

A. Preliminaries: Table Operators

Suppose we have two tables, T_1, T_2 , that share common columns C , and are in their minimal forms in which there are no duplicates and no tuples that can be subsumed or complemented. We show that for each pairwise table operator, Inner Union, Inner Join, Left Join, Outer Join, Cross Product, there exists an equivalent query consisting of Outer Union and/or unary operators. (SP of SPJU queries are accounted for by the unary operators).

Lemma 11 (Inner Union): Inner Union(\cup): it is known that if the schemas of two tables are equal, then Inner Union = Outer Union

Lemma 12 (Inner Join): Inner Join (\bowtie):

$$T_1 \bowtie T_2 = \sigma(T_1.C = T_2.C \neq \perp, \beta(\kappa(T_1 \uplus T_2))) \quad (7)$$

Lemma 13 (Left Join): Left Join (\ltimes) [32]:

$$T_1 \ltimes T_2 = \beta((T_1 \bowtie T_2) \uplus T_1) \quad (8)$$

Lemma 14 (Outer Join): Full Outer Join (\bowtie) [32]:

$$T_1 \bowtie T_2 = \beta(\beta((T_1 \bowtie T_2) \uplus T_1) \uplus T_2) \quad (9)$$

Lemma 15 (Cross Product): Cross Product(\times): We denote columns in T_1 and T_2 as $T_1.C$ and $T_2.C$, respectively. Consider a constant column c .

$$T_1 \times T_2 = \kappa(\pi((T_1.C, c), T_1) \uplus \pi((T_2.C, c), T_2)) \quad (10)$$

Thus, $\uplus, \sigma, \pi, \kappa, \beta$ operators form queries that are equivalent to all SPJU queries.

1) **Proof of Lemma [12] Inner Join:** Given two tables T_1, T_2 that join on a set of common columns C , such that T_1, T_2 are in their minimal forms in which they contain no duplicate tuples and no tuples can be subsumed or complemented, $T_1 \bowtie T_2$ can be expressed by an equivalent query containing Outer Union, complementation, and subsumption. Specifically, $T_1 \bowtie T_2$ is equivalent to query $\sigma(T_1.C = T_2.C \neq \perp, \beta(\kappa(T_1 \uplus T_2)))$.

Proof: We first prove that all tuples in $T_1 \bowtie T_2$ are contained in $\sigma(T_1.C = T_2.C \neq \perp, \beta(\kappa(T_1 \uplus T_2)))$. Let tuple $t \in T_1 \bowtie T_2$, such that join columns C 's values in t appear in both $T_1.C$ and $T_2.C$, and are non-null: $t.C \in T_1.C \cap T_2.C$ s.t. $t.C \neq \perp$.

When applying $\beta(\kappa(T_1 \uplus T_2))$, only tuples with common non-null values $T_1.C_i = T_2.C_i \neq \perp$ in same column(s) i are complemented and subsumed. This is similar to tuple t , which is formed by joining on $T_1.C_i = T_2.C_i$. Thus, tuple t is derived by selecting on tuples from $\beta(\kappa(T_1 \uplus T_2))$ with non-null C values in both $T_1.C$ and $T_2.C$, so $t \in \sigma(T_1.C = T_2.C \neq \perp, \beta(\kappa(T_1 \uplus T_2)))$.

Next, we show that all tuples in $\sigma(T_1.C = T_2.C \neq \perp, \beta(\kappa(T_1 \uplus T_2)))$ are found in $T_1 \bowtie T_2$. Let tuple $t' \in \sigma(T_1.C = T_2.C \neq \perp, \beta(\kappa(T_1 \uplus T_2)))$. Here, all C values in t' are non-null values found in both $T_1.C$ and $T_2.C$ as a result of selection. From $\beta(\kappa(T_1 \uplus T_2))$, t' contains all values from all columns in T_1 and T_2 in a single tuple, formed by

complementing and subsuming based on common C values. Thus, $t' \in T_1 \bowtie T_2$.

