

# Subset Sampling over Joins\*

ARYAN ESMAILPOUR, University of Illinois Chicago, USA

XIAO HU, University of Waterloo, Canada

JINCHAO HUANG, The Chinese University of Hong Kong, Hong Kong SAR

STAVROS SINTOS, University of Illinois Chicago, USA

Subset sampling (also known as Poisson sampling), where the decision to include any specific element in the sample is made independently of all others, is a fundamental primitive in data analytics, enabling efficient approximation by processing representative subsets rather than massive datasets. While sampling from explicit lists is well-understood, modern applications—such as machine learning over relational data—often require sampling from a set defined implicitly by a relational join. In this paper, we study the problem of *subset sampling over joins*: drawing a random subset from the join results, where each join result is included independently with some probability. We address the general setting where the probability is derived from input tuple weights via decomposable functions (e.g., product, sum, min, max). Since the join size can be exponentially larger than the input, the naive approach of materializing all join results to perform subset sampling is computationally infeasible. We propose the first efficient algorithms for subset sampling over acyclic joins: (1) a *static index* for generating multiple (independent) subset samples over joins; (2) a *one-shot* algorithm for generating a single subset sample over joins; (3) a *dynamic index* that can support tuple insertions, while maintaining a one-shot sample or generating multiple (independent) samples. Our techniques achieve near-optimal time and space complexity with respect to the input size and the expected sample size.

CCS Concepts: • **Theory of computation** → **Sketching and sampling**.

Additional Key Words and Phrases: Join, Subset Sampling, Direct Access

## ACM Reference Format:

Aryan Esmailpour, Xiao Hu, Jinchao Huang, and Stavros Sintos. 2026. Subset Sampling over Joins. *Proc. ACM Manag. Data* 4, 2 (PODS), Article 117 (May 2026), 26 pages. <https://doi.org/10.1145/3801913>

## 1 Introduction

Sampling is a cornerstone technique in modern data analytics and machine learning, allowing systems to trade a small amount of accuracy for significant gains in performance by processing a representative subset of data rather than the entire dataset [21, 44]. A particularly powerful primitive is *subset sampling* (also known as *Poisson sampling*), where the decision to include any specific element in the sample is made independently of all others, based on a specific probability or weight. This independence is often crucial for theoretical guarantees in randomized algorithms and allows for uncoordinated execution in distributed systems [25, 42].

While subset sampling from explicit lists of items is well-understood, a new challenge arises in modern data pipelines where the dataset of interest is not stored in a single table but is defined

\*This work was partially supported by NSF grant IIS-2348919, and NSERC Discovery Grant.

Authors' Contact Information: Aryan Esmailpour, University of Illinois Chicago, Chicago, USA, [aesmai2@uic.edu](mailto:aesmai2@uic.edu); Xiao Hu, [xiaohu@uwaterloo.ca](mailto:xiaohu@uwaterloo.ca), University of Waterloo, Waterloo, Canada, [xiaohu@uwaterloo.ca](mailto:xiaohu@uwaterloo.ca); Jinchao Huang, [jchuang@se.cuhk.edu.hk](mailto:jchuang@se.cuhk.edu.hk), The Chinese University of Hong Kong, Shatin, Hong Kong SAR; Stavros Sintos, University of Illinois Chicago, Chicago, USA, [stavros@uic.edu](mailto:stavros@uic.edu).



This work is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License.

© 2026 Copyright held by the owner/author(s).

ACM 2836-6573/2026/5-ART117

<https://doi.org/10.1145/3801913>

*implicitly* as the result of a relational join. In this paper, we study the problem of Subset Sampling over Joins: efficiently drawing a random subset from the join results, where every join result is associated with a weight and is included in the sample independently with probability equal to its weight. Crucially, we focus on the standard setting where the weight of a join result is not arbitrary, but is derived from weights assigned to the input tuples via an aggregation function—such as the product, minimum, maximum, or summation of the component weights. A compelling motivation for this problem lies in dataset condensation (or coreset construction) for machine learning over multi-relational data [37, 46].

**Example 1.1.** Consider a high-stakes scientific computing task—virtual screening in drug discovery—where a researcher wishes to train a graph neural network to predict the efficacy of complex molecular interactions. The training data is defined by joining massive tables of *chemical compounds*, *protein targets*, and *metabolic pathways*. The full set of valid interaction subgraphs (the join result) is combinatorially explosive—often exceeding petabytes in size—making it impossible to materialize or exhaustively simulate. To make training feasible, the system must generate a condensed, representative dataset [6, 27]. Here, input tuples carry crucial domain-specific weights: a compound has a binding affinity score, and a protein target has a clinical relevance score. The importance of a resulting interaction tuple is naturally derived from these components (e.g., the product of affinity and relevance, representing the joint probability of success; or the minimum, representing a bottleneck constraint). By performing subset sampling with these derived weights, we create a manageable, high-quality summary of the chemical space.

However, performing this sampling efficiently is algorithmically difficult due to the implicit nature of the data. A naive approach would be to first compute the join results as well as their aggregated weights (e.g., the product of its components), and then apply the classic subset sampling algorithm. This strategy is prohibitively costly because of the “materialization bottleneck”: the size of the join result can be exponentially larger than the input relations [5], which is very costly to compute and store. In this naive process, the algorithm would spend the vast majority of its time generating join results that are immediately discarded by the coin flips, wasting massive amounts of computation and memory. To the best of our knowledge, no prior work has addressed the problem of relational subset sampling except this baseline. In this paper, we aim to fill this gap by providing the first efficient subset sampling algorithms that operate directly on the input relations. By exploiting the join structure of the query and the decomposable nature of the weight functions, we can effectively “skip” over the vast sea of rejected tuples and directly generate the successful samples. Our collective results establish subset sampling over joins as a tractable database primitive, which will pave the way for scalable, interactive analytics over massive multi-relational datasets.

## 1.1 Problem Definition

**Subset Sampling.** Given a set of elements  $S = \{e_1, e_2, \dots, e_n\}$  and a function  $p$  that assigns each element  $e \in S$  a probability  $p(e)$ , a *subset sampling* query asks for a *subset sample* of  $S$ . A subset sample of  $S$  is a random subset  $X \in 2^S$ , where each element  $e \in S$  is sampled into  $X$  independently with probability  $p(e)$ . More formally, the distribution of  $X$  is:  $\Pr[X = Y] = (\prod_{e \in Y} p(e)) (\prod_{e \in S \setminus Y} (1 - p(e)))$ ,  $Y \subseteq S$ . We refer to  $\Psi = \langle S, p \rangle$  as a *subset sampling problem instance*. We denote the expected size of  $X$ , i.e.,  $E[|X|] = \sum_{e \in S} p(e)$ , as  $\mu_\Psi$ .

**(Natural) Joins.** Let **att** be a set where each element is called an *attribute*, and **dom** be another set where each element is called a *value*. A *tuple* over a set  $U \subseteq \mathbf{att}$  is a function  $\mathbf{u} : U \rightarrow \mathbf{dom}$ . For any subset  $U' \subseteq U$ , define  $\mathbf{u}[U']$  as the tuple  $\mathbf{u}'$  over  $U'$  such that  $\mathbf{u}'(A) = \mathbf{u}(A)$  for every  $A \in U'$ . We also call  $U$  the support attributes of  $\mathbf{u}$ . A *relation* is a set  $R$  of tuples over the same set  $U$  of attributes; we call  $U$  the *schema* of  $R$ , a fact denoted as  $\text{schema}(R) = U$ . Given a subset

$U \subseteq \text{schema}(R)$ , the *projection of  $R$  on  $U$* —denoted as  $R[U]$ —is a relation with schema  $U$  defined as  $R[U] = \{\text{tuple } \mathbf{u} \text{ over } U \mid \exists \text{ tuple } \mathbf{v} \in R \text{ s.t. } \mathbf{v}[U] = \mathbf{u}\}$ . We represent a *natural join* (henceforth, simply a “join”) as a set  $Q$  of relations. We denote the number of relations in  $Q$  as  $k = |Q|$ . The input size of  $Q$  is defined as the total number of tuples in  $Q$ , denoted as  $N = \sum_{R \in Q} |R|$ . Define  $\text{attset}(Q) = \bigcup_{R \in Q} \text{schema}(R)$ . The join result is the following relation over  $\text{attset}(Q)$ :

$$\text{Join}(Q) = \{\text{tuple } \mathbf{u} \text{ over } \text{attset}(Q) \mid \forall R \in Q, \mathbf{u}[\text{schema}(R)] \in R\}.$$

The join  $Q$  can be characterized by a *schema graph*  $G_Q = (V, E)$ , which is a hypergraph where each vertex in  $V$  is a distinct attribute in  $\text{attset}(Q)$ , and each edge in  $E$  is the schema of a distinct relation in  $Q$ . The set  $E$  may contain identical edges because two (or more) relations in  $Q$  can have the same schema. The term “hyper” suggests that an edge can have more than two attributes. The join  $Q$  is *acyclic* [50] if  $G_Q$  is acyclic. A hypergraph is acyclic if there exists a tree  $\mathcal{T}$ , called a *join tree*, whose nodes have a one-to-one correspondence with the relations in  $Q$  (i.e., hyperedges in  $E$ ), such that for any attribute  $A \in V$ , the set of nodes corresponding to relations containing  $A$  forms a connected subtree in  $\mathcal{T}$ . A join that is not acyclic is called a *cyclic join*.

As a variant of join, the semi-join between  $R_1$  and  $R_2$  is defined as  $R_1 \bowtie R_2 = R_1 \bowtie (R_2[\text{schema}(R_1) \cap \text{schema}(R_2)])$ , which returns all tuples in  $R_1$  that can be joined with some tuple from  $R_2$ . When the second operand is a single tuple  $\mathbf{v}$ ,  $R_1 \bowtie \mathbf{v}$  is a shorthand for  $R_1 \bowtie \{\mathbf{v}\}$ . As a special case,  $R \bowtie \emptyset = R$ .

**Subset Sampling over Joins.** Let  $Q = \{R_1, R_2, \dots, R_k\}$  be a join as defined above. Each relation  $R_j \in Q$  is associated with a function  $p_j : R_j \rightarrow [0, 1]$  that assigns a weight from  $[0, 1]$  to each tuple in it. For each join result  $\mathbf{u} \in \text{Join}(Q)$ , its weight is a function  $\mathcal{F}$  defined on the weights of these tuples that constitute this join result:

$$p(\mathbf{u}) = \mathcal{F}(p_1(\mathbf{u}[\text{schema}(R_1)]), p_2(\mathbf{u}[\text{schema}(R_2)]), \dots, p_k(\mathbf{u}[\text{schema}(R_k)])), \quad (1)$$

where  $\mathcal{F}$  can be the MAX, MIN, PRODUCT, or SUM function. It is required that  $p(\mathbf{u}) \in [0, 1]$ . Given a join instance  $Q$  and a set of associated weight functions  $\{p_j\}_{j \in \{1, 2, \dots, k\}}$ , a *relational subset sampling query* is a subset sampling query to the subset sampling instance  $\Psi = \langle \text{Join}(Q), p \rangle$ , i.e., asks to sample a random subset  $X$  of  $\text{Join}(Q)$ , where each join result  $\mathbf{u}$  is sampled into  $X$  independently with probability defined by its weight  $p(\mathbf{u})$ . More formally, the distribution of  $X$  is

$$\Pr[X = Y] = \left( \prod_{\mathbf{u} \in Y} p(\mathbf{u}) \right) \left( \prod_{\mathbf{u} \in \text{Join}(Q) \setminus Y} (1 - p(\mathbf{u})) \right), Y \subseteq \text{Join}(Q). \quad (2)$$

We denote the expected size of  $X$ , i.e.,  $E[|X|] = \sum_{\mathbf{u} \in \text{Join}(Q)} p(\mathbf{u})$ , as  $\mu_\Psi$ .

The first problem we study is an *indexing* (data structure) problem, where we wish to have an index that supports drawing multiple independent subset samples from join results efficiently:

**Problem 1.2** (Indexed Subset Sampling over Joins). Given a join instance  $Q$  and a set of associated weight functions  $\{p_i\}_{i \in \{1, 2, \dots, k\}}$ , the goal is to construct an index for answering subset sampling queries to the instance  $\langle \text{Join}(Q), p \rangle$ . Additionally, the random subsets returned by distinct queries must be independent.

We are interested in the preprocessing time  $t_p$  and space usage of the index constructed, as well as the time  $t_s$  for answering a relational subset sampling query from the index.

The second problem we study is a *one-shot* version of relational subset sampling where we wish to output one subset sample from join results efficiently:

**Problem 1.3** (One-shot Subset Sampling over Joins). Given a join instance  $Q$  and a set of associated weight functions  $\{p_i\}_{i \in \{1, 2, \dots, k\}}$ , the goal is to compute one subset sample from the instance  $\langle \text{Join}(Q), p \rangle$ .

In this case, we are interested in the total answering time. Although any solution to Problem 1.2 gives a solution to Problem 1.3 with total time  $O(t_p + t_s)$ , this may not be efficient enough.

**Dynamic Subset Sampling over Joins.** We also consider the dynamic setting with tuple insertions<sup>1</sup>. Now, each insertion is a triple  $\langle \mathbf{u}, R_h, \eta, p \rangle$  for  $\eta \in \mathbb{Z}^+$ ,  $p \in [0, 1]$  and  $R_h \in Q$ , indicating that tuple  $\mathbf{u}$  is inserted into relation  $R_h$  at timestamp  $\eta$  with weight  $p$ . We follow the set semantics, so inserting a tuple into a relation that already has it has no effect. Consider the stream of input tuples, ordered by their timestamp. Let  $Q^\eta$  be the join defined by the first  $\eta$  tuples of the stream, and set  $Q^0 = \emptyset$ . Also, let  $\{p_i^\eta\}_{i \in \{1, 2, \dots, k\}}$  be the associated weight functions for  $Q^\eta$ . Each join result  $\mathbf{u} \in Q^\eta$  has its probability  $p^\eta(\mathbf{u})$  defined similar to (1). We use  $N$  to denote the length of the stream, which is only used in the analysis. The algorithms will not need the knowledge of  $N$ , so they work over an unbounded stream. We have the corresponding versions of both problems in the dynamic setting:

**Problem 1.4** (Dynamic Indexed Subset Sampling over Joins). Suppose tuples come in a streaming fashion. The goal is to maintain an index for answering subset sampling queries to the instance  $\langle \text{Join}(Q^\eta), p^\eta \rangle$  for every timestamp  $\eta \in \mathbb{Z}^+$ . Additionally, the random subsets returned by distinct queries must be independent.

**Problem 1.5** (Dynamic One-shot Subset Sampling over Joins). Suppose tuples come in a streaming fashion. The goal is to maintain one subset sample from the instance  $\langle \text{Join}(Q^\eta), p^\eta \rangle$  for every timestamp  $\eta \in \mathbb{Z}^+$ .

For the dynamic index, we are interested in the maintenance time  $t_u$  of the index when a tuple is inserted, and the time  $t_s$  for answering a relational subset sampling query from the index. For the dynamic one-shot problem, we care about the total running time. Again, any solution to Problem 1.4 yields a solution to Problem 1.5 with total time  $O(t_u \cdot N + t_s)$ .

For all versions of the relational subset sampling problem, we study the data complexity [1] and analyze the complexity in terms of the data-dependent quantities (such as input size  $N$  and expected size of subset samples  $\mu_\Psi$ ), while taking the schema size of  $Q$  (i.e.,  $|V|$  and  $|E|$ ) as a constant.

**Computation Model.** We discuss our algorithms in the *real RAM* model of computation [10, 45]. In particular, we assume that the following operations take constant time: (i) accessing a memory location; (ii) generating a random value from the standard uniform distribution  $\text{Uniform}(0, 1)$ ; and (iii) performing basic arithmetical operations involving real numbers, such as addition, multiplication, division, comparison, truncation, and evaluating fundamental functions like  $\log$ . We denote by  $\text{Geometric}(p)$  the geometric distribution over  $\{0, 1, \dots\}$  with probability mass function  $\Pr[X = k] = (1 - p)^k p$ . We denote by  $\text{TruncatedGeometric}(p, n)$  the geometric distribution conditioned on the value being strictly less than  $n$ , i.e., with support  $\{0, 1, \dots, n - 1\}$ . Under the RAM model, we can generate a random value  $x$  from  $\text{Geometric}(p)$  in  $O(1)$  time [11] by setting  $x = \lfloor \frac{\log \text{Uniform}(0, 1)}{\log(1-p)} \rfloor$ . Additionally, we can generate a random value  $x$  from  $\text{TruncatedGeometric}(p, n)$  in  $O(1)$  time [11] by setting  $x = \lfloor \frac{\log(1-q \cdot \text{Uniform}(0, 1))}{\log(1-p)} \rfloor$ , where  $q = 1 - (1 - p)^n$ .

**Math Conventions.** For an integer  $x \in \mathbb{Z}^+$ , the notation  $[x]$  denotes the set  $\{1, 2, \dots, x\}$ , and  $\llbracket x \rrbracket$  denotes the set  $\{0, 1, \dots, x\}$ . The notation  $\sqcup$  denotes the disjoint union of sets. Given an ordered set of elements  $X$  and a pair of elements  $x, x' \in X$ , we use  $x \prec x'$  to indicate that  $x$  is smaller than  $x'$ , and  $x \preceq x'$  to indicate that  $x$  is no larger than  $x'$ . We use  $\log x$  to denote  $\log_2 x$ .

