ORACLE

MySQL HeatWave: A Deep Dive Into Optimizations Robust OLAP Query Performance At Scale

Kaan Kara, PhD Principal Member of Technical Staf

MySQL HeatWave

OLTP, real-time analytics, and ML in one cloud database service – without ETL



How do we scale-out query processing to hundreds of nodes and thousands of cores efficiently?

- 1. Main Principal in Execution: Data Partitioning
- 2. Query Optimization: A Holistic Approach

MySQL HeatWave Query Performance

Significant Advantage in Perf and Price/Perf



https://github.com/oracle/heatwave-tpch

0

In-Memory Hybrid Columnar Format

OLAP Performance Fundamentals

chunk M chunk I vector 1 vector N vector N vector N vector N vector N			vector	Multi-core scalability				izontal scalability			
E Contraction of the second se	→ tile 1	vector N	•••		vector 3		vector 2		vector 1	chunk 1	
			:							chunk M	
column 1 column 2 column 3 colum	nn N	column N	:	column 3		column 2		column 1			

High fanout partitioning is the core aspect of high performance.

Distributed Relational Join Processing



How much do we need to partition?



Back of the envelope

- Orders Table: key column is 8B integers, 768 Billion rows: 6.1TB
- 6.1TB / 64KB = ~96 Million Partitions
- Next power of two: 2^27, a partitioning fanout of 134 Million

Data Partitioning Problem: Where is the bottleneck?

foreach tuple t in relation R:

p = hash(t.key)
partitions[p][counts[p]++] = t

- p is random and will cause heavy random access.
- In case the range of p is high (high fanout), chances of TLB misses are high.

TLB (Translation Lookaside Buffer)

- A specialized cache for virtual-to-physical address translation.
- If a virtual address is not found, expensive page walk occurs.
 - Can be even more expensive than simple cache miss, because page walk might perform multiple memory accesses.
- Typical TLB has around 64 entries.
- Any partitioning fanout larger than 64 is likely to use more than 64 pages, causing a huge performance impact.





Reducing the TLB bottleneck with SW-managed buffers [2]

foreach tuple t in relation R:



- Maintain smaller (i.e. cacheline size) buffers per partition.
- Most of the writes happen to a small number of pages; avoiding TLB misses.

Back of the envelope

- 4KB page / 64B cacheline = 1024 partitions per page
- Assume 32 threads do 2048 fanout partitioning => this will consume 64 TLB entries.
 - SW-managed buffers are a good optimization for the TLB bottleneck, yet higher partitioning fanout can still cause a problem.

Multi Round Partitioning



MySQL HeatWave: Query Optimization

Query Optimization: Holistic Optimization

TPC-H Query 21



0

Query Optimization: A Detailed Overview

TPC-H Query 21



Logical Optimization: Pushdown Transformation



- Logical equivalence: We can perform the join between T1.L_PARTKEY with P_PARTKEY before or after the group-by. Which one is cheaper though?
 - Intuition: Selective join, especially with filtered build-side, is usually cheap.
- Without pushdown, group-by has to process 180 Billion rows from LINEITEM.
- After the pushdown, group-by processes 180 Million rows (1000x reduction).



- False positives are possible (due to collisions at insert)
- False positive rate estimate: $E=(1-exp(-kN/m))^k$ (m: #bits, N: #inserts, k: #hashes)
 - For a low false positive rate:
 - More bits

٠

• Less inserts

Bloom Filter Application: TPC-H q08



Observations

- **False-Positive-Rate:** More inserts require larger BF. Trade-off between low FPR and cache-locality.
- **Expected filtering** does not only depend on FPR, but also join output cardinality. Advantage only for high selective joins.
- **Runtime improvement:** Depends on the filtering rate, however BF lookup itself turns out to be costly. Apply a BF only for good filtering rate.
- Join pattern: System broadcast on the build side. No need to share local BFs across the cluster.

Bloom Filter Application: TPC-H q21



Observations

- **Runtime improvement** is great for this query, since we can push down the BF filtering through many operators.
- Join pattern: System partitioning on the build side. Costly, since we need to share local BFs across the network. Cost-based decision for these type of queries.

Bloom Filter Implementation

Standard vs. Blocked

Standard



- Set k bits in the entire bloom filter buffer.
- Each insert/lookup requires accessing k cachelines.

- - k random accesses per single lookup.
- - Hard to parallelize with AVX (row dimension rather than hash dimension).
- + Better false-positive-rate.



Blocked

- Set 1 bit in each sector in a block, for each insert. (B/S bits per insert)
- Each insert/lookup requires accessing a single block.

Choose to accommodate AVX2: B: 256 bits, S: 32 bits => 8 hashes (close to the optimal 7) A single lookup can happen in few 256-bit AVX instructions.

- + Single access per lookup.
- + Suitable for AVX.
- Worse false-positive-rate.

Bloom Filter Implementation

Standard vs. Blocked



Vanilla-BF vs Optimized-BF



a10

q11

Queries

q12

q13

q14

q15

q17

a18

C

a21

Summary

MySQL HeatWave: A Deep Dive into Optimizations

- Data partitioning is a first principle
 - CPU optimizations
- A Holistic Optimization Approach
 - Inform MySQL Hypergraph optimizer with HeatWave query properties (logical and physical).
- Logical optimizations save lots of work
 - Pushdown transformations
 - Cardinality results from stats cache
 - ...
- Physical optimizations needed to achieve good scalability
 - Bloom filter applications
 - Join patterns (partitioned vs. broadcast)
 - Group-by patterns (high-NDV, medium-NDB, low-NDV, approx. low-NDV)
 - Data placement
 - Late decompression
 - ...

