

Accelerating Joins with Filters

Nicholas Corrado

Xiating Ouyang

Abstract

In query optimization on star schemas, lookahead information passing (LIP) is a strategy exploiting the efficiency of probing succinct filters to eliminate practically all facts that do not appear in the final join results before performing the actual join. Assuming data independency across all columns in the fact table, LIP achieves efficient and robust query optimization. We present LIP-k, a variant of LIP that only remembers the hit/miss statistics for the previous k batches, achieving empirically efficient query execution on fact table with correlated and even adversarial data columns. We implemented LIP and LIP-k on a skeleton database on top of Apache Arrow and analyze the performance of each variant of LIP using the notion of competitive ratio in online algorithms.

1 Introduction

Performing join operations in database management systems is a fundamental and prevalent task in the modern data industry. Continuous efforts have been spent on building a reliable query optimizer over the last few decades. However, the current optimizers may still produce disastrously inefficient query plans which involves processing unnecessarily gigantic intermediate tables [?, ?]. The *Lookahead Information Passing (LIP)* strategy aggressively uses Bloom Filters to filter the fact tables, effectively reducing the sizes of the intermediate tables. LIP is provably as efficient and robust as computing the join using the optimal query plan [?]. The key idea behind LIP is to estimate the filter selectivity of each dimension table and adaptively reorder the sequence of applying the filters to the fact table.

The filtering process can be modeled as the LIP problem in an online setting: Suppose we fix n filters, and the tuples in the fact table arrives in an online fashion. Upon arrival of each tuple, one has to decide a sequence of filters to probe the tuple, with an objective of minimizing the number of probes needed to decide whether to accept the tuple and forward it to the hash join phase, or to eliminate it. A mechanism deciding the sequence of applying the filters is thus crucial to the success of LIP. If a tuple passes all filters, *all* mechanisms must probe the tuple to all n filters to confirm its passage. If a tuple is eliminated by at least one filter, then the *optimal* mechanism would apply any filter that rejects the tuple first, using only one probe. Thus, given any fact table F , the number of probes that an optimal mechanism requires to process all tuples in the fact table can be readily computed:

$$\text{OPT}(F) = n|F_{\text{pass}}| + |F_{\text{reject}}|,$$

where $|F_{\text{pass}}|$ and $|F_{\text{reject}}|$ are the number of tuples in F that pass all filters and are rejected in F respectively. For any mechanism \mathcal{M} , denoted by $\text{ALG}_{\mathcal{M}}(F)$ the number of probes that \mathcal{M} performed to process all tuples in the fact table. The performance of any mechanism \mathcal{M} can thus be measured by

$$\max_F \frac{\text{ALG}_{\mathcal{M}}(F)}{\text{OPT}(F)},$$

called the *competitive ratio* of \mathcal{M} . The competitive ratio is always at least 1 by definition, and in this problem the competitive ratio is at most n , the number of filters, since one mechanism can probe each tuple to at most n filters.

This project aims at designing efficient LIP mechanisms and measure their performance in terms of their overall running time and competitive ratio. We will build a skeleton database system on top of Apache Arrow supporting LIP and hash-joins to conduct our experiments and test the performance of LIP- k —a variant mechanism of LIP that only remembers the statistics for the previous k batches—against the hash-join and the original LIP. We also present a theoretical result showing that no deterministic mechanism can have a competitive ratio better than n , and discuss possible extensions of LIP incorporating randomness to design a mechanism with better theoretical guarantee.

2 Lookahead Information Passing (LIP)

In this section, we first present the LIP strategy in [?]. We then discuss our variant LIP- k , designed to respond to local skewness more quickly than LIP. Finally, we discuss the competitive ratios of all deterministic mechanisms and provide proof on its lower bound. Possible extensions of LIP using randomness is also discussed.

2.1 LIP

The LIP strategy has three stages: (1) Build a hash table and a filter for each dimension table, (2) probe each fact tuple on the filters, producing a set of fact tuples with false positives, and (3) probe the hash table of each dimension table to eliminate the false positives. In what follows, we mainly discuss stage (2), since stages (1) and (3) are readily implemented by the filter constructors and database engine, respectively.

