

Towards Output-Optimal Uniform Sampling and Approximate Counting for Join-Project Queries

Xiao Hu¹ Jinchao Huang²

¹University of Waterloo ²The Chinese University of Hong Kong

PODS 2026

June 2026

▪ Bengaluru, India

Why Sampling for Join-Project Queries?

Modern analytics queries are not bare joins. They look like

$$\pi_{\mathbf{y}} (R_1 \bowtie R_2 \bowtie \cdots \bowtie R_k).$$

Many applications need a *sample*, not the full result

- Approximate query processing & visualization
- Cardinality / selectivity estimation for the optimizer
- Training ML models without materializing the join
- Interactive exploration of warehouse-scale data

Prior breakthroughs: optimal sampling for **full joins** (Zhao'18; Deng-Lu-Tao'23; Kim et al.'23).

Reality: real queries **nearly always** project away non-output attributes.

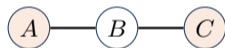
Join-Project Queries

Definition

$Q = (\mathcal{V}, \mathcal{E}, \mathbf{y})$ with output attributes $\mathbf{y} \subseteq \mathcal{V}$; $Q(\mathcal{R}) = \pi_{\mathbf{y}}\left(\bigbowtie_{e \in \mathcal{E}} R_e\right)$.

Matrix Q_{matrix}

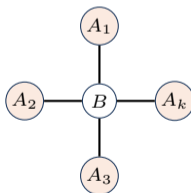
$\pi_{A,C}(R_1(A, B) \bowtie R_2(B, C))$



red = output attr.

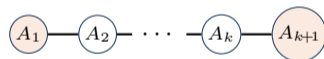
k-Star Q_{star}

$\pi_{A_1, \dots, A_k}(\bigbowtie_i R_i(A_i, B))$



k-Chain Q_{chain}

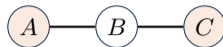
$\pi_{A_1, A_{k+1}}(\bigbowtie_i R_i(A_i, A_{i+1}))$



Notation: N = input size; $\text{OUT} = |Q(\mathcal{R})|$; OUT_{\bowtie} = full-join size; ρ^* = fractional edge-cover number (of the underlying full join query).

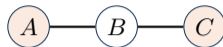
Prior Framework: “Propose-and-Verify” [Chen-Yi’20]

1. **Propose:** draw a candidate t from the *projection-induced full join*
 $\pi_{\mathbf{y}}(R_{e_1}) \bowtie \pi_{\mathbf{y}}(R_{e_2}) \bowtie \dots$
2. **Verify:** check whether t extends to a real full-join tuple.



Prior Framework: “Propose-and-Verify” [Chen-Yi’20]

1. **Propose:** draw a candidate t from the *projection-induced full join*
 $\pi_{\mathbf{y}}(R_{e_1}) \bowtie \pi_{\mathbf{y}}(R_{e_2}) \bowtie \dots$
2. **Verify:** check whether t extends to a real full-join tuple.



Two fundamental bottlenecks

1. **Low acceptance probability:** the projection-induced join is much larger than OUT.
2. **High verification cost:** essentially evaluating the *residual join* on non-output attrs.

Resulting sampling complexity: $\tilde{O}\left(\frac{N^{\rho^*}}{\text{OUT}}\right)$.

Specializing:

$$Q_{\text{matrix}} : \tilde{O}(N^2/\text{OUT}), \quad Q_{\text{star}} : \tilde{O}(N^k/\text{OUT}), \quad Q_{\text{chain}} : \tilde{O}\left(N^{\lceil (k+1)/2 \rceil} / \text{OUT}\right).$$

Are these optimal? — No.