<sup>1</sup>In the fully dynamic case with both insertions and deletions, the subset sampling problem (with all weights as 1) is at least as hard as maintaining the join results. For general (more precisely, a non-hierarchical) join query, conditioned on the Online Matrix-Vector (OMv) multiplication conjecture, the update time must be  $\Omega(N^{0.5-\epsilon})$  for every  $\epsilon > 0$  just to maintain the Boolean answer, when both insertions and deletions are allowed [8]. Even for this Boolean problem, establishing a matching upper bound is still open, except for some specific joins.

To help understanding, we list the main notations used by this paper in Table 2.

## 1.2 Related Work

**Uniform Sampling over Joins.** Chaudhuri et al. [14] and Olken [44] first established the uniform sampling over join problem, which returns a uniform sample from the join results. Zhao et al. [52] introduced a linear space index for drawing uniform independent samples from acyclic joins. Later, Chen and Yi [16] investigated the cyclic joins and their results were improved to be conditionally optimal by two independent works [24, 36]. Huang et al. [33] investigated uniform sampling for acyclic joins with selection predicates. Dai et al. [23] proposed a reservoir sampling framework for maintaining uniform samples over joins in streaming settings.

**Non-uniform Sampling over Joins.** Non-uniform sampling has been extensively studied in the context of online aggregation (OLA). The ripple join [30] and its variants generalize nested-loop joins to incrementally estimate aggregates with confidence intervals. Similarly, wander join [38] uses random walks over the join graph to provide unbiased estimators for aggregates like SUM or COUNT. These approaches differ fundamentally from ours: they are designed to estimate scalar statistics using algorithmic probabilities that minimize variance, whereas our goal is to materialize a concrete subset of independent samples according to user-defined importance weights.

Despite the substantial body of work on sampling over joins, we are not aware of any existing approach that can be adapted to support subset sampling over joins, indicating the need for novel techniques. A straightforward approach to subset sampling over joins is to first materialize the join result and then build a standard index over all output tuples. However, the size of the join result can be orders of magnitude larger than that of the input database, making this approach prohibitively inefficient. More broadly, there is a growing line of work on solving optimization problems on join results without fully materializing the join [2–4, 13, 15, 17, 18, 22, 26, 34, 35, 37, 40, 41, 46–49].

**Post-Acceptance Note.** Concurrent to our work, a nice work [7] independently studied the same core statistical problem in slightly different settings. Their problem formulation assumes the sampling probability of a join result is dictated by a single attribute originating from a single base relation. In contrast, our problem definition is broader, supporting probabilities derived from the aggregation of tuple weights across multiple relations via decomposable functions (e.g., PRODUCT, MIN, MAX, SUM). Algorithmically, their work achieves an  $O(N)$  index construction time and  $O(\kappa \log N)$  sampling time (where  $\kappa$  is the sample size), with a focus on practical implementations in main-memory column stores using the Shredded Yannakakis framework.

## 1.3 Our Results

Our main results are summarized in Table 1. In this paper, we focus on acyclic joins, and all of our algorithms can be extended to cyclic joins using the standard tree decomposition approach [29]. In this transformation, the effective input size scales from  $N$  to  $N^w$ , where  $w \geq 1$  represents the fractional hypertree width (or more complex measures such as #submodular width [39]). Consequently, in our complexity results, linear dependencies on  $N$  are replaced by  $N^w$ , while logarithmic dependencies (e.g.,  $\log N$ ,  $\log \log N$ ) remain unchanged, as the width  $w$  is treated as a constant under data complexity. We address three primary scenarios:

- **(Section 3)** For the static setting, we design an index that can be constructed in  $O(N \log N \log \log N)$  time and  $O(N \log N)$  space, supporting subset sampling queries in expected  $O(1 + \mu_\Psi \log N)$  time, where  $N$  is the input size and  $\mu_\Psi$  is the expected sample size.
- **(Section 4)** Then, we present a one-shot algorithm that bypasses index construction to generate a single subset sample in  $O(N \log^2 N + \mu_\Psi)$  expected time.

Static	Method	Preprocessing Time	Sampling Time	Space Usage
Index	Baseline	$O(N +  \text{Join}(Q) )$	$O(1 + \mu_\Psi)$	$O( \text{Join}(Q) )$
	Our Algorithm	$O(N \log N \log \log N)$	$O(1 + \mu_\Psi \log N)$	$O(N \log N)$
One-shot	Baseline	—	$O(N +  \text{Join}(Q) )$	$O( \text{Join}(Q) )$
	Our Algorithm	—	$O(N \log^2 N + \mu_\Psi)$	$O(N \log^2 N + \mu_\Psi)$
Dynamic	Method	Update Time	Sampling Time	Space Usage
Index	Our Algorithm	$O(\log^3 N \log \log N)$	$O(\mu_\Psi \log N)$	$O(N \log N)$
One-shot	Our Algorithm	—	$O(N \log^3 N \log \log N + \mu_\Psi \log N)$	$O(N \log N)$

Table 1. Complexity of subset sampling on acyclic joins (extendable to cyclic joins with  $N$  increased to  $N^w$ ) for product function.  $N$  is the input size,  $|\text{Join}(Q)|$  is the join size, and  $\mu_\Psi$  is the expected output size.

- (Section 5) Finally, we extend our framework to the dynamic setting with insertions, by applying and adapting the dynamic direct access index for the acyclic join in [23].

**Remark 1.** In the main text, we assume the function  $\mathcal{F}$  to be the product of the input weights. In Appendix C, we show how other functions can be supported by simply adapting our algorithms.

#### 1.4 Prior Results as Preliminaries

**Subset Sampling.** The classic subset sampling problem with a given set of elements has been well studied [9, 11, 31, 51]. A naive approach is to iterate through all  $n$  elements and flip a biased coin for each, taking  $O(n)$  time. So, the goal is to design algorithms with query time proportional to the expected output size  $\mu_\Psi$ , which can be significantly smaller than the total number of elements. In the dynamic setting, where elements and probabilities can change, any algorithm requires  $\Omega(1 + \mu_\Psi)$  query time and  $\Omega(1)$  update time [51]. In the static setting, [11] provided the first index with an optimal query time of  $O(1 + \mu_\Psi)$ . However, their index requires  $O(\log^2 n)$  update time, which is suboptimal in the dynamic setting. Recently, [9, 31, 51] optimally solved this dynamic problem with indexes that take expected query time  $O(1 + \mu_\Psi)$ ,  $O(n)$  space and  $O(1)$  update time. [28, 32] explored weighted subset sampling where the probabilities are normalized by the total weight.

**Worst-Case Optimal Joins.** The AGM bound states that for any join instance  $Q$  of size  $N$ , the maximum number of join results is  $\Theta(N^{\rho^*})$  [5], where  $\rho^*$  is the fractional edge covering number of the schema graph  $G_Q$ . The worst-case optimal join algorithms can compute any join instance  $Q$  of size  $N$  within  $O(N^{\rho^*})$  time [43].

**Direct Access for Joins.** The *direct access* problem for joins was first studied by Carmeli et al. [12], which assumes a fixed ordering on the join results in  $\text{Join}(Q)$  and asks for an index to return the join result in some specific position. For any acyclic join  $Q$ , an index can be built in  $O(N)$  time such that each direct access query can be answered in  $O(\log N)$  time [12, 52].

## 2 Classic Subset Sampling Revisited

In this section, we revisit the classic subset sampling problem, which serves as the foundation for our relational algorithms. We rely on the `DirectAccess` oracle for the input set  $S$  (under some fixed ordering<sup>2</sup>). Given an integer  $i \in [|S|]$ , this oracle returns the  $i$ -th element of  $S$  in  $O(1)$  time.

### 2.1 Perfect-Uniform Subset Sampling

A subset sampling problem instance  $\Psi = \langle S, p \rangle$  is *perfect-uniform* if  $p(e) = p$  for any  $e \in S$ , which is also simplified as  $\Psi = \langle S, p \rangle$ . We have the following folklore result for the perfect-uniform scenario:

<sup>2</sup>The ordering can be arbitrary but must remain consistent across multiple invocations of this oracle.

**Lemma 2.1** (Perfect-Uniform Subset Sampling). *For a perfect-uniform subset sampling instance  $\Psi = \langle S, p \rangle$ , there is an index  $\mathcal{D}$  that can be built within  $O(|S|)$  space and time, using which each subset sampling query can be answered in  $O(1 + \mu_\Psi)$  expected time.*

The index  $\mathcal{D}$  consists of an array storing the elements of  $S$ , which allows us to retrieve the  $i$ -th element in  $O(1)$  time. Additionally, we precompute the probability  $q = 1 - (1 - p)^{|S|}$  that the resulting sample is non-empty. The construction takes  $O(|S|)$  time and space. Given this index, there are two algorithms to achieve the query time complexity. Algorithm 1 represents the standard approach using geometric jumps to skip over rejected elements. Algorithm 2, however, utilizes the precomputed probability  $q$  to introduce a preliminary step: it first tosses a coin to decide whether to sample at least one element from  $S$ . While both algorithms are optimal for a single instance, Algorithm 2 offers a special benefit when handling multiple perfect-uniform instances. Instead of invoking the sampling procedure for every instance, this structure allows us to first sample the subset of instances that yield non-empty results, and then perform subset sampling only within those. This strategy significantly improves the total time complexity in the batched setting, as will be detailed in Section 2.3.

---

**Algorithm 1: `uss-vanilla`( $S, p$ )**


---

```

1  $X \leftarrow \emptyset; i \leftarrow 0;$ 
2 while  $i < |S|$  do
3    $i \leftarrow i + 1 + \text{Geometric}(p);$ 
4   Add DirectAccess ( $S, i$ ) to  $X;$ 
5 return  $X;$ 

```

---



---

**Algorithm 2: `uss-advanced`( $S, p$ )**


---

```

1  $X \leftarrow \emptyset; q \leftarrow 1 - (1 - p)^{|S|};$ 
2 if  $\text{Uniform}(0, 1) \leq q$  then
3    $i \leftarrow$ 
4      $1 + \text{TruncatedGeometric}(p, |S|);$ 
5   Add DirectAccess ( $S, i$ ) to  $X;$ 
6   while  $i < |S|$  do
7      $i \leftarrow i + 1 + \text{Geometric}(p);$ 
8     Add DirectAccess ( $S, i$ ) to  $X;$ 
8 return  $X;$ 

```

---

## 2.2 Rejection-based Subset Sampling

The perfect-uniform condition is too restrictive for general applications. We can relax this requirement using a rejection-based strategy without affecting the asymptotic running time. Given a general instance  $\Psi = \langle S, p \rangle$ , we define an upper bound  $p_{\max} = \max_{e \in S} p(e)$ . We construct the index described in Lemma 2.1 for the bounding perfect-uniform instance  $\Psi^+ = \langle S, p_{\max} \rangle$ . This preprocessing takes  $O(|S|)$  time.

The query procedure first retrieves a subset sample  $X$  from  $\Psi^+$  using the index (via Algorithm 1 or Algorithm 2). Then, for each element  $e' \in X$ , we retain it in the final sample with probability  $p(e')/p_{\max}$ . The overall runtime is proportional to the size of the intermediate sample  $X$ . There are two conditions under which this runtime is efficiently bounded:

**Lemma 2.2** ( $\beta$ -uniform Subset Sampling). *Given a subset sampling instance  $\Psi = \langle S, p \rangle$ , if  $\max_{e \in S} p(e) \leq \beta \cdot \min_{e \in S} p(e)$  for some parameter  $\beta \geq 1$ , there is an index  $\mathcal{D}$  that can be built within  $O(|S|)$  space and time, using which each subset sampling query can be answered in  $O(1 + \beta \cdot \mu_\Psi)$  expected time.*

**Lemma 2.3** (Light Subset Sampling). *Given a subset sampling instance  $\Psi = \langle S, p \rangle$ , if  $\max_{e \in S} p(e) \leq 1/|S|$ , there is an index  $\mathcal{D}$  that can be built within  $O(|S|)$  space and time, using which each subset sampling query can be answered in  $O(1)$  expected time.*

## 2.3 Batched Rejection-based Subset Sampling

Consider a general subset sampling problem instance  $\Psi = \langle S, p \rangle$  that is partitioned into a set of  $m$  disjoint sub-instances  $\{\Psi_i = \langle S_i, p_i \rangle\}_{i \in [m]}$ , where  $S = \bigsqcup_{i \in [m]} S_i$ . Each sub-instance has an upper bound  $p_i^+$  on the maximum probability (i.e.,  $\max_{e \in S_i} p(e) \leq p_i^+$ ) and is either light or  $\beta$ -uniform.

**Algorithm 3: ss-rejected-batch**( $\{\Psi_i\}_{i \in [m]}$ )**Input:** Implicit access to the meta-index for  $q$  and sub-indexes for  $\{\Psi_i\}$ .**Output:** A subset sample  $X$ .

---

```

1  $I \leftarrow \text{ss-query}(\Psi_{\text{meta}})$ ; // Query the meta-index
2 foreach  $i \in I$  do
3    $X_i \leftarrow \emptyset$ ;
4    $j \leftarrow 1 + \text{TruncatedGeometric}(p_i^+, |S_i|)$ ; // Simulate Algorithm 2 using sub-index
5   Add DirectAccess ( $S_i, j$ ) to  $X_i$ ;
6   while  $j < |S_i|$  do
7      $j \leftarrow j + 1 + \text{Geometric}(p_i^+)$ ;
8     if  $j \leq |S_i|$  then Add DirectAccess ( $S_i, j$ ) to  $X_i$ ;
9   foreach  $e \in X_i$  do Remove  $e$  from  $X_i$  with probability  $1 - \frac{p_i(e)}{p_i^+}$ ;
10 return  $\bigcup_{i \in I} X_i$ ;

```

---

A naive approach would be to construct the index from Section 2.2 for each sub-instance and invoke the query procedure for each one. However, this would take  $\Omega(m)$  time per query, as we must spend at least  $O(1)$  time to check every sub-instance. To improve efficiency, we construct a composite index that supports a two-stage sampling process:

- **Preprocessing phase:** First, for each sub-instance  $\Psi_i$ , we build the standard index for  $\langle S_i, p_i^+ \rangle$  as described in Lemma 2.1. Second, we construct a meta-index to identify which sub-instances are likely to contribute to the final sample. We define a probability function  $q : [m] \rightarrow [0, 1]$  where  $q(i) = 1 - (1 - p_i^+)^{|S_i|}$ . Note that  $q(i)$  is the probability that sampling from  $S_i$  with uniform probability  $p_i^+$  yields a non-empty set. We build an optimal subset sampling index (e.g., [9]) for the meta-instance  $\Psi_{\text{meta}} = \langle [m], q \rangle$ . The total preprocessing time is linear in  $|S|$ .
- **Query phase:** Algorithm 3 proceeds in two stages. First, it queries the meta-index to obtain a random set of indices  $I \subseteq [m]$ . Then, only for the selected indices  $i \in I$ , it invokes the query procedure on the sub-instance  $\Psi_i$  (using Algorithm 2 logic) and performs rejection sampling.

**Lemma 2.4** (Batched Subset Sampling Index). *Given a partitioned instance as defined above, there is an index that can be built in  $O(|S|)$  time and space. Using this index, Algorithm 3 returns a valid subset sample in  $O(1 + \sum_{i=1}^m |S_i| \cdot p_i^+)$  expected time.*

**PROOF OF LEMMA 2.4.** The preprocessing phase constructs two levels of indexes: standard subset sampling indexes for each sub-instance  $\Psi_i$  (taking  $O(\sum |S_i|) = O(|S|)$  time) and a meta-index for  $\Psi_{\text{meta}} = \langle [m], q \rangle$  (taking  $O(m) \leq O(|S|)$  time). The space complexity is clearly  $O(|S|)$ .

We now analyze the expected query time of Algorithm 3. The running time consists of two parts:

- (1) *Meta-sampling.* Using the meta-index, generating the set of indices  $I$  takes expected time  $O(1 + \mu_{\Psi_{\text{meta}}})$ , where  $\mu_{\Psi_{\text{meta}}} = \sum_{i=1}^m q(i)$ .
- (2) *Sub-sampling.* For each selected index  $i \in I$ , we generate an intermediate sample  $X'_i$  using the geometric jump procedure (simulating **uss-advanced**). The cost for a specific  $i$  is proportional to the size of this intermediate sample. The total expected cost is:  $\mathbf{E}[\text{Cost}_{\text{sub}}] = \sum_{i=1}^m q(i) \cdot O(\mathbf{E}[|X'_i| \mid X'_i \neq \emptyset])$ . Note that  $\mathbf{E}[|X'_i| \mid X'_i \neq \emptyset] = \frac{\mathbf{E}[|X'_i|]}{\Pr[X'_i \neq \emptyset]} = \frac{|S_i| p_i^+}{q(i)}$ . Substituting this back into the summation, we have:

$$\mathbf{E}[\text{Cost}_{\text{sub}}] = O\left(\sum_{i=1}^m q(i) \cdot \frac{|S_i| p_i^+}{q(i)}\right) = O\left(\sum_{i=1}^m |S_i| p_i^+\right).$$

Combining both parts, the total expected time is  $O(1 + \sum_{i=1}^m q(i) + \sum_{i=1}^m |S_i| p_i^+)$ . Since  $q(i) \leq |S_i| p_i^+$ , the total complexity simplifies to  $O(1 + \sum_{i=1}^m |S_i| \cdot p_i^+)$ .  $\square$

### 3 Indexed Subset Sampling over Joins

We start with Problem 1.2 for the join instance  $Q$  and the associated weight functions  $\{p_i\}_{i \in [k]}$ . We assume tuples in each relation  $R_i \in Q$  are in some arbitrary fixed ordering, and relations in  $Q$  are also in some arbitrary fixed ordering. All missing proofs of this section are given in Appendix B.

