

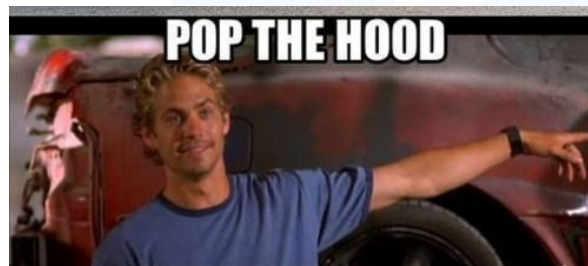
WWW.PICODATA.IO



PICODATA

Picodata SQL internals

What it is?



PICODATA SQL



Plan

1. Overview
2. Read query 1
3. Read query 2
4. DML
5. DDL

Overview

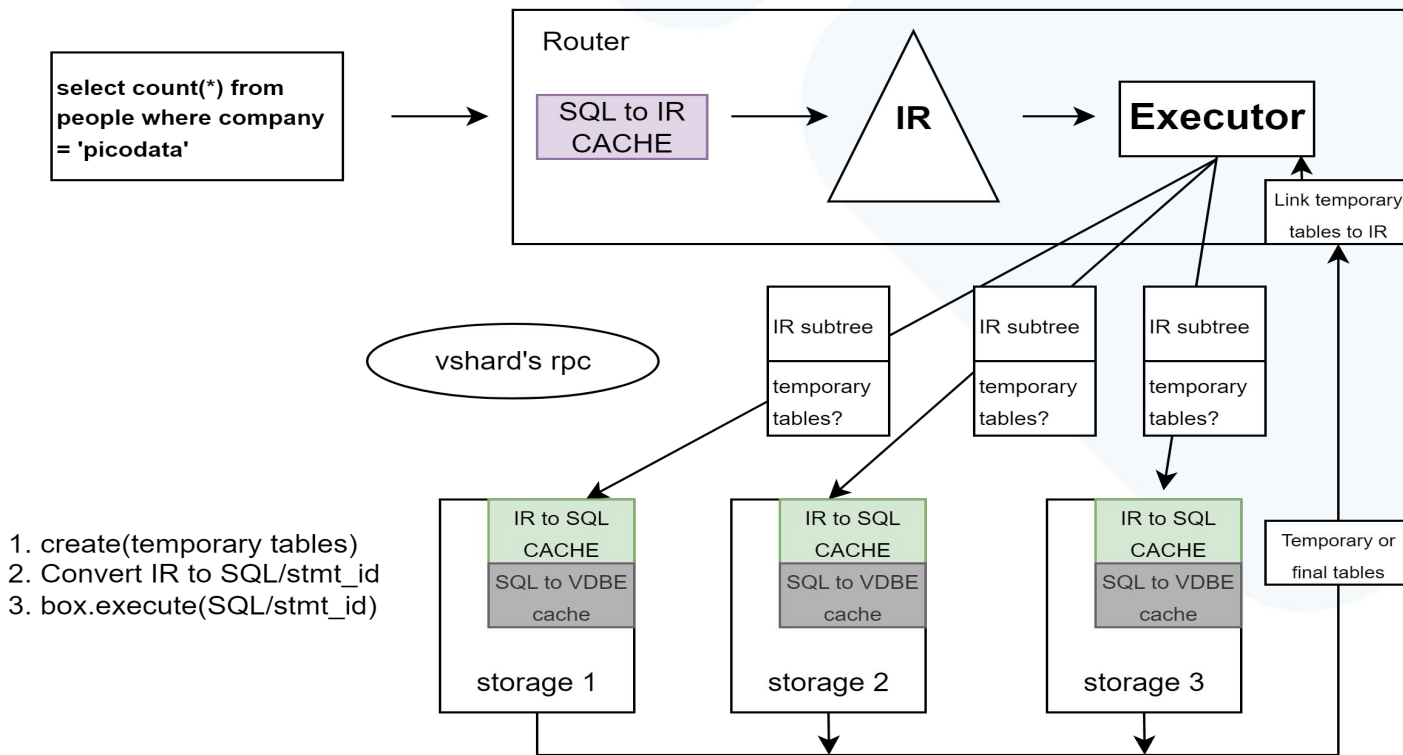
Что такое Picodata SQL



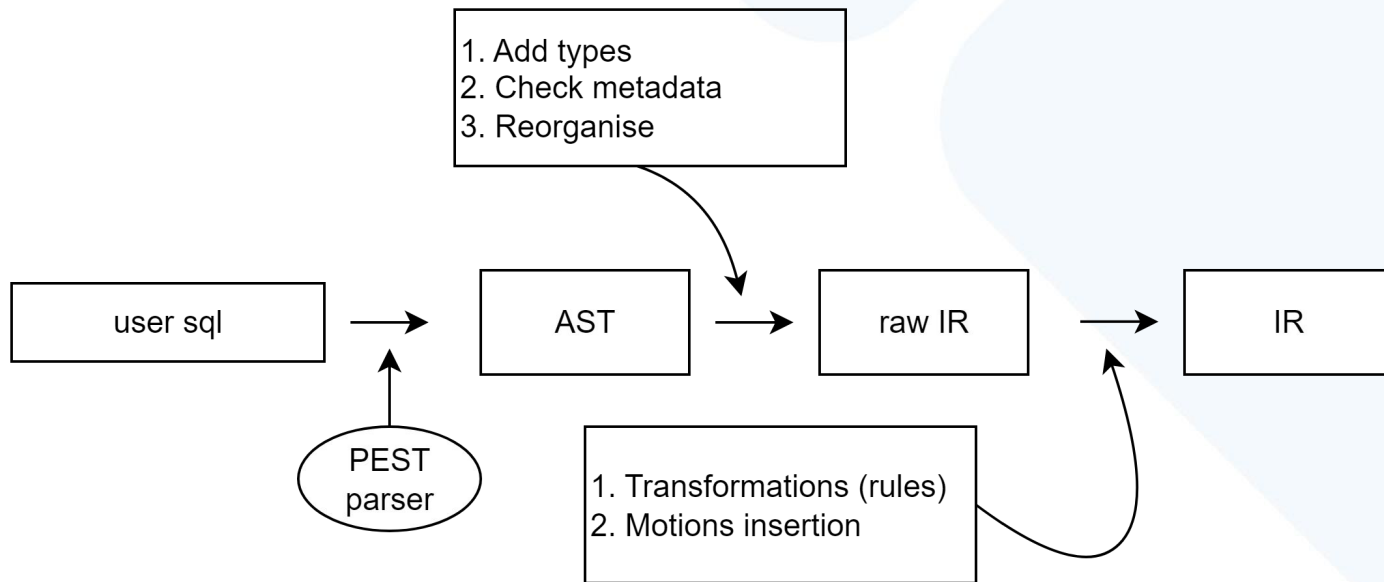
1. Library for distributed SQL
2. Uses vshard for queries dispatch
3. Needs several traits to be implemented before usage (e.g DDL)
4. For executing on single node uses SQL Tarantool (read queries), tarantool-module (update/insert/delete)

Overview: big picture

IR == Intermediate representation == Plan

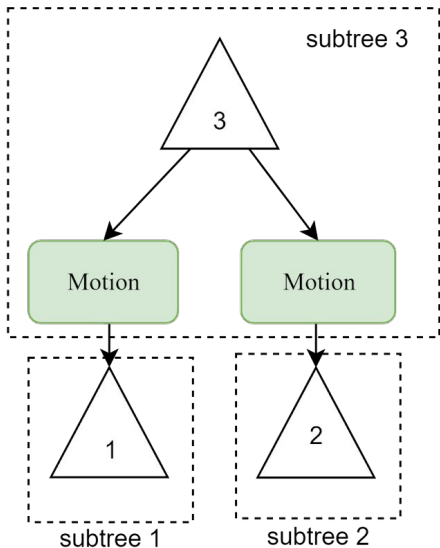


Overview: IR creation

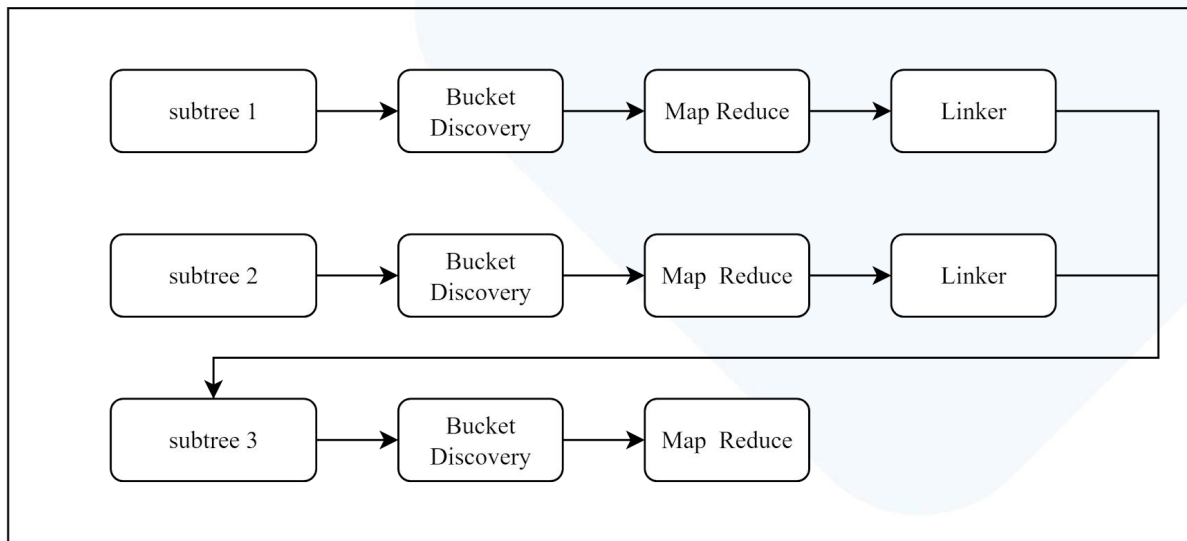


Overview: Executor

IR



Executor



Example 1

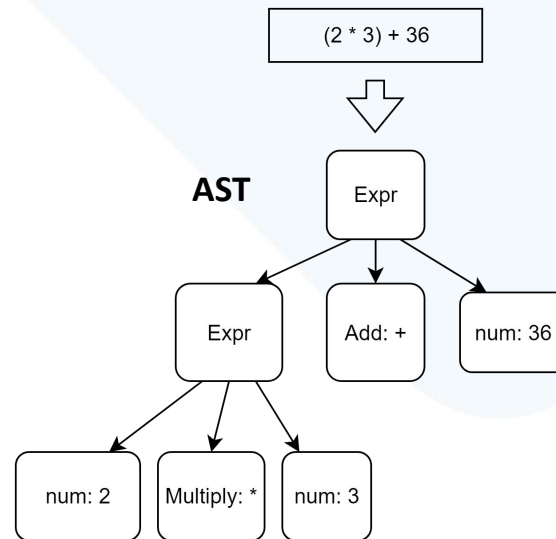
Query

```
SELECT b FROM t WHERE a = 1
```

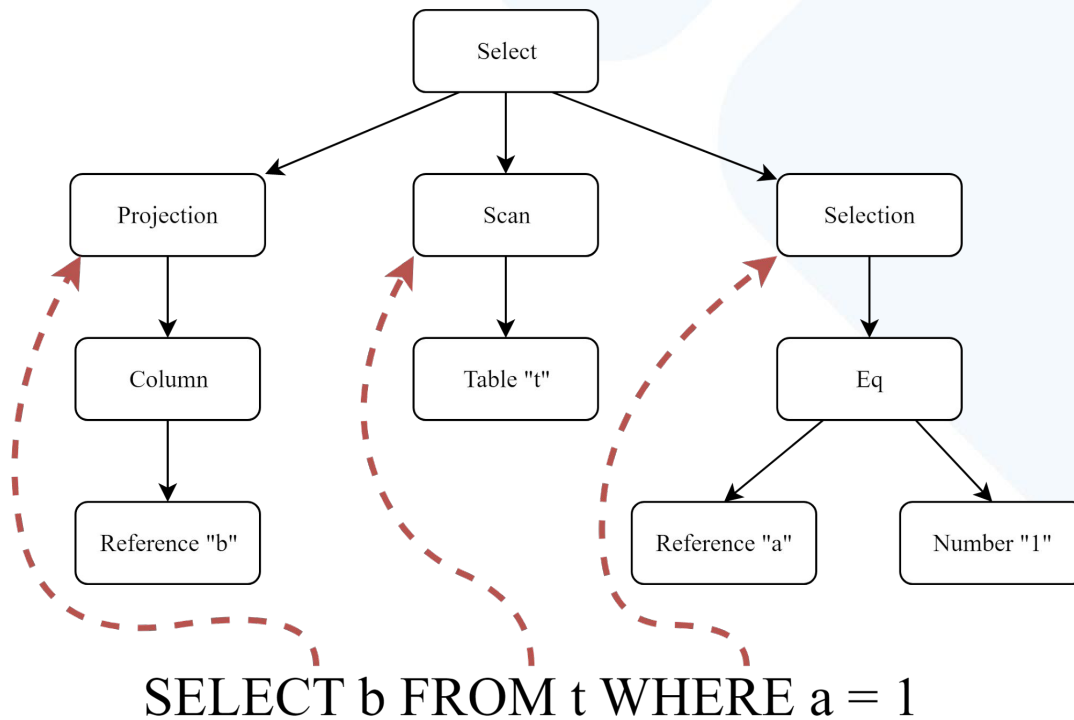
Grammar

- Grammar - a way to describe formal language
- AST (abstract syntax tree) - representation of sentence (that satisfies grammar) in form of a tree.
- PEST - external library for building AST using PEG (Parsing Expression Grammar)

```
num = @{ ("+" | "-")? ~ ASCII_DIGIT+ }  
operation = _{ add | subtract | multiply | divide }  
    add      = { "+" }  
    subtract = { "-" }  
    multiply  = { "*" }  
    divide   = { "/" }  
expr = { term ~ (operation ~ term)* }  
term  = _{ num | "(" ~ expr ~ ")" }  
calculation = _{ SOI ~ expr ~ EOI }
```



