Paper Presentation: Reservoir Sampling over Joins

Yuxin Kang



Background

(1) Fully materializing the join can be extremely expensive in both time and space. Uniform sample of the join results would suffice for many purposes, (eg. answer analytical queries, train ML models.

(2) Real-world data arrives continuously, so the database is always growing. We need a method that keeps a uniform random sample of the join results as new tuples stream in, rather than sampling after everything is collected.

Previous Results

(1) Static Join Sampling:

Acyclic joins: Build an index in O(N), draw one sample in O(1).

Cyclic joins: Build an index in O(N), draw one sample in $O\left(\frac{N^{\rho^*}}{|Q(\mathcal{R})|}\right)$

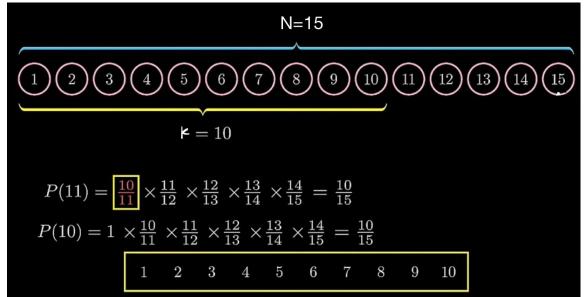
N: the total size of all input relations

|Q(R)|: the number of tuples in the join result.

 ρ^* : the fractional edge cover number

(2) Reservoir Sampling (Classical):

The total number of elements in the stream is unknown; Each element can only be accessed once. We want each element has an equal probability of being sampled.



For item number t: If $t \le k$: store it directly in the sample.

If t > k:Draw i uniformly at random from $\{0, 1, ..., t-1\}$.

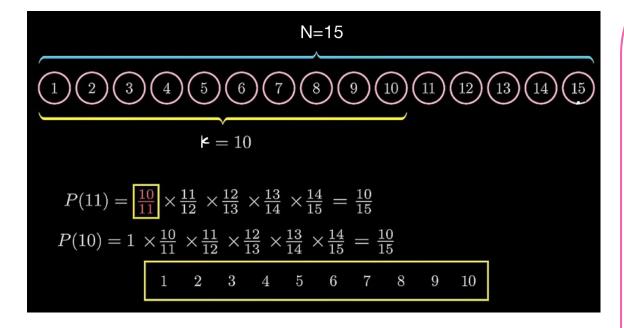
If i < k, replace sample[i] with the new item.

Otherwise, do nothing.

Each item has the probability:

$$rac{k}{k+1} imesrac{k+1}{k+2} imesrac{k+2}{k+3} imes\cdots imesrac{N-2}{N-1} imesrac{N-1}{N}\;=\;rac{k}{N}$$

Reservoir Sampling (Classical), cont.



Without skip optimization: O(N)

With skip optimization: $O\left(k\log\frac{N}{k}\right)$

The skip optimization:

Suppose we have N > 100,000 and k = 10. The probability that the t = 100,000 th item will be accepted is: 0.0001.

99.99% of items will be rejected. So, checking each one is a waste of time.

The probability of keeping the t-th item in the reservoir is $p_t = \frac{k}{t}$.

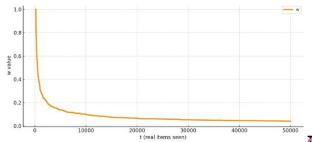
So, Pr(skip exactly q items) =
$$(1 - p_t)^q \cdot p_t$$
 —-> q ~Geom (pt)

The skip-based algorithm does not examine every item, so we can not compute t.

So the algorithm maintains a value of w such that

$$w \leftarrow w \cdot (u')^{1/k} \quad \text{where } u' \sim \text{Unif}(0,1)$$

And w behave like $w \approx \frac{k}{t}$



Reservoir Sampling (Classical) — the skip optimization, cont.

Algorithm:

Initialize reservoir S with the first k items.

Initialize w after filling the reservoir.