Our Results at a Glance

		Matrix	k -Star	k -Chain	
				$k = 3$	$k \geq 4$
Uniform Sampling	Prior	$\tilde{O}(N^2/\text{OUT})$	$\tilde{O}(N^k/\text{OUT})$	$\tilde{O}(N^2/\text{OUT})$	$\tilde{O}(N^{\lceil \frac{k+1}{2} \rceil} / \text{OUT})$
	Our UB	$\tilde{O}(N/\sqrt{\text{OUT}})$	$\tilde{O}(N/\text{OUT}^{1/k})$	$\tilde{O}(\min\{N, N^2/\text{OUT}\})$	$\tilde{O}(N)$
	Our LB	$\tilde{\Omega}(N/\sqrt{\text{OUT}})$	$\tilde{\Omega}(N/\text{OUT}^{1/k})$	$\tilde{\Omega}(\min\{N, N^2/\text{OUT}\})$	$\tilde{\Omega}(N)$
Approx. Counting	Prior	$\tilde{O}(N^2/\text{OUT})$	$\tilde{O}(N^k/\text{OUT})$	$\tilde{O}(N^2/\text{OUT})$	$\tilde{O}(N^{\lceil \frac{k+1}{2} \rceil} / \text{OUT})$
	Our UB	$\tilde{O}(N/\sqrt{\text{OUT}})$	$\tilde{O}(N/\text{OUT}^{1/k})$	$\tilde{O}(\min\{N, N^2/\text{OUT}\})$	$\tilde{O}(N)$
	Our LB	$\tilde{\Omega}(N/\sqrt{\text{OUT}})$	$\tilde{\Omega}(N/\text{OUT}^{1/k})$	$\tilde{\Omega}(\min\{N, N^2/\text{OUT}\})$	$\tilde{\Omega}(N)$

- **Polynomial** speed-ups across all three query classes.
- Counting LBs via **communication complexity** — rule out *any* algorithm, incl. fast matrix mult.

Matrix Query: Two-Step Framework

$$Q_{\text{matrix}} = \pi_{A,C}(R_1(A,B) \bowtie R_2(B,C))$$

SampleMatrix — repeat until acceptance

1. **Step 1.** Draw $(a, b, c) \in R_1 \bowtie R_2$ uniformly at random.
2. **Step 2.** Accept (a, c) as the final sample with probability $\frac{1}{\text{deg}(a, c)}$,
where $\text{deg}(a, c) = |\pi_B(R_1 \bowtie a) \cap \pi_B(R_2 \bowtie c)|$.

Matrix Query: Two-Step Framework

$$Q_{\text{matrix}} = \pi_{A,C}(R_1(A,B) \bowtie R_2(B,C))$$

SampleMatrix — repeat until acceptance

1. **Step 1.** Draw $(a, b, c) \in R_1 \bowtie R_2$ uniformly at random.
2. **Step 2.** Accept (a, c) as the final sample with probability $\frac{1}{\text{deg}(a, c)}$,
where $\text{deg}(a, c) = |\pi_B(R_1 \bowtie a) \cap \pi_B(R_2 \bowtie c)|$.

Why uniform? A query result (a, c) appears in the full-join sample with probability $\frac{\text{deg}(a, c)}{\text{OUT}_{\bowtie}}$, hence

$$\Pr[\text{output } (a, c)] = \underbrace{\frac{\text{deg}(a, c)}{\text{OUT}_{\bowtie}}}_{\text{Step 1}} \cdot \underbrace{\frac{1}{\text{deg}(a, c)}}_{\text{Step 2}} = \frac{1}{\text{OUT}_{\bowtie}} \quad \text{— independent of } (a, c). \quad \checkmark$$

Matrix Query: Two-Step Framework

$$Q_{\text{matrix}} = \pi_{A,C}(R_1(A,B) \bowtie R_2(B,C))$$

SampleMatrix — repeat until acceptance

1. **Step 1.** Draw $(a,b,c) \in R_1 \bowtie R_2$ uniformly at random.
2. **Step 2.** Accept (a,c) as the final sample with probability $\frac{1}{\text{deg}(a,c)}$,
where $\text{deg}(a,c) = |\pi_B(R_1 \bowtie a) \cap \pi_B(R_2 \bowtie c)|$.