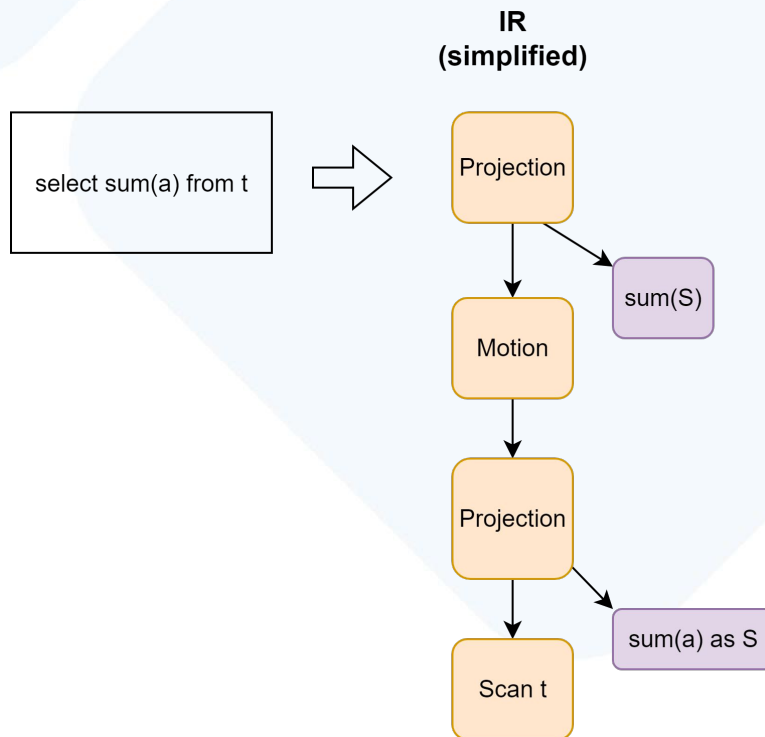
AST



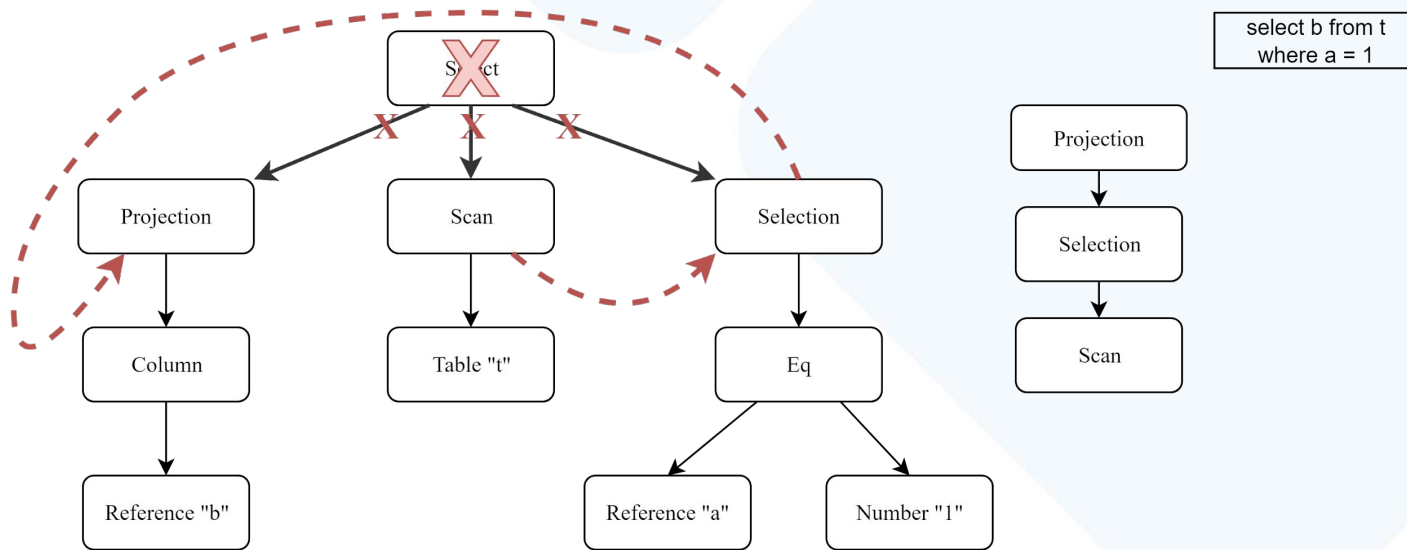
select b from t
where a = 1

IR

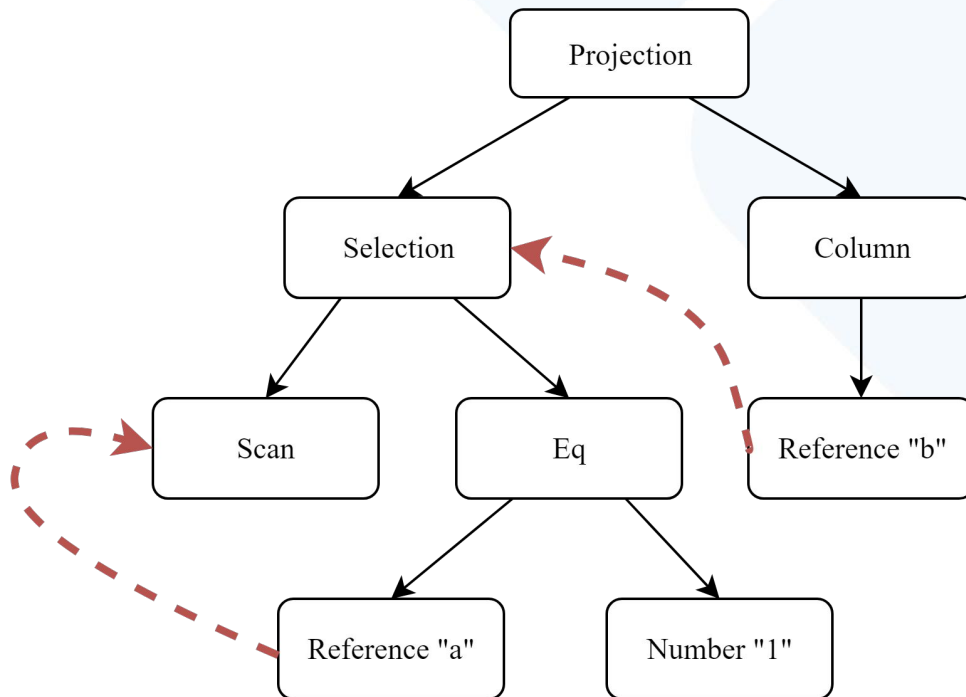
1. IR (intermediate representation) - represents a plan for distributed query execution
2. It is a tree consisting of Relational and Expression nodes
3. Motion - special Relational node, that materializes subtree below it
4. Each subtree below Motion can be converted to valid tarantool's SQL



AST: transform select



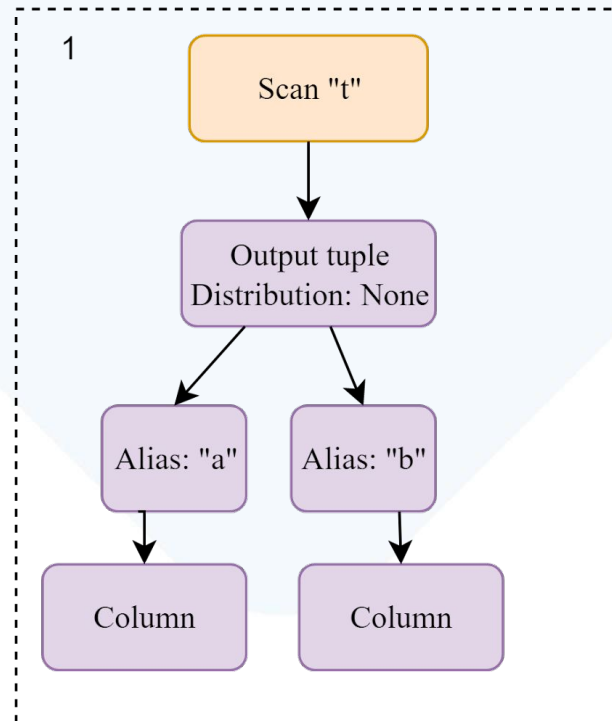
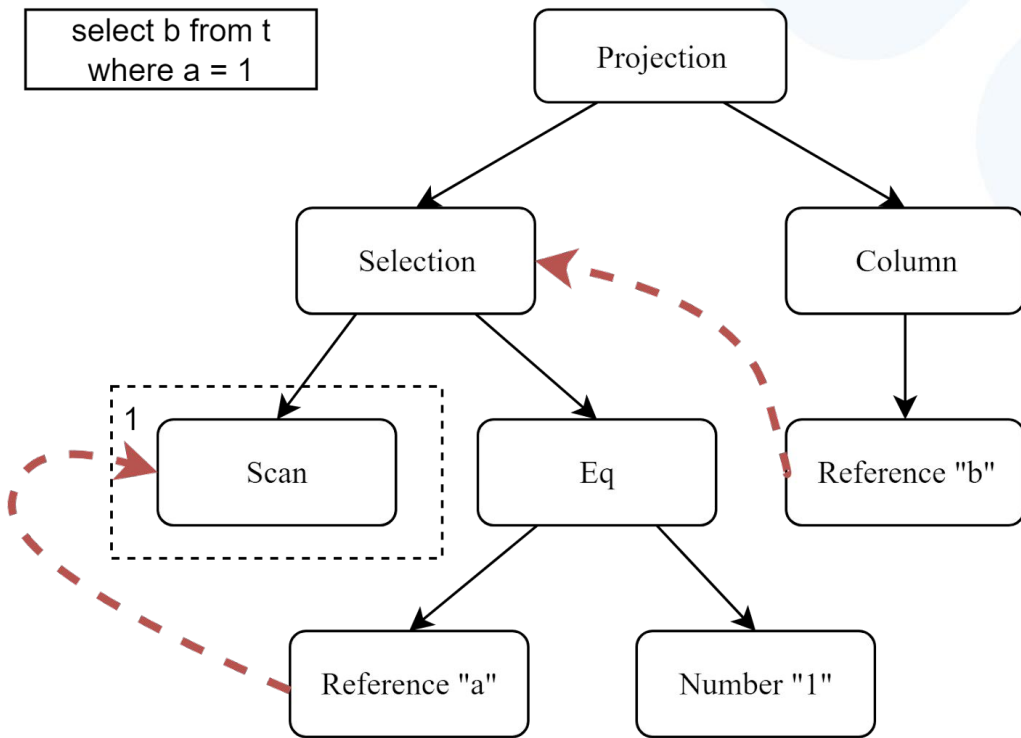
AST: bind references



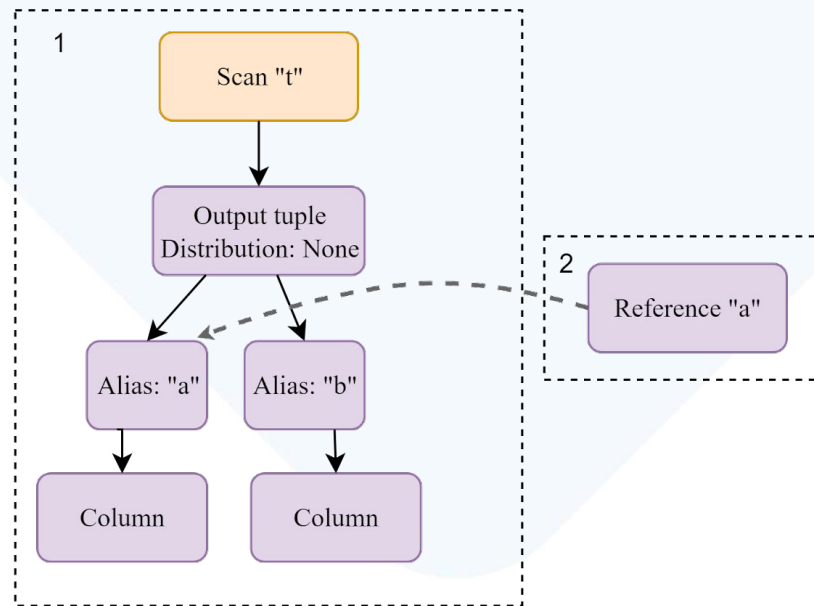
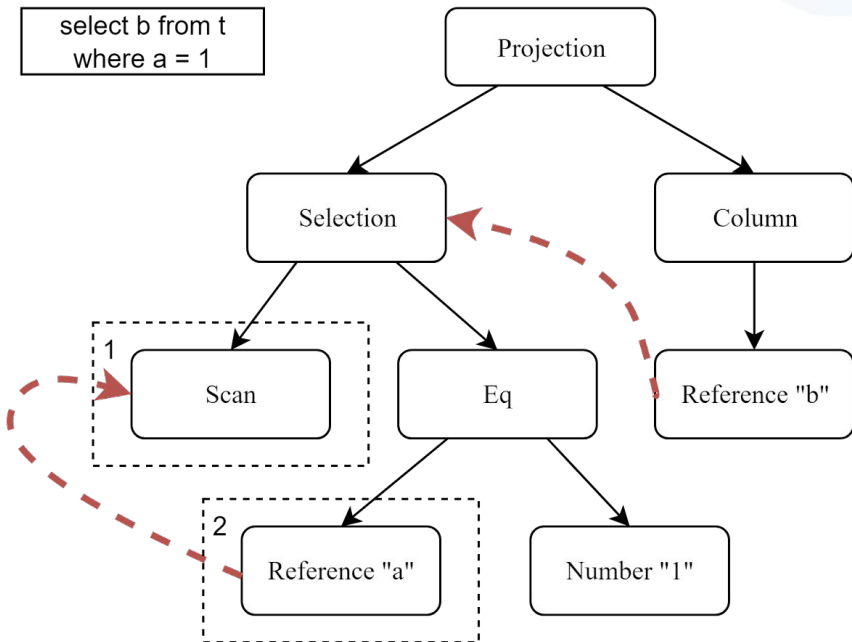
select b from t
where a = 1

AST to raw IR

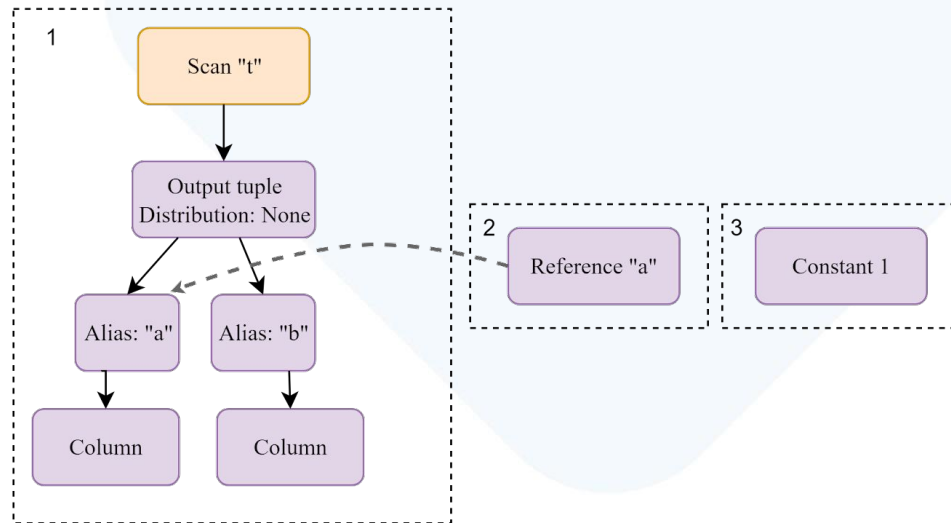
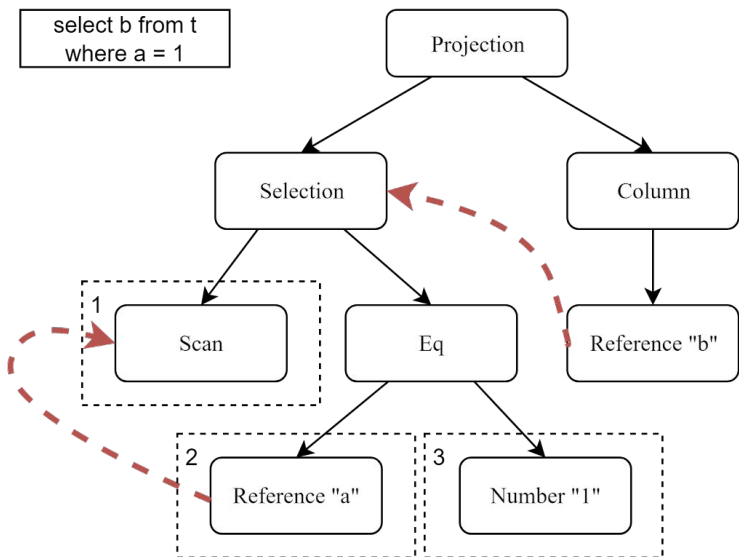
select b from t
where a = 1



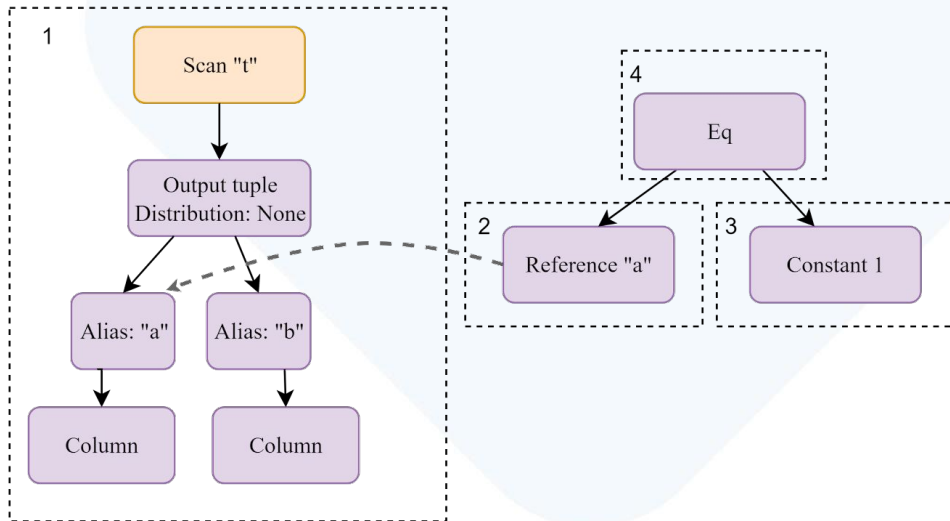
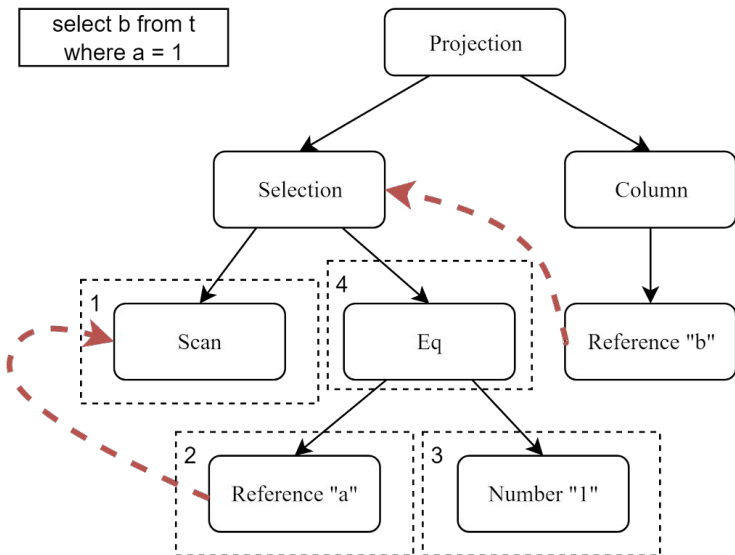
AST to raw IR



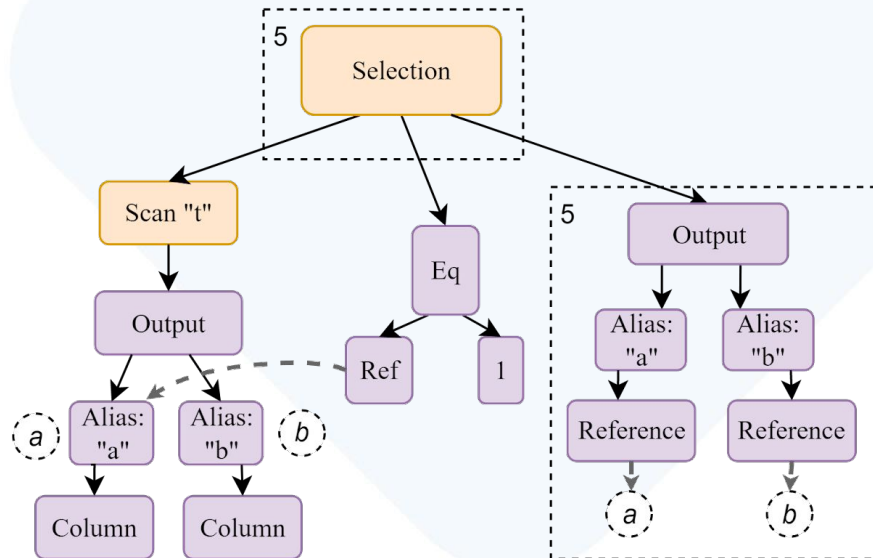
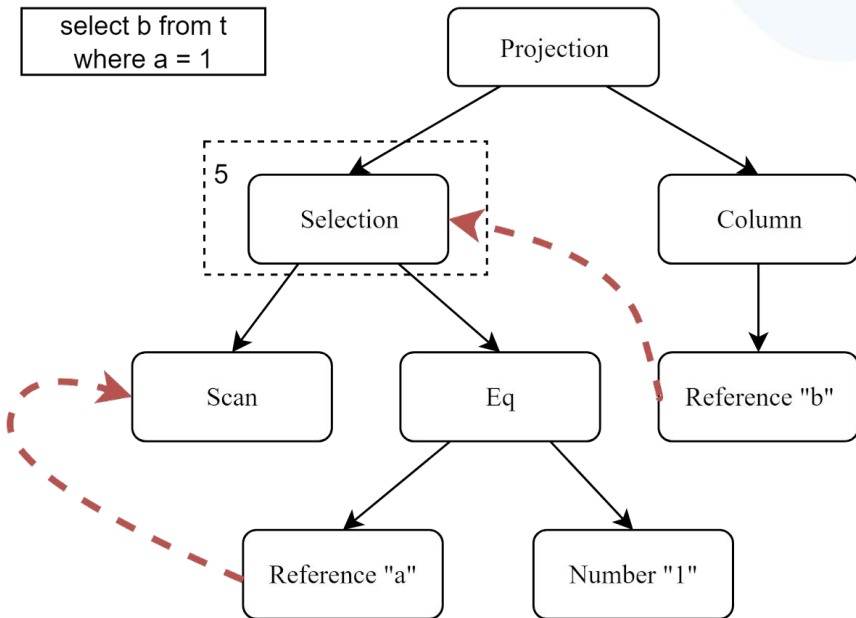
AST to raw IR



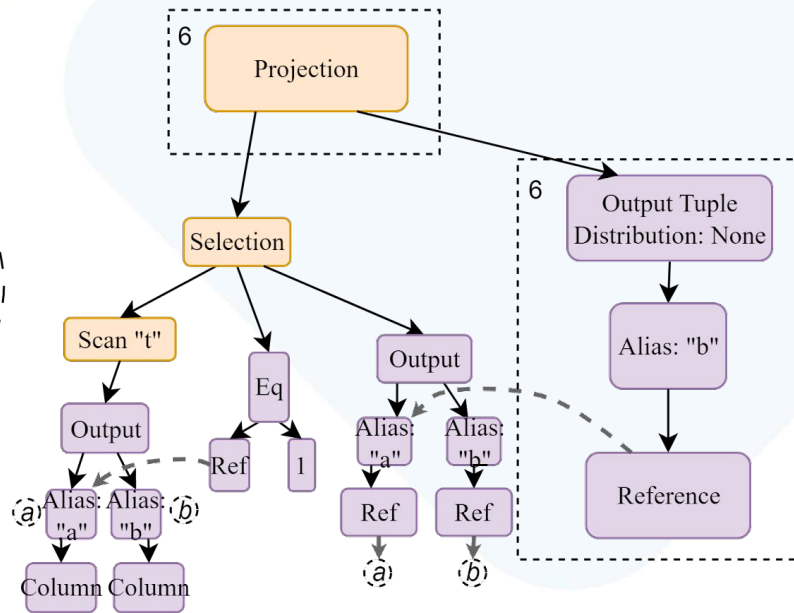
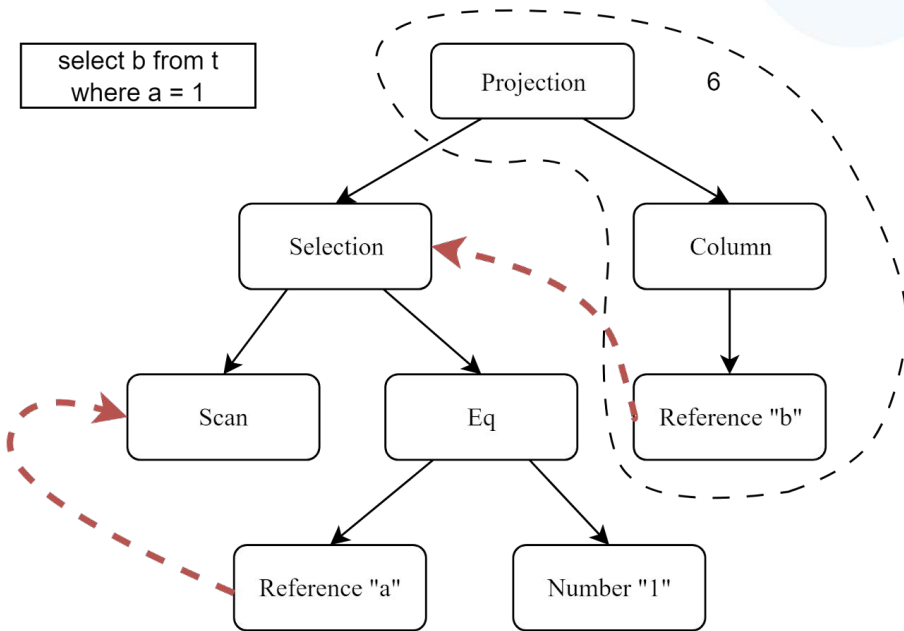
AST to raw IR



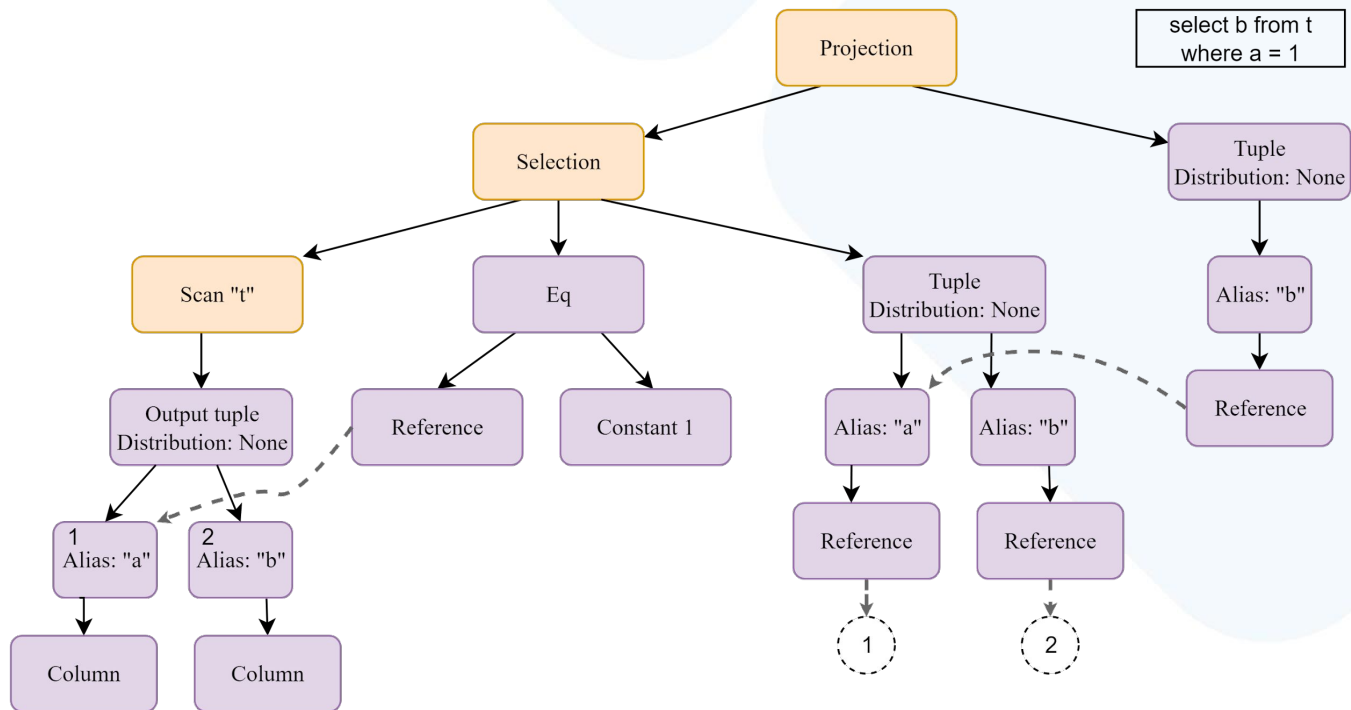
AST to raw IR



AST to raw IR



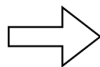
raw IR



Sharded tables

t sharded by a

a	b	bucket_id
1	1	100
1	2	100
150	1	2400



Storage 1
buckets: 1 - 1000

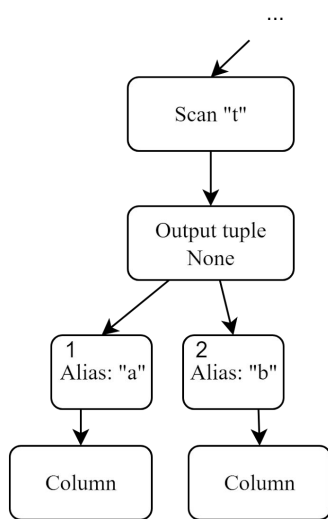
a	b	bucket_id
1	1	100
1	2	100

Storage 2
buckets: 1001 - 3000

a	b	bucket_id
150	1	2400

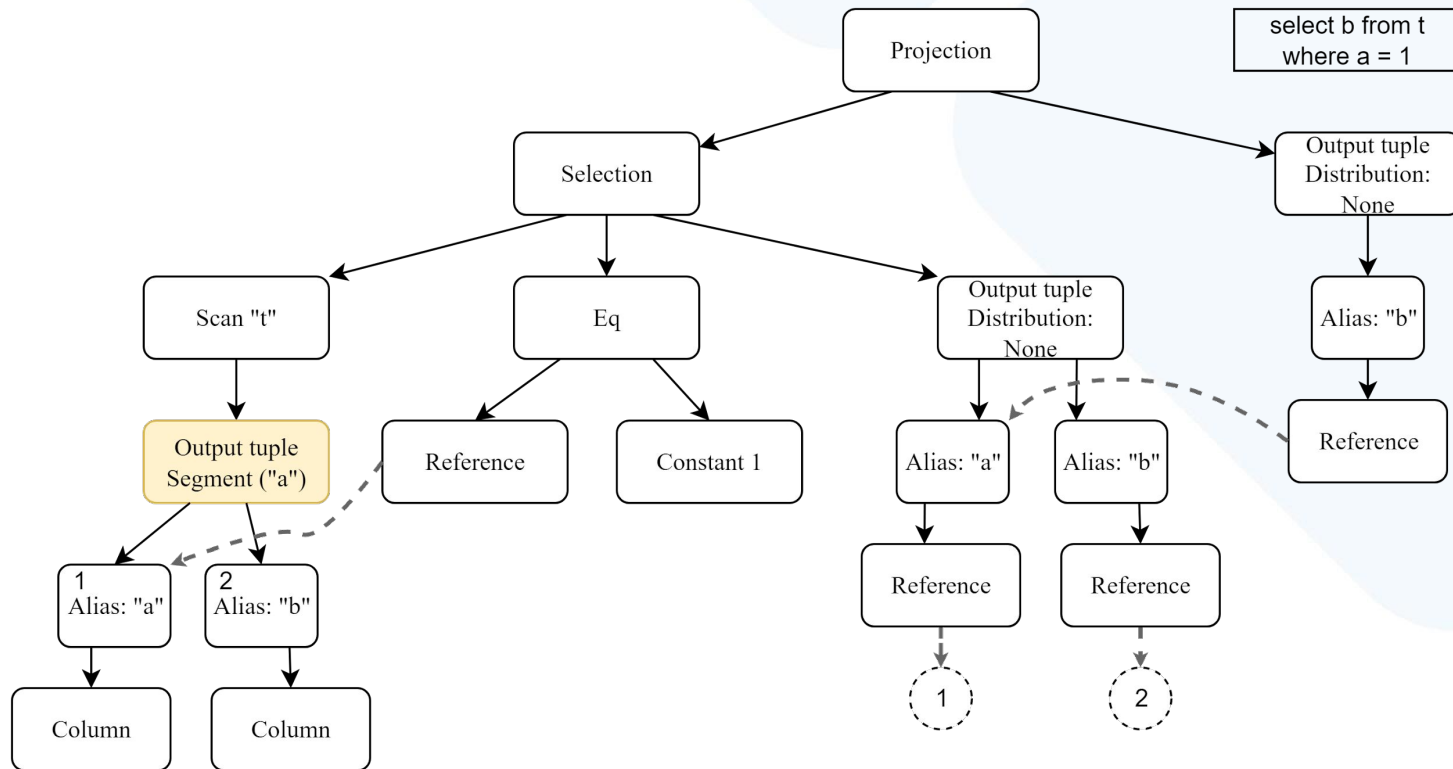
Distributions

- How during IR build, we understand that Motion must be inserted?
- Distribution - properties of output table, that must be true for Relational node output's table after subtree with its node was dispatched and materialized on storages

