# QJoin: Transformation-aware Joinable Data Discovery Using Reinforcement Learning

Ning Wang
Cornell University
Ithaca, NY, USA
nw366@cornell.edu

Sainyam Galhotra
Cornell University
Ithaca, NY, USA
sg@cs.cornell.edu

## Abstract

Discovering which tables in large, heterogeneous repositories can be joined and by what transformations is a central challenge in data integration and data discovery. Traditional join discovery methods are largely designed for equi-joins, which assume that join keys match exactly or nearly so. These techniques, while efficient in clean, well-normalized databases, fail in open or federated settings where identifiers are inconsistently formatted, embedded, or split across multiple columns. Approximate or fuzzy joins alleviate minor string variations but cannot capture systematic transformations. We introduce QJoin, a reinforcement-learning framework that learns and reuses transformation strategies across join tasks. QJoin trains an agent under a uniqueness-aware reward that balances similarity with key distinctiveness, enabling it to explore concise, high-value transformation chains. To accelerate new joins, we introduce two reuse mechanisms: (i) agent transfer, which initializes new policies from pretrained agents, and (ii) transformation reuse, which caches successful operator sequences for similar column clusters. On the AutoJoin Web benchmark (31 table pairs), QJoin achieves an average F1-score of 91.0%. For 19,990 join tasks in NYC+Chicago open datasets, Qjoin reduces runtime by up to 7.4% (13,747 s) by using reusing. These results demonstrate that transformation learning and reuse can make join discovery both more accurate and more efficient.

## 1 Introduction

Joining tables is the foundation of data integration and data discovery. Analysts routinely need to combine heterogeneous sources to answer questions that no single dataset can resolve. At web scale, systems such as Google Goods [13] and subsequent efforts to organize and navigate data lakes highlight that discovery increasingly entails identifying joinable attributes across thousands of decentralized repositories [9, 10, 22, 33].

However, classical data pipelines often assume equi-joins or lightly normalized schemas. In reality, repositories are far more heterogeneous as identifiers may be embedded or inconsistently formatted (e.g., emails vs. names, codes vs. human labels, or composite keys split across multiple columns). Discovering the correct join thus frequently requires learning transformations that normalize values before matching. A recent work, AutoJoin [34] formalized this challenge as a search problem over transformation operators, demonstrating that transformation-based joins can bridge heterogeneous schemas. Yet applying AutoJoin-style search to large repositories exposes three fundamental bottlenecks:

| | Equi-join | Transformation join |
|---|---|---|
| **Small scale repositories** | [8, 28, 29] | [34] |
| **Large scale repositories** (Data discovery) | [13, 22] | **This work** |

**Table 1: Landscape of join discovery by join type and scale.**

(i) **Limited Expressiveness leading to lower recall:** AutoJoin is designed to identify transformations over unary columns, which does not allow combining multiple columns to identify the join column. (ii) **Combinatorial search cost:** the space of operator chains is computationally expensive and needs to be performed for every candidate pair; (iii) **Higher false positives:** the similarity objective used to identify join can lead to spurious matches.

We demonstrate these with the following example.

*Example 1.* Consider the NYC Open Data election datasets [25, 26], where candidate names appear in split columns (CANDLAST, CANDFIRST, CANDMI) in one table and as a combined field CANDNAME in another. Table 2 contains a snapshot of these tables.

To align them, the key must be composed as:

$$\text{CANDLAST} + \text{', '} + \text{CANDFIRST} [+ \text{' '} + \text{CANDMI}]$$

Traditional join discovery systems including those based on equi-joins or fuzzy similarity cannot infer such multi-attribute transformations automatically. Even specialized transformation-based approaches like AutoJoin [34] fall short because of following reasons. (1) Unary scope: AutoJoin searches only unary operator chains between one source column and one target column. (2) False positives: If AutoJoin tries to join CANDLAST column with CANDNAME, it can lead to several mistakes because many individuals have the same last name.

These issues are not limited to a single dataset. Similar composite-key or embedded-identifier patterns are pervasive in open-data and enterprise repositories.

| CANDLAST | CANDFIRST | CANDMI | | CANDNAME |
|---|---|---|---|---|
| de Blasio | Bill | | | de Blasio, Bill |
| Chen | Ethel | T | | Chen, Ethel T |
| Perkins | Bill | | | Perkins, Bill |
| Chen | Hailing | | | Chen, Hailing |
| Chen | Jin Liang | | | Chen, Jin Liang |
| Qiu | Helen | J | | Qiu, Helen J |
| Sears | Helen | | | Sears, Helen |

**Campaign Expenditures.**       **Funds Payments.**

**Table 2: NYC Open Data example: join requires composing CANDLAST, CANDFIRST, and optional CANDMI into CANDNAME (multi→1). Datasets: *Campaign Expenditures* [25] and *Campaign Public Funds Payments* [26].**

We present QJoin, a transformation-aware, reuse-centric join-discovery framework that amortizes the cost of transformation search across datasets. QJoin models operator selection as a Markov Decision Process (MDP) and trains a Q-learning agent that favors concise, high-reward transformation sequences. Our key intuition is that transformation behaviors exhibit strong regularities across structurally similar columns and thus can be learned once and reused. Reinforcement learning provides a natural mechanism to capture and transfer these recurring operator patterns, allowing QJoin to generalize from prior joins instead of restarting search from scratch. QJoin allows the user to initialize the system with any complex operator that may combine several different columns. Note that plugging in standard notions of reward like jaccard similarity or edit distance between the considered column pairs into the Q-learning approach is not sufficient and yield suboptimal results. Therefore, we develop a novel Longest Continuous Substring based scoring mechanism that is designed for join identification. Specifically, the pipeline proceeds in three intuitive stages:

(i) Candidate Screening and Clustering. Column pairs are first scored using Adjusted Longest Continuous Substring (ALCS) similarity, a contiguous-match metric robust to tokenization noise. Highly aligned pairs are clustered so that families likely to share transformations are discovered together.

(ii) Transformation Learning. Within each cluster, the agent explores transformation chains under a uniqueness-aware reward that jointly optimizes similarity (ALCS) and key distinctiveness, preventing over-aggregation.

(iii) Transformation Reuse. Successful chains are cached and reused within or across clusters, eliminating redundant exploration and accelerating subsequent joins.

This design directly resolves the motivating cases: QJoin learns to compose multi-column keys (e.g., last + ', ' + first + mi) and automatically reuses those operators in similar contexts—achieving both higher recall and greater efficiency.

We make the following contributions.

- We extend the problem of join discovery to consider data transformations.
- We formulate the problem as a a reinforcement-learning formulation for transformation-based join discovery, enabling policy transfer and transformation reuse across candidate pairs.
- Uniqueness-aware reward. We design a reward that couples alignment quality (ALCS) with key distinctiveness, guiding learning toward semantically valid joins.
- Empirical validation: We evaluate QJoin on several benchmarks like AutoJoin Web, NYC Open Data, and Chicago Open Data repositories. QJoin achieves the highest F1 score, while providing over 7.4 % runtime savings, demonstrating that transformation-learning and reuse substantially outperform equi-join based discovery and other baselines.

## 2 Preliminaries

In this section, we define the framework for discovering joins with transformations. We first define a repository of datasets, which acts as a search space for joins and then define join operations, and the formal problem statement.

**Data Repository.** Modern data repositories consist of numerous datasets maintained by different teams, leading to significant heterogeneity in schemas, formats, and data quality. We begin by formalizing the repository structure and the relationships between datasets.

DEFINITION 1 (DATA REPOSITORY). *A data repository $\mathcal{R} = \{D_1, \ldots, D_k\}$ is a collection of datasets, where each dataset $D_i$ has columns $\text{COL}(D_i)$ $= \{c_1^i, c_1^i, \ldots, c_{n_i}^i\}$.*

**Join Operations.** We now define the fundamental join operations that form the basis of our discovery framework. We first formalize the traditional notions of equi-join and approximate fuzzy joins, followed by transformation based joins.

DEFINITION 2 (EQUI-JOIN). *Let $D_a$ and $D_b$ be two different datasets, and let $c_a$ (resp. $c_b$) denote a designated join column in $D_a$ (resp. $D_b$). The equi-join of $D_a$ and $D_b$ on these columns is a combined dataset where rows having same value of $c_a$ and $c_b$ are combined together.*

$$D_a \bowtie^{\text{eq}} D_b = \{(x, y) \in D_a \times D_b \mid x[c_a] = y[c_b]\}.$$

Equi-joins look for exact matches and are easy to discover with LSH based indexing techniques [5]. However, they fail even with minor formatting differences or typos. For example, "Barack Obama" would not join with "Barack J. Obama" in another dataset. Fuzzy join is an approximate join that is robust to such noise.

DEFINITION 3 (FUZZY JOIN). *Given a threshold $0 \leq \delta \leq 1$, the fuzzy join includes all pairs $(x, y)$ such that $\text{sim}(x[c_a], y[c_b]) \geq \delta$, where sim(r,s) denotes the similarity between r and s.*

$$D_a \bowtie^{\text{fuz}} D_b = \{(x, y) \in D_a \times D_b \mid \text{sim}(x[(c_a)], y[(c_b)]) \geq \delta\}.$$

**Fuzzy joins** use approximate string-matching algorithms (e.g., *Levenshtein distance*, n-gram overlap) to score how "close" two values are. This approach handles minor mistakes and variations ("New Yrok" ↔ "New York"), abbreviations ("N.Y." ↔ "NY"), spelling variants ("color" ↔ "colour") that equi-joins would miss. However, confidence in a fuzzy match is directly tied to the similarity score: when that score is low (as in "NY" vs. "New York"), the confidence drops significantly, increasing the risk of false positives or negatives. Additionally, fuzzy joins cannot identify cases where the correspondence between two values is not due to surface-level similarity but rather due to a *systematic transformation*.

For example, fuzzy matching fails to detect equivalences such as CA" ↔ California", IBM" ↔ International Business Machines", or 123-45-6789" ↔ 123456789", since these pairs are not lexically similar but are connected through abbreviation, expansion, normalization, or semantic substitution. In such cases, similarity metrics provide little guidance, and rule-based or learned transformations are needed to bridge the gap.

This limitation motivates the class of **transformation-based joins**, which explicitly model how one value can be transformed into another through a sequence of operations or learned mappings. Rather than relying solely on character-level proximity, these joins reason over structural, linguistic, or semantic transformations to capture deeper correspondences between data values. A transformation operator $\omega$ produces a normalized form of its input. We distinguish two main classes:

**Unary operators** Apply a simple, per-value normalization to a single column. Examples include: lower casing, stripping whitespace, removing special characters, extracting substrings, standardizing date formats, and expanding abbreviations. These operators are computationally efficient but may leave multiple rows ambiguous if raw values remain too similar after transformation.

**Concatenated operators** Merge two or more columns into one composite key. For instance, concatenating the first initial of a first name with the full last name can create unique join keys when direct normalization fails. These operators can resolve ambiguities but require careful selection of columns to combine.

The power of transformation operators lies in their composability. By chaining multiple operators, we can handle complex normalization scenarios that arise in practice.

*Example 2.* Consider three people with first names {Patrick, Audra, Andy}, last names {Zhang, Zhang, Wang}, and emails {PZhang@..., AZhang@..., AWang@...}. Using unary operators alone (e.g., lowercasing, stripping domains) can cause "zhang" to match both "PZhang" and "AZhang," and "a" to match both "Audra" and "Andy." Concatenating the first initial with the full last name yields distinct keys: "pzhang," "azhang," and "awang," resolving the ambiguity.

**Transformation-Based Join Framework.** We now formalize how operator sequences are applied and how joins are defined over transformed columns.

DEFINITION 4 (TRANSFORMATION-BASED JOIN). *Let $D_a$ and $D_b$ be two different datasets with designated join columns $c_a$ and $c_b$, respectively. Let*

$$\Omega = \{\omega_1, \omega_2, \ldots, \omega_m\}$$

*be a finite set of transformation operators. For a fixed maximum length $L_{max} \in \mathbb{N}$, define*

$$\Omega^* = \bigcup_{\ell=0}^{L_{max}} \left\{ \omega_{i_1} \circ \omega_{i_2} \circ \cdots \circ \omega_{i_\ell} \mid (i_1, \ldots, i_\ell) \in \{1, \ldots, m\}^\ell \right\}.$$

*The datasets $D_a$ and $D_b$ are considered to join if $\exists (\Omega_a, \Omega_b) \in \Omega^* \times \Omega^*$ such that $\Omega_a(D_a) \bowtie^{eq} \Omega_b(D_b)$, often denoted as $D_a \bowtie^{eq}_{\Omega_a, \Omega_b} D_b$ and the matched rows are valid join pairs.*

In this definition, $l = 0$ means that there is no transformation and the definition identifies equi-joins. Note that repetition of $\omega_i$ are allowed to capture any complex transformation.

Note that a common challenge of all prior definitions equi-join, fuzzy and transformation aware join is that the join result may not be semantically meaningful. While semantic validity is difficult to measure, we evaluate validity of transformation in terms of uniqueness of joins. Data discovery relies on this common practice to use computable statistics to identify joins, e.g. jaccard similarity for identifying equi-joins [5].

