

How to Write Fast Flink SQL

贺小令 | 阿里巴巴高级技术专家、Apache Flink Committer

01 Flink SQL Insight

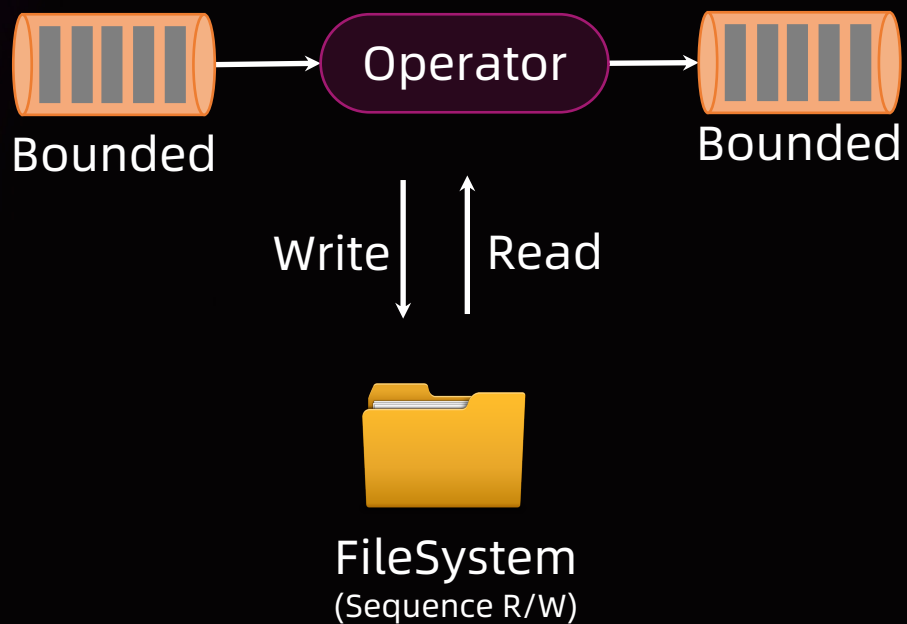
02 Best Practices

03 Future Works

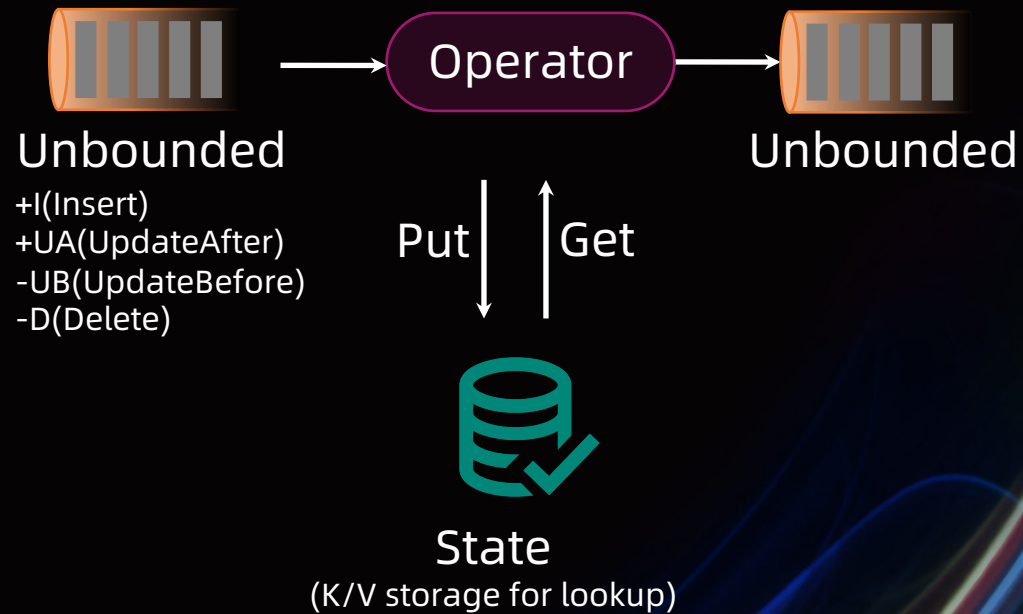
01 Flink SQL Insight

Flink Operator Execution Mode

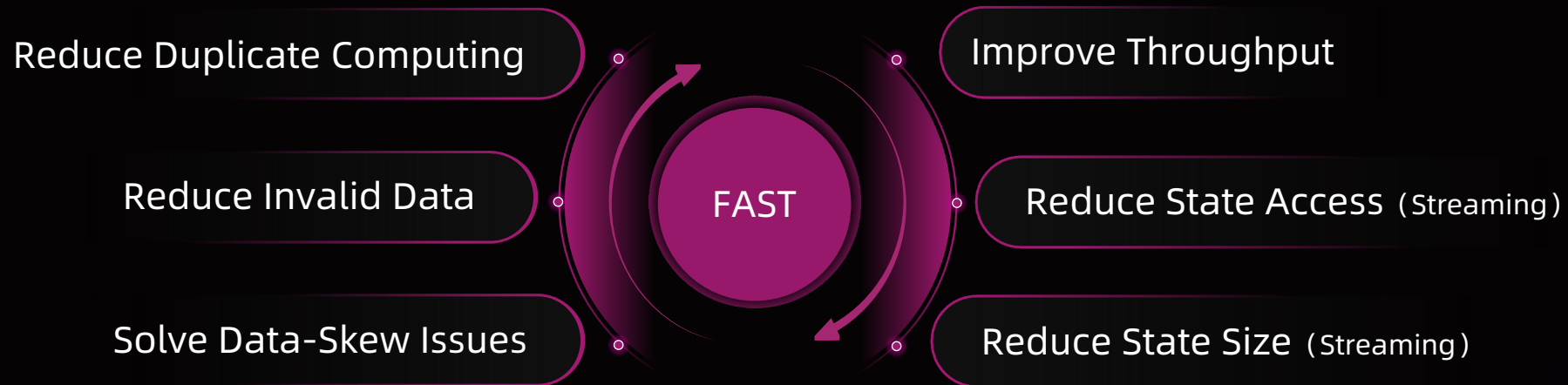
Batch



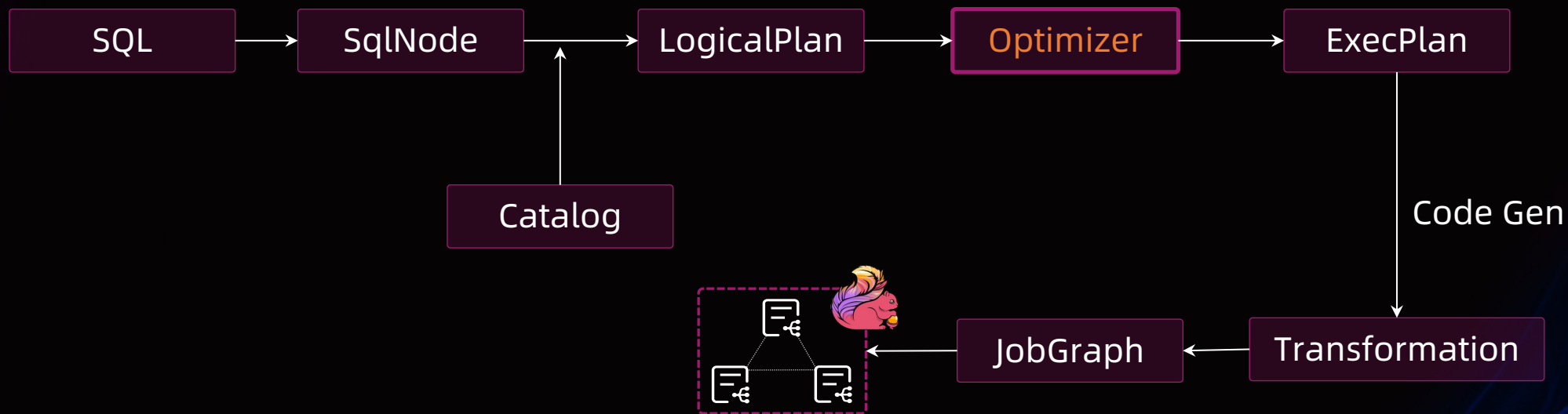
Streaming



How To Make Flink Job Run Fast

