



StarRocks Query Optimizer

Kaisen Kang



•CelerData Query Team Leader

•StarRocks PMC & Committer

•Apache Kylin PMC & Committer (Inactive)

•Apache Doris PMC & Committer (Retired)



StarRocks Query Optimizer





Development History of StarRock



DataLake Analytic

Resource Isolation







KS			StarRocks
2.x on and Unity		3 Lake	X House
	Multi Table MV	Shared Data	Data Lake Write
	Global Dict	Semi Structured	Spill To Disk
c	Unified Catalog	MV For Lake	Flat Json
on	Primary Key Model	Full Text Search	Vector Search

From 2020 To 2025

StarRocks Architecture

HTTP

In-Memory Metadata Java **Query Optimizer Query Scheduler**





StarRocks Optimizer Overview 1

Cascades Search Framwork



Based on the Cascades and Columbia Papers



StarRocks Optimizer Overview 2





7

StarRocks Optimizer Introduction



Representative Optimizations

Multi Left Join Colocate

Global Dict Optimization

Partitioned MV Union Rewrite





StarRocks Colocate Join



Shuffle Join





Colocate Join

Enforce For Shuffle and Colocate

SELECT * FROM T1 INNER JOIN T2 ON a = b

Shuffle Join Enforce



T1 and T2 are normal tables



Colocate Join Enforce

T1.a Colocate With T2.b

Scan satisfy the required property

arRock	S
--------	---





Multi Inner Join Colocate

SELECT * FROM T1 JOIN T2 ON t1.c1 = t2.c1 Join T3 On T2.c1 = T3.c1

Multi Inner Join Colocate Enforce



Inner Join won't procude NULL So Multi Inner Join Colocate is same as One Inner Join



T1.c1, T2.c1, T3.c1 are Colocated

Left Outer Join will Produce NULL

SELECT * FROM T1 Left JOIN T2 ON T1.c1 = T2.c1 Left Join T3 On T2.c1 = T3.c1

T1













C2
a
b

C1	C2
1	a
3	b

T1 Left Join T2

1	T2.C1
	1
	NULL

The right table of left outer Join will produce NULL

Left Outer Join will Produce NULL

SELECT * FROM T1 Left JOIN T2 ON T1.c1 = T2.c1 Left Join T3 On T2.c1 = T3.c1



The NULL generated by left join makes property enforce unable to satisfy





The Result of Colocate NULL and Shuffle NULL are the same

SELECT * FROM T1 Left JOIN T2 ON T1.c1 = T2.c1 Left Join T3 On T2.c1 = T3.c1



T1.C1	T1.C2	T2.C1	T2.C2	T3.C
1	a	1	a	1
2	b	NULL	NULL	NUL

The NULL distribution of T2 does not affect the final result So T2 and T3 still could do Colocate Join We could add Null strict and Null relax mode to Distribution Property



Add Null strict and Null relax mode in DistributionProperty







Add Null strict and Null relax mode in DistributionProperty

SELECT * FROM T1 Left JOIN T2 ON T1.c1 = T2.c1 Left Join T3 On T2.c1 = T3.c1

Multi Left Join Colocate Enforce



T1.c1, T2.c1, T3.c1 are Colocated

For the dirstribution column, when Null relax mode, could ignore NULL value when enforce



StarRocks Optimizer Introduction



Representative Optimizations

Multi Left Join Colocate

Global Dict Optimization

Partitioned MV Union Rewrite





Dict Optimization (Operations on Encoded Data)

String Column With Dict Encode



Local Dict in Storage Engine





Int Compare is Very Faster Than String



StarRocks Low Cardinality Global Dict Optimization

Aggregate Hash Table

Key(Binary) AndroidBeijing AndroidShanghai IosBeijing IosShanghai G An

3X Performance Improvement For Aggregate

- Scan
- Filter
- Agg
- Sort
- Join
- String Functions



Select Sum(PV) From Table Group By City, Platform

Aggregate Hash Table With Encoded

				_		
Sum					Key(Int)	Sum
40000		String Encode To Int			11	40000
30000					12	30000
20000		Faster Hash Faster Equal		21	20000	
10000		Faster Memcpy		22	10000	
Platform lobal Dictionary			Ci Global D	ity Dictionary		
alue	Id		Value	Id		
droid	1		Beijing	1		
los	2		Shanghai	2		

Low Cardinality Global Dict Rewrite Overview





Column Statistic Cache In FE Memory

Low Cardinality Global Dict Rewrite

- Tree Based Rewrite
- Add Decode Node



- Attach Global Dict to Exchange
- String and Array String Type



StarRocks MV

- Async and Sync
- Single Table and Multi Tables
- Auto Rewrite: No need to change the query
- Partition Refresh
- MV TTL

Year	City	Cost(sum)	
2016	beijing	10	
2016	shanghai	30	K
2017	beijing	20	
2017	shnaghai	40	

Base Table





MV (Year, Cost)

	Year	Cost(sum)	
Materialized	2016	40	select sum(cost) from ta Auto Rewrite where year = 2016
	2017	60	



ıble

StarRocks MV Auto Rewrite

- Agg Rewrite
- Join Rewrite
- Union Rewrite
- Nested MV Rewrite
- Text-Based Rewrite
- View-Based Rewrite







Why Need Auto MV Union Rewrite

MV For DataLake



MV speed up the recent data of Data Lake

We need to union the MV and Historical Table Data



MV For RealTime



MV speed up the historical data

We need to union the MV and RealTime Table Data

MV Union Rewrite 1: MV is up-to-date





MV Union Rewrite 2: Partial Partitions Refreshed

MV:

select count(*) from t group by a

Query:

select count(*) from t group by a

MV and Base Table have four partitions:







MV Union Rewrite 3: Partitioned MV With Predicate

MV:

Select * from t where a > 20

Query:

Select * from t where a > 10

select * from t

where **a** > 10 **and a** <= 20

union all

select * from "MV"

MV and Base Table have two partitions:





select * from "MV"

=>

select * from t where a > 20 and

dt='20250202'

union all

select * from MV

Step2: Compensating Partition Predicates

Consider MV Partitions refresh status select * from t where a > 10 and

union all

(select * from t where a > 20 and

dt='20250202'

union all

select * from mv)

Step3: Get the final result

Merge the result from step1 and step2

_	_	_	_	_	-
a	<	(=	=	2	0
_	_	_	_	_	_'
_	_	_		_	-,
					į
					į
					į
					i
					İ
_	_	_	_	_	_1





Challenges Of CBO

CBO Challenges

No Statistics

Inaccurate Sampling Statistics

Statistics Estimation Errors

Cost Estimation Errors

Data Skew

Query Plan Jitter



StarRocks Solutions

Auto Analyze

Table Sample

Full Analyze Predicate Columns

Query Feedback Adaptive Execution

Histogram Skew-Awared Shuffle Join

SQL Plan Management Hint

StarRocks Statistics Overview

- •Count
- •Max, Min
- •Null Fraction
- •Distinct Count
- •Avg Row Size

- •Histogram
 - •MCV
 - Bucket Statistics





Optimizes Performance and Accuracy of Statistics Collection







StarRocks Table Sample 1





StarRocks Distributed Table Sampler

StarRocks Table Sample 2







StarRocks Predicate Columns

- Where Filter Columns
- Group By Columns
- Distinct Columns
- Join On Columns





Column Statistics need to be collected Only for Predicate (Dimension) columns

StarRocks Optimizer Introduction



Representative Optimizations

Multi Left Join Colocate

Global Dict Optimization

Partitioned MV Union Rewrite





Query Feedback: Overview



The result of Query Feedback is the tuning guide



Query Feedback: Use Case & What's the Tuning Guide

- Join tuning guide
 - Optimize the left and right order of Join
 - Optimize the distributed execution mode of Join (eg: broadcast -> shuffle)
- Agg tuning guide
 - and aggregate more data in first phase.