Let F be the fact table and D_i the dimension tables for $1 \leq i \leq n$. We denote the number of facts in F and each D_i as $|F|$ and $|D_i|$. Each LIP filter on D_i is a Bloom filter with false positive rate ϵ . The true selectivity σ_i of D_i on fact table F is given by $\sigma_i = |D_i \bowtie_{pk_i=fk_i} F|/|F|$, where pk_i is the primary key of D_i and fk_i is the foreign key of D_i in F . The LIP-join algorithm, depicted in Figure 1, computes the indices of tuples in F that pass all LIP filters. Note that because of the false positive rate ϵ associated with each filter, the set of indices is a superset of the true set of indices of tuples appearing in the final join result.

The partition in [?] satisfies that $|F_{t+1}| = 2|F_t|$ at line 5, and the algorithm approximates the true selectiveness σ_i of each dimension D_i using $\text{pass}[i]/\text{count}[i]$, the aggregated selectiveness since the beginning.

2.2 LIP- k

The LIP strategy in Figure 1 estimates the selectivity of each filter using statistics from all previous batches, which can be inefficient for certain foreign key distributions in the fact table. Consider the case where some filter f is very selective for the first T batches and not selective for the remaining batches. (For example, a filter f filtering for $\text{year} \geq 2017$ and the Date table is sorted in year.) In this case, LIP would obtain a good estimate of the selectivity of f during the first T iterations, and thus tend to apply f early in the remaining iterations. However, it is more efficient to postpone applying f in the remaining iterations, despite f has good selectivity in the first T iterations. One remedy to this is to only “remember” the hit/miss statistics of each filter over the previous k batches.

Empirical data shows that for the original SSB dataset and certain queries, LIP- k is as fast as LIP, and for certain datasets and queries LIP- k is faster than LIP. Detailed empirical data are presented and analyzed in Section 4.

```

PROCEDURE: LIP-join
INPUT: a fact table F and a set of n Bloom filters  $f_i$  for each  $D_i$  with  $1 \leq i \leq n$ 
OUTPUT: Indices of tuples in F that pass the filtering

1. Initialize  $I = \emptyset$ 
2. foreach filter  $f$  do
3.    $\text{count}[f] \leftarrow 0$ 
4.    $\text{pass}[f] \leftarrow 0$ 
5. Partition  $F = \bigcup_{1 \leq t \leq T} F_t$ .
6. foreach fact block  $F_t$  do
7.   foreach filter  $f$  in order do
8.     foreach index  $j \in F_t$  do
9.        $\text{count}[f] \leftarrow \text{count}[f] + 1$ 
10.      if  $f$  contains  $F_t[j]$ 
11.         $I \leftarrow I \cup \{j\}$ 
12.       $\text{pass}[f] \leftarrow \text{pass}[f] + 1$ 
13. sort filters  $f$  in nondesending order of  $\text{pass}[f]/\text{count}[f]$ 
14. return  $I$ 

```

Figure 1: The LIP algorithm for computing the joins.

2.3 Competitive Ratio Analysis

The LIP strategy and its variant LIP-k depicted in Figures 1 and 2 are *deterministic*, i.e. multiple executions over the same fact table would produce the same result. Experimental results show that LIP has faster execution time compared to hash join in the optimal sequence [?] on the benchmark dataset, in which the keys are distributed almost uniformly. However, it can be shown that any deterministic mechanism in the worst case can never achieve a competitive ratio less than n when played against an adversary producing an adversarial dataset.

