



Refining Labeling Functions with Limited Labeled Data

Chenjie Li
University of Illinois Chicago
Chicago, IL, USA
cl206@uic.edu

Amir Gilad
Hebrew University
Jerusalem, Israel
amirg@cs.huji.ac.il

Boris Glavic
University of Illinois Chicago
Chicago, IL, USA
bglavic@uic.edu

Zhengjie Miao
Simon Fraser University
Burnaby, BC, Canada
zhengjie@sfu.ca

Sudeepa Roy
Duke University
Durham, NC, USA
sudeepa@cs.duke.edu

Abstract

Programmatic weak supervision (PWS) significantly reduces human effort for labeling data by combining the outputs of user-provided labeling functions (LFs) on unlabeled datapoints. However, the quality of the generated labels depends directly on the accuracy of the LFs. In this work, we study the problem of fixing LFs based on a small set of labeled examples. Towards this goal, we develop novel techniques for repairing a set of LFs by minimally changing their results on the labeled examples such that the fixed LFs ensure that (i) there is sufficient evidence for the correct label of each labeled datapoint and (ii) the accuracy of each repaired LF is sufficiently high. We model LFs as conditional rules, which enables us to refine them, i.e., to selectively change their output for some inputs. We demonstrate experimentally that our system improves the quality of LFs based on surprisingly small sets of labeled datapoints.

CCS Concepts

• Information systems → Data cleaning.

Keywords

weak supervision; labeling functions; label repair; rule refinement; label quality

ACM Reference Format:

Chenjie Li, Amir Gilad, Boris Glavic, Zhengjie Miao, and Sudeepa Roy. 2025. Refining Labeling Functions with Limited Labeled Data. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V.2 (KDD '25)*, August 3–7, 2025, Toronto, ON, Canada. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3711896.3737102>

KDD Availability Link:

The source code of this paper has been made publicly available at <https://doi.org/10.5281/zenodo.15558280>.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
KDD '25, Toronto, ON, Canada

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 979-8-4007-1454-2/2025/08
<https://doi.org/10.1145/3711896.3737102>

```
def key_word_star(v): #LF-1
    words = ['star', 'stars']
    return POS if words.intersection(v) else ABSTAIN

def key_word_waste(v): #LF-2
    return NEG if 'waste' in v else ABSTAIN

def key_word_poor(v): #LF-3
    words = ['poorly', 'useless', 'horrible', 'money']
    return NEG if words.intersection(v) else ABSTAIN
```

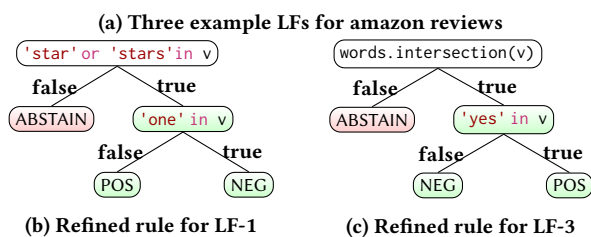


Figure 1: LFs before / after refinement by RULECLEANER

1 Introduction

Programmatic weak supervision (PWS) [25, 34] is a powerful technique for creating training data. Unlike manual labeling, where labels are painstakingly assigned by hand to each training datapoint, data programming assigns labels by combining the outputs of labeling functions (LFs) — heuristics that take a datapoint as input and output a label — using a model. This approach dramatically reduces the human effort required to label data. To push this reduction even further, recent approaches automate the generation of LFs [4, 7, 11, 30]. For example, Witan [7] creates LFs from simple predicates that are effective in differentiating datapoints, subsequently guiding users to select and refine sensible LFs. Guan et al. [11] employ large language models (LLMs) to derive LFs based on a small amount of labeled data, further reducing the dependency on human intervention. One advantage of PWS over weak supervision with a black box model is that LFs are inherently interpretable.

Regardless of whether LFs are manually crafted by domain experts or generated by automated techniques, users face significant challenges when it comes to repairing these LFs to correct issues with the resulting labeled data. The black-box nature of the model that combines LFs results obscures which specific LFs are responsible for mislabeling a datapoint, and large training datasets make it difficult for users to manually identify effective repairs. While explanation techniques for PWS [12, 35] can identify LFs responsible

id	text	true label	old predicted label by Snorkel	new predicted label by Snorkel	LF labels: old (new)		
					1	2	3
0	five stars. product works fine	P	P	P	P	-	-
1	one star. rather poorly written needs more content and an editor	N	P	N	P (N)	-	N
2	five stars. awesome for the price lightweight and sturdy	P	P	P	P	-	-
3	one star. not my subject of interest, too dark	N	P	N	P (N)	-	-
4	yes, get it! the best money on a pool that we have ever spent. really cute and holds up well with kids constantly playing in it	P	N	P	-	-	N (P)

Table 1: Products reviews with ground truth labels ("P" ositive or "N" egative), predicted labels by Snorkel [25] (before and after rule repair), and the results of the LFs from Figure 1 ("-" means ABSTAIN). Results for repaired rules are shown in blue.

for erroneous labels, determining how to repair the LFs to fix these errors remains a significant challenge.

In this work, we tackle the challenge of automatically suggesting repairs for a set of LFs based on a small set of labeled datapoints. Our approach refines an LF by locally overriding its outputs to align with expectations for specific datapoints. Rather than replacing human domain expertise or existing automated LF generation, our method, RULECLEANER, improves an existing set of LFs. Our approach is versatile, supporting arbitrarily complex LFs, and various black box models that combine them such as Snorkel [25], FlyingSquid [8] or simpler models like majority voting. RULECLEANER is agnostic to the source of LFs, enabling the repair of LFs generated by tools like Witan [7], LLMs [11], and those created by domain experts.

To address the challenge of refining LFs expressed in a general-purpose programming language, we model LFs as *rules*, represented as trees. In these trees, inner nodes are *predicates*, Boolean conditions evaluated on datapoints, and leaf nodes correspond to labels. Such a tree encodes a cascading series of conditions, starting at the root, each predicate directs navigation to a *true* or *false* child until a leaf node is reached, which assigns the label to the input datapoint. This model can represent any LF as a rule by creating predicates that match the result of the LF against every possible label (see [18]).

EXAMPLE 1. Consider the Amazon Review Dataset from [7, 14] which contains reviews for products bought from Amazon and the task of labeling the reviews as POS or P (positive), or NEG or N (negative). A subset of LFs generated by the Witan system [7] for this task are shown in Figure 1a. `key_word_star` labels reviews as POS that contain either "star" or "stars" and otherwise returns ABSTAIN (the function cannot make a prediction). Some reviews with their ground truth labels (unknown to the user) and the labels predicted by Snorkel [26] are shown in Table 1, which also shows the three LFs from Figure 1a. Reviews 1, 3, and 4 are mispredicted by the model trained by Snorkel over the LF outputs. Our goal is to reduce such misclassifications by refining the LFs. We treat systems like Snorkel as a blackbox that can use any algorithm or model to generate labels.

Suppose that for a small set of reviews, the true label is known (Table 1). RULECLEANER uses these ground truth labels to generate a set of repairs for the LFs by refining LFs to align them with the ground truth. Table 1 also shows the labels produced by the repaired LFs (updated labels shown in blue), and the predictions generated by Snorkel before and after LF repair. RULECLEANER repairs LF-1 and LF-3 from Figure 1a by adding new predicates (refinement). Figure 1b and 1c show the refined rules in tree form with new nodes highlighted in green. Consider the repair for LF-1. Intuitively, this repair is sensible: a review mentioning "one" and "star(s)" is likely negative.

Our RULECLEANER system produces repairs as shown in the example above. We make the following contributions.

- **The PWS Repair Problem.** We introduce a general model for PWS as programmatic weak supervision systems (PWSSs), where interpretable LFs (rules) are combined to predict labels for a dataset \mathcal{X} (Section 2). We formalize the problem of repairing LFs in PWSS through refinement, proving the problem to be NP-hard. To avoid overfitting, we (i) minimize the changes to the outputs of the original LFs and (ii) allow some LFs to return incorrect labels for some labeled datapoints.
- **Efficient Rule Repair Algorithm.** In spite of its hardness, in practice it is feasible to solve the repair problem exactly as the number of labeled examples is typically small. We formalize this problem as a mixed-integer linear program (MILP) that determines changes to the LF output on the labeled examples (Section 3). To implement these changes, we refine individual rules to match the desired outputs (Section 4). To further decrease the likelihood of overfitting and limit the complexity of the fixed rules, we want to minimize the number of new predicates that are added. This problem is also NP-hard. We propose a PTIME information-theoretic heuristic algorithm (Section 4.2).
- **Comprehensive Experimental Evaluation.** We conduct experiments on 11 real datasets using Snorkel [25] over LFs generated by Witan [7] or LLMs [11] (Section 5). Furthermore, we compare against using LLMs directly for labeling and for repairing LFs. RULECLEANER significantly improves labeling accuracy using a small number of labeled examples. While direct labeling with LLMs achieves impressive accuracy for advanced models like GPT-4o, it is also prohibitively expensive.

2 The RULECLEANER Framework

As shown in Figure 2, we assume as input a set of LFs modeled as rules \mathcal{R} , the corresponding labels produced by an PWSS for an unlabeled dataset \mathcal{X} , and a small subset of labeled datapoints $\mathcal{X}^* \subset \mathcal{X}$. RULECLEANER refines specific LFs based on this input, generating an updated set of rules \mathcal{R}^* . Finally, the PWSS applies these revised rules to re-label the dataset.

