

PUER: Boosting Few-shot Positive-Unlabeled Entity Resolution with Reinforcement Learning

Yaoshu Wang¹, Mengyi Yan^{2*}, Wei Wang³

¹Shenzhen Institute of Computing Sciences, ²Shandong University

³The Hong Kong University of Science and Technology, Guangzhou

¹yaoshuw@sics.ac.cn, ²yanmy@sdu.edu.cn, ³weiwcs@ust.hk

Abstract

Entity resolution is a fundamental problem in data management that aims to identify all duplicate entries within collections of multi-attribute tuples. Most existing works focus on supervised learning, relying on large amounts of high-quality labeled data, including both positive and negative tuple pairs that are meticulously prepared. However, in reality, the manual annotation process is labor-intensive; in particular, selecting high-quality negative data for labeling is both important and challenging. In this paper, we propose an end-to-end ER solution, PUER, to address low-resource entity resolution (ER) by leveraging Large Language Models (LLMs) in a Positive-Unlabeled (PU) learning setting, where only a small number of positively labeled examples, *e.g.*, 50, and unlabeled data are provided. Unlike directly fine-tuning LLMs in a supervised manner, we solve the entity matching task using reinforcement learning and propose a self-adaptive reward function in the process of RL. To enhance performance, we design an iterative workflow based on the co-training mechanism that fully utilizes entity blocking component to assist the entity matching. This workflow aims to improve the robustness and quality of pseudo-labels so that the performance of entity matching is improved. Comprehensive experimental results on various benchmark datasets demonstrate the superiority of PUER. Full version and code are available¹.

1 Introduction

Entity resolution (ER) aims to identifying all tuple pairs from two relational tables that refer to the same entities, making it a key components of data cleaning with the goal of deduplicating records in datasets. Traditionally, the ER task consists of two components, entity blocking (EB) and entity matching (EM). Entity blocking efficiently retrieves potentially matched tuple pairs, while entity matching

verifies whether these tuple pairs refer to the same entities.

Traditionally, the entity resolution has been extensively studied and most of solutions reply on a sufficient number of annotated tuple pairs to achieve good performance. However manual annotation is costly, as demonstrated by methods like Ditto (Li et al., 2020b) and DeepMatcher (Mudgal et al., 2018). To address it, a few existing EM approaches focus on unsupervised learning, semi-supervised learning and active learning. For instance, TDmatch (Ahmadi et al., 2022) is an unsupervised ER approach based on graph creation and random walk, while only relying on the data distribution cannot have very high accuracy due to the extreme class imbalance. PromptEM (Wang et al., 2022) generates pseudo-labels for low-resource ER in the semi-supervised learning. While it partially alleviates the annotation cost, selecting and labeling high-quality positive and negative tuple pairs from large datasets remains challenging. Active learning approaches, *e.g.*, (Arasu et al., 2010), select ambiguous tuple pairs for user labeling, but this also incurs significant manual annotation costs. In this work, we focus on the few-shot Positive-Unlabeled (PU) learning, where only a small number of labeled positive tuple pairs are provided along with two relational tables. To the best of our knowledge, we are the first to explore the ER task in the few-shot PU learning context.

Entity resolution, which typically replies on both positive and negative training data, is particularly relevant to few-shot PU learning (Bekker and Davis, 2020). In practice, only users who detect duplicate issues in their datasets often invest resources to find and integrate these duplicates, so that these detected duplicates instances are treated as labeled positive data. Negative tuple pairs are generally not provided. Additionally, in a search engine scenario, users might ask a question and receive multiple semantically identical re-

* Corresponding author

¹<https://github.com/authurlord/PUER>

sponses (Niu et al., 2016). These responses are considered positive data. While negative tuple pairs can be relatively easy to collect, acquiring high-quality negative pairs in large-scale datasets is non-trivial and requires manual annotation and verification, presenting a challenge for annotators (Wang et al., 2024b). Building on these observations, we investigate the ER task within the framework of few-shot PU learning to address these challenges effectively with minimal labelling cost, while aligning with human preference.

However, existing ER methods based on pre-trained Language Models (e.g., RoBERTa) primarily learn distribution and decision boundaries from large sets of annotated samples. This approach results in inefficiencies in labeled data utilization, placing these methods at a disadvantage in few-shot scenarios. Furthermore, they lack the capability to generalize well, and cannot achieve high-performance EM tasks based solely on a limited number of positive sample annotations. Recently, as the advanced performance of large language models (LLMs), they have been explored in the entity matching task. JellyFish (Zhang et al., 2024) addresses various data pre-processing tasks, including entity matching, by leveraging LLMs in the instruction-tuning and reasoning manner. TableGPT (Li et al., 2024b) employs and fine-tunes on-line GPT-3.5 on various data pre-processing tasks. Considering the data privacy concerns, we focus on using local LLMs that are open-sourced and can be fine-tuned in local environments.

Motivated by the above considerations, we explore the ER problem within the framework of few-shot PU learning by harnessing the capabilities of local LLMs. Our objective is to fully utilize entity blocking to assist the entity matching process and to develop an LLM-based model that can efficiently and effectively retrieve all matched tuple pairs from limited labeled data. Unlike traditional methods that treat entity matching purely as a binary classification task, our approach is the first to formulate entity matching as a reinforcement learning problem while simultaneously fine-tuning the model using both Supervised Fine-Tuning (SFT) and reinforcement learning. Furthermore, to integrate entity blocking with entity matching, we introduce an iterative workflow that progressively generates high-quality pseudo-labels, facilitating mutual learning among these components.

Our contributions are as follows:

- Beyond binary classification, we are the first to employ **reinforcement learning** to solve the entity matching, and design a self-adaptive reward function to enhance the convergence speed.
- We propose an **end-to-end** entity resolution workflow that iteratively make full use of the entity blocking model to select high-quality training data and jointly fine-tunes two entity matching models through a co-training mechanism.
- We conduct **comprehensive experiments** to evaluate the efficiency and effectiveness of our approach, demonstrating its superiority over existing methods.

2 Related Work

We classify ER into entity blocking and matching.

Entity blocking. We classify entity blocking as (1) rules, *e.g.*, handcrafted rules (Papadakis et al., 2020, 2014; Fan et al., 2009; Kejriwal and Miranker, 2015), and learned rules (Michelson and Knoblock, 2006; Kejriwal and Miranker, 2015; Singh et al., 2017a; Paulsen et al., 2023), (2) traditional ML, *e.g.*, (C. et al., 2018; Efthymiou et al., 2015), and (3) deep learning, *e.g.*, (Thirumuranathan et al., 2021; Brinkmann et al., 2024; Wang et al., 2023; Reimers and Gurevych, 2019; Wu et al., 2023; Wang et al., 2024a), which retrieve potentially matched tuple pairs from large-scale datasets.

Entity matching. There are host of works on entity matching, including rule-based methods (Guo et al., 2010; Fan et al., 2011; Whang and Garcia-Molina, 2013; Singh et al., 2017b), ML-based methods (Konda et al., 2016; Bilenko and Mooney, 2003; Wu et al., 2020) and deep learning-based methods (Li et al., 2020b; Mudgal et al., 2018; Ebraheem et al., 2018; Zhao and He, 2019; Li et al., 2020a; Fu et al., 2019). Recently low resource entity matching based on deep learning models has been paid attention, including (1) active learning ER, *e.g.*, (Qian et al., 2017; Meduri et al., 2020; Kasai et al., 2019; Nafa et al., 2022), (2) data augmentation ER, *e.g.*, Rotom (Miao et al., 2021), (3) unsupervised learning ER, *e.g.*, (Zeng et al., 2024; Ahmadi et al., 2022), (4) transfer learning ER, *e.g.*, (Kirielle et al., 2022; Tu et al., 2022; Sun et al., 2024; Loster et al., 2021), (5) semi-supervised learning ER, *e.g.*, PromptEM (Wang et al., 2022), (6) multi-task learning, *e.g.*, Unicorn (Fan et al., 2024a). and (7) information fusion (Yao et al.,

2021). There are also works to combine the entity blocking and matching models for mutual learning, *e.g.*, (Wu et al., 2023; Wang et al., 2023; Li et al., 2021), and works by leveraging local LLMs (Zhang et al., 2024; Wadhwa et al., 2024) and online LLMs (Li et al., 2024b; Wang et al., 2025; Li et al., 2024a; Fan et al., 2024b). However, none of the above works address few-shot the Positive-Unlabeled setting, such that only as small number of positive instances are given, which is more practical in real-life.

3 Preliminaries

In this section, we first present the ER problem, and then introduce the entity blocking and matching.

3.1 Problem Formulation

Given two relational tables of multi-attribute tuples, the goal of entity resolution (ER) is to identify pairs of tuples that refer to the same entity. The ER task generally consists of two main components: entity blocking and entity matching. The entity blocking component efficiently retrieves a candidate set of potentially matching tuple from large tables, thereby avoiding the quadratic time complexity of comparing all tuple pairs between relational tables. The entity matching component then predicts whether tuple pairs in the candidate set are matches.

Definition 1: (ER under the few-shot positive-unlabeled setting.) Given two relational tables of multi-attribute tuples \mathcal{R}_l and \mathcal{R}_r , and a set \mathcal{P} consists of a small number of positive tuple pairs, the objective of few shot PU entity resolution (ER) is to identify all matching tuple pairs from $\mathcal{R}_l \times \mathcal{R}_r$. \square

In this paper, we mainly focus on addressing entity matching task of ER. We fully utilize existing entity blocking techniques to enhance the efficiency and effectiveness of the matching process.

3.2 Entity Resolution

Following previous work (Wu et al., 2023), we decompose entity resolution into entity blocking and entity matching.

Entity blocking. The entity blocking Blocker, primarily utilizes the SentenceBert model, denoted as \mathcal{F}_{RAG} , to transform each tuple t into an embedding vector (Thirumuruganathan et al., 2021; Wang et al., 2023; Wu et al., 2023; Li et al., 2020b). Given

the training data $(t, \mathcal{P}_t, \mathcal{N}_t)$, where \mathcal{P}_t and \mathcal{N}_t represent the positive and negative sets of tuples that match and mismatch with t , we fine-tune \mathcal{F}_{RAG} via contrastive learning (Oord et al., 2018).

Entity matching. Previous work (Wang et al., 2025) formulates the entity matching into two sub-tasks: Matcher $\mathcal{F}_{\text{EM}}^{\text{M}}$ and Selector $\mathcal{F}_{\text{EM}}^{\text{S}}$.

Matcher. Given a tuple pair (t, s) and a domain-specific prompt pt_m as EM instruction, we could query LLM to transform (t, s) to a binary decision $p_m \in \{\text{Yes}, \text{No}\}$, *s.t.* $p_m = \text{LLM}(\text{pt}_m, (t, s))$. The Matcher $\mathcal{F}_{\text{EM}}^{\text{M}}$ is consistent with all existing entity matching works, *e.g.*, JellyFish (Zhang et al., 2024), Ditto (Li et al., 2020b), aiming to identifying whether a tuple pair is matched or not. Matcher is supervised fine-tuned (SFT) with LoRA and aims to provide domain-specific decision boundary.

Selector. Selector takes a pivotal tuple t , a list of candidate tuples $\mathcal{C}_s(t) = \{s_1, \dots, s_{|\mathcal{C}_s|}\}$, and a prompt pt_s as inputs, and outputs a list of positive ones in $\mathcal{C}_s(t)$. It lets LLMs check more examples so that they make correct decision. Selector targets at re-ranking \mathcal{C}_s by simulating human preference.

The Matcher subtask is mainly used in most EM approaches, *e.g.*, JellyFish, while Selector has not been as extensively studied. Although (Wang et al., 2025) introduced it to address the EM, they did not further fine-tune it to improve its performance.

4 RL-based Entity Matching

As discussed above, the entity matching task is formulated as two sub-tasks, Matcher $\mathcal{F}_{\text{EM}}^{\text{M}}$ and Selector $\mathcal{F}_{\text{EM}}^{\text{S}}$. In this section, we focus on how to fine-tune the Selector to make policies from a list of candidate tuples. Here we assume that $\mathcal{F}_{\text{EM}}^{\text{S}}$ is fine-tuned using a (pseudo-)labeled training dataset $\mathcal{D}_{\text{train}}^{\text{S}} = \{(t, \mathcal{C}_s(t), \mathbf{L}_t)\}$, where \mathbf{L}_t is the label of the pivotal tuple t . The process of generating $\mathcal{D}_{\text{train}}^{\text{S}}$ will be discussed in Section 5.

RL-based Selector. To select matching tuples from a pivotal tuple, we employ the Group Relative Policy Optimization (GRPO) (DeepSeek-AI et al., 2025) to fine-tune the Selector so that it can better adapt to the dynamic changes and improve its accuracy. However, the number of positive tuples in the candidate list is very small, resulting in a continuously low reward value during the learning process of GRPO. To address this issue, we design a self-adaptive reward model.