Repeat for the rest of the stream:

Draw $q \sim \text{Geom}(w)$.

Skip the next q items in the stream.

Let x be the next item.

Draw $u \sim \text{Unif}(0,1)$.

If u < w, then:

replace a random item in S with x;

draw $u' \sim \text{Unif}(0,1)$ and update $w \leftarrow w \cdot (u')^{1/k}$.

drawing u from a uniform distribution and comparing it to w gives: Pr[u<w] = w.

The expectation of the acceptance of the t-th item is:

$$0 \cdot \frac{t-k}{t} + 1 \cdot \frac{k}{t} = \frac{k}{t}$$

The expected number of acceptances over the entire stream:

$$\sum_{t=k+1}^{N} \frac{k}{t} = k \sum_{t=k+1}^{N} \frac{1}{t} = k (H_N - H_k)$$

$$\approx k (\log N - \log k)$$

$$= k \log \left(\frac{N}{k}\right).$$

Each skip and replacement both take O(1) time. So the total expected time is:

$$O\left(k\log\frac{N}{k}\right)$$

Previous Results

(3) Attempts of Applying Reservoir Sampling to Joins: Zhao et al. investigated the problem over acyclic joins. But their solution takes $O(N^2)$ time in the worst case since their index needs costly maintenance when new tuples arrive. This is essentially no better than a naive approach that re-build the static join sampling index and re-draw the samples after each tuple has arrived $(N *O(N) = O(N^2))$.

Motivation

Can we build a new reservoir sampling algorithm that maintains a sample over joins with a near-linear running time?

Reservoir Sampling with Predicate

Goal: We want to maintain a reservoir sample of size k only over items that satisfy a predicate θ . (Items that satisfy θ are called real items. Items that do not satisfy θ are called dummy items.)

Key changes: The skip counts real items only. Dummy items do not count.

Repeat for the rest of the stream: Draw a skip length $q \sim \text{Geom}(w)$. Skip the next q items in the stream. Let x be the next item. If $\theta(x)$ is false, then continue to the next iteration. Draw $u \sim \text{Unif}(0,1)$. If u < w, then: Choose a random position in S and replace its item with x. Draw $u' \sim \text{Unif}(0,1)$ and update $w \leftarrow w \cdot (u')^{1/k}$.

Reservoir Sampling with Predicate

Why It Remains Correct:

We never adjust w when seeing dummy items, So,

$$w \approx \frac{t}{k} \implies w \approx \frac{t}{r}$$

where r is the counts of only the real items.

Each real item is accepted with probability $\frac{k}{r}$, ensuring every real item is equally likely to appear in the reservoir.

Upper bound: Reservoir Sampling with Predicate

Let: r_i = number of real items among the first (i-1) items

- If x_i is not dummy, then: It cannot enter the reservoir; It does not update w. So, the cost is o.
- If x_i is real, it could enter the reservoir. The probability that a real item enters the sample at this moment is:

$$\Pr(ext{real } x_i ext{ enters}) = rac{k}{r_{i+1}}$$

So, expected work at step i is:

$$\mathrm{Cost}_i = \min\left(1, rac{k}{r_{i+1}}
ight)$$

(If we haven't filled the reservoir yet, cost is 1; once full, it becomes $\frac{k}{r_{i+1}}$)
Sum over all items:

$$ext{Total Time} = O\!\left(\sum_{i=1}^{N} \min\!\left(1, rac{k}{r_i + 1}
ight)
ight)$$

If every item is real, then $r_i = i - 1$:

$$\sum_{i=k}^{N} \frac{k}{i} = O\left(k \log \frac{N}{k}\right)$$



Lower bound: Reservoir Sampling with Predicate

To output a uniform sample of size k over only the real items, a real item at position i must be given probability $\frac{k}{r_{i+1}}$ to enter the reservoir.