2.1 Rules and PWSSs

To be able to repair a LF by selectively overriding its output based on conditions that hold for an input datapoint, we model LFs as a set of cascading conditions. A rule r is a *tree* where leaf nodes represent labels from a set of labels \mathcal{Y} and the non-leaf nodes are labeled with Boolean predicates from a space of predicates \mathcal{P} . Each non-leaf node has two outgoing edges labeled with **true** and **false**. A rule r assigns a label $r(x)$ to an input datapoint x by evaluating

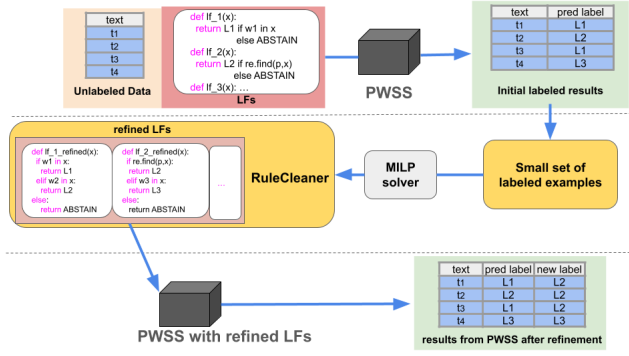


Figure 2: The RULECLEANER framework for repairing LFs (rules \mathcal{R}). After running an PWSS (e.g., Snorkel) on the rule outputs, RULECLEANER fixes the rules \mathcal{R} using labeled examples \mathcal{X}^* . Finally, the PWSS is rerun on the output of the refined rules \mathcal{R}^* on the whole dataset \mathcal{X} to produce the repaired labels.

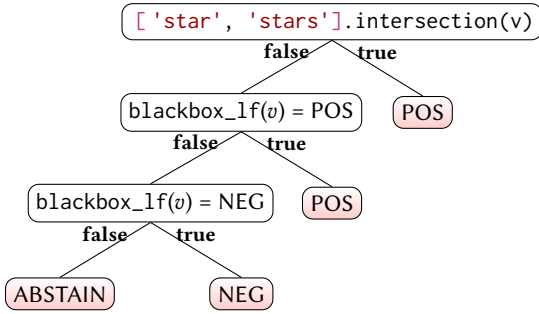


Figure 3: Translating a LF wrapping parts into a blackbox function

the predicate at the root, following the outgoing edge **true** if the predicate evaluates to true and the **false** edge otherwise. Then the predicate of the node at the end of the edge is evaluated. This process is repeated until a leaf node is reached. The label of the leaf node is the label assigned by r to x .

EXAMPLE 2. Figure 4 shows the rule for a LF that returns NEG if the review contains the word ‘waste’ and returns ABSTAIN otherwise. In [18], we show how to translate any LF written in a general-purpose programming language into a rule in PTIME.

To demonstrate that RULECLEANER can support arbitrary LFs, including those with complex components, consider an LF that first checks whether the input text contains the keywords ‘star’ or ‘stars’. If so, it returns the label POSITIVE. Otherwise, it calls an external sentiment analysis function, returning POSITIVE only if the sentiment score exceeds 0.7, and ABSTAIN otherwise. The first condition is directly translated into a predicate in the rule tree, while the sentiment analysis branch is wrapped as a blackbox component `blackbox_lf`. The translated tree rule is shown in Figure 3. Any black-box LF can be translated into a rule tree using our Translate-BBox algorithm [18]. From here on, we will use the terms LF and rule interchangeably.

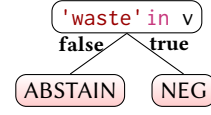


Figure 4: Rule form of the LF `keyword_word_waste` (Figure 1a)

Consider a set of input *datapoints* \mathcal{X} and a set of discrete *labels* \mathcal{Y} . For a datapoint $x \in \mathcal{X}$, y_x^* denotes the datapoint’s (unknown) true label. A PWSS takes \mathcal{X} , the labels \mathcal{Y} , and a set of rules \mathcal{R} as input and produces a *model* $\mathcal{M}_{\mathcal{R},\mathcal{X}}$ as the output that maps each datapoint in \mathcal{X} to a label in \mathcal{Y} . Without loss of generality, we assume the presence of an *abstain label* $y_0 \in \mathcal{Y}$ that is used by the PWSS or a rule to abstain from providing a label to some datapoints.

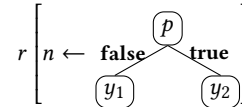
DEFINITION 1 (PWSS). Given a set of datapoints \mathcal{X} , a set of labels \mathcal{Y} , and a set of rules \mathcal{R} , a PWSS takes $\mathcal{R}(\mathcal{X})$ as input and produces a model $\mathcal{M}_{\mathcal{R},\mathcal{X}} : \mathcal{X} \rightarrow \mathcal{Y}$ that maps datapoints $x \in \mathcal{X}$ to labels

$$\mathcal{M}_{\mathcal{R},\mathcal{X}}(x) = y_x.$$

In the following, we will often drop \mathcal{R} and \mathcal{X} from $\mathcal{M}_{\mathcal{R},\mathcal{X}}$ when they are irrelevant to the discussion.

2.2 The Rule Repair Problem

Rule Refinement Repairs. We model a repair of a set of rules \mathcal{R} as a **repair sequence** $\Phi = \phi_1, \dots, \phi_k$ of refinement steps ϕ_i and use $\mathcal{R}^* = \{r'_1, \dots, r'_m\}$ to denote $\Phi(\mathcal{R})$. We repair rules by *refining* them by replacing a leaf node with a new predicate to achieve a desired change to the rule’s result on some datapoints. Consider a rule r , a path P ending in a node n , and a predicate p and two labels y_1 and y_2 . The refinement $\text{refine}(r, P, p, y_1, y_2)$ of r replaces n with a new node labeled p and adds the new leaf nodes for y_1, y_2 :



For example, Figure 1b shows the result of refinement where a leaf POS was replaced with the subtree highlighted in green.

Desiderata. Given the labeled training data \mathcal{X}^* , we would like the repaired rules to provide sufficient information about the true labels for datapoints in \mathcal{X}^* to the PWSS without overfitting to the small number of labeled datapoints in \mathcal{X}^* . Specifically, we want the repair to fulfill the following desiderata:

Datapoint Evidence. We define the *evidence* for a datapoint x_i as the fraction of non-abstain labels ($\neq y_0$) the datapoint receives from the m rules in \mathcal{R} . The repaired rules should provide sufficient evidence for each datapoint x_i , such that the PWSS can make an informed decision about x_i ’s label.

$$\text{Evidence}(x_i) = \frac{\sum_j \mathbb{1}[r'_j(x_i) \neq y_0]}{m}$$

Datapoint Accuracy. The *accuracy* for a datapoint x_i is defined below. The accuracy of the repaired rules that do not abstain on x_i should be high.

$$\text{Acc}(x_i) = \frac{\sum_{j: r'_j(x_i) \neq y_0} \mathbb{1}[r'_j(x_i) = y_{x_i}^*]}{m}$$

Rule Accuracy. In addition, the rules should have high accuracy. The *accuracy* of a rule $r_j \in \mathcal{R}^*$ is defined as the fraction of the n datapoints in \mathcal{X}^* on which it returns the ground truth label.

$$\text{Acc}(r'_j) = \frac{\sum_{i: r'_j(x_i) \neq y_0} \mathbb{1}[r'_j(x_i) = y_{x_i}^*]}{n}$$

Repair Cost. For a repair sequence Φ and $\mathcal{R}^* = \Phi(\mathcal{R})$ we define its cost as the number of labels that differ between the results of $\mathcal{R} = \{r_1, \dots, r_m\}$ and $\mathcal{R}^* = \{r'_1, \dots, r'_m\}$ on \mathcal{X}^* . Optimizing for low repair cost avoids overfitting to \mathcal{X}^* and preserves rule semantics where feasible.

$$\text{cost}(\Phi) = \sum \mathbb{1}[r_j(x_i) \neq r'_j(x_i)]$$

We state the rule repair problem as an optimization problem: minimize the number of changes to labeling function results ($\text{cost}(\Phi)$) while ensuring the desiderata enforced by thresholds τ_E (evidence), τ_{acc} (accuracy), and τ_{racc} (rule accuracy).

DEFINITION 2 (RULE REPAIR PROBLEM). Consider a black-box model $\mathcal{M}_{\mathcal{R}, \mathcal{X}}$ that uses a set of m rules \mathcal{R} , a dataset of n datapoints \mathcal{X} , output labels \mathcal{Y} , and ground truth labels for a subset of datapoints \mathcal{X}^* . Given thresholds $\tau_{acc} \in [0, 1]$, $\tau_E \in [0, 1]$, and $\tau_{racc} \in [0, 1]$, the rule repair problem is to find a repair sequence Φ such that:

$$\begin{aligned} & \mathbf{argmin}_{\Phi} \quad \text{cost}(\Phi) \\ & \mathbf{subject\ to} \quad \forall i \in [1, n] : \text{Acc}(x_i) \geq \tau_{acc} \wedge \text{EVIDENCE}(x_i) \geq \tau_E \\ & \quad \quad \quad \forall j \in [1, m] : \text{Acc}(r'_j) \geq \tau_{racc} \end{aligned}$$

Note that since we treat the PWSS as a black box, we can, in general, not guarantee that the PWSS's performance on the unlabeled dataset \mathcal{X} will improve. Nonetheless, we will demonstrate experimentally in Section 5 that significant improvements in the accuracy of rules on \mathcal{X} can be achieved based on 10s of training examples. This is due to the use of predicates in rule repairs that generalize beyond \mathcal{X}^* . While finding an optimal repair is NP-hard, we can still solve this problem exactly as \mathcal{X}^* is expected to be small.

THEOREM 1. The rule repair problem is NP-hard in the size of \mathcal{R} .

3 Ruleset Repair Algorithm

We now present an algorithm that solves the rule repair problem in two steps. In the first step, we use an MILP to determine desired changes to the outputs of rules, and in the second step, described in Section 4, we implement these changes by refining individual rules to return the desired output on \mathcal{X}^* .

