

Introduction

Problem statement:

Computing fine grained provenance at interactive speeds

Q : SELECT * from Personal_Info WHERE age < 30 ORDER BY age DESC; LQ : Lineage(Q, 0);

Personