

# Aggregation Consistency Errors in Semantic Layers and How to Avoid Them

Ze Zhou Huang  
zh2408@columbia.edu  
Columbia University

Pavan Kalyan Damalapati  
pd2720@columbia.edu  
Columbia University

Eugene Wu  
ewu@cs.columbia.edu  
DSI, Columbia University

## ABSTRACT

Analysts often struggle with analyzing data from multiple tables in a database due to their lack of knowledge on how to join and aggregate the data. To address this, data engineers pre-specify "semantic layers" which include the join conditions and "metrics" of interest with aggregation functions and expressions. However, joins can cause "aggregation consistency issues". For example, analysts may observe inflated total revenue caused by double counting from join fanouts. Existing BI tools rely on heuristics for deduplication, resulting in imprecise and challenging-to-understand outcomes. To overcome these challenges, we propose "weighing" as a core primitive to counteract join fanouts. "Weighing" has been used in various areas, such as market attribution and order management, ensuring metrics consistency (e.g., total revenue remains the same) even for many-to-many joins. The idea is to assign equal weight to each join key group (rather than each tuple) and then distribute the weights among tuples. Implementing weighing techniques necessitates user input; therefore, we recommend a human-in-the-loop framework that enables users to iteratively explore different strategies and visualize the results.

## 1 INTRODUCTION

Cloud technology has enabled enterprises to store unlimited amounts of tables in data warehouses. However, analysts with domain knowledge of the business but with limited expertise in data management face challenges when analyzing metrics of interest across multiple tables. In particular, they may find it difficult to determine which tables to join, what the join conditions are, and how the attributes of these tables relate to the metrics of interest [17, 36].

To bridge this knowledge gap, one popular and easy approach is to denormalize (join) tables into a wide table [10, 20], which is usually carried out by data engineers offline. This simplifies the data analysis process for analysts, who only need to work with one table as the single source of truth. However, join causes spurious duplications, and drops rows due to the lack of matching data or missing values [11, 12, 22], which are difficult to detect, understand, and ameliorate. We illustrate these with the following example:

**EXAMPLE 1.1.** *The example retail business database Table 1 and join graph Figure 1 contain information on the user ad view and purchase history. There are two fact tables: Ad\_view and Purchase. The standard approach of denormalization [10, 20] can lead to incorrect results. The denormalized relations are defined as follows:*

```
CREATE VIEW Denormalized_Ad_View AS SELECT *
FROM V JOIN U ON V.uid = U.uid
      JOIN A ON V.aid = A.aid;
CREATE VIEW Denormalized_Purchase AS SELECT *
FROM H JOIN U ON H.uid = U.uid
      JOIN I ON H.iid = I.iid
      JOIN P ON H.pid = P.pid;
```

User		Purchase History			Item		
uid	name	uid	iid	pid	iid	size	price
1	Joe	1	1	1	1	1	20
2	Mary	2	1	1	2	N	30
		2	2	N	3	5	35
		2	4	2			

Payment		Ad			Ad View	
pid	name	aid	source	cost	uid	aid
1	Paypal	1	Google	500	1	1
2	Visa	2	Facebook	600	1	2
					2	1

Table 1: Retail Business Database. N is NULL.

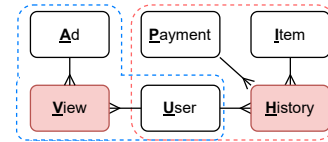


Figure 1: The join graph for a retail business database. The crow's foot represents the "many" relationships. Within this graph, the fact tables are in red. The sets of relations to join for denormalization are enclosed within dotted borders.

Let us consider three simple questions and their SQL queries:

**Q1: What is the total cost of ads from all sources?**

```
SELECT SUM(A.cost) FROM Denormalized_Ad_View;
```

This query above is incorrect because Ad\_view duplicates the cost for each view. In the table below, Google's cost is double-counted in red:

V.uid	aid	U.uid	U.name	A.aid	source	cost
1	1	1	Joe	1	Google	500
1	2	1	Joe	2	Facebook	600
2	1	2	Mary	1	Google	500

We should instead aggregate on only A: `SELECT SUM(A.cost) FROM A`

**Q2: What is the total revenue from the purchased items?**

```
SELECT SUM(I.price) FROM Denormalized_Purchase;
```

Unlike ad cost, item price should be duplicated for each purchase in order to compute the total revenue. The above query is incorrect because there may be missing payments, and joining on NULL values removes the tuple. One fix is to use outer joins.

**Q3: What is the total revenue from different ad sources?**

```
SELECT A.source, SUM(I.price)
FROM V FULL OUTER JOIN U ON V.uid = U.uid
      FULL OUTER JOIN A ON V.aid = A.aid
      FULL OUTER JOIN H ON H.uid = U.uid
      FULL OUTER JOIN I ON H.iid = I.iid
```

```
FULL OUTER JOIN P ON H.pid = P.pid
GROUP BY A.source;
```

The query above is incorrect because of the one-to-many relationship between a user’s revenue (Q2) and their Ad Views; a full outer join would result in duplications and an increase in the total revenue.