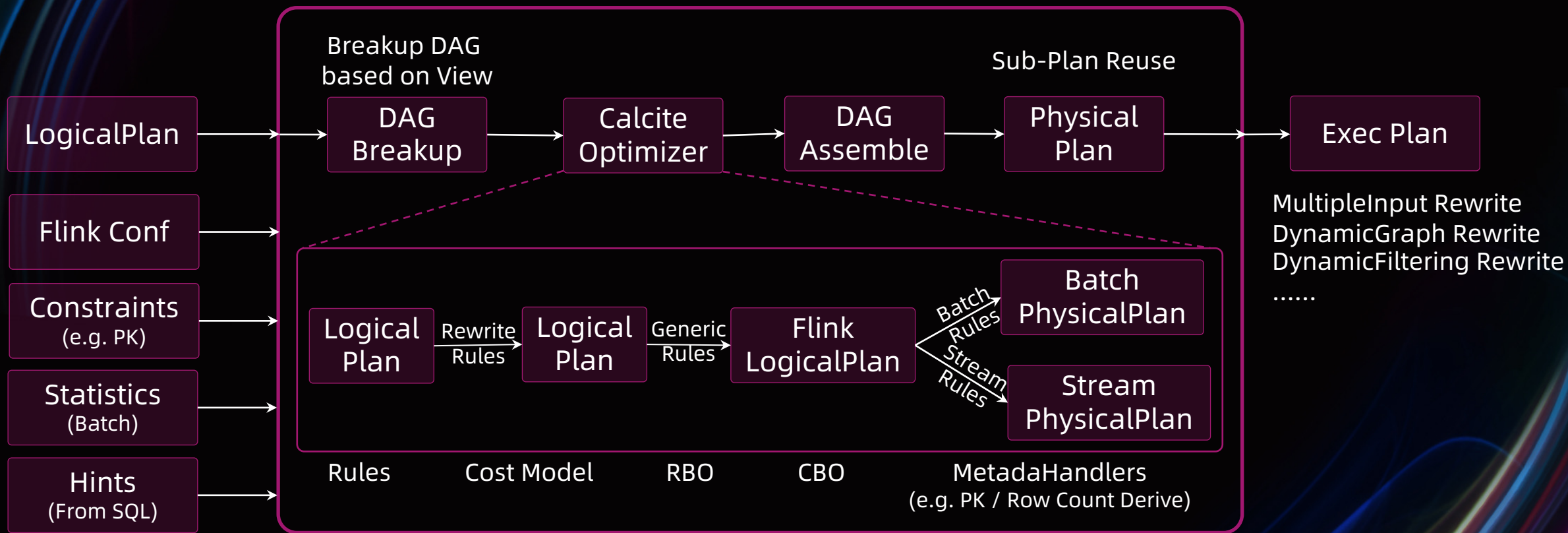


Flink SQL Workflow



Insight Optimizer

DAG Optimizer

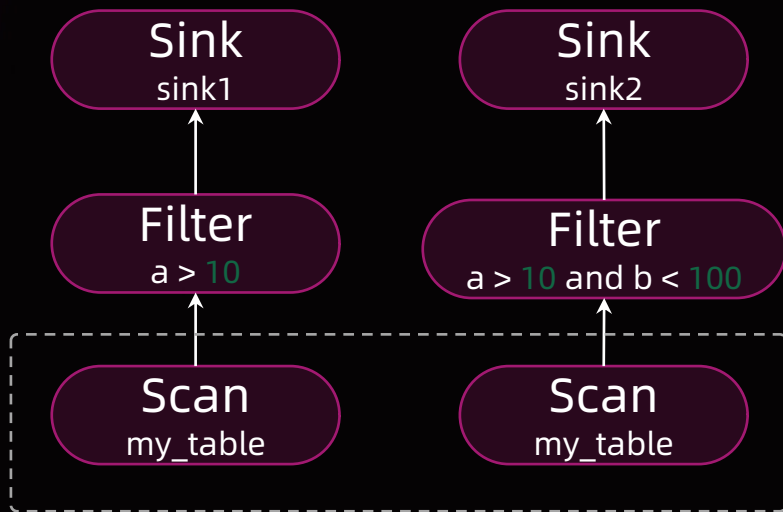


02 Best Practices

Sub-Plan Reuse

```
insert into sink1 select * from my_table where a > 10;  
insert into sink2 select * from my_table where a > 10 and b < 100;
```

Duplicate
Computing



Optimize

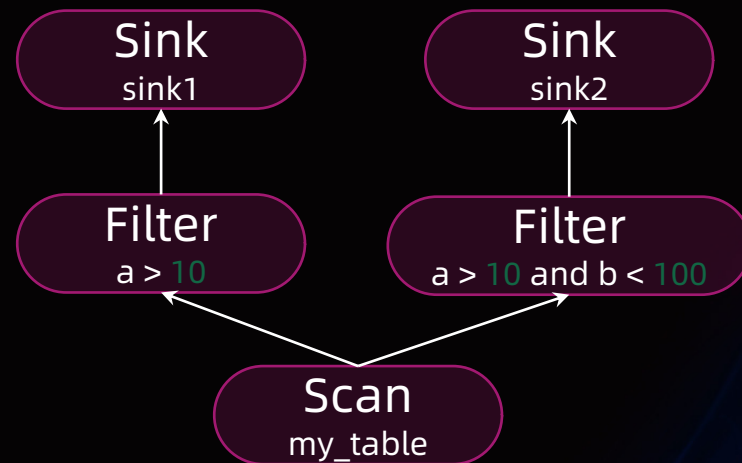
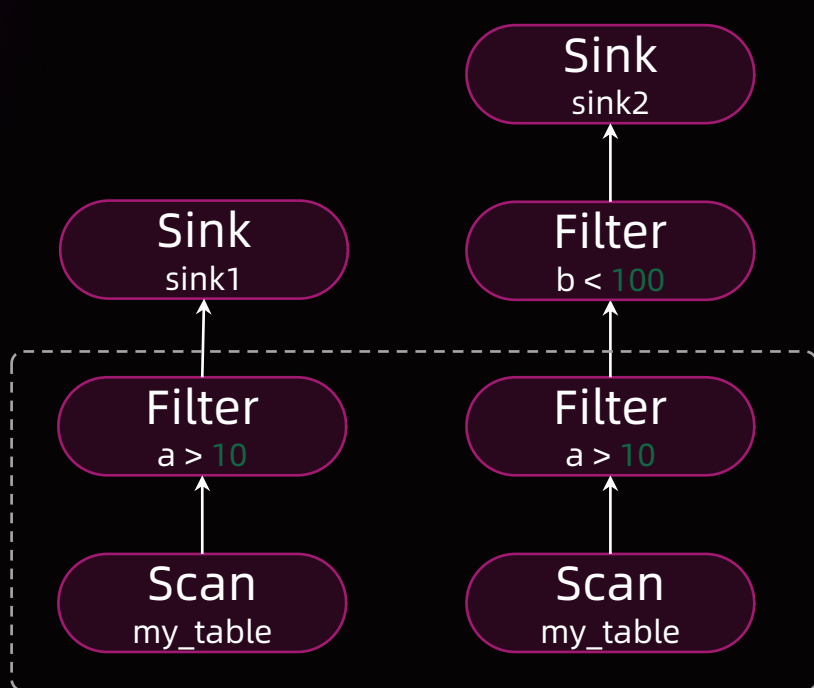