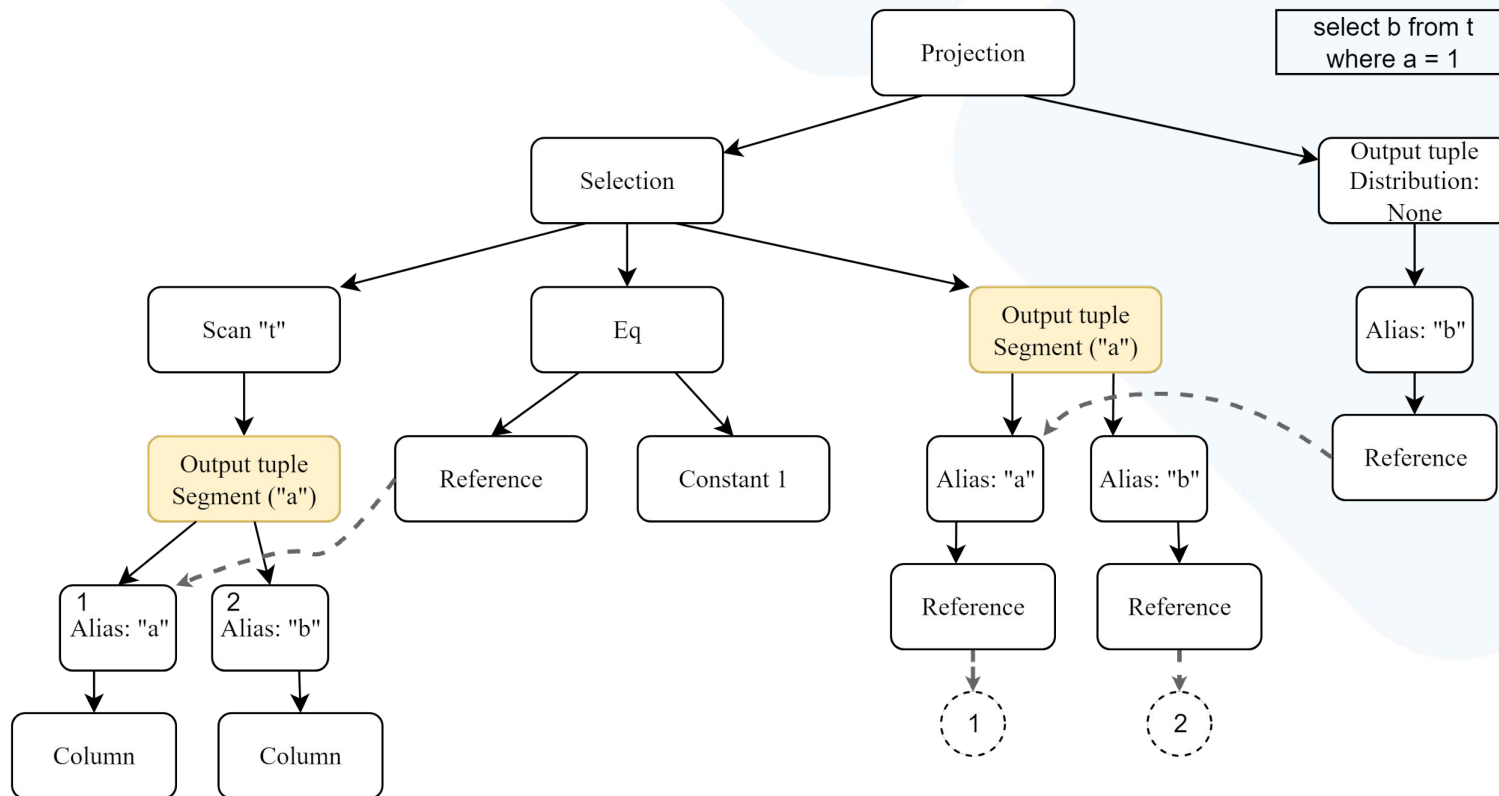


select b from t
where a = 1

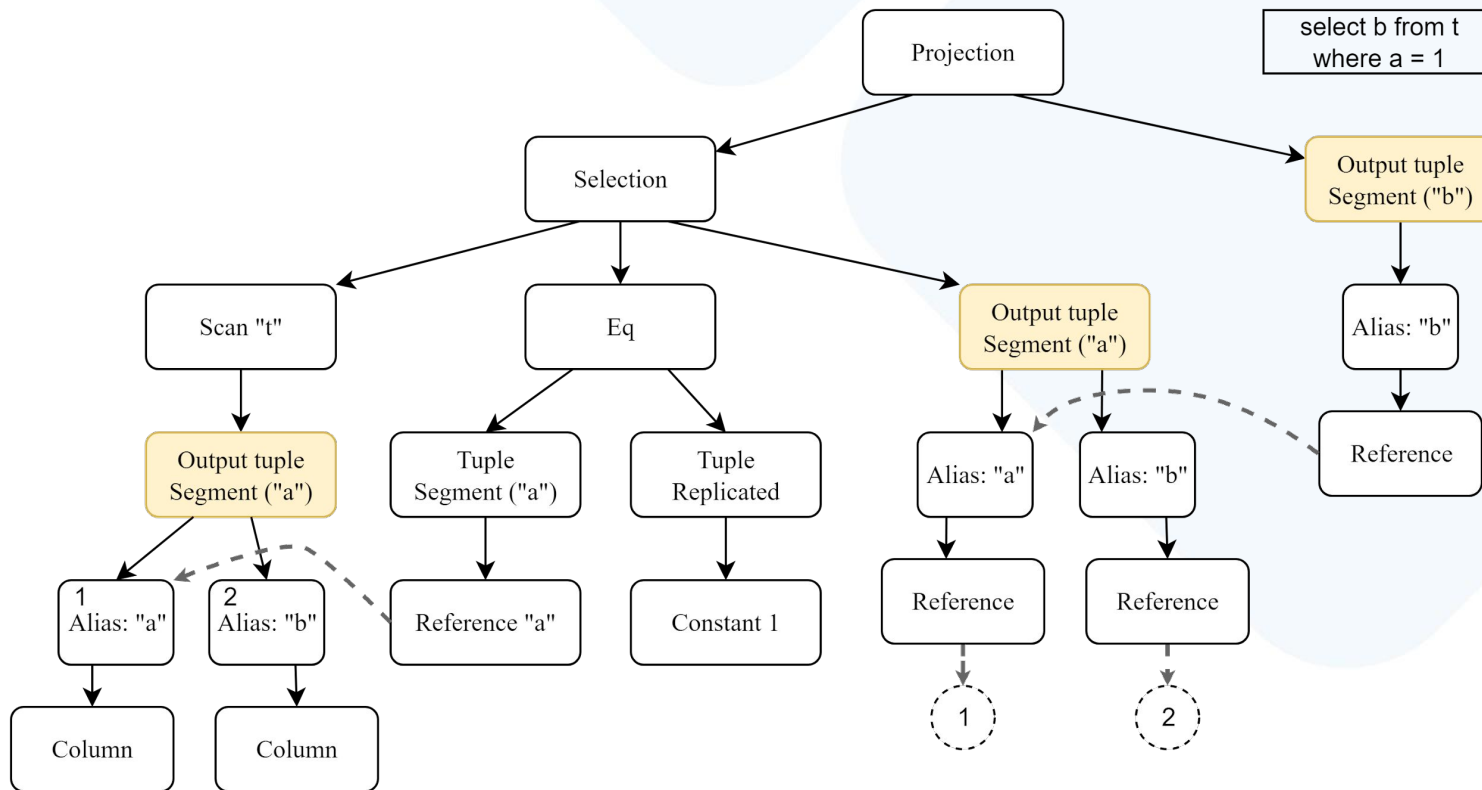
Raw IR to IR



Raw IR to IR

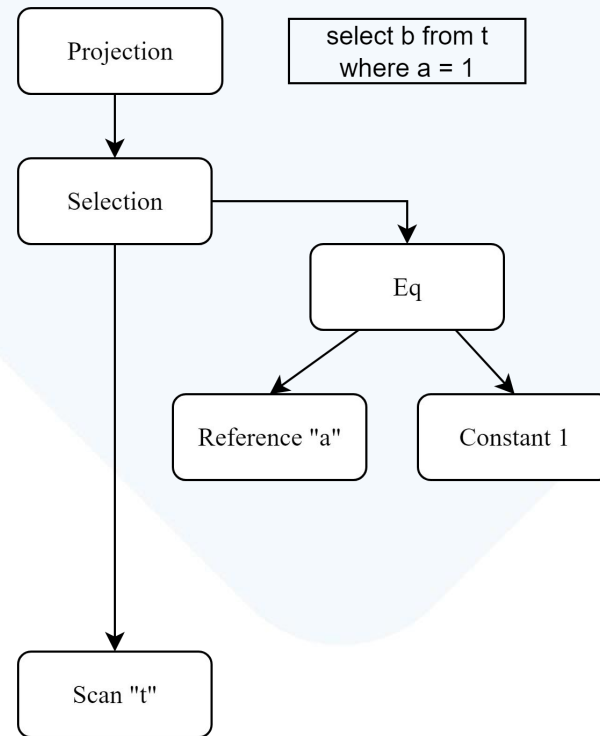


Raw IR to IR

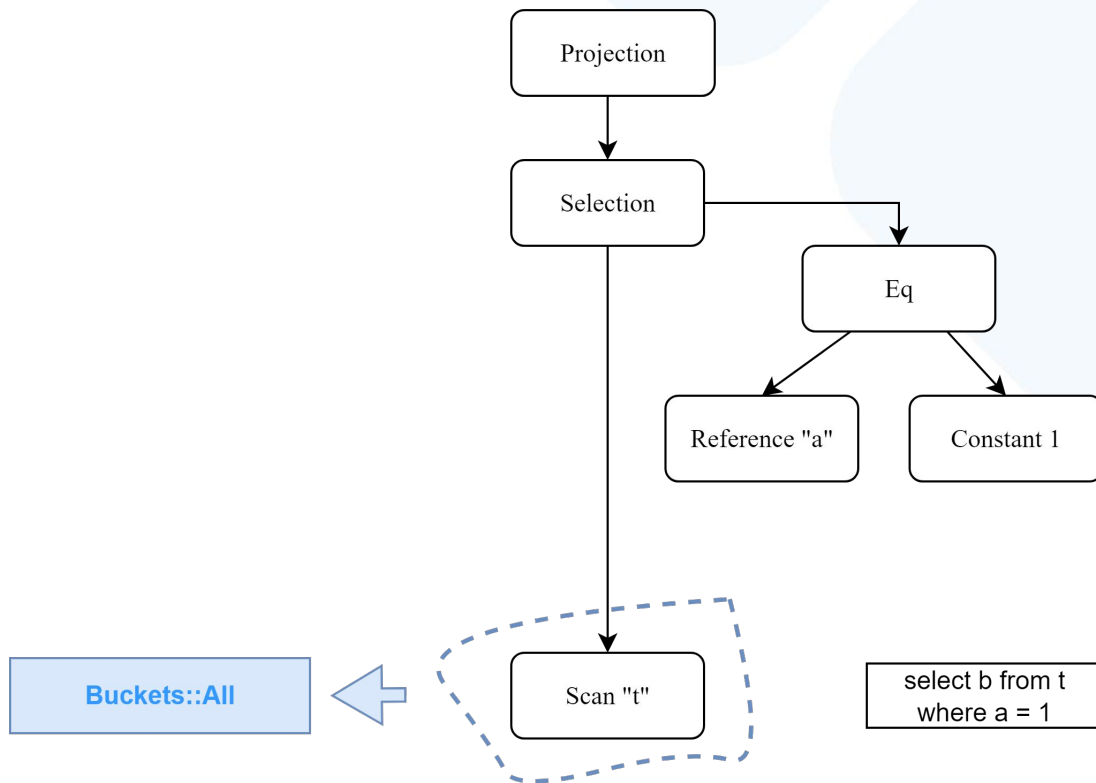


Bucket discovery

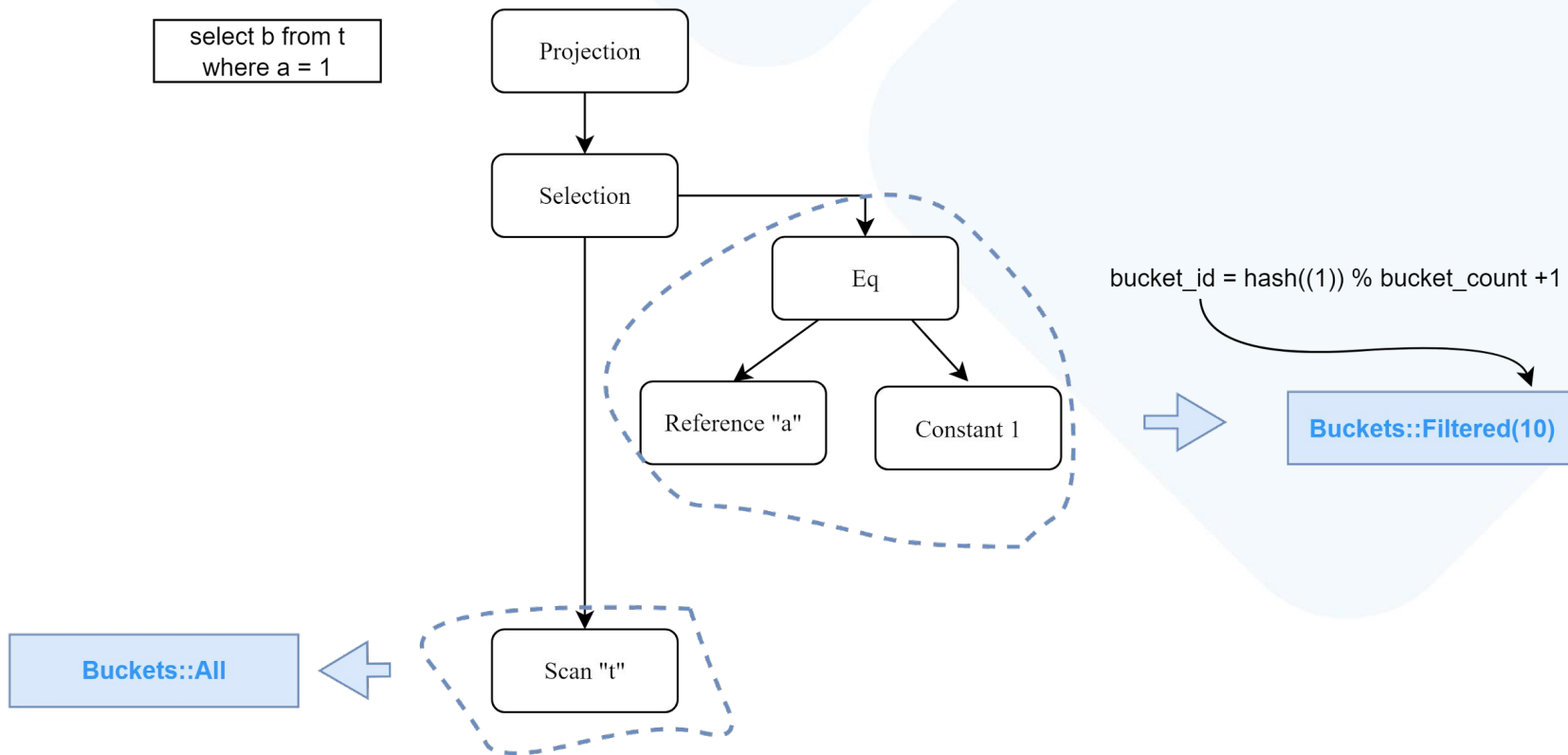
- Bucket discovery - calculation of buckets, where IR subtree must be executed.