Therefore, we consider $D_a$ and $D_b$ as transformation aware joins only when the columns equi-join after the transformation and the joined pairs are semantically valid. The existence of false positives can be exacerbated under transformations as certain transformations may collapse distinct values into the same form. For instance, consider two ID columns from different tables:

$$ID_1 = \{0123, 0234, 0345\} \quad \text{and} \quad ID_2 = \{0123A, 0234A, 0356A\}.$$

If we apply the single-operator sequence "keep first character" to each column, both sets reduce to {0, 0, 0}, so that every pair of rows appears identical under both exact and fuzzy matching. In reality, only the first two IDs share a common prefix; the third pair (0345 vs. 0356A) is unrelated. This collapse inflates confidence in both equi-join and fuzzy-join, yet yields a cascade of invalid matches, a clear sign of operator overfitting. Hence, transformation-based joins must jointly optimize similarity and *key uniqueness* to ensure semantic correctness.

Given this formalism, we define the problem of transformation-aware join discovery.

PROBLEM 1 (TRANSFORMATION-BASED JOIN DISCOVERY). *Given a repository of tables, a library of operators $\Omega$, and a similarity threshold $\tau$, find for each candidate column pair $(A_i, A_j)$ a transformation chain $\Omega$ such that $A_i$ and $A_j$ join under the given transformation.*
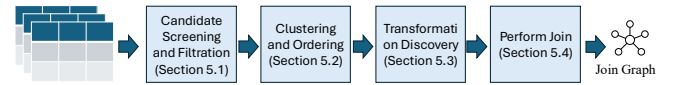

**Figure 1: QJoin Architecture**

## 3 QJoin: Reinforcement-Learning-Based Join Discovery

We now present QJoin, a transformation-aware framework for discovering joins in heterogeneous repositories. QJoin formulates transformation search as a *reinforcement learning* problem, in which an agent learns to select and compose transformation operators that yield semantically valid and high-quality joins.

### 3.1 Our Intuition

Transformation patterns often recur across repositories. For example, many "name" columns require similar normalization pipelines (e.g., lowercase → strip → concat) across datasets in education or civic domains. Instead of re-exploring the operator space for every new column pair, QJoin uses reinforcement learning to *learn reusable operator policies* that generalize across structurally similar columns.

The key intuition is that transformation search exhibits regularities: operator sequences that improve similarity and preserve key uniqueness in one setting are likely to succeed in others. By representing each transformation step as an action in a Markov Decision Process (MDP), QJoin learns these patterns through trial and reward, then transfers the knowledge to new join tasks.

### 3.2 System Overview

QJoin consists of the following five steps.
**Stage 1: Candidate Screening and Filtration.** Given a repository $\mathcal{R} = \{D_1, \ldots, D_k\}$, the first stage enumerates cross-dataset column pairs and computes lightweight similarity proxies. Two complementary metrics are used: *(i) Jaccard q-gram overlap*, a fast lexical indicator, and *(ii) Adjusted Longest Continuous Substring (ALCS)* similarity, which captures contiguous token alignment even under noise. The combination of the two provides both recall and structural awareness. Column pairs whose scores exceed a global threshold $\delta$ and rank within the top-$k$ per table pair are retained for further analysis (details in Section 5.1).

This pruning step reduces the candidate space by several orders of magnitude while ensuring that potentially joinable pairs—including those requiring transformations—are preserved.

**Stage 2: Clustering and Ordering** Many surviving pairs share similar syntactic or semantic patterns (e.g., names, dates, locations). We embed each pair using its similarity descriptors (Jaccard, ALCS, entropy, length ratio) and perform hierarchical clustering (Section 5.2) to obtain coherent groups likely to share transformation behaviors. Within each cluster, pairs are ordered by their pre-scores, so that successful transformations learned on simple pairs can bootstrap harder ones. This structure establishes the units of reuse: all subsequent learning and transfer operate at the cluster level.

**Stage 3: Reinforcement-Learning-Based Transformation Discovery** For each cluster, QJoin formulates transformation discovery as a Markov Decision Process $\mathcal{M} = \langle S, A, P, R, \gamma \rangle$. Each *state* encodes the current transformation context (ALCS score, uniqueness ratio, operator depth); each *action* corresponds to applying an operator (e.g., lowercase, strip, concat); the *reward* integrates alignment improvement and key distinctiveness. An episode terminates once further transformations yield negligible gain. The agent's Q-table is updated online using the reward defined in Section 4.1.

**Stage 4: Adaptive Join Execution** Once optimal transformation chains are learned, the transformed columns are joined using an *adaptive fuzzy-join operator* based on ALCS (Section 5.4). The threshold dynamically adjusts to token length and data variability—tight for short identifiers, relaxed for long phrases— to maintain both precision and recall.

**Stage 5: Validation and Robustness** The system consolidates learned transformations using an *update-and-selection wrapper* (Section 5.5) that validates each operator chain across multiple directions and retains only those with consistent reward improvements. Validated chains are stored in a *Transformation Reuse Library*, indexed by column metadata and cluster ID. When new column pairs arrive, QJoin first attempts zero-shot reuse; if no exact match exists, the closest prior agent's Q-table is used to warm-start training.

Overall, the workflow can be viewed as a closed-loop control system: pre-scores supply priors, reinforcement learning refines beliefs via reward feedback, and the reuse library embodies long-term memory. We now discuss the RL based formulation and then provide the details of the other components.

## 4 MDP Formulation

QJoin represents the transformation-discovery process as a Markov Decision Process (MDP) defined by the tuple: $\mathcal{M} = \langle S, \mathcal{A}, \mathcal{P}, R, \gamma \rangle$, where $S$ is the state space, $\mathcal{A}$ is the set of actions (transformation operators), $\mathcal{P}$ is the transition function, $R$ is the reward, and $\gamma$ is the discount factor.

**State.** A state $s_t \in S$ encodes the current transformation context for a column pair, represented by features such as: (1) token-level similarity statistics (ALCS, Jaccard), (2) uniqueness ratio of transformed keys, and (3) length and composition of the current operator chain. States are continuous and capture the evolving join quality as transformations are applied.

**Action.** Each action $a_t \in \mathcal{A}$ corresponds to applying a transformation operator $\omega \in \Omega$ (e.g., lowercase, trim, concat). Applying an

operator transitions the agent to a new state $s_{t+1}$ that reflects the updated column pair.

**Transition.** The transition function $\mathcal{P}(s_{t+1} \mid s_t, a_t)$ is deterministic given the operator semantics and current column values, though the resulting join score is stochastic due to data variability.

**Reward.** At each step, the agent receives a reward $R(s_t, a_t)$ based on both similarity improvement and key uniqueness. Purely similarity-based measures (e.g., Jaccard) are insufficient because they reward non-unique transformations. QJoin's reward function explicitly balances alignment quality and distinctiveness, as detailed in Section 4.1.

**Episode Termination.** An episode terminates when (i) the agent reaches a maximum operator depth $L_{\max}$, or (ii) no further improvement in reward is observed. The final transformation chain $F^*$ is taken as the composition of operators along the highest-reward trajectory.

*4.0.1 Learning and Reuse Mechanisms.* QJoin trains the agent using Q-learning:

$$Q(s_t, a_t) \leftarrow R(s_t, a_t) + \gamma \max_{a'} Q(s_{t+1}, a'),$$

where $\gamma$ is the discount factor.

This choice accelerates adaptation across related tables by immediately propagating successful outcomes:

- **Agent Transfer.** When a new cluster is encountered, its Q-table is initialized with parameters learned from the most similar prior cluster (measured by schema and token overlap). This warm start accelerates convergence by exploiting cross-cluster regularities.
- **Transformation Cache.** High-performing operator chains are stored in a cache indexed by feature signatures of column pairs. When new pairs resemble cached patterns, QJoin reuses these sequences directly, bypassing the RL exploration phase.

Together, these mechanisms allow QJoin to generalize beyond individual column pairs, making transformation learning scalable and transferable across the repository.

## 4.1 Reward Design

The reward function lies at the core of QJoin. It guides the agent to construct transformation chains that both (i) increase join alignment and (ii) preserve key distinctiveness. A purely similarity-driven reward would encourage degenerate transformations that collapse multiple values into identical strings, while a purely uniqueness-driven reward would reject legitimate normalizations (e.g., removing punctuation). QJoin therefore integrates both through a *composite, uniqueness-aware reward.*

**Design Goals and Pitfalls.** We design the reward $R(s_t, a_t)$ with three goals:

(i) **Balance similarity and uniqueness.** The agent should increase column alignment while avoiding transformations that collapse distinct keys.

(ii) **Encourage efficient exploration.** Simpler or cheaper operators should be preferred unless complex ones (e.g., concat) yield substantial gain.

(iii) **Enable transfer and reuse.** The same reward formulation should generalize across datasets, producing policies that can be reused by similar clusters.

A naive similarity metric such as Jaccard or cosine, when used directly as $R(s_t, a_t)$, violates these goals:

- It rewards transformations that trivially increase overlap (e.g., mapping every token to "a"), producing false joins.
- It fails to differentiate operator effects (e.g., `strip` vs. `concat`) or the structure of the column.

To overcome these issues, QJoin integrates three interacting reward components. Specifically, we (a) mix similarity *gains* with uniqueness preservation, (b) gate rewards by the fraction of rows improved, and (c) differentiate operator classes (direct vs. concatenation) with sensible weighting.

*4.1.1 Similarity Component: ALCS Gain.* The first component measures how much a transformation improves column alignment. In our MDP formulation, transformations are discovered incrementally through a sequence of operator applications. This sequential nature imposes specific requirements on our similarity metric, it must provide meaningful feedback not only for complete transformations but also for intermediate states. Traditional similarity metrics like edit distance or Jaccard similarity often fail to capture the nuanced progress made by partial transformations, potentially misleading the learning process.

- **Monotonic, Parameter-Free Feedback.** Fixed hyperparameters (e.g. $q$ in some similarity measures) may be optimal for one dataset but fail on another. A monotonic, parameter-free function guarantees that any extension of the core match yields a higher score. Then the agent immediately recognizes and reinforces partial gains, preventing stalls due to mis-tuned thresholds.
- **Focus on Core Substring Growth.** Real-world join keys (e.g. IDs, email local-parts) typically hinge on an uninterrupted block of text. Emphasizing the longest continuous overlap isolates this true key from peripheral noise. By rewarding the growth of the main substring, the agent learns to prioritize transformations that reveal the true join key, accelerating convergence.
- **Tolerance to Peripheral Changes.** Transformations may insert, delete, or reorder characters outside the core block. Penalizing these harmless edits misleads learning. By ignoring extraneous changes and scoring only the core overlap, the function prevents negative feedback for correct matches, enabling robust exploration.

Therefore, we use Adjusted Longest Continuous Substring (**ALCS**):

$$\text{ALCS}(S_1, S_2) = \frac{\text{LCS}(S_1, S_2)}{\frac{1}{2}\left(|S_1| + |S_2|\right)},$$

where LCS is the longest common *substring* (length $\geq n$ for significance). For a source column $c_a$ with values $S_i$ and a target $c_b$ with values $T_j$, define (mean-max style):

$$A_{ij}^{\text{prev}} = \text{ALCS}(S_i^{\text{prev}}, T_j), \quad A_{ij}^{\text{new}} = \text{ALCS}(S_i^{\text{new}}, T_j),$$

where $A_{ij}^{\text{prev}}$ denotes the ALCS before applying the operator and $A_{ij}^{\text{new}}$ denites the ALCS after applying the operator on $c_a$. Denote

$$\Delta_i = \max_j A_{ij}^{\text{new}} - \max_j A_{ij}^{\text{prev}}.$$

The aggregate improvement is $\Delta\text{ALCS} = \sum_i \Delta_i$. We grant the *ALCS reward* $R_{\text{ALCS}} = \alpha_{\text{ALCS}} \cdot \Delta\text{ALCS}$ *only if* a sufficient fraction of rows

benefit:

$$p_{\text{ALCS}} = \frac{1}{m}\left|\{\, i : \Delta_i > 0 \,\}\right| \geq p_{\text{ALCS}}^{\text{min}}.$$

This proportional-impact based gating blocks outlier-only gains and provides monotone, parameter-free feedback as contiguous blocks are revealed by partial transformations.

*4.1.2 Uniqueness Component: Duplicate Penalty.* Similarity alone is insufficient—some transformations increase overlap by erasing distinctions. To preserve discriminative power, QJoin penalizes any drop in key uniqueness.

Let $x_i$ be the best-matching target token chosen for row $i$ (via $\max_j A_{ij}$ with an LCS tie-break), and let $f(x)$ be the frequency of token $x$ among $\{x_i\}$. The duplication score is

$$\phi = \sum_{i=1}^{m} \max\left(0, f(x_i) - 1\right).$$

With $\phi^{\text{prev}}$ and $\phi^{\text{new}}$ before/after a step, define

$$\Delta\text{dup} = \begin{cases} \dfrac{\phi^{\text{new}} - \phi^{\text{prev}}}{\phi^{\text{prev}}}, & \phi^{\text{prev}} > 0, \\ 0, & \phi^{\text{prev}} = 0. \end{cases} \quad \text{and} \quad R_{\text{uniq}} = -\alpha_{\text{uniq}} \cdot \Delta\text{dup}.$$

Thus, increases in duplication (loss of distinctiveness) incur negative reward. As with ALCS, we gate by the fraction of rows whose duplicates do *not* worsen: $p_{\text{uniq}} \geq p_{\text{uniq}}^{\text{min}}$.