Why uniform? A query result (a,c) appears in the full-join sample with probability $\frac{\text{deg}(a,c)}{\text{OUT}_{\bowtie}}$, hence

$$\Pr[\text{output } (a,c)] = \underbrace{\frac{\text{deg}(a,c)}{\text{OUT}_{\bowtie}}}_{\text{Step 1}} \cdot \underbrace{\frac{1}{\text{deg}(a,c)}}_{\text{Step 2}} = \frac{1}{\text{OUT}_{\bowtie}} \quad \text{— independent of } (a,c). \quad \checkmark$$

(i) Step 1: Easy. (ii) Step 2 *without* computing $\text{deg}(a,c)$? **The challenge.**

Step 1: Sampling a Full-Join Tuple in $O(1)$

Auxiliary index (built in $O(N)$): for each $b \in \text{adom}(B)$, store

$$W_b = |R_1 \times b| \cdot |R_2 \times b|, \quad W = \sum_{b \in \text{adom}(B)} W_b = \text{OUT}_{\bowtie}.$$

JoinSampleMatrix-H

1. Sample $b \in \text{adom}(B)$ with probability W_b/W .
2. $a \leftarrow$ random A -value in $R_1 \times b$.
3. $c \leftarrow$ random C -value in $R_2 \times b$.
4. Return (a, b, c) .

Correctness:

$$\Pr[(a, b, c)] = \frac{W_b}{W} \cdot \frac{1}{|R_1 \times b|} \cdot \frac{1}{|R_2 \times b|} = \frac{1}{W} = \frac{1}{\text{OUT}_{\bowtie}}.$$

\Rightarrow uniform sample in $O(1)$ time. \checkmark

Step 2: Acceptance *without* Computing $\deg(a, c)$

Goal: accept (a, c) with probability $1/\deg(a, c)$ without enumerating all common B -witnesses.

Simulation Lemma

If $Y \in [M]$ and U is uniform on $[M]$, then

$$\Pr[U \leq Y] = \frac{\mathbf{E}[Y]}{M}.$$

Step 2: Acceptance *without* Computing $\deg(a, c)$

Goal: accept (a, c) with probability $1/\deg(a, c)$ without enumerating all common B -witnesses.

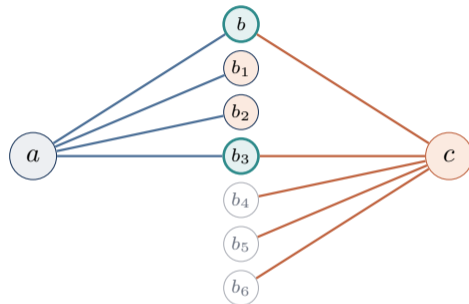
Simulation Lemma

If $Y \in [M]$ and U is uniform on $[M]$, then

$$\Pr[U \leq Y] = \frac{\mathbf{E}[Y]}{M}.$$

Constructing Proxy r.v.

Sample w/o replacement the smaller neighbor list until the next overlap or exhaust. Let F be the number of non-overlap values before stop.



Step 2: Acceptance *without* Computing $\deg(a, c)$

Goal: accept (a, c) with probability $1/\deg(a, c)$ without enumerating all common B -witnesses.

Simulation Lemma

If $Y \in [M]$ and U is uniform on $[M]$, then

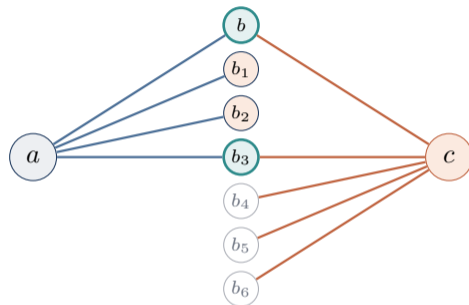
$$\Pr[U \leq Y] = \frac{\mathbf{E}[Y]}{M}.$$

Constructing Proxy r.v.

