

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 using Star Schemas 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 [1, 4]. The *lookahead Information Passing (LIP)* strategy aggressively uses Bloom Filters to filter the fact tables to effectively reduce the sizes of the intermediate tables, provably as efficient and robust as computing the join using the optimal query plan [5]. 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 have to probe the tuple to all filters to confirm its passage; and if a tuple is eliminated by the filters, the *optimal* mechanism would apply any filter that rejects the tuple in the first place, 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 our variant LIP mechanisms 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 to use randomness to design a better mechanism.

2 Lookahead Information Passing (LIP)

In this section, we first present the LIP strategy in [5], and 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. Some possible extensions of LIP using randomness is also discussed.

2.1 LIP

The LIP strategy has three stages: (1) Building 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 either the database engine or the filter constructors.

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 $|D_i|$. A LIP filter on D_i is implemented using 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 the filtering of each Bloom filter of D_i . Note that there is an innate false positive rate ε associated with each Bloom Filter, and thus the set of indices is a superset of the true set of indices of tuples appearing in the final join result.

The partition in [5] 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

LIP strategy in Figure 1 estimates the selectivity of each filter using statistics from the beginning batch, which is inefficient for certain distribution or physical layout of data. Consider some filter f that is very selective for the first t_0 iterations at line 6 and not selective for the remaining iterations. (For example, a filter f filtering for year ≥ 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_0 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_0 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 Figure 1 and 2 are *deterministic*, i.e. multiple executions over the same fact table would produce the same result. Experimental results show that LIP performs almost optimally compared to the performance of hash join in the optimal sequence [5] on the benchmark dataset, in which the keys are distributed almost uniformly. However, it can be shown that 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 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 consider n tuples t_1, t_2, \dots, t_n , where $t_i \notin f_i$ but $t_i \in f_j$ for any $i \neq j$.

The adversary proceeds as follows: It first observes the sequence of filters σ_t at any step t set by the mechanism \mathcal{M} , and produce the input $f_{\sigma_t(n)}$ to the mechanism \mathcal{M} . Thus the mechanism \mathcal{M} would require n filter probes to eliminate $f_{\sigma_t(n)}$ at each step t , whereas the optimal sequence is to apply $\sigma_t(n)$ at the first place. Thus it yields a competitive ratio of n . \square

It might be possible to design a randomized mechanism $\mathcal{M}_{\sqrt{\cdot}}$ 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.

```

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.      Dequeue one element from both  $\text{count\_queue}[f]$  and  $\text{pass\_queue}[f]$ 
14.      Enqueue  $\text{count}[f]$  and  $\text{pass}[f]$  to  $\text{count\_queue}[f]$  and  $\text{pass\_queue}[f]$  respectively
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 \mathcal{M} , 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 join/select operations on star schemas sufficient to benchmark the performance of LIP and its variants on top of Apache Arrow. We assume that the fact table schema contains foreignkeys 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 .

The supported primitive selectors are comparison with integer/string values ($=, \leq, \geq, <, >$) and between, where the semantic of BETWEEN (ℓ, h) is to select all x with $\ell \leq x \leq h$. The selectors for each dimension can be either a primitive selector, or a composition (logical AND and OR) of multiple primitive/composite selectors. This is implemented using the Composite design pattern.

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.

The succinct filter structure we choose is the Bloom filters. The default false-positive rate is set to 0.001 and the default number of inserts to the filter is set to 50,000. The hash function for the Bloom Filter is Knuth’s Multiplicative hash function, extended to accept a 64-bit integer as a seed.

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 dataset we used and how we manually apply skewness to generate skewed and adversarial datasets; in Section 4.2 we present the running time of multiple strategies and discuss their performance; and in Section 4.3 we discuss how k affects the competitive ratio of LIP-k empirically on our dataset.

4.1 Skewed and Adversarial Datasets

The dataset for testing is obtained from the Star Schema Benchmark [3]. We will hard-code each queries in [3] to measure the join processing time.

4.2 Execution Time

4.3 Competitive Ratio

5 Concluding remarks

Acknowledgement

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.

References

- [1] Viktor Leis, Andrey Gubichev, Atanas Mirchev, Peter Boncz, Alfons Kemper, and Thomas Neumann. How good are query optimizers, really? *Proceedings of the VLDB Endowment*, 9(3):204–215, 2015.
- [2] Nick Littlestone and Manfred K Warmuth. The weighted majority algorithm. *Information and computation*, 108(2):212–261, 1994.
- [3] Patrick O’Neil, Elizabeth O’Neil, Xuedong Chen, and Stephen Revilak. The star schema benchmark and augmented fact table indexing. In *Technology Conference on Performance Evaluation and Benchmarking*, pages 237–252. Springer, 2009.

- [4] Tilmann Rabl, Meikel Poess, Hans-Arno Jacobsen, Patrick O’Neil, and Elizabeth O’Neil. Variations of the star schema benchmark to test the effects of data skew on query performance. In *Proceedings of the 4th ACM/SPEC International Conference on Performance Engineering*, pages 361–372. ACM, 2013.
- [5] Jianqiao Zhu, Navneet Potti, Saket Saurabh, and Jignesh M Patel. Looking ahead makes query plans robust: Making the initial case with in-memory star schema data warehouse workloads. *Proceedings of the VLDB Endowment*, 10(8):889–900, 2017.