We have thus shown that all tuples from $T_1 \bowtie T_2$ are found in $\sigma(T_1.C = T_2.C \neq \perp, \beta(\kappa(T_1 \uplus T_2)))$ and vice versa, and so $T_1 \bowtie T_2$ is an equivalent query to $\sigma(T_1.C = T_2.C \neq \perp, \beta(\kappa(T_1 \uplus T_2)))$. ■

2) **Proof of Lemma [13] Left Join:** Given two tables T_1, T_2 that join on a set of common columns C , such that T_1, T_2 are in their minimal forms in which there are no duplicates and no tuples can be subsumed or complemented, $T_1 \ltimes T_2$ can be expressed by an equivalent query containing Outer Union and subsumption. Specifically, $T_1 \ltimes T_2$ is equivalent to query $\beta((T_1 \bowtie T_2) \uplus T_1)$.

Proof: We first prove that the resulting table of $T_1 \ltimes T_2$ is contained in the resulting table of $\beta((T_1 \bowtie T_2) \uplus T_1)$:

Let tuple $t \in T_1 \ltimes T_2$. There are two cases for join column C 's values in tuple t : $t.C \in T_1.C \cap T_2.C$ (i.e., $t.C$ values are in both $T_1.C$ and in $T_2.C$) and $t.C \in T_1.C \setminus T_2.C$ (i.e., $t.C$ values are only in $T_1.C$ and not in $T_2.C$). Since we are performing left join on T_1 and T_2 , $t.C \notin T_2.C \setminus T_1.C$.

- 1) $t.C \in T_1.C \cap T_2.C \implies t \in (T_1 \bowtie T_2)$. Since t is in the inner join result and contains more non-Null values than other tuples with C values only in T_1 or T_2 , it would not be subsumed when applying $\beta((T_1 \bowtie T_2) \uplus T_1)$.
- 2) $t.C \in (T_1.C \setminus T_2.C) \implies t \in \beta((T_1 \bowtie T_2) \uplus T_1)$. Since T_1 is in its minimal form, and t does not share any C values with any tuple in T_2 , it is not subsumed when applying β to $(T_1 \bowtie T_2) \uplus T_2$, and thus appear as is in $\beta((T_1 \bowtie T_2) \uplus T_1)$.

Thus, all tuples from $T_1 \ltimes T_2$ are contained in the resulting table of $\beta((T_1 \bowtie T_2) \uplus T_1)$.

Next, we show that the resulting tuples of $\beta((T_1 \bowtie T_2) \uplus T_1)$ are contained in the resulting table of $T_1 \ltimes T_2$.

Let's consider tuple $t' \in \beta((T_1 \bowtie T_2) \uplus T_1)$. There are two cases for C values in tuple t' : $t'.C \in T_1.C \cap T_2.C$ and $t'.C \notin T_1.C \cap T_2.C$.

- 1) $t'.C \in (T_1.C \cap T_2.C) \implies t' \in (T_1 \bowtie T_2)$. Since $(T_1 \bowtie T_2) \subseteq (T_1 \ltimes T_2)$, $t' \in (T_1 \ltimes T_2)$.
- 2) All tuples in $((T_1 \bowtie T_2) \uplus T_1)$ are either subsumed by tuples from $(T_1 \bowtie T_2)$, or are in $T_1 \setminus (T_1 \bowtie T_2)$. Thus, $t'.C \notin T_1.C \cap T_2.C \implies t' \in T_1 \setminus (T_1 \bowtie T_2) \implies t' \in T_1 \ltimes T_2$.

Thus, all tuples from $\beta((T_1 \bowtie T_2) \uplus T_1)$ are contained in the resulting table of $T_1 \ltimes T_2$.

Now that we have shown that tuples from $T_1 \ltimes T_2$ are contained in the resulting table of $\beta((T_1 \bowtie T_2) \uplus T_1)$ and vice versa, we have shown that $\beta((T_1 \bowtie T_2) \uplus T_1)$ is an equivalent query to $T_1 \ltimes T_2$. ■

3) **Proof of Lemma [14] Outer Join:** Given two tables T_1, T_2 that join on a set of common columns C , such that T_1, T_2 are in their minimal forms in which there are no duplicates and no tuples can be subsumed or complemented, $T_1 \bowtie T_2$ can be expressed by an equivalent query containing Outer Union and subsumption. Specifically, $T_1 \bowtie T_2$ is equivalent to query $\beta(\beta((T_1 \bowtie T_2) \uplus T_1) \uplus T_2)$.