3.1 MILP Formulation

In the MILP, we use an integer variable o_{ij} for each datapoint $x_i \in \mathcal{X}^*$ and rule r_j that stores the label that the repaired rule r'_j should assign to x_i . That is, in combination these variables store the desired changes to the results of rules that we then have to implement by refining each rule r_j to a rule r'_j . We restrict these variables to take values in $[0, |\mathcal{Y}| - 1]$ where value i represents the label $y_i \in \mathcal{Y}$ with 0 encoding y_0 . To encode the objective (minimizing the changes to the outputs of rules on \mathcal{X}^*) we use a Boolean variable m_{ij} for each rule r_j and datapoint x_i that is 1 iff

$o_{ij} \neq r_j(x_i)$ (the output of r'_j on x_i is different from $r_j(x_i)$). The objective is then to minimize the sum of these indicators m_{ij} .

To encode the side constraints of the rule repair problem, we introduce additional indicators: c_{ij} is 1 if $o_{ij} = y_{x_i}^*$, and e_{ij} is 1 if $o_{ij} \neq y_0$. To ensure that the accuracy for each datapoint x_i is above τ_{acc} , we have to ensure that out of rules that do not return y_0 on x_i , i.e., all $j \in [1, m]$ where $e_{ij} = 1$, at least a fraction of τ_{acc} have the correct label ($c_{ij}=1$). This can be enforced if $\sum_j c_{ij} - \sum_j e_{ij} \cdot \tau_{acc} \leq 0$ or equivalently $\sum_j c_{ij} \geq \sum_j e_{ij} \cdot \tau_{acc}$. A symmetric condition is used to ensure LF accuracy using the threshold τ_{racc} and summing up over all datapoints instead of over all rules. Finally, we need to ensure that each datapoint x_i receives a sufficient number of labels $\neq y_0$. Recall that e_{ij} encodes whether LF r'_j returns a non-abstain label. Thus, for m rules we have to enforce: $\forall i \in [1, n] : \sum_j e_{ij} \geq m \cdot \tau_E$. The full MILP is shown below. The non-linear constraints for indicator variables can be translated into linear constraints using the so-called Big M technique [10].

minimize $\sum_i \sum_j m_{ij}$ **subject to**

$$\begin{aligned} \forall i \in [1, n], j \in [1, m] : & \quad \forall i \in [1, n] : \sum_j c_{ij} \geq \sum_j e_{ij} \cdot \tau_{acc} \\ o_{ij} \in [0, |\mathcal{Y}| - 1] & \quad \forall i \in [1, n] : \sum_j e_{ij} \geq m \cdot \tau_E \\ m_{ij} = \mathbb{1}[o_{ij} \neq r_j(x_i)] & \quad \forall j \in [1, m] : \sum_i c_{ij} \geq \sum_i e_{ij} \cdot \tau_{racc} \\ c_{ij} = \mathbb{1}[o_{ij} = y_{x_i}^*] & \\ e_{ij} = \mathbb{1}[o_{ij} > 0] & \end{aligned}$$

As we show next, the solution of the MILP is a solution for the rule repair problem as long as the expected changes to the LF results on \mathcal{X}^* encoded in the variables o_{ij} can be implemented as a repair sequence Φ . As we will show in Section 4 such a repair sequence is guaranteed to exist as long as we choose the space of predicates to use in refinements carefully.

PROPOSITION 1. Consider rules \mathcal{R} , \mathcal{X}^* , and the output o_{ij} produced as a solution to the MILP. If there exists a repair sequence Φ such that for $\mathcal{R}^* = \Phi(\mathcal{R})$ the output on \mathcal{X}^* is equal to o_{ij} for all $i \in [1, n]$ and $j \in [1, m]$, then Φ is a solution to the rule repair problem.

MILP Size. The number of constraints and variables in the MILP is both in $O(n \cdot m)$ where $n = |\mathcal{X}^*|$ and $m = |\mathcal{R}|$. While solving MILPs is hard in general, we demonstrate experimentally that the runtime is acceptable for $|\mathcal{X}^*| \leq 200$.

EXAMPLE 3. Consider a set of 3 datapoints $\mathcal{X}^* = \{x_1, x_2, x_3\}$ with ground truth labels $y_{x_1}^* = 2$, $y_{x_2}^* = 1$, $y_{x_3}^* = 2$, and three rules r_1 to r_3 labels $\mathcal{Y} = \{0, 1, 2\}$ where $y_0 = 0$ and assume that these rules return the results on \mathcal{X}^* shown below on the left where abstain (incorrect) labels are highlighted in blue (red). Assume that all thresholds are set to 50%. That is, each datapoint should receive at least two labels $\neq y_0$, and the accuracy for datapoints and rules is at least 50% (1 correct label if 2 non-abstain labels are returned and 2 correct labels for no abstain label). The minimum number of changes required to fulfill these constraints is 4. One possible solution for the MILP is shown below on the right with modified cells (with correct labels) shown with a black background.

	r_1	r_2	r_3
x_1	1	1	2
x_2	0	1	0
x_3	0	1	0

o_{ij}	$i = 1$	$i = 2$	$i = 3$
$j = 1$	2	1	2
$j = 2$	0	1	1
$j = 3$	2	2	0

Given the outputs o_{ij} of the MILP, we need to find a repair sequence Φ such that for $\mathcal{R}^* = \Phi(\mathcal{R}) = \{r'_1, \dots, r'_m\}$ we have $r'_j(x_i) = o_{ij}$ for all $i \in [1, n]$ and $j \in [1, m]$. An important observation regarding this goal is that as rules operate independently of each other, we can solve this problem one rule at a time.

4 Single Rule Refinement

We now detail our approach for refining a single rule.

4.1 Rule Repair

We now formalize the problem of generating a sequence of refinement steps Φ of minimal size for a rule r_j such that for a set of datapoints and labels $\mathcal{Z} = \{(x_i, y_i)\}_{i=1}^n$, $r'_j = \Phi(r_j)$ we have $r'_j(x_i) = y_i$ for all $(x_i, y_i) \in \mathcal{Z}$. We can use this algorithm to implement the changes to rule outputs computed by the MILP from the previous section using: $\mathcal{Z} = \{(x_i, o_{ij}) \mid x_i \in \mathcal{X}^*\}$. For a sequence of refinement repairs Φ we define its cost as $rcost(\Phi) = |\Phi|$.

DEFINITION 3 (THE SINGLE RULE REFINEMENT PROBLEM). *Given a rule r , a set of datapoints with desired labels \mathcal{Z} , and a set of allowable predicates \mathcal{P} , find a sequence of refinements Φ_{min} using predicates from \mathcal{P} such that for $r_{fix} = \Phi(r)$ we have:*

$$\Phi_{min} = \underset{\Phi}{\operatorname{argmin}} rcost(\Phi) \quad \textbf{subject to} \quad \forall (x, y) \in \mathcal{Z} : r_{fix}(x) = y$$

Let P_{fix} denote the set of paths (from the root to a label on a leaf) in rule r that are taken by the datapoints from \mathcal{X} . For $P \in P_{fix}$, \mathcal{X}_P denotes all datapoints from \mathcal{X} for which the path is P , hence also $\mathcal{X} = \bigcup_{P \in P_{fix}} \mathcal{X}_P$. Similarly, \mathcal{Z}_P denotes the subset of \mathcal{Z} for datapoints x with path P in rule r . The algorithm for solving the single rule repair problem we will present in the following exploits two important properties of this problem.

Independence of path repairs. As any refinement in a minimal repair will only extend paths in P_{fix} (any other refinement does not affect the labels for \mathcal{X}) and refinements at any path P_1 do not affect the labels of datapoints in \mathcal{X}_{P_2} for a path $P_2 \neq P_1$, a solution to the single rule repair problem can be constructed one path at a time (see [18] for the formal proof).

Existence of path repairs. In [18], we show that path repairs with a cost of at most $|\mathcal{X}_P|$ are guaranteed to exist as long as the space of predicates is partitioning. That is, for any two datapoints x_1 and x_2 we can find a predicate p such that $p(x_1) \neq p(x_2)$. Note that for textual data, even a simple predicate space that only contains predicates of the form $w \in x$ where w is a word is partitioning as long as any two datapoints (documents in the case of text data) will differ in at least one word. Intuitively, this guarantees the existence of a repair as for any two datapoints x_1 and x_2 with $\mathcal{Z}(x_1) \neq \mathcal{Z}(x_2)$ that share the same path (and, thus, also label) in a rule r we can refine r using an appropriate predicate p with $p(x_1) \neq p(x_2)$ to assign the desired labels to x_1 and x_2 .

The pseudocode for `SingleRuleRefine` is given in Algorithm 1. Given a single rule r , this algorithm determines a refinement-based repair Φ_{min} for r such that $\Phi_{min}(r)$ returns the designed label $\mathcal{Z}(x)$ for all datapoints specified in \mathcal{Z} by refining one path at a time using a function `RefinePath`. The problem solved by `RefinePath` is NP-hard. Next, we introduce an algorithm implementing `RefinePath`

Algorithm 1: SingleRuleRefine

Input : Rule r , Labeled datapoints \mathcal{Z} .
Output : Repair sequence Φ such that $\Phi(r)$ fixes \mathcal{Z}

```

1  $Y \leftarrow \emptyset, \Phi \leftarrow \emptyset$ 
2  $P_{fix} \leftarrow \{P[r, x] \mid x \in \mathcal{X}\}$ 
3  $r_{cur} \leftarrow r$ 
4 foreach  $P \in P_{fix}$  do /* Fix one path at a time */
5     /* Fix path  $P$  to return correctly labels on  $\mathcal{Z}$  */
6      $\mathcal{Z}_P \leftarrow \{(x, y) \mid (x, y) \in \mathcal{Z} \wedge P[r, x] = P\}$ 
7      $\phi \leftarrow \text{RefinePath}(r_{cur}, P, \mathcal{Z}_P)$ 
8      $r_{cur} \leftarrow \phi(r_{cur})$ 
9      $\Phi \leftarrow \Phi.append(\phi)$ 
10 return  $\Phi$ 