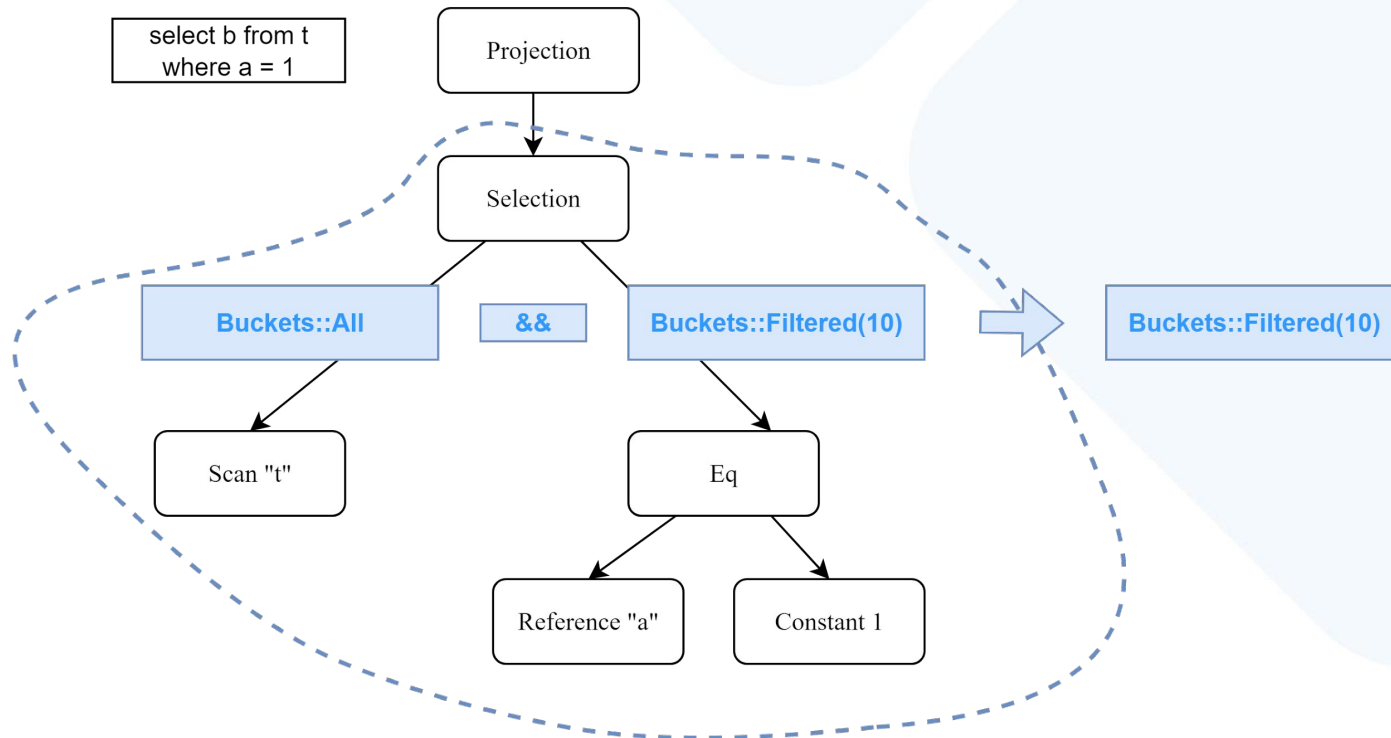
Execution: bucket discovery



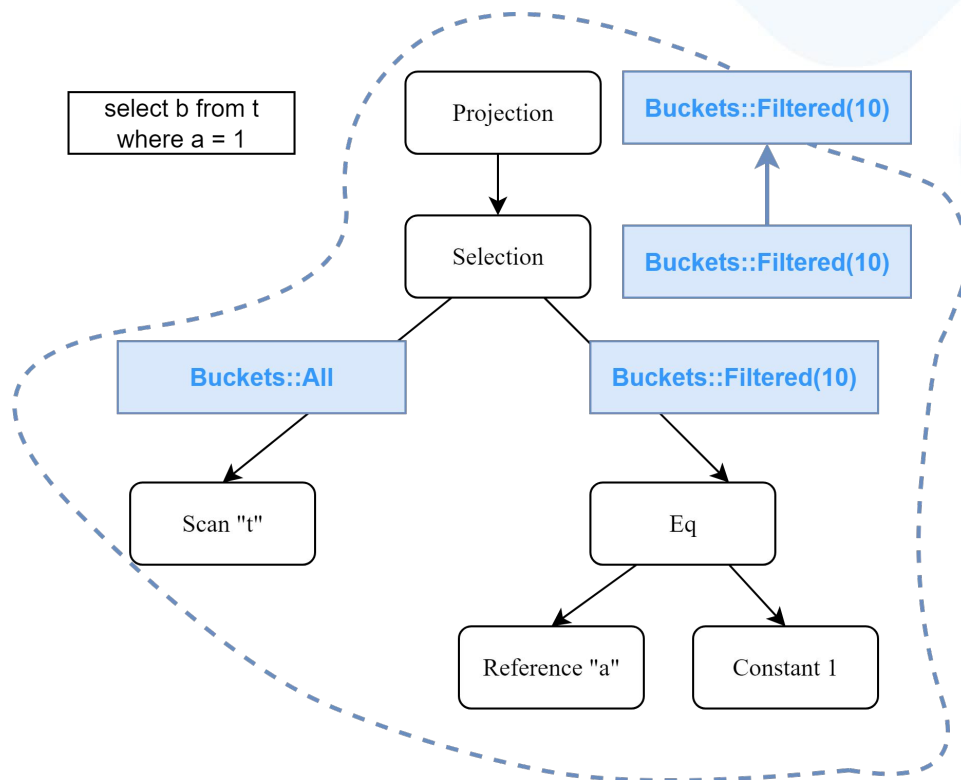
Execution: bucket discovery



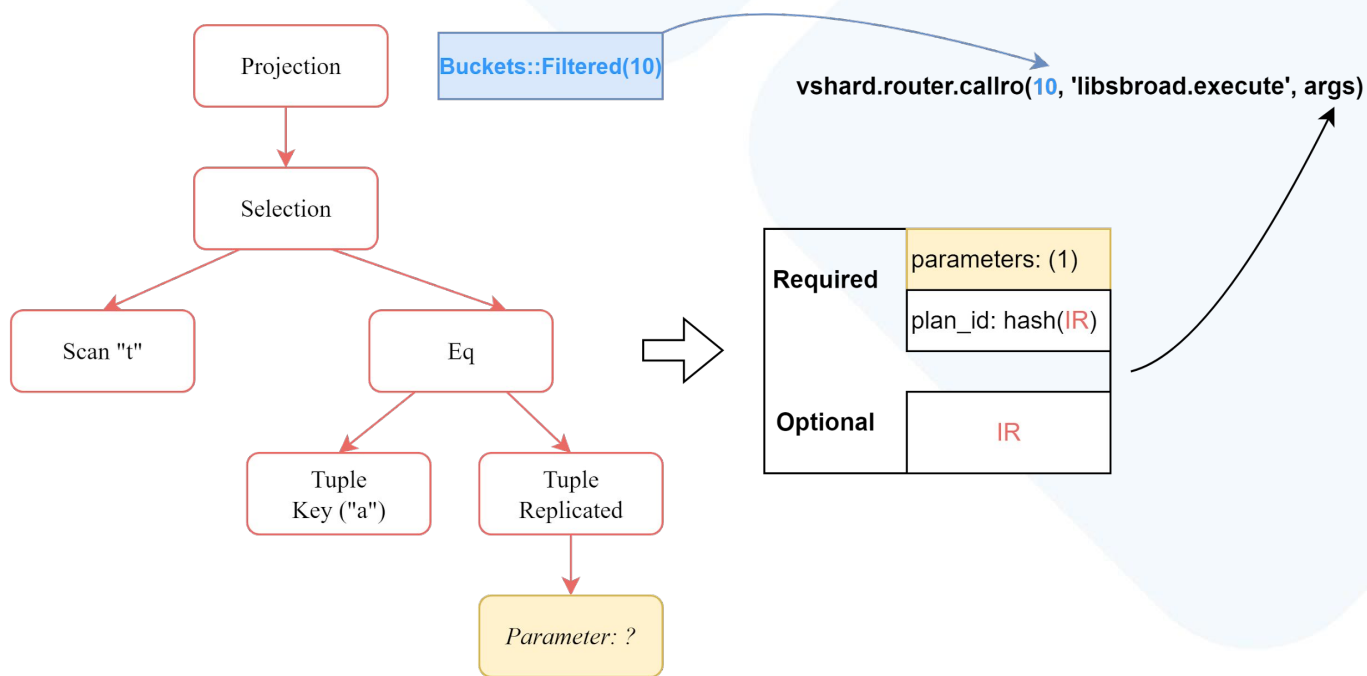
Execution: bucket discovery



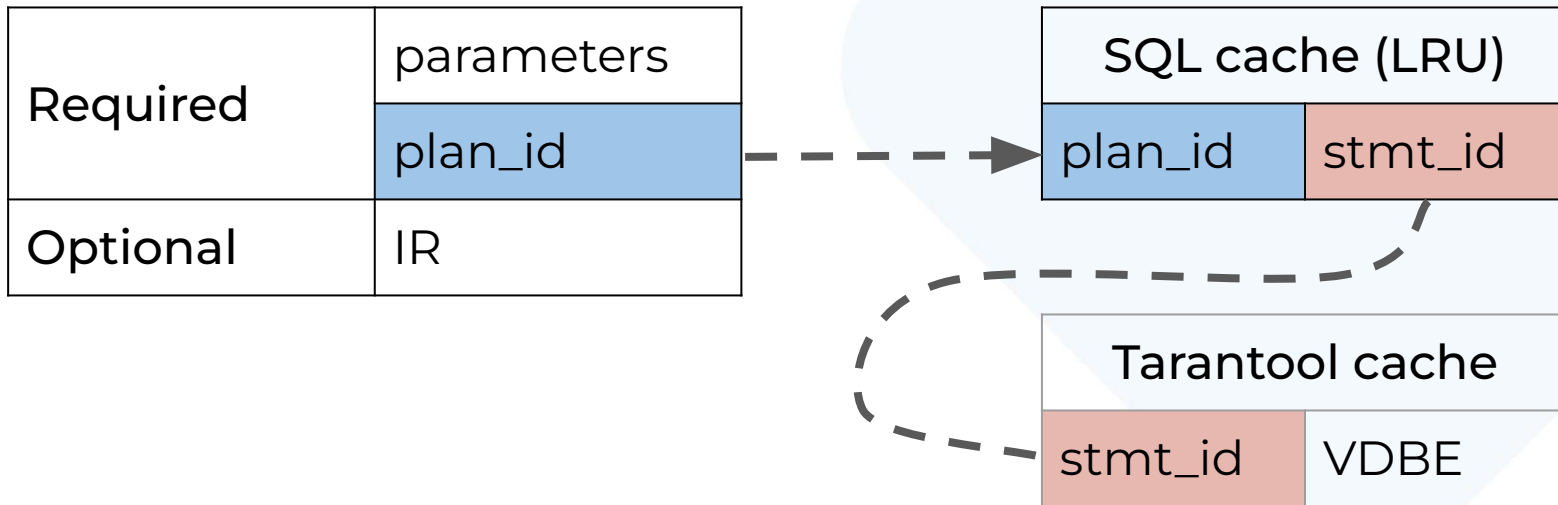
Execution: bucket discovery



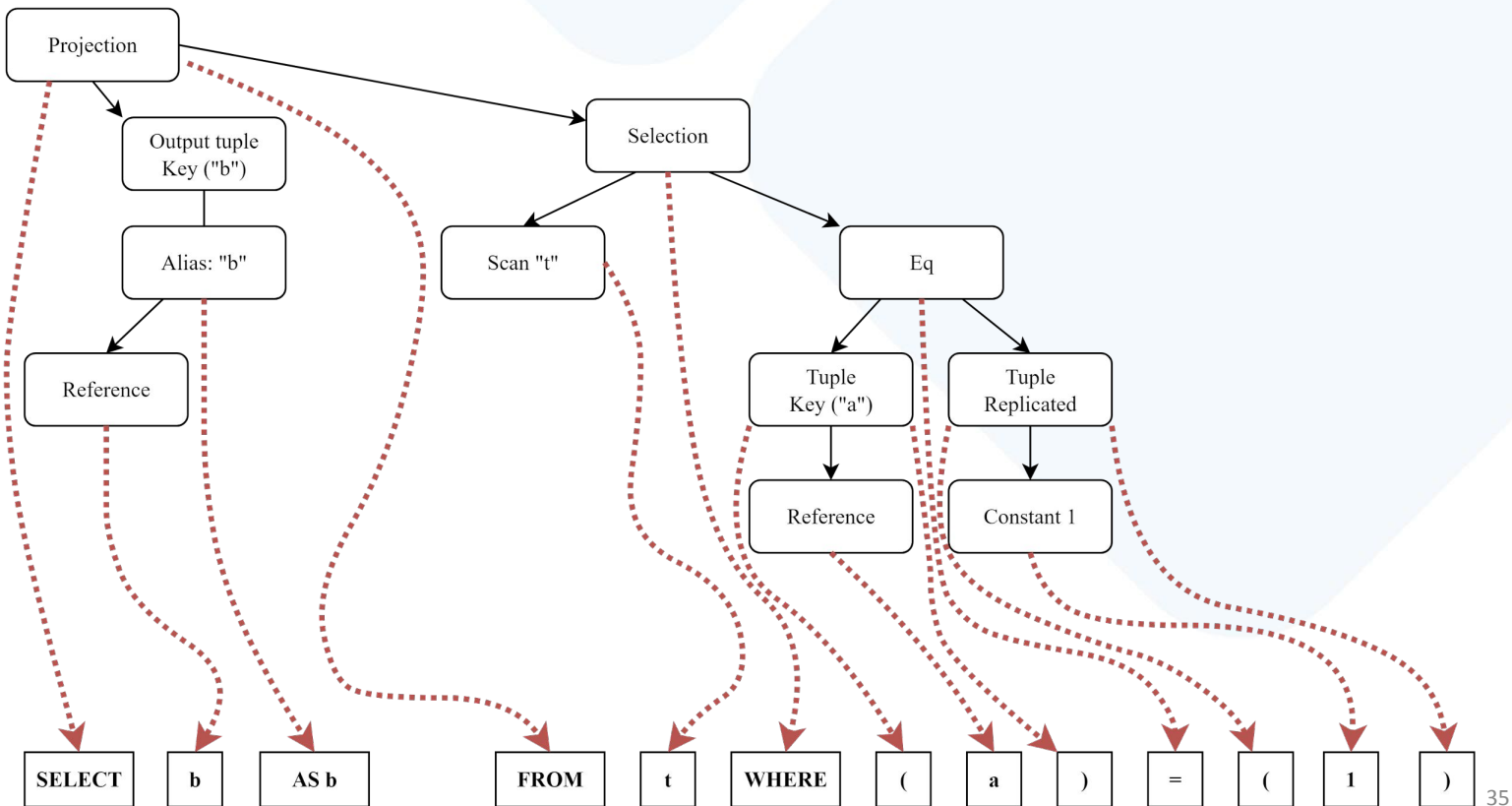
Execution: vshard rpc



Storage cache



IR to SQL



Tarantool SQL: VDBE



```
box.execute([[
    explain select "b" from "t"
    where "a" = 1
]])
```

idx	opcode	p1	p2	p3	p4	p5
0	Init	0	1	0		00
1	IteratorOpen	1	0	0	space<name=t>	02
2	Integer	1	1	0		00
3	IsNull	1	5	0		00
4	MustBeInt	1	10	0		00
5	SeekGE	1	10	1	1	00
6	IdxGT	1	10	1	1	00
7	Column	1	1	2		00
8	ResultRow	2	1	0		00
9	Next	1	6	0		00
10	Halt	0	0	0		00

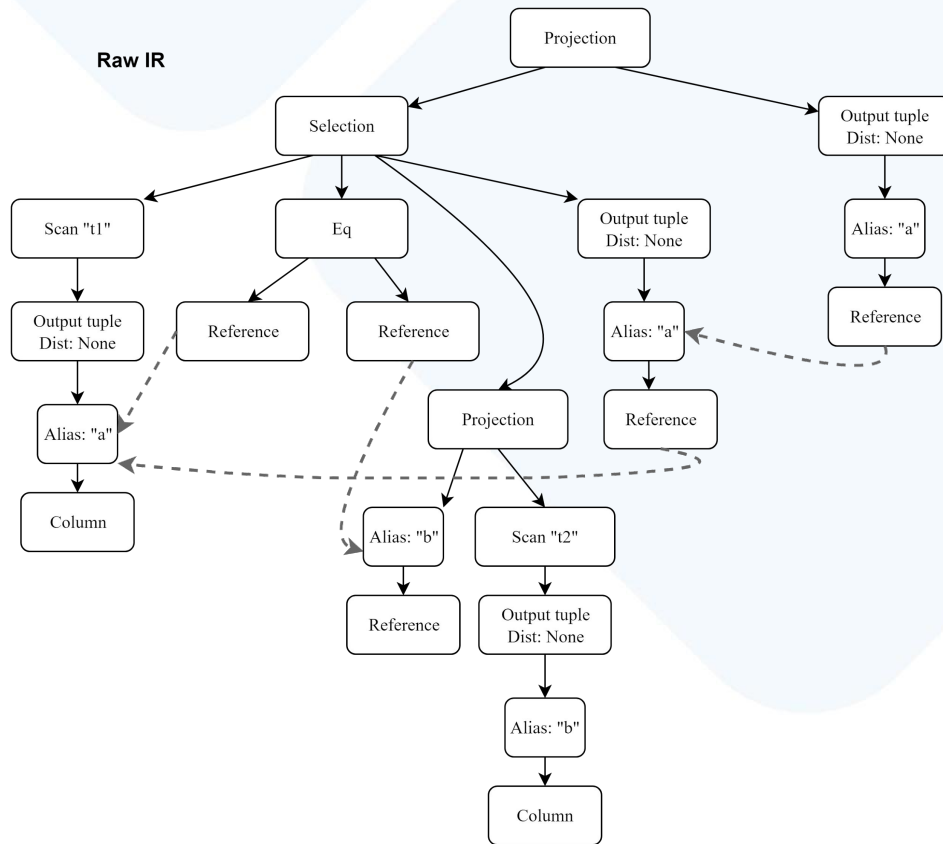
Example 2

Query

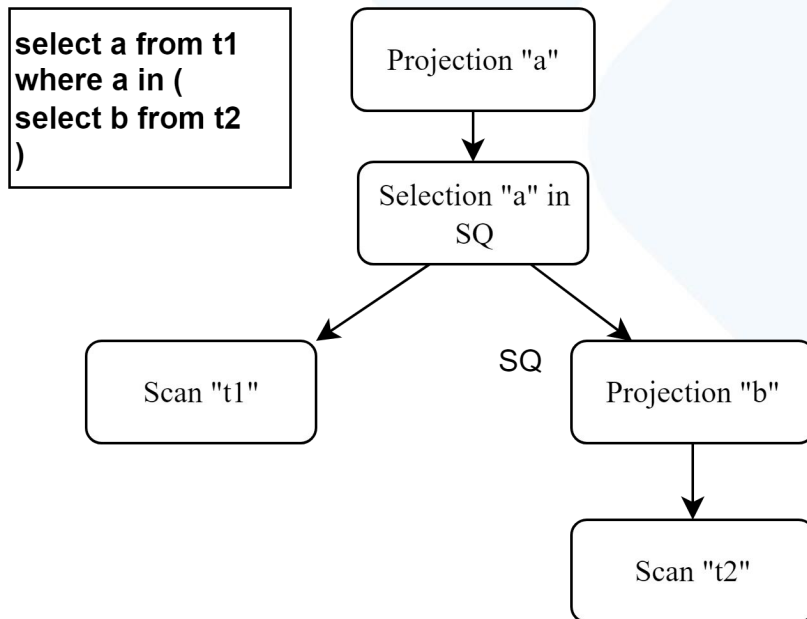
```
select a from t1 where a in (select b from t2)
```