#### 3.1 First Attempt

So far, we first assume that a DirectAccess oracle is available for the input join instance  $Q$  (under some fixed ordering), such that it receives an integer  $i \in [|\text{Join}(Q)|]$ , and returns the  $i$ -th element of  $\text{Join}(Q)$ . Similar to before, the ordering can be arbitrary but must remain consistent across multiple invocations of this oracle. Let  $L = \lceil 2\rho^* \log N \rceil$ , where  $\rho^*$  is the fractional edge covering number of the schema graph  $G_Q$ . Note that  $2^L \geq |\text{Join}(Q)|$ .

**Step 1: Partition input relations.** For each relation  $R_i \in Q$ , we partition tuples into  $L + 1$  sub-relations  $R_i^{(0)}, R_i^{(1)}, \dots, R_i^{(L)}$  based on their probabilities, where  $R_i^{(j)} = \{\mathbf{u} \in R_i \mid 2^{-j-1} < p_i(\mathbf{u}) \leq 2^{-j}\}$  for  $j \in \llbracket L-1 \rrbracket$  and  $R_i^{(L)} = \{\mathbf{u} \in R_i \mid p_i(\mathbf{u}) \leq 2^{-L}\}$ . Then, each combination  $\mathbf{j} = (j_1, j_2, \dots, j_k) \in \llbracket L \rrbracket^k$  defines a sub-join of  $Q$ , denoted as:  $Q_{\mathbf{j}} = \{R_1^{(j_1)}, R_2^{(j_2)}, \dots, R_k^{(j_k)}\}$ . The join results of  $Q_{\mathbf{j}}$  are a disjoint partition of the join results of  $Q$ , i.e.,  $\text{Join}(Q) = \bigsqcup_{\mathbf{j} \in \llbracket L \rrbracket^k} \text{Join}(Q_{\mathbf{j}})$ . To answer a query to the subset sampling instance  $\Psi = \langle \text{Join}(Q), p \rangle$ , we can return the disjoint union of the subsets sampled from all sub-instances  $\Psi_{\mathbf{j}} = \langle \text{Join}(Q_{\mathbf{j}}), p \rangle$  for  $\mathbf{j} \in \llbracket L \rrbracket^k$ . There are  $(L + 1)^k$  sub-instances in total.

**Step 2: Apply subset sampling to all sub-instances.** We first point out the following observation on each sub-instance in such a partition:

**Lemma 3.1** (Either Light or Near-uniform Sub-instance). *Let  $\mathbf{j} = (j_1, j_2, \dots, j_k)$ . If  $\sum_{i \in [k]} j_i \geq L$ , the instance  $\Psi_{\mathbf{j}}$  is light, and otherwise, it is  $2^k$ -uniform.*

**PROOF OF LEMMA 3.1.** For the largest probability, we always have  $\max_{\mathbf{u} \in \text{Join}(Q_{\mathbf{j}})} p(\mathbf{u}) \leq 2^{-\sum_{i \in [k]} j_i}$ . If  $\sum_{i \in [k]} j_i \geq L$ , we have  $\max_{\mathbf{u} \in \text{Join}(Q_{\mathbf{j}})} p(\mathbf{u}) \leq 2^{-\sum_{i \in [k]} j_i} \leq 2^{-L} \leq 1/|\text{Join}(Q)|$ , following our assumption on  $L$ . Hence,  $\Psi_{\mathbf{j}}$  is light. Otherwise, we must have  $L \notin \{j_1, j_2, \dots, j_k\}$ . In this case, for the smallest probability, we have  $2^{-k - \sum_{i \in [k]} j_i} \leq \min_{\mathbf{u} \in \text{Join}(Q_{\mathbf{j}})} p(\mathbf{u})$ . Hence,  $\Psi_{\mathbf{j}}$  is  $2^k$ -uniform since  $\max_{\mathbf{u} \in \text{Join}(Q_{\mathbf{j}})} p(\mathbf{u}) \leq \min_{\mathbf{u} \in \text{Join}(Q_{\mathbf{j}})} p(\mathbf{u})$ .  $\square$

To efficiently sample from these sub-instances, we adopt the composite index strategy from Section 2.3. We treat the sub-joins  $\{Q_{\mathbf{j}}\}_{\mathbf{j} \in \llbracket L \rrbracket^k}$  as the disjoint sub-instances. We define the meta-probability  $q(\mathbf{j}) = 1 - (1 - 2^{-\sum_{i \in [k]} j_i})^{|\text{Join}(Q_{\mathbf{j}})|}$  and construct the meta-index for the instance  $\langle \llbracket L \rrbracket^k, q \rangle$ . This allows us to skip empty or non-selected sub-instances efficiently. To put it formally:

- **Preprocessing phase:** We partition input tuples by weights. For each sub-instance  $Q_{\mathbf{j}}$ , we build a DirectAccess oracle (the sub-index) and compute the join size (e.g., [52]). Additionally, we build the meta-index for  $q$  as described in Section 2.3.
- **Query phase:** We invoke Algorithm 3 using the constructed meta-index and sub-indexes. The algorithm first samples a set of indices  $\mathbf{I}$  using the meta-index, and then retrieves tuples from the selected sub-instances  $Q_{\mathbf{j}}$  via their DirectAccess oracles.

**Theorem 3.2.** *Given an acyclic join instance  $Q$  consisting of  $k$  relations, and a set of associated weight functions  $\{p_j\}_{R_j \in Q}$ , there is an index  $\mathcal{D}$  that can be built within  $O(N \log^{k-1} N)$  space and time, using which each subset sampling query can be answered in  $O(1 + \mu_{\Psi} \log N)$  expected time.*

PROOF OF THEOREM 3.2. The preprocessing phase constructs the index defined in Lemma 2.4.

- (1) *Partitioning and Sub-indexes.* We partition each relation  $R_i \in Q$  into  $L + 1$  sub-relations, defining  $(L + 1)^k$  sub-instances  $Q_j$ . For each sub-instance, we construct a DirectAccess oracle (sub-index). According to [52], building a DirectAccess oracle for an acyclic join takes time linear in its input size. Since each tuple participates in  $(L + 1)^{k-1}$  sub-instances and  $L = O(\log N)$ , the total time and space are  $O(N \log^{k-1} N)$ .
- (2) *Meta-index.* We build the meta-index for  $\Psi' = \langle \llbracket L \rrbracket^k, q \rangle$ . As the universe size is  $(L + 1)^k = O(\log^k N)$ , this index can be built in  $O(\log^k N)$  time and space, which is subsumed by the cost of building the sub-indexes.

The query proceeds according to Algorithm 3.

- (1) *Meta-sampling.* We query the meta-index to obtain indices  $I$ . The expected time is  $O(1 + \mu_{\Psi'})$ , where  $\mu_{\Psi'} = \sum_j q(j)$ . As shown in the main text,  $\mu_{\Psi'} \leq \mu_{\Psi} + 1$ .
- (2) *Sub-sampling.* For each  $j \in I$ , we draw samples from  $Q_j$ . Since each sub-instance is either light or  $2^k$ -uniform, the expected cost per selected instance is proportional to the number of samples times the DirectAccess access cost  $O(\log N)$ . Following the analysis in Lemma 2.4, the total expected cost is  $O((1 + \mu_{\Psi}) \log N)$ .

Summing these gives the total expected query time  $O(1 + \mu_{\Psi} \log N)$ .  $\square$

Next, we improve upon this initial result by showing an optimized index with much smaller space and preprocessing time, while maintaining the same query time.

### 3.2 Optimized Index – Framework

The framework in Section 3.1 establishes feasibility but is not space-efficient because it treats every combination of weight buckets  $j \in \llbracket L \rrbracket^k$  as a distinct sub-instance. We now present an optimized index that reduces space usage to  $O(N \log N)$ .

**Merging Sub-instances by Score.** Recall from Lemma 3.1 that the classification of a sub-instance  $Q_j$  (as either light or near-uniform) depends solely on the sum of indices  $\sum_{i=1}^k j_i$ . This observation allows us to merge sub-instances. Instead of maintaining  $(L + 1)^k$  disjoint partitions, we group join results based on their total weight “score.”

For any tuple  $\mathbf{u}$  in a relation  $R_i$ , let its *score* be  $\phi(\mathbf{u}) = \lfloor -\log p_i(\mathbf{u}) \rfloor$ . For a join result  $\mathbf{u} \in \text{Join}(Q)$ , its *score* is the sum of component scores:  $\bar{\phi}(\mathbf{u}) = \sum_{i=1}^k \phi(\mathbf{u}[\text{schema}(R_i)])$ . We partition the join results  $\text{Join}(Q)$  into disjoint buckets based on this score. Let  $\mathcal{B}_\ell = \{\mathbf{u} \in \text{Join}(Q) \mid \bar{\phi}(\mathbf{u}) = \ell\}$  denote the set of join results with score  $\ell$ . We treat these buckets as the disjoint sub-instances required by the batched framework in Section 2.3. Specifically:

- (1) **Buckets  $\ell < L$ :** Each bucket  $\mathcal{B}_\ell$  contains join results with probabilities in the range  $(2^{-\ell-1}, 2^{-\ell}]$ . Thus,  $\mathcal{B}_\ell$  forms a 2-uniform sub-instance with upper bound  $p_\ell^+ = 2^{-\ell}$ .
- (2) **Tail Bucket  $\ell \geq L$ :** We merge all join results with score  $\ell \geq L$  into a single tail bucket  $\mathcal{B}_{\geq L} = \bigsqcup_{\ell \geq L} \mathcal{B}_\ell$ . Since any  $\mathbf{u} \in \mathcal{B}_{\geq L}$  has  $p(\mathbf{u}) \leq 2^{-L} \leq 1/|\text{Join}(Q)|$ , this combined bucket is a *light* sub-instance (Lemma 2.3) with upper bound  $p_{\geq L}^+ = 2^{-L}$ .

**Algorithm Overview.** We apply Algorithm 3 to these  $L + 1$  sub-instances.

- **Preprocessing phase:** We compute the size  $|\mathcal{B}_\ell|$  for each  $\ell < L$  to build the meta-index. For the tail bucket  $\mathcal{B}_{\geq L}$ , we do not maintain an exact count or index; we simply assign it a weight proxy of  $N^{\rho^*}$  (an upper bound on join size) for the meta-index, or handle it via a fallback mechanism since its selection probability is negligible.
- **Query phase:** We invoke Algorithm 3 with inputs  $\{\langle \mathcal{B}_\ell, p \rangle \mid \ell < L\}$  and the tail bucket.

- If the meta-sampling selects a bucket  $\mathcal{B}_\ell$  with  $\ell < L$ , we use a specialized DirectAccess oracle (described below) to retrieve a given index in  $\mathcal{B}_\ell$  in  $O(\log N)$  time.
- If the tail bucket  $\mathcal{B}_{\geq L}$  is selected (which happens with very low probability), we materialize the necessary results from  $\mathcal{B}_{\geq L}$  on the fly to support the retrieval.

The challenge is to efficiently support DirectAccess for the buckets  $\mathcal{B}_\ell$  ( $\ell < L$ ) without materializing them. We next describe how this can be achieved.

**Join Tree with Notations.** As mentioned, any acyclic join  $Q$  has a join tree  $\mathcal{T}$  such that (1) there is a one-to-one correspondence between relations in  $Q$  and nodes in  $\mathcal{T}$ , and (2) for each attribute  $x$ , the set of nodes containing  $x$  forms a connected subtree. With a slight abuse of notation, we also use  $R_i$  to denote the corresponding node in  $\mathcal{T}$  for relation  $R_i \in Q$ . For simplicity, we assume the nodes in  $\mathcal{T}$  are ordered by the in-order traversal. For node  $R_i$ , let  $C_i$  be the child nodes of  $R_i$  if  $R_i$  is not a leaf. For a node  $R_j \in C_i$ , let  $R_{\text{next}(j)}$  be the immediately next child node after  $R_j$  in  $C_i$ . For the largest child node  $R_{j^*}$ , we set  $\text{next}(j^*) = \text{null}$ . For node  $R_i$ , we use  $\text{parent}(i)$  to denote the (unique) parent node of  $R_i$  if  $R_i$  is not the root. Let  $\text{key}(i) = \text{schema}(R_i) \cap \text{schema}(R_{\text{parent}(i)})$  be the (common) join attributes between relation  $R_i$  and its parent if  $R_i$  is not the root. For the root node  $r$ ,  $\text{key}(r) = \emptyset$ . For a node  $R_i$ , let  $\mathcal{T}_i$  be the subtree rooted at node  $R_i$ . For an arbitrary child  $R_j \in C_i$ , we define  $\mathcal{T}_i^j$  to be the subtree of  $\mathcal{T}_i$  that excludes any subtrees rooted at the child nodes in  $C_i$  coming before  $R_j$ , i.e.,  $\mathcal{T}_i^j = \mathcal{T}_i \setminus \cup_{R_{j'} \in C_i: R_{j'} \prec R_j} \mathcal{T}_{j'}$ . For completeness, we also define  $\mathcal{T}_i^0 = \mathcal{T}_i$ .

**Data Structure.** We first remove all dangling tuples from the database instance. For each relation  $R_i \in Q$ , we store input tuples in a hash table, so that for any subset of attributes  $U \subseteq \text{schema}(R_i)$  and a tuple  $t \in \text{dom}(U)$ , we can get the list of tuples  $R_i \bowtie t$  in  $O(1)$  time. We take an arbitrary join tree  $\mathcal{T}$  for  $Q$  and store the following statistics to augment  $\mathcal{T}$ .

**W-values.** Consider any internal node  $R_i \in Q$ . For each tuple  $\mathbf{u} \in R_i$ , we store  $L$  counters  $W_{i,\mathbf{u}}^j(\ell)$  with  $\ell \in \llbracket L-1 \rrbracket$  specifically for each child node  $R_j \in C_i$ , where  $W_{i,\mathbf{u}}^j(\ell)$  stores the number of join results produced by relations in the subtree  $\mathcal{T}_i^j$  that is also participated by  $\mathbf{u}$ , with score  $\ell$ , i.e.,

$$W_{i,\mathbf{u}}^j(\ell) = \left| \left\{ t \in (\bowtie_{R_h \in \mathcal{T}_i^j} R_h) \bowtie \mathbf{u} : \sum_{R_h \in \mathcal{T}_i^j} \phi(t[\text{schema}(R_h)]) = \ell \right\} \right|. \quad (3)$$

Furthermore, we define and store  $W_{i,\mathbf{u}}^0(\ell)$ , the number of join results produced by relations in the subtree  $\mathcal{T}_i$  that is also participated by  $\mathbf{u}$ , with score  $\ell$ , by replacing  $\mathcal{T}_i^j$  in (3) with  $\mathcal{T}_i^0$ . For completeness, for the largest child node  $R_{j^*} \in C_i$ , we also define  $W_{i,\mathbf{u}}^{\text{next}(j^*)}(0) = 1$  and  $W_{i,\mathbf{u}}^{\text{next}(j^*)}(\ell) = 0$  for every  $\ell \in [L]$ . Moreover, for each score  $\ell \in \llbracket L-1 \rrbracket$ , and each child node  $R_j \in C_i$ , we build a prefix-sum array on  $\langle W_{i,\mathbf{u}}^j(\ell) : \mathbf{u} \in R_i \rangle$  under the pre-determined ordering of tuples in  $R_i$ . For each leaf node  $R_i \in Q$ , we only store the counter  $W_{i,\mathbf{u}}^0(\ell)$  for each tuple  $\mathbf{u} \in R_i$ .

**M-values.** For a non-root node  $R_i$ , we compute the projection of  $R_i$  onto  $\text{key}(i)$ , as  $R_i[\text{key}(i)]$ . For each tuple  $\mathbf{v} \in R_i[\text{key}(i)]$ , we also compute  $M_{i,\mathbf{v}}(0), M_{i,\mathbf{v}}(1), \dots, M_{i,\mathbf{v}}(L-1)$ : for  $\ell \in \llbracket L-1 \rrbracket$ ,  $M_{i,\mathbf{v}}(\ell)$  stores the sum of counters  $W_{i,\mathbf{u}}^0(\ell)$  over all tuples  $\mathbf{u} \in R_i \bowtie \mathbf{v}$ , i.e.,

$$M_{i,\mathbf{v}}(\ell) = \sum_{\mathbf{u} \in R_i \bowtie \mathbf{v}} W_{i,\mathbf{u}}^0(\ell). \quad (4)$$

### 3.3 Optimized Index – Preprocessing

We will show how to compute the statistics to build the desired index. The complete pseudocode is provided in Algorithm 4. Below, we show the intuition in detail. Our preprocessing phase performs computation in a bottom-up way. We note that  $M$ -values are defined on top of  $W$ -values, hence for

each node  $R_i \in Q$ ,  $M$ -values will be computed according to (4) straightforwardly once  $W$ -values are well computed. Below, we focus on computing  $W$ -values.

Consider a leaf node  $R_i$ . By definition of (3), we have  $W_{i,\mathbf{u}}^0(\phi(\mathbf{u})) = 1$  and  $W_{i,\mathbf{u}}^0(\ell) = 0$  for any  $\ell \in \llbracket L-1 \rrbracket \setminus \{\phi(\mathbf{u})\}$ . Consider an internal node  $R_i$ . Suppose  $W$ -values and  $M$ -values are well defined for each child  $R_j \in C_i$ . We next compute  $W$ -values for each tuple  $\mathbf{u} \in R_i$  and each child node  $R_j \in C_i$ . More specifically, we compute  $W_{i,\mathbf{u}}^j$  in a decreasing ordering of nodes in  $C_i$ . Recall that for the largest child node  $R_{j^*} \in C_i$ ,  $W_{i,\mathbf{u}}^{j^*}(0) = 1$  and  $W_{i,\mathbf{u}}^{j^*}(\ell) = 0$  for  $\ell \in \llbracket L-1 \rrbracket$ . We next compute  $W_{i,\mathbf{u}}^j(\cdot)$ , assuming  $W_{i,\mathbf{u}}^{\text{next}(j)}(\cdot)$  has been computed. We distinguish the following two cases on the score  $\ell$  (recall that  $\ell \in \llbracket L-1 \rrbracket$ ):

- **Case 1:**  $\ell < \phi(\mathbf{u})$ . In this case,  $W_{i,\mathbf{u}}^j(\ell) = 0$  since any join result participated by  $\mathbf{u}$  has its score at least  $\phi(\mathbf{u})$ ;
- **Case 2:**  $\phi(\mathbf{u}) \leq \ell < L$ . First, the Cartesian product of join result in each subtree  $\mathcal{T}_{j'}$  rooted at any child node  $R_{j'} \in C_i$  coming no earlier than  $R_j$ , that can be joined with  $\mathbf{u}$ , is essentially the join result in  $\mathcal{T}_i^j$  that can be joined with  $\mathbf{u}$ . We also need to take the scores into consideration—all possible combinations of  $\langle \ell_{j'} \in \llbracket L-1 \rrbracket : R_{j'} \in C_i, R_j \preceq R_{j'} \rangle$  with  $\sum_{R_{j'} \in C_i} \ell_{j'} = \ell - \phi(\mathbf{u})$ . Each combination will contribute  $\prod_{R_{j'} \in C_i: R_j \preceq R_{j'}} M_{j',\mathbf{u}[\text{key}(j')]}(\ell_{j'})$  to  $W_{i,\mathbf{u}}^j(\ell)$ . However, the number of such combinations is as large as  $O(L^k)$ , which translates to  $O(\log^k N)$  time even for computing a single  $W$ -value. A critical observation is that the products of  $M$ -values over  $R_{j'}$  with  $R_j \prec R_{j'}$  is actually captured by  $W_{i,\mathbf{u}}^{\text{next}(j)}$ , where  $R_{\text{next}(j)}$  is the immediately next node of  $R_j$  in  $C_i$ . Putting these observations together, we obtain:

$$\begin{aligned}
W_{i,\mathbf{u}}^j(\ell) &= \sum_{\substack{\langle \ell_{j'} \in \llbracket L-1 \rrbracket : R_{j'} \in C_i, R_j \preceq R_{j'} \rangle \\ \ell - \phi(\mathbf{u}) = \sum_{R_{j'} \in C_i} \ell_{j'}} \prod_{R_{j'} \in C_i: R_j \preceq R_{j'}} M_{j',\mathbf{u}[\text{key}(j')]}(\ell_{j'}) \\
&= \sum_{\substack{(\ell_1, \ell_2) \in \llbracket \ell - \phi(\mathbf{u}) \rrbracket : \\ \ell_1 + \ell_2 = \ell - \phi(\mathbf{u})}} M_{j,\mathbf{u}[\text{key}(j)]}(\ell_1) \cdot \sum_{\substack{\langle \ell_{j'} \in \llbracket L-1 \rrbracket : R_{j'} \in C_i, R_j \prec R_{j'} \rangle \\ \ell_2 = \sum_{R_{j'} \in C_i} \ell_{j'}} \prod_{R_{j'} \in C_i: R_j \prec R_{j'}} M_{j',\mathbf{u}[\text{key}(j')]}(\ell_{j'}) \\
&= \sum_{\substack{(\ell_1, \ell_2) \in \llbracket \ell - \phi(\mathbf{u}) \rrbracket : \\ \ell_1 + \ell_2 = \ell - \phi(\mathbf{u})}} M_{j,\mathbf{u}[\text{key}(j)]}(\ell_1) \cdot W_{i,\mathbf{u}}^{\text{next}(j)}(\ell_2). \tag{5}
\end{aligned}$$

Hence  $W_{i,\mathbf{u}}^j(\ell)$  can be computed by (5). At that point, the  $M$ -values for the child node  $R_j$  and the  $W$ -values for the subsequent node  $R_{\text{next}(j)}$  have already been computed. Note that  $W_{i,\mathbf{u}}^0(\ell)$  can be computed similarly by (5):

$$W_{i,\mathbf{u}}^0(\ell) = \sum_{\substack{(\ell_1, \ell_2) \in \llbracket \ell - \phi(\mathbf{u}) \rrbracket : \\ \ell_1 + \ell_2 = \ell - \phi(\mathbf{u})}} M_{j^\circ, \mathbf{u}[\text{key}(j^\circ)]}(\ell_1) \cdot W_{i,\mathbf{u}}^{\text{next}(j^\circ)}(\ell_2),$$

where  $R_{j^\circ}$  is the first node in  $C_i$ .

To speed up the process, let's fix a node  $R_i \in Q$ , a tuple  $\mathbf{u} \in R_i$ , and a child node  $R_j \in C_i$  as an example. Now, we need to compute  $W_{i,\mathbf{u}}^0(\cdot)$  for  $L$  different scores. There is a convolution structure between  $W_{i,\mathbf{u}}^j(\cdot)$ ,  $M_{j,\mathbf{u}[\text{key}(j)]}(\cdot)$  and  $W_{i,\mathbf{u}}^{\text{next}(j)}(\cdot)$ , so we apply the Fast Fourier Transform for computing these  $L$  scores together. As we will show in the analysis of the preprocessing step, this observation improves the preprocessing time from  $O(NL^2)$  to  $O(NL \log L)$  (Lemma 3.3).

**Algorithm 4: Preprocess( $Q, \mathcal{T}$ )**


---

**Input:** Join tree  $\mathcal{T}$  for  $Q$ , Relations  $\{R_i\}_{R_i \in Q}$ .

```

1 foreach node  $R_i \in \mathcal{T}$  in bottom-up order do
2   if  $R_i$  is a leaf node then
3     foreach tuple  $\mathbf{u} \in R_i$  do
4       Set  $W_{i,\mathbf{u}}^0(\phi(\mathbf{u})) \leftarrow 1$  and  $W_{i,\mathbf{u}}^0(\ell) \leftarrow 0$  for  $\ell \neq \phi(\mathbf{u})$ ;
5        $M_{i,\mathbf{u}[\text{key}(i)]} \leftarrow M_{i,\mathbf{u}[\text{key}(i)]} + W_{i,\mathbf{u}}^0$ ; // element-wise vector addition
6   else //  $R_i$  is an internal node
7     foreach tuple  $\mathbf{u} \in R_i$  do
8       Let  $R_{j^*}$  be the last child in  $C_i$ ;
9       Set  $W_{i,\mathbf{u}}^{\text{next}(j^*)}(0) \leftarrow 1$  and  $W_{i,\mathbf{u}}^{\text{next}(j^*)}(\ell) \leftarrow 0$  for  $\ell > 0$ ;
10      foreach child  $R_j \in C_i$  in decreasing order do
11        Let  $R_{\text{next}(j)}$  be the child immediately following  $R_j$ ;
12        foreach score  $\ell \in \llbracket L-1 \rrbracket$  do
13          if  $\ell < \phi(\mathbf{u})$  then  $W_{i,\mathbf{u}}^j(\ell) \leftarrow 0$ ;
14          else  $W_{i,\mathbf{u}}^j(\ell) \leftarrow \sum_{k=0}^{\ell-\phi(\mathbf{u})} M_{j,\mathbf{u}[\text{key}(j)]}(k) \cdot W_{i,\mathbf{u}}^{\text{next}(j)}(\ell - \phi(\mathbf{u}) - k)$ ;
15         $M_{i,\mathbf{u}[\text{key}(i)]} \leftarrow M_{i,\mathbf{u}[\text{key}(i)]} + W_{i,\mathbf{u}}^0$ ; // element-wise vector addition
16      foreach child  $R_j \in C_i$  do
17        foreach  $\mathbf{v} \in R_j[\text{key}(i)]$  and  $\ell \in \llbracket L-1 \rrbracket$  do
18          Build prefix-sum array on  $\langle W_{i,\mathbf{u}}^j(\ell) : \mathbf{u} \in R_i \bowtie \mathbf{v} \rangle$ ;
```

---

After all  $W$ -values are computed, for each internal node  $R_i \in Q$  with each child node  $R_j \in C_i$ , each value  $\mathbf{v} \in R_j[\text{key}(i)]$ , and each score  $\ell \in \llbracket L-1 \rrbracket$ , we build a prefix-sum array on  $\langle W_{i,\mathbf{u}}^j(\ell) : \mathbf{u} \in R_i \bowtie \mathbf{v} \rangle$  under the pre-determined ordering.

Lastly, we compute the bucket size  $|\mathcal{B}_\ell|$  for each score  $\ell \in \llbracket L-1 \rrbracket$ , which by definition can be rewritten as  $|\mathcal{B}_\ell| = \sum_{\mathbf{u} \in R_r} W_{r,\mathbf{u}}^0(\ell)$  for the root node  $r$ .

We briefly describe how to preprocess the data structures (see Appendix B for complete details). We compute  $W/M$  values in a bottom-up way.  $M$ -values will be computed according to (4) straightforwardly once  $W$ -values are well computed. Consider a leaf node  $R_i$ , we have  $W_{i,\mathbf{u}}^0(\phi(\mathbf{u})) = 1$  and  $W_{i,\mathbf{u}}^0(j) = 0$  for any  $j \in \llbracket L-1 \rrbracket \setminus \{\phi(\mathbf{u})\}$ . Consider an internal node  $R_i$ . Suppose  $W$ -values and  $M$ -values are well defined for each child  $R_j \in C_i$ . We compute  $W_{i,\mathbf{u}}^j$  in a decreasing ordering of nodes in  $C_i$  and follow a recursive structure. If  $\ell < \phi(\mathbf{u})$ ,  $W_{i,\mathbf{u}}^j(\ell) = 0$ , and otherwise,

$$W_{i,\mathbf{u}}^j(\ell) = \sum_{(\ell_1, \ell_2) \in \llbracket \ell - \phi(\mathbf{u}) \rrbracket : \ell_1 + \ell_2 = \ell - \phi(\mathbf{u})} M_{j,\mathbf{u}[\text{key}(j)]}(\ell_1) \cdot W_{i,\mathbf{u}}^{\text{next}(j)}(\ell_2). \quad (6)$$

Note that  $W_{i,\mathbf{u}}^0(\ell)$  can be computed similarly by resorting to the first child node. Importantly, the computation of (6) can be sped up by Fast Fourier Transform [19, 20], and the whole computation only takes  $O(NL \log L)$  time. Lastly, the bucket size  $|\mathcal{B}_\ell|$  for each score  $\ell \in \llbracket L-1 \rrbracket$  can be rewritten as  $|\mathcal{B}_\ell| = \sum_{\mathbf{u} \in R_r} W_{r,\mathbf{u}}^0(\ell)$  for the root node  $r$ .

**Lemma 3.3.** *Algorithm 4 runs  $O(N \log N \log \log N)$  time, and uses  $O(N \log N)$  space.*

**Algorithm 5: RecursiveAccess**( $i, j, v, \ell, \tau$ )

**Input:**  $i, j \in [k]$  such that  $R_j \in C_i$ , tuple  $v$ , score  $\ell \in \llbracket L-1 \rrbracket$ , integer  $\tau \in \mathbb{Z}^+$ .

**Output:** The  $\tau$ -th tuple in the join result of  $\mathcal{T}_i^j$ , with score  $\ell$ , that can be joined with  $v$ .

- 1  $\mathbf{u} \leftarrow$  the smallest tuple in  $R_i \times v$  such that  $\sum_{\mathbf{u}' \in R_i \times v: \mathbf{u}' \prec \mathbf{u}} W_{i, \mathbf{u}'}^j(\ell) \geq \tau$ ;
- 2 **if**  $R_i$  is a leaf node **then return**  $\mathbf{u}$ ;
- 3  $\tau \leftarrow \tau - \sum_{\mathbf{u}' \in R_i \times v: \mathbf{u}' \prec \mathbf{u}} W_{i, \mathbf{u}'}^j(\ell)$ ;
- 4  $\Phi \leftarrow \{(\ell_1, \ell_2) \in \llbracket \ell \rrbracket^2 : \ell_1 + \ell_2 = \ell - \phi(\mathbf{u})\}$  in lexicographical ordering by  $(\ell_1, \ell_2)$ ;
- 5  $(\ell_1, \ell_2) \leftarrow$  the smallest pair in  $\Phi$  that  $\sum_{(\ell'_1, \ell'_2) \in \Phi: (\ell'_1, \ell'_2) \preceq (\ell_1, \ell_2)} M_{j, \mathbf{u}[\text{key}(j)]}(\ell'_1) \cdot W_{i, \mathbf{u}}^{\text{next}(j)}(\ell'_2) \geq \tau$ ;
- 6  $\tau \leftarrow \tau - \sum_{(\ell'_1, \ell'_2) \in \Phi: (\ell'_1, \ell'_2) \prec (\ell_1, \ell_2)} M_{j, \mathbf{u}[\text{key}(j)]}(\ell'_1) \cdot W_{i, \mathbf{u}}^{\text{next}(j)}(\ell'_2)$ ;
- 7  $\tau_1 \leftarrow \left\lfloor \frac{\tau}{W_{i, \mathbf{u}}^{\text{next}(j)}(\ell_2)} \right\rfloor$ ,  $\tau_2 \leftarrow ((\tau - 1) \bmod W_{i, \mathbf{u}}^{\text{next}(j)}(\ell_2)) + 1$ ;
- 8  $\mathbf{u}_1 \leftarrow \text{RecursiveAccess}(j, \emptyset, \mathbf{u}, \ell_1, \tau_1)$ ;
- 9  $\mathbf{u}_2 \leftarrow \text{RecursiveAccess}(i, \text{next}(j), \mathbf{u}, \ell_2, \tau_2)$ ;
- 10 **return**  $\mathbf{u}_1 \bowtie \mathbf{u} \bowtie \mathbf{u}_2$ ;

PROOF OF LEMMA 3.3. For any relation  $R_i$ , and for every tuple  $\mathbf{u} \in R_i$ , we store  $O(|C_i| \cdot L) = O(L)$  different  $W$  values. So, in total, the space needed for storing all the  $W$  values is  $O(NL)$ . Moreover, for each  $v \in R_i[\text{key}(i)]$ , we store  $O(L)$  different  $M$  values. So, in total, the space needed for storing the  $M$  values is  $O(NL)$ . Hence, the overall space used is  $O(NL) = O(N \log N)$ .

To compute each  $W$  value, we use Equation (6). Calculating each  $W$  value by (6), needs a summation of  $O(L)$  different terms, leading to an  $O(NL^2)$  overall running time for calculating all the  $W$  values naively. To speed up the process, let's fix a node  $R_i \in Q$ , a tuple  $\mathbf{u} \in R_i$ , and a child node  $R_j \in C_i$  as an example. Now, we need to compute  $W_{i, \mathbf{u}}^j(\cdot)$  for  $L$  different scores. There is a convolution structure between  $W_{i, \mathbf{u}}^j(\cdot)$ ,  $M_{j, \mathbf{u}[\text{key}(j)]}(\cdot)$  and  $W_{i, \mathbf{u}}^{\text{next}(j)}(\cdot)$  in Equation (6), so we can apply the Fast Fourier Transform for computing these  $L$  scores together. This can improve the  $O(L^2)$  computations to  $O(L \log L)$ . More formally, consider the multiplication of two degree  $L$  polynomials:

$$p(x) = M_{j, \mathbf{u}[\text{key}(j)]}(L-1)x^{L-1} + M_{j, \mathbf{u}[\text{key}(j)]}(L-2)x^{L-2} + \dots + M_{j, \mathbf{u}[\text{key}(j)]}(0)x^0,$$

and,

$$q(x) = W_{i, \mathbf{u}}^{\text{next}(j)}(L-1)x^{L-1} + W_{i, \mathbf{u}}^{\text{next}(j)}(L-2)x^{L-2} + \dots + W_{i, \mathbf{u}}^{\text{next}(j)}(0)x^0.$$

By Equation (6), notice that for all  $\ell \in \llbracket L-1 \rrbracket$ , the value  $W_{i, \mathbf{u}}^j(\ell)$  is equal to the coefficient of  $x^{\ell - \phi(\mathbf{u})}$  in the polynomial  $p(x) \cdot q(x)$ . Using the Fast Fourier Transform, we can calculate all the coefficients of the degree  $2L$  polynomial  $p(x) \cdot q(x)$  in  $O(L \log L)$  time. Hence, for a fixed node  $R_i \in Q$ , a tuple  $\mathbf{u} \in R_i$ , and a child node  $R_j \in C_i$ , we can find all the  $L$  different  $W_{i, \mathbf{u}}^j(\cdot)$  values in  $O(L \log L)$  time. Putting everything together, calculating all the  $W$  values takes  $O(NL \log L) = O(N \log N \log \log N)$ . Having  $W$ -values, then calculating  $M$ -values is straightforward using Equation (4) in  $O(NL)$  total time. Therefore, the overall running time of the preprocessing step is  $O(N \log N \log \log N)$ .  $\square$

### 3.4 Optimized Index – Direct Access

Our goal is to retrieve the  $\tau$ -th tuple from bucket  $\mathcal{B}_\ell$  given score  $\ell \in \llbracket L-1 \rrbracket$  and rank  $\tau \in \llbracket |\mathcal{B}_\ell| \rrbracket$ . To ensure the “ $\tau$ -th” tuple is well-defined, we rely on the canonical lexicographic ordering of join

results. This ordering is naturally induced by the fixed traversal order of relations in the join tree  $\mathcal{T}$  and the total ordering of tuples within each relation  $R_i$ .