table.optimizer.reuse-sub-plan-enabled: true

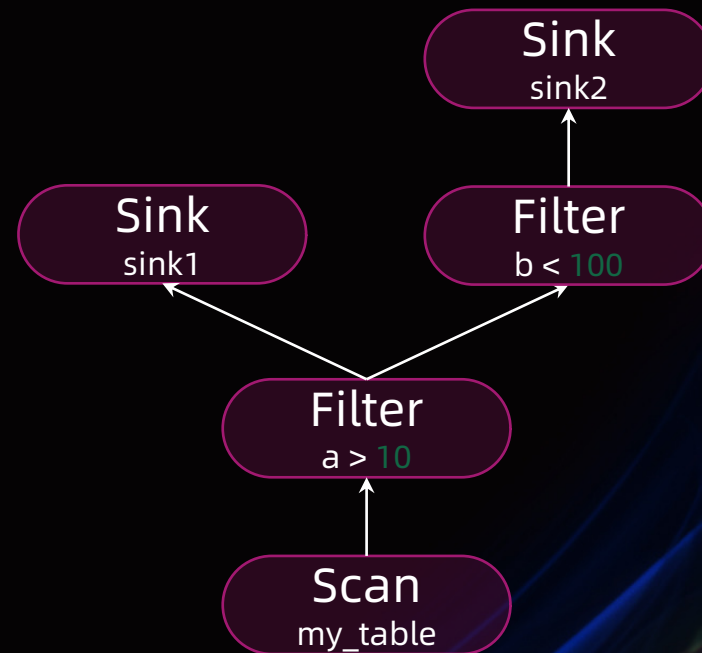
Sub-Plan Reuse

```
create temporary view v1 select * from my_table where a > 10;  
insert into sink1 select * from v1;  
insert into sink2 select * from v1 where b < 100;
```