Proof: We first prove that the resulting table of $T_1 \bowtie T_2$ is contained in the resulting table of $\beta(\beta((T_1 \bowtie T_2) \uplus T_1) \uplus T_2)$:

Let tuple $t \in T_1 \bowtie T_2$. There are three cases for join column C 's values in tuple t : $t.C \in T_1.C \cap T_2.C$ (i.e., $t.C$ values are in both $T_1.C$ and in $T_2.C$), $t.C \in T_1.C \setminus T_2.C$ (i.e., $t.C$ values are only in $T_1.C$ and not in $T_2.C$), and $t.C \in T_2.C \setminus T_1.C$ (i.e., $t.C$ values are only in $T_2.C$ and not in $T_1.C$).

- 1) $t.C \in T_1.C \cap T_2.C \implies t \in T_1 \bowtie T_2$. Tuple t is a result of inner joining two tuples from T_1, T_2 on shared values in common columns C . This is similar to taking $T_1 \uplus T_2$, and applying subsumption and complementation on tuples with shared values in C (Lemma 12) to get t . Since t does not share any values in C with other tuples, it cannot be subsumed. Thus, $t \in \beta(\beta((T_1 \bowtie T_2) \uplus T_1) \uplus T_2)$.
- 2) $t.C \in T_1.C \setminus T_2.C \implies t \in (T_1 \bowtie T_2) \setminus (T_1 \bowtie T_2)$. When we take $(T_1 \bowtie T_1) \uplus T_1$, we append all tuples from T_1 to $T_1 \bowtie T_2$. After applying subsumption, all tuples from T_1 that are used in $T_1 \bowtie T_2$ are subsumed by tuples from $T_1 \bowtie T_2$ on shared values in C . Thus, the only tuples remaining are tuples like t in $(T_1 \bowtie T_2) \setminus (T_1 \bowtie T_2)$. Since t does not share any common values with any tuple in T_2 , it is not subsumed when taking $\beta(\beta((T_1 \bowtie T_2) \uplus T_1) \uplus T_2)$, and so $t \in \beta(\beta((T_1 \bowtie T_2) \uplus T_1) \uplus T_2)$.
- 3) $t.C \in T_2.C \setminus T_1.C \implies t \in (T_2 \bowtie T_1) \setminus (T_1 \bowtie T_2)$. Taking the subsumption of $\beta((T_1 \bowtie T_2) \uplus T_1) \uplus T_2$ removes all tuples from T_2 that are subsumed by tuples in $T_1 \bowtie T_2$. Since the remaining tuples in T_2 cannot be subsumed by any tuple from T_1 not in $T_1 \bowtie T_2$, $t \in (T_2 \bowtie T_1) \setminus (T_1 \bowtie T_2)$. Thus, $t \in \beta(\beta((T_1 \bowtie T_2) \uplus T_1) \uplus T_2)$.

Thus, all tuples from $T_1 \bowtie T_2$ are contained in the resulting table of $\beta(\beta((T_1 \bowtie T_2) \uplus T_1) \uplus T_2)$.

Next, we show that all tuples in $\beta(\beta((T_1 \bowtie T_2) \uplus T_1) \uplus T_2)$ are contained in the resulting table of $T_1 \bowtie T_2$. Let's consider tuple $t' \in \beta(\beta((T_1 \bowtie T_2) \uplus T_1) \uplus T_2)$. There are two cases for C values in tuple t' : $t'.C \in T_1.C \cap T_2.C$ and $t'.C \notin T_1.C \cap T_2.C$.