One approach widely adopted by the marketing domain [25] is to weigh each ad view based on its “importance”, and ensure that the sum of weights for each join key value (uid) is 1 to counteract join fanouts. The choice of weights is necessarily based on the analyst’s domain knowledge: one analyst may believe that the first ad view is the most important, while another prioritizes the last one [23]. The following illustrates weights where all ad views are equally important:

User		uniformly distributed	Ad View			
uid	revenue		uid	aid	weight	revenue
1	20	→	1	1	0.5	20 × 0.5
2	50		1	2	0.5	20 × 0.5
			2	1	1	50 × 1

For each user, the weights are uniform among ad views. For instance, uid=1 has weights of 1/2=0.5 as there are 2 views. Note that the total revenue from User and Ad View are the same (70). Therefore, analysts can examine how each ad source contributes to revenue based on their assumptions, despite the one-to-many relationship.

Such an idea of “weighing” has been used across multiple domains, as summarized by Kimball and Ross [34]: In order management, freight charges are allocated (or weighed) to a line’s products based on their sizes. In financial services, personal incomes are weighed across individual accounts. In accounting, payments are weighed across organizations according to ownership.

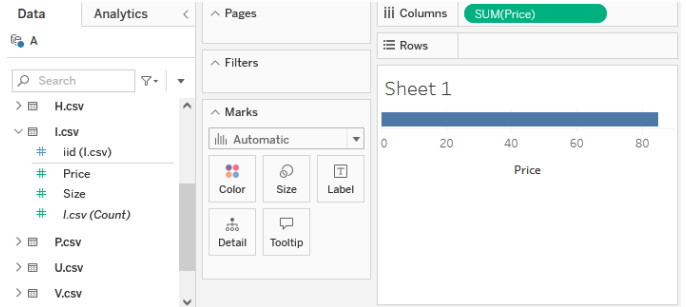
From the above example, we see that the correctness of the aggregates hinges upon the selection of tables to join, deduplication methods, null handling, and weighing designs. Unfortunately, the choices depend on the specific analysis query and the analyst’s understanding of the issue at hand, rendering a singular static interpretation unsuitable for all analyses. Nevertheless, the denormalized wide table is still considered the conventional method.

To overcome denormalization’s limits and tailor the decisions to specific queries, the industry has developed the notion of a “semantic layer” [5, 7, 9, 14, 15, 26] as a way of decoupling the needs and expertise of data engineers and analysts. In the offline phase, a data engineer designs a wide table in the form of a **join graph**, along with appropriate **metrics** (i.e., aggregation functions over expressions) pertinent to a given business problem (by discussing the requirements from analysts). In the online exploration phase, analysts use BI tools to specify **exploratory queries** over the metrics with grouping and filtering. Behind the scenes, the semantic layer employs heuristics to determine the tables to join, deduplication methods, and null handling specifically for each query.

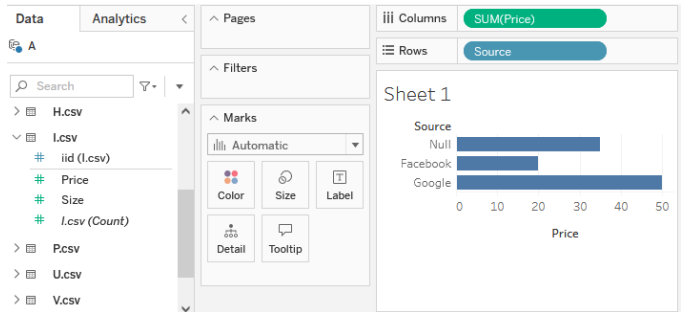
For example, Tableau “relationships” [8] consist of a join graph (Figure 2), and analysts treat it as a base table and create visualizations by dragging and dropping attributes onto the canvas to explore different metrics, such as Q2-3 in Figure 3. Unfortunately, to decide the tables to join for each query, Tableau arbitrarily selects those in the minimum sub-tree of the join graph that cover all referenced attributes, which leads to subtle errors. For example, Q2 only references attribute I.price. Tableau, therefore, sums the



Figure 2: Tableau data relationships.



(a) Q2: SUM(I.price).



(b) Q2: SUM(I.price) GROUP BY A.source.

Figure 3: Tableau sheet, where analysts drag the attributes to visualize without worrying about the underlying join relationships. However, Tableau gives wrong results.

prices over I and computes the total revenue as SUM(I.price)=85 (Figure 3a), even though I.price should be duplicated for each purchase (by joining with H) to compute total revenue correctly. For Q3, Tableau generates results “Google” 50 and “FaceBook” 20, which internally applies specialized deduplication [8] to handle the join fanout. However, the deduplication mechanism is not surfaced to the analysts, even though it impacts query results and may not align with the analysts’ intentions.

Tableau is not alone in making incorrect heuristics. Of the 5 BI tools and semantic layers we surveyed, 2 of them determine the tables to join in a heuristic manner that’s not disclosed to the analyst, which leads to incorrect outcomes for Q2. Additionally, for Q3 with many-to-many relationships, 2 of these BI tools don’t support it, and the remaining 3 apply arbitrary deduplication rules solely based on heuristics. They are therefore all susceptible to correctness errors, and hide from the analyst the ability to interpret and control how the final metrics in a query are decided upon.

How do analysts make informed decisions about how to handle duplication problems that arise from joins? How can analysts best decide to include a table in a query, duplicate or deduplicate tuples, or reweigh duplicates, to correctly compute their desired business metrics? This paper formalizes these issues and proposes “weighing” as the core primitive to address them, whose idea is as

follows: Naive query treats every tuple as having an equal "weight" when aggregated. However, with join fanout, the "weights" of join key groups (calculated as the sum of the "weights" of their tuples) can be amplified, which may inadvertently bias the join results. To overcome this, we "weigh" tables by assigning equal "weight" to each join group (instead of each tuple) and then distributing the weight among the tuples within the group. Consequently, tuple values are aggregated based on these weights to counteract join fanouts. Previous BI tools use deduplication, which can be considered as a special "weighing" of 1/0 without fractional numbers. We demonstrate how reweighing can resolve fanout-induced errors, and how instances of this formulation have been used in existing problems such as market attribution, unbiased sampling, and ML fairness.