We present a general procedure **RecursiveAccess** $(i, j, v, \ell, \tau)$ , that returns the  $\tau$ -th tuple in the join results of the relations in the subtree  $\mathcal{T}_i^j$ , with score  $\ell$ , that can be joined with  $v$  on its support attributes. The original problem can be solved by invoking **RecursiveAccess** $(r, \emptyset, \emptyset, \ell, \tau)$ .

Algorithm 5 begins by identifying the specific tuple  $u \in R_i$  that participates in the target join result. Line 1 performs a binary search on the prefix-sum stored over  $W_{i,*}^j(\ell)$  to locate the first tuple  $u$  where the cumulative count of join results in the subtree  $\mathcal{T}_j^i$  (with score  $\ell$ ) covers the rank threshold  $\tau$ . Line 3 then subtracts the contribution of tuples preceding  $u$  from  $\tau$ . This reduces the problem to retrieving the (updated)  $\tau$ -th join result within the subtree  $\mathcal{T}_j^i$  specifically involving  $u$ . We then distinguish two cases. If  $R_i$  is a leaf node, we simply return  $u$ . In the general case, we decompose the retrieval into  $\mathcal{T}_j$  and  $\mathcal{T}_i^{\text{next}(j)}$  based on the fact that  $\mathcal{T}_i^j = \mathcal{T}_j \sqcup \mathcal{T}_i^{\text{next}(j)}$ . The objective is to determine the local indexes to query within these sub-problems. Recall that  $W_{i,u}^j(\ell)$  aggregates disjoint contributions from score pairs  $(\ell_1, \ell_2)$  (Equation 6). The algorithm iterates through these pairs in lexicographical order (Line 4) to find the specific pair  $(\ell_1, \ell_2)$  that covers the current rank threshold  $\tau$  (Line 5). After further adjusting  $\tau$  to account for preceding pairs (Line 6), we effectively map the remaining rank to indexes  $\tau_1$  and  $\tau_2$  within the Cartesian product of valid results from  $\mathcal{T}_j$  and  $\mathcal{T}_i^{\text{next}(j)}$  (Line 7). Finally, we invoke the procedure recursively on  $\mathcal{T}_j$  and  $\mathcal{T}_i^{\text{next}(j)}$  and combine the returned results.

**Correctness.** To establish the correctness of Algorithm 5, we show that for **DirectAccess** $(S_\ell, j)$  in Algorithm 3, Algorithm 5 returns the  $\tau$ -th tuple in  $\mathcal{B}_\ell$ , according to a fixed total ordering of the tuples in each bucket  $\mathcal{B}_\ell$ . Recall from Section 2 that the actual ordering of the tuples is irrelevant, as long as it is fixed. Since Algorithm 5 is deterministic, it suffices to prove that for any two distinct calls of **DirectAccess** $(S_{\hat{\ell}}, \hat{\tau})$  and **DirectAccess** $(S_{\bar{\ell}}, \bar{\tau})$  in Algorithm 3, Algorithm 5 returns two distinct tuples  $\hat{w}$  and  $\bar{w}$  from  $\text{Join}(Q)$ . This is exactly proved by Lemma B.1 in Appendix B.

**Complexity.** We will analyze its cost step by step. Line 1 performs a binary search on the prefix-sum tree for  $N$  tuples, which takes  $O(\log N)$  time. Line 3 can also be efficiently done through the prefix-sum tree in  $O(\log N)$  time, since any prefix sum can be decomposed into  $O(\log N)$  canonical nodes in this tree. Line 2 takes  $O(1)$  time. Lines 5 - 6 take  $O(L) = O(\log N)$  time, since there are at most  $L$  pairs in  $\Phi$  to be explored. Line 7 takes  $O(1)$  time. Hence, all lines before the recursion take  $O(\log N)$  time. As there are at most  $k^2$  different combinations of  $(i, j)$ , and the recursion is only invoked once for each combination of  $(i, j)$ , the total number of recursions is  $O(k^2)$ . As  $k$  is a constant, the total running time of Algorithm 5 is  $O(\log N)$ .

Now, we are ready to show the following theorem for our optimized index:

**Theorem 3.4.** *Given an acyclic join instance  $Q$  of  $k$  relations, and a set of associated weight functions  $\{p_j\}_{R_j \in Q}$ , there exists an index  $\mathcal{D}$  of size  $O(N \log N)$  that can be built in  $O(N \log N \log \log N)$  time, such that each subset sampling query can be answered in  $O(1 + \mu_\Psi \log N)$  expected time.*

**PROOF OF THEOREM 3.4.** The preprocessing cost follows from Lemma 3.3. We focus on the query answering time. In answering a subset sampling query, whenever some join result needs to be retrieved from  $\mathcal{B}_\ell$ , we distinguish the following two cases. If  $\ell = L$ , we evaluate the full join results in  $\mathcal{B}_{\geq L}$  (if these join results have not been computed before) to support the direct access. Otherwise, we call the **RecursiveAccess** procedure accordingly. For the first case, the expected cost is  $O(1)$  because there are at most  $O(N^{p^*})$  join results in  $\mathcal{B}_{\geq L}$  that takes  $O(N^{p^*})$  time to materialize [5], and the probability that at least a join result is retrieved from  $\mathcal{B}_{\geq L}$  is at most  $1 - (1 - p_L^+)^{|\mathcal{B}_{\geq L}|} \leq$

$1 - (1 - \frac{1}{N^{2\rho^*}})^{N\rho^*} \leq \frac{1}{N\rho^*}$ , since  $|\mathcal{B}_{\geq L}| \leq N\rho^*$ , and  $p_L^+ \leq \frac{1}{2L} \leq \frac{1}{N^{2\rho^*}}$ . Following Lemma 2.2 and the runtime of Algorithm 5, the expected query answering time is  $O(1 + \mu_\Psi \log N)$ .  $\square$

#### 4 One-shot Subset Sampling over Joins

We observe that the data structure from Theorem 3.4 provides a baseline solution for the one-shot problem. By constructing the index and issuing a single query, we can generate a subset sample in  $O(N \log N \log \log N + \mu_\Psi \log N)$  expected time. In this section, we present a specialized algorithm for the one-shot setting that runs in  $O(N \log^2 N + \mu_\Psi)$  expected time. This approach improves upon the baseline solution when the expected sample size is large (i.e.,  $\mu_\Psi \gg N$ ), as it eliminates the logarithmic factor associated with  $\mu_\Psi$ .

**High-level Idea.** We first keep all statistics computed in the preprocessing phase as in Section 3.2. We conceptually partition the join results into buckets. Recall that the high-level idea of answering one subset sampling query is to invoke Algorithm 3 on instances  $\{(\mathcal{B}_\ell, p) : \ell \in \llbracket L-1 \rrbracket\}$  with  $p_\ell^+ = 2^{-\ell}$ , and  $(\bigsqcup_{\ell \geq L} \mathcal{B}_\ell, p)$  with  $p_L^+ = 2^{-L}$ . Whenever a join result needs to be retrieved from  $\mathcal{B}_\ell$ , if  $\ell < L$ , we use the DirectAccess oracle to support the retrieval, and otherwise, we compute all the results in  $\bigsqcup_{\ell \geq L} \mathcal{B}_\ell$  to support the retrieval. In the first case, whenever the DirectAccess oracle on bucket  $\mathcal{B}_\ell$  with  $\ell < L$  is invoked, it pays  $O(\log N)$  time for each retrieval. To overcome this barrier in the query time, we need to compute additional statistics. Careful inspection of Algorithm 5 reveals that the  $O(\log N)$  cost comes from the binary search in Line 1 and the exhaustive search in Line 5. A natural idea is to precompute all these statistics. Consider an arbitrary relation  $R_i \in Q$  with an arbitrary child node  $R_j \in C_i$ . For every score  $\ell \in \llbracket L-1 \rrbracket$ , and every tuple  $\mathbf{v} \in R_i[\text{key}(i)]$ , we compute an array  $X_{i,j,\mathbf{v},\ell}$  of size  $|R_i \bowtie \mathbf{v}|$ , where for every tuple  $\mathbf{u} \in R_i \bowtie \mathbf{v}$ ,

$$X_{i,j,\mathbf{v},\ell}(\mathbf{u}) = \sum_{\mathbf{u}' \in R_i \bowtie \mathbf{v}, \mathbf{u}' \preceq \mathbf{u}} W_{i,\mathbf{u}'}^j(\ell).$$

Next, for every score  $\ell \in \llbracket L-1 \rrbracket$  and every tuple  $\mathbf{u} \in R_i$ , we compute an array  $Y_{i,j,\mathbf{u},\ell}$  of size at most  $L$ , where for every pair  $(\ell_1, \ell_2) \in \llbracket \ell \rrbracket^2$  with  $\ell_1 + \ell_2 = \ell - \phi(\mathbf{u})$ , in lexicographically sorted order,

$$Y_{i,j,\mathbf{u},\ell}(\ell_1, \ell_2) = \sum_{(\ell'_1, \ell'_2) \in \llbracket \ell \rrbracket : (\ell'_1, \ell'_2) \preceq (\ell_1, \ell_2), \ell'_1 + \ell'_2 = \ell - \phi(\mathbf{u})} M_{j,\mathbf{u}[\text{key}(j)]}(\ell'_1) \cdot W_{i,\mathbf{u}}^{\text{next}(j)}(\ell'_2).$$

We note that by constructing  $X$  and  $Y$  arrays, we do not define a new index for subset sampling over joins. Instead, these statistics are computed when answering one-shot subset sampling query.

Each invocation of Algorithm 5 is identified by a quintuple  $(i, j, \mathbf{v}, \ell, \tau)$  as its input. The high-level idea is to maintain the list of quintuples  $(i, j, \mathbf{v}, \ell, \tau)$  on which we will invoke Algorithm 5, and run multiple invocations of Algorithm 5 together by merging certain operations. We group these quintuples by  $(i, j, \mathbf{v}, \ell)$ , where quintuples inside each group only have different  $\tau$ -values sorted by  $\tau$  in increasing order. We maintain the set of quintuples  $\mathcal{P}_{i,j}$  we should satisfy in each call of Algorithm 5. For a node  $R_i$  and  $R_j \in C_i \cup \{\emptyset\}$ ,  $\mathcal{P}_{i,j}$  is defined as the parameter set of all invocations of Algorithm 5 on  $(i, j)$ . Initially, for each request DirectAccess  $(S_\ell, \tau)$  issued in Algorithm 3, we add a quintuple  $(r, \emptyset, \emptyset, \ell, \tau)$  to  $\mathcal{P}_{r,\emptyset}$ , where  $R_r$  is the root node. Below, we describe a general procedure **BatchRecursiveAccess** $(i, j, \mathcal{P}_{i,j})$ . The new algorithm executes the same steps as in Algorithm 5, but it computes all smallest tuples  $\mathbf{u}$  and pairs  $(\ell_1, \ell_2)$  for all quintuples simultaneously.

**BatchRecursiveAccess** $(i, j, \mathcal{P}_{i,j})$ : This procedure takes as input a node  $R_i \in Q$  with a child node  $R_j \in C_i$ , and the set of quintuples  $\mathcal{P}_{i,j}$ . The output will be exactly the collection of tuples returned by Algorithm 5 for each invocation parameterized by quintuple  $(i, j, \mathbf{v}, \ell, \tau) \in \mathcal{P}_{i,j}$ . The pseudocode is shown in Algorithm 6.

We first sort  $\mathcal{P}_{i,j}$  based on  $\tau$  using the Radix sort algorithm [20]. Since  $\tau$  is polynomially bounded in  $N$ , using radix sort saves a factor of  $\log N$  in the sorting time. Then we group  $\mathcal{P}_{i,j}$  by the

**Algorithm 6: BatchRecursiveAccess**( $i, j, \mathcal{P}_{i,j}$ )

---

```

1  $\mathcal{U} \leftarrow \emptyset$ ;
2 Sort  $\mathcal{P}_{i,j}$  by  $\tau$  using Radix sort;
3 Group  $\mathcal{P}_{i,j}$  further by  $(i, j, v, \ell)$  creating groups  $\mathcal{P}_{i,j,v,\ell}$  using the sorted order by  $\mathcal{P}_{i,j}$ ;
4 for  $\ell \in \llbracket L - 1 \rrbracket$  do
5   for each  $(i, j, v, \ell, \tau) \in \mathcal{P}_{i,j,v,\ell}$  in sorted order do
6      $\mathbf{u} \leftarrow$  first tuple in  $X_{i,j,v,\ell}$ ;
7     while  $X_{i,j,v,\ell}(\mathbf{u}) < \tau$  do  $\mathbf{u} \leftarrow$  next tuple in  $X_{i,j,v,\ell}$ ;
8     if  $R_i$  is a leaf node then Add  $\mathbf{u}$  as  $\vec{\mathbf{u}}(i, j, v, \ell, \tau)$  to  $\mathcal{U}$ ;
9      $\mathcal{P}_{i,j} \leftarrow (\mathcal{P}_{i,j} - \{(i, j, v, \ell, \tau)\}) \cup \{(i, j, \mathbf{u}, \ell, \tau - X_{i,j,v,\ell}(\text{prev}(\mathbf{u}))\}$ ;
10 if  $R_i$  is a leaf node then return  $\mathcal{U}$ ;
11 Sort  $\mathcal{P}_{i,j}$  by  $\tau$  using Radix sort;
12 Group  $\mathcal{P}_{i,j}$  further by  $(i, j, v, \ell)$  creating groups  $\mathcal{P}_{i,j,v,\ell}$  using the sorted order by  $\mathcal{P}_{i,j}$ ;
13  $\mathcal{P}_{j,\emptyset} \leftarrow \emptyset, \mathcal{P}_{i,\text{next}(j)} \leftarrow \emptyset$ ;
14 for  $\ell \in \llbracket L - 1 \rrbracket$  do
15   for each  $(i, j, \mathbf{u}, \ell, \tau) \in \mathcal{P}_{i,j,\mathbf{u},\ell}$  in sorted order do
16      $(\ell_1, \ell_2) \leftarrow$  first pair in  $Y_{i,j,\mathbf{u},\ell}$ ;
17     while  $Y_{i,j,\mathbf{u},\ell}((\ell_1, \ell_2)) < \tau$  do  $(\ell_1, \ell_2) \leftarrow$  the next pair in  $Y_{i,j,\mathbf{u},\ell}$ ;
18     Compute  $\tau_1, \tau_2$  based on  $\ell_1, \ell_2$  according to Line 7 of Algorithm 5;
19      $\mathcal{P}_{j,\emptyset} \leftarrow \mathcal{P}_{j,\emptyset} \cup \{(j, \emptyset, \mathbf{u}, \ell_1, \tau_1)\}$ ;
20      $\mathcal{P}_{i,\text{next}(j)} \leftarrow \mathcal{P}_{i,\text{next}(j)} \cup \{(i, \text{next}(j), \mathbf{u}, \ell_2, \tau_2)\}$ ;
21  $\mathcal{U} \leftarrow \mathcal{U} \cup \text{BatchRecursiveAccess}(j, \emptyset, \mathcal{P}_{j,\emptyset})$ ;
22  $\mathcal{U} \leftarrow \mathcal{U} \cup \text{BatchRecursiveAccess}(i, \text{next}(j), \mathcal{P}_{i,\text{next}(j)})$ ;
23 for  $(i, j, \mathbf{u}, \ell, \tau) \in \mathcal{P}_{i,j}$  do
24    $\vec{\mathbf{u}}(i, j, \mathbf{u}, \ell, \tau) \leftarrow \vec{\mathbf{u}}(j, \emptyset, \mathbf{u}, \ell_1, \tau_1) \bowtie \mathbf{u} \bowtie \vec{\mathbf{u}}(i, \text{next}(j), \mathbf{u}, \ell_2, \tau_2)$ ;
25   //  $\ell_1, \ell_2, \tau_1, \tau_2$  are computed in Lines 17 and 18 above
26   Add  $\vec{\mathbf{u}}(i, j, \mathbf{u}, \ell, \tau)$  to  $\mathcal{U}$ ;
27 return  $\mathcal{U}$ ;

```

---

different  $(i, j, v, \ell)$  creating the groups  $\mathcal{P}_{i,j,v,\ell}$ . We notice that all quintuples in each  $\mathcal{P}_{i,j,v,\ell}$  are sorted by  $\tau$  using the sorted ordering of  $\mathcal{P}_{i,j}$ . The goal is to compute the smallest tuple  $\mathbf{u}$  for each quintuple, as in Line 1 of Algorithm 5. For every score  $\ell \in \llbracket L - 1 \rrbracket$ , we consider every quintuple  $(i, j, v, \ell, \tau) \in \mathcal{P}_{i,j,v,\ell}$  following the sorted order. We traverse table  $X_{i,j,v,\ell}$  until we get the first tuple  $\mathbf{u}$  such that  $X_{i,j,v,\ell}(\mathbf{u}) \geq \tau$ . We also update the quintuple  $(i, j, v, \ell, \tau)$  to  $(i, j, v, \ell, \tau - X_{i,j,v,\ell}(\text{prev}(\mathbf{u})))$  as in Line 3 of Algorithm 5. If  $R_i$  is a leaf node we return all smallest tuples  $\mathbf{u}$  we computed from every quintuple in  $\mathcal{P}_{i,j}$ . For later reference, we denote the returned tuple  $\mathbf{u}$  for the quintuple  $(i, j, v, \ell, \tau)$  as  $\vec{\mathbf{u}}(i, j, v, \ell, \tau)$ .