- 1) $t'.C \in (T_1.C \cap T_2.C) \implies t' \in (T_1 \bowtie T_2)$. Since $(T_1 \bowtie T_2) \subseteq (T_1 \bowtie T_2)$, $t' \in (T_1 \bowtie T_2)$.
- 2) All tuples in $((T_1 \bowtie T_2) \uplus T_1) \uplus T_2$ are either subsumed by tuples from $(T_1 \bowtie T_2)$, are in $T_1 \setminus (T_1 \bowtie T_2)$, or are in $T_2 \setminus (T_1 \bowtie T_2)$. Thus, $t'.C \notin T_1.C \cap T_2.C \implies t' \in (T_1 \uplus T_2) \setminus (T_1 \bowtie T_2) \implies t' \in T_1 \bowtie T_2$.

Thus, all tuples from $\beta(\beta((T_1 \bowtie T_2) \uplus T_1) \uplus T_2)$ are contained in the resulting table of $T_1 \bowtie T_2$.

Now that we have shown that tuples from $T_1 \bowtie T_2$ are contained in the resulting table of $\beta(\beta((T_1 \bowtie T_2) \uplus T_1) \uplus T_2)$ and vice versa, we have shown that $\beta(\beta((T_1 \bowtie T_2) \uplus T_1) \uplus T_2)$ is an equivalent query to $T_1 \bowtie T_2$. ■

4) *Proof of Lemma 15 [Cross Product]:* Given two tables T_1, T_2 , each with columns C_{T_1}, C_{T_2} respectively and do not share any columns, and a constant column c , $T_1 \times T_2$ can be expressed by an equivalent query containing Outer Union,

projection, and complementation. Specifically, $T_1 \times T_2$ is equivalent to query $\kappa(\pi((C_{T_1}, c), T_1) \uplus \pi((C_{T_2}, c), T_2))$.

Proof: Since T_1 and T_2 do not share any columns, the complementation operator cannot be applied to $T_1 \uplus T_2$. Thus, we project on all columns C_{T_1} and constant column c in T_1 , and columns C_{T_2}, c in T_2 . This way, T_1, T_2 now share all values in c and we can apply complementation on $\pi((C_{T_1}, c), T_1) \uplus \pi((C_{T_2}, c), T_2)$ since T_1, T_2 . Thus, we iteratively apply complementation on all tuples from T_1 on all tuples from T_2 to form all tuples in $T_1 \times T_2$. Recall that in every tuple in $T_1 \times T_2$, every value in $t.C_{T_1}$ is from T_1 and every value in $t.C_{T_2}$ is from T_2 . Therefore, every tuple in $T_1 \times T_2$ is contained in $\kappa(\pi((C_{T_1}, c), T_1) \uplus \pi((C_{T_2}, c), T_2))$ and every tuple in $\kappa(\pi((C_{T_1}, c), T_1) \uplus \pi((C_{T_2}, c), T_2))$ is contained in $T_1 \times T_2$, and so $\kappa(\pi((C_{T_1}, c), T_1) \uplus \pi((C_{T_2}, c), T_2))$ is an equivalent query to $T_1 \times T_2$. ■

B. Set Similarity

Algorithm 3: Set Similarity

```

1 Input:  $\mathcal{T} = \{T_1, T_2, \dots, T_n\}$ : set of data lake tables;  $S$ : Source Table;  $\tau$ : Similarity Threshold
2 Output:  $\mathcal{T}' = \{T_1, T_2, \dots, T_n\}$ : a set of candidate tables with high syntactic overlap with  $S$ 
3  $\mathcal{T}'_{\text{scores}} \leftarrow \{\}$  //Store a list of scores for each candidate table
4 for all  $S$  columns  $c \in C$  do
5    $\mathcal{T}_C$ , overlapScores  $\leftarrow$  SetOverlap( $\mathcal{T}, c, \tau$ )
6    $\mathcal{T}_C$ , diverseOverlapScores  $\leftarrow$  diversifyCandidates( $\mathcal{T}_C, c, \tau$ )
7   for all tables  $T \in \mathcal{T}_C$  do
8      $\mathcal{T}'_{\text{scores}}[T] \leftarrow$  diverseOverlapScores[ $T$ ]
9 Order  $\mathcal{T}'_{\text{scores}}$  by average diverseOverlapScores, in descending order
10  $\mathcal{T}' \leftarrow$  keys( $\mathcal{T}'_{\text{scores}}$ )
11 for all tables  $T \in \mathcal{T}'$  do
12   alignedTuples  $\leftarrow$  tuples in  $T$  that contain  $S$ 's column values
13   if set overlap of  $T$  values in alignedTuples with  $S < \tau$  then
14     Discard  $T$ ;
15     Remove  $T$  if its values are contained in another table
16      $T' \in \mathcal{T}'$ 
17     Rename  $T$  columns to aligned  $S$  columns
18 return  $\mathcal{T}'$ ;