Therefore, any correct algorithm must stop and examine position i with probability at least $\min\left(1, \frac{k}{r_i+1}\right)$

If it stops less often than this, it would miss some real items that should have been considered. Those items would have too small a chance of entering the reservoir. The resulting sample would not be uniform anymore.

$$\Omega\!\left(\sum_{i=1}^N \min\left(1,rac{k}{r_i+1}
ight)
ight)$$

Goal: We want to build a data structure that allows us to:

- Incrementally update the join when new tuples arrive
- Efficiently retrieve or skip join results later. b.

(1) Two-table join $R_1(X,Y) \bowtie R_2(Y,Z)$

a. We maintain one bucket per value of Y

b. Each tuple is stored exactly once in its bucket. So, it uses linear space: $\sum (|R_1 \bowtie b| + |R_2 \bowtie b|) = O(N)$

 R_1

Example:

X	Υ	
а	2	
b	3	

$$R_2$$

Y	Z
2	Р
2	Q
3	R

$$\sum_b (|R_1 owtie b| + |R_2 owtie b|) = O(N)$$

(1) Two-table join —cont

- c. If a new tuple t arrives:
- If $t \in R_1$, insert it into bucket $R_1 \bowtie t.Y$
- If $t \in R_2$, insert it into bucket $R_2 \bowtie t.Y$

This is just an append \rightarrow O(1) time.

d. If $t \in R_1$, then new join results are:

$$\Delta J = R2 \bowtie t.Y$$

because the lists are directly stored, accessing any joined result is O(1)

Example:

y value
 List
$$R_1 \bowtie y$$
 List $R_2 \bowtie y$

 2
 $\{(a,2)\}$
 $\{(2,P),(2,Q)\}$

 3
 $\{(b,3)\}$
 $\{(3,R)\}$

$$t=(c,2)\in R_1$$

$$\Delta J = \{(c,2,P),\; (c,2,Q)\}$$



(2) Line-3 Join

a. Index Rebuild: $O(\log N)$

Initialization:

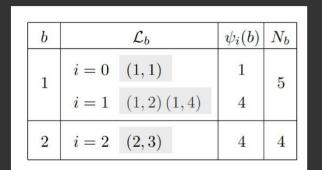
$$R_1(X, Y) \bowtie R_2(Y, Z) \bowtie R_3(Z, W)$$

X	Y
2	5
4	4
5	5
5	1
5	6
6	3
6	6
7	4
7	5
7	6

\mathbf{Y}	Z
1	1
1	2
1	4
2	3
6	5
7	7

L
3
1
1
5
6
7
1
5
6

	b	\mathcal{L}_b	$ \psi_i(b) N_b$	For each tuple $(b,c) \in \mathbb{R}_2$, we converte: $(\widetilde{nt}(C) = 2^{\lceil \log_2 cnt(C) \rceil}, C\widetilde{nt}(C) = 2^{\widetilde{l}} $
	1	i = 0 (1, 1)	1 5	where $Cnt(c) = R>PACI is the number of match in Rs$
	2	i = 1 $(1, 2) (1, 4)i = 2$ $(2, 3)$	4 4	For $b=1$: $(1,1): C=1$, $\widetilde{Cnt}(1)=1=2^{\circ} \Rightarrow bncket i=0$ $(1,2): C=2$, $\widetilde{Cnt}(2)=2=2^{\circ} \Rightarrow bncket i=1$
b: f	A VIV	we from attribute the from attribute	Yin Rel Zin Rel	$(1,4): C=4 / Cnt(4)=2=2' \Rightarrow bucket '='$ (5)



Lb: the list of non-empty buckets for a given be 24(R2)

(Filb): i-th sub-bucket inside bucket b: 41(b)=21. | Lb,il

Nb: Total number of items in bucket b

$$\psi_{0}(1) = 2^{0} \cdot | = 1$$
 $\psi_{1}(1) = 2^{1} \cdot 2 = 4$
 $\psi_{2}(2) = 2^{2} \cdot 1 = 4$
 $N_{1} = 1 + 4 = 5$
 $N_{2} = 4$