Sample w/o replacement the smaller neighbor list until the next overlap or exhaust. Let F be the number of non-overlap values before stop.

Key calculation: $F \sim \text{NegHypGeo}(1, |S| - 1, \deg(a, c) - 1)$, so $\mathbf{E}[F + 1] = \frac{|S|}{\deg(a, c)}$.

$$\Pr[\text{RandInt}(1, |S|) \leq F + 1] = \frac{\mathbf{E}[F + 1]}{|S|} = \frac{1}{\deg(a, c)}. \quad \checkmark$$



Matrix Query: $\tilde{O}(N/\sqrt{\text{OUT}})$ Analysis

Step-2 cost per attempt: $\mathbf{E}[F + 1] = \frac{|S|}{\text{deg}(a, c)} = \frac{\min\{|R_1 \times a|, |R_2 \times c|\}}{\text{deg}(a, c)}.$

Per (Step 1 + Step 2) iteration:

$$1 + \sum_{(a,c) \in \mathcal{Q}_{\text{matrix}}(\mathcal{R})} \frac{\text{deg}(a, c)}{\text{OUT}_{\boxtimes}} \cdot \frac{\min\{|R_1 \times a|, |R_2 \times c|\}}{\text{deg}(a, c)} \leq 1 + \frac{N\sqrt{\text{OUT}}}{\text{OUT}_{\boxtimes}}.$$

(using $\min\{x, y\} \leq \sqrt{xy}$ and Cauchy–Schwarz.)

Matrix Query: $\tilde{O}(N/\sqrt{\text{OUT}})$ Analysis

Step-2 cost per attempt: $\mathbf{E}[F + 1] = \frac{|S|}{\text{deg}(a, c)} = \frac{\min\{|R_1 \times a|, |R_2 \times c|\}}{\text{deg}(a, c)}.$

Per (Step 1 + Step 2) iteration:

$$1 + \sum_{(a,c) \in \mathcal{Q}_{\text{matrix}}(\mathcal{R})} \frac{\text{deg}(a, c)}{\text{OUT}_{\boxtimes}} \cdot \frac{\min\{|R_1 \times a|, |R_2 \times c|\}}{\text{deg}(a, c)} \leq 1 + \frac{N\sqrt{\text{OUT}}}{\text{OUT}_{\boxtimes}}.$$

(using $\min\{x, y\} \leq \sqrt{xy}$ and Cauchy–Schwarz.) **Iterations needed** (acceptance prob. = $\text{OUT}/\text{OUT}_{\boxtimes}$): $\text{OUT}_{\boxtimes}/\text{OUT}.$

Total expected time:

$$\frac{\text{OUT}_{\boxtimes}}{\text{OUT}} \cdot \left(1 + \frac{N\sqrt{\text{OUT}}}{\text{OUT}_{\boxtimes}}\right) = O\left(\frac{\text{OUT}_{\boxtimes}}{\text{OUT}} + \frac{N}{\sqrt{\text{OUT}}}\right) = O\left(\frac{N}{\sqrt{\text{OUT}}}\right).$$

(using $\text{OUT}_{\boxtimes} \leq N\sqrt{\text{OUT}}$ for $\mathcal{Q}_{\text{matrix}}$, Amossen–Pagh 2009.)

Theorem

For $\mathcal{Q}_{\text{matrix}}$: $O(N)$ preprocessing, $O(N/\sqrt{\text{OUT}})$ expected sampling time.

Matrix: Matching Lower Bounds

Thm A: Counting LB — info-theoretic

Any algorithm for $O(1)$ -approx counting on $\mathcal{Q}_{\text{matrix}}$ needs $\tilde{\Omega}(N/\sqrt{\text{OUT}})$ time.

Reduction from 2-party Set Disjointness.

Alice has S_A , Bob has $S_B \subseteq [m]$; decide $S_A \cap S_B \stackrel{?}{=} \emptyset$
(needs $\Omega(m)$ comm.).