```

Algorithm 2: EntropyPathRepair

Input : Rule r , Path P_{in} , Ground truth labels $\mathcal{Z}_{P_{in}}$
Output : Repair sequence Φ which fixes r wrt. $\mathcal{Z}_{P_{in}}$

```

1  $todo \leftarrow [(P_{in}, \mathcal{Z}_{P_{in}})]$ 
2  $\Phi \leftarrow []$ 
3  $r_{cur} \leftarrow r$ 
4  $\mathcal{P}_{all} \leftarrow \text{GetAllCandPredicates}(P_{in}, \mathcal{Z}_{P_{in}})$ 
5 while  $todo \neq \emptyset$  do
6      $(P, \mathcal{Z}_P) \leftarrow pop(todo)$ 
7      $p_{new} \leftarrow \operatorname{argmin}_{p \in \mathcal{P}_{all}} I_G(\mathcal{Z}_P, p)$ 
8      $\mathcal{Z}_{false} \leftarrow \{(x, y) \mid (x, y) \in \mathcal{Z}_P \wedge \neg p(x)\}$ 
9      $\mathcal{Z}_{true} \leftarrow \{(x, y) \mid (x, y) \in \mathcal{Z}_P \wedge p(x)\}$ 
10     $y_{max} \leftarrow \operatorname{argmax}_{y \in \mathcal{Y}} |\{x \mid \mathcal{Z}_{true}(x) = y\}|$ 
11     $\phi_{new} \leftarrow \text{refine}(r_{cur}, P, p, Y[P], y_{max})$ 
12     $r_{cur} \leftarrow \phi_{new}(r_{cur})$ 
13     $\Phi \leftarrow \Phi.append(\phi_{new})$ 
14    if  $|\mathcal{Y}_{\mathcal{Z}_{false}}| > 1$  then
15         $todo.push((P[r_{cur}, \mathcal{Z}_{false}], \mathcal{Z}_{false}))$ 
16    if  $|\mathcal{Y}_{\mathcal{Z}_{true}}| > 1$  then
17         $todo.push((P[r_{cur}, \mathcal{Z}_{true}], \mathcal{Z}_{true}))$ 
18 return  $\Phi$ 

```

that utilizes an information-theoretic heuristic that does not guarantee that the returned repair is minimal but works well in practice.

4.2 Path Repair: EntropyPathRepair

Given a rule r , a path $P_{in} \in P_{fix}$, and the datapoints and desired labels for this path $\mathcal{Z}_{P_{in}}$, our algorithm `EntropyPathRepair` avoids the exponential runtime of an optimal brute force algorithm `BruteForcePathRepair` that enumerates all possible refinements (see [18]). We achieve this by greedily selecting predicates that best separate datapoints with different labels at each step. To measure the quality of a split, we employ the entropy-based *Gini impurity score* I_G [16]. Given a candidate predicate p for splitting a set of datapoints and their labels at path P (\mathcal{Z}_P), we denote the subsets of \mathcal{Z}_P generated

by splitting \mathcal{Z}_P based on p :

$$\mathcal{Z}_{\text{false}} = \{(x, y) \mid (x, y) \in \mathcal{Z}_P \wedge \neg p(x)\}$$

$$\mathcal{Z}_{\text{true}} = \{(x, y) \mid (x, y) \in \mathcal{Z}_P \wedge p(x)\}$$

Using $\mathcal{Z}_{\text{false}}$ and $\mathcal{Z}_{\text{true}}$ we define the score $I_G(\mathcal{Z}_P, p)$ for p :

$$I_G(\mathcal{Z}_P, p) = \frac{|\mathcal{Z}_{\text{false}}| \cdot I_G(\mathcal{Z}_{\text{false}}) + |\mathcal{Z}_{\text{true}}| \cdot I_G(\mathcal{Z}_{\text{true}})}{|\mathcal{Z}_P|}$$

$$I_G(Z) = 1 - \sum_{y \in \mathcal{Y}_Z} p(y)^2 \quad p(y) = \frac{|\{x \mid Z(x) = y\}|}{|Z|}$$

For a set of ground truth labels Z , $I_G(Z)$ is minimal if $\mathcal{Y}_Z = \{y \mid \exists x : (x, y) \in Z\}$ contains a single label. Intuitively, we want to select predicates such that all datapoints that reach a particular leaf node are assigned the same label. At each step, the best separation is achieved by selecting a predicate p that minimizes $I_G(\mathcal{Z}_P, p)$.

Algorithm 2 first determines all candidate predicates using function `GetAllCandPredicates`. Then, it iteratively selects predicates until all datapoints are assigned the expected label by the rule. For that, we maintain a queue of paths paired with a set \mathcal{Z}_P of datapoints with expected labels that still need to be processed. In each iteration of the algorithm’s main loop, we pop one pair of a path P and datapoints with labels \mathcal{Z}_P from the queue. We then determine the predicate p that minimizes the entropy of \mathcal{Z}_P . Afterward, we determine two subsets of datapoints from \mathcal{X}_P : datapoints fulfilling p and those that do not. We then generate a refinement repair step ϕ_{new} for the current version of the rule (r_{cur}) that replaces the last element on P_{cur} with predicate p ($Y[P]$ denotes the label of the node at the end of P). The child at the **true** edge of the node for p is then assigned the most prevalent label y_{max} for the datapoints at this node (the datapoints from $\mathcal{Z}_{\text{true}}$). Finally, unless they only contain one label, new entries for $\mathcal{Z}_{\text{false}}$ and $\mathcal{Z}_{\text{true}}$ are appended to the queue. As shown below, `EntropyPathRepair` is correct (the proof is shown in [18]).

THEOREM 2 (CORRECTNESS). *Consider a rule r , ground-truth labels of a set of datapoints \mathcal{Z}_P , and partitioning space of predicates \mathcal{P} . Let Φ be the repair sequence produced by `EntropyPathRepair` for path P . Then we have:*

$$\forall (x, y) \in \mathcal{Z}_P : \Phi(r)(x) = y$$

5 Experiments

We evaluate the runtime of `RULECLEANER` and its effectiveness in improving the accuracy of rules produced by Witan [7] and LLMs. Additionally, we analyze the trade-offs introduced by the three path repair algorithms proposed in this work. Our experiments use *Snorkel* [25] as the default PWSS. To demonstrate that `RULECLEANER` is agnostic to the choice of PWSS, we also test it with alternative PWSSs from [36], measuring improvements in global accuracy. We assess both the runtime and the quality of the refinements produced by `RULECLEANER` across several parameters. `RULECLEANER` is implemented in Python, and all experiments are conducted on Oracle Linux Server 7.9 with 2 x AMD EPYC 7742 CPUs and 128GB RAM.

Datasets and rules. The datasets used in the experiments are listed in Table 2. Note that because of the complexity of and the nature of multi-class labels, the LFs we used for *CmPt* are all from [33], which has 26 LFs. We give a brief description of each dataset:

Dataset	#row	#word	\mathcal{Y}	#LFs _{witan}	#LFs _{llm}
<i>Amazon</i>	200000	68.9	pos/neg	15	23
<i>AGnews</i>	60000	37.7	busi/tech	9	21
<i>PP</i>	54476	55.8	physician/prof	18	20
<i>IMDB</i>	50000	230.7	pos/neg	7	20
<i>FNews</i>	44898	405.9	true/false	11	20
<i>Yelp</i>	38000	133.6	neg/pos	8	20
<i>PT</i>	24588	62.2	prof/teacher	7	19
<i>CmPt</i>	16075	27.9	10 relations	-	-
<i>PA</i>	12236	62.6	painter/architect	10	18
<i>Tweets</i>	11541	18.5	pos/neg	16	18
<i>SMS</i>	5572	15.6	spam/ham	17	16
<i>MGenre</i>	1945	26.5	action/romance	10	14

Table 2: LF dataset statistics.

Amazon: product reviews from Amazon and their sentiment label [14]. *AGnews*: categorized news articles from AG’s corpus of news articles. For this dataset, we chose a binary class version from [7]. *PP*: descriptions of biographies, each labeled as a physician or a professor [6]. *IMDB*: IMDB movie reviews [21]. *FNews*: Fake news identification [1]. *Yelp*: Yelp reviews [38]. *PT*: descriptions of individuals, each labeled as a professor or a teacher.[6]. *PA*: descriptions of individuals, each labeled a painter or an architect. [6]. *Tweets*: classification of tweets on disasters [22]. *SMS*: classification of SMS messages [2]. *MGenre*: movie genre classification based on plots [30]. *CmPt*: chemical-protein relationship classification from [17]. Unless stated otherwise, the experiments in this section are run with `EntropyPathRepair`. We present a detailed evaluation of all path repair algorithms in Section 5.4.

5.1 Refining labelling functions

In this experiment, we investigate the effects of several parameters on the performance and quality of the rules repaired with `RULECLEANER` for several datasets.

Varying the number of labeled examples. We evaluate how the size of \mathcal{X}^* affects global accuracy. Global accuracy is defined as the accuracy of the labels predicted by the trained PWSS using the LFs compared to the ground truth labels. Given the limited scalability of MILP solvers in the number of variables, we used at most $|\mathcal{X}^*| = 150$ datapoints. The labeled datapoints are randomly sampled from \mathcal{X} , with 50% correct predictions by PWSS and 50% wrong predictions within each sample. The reason for sampling in this manner is to provide sufficient evidence for correct predictions and predictions that need to be adjusted. Even if we have no control over the creation of \mathcal{X}^* , we can achieve this by sampling from a larger set of labeled examples. Figure 6a shows the global accuracy after retraining a *Snorkel* (PWSS) model with the rules refined by `RULECLEANER`. The repairs improve the global accuracy on 8 out of 9 datasets, even for very small sample sizes. The variance of the new global accuracy also decreases as the amount of labeled examples increases.

Varying thresholds. We evaluate the relationships between τ_{acc} , τ_E , τ_{racc} and new global accuracy. We used *Tweets* with 20 labeled examples. The details of the experiments and analysis are shown in Appendix B.2. Based on the experiments, we recommend setting all of the thresholds to ~ 0.7 .