Raw IR

select a from t1
where a in
(select b from t2)



Simplified raw IR



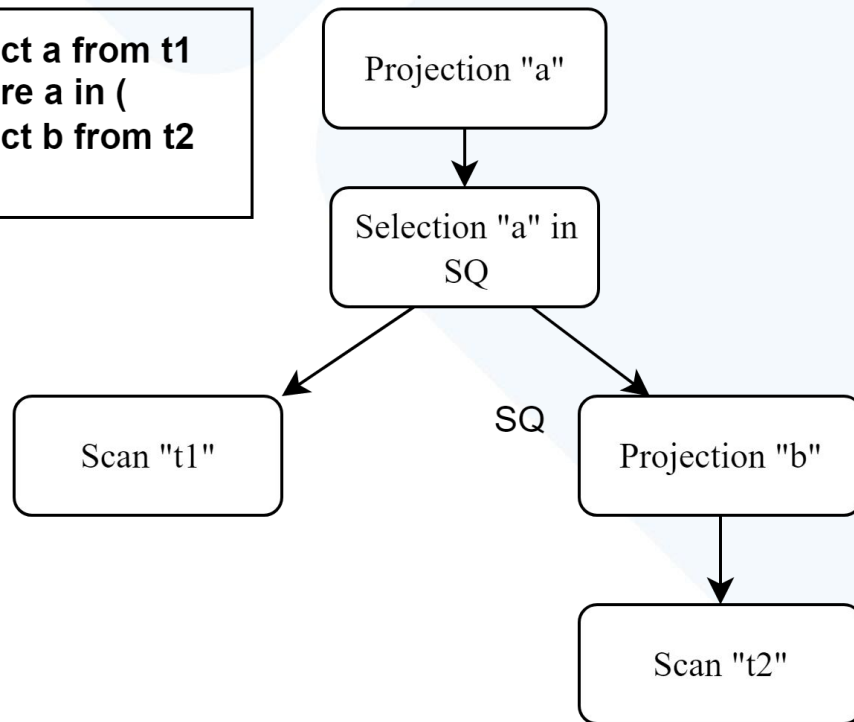
Case 1: $t_1(a)$, $t_2(b)$

Initial conditions

t1 sharded by "a"

t2 sharded by "b"

```
select a from t1
where a in (
select b from t2
)
```

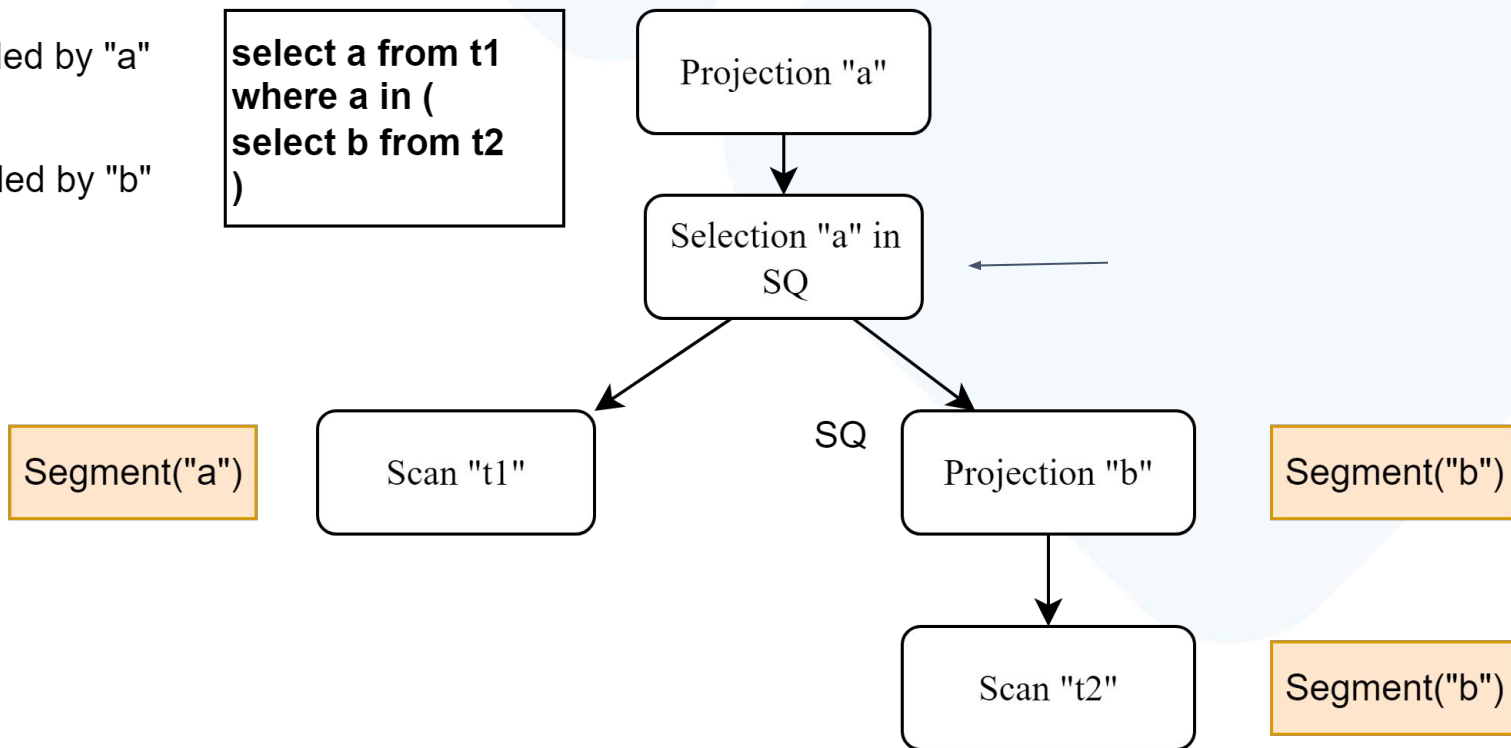


Distribution conflict?

t1 sharded by "a"

t2 sharded by "b"

**select a from t1
where a in (
select b from t2
)**

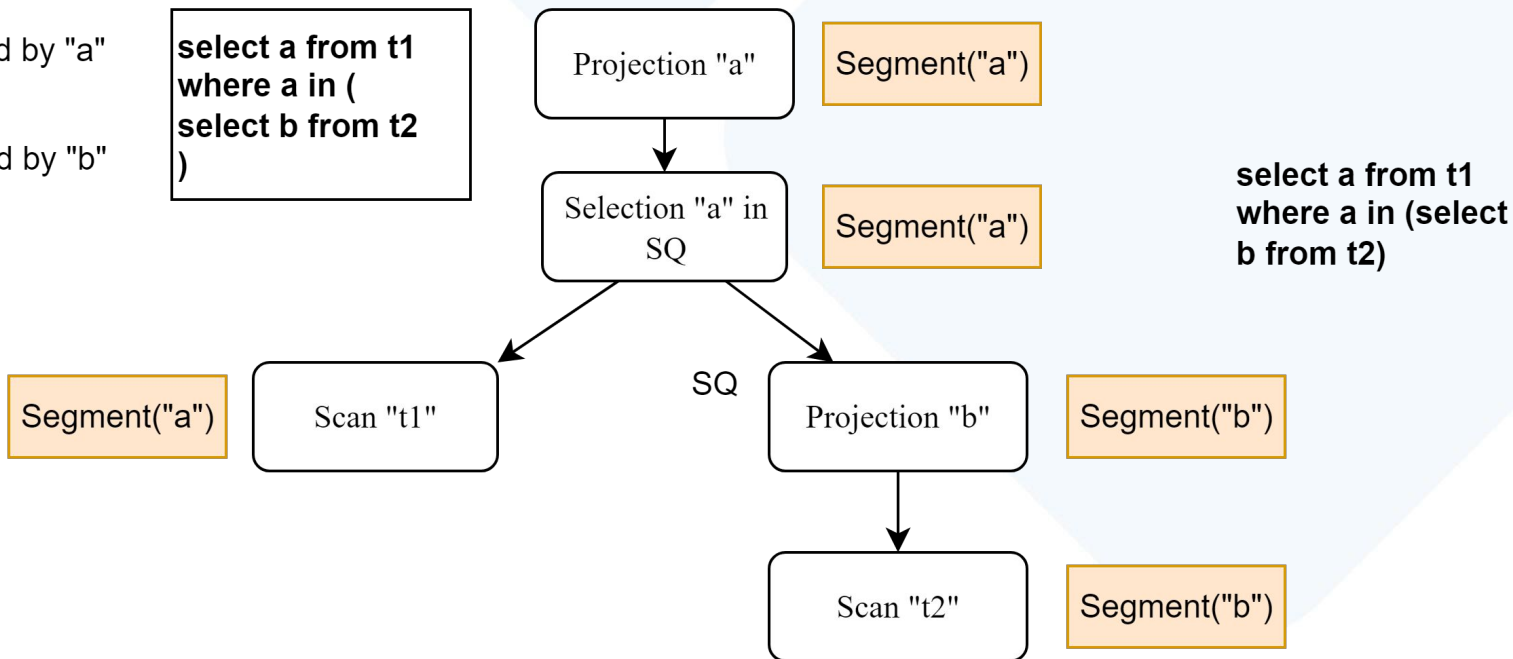


Answer

t1 sharded by "a"

t2 sharded by "b"

**select a from t1
where a in (
select b from t2
)**



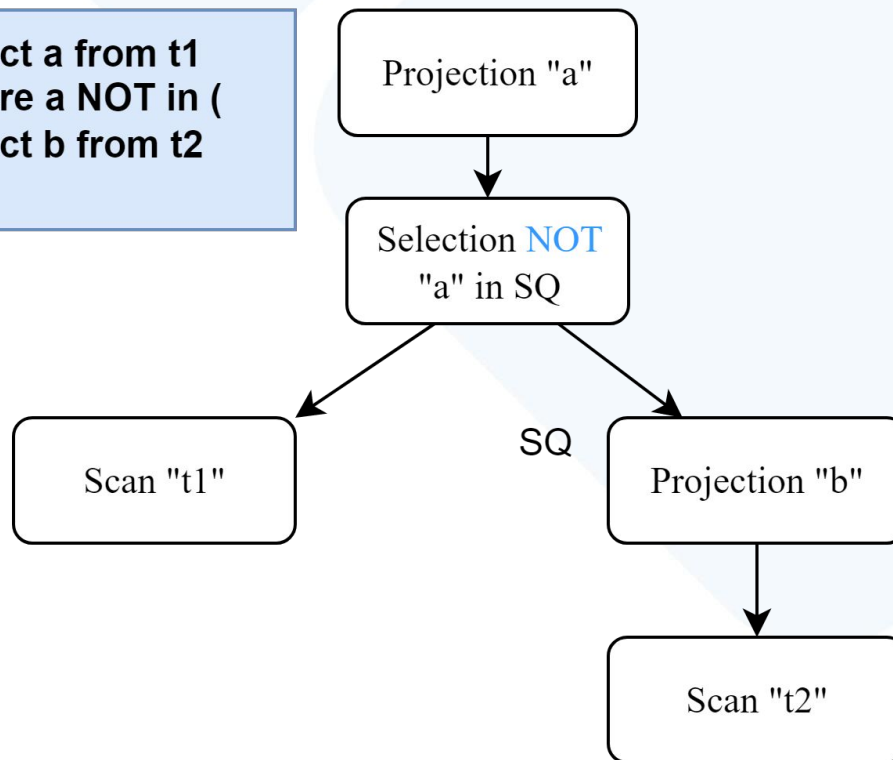
Case 2: different filter

Different filter

t1 sharded by "a"

t2 sharded by "b"

```
select a from t1
where a NOT in (
  select b from t2
)
```

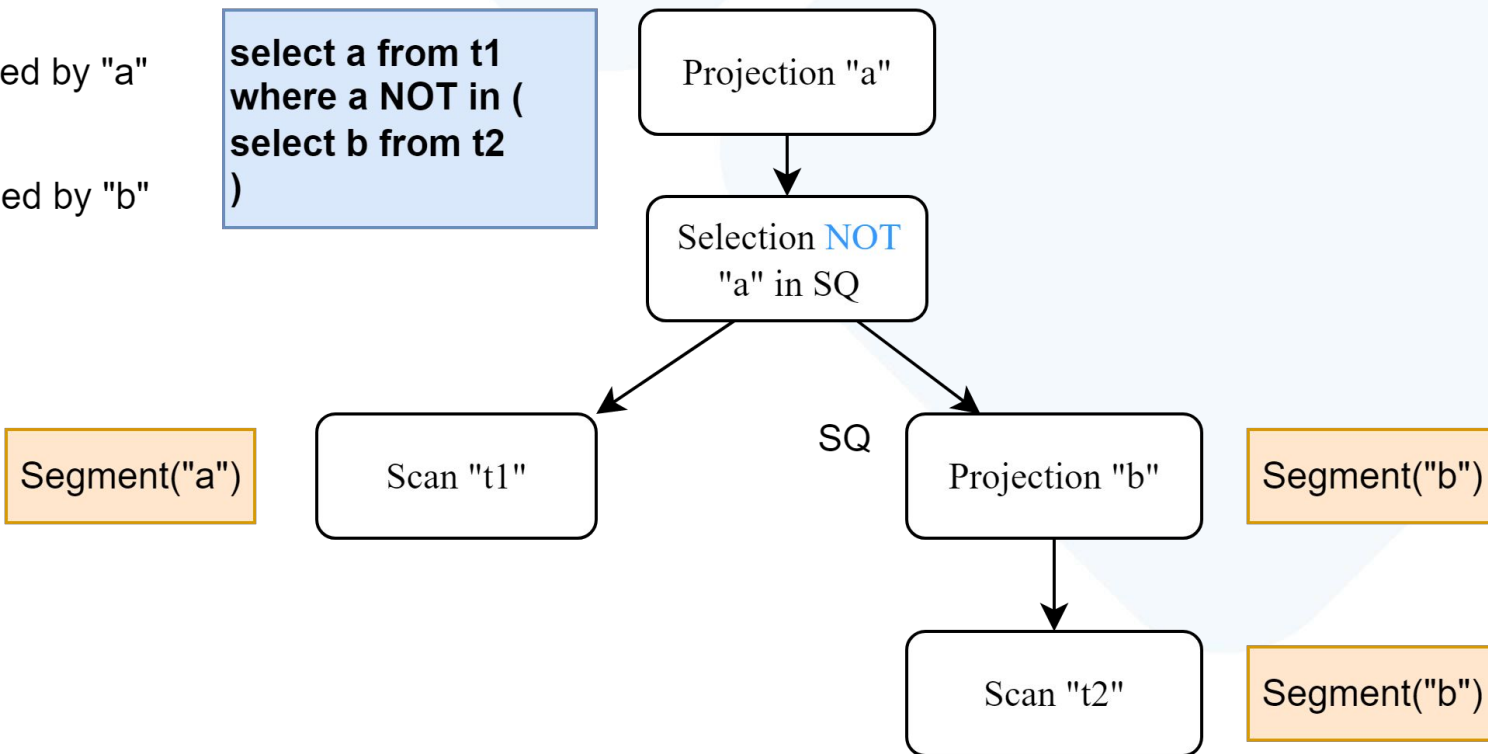


Different filter

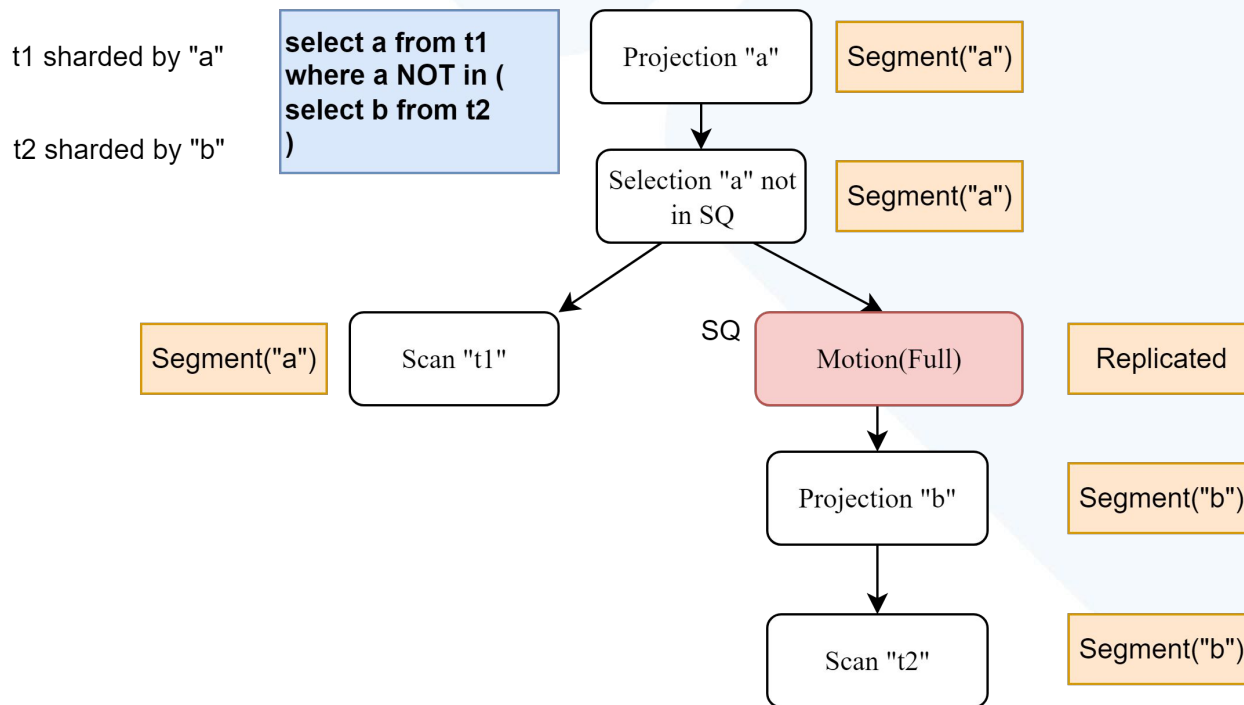
t1 sharded by "a"

t2 sharded by "b"