Theorem 2.1. *Let n be the number of filters in the LIP problem. There is no deterministic mechanism \mathcal{M} achieving a competitive ratio less than n for the LIP problem.*

Proof. We present an adversary to the arbitrary mechanism \mathcal{M} such that \mathcal{M} only achieves a competitive ratio of n in the worst case. Let the n filters be f_1, f_2, \dots, f_n , and let S_k denote the filter sequence that will be used to filter batch k . Let each batch contain m tuples. At each iteration k , prior to LIP's probe phase, the adversary observes S_k and produces a batch of tuples $\{t_1, t_2, \dots, t_m\}$, where $t_i \in f_n$ but $t_i \notin f_j$ for any $j < n$. The adversary then feeds this batch into mechanism \mathcal{M} . Thus, \mathcal{M} performs n filter probes to eliminate each tuple, whereas the optimal sequence is to apply f_n first. Thus, \mathcal{M} achieves a competitive ratio of n . \square

It might be possible to design a randomized mechanism that can achieve a better competitive ratio than n . The randomized mechanism would, at the end of each batch, select a sequence of applying the filters from a distribution of all filter permutations, based on the estimated selectivities. However, we have not obtained any algorithmic upper bound on the competitive ratio. Furthermore, a randomized approach creates a new concern of being too computationally heavy.

```

PROCEDURE: LIP-k
INPUT: a fact table F and a set of n Bloom filters  $f_i$  for each  $D_i$  with  $1 \leq i \leq n$ 
OUTPUT: Indices of tuples in F that pass the filtering

1. Initialize  $I = \emptyset$ 
2. foreach filter  $f$  do
3.   Initialize  $\text{count}[f] \leftarrow 0$ ,  $\text{pass}[f] \leftarrow 0$ 
4.   Initialize  $\text{count\_queue}[f]$  with k zeros and  $\text{pass\_queue}[f]$  with k zeros.
5. Partition  $F = \bigcup_{1 \leq t \leq T} F_t$ .
6. foreach fact block  $F_t$  do
7.   foreach filter  $f$  in order do
8.     foreach index  $j \in F_t$  do
9.        $\text{count}[f] \leftarrow \text{count}[f] + 1$ 
10.      if  $f$  contains  $F_t[j]$ 
11.         $I \leftarrow I \cup \{j\}$ 
12.         $\text{pass}[f] \leftarrow \text{pass}[f] + 1$ 
13.       $\text{count\_queue}[f].\text{dequeue}()$  and  $\text{pass\_queue}[f].\text{dequeue}()$ 
14.       $\text{count\_queue}[f].\text{enqueue}(\text{count}[f])$  and  $\text{pass\_queue}[f].\text{enqueue}(\text{pass}[f])$ 
15.      Reset  $\text{count}[f] \leftarrow 0$ ,  $\text{pass}[f] \leftarrow 0$ 
16. sort filters  $f$  in nondesending order of  $\text{sum}(\text{pass\_queue}[f])/\text{sum}(\text{count\_queue}[f])$ 
17. return  $I$ 

```

Figure 2: The LIP algorithm for computing the joins.

In the practical perspective however, one wishes to minimize the total running time of the mechanism, which is effectively the sum of the running time of the mechanism and the running time of building the filters and performing the probes. A trade-off between having a near optimal mechanism that consumes much time and allowing many failed probes to eliminate each non-participating tuple is therefore of much interest.

3 Database Implementation

We have developed a prototype database system supporting basic select and join operations on top of Apache Arrow [?], a column-store format. This minimal prototype is sufficient to benchmark the performance of our Hash join, LIP, and LIP-k.

Our implementations of LIP and LIP-k only support left-deep join tree plans where the fact table is the “outer table” in every join. We assume that the fact table schema contains foreign keys to all dimension tables, and each dimension table is single-key. Given a star schema fact table F and dimension tables D_i for $1 \leq i \leq n$, a join query in our system specifies selectors σ_F for F and σ_i for each D_i , and executing that query will return

$$\sigma_F(F) \bowtie \sigma_1(D_1) \bowtie \dots \bowtie \sigma_n(D_n)$$

projected on the schema of F , *i.e.* we output the tuples in F that can be joined with each D_i . The supported

primitive selectors allow for selection ($=, \leq, \geq, <, >$) on scalar values and ranges. Range selections are executed using BETWEEN (ℓ, h), which selects all x with $\ell \leq x \leq h$.