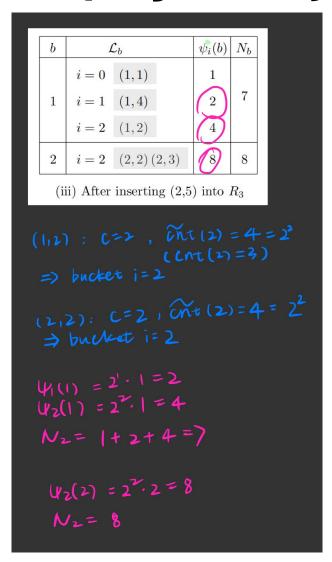
Inserting (2,2) into R2:

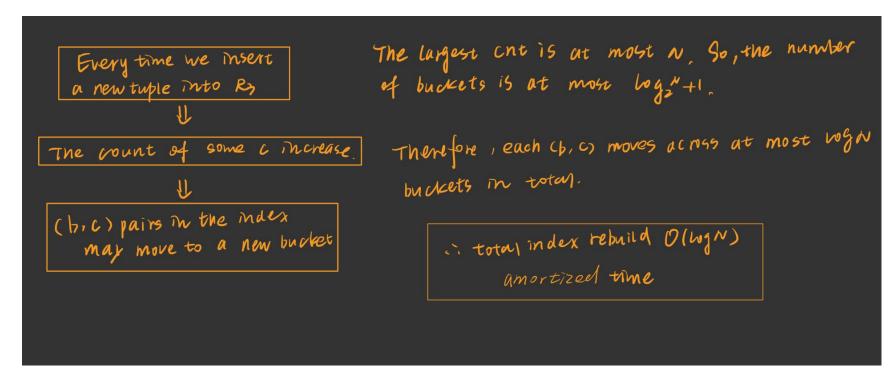
b	\mathcal{L}_b	$\psi_i(b)$	N_b
1	i = 0 (1,1)	1	5
1	i = 1 (1,2) (1,4)	4	0
2	i = 1 (2,2)	2	6
2	i = 2 (2,3)	4	0

(ii) After inserting (2,2) into R_2

(2,2):
$$C=2$$
, $Cnt(2)=2=2^{\circ}$ bucket i= 1
(doesn't affect cut (67))
 $V_1(2)=2^{\circ}\cdot 1=2$
 $N_2=2+4=6$







b. Get a sample:

Step 1: Sample (a,b) from R1 (but not uniformly)

We prefer (a, b) pairs that lead to more join results.

Step 2 : Sample (b, c) from $R_2 \bowtie b$

Binary search over cumulative weights

Step 3: Sample (c,d) uniformly from $R_3 \bowtie c$

Total: $O(\log N)$

 $\longrightarrow O(\log N)$



Reservoir Sampling over Joins

For acyclic:

(1) Update the join index

Cost per tuple: $O(\log N)$

Total over N tuples: $O(N \log N)$

(2)Possibly update the reservoir sample

Retrieve a random join result via the index: $O(\log N)$

Number of times we perform a reservoir replacement: $k \log \frac{N}{k}$

Total reservoir sampling cost: $O\left(k \log N \log \frac{N}{k}\right)$

$$oxed{O(N \log N)} \ + \ \overline{O(k \log N \log(N/k))} = \overline{O(N \log N + k \log N \log(N/k))}$$

Reservoir Sampling over Joins

For cyclic:

Cyclic joins can create many more join results, much larger than N

By the Generalized Hypertree Decomposition:

Fractional Hypertree Width w(Q) tells us how cyclic the query is,

If w = 1, the join is acyclic.

If w > 1, the join is cyclic.

The worst -case join size is $O(N^{w(Q)})$

 $O(N^{w(Q)} \log N + k \log N \log(N/k))$

So the index maintenance time become: $O(N^{w(Q)} \log N)$ And the reservoir sampling part does not change.