5.2 Runtime

Runtime breakdowns for a subset of the experiments from Section 5.1 are shown in Figure 7. For the breakdown of the other datasets, please refer to Appendix B.1. The total runtime increases as we increase the amount of labeled examples. The runtime of the refinement step is strongly correlated with the average length of the texts in the input dataset, i.e., the longer the average text length (as presented in *average # words* in Table 2), the more time is required to select the best predicate using EntropyPathRepair.

It is important to note that the runtime changes for both *snorkel run after refinement* and *MILP* do not exhibit a strictly linear pattern. The reason for such non-linearity arises from the fact that the labeled datapoints are randomly sampled from \mathcal{X} , and the complexity of solving the MILP problem depends on the sparsity of the solution space. The same reason applies for retraining with Snorkel using the refined rules. Some sets of labeled examples result in more complex rules even when the sample size is small, increasing the time required for Snorkel to fit a model.

5.3 LLMs vs RULECLEANER

In this section, we compare our approach against LLMs. We consider three setups: (i) using the LLM as a labeler (without any use of LFs) and (ii) using the LLM to generate LFs based on with labeled examples; and (iii) using the LLM to repair labeling functions (refer to Appendix C). For (ii), we then investigate whether RULECLEANER can successfully improve the LFs generated by the LLM.

The LLM as a labeler. We compare the performance and quality of RULECLEANER with Snorkel and LLMs as a standalone labeler. In this experiment, we use GPT-4o and Llama-3-8B-instruct (Llama 3 8B), using a zero-shot prompt to describe the task. Both LLMs receive the possible labels along with the sentences, but not the LFs. To optimize API usage, we batch 10 datapoints per call for GPT-4o, whereas for Llama-3-8B-instruct, we label one sentence per call to maintain response validity. Experiments were conducted on a Mac Studio (Apple M2 Max, 12-core CPU, 64GB unified memory, SSD storage). The setup with RULECLEANER (RC for short shown in the plot) is the default setup from Section 5.1. The runtime and quality comparisons are presented in Figure 5. We did set a 48-hour time limit, including only datasets where all three competitors completed the task within this time limit. GPT-4o achieves the highest accuracy in 6 out of 7 datasets. However, for dataset *CmPt*, RULECLEANER with Snorkel (37.9%) outperforms GPT-4o (27.4%). Upon further analysis, we speculate that GPT-4o’s poor performance could stem from the specialized terminology and complex domain knowledge required for labeling in ChemProt. Unlike general-purpose datasets where LLMs excel, ChemProt contains highly domain-specific biomedical entity interactions, which may be challenging for zero-shot prompting. Using LLMs as an end-to-end labeler comes at an unacceptable computational/monetary cost. The experiments with GPT-4o did cost \$254.42 for API usage. While we do not know the precise computational resources that were required, our local experiments with Llama-3-8B-instruct, a significantly smaller model that also cannot compete with RULECLEANER in terms of accuracy on most datasets, demonstrate the high computational cost of using an LLM for this purpose. In fact, RULECLEANER is ~ 2 to ~ 4 orders of magnitude faster than Llama-3-8B-instruct. In

summary, while large models like GPT-4o, but not smaller models like Llama-3-8B-instruct, can achieve high accuracy as labelers, this comes at a prohibitively high computational cost. RULECLEANER outperforms Llama-3-8B-instruct in terms of accuracy on most datasets and between 56x to 1,312x in terms of runtime. Furthermore, LFs have the additional advantage of being inherently interpretable which is not the case for labeling with LLMs.

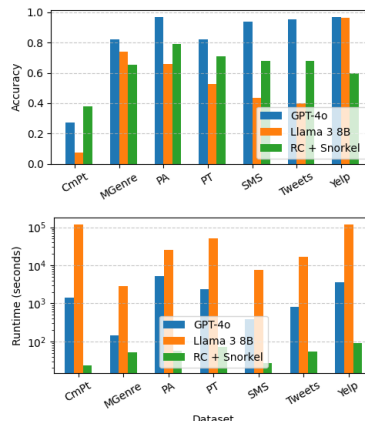


Figure 5: Comparison of 3 labelers

The LLM as a LF generator. We investigate whether LLMs can generate effective LFs and whether RULECLEANER can fix such functions to improve their accuracy. The use of LLMs for generating LFs has been explored in recent years [11, 19]. We adapted these existing methods to generate LFs. For an example prompt, see [18]. In each prompt, we sample a small set of sentences along with their ground truth labels and provide two LF templates based on keywords and regular expressions. Since the number of LFs is determined by how many times we query the LLM, we scale the number of LFs logarithmically in the dataset size. The number of LFs produced for each dataset is reported as $\#LFs_{llm}$ in Table 2.

To evaluate the quality of generated LFs LFs_{llm} , we use the same experimental setup described in Section 5.1. The results are shown in Figure 6b. Out of the nine datasets used in this experiment, three exhibit higher original global accuracy after training with Snorkel compared to using LFs_{witan} . After refinement using labeled datapoints, we observe improvements in eight out of nine datasets.

Two datasets, *SMS* and *Tweets*, show significant improvement. Upon further inspection, we observed that some of the LFs generated by the LLM assign incorrect labels. For example, in *SMS*, one LF is defined as follows: `return HAM if any(x in text for x in ['sorry', 'please', 'home', 'call', 'message', 'buy', 'talk', 'problem', 'help', 'ask']) else ABSTAIN`. Some of these keywords, such as “call” and “message”, frequently appear in spam messages. RULECLEANER successfully refines this LF by correcting its label assignment, thereby improving accuracy.

5.4 Path Repair Algorithms

Next, we compare the three path repair algorithms discussed in Section 4.2. In this experiment, we used the *Tweets* dataset and randomly selected between 2 and 10 labeled examples. We show *cost*, the repaired rule size in terms of number of nodes and the

runtime in Figure 8. Note that for 10 labeled examples, the repair runtime for BruteForcePathRepair exceeded the time limit we set for this experiment (600 secs) and, thus, is absent from the plot. The runtime of BruteForcePathRepair is prohibitory large even for just 8 datapoints. EntropyPathRepair achieves almost the same repair cost as BruteForcePathRepair while being significantly faster. While GreedyPathRepair is the fastest algorithm, this comes at the cost of a significantly higher repair cost.

It is obvious that the runtime for BruteForcePathRepair is significantly higher than the other 2 algorithms. GreedyPathRepair is the fastest since it picks the first available predicate without any additional computation. In terms of rule sizes after the repairs, BruteForcePathRepair generates the smallest rules since it will exhaustively enumerate all the possible solutions and is guaranteed to find the smallest possible solution. It is worth noting that EntropyPathRepair is only slightly worse than BruteForcePathRepair while being significantly faster.

5.5 Complexity of Refined LFs

To analyze the evolution of rule complexity during refinement, we conducted experiments using labeling functions generated by GPT-4o. We varied the number of labeled data points (20 and 40) and evaluated the resulting rule trees across multiple random samples. Table 3 reports the average tree depth and node count for three representative datasets. We observe a consistent increase in rule complexity (depth) with more input data, while the tree size (in terms of both depth and node count) grows sublinearly with respect to the input size. This suggests that the refinement process increases expressiveness efficiently without leading to overfitting.

Dataset	Input Size	Depth	Node Count
AGnews	20	7.52	16.48
AGnews	40	11.54	28.06
IMDB	20	4.95	11.67
IMDB	40	7.48	20.21
SMS	20	7.08	14.24
SMS	40	11.08	25.47

Table 3: Average depth and node count of refined LFs.

5.6 Other PWSSs

In addition to using Snorkel as PWSS, we also tested models from [36] using the datasets from Table 2. The results are shown in Table 4. RULECLEANER consistently improves global accuracy across all PWSSs. The most substantial relative gains are observed for MetaL (+15.5%), DawidSkene (+10.5%), and FlyingSquid (+6.9%), while Majority Voting sees a modest improvement of +0.7%.

Model	Before	After	Rel. Gain
MetaL	0.538	0.693	+15.5%
DawidSkene	0.566	0.672	+10.5%
FlyingSquid	0.619	0.688	+6.9%
Majority	0.685	0.690	+0.7%

Table 4: Global accuracy improvements of PWSSs after refinement.

6 Related Work

We next survey related work on tasks that can be modeled as PWSSs as well as discuss approaches for automatically generating rules for PWSSs and improving a given rule set.

Programmatic weak supervision (PWS). Weak supervision is a general technique of learning from noisy supervision signals, widely applied for data labeling to generate training data [25, 26, 30] (the main use case we target in this work), data repair [27], and entity matching [24]. Its main advantage is reducing the effort of creating training data from unlabeled data. The programmatic weak supervision paradigm pioneered in Snorkel [25] has the additional advantage that the labeling functions are interpretable. However, as such LFs are typically noisy heuristics, systems like Snorkel combine the output of LFs using a model.

Automatic generation and fixing labeling functions. While PWS proves effective, asking human annotators to create a large set of high-quality labeling functions requires domain knowledge, programming skills, and time. As a result, the automatic generation or improvement of labeling heuristics has received much attention from the research community. Some existing methods demand interactive user feedback in creating labeling functions [4, 9]. *Witan* [7] asks a domain expert to select the automatically generated LFs and assign labels to the LFs. While the LFs produced by Witan are certainly useful, we demonstrate in our experimental evaluation that applying RULECLEANER to Witan LFs can significantly improve accuracy. Other methods generate LFs without requiring user annotations. *Snuba* [30] fits classification models, such as decision trees and logistic regressions, as LFs on a small labeled training set, followed by a pruning process to determine the final set of LFs. *Datasculpt* [11] prompts a large language model (LLM) with a small set of labeled training data and keyword- or pattern-based LFs as in-context examples. The LLM then generates LFs for unlabeled examples based on this input. *Evaporate* [3] uses an LLM to generate data extraction functions, and then applies weak supervision to filter out low-quality functions and aggregate the results.

Hsieh et al. [15] propose *Nemo*, a framework for selecting data to guide users in developing LFs. It estimates the likelihood of users proposing specific LFs using a utility metric for LFs and a model of user behavior. *Nemo* tailors LFs to the neighborhood of the data, assuming that user-developed LFs are more accurate for data similar to those used for LFs creation. However, unlike RULECLEANER, *Nemo* lacks a mechanism for the user to provide feedback on the labeling results, preventing the automatic deletion and refinement of LFs. *ULF* [29] is an unsupervised system for adjusting LFs for unlabeled samples (instead of repairing them) using k-fold cross-validation, extending previous approaches addressing labeling errors [23, 31].

Explanations for weakly supervised systems. There is a large body of work on explaining the results of weak-supervised systems that target improving the final model or better involving human annotators [4, 5, 12, 32, 35, 35, 37]. For instance, [35] uses influence function to identify LFs responsible for erroneous labels; WeShap [12] measures the Shapley value of LFs to rank and prune LFs. However, most of this work has stopped short of repairing the rules in a PWSS and, thus, are orthogonal to our work. Still,

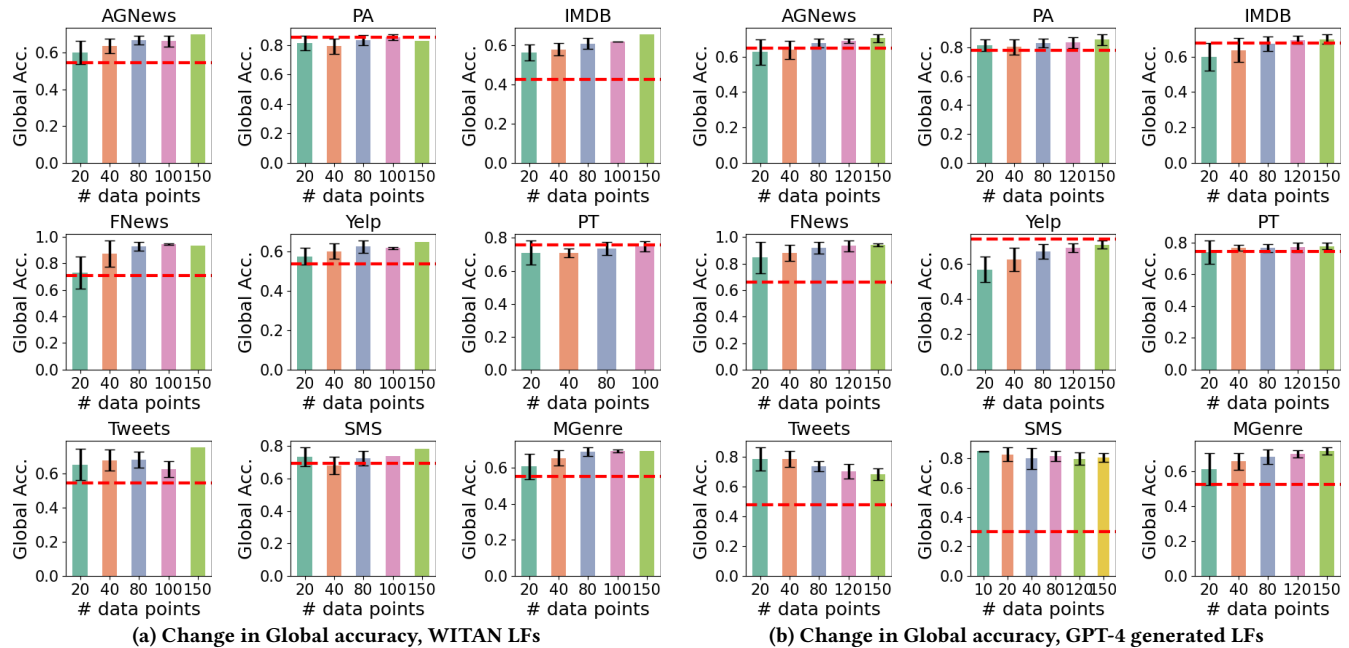


Figure 6: Impact of repairs on global accuracy (the red dotted line is accuracy before the repair) for LFs generated by Witan and GPT-4. We vary the number of labeled examples \mathcal{X}^* .

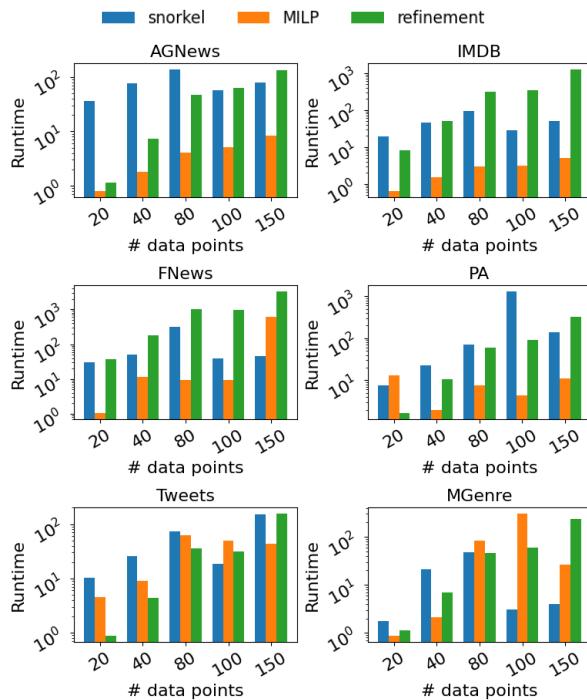


Figure 7: Runtime, varying the size of \mathcal{X}^* .

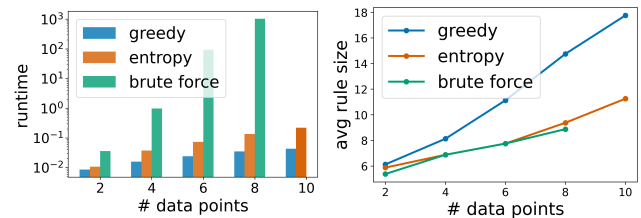


Figure 8: Comparing path repair algorithms

7 Conclusions and Future Work

We study repairs for LFs in PWS based on a small set of labeled examples. Our algorithm is highly effective in improving the accuracy of PWSSs by improving rules created by a human expert or automatically discovered by a system like Witan [7]. In future work, we will explore the application of our rule repair algorithms to other tasks that can be modeled as PWSS, e.g., information extraction based on user-provided rules [20, 28].

8 Acknowledgments

This work is supported in part by NSF Awards IIS-2420577, IIS-2420691, and IIS-2147061, and by the Natural Sciences and Engineering Research Council of Canada (NSERC) under award numbers RGPIN-2025-04724 and DGECR-2025-00373. The work of Amir Gilad was funded by the Israel Science Foundation (ISF) under grant 1702/24, the Scharf-Ullman Endowment, and the Alon Scholarship.

explanations provided by such systems might guide users in selecting what datapoints to label. People have also studied using human-annotated natural language explanations to build LFs [13].