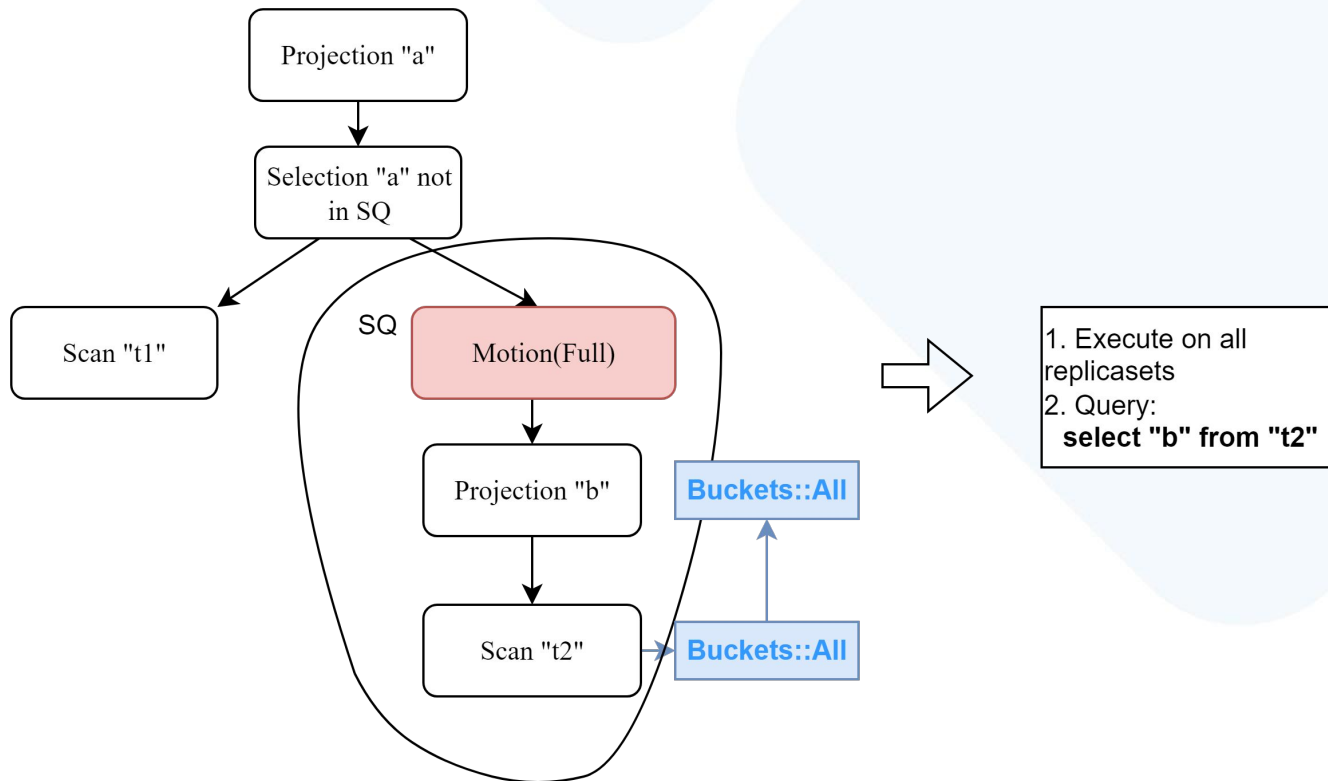
**select a from t1
where a NOT in (
select b from t2
)**