Then we proceed in a similar manner to compute the pairs  $(\ell_1, \ell_2)$ . We sort (the updated)  $\mathcal{P}_{i,j}$  by  $\tau$  using the Radix sort algorithm. Then we group  $\mathcal{P}_{i,j}$  by the different  $(i, j, v, \ell)$  creating the groups  $\mathcal{P}_{i,j,v,\ell}$ . The goal is to compute the smallest pair  $(\ell_1, \ell_2)$  for each quintuple, as in Line 5 of Algorithm 5. For every score  $\ell \in \llbracket L - 1 \rrbracket$  we consider every quintuple  $(i, j, v, \ell, \tau) \in \mathcal{P}_{i,j,v,\ell}$  following the sorted order. We traverse table  $Y_{i,j,v,\ell}$  until we get the first pair  $(\ell_1, \ell_2)$  such that  $Y_{i,j,v,\ell}(\ell_1, \ell_2) \geq \tau$ . We compute  $\tau_1, \tau_2$  as in Line 7 of Algorithm 5. Furthermore, we add the quintuple  $(j, \emptyset, \mathbf{u}, \ell_1, \tau_1)$  to

the (initially empty) set  $\mathcal{P}_{j,\emptyset}$  and add the quintuple  $(i, \text{next}(j), \mathbf{u}, \ell_2, \tau_2)$  in the (initially empty) set  $\mathcal{P}_{i,\text{next}(j)}$ . Notice that these sets are used as input to the next calls **BatchRecursiveAccess** $(j, \emptyset, \mathcal{P}_{j,\emptyset})$  and **BatchRecursiveAccess** $(i, \text{next}(j), \mathcal{P}_{i,\text{next}(j)})$  as in Lines 8 and 9 of Algorithm 5.

Now, we will generate tuples to be returned for quintuples in  $\mathcal{P}_{i,j}$ . Suppose the sets of tuples returned by the two recursive invocations are available. For every quintuple  $(i, j, \mathbf{v}, \ell, \tau) \in \mathcal{P}_{i,j}$ , let  $\vec{\mathbf{u}}(j, \emptyset, \mathbf{u}, \ell_1, \tau_1)$  be the tuple returned for the quintuple  $(j, \emptyset, \mathbf{u}, \ell_1, \tau_1) \in \mathcal{P}_{j,\emptyset}$ , and  $\vec{\mathbf{u}}(i, \text{next}(j), \mathbf{u}, \ell_2, \tau_2)$  be the tuple returned for the quintuple  $(i, \text{next}(j), \mathbf{u}, \ell_2, \tau_2) \in \mathcal{P}_{i,\text{next}(j)}$ . We simply join tuple  $\vec{\mathbf{u}}(j, \emptyset, \mathbf{u}, \ell_1, \tau_1)$ , and tuple  $\vec{\mathbf{u}}(i, \text{next}(j), \mathbf{u}, \ell_2, \tau_2)$  together with tuple  $\mathbf{u}$  to form the tuple  $\vec{\mathbf{u}}(i, j, \mathbf{u}, \ell, \tau)$ . Finally, we will return all tuples for each quintuple  $(i, j, \mathbf{v}, \ell, \tau) \in \mathcal{P}_{i,j}$ .

**Theorem 4.1.** *For an acyclic join instance  $Q$  of size  $N$ , and a set of associated weight functions  $\{p_j\}_{R_j \in Q}$ , there exists an algorithm that returns one subset sample of the instance  $\Psi = \langle \text{Join}(Q), p \rangle$  in  $O(\min\{N \log N \log \log N + \mu_\Psi \log N, N \log^2 N + \mu_\Psi\})$  expected time.*

**PROOF OF THEOREM 4.1.** The correctness follows from that of Algorithm 5, since we exactly simulate the execution of all invocations of Algorithm 5. The data statistics inherited from Theorem 3.4 can be computed in  $O(N \log N \log \log N)$  time. All  $X$ -arrays can be computed in  $O(N \log N)$  time since for every  $\mathbf{u} \in R_i$  there is a unique  $\mathbf{v} \in R_i[\text{key}(i)]$  such that  $\mathbf{u} \in R_i \times \mathbf{v}$ . All  $Y$ -arrays can be computed in  $O(N \log^2 N)$  time, since there are  $O(N \log N)$  arrays and each array has size  $O(\log N)$ . Using the precomputed  $M$  and  $W$  statistics, every value in every  $Y$  array is computed in  $O(1)$  time.

We fix a table  $R_i$  and one of its children  $R_j$  in  $\mathcal{T}$ . Since  $|\mathcal{P}_{i,j}| = O(1 + \mu_\Psi)$ , Radix sort runs in  $O(N + \mu_\Psi)$  expected time. To compute the smallest tuple  $\mathbf{u}$  for each quintuple in  $\mathcal{P}_{i,j,\mathbf{v},\ell}$ , we iterate over all quintuples in  $\mathcal{P}_{i,j,\mathbf{v},\ell}$  and traverse  $X_{i,j,\mathbf{v},\ell}$  once, in the worst case. Since for every  $\mathbf{u} \in R_i$  there is a unique  $\mathbf{v}$  with  $\mathbf{u} \in R_i \times \mathbf{v}$  and every quintuple in  $\mathcal{P}_{i,j}$  belongs to a unique group  $\mathcal{P}_{i,j,\mathbf{v},\ell}$ , we have  $\sum_\ell (|X_{i,j,\mathbf{v},\ell}| + |\mathcal{P}_{i,j,\mathbf{v},\ell}|) = O(N \log N + \mu_\Psi)$ . Similarly, computing all pairs  $(\ell_1, \ell_2)$  requires a single pass over all quintuples in each  $\mathcal{P}_{i,j,\mathbf{u},\ell}$  and all  $Y$  tables. Since there are  $O(N \log N)$  tables  $Y$  of size  $O(\log N)$ , the total cost is  $O(N \log^2 N + \mu_\Psi)$  expected time. The number of nodes in  $\mathcal{T}$  is  $k = O(1)$  and every node has  $O(1)$  children, so, overall, the running time of the one-shot algorithm is  $O(N \log^2 N + \mu_\Psi)$  expected time.  $\square$

## 5 Dynamic Subset Sampling over Joins

In this section, we move to the dynamic setting, where tuples are inserted one by one. We first adapt the techniques of [23] to directly extend our non-optimal (static) index from Theorem 3.2 to a dynamic variant that maintains an approximate **DirectAccess** index. We then introduce new techniques and ideas to extend our optimized static index from Theorem 3.4 to the dynamic setting. Throughout, we assume that, for each relation, tuples are ordered by their insertion timestamps.

### 5.1 Basic Index

The high-level idea is to approximately maintain the static index presented in Section 3. However, as pointed out by [23], even maintaining the join size for line-3 join  $Q = \{R_1(X, Y), R_2(Y, Z), R_3(Z, W)\}$  requires at least  $\Omega(\sqrt{N})$  time for each insertion, which is too costly in the streaming setting. Hence, we follow the approach proposed by [23] by maintaining an approximate **DirectAccess** index more efficiently while not sacrificing the overall sampling efficiency asymptotically. We first define the delta query  $\Delta \text{Join}(Q, \mathbf{u})$  as the set of new join results generated when tuple  $\mathbf{u}$  is inserted into a relation of  $Q$ . Formally, if  $\mathbf{u}$  is inserted into  $R_i$ ,  $\Delta \text{Join}(Q, \mathbf{u}) = \text{Join}(Q \setminus \{R_i\} \cup \{R_i \cup \{\mathbf{u}\}\}) \setminus \text{Join}(Q)$ .

**Theorem 5.1** ([23]). *Given an initially empty acyclic join  $Q$ , we can maintain an index  $\mathcal{L}$  on  $Q$  that supports the following operations:*

- When a tuple  $t$  is added to  $Q$ ,  $\mathcal{L}$  can be updated in amortized  $O(\log N)$  time.

- The index defines an array  $J \supseteq \text{Join}(Q)$  where the tuples in  $\text{Join}(Q)$  are the real tuples and the others are dummy. The index can return  $|J|$  in  $O(1)$  time. For any given  $j \in \llbracket L \rrbracket$ , it returns  $J[j]$  in  $O(\log N)$  time. Furthermore,  $J$  is guaranteed to be  $\lambda$ -dense<sup>3</sup> for some constant probability  $0 < \gamma \leq 1$ .
- The above also holds when  $\text{Join}(Q)$  is replaced by the delta query  $\Delta\text{Join}(Q, \mathbf{u})$  for any  $\mathbf{u} \notin Q$ .

where  $N$  is the total number of insertions, but the index does not need the knowledge of  $N$  in advance.

Using Theorem 5.1, we derive a simple dynamic algorithm by adapting the framework from Section 3.1. As in Section 3.1, we define  $(L + 1)^k$  sub-instances  $Q_j$  for  $j \in \llbracket L \rrbracket^k$ . In the dynamic setting, we maintain a dynamic index  $\mathcal{L}_j$  (as described in Theorem 5.1) for each sub-instance  $Q_j$ . When a tuple  $t$  is inserted into relation  $R_i$ , we first determine its bucket index  $a$  such that  $t \in R_i^{(a)}$ . This tuple  $t$  participates in  $(L + 1)^{k-1}$  sub-instances. We update the corresponding  $(L + 1)^{k-1}$  indexes. Since each update takes  $O(\log N)$  time amortized, the total amortized update time is  $O(L^{k-1} \log N)$ . To answer a subset sampling query, we use the batched rejection sampling strategy. For each sub-instance  $Q_j$ , the index  $\mathcal{L}_j$  provides the size of the implicit array  $|J_j|$ . We use  $|J_j|$  as an upper bound on the size of the sub-join to define the meta-sampling probability  $q(j)$ . Suppose a sub-instance  $Q_j$  is selected and the algorithm needs to retrieve the tuple  $u = J_j[z]$ . If  $u$  is a dummy tuple (which happens with probability at most  $1 - \gamma$ ) we reject it. If  $u$  is real, we proceed with the standard rejection check based on its weight. Since  $\gamma$  is a constant, the expected query time remains asymptotically the same as the static case.

**Corollary 5.2.** *Given an acyclic join instance  $Q$  consisting of  $k$  relations, and a set of associated weight functions  $\{p_j\}_{R_j \in Q}$ , we can maintain a dynamic index under tuple insertions using  $O(N \log^{k-1} N)$  space, with amortized update time  $O(\log^k N)$ , while supporting answering the subset sampling query in  $O(1 + \mu_\Psi \log N)$  expected time, where  $N$  is the total number of insertions.*

## 5.2 Optimized Dynamic Index for Subset Sampling over Joins

To improve both space and update time, we adapt the optimized index from Section 3.2 to the dynamic setting. This adaptation is nontrivial, since the index in Section 3.2 is inherently static and therefore requires new ideas. Rather than maintaining disjoint indexes for every bucket combination, we maintain the aggregated statistics (the  $W$  and  $M$  values) directly. Because maintaining exact counts is too costly, we instead store *approximate* statistics, rounded up to the next power of 2.

For each relation  $R_i \in Q$ , each child node  $R_j \in C_i$ , and each tuple  $\mathbf{u} \in R_i$ , we maintain an upper bound  $\tilde{W}_{i,\mathbf{u}}^j(\ell)$  on  $W_{i,\mathbf{u}}^j(\ell)$ . Similarly, for each non-root node  $R_i$ , tuple  $\mathbf{v} \in R_i[\text{key}(i)]$ , and score  $\ell \in \llbracket L - 1 \rrbracket$ , we maintain  $\hat{M}_{i,\mathbf{v}}(\ell)$  as the sum of the approximate counters:

$$\hat{M}_{i,\mathbf{v}}(\ell) = \sum_{\mathbf{u} \in R_i \bowtie \mathbf{v}} \tilde{W}_{i,\mathbf{u}}^0(\ell), \quad (7)$$

and use it to compute the upper bound  $\tilde{M}_{i,\mathbf{v}}(\ell) = 2^{\lceil \log \hat{M}_{i,\mathbf{v}}(\ell) \rceil}$ . Consider a leaf node  $R_i$ . As  $C_i = \emptyset$ , we have  $\tilde{W}_{i,\mathbf{u}}^0(\ell) = 1$  if  $\ell = \phi(\mathbf{u})$ , and  $\tilde{W}_{i,\mathbf{u}}^0(\ell) = 0$  otherwise. Consider an internal node  $R_i$ . Suppose  $\tilde{W}$  and  $\tilde{M}$ -values have been computed for each child node  $R_j \in C_i$ . We compute  $\tilde{W}_{i,\mathbf{u}}^j$  in a decreasing ordering of nodes in  $C_i$ . Let  $R_{\text{next}(j)} \in C_i$  be the immediate next sibling of  $R_j$ .

$$\tilde{W}_{i,\mathbf{u}}^j(\ell) = \sum_{(\ell_1, \ell_2) \in \llbracket \ell - \phi(\mathbf{u}) \rrbracket : \ell_1 + \ell_2 = \ell - \phi(\mathbf{u})} \tilde{M}_{j,\mathbf{u}[\text{key}(j)]}(\ell_1) \cdot \tilde{W}_{i,\mathbf{u}}^{\text{next}(j)}(\ell_2). \quad (8)$$

Also, we define  $\tilde{W}_{i,\mathbf{u}}^{\text{next}(j^*)}(\ell) = 1$  and  $\tilde{W}_{i,\mathbf{u}}^{\text{next}(j^*)}(\ell) = 0$  for  $\ell > 0$  for the largest child node  $R_{j^*} \in C_i$ .

<sup>3</sup>A stream  $S = \langle x_1, x_2, \dots, x_n \rangle$  is  $\lambda$ -dense for  $0 < \lambda \leq 1$ , if  $r_i \geq \lambda \cdot (i - 1)$  for all  $i$ , where  $r_i$  is the number of real (non-dummy) items in the first  $i - 1$  items.

**Algorithm 7: Update**( $i, v, \ell, \Delta$ )

---

**Input:**  $i \in [k]$ , tuple  $v \in \mathbf{dom}(\text{key}(i))$ , score  $\ell$ , and  $\tilde{W}_{i,u}^0(\ell)$  increased by  $\Delta$  for  $u \in R_i \times v$ .

- 1 Update the prefix-sum tree built for  $\{\tilde{W}_{i,u'}^0(\ell) : u' \in R_i, u'[\text{key}(i)] = u[\text{key}(i)]\}$ ;
- 2  $\hat{M}_{i,v}(\ell) \leftarrow \hat{M}_{i,v}(\ell) + \Delta$ ;
- 3  $\tilde{M}_{i,v}(\ell) \leftarrow 2^{\lceil \log \hat{M}_{i,v}(\ell) \rceil}$ ;
- 4 **if**  $\tilde{M}_{i,v}(\ell)$  changes and  $\text{parent}(i) \neq \text{null}$  **then**
- 5      $p \leftarrow \text{parent}(i)$ ;
- 6     **foreach** tuple  $u \in R_p \times v$  **do**
- 7         Compute  $\tilde{W}_{p,u}^j(\ell')$  for each score  $\ell' \in \llbracket L-1 \rrbracket$  and  $R_j \in C_p$  according to (8);
- 8         **foreach**  $\ell' \in \llbracket L-1 \rrbracket$  **do**
- 9             **if**  $\tilde{W}_{p,u}^0(\ell')$  is increased (by  $\Delta'$ ) **then** **Update**( $p, u[\text{key}(p)], \ell', \Delta'$ );

---

We point out that  $\tilde{W}_{i,u}^j(\ell)$  (resp.  $\tilde{M}_{i,v}(\ell)$ ) is a constant-approximation of  $W_{i,u}^j(\ell)$  (resp.  $M_{i,v}(\ell)$ ). For each score  $\ell \in \llbracket L-1 \rrbracket$ , and each child node  $R_j \in C_i$ , we maintain a dynamic prefix-sum tree on the values  $\tilde{W}_{i,u}^j(\ell)$  over all tuples  $u \in R_i$  under the pre-determined ordering. This allows us to support point updates and prefix sum queries in  $O(\log N)$  time.

**Update Procedure.** When a new tuple  $u$  is inserted into relation  $R_i$ , we first initialize its local statistics. If  $R_i$  is a leaf node,  $\tilde{W}_{i,u}^0(\ell) = 1$  if  $\ell = \phi(u)$  and  $\tilde{W}_{i,u}^0(\ell) = 0$  otherwise. If  $R_i$  is an internal node, we compute  $\tilde{W}_{i,u}^0(\ell)$  for all score  $\ell \in \llbracket L-1 \rrbracket$  via Equation (8), which takes  $O(L \log L)$  time using FFT. After initializing  $u$ , we must update the statistics in the parent node. Specifically, inserting  $u$  increases  $\hat{M}_{i,v}(\ell) = \sum_{u' \in R_i \times v} \tilde{W}_{i,u'}^0(\ell)$  for  $v = u[\text{key}(i)]$ . We update  $\hat{M}_{i,v}(\ell)$  and check if  $\tilde{M}_{i,v}(\ell)$  changes. If  $\tilde{M}_{i,v}(\ell)$  changes, we must propagate this change to the parent relation  $R_{\text{parent}(i)}$ . This propagation is described in Algorithm 7. When a tuple  $u$  is inserted into a leaf  $R_i$ , we just invoke  $\text{UPDATE}(i, u[\text{key}(i)], \phi(u), 1)$ . If it is inserted in an inner node we invoke  $\text{UPDATE}(i, u[\text{key}(i)], \ell, \tilde{W}_{i,u}^0(\ell))$  for every  $\ell \in \llbracket L-1 \rrbracket$ .