Dataset:

Graph dataset: Epinions;

Relational dataset: TPC-DS for decision-support joins;

Relational dataset: LDBC-SNB for join-heavy complex queries.

Join queries tested:

Multiple acyclic joins (e.g., line-3, line-4, line-5, star joins).

Some cyclic joins to evaluate the extension that uses fractional hypertree width.

Baseline:

Zhao et al.



Key Results:

(1) The proposed (RSJoin) update time is stable and stays around 10 microseconds, while the method from Zhao et al. (SJoin) update time is unpredictable and can jump to hundreds of milliseconds, which makes it unreliable in streaming environments.

(2) Although the number of join results grows rapidly as more data is processed, RSJoin's runtime increases almost linearly with the input, while SJoin's runtime grows roughly in proportion to the join size.

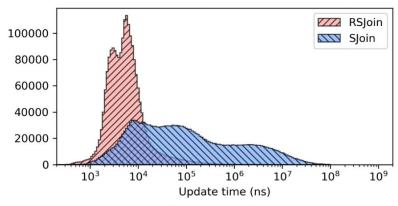


Fig. 6. Update time distribution

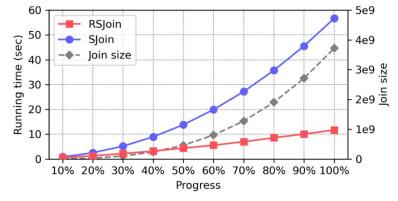


Fig. 7. Running time v.s. input size and join size



Key Results:

(3) RSJoin achieves 4.6×-147.6× speedups over prior methods and successfully completes all join queries, including cyclic ones where SJoin may fail. It also does not depend on foreign-key constraints, making it more scalable and robust across both graph and relational workloads.

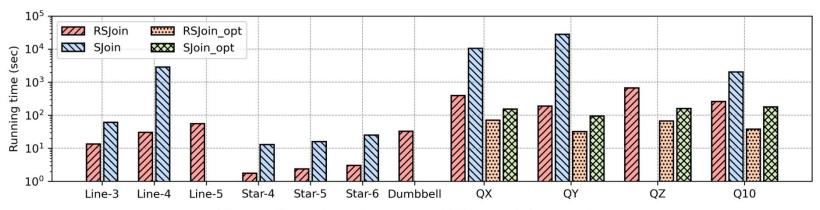


Fig. 5. Running time over different join queries

Key Results:

(4) RSJoin uses much less memory than SJoin as input grows. For complex relational joins (e.g., Q10), SJoin's memory usage increases to tens of GB, while RSJoin remains low and stable.

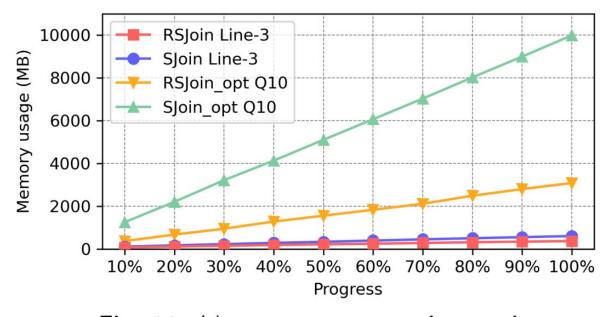


Fig. 11. Memory usage v.s. input size



Conclusion

(1) Introduce reservoir sampling with a predicate, and prove an instance-optimal running time of

$$O\!\left(\sum_{i=1}^N \min\!\left(1,rac{k}{r_i+1}
ight)
ight)$$

- (2) Present a dynamic index for acyclic joins supporting O(log N) amortized updates and O(log N) sampling.
- (3) Combine both to solve reservoir sampling over joins in $O(N \log N + k \log N \log \frac{N}{k})$.
- (4) Extend to cyclic joins with fractional hypertree width , achieving $O(N^w \log N + k \log N \log \frac{N}{k})$

References

https://arxiv.org/abs/2404.03194

http://xhslink.com/o/47rnsirYAcF