A key challenge that necessitates a human-in-the-loop approach is that the appropriate weighing cannot be determined offline, but is dependent on the specific wide table, query, and even analyst needs. For instance, the weights for **Q3** may vary depending on different analysts' belief in the importance of different ad views—Shannon may only care about the ad user viewed last, while Erika may care about all ad views equally. To help users specify the weights and visualize the outcomes, we propose a human-in-the-loop framework. We observe that visualizing each table one-by-one makes it difficult to contextualize the aggregates across joins, and presenting the full join can be overwhelming. Our framework addresses these issues by (1) enabling users to decide the weights iteratively along the join path, (2) visualizing partial aggregates that summarize the current aggregates for the decided weights (instead of computing the full join), and (3) providing an interface with common weighing options to declaratively specify the weights.

## 2 BACKGROUND

This section surveys the existing academic literature for pitfalls that can arise when aggregating over join graphs, and identifies the missing gaps and limitations of existing approaches. We further survey prominent BI tools for how they address these pitfalls.

### 2.1 Pitfalls in Join Aggregation Query

"Summarizability" [24, 37, 43] studies the correctness of querying aggregated fine-grained values at a coarser level. Within snowflake-schema multidimensional models, fact values are usually fine-grained, while aggregation queries that join and group dimensional attributes are considered coarse-grained. To address the fanout issues in **Q1** and **Q2**, they require users to specify the "level of detail" of metrics, which determines what should be joined for duplications before the values are aggregated. They cover other additive issues not directly related to join fanouts. For example, age should not be added, and population can be added over geographical areas but not over time. These are complementary to this work.

However, "summarizability" is too strict for practical exploratory queries. For example, when there are missing join keys in dimension tables, "roll up" and "drill down" are considered "incomplete" and therefore not summarizable, which can be addressed by outer join in practice [8]. Besides, many-to-many joins are considered "non-strict" and also not summarizable, but are important in applications

like market attribution [25] and order management [34]. In such cases, weighing can be used (Section 3.3) to counteract join fanouts.

The same join-induced issues arise in domains beyond BI. Previous analytics works over joins, such as ML [41] and sampling [42], are susceptible to bias [38] and correctness errors when computed over the materialized join result. Consider the following:

**EXAMPLE 2.1.** Consider the database shown in Table 2. Suppose analysts aim to train a model to predict user purchases or sample tuples for insights. They want to include features in the join of all three tables  $A \bowtie U \bowtie H$  for enrichment. Despite having an even gender distribution in  $U$  (1 for each), the full join produces 6× more tuples for female due to fanout. Consequently, this creates a potential imbalance in the ML training and bias for sampled data.

Ad View		User		Purchase History	
aid	uid	uid	gender	uid	iid
1	1	1	male	1	1
2	2	2	female	2	1
2	2	2	female	2	2
				2	4

**Table 2: Example databases where gender is evenly distributed in User, but joining it with Purchase results in bias.**

A common approach to address data imbalance in the single table regime is via weighing [28]. We extend weighing to analytics over joins. Another approach to address join fanout is to pre-aggregate (e.g., averaging) [27, 31] to reduce a N-M relationship into a N-1 relationship. However, this introduces errors for long join paths (e.g., the average of averages is not the total average), and so BI tools that apply pre-aggregation only support a single join [4].

### 2.2 Current Business Intelligence Tools

Modern BI tools offer different ways to pre-define join graphs and metrics offline, and then allow for online exploratory analysis. For join graphs and metrics, some automatically infer default joins from the database schema, while others offer GUI, SQL, or custom languages for user inputs. During online exploration, analysts query the join graphs as a wide table, and the system dynamically determines the appropriate joins to avoid logical pitfalls. In all cases,

Tool	Method	Source	Q1	Q2	Q3
Tableau	Blend		A	I	NA
	Relationship		A	I	WRG
PowerBI	Relationship		A	I	ERR
Looker	Blend		A	I	NA
	Model	V/H	V $\bowtie$ A	H $\bowtie$ I	WRG
	Model	A/I	A	I	
Malloy	Model	V/H	V $\bowtie$ A	H $\bowtie$ I	WRG
	Model	A/I	A	I	
Sigma	Lookup	V/H	V $\bowtie$ A	H $\bowtie$ I	NA
	Lookup	A/I	A	I	NA

**Table 3: The joins current BI tools use to execute queries.  $\bowtie$  is left outer join. Some methods allow source tables to specify the duplication of measurements. Red are wrong results. NA means the queries can't be specified. ERR means the queries can be specified but errors are thrown. WRG means that the results are based on unintended deduplication mechanisms.**

they vary in how they handle different fanout conditions (1-1, 1-N, N-M) and NULLs. We now examine how popular BI tools address the above pitfalls for **Q1-3**. Our findings are summarized in Table 3.