• For cases where first-phase aggregation works well, disable streaming aggregation

Query Feedback: Use Case







Optimize the left and right order of Join nodes

Query Feedback: Use Case



Optimize the distributed execution mode of Join nodes: Shuffle -> Broadcast





StarRocks Optimizer Introduction



Representative Optimizations

Multi Left Join Colocate

Global Dict Optimization

Partitioned MV Union Rewrite





Adaptive Streaming Aggregate 1







ar]	R	0	C	ks	5
	Γ	U	C	K)

Adaptive Streaming Aggregate 2

Distributed One-Phase Aggregation

Local Agg has no aggregation effect:

- Reduces the build of useless HashTable
- Reduces CPU and Memory resources





Fragment Instance

Adaptive Streaming Aggregate 3

- Prefer two-phase aggregation
- Runtime adaptive
- Good aggregation effect: use hash table
- Bad aggregation effect: directly shuffle data

Do not rely on error-prone cardinality estimates Rely on real runtime statistics





Adaptive Runtime Filters Selection 1

- Support Join, TopN, Agg Runtime Filter
- Support Local And Global Runtime Filter
- Shuffle Aware
- Push down Max/Min, In Filter To Storage Engine
- Support Runtime Filter Cache
- Push Runtime Filter To Two Sides
- SIMD Bloom Filter
- Adaptive Join Runtime Filters Selection

1B Rows
Fact Scan

Reduce Disk IO, Reduce Network Transport, Reduce Join Probe Rows

10 x ~ 100x Performance Improvement For Complex Join Query



Select * From Fact Join Dim On Fact.DimId = Dim.Id Where Dim.Price > 100





Adaptive Runtime Filters Selection 2

Issues:	dou	uble	sel
	if	(se]	lect
• Some filters are ineffective		if	(se
			_S
• Filters with low selectivity are applied las	st		_S
			ch
			re
		}	
		//	Onl
Filter compute has cost,		if	(_s
			_S
we should use most selective filters		} (else
		1	

Rely on real runtime statistics



```
.ectivity = true_count * 1.0 / chunk_size;
ivity <= 0.5) { // useful filter</pre>
lectivity < 0.05) { // very useful filter, could early return</pre>
electivity.clear();
electivity.emplace(selectivity, rf_desc);
unk->filter(new_selection);
turn;
```

y choose three most selective runtime filters electivity.size() < 3) {</pre> electivity.emplace(selectivity, rf_desc); { . .. -

Do not rely on error-prone selective estimates



StarRocks Optimizer Introduction



Representative Optimizations

Multi Left Join Colocate

Global Dict Optimization

Partitioned MV Union Rewrite





The Goal of SQL Plan Management

- •Ensure no performance regression after upgrading
 - It is difficult to ensure that the changes to the optimizer are 100% positive optimization.

- •Ensure that the query plan for online business remains unchanged
 - •Plan changes may lead to big queries, which in turn may cause production accidents



SQL Plan Management 1: How to Generate Baseline Plan

select * from t0, t1, t2 where t0.v1 = t1.v1and t0.v1 = t2.v1and t0.v2 > 200

Bind SQL

Bind SQL With Hint

select *

from t0 join **[SHUFFLE]** t2 on t0.v1 = t2.v1 join **[BUCKET]** t1 on t0.v1 = t1.v1

where $t0.v2 > _spm_const_var(0)$

Plan Storage





SQL Plan Management 2: How to Use Baseline Plan

select * from t0, t1, t2 where t0.v1 = t1.v1 and t0.v1 = t2.v1 and t0.v2 > 1000

User Query

Removing Parameters

select *

from t0 join [**SHUFFLE**] t2 on t0.v1 = t2.v1 join [**BUCKET**] t1 on t0.v1 = t1.v1

where t0.v2 > **1000**

Normal Query Flow



```
select *
from t0, t1, t2
where t0.v1 = t1.v1
and t0.v1 = t2.v1
and t0.v2 >spm_const_var(0)
```



Lessons From StarRocks Query Optimizer

- Errors in optimizer cost estimation are unavoidable, and the executor needs to be able to make autonomous decisions and provide timely feedback.
- In engineering, the optimizer testing system is as important as the optimizer itself. The optimizer requires correctness testing, performance testing, plan quality testing, etc. • Null and Nullable is Interesting and Annoying
- - For Performance, we need to special handle Null and nullable
 - For Correctness, we need to take care of Null and nullable in Optimizer and Executor
- Integrated optimizer could be more powerful than Standalone optimizer





Thanks kangkaisen.com