GRPO. Given each pivotal tuple t from $\mathcal{D}_{\text{train}}^S$ that follows the distribution P , we adopt GRPO (DeepSeek-AI et al., 2025) with the following loss function.

$$\mathcal{J}_{\text{GRPO}} = \mathbb{E}[t \sim P(\mathcal{D}_{\text{train}}^S), \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(O|t)] \frac{1}{G} \sum_{i=1}^G \left(\min \left(\frac{\pi_{\theta}(o_i|t)}{\pi_{\theta_{\text{old}}}(o_i|t)} A_i, \text{clip} \left(\frac{\pi_{\theta}(o_i|t)}{\pi_{\theta_{\text{old}}}(o_i|t)}, 1 - \epsilon, 1 + \epsilon \right) A_i \right) - \eta \mathbb{D}_{\text{KL}}(\pi_{\theta} || \pi_{\text{ref}}) \right)$$

where A_i is the advantage function computed within a group of rewards. The Selector $\mathcal{F}_{\text{EM}}^S$ is fine-tuned in two stages. First, to address the cold start problem, we initialize $\mathcal{F}_{\text{EM}}^S$ using SFT, which enables it to adapt to the selection task. Subsequently, $\mathcal{F}_{\text{EM}}^S$ is further fine-tuned using GRPO with the $\mathcal{J}_{\text{GRPO}}$ loss function.

Input and Output Formats of $\mathcal{F}_{\text{EM}}^S$. Given a handcrafted instruction pt_s , a tuple t , and a candidate list $\mathcal{C}_s(t)$, we first formulate them into a final prompt following (Wang et al., 2025). We then feed this prompt into the function $\mathcal{F}_{\text{EM}}^S$. Then $\mathcal{F}_{\text{EM}}^S$ subsequently produces the response r .

$r = \langle \text{positive} \rangle [\dots] / \langle \text{positive} \rangle \langle \text{negative} \rangle [\dots] / \langle \text{negative} \rangle$

Here, the lists of positive and negative tuple IDs from $\mathcal{C}_s(t)$ are enclosed within the markers $\langle \text{positive} \rangle$ and $\langle \text{negative} \rangle$, respectively. Here notice that we also let $\mathcal{F}_{\text{EM}}^S$ return negative tuples to make sure that it also focuses on negative ones.

Self-adaptive Reward Function. Given a tuple t , the candidates $\mathcal{C}_s(t)$, and the label vector \mathbf{L}_t , we design a reward function R that returns a scalar reward value for RL, where $\mathbf{L}_t \in \{0, 1\}^{|\mathcal{C}_s(t)|}$ is a binary vector indicating whether each candidate in $\mathcal{C}_s(t)$, e.g., the i -th element in $\mathcal{C}_s(t)$, is a true match ($\mathbf{L}_t[i] = 1$) or not ($\mathbf{L}_t[i] = 0$). The reward function R contains the following steps.

(1) Step 1: Answer Extraction. We handcraft the regular expression to extract the list L_{pos} of positive tuple IDs and the list L_{neg} of negative ones from the response r of $\mathcal{F}_{\text{EM}}^S$. We return a zero reward if L_{pos} or L_{neg} cannot be parsed successfully from r , if the tuple IDs in L_{pos} or L_{neg} are not within the range $[1, |\mathcal{C}_s(t)|]$, or if $|L_{\text{pos}}| + |L_{\text{neg}}| \neq |\mathcal{C}_s(t)|$. In other words, the answer extracted from the response must be valid. If this condition is met, we proceed to Step 2; otherwise, a zero reward is returned.

Input: a collection of training data $\mathcal{D}_{\text{train}}^S = \{(t, \mathcal{C}_s(t), \mathbf{L}_t)\}$, the number of iteration iter_{max} , the smoothing factor α .

Output: the policy π_{EM} .

1. Split $\mathcal{D}_{\text{train}}^S$ into training data $\mathcal{D}_{\text{train}}^S$ and validation data $\mathcal{D}_{\text{valid}}^S$;
2. SFT $\mathcal{F}_{\text{EM}}^S$ in $\mathcal{D}_{\text{train}}^S$ as the cold start.
3. $\text{iter} := 0, w_{\text{pos}}^{(0)} = 1, w_{\text{neg}}^{(0)} = 1$;
4. **while** $\text{iter} \leq \text{iter}_{\text{max}}$ **do**
5. The reward $R = \mathcal{H}_w(\mathbf{P}_t, \mathbf{L}_t, w_{\text{pos}}^{(\text{iter})}, w_{\text{neg}}^{(\text{iter})})$;
6. Fine-tune $\mathcal{F}_{\text{EM}}^S$ in $\mathcal{D}_{\text{train}}^S$ with GRPO using the reward R ;
7. Compute the prediction $\mathbf{P}_{\text{valid}} = \mathcal{F}_{\text{EM}}^S(\mathcal{D}_{\text{valid}}^S)$;
8. Compute FN and FP by comparing $\mathbf{P}_{\text{valid}}$ and $\mathbf{L}_{\text{valid}}$;
9. $w_{\text{pos}} = \frac{\text{FP} + \epsilon}{\text{FN} + \text{FP} + \epsilon}, w_{\text{neg}} = \frac{\text{FN} + \epsilon}{\text{FN} + \text{FP} + \epsilon}$;
10. $w_{\text{pos}}^{(\text{iter}+1)} := (1 - \alpha)w_{\text{pos}}^{(\text{iter})} + \alpha w_{\text{pos}}$;
11. $w_{\text{neg}}^{(\text{iter}+1)} := (1 - \alpha)w_{\text{neg}}^{(\text{iter})} + \alpha w_{\text{neg}}$;
12. $\text{iter} := \text{iter} + 1$;
13. **return** $\mathcal{F}_{\text{EM}}^S$;

Figure 1: RL-based Selector

(2) Step 2: Similarity reward computation. Intuitively, we aim to measure the similarity between the answer and the ground truth. Hamming similarity is a good option.

$$R(t, \mathcal{C}_s(t), \mathbf{L}_t) = \mathcal{H} \left(\text{Enc}(L_{\text{pos}}, L_{\text{neg}}), \mathbf{L}_t \right)$$

where Enc is a handcrafted encoding function that transforms L_{pos} and L_{neg} into a binary vector \mathbf{P}_t of the same dimension as \mathbf{L}_t . \mathcal{H} represents the Hamming similarity, s.t. $\mathcal{H}(\mathbf{P}_t, \mathbf{L}_t) = \frac{\sum_{i=1}^{|\mathcal{C}_s(t)|} \mathbf{1}_{\mathbf{P}_t[i] = \mathbf{L}_t[i]}}{|\mathcal{C}_s(t)|}$. A higher $\mathcal{H}(\mathbf{P}_t, \mathbf{L}_t)$ indicates better performance of $\mathcal{F}_{\text{EM}}^S$, while a lower value suggests suboptimal performance.

However, directly using \mathcal{H} has the following drawbacks. First, the ratio of positive tuples in \mathcal{C}_s is very small, and the rewards from negative tuples would dominate the exploration process, causing the feedback from the reward function to remain at a very low value. Second, the goal of entity matching is to reduce both false positives (FPs) and false negatives (FNs). When the number of FPs increases, we expect $\mathcal{F}_{\text{EM}}^S$ to focus on reducing FPs; otherwise, it should focus on reducing FNs. The current \mathcal{H} does not encourage this behavior in $\mathcal{F}_{\text{EM}}^S$, causing it to spend a large number of iterations exploring unseen regions.

To address these issues, we design a weighted Hamming similarity \mathcal{H}_w :

$$\mathcal{H}_w(\mathbf{P}_t, \mathbf{L}_t, w_{\text{pos}}, w_{\text{neg}}) = \frac{\sum_{i=1}^{|\mathcal{C}_s(t)|} \mathbf{1}_{\mathbf{P}_t[i] = \mathbf{L}_t[i] = 1} w_{\text{pos}} + \mathbf{1}_{\mathbf{P}_t[i] = \mathbf{L}_t[i] = 0} w_{\text{neg}}}{\sum_{i=1}^{|\mathcal{C}_s(t)|} \mathbf{1}_{\mathbf{L}_t[i] = 1} w_{\text{pos}} + \mathbf{1}_{\mathbf{L}_t[i] = 0} w_{\text{neg}}}$$

For the positive tuples in \mathbf{L}_t , we assign a weight w_{pos} , and for the negative tuples, we assign a

weight w_{neg} . By integrating these weights into the reward function, $\mathcal{F}_{\text{EM}}^S$ is more inclined to focus on one side, *i.e.*, either positive or negative data.

The final problem is how to set the values of w_{pos} and w_{neg} . Our idea is that if $\mathcal{F}_{\text{EM}}^S$ has an increasing number of false positives, we should increase the value of w_{pos} so that GRPO focuses on reducing the false positives. Otherwise, we encourage GRPO to find true positives from the negative tuples. To achieve this, we split $\mathcal{D}_{\text{train}}^S$ into validation data $\mathcal{D}_{\text{valid}}^S$ and compute the false positives (FPs) and false negatives (FNs) in each iteration. Let $w_{\text{pos}}^{(i)}$ and $w_{\text{neg}}^{(i)}$ be the weights for the i -th iteration. We use $\mathcal{F}_{\text{EM}}^S$ to make predictions $\mathbf{P}_{\text{valid}}$ on $\mathcal{D}_{\text{valid}}^S$, and then compute FPs and FNs. The current weights w_{pos} and w_{neg} are set to the percentages of FPs and FNs, respectively. However, resetting these weights in each iteration would lead to an unstable reward function. To gradually change these values, we introduce a smoothness factor α to update the weights with small adjustments. Specifically, we set $w_{\text{pos}}^{(i+1)} := (1 - \alpha)w_{\text{pos}}^{(i)} + \alpha w_{\text{pos}}$ and $w_{\text{neg}}^{(i+1)} := (1 - \alpha)w_{\text{neg}}^{(i)} + \alpha w_{\text{neg}}$.

To further reinforce the impact of positive tuples in \mathbf{L}_t , we incorporate semantic similarity into the reward function as prior knowledge. This guides the RL process to find a good direction. Our final reward is as follow.

$$R(t, \mathcal{C}_s(t), \mathbf{L}_t) = \mathcal{H}_w(\mathbf{P}_t, \mathbf{L}_t, w_{\text{pos}}, w_{\text{neg}}) + \beta \cdot \frac{1}{|S|} \sum_{s \in S} \text{Sim}_{\text{cos}}(\text{vec}(t), \text{vec}(s))$$

where β is a hyper-parameter and 0.2 by default, S is the set of true positives in the prediction of $\mathcal{F}_{\text{EM}}^S$, Sim_{cos} is the cosine similarity between two vectors, and $\text{vec}()$ is the embedding returned by \mathcal{F}_{RAG} .

Figure 1 illustrates the RL process of $\mathcal{F}_{\text{EM}}^S$. In addition to $\mathcal{D}_{\text{train}}^S$, the number of iterations iter_{max} and the smoothing factor α are added as inputs. Initially, we set $w_{\text{pos}}^{(0)}$ and $w_{\text{neg}}^{(0)}$ to 1, indicating the normal Hamming similarity (line 3). In each iteration, we re-formulate the reward function using \mathcal{H}_w (line 5) and fine-tune $\mathcal{F}_{\text{EM}}^S$ using GRPO with the reward function R (line 6). We then update the weights of positive and negative tuples using the gradual update rules (lines 7-11).

5 An ER Workflow

In this section, we present an ER workflow to train the entity matching models \mathcal{F}_{EM} with the assistance

of an entity blocking model \mathcal{F}_{RAG} . The workflow takes as input two relational tables, \mathcal{R}_l and \mathcal{R}_r , and a set \mathcal{P} of positive tuple pairs. As shown in Figure 2, the workflow consists of three main steps: data enrichment, entity blocking, and entity matching.

Step 1: Data enrichment. Enriching tuples in \mathcal{R}_l and \mathcal{R}_r with additional attributes, denoted as \bar{B} , is a common and effective method. Due to the uncertainty (Farquhar et al., 2024) inherent in LLMs, we observe that **the values of \bar{B} imputed for a tuple can vary depending on its paired tuples (*w.r.t.* context)**. For each tuple $t \in \mathcal{R}_l$ and a set $S_t \subset \mathcal{R}_r$, we generate $|S_t|$ tuple pairs, *i.e.*, $P_t = \{(t, s_1), \dots, (t, s_{|S_t|})\}$, where $s_i \in S_t$. Given each pair $(t, s_i) \in P_t$, we query the LLM to impute the values of \bar{B} as $a_i = \text{LLM}((t, s_i), \bar{B}, \text{pt}_{\text{enr}})$. Due to LLM uncertainty, the imputed values $a_1, \dots, a_{|P_t|}$ may differ across pairs. To leverage these variations, we enumerate all imputations and augment each tuple pair with multiple enriched versions.