- **Blend** is adopted by both Tableau and Looker [4]. The user specifies the join condition between two tables and query, and the BI tool executes it using pre-aggregation before the outer join. Blend is limited in two ways: (1). It only supports queries of join path with a length of 2 (and not general join graphs), therefore doesn't support **Q3** (with a join path of 5 tables). (2). It doesn't require an explicit definition of the metrics and the decision to join is based on whether the attributes are referenced. For **Q2**, it avoids the join since the referenced attributes ( $I.price$ ) are not referenced in  $H$ , which is incorrect because the join is necessary to duplicate prices for the total revenue metric.
- **Relationship** is supported by Tableau and PowerBI, where users specify the join graph (tables and join conditions) and can query over it directly. However, both tools make the same error as Blend for **Q2** as they don't require explicit metric definitions. For **Q3**, PowerBI generates errors as it doesn't support many-to-many relationships, while Tableau applies arbitrary deduplication that impacts query results but may deviate from the analyst's intentions without even notifying the analyst.
- Looker [5] and Malloy [16] use a SQL-based declarative language to **model** a source table to query against, which defines both join graph and metrics. For instance, the following is used for **Q2**, which defines the join condition with  $I$  to construct a join graph, along with the total revenue as metric ( $\text{sum}(I.price)$ ).

```
source: H is table('PurchaseHistory') {
  join_one: I is table('Items') on iid = I.iid
  measure: total_revenue is sum(I.price)
}
```

Note that there are different options for selecting the source table. In the example,  $H$  is the source, but  $I$  could also be the source. The choice of the source table determines the metric duplication level: In this example, because the total revenue metric is sourced from  $H$ , it is aggregated after  $H \bowtie I$ . For **Q1** and **Q2**, the correctness depends on if the correct source table is used for duplication ( $A$  for **Q1** and  $H$  for **Q2** in the **Source** column). For **Q3**, choosing different source tables yields different results, and both BI tools apply arbitrary deduplication without user interactions.

- Finally, Sigma Computing supports queries over attributes across two tables through **lookup** join [6], which employs preaggregations and is not applicable to **Q3** similar to **blend**. Unlike **blend** but similar to **model**, users can specify explicitly whether the aggregation is being performed before or after join for **Q1-2**.

In summary, current BI tools frequently depend on implicit assumptions for deduplication, which can be dangerous if they are incorrect and challenging to identify the errors.

### 3 WEIGHING TO THE RESCUE

In this work, we argue that users should be able to easily specify reweighing policies in order to avoid correctness errors when performing analytics over joins. This section first defines the consistency errors in the context of semantic layers that store pre-defined join graphs and associated metrics. Our solution then builds on

semi-ring aggregation. Although semi-ring aggregation is traditionally used to accelerate join-aggregation queries by pushing the aggregations through joins, this push-down computes partial aggregates that are also useful for reasons about the appropriate reweighing decisions.

#### 3.1 Problem

We consider the setting where a data engineer has pre-defined (1) the metric as an aggregation function and (2) the duplication as an acyclic join; the aggregation and join together constitute a *base query*  $Q_{base} = \gamma(R_1 \bowtie \dots \bowtie R_n)$  (without group-by and selection). For example, the base query for **Q2** can be represented as  $Q_{base} = \gamma_{SUM(I.price)}(H \bowtie I)$ . The analyst then composes an SPJA query  $Q$  that uses the same metric but may also include selection and group-by expressions that reference attributes in tables that require left outer joining with additional tables. For instance, to answer **Q3**, they may issue  $Q = \gamma_{A.source, SUM(I.price)}(H \bowtie I \bowtie U \bowtie V \bowtie A)$ .

We next formalize the "consistency" error. To facilitate understanding, we consider "sum" as the aggregation in this section, and we will extend it to arbitrary aggregation when introducing the solution framework. The "consistency" error refers to the inequality between the sum of the groupby<sup>1</sup> results for  $Q$  and  $Q_{base}$ :  $Q$  includes additional joins for enriched group-bys and selections, and the total measurement (e.g., revenue or expense) can be amplified compared to  $Q_{base}$ . Such inconsistently larger results have been complained about across different BI tools [1–3, 12, 13, 17, 36], and we want to help analysts understand and avoid them.

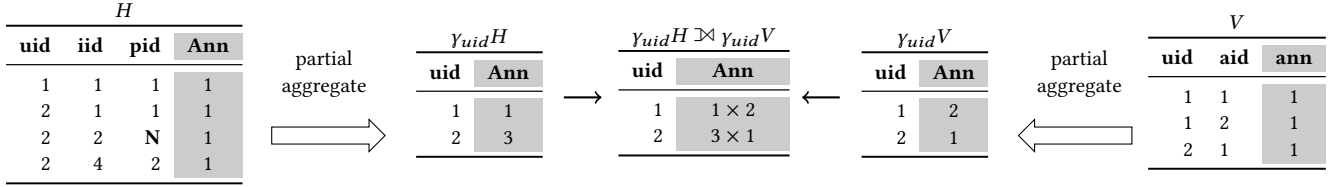
Given the base query  $Q_{base} = \gamma_{SUM}(R_1 \bowtie \dots \bowtie R_k)$ , and the exploration query  $Q = \gamma_{gb, SUM}(\sigma(R_1 \bowtie \dots \bowtie R_k \bowtie \dots \bowtie R_j))$ , that includes additional joins, selections, and group-bys (in red), the objective is to find the re-weighing function  $W$  that weighs (multiplies) the sum when joined (e.g., the revenue is weighed for **Q3** in Example 1.1), such that  $Q^* = \gamma_{gb}(\sigma(R_1 \bowtie \dots \bowtie R_k \bowtie W(R_{k+1}) \dots W(R_j)))$  is:

- **Consistent:** If  $Q$  doesn't have a selection  $\sigma$ , then  $\gamma(Q^*) = Q_{base}$ . If  $Q$  has a selection  $\sigma$ , let  $\neg\sigma$  be its negation and  $\gamma(Q_{\neg}^*) = \gamma_{gb}(\neg\sigma(W(R_1) \bowtie \dots \bowtie W(R_j)))$ . Then  $\gamma(Q^*) + \gamma(Q_{\neg}^*) = Q_{base}$ .
- **User-Directed:** There could be various valid weighing strategies. For example, in **Q3**, different analysts may assign different weights based on their opinions about the importance of the ad views. The weighing should be transparent and understandable to analysts, who should be able to guide based on their interpretation of the data and domain knowledge.

We observe that the re-weighing function  $W$  is applied to relations only in the exploratory query ( $R_{j+1} \bowtie \dots \bowtie R_k$ ) to mitigate fanout effects, but not to those in the base query. The duplication in the base query is assumed intentional to compute the total metric correctly. Although there might be other reasons to weigh the base query relations, such as addressing an imbalanced distribution [38], we consider these as a preprocessing step, separate from  $W$ .

**Scope.** Many previous works [19, 30, 33, 40] studied syntax and semantic issues related to SQL queries for correctness, which remains challenging. We specifically focus on the amplified or reduced aggregate result caused by the join fanout and ensure "consistency" as

<sup>1</sup>Selection removes data and the inequality is expected. For the problem definition purposes, we regard "selection" as a groupby based on whether the selection predicate is satisfied and post-process it to obtain the final selected value.



**Table 4: Partial aggregations  $\gamma(H \bowtie V) = \gamma(Y_{uid}H \bowtie Y_{uid}V)$  reduce many-to-many joins to one-to-many joins.**

formally defined above. We hope that, ensuring consistency helps users identify semantic errors and achieve final query correctness.

### 3.2 Semi-ring Aggregation Preliminary

Semi-ring aggregation breaks down the aggregation into two fundamental operations: addition (e.g., to aggregate values) and multiplication (e.g., to weigh values). It is highly expressive and can express nearly all common aggregations with the benefit of efficient partial aggregation, such as sum, count, average, and max. Other aggregations, such as median, can also be expressed using semi-ring aggregation (by exhaustively tracking all values in a long list), but don't derive as much benefit from partial aggregation for both performance and interpretability perspectives.

**Data Model.** We use the traditional relational data model: Given relation  $R$ , let  $A$  be an attribute,  $dom(A)$  be its domain,  $S_R = [A_1, \dots, A_n]$  be its schema,  $t \in R$  be a tuple of  $R$ , and  $t[A]$  be the value of attribute  $A$  in tuple  $t$ . The domain of  $R$  is then the Cartesian product of attribute domains, i.e.,  $dom(R) = dom(A_1) \times \dots \times dom(A_n)$ .

**Semi-ring Aggregation Query.** We begin by extending the relation table with annotations, [29, 32, 39], which maps  $t \in R$  to a commutative semi-ring  $(D, +, \times, 0, 1)$ , where  $D$  is a set,  $+$  and  $\times$  are commutative binary operators closed over  $D$ , and  $0/1$  are the zero/unit elements. Annotations are useful for query optimizations based on algebraic manipulation. Different semi-ring definitions support various aggregation functions, from standard statistical functions to machine learning models. For example, the natural numbers semi-ring  $(\mathbb{N}, +, \times, 0, 1)$  allows for integer addition and multiplication, and supports the *COUNT* aggregate. For an annotated relation  $R$ , let  $R(t)$  represent the annotation of tuple  $t$ . Tuples in the domain but not in the table are assumed to have 0 as annotations.

Aggregation queries can now be redefined over annotated relations by translating group-by and join operations into  $+$  and  $\times$  operations over the semi-ring annotations, respectively:

$$(\gamma_A R)(t) = \sum \{R(t_1) \mid t_1 \in R, t = \pi_A(t_1)\} \quad (1)$$

$$(R \bowtie T)(t) = R(\pi_{S_R}(t)) \times T(\pi_{S_T}(t)) \quad (2)$$

(1) The annotation for each group-by result in  $\gamma_A R$  is the sum of the annotations of all tuples in its input group. (2) The annotation for each join result in  $R \bowtie T$  is the product of semi-ring annotations from its contributing tuples in  $R$  and  $T$ . The  $\bowtie$  can be extended for outer join, where the non-matching tuple retains its original annotation and has the rest of the attributes as NULL.

To translate an aggregation function like SUM into semi-ring aggregation, we need to determine (1) the semi-ring and (2) the annotation for initial tables. For example, in the case of  $\text{SUM}(I.\text{price})$ , the semi-ring is the real number following standard mathematics.

For the annotation, only  $I$  has each tuple  $t$  annotated with  $t[\text{price}]$ ; all other tables have all annotations of 1.

**Partial Aggregation.** The key optimization in factorized query execution [18, 41] is to distribute aggregations (additions) through joins (multiplications). Consider the  $\gamma(H \bowtie V)$  for count semi-ring, which is a many-to-many join. We can push down part of aggregations before join for one-to-many as illustrated in Table 4.