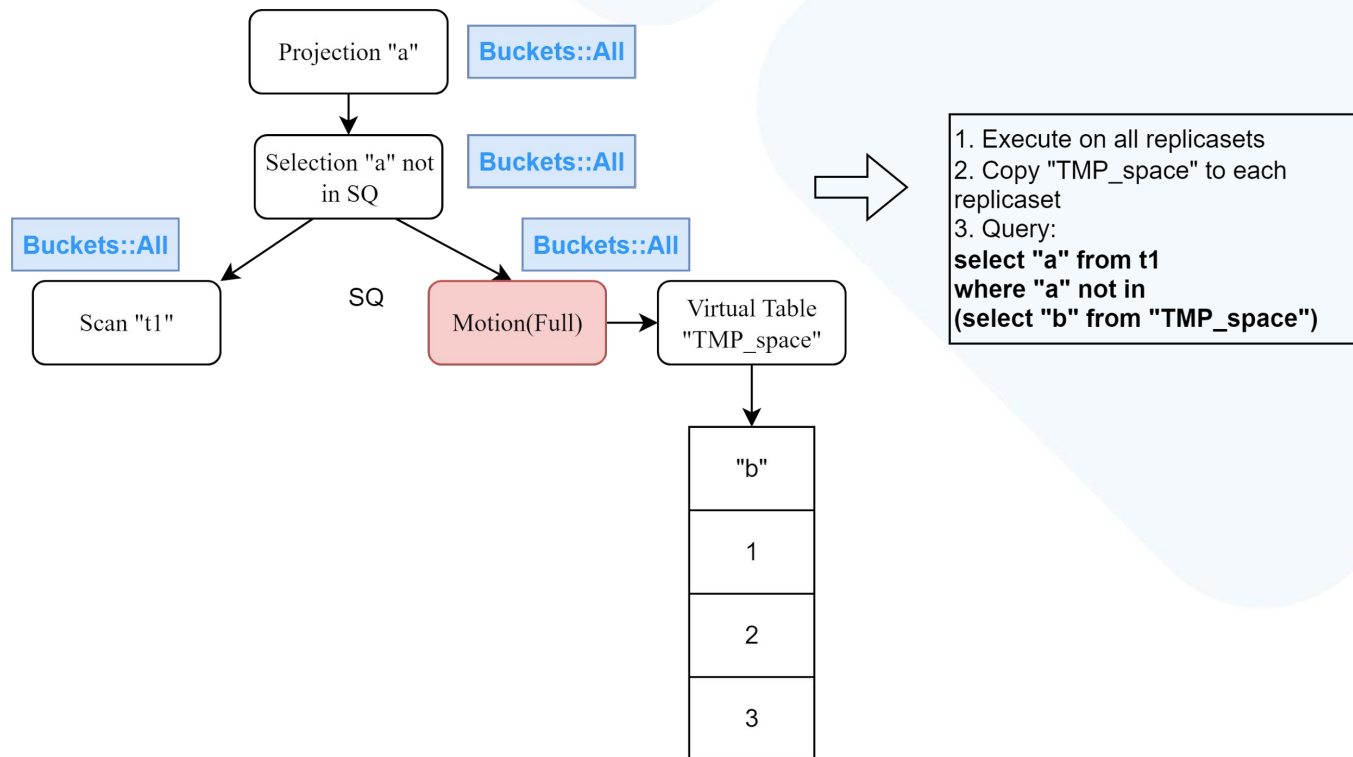
Motion with Full policy



Motion with Full policy

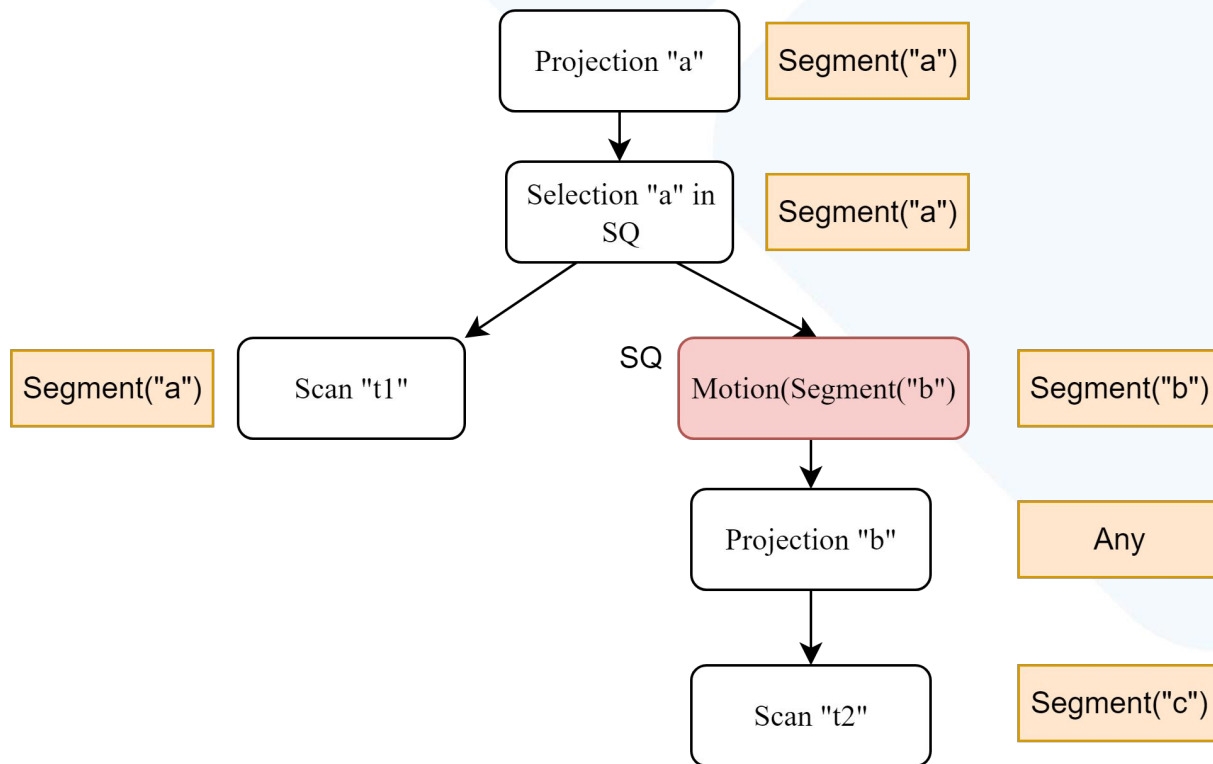


Motion with Full policy

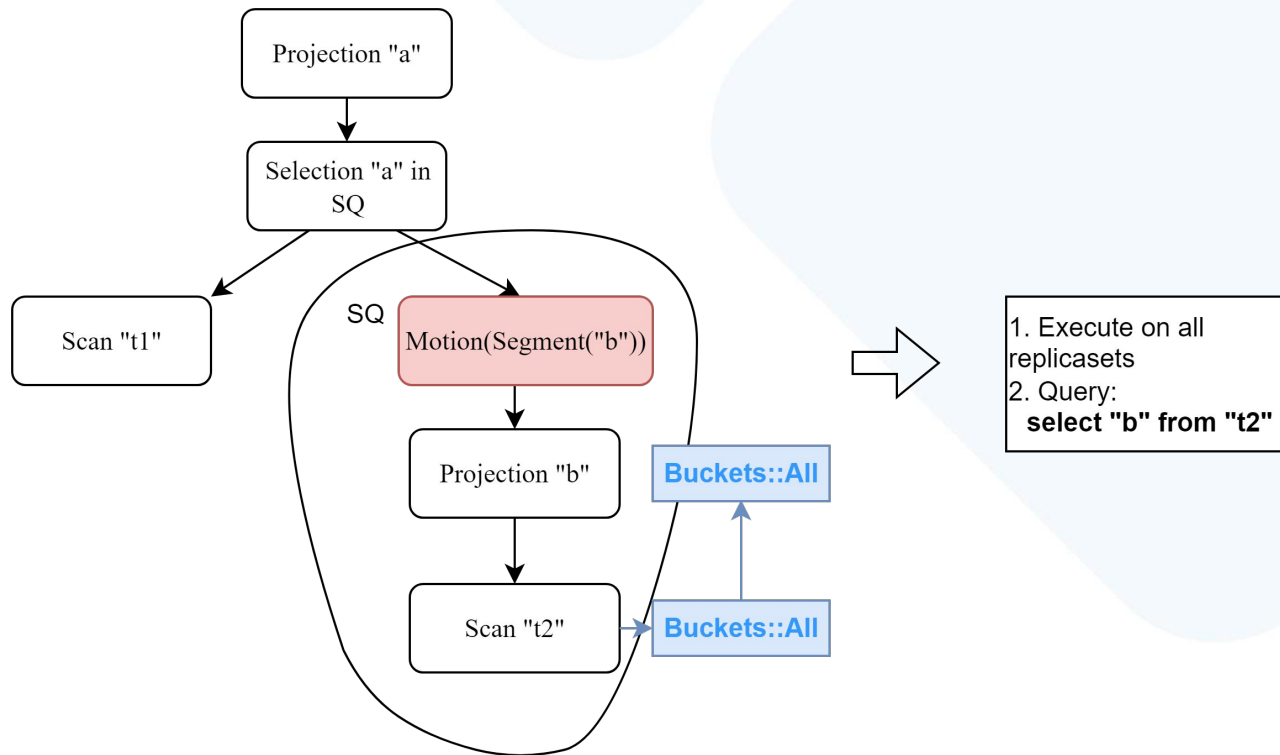


Case 3: $t_1(a)$, $t_2(c)$

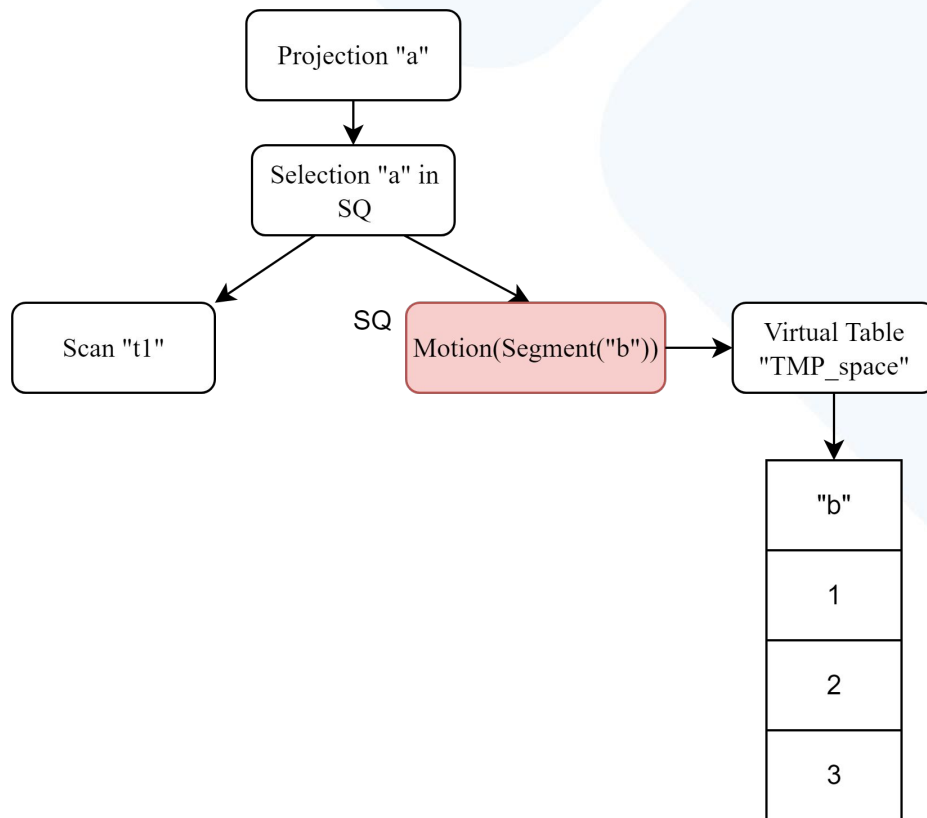
Motion with Segment policy



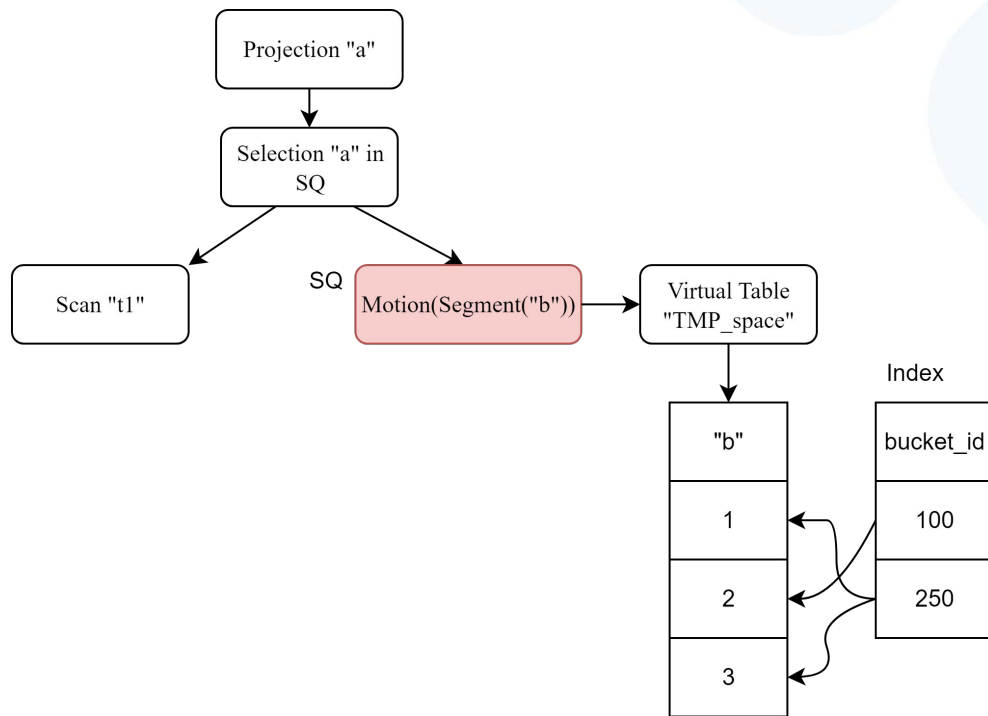
Motion(Segment): execution



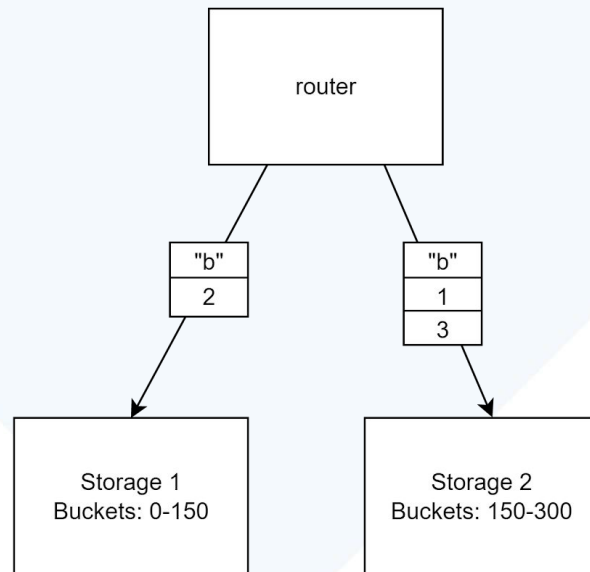
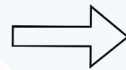
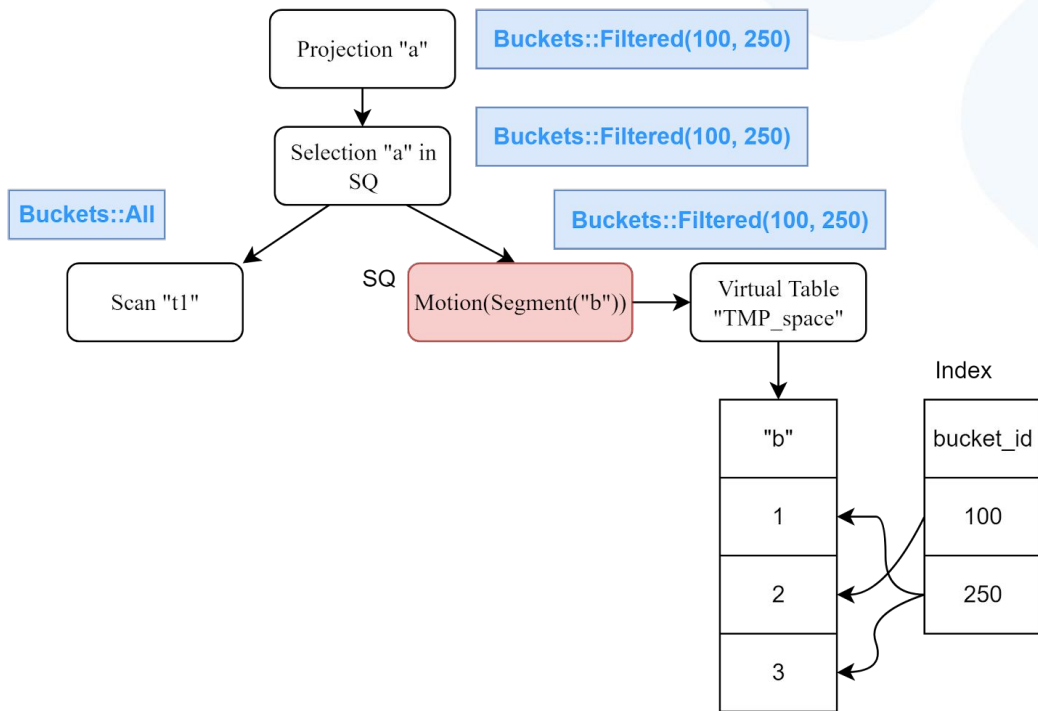
Motion(Segment): Linking



Motion(Segment): Index



Motion(Segment): sending



The background is a solid red color with a complex, abstract pattern of overlapping, rounded geometric shapes in various shades of red, creating a textured, layered effect.

DML query

DML

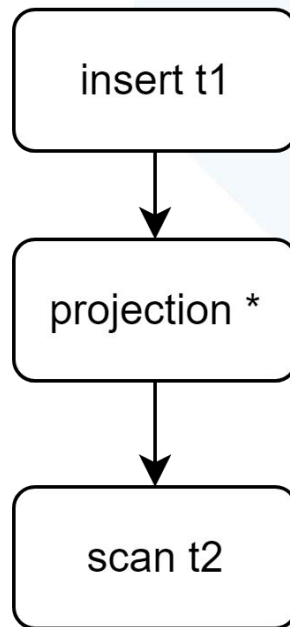
insert into "t1"
select * from "t2"

t1

a	b	bucket_id
---	---	-----------

t2

c	d	bucket_id
---	---	-----------



Do we need a motion?

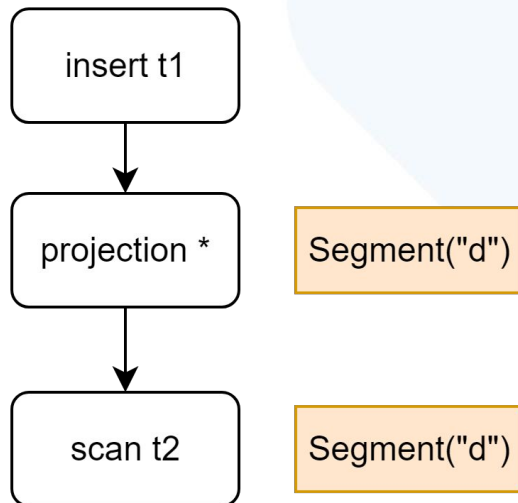
insert into "t1"
select * from "t2"

t1(a)

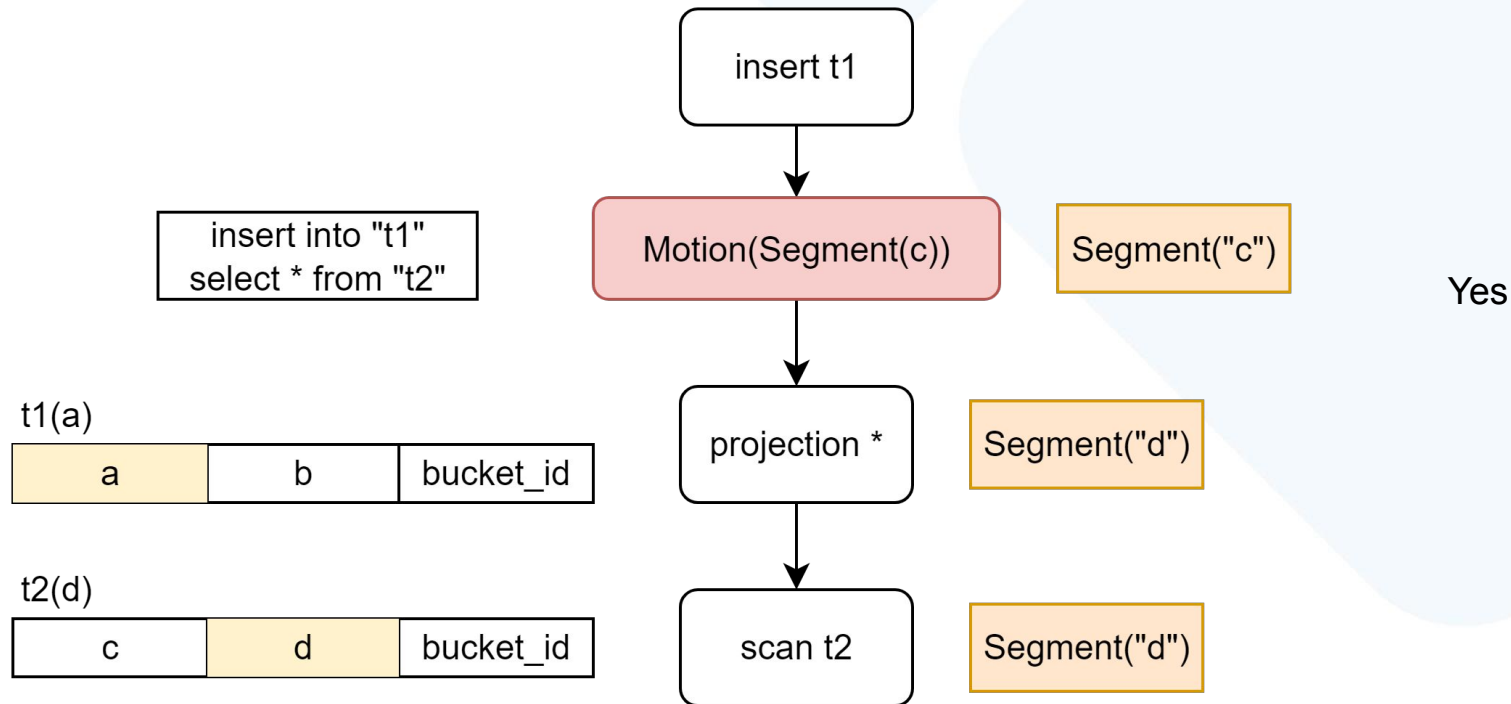
a	b	bucket_id
---	---	-----------

t2(d)

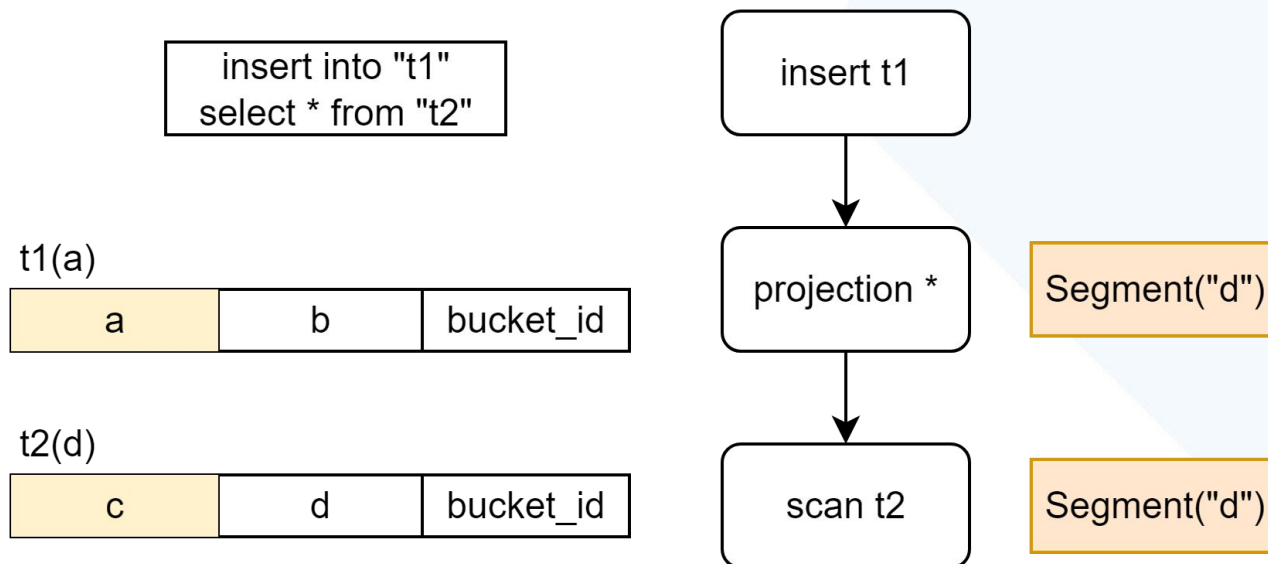
c	d	bucket_id
---	---	-----------



Motion with Segment policy

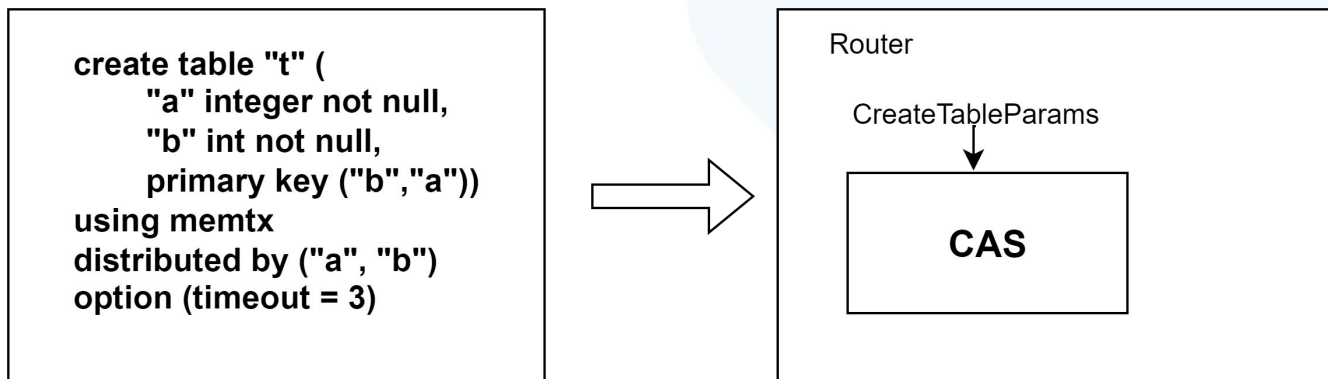


Motion with Segment policy



The background is a solid red color with a complex, abstract pattern of overlapping, rounded geometric shapes in various shades of red, creating a sense of depth and movement.

DDL query



**Thanks for your
attention!**

Glossary



1. IR (intermediate representation) - plan of a query in picosql library.
2. Executor - module responsible for plan (IR) execution in picosql
3. VDBE - internal representation of sql query in Tarantool
4. Virtual Table - temporary table during distributed query execution
5. Router - node where IR is built and where Executor works.
6. Storage - node where data is stored and where execution happens.