Step 2: Entity Blocking. After data enrichment, we enrich the positive set \mathcal{P} into an augmented set \mathcal{P}_{enr} . We then fine-tune our entity blocking model \mathcal{F}_{RAG} using contrastive learning with a randomly sampled negative set, following the approach in (Wang et al., 2024a). The final output of this step is the fine-tuned \mathcal{F}_{RAG} .

Step 3: An iterative EM workflow. Given \mathcal{F}_{RAG} , \mathcal{P}_{enr} , \mathcal{R}_l and \mathcal{R}_r , we propose a progressive training workflow that fine-tunes $\mathcal{F}_{\text{EM}}^M$ and $\mathcal{F}_{\text{EM}}^S$.

Overview. We show the EM workflow. Given a tuple $t \in \mathcal{R}_l$, \mathcal{F}_{RAG} conducts similarity search by retrieving its K nearest neighbors $\text{NN}_K(t)$, which forms the candidate list for t . We define two pointers: ptr_s and ptr_e , where ptr_s indicates the boundary separating positive tuples from the rest, such that all tuples in the range $[1, \text{ptr}_s]$ are considered positive, while tuples in the range $[\text{ptr}_e, K]$ are considered negative. Specifically, $(t, \text{NN}_K(t)[i])$ are treated as positive tuple pairs for $i \in [1, \text{ptr}_s]$ and $(t, \text{NN}_K(t)[j])$ are negative ones for $j \in [\text{ptr}_e, K]$. $[\text{ptr}_s, \text{ptr}_e]$ are ambiguous pairs.

In the beginning of the training procedure, ptr_s is set to 1 and ptr_e is set to K , and \mathcal{F}_{RAG} generates the potentially positive and negative tuple pairs \mathcal{P}_{RAG} and \mathcal{N}_{RAG} within $[1, \text{ptr}_s]$ and $[\text{ptr}_e, K]$, respectively. Next \mathcal{F}_{RAG} sends them to \mathcal{F}_{EM} , which processes them using the co-training strategy. In the next iteration, \mathcal{F}_{RAG} retrieves the new $\text{NN}_K(t)$ for each tuple t and adjust ptr_s and ptr_e by a step

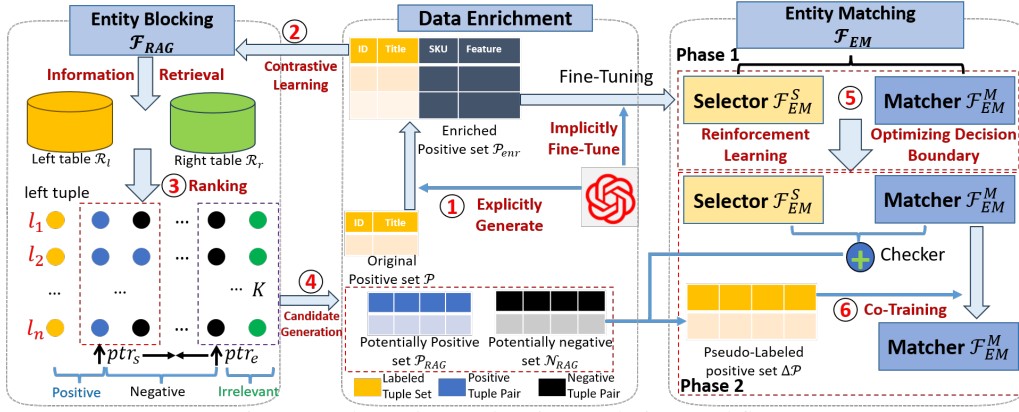


Figure 2: The end-to-end entity resolution workflow

size δ , updating the pointers as $\text{ptr}_s = \text{ptr}_s + \delta$ and $\text{ptr}_e = \text{ptr}_e - \delta$. The iterative process continues until ptr_s is no longer less than ptr_e .

Co-training strategy. Given potentially positive tuple pairs \mathcal{P}_{RAG} and negative tuple pairs \mathcal{N}_{RAG} , and the augmented set \mathcal{P}_{enr} , we simultaneously learn Matcher $\mathcal{F}_{\text{EM}}^{\text{M}}$ and Selector $\mathcal{F}_{\text{EM}}^{\text{S}}$. Considering the extremely low ratio of positive to negative tuple pairs in \mathcal{R}_l and \mathcal{R}_r , and the fact that existing methods, *e.g.*, (Thirumuruganathan et al., 2021), rely on random sampling for negative tuple pairs, we assume that \mathcal{N}_{RAG} are more likely to be correct in the first few iteration. Thus, we introduce a warmup period during which \mathcal{N}_{RAG} are initially treated as ground truth negatives and combined with \mathcal{P}_{enr} to fine-tune $\mathcal{F}_{\text{EM}}^{\text{M}}$ for the first λ iterations. As ptr_e approaches ptr_s after λ iterations, $\mathcal{F}_{\text{EM}}^{\text{M}}$ is then responsible for selecting which negative tuple pairs should be included in the training data.

Specifically, we design a two-phase learning method to simultaneously train $\mathcal{F}_{\text{EM}}^{\text{M}}$ and $\mathcal{F}_{\text{EM}}^{\text{S}}$.

Phase 1. In the first phase, we add \mathcal{P}_{enr} into the training data $\mathcal{D}_{\text{train}}$. If the current iteration is less than λ , we generate the training data $\mathcal{D}_{\text{train}}$ as $\mathcal{D}_{\text{train}} = \mathcal{P}_{\text{enr}} \cup \mathcal{N}_{\text{RAG}}$. For iterations beyond λ , we have a checker step by using $\mathcal{F}_{\text{EM}}^{\text{M}}$ to verify whether the labels of tuple pairs in \mathcal{N}_{RAG} are consistent with its predictions. The training data $\mathcal{D}_{\text{train}}$ is then updated to $\mathcal{D}_{\text{train}} = \mathcal{P}_{\text{enr}} \cup \{(t, s) | \mathcal{F}_{\text{EM}}^{\text{M}}(t, s) = \text{No}, (t, s) \in \mathcal{N}_{\text{RAG}}\}$.

After generating $\mathcal{D}_{\text{train}}$, we proceed to generate the training data $\mathcal{D}_{\text{train}}^{\text{S}}$ for the Selector. For each pivot tuple t , we retrieve all tuple pairs $(t, s_1), \dots, (t, s_L)$ from $\mathcal{D}_{\text{train}}$ where t appears on the left-hand side of the tuple pairs. We then define $\mathcal{C}_s(t) = \{s_1, \dots, s_L\}$, and the label \mathbf{L}_t is an L -dimensional vector, where $\mathbf{L}_t[i] = 1$ if (t, s_i) is a matched tuple, and $\mathbf{L}_t[i] = 0$ otherwise. We gen-

erate a triplet $(t, \mathcal{C}_s(t), \mathbf{L}_t)$ as an element of $\mathcal{D}_{\text{train}}^{\text{S}}$. Finally, we fine-tune $\mathcal{F}_{\text{EM}}^{\text{M}}$ using SFT on $\mathcal{D}_{\text{train}}$ and $\mathcal{F}_{\text{EM}}^{\text{S}}$ using GRPO on $\mathcal{D}_{\text{train}}^{\text{S}}$, respectively.

Phase 2. Once $\mathcal{F}_{\text{EM}}^{\text{S}}$ and $\mathcal{F}_{\text{EM}}^{\text{M}}$ are well fine-tuned, we adopt $\mathcal{F}_{\text{EM}}^{\text{S}}$ to generate more positive training data with pseudo-labels. For each tuple $t \in \mathcal{R}_l$, we first generate m augmented tuples, denoted by t_1, \dots, t_m , and query $\mathcal{F}_{\text{EM}}^{\text{S}}$ with the list $\mathcal{C}_s(t_i)$ for t_i for $i \in [1, m]$. This process yields m lists of positive tuples R_1, \dots, R_m . To enhance the accuracy of these positive instances, we use $\mathcal{F}_{\text{EM}}^{\text{M}}$ to make inferences, obtaining additional positive training data as $\Delta\mathcal{P} = \{(t, s) | \mathcal{F}_{\text{EM}}^{\text{M}}(t, s) = \text{Yes}, (t, s) \in R_1 \cup \dots \cup R_m\}$. Finally, we incorporate $\Delta\mathcal{P}$ into $\mathcal{D}_{\text{train}}$, and $\mathcal{F}_{\text{EM}}^{\text{M}}$ is continuously fine-tuned using SFT to further enhance its performance.

The ER workflow that relies on RL training can be time-consuming; nevertheless, they are usually carried out in an offline preprocessing phase. For example, deduplication of corpora executed offline ensures that search engines do not return duplicate or valueless results when an online query is given. Similarly, when a new data source is introduced, data integration is usually executed offline, allowing the newly integrated data to be used for downstream online business processes. Consequently, users mainly focus on the accuracy of detecting duplicates, and the workflow remains worthwhile as long as it does not introduce prohibitive delays.

6 Experimental Results

In this section, we empirically evaluated our method, PUER using benchmark datasets on (1) the effectiveness and efficiency of entity matching, (2) the ablation study, and (3) a comparative analysis with online LLMs, particular ComEM (Wang et al., 2025) with GPT-4o-mini. More experimental results are provided in the supplemental material.

Experimental settings. We start with our settings.

Datasets. We conducted experiments using 11 benchmark datasets from the ER Benchmark datasets (Köpcke et al., 2010), the Magellan data repository (The Magellan Data Repository) and WDC product data corpus (Primpeli et al., 2019) used for evaluating Ditto (Li et al., 2020b). These datasets include Amazon-Google (AG), Walmart-Amazon (WA), Abt-Buy (AB), DBLP-ACM (DA), DBLP-Scholar (DS), Company (CO), Cameras (CA), Computers (COM), Shoes (SH), Watch (WAT) and WDC-All-Small (WS). Following Ditto, we randomly sample 50 positive tuple pairs as labeled training data and retained the left and right relational tables as all unlabeled data. The statistics of all datasets are summarized in the supplementary material.

Baselines. We implemented PUER in Python and used the following baselines. (1) Ditto (Li et al., 2020b), an entity matching model based on BERT; (2) Rotom (Miao et al., 2021), an entity matching model leveraging language models and data augmentation through RL; (3) PromptEM (Wang et al., 2022), a prompt-tuning model based on pretrained language models; (4) Unicorn (Fan et al., 2024a), a multi-task data matching model using a mixture of experts; (5) CLER (Wu et al., 2023), a low-resource entity resolution model that integrates entity blocking and matching; (6) JellyFish (Zhang et al., 2024), an LLM based entity matching model using LoRA-based instruction-tuning; (7) Sudowoodo (Wang et al., 2023), an entity resolution framework based on contrastive representation learning. We also compared with the following online LLMs: (8) BatchER (Fan et al., 2024b), a cost-effective batch prompting to ER based on online LLMs, and (9) ComEM (Wang et al., 2025), an LLM-based ER model using Matcher, Comparer and Selector.

For a fair comparison, all baselines are provided with the same set of 50 positive tuple pairs (denoted as \mathcal{P}), all labeled negative pairs from the training set of the benchmark, and the same relational left and right tables. In contrast, PUER is provided with \mathcal{P} and relational left and right tables but without the labeled negative pairs, representing a few-shot PU (positive-unlabeled) setting. Notably, Unicorn, and JellyFish were also pretrained on additional labeled entity matching corpus.

Measures. We report precision (P), recall (R) and F1 (F) score for entity matching following (Li et al.,

2020b). All results are reported in 100-scale.

Configuration. We select Qwen-2.5-7B-instruct (Yang et al., 2024) as the backbone of \mathcal{F}_{EM} and bge-large-en as the pre-trained model for \mathcal{F}_{RAG} . We set K as 20, δ as 5, and τ as 0.02 by default and adopt the AdamW optimizer with the learning rate of $1e-4$ and $1e-5$ for \mathcal{F}_{EM}^M and \mathcal{F}_{RAG} , respectively. In GRPO of \mathcal{F}_{EM}^S , we set the training batch size as 16, the length of input prompt as 1024, the maximum output length as 64, the mini-batch as 16, the learning rate of the actor model as $1e-6$, and the coefficient η of KL loss as 0.001. We adopt verl (Sheng et al., 2025), a RL training framework to fine-tune \mathcal{F}_{EM}^S and remains other hyper-parameters by default. For all baselines, we use their default settings. We conduct our experiment on a single machine powered by 1.5TB RAM and 128 processors with Intel(R) Xeon(R) Platinum 8358 CPU @2.60GHz and 4 NVIDIA A800 GPUs. Each experiment was conducted twice, averaging the results reported here.

Experimental results. We next report our findings.