**Complexity Analysis.** The update cost is dominated by the propagation of changes. Note that  $\tilde{M}_{i,v}(\ell)$  is an upper bound that increases at most  $O(\log N)$  times for any fixed  $i, v, \ell$ , since  $\hat{M}_{i,v}(\ell)$  only increases if  $\hat{M}_{i,v}(\ell)$  gets doubled. When  $\tilde{M}_{i,v}(\ell)$  changes, we perform convolution operations for tuples in the parent node that can be joined with  $v$ . Since  $L = O(\log N)$ , the convolution takes  $O(L \log L)$  time. We show that the total number of times  $\tilde{M}$  values change across the entire stream of length  $N$  is bounded as  $O(NL \log N)$ . So, the amortized update time is  $O(L^2 \log N \log L) = O(\log^3 N \log \log N)$ . The space complexity remains  $O(NL) = O(N \log N)$ .

**Query Procedure.** The query procedure follows the optimized strategy from Section 3.2, invoking Algorithm 3 on the score buckets  $\mathcal{B}_\ell$  (for  $\ell < L$ ) and the tail bucket  $\mathcal{B}_{\geq L}$ . To retrieve the  $\tau$ -th tuple from a selected bucket, we employ a modified version of Algorithm 5 that traverses the approximate statistics  $\tilde{W}$  and  $\tilde{M}$ . As discussed in Section 5.1, since these statistics are upper bounds, the index implicitly defines a superset of the join results containing “dummy” tuples. If the retrieved tuple is a dummy, we reject it. Since  $\tilde{W}$  is a constant-factor approximation of  $W$ , the probability of drawing a real tuple is at least  $2^{-c}$  for some constant  $c$ . Consequently, sampling  $\chi$  tuples takes  $O(\chi)$  expected time, and the overall expected query time remains  $O(1 + \mu_\Psi \log N)$ .

**Theorem 5.3.** *Given an acyclic join instance  $Q$  consisting of  $k$  relations, and a set of associated weight functions  $\{p_j\}_{R_j \in Q}$ , we can maintain a dynamic index under tuple insertions using  $O(N \log N)$  space,*

with amortized update time of  $O(\log^3 N \log \log N)$ , while supporting answering the subset sampling query in  $O(1 + \mu_\Psi \log N)$  expected time, where  $N$  is the total number of insertions.

**PROOF OF THEOREM 5.3.** We first show that all values of  $\tilde{W}$ ,  $\hat{M}$ , and  $\tilde{M}$  are correctly updated. Wlog, assume that all values of  $\tilde{W}$ ,  $\tilde{M}$ , and  $\hat{M}$  are correct in the subtree  $\mathcal{T}_i$  up to the moment when  $\text{UPDATE}(i, v, \ell, \Delta)$  is invoked for the first time at a node  $R_i$  located at level  $h$  of  $\mathcal{T}$ . Let  $\Delta$  denote the correct increment to  $\hat{M}_{i,v}(\ell)$ . We distinguish two cases. If  $\hat{M}_{i,v}(\ell)$  does not change, then by Equation (8), no value  $\tilde{W}$  needs to be updated in any ancestor of  $R_i$ . Since  $\tilde{W}$  remains unchanged, all corresponding  $\hat{M}$  values on the ancestors of  $R_i$  also remain unchanged. Therefore, no further updates are required. Otherwise, suppose that  $\hat{M}_{i,v}(\ell)$  changes. We then recompute all values  $\tilde{W}_{p,u}^j(\ell')$  according to Equation (8) (Algorithm 7, line 7). If a value  $\tilde{W}_{p,u}^j(\ell')$  with  $j \neq \emptyset$  increases, then no further  $\tilde{W}$  or  $\hat{M}$  values in the ancestors of  $R_i$  need to be updated. If a value  $\tilde{W}_{p,u}^0(\ell')$  does not increase, again no changes are required in the ancestors. Finally, if  $\tilde{W}_{p,u}^0(\ell')$  increases by a nonzero amount  $\Delta'$ , then by Equation (7), the value  $\hat{M}_{i,v}(\ell')$  must be increased by  $\Delta'$ . Consequently, the recursive call  $\text{UPDATE}(p, \mathbf{u}[\text{key}(p)], \ell', \Delta')$  is triggered, and  $\hat{M}_{i,v}(\ell')$  is updated at line 2 of Algorithm 7. Applying the argument recursively establishes the correctness of the entire process.

We next show that the amortized update time is  $O(\log^3 N \log \log N)$ . The key observation is that for each relation  $R_i$  and its parent  $R_p$ , we update the values  $\tilde{W}_{p,u}^j(\ell')$  if and only if  $\hat{M}_{i,v}(\ell)$  changes, where  $\mathbf{u} \in R_p \times v$  (note that for each  $\mathbf{u}$  there exists a unique  $v$  such that  $\mathbf{u} \in R_p \times v$ ). For any fixed  $i, v$ , and  $\ell$ , the value  $\hat{M}_{i,v}(\ell)$  changes at most  $O(\log N)$  times. Whenever an update occurs, we recompute all values  $\tilde{W}_{p,u}^j(\ell')$  for every  $\ell' \in \llbracket L - 1 \rrbracket$  using the FFT algorithm, which takes  $O(L \log L)$  time. Since there are  $L$  possible scores, for fixed  $i$  and  $v$  there are  $O(L)$  distinct values  $\hat{M}_{i,v}(\ell)$ . Hence, the values  $\tilde{W}_{p,u}^j(\ell')$  are updated at most  $O(L \log N)$  times in total, each costing  $O(L \log L)$  time. This yields an amortized update time of  $O(\log^3 N \log \log N)$ . Finally, after every  $N$  insertions, we rebuild the dynamic index from scratch and update  $N \leftarrow 2N$ . This rebuilding step preserves the same amortized update bound of  $O(\log^3 N \log \log N)$ .

By definition, there exists a constant  $c^*$ , depending only on  $|\mathcal{T}|$ , such that  $W_{i,u}^j \leq \tilde{W}_{i,u}^j \leq c^* W_{i,u}^j$  and  $M_{i,u} \leq \tilde{M}_{i,u} \leq c^* M_{i,u}$ , for all  $i, j$ , and  $\mathbf{u}$ . Therefore, the index implicitly represents a superset of the true join results that may include dummy tuples. Since the number of non-dummy tuples is a constant fraction of this superset, the probability of sampling a non-dummy tuple using Algorithm 5 is at least  $2^{-c}$ , where  $c$  is a constant depending only on  $|\mathcal{T}|$ , similarly to [23]. Consequently, the expected query time remains  $O(1 + \mu_\Psi \log N)$ .  $\square$

This index can be used to generate a one-shot subset sample with the following running time:

**Corollary 5.4.** *Given an acyclic join instance  $Q$  consisting of  $k$  relations, and a set of associated weight functions  $\{p_j\}_{R_j \in Q}$ , we can maintain a one-shot subset sample in  $O(N \log^3 N \log \log N + \mu_\Psi \log N)$  expected time, where  $N$  is the total number of insertions.*

## 6 Conclusion

We presented the first efficient framework for subset sampling over joins without full materialization. By exploiting the join structure of the query and decomposability of weight functions, our algorithms achieve near-optimal complexity for static indexing, one-shot sampling and dynamic maintenance under insertions. We aim to extend our theoretical framework to support a broader class of weight definitions. While this paper addressed standard aggregation functions, efficient sampling under complex, non-monotonic, or holistic aggregation functions remains open. Furthermore, we plan to bridge the gap between theory and practice by implementing and evaluating our indexes on real-world machine learning workloads over relational data.

## References

- [1] S. Abiteboul, R. Hull, and V. Vianu. *Foundations of databases*, volume 8. Addison-Wesley Reading, 1995.
- [2] M. Abo-Khamis, S. Im, B. Moseley, K. Pruhs, and A. Samadian. A relational gradient descent algorithm for support vector machine training. In *Symposium on Algorithmic Principles of Computer Systems (APOCS)*, pages 100–113. SIAM, 2021.
- [3] P. K. Agarwal, A. Esmailpour, X. Hu, S. Sintos, and J. Yang. Computing a well-representative summary of conjunctive query results. *Proceedings of the ACM on Management of Data*, 2(5):1–27, 2024.
- [4] M. Arenas, T. C. Merkl, R. Pichler, and C. Riveros. Towards tractability of the diversity of query answers: Ultrametrics to the rescue. *Proceedings of the ACM on Management of Data*, 2(5):1–26, 2024.
- [5] A. Atserias, M. Grohe, and D. Marx. Size Bounds and Query Plans for Relational Joins. *SIAM J. of Comp.*, 42(4):1737–1767, 2013.
- [6] O. Bachem, M. Lucic, and A. Krause. Practical coresets constructions for machine learning. *arXiv preprint arXiv:1703.06476*, 2017.
- [7] L. Bekkers, F. Neven, L. Pantelis, and S. Vansummeren. Poisson sampling over acyclic joins. *Proceedings of the ACM on Management of Data*, 4(2), 2026.
- [8] C. Berkholz, J. Keppeler, and N. Schweikardt. Answering conjunctive queries under updates. In *Proceedings of the 36th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*, PODS '17, page 303–318, New York, NY, USA, 2017. Association for Computing Machinery. ISBN 9781450341981.
- [9] S. Bhattacharya, P. Kiss, A. Sidford, and D. Wajc. Near-Optimal Dynamic Rounding of Fractional Matchings in Bipartite Graphs. In *STOC*, pages 59–70. ACM, 2024.
- [10] A. Borodin and J. I. Munro. *The computational complexity of algebraic and numeric problems*, volume 1. Elsevier, 1975.
- [11] K. Bringmann and K. Panagiotou. Efficient sampling methods for discrete distributions. In *ICALP*, pages 133–144, 2012.
- [12] N. Carmeli, S. Zeevi, C. Berkholz, B. Kimelfeld, and N. Schweikardt. Answering (unions of) conjunctive queries using random access and random-order enumeration. In *Proceedings of the 39th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*, pages 393–409, 2020.
- [13] N. Carmeli, N. Tziavelis, W. Gatterbauer, B. Kimelfeld, and M. Riedewald. Tractable orders for direct access to ranked answers of conjunctive queries. *ACM Transactions on Database Systems*, 48(1):1–45, 2023.
- [14] S. Chaudhuri, R. Motwani, and V. Narasayya. On random sampling over joins. In *Proceedings of the 1999 ACM SIGMOD International Conference on Management of Data*, SIGMOD '99, page 263–274, New York, NY, USA, 1999. Association for Computing Machinery. ISBN 1581130848.
- [15] J. Chen, Q. Yang, R. Huang, and H. Ding. Coresets for relational data and the applications. *Advances in Neural Information Processing Systems*, 35:434–448, 2022.
- [16] Y. Chen and K. Yi. Random Sampling and Size Estimation Over Cyclic Joins. In C. Lutz and J. C. Jung, editors, *23rd International Conference on Database Theory (ICDT 2020)*, volume 155 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 7:1–7:18, Dagstuhl, Germany, 2020. Schloss Dagstuhl – Leibniz-Zentrum für Informatik. ISBN 978-3-95977-139-9.
- [17] Z. Cheng and N. Koudas. Nonlinear models over normalized data. In *2019 IEEE 35th International Conference on Data Engineering (ICDE)*, pages 1574–1577. IEEE, 2019.
- [18] Z. Cheng, N. Koudas, Z. Zhang, and X. Yu. Efficient construction of nonlinear models over normalized data. In *2021 IEEE 37th International Conference on Data Engineering (ICDE)*, pages 1140–1151. IEEE, 2021.
- [19] J. W. Cooley and J. W. Tukey. An algorithm for the machine calculation of complex fourier series. *Mathematics of computation*, 19(90):297–301, 1965.
- [20] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein. *Introduction to algorithms*. MIT press, 2022.
- [21] G. Cormode, M. Garofalakis, P. J. Haas, and C. Jermaine. Synopses for massive data: Samples, histograms, wavelets, sketches. *Foundations and Trends® in Databases*, 4(1–3):1–294, 2012.
- [22] R. Curtin, B. Moseley, H. Ngo, X. Nguyen, D. Olteanu, and M. Schleich. Rk-means: Fast clustering for relational data. In *International Conference on Artificial Intelligence and Statistics*, pages 2742–2752. PMLR, 2020.
- [23] B. Dai, X. Hu, and K. Yi. Reservoir sampling over joins. In *SIGMOD*, pages 1–26, 2024.
- [24] S. Deng, S. Lu, and Y. Tao. On join sampling and the hardness of combinatorial output-sensitive join algorithms. In *Proceedings of the 42nd ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*, PODS '23, page 99–111, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400701276.
- [25] N. Duffield, C. Lund, and M. Thorup. Priority sampling for estimation of arbitrary subset sums. *JACM*, 54(6):32, 2007.
- [26] A. Esmailpour and S. Sintos. Improved approximation algorithms for relational clustering. *Proceedings of the ACM on Management of Data*, 2(5):1–27, 2024.
- [27] D. Feldman and M. Langberg. A unified framework for approximating and clustering data. In *STOC*, pages 569–578, 2011.

- [28] J. Gan, S. W. Umboh, H. Wang, A. Wirth, and Z. Zhang. Optimal dynamic parameterized subset sampling. *Proceedings of the ACM on Management of Data*, 2(5):1–26, 2024.
- [29] G. Gottlob, G. Greco, F. Scarcello, et al. Treewidth and hypertree width. *Tractability: Practical Approaches to Hard Problems*, 1:20, 2014.
- [30] P. J. Haas and J. M. Hellerstein. Ripple joins for online aggregation. In *Proceedings of the 1999 ACM SIGMOD International Conference on Management of Data*, SIGMOD '99, page 287–298, New York, NY, USA, 1999. Association for Computing Machinery. ISBN 1581130848.
- [31] J. Huang and S. Wang. Subset Sampling and Its Extensions. *arXiv preprint arXiv:2003.01075*, 2023.
- [32] J. Huang and S. Wang. Dips: Optimal dynamic index for poisson  $\pi$ ps sampling. In *SIGKDD*, pages 520–531, 2025.
- [33] J. Huang, Y. Tao, and S. Wang. Acyclic join sampling under selections: Dichotomy, union sampling, and enumeration. In *ICDT*, 2026.
- [34] A. Kara, M. Nkolik, D. Olteanu, and H. Zhang. F-ivm: analytics over relational databases under updates. *The VLDB Journal*, 33(4):903–929, 2024.
- [35] M. A. Khamis, H. Q. Ngo, X. Nguyen, D. Olteanu, and M. Schleich. Ac/dc: in-database learning thunderstruck. In *Proceedings of the second workshop on data management for end-to-end machine learning*, pages 1–10, 2018.
- [36] K. Kim, J. Ha, G. Fletcher, and W.-S. Han. Guaranteeing the  $\tilde{O}(\text{agm}/\text{out})$  runtime for uniform sampling and size estimation over joins. In *Proceedings of the 42nd ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*, PODS '23, page 113–125, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400701276.
- [37] A. Kumar, J. Naughton, and J. M. Patel. Learning generalized linear models over normalized data. In *SIGMOD*, pages 1969–1984, 2015.
- [38] F. Li, B. Wu, K. Yi, and Z. Zhao. Wander join: Online aggregation via random walks. In *Proceedings of the 2016 International Conference on Management of Data*, pages 615–629, 2016.
- [39] D. Marx. Tractable hypergraph properties for constraint satisfaction and conjunctive queries. *Journal of the ACM (JACM)*, 60(6):1–51, 2013.
- [40] T. C. Merkl, R. Pichler, and S. Skritek. Diversity of answers to conjunctive queries. *Logical Methods in Computer Science*, 21, 2025.
- [41] B. Moseley, K. Pruhs, A. Samadian, and Y. Wang. Relational algorithms for k-means clustering. In *International Colloquium on Automata, Languages, and Programming*, 2021.
- [42] R. Motwani and P. Raghavan. *Randomized algorithms*. Cambridge University Press, 1995.
- [43] H. Q. Ngo. Worst-case optimal join algorithms: Techniques, results, and open problems. In *PODS*, pages 111–124. ACM, 2018.
- [44] F. Olken. *Random sampling from databases*. PhD thesis, University of California, Berkeley, 1993.
- [45] F. P. Preparata and M. I. Shamos. *Computational geometry. texts and monographs in computer science. Berlin, Springer-Verlag*, 1985.
- [46] M. Schleich, D. Olteanu, and R. Ciucanu. Learning linear regression models over factorized joins. In *SIGMOD*, pages 3–18, 2016.
- [47] M. Schleich, D. Olteanu, M. Abo-Khamis, H. Q. Ngo, and X. Nguyen. Learning models over relational data: A brief tutorial. In *Scalable Uncertainty Management: 13th International Conference, SUM 2019*, pages 423–432. Springer, 2019.
- [48] V. Surianarayanan, N. Kumar, and S. Sintos. Clustering with set outliers and applications in relational clustering. *Proc. ACM Manag. Data*, 3(5), 2025. doi: 10.1145/3767712.
- [49] K. Yang, Y. Gao, L. Liang, B. Yao, S. Wen, and G. Chen. Towards factorized svm with gaussian kernels over normalized data. In *2020 IEEE 36th International Conference on Data Engineering (ICDE)*, pages 1453–1464. IEEE, 2020.
- [50] M. Yannakakis. Algorithms for acyclic database schemes. In *Proceedings of the Seventh International Conference on Very Large Data Bases - Volume 7*, VLDB '81, page 82–94. VLDB Endowment, 1981.
- [51] L. Yi, H. Wang, and Z. Wei. Optimal Dynamic Subset Sampling: Theory and Applications. In *SIGKDD*, pages 3116–3127. ACM, 2023.
- [52] Z. Zhao, R. Christensen, F. Li, X. Hu, and K. Yi. Random sampling over joins revisited. In *SIGMOD*, 2018.

## A Full Notation Table

Please find our main notations in Table 2.