Duplicate
Computing

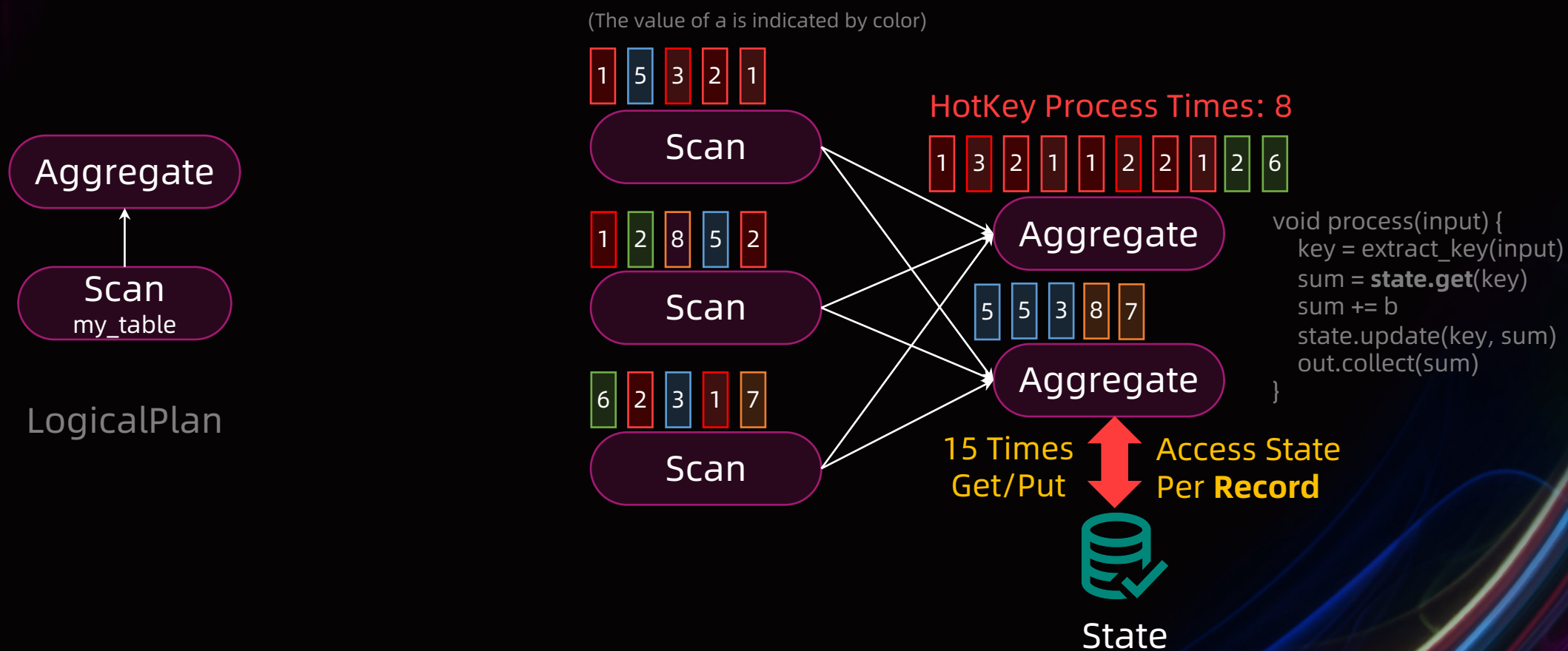


Optimize



Fast Aggregation

`select a, sum(b) from my_table group by a;`



Fast Aggregation

`select a, sum(b) from my_table group by a;`

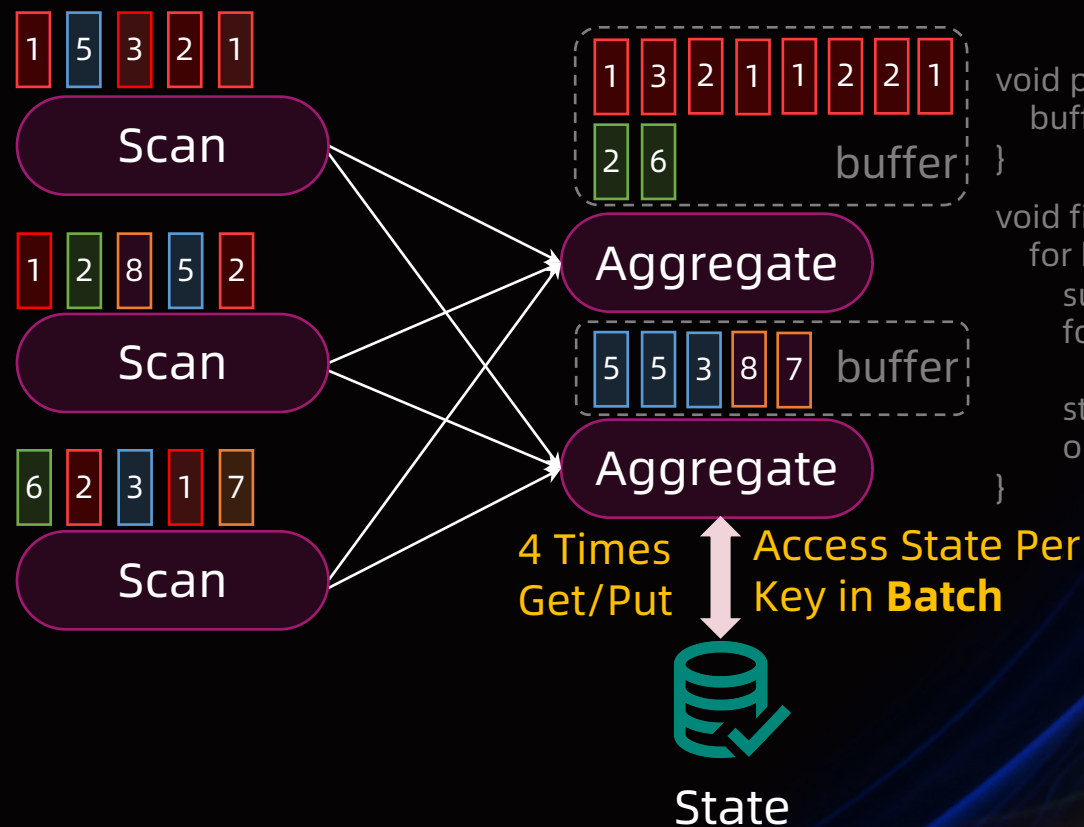
Enable MiniBatch:

- `table.exec.mini-batch.enabled: true`
- `table.exec.mini-batch.allow-latency: 5s`

Suitable scenarios:

- Low requirements for job delays (delays in collecting a batch)
- Insufficient state processing capability (Reduce state access)
- Insufficient downstream processing capability (Output less records)

(The value of a is indicated by color)



```
void process(input) {
    buffer.add(input)
}

void finish_buffer() {
    for key in buffer:
        sum = state.get(key)
        for v in buffer.get(key):
            sum += v
        state.update(key, sum)
    out.collect(sum)
}
```


Fast Aggregation

`select a, sum(b) from my_table group by a;`

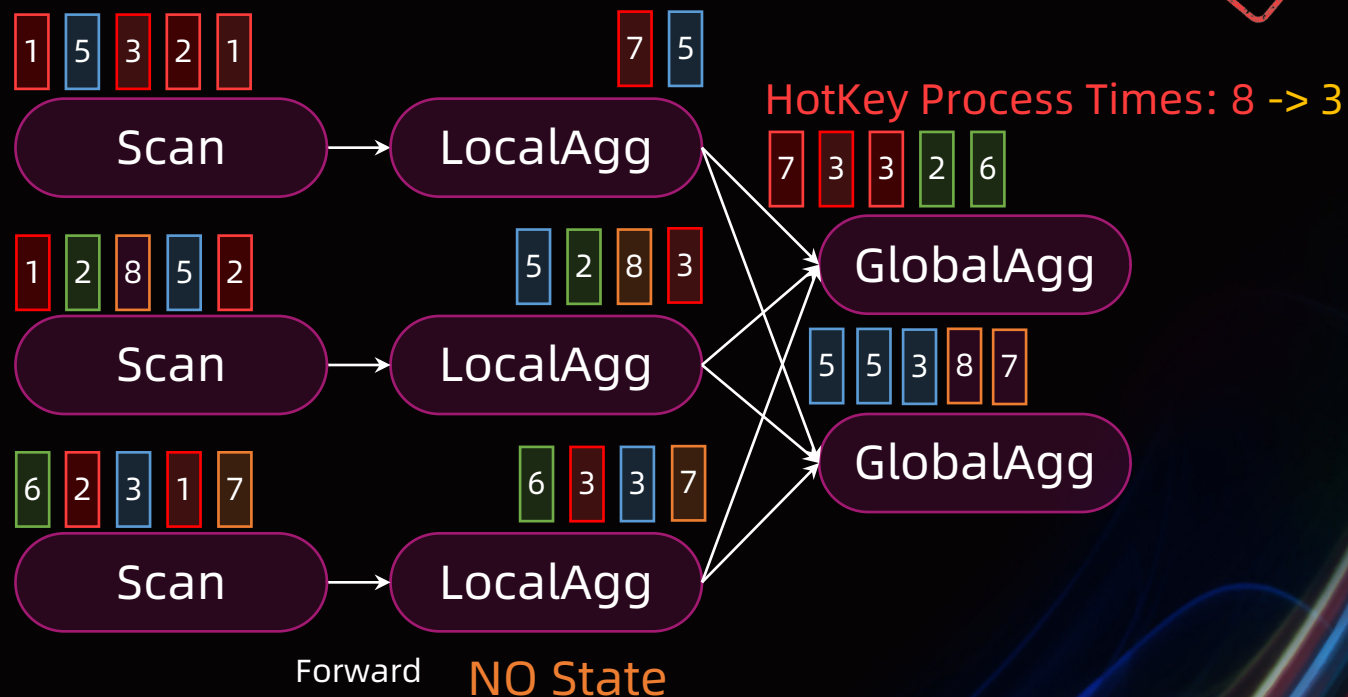
Enable Local/Global:

- `table.exec.mini-batch.enabled: true`
- `table.exec.mini-batch.allow-latency: 5s`
- `table.optimizer.agg-phase-strategy: TWO_PHASE/AUTO`

Suitable scenarios:

- All Agg Functions must implement the *merge* method
- Data skew after shuffle

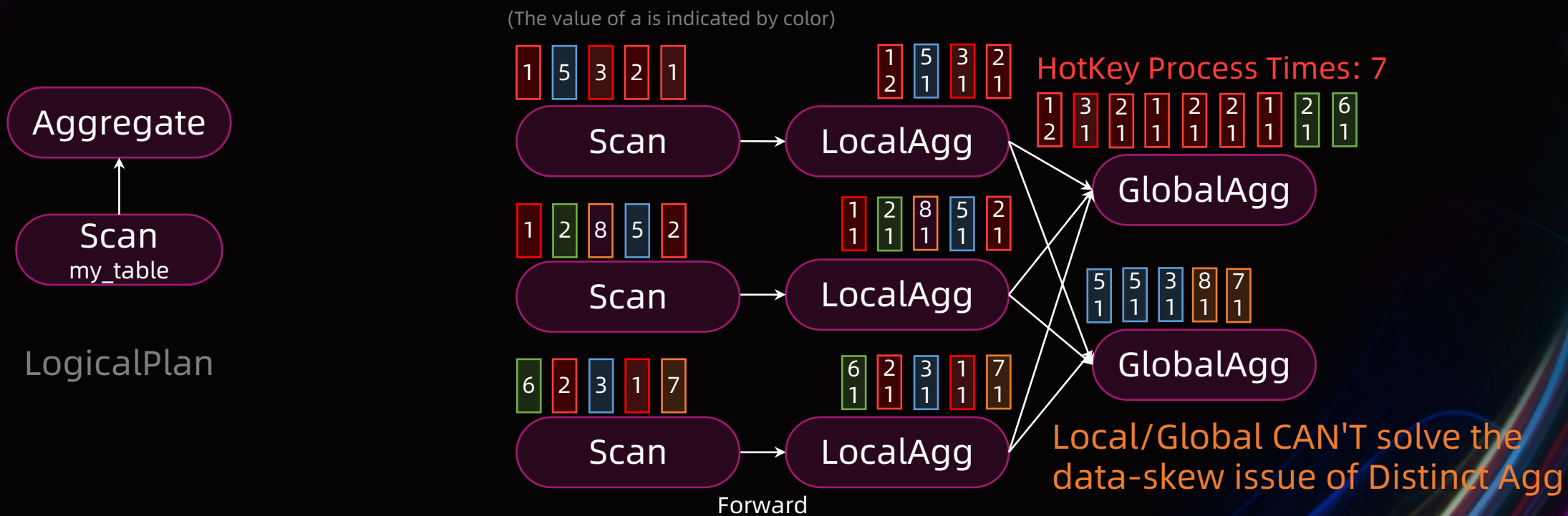
(The value of a is indicated by color)



Solve
Data Skew

Fast Aggregation

`select a, count(distinct b) from my_table group by a;`



Fast Aggregation

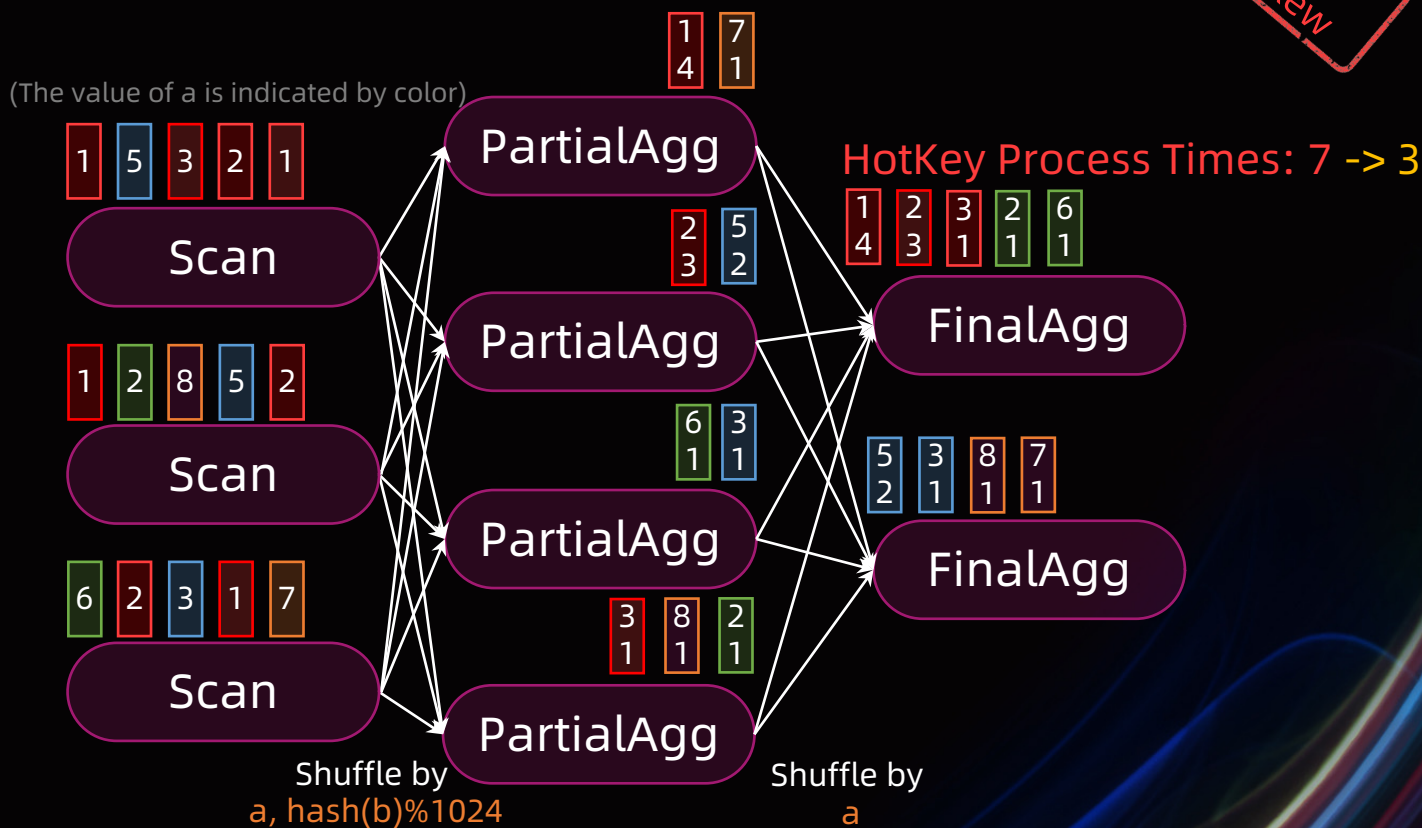
`select a, count(distinct b) from my_table group by a;`

Enable Partial/Final:

- `table.optimizer.distinct-agg.split.enabled: true`
- `table.optimizer.distinct-agg.split.bucket-num: 1024`

Suitable scenarios:

- Data skew in distinct function
- Only some built-in UDAFs are supported (e.g. sum, max)
- Large data set is beneficial (additional Shuffle is introduced)
- Note: Partial Agg will add computation and state