```

We find candidate tables with values that have high set overlap with those in a Source Table. As shown in Algorithm 3 we perform Set Similarity with an input set of data lake tables \mathcal{T} , the Source Table S , and a similarity threshold τ (Line 1), and output a set of candidate tables (Line 2). We first find a set of candidate tables, where each table contains a column whose set overlap with a column from S (overlapScore) is above a specified threshold (Lines 4-8). This can be done efficiently with a system like JOSIE [10] that computes exact set containment or MATE [44] that supports multi-attribute joins. In addition, when finding tables with columns that have a high set overlap with columns in S , we call diversifyCandidates() (Line 6) to ensure that each candidate table not only has a high overlap with S , but also has minimal overlap with the previous candidates, shown in Diversify Candidates Algorithm 4

Algorithm 4: Diversify Candidates

```

1 Input:  $c$ : column from Source Table;  $\mathcal{T}_C = \{T_1, T_2, \dots, T_n\}$ :
  set of candidate tables with columns having high overlap
  with  $c$ ;  $\tau$ : Similarity Threshold
2 Output:  $\mathcal{T}'_C = \{T_1, T_2, \dots, T_n\}$ : a set of diverse candidate
  tables
3  $\mathcal{T}_{\text{scores}} \leftarrow \{\}$ 
4 for all tables  $T \in \mathcal{T}_C$  do
5    $C \leftarrow$  column from  $T$  with highest set overlap with  $c$ 
6    $\text{Ind}_T \leftarrow$  index of  $T$  in  $\mathcal{T}_C$ 
7   if  $\text{Ind}_T = 0$  then
8     Continue;
9    $C_{\text{prev}} \leftarrow$  column from  $\mathcal{T}_C[\text{Ind}_T - 1]$  with highest set
    overlap with  $c$  //Get column from previous candidate
    table with high overlap with  $c$ 
10   $\text{prevColOverlap} \leftarrow (C \cap C_{\text{prev}})/|C|$  //Set overlap with
    previous column
11   $\text{sourceColOverlap} \leftarrow (C \cap c)/|c|$  //Set overlap with
    column from Source table
12  if  $\text{sourceColOverlap} < \tau$  then
13    Continue;
14   $\text{overlapScore} \leftarrow \text{sourceColOverlap} - \text{prevColOverlap}$ 
15   $\mathcal{T}_{\text{scores}}[T] \leftarrow \text{overlapScore}$ 
16  Order  $\mathcal{T}_{\text{scores}}$  by values in descending order
17   $\mathcal{T}'_C \leftarrow \mathcal{T}_{\text{scores}}.\text{keys}$ 
18 return  $\mathcal{T}'_C$ ;

```

After we find candidate tables for each column in the Source Table, we average over all overlap scores such that each is for a Source Table's column with which they share many values, and rank them in descending order of averaged scores (Line 9). With a set of candidate tables, we find tuples in each candidate table that contain column values from S . Within these aligned tuples, we check if each aligned column in a candidate table, with respect to a column in S , still has high set overlap (above threshold τ). If not, we remove them (Line 14). Next, we remove any subsumed candidate table, whose columns and column values are all contained in another candidate table (Line 15). We then rename each candidate tables' columns to the names of S 's columns with which they align (Line 16), thus implicitly performing schema matching between S 's columns and the columns from the candidate tables that have overlapping values with S 's columns. Finally, we return the set of candidate tables.