The selectors for each dimension can be either a primitive selector consisting of simple predicate (e.g. ORDER DATE = 1997) or a composition (logical AND/OR) of multiple primitive/composite selectors. This is implemented using the Composite design pattern. Apache Arrow does not yet support vectorized string comparison operations nor vectorized range comparison operations. For queries involving string and/or range selections, we instead scan along the column, checking which rows satisfy the selection predicate. Such selections are inherently slower than the supported vectorized selections. Because all of our join implementations must implement row-wise selection for string and range predicates, all algorithms suffer the same slowdown. Thus, this implementation caveat does not preclude us from studying the relative performance of LIP and LIP-k.

We only support string and numeric data types. All numeric data is stored as 64-bit integers by default.

The hash join algorithm first produces a hash table T_i for each $\sigma_i(D_i)$, projected on the k_i , and then probe each tuple in the fact table against all T_i . We used Sparseepp (accessible at <https://github.com/greg7mdp/sparsepp>) as our implementation of the hash table, in which the sparsehash by Google (accessible at <https://github.com/sparsehash/sparsehash>) is used as the underlying hash function. All primary keys are regarded as 64-bit integers.

For LIP and LIP-k, the succinct filter structure we choose is the Bloom filters. The default false-positive rate is set to 0.001, which requires 10 hash functions. We use Knuth’s Multiplicative hash function, extended to accept a 64-bit integer as a seed. Our minimal implementation does not support selectivity estimation, so we initialize each Bloom filter assuming $\sigma(D_i) = \frac{1}{2}$, i.e. assuming half of the keys in each dimension table will be inserted into the filter.

Our code is available at <https://github.com/NicholasCorrado/CS764>.

4 Empirical Results

In this section we present the empirical results obtained from running Hash-join, LIP and LIP-k on several datasets. In Section 4.1, we present the datasets we use and describe how we generate skewed and adversarial datasets. In Section 4.2 we present the running time of multiple strategies and discuss their performance. In Section 4.3 we discuss how k affects the competitive ratio of LIP-k empirically on our datasets.

4.1 Datasets

SSB query groups 1 and 2 do not select on DATE.YEAR and are thus unaffected by this skew. Query groups 3 and 4 do select on DATE.YEAR, so this skew can potentially affect the execution of LIP and LIP-k on these queries.

4.1.1 Skew Datasets

We modify the normal LINEORDER file produced by the SSB data generator with SF = 1 to produce four skewed LINEORDER tables:

- LINEORDER-DATE-50-50: The first 50 batches have $\sigma_{\text{SKEW PRED}} = 1$ while the next 50 batches have $\sigma_{\text{SKEW PRED}} = 0$, and so on. In other words, The first 50 batches contain only keys satisfying SKEW PRED, while the next 50 batches contain only keys not satisfying SKEW PRED, and so on.

- **LINEORDER-DATE-FIRST-HALF:** Let N be the total number of batches in **LINEORDER**. The first $N/2$ batches have $\sigma_{\text{SKEW PRED}} = 1$ while the remaining $N/2$ batches have $\sigma_{\text{SKEW PRED}} = 0$.
- **LINEORDER-DATE-LINEAR:** $\sigma_{\text{SKEW PRED}}$ increases linearly from 0 to 1 across the **LINEORDER** table, *i.e.* batch k has $\sigma_{\text{SKEW PRED}} = k/N$.

LINEORDER-DATE-PART-ADVERSARY.

For simplicity in generating skew datasets, we skew the **ORDER DATE** foreign key in the **LINEORDER** table, corresponding to the primary key of the **DATE** table. We alter the distribution of **ORDER DATE** foreign keys satisfying **DATE.YEAR** = 1997 OR **DATE.YEAR** = 1998, which we will shorten to **SKEW PRED**. Let

Let f_i^A denote the selectivity of a Bloom filter on key A after processing the i^{th} batch. Let σ_i^A denote the selectivity of a Bloom filter on key A on the i^{th} batch alone.

Skewed datasets are generated by first generating the SSB dataset with $SF = 1$. We then edit the columns we want to skew.