References

- [1] Hadeer Ahmed, Issa Traoré, and Sherif Saad. 2018. Detecting opinion spams and fake news using text classification. *Secur. Priv.* 1, 1 (2018).
- [2] Tiago A. Almeida, José María Gómez Hidalgo, and Akebo Yamakami. 2011. Contributions to the study of SMS spam filtering: new collection and results. In *ACM Symposium on Document Engineering*. ACM, 259–262.
- [3] Simran Arora, Brandon Yang, Sabri Eyuboglu, Avani Narayan, Andrew Hojel, Immanuel Trummer, and Christopher Ré. 2023. Language Models Enable Simple Systems for Generating Structured Views of Heterogeneous Data Lakes. *Proceedings of the VLDB Endowment* 17, 2 (2023), 92–105.
- [4] Benedikt Boecking, Willie Neiswanger, Eric Xing, and Artur Dubrawski. 2021. Interactive Weak Supervision: Learning Useful Heuristics for Data Labeling. In *International Conference on Learning Representations*.
- [5] Bradley Butcher, Miri Zilka, Darren Cook, Jiri Hron, and Adrian Weller. 2023. Optimising Human-Machine Collaboration for Efficient High-Precision Information Extraction from Text Documents. *arXiv preprint arXiv:2302.09324* (2023).
- [6] Maria De-Arteaga, Alexey Romanov, Hanna M. Wallach, Jennifer T. Chayes, Christian Borgs, Alexandra Chouldechova, Sahin Cem Geyik, Krishnaram Kenchadadi, and Adam Tauman Kalai. 2019. Bias in Bios: A Case Study of Semantic Representation Bias in a High-Stakes Setting. In *FAT*, danah boyd and Jamie H. Morgenstern (Eds.), 120–128.
- [7] Benjamin Denham, Edmund M.-K. Lai, Roopak Sinha, and M. Asif Naeem. 2022. Witan: Unsupervised Labelling Function Generation for Assisted Data Programming. *PVLDB* 15, 11 (2022), 2334–2347.
- [8] Daniel Fu, Mayee Chen, Frederic Sala, Sarah Hooper, Kayvon Fatahalian, and Christopher Ré. 2020. Fast and three-rious: Speeding up weak supervision with triplet methods. In *International conference on machine learning*. PMLR, 3280–3291.
- [9] Sainyam Galhotra, Behzad Golshan, and Wang-Chiew Tan. 2021. Adaptive Rule Discovery for Labeling Text Data. In *SIGMOD*. 2217–2225.
- [10] Igor Griva, Stephen G Nash, and Ariela Sofer. 2008. *Linear and Nonlinear Optimization 2nd Edition*. SIAM.
- [11] Naiqing Guan, Kaiwen Chen, and Nick Koudas. 2025. DataSculpt: Cost-Efficient Label Function Design via Prompting Large Language Models. In *Proceedings 28th International Conference on Extending Database Technology, EDBT 2025, Barcelona, Spain, March 25–28, 2025*. OpenProceedings.org, 226–232.
- [12] Naiqing Guan and Nick Koudas. 2024. Weshap: Weak Supervision Source Evaluation With Shapley Values. *CoRR abs/2406.11010* (2024). arXiv:2406.11010
- [13] Braden Hancock, Martin Bringmann, Paroma Varma, Percy Liang, Stephanie Wang, and Christopher Ré. 2018. Training classifiers with natural language explanations. In *ACL*, Vol. 2018. 1884.
- [14] Ruining He and Julian J. McAuley. 2016. Ups and Downs: Modeling the Visual Evolution of Fashion Trends with One-Class Collaborative Filtering. In *WWW*, Jacqueline Bourdeau, Jim Hendler, Roger Nkambou, Ian Horrocks, and Ben Y. Zhao (Eds.), 507–517.
- [15] Cheng-Yu Hsieh, Jieyu Zhang, and Alexander J. Ratner. 2022. Nemo: Guiding and Contextualizing Weak Supervision for Interactive Data Programming. *PVLDB* 15, 13 (2022), 4093–4105.
- [16] Sotiris B. Kotsiantis. 2013. Decision Trees: a Recent Overview. *Artif. Intell. Rev.* 39, 4 (2013), 261–283.
- [17] Martin Krallinger, Obdulia Rabal, Saber A Akhondi, Martín Pérez Pérez, Jesús Santamaría, Gael Pérez Rodríguez, Georgios Tsatsaronis, Ander Intxaurreondo, José Antonio López, Umesh Nandal, et al. 2017. Overview of the BioCreative VI chemical-protein interaction Track. In *Proceedings of the sixth BioCreative challenge evaluation workshop*, Vol. 1. 141–146.
- [18] Chenjie Li, Amir Gilad, Boris Glavic, Zhengjie Miao, and Sudeepa Roy. 2025. Refining Labeling Functions with Limited Labeled Data. *arXiv:2505.23470 [cs.LG]* <https://arxiv.org/abs/2505.23470>
- [19] Chenjie Li, Dan Zhang, and Jin Wang. 2024. LLM-assisted Labeling Function Generation for Semantic Type Detection. In *Proceedings of Workshops at the 50th International Conference on Very Large Data Bases, VLDB 2024, Guangzhou, China, August 26–30, 2024*.
- [20] B. Liu, L. Chiticariu, V. Chu, HV Jagadish, and F.R. Reiss. 2010. Automatic Rule Refinement for Information Extraction. *PVLDB* 3, 1 (2010).
- [21] Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011. Learning Word Vectors for Sentiment Analysis. In *Association for Computational Linguistics: Human Language Technologies*, Dekang Lin, Yuji Matsumoto, and Rada Mihalcea (Eds.), 142–150.
- [22] Hussein Mouzannar, Yara Rizk, and Mariette Awad. 2018. Damage Identification in Social Media Posts using Multimodal Deep Learning. In *International Conference on Information Systems for Crisis Response and Management*, Kees Boersma and Brian M. Tomaszewski (Eds.), ISCRAM Association.
- [23] Curtis G. Northcutt, Lu Jiang, and Isaac L. Chuang. 2021. Confident Learning: Estimating Uncertainty in Dataset Labels. *J. Artif. Intell. Res.* 70 (2021), 1373–1411.
- [24] Fatemah Panahi, Wentao Wu, AnHui Doan, and Jeffrey F. Naughton. 2017. Towards Interactive Debugging of Rule-based Entity Matching.. In *EDBT*. 354–365.
- [25] Alexander Ratner, Stephen H. Bach, Henry R. Ehrenberg, Jason A. Fries, Sen Wu, and Christopher Ré. 2020. Snorkel: rapid training data creation with weak supervision. *VLDBJ* 29, 2-3 (2020), 709–730.
- [26] Alexander J. Ratner, Christopher De Sa, Sen Wu, Daniel Selsam, and Christopher Ré. 2016. Data Programming: Creating Large Training Sets, Quickly. In *NIPS*. 3567–3575.
- [27] Theodoros Rekatsinas, Xu Chu, Ihab F. Ilyas, and Christopher Ré. 2017. HoloClean: Holistic Data Repairs with Probabilistic Inference. *PVLDB* 10, 11 (2017), 1190–1201.
- [28] Sudeepa Roy, Laura Chiticariu, Vitaly Feldman, Frederick R Reiss, and Huaiyu Zhu. 2013. Provenance-based dictionary refinement in information extraction. In *SIGMOD*. 457–468.
- [29] Anastasiya Sedova and Benjamin Roth. 2022. ULF: Unsupervised Labeling Function Correction using Cross-Validation for Weak Supervision. *CoRR abs/2204.06863* (2022).
- [30] Paroma Varma and Christopher Ré. 2018. Snuba: Automating Weak Supervision to Label Training Data. *PVLDB* 12, 3 (2018), 223–236.
- [31] Zihan Wang, Jingbo Shang, Liyuan Liu, Lihao Lu, Jiacheng Liu, and Jiawei Han. 2019. CrossWeigh: Training Named Entity Tagger from Imperfect Annotations. In *EMNLP-IJCNLP*. 5153–5162.
- [32] Peilin Yu and Stephen Bach. 2023. Alfred: A System for Prompted Weak Supervision. *arXiv preprint arXiv:2305.18623* (2023).
- [33] Yue Yu, Simiao Zuo, Haoming Jiang, Wendi Ren, Tuo Zhao, and Chao Zhang. 2020. Fine-tuning pre-trained language model with weak supervision: A contrastive-regularized self-training approach. *arXiv preprint arXiv:2010.07835* (2020).
- [34] Jieyu Zhang, Cheng-Yu Hsieh, Yue Yu, Chao Zhang, and Alexander Ratner. 2022. A Survey on Programmatic Weak Supervision. *CoRR abs/2202.05433* (2022). arXiv:2202.05433
- [35] Jieyu Zhang, Haonan Wang, Cheng-Yu Hsieh, and Alexander J. Ratner. 2022. Understanding Programmatic Weak Supervision via Source-aware Influence Function. In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.
- [36] Jieyu Zhang, Yue Yu, Yinghao Li, Yujing Wang, Yaming Yang, Mao Yang, and Alexander Ratner. 2021. WRENCH: A comprehensive benchmark for weak supervision. *arXiv preprint arXiv:2109.11377* (2021).
- [37] Xiaoyu Zhang, Xiwei Xuan, Alden Dima, Thurston Sexton, and Kwan-Liu Ma. 2023. LabelVizier: Interactive Validation and Relabeling for Technical Text Annotations. In *2023 IEEE 16th Pacific Visualization Symposium (PacificVis)*. IEEE, 167–176.
- [38] Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015. Character-level Convolutional Networks for Text Classification. In *NeurIPS*, Corinna Cortes, Neil D. Lawrence, Daniel D. Lee, Masashi Sugiyama, and Roman Garnett (Eds.), 649–657.

A Path Refinement Repairs - Proofs and Additional Details

A.1 GreedyPathRepair

The function GreedyPathRepair is shown Algorithm 3. This algorithm maintains a list of pairs of paths and datapoints at these paths to be processed. This list is initialized with all datapoints \mathcal{X}_P from \mathcal{Z}_P and the path P provided as input to the algorithm. In each iteration, the algorithm picks two datapoints x_1 and x_2 from the current set and selects a predicate p such that $p(x_1) \neq p(x_2)$. It then refines the rule with p and appends $\mathcal{X}_1 = \{x \mid x \in \mathcal{X}_P \wedge p(x)\}$ and $\mathcal{X}_2 = \{x \mid x \in \mathcal{X}_P \wedge \neg p(x)\}$ with their respective paths to the list. As shown in [18], this algorithm terminates after adding at most $|\mathcal{X}_P|$ new predicates.

To ensure that all datapoints ending in path P get assigned the desired label based on \mathcal{Z}_P , we need to add predicates to the end of P to “reroute” each datapoint to a leaf node with the desired label. As mentioned above, this algorithm implements the approach from [18]: for a set of datapoints taking a path with prefix P ending in a leaf node that is not pure (not all datapoints in the set have the same expected label), we pick a predicate that “separates” the datapoints, i.e., that evaluate to true on one of the datapoints and false on the other. Our algorithm applies this step until all leaf nodes are pure wrt. the datapoints from \mathcal{X}_P . For that, we maintain a queue of path

Algorithm 4: BruteForcePathRepair

```

Input : Rule  $r$ 
         Path  $P$ 
         Datapoints to fix  $\mathcal{X}_P$ 
         Expected labels for datapoints  $\mathcal{Z}_P$ 
Output: Repair sequence  $\Phi$  which fixes  $r$  wrt.  $\mathcal{Z}_P$ 

1  $todo \leftarrow [(r, \emptyset)]$ 
2  $\mathcal{P}_{all} = \text{GetAllCandPredicates}(P, \mathcal{X}_P, \mathcal{Z}_P)$ 
3 while  $todo \neq \emptyset$  do
4    $(r_{cur}, \Phi_{cur}) \leftarrow \text{pop}(todo)$ 
5   foreach  $P_{cur} \in \text{leafpaths}(r_{cur}, P)$  do
6     foreach  $p \in \mathcal{P}_{all} - \mathcal{P}_{r_{cur}}$  do
7       foreach  $y_1 \in \mathcal{Y} \wedge y_1 \neq \text{last}(P_{cur})$  do
8          $\phi_{new} \leftarrow \text{refine}(r_{cur}, P_{cur}, y_1, p, \text{true})$ 
9          $r_{new} \leftarrow \phi_{new}(r_{cur})$ 
10         $\Phi_{new} \leftarrow \Phi_{cur}, \phi_{cur}$ 
11        if  $\text{Acc}(r_{new}, \mathcal{Z}_P) = 1$  then
12          return  $\Phi_{new}$ 
13        else
14           $todo.push((r_{new}, \Phi_{new}))$ 