### 3.3 Consistency through Weighing

There are many aggregations like min/max [37] and other specialized probabilistic databases [18, 35] that provides consistent results, even with many-to-many joins. We highlight the key property they have that ensures consistency, and generalize it to other aggregations through reweighing.

Given  $Q_{base} = \gamma(R_1 \bowtie \dots \bowtie R_k)$ , consider the query with an additional join  $\gamma(Q) = \gamma(R_1 \bowtie \dots \bowtie R_k \bowtie R_{k+1}) = \gamma(\gamma_J(R_1 \bowtie \dots \bowtie R_k) \bowtie \gamma_J(R_{k+1}))$ , where  $J$  is the join key between  $R_{k+1}$  and the rest of the relations. Then, the key sufficient condition for consistency is that:

$$\forall j \in Dom(J), \gamma_J(R_{k+1})(j) = 1$$

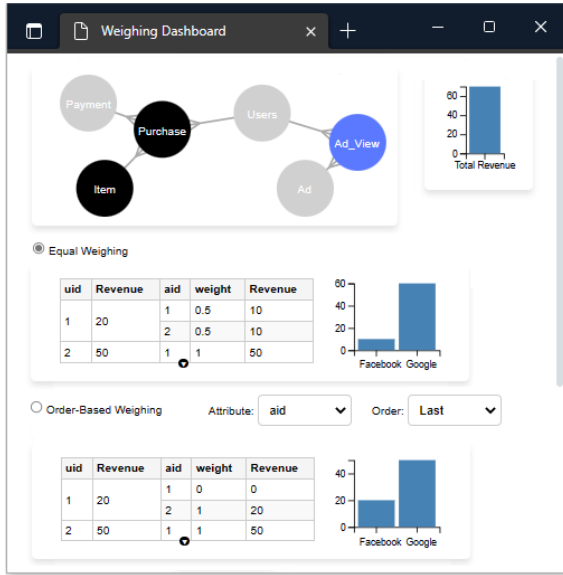
The consistency  $\gamma(Q) = Q_{base}$  is ensured because of the multiplicative identity property of 1 element in semi-ring, and can be recursively applied to an arbitrary number of joins. This is also satisfied with additional group-bys (because the groupings are summed) and selections (as special group-bys of whether selected or not).

Previous applications that support many-to-many joins satisfy this property. For probabilistic tables, it is ensured that each conditional probability table has the sum of the probabilities conditioned (selected) on the join key as 1. In the case of min/max aggregation, the relation without the attribute to maximize have all tuples annotated with 1, and the min/max semi-ring has its addition operator  $\oplus$  as max that ensures  $1 \oplus \dots \oplus 1 = \max(1, \dots, 1) = 1$ .

For other aggregation queries such as sum and count, such a property is not natively satisfied and it's necessary to weigh relations. For instance, in the context of market attribution of **Q3**, we can assign equal weights to all tuples in  $V$  with the same `uid`. However, assigning weights is not one-time, as different analysts may have different opinions on how to weigh the ad views. Therefore, we next introduce a framework to assist users in assigning weights.

### 3.4 Human-in-the-loop Weighing

To make informed decisions, users need to contextualize aggregation within the larger join graph and be capable of specifying weights for the joined table. However, visualizing each table individually makes it challenging to contextualize aggregates across joins, while presenting the full join can be overwhelming. A key design decision for understanding the aggregation relationship is utilizing partial aggregates. These partial aggregates progressively join with



**Figure 4: The weighing interface screenshot displays the join graph and base query result visualizations at the top. The tabular view of the weighed relations and visualizations are displayed for different weighing options.**

new tables, compress the results by aggregating out attributes, not in use and retaining only the necessary attributes for users to comprehend and design weights. The challenge in specifying weights lies in the worst-case scenario, where users need to assign a weight for each tuple, which is impractical. Instead, we identify common and efficient special cases that simplify the process for the user. We aim to determine the minimum information required, given reasonable assumptions for specific settings that are clearly stated and easier for the analyst to specify. This section delves into these details and introduces an interface.

**Interface.** Our interface comprises two panels illustrated in Figure 4 to solicit weighings from users, and visualize outcomes.

**3.4.1 Weighing.** For exploratory query  $Q$ , we request weighing only the necessary relations in a depth-first fashion, where the relations in  $Q_{base}$  are the roots. Some relations, such as those in the  $Q_{base}$  don't require weighing (Section 3.1). Other relations, with one-to-many or one-to-one relationships (to the parent relation in the depth-first search), have one tuple per join key and are directly annotated with a weight of 1 without user request. We only request weights to relations that (1) have many-to-one or many-to-many relationships and (2) are part of any paths from the  $Q_{base}$  relations to the  $Q$  relations, excluding the  $Q_{base}$  relations themselves.

Based on the sophistication of users, we offer parameterized options for common cases to streamline the process such as:

- *Equal Weighing:* Tuples within the group have the same weight.
- *Order-Based Weighing:* The first/last tuple ordered by some user-specified attribute has weight 1, while the rest have 0.
- *Position-Based Weighing:* Similar to order-based one, but users can additionally determine the allocation percentages for the first/last tuple and the rest. E.g., the first and last have weights of 0.4, while the remaining 0.2 is distributed evenly for the rest.

- *Proportional Weighing:* The weights are distributed proportionally to some user-specified attribute (e.g., distributing freight charges based on the item sizes [34]).