The aim of the adversarial dataset is to force LIP to perform the maximum number of filter probes possible, mimicking the proof of Theorem 2.1. We now show how to construct such a dataset where $n = 2$. Let subscripts denote the batch index and superscripts denote the key column (A or B). We start with the first batch having $\sigma_1^A = \frac{1}{2} - \epsilon$ and $\sigma_1^B = \frac{1}{2} + \epsilon$ where $0 < \epsilon < \frac{1}{2}$. Then for all $j > 1$, we let

$$\sigma_j^A = \begin{cases} 1 & \text{if } j \text{ is even,} \\ 0 & \text{if } j \text{ is odd} \end{cases} \quad \text{and} \quad \sigma_j^B = \begin{cases} 0 & \text{if } j \text{ is even,} \\ 1 & \text{if } j \text{ is odd.} \end{cases}$$

Thus, the optimal filter sequence S^{OPT} is

$$S^{\text{OPT}} = \begin{cases} (B, A) & \text{if } j \text{ is even,} \\ (A, B) & \text{if } j \text{ is odd.} \end{cases}$$

After processing batch j , we have $f_j^A = \frac{\frac{1}{2} - \epsilon + \lfloor \frac{j}{2} \rfloor}{j}$ and $f_j^B = \frac{\frac{1}{2} + \epsilon + \lfloor \frac{j}{2} \rfloor}{j}$ which can be rewritten as

$$f_j^A = \begin{cases} \frac{1}{2} + \frac{\frac{1}{2} - \epsilon}{j} & \text{if } j \text{ is even,} \\ \frac{1}{2} - \frac{\epsilon}{j} & \text{if } j \text{ is odd;} \end{cases} \quad \text{and} \quad f_j^B = \begin{cases} \frac{1}{2} - \frac{\frac{1}{2} - \epsilon}{j} & \text{if } j \text{ is even,} \\ \frac{1}{2} + \frac{\epsilon}{j} & \text{if } j \text{ is odd,} \end{cases}$$

and thus, LIP's filter ordering will be

$$S = \begin{cases} (A, B) & \text{if } j \text{ is even,} \\ (B, A) & \text{if } j \text{ is odd,} \end{cases}$$

which is the reverse of S^{OPT} . Hence, after the first batch has been processed, LIP will have worst-case performance on all remaining batches.

We summarize the skew datasets studied in this paper in Table ??

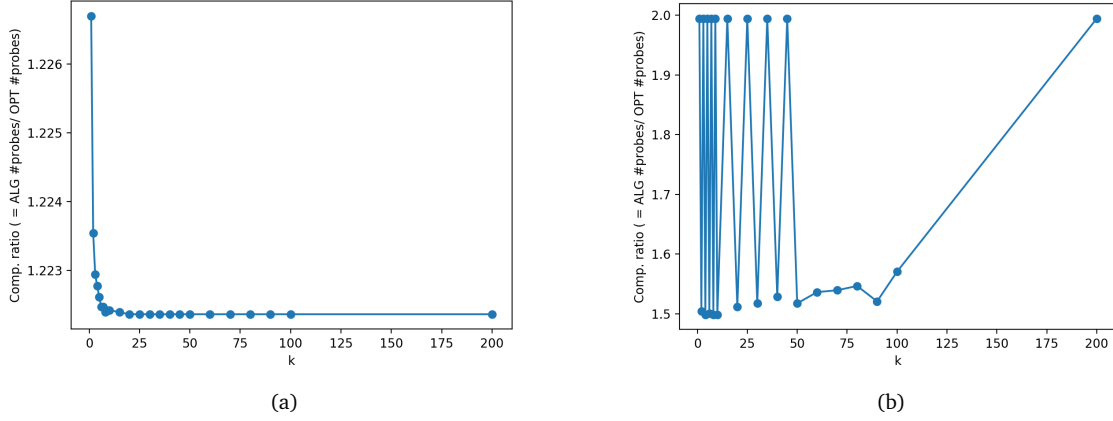


Figure 3: The competitive ratios of LIP-k against different k values. We ran LIP-k on uniform data and skewed (and adversarial) dataset to produce (a) and (b) respectively. The data point at $k = 200$ represents LIP (which is essentially LIP- ∞).