Solve
Data Skew

Fast Aggregation

`select a, count(distinct b) from my_table group by a;`

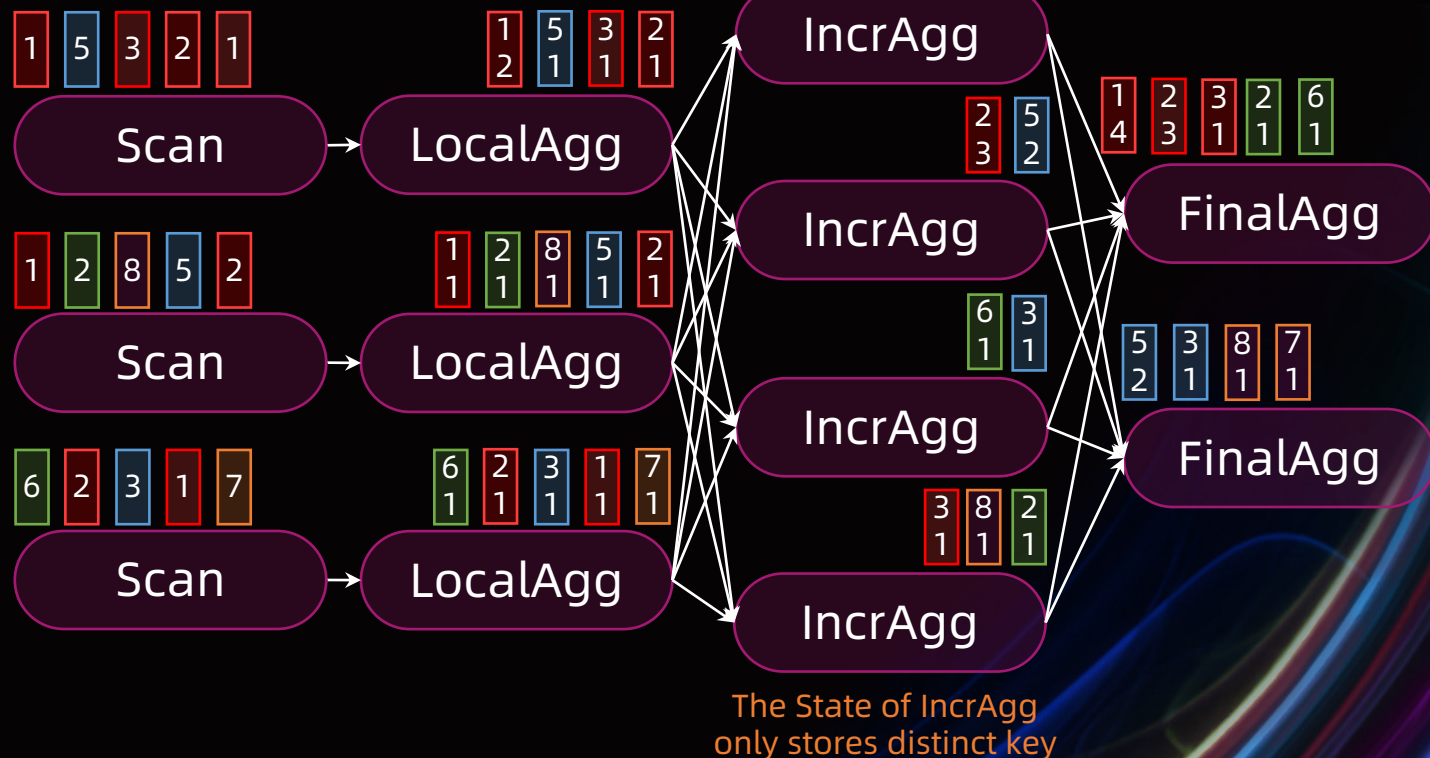
Enable Incremental:

- `table.exec.mini-batch.enabled: true`
- `table.exec.mini-batch.allow-latency: 5s`
- `table.optimizer.agg-phase-strategy: TWO_PHASE/AUTO`
- `table.optimizer.distinct-agg.split.enabled: true`
- `table.optimizer.distinct-agg.split.bucket-num: 1024`
- `table.optimizer.incremental-agg.enabled: true`

Suitable scenarios:

- The state size of partial aggregate is too large

(The value of a is indicated by color)



Fast Aggregation

Reduce
State Size

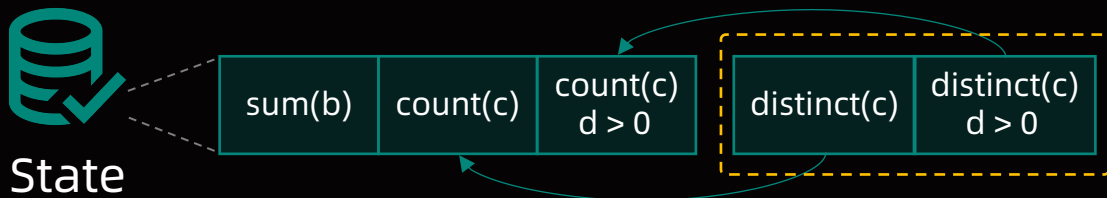
Aggregate Function with CASE WHEN

```
select
  a,
  sum(b),
  count(distinct c) as c1,
  count(distinct
    case when d > 0 then c
    else null end) as c2
from my_table
group by a;
```

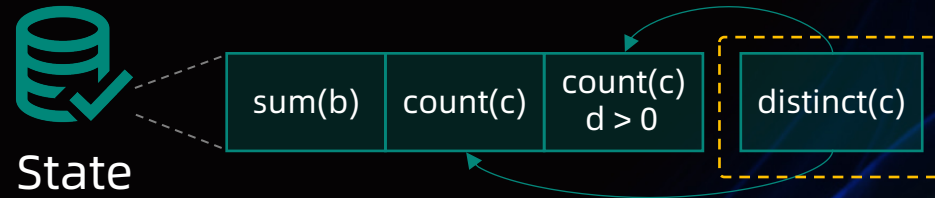
Rewrite

Aggregate Function with FILTER

```
select
  a,
  sum(b),
  count(distinct c) as c1,
  count(distinct c) filter (
    where d > 0) as c2
from my_table
group by a;
```



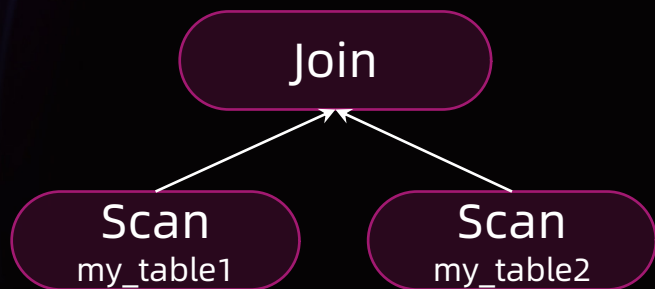
The State corresponding to the same function with same distinct field is stored **independently**



The State corresponding to the same function with same distinct field is stored **shared**

Fast Join

```
select * from my_table1 t1 join my_table2 t2 where t1.a = t2.c;
```



LogicalPlan

(The value of a is indicated by color)

1 5 3 2 1

Scan
my_table1

1 2 8 5 2

Scan
my_table1

(The value of c is indicated by color)

6 2 3 1 7

Scan
my_table2

Left State

1 5 3 2 1
1 2 8 5 2

Left State

Join

Right State

The State of the regular Join stores **all the input data**, which cannot handle large data sets

Fast Join

```
select * from my_table1 t1 join my_table2 t2 where t1.a = t2.c;
```

Reduce
State Size

The suggestions to reduce the State of regular Join:

- Join Key contains PK OR Join Input has PK
- Only keep the necessary fields before Join

Join Key => Map<Row, Tuple2<Int, Int>>
(The Map key is the input row, the value is the pair of the number of occurrences and the number of associations)

Join Key => Map<PK, Tuple2<Row, Int>>
(The Map key is PK, the value is the input row and its number of associations)

Join Key => Tuple2<Row, Int>
(The tuple2 is the input row and its number of associations)