For advanced users, we provide a SQL-based interface that allows for customized weighing specifications. The SQL query is intended to create a weight column. Since relations are unordered, we ask users to create a weight table  $W[rowid, weight]$  and then join it with the relation to append the weight column. For instance, the uniform weighing (linear attribution) in Example 1.1 Q3 can be specified using the following SQL queries:

```
SELECT rowid, 1/COUNT(*)
      OVER (PARTITION BY uid) AS weight
FROM V;
```

Users can then modify SQL queries to customize the weights as per their requirements. The weighing can also be based on multiple attributes and even other tables in SQL. Once the weight column is specified, we perform sanity checks: For exploratory queries, we verify whether the weights grouped by the join key are all 1.

**3.4.2 Visualization.** We have developed two types of visualizations for the weighing interface: (1) Tabular views of weights for detailed weights, (2) Visualizations of query results (e.g.,  $Q_{base}$ ,  $Q$ ) over the whole join along with the join graph for an overview view. During the weighing process, we provide users with a tabular view of the relation with weights for detailed inspection. It is necessary to visualize both the relation being weighed and the relation it joins, to understand the distribution of the aggregation. But these two relations have potentially many-to-many relationships, which is hard to visualize. To address this, we implement partial aggregation to summarize aggregates of the parent table from the depth-first search grouped by its join key and present it as a one-to-many join using nested table layout [21]. For instance, in Figure 4, we aggregate the revenue of each user by the join key "uid". The relation could be large, and we only sample  $n$  ( $= 100$  by default) join key groups to display, but the table can be expanded to display more rows (by clicking on the button at the bottom of the table).

For visualizations, we display  $Q_{base}$  and  $Q$  by default but offer a library for users to build extra visualizations. Before weighing, we use the "equal weighing" by default to compute the query results for visualizations, which are progressively updated as users specify weighing. We place  $Q_{base}$  at the top since it remains unaffected by weighing and show other possible visualizations for various weighing options to assist users in understanding potential outcomes. Additionally, we present the join graph, which depicts the relationships (e.g., many-to-many, one-to-one), highlighting the  $Q_{base}$  relations (in black), the relation being weighed (in blue).

### 3.5 Limitations and Future works

Our existing interface evaluates each metric independently, necessitating users to examine and navigate various metrics individually. However, users might prefer to compare and visualize multiple metrics collectively and further generate metrics based on the existing ones. In our future work, we plan to create a unified interface that not only displays multiple metrics but also enables users to seamlessly comprehend weighed outcomes across diverse metrics and weighing options, in order to fully actualize the semantic layer.