Dataset Name	Skewed Foreign Keys	Description of Skew
lineorder-date-50-50	ORDER DATE	First 50 batches have $\sigma_{\text{SKEW PRED}} = 1$ Next 50 batches have $\sigma_{\text{SKEW PRED}} = 0$
lineorder-date-first-half	ORDER DATE	First $N/2$ batches have $\sigma_{\text{SKEW PRED}} = 1$ Next $N/2$ batches have $\sigma_{\text{SKEW PRED}} = 0$, where N is the total number of batches in LINEORDER
lineorder-date-linear	ORDER DATE	$\sigma_{\text{SKEW PRED}}$ increases linearly from 0 to 1
lineorder-date-part-adversary	ORDER DATE	See Section 4.3

The adversarial dataset is constructed to make SSB query 4.2 achieve near worst-case performance for LIP.

4.2 Execution Time

4.3 Competitive Ratio

We ran LIP-k on each skewed dataset we produced and compute the competitive ratio of each LIP-k by taking the maximum of the competitive ratios achieved across all queries in all datasets. The results are depicted in Figure 7.

When the keys in the fact table columns are distributed uniformly, the filters need not react to the local changes. Hence for LIP with higher k , it remembers more batches in the uniform data, therefore may produce relatively more accurate estimate of the selectivities than the LIP with lower k . Hence the competitive ratio would decrease (slightly) when k increases, as depicted in Figure 7.

Figure 7 (b) displays how an adversarial dataset can make LIP-k and LIP perform inefficiently. In fact, each LIP-k with even k achieves an approximation ratio less than 2; LIP and LIP-k with odd k achieve an approximation ratio of almost 2 precisely at Query 3.2 in dataset date-part-adversary. Query 3.2 has two joins, and thus the performance of LIP-k with even k matches the worst case competitive ratio. By construction of date-part-adversary,

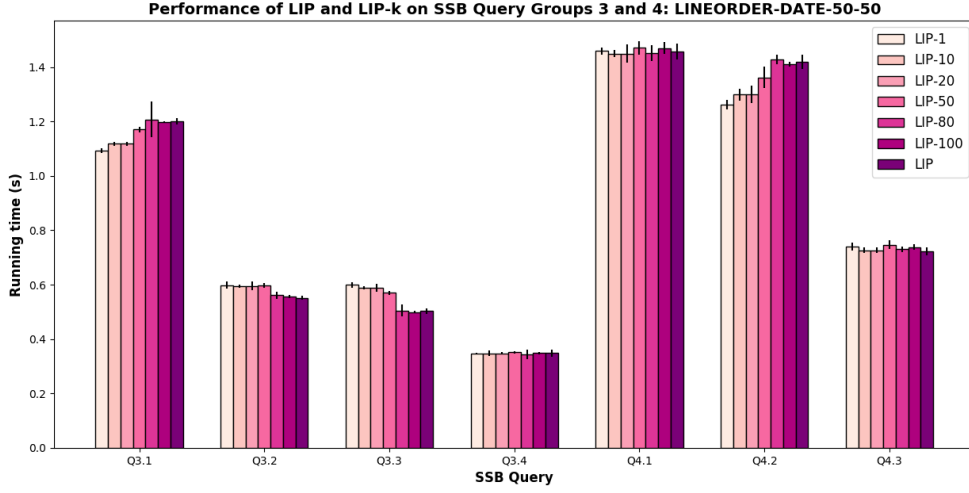


Figure 4

the first batch has $\sigma_1^1 = 1/2 - \varepsilon$ and $\sigma_2^1 = 1/2 + \varepsilon$, and each batch B_{2i+1} is an *odd* batch, with $\sigma_1^{2i} = 0$ and $\sigma_2^{2i} = 1$; and each batch B_{2i} is an *even* batch, with $\sigma_1^{2i} = 1$ and $\sigma_2^{2i} = 0$.