State

Join Key contains PK



State

Join Input has PK

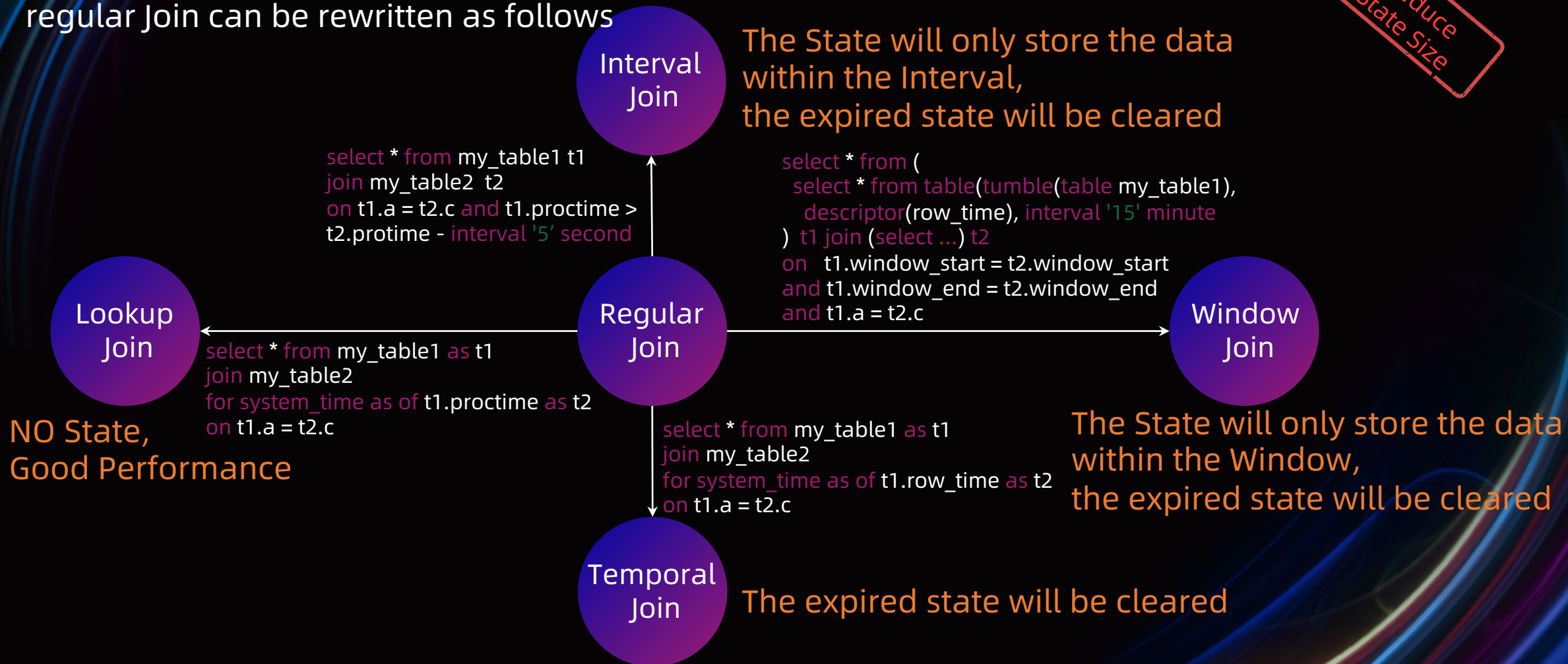


State

No PK

Fast Join

When meeting business requirements,
regular Join can be rewritten as follows



Reduce
State Size

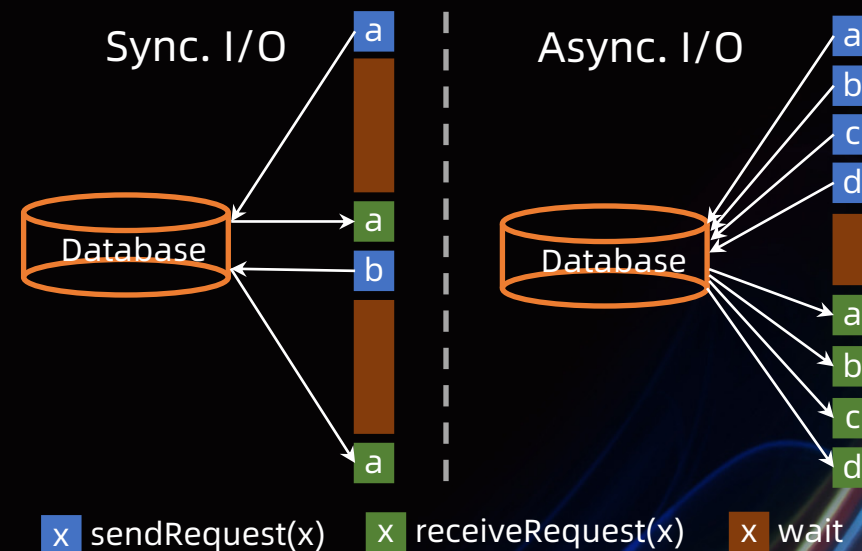
Fast Lookup Join

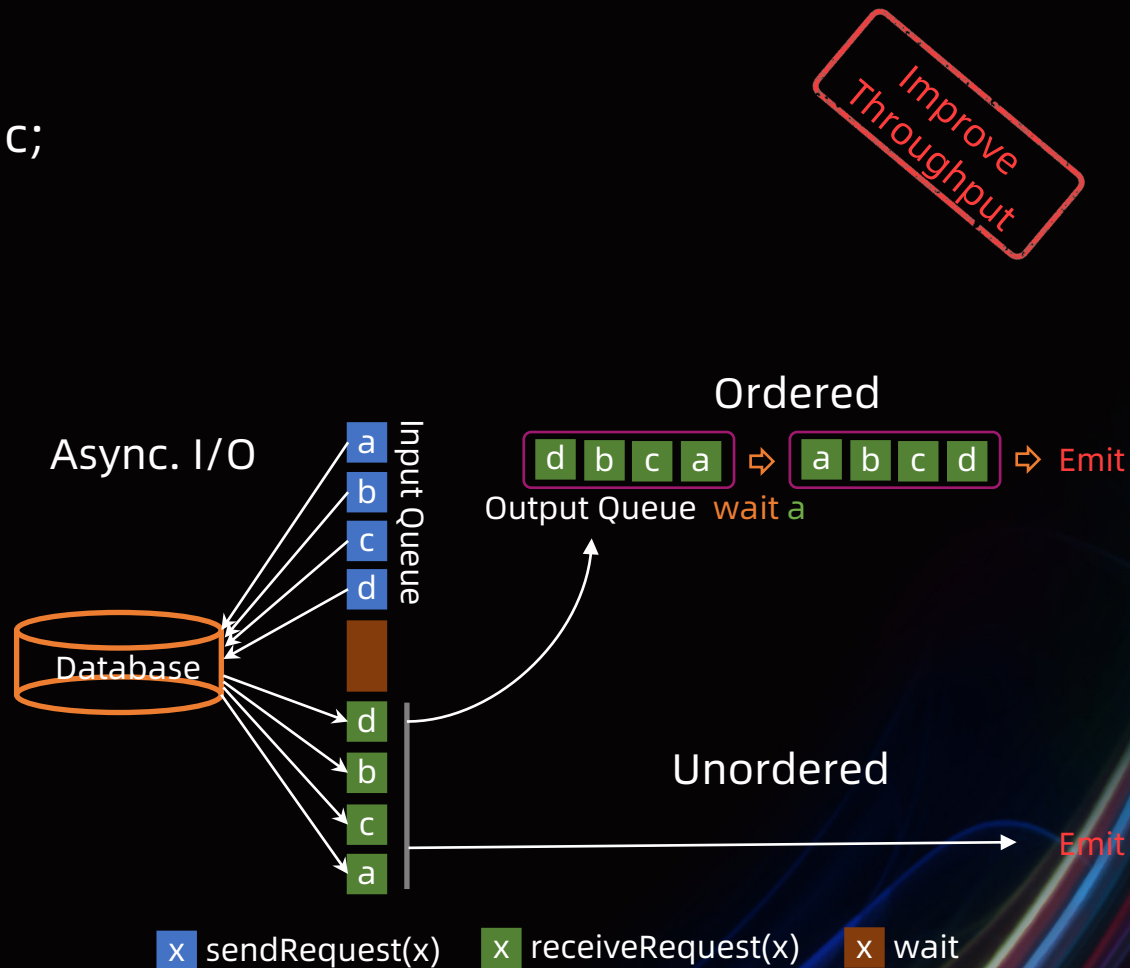
```
select * from my_table1 as t1 join my_table2
for system_time as of t1.proctime as t2 on t1.a = t2.c;
```