Exp-1: Entity matching. We evaluate the effectiveness of PUER in comparison to other baselines with aspect to entity matching. Table 1 shows the performance of all baselines. PUER consistently outperforms all other baselines across all 11 datasets in terms of precision, recall and F1-score, achieving average improvements of 26.31%, 36.87% and 40.19%, respectively, and up to 63.21%, 50.09% and 53.42%. This verifies that the co-training strategy between \mathcal{F}_{EM}^S and \mathcal{F}_{EM}^M and interaction between \mathcal{F}_{EM} and \mathcal{F}_{RAG} are effective, where the Selector and the RAG blocker enhance the performance of the entity matching component. Furthermore, PUER exhibits greater robustness compared to other baselines and is less not sensitive with data distribution. Specifically, PUER shows superior performance in 10 out of 11 datasets in terms of F1-score, e.g., at least 23.01%, 14.82% and 47.74% improvement in WS, DS and Company dataset across different domains. This highlights the stability of PUER and its effectiveness independent of specific data distribution.

Compared with pre-trained baselines, PUER also outperforms JellyFish, an LLM-based entity matching model using handcraft prompts and LoRA tuning, e.g., 16.63% F1-score improvement on average. This underscores the effectiveness of

Datasets	PUER (Ours)			Ditto			Rotom			Unicorn			PromptEM			JellyFish			Sudowoodo		
	P	R	F	P	R	F	P	R	F	P	R	F	P	R	F	P	R	F	P	R	F
AG	84.83	76.49	80.45	39.28	4.70	8.39	19.50	59.01	29.30	90.66	11.53	20.14	64.62	17.95	28.09	92.03	44.44	59.94	55.85	52.99	54.88
WA	93.95	88.60	91.20	78.43	20.72	32.78	11.80	73.10	20.30	89.99	60.62	72.44	93.55	30.05	45.49	80.09	93.78	86.39	46.41	50.25	48.25
AB	90.04	87.86	88.94	97.24	51.45	67.3	14.60	42.70	21.70	97.11	49.02	65.16	98.04	48.54	64.94	99.38	78.15	87.5	42.30	32.03	36.46
DA	95.47	99.77	97.57	99.5	89.63	94.31	80.9	97.1	88.2	99.29	95.27	97.24	100	86.04	92.49	99.76	97.07	98.40	78.55	98.19	87.28
DS	99.31	94.85	97.03	98.2	30.65	46.72	65.6	94.7	77.5	98.81	70.18	82.07	98.73	58.04	73.1	99.7	64.01	77.97	73.33	93.55	82.21
CO	98.56	79.18	87.82	25.06	100	40.08	n/a	n/a	n/a	87.49	3.84	7.37	n/a	n/a	n/a	96.43	22.07	35.92	30.62	25.28	27.69
CA	93.20	100.00	96.48	70.53	27.43	39.5	40.1	35.8	37.8	96.29	36.11	52.52	89.02	25.35	39.46	96.07	68.05	79.67	57.58	44.79	50.39
COM	98.67	99.33	99.00	65.64	28.76	39.99	29.3	50.2	37	92.82	69.23	79.31	80.57	74.58	77.57	97.5	78.26	86.82	37.08	48.49	42.02
SH	98.48	88.13	93.02	74.7	43.05	54.62	26.7	55.9	36.1	80.47	58.64	67.84	50.00	1.00	1.97	94.44	51.86	66.95	34.38	36.94	35.62
WAT	97.89	85.03	91.01	27.36	27.09	27.22	26.9	65.9	38.2	93.71	49.83	65.06	42.86	1.00	1.96	89.45	77.37	82.97	31.49	51.50	39.08
WS	87.67	95.57	91.45	71.89	20.28	31.64	27.4	75.00	40.2	95.97	43.73	60.09	91.55	28.05	42.94	96.64	52.92	68.39	40.99	63.10	49.70
Average	94.37	90.43	92.17	67.98	40.34	43.87	31.16	59.04	38.75	92.96	49.82	60.84	73.54	51.69	42.55	94.68	66.18	75.54	48.05	54.28	50.33

Table 1: Entity matching performance in comparison to baselines, n/a means the method cannot be terminated within 10 hours

Methods/Model	AB	AG	DA	DS	WA
PUER (Ours)	88.94	80.45	97.57	97.03	91.20
CLER	75.86	47.56	80.04	55.96	70.02
BatchER (GPT-4)	85.22	64.06	96.04	89.48	81.22
ComEM (GPT-3.5-turbo)	87.62	69.63	90.85	84.68	86.37
ComEM (GPT-4o-mini)	88.24	71.47	90.58	87.84	88.56

Table 2: Comparison with Online Model (F1-score)

Datasets	AG		AB		WA	
	Train	Predict	Train	Predict	Train	Predict
Ditto	235	23	208	20	215	22
Rotom	522	23	435	20	487	23
Unicorn	725	17	602	13	654	14
PromptEM	1420	65	1533	42	1365	55
JellyFish	1243	45	1010	30	1190	42
PUER (Ours)	3561	208	4323	255	4555	339

Table 3: The Efficiency of Entity Matching (in seconds)

our end-to-end iterative framework. While Unicorn achieves relatively good performance among the baselines due to its mixture-of-expert architecture, it struggles to attain high recall because of the small set of positive instances.

To evaluate the ER framework that integrates entity blocking and matching, we compare with CLER and Sudowoodo in Table 1 and 2. PUER is 10.85%, 29.73% and 23.63% more accuracy in aspects of precision, recall and F1-score than CLER on average (Table 2), which indicates that the proposed workflow that interacts \mathcal{F}_{RAG} and \mathcal{F}_{EM} could be beneficial to both of them. Due to the few-shot PU setting, Sudowoodo struggles to generate high-quality pseudo-labels relying solely on similarity threshold and positive ratio, and the number of high-quality pseudo-labels it generates is insufficient (Table 1). In contrast, PUER employs the GRPO algorithm with a carefully designed reward function, which enables it to tolerate noisy data with better generalization ability.

In Table 2, we compared PUER with the online BatchER (Fan et al., 2024b) and ComEM (Wang et al., 2025) using GPT-3.5, GPT-4 and GPT-4o-mini as the backbones. The result shows that PUER achieves up to 15.5% higher F1-score, indicating the effectiveness of the fine-tuning workflow.

Methods	AG			WA			AB		
	P	R	F	P	R	F	P	R	F
PUER	84.83	76.49	80.45	93.95	88.60	91.20	90.04	87.86	88.94
w.o. Selector	33.43	97.86	49.83	11.42	100.00	20.51	42.20	94.66	58.38
w.o. enrich	65.92	76.06	70.63	81.42	88.60	84.86	88.62	90.77	89.68
w.o. co-train	63.66	84.61	72.66	75.73	93.78	83.79	85.30	87.37	86.33

Table 4: Ablation Study for Entity Matching

Exp-2: Efficiency of entity matching. We present the fine-tuning time (Train) and inference time (Predict) of \mathcal{F}_{EM} in Table 3. Since \mathcal{F}_{EM} employs co-training of the Selector and Matcher, as well as an iterative workflow to gradually generate more training data with pseudo-labels, PUER requires significantly more time to fine-tune its entity matching models and perform inferences compared to other methods. Although PUER is notably slower, its high accuracy in entity matching, as shown in Table 1, justifies the use of RL and co-training strategy despite the increased computational cost.

The training cost of PUER is bounded by the number of training data. Despite its higher training cost, it does not result in prohibitively high expenses in few-shot setting. Considering the significant improvement in the F1-score of PUER (e.g., at least 16.63% F1 improvement over baselines) and the unique background of the positive-unlabeled setting and ER task, the higher cost can be tolerated and is relatively negligible.

Exp-3: Quality of pseudo-labels. To further demonstrate the effectiveness of our pseudo-labeling mechanism in PUER, we conducted experiments to measure the precision and recall of pseudo-positive and pseudo-negative examples under varying values of K in Table 5 following (Wang et al., 2023). Notice that we discard tuple pairs of pseudo-labels that are not included in the ground truths provided by benchmarks so that the evaluations are accurate and totally based on benchmarks.

By varying K from 5 to 9, the average precision and recall are 88.28% and 81.53%, respectively. These results indicate the high quality of

Dataset	K=5		K=7		K=9	
	P	R	P	R	P	R
AG	78.33	67.52	77.24	67.52	76.17	67.09
WA	95.04	87.63	95.23	87.21	95.05	87.83
AB	92.96	89.88	93.01	89.30	91.47	89.78

Table 5: Quality of pseudo-labels under different K

our generated pseudo-labels. Furthermore, as K increases, the quality of the pseudo-labels only decreases slightly. Although it may involve more tuples for checking, our PUER is robust enough to generate correct pseudo-labels, which is critical for the further fine-tuning of the selector and matcher.

Exp-4: Flexibility of PUER. Given that our ER framework does not rely on the matcher, we adopted the EM model in Unicorn (Fan et al., 2024a), a smaller model as the matcher of PUER. Table 6 and 7 present the effectiveness and efficiency of PUER_{unicorn}, which utilizes the EM model of Unicorn as the matcher, respectively.

Dataset	PUER			PUER _{unicorn}		
	P	R	F	P	R	F
AG	84.83	76.49	80.45	80.18	76.07	78.07
WA	93.95	88.60	91.20	76.19	82.90	79.40
AB	90.04	87.86	88.94	90.28	76.69	82.93

Table 6: Effectiveness of PUER v.s. PUER_{unicorn}.

Dataset	Train		Predict	
	PUER	PUER _{unicorn}	PUER	PUER _{unicorn}
AG	3561	3027 (-534)	208	160 (-48)
WA	4555	3761 (-794)	339	205 (-134)
AB	4323	3618 (-705)	255	202 (-53)

Table 7: Efficiency of PUER vs. PUER_{unicorn} (in seconds).

As evidenced by Table 6 and 7, the training time is reduced to approx. 80% of the original time, but the accuracy remains largely unaffected, which indicates the effectiveness of the PUER framework, regardless of whether the Matcher is a powerful LLM or a simpler model. Under the few-shot PU setting, considering the superior accuracy achieved, the training time is deemed acceptable.

Exp-5: Ablation study. In Table 4, we present the ablation study of PUER and its variants, namely PUER without the Selector, without enrichment and PUER without co-training that the selector and matcher are trained independently. The results demonstrate that each variant achieves lower accuracy compared to the complete PUER model.

Specifically, PUER without enrichment and without the Matcher shows only a relatively small

performance drop. This suggests that LLMs already possess substantial prior knowledge, and the Matcher does not contribute much in the few-shot data setting. In contrast, PUER without the Selector experiences a significant drop in accuracy, such as a 31% decrease on the AG dataset. This finding indicates that the RL component in the Selector plays a crucial role in generalization, particularly in few-shot scenarios. It enables the model to collect sufficient data from the environment to achieve robust performance.

To determine whether the pre-training process of the LLM backbone includes the benchmark datasets, we conducted an experiment of the few-shot prompting using Qwen-2.5-7B-instruct.

	CO	WA	DS	AB	AG
P	99.35	74.26	92.02	77.63	46.74
R	10.92	52.33	62.52	60.67	82.90
F	19.68 (-68.14)	61.40 (-29.80)	74.45 (-22.58)	68.11 (-20.83)	59.78 (-20.67)

Table 8: The effectiveness of Qwen-2.5-7B-instruct (P/R/F)

As illustrated in Table 8, the F1-score of PUER is at least 20.67 points higher than that of few-shot prompting and 68.14 points at most, demonstrating that only LLM backbone is not sufficient and PUER is necessary to make more accurate predictions.

Dataset	Methods	$ \mathcal{P} =10$	20	30	40	50	100
AB	PUER	55.15	74.01	73.11	79.49	88.92	89.26
	Unicorn	48.17	49.08	59.66	71.42	65.16	84.57
	DITTO	17.39	21.58	34.04	45.66	67.3	82.35
WS	PUER	65.82	68.96	65.01	82.29	91.45	92.12
	Unicorn	23.12	35.02	65.22	54.24	40.20	70.19
	DITTO	13.25	20.57	27.45	35.69	31.64	66.10

Table 9: Performance Vary Labelling Budget $|\mathcal{P}|$ (F1)

Exp-6: Hyper-parameter study. We vary the labeling budget $|\mathcal{P}|$ from 10 to 50 in Table 9. PUER demonstrates robustness with respect to the number of positive tuples, *e.g.*, only 6.9% drop when $|\mathcal{P}|$ decreases from 50 to 10 in WS.

7 Conclusion

In this paper, we propose, PUER, an end-to-end ER solution for few-shot PU learning. We adopt the reinforcement learning method to solve the entity matching task, and design a self-adaptive reward function. Furthermore, we introduce an iterative training workflow that fully utilizes the entity blocking model to assist the entity matching via a co-training mechanism. Finally comprehensive experiments across 11 benchmarks demonstrate the superior performance of PUER.