C. Expanding Candidate Tables

In order to represent candidate tables as matrices, their tuples need to align with those in the Source Table. However, not all candidate tables may share a key column with the Source Table. Thus, we need to join a given candidate table that does not share a key column with the Source Table with those that do. This way, tuples from all candidate tables can be aligned with tuples from a Source Table using key values.

As illustrated in Expand Algorithm 5, we traverse a graph that consists of candidate tables as nodes and we find a join path between candidate tables that do not have a Source Table's key (start nodes), and candidate tables that do (potential end nodes). If two candidate tables can join on common columns, their nodes are connected with an edge. For each edge, we use standard join-size cardinality estimation to find edge weights [75].

After we find a path from a start node to an end node, we iteratively join all tables in the path, resulting in a table that shares a key column with the Source Table. This way, all candidate tables share a key column with the Source Table.

Algorithm 5: Expand

```

1 Input:  $\mathcal{T} = \{T_1, T_2, \dots, T_n\}$ : set of candidate tables;  $k$ :
  Source Table's key column(s)
2 Output:  $\mathcal{T}_k = \{T_1, T_2, \dots, T_n\}$ : set of candidate tables that
  all now contain Source Table's key column
3 //Initialize Graph
4 nodes = candidate tables
5 edges = tables that have joinable columns
6 edge Weights = value overlap of joinable columns
7  $\text{start\_nodes} = \{\text{candidate tables that do not contain } k\}$ 
8  $\text{end\_nodes} = \{\text{candidate tables that contain } k\}$ 
9 for each start in  $\text{start\_nodes}$  do
10  //Initialize sets and dictionaries
11   $\text{visited} \leftarrow \text{set}()$  //visited nodes
12   $\text{node\_weights} \leftarrow \{\}$  //maximum weights before each
    node
13   $\text{descendant} \leftarrow \{\}$  //best child for each node
14   $\text{max\_weight} \leftarrow 0$  //Maximum weight found so far
15   $\text{end\_node} \leftarrow \text{None}$  //end node for a given start node
16  //Initialize stack for DFS
17   $\text{stack} \leftarrow \text{stack} + \text{start}$ 
18   $\text{visited} \leftarrow \text{visited} + \text{start}$ 
19  while  $\text{stack}$  is not empty do
20     $\text{node} \leftarrow \text{stack.pop}()$ 
21     $\text{unvisited\_children} \leftarrow$  children of  $\text{node}$  not in
      visited
22    //Current child's weight is the weight of the path so
      far, including the edge weight between node and
      current child
23    for each child in  $\text{unvisited\_children}$  do
24       $\text{child\_weight} \leftarrow$ 
25         $\text{node\_weights}[\text{node}] + \text{weight}(\text{node}, \text{child})$ 
26      //update descendant if it contains the maximum
27      sum of weights so far
28      if  $\text{child\_weight} > \text{node\_weights}[\text{child}]$  then
29         $\text{node\_weights}[\text{child}] \leftarrow \text{child\_weight}$ 
30         $\text{descendant}[\text{child}] \leftarrow \text{node}$ 
31        if child is in  $\text{end\_nodes}$  then
32          if  $\text{child\_weight} > \text{max\_weight}$  then
33            //child has  $k$  and the maximum
34            weighted path so far
35             $\text{max\_weight} \leftarrow \text{child\_weight}$ 
36             $\text{end\_node} \leftarrow \text{child}$ 
37             $\text{stack} \leftarrow \text{stack} + \text{child}$ 
38             $\text{visited} \leftarrow \text{visited} + \text{child}$ 
39    if  $\text{end\_node}$  is not null then
40      //reconstruct path with maximum sum of weights by
41      reversing path, starting with found end node
42       $\text{path} \leftarrow []$ 
43       $\text{current\_node} \leftarrow \text{end\_node}$ 
44      while  $\text{current\_node}$  is in  $\text{descendant}$  do
45         $\text{path} \leftarrow \text{path} + \text{current\_node}$ 
46         $\text{current\_node} \leftarrow \text{descendant}[\text{current\_node}]$ 
47       $\text{path.reverse}()$ 
48       $\text{table} \leftarrow$  first node in  $\text{path}$ 
49      for each join\_table in  $\text{path}[1 : ]$  do
50         $\text{table} \leftarrow \text{join}(\text{table}, \text{join\_table})$ 
51       $\mathcal{T}_k \leftarrow \mathcal{T}_k + \text{table}$ 