Hard instance. $|\text{adom}(A)| = |\text{adom}(C)| = \sqrt{K}$;

- Alice: $R_1 = \text{adom}(A) \times S_A$
- Bob: $R_2 = S_B \times \text{adom}(C)$

$\text{OUT} = K$ iff $S_A \cap S_B \neq \emptyset$; else 0. $N = \sqrt{K} m/2$.

$O(1)$ -approx count \Rightarrow decides disjointness \Rightarrow
 $\Omega(m) = \Omega(N/\sqrt{\text{OUT}})$.

\Rightarrow **Info-theoretic:** rules out *any* algorithm, incl. fast matrix mult.

Matrix: Matching Lower Bounds

Thm A: Counting LB — info-theoretic

Any algorithm for $O(1)$ -approx counting on $\mathcal{Q}_{\text{matrix}}$ needs $\tilde{\Omega}(N/\sqrt{\text{OUT}})$ time.

Reduction from 2-party Set Disjointness.

Alice has S_A , Bob has $S_B \subseteq [m]$; decide $S_A \cap S_B \stackrel{?}{=} \emptyset$ (needs $\Omega(m)$ comm.).

Hard instance. $|\text{adom}(A)| = |\text{adom}(C)| = \sqrt{K}$;

- Alice: $R_1 = \text{adom}(A) \times S_A$
- Bob: $R_2 = S_B \times \text{adom}(C)$

$\text{OUT} = K$ iff $S_A \cap S_B \neq \emptyset$; else 0. $N = \sqrt{K} m/2$.

$O(1)$ -approx count \Rightarrow decides disjointness \Rightarrow

$\Omega(m) = \Omega(N/\sqrt{\text{OUT}})$.

\Rightarrow **Info-theoretic:** rules out *any* algorithm, incl. fast matrix mult.

The info-theoretic counting LB also implies a sampling LB via reduction (Chen-Yi'20), but it only rules out simultaneously sublinear preprocessing and sublinear sampling — which does not cover our $O(N)$ -preprocessing algorithm. The matching sampling LB therefore requires the separate combinatorial-hardness argument.

Thm B: Sampling LB — conditional on comb. BMM

Assuming comb. hardness of Boolean mat. mult., no algorithm preprocesses in $O(N \cdot \text{OUT}^{1/2-\gamma})$ and samples in $O(N/\text{OUT}^{1/2+\gamma})$ for any $\gamma > 0$.

Reduction from listing. Combinatorial listing of $\mathcal{Q}_{\text{matrix}}(\mathcal{R})$ needs $\Omega(N\sqrt{\text{OUT}})$ time (Amossen–Pagh 2009).

Idea. A fast sampler + coupon-collector ($\tilde{O}(\text{OUT})$ samples) lists $\mathcal{Q}_{\text{matrix}}(\mathcal{R})$ in

$$\tilde{O}\left(N\text{OUT}^{\frac{1}{2}-\gamma} + \text{OUT} \cdot \frac{N}{\text{OUT}^{\frac{1}{2}+\gamma}}\right) = \tilde{O}(N\text{OUT}^{\frac{1}{2}-\gamma}),$$

beating $\Omega(N\sqrt{\text{OUT}})$. Contradiction.

\Rightarrow **Conditional:** matches our UB *only* against combinatorial algorithms.

Star and Chain: Extensions and Dichotomy

Star $Q_{\text{star}} = \pi_{A_1, \dots, A_k} (\bowtie_i R_i(A_i, B))$

Theorem (tight bound)

$\Theta(N/\text{OUT}^{1/k})$ sampling *and* counting time.

Algorithm. Same two-step framework: weight b by $W_b = \prod_i |R_i \times b|$, then accept via the simulation lemma on $\pi_B(R_{j^*} \times a_{j^*})$ for $j^* = \arg \min_i |R_i \times a_i|$.