```

Algorithm 3: GreedyPathRepair

```

Input : Rule  $r$ 
         Path  $P$ 
         datapoints to fix  $\mathcal{X}_P$ 
         Expected labels for assignments  $\mathcal{Z}_P$ 
Output: Repair sequence  $\Phi$  which fixes  $r$  wrt.  $\mathcal{Z}_P$ 

1  $todo \leftarrow [(P, \mathcal{Z}_P)]$ 
2  $\Phi = []$ 
3 while  $todo \neq \emptyset$  do
4    $(P, \mathcal{Z}_P) \leftarrow \text{pop}(todo)$ 
5   if  $\exists x_1, x_2 \in \mathcal{X}_P : \mathcal{Z}_P(x_1) \neq \mathcal{Z}_P(x_2)$  then
6     /* Determine predicates that distinguish assignments
7       that should receive different labels for a path */
8      $p \leftarrow \text{GetSeperatorPred}(x_1, x_2)$ 
9      $y_1 \leftarrow \mathcal{Z}_P(x_1)$ 
10     $\phi \leftarrow \text{refine}(r_{cur}, P, y_1, p, \text{true})$ 
11     $\mathcal{X}_1 \leftarrow \{x \mid x \in \mathcal{X}_P \wedge p(x)\}$ 
12     $\mathcal{X}_2 \leftarrow \{x \mid x \in \mathcal{X}_P \wedge \neg p(x)\}$ 
13     $todo.push((P[r_{cur}, x_1], \mathcal{X}_1))$ 
14     $todo.push((P[r_{cur}, x_2], \mathcal{X}_2))$ 
15  else
16     $\phi \leftarrow \text{refine}(r_{cur}, P, \mathcal{Z}_P(x))$ 
17     $r_{cur} \leftarrow \phi(r_{cur})$ 
18     $\Phi.append(\phi)$ 
19 return  $\Phi$ 

```

and datapoint set pairs which tracks which combination of paths

and datapoint sets still have to be fixed. This queue is initialized with P and all datapoints \mathcal{X}_P for P . The algorithm processes sets of datapoints until the todo queue is empty. In each iteration, the algorithm greedily selects a pair of datapoints x_1 and x_2 ending in this path that should be assigned different labels (line 5). It then calls method `GetSeperatorPred` (line 7) to determine a predicate p which evaluates to true on x_1 and false on x_2 (or vice versa). If we extend path P with p , then x_1 will follow the **true** edge of p and x_2 will follow the **false** edge (or vice versa). This effectively partitions the set of datapoints for path P into two sets \mathcal{X}_1 and \mathcal{X}_2 where \mathcal{X}_1 contains x_1 and \mathcal{X}_2 contains x_2 . We then have to continue to refine the paths ending in the two children of p wrt. these sets of datapoints. This is ensured by adding these sets of datapoints with their new paths to the todo queue (lines 12 and 13). If the current set of datapoints does not contain two datapoints with different labels, then we know that all remaining datapoints should receive the same label. The algorithm picks one of these datapoints x (line 14) and changes the current leaf node's label to $\mathcal{Z}_P(x)$.

Generating Predicates. The implementation of `GetCoveringPred` is specific to the type of PWSS. In [18] we present implementations of this procedure for weak supervised labeling that exploit the properties of these two application domains. However, note that, as we have shown in [18], as long as the space of predicates for an application domain contains equality and inequality comparisons for the atomic elements of datapoints, it is always possible to generate a predicate for two datapoints such that only one of these two datapoints fulfills the predicate. The algorithm splits the datapoint set \mathcal{X}_P processed in the current iteration into two subsets, which each are strictly smaller than \mathcal{X}_P . Thus, the algorithm is guaranteed to terminate and by construction assigns each datapoints x in \mathcal{X}_P its desired label $\mathcal{Z}_P(x)$.

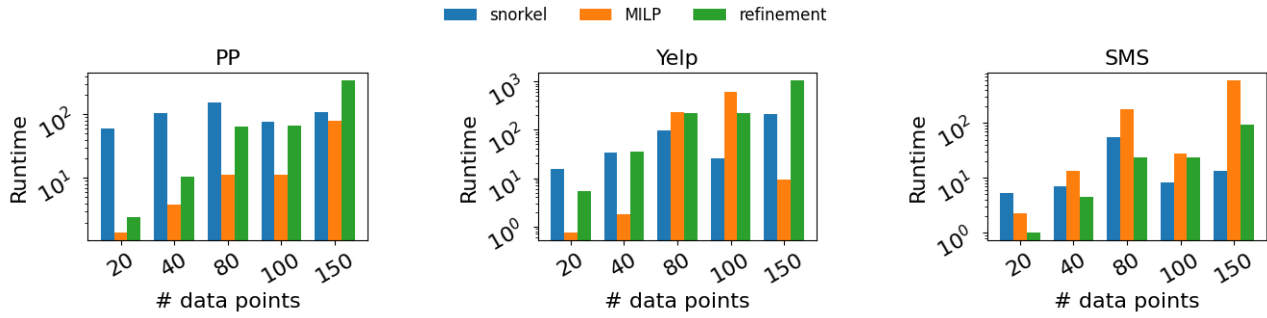
A.2 BruteForcePathRepair

The brute-force algorithm (Algorithm 4) is optimal, i.e., it returns a refinement of minimal cost (number of new predicates added). This algorithm enumerates all possible refinement repairs for a path P . Each such repair corresponds to replacing the last element on P with some rule tree. We enumerate such trees in increasing order of their size and pick the smallest one that achieves perfect accuracy on \mathcal{Z}_P . We first determine all predicates that can be used in the candidate repairs. As shown in [18], there are only finitely many distinct predicates (up to equivalence) for a given set \mathcal{X}_P . We then process a queue of candidate rules, each paired with the repair sequence that generated the rule. In each iteration, we process one rule from the queue and extend it in all possible ways by replacing one leaf node, and selecting the refined rule with minimum cost that satisfies all assignments. As we generate subtrees in increasing size, as shown in [18], the algorithm will terminate and its worst-case runtime is exponential in $n = |\mathcal{X}_P|$ as it may generate all subtrees of size up to n .

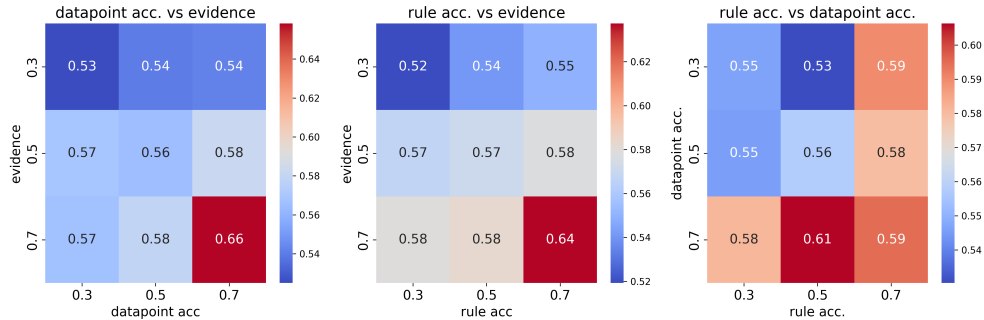
B Additional Experiment Details

B.1 runtime breakdown

The runtime breakdown for the remaining datasets from Section 5.2 are shown in Figure 9a.



(a) Runtime, varying the size of λ^* .



(b) Pairwise interaction heatmaps for New Global Acc.

Figure 9: Additional experimental results

dataset	repairer	fix%	preserv%	global acc.	new global acc.
<i>FNews</i>	RC	1	1	0.71	0.92
<i>FNews</i>	LLM	0.9	0.65	0.71	0.81
<i>Amazon</i>	RC	1	1	0.6	0.77
<i>Amazon</i>	LLM	0.9	0.5	0.6	0.7

Table 5: LLM vs RULECLEANER (RC) quality rule refinement comparison

B.2 MILP thresholds

The effects on global accuracy of the pairwise relationships of τ_{acc} , τ_E , and τ_{racc} are shown in Figure 9b. The color of a square represents global accuracy after the repair. Based on these results, it is generally preferable to set the thresholds higher as discussed in Section 2.2. However, larger thresholds reduce the amount of viable solutions to the MILP and, thus, can significantly increase the runtime of solving the MILP and lead to overfitting to λ^* .

C Repairing LFs with LLMs

In this section, we compare RULECLEANER against a baseline using a large language model (LLM). We designed a prompt instructing the LLM to act as an assistant and refine the LFs given a set of labeled datapoints. In this experiment, we used the *FNews* dataset with 40 labeled datapoints and *Amazon* dataset with 20 labeled examples.

We used GPT-4-turbo as the LLM. A detailed description of the prompt and responses from the LLM are presented in the [18].

We manually inspected the rules returned by the LLM to ensure that they are semantically meaningful. The quality of results after running Snorkel with the refined LFs from LLM and RULECLEANER are shown in Table 5. *fix%* measures the percentage of the wrong predictions by Snorkel that are fixed after retraining Snorkel with the refined rules. *preserv%* measures the percentage inputs correctly predicted by Snorkel that remain valid after retraining with the refined rules. RULECLEANER outperforms the LLM in both global accuracy and accuracy on labeled input data. We observe that the LLM tends to preserve the semantic meaning of the original LFs in the repairs it produces. For example, in one of the rules from *FNews*, the original rule is `if "talks" in text: REAL else ABSTAIN` and the repaired rule was `if any(x in text for x in ['discussions', 'negotiations', 'talks']):`. In one of the rules from *Amazon*, the original rule is `if any(x in text for x in ['junk', 'disappointed', 'useless']):NEGATIVE else ABSTAIN` mainly covers negative reviews. The LLM did add more negative words such as 'defective' whereas RULECLEANER could possibly add opposite sentiment conditions based on the solutions provided by the MILP. For example, it is possible for RULECLEANER to refine a rule with negative sentiment by adding `else if 'great' in text: POSITIVE else ABSTAIN`. The returned refined functions for *FNews* can be found in [18].