```

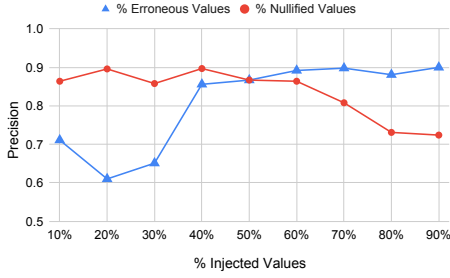


Fig. 8: Gen-T Precision as TP-TR Med has different % of Erroneous (blue triangles) and Nullified Values (red circles).

D. Two-Valued vs. Three-Valued Matrix Representations

After aligning tuples in candidate tables to Source Table’s tuples that share the same key values, we can represent candidate tables as matrices to show how similar the values in candidate tables are to those in the Source Table. These matrices have the same dimensions and indices as the Source Table. First, we consider matrices that consist of binary values, where a 0 in tuple i , column j represents a value at index (i, j) in the candidate table that is different from the value in the same position in the Source Table, and a 1 represent common values in the same indices.

However, populating alignment matrices with binary values does not fully encode how many values in candidate tables differ from those in the Source Table. Specifically, this representation does not distinguish between nullified values in candidate tables (null values in candidate tables at index (i, j) where there is a non-null value in the Source Table at the same position), and erroneous values (different non-null values in candidate tables at index (i, j) from those in Source Tables at the same position).

Instead, we encode matrix representations using three values, where at a given index in an aligned tuple, there is a 1 for a value shared between a candidate table and the Source Table, 0 for a null value in the candidate table where there is a non-null value in the Source Table, and -1 for a non-null value in the candidate table that differs from the value in the Source Table.

E. Metrics

Conditional KL-divergence: Given column C shared between a Source Table and a reclaimed table T , suppose we have probability distributions, \mathcal{P} for C in the Source Table and \mathcal{Q} for C in the reclaimed table. We condition on the key values in key column K . The conditional KL-divergence (or conditional relative entropy) between \mathcal{P} and \mathcal{Q} of sample space X of column C conditioned on key K is as follows:

$$D_{KL}(\mathcal{Q}||\mathcal{P}) = - \sum_{x \in X, k \in K} \mathcal{P}(x|k) \log \left(\frac{\mathcal{Q}(x|k)(1 - \mathcal{Q}(\neg x|k))}{\mathcal{P}(x|k)} \right) \quad (11)$$

Given n non-key columns \mathcal{C} in a Source Table we take the average D_{KL} for each column divided by the probability of a

key value in T matching a key value from the Source Table ($\mathcal{Q}(K)$) and the number of non-key columns (n). Then, the conditional KL-divergence of the reclaimed table is as follows:

$$D_{KL}(T) = \frac{D_{KL}(\mathcal{Q}_1||\mathcal{P}_1) + D_{KL}(\mathcal{Q}_2||\mathcal{P}_2) + \dots + D_{KL}(\mathcal{Q}_n||\mathcal{P}_n)}{\mathcal{Q}(K) * n} \quad (12)$$

The conditional KL-divergence of the reclaimed table is a score $\in [0, \infty)$, with 0 being the ideal score. There is no upper limit on this metric since it naturally approaches ∞ when no key value from the Source Table is found in the reclaimed table.