## REFERENCES

- [1] 2012. Joins cause incorrect values with aggregated measures. <https://community.tableau.com/s/question/0D54T00000C5aT5SAJ/joins-cause-incorrect-values-with-aggregated-measures-can-i-specify-aggregation-before-i-calculate-a-sum>.
- [2] 2014. Aggregate after join without duplicates. <https://stackoverflow.com/questions/26799813/aggregate-after-join-without-duplicates>.
- [3] 2017. Merge two tables while summing values from subtable. <https://community.powerbi.com/t5/Desktop/Merge-two-tables-while-summing-values-from-subtable/m-p/201112>.
- [4] 2020. Blend Your Data. [https://help.tableau.com/current/pro/desktop/en-us/multiple\\_connections.htm](https://help.tableau.com/current/pro/desktop/en-us/multiple_connections.htm).
- [5] 2020. Looker data modeling. <https://www.looker.com/platform/data-modeling/>.
- [6] 2020. Lookup Join. <https://help.sigmacomputing.com/hc/en-us/articles/4408784908819-Lookup-Join>.
- [7] 2020. Microstrategy Semantic Graph. <https://www.microstrategy.com/pt/resources/video/what-is-a-semantic-graph>.
- [8] 2020. Relationships: Data modeling in Tableau. <https://www.tableau.com/blog/relationships-tableau-data-model>.
- [9] 2020. The Tableau Data Model. [https://help.tableau.com/current/online/en-us/datasource\\_datamodel.htm](https://help.tableau.com/current/online/en-us/datasource_datamodel.htm).
- [10] 2021. Kimball in the context of the modern data warehouse: what's worth keeping, and what's not. <https://www.youtube.com/watch?v=3OcS2TMXELU>.
- [11] 2021. Left join not operating correctly, appears to be working as an inner join. <https://community.looker.com/lookml-5/left-join-not-operating-correctly-appears-to-be-working-as-an-inner-join-18211>.
- [12] 2022. Aggregate functions gone bad and the joins who made them that way. <https://community.looker.com/blog-archives-1027/aggregate-functions-gone-bad-and-the-joins-who-made-them-that-way-29370?postid=54109#post54109>.
- [13] 2022. Tackling the complexity of joining snapshots. <https://docs.getdbt.com/blog/joining-snapshot-complexity>.
- [14] 2023. BI solution architecture in the Center of Excellence. <https://learn.microsoft.com/en-us/power-bi/guidance/center-of-excellence-business-intelligence-solution-architecture>.
- [15] 2023. The dbt Semantic Layer. <https://www.getdbt.com/product/semantic-layer/>.
- [16] 2023. Malloy. <https://github.com/malloymalloy/malloy>.
- [17] Dr Saif Abed. 2020. Tweet. [https://twitter.com/Saif\\_Abed/status/1328488315771817988](https://twitter.com/Saif_Abed/status/1328488315771817988).
- [18] Mahmoud Abo Khamis, Hung Q Ngo, and Atri Rudra. 2016. FAQ: questions asked frequently. In *Proceedings of the 35th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*. 13–28.
- [19] Alireza Ahadi, Julia Prior, Vahid Behbood, and Raymond Lister. 2016. Students' semantic mistakes in writing seven different types of SQL queries. In *Proceedings of the 2016 ACM Conference on Innovation and Technology in Computer Science Education*. 272–277.
- [20] Gleb B. Alex B. 2020. <https://rubygarage.org/blog/database-denormalization-with-examples>. <https://rubygarage.org/blog/database-denormalization-with-examples>.
- [21] Eirik Bakke and David R Karger. 2016. Expressive query construction through direct manipulation of nested relational results. In *Proceedings of the 2016 International Conference on Management of Data*. 1377–1392.
- [22] Itzik Ben-Gan. 2019. T-SQL bugs, pitfalls, and best practices – joins. <https://sqlperformance.com/2019/06/sql-performance/t-sql-bugs-pitfalls-best-practices-joins>.
- [23] Ron Berman. 2018. Beyond the last touch: Attribution in online advertising. *Marketing Science* 37, 5 (2018), 771–792.
- [24] Mónica Caniupán, Loreto Bravo, and Carlos A Hurtado. 2012. Repairing inconsistent dimensions in data warehouses. *Data & Knowledge Engineering* 79 (2012), 17–39.
- [25] Claire Carroll. 2020. Modeling marketing attribution. <https://www.getdbt.com/blog/modeling-marketing-attribution/>.
- [26] Damianos Chatziantoniou and Verena Kantere. 2020. Data Virtual Machines: Data-Driven Conceptual Modeling of Big Data Infrastructures.. In *EDBT/ICDT Workshops*.
- [27] Richard Frank, Flavia Moser, and Martin Ester. 2007. A method for multi-relational classification using single and multi-feature aggregation functions. In *Knowledge Discovery in Databases: PKDD 2007: 11th European Conference on Principles and Practice of Knowledge Discovery in Databases, Warsaw, Poland, September 17-21, 2007. Proceedings 11*. Springer, 430–437.
- [28] Vaishali Ganganwar. 2012. An overview of classification algorithms for imbalanced datasets. *International Journal of Emerging Technology and Advanced Engineering* 2, 4 (2012), 42–47.
- [29] Todd J Green, Grigoris Karvounarakis, and Val Tannen. 2007. Provenance semirings. In *Proceedings of the twenty-sixth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*. 31–40.
- [30] Torsten Grust, Fabian Kliebhan, Jan Rittinger, and Tom Schreiber. 2011. True language-level SQL debugging. In *Proceedings of the 14th International Conference on Extending Database Technology*. 562–565.
- [31] Hongyu Guo and Herna L Viktor. 2008. Multirelational classification: a multiple view approach. *Knowledge and Information Systems* 17 (2008), 287–312.
- [32] Manas Joglekar, Rohan Puttagunta, and Christopher Ré. 2015. Aggregations over generalized hypertree decompositions. *arXiv preprint arXiv:1508.07532* (2015).
- [33] R Kearns, Stephen Shead, and Alan Fekete. 1997. A teaching system for SQL. In *Proceedings of the 2nd Australasian conference on Computer science education*. 224–231.
- [34] Ralph Kimball and Margy Ross. 2011. *The data warehouse toolkit: the complete guide to dimensional modeling*. John Wiley & Sons.
- [35] Daphne Koller and Nir Friedman. 2009. *Probabilistic graphical models: principles and techniques*. MIT press.
- [36] lloyd tabb. 2022. Tweet. <https://twitter.com/lloydtabb/status/1557032026976313345>.
- [37] Jose-Norberto Mazón, Jens Lechtenböcker, and Juan Trujillo. 2009. A survey on summarizability issues in multidimensional modeling. *Data & Knowledge Engineering* 68, 12 (2009), 1452–1469.
- [38] Fatemeh Nargesian, Abolfazl Asudeh, and HV Jagadish. 2021. Tailoring data source distributions for fairness-aware data integration. *Proceedings of the VLDB Endowment* 14, 11 (2021), 2519–2532.
- [39] Milos Nikolic and Dan Olteanu. 2018. Incremental view maintenance with triple lock factorization benefits. In *Proceedings of the 2018 International Conference on Management of Data*. 365–380.
- [40] Kai Presler-Marshall, Sarah Heckman, and Kathryn Stolee. 2021. SQLRepair: Identifying and repairing mistakes in student-authored SQL queries. In *2021 IEEE/ACM 43rd International Conference on Software Engineering: Software Engineering Education and Training (ICSE-SEET)*. IEEE, 199–210.
- [41] Maximilian Schleich, Dan Olteanu, and Radu Ciucanu. 2016. Learning linear regression models over factorized joins. In *Proceedings of the 2016 International Conference on Management of Data*. 3–18.
- [42] Ali Mohammadi Shanghooshabad, Meghdad Kurmanji, Qingzhi Ma, Michael Shekelyan, Mehrdad Almasi, and Peter Triantafillou. 2021. PGMJoins: Random Join Sampling with Graphical Models. In *Proceedings of the 2021 International Conference on Management of Data*. 1610–1622.
- [43] Eric Simon, Bernd Amann, Rutian Liu, and Stéphane Gançarski. 2022. Controlling the Correctness of Aggregation Operations During Sessions of Interactive Analytic Queries. *ACM Journal of Data and Information Quality* (2022).