When $i \leq k$, LIP and LIP- k execute identically, since LIP- k has not yet “forgotten” any previous batches. We now consider the case where $i > k$, *i.e.* where LIP- k has forgotten at least the first batch. For odd k , LIP- k ’s selectivity estimates always contains one more odd (or even) batch than the other, and thus the estimated selectivity and filtering sequence are in favor of the majority batch type. This yields a false prediction of the next batch, yielding a competitive ratio of 2. For even k , LIP- k ’s selectivity estimates contain an equal amount of even an odd batches, and thus estimated selectivities remain the same ($1/2$ by construction) throughout the execution, and the filter sequence does not change. Thus for half batches it is optimal, and for the other half it is worse, resulting in a competitive ratio of

$$\frac{1 \times 1/2 + 2 \times 1/2}{1} = 1.5,$$

as depicted in Figure 7. For LIP, since it remembers the statistics from the beginning, by construction it would also make the same decision as LIP- k with odd k , losing at every batch.

LIP performs poorly because it never forgets the first batch. ¹

$$\varepsilon_i = \begin{cases} \frac{\varepsilon}{2k+1} & \text{for odd } i, \text{ where } i = 2k + 1 \\ \frac{\frac{1}{2} - \varepsilon}{2k} & \text{for even } i, \text{ where } i = 2k \end{cases}$$

If k is even, then LIP- k will (after the first k batches have been processed) have $\sigma_i^A = \sigma_i^B = \frac{1}{2}$. Thus,

¹An astute observer might notice that for $k = 200$, the worst-case competitive ratio of 2 is achieved even though k is not odd. This is explained by recognizing that LIP-200 is essentially an approximation of LIP. LIP-200 performs just as poorly as LIP up until the 201st batch, after which it estimates both selectivities as $1/2$.