F. Effectiveness: Ablation Study

Tuning % Erroneous vs. Nullified Values: We further analyze Gen-T’s performance on data lake tables with different number of erroneous and nullified values in TP-TR Med tables (Figure 8). So far, TP-TR Med tables have 50% erroneous values in erroneous versions and 50% nulls in nullified versions (intersection point on the graph where Gen-T has 0.867 Precision). Now, we tune the percentage of values replaced with non-null, random strings (blue line in Figure 8) in erroneous versions, while the nullified versions always contain 50% nulls. Similarly, we tune the percentage of values replaced with nulls (red line in Figure 8) while holding the erroneous versions constant. For Gen-T to produce a perfect reclamation of a Source Table, it should only have originating tables with injected nulls so that these nulls can be replaced with correct values during table integration.

As data lake tables have more erroneous values (blue line), Gen-T is more likely to contain tables with nullified tuples in its set of originating tables, which results in an integrated table with higher precision. On the other hand, as we tune the percentage of values replaced with nulls (red line), precision decreases. As more nulls are injected, these tables also have fewer correct values. Gen-T is thus more inclined to have originating tables with 50% erroneous values, or 50% correct values, leading to a final integration with lower precision.

Effectiveness of Pruning in Gen-T: For a more detailed analysis, we analyze Recall, Precision, and F1-Scores of Gen-T and baseline, ALITE-PS, on TP-TR Med for each of its 26 Source Tables (Figure 9). Note that ALITE-PS directly integrates a set of candidate tables, whereas Gen-T first prunes the set of candidate tables to a set of originating tables before performing table integration. Gen-T outperforms ALITE-PS in Precision for all Source Tables, and outperforms ALITE-PS in Recall for 24 of 26 total Source Tables. This shows that ALITE-PS, which directly integrates candidate tables without pruning, reclaims more Source Tuples than Gen-T, which does include a pruning step, for only a few Source Tables. Also, Gen-T outperforms ALITE-PS in F1-Score for all Source Tables (Figure 9(c)), showing that even if ALITE-PS outperforms Gen-T in Recall, it does not impact the F1-Score.

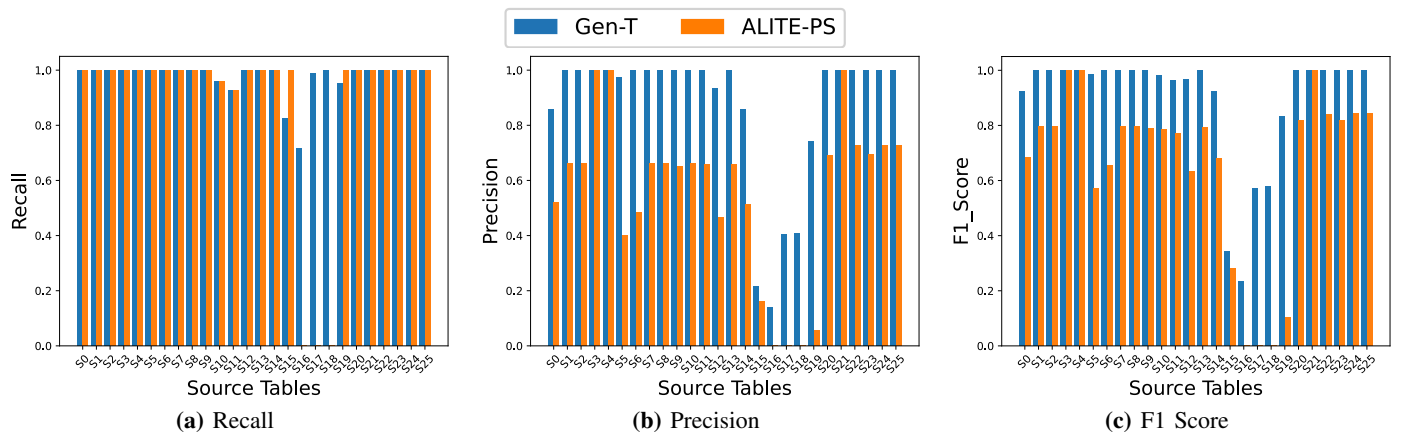


Fig. 9: Recall, Precision, and F1 Scores of Gen-T and ALITE-PS for each Source Table in TP-TR Med benchmark.