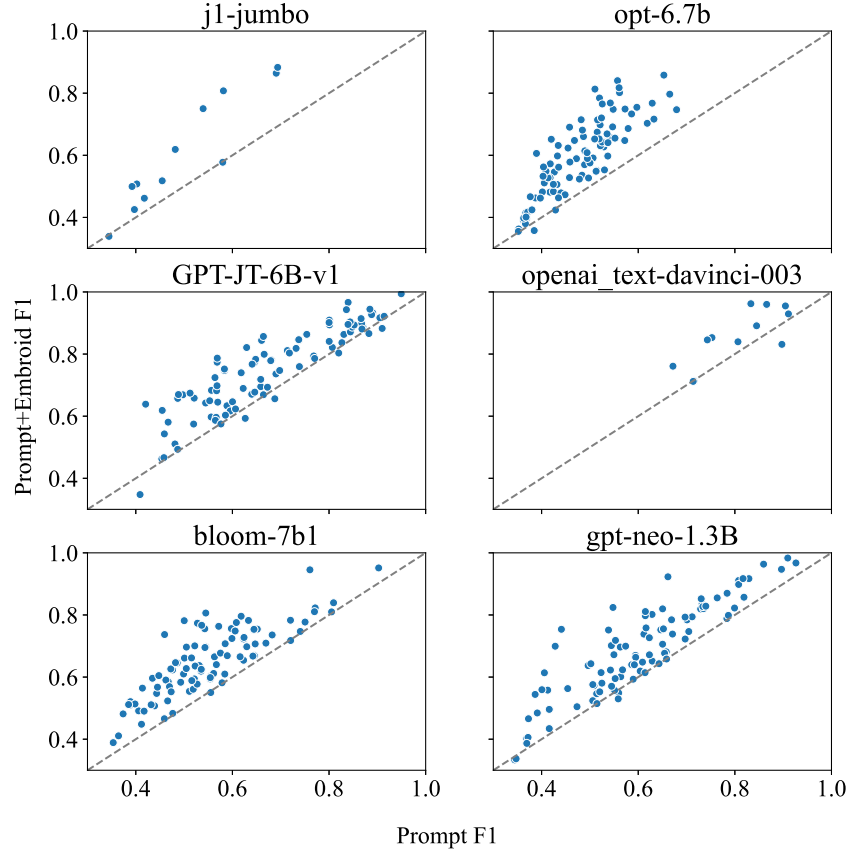


Figure 7: We visualize the improvement from EMBROID over AMA [2] for each model’s tasks. Each dot corresponds to a task. The x-axis measures the average macro F1 of AMA, and the y-axis measures the average macro F1 of EMBROID. Because GPT-3.5 and J1-Jumbo are studied on a subset of 12 tasks, there are fewer dots in the plot.