*4.1.3 Operator-Aware Composite Rewards.* Certain operators inherently alter column semantics more drastically than others. QJoin incorporates an *operator cost* term, $C_{\text{op}}(a_t)$, which regularizes exploration by penalizing complex, multi-column, or high-similarity operations. For example, Concatenation operators have the inherent advantage of increasing the uniqueness by continuously concatenating new columns. Then we should weigh the similarity reward factor to constrain it. For unary operators, extraction often leads to losing the uniqueness. Then we add the gated uniqueness and ALCS-based rewards for unary operators.

We distinguish *direct* (unary, potentially lossy) vs. *concatenation* (multi-column, potentially information-preserving) operators.

*Unary operators.* We use a conservative sum with gates:

$$R_{\text{direct}} = \underbrace{\tilde{R}_{\text{ALCS}}}_{\text{gated}} + \underbrace{\tilde{R}_{\text{uniq}}}_{\text{gated}},$$

granting positive reward only when alignment improves for enough rows and duplicates do not materially increase.

*Concatenation operators (adaptive weighting).* Concatenation can unlock one-to-one keys; we adapt weights based on the similarity ($\sigma$) landscape over candidate pairs $\mathcal{P}$. Let

$$\sigma_{\max} = \max_{(c_i, c_j) \in \mathcal{P}} \sigma(c_i, c_j), \quad N_1 = \left|\{(c_i, c_j) : \sigma(c_i, c_j) > \tau_{\text{high}}\}\right|,$$

$$N_2 = \left|\{(c_i, c_j) : \sigma_{\max} - \sigma(c_i, c_j) > \tau_{\text{diff}}\}\right|.$$

If $N_{\text{high}} = \max(N_1, N_2) = 1$ (a single dominant candidate), we boost alignment weight and relax uniqueness slightly:

$$\alpha_{\text{ALCS}} = \alpha_{\text{ALCS}}^{\text{high}}, \quad \alpha_{\text{uniq}} = \alpha_{\text{uniq}}^{\text{low}};$$

otherwise use default balanced weights ($\alpha_{\text{ALCS}}^{\text{def}}, \alpha_{\text{uniq}}^{\text{def}}$). The concatenation reward mirrors $R_{\text{direct}}$ but with these adaptive coefficients.

*4.1.4 Putting It Together.* For a chain of steps $F$ (length $|F|$), the reward at $i$th steo is

$$R(F_i) = R - \lambda * i,$$

penalizing longer sequences. This shaping prevents degenerate similarity hacks, preserves distinctiveness, and favors short, reusable chains.

At each iteration, the RL Engine evaluates $R(s_t, a_t)$, updates its Q-table:

$$Q(s_t, a_t) \leftarrow (1 - \alpha)Q(s_t, a_t) + \alpha\big[R(s_t, a_t) + \gamma \max_{a'} Q(s_{t+1}, a')\big],$$

and logs the best transformation chains into the Reuse Library.

## 5  QJoin Implementation

QJoin proceeds as a sequence of lightweight, progressively more informed decisions. Each stage prunes the search space, so that by the time reinforcement learning begins, the agent faces a compact, high-value subset of candidates. The system's power lies in the synergy between statistical pre-scoring, adaptive reward design, and cross-pair transformation reuse. Sec. 5.1 describes efficient pre-scoring and filtering; Sec. 5.2 presents the clustering and ordering; Sec. 5.3 introduces the RL formulation; and Sec. 5.4–5.5 cover join execution, selection, and reuse.

### 5.1  Compute Pre-Scores

Searching for joins across all possible column pairs in a repository is computationally prohibitive. Even a moderate repository with thousands of columns yields millions of pairs, and only a tiny fraction are meaningful. Hence, we begin with a lightweight *pre-scoring phase* that estimates the likelihood of joinability to run at repository scale, yet expressive enough to surface hidden matches that simple token overlap might miss.

**Approach.** Given a repository $\mathcal{R} = \{D_1, \ldots, D_K\}$, where each dataset $D_i = \{C_{i1}, C_{i2}, \ldots, C_{in_i}\}$, we compute pre-scores for all cross-dataset column pairs $(C_{ip}, C_{jq})$ with $i \neq j$. The pre-score $s_{ij}$ quantifies how promising the pair is for further transformation learning. We implement two complementary scoring strategies—one surface-level and one transformation-aware:

**(a) Jaccard pre-score (surface-level).** We first compute a *Jaccard similarity* on $q$-grams as a rapid, structural proxy:

$$s_{ij}^J = \frac{1}{m} \sum_{u=1}^{m} \max_v \text{Jaccard}_{q\text{-gram}}(S_u, T_v),$$

where a proportion $p$ of source items $S_u \in C_{ip}$ are sampled, $q$-grams are extracted for both source and target columns, and $m = p|C_{ip}|$ denotes the sample size. Indexing techniques help to efficiently implement this low-cost filter.

**(b) ALCS-based direct-operator pre-score (transformation-aware).** When simple token overlap fails (e.g., due to prefixes, concatenations, or embedded keys), we evaluate whether lightweight transformations already unlock alignment. We sample the same proportion $p$ from $C_{ip}$, apply a small set of direct operators $\Omega_i, \Omega_j$ (e.g., `lower`, `strip`, `concat`), and compute:

$$s_{ij}^A = \frac{1}{m} \sum_{u=1}^{m} \max_v \text{ALCS}(S_u, T_v), \quad s_{ij}'^A = \frac{1}{m} \sum_{u=1}^{m} \max_v \text{ALCS}(S_u', T_v'),$$

where $S_u' \in \Omega_i(C_{ip})$ and $T_v' \in \Omega_j(C_{jq})$. The improvement $\Delta s_{ij}^A = s_{ij}'^A - s_{ij}^A$ reflects how much transformation helps, distinguishing easy pairs (high $s_{ij}^A$) from structurally misaligned pairs (low $s_{ij}^A$ but large $\Delta s_{ij}^A$).

*Interpretation.* Jaccard pre-scoring offers broad recall at negligible cost, which is ideal for large repositories while ALCS pre-scoring probes deeper structural similarity, highlighting columns that *could* align after normalization. Together, they form a two-tier sieve: the former detects direct matches, the latter reveals hidden ones. Subsequent stages focus computational effort only on these promising regions of the search space.

Each column pair now has a concise descriptor $(s_{ij}^J, s_{ij}^A, \Delta s_{ij}^A)$ that summarizes its join potential. These descriptors feed directly into the next stage's filtering logic.

*5.1.1  Filter.* Even after pre-scoring, many column pairs remain uninformative or redundant. Therefore, this step performs filtration to prune uninformative candidates. This step is based on ideas from blocking in entity resolution [27].

**(1) Absolute thresholding.** We first remove pairs whose pre-score (either $s_{ij}^J$ or $s_{ij}^A$) falls below a global threshold $\delta$. Formally:

$$\text{Keep}(C_{ip}, C_{jq}) \text{ iff } \max(s_{ij}^J, s_{ij}^A) \geq \delta.$$

This step discards clearly incompatible pairs, those with near-zero overlap—even if they appear in large tables or popular domains.

**(2) Relative top-$k$ per table pair.** Similarity scores can vary widely across domains. For each table pair $(D_i, D_j)$, we retain only the top-$k$ column pairs by score:

$$\mathcal{S}_{ij}^{\text{filtered}} = \text{TopK}_k\Big\{ (s_{ip,jq}, \max(s_{ij}^J, s_{ij}^A)) : p \in D_i, q \in D_j \Big\}.$$

This relative selection preserves diversity—ensuring each table contributes candidates—while controlling per-domain imbalance.

*Intuition.* Absolute thresholding eliminates uninformative pairs regardless of scale, while per-table top-$k$ selection ensures fairness: highly similar domains (e.g., multiple product tables) cannot crowd out rare but valid cross-domain joins (e.g., customer–region). Together, they adapt to both sparse and dense repositories, scaling computational cost with join density rather than repository size.

often one to two orders of magnitude smaller than the Cartesian product. Each pair retains its pre-score features, ready for clustering and ordering in the next step.

### 5.2  Cluster and Order

Many candidate join pairs share similar syntactic or semantic characteristics e.g., name fields across datasets, or location strings differing only by formatting. The key insight of this component is that column pairs with similar similarity profiles tend to admit *similar transformation sequences*. By clustering such pairs together and ordering them strategically, we both enhance reusability and accelerate convergence of transformation learning.

*5.2.1  ClusterPairs.* Each surviving column pair $(C_{ip}, C_{jq})$ from the filtered set is represented by a feature vector $x_i \in \mathbb{R}^F$ summarizing its similarity profile (e.g., pre-scores, ALCS statistics, text-length

ratios, token entropy, and domain hints). Let $X = [x_1, \dots, x_n]^\top$ denote the resulting matrix for all $n$ pairs.

We compute pairwise Euclidean distances:

$$d_{ij} = \|x_i - x_j\|_2, \quad 1 \le i < j \le n,$$

and perform average-linkage hierarchical clustering:

$$Z = \texttt{linkage}(D, \texttt{average}),$$

where $D = \{d_{ij}\}$ is the condensed distance matrix. Clusters are extracted via:

$$\ell_i = \texttt{fcluster}(Z, t, \texttt{criterion} = \texttt{distance}) - 1, \quad i = 1, \dots, n.$$

Each cluster $C_k = \{\, i : \ell_i = k \,\}$ groups pairs with similar transformation behavior, and the cluster centroid is:

$$\mu_k = \frac{1}{|C_k|} \sum_{i \in C_k} x_i.$$

Outcome. The repository is now partitioned into $K$ coherent clusters of column pairs, each representing a distinct structural or semantic relation type (e.g., names, dates, codes). This structure forms the basis for the next step: determining an efficient processing order within and across clusters.

*5.2.2 DetermineClusterOrder.* Within each cluster, some pairs are easier to align (high pre-score, near-duplicates), while others are harder (lower similarity, multiple transformations needed). Starting from easy pairs allows us to learn transformations quickly and reuse them downstream. Thus, we order pairs to *front-load the most promising and reusable cases.*

We compute for each column $C_i$ a cumulative similarity score:

$$\text{MaxSim}(C_i) = \max \{\, \text{sim}(C_i, C_k),\ \text{sim}(C_i, C_\ell), \dots \},$$

and for each table pair $(T_L, T_K)$, we aggregate:

$$\text{TotalSim}(T_L, T_K) = \sum_{C_i \in T_L} \text{MaxSim}(C_i).$$

Pairs are then globally ranked by TotalSim, with priority given to those satisfying: (i) at least one column has been successfully transformed earlier, and (ii) the ALCS score exceeds a per-cluster percentile threshold.

## 5.3 RL: Process Each Pair in Order

Now, we cast transformation discovery as a *Markov Decision Process (MDP)* and train a reinforcement-learning (RL) agent via Q-learning to efficiently explore and optimize operator chains. Each state in the MDP encodes the current configuration of transformed data; each action applies one transformation operator; and the reward (Sec. 4.1) measures improvement in similarity and uniqueness. Through repeated interaction, the agent learns which transformation sequences maximize join quality while preserving discriminative structure.

*5.3.1 Agent Initialization and Setup.* We begin by defining the environment and learning parameters.

**State Space.** A state $s_t$ summarizes the current transformation context, including: (i) the active set of transformed columns $C_{\text{set}}$, (ii) their cumulative similarity and uniqueness statistics, and (iii) metadata such as transformation depth.

**Action Space.** Actions correspond to transformation operators $\Omega = \{\omega_1, \omega_2, \dots, \omega_K\}$ (e.g., `lower`, `strip`, `concat`, `substr`). At each step, the agent selects an operator $\omega_i$ and applies it to one column in $C_{\text{set}}$.

**Initialization.** Each operator is assigned an initial uniform probability $Pr(\omega_i) = 1/K$. The agent's learning rate $\alpha$ determines how quickly new experience overrides past knowledge, while the exploration parameter $\epsilon$ controls the probability of trying random transformations rather than exploiting the best-known ones. Reward coefficients $(\lambda_1, \lambda_2, \lambda_3)$ weight ALCS gain, uniqueness preservation, and operator cost respectively, directly tying this stage to the reward formulation in Sec. 4.1.

This setup encourages data-efficient learning: the agent does not require full exploration of the exponential space of operator chains. Instead, Q-learning's incremental updates bias it toward transformations that repeatedly yield positive reward, pruning unproductive sequences early.

*5.3.2 Representative Sample Selection.* Training directly on all rows is unnecessary and costly. Moreover, the distribution of similarities is typically skewed—some value pairs are trivially similar, others nearly disjoint. Balanced sampling ensures that the agent encounters both easy and hard examples.

We compute the ALCS similarity matrix $\{\text{ALCS}(r_i, r_j) \mid r_i \in C_i,\ r_j \in C_j\}$ and, for each source record $r_i$, extract the best target match $\max_j \text{ALCS}(r_i, r_j)$. Using these maxima as features, we apply $k$-means clustering ($k = 3$) to partition the data into low-, medium-, and high-similarity strata. From each cluster, we sample proportions $(p_1, p_2, p_3)$ respectively, ensuring representation across difficulty levels. For each sampled $r_i$, we retain its top-$k$ matching targets to form an ALCS submatrix that guides the RL process. Stratified sampling prevents the agent from overfitting to high-similarity cases. By exposing it to challenging examples early, the agent learns robust transformation patterns that generalize beyond the sampled subset.

*5.3.3 Iterative Learning Process.* Learning proceeds in alternating exploration and exploitation phases. During exploration (probability $\epsilon$), the agent randomly selects a column $C_{\text{chosen}} \in C_{\text{set}}$ and applies a transformation $\omega_i$ sampled from $Pr(\omega_i | C_{\text{chosen}})$. This randomness enables discovery of non-obvious yet beneficial transformations. During exploitation (probability $1 - \epsilon$), the agent simulates each operator on all columns, computes the expected reward

$$R(s_t, a_t) = \lambda_1 \Delta_{\text{ALCS}}(s_t, a_t) - \lambda_2 P_{\text{dup}}(s_t, a_t) - \lambda_3 C_{\text{op}}(a_t),$$

and applies the transformation yielding the highest positive gain. If no operator yields positive reward, the current configuration is reset, preventing accumulation of detrimental changes.

*5.3.4 Transformation Application and Evaluation.* Each iteration maintains two working column sets, $C_{\text{set},i}$ and $C_{\text{set},j}$, including both raw and concatenated variants. Transformation proceeds in three steps:

(i) Apply the chosen operator to update $C_i \leftarrow \omega(C_i)$.
(ii) Concatenate relevant columns into a composite representation $col' = \texttt{concatenate}(C_{\text{set}})$.
(iii) Evaluate the reward increment $\Delta R = R(col_i', col_j) - R(col_i, col_j)$.

Operator probabilities are updated accordingly:

$$Pr(\omega_i|C_{\text{chosen}}) \leftarrow \begin{cases} Pr(\omega_i|C_{\text{chosen}}) + \alpha(1 - Pr(\omega_i|C_{\text{chosen}})), \text{if } R > 0, \\ Pr(\omega_i|C_{\text{chosen}}) - \alpha Pr(\omega_i|C_{\text{chosen}}), \text{otherwise.} \end{cases}$$

Probabilities are renormalized to sum upto 1.

**Convergence and Termination**

Training terminates when one of the following holds: (i) the ALCS similarity exceeds a threshold $\tau_{\text{sim}}$; (ii) reward improvement falls below $\epsilon_{\text{tol}}$ for $n$ consecutive iterations; or (iii) a maximum iteration limit $T_{\text{max}}$ is reached.

## 5.4 Perform Join

Even with well-learned transformations, perfect normalization across real-world datasets is rare. Identifiers may contain typos, abbreviations, or minor structural differences that strict equi-joins. To effectively perform a fuzzy join, we design an *adaptive* strategy based on the Adjusted Longest Common Substring (ALCS) metric. It dynamically adjusts its similarity threshold to the length and variability of the data, ensuring high precision on short tokens and adequate tolerance on long, noisy text.

Given two transformed columns $col_a$ and $col_b$ from datasets $D_a$ and $D_b$ ($a \neq b$), we evaluate all record pairs $(r_i, r_j)$ drawn from the transformed value sets $\Omega_a(col_a)$ and $\Omega_b(col_b)$ respectively. The base similarity between two values is

$$\text{ALCS}(r_i, r_j) = \frac{\text{LCS}(r_i, r_j)}{\frac{1}{2}(|r_i| + |r_j|)}.$$

For each $r_i$ in the source column, we identify its best target match:

$$\text{BestMatch}(r_i) = \max_j \text{ALCS}(r_i, r_j),$$

and declare a candidate join when the similarity exceeds an adaptive threshold (defined below).

*Length-Aware Adaptivity.* A single fixed threshold is insufficient: short identifiers (e.g., "NY") demand near-perfect matching, whereas longer phrases ("New York City, NY") should tolerate mild variations. We therefore compute the minimum average text length across both columns:

$$l_{\text{min}} = \min\left(\frac{1}{m}\sum_{i=1}^{m}|r_i|, \frac{1}{n}\sum_{j=1}^{n}|r_j|\right),$$

and set the pair-specific distance tolerance $d$ and similarity weight $\alpha_{\text{sim}}$.

*Robust Threshold Computation.* To avoid sensitivity to outliers, we combine two complementary statistics:

$$\text{ALCS}_{\text{mean}} = \frac{1}{m}\sum_{i=1}^{m}\max_j \text{ALCS}(S_i, T_j),$$

$$\text{ALCS}_{\text{median}} = \text{mean}\big(\text{cluster}_{\text{middle}}(\{\max_j \text{ALCS}(S_i, T_j)\})\big),$$

where the median term is computed over the "middle" similarity cluster obtained via $k$-means on best-match scores. The final join threshold is $\text{thr}_{\text{join}} = \max(\text{ALCS}_{\text{mean}}, \text{ALCS}_{\text{median}})$. A join between rows $S_i$ and $T_j$ is executed if $\text{ALCS}(S_i, T_j) \geq \text{thr}_{\text{join}} - d$.

This adaptive rule captures three key intuitions: (i) longer strings invite greater tolerance since local edits have smaller proportional impact; (ii) averaging and median statistics jointly stabilize the threshold against noise and skewed distributions; and (iii) using the maximum of both metrics biases the system toward precision when data are clean and toward recall when noise is high. In effect, the algorithm emulates how an analyst would tune a fuzzy join threshold after inspecting a few sample matches—except that it does so automatically for every pair.

## 5.5 Robustness with an update and Selection Mechanism

Even with adaptive learning, the initial choice of column pairs may mislead the agent. Some pairs may achieve high reward due to spurious correlations rather than genuine joinability. To guard against this, QJoin includes an *update-and-selection wrapper* that validates and consolidates learned transformations across multiple candidate pairs. The goal is to retain only those transformation sequences that generalize across directions (A→B and B→A) and yield consistently positive reward.

*Approach.* Let $\mathcal{P} = [p_1, p_2, \ldots, p_n]$ denote the list of transformed column pairs obtained from the RL process. The wrapper compares these candidates sequentially and selects the one that maximizes the composite reward.

This wrapper acts as a validation layer, enforcing monotonic improvement in composite reward. It ensures that transformations adopted by the system are not merely pair-specific optimizations but produce genuine cross-pair gains. Whenever the reward of a new pair exceeds that of the current best, the system updates its "canonical" transformation. Otherwise, it conservatively retains the previous one—preventing reward drift and instability.

*5.5.1 Transformation Reuse Across Clusters.* In large repositories, structurally similar column pairs recur frequently: different datasets often encode the same concept (e.g., `name`, `vendor`, `region`) using analogous formats. Relearning transformations for each pair would be redundant. To avoid this, QJoin maintains a *Transformation Reuse Library* that stores successful operator sequences, together with their contextual metadata (cluster membership, operator probabilities, and reward traces). When a new pair is encountered, the system first attempts zero-shot reuse or warm-starts the RL agent from the most similar stored case.

*Approach.* Let $D_A$ and $D_B$ be the current datasets on sides A and B, and $C$ the clustering of column pairs (Sec. 5.2). For each cluster $C_k \subseteq \{(c_i, c_j)\}$, we maintain stored transformation blocks $\Omega_{\text{stored},i}$ and agent states $A_i$. When a new pair $(c'_a, c'_b) \in C_k$ is processed:

(i) **Check direct compatibility:** If all column references in $\Omega_{\text{stored},i}$ exist in the new schema, apply the stored transformation directly and evaluate reward. If the gain remains positive, the transformation is accepted as-is.

(ii) **Find equivalent replacements:** If missing keys exist, find cluster-equivalent replacements by substituting semantically similar columns from the same cluster.

(iii) **Generate candidate reuses:** Enumerate all combinations of replacements (including a "no-change" $\perp$ option) For each candidate, evaluate its reward on sampled rows.

(iv) **Select and update:** Retain only candidates with positive reward; the highest-rewarding sequence becomes the adopted

Table 3: Datasets used in our evaluation.

| Data | #Tbl | #Total Rows | Size |
|---|---|---|---|
| Auto-join web [34] | 62 | 9,220 | 1.07 MB |
| NYC Open [9] | 1,614 | 49.8M | 8.18 GB |
| Chicago Open [30] | 802 | 2,.9M | 0.868 GB |
| NYC+Chicago [9, 30] | 2,416 | 52.8M | 9.05 GB |

transformation for $(c'_a, c'_b)$. Append it to the library for future reuse.

Once reuse is complete, the library maintains a set of transferable transformations annotated by cluster metadata and empirical rewards. Subsequent runs over new datasets immediately benefit from this accumulated experience: they can skip pre-scoring and filtering stages for familiar clusters and proceed directly to RL refinement. Over time, the library evolves into a rich repository of canonical transformation templates, yielding logarithmic amortization of join discovery cost across repository growth.

## 6 Experiments

We evaluate QJoin along three questions:

**RQ1 (Effectiveness vs. baselines).** Does QJoin improve join quality (F1) compared to AutoJoin and LLM-based methods?
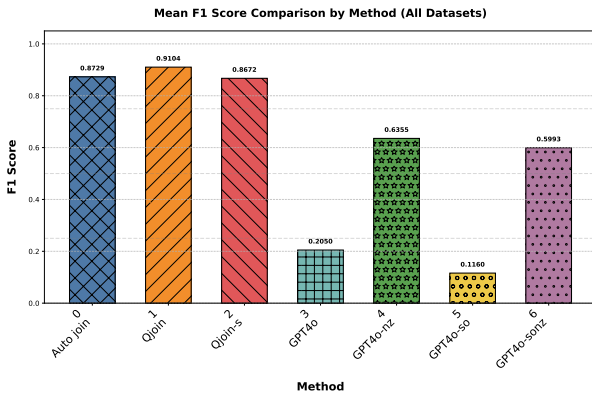
**RQ2 (Efficiency via reuse).** Does reusing transformations yield meaningful runtime savings without degrading F1?

**RQ3 (Scalability).** How does QJoin behave at repository scale, and how much do reuse strategies reduce learning-time overhead as the number of join tasks grows?

**Experimental settings.** All experiments run on a workstation with an Intel i9-13900KF CPU and 64 GB RAM, Python 3.11. For ChatGPT-based baselines we use the Azure-hosted GPT-4o API. We report *Precision*, *Recall*, and *F1* for identified joins; averages are over five runs unless noted.

**Datasets.** We consider the following datasets for evaluation. ● **Auto-join web benchmark** [34]: Contains 31 pairs of tables.62 tables in total. For each pair, it provides a source table and a target table.

● **NYC open datasets** [9]: A diverse collection of CSV tables drawn from New York City's Open Data portal, covering topics such as taxi trips, 311 service requests, and public school performance. Containing 1614 tables in total.

**Figure 2: Web benchmark average F1 score comparison**

● **Chicago datasets** [30]: A set of CSV tables extracted from the City of Chicago's Open Data portal—including crime incidents, building permits, and 311 service requests—preprocessed and stored in the Pneuma 'data_src/tables' directory. Containing 802 tables in total. These datasets are widely used in the literature on transformation-aware joins and open-data integration [9, 30, 34]. We compare the following methods.

- **AutoJoin** [34]: original implementation.
- **QJoin** (reward wrapper): selects the join using the composite reward (Sec. 5.5); discount $\gamma = 1$, learning rate $\alpha = 0.1$.
- **QJoin-S** (similarity wrapper): selects the join sequence by highest ALCS; $\gamma = 1$, $\alpha = 0.1$.
- **GPT-4o**: LLM prompted to (i) produce Python code for the join and (ii) output the merged Pandas DataFrame.
- **GPT-4o-SO**: GPT-4o constrained to the same operator set as AutoJoin (e.g., `concatenate_front/back`, `auto_split_by_operator`, `substring_operator`, ...).

For completeness we also report **GPT-4o-NZ** and **GPT-4o-SO-NZ** which exclude runs with F1 = 0 (to separate execution failures from modeling quality).

### 6.1 AutoJoin Web Benchmark: Results

Figures 2 show that QJoin (reward) achieves the best average F1 ($\approx 91.0\%$), surpassing AutoJoin (87.3%) and QJoin-S (86.7%). GPT-4o methods underperform markedly: GPT-4o averages 20.5% and GPT-4o-SO 11.6%. Even after excluding F1 = 0 cases, GPT-4o reaches only 63.6% and GPT-4o-SO 59.9%. The primary failure modes are: non-executable code (8 datasets for GPT-4o; 15 for GPT-4o-SO), incorrect join-column selection, and weak transformation strategies.
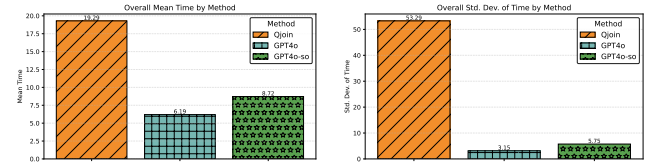
**Figure 3: Web benchmark time**

Figure 4 shows the variation in F1 score for different datasets. The benchmark consisted of 62 tables but due to space restrictions, we show the variations for a random sample of 10 datasets. We observe that QJoin achieves higher or comparable F1 score to AutoJoin across all datasets. The only dataset where AutoJoin performs much better is NY govermnent dataset, where there's no concatenation and uniqueness need, and the vice president column overlaps highly with the president column. Since the vice president column has both high similarity and uniqueness with the president column, Qjoin sometimes gets trapped by trying to match vice presidents to presidents.