Limitations

Our work, PUER, introduces an end-to-end entity resolution solution tailored for few-shot positive-unlabeled (PU) learning scenarios by leveraging Large Language Models (LLMs) and reinforcement learning. We propose an iterative co-training mechanism that integrates entity blocking and entity matching, including a novel self-adaptive reward function for the reinforcement learning component, to enhance performance with minimal labeled positive data.

Despite the promising results, our approach has several limitations. Firstly, the reinforcement learning (RL) component, particularly the Selector fine-tuned with Group Relative Policy Optimization (GRPO), inherently introduces a higher level of complexity in terms of training and hyperparameter tuning compared to simpler supervised methods. Secondly, while LLMs offer powerful generative capabilities, the quality and consistency of the generated outputs (e.g., enriched attributes or pseudo-labels) can be uncertain and may occasionally require careful validation. Lastly, as indicated by our efficiency experiments, the proposed PUER framework, with its iterative workflow and co-training of multiple components including RL-based Selector and Matcher, exhibits a higher computational complexity during both training and inference compared to some traditional entity resolution methods. This increased cost is a trade-off for the achieved accuracy in low-resource settings.

Ethics Statement

The experiments were conducted on publicly available benchmark datasets and models, eliminating any data privacy concerns. To the best of our knowledge, there is no negative societal impact in this research.

Acknowledgments

This work was supported by Guangdong Provincial Key Lab of Integrated Communication, Sensing and Computation for Ubiquitous Internet of Things (No.2023B1212010007, SL2023A03J00934), Guangzhou Municipal Science and Technology Project (No. 2023A03J0003, 2023A03J0013 and 2024A03J0621).

References

- Naser Ahmadi, Hansjörg Sand, and Paolo Papotti. 2022. Unsupervised matching of data and text. In *2022 IEEE 38th International Conference on Data Engineering (ICDE)*, pages 1058–1070. IEEE.
- Arvind Arasu, Michaela Götz, and Raghav Kaushik. 2010. On active learning of record matching packages. In *SIGMOD*, pages 783–794.
- Jessa Bekker and Jesse Davis. 2020. Learning from positive and unlabeled data: a survey. *Mach. Learn.*
- Mikhail Bilenko and Raymond J. Mooney. 2003. Adaptive duplicate detection using learnable string similarity measures. In *Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Washington, DC, USA, August 24 - 27, 2003*.
- Alexander Brinkmann, Roei Shraga, and Christina Bizer. 2024. Sc-block: Supervised contrastive blocking within entity resolution pipelines. In *ESWC*.
- Paul Suganthan G. C., Adel Ardalan, AnHai Doan, and Aditya Akella. 2018. Smurf: Self-service string matching using random forests. *Proc. VLDB Endow.*
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jia Shi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojuan Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, and S. S. Li. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *CoRR*, abs/2501.12948.
- Muhammad Ebraheem, Saravanan Thirumuruganathan, Shafiq R. Joty, Mourad Ouzzani, and Nan Tang. 2018. Distributed representations of tuples for entity resolution. *PVLDB*, 16(8):1944–1957.
- Vasilis Efthymiou, George Papadakis, George Papastefanatos, Kostas Stefanidis, and Themis Palpanas.

2015. Parallel meta-blocking: Realizing scalable entity resolution over large, heterogeneous data. In *IEEE BigData*.
- Ju Fan, Jianhong Tu, Guoliang Li, Peng Wang, Xiaoyong Du, Xiaofeng Jia, Song Gao, and Nan Tang. 2024a. Unicorn: A unified multi-tasking matching model. *SIGMOD Rec.*
- Meihao Fan, Xiaoyue Han, Ju Fan, Chengliang Chai, Nan Tang, Guoliang Li, and Xiaoyong Du. 2024b. Cost-effective in-context learning for entity resolution: A design space exploration. In *40th IEEE International Conference on Data Engineering, ICDE 2024, Utrecht, The Netherlands, May 13-16, 2024*.
- Wenfei Fan, Hong Gao, Xibei Jia, Jianzhong Li, and Shuai Ma. 2011. Dynamic constraints for record matching. *VLDB J.*, 20(4):495–520.
- Wenfei Fan, Xibei Jia, Jianzhong Li, and Shuai Ma. 2009. Reasoning about record matching rules. *PVLDB*, 2(1):407–418.
- Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and Yarin Gal. 2024. Detecting hallucinations in large language models using semantic entropy. *Nature*, 630(8017):625–630.
- Cheng Fu, Xianpei Han, Le Sun, Bo Chen, Wei Zhang, Suhui Wu, and Hao Kong. 2019. End-to-end multi-perspective matching for entity resolution. In *IJCAI*, pages 4961–4967.
- Songtao Guo, Xin Luna Dong, Divesh Srivastava, and Remi Zajac. 2010. Record linkage with uniqueness constraints and erroneous values. *PVLDB*, 3(1):417–428.
- Jungo Kasai, Kun Qian, Sairam Gurajada, Yunyao Li, and Lucian Popa. 2019. Low-resource deep entity resolution with transfer and active learning. In *ACL*, pages 5851–5861.
- Mayank Kejriwal and Daniel P. Miranker. 2015. A DNF blocking scheme learner for heterogeneous datasets. *CoRR*, abs/1501.01694.
- Nishadi Kirielle, Peter Christen, and Thilina Ranbaduge. 2022. Transer: Homogeneous transfer learning for entity resolution. In *EDBT*.
- Pradap Konda, Sanjib Das, Paul Suganthan G. C., An-Hai Doan, Adel Ardalani, Jeffrey R. Ballard, Han Li, Fatemah Panahi, Haojun Zhang, Jeffrey F. Naughton, Shishir Prasad, Ganesh Krishnan, Rohit Deep, and Vijay Raghavendra. 2016. Magellan: Toward building entity matching management systems. *PVLDB*, 9(12):1197–1208.
- Hanna Köpcke, Andreas Thor, and Erhard Rahm. 2010. Evaluation of entity resolution approaches on real-world match problems. *Proc. VLDB Endow.*, 3(1):484–493.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.
- Bing Li, Yukai Miao, Yaoshu Wang, Yifang Sun, and Wei Wang. 2021. Improving the efficiency and effectiveness for bert-based entity resolution. In *AAAI*.
- Bing Li, Wei Wang, Yifang Sun, Linhan Zhang, Muhammad Asif Ali, and Yi Wang. 2020a. Grapher: Token-centric entity resolution with graph convolutional neural networks. In *AAAI*, pages 8172–8179.
- Huahan Li, Shuangyin Li, Fei Hao, Chen Jason Zhang, Yuanfeng Song, and Lei Chen. 2024a. Booster: leveraging large language models for enhancing entity resolution. In *WWW*, pages 1043–1046.
- Peng Li, Yeye He, Dror Yashar, Weiwei Cui, Song Ge, Haidong Zhang, Danielle Rifinski Fainman, Dongmei Zhang, and Surajit Chaudhuri. 2024b. Table-gpt: Table fine-tuned GPT for diverse table tasks. *Proc. ACM Manag. Data*.
- Yuliang Li, Jinfeng Li, Yoshihiko Suhara, AnHai Doan, and Wang-Chiew Tan. 2020b. Deep entity matching with pre-trained language models. *PVLDB*, 14(1):50–60.
- Michael Loster, Ioannis K. Koumarelas, and Felix Naumann. 2021. Knowledge transfer for entity resolution with siamese neural networks. *ACM J. Data Inf. Qual.*
- Venkata Vamsikrishna Meduri, Lucian Popa, Prithviraj Sen, and Mohamed Sarwat. 2020. A comprehensive benchmark framework for active learning methods in entity matching. In *SIGMOD*.
- Zhengjie Miao, Yuliang Li, and Xiaolan Wang. 2021. Rotom: A meta-learned data augmentation framework for entity matching, data cleaning, text classification, and beyond. In *SIGMOD*, pages 1303–1316. ACM.
- Matthew Michelson and Craig A. Knoblock. 2006. Learning blocking schemes for record linkage. In *AAAI*.
- Sidharth Mudgal, Han Li, Theodoros Rekatsinas, An-Hai Doan, Youngchoon Park, Ganesh Krishnan, Rohit Deep, Esteban Arcaute, and Vijay Raghavendra. 2018. Deep learning for entity matching: A design space exploration. In *SIGMOD*, pages 19–34.
- Youcef Nafa, Qun Chen, Zhaoqiang Chen, Xingyu Lu, Haiyang He, Tianyi Duan, and Zhanhuai Li. 2022. Active deep learning on entity resolution by risk sampling. *Knowl. Based Syst.*

- Gang Niu, Marthinus Christoffel du Plessis, Tomoya Sakai, Yao Ma, and Masashi Sugiyama. 2016. Theoretical comparisons of positive-unlabeled learning against positive-negative learning. In *NeurIPS*.
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*.
- George Papadakis, Georgia Koutrika, Themis Palpanas, and Wolfgang Nejdl. 2014. Meta-blocking: Taking entity resolution to the next level. *IEEE Trans. Knowl. Data Eng.*
- George Papadakis, Dimitrios Skoutas, Emmanouil Thanos, and Themis Palpanas. 2020. Blocking and filtering techniques for entity resolution: A survey. *ACM Computing Surveys (CSUR)*, 53(2):1–42.
- Derek Paulsen, Yash Govind, and AnHai Doan. 2023. Sparkly: A simple yet surprisingly strong TF/IDF blocker for entity matching. *Proc. VLDB Endow.*
- Anna Primpeli, Ralph Peeters, and Christian Bizer. 2019. The WDC training dataset and gold standard for large-scale product matching. In *WWW*.
- Kun Qian, Lucian Popa, and Prithviraj Sen. 2017. Active learning for large-scale entity resolution. In *CIKM*, pages 1379–1388.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In *EMNLP*.
- Guangming Sheng, Chi Zhang, Zilinfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. 2025. Hybridflow: A flexible and efficient RLHF framework. In *Proceedings of the Twentieth European Conference on Computer Systems, EuroSys 2025, Rotterdam, The Netherlands, 30 March 2025 - 3 April 2025*.
- Rohit Singh, Venkata Vamsikrishna Meduri, Ahmed K. Elmagarmid, Samuel Madden, Paolo Papotti, Jorge-Arnulfo Quiané-Ruiz, Armando Solar-Lezama, and Nan Tang. 2017a. Synthesizing entity matching rules by examples. *Proc. VLDB Endow.*
- Rohit Singh, Venkata Vamsikrishna Meduri, Ahmed K. Elmagarmid, Samuel Madden, Paolo Papotti, Jorge-Arnulfo Quiané-Ruiz, Armando Solar-Lezama, and Nan Tang. 2017b. Synthesizing entity matching rules by examples. *PVLDB*.
- Chenchen Sun, Yang Xu, Derong Shen, and Tiezheng Nie. 2024. Matching feature separation network for domain adaptation in entity matching. In *WWW*, pages 1975–1985. ACM.
- The Magellan Data Repository. Sanjib das, anhai doan, paul suganthan g. c., chaitanya gokhale, pradap konda, yash govind, and derek paulsen. <https://sites.google.com/site/anhaidgroup/projects/data>.
- Saravanan Thirumuruganathan, Han Li, Nan Tang, Mourad Ouzzani, Yash Govind, Derek Paulsen, Glenn Fung, and AnHai Doan. 2021. Deep learning for blocking in entity matching: A design space exploration. *Proc. VLDB Endow.*, 14(11):2459–2472.
- Jianhong Tu, Xiaoyue Han, Ju Fan, Nan Tang, Chengliang Chai, Guoliang Li, and Xiaoyong Du. 2022. DADER: hands-off entity resolution with domain adaptation. *PVLDB*, 15(12):3666–3669.
- Somin Wadhwa, Adit Krishnan, Runhui Wang, Byron C. Wallace, and Luyang Kong. 2024. Learning from natural language explanations for generalizable entity matching. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*. Association for Computational Linguistics.
- Pengfei Wang, Xiaocan Zeng, Lu Chen, Fan Ye, Yuren Mao, Junhao Zhu, and Yunjun Gao. 2022. Promptem: Prompt-tuning for low-resource generalized entity matching. *PVLDB*.
- Runhui Wang, Yuliang Li, and Jin Wang. 2023. Sudooodo: Contrastive self-supervised learning for multi-purpose data integration and preparation. In *ICDE*.
- Tianshu Wang, Xiaoyang Chen, Hongyu Lin, Xuanang Chen, Xianpei Han, Le Sun, Hao Wang, and Zhenyu Zeng. 2025. Match, compare, or select? an investigation of large language models for entity matching. In *ACL*, pages 96–109.
- Tianshu Wang, Hongyu Lin, Xianpei Han, Xiaoyang Chen, Boxi Cao, and Le Sun. 2024a. Towards universal dense blocking for entity resolution. *CoRR*, abs/2404.14831.
- Ye Wang, Huazheng Pan, Tao Zhang, Wen Wu, and Wenxin Hu. 2024b. A positive-unlabeled metric learning framework for document-level relation extraction with incomplete labeling. In *AAAI*.
- Steven Euijong Whang and Hector Garcia-Molina. 2013. Joint entity resolution on multiple datasets. *VLDB J.*, 22(6):773–795.
- Brandon T Willard and Rémi Louf. 2023. Efficient guided generation for llms. *arXiv preprint arXiv:2307.09702*.
- Renzhi Wu, Sanya Chaba, Saurabh Sawlani, Xu Chu, and Saravanan Thirumuruganathan. 2020. ZeroER: Entity resolution using zero labeled examples. In *SIGMOD*, pages 1149–1164.
- Shiwen Wu, Qiyu Wu, Honghua Dong, Wen Hua, and Xiaofang Zhou. 2023. Blocker and matcher can mutually benefit: A co-learning framework for low-resource entity resolution. *Proc. VLDB Endow.*
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian

Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. 2024. Qwen2.5 technical report. *CoRR*, abs/2412.15115.