	Notation	Description
General	$S$	A set of elements.
	$p$	A probability function assigning a probability to each element in $S$ .
	$\Psi = \langle S, p \rangle$	A subset sampling problem instance.
	$X$	A random subset sample from a subset sampling instance.
	$\mu_\Psi$	The expected size of a subset sample $X$ from $\Psi$ .
	$[x]$	The set of integers $\{1, 2, \dots, x\}$ .
	$\llbracket x \rrbracket$	The set of integers $\{0, 1, 2, \dots, x\}$ .
	DirectAccess	A Direct Access oracle.
Subset Sampling over Joins	<b>att, dom</b>	The set of all attributes and the domain of all values.
	$\mathbf{u}, \mathbf{v}$	Tuples.
	$R_i$	The $i$ -th relation in a join query.
	schema( $R_i$ )	The schema (set of attributes) of relation $R_i$ .
	$Q = \{R_1, \dots, R_k\}$	A join query consisting of $k$ relations.
	$N$	The input size of $Q$ , i.e., $\sum_{i=1}^k  R_i $ .
	Join( $Q$ )	The set of result tuples from the join query $Q$ .
	$G = (V, E)$	The schema hypergraph of a join query $Q$ .
	$\rho^*$	The fractional edge covering number of schema graph $G$ .
	$p_i(\mathbf{u})$	The probability associated with a tuple $\mathbf{u} \in R_i$ .
	$p(\mathbf{u})$	Prob. of join result $\mathbf{u}$ , defined as $\prod_{i=1}^k p_i(\mathbf{u}[\text{schema}(R_i)])$ .
	$L$	The number of partitions/buckets, defined as $\lceil \rho^* \log N \rceil$ .
	$\phi(\mathbf{u})$	The score of a tuple $\mathbf{u}$ .
$\mathcal{B}_\ell$	The set (bucket) of join results with score $\ell$ .	
Optimized Index	$\mathcal{T}$	Join tree for an acyclic join.
	$\mathcal{T}_i$	The subtree of $\mathcal{T}$ rooted at relation $R_i$ .
	$C_i$	The set of child nodes of relation $R_i$ in $\mathcal{T}$ .
	$\mathcal{T}_i^j$	The partial subtree of $\mathcal{T}_i$ excluding children preceding $R_j$ .
	key( $i$ )	Common join attributes between $R_i$ and its parent in $\mathcal{T}$ .
	$W_{i,\mathbf{u}}^j(\ell)$	# of join results with score $\ell$ in $\mathcal{T}_i^j$ involving $\mathbf{u}$ .
	$M_{i,\mathbf{v}}(\ell)$	An aggregated count of $W$ -values for tuples projecting to $\mathbf{v}$ .
Dynamic	$\eta$	A timestamp in the streaming setting.
	$Q^\eta$	The join defined by the first $\eta$ tuples in the stream.
	$p^\eta$	The weight functions associated with $Q^\eta$ .
	$\tilde{W}_{i,\mathbf{u}}^j(\ell)$	An approximated upper bound of $W_{i,\mathbf{u}}^j(\ell)$ .
	$\tilde{M}_{i,\mathbf{v}}(\ell)$	An approximated upper bound of $M_{i,\mathbf{v}}(\ell)$ .

Table 2. Table of Notations

## B Missing Materials in Section 3

Formally, we prove the correctness of Algorithm 5 using the following lemma.

**Lemma B.1.** *Let  $\text{DirectAccess}(S_{\hat{\tau}}, \hat{\tau})$  and  $\text{DirectAccess}(S_{\bar{\tau}}, \bar{\tau})$  be two distinct calls in Algorithm 3. Let  $\hat{\mathbf{w}} \in \text{Join}(Q)$  be the tuple returned by Algorithm 5 on input  $(S_{\hat{\tau}}, \hat{\tau})$ , and let  $\bar{\mathbf{w}} \in \text{Join}(Q)$  be the tuple returned on input  $(S_{\bar{\tau}}, \bar{\tau})$ . Then  $\hat{\mathbf{w}} \neq \bar{\mathbf{w}}$ .*

**PROOF.** We assume that all the  $W$  and  $M$ -values are correctly computed. Moreover,  $\hat{\ell} \neq \bar{\ell}$ . We will show that  $\hat{\mathbf{w}} \in \mathcal{B}_{\hat{\tau}}$  and  $\bar{\mathbf{w}} \in \mathcal{B}_{\bar{\tau}}$ . Consider any call of  $\text{DirectAccess}(S_{\ell'}, \tau')$  in Algorithm 3. We claim that Algorithm 5 always returns a join result from  $\mathcal{B}_{\ell'}$ . We prove this by strong induction on the height of the join tree  $\mathcal{T}$ . If the join tree has a single leaf node  $R_i$  (height 1), then in Line 1 the algorithm only considers tuples of  $R_i$  with score  $\ell'$ , so the returned tuple belongs to  $\mathcal{B}_{\ell'}$ . Assume that the claim holds for all join trees of height at most  $h$ . Let  $R_i$  be the root of a join tree of height  $h + 1$ . In Line 4, the algorithm considers pairs  $(\ell_1, \ell_2)$  such that  $\ell_1 + \ell_2 + \phi(\mathbf{u}) = \ell'$ . By the induction hypothesis, we have  $\phi(\mathbf{u}_1) = \ell_1$  and  $\phi(\mathbf{u}_2) = \ell_2$ , and therefore  $\phi(\mathbf{u}_1 \bowtie \mathbf{u} \bowtie \mathbf{u}_2) = \ell'$ .

Applying this argument to  $\text{DirectAccess}(S_{\hat{\tau}}, \hat{\tau})$  and  $\text{DirectAccess}(S_{\bar{\tau}}, \bar{\tau})$ , we obtain  $\hat{\mathbf{w}} \in \mathcal{B}_{\hat{\tau}}$  and  $\bar{\mathbf{w}} \in \mathcal{B}_{\bar{\tau}}$ . Since  $\hat{\ell} \neq \bar{\ell}$ , we have  $\mathcal{B}_{\hat{\tau}} \cap \mathcal{B}_{\bar{\tau}} = \emptyset$ , and thus  $\hat{\mathbf{w}} \neq \bar{\mathbf{w}}$ . We now assume that  $\hat{\ell} = \bar{\ell} = \ell^*$  and  $\hat{\tau} \neq \bar{\tau}$ . Wlog, suppose that  $\hat{\tau} > \bar{\tau}$ . We prove the claim by strong induction on the height of the join tree. If the join tree consists of a single leaf node  $R_i$  (height 1), let  $\hat{\mathbf{u}}$  (resp.,  $\bar{\mathbf{u}}$ ) be the smallest tuple found in Line 1 for the pair  $(\hat{\ell}, \hat{\tau})$  (resp.,  $(\bar{\ell}, \bar{\tau})$ ). For every  $\mathbf{u}' \in R_i$ , the value  $W_{i,\mathbf{u}'}^j(\ell^*)$  is equal to 1 if  $\phi(\mathbf{u}') = \ell^*$  and 0 otherwise. Hence, the sum  $\sum_{\mathbf{u}' \in R_i: \mathbf{u}' \preceq \mathbf{u}} W_{i,\mathbf{u}'}^j(\ell^*)$  increases by one for each new qualifying tuple. Since  $\hat{\tau} \neq \bar{\tau}$ , we conclude that  $\hat{\mathbf{u}} \neq \bar{\mathbf{u}}$  and therefore  $\hat{\mathbf{w}} \neq \bar{\mathbf{w}}$ .

Assume that the claim holds for all join trees of height at most  $h$ . Let  $R_i$  be the root of a join tree of height  $h + 1$ . If  $\hat{\mathbf{u}} \neq \bar{\mathbf{u}}$  in Line 1, then clearly  $\hat{\mathbf{w}} \neq \bar{\mathbf{w}}$ . Hence, we assume that  $\hat{\mathbf{u}} = \bar{\mathbf{u}} = \mathbf{u}$ . The sum  $\sum_{\mathbf{u}' \in R_i: \mathbf{u}' \preceq \mathbf{u}} W_{i,\mathbf{u}'}^j(\ell^*)$  is identical for both instances, so  $\hat{\tau}$  and  $\bar{\tau}$  are reduced by the same amount in Line 2. By a slight abuse of notation, we again denote the updated values by  $\hat{\tau}$  and  $\bar{\tau}$ . Let  $(\hat{\ell}_1, \hat{\ell}_2)$  and  $(\bar{\ell}_1, \bar{\ell}_2)$  be the smallest pairs found in Line 5 for the inputs  $(\hat{\ell}, \hat{\tau})$  and  $(\bar{\ell}, \bar{\tau})$ , respectively. If  $\hat{\ell}_1 \neq \bar{\ell}_1$  or  $\hat{\ell}_2 \neq \bar{\ell}_2$ , then by the first part of the proof we immediately obtain  $\hat{\mathbf{w}} \neq \bar{\mathbf{w}}$ . Hence, we assume that  $\hat{\ell}_1 = \bar{\ell}_1 = \ell_1$  and  $\hat{\ell}_2 = \bar{\ell}_2 = \ell_2$ . The sum  $\sum_{(\ell'_1, \ell'_2) \in \Phi: (\ell'_1, \ell'_2) \prec (\ell_1, \ell_2)} M_{j,\mathbf{u}[\text{key}(j)]}(\ell'_1) \cdot W_{i,\mathbf{u}}^{\text{next}(j)}(\ell'_2)$  is identical for both instances, so  $\hat{\tau}$  and  $\bar{\tau}$  are again reduced by the same amount in Line 6.

We now reach the final stage of the algorithm, where the smallest tuples and score pairs coincide, but  $\hat{\tau} > \bar{\tau}$ . We compute  $\hat{\tau}_1, \hat{\tau}_2$  and  $\bar{\tau}_1, \bar{\tau}_2$  as in Lines 7 and 8. The denominator  $W_{i,\mathbf{u}}^{\text{next}(j)}(\ell_2)$  in Line 7 is the same in both executions. If  $\hat{\tau}_1 \neq \bar{\tau}_1$ , then by the induction hypothesis  $\hat{\mathbf{u}}_1 \neq \bar{\mathbf{u}}_1$ , which implies  $\hat{\mathbf{w}} \neq \bar{\mathbf{w}}$ . Otherwise, assume that  $\hat{\tau}_1 = \bar{\tau}_1 = \tau_1$ . There are two cases. If neither  $\hat{\tau}$  nor  $\bar{\tau}$  is divisible by  $W_{i,\mathbf{u}}^{\text{next}(j)}(\ell_2)$ , then  $\hat{\tau}_2 = ((\hat{\tau} - 1) \bmod W_{i,\mathbf{u}}^{\text{next}(j)}(\ell_2)) + 1 = \hat{\tau} \bmod W_{i,\mathbf{u}}^{\text{next}(j)}(\ell_2) > 0$  and  $\bar{\tau}_2 = ((\bar{\tau} - 1) \bmod W_{i,\mathbf{u}}^{\text{next}(j)}(\ell_2)) + 1 = \bar{\tau} \bmod W_{i,\mathbf{u}}^{\text{next}(j)}(\ell_2) > 0$ . Since  $\hat{\tau} > \bar{\tau}$ , it follows that  $\hat{\tau}_2 > \bar{\tau}_2$ , and by the induction hypothesis  $\hat{\mathbf{u}}_2 \neq \bar{\mathbf{u}}_2$ , again implying  $\hat{\mathbf{w}} \neq \bar{\mathbf{w}}$ . In the second case,  $\hat{\tau}$  is divisible by  $W_{i,\mathbf{u}}^{\text{next}(j)}(\ell_2)$  but  $\bar{\tau}$  is not. Then  $\hat{\tau}_2 = ((\hat{\tau} - 1) \bmod W_{i,\mathbf{u}}^{\text{next}(j)}(\ell_2)) + 1 = W_{i,\mathbf{u}}^{\text{next}(j)}(\ell_2)$  while  $\bar{\tau}_2 = ((\bar{\tau} - 1) \bmod W_{i,\mathbf{u}}^{\text{next}(j)}(\ell_2)) + 1 = \bar{\tau} \bmod W_{i,\mathbf{u}}^{\text{next}(j)}(\ell_2) < W_{i,\mathbf{u}}^{\text{next}(j)}(\ell_2)$ . Hence  $\hat{\tau}_2 > \bar{\tau}_2$ , and by the induction hypothesis  $\hat{\mathbf{u}}_2 \neq \bar{\mathbf{u}}_2$ , which concludes that  $\hat{\mathbf{w}} \neq \bar{\mathbf{w}}$ .  $\square$

## C Discussion On Other Functions

We next discuss the extension to other aggregation functions. For our algorithms, the results can be extended to support MIN, MAX, and SUM. We will show how to adapt the high-level idea in Section 3.1 to these functions.

**MIN and MAX.** Without loss of generality, we consider the function MIN. We assume that a  $\text{DirectAccess}$  oracle is available for the input join  $Q$  (under some fixed ordering), such that it receives an integer  $i \in [|\text{Join}(Q)|]$ , and returns the  $i$ -th element of  $\text{Join}(Q)$ . Similar to before, the ordering can be arbitrary but must remain consistent across multiple invocations of this oracle. Let

$L = \lceil 2\rho^* \log N \rceil$ . The first step of partitioning join result is exactly the same as before. The second step of applying subset sampling to all sub-instances also follows. Let  $\mathbf{j} = (j_1, j_2, \dots, j_k)$ . Every join result in the sub-instance  $Q_j$  has the probability  $\min_{i \in [k]} 2^{-j_i} = 2^{-\max_{i \in [k]} j_i}$  to be sampled. We first point out the following observation on each sub-instance in such a partition:

**Lemma C.1** (Either Light or Near-uniform Sub-instance). *Let  $\mathbf{j} = (j_1, j_2, \dots, j_k)$ . If  $\max_{i \in [k]} j_i \geq L$ , the instance  $\Psi_j$  is light, and otherwise, it is 2-uniform.*

Intuitively, if  $\max_{i \in [k]} j_i \geq L$ , every join result has its probability at most  $\frac{1}{|\text{Join}(Q)|}$ , so by definition this sub-instance is light. Otherwise,  $\max_{i \in [k]} j_i < L$ . Now, the ratio between the maximum and the minimum probability is at most 2, so by definition this sub-instance is 2-uniform. Hence,

- **Preprocessing phase:** We partition input tuples by weights, as well as the join instance  $Q$ . For each sub-instance  $Q_j$  with  $\mathbf{j} \in \llbracket L \rrbracket^k$ , we build DirectAccess oracles and compute the join size.
- **Query phase:** We invoke Algorithm 3 with input  $\{\langle Q_j, p \rangle : \mathbf{j} \in \llbracket L \rrbracket^k\}$  with  $p_j^+ = 2^{-\max_{i \in [k]} j_i}$ . Whenever a join result is sampled from  $Q_j$ , we use the DirectAccess oracle to retrieve it.

For the optimized index, for every  $\mathbf{u}$  in  $R_i$  we define the score  $\phi(\mathbf{u}) = \lfloor -\log p_i(\mathbf{u}) \rfloor$ . For a join result  $\mathbf{u} \in \text{Join}(Q)$ , its score is  $\bar{\phi}(\mathbf{u}) = \max_{i=1}^k \phi(\mathbf{u}[\text{schema}(R_i)])$ . The buckets  $\mathcal{B}_\ell = \{\mathbf{u} \in \text{Join}(Q) \mid \bar{\phi}(\mathbf{u}) = \ell\}$  are defined in a similar way, as for the product function. We note that for any  $\mathbf{u} \in \mathcal{B}_\ell$ ,  $2^{-\ell-1} < p(\mathbf{u}) \leq 2^{-\ell}$ , hence the ratio between the maximum and minimum probability in each bucket  $\mathcal{B}_\ell$  is at most 2. The index is similar to the one we constructed in Section 3.2. During the query phase, we run an algorithm similar to Algorithm 5. In Steps 4 and 5 the number of pairs we try is not  $L = O(\log N)$ , as in the case of the product function, but  $O(L^2) = O(\log^2 N)$ . Additionally, we do not apply the FFT in the preprocessing phase. Overall, our results extend to the MIN (and MAX) function with the same guarantees up to a  $\log N$  factor.

**SUM.** Similarly, our indexing framework naturally extends to the SUM function. Interestingly, SUM can be handled in exactly the same way as MAX: we define identical scores and buckets for the SUM function as those used for the MAX function. Namely, for every  $\mathbf{u}$  in  $R_i$  we define the score  $\phi(\mathbf{u}) = \lfloor -\log p_i(\mathbf{u}) \rfloor$ . For a join result  $\mathbf{u} \in \text{Join}(Q)$ , its score is  $\bar{\phi}(\mathbf{u}) = \min_{i=1}^k \phi(\mathbf{u}[\text{schema}(R_i)])$ . For every  $\ell$  we define the bucket  $\mathcal{B}_\ell = \{\mathbf{u} \in \text{Join}(Q) \mid \bar{\phi}(\mathbf{u}) = \ell\}$ . We note that for any  $\mathbf{u} \in \mathcal{B}_\ell$ ,  $2^{-\ell-1} < p(\mathbf{u}) \leq k \cdot 2^{-\ell}$ , hence the ratio between the maximum and minimum probability in each bucket  $\mathcal{B}_\ell$  is at most  $2k = O(1)$ . The index is similar to the one we constructed in Section 3.2. During the query phase, we run an algorithm similar to Algorithm 5. In Steps 4 and 5 the number of pairs we try is not  $L = O(\log N)$ , as in the case of the product function, but  $O(L^2) = O(\log^2 N)$ . Additionally, we do not apply the FFT in the preprocessing phase. Overall, our results extend to the SUM function with the same guarantees up to a  $\log N$  factor.

Received December 2025; accepted February 2026