**Efficiency.** As shown in Fig. 3, QJoin averages 19.3 s per pair across 31 pairs; GPT-4o averages 6.2 s and GPT-4o-SO 8.7 s. Small heterogeneous tables yield larger variance for QJoin due to differing learned chain depths; nevertheless, QJoin's higher F1 at competitive runtime demonstrates favorable quality–cost trade-offs.
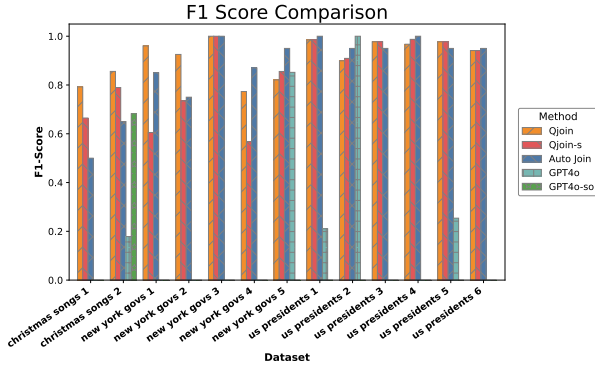
Figure 4: Web benchmark F1 score comparison

**AutoJoin Benchmark: Summary**

QJoin (reward wrapper) attains **91.0% F1**, outperforming AutoJoin (87.3%), QJoin-S (86.7%), GPT-4o (20.5%; 63.6% excl. zeros), and GPT-4o-SO (11.6%; 59.9% excl. zeros), at an average runtime of 19.3 s across 31 pairs.



Figure 5: Time comparison for auto-join folder.

## 6.2 Benefits of Reuse

We assess transformation reuse on individual folders of related datasets derived from the AutoJoin web benchmark. Results are averaged over five runs. We consider the following variants.

- **Direct Ops**: reuse only the unary transformation operators.
- **Agents**: reuse only trained agent policies.
- **All**: reuse every stored operator.
- **QJoin**: reuse both unary operators and agent policies.
- **All+Agents**: reuse all operators plus policies.

*Time savings (Fig. 5).* In the Christmas Songs folder, all reusing methods complete in about 27 s versus 37 s for the standard pipeline, yielding roughly $\frac{37-27}{37} \approx 27\%$ savings. In the New York Gov datasets, the standard method takes $\approx 27$ s, while reuse variants average $\approx 21$ s—around a 22% reduction in runtime.

*F1-score comparison.* Figure 6 shows that, for the Christmas Songs 1 dataset, all reuse variants outperform the standard Ada Join with the similarity-based selection, and their F1 scores closely match the reward-based standard pipeline. In Christmas Songs 2,

*Direct Ops* and *Qjoin* are the closest to the similarity-wrapper baseline. For the New York Gov datasets On Gov 3 and Gov 5 every reuse scheme matches the standard + similarity-wrapper performance, while on the other datasets *Qjoin* most closely approaches this baseline.

Reusing only agent policies can lead to suboptimal outcomes: the agent will accept any positively rewarded transformation, which accelerates convergence but sometimes sacrifices final accuracy. Conversely, reusing all operators risks poor choices (e.g., wrong column mappings for concatenation) instead of allowing the agent to adapt. The Qjoin scheme best balances convergence speed and result quality by combining the reliability of direct operators with the adaptivity of learned policies.

**Conclusion**

Reusing transformations cuts runtime by about 27% on the Christmas Songs datasets and 22% on New York Gov, while matching Qjoin-non reusing's F1 scores. The *Qjoin* variant achieves the best balance of speed and accuracy, converging quickly without sacrificing join quality.
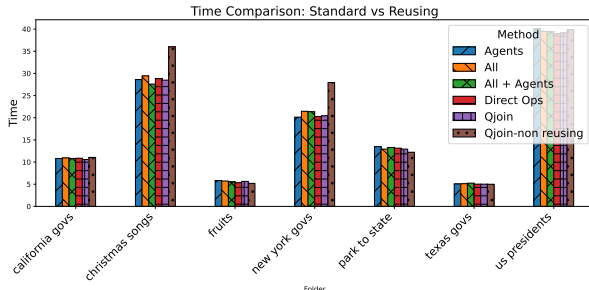
## 6.3 Data Discovery at Repository Scale

Since, large scale data repositories do not have any ground truth, we evaluate the efficiency of QJoin to identify transformation based join and the gain of reusing transformation across varying repository scale.

**Datasets Setup.** First, for each column-pair we estimate the q-gram (q=1–3) Jaccard similarity via LSH, retaining only those with $\hat{J}_q \geq$ 0.6 among non-numeric columns (or when the column names match or are date-related). We then prune any remaining pairs whose full-string MinHash Jaccard also exceeds 0.6 (since no transformation is needed). Next, for each table-pair we keep only the column-pair with the highest $\hat{J}_q$.

**Automatic Folder Construction.** We group join tasks into three folder types: (i) **Same Column Names** (exact name matches), (ii) **Date Column Names** (name contains date/time/year/month), and (iii) **Else Column Names** (remaining tasks clustered by $q$-gram similarity via $k$-means). Within each folder, tables are ordered by frequency across tasks; tasks are then assigned greedily to avoid redundancy.

**Performance metrics** We focus on learning-time savings; join execution time $t_{\text{join}}$ is identical across methods and counted only in total-time savings. We use these time savings into percentage to compare the effectiveness of QJoin.

**Methods.** We consider the standard QJoin along with two variants. **One-Shot**: apply all reused transforms at once; accept if reward increases; **Sequential**: apply reused transforms one by one, accepting only if reward increases (stop at first non-improving step).

*6.3.1 Effect from sample size.* **NYC.** Figure 7 shows total savings vs. sample size. With only 10 joins in the repository, neither one-shot nor sequential reuse provides a measurable benefit: the stored transformation library is nearly empty, so reuse attempts produce zero hits and only add overhead. At 50 joins, early reuse starts to appear in the *date-related folders*, where column names or values
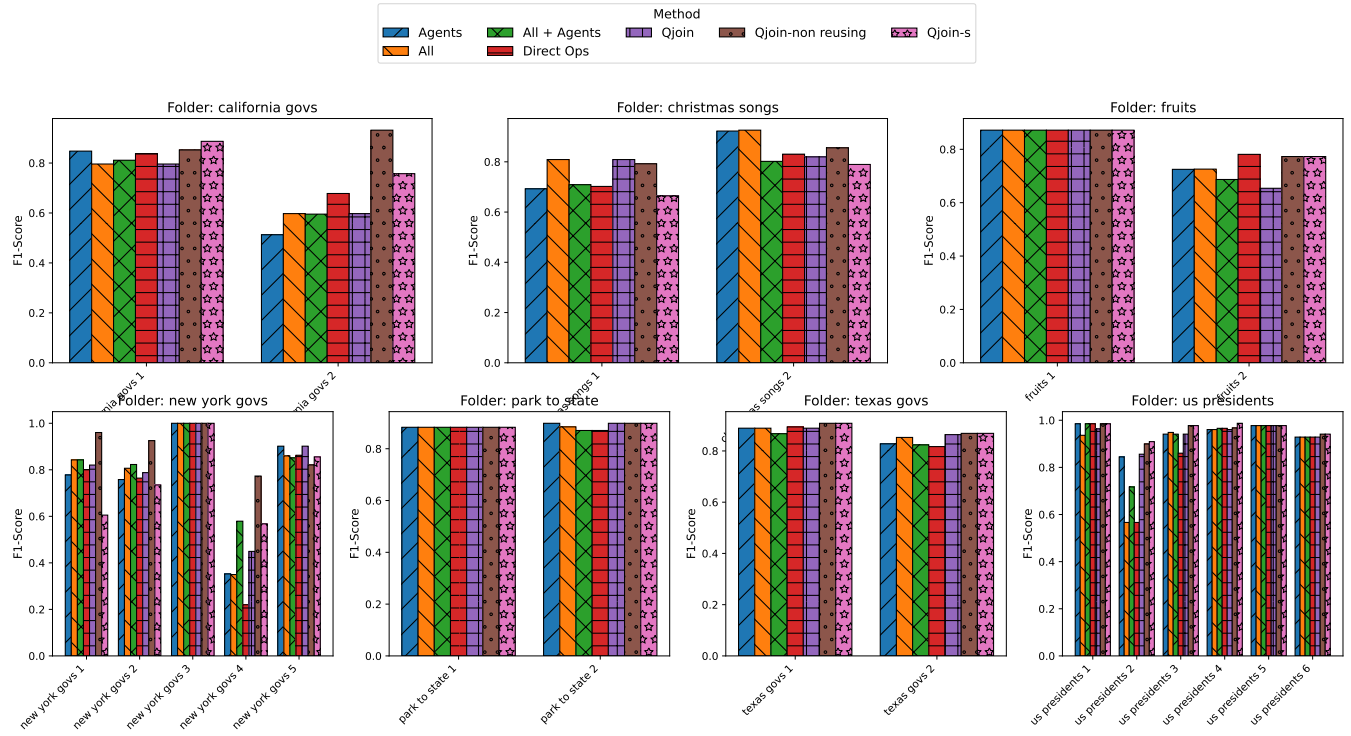
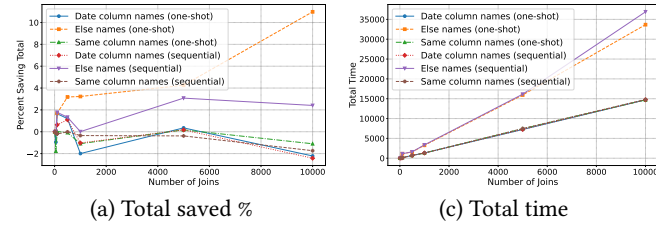**Figure 6: F1 comparison across all folders.**



(a) Total saved %

(c) Total time

**Figure 7: Time saving across all folders for different sample sizes**
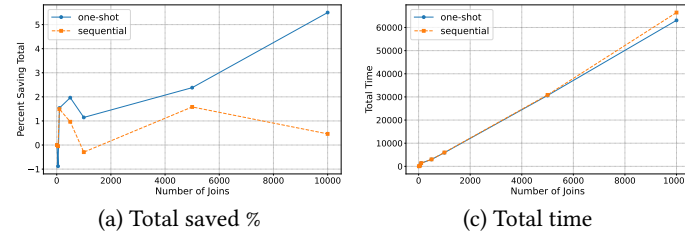


(a) Total saved %

(c) Total time

**Figure 8: Total Time saving for different sample sizes for NYC open datasets**

(e.g., date, year, time) repeat but the overhead of loading and validation still dominates, resulting in a minor slowdown (+0.88% total time).

At 100–500 joins, reuse begins to pay off consistently. The date-related folders show 1.65% and 1.20% time savings, while the remaining non-date folders—containing diverse text attributes like names, addresses, and identifiers—achieve stronger gains of 1.69%

and 3.18%. Overall savings rise to 1.54% (21 s) and 1.97% (59 s). This behavior reflects a growing *reuse density*: as more joins are processed, the likelihood that a new task overlaps with an existing transformation chain increases sharply.

At 1,000–5,000 joins, date-related folders begin to plateau since their simple normalization chains are quickly learned and reused. In contrast, non-date folders continue to show steady improvements—3.23% and 4.26%—bringing total savings to 1.15% and 2.38% (68 s). By 10,000 joins, accumulated reuse yields substantial efficiency: non-date folders achieve 11% (≈ 4,151 s) savings, driving an overall reduction of 5.50% (≈ 3,672 s). These results confirm that reuse benefits grow superlinearly with repository scale once transformation diversity and repetition are sufficient.

Sequential reuse remains nearly cost-free: at small scales, it introduces negligible overhead (0% at 10 joins; 0.05% at 50), then gradually improves (1.49% at 100 joins, 0.97% at 500). A transient dip occurs at 1,000 (−0.30%) due to marginal hit rates, followed by recovery to 1.58% (≈496 s) at 5,000 and 0.46% (≈309 s) at 10,000. In essence, sequential reuse trades peak performance for robustness—minimizing risk when few reusable chains exist.

*Chicago (Fig. 10).* Chicago's repository exhibits fewer recurring patterns and shorter textual values, so reuse takes longer to manifest. With 10 joins, savings are negligible. At 50 joins, both approaches incur small overheads (one-shot −0.24%; sequential −0.13%). At 100 joins, one-shot begins to help (+2.50%), while sequential remains slightly negative (−0.55%). By 500, one-shot stabilizes (+0.20%) and sequential drops (−1.66%). Both peak near 1,000 joins (one-shot +1.75%; sequential +1.21%) before declining again (5,000:
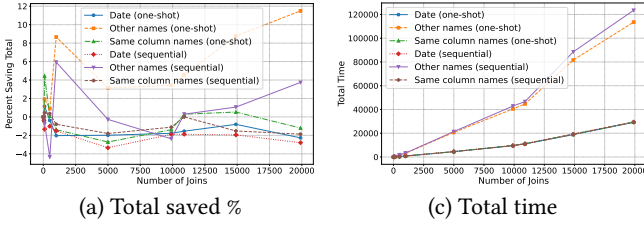
(a) Total saved %   (c) Total time

**Figure 9: Time saving across all folders for different sample sizes for Chicago datasets**



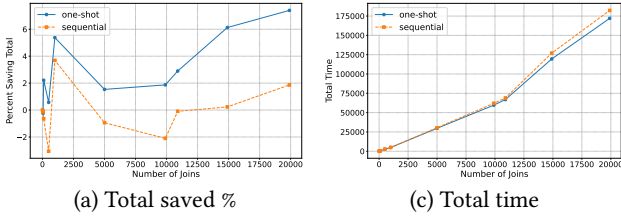(a) Total saved %   (c) Total time

**Figure 10: Total Time saving for different sample sizes for NYC and Chicago merged datasets**

$-0.52\%/-1.80\%$; 9,900: $+0.09\%/-1.79\%$). Most of these improvements come from non-date folders—where irregular text fields benefit from string splitting, concatenation, and normalization—while date-related folders tend to show neutral or slight slowdowns.