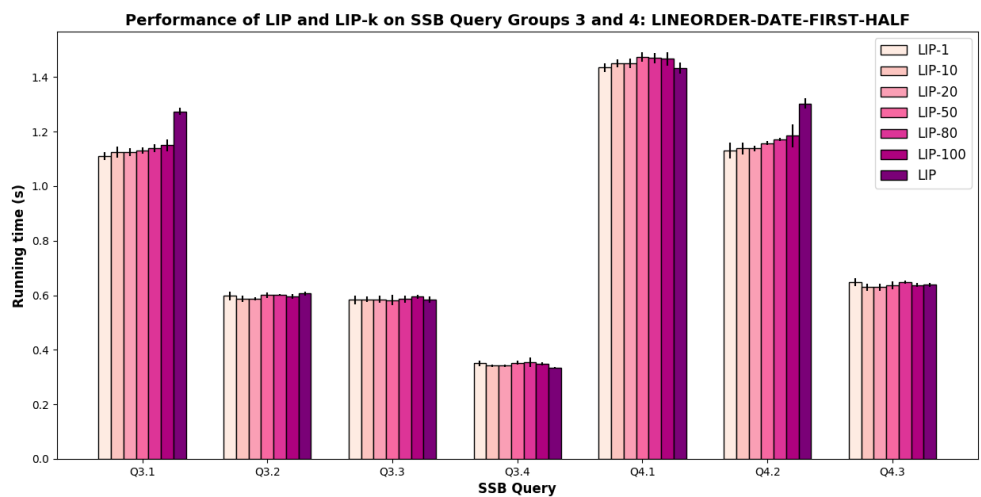


Figure 5

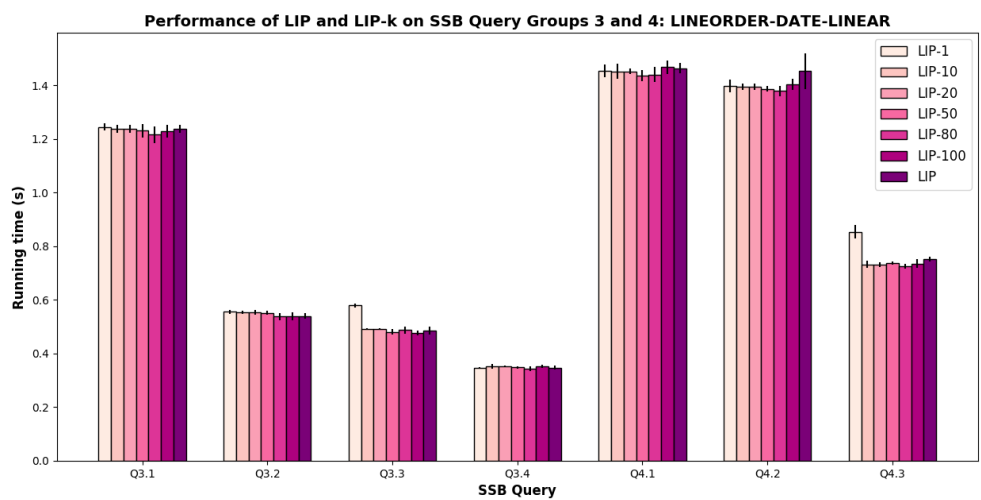


Figure 6

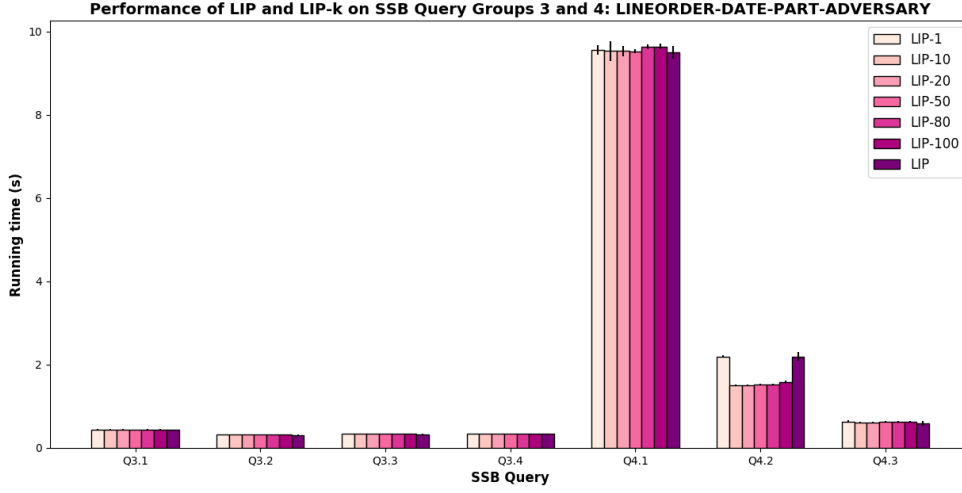


Figure 7

5 Future Works

A huge take-away from this project is that the budget allowed for filtering each filter is very limited, in the magnitude of near 100 CPU cycles. Hence maintaining the statistics that require heavy computation or using randomness that requires much computation cost are not practically beneficial. Given this, we will investigate dynamically start sampling a segment of the fact table to have true estimates of each filters, and then stick to this estimate for the next few batches. This has low maintenance overhead, and once the adaptive policy for dynamical sampling is established, we expect that to have better practical performance than LIP and LIP-k. Note that this strategy is still deterministic, still prone to the worst case n competitive ratio bound established by Theorem 2.1.

6 Concluding Remarks

The contributions of this project come in two-folds: We implemented Hash-join and LIP on top of Apache Arrow and thus provided an interface for future integration of Hustle on Apache Arrow; and we proposed LIP-k, a variant of LIP, which opens up an area for improving the performance LIP. It would also be interesting to study the problem of LIP in the online algorithmic setting to design an efficient mechanism utilizing randomness with better worst case guarantees.

Acknowledgements

The authors wish to thank Prof. Jignesh Patel for constant feedbacks on this project and Kevin Gaffney for helping us with Apache Arrow specifics. The second author wishes to thank Prof. Paris Koutris for the suggestion of working on a practical project when the lemma production pipe is jammed. It works.