Zijun Yao, Chengjiang Li, Tiansi Dong, Xin Lv, Jifan Yu, Lei Hou, Juanzi Li, Yichi Zhang, and Zelin Dai. 2021. Interpretable and low-resource entity matching via decoupling feature learning from decision making. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*.

Xiaocan Zeng, Pengfei Wang, Yuren Mao, Lu Chen, Xiaozhe Liu, and Yunjun Gao. 2024. Multiem: Efficient and effective unsupervised multi-table entity matching. In *40th IEEE International Conference on Data Engineering, ICDE 2024, Utrecht, The Netherlands, May 13-16, 2024*.

Haochen Zhang, Yuyang Dong, Chuan Xiao, and Masafumi Oyamada. 2024. [Jellyfish: Instruction-tuning local large language models for data preprocessing](#). In *EMNLP*, pages 8754–8782, Miami, Florida, USA. Association for Computational Linguistics.

Peitian Zhang, Shitao Xiao, Zheng Liu, Zhicheng Dou, and Jian-Yun Nie. 2023. Retrieve anything to augment large language models. *CoRR*.

Chen Zhao and Yeye He. 2019. Auto-EM: End-to-end fuzzy entity-matching using pre-trained deep models and transfer learning. In *WWW*, pages 2413–2424.

Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma. 2024. [Llamafactory: Unified efficient fine-tuning of 100+ language models](#). In *ACL*, Bangkok, Thailand. Association for Computational Linguistics.

A Overview

In the supplementary material, we mainly provide (1) a running example of entity resolution; (2) an example of data enrichment using LLMs along with our observation; (3) detailed hyper-parameter configurations used in our experiment; (4) the notation table from our full paper; (5) comprehensive descriptions of the datasets used in our experiments; (6) experimental results for entity blocking across all datasets, comparing with other baselines in various settings; (7) a detailed comparison with on-line model, *e.g.*, ChatGPT, under different parameter size for LLMs, and (8) detailed information of prompts used in PUER.

B Examples of Entity Resolution

We present an example of entity resolution in Table 11 (left table, denoted as \mathcal{R}_l) and Table 12 (right table, denoted as \mathcal{R}_r). Both tables contain multiple tuples with three fundamental attributes, specifically $\bar{A} = \{\text{Manufacturer, price, title}\}$. Among these, the pairs l_{1162} and r_{2109} , as well as l_{587} and r_{2816} , represent the same real-world entities. The objective of the entity resolution task is to efficiently and effectively identify all matching tuple pairs between \mathcal{R}_l and \mathcal{R}_r .

C Examples of Data Enrichment

C.1 Schema Enrichment

In Table 13 and Table 14, we first use LLMs to enrich the data with 6 additional attributes, denoted as \bar{B} , where $\bar{B} = \{\text{category, subcategory, platform, edition, type, modelno}\}$. We then query LLMs for each tuple using the attributes in \bar{B} . The values of the columns highlighted in blue are imputed by LLMs. These imputed values provide PUER more valuable information to identify matched tuple pairs.

C.2 Our Observation

In Table 13, we demonstrate that l_{1162} is imputed with different values for \bar{B} , specifically l_{1162}^1 and l_{1162}^3 , depending on whether it is paired with r_{2109} or r_{2816} (in Table 14). Based on these observations, we can enhance the data quality of training data in our Positive-Unlabeled (PU) setting. Such **pair-wise enrich** problem can activate the ability of LLM generating different information from various perspective and context.

D Hyper-parameter Configuration

Please check Table 19 for detailed hyper-parameter configuration and corresponding explanations. Our code is available at <https://anonymous.4open.science/r/PUER-CB71>.

We apply LLaMA-Factory (Zheng et al., 2024) for training, and apply vLLM (Kwon et al., 2023) for efficient inference. We use verl (Sheng et al., 2025) for training RL-based Selector. To stabilize the output result for Matcher and Selector, during inference with vLLM, for querying LLM, we set the temperature to 0, and top-p to 1 for deterministic output.

We also incorporate outlines (Willard and Louf, 2023) to fix the output of LLMs to JSON format.

E Notation Table

In Table 15, we provide notation table and their corresponding descriptions.

F The dataset descriptions

In Table 16, we provide descriptions of all benchmark datasets used in this paper. Following our PU learning setting, for each dataset we only use 50 random sampled positive samples for PUER.

- The column *# Dataset* lists all datasets used in this paper along with their abbreviation. For the WDC dataset (*w.r.t.* WS, COM, CA, SH, WAT), we sampled 50 positive tuple pairs within the small size (1/20 of all pairs, following (Mudgal et al., 2018; Li et al., 2020b)) of each dataset.
- The column *# All* provides the total number of labeled examples for each dataset, and the column *Match* specifies the number of matched examples for each dataset, including the train/valid/test splits.
- The column *# of Original Attr* shows the number of attributes in each original dataset, and the column *# of Enriched Attr* displays the number of attributes for enrichment by PUER, *i.e.*, $|\bar{A}| + |\bar{B}|$. A detailed example for dataset is provided in Figure 13 and 14.
- The column *# of $|\mathcal{R}_l|, |\mathcal{R}_r|$* indicates the sizes of left and right tables for each dataset, respectively.
- The column *Proportion of PU* represents the ratio of PU positive sample (*e.g.*, 50) to all labeled training samples in benchmark datasets, while

Methods/Model	AB	AG	DA	DS	WA	Avg
PUER (Qwen-2.5-7B) w.o. SFT	39.41	66.81	89.25	77.04	49.62	64.43
PUER (Qwen-2.5-0.5B)	12.54	17.40	99.32	96.32	14.80	48.08
PUER (Qwen-2.5-1.5B)	73.03	81.00	96.08	97.72	50.19	79.60
PUER (Qwen-2.5-3B)	82.72	79.30	95.95	97.70	87.07	88.55
PUER (Qwen-2.5-7B)	88.94	80.45	97.57	97.03	91.20	91.04
ComEM (Mistral-7B)	40.70	37.77	24.68	28.89	55.96	37.60
ComEM (Qwen-2.5-7B)	72.39	61.03	81.49	76.57	72.96	72.89
ComEM (LLAMA3-8B)	74.37	49.50	78.91	68.79	42.33	62.78
ComEM (Mixtral-8×7B)	77.67	34.76	67.20	60.09	50.57	58.06
BatchER (GPT-4)	85.22	64.06	96.04	89.48	81.22	83.20
ComEM (GPT-3.5-turbo)	87.62	69.63	90.85	84.68	86.37	83.83
ComEM (GPT-4o-mini)	88.24	71.47	90.58	87.84	88.56	85.34

Table 10: Comparison with Online Model (F1 Score). For our method PUER, we fix the RL-based Selector model as Qwen-2.5-7B, and only change Matcher with different backbone model.

id	title	Manufacturer	price
l_{1162}	motu digital performer 5 digital audio software competitive upgrade (mac only)	motu	395.0
l_{587}	microsoft word 2007 version upgrade	microsoft	109.95

Table 11: Examples of Amazon dataset with basic attributes \bar{A}

the column *Proportion of Positive Samples* represents the ratio of PU positive sample(*e.g.*, 50) to all labeled positive training samples.

G Full Version of Blocking Result

Table 17 shows full version of blocking result in all datasets. PUER shows the superior performance, *i.e.*, highest values of PC and PQ and the smallest value of K in most cases.

H Detailed Comparison with Online Model

Table 10 shows the comparison of PUER and other offline and online models without SFT (Supervised Fine-Tuning). The results demonstrate the effectiveness of co-training of Matcher and Selector subtasks in our proposed method PUER. All performances are evaluated under the same In-Context Learning settings, using an equal number of positive and negative samples for demonstration.

In Table 10, the upper section includes our methods PUER and PUER without SFT, while the lower section follows the setting from (Wang et al., 2025).

We also compare the matching result of PUER with backbone LLMs under different parameter size for Matcher, *e.g.*, 0.5B, 1.5B, 3B and 7B in Table 10. We can observe that larger backbone LLMs

are not consistently stronger in entity matching performance. However, the 7B model achieve the highest average performance among different scenarios.

I PC/PQ Curve for Blocking Experiment

Figure 7 shows the performance of our RAG blocker in PUER, in terms of PC (Top- K Recall) for different values of K (candidate set size). A curve approaching the **upper left corner** of the figure indicates better performance. The results show that the RAG blocker of PUER is highly effective, capable of retrieving the smallest number of candidate tuples while achieving the highest recall.

J CSSR Curve for Blocking Experiment varying K

Figure 8 provide the Blocker performance(in PC, *w.r.t.* Top- K Recall) under different K , following the setting of DeepBlocker (Thirumuruganathan et al., 2021). The curve approaching the **lower right corner** of the figure indicates better performance. PUER also performs the best among all baseliens in most datasets.

id	title	Manufacturer	price
r_{2816}	microsoft word 2007 upgrade (pc)	null	109.95
r_{2109}	motu digital performer dp5 software music production software	null	319.95

Table 12: Examples of Google dataset with basic attributes \bar{A}

id	title	Ma.	price	category	sub-category	platform	edition	type	modelno
l_{1162}^1	motu digital performer 5 digital audio software competitive upgrade (mac only)	motu	395.0	Audio Production	DAWs	Mac	Competitive Upgrade	Software	DP5
l_{1162}^2	motu digital performer 5 digital audio software competitive upgrade (mac only)	motu	395.0	Audio & Music Software	Audio Editing & Production	Mac	Standard	Software	5
l_{1162}^3	motu digital performer 5 digital audio software competitive upgrade (mac only)	motu	395.0	Audio Editing Software	DAW (Digital Audio Workstation)	Mac	Upgrade	Software	5
l_{587}^1	microsoft word 2007 version upgrade	microsoft	109.95	Productivity Software	Office Suites	Windows	Standard	Upgrade	2007
l_{587}^2	microsoft word 2007 version upgrade	microsoft	109.95	Productivity Software	Word Processing	Windows	home	Upgrade	2007
l_{587}^3	microsoft word 2007 version upgrade	microsoft	109.95	software	office	Windows	ultimate	Upgrade	2007

Table 13: Examples for Amazon dataset (left table for Amazon-Google dataset). Grey columns are original attributes(*w.r.t.* \bar{A}), and blue columns are enriched attributes(*w.r.t.* \bar{B}). For each entity (*e.g.*, l_{1162} , l_{587}), we report three different enrichment outputs, to demonstrate the uncertainty of our proposed data enrichment methods. Ma. is short for attribute Manufacturer.

K Example of the Schema Enrichment Prompt pt_{SE}

In Prompt Template 3, **Entity 1** l_{1162} is from Table 13(left table *Amazon*), and **Entity 2** r_{2109} is from Table 14 (right table *Google*).

For each dataset, we query LLM using the same prompt pt_{SE} with varying different **Entity 1** and **Entity 2** multiple times. We then apply majority voting to the different generated attributes to determine \bar{B} .

L Example of the Data Enrichment Prompt pt_{enr}

Prompt Template 5 provides an example of pt_{enr} using the Amazon-Google dataset. The enriched attribute set \bar{B} is obtained from the previous step using pt_{SE} .

pt_{enr} is queried with different entity pairs, *e.g.*, (l_{1162}, l_{587}) , (l_{1162}, r_{2109}) , (l_{587}, r_{2816}) , (l_{587}, r_{2109}) to generate different values of \bar{B} .

M Example of the Subtask Matcher Prompt pt_m

Prompt Template 4 provides an example for the Matcher subtask using the DBLP-Scholar (DS) dataset. **Paper 1** and **Paper 2** both contain enriched attributes that are extracted in the previous step using the prompt pt_{enr} .

N Example of the Subtask Selector Prompt pt_s

Prompt Template 6 provides an example for the Selector subtask using the Amazon-Google (AG) dataset. **Entity 1** and **Candidate** already contain enriched attributes extracted in the previous step using the enrichment prompt pt_{enr} . Additionally, **Candidate** entities are also retrieved and ranked using the preceding Blocker component, *i.e.*, \mathcal{F}_{RAG} .

O Training curve for Selector

We also list the curve for training GRPO-based Selector in Fig. 9.

id	title	Ma.	price	category	sub-category	platform	edition	type	modelno
r_{2816}^1	microsoft word 2007 upgrade (pc)	null	109.95	Productivity Software	Office Suites	Windows	Standard	Upgrade	2007
r_{2816}^2	microsoft word 2007 upgrade (pc)	null	109.95	Software	Office Suites	Windows	Upgrade	Desktop Software	2007
r_{2816}^3	microsoft word 2007 upgrade (pc)	null	109.95	Productivity Software	Word Processors	PC	Upgrade	Desktop Software	WORD2007UPG
r_{2109}^1	motu digital performer dp5 software music production software	null	319.95	Audio Production	DAWs			Software	DP5
r_{2109}^2	motu digital performer dp5 software music production software	null	319.95	Audio Production	DAWs	Mac	Pro	Software	DP5

Table 14: Examples of Google dataset (right table for Amazon-Google dataset). Grey columns are original attributes(w.r.t. \bar{A}), and blue columns are enriched attributes(w.r.t. \bar{B}). For each entity (e.g., r_{2816}, r_{2109}), we report three different enrichment outputs, to demonstrate the uncertainty of our proposed data enrichment methods. Ma. is short for attribute Manufacturer.

Symbol	Description
$t, \{A_1, \dots, A_m\}$	tuple t with multi-attributes $\{A_1, \dots, A_m\}$
$\mathcal{P}, \mathcal{P}_{\text{enr}}$	the labeled positive training dataset, and its enriched version
$\mathcal{P}_{\text{RAG}}, \mathcal{N}_{\text{RAG}}$	the set of potentially positive and negative tuple pairs by the RAG blocker
$\mathcal{R}_l, \mathcal{R}_r$	the left and right relational tables of multi-attribute tuples
\bar{B}, m	the set of enriched attributes, the number of enriched attributes
K	the top- K most similar tuples to retrieve by the blocker
$\text{NN}_K(t)$	the set of top- K most similar tuples with the tuple t
\mathcal{F}_{RAG}	the entity blocking model of PUER
\mathcal{F}_{EM}	the entity matching model of PUER
$\mathcal{F}_{\text{EM}}^M$	the Matcher subtask in \mathcal{F}_{EM}
$\mathcal{F}_{\text{EM}}^S$	the Selector subtask in \mathcal{F}_{EM}
$\mathcal{C}_s(t)$	the candidate list of the tuple t in $\mathcal{F}_{\text{EM}}^S$
$\mathcal{F}_{\text{label}}$	the labeler of the Selector
$\mathcal{D}_{\text{train}}$	the generated training data to fine-tune \mathcal{F}_{EM} , including labeled and pseudo-labeled training instances
pt_m, pt_s	the prompts of Matcher and Selector
pt_{enr}	the prompt of data enrichment by LLMs
pt_{SE}	the prompt of enriching more attributes by LLMs
λ	the warmup iteration
$\mathcal{M}_{\text{embed}}$	Embedding model for Blocker
S_t	pairwise enriched tuple set in right table \mathcal{R}_r for t
$p_m(s, t)$	query for LLM-based Matcher, to determine whether tuple pair s, t is match or mismatch

Table 15: General notations with corresponding descriptions.

Dataset	Domain	# All	# Match	# of Original Attr	# of Enriched Attr	# $ \mathcal{R}_l , \mathcal{R}_r $	Proportion of PU	Proportion of Positive Samples
Abt-Buy (AB)	Product	9,575	1,028	3	8	1081, 1092	0.87%	8.11%
Walmart-Amazon (WA)	Electronic	10,242	962	5	9	2554, 22074	0.81%	8.68%
Amazon-Google (AG)	Electronic	11,460	1,300	3	9	1363, 3226	0.72%	7.15%
DBLP-ACM (DA)	Citation	12,363	2,224	4	6	2616, 2294	0.67%	3.75%
DBLP-Scholar (DS)	Citation	28,707	5,347	4	6	2616, 64263	0.29%	1.56%
Company(CO)	Company	112,632	28,200	1	3	28200, 28200	0.07%	0.29%
WDC-All-Small(WS)	Product	13,436	3,516	1	6	7437, 8091	0.77%	2.69%
Computer(COM)	Electronic	3,865	1,005	1	7	2204, 2443	2.24%	8.98%
Camera(CA)	Product	2,858	752	1	7	1561, 1743	3.54%	13.62%
Shoes(SH)	Product	3,099	812	1	8	1600, 1767	3.10%	11.85%
Watch(WAT)	Product	3,181	831	1	9	1821, 1991	2.84%	10.98%

Table 16: Datasets used in our experiments, # means Number of, # Attr provide the original/enriched attribute number, Proportion of PU means the number of labeled samples divide **all train samples** in benchmark; Proportion of Positive Samples means the number of labeled positive samples in PU settings divide **all positive samples** in training samples.

	AG	AB	WA	DA	DS	WS	COM	CA	WAT	SH
DeepBlocker	85.69 / 3.67 / 20	75.19 / 3.57 / 20	90.12 / 3.39 / 10	97.21 / 82.49 / 1	90.14 / 18.43 / 10	55.00 / 0.72 / 20	61.59 / 0.80 / 20	60.37 / 0.79 / 20	28.03 / 0.30 / 20	25.86 / 0.27 / 20
Sudowoodo	90.06 / 9.80 / 8	90.37 / 28.73 / 3	90.54 / 8.53 / 4	98.92 / 84.92 / 1	90.24 / 12.30 / 15	53.04 / 0.92 / 20	68.55 / 1.12 / 20	61.96 / 1.01 / 20	26.83 / 0.35 / 20	26.23 / 0.34 / 20
STransformer	91.60 / 15.69 / 5	74.32 / 3.53 / 20	86.38 / 1.63 / 20	97.03 / 82.34 / 1	91.17 / 26.62 / 7	57.39 / 0.65 / 20	52.73 / 0.68 / 20	71.41 / 0.94 / 20	59.08 / 0.77 / 20	49.51 / 0.65 / 20
CLER	90.59 / 21.25 / 4	94.96 / 48.88 / 2	92.14 / 13.47 / 3	98.04 / 84.40 / 1	90.72 / 30.14 / 6	63.68 / 0.91 / 20	74.91 / 1.07 / 20	60.00 / 0.95 / 20	33.21 / 0.47 / 20	30.84 / 0.44 / 20
PUER	95.80 / 27.34 / 3	94.06 / 89.45 / 1	93.76 / 17.65 / 2	99.72 / 84.64 / 1	92.79 / 31.61 / 6	90.35 / 1.39 / 17	90.84 / 2.15 / 11	90.55 / 2.97 / 8	90.49 / 1.68 / 14	90.51 / 1.39 / 17

Table 17: Performance Evaluation. Following UniBlocker (Wang et al., 2024a), we report the first results (in order of PC/PQ/ K , also known as Top- K recall/precision) of baselines when their PC exceeds the threshold (90%). If both methods have larger PC than the threshold, we evaluate K , otherwise we evaluate their PC. If their K are the same, we evaluate their PC and PQ.

Dataset	Original Attribute \bar{A}	Enriched Attribute \bar{B}
Amazon-Google (AG)	title, manufacturer, price	category, subcategory, platform, edition, type, modelno
Abt-Buy (AB)	name, description, price	category, sku, brand, modelno, keyfeatures
Walmart-Amazon (WA)	title, category, brand, modelno, price	subcategory, key-features, sku, color
DBLP-ACM (DA)	title, authors, venue, year	keywords
DBLP-Scholar (DS)	title, authors, venue, year	keywords, research-area
Company	Description	CompanyName, CompanyType, ShortDescription
WDC-All-Small (WS)	title	category, subcategory, brand, modelno, key-features
WDC-Computer (COM)	title	category, subcategory, brand, modelno, sku, edition
WDC-Camera (CA)	title	category, subcategory, brand, modelno, sku, key-features
WDC-Shoes (SH)	title	category, sku, brand, modelno, colorway, type, edition
WDC-Watch (WAT)	title	brand, sku, gender, modelno, diameter, type, colorway, price

Table 18: Original and enriched attribute for all datasets

Hyper-Parameter	Value	Description(Optional)
Backbone model of \mathcal{F}_{EM}	Qwen-2.5-7B (Yang et al., 2024)	Applied for both Enrichment, Matcher and Selector
Backbone Model of \mathcal{F}_{RAG}	bge-large-en-1.5 (Zhang et al., 2023)	Applied for Blocker \mathcal{M}_{embed}
Learning Rate for \mathcal{F}_{EM}	1e-4	
Learning Rate for \mathcal{F}_{RAG}	1e-5	
τ	0.02	Temperature parameter for contrastive learning of \mathcal{M}_{embed}
K	20	Range of default NN search for Blocker, controlled by pointer $\text{ptr}_s, \text{ptr}_e$
δ	5	Step length for each iteration of pointer $\text{ptr}_s, \text{ptr}_e$
λ	2	iteration of co-training
n	6	number of candidate set for Selector during DPO phase
Max Input Length of \mathcal{F}_{EM}	2048	
Max Input Length of \mathcal{F}_{RAG}	256	
Lora-rank	16	Lora-Rank for fine-tune \mathcal{F}_{EM}
Training epoch	3	Epoch for fine-tune $\mathcal{F}_{EM}, \mathcal{F}_{RAG}$

Table 19: Hyper-Parameter List

Instruction for pt_{SE}

(system message) You are an AI assistant that follows instruction extremely well. User will give you a question. Your task is to answer as faithfully as you can.

(task description) Your task is to determine additional attributes for dataset Amazon-Google. By adding these attributes, you will be leaded to a more clear justification on whether Entity 1 and Entity 2 are the same entity or not.

(instruction) Your output should be in JSON format, only contain the set of enriched attributes. You should take the following Incomplete Entity 1 and Entity 2 as reference.

(input)

Entity 1: { 'title': 'motu digital performer dp5 software music production software', 'manufacturer': '', 'price': 319.95 }

Entity 2: { 'title': 'motu digital performer 5 digital audio software competitive upgrade (mac only)', 'manufacturer': 'motu', 'price': 395.0 }

(output format) Enriched Attributes:

{ Attribute 1:',Attribute 2:' }

Figure 3: Schema Enrichment Prompt pt_{SE}

Instruction for pt_m

(system message) You are an AI assistant that follows instruction extremely well. User will give you a question. Your task is to answer as faithfully as you can.

(task description) You are an expert in computer science and database.

Judge whether record Paper 1 from DBLP, and record Paper 2 from Google Scholar are match or mismatch (refer to the same paper or not), and choose within the given Options.

(input)

Paper 1:

{title: fast algorithms for mining association rules in large databases, authors: R Agrawal, R Srikant, venue: VLDB, year: 1994, keywords: [association rules, large databases, data mining, algorithms, Apriori algorithm, FP-growth algorithm]}

Paper 2:

{title: an efficient algorithm for mining association rules in large databases, authors: a savasere , e omiecinski , s navathe, venue: , year: 1995 }

(output format)

Options: [match,mismatch]

Output format example:{ Output: }

Figure 4: The Prompt pt_m for the Matcher Subtask

Instruction for pt_{enr}

(task description) You are an expert in e-commerce, and you are well known to various goods in Amazon platform. Enrich Entity 1 and Entity 2 with attributes: category/subcategory/platform/edition/type/modelno.

(instruction) Your output should be in JSON format, only contain the value of enriched attributes. You should take the following Incomplete Entity 1 and Entity 2 as reference.

(input)

```
1 Entity 1:{\'title\': \'microsoft visio standard 2007 version upgrade\', \'
  manufacturer\': \'microsoft\', \'price\': 129.95}\n
2 Entity 2:{\'title\': \'adobe cs3 design standard upgrade\', \'manufacturer\': \'\',
  \'price\': 413.99}
```

(output format)

```
1 {"Entity 1": {"title": "", "manufacturer": "", "price": "", "category": "", "
  subcategory": "", "platform": "", "edition": "", "type": "", "modelno": ""},
2 "Entity 2": {"title": "", "manufacturer": "", "price": "", "category": "", "
  subcategory": "", "platform": "", "edition": "", "type": "", "modelno": ""}}
```

Figure 5: Data Imputation (Enrichment) Prompt pt_{enr}

Instruction for pt_s

Task: Entity Matching.

Objective: For the given Entity 1, determine which of the numbered Entity 2 candidates refer to the same real-world entity.

Instructions for your response:

1. Specify the id of candidates that match Entity 1 within `<positive>...</positive>` tags, return in list format.
2. Specify the id of candidates that DO NOT match Entity 1 within `<negative>...</negative>` tags, return in list format.
3. Ensure all candidate indices are covered in either the positive or negative set.

(input)

```
1 Entity 1: {'id': 574, 'title': 'microsoft mappoint 2006 with gps', 'manufacturer'
  : 'microsoft', 'price': 349.0}
```

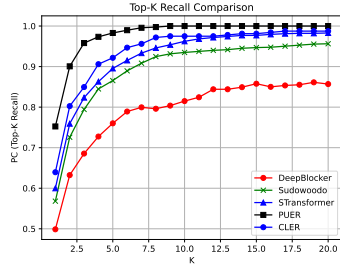
Candidate Options:

```
1 Entity 2 Candidates:
2 {'id': 3029, 'title': 'microsoft mappoint 2006 with gps locator ( pc )', '
  manufacturer': '', 'price': 349.99}
3 {'id': 3190, 'title': 'microsoft ( r ) mappoint ( r ) 2006', 'manufacturer': '', '
  price': 249.99}\n
4 {'id': 2480, 'title': 'microsoft b21-00806 ae mappoint 2006 cd', 'manufacturer': '',
  'price': 50.39}
5 {'id': 1623, 'title': 'language guide for nuvi 350', 'manufacturer': '', 'price':
  79.95}
```

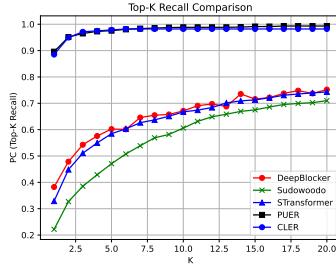
(output format)

`<think> ... </think><positive> [1,2] </positive><negative> [3] </negative>`

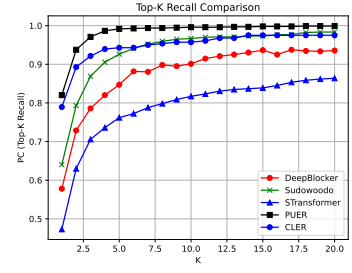
Figure 6: The Prompt pt_s for the Selector Subtask



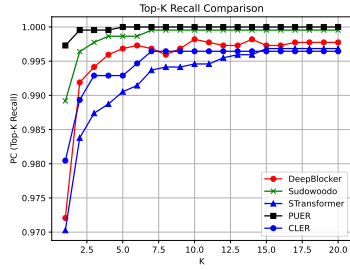
(a) Amazon-Google (varying K)



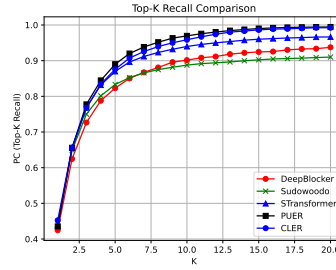
(b) Abt-Buy (varying K)



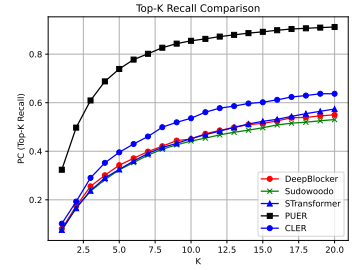
(c) Walmart-Amazon (varying K)



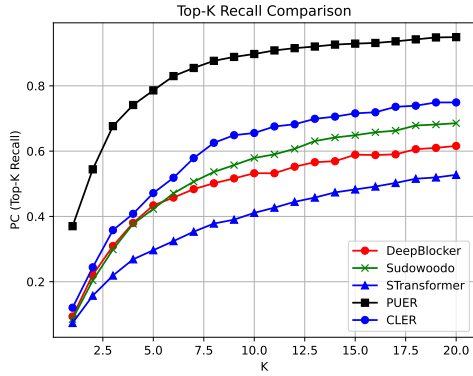
(d) DBLP-ACM (varying K)



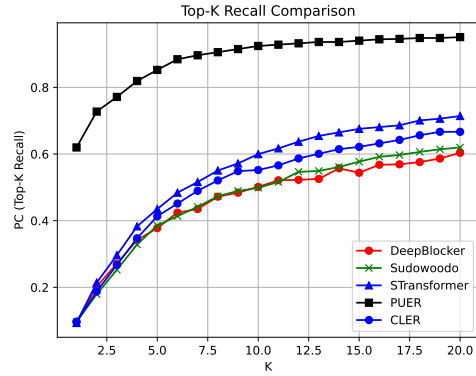
(e) DBLP-Scholar (varying K)



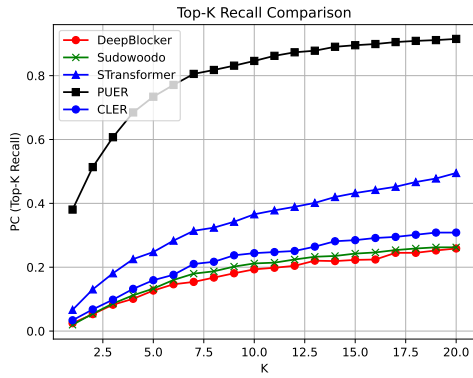
(f) WDC-All-Small (varying K)



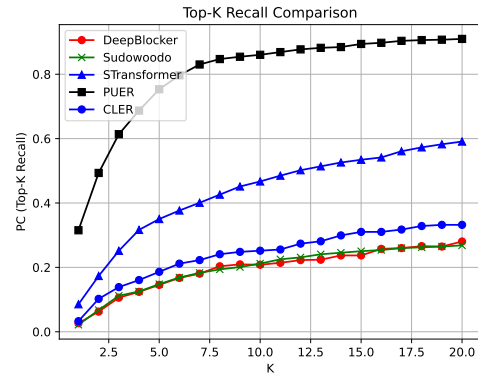
(g) WDC-Computer (varying K)



(h) WDC-Camera (varying K)

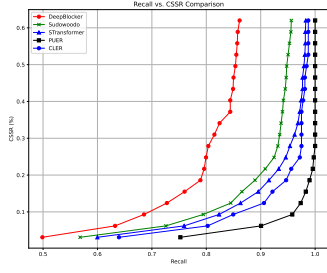


(i) WDC-Shoes (varying K)

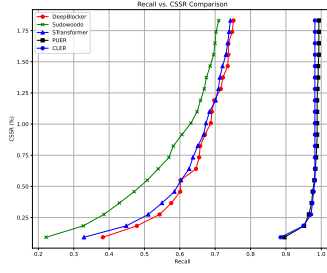


(j) WDC-Watch (varying K)

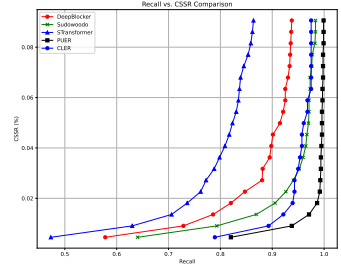
Figure 7: Effectiveness evaluation for Blocker vary K . The curve approaching the **upper left corner** of the figure indicates better performance



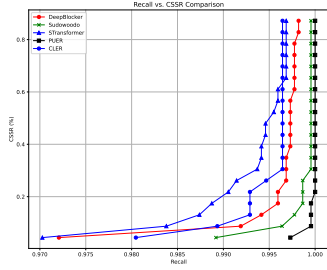
(a) Amazon-Google (varying PC)



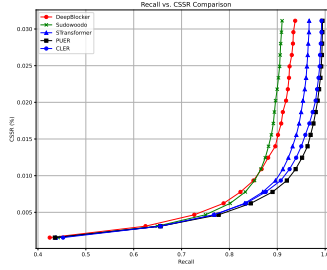
(b) Abt-Buy (varying PC)



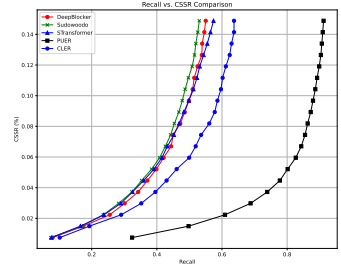
(c) Walmart-Amazon (varying PC)



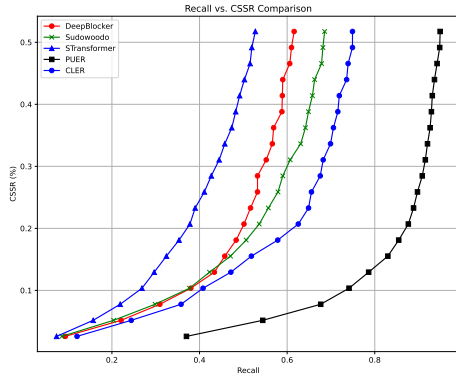
(d) DBLP-ACM (varying PC)



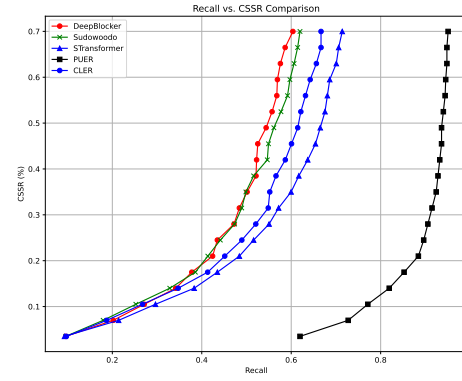
(e) DBLP-Scholar (varying PC)



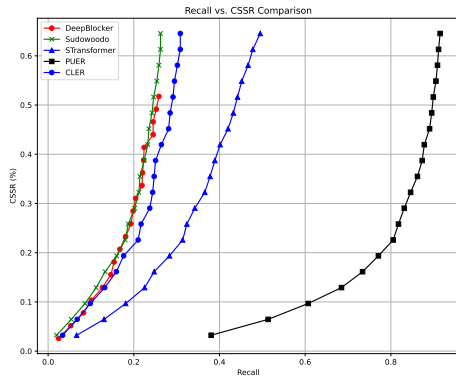
(f) WDC-All-Small (varying PC)



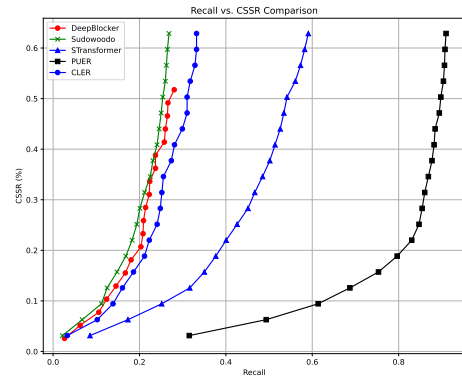
(g) WDC-Computer (varying PC)



(h) WDC-Camera (varying PC)

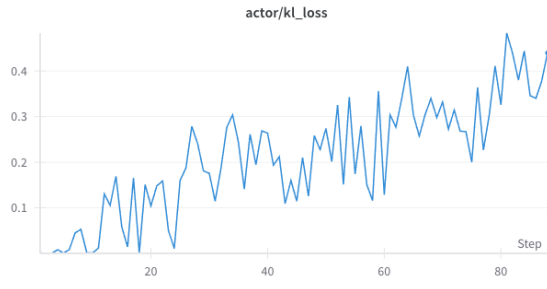


(i) WDC-Shoes (varying PC)

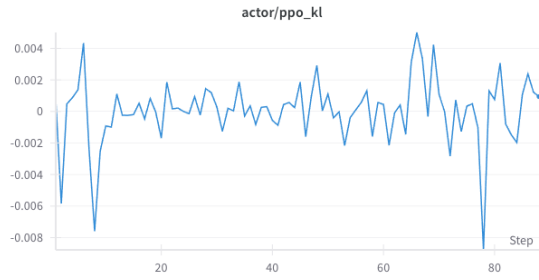


(j) WDC-Watch (varying PC)

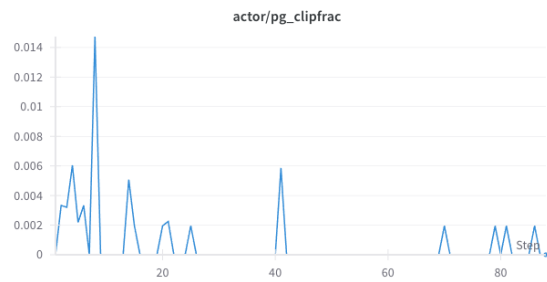
Figure 8: Effectiveness evaluation for Blocker vary PC(*w.r.t.* Recall in figure). The curve approaching the **lower right corner** of the figure indicates better performance



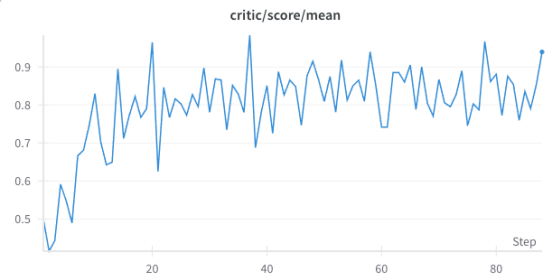
(a) KL Loss Curve for Actor



(b) The curve of estimated KL divergence between old and new policies for Actor



(c) The curve for fraction of policy gradient loss being clipped for Actor



(d) Mean reward score curve for Critic



(e) Mean validation reward score curve for Critic

Figure 9: The training curve for Selector model with dataset AG.