**NYC+Chicago combined.** When we merge the repositories, reuse opportunities multiply. Starting from 9,900 Chicago-only joins, adding cross-repository pairs of size $n_2 \in \{1,000, 5,000, 10,000\}$ produces total sizes $n_{total} \in \{10,900, 14,900, 19,900\}$. One-shot reuse achieves $+1.10\%$, $+2.83\%$, and $+2.67\%$ savings respectively, while sequential reuse remains slightly negative ($-0.53\%$, $-0.80\%$, $-0.31\%$). These gains stem primarily from non-date folders, which now include cross-city attributes (e.g., building IDs, agency codes, zip codes) that share syntactic regularities across repositories. Peak folder-level improvements reach $+11.49\%$ at 19,900 joins, whereas date-related folders remain near break-even due to already saturated patterns.

> (i) Cross-repository reuse *amplifies* value by enlarging the pool of reusable chains and the likelihood of immediate matches. (ii) Larger, mixed workloads stabilize performance. Combining NYC with Chicago dataset raises one-shot savings to **2.67%–2.83%** at 15–20k joins, with folder-level peaks of **11.49%** for certain clusters.

## 7 Related Work

**Transformation-Based Data Integration and Joins.** Program synthesis techniques have revolutionized data transformation and integration. FlashFill [12] automates string processing in spreadsheets using input-output examples, with extensions like FlashExtract [16] providing frameworks for extracting data from semistructured documents. FlashRelate [3] addresses extracting relational data from spreadsheets using novel extraction languages. Interactive approaches such as Wrangler [15] combine direct manipulation of visualized data with automatic inference of transforms, while DataXFormer [1] provides robust transformation discovery through systematic exploration of the transformation space. Most relevant to our work, *Auto-Join* [34] automatically discovers transformations to enable joins between disparate tables, though it treats each join problem independently.

**Join Discovery in Data Lakes.** Modern data repositories present unique challenges for join discovery and data exploration. Data lake systems like those at Google [13] organize datasets through comprehensive metadata management. JOSIE [33] formulates joinable table discovery as overlap set similarity search, minimizing the cost of set reads and inverted index probes. The LSH Ensemble [35] addresses Internet-scale domain search using Jaccard set containment, particularly suitable for Open Data. Recent advances include Nexus [10], which enables correlation discovery over spatiotemporal tabular data, Metam [9] for goal-oriented data discovery, and industrial systems like DataHub [17] offering generalized metadata search tools. Table union search [23] and data lake organization [22] use probabilistic models to enhance data discovery and navigation. Schema-level signals such as data profiling and inclusion dependency (IND) discovery are practical precursors for equi-join feasibility and quality control [8, 28, 29].

**Traditional Join Algorithms and Similarity Search.** The foundation of join processing has been extensively studied [11], with techniques evolving from exact to approximate string matching [24]. Set similarity joins have emerged as fundamental for data cleaning [6], with efficient algorithms [31, 32] achieving significant speedups through prefix filtering. Scalability has been addressed through approaches like scaling up all pairs similarity search [4], while Arasu et al. [2] proposed exact algorithms with precise performance guarantees.

**Machine Learning for Schema Matching.** Machine learning has transformed schema matching since LSD [7] introduced multilearner approaches for semantic mappings. Cupid [18] combined multiple matching techniques, while recent deep learning approaches [21] explore various architectures for entity matching. Sherlock [14] uses multi-input deep neural networks for semantic data type detection. In query optimization, reinforcement learning has shown promise with Neo [19] learning from feedback and ReJOIN [20] applying deep RL to join order enumeration.

Our work builds upon these foundations by introducing a reinforcement learning approach that learns and reuses transformation strategies across multiple join tasks. Unlike previous methods that treat each join problem independently, Qjoin maintains a memory of successful transformations and applies them intelligently to new joining scenarios, significantly reducing the computational overhead of join discovery in large-scale data repositories.

## References

[1] Ziawasch Abedjan, John Morcos, Ihab F. Ilyas, Mourad Ouzzani, Paolo Papotti, and Michael Stonebraker. 2016. DataXFormer: A robust transformation discovery

system. In *2016 IEEE 32nd International Conference on Data Engineering (ICDE)*. 1134–1145. doi:10.1109/ICDE.2016.7498319

[2] Arvind Arasu, Venkatesh Ganti, and Raghav Kaushik. 2006. Efficient exact set-similarity joins. In *Proceedings of the 32nd International Conference on Very Large Data Bases* (Seoul, Korea) *(VLDB '06)*. VLDB Endowment, 918–929.

[3] Daniel W. Barowy, Sumit Gulwani, Ted Hart, and Benjamin Zorn. 2015. FlashRelate: extracting relational data from semi-structured spreadsheets using examples. In *Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation* (Portland, OR, USA) *(PLDI '15)*. Association for Computing Machinery, New York, NY, USA, 218–228. doi:10.1145/2737924.2737952

[4] Roberto J. Bayardo, Yiming Ma, and Ramakrishnan Srikant. 2007. Scaling up all pairs similarity search. In *Proceedings of the 16th International Conference on World Wide Web* (Banff, Alberta, Canada) *(WWW '07)*. Association for Computing Machinery, New York, NY, USA, 131–140. doi:10.1145/1242572.1242591

[5] Raul Castro Fernandez, Ziawasch Abedjan, Famien Koko, Gina Yuan, Samuel Madden, and Michael Stonebraker. 2018. Aurum: A Data Discovery System. In *2018 IEEE 34th International Conference on Data Engineering (ICDE)*. 1001–1012. doi:10.1109/ICDE.2018.00094

[6] S. Chaudhuri, V. Ganti, and R. Kaushik. 2006. A Primitive Operator for Similarity Joins in Data Cleaning. In *22nd International Conference on Data Engineering (ICDE'06)*. 5–5. doi:10.1109/ICDE.2006.9

[7] AnHai Doan, Pedro Domingos, and Alon Y. Halevy. 2001. Reconciling schemas of disparate data sources: a machine-learning approach. In *Proceedings of the 2001 ACM SIGMOD International Conference on Management of Data* (Santa Barbara, California, USA) *(SIGMOD '01)*. Association for Computing Machinery, New York, NY, USA, 509–520. doi:10.1145/375663.375731

[8] Falco Dürsch, Axel Stebner, Fabian Windheuser, Maxi Fischer, Tim Friedrich, Nils Strelow, Tobias Bleifuß, Hazar Harmouch, Lan Jiang, Thorsten Papenbrock, and Felix Naumann. 2019. Inclusion Dependency Discovery: An Experimental Evaluation of Thirteen Algorithms. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management* (Beijing, China) *(CIKM '19)*. Association for Computing Machinery, New York, NY, USA, 219–228. doi:10.1145/3357384.3357916

[9] Sainyam Galhotra, Yue Gong, and Raul Castro Fernandez. 2023. Metam: Goal-Oriented Data Discovery. In *39th IEEE International Conference on Data Engineering, ICDE 2023, Anaheim, CA, USA, April 3-7, 2023*. IEEE, 2780–2793. doi:10.1109/ICDE55515.2023.00213

[10] Yue Gong, Sainyam Galhotra, and Raul Castro Fernandez. 2024. Nexus: Correlation Discovery over Collections of Spatio-Temporal Tabular Data. *Proc. ACM Manag. Data* 2, 3 (2024), 154. doi:10.1145/3654957

[11] Goetz Graefe. 1993. Query evaluation techniques for large databases. *ACM Comput. Surv.* 25, 2 (June 1993), 73–169. doi:10.1145/152610.152611

[12] Sumit Gulwani. 2011. Automating string processing in spreadsheets using input-output examples. In *Proceedings of the 38th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (Austin, Texas, USA) *(POPL '11)*. Association for Computing Machinery, New York, NY, USA, 317–330. doi:10.1145/1926385.1926423

[13] Alon Halevy, Flip Korn, Natalya F Noy, Christopher Olston, Neoklis Polyzotis, Sudip Roy, and Steven Euijong Whang. 2016. Goods: Organizing Google's Datasets. In *Proceedings of the 2016 International Conference on Management of Data (SIGMOD)*. ACM, 795–806. https://static.googleusercontent.com/media/research.google.com/en//pubs/archive/45390.pdf

[14] Madelon Hulsebos, Kevin Hu, Michiel Bakker, Emanuel Zgraggen, Arvind Satyanarayan, Tim Kraska, Çağatay Demiralp, and César Hidalgo. 2019. Sherlock: A Deep Learning Approach to Semantic Data Type Detection. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (Anchorage, AK, USA) *(KDD '19)*. Association for Computing Machinery, New York, NY, USA, 1500–1508. doi:10.1145/3292500.3330993

[15] Sean Kandel, Andreas Paepcke, Joseph Hellerstein, and Jeffrey Heer. 2011. Wrangler: interactive visual specification of data transformation scripts. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (Vancouver, BC, Canada) *(CHI '11)*. Association for Computing Machinery, New York, NY, USA, 3363–3372. doi:10.1145/1978942.1979444

[16] Vu Le and Sumit Gulwani. 2014. FlashExtract: a framework for data extraction by examples. *SIGPLAN Not.* 49, 6 (June 2014), 542–553. doi:10.1145/2666356.2594333

[17] LinkedIn Corporation. 2020. DataHub: A Generalized Metadata Search & Discovery Tool. https://github.com/linkedin/datahub

[18] Jayant Madhavan, Philip A. Bernstein, and Erhard Rahm. 2001. Generic Schema Matching with Cupid. In *Proceedings of the 27th International Conference on Very Large Data Bases (VLDB '01)*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 49–58.

[19] Ryan Marcus, Parimarjan Negi, Hongzi Mao, Chi Zhang, Mohammad Alizadeh, Tim Kraska, Olga Papaemmanouil, and Nesime Tatbul. 2019. Neo: A Learned Query Optimizer. In *Proceedings of the VLDB Endowment*, Vol. 12. VLDB Endowment, 1705–1718. http://www.vldb.org/pvldb/vol12/p1705-marcus.pdf

[20] Ryan Marcus and Olga Papaemmanouil. 2018. Deep Reinforcement Learning for Join Order Enumeration. In *Proceedings of the First International Workshop on Exploiting Artificial Intelligence Techniques for Data Management* (Houston, TX, USA) *(aiDM'18)*. Association for Computing Machinery, New York, NY, USA, Article 3, 4 pages. doi:10.1145/3211954.3211957

[21] Sidharth Mudgal, Han Li, Theodoros Rekatsinas, AnHai Doan, Youngchoon Park, Ganesh Krishnan, Rohit Deep, Esteban Arcaute, and Vijay Raghavendra. 2018. Deep Learning for Entity Matching: A Design Space Exploration. In *Proceedings of the 2018 International Conference on Management of Data* (Houston, TX, USA) *(SIGMOD '18)*. Association for Computing Machinery, New York, NY, USA, 19–34. doi:10.1145/3183713.3196926

[22] Fatemeh Nargesian, Ken Q. Pu, Erkang Zhu, Bahar Ghadiri Bashardoost, and Renée J. Miller. 2020. Organizing Data Lakes for Navigation. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data* (Portland, OR, USA) *(SIGMOD '20)*. Association for Computing Machinery, New York, NY, USA, 1939–1950. doi:10.1145/3318464.3380605

[23] Fatemeh Nargesian, Erkang Zhu, Ken Q. Pu, and Renée J. Miller. 2018. Table union search on open data. *Proc. VLDB Endow.* 11, 7 (March 2018), 813–825. doi:10.14778/3192965.3192973

[24] Gonzalo Navarro. 2001. A guided tour to approximate string matching. *ACM Comput. Surv.* 33, 1 (March 2001), 31–88. doi:10.1145/375360.375365

[25] New York City Campaign Finance Board. 2025. Campaign Expenditures. https://data.cityofnewyork.us/City-Government/Campaign-Expenditures/qxzj-vkn2. NYC Open Data.

[26] New York City Campaign Finance Board. 2025. Campaign Public Funds Payments. https://data.cityofnewyork.us/City-Government/Campaign-Public-Funds-Payments/u69g-mvrb. NYC Open Data.