Lower bounds.

- Counting: k -party set disjointness (info-theoretic)
- Sampling: combinatorial hardness of star listing (conditional)

Star and Chain: Extensions and Dichotomy

Star $Q_{\text{star}} = \pi_{A_1, \dots, A_k} (\boxtimes_i R_i(A_i, B))$

Theorem (tight bound)

$\Theta(N/\text{OUT}^{1/k})$ sampling *and* counting time.

Algorithm. Same two-step framework: weight b by $W_b = \prod_i |R_i \times b|$, then accept via the simulation lemma on $\pi_B(R_{j^*} \times a_{j^*})$ for $j^* = \arg \min_i |R_i \times a_i|$.

Lower bounds.

- Counting: k -party set disjointness (info-theoretic)
- Sampling: combinatorial hardness of star listing (conditional)

Chain $Q_{\text{chain}} = \pi_{A_1, A_{k+1}} (\boxtimes_i R_i(A_i, A_{i+1}))$

Theorem (sharp phase transition)

$$k = 3 : \Theta(\min\{N, N^2/\text{OUT}\})$$

$$k \geq 4 : \Theta(N)$$

Algorithm. For $k = 3$: combine prior $\tilde{O}(N^2/\text{OUT})$ with new $O(N)$ via the k -minimum-values (KMV) summary [Cohen 2008].

Both LBs are info-theoretic.

Approximate Counting: Heavy/Light Hybrid Reduction

Naive sampling-to-counting reduction [Chen-Yi'20] gives $\tilde{O}(N + N^{\rho^*}/\text{OUT}) = \tilde{O}(N)$ for $\mathcal{Q}_{\text{matrix}}$ — not sublinear in OUT .

Our refinement. Geometric guess $\Lambda \in [\text{OUT}, 2\text{OUT}]$ and partition:

$$B^{\text{heavy}} = \{b : |R_2 \times b| \geq \sqrt{\Lambda}\}, \quad B^{\text{light}} = \text{adom}(B) \setminus B^{\text{heavy}}.$$

Three estimators, then combine

- **Heavy:** SAMPLEMATRIX-H — efficient because $|B^{\text{heavy}}| = O(N/\sqrt{\Lambda})$.
- **Light:** SAMPLEMATRIX-L with rejection threshold $\Delta = \sqrt{\Lambda}$.
- **Intersection:** estimate $|R^{\text{heavy}} \cap R^{\text{light}}|$ via sampling.
- **Combine:** $\widehat{\text{OUT}} = s^{\text{heavy}} + s^{\text{light}} - \min\{s^{\text{heavy}}, \beta s^{\text{light}}\}$.

Theorem

$\tilde{O}(N/\sqrt{\text{OUT}})$ for ϵ -approx counting on $\mathcal{Q}_{\text{matrix}}$ — matches the info-theoretic LB.

Analogous partition (threshold $\text{OUT}^{1/k}$) gives $\tilde{O}(N/\text{OUT}^{1/k})$ for $\mathcal{Q}_{\text{star}}$; chain counting uses KMV.

Summary & Open Problems

Query	Sampling & Counting	Improvement	Lower-bound type
Matrix	$\Theta(N/\sqrt{\text{OUT}})$	poly.	Counting: info-th.; Sampling: comb.
k -Star	$\Theta(N/\text{OUT}^{1/k})$	poly.	Counting: info-th.; Sampling: comb.
3-Chain	$\Theta(\min\{N, N^2/\text{OUT}\})$	tight	Both: info-theoretic
k -Chain ($k \geq 4$)	$\Theta(N)$	tight (hard)	Both: info-theoretic

Two reusable techniques:

- **Rejection + simulation lemma** for output-optimal sampling.
- **Heavy/light hybrid reduction** for output-optimal counting.

Open problems: combining $\sigma + \pi$; other join-project query classes; other operators; dynamic settings under updates.

Thank you — questions?