Sync vs Async

- Sync:
select /*+ LOOKUP('table'='my_table2', 'async'='false') */ *
from my_table1 as t1 join my_table2
for system_time as of t1.proctime as t2 on t1.a = t2.c;
- Async (Improve throughput):
select /*+ LOOKUP('table'='my_table2', 'async'='true') */ *
from my_table1 as t1 join my_table2
for system_time as of t1.proctime as t2 on t1.a = t2.c;

Improve
Throughput





Fast Lookup Join

```
select * from my_table1 as t1 join my_table2  
for system_time as of t1.proctime as t2 on t1.a = t2.c;
```

Cache (Memory Lookup):

- Full Caching (Small data set):
 - 'lookup.cache': 'FULL'
 - 'lookup.full-cache.reload-strategy': 'PERIODIC'
 - 'lookup.full-cache.periodic-reload.interval': '1h'
 - 'lookup.full-cache.periodic-reload.schedule-mode': 'FIXED_DELAY'
- Partial Caching (Large data set):
 - 'lookup.cache': 'PARTIAL'
 - 'lookup.partial-cache.max-rows': 10000
 - 'lookup.partial-cache.expire-after-access': '1h'
 - 'lookup.partial-cache.expire-after-write': '1h'
 - 'lookup.partial-cache.cache-missing-key': 'true'
- No Caching (Disable cache):
 - 'lookup.cache': 'NONE'

Improve
Throughput

Fast Deduplication

Find the FIRST row per key

```
select
  a,
  first_value(b),
  first_value(c)
from my_table
group by a;
```

Find the LAST row per key

```
select
  a,
  last_value(b),
  last_value(c)
from my_table
group by a;
```

BAD Practices

- Large State
- Incomplete semantics (can't handle *null* values)

Fast Deduplication

Reduce
State Size

Find the FIRST row per key

```
select a, b, c from (  
  select a, b, c,  
    row_number() over (  
      partition by a  
      over by time_attr asc) as rn  
  from my_table)  
where rn = 1;
```

The State will only store Keys

Find the LAST row per key

```
select a, b, c from (  
  select a, b, c,  
    row_number() over (  
      partition by a  
      over by time_attr desc) as rn  
  from my_table)  
where rn = 1;
```

The State will only store the last row of each Key
(NOT supported when upstream is changelog
and time attribute is row_time)

Fast TopN

TopN (Rank) implementations

(Performance decreases sequentially)

- AppendRank
 - Input produces insert-only changes
 - (The State stores the records in TopN per key)
- UpdateFastRank
 - Input produces update changes
 - Upsert key contains partition key
 - Order-by fields are monotonic, and the monotonic direction is opposite to the order-by direction
 - (The State stores a Map, its key is order key and its value is the record and the order number)
- RetractRank
 - NO required for input
 - (The State stores all input data)

Reduce State Size

```
select a, b, c from (  
  select a, b, c,  
    row_number() over (  
      -- The upsert key from upstream contains partition key  
      partition by a, b  
      -- c is monotonically increasing, while the order  
      -- direction is monotonically decreasing  
      order by c desc) as rn  
  from (  
    select a, b,  
      -- Declare the argument of sum to be a positive number,  
      -- So that the result of sum is monotonically increasing  
      sum(c) filter (where c >= 0) as c  
    from my_table  
    -- Produce upsert stream  
    group by a, b))  
where rn < 10;
```

Fast TopN

TopN optimization methods

- DO NOT output row_number field
 - Greatly reduces the amount of data in the result table
 - (the data can be sorted in front-end if needed)
- Increase Cache (LRU memory) size
 - table.exec.rank.topn-cache-size: 10000
 - $\text{cache_hit} = \text{cache_size} * \text{parallelism} / \text{top_n_num} / \text{partition_key_num}$
 - (NOTE: The TM memory needs to be increased accordingly)
- Partition fields are best related to time
 - Otherwise, state TTL will leads to wrong result

Reduce
State Access

Efficient User Defined Connector

Reduce
Invalid Data

Implement the following interfaces to improve execution efficiency:

- SupportsFilterPushDown
 - Reduce Scan I/O
- SupportsProjectionPushDown
 - Reduce invalid field reads
- SupportsPartitionPushDown
 - Reduce the invalid Partitions (static optimization)
- SupportsDynamicFiltering
 - Reduce the invalid Partitions (dynamic optimization)
- SupportsLimitPushDown
 - Reduce Scan I/O
- SupportsAggregatePushDown
 - Reduce Scan I/O, output less data
- SupportsStatisticReport
 - Report statistics, the optimizer can generate a better execution plan

Use Hints Well

- Table Hints: Change Table Options
 - e.g. Change the cache strategy of Lookup Table
 - `/*+ OPTIONS('lookup.cache'='FULL') */`
- Query Hints:
 - LOOKUP: Change Lookup Join Strategy
 - `/*+ LOOKUP('table'='my_table2', 'async'='true') */`
 - BROADCAST: Suggest the Optimizer chooses Broadcast Hash Join (batch only)
 - `/*+ BROADCAST(t1)*/`
 - SHUFFLE_HASH: Suggest the Optimizer chooses Shuffle Hash Join (batch only)
 - `/*+ SHUFFLE_HASH(t1)*/`
 - SHUFFLE_MERGE: Suggest the Optimizer chooses Sort Merge Join (batch only)
 - `/*+ SHUFFLE_MERGE(t1)*/`
 - NEST_LOOP: Suggest the Optimizer chooses Nested Loop Join (batch only)
 - `/*+ NEST_LOOP(t1)*/`

03 Future Works

Future Works

Deeper Optimization

Support more advanced
optimizations

Richer Optimization

Support more scenarios

Smarter Optimization

Support dynamic
optimization

THANK YOU

谢 谢 观 看