[27] George Papadakis, Dimitrios Skoutas, Emmanouil Thanos, and Themis Palpanas. 2020. Blocking and Filtering Techniques for Entity Resolution: A Survey. *ACM Comput. Surv.* 53, 2, Article 31 (March 2020), 42 pages. doi:10.1145/3377455

[28] Thorsten Papenbrock, Tanja Bergmann, Moritz Finke, Jakob Zwiener, and Felix Naumann. 2015. Data profiling with metanome. *Proc. VLDB Endow.* 8, 12 (Aug. 2015), 1860–1863. doi:10.14778/2824032.2824086

[29] Thorsten Papenbrock, Sebastian Kruse, Jorge-Arnulfo Quiané-Ruiz, and Felix Naumann. 2015. Divide & conquer-based inclusion dependency discovery. *Proc. VLDB Endow.* 8, 7 (Feb. 2015), 774–785. doi:10.14778/2752939.2752946

[30] TheDataStation. 2025. LLM-Powered Data Discovery System for Tabular Data. GitHub repository. https://github.com/TheDataStation/pneuma/tree/main/data_src/tables

[31] Chuan Xiao, Wei Wang, Xuemin Lin, and Jeffrey Xu Yu. 2008. Efficient similarity joins for near duplicate detection. In *Proceedings of the 17th International Conference on World Wide Web* (Beijing, China) *(WWW '08)*. Association for Computing Machinery, New York, NY, USA, 131–140. doi:10.1145/1367497.1367516

[32] Chuan Xiao, Wei Wang, Xuemin Lin, Jeffrey Xu Yu, and Guoren Wang. 2011. Efficient similarity joins for near-duplicate detection. *ACM Trans. Database Syst.* 36, 3, Article 15 (Aug. 2011), 41 pages. doi:10.1145/2000824.2000825

[33] Erkang Zhu, Dong Deng, Fatemeh Nargesian, and Renée J. Miller. 2019. JOSIE: Overlap Set Similarity Search for Finding Joinable Tables in Data Lakes. In *Proceedings of the 2019 International Conference on Management of Data* (Amsterdam, Netherlands) *(SIGMOD '19)*. Association for Computing Machinery, New York, NY, USA, 847–864. doi:10.1145/3299869.3300065

[34] Erkang Zhu, Yeye He, and Surajit Chaudhuri. 2017. Auto-join: joining tables by leveraging transformations. *Proc. VLDB Endow.* 10, 10 (June 2017), 1034–1045. doi:10.14778/3115404.3115409

[35] Erkang Zhu, Fatemeh Nargesian, Ken Q. Pu, and Renée J. Miller. 2016. LSH ensemble: internet-scale domain search. *Proc. VLDB Endow.* 9, 12 (Aug. 2016), 1185–1196. doi:10.14778/2994509.2994534

# A  Appendix

## A.1  ALCS vs other metrics

### A.1.1  Comparison with Jaccard Similarity.
We first contrast ALCS with the widely used Jaccard similarity on $q$-grams, which measures token overlap within fixed-length windows.

*Jaccard on $q$-grams.* For two strings $S_1, S_2$, the Jaccard similarity is defined as:

$$\text{Jaccard}_q(S_1, S_2) = \frac{|\text{QGrams}_q(S_1) \cap \text{QGrams}_q(S_2)|}{|\text{QGrams}_q(S_1) \cup \text{QGrams}_q(S_2)|},$$

where $\text{QGrams}_q(S)$ is the multiset (or set) of all contiguous substrings of $S$ of length $q$.

**Sensitivity to a Single $q$.** Jaccard's dependence on fixed-length windows makes it brittle under small shifts or boundary misalignment.

LEMMA 1 (SENSITIVITY OF JACCARD TO MISALIGNMENT). *For any fixed $q$, there exist strings $S_1, S_2$ that share a long contiguous block but exhibit a small intersection of $q$-grams due to slight positional shifts. In contrast, $\text{ALCS}(S_1, S_2)$ captures nearly the entire block.*

PROOF 1. *Let $S_1 = uXv$ and $S_2 = u'Xv'$, where $X$ is a common block of length $k \gg q$. Assume $X$ starts at position $p_1$ in $S_1$ and $p_2$ in $S_2$ with $|p_1 - p_2| = 1$. Then the $q$-grams covering $X$ in $S_1$ run from indices $p_1$ to $p_1 + k - q$, while in $S_2$ they run from $p_2$ to $p_2 + k - q$. A one-character shift prevents most $q$-grams from aligning exactly, so $|\text{QGrams}_q(S_1) \cap \text{QGrams}_q(S_2)|$ is small even though $X$ is nearly identical in both strings. By contrast, ALCS identifies $X$ as the longest shared substring of length $k$, producing a high similarity value $\approx k / (\frac{1}{2}(|S_1| + |S_2|))$.*

Thus, ALCS remains robust to local shifts or token misalignments, while $\text{Jaccard}_q$ can underestimate similarity for otherwise well-aligned blocks.

*Effect of Transformations.* Under transformations $\Omega_a$ and $\Omega_b$, columns $col_a$ and $col_b$ may be reordered or structurally modified. When two disjoint common regions become a *single* contiguous overlap, ALCS strictly increases since the merged block is longer. By contrast, $\text{Jaccard}_q$ improves only if the merge creates perfectly aligned $q$-grams—an unlikely condition when transformations shift or reorder text boundaries.

### A.1.2  Comparison with Cosine Similarity.
Cosine similarity, when applied to bag-of-tokens or embedding vectors, ignores token order entirely.

*Cosine Similarity on Token Vectors.* Let $\vec{v}(S)$ denote a token-count or embedding vector. Then:

$$\text{Cosine}(S_1, S_2) = \frac{\vec{v}(S_1) \cdot \vec{v}(S_2)}{\|\vec{v}(S_1)\| \, \|\vec{v}(S_2)\|}.$$

*Ignoring Contiguity and Order.* Because this representation discards sequential structure, it cannot distinguish between identical token sets appearing in different orders.

LEMMA 2 (COSINE FAILS TO CAPTURE CONTINUITY). *If $S_1$ and $S_2$ contain the same multiset of tokens in different orders, then $\text{Cosine}(S_1, S_2) = 1$, while $\text{ALCS}(S_1, S_2)$ can be near $0$ if no contiguous substring longer than the threshold $n$ is shared.*

PROOF 2. *Let $S_1$ be a permutation of tokens $\{t_1, t_2, \ldots, t_m\}$ and $S_2$ another permutation of the same multiset. Since $\vec{v}(S_1) = \vec{v}(S_2)$, their dot product equals the product of their norms, yielding $\text{Cosine}(S_1, S_2) = 1$. However, if the orderings are disjoint, the longest shared contiguous substring may have length $\leq 1$, so $\text{ALCS}(S_1, S_2)$ is small.*

Thus, ALCS distinguishes between re-ordered and truly aligned text segments, whereas Cosine similarity cannot.

*Impact of Transformations.* When transformations $\Omega_a$ and $\Omega_b$ reorder tokens to improve contiguity, Cosine often remains unchanged (token counts are constant). ALCS, however, increases proportionally to the newly formed contiguous overlap.

### A.1.3  Comparison with Edit Distance.
Edit distance measures character-level transformations but penalizes large block moves heavily.

*Edit (Levenshtein) Distance.*
$$\text{ED}(S_1, S_2) = \min\{\#\text{ins, del, or subs } S_1 \rightarrow S_2\},$$

$$\text{Sim}_{\text{edit}}(S_1, S_2) = 1 - \frac{\text{ED}(S_1, S_2)}{\max(|S_1|, |S_2|)}.$$

*Block Reordering Cost.* Standard edit distance does not include a "block move" operation, so reordering contiguous substrings is expensive.

LEMMA 3 (COST OF BLOCK SWAPPING IN EDIT DISTANCE). *If $S_1$ can be transformed to $S_2$ only by swapping two blocks of length $m$, then $\text{ED}(S_1, S_2) \geq m$ unless block-move operations are given negligible cost.*

PROOF 3. *Under the standard model (insertions, deletions, substitutions), swapping two blocks of length $m$ requires deleting $m$ characters and re-inserting them at the new location. Thus the total edit cost grows linearly with $m$, and $\text{ED}(S_1, S_2) \geq m$ unless augmented with a free block-swap operation.*

Therefore, $\text{Sim}_{\text{edit}}$ penalizes even semantically equivalent reorderings. In contrast, ALCS focuses only on the final aligned substrings, ignoring the number of edits required to reach them.

### A.1.4  Monotonicity Property of ALCS.
Finally, ALCS satisfies a key structural property—*monotonicity under merges of disjoint blocks*—which guarantees that transformations improving contiguity always increase similarity.

PROPOSITION 1 (MONOTONICITY OF ALCS UNDER MERGING). *Let $(S_1', S_2')$ have a set of significant common substrings $L'$. Suppose transformations $\Omega_a, \Omega_b$ produce new strings $(S_1'', S_2'')$ in which two disjoint matched blocks $s_1, s_2 \in L'$ become adjacent, forming a longer block $s_{12}$. Then:*
$$\text{ALCS}(S_1'', S_2'') > \text{ALCS}(S_1', S_2').$$

PROOF 4. *Since $s_{12} = s_1 \| s_2$ is contiguous in both $S_1''$ and $S_2''$, its length is $|s_{12}| = |s_1| + |s_2| > \max(|s_1|, |s_2|)$. The longest common substring length therefore strictly increases. Because $\text{ALCS}(S_1, S_2) = \frac{\max_{s \in L} |s|}{\frac{1}{2}(|S_1| + |S_2|)}$ and transformations preserve total string lengths, the denominator is unchanged. Hence $\text{ALCS}(S_1'', S_2'') > \text{ALCS}(S_1', S_2')$.*

Thus, ALCS naturally rewards transformations that merge scattered matches into fewer, longer contiguous blocks. Other similarity measures may fail to improve ($\text{Jaccard}_q$, due to misalignment), remain unchanged (Cosine, due to token order invariance), or even penalize such merges (Edit Distance, due to high reordering cost).

## A.2 Replacements pseudo code

---

**Algorithm 1** Select Best Column Pair by Reward

---

**Require:** List of transformed pairs $\mathcal{P} = [p_1, p_2, \ldots, p_n]$
**Ensure:** Best pair $p_{\text{best}}$
1:   $p_{\text{best}} \leftarrow p_1$
2:   **for** $i \leftarrow 2$ to $|\mathcal{P}|$ **do**
3:      $p_{\text{new}} \leftarrow p_i$
4:      $R_{\text{alcs}} \leftarrow \text{ComputeALCSReward}(p_{\text{best}}, p_{\text{new}})$
5:      $R_{\text{uniq}} \leftarrow \text{ComputeUniquenessReward}(p_{\text{best}}, p_{\text{new}})$
6:      $R_{\text{final}} \leftarrow \lambda_1 R_{\text{alcs}} + \lambda_2 R_{\text{uniq}}$
7:      **if** $R_{\text{final}} > 0$ **then**
8:         $p_{\text{best}} \leftarrow p_{\text{new}}$
9:   **return** $p_{\text{best}}$

---

**Algorithm 2** Find Equivalent Replacements (Single Side)

---

**Require:** Cluster map $C$, learning pair $L$, stored pair $S = (p_a, p_b)$, missing column $x$
**Ensure:** Set of replacements Reps
1:   **if** $x_{\text{table}} = p_a^{\text{table}}$ **then**
2:      $p_{\text{other}} \leftarrow p_b$
3:   **else**
4:      $p_{\text{other}} \leftarrow p_a$
5:   $k \leftarrow \textsc{find\_cluster\_key}(C, (p_{\text{other}}, x))$
6:   Reps $\leftarrow \{ q : (L, q) \in C[k] \} \cup \{ p : (p, L) \in C[k] \}$
7:   **return** Reps

---

*Interpretation.* This reuse mechanism serves two complementary purposes: (i) it dramatically reduces search cost by leveraging past experience, and (ii) it promotes consistency of transformations across datasets sharing related semantics. Conceptually, it functions like a "policy cache" in reinforcement learning—providing prior knowledge that biases exploration toward successful operator patterns.

## A.3 Scalable Join-Discovery Workflow & Optimizations

**Overview.** Given $n$ tables, there are $\binom{n}{2}$ possible column-pairs. We apply a three-stage pipeline to reduce $O(n^2)$ costs:

(1) *Discover via q-gram LSH* (Alg. 3).
(2) *Prune trivial high-sim pairs* using full-string LSH (Alg. 4).
(3) *Select final joins* as a maximum-spanning forest (Alg. 5).

---

**Algorithm 3** Discover Joinable Pairs via Q-gram LSH

---

**Require:** $T = \{t \mapsto D_t\}$, $q$, $p$, $s$, $\theta$, $w$, $e$
**Ensure:** $\mathcal{P}_0 = \{(t_a, c_a, t_b, c_b, \hat{J}_q) \mid \hat{J}_q \geq \theta\}$
1:   Build $\tau = \{(k, \mathbf{v}_k, p)\}$, where $k = t.c$, $\mathbf{v}_k = \text{sample}(D_t[c], s)$
2:   **for all** $(k, \mathbf{v}_k, p) \in \tau$ in parallel $(w)$ **do**
3:      $\ell_k \leftarrow \text{MinHash}(p)$
4:      **for all** $v \in \mathbf{v}_k$ **do**
5:         **for all** $g \in \text{qgrams}(v, q)$ **do**
6:             $\ell_k.\text{update}(g)$
7:   Index $\{(k, \ell_k, |\mathbf{v}_k|)\}$ in LSHEnsemble$(\theta, p, e)$
8:   $\mathcal{P}_0 \leftarrow \emptyset$
9:   **for all** $k$ **do**
10:     **for all** $k' \in \text{ensemble.query}(\ell_k, |\mathbf{v}_k|)$ **do**
11:        **if** $k' \neq k$ **then**
12:           $\hat{J}_q \leftarrow \ell_k.\text{jaccard}(\ell_{k'})$
13:           **if** $\hat{J}_q \geq \theta$ and $\neg\texttt{same-table}(k, k')$ **then**
14:              parse $(t_a, c_a), (t_b, c_b) \leftarrow \text{sort}(k, k')$
15:              add $(t_a, c_a, t_b, c_b, \hat{J}_q)$ to $\mathcal{P}_0$
16:   **return** $\mathcal{P}_0$

---

**Algorithm 4** Prune Trivial Pairs via Full-String LSH

---

**Require:** $\mathcal{P}_0$, $T$, $p$, $w$, $e$, $\theta$
**Ensure:** $\mathcal{R} = \{(t_a, c_a, t_b, c_b) \mid \hat{J} \geq \theta\}$
1:   Extract keys $\{k \mid (k, \ldots) \in \mathcal{P}_0\}$
2:   Build $\tau' = \{(k, D_t[c], p)\}$ on full values
3:   **for all** $(k, \mathbf{v}_k, p) \in \tau'$ in parallel **do**
4:      $\ell'_k \leftarrow \text{MinHash}(p)$ over each $v \in \mathbf{v}_k$
5:   Index $\{(k, \ell'_k, |\mathbf{v}_k|)\}$ in LSHEnsemble$(\theta, p, e)$
6:   $\mathcal{R} \leftarrow \emptyset$
7:   **for all** $k$ **do**
8:      **for all** $k' \in \text{ensemble.query}(\ell'_k, |\mathbf{v}_k|)$ **do**
9:        **if** $k' \neq k$ **then**
10:       $\hat{J} \leftarrow \ell'_k.\text{jaccard}(\ell'_{k'})$
11:       **if** $\hat{J} \geq \theta$ and $\neg\texttt{same-table}(k, k')$ **then**
12:         parse $(t_a, c_a), (t_b, c_b) \leftarrow \text{sort}(k, k')$
13:         add $(t_a, c_a, t_b, c_b)$ to $\mathcal{R}$
14:   **return** unique rows of $\mathcal{R}$

---

**Algorithm 5** Final Join Tasks via Maximum-Spanning Forest

---

**Require:** $\mathcal{P}_1 = \mathcal{P}_0 \setminus \mathcal{R}$
**Ensure:** $\mathcal{L}$: MST edges $(t_a, c_a, t_b, c_b, \hat{J}_q)$
1:   Sort $\mathcal{P}_1$ by descending $\hat{J}_q$
2:   Init union-find on tables
3:   $\mathcal{L} \leftarrow []$
4:   **for all** edge $(t_a, c_a, t_b, c_b, \hat{J}_q) \in \mathcal{P}_1$ **do**
5:      **if** $\text{find}(t_a) \neq \text{find}(t_b)$ **then**
6:         append $(t_a, c_a, t_b, c_b, \hat{J}_q)$ to $\mathcal{L}$
7:         union$(t_a, t_b)$
8:   **return** $\mathcal{L}$

---

**Time-Efficiency Optimizations.** Let each cluster $C$ be down–sampled:

$$S_C \subseteq C, \quad |S_C| = \min\{\max\{|C|, 10\}, 20\}.$$

Then:

(i) *Pruned ALCS Search.* Compute q-gram Jaccard $J_q(i, j)$ on $S_C$. Let $\tau = \text{Quantile}_{0.9}\{J_q\}$. Only pairs with $J_q \geq \tau$ incur full ALCS cost.

**Reuse-Specific Optimizations.**

  (a) *Early Termination.* Stop replacement trials once reward $> 0$.

  (b) *Cluster-Sampling.* Apply the same $|S_C|$-rule when clustering for reuse.

## A.4 Optimization of Concatenation Operators

Since Concatenated Operators will lengthen the values and will cause the training time to be longer. We do the optimizations as below.

Consider two sequences of strings defined as:

$$\text{str}_1 = \text{str}_i \ldots \text{str}_m, \tag{1}$$

$$\text{str}_2 = \text{str}_j \ldots \text{str}_n. \tag{2}$$

Our goal is to iteratively concatenate substrings either at the front or back positions of , and substrings similarly to . After each concatenation, we evaluate the resulting strings using the ALCS.

Direct evaluation using naive methods, attempting all possible concatenation positions and recalculating ALCS, is computationally expensive. We propose optimization by excluding redundant computations based on the following observations:

  (1) Concatenating to the back of produces an equivalent result as concatenating to the front of . Thus, concatenating at the front of can be excluded if concatenation at the back of yields a negative reward (we only select actions with the highest positive reward).

  (2) This equivalence holds true provided the overlapping part remains unchanged, primarily since direct operators typically shorten strings by removing irrelevant parts, rarely affecting the overlapping portion. Under this assumption, removing actions do not modify the overlapping segments.

To systematically exploit these observations, we maintain two dictionaries to record potential concatenations that can safely be excluded from consideration. Specifically, whenever a concatenation action yields a negative reward, we update these dictionaries accordingly. Moreover, once a positive-reward concatenation is executed:

  - We reset all excluded concatenation choices for after concatenating at the back of .
  - Similarly, we reset concatenation choices for when concatenation occurs for , and vice versa, given their reversibility.

This structured optimization significantly reduces computational overhead by eliminating redundant ALCS calculations without compromising performance.

### A.4.1 MST joins for data repositories.

*Setup.* Following the same standardized dataset preparation, we built a maximum spanning forest over the candidate join graph of the NYC Open Data repository, which produced 905 final join tasks.

*Distribution of Time Differences.* Figures 11 and 12 display histograms of $\Delta t$ for one-shot and sequential rewards, respectively. In both cases:

  - The distributions are approximately Gaussian, centered around $\Delta t \approx 0$.
  - This reflects (a) relatively short transformation chains, so reuse only marginally reduces exploration, and (b) small absolute times (mean learning time $\approx 3\,\text{s}$).

## A.5 NYC MST joins

*One-Shot Reward.* Figure 11 shows the distribution of learning-time savings under the one-shot reward scheme, while Figure 13 reports both percentage and absolute savings for each dataset folder. In the *Date Names* folder, the 8 successful reuse attempts yield a learning-time reduction of approximately 18% (and an 11% reduction in total time), whereas across all 22 trials the savings are 15% and 9%, respectively—implying a reuse success rate of about 36%. The decline in average savings from successful to all attempts indicates that date-parsing tasks benefit most from direct transformation reuse; when direct reuse fails, the subsequent agent fallback does not recover enough time to offset the overhead of testing stored transformations.

In *Else Column Names* folder above 62%—achieve absolute learning-time reductions of roughly 213 s ($\approx$ 27%) for the 56 successful attempts and 247 s ($\approx$ 26%) when considering all 90 trials. This corresponds to a reuse success rate of nearly 62%, demonstrating that high-similarity column pairs yield consistently reliable direct transformations, and that agents effectively compensate when reuse does not succeed.

In the *Same Column Names* folder, 8 successful reuse attempts actually incur a modest 21% learning-time increase and 14% for all 24 trials, suggesting that trivial name matches rarely correspond to non-trivial transformations.

*Sequential Reward.* Figure 12 shows the distribution of learning-time savings under the sequential reward scheme, while Figure 14 reports both percentage and absolute savings for each dataset folder. Compared with the one-shot results, the sequential scheme achieves far fewer successful direct-transformation reuses—just one in the *Same Column Names* folder, one in the *Date Names* folder, and 6 in the *Other Similarity* folder. Because successful reuses are so rare, we focus here on the *all attempts* metrics.

Across all reuse trials, the *Other* folder achieves roughly a 10 % reduction in both learning time and total execution time. In the *Date Names folder*, we observe about a 10 % savings in learning time and a 7 % savings in total time. This confirms that, even when direct reuse seldom succeeds, the agent's recovery mechanism still expedites the search for effective transformations relative to training from scratch.

In the *Same Column Names* folder, sequential reward delivers about both 1.9% savings in learning time and in total time over all trials. This improvement over one-shot reward arises because the sequential scheme aborts the reuse process immediately upon encountering a transformation step that does not increase reward, thereby avoiding unnecessary checks and reducing overall learning overhead.

*Final contrast.* By including the non-reusable trials, Table 4 and Table 4 highlight the differing strengths of the two reuse schemes. Under One-Shot reward, the *Date Names* folder yields a 1.34 % learning-time saving (down to 2.60 s) and a 1.00 % total-time saving (down to 3.50 s), whereas Sequential reward produces 1.15 % and 0.86 % savings (to 2.60 s and 3.51 s). For the *Other Similarity* folder, One-Shot achieves 12.17 % learning-time and 10.73 % total-time
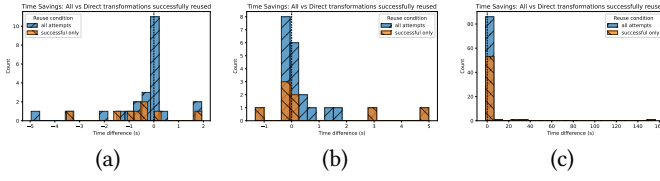
(a)                          (b)                          (c)

**Figure 11: Time-difference distributions for one-shot reward reusing: (a) Same Column names, (b) Date Column names, and (c) Else Column names.**



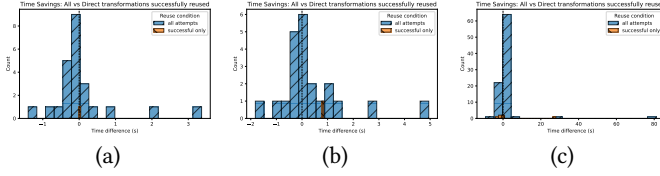(a)                          (b)                          (c)

**Figure 12: Time-difference distributions for sequential reward reusing: (a) Same Column names, (b) Date Column names, and (c) Else Column names.**
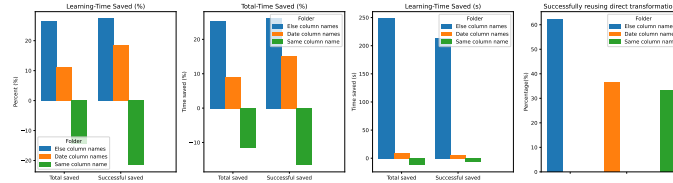


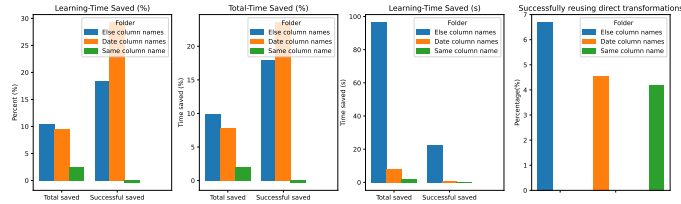**Figure 13: All time info for one-shot reward reusing.**



**Figure 14: All time info for sequential reward reusing.**

savings (to 4.69 s and 5.40 s), compared to 4.73 % and 4.18 % under Sequential (to 5.09 s and 5.80 s). In contrast, *Same Column Names* incurs a 1.75 % increase in learning time (to 2.50 s) and a 1.30 % increase in total time (to 3.37 s) under One-Shot, while Sequential

still delivers modest savings of 0.30 % (to 2.45 s) and 0.22 % (to 3.32 s) by aborting non-rewarding steps early. Aggregated across all three folders, One-Shot reward saves a total of 245.49 s—equivalent to 7.29 % of learning time and 5.98 % of total time—whereas Sequential reward saves 106.21 s (3.15 % learning-time, 2.59 % total-time). Overall, One-Shot reward maximizes benefits when stored transformations closely match new tasks, whereas Sequential reward minimizes wasted checks when reuse potential is low.

> **Conclusion**
>
> The one-shot scheme achieves up to 7.29% learning time savings (5.98% for total saved time, 245.5 s). The sequential scheme delivers up to 3.15% learning savings (2.59% for total saved time, 106.2 s).

**Table 4: Learning-time and total-time savings by folder and reuse scheme**

| Folder | Scheme | % Learning-s | % Total Saved |
|---|---|---|---|
| Date Column Names | One-Shot | 1.34 | 1.00 |
| Same Column Names | One-Shot | -1.75 | -1.30 |
| Else Column Names | One-Shot | 12.17 | 10.73 |
| Date Column Names | Sequential | 1.15 | 0.86 |
| Same Column Names | Sequential | 0.30 | 0.22 |
| Else Column Names | Sequential | 4.73 | 4.18 |

**Table 5: Average Learning-time and Average total-time by folder and reuse scheme**

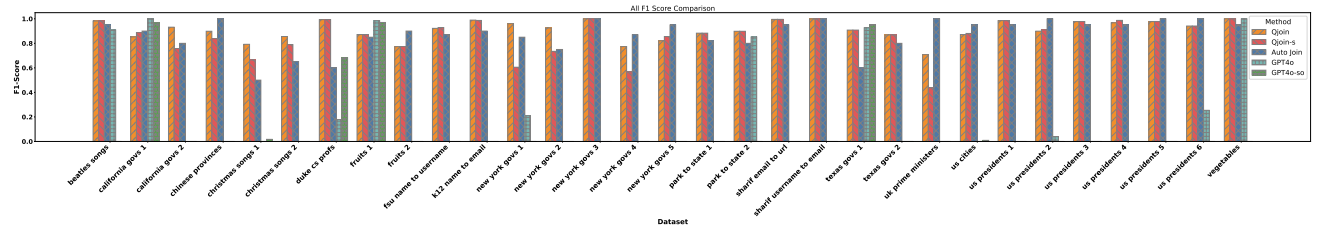| Folder | Scheme | Avg Learning(s) | Avg Total(s) |
|---|---|---|---|
| Date Column Names | One-Shot | 2.60 | 3.50 |
| Same Column Names | One-Shot | 2.50 | 3.37 |
| Else Column Names | One-Shot | 4.69 | 5.40 |
| Date Column Names | Sequential | 2.60 | 3.51 |
| Same Column Names | Sequential | 2.45 | 3.32 |
| Else Column Names | Sequential | 5.09 | 5.80 |

## A.6    Auto-join Benchmark comparison

**Figure 15: Auto-join Benchmark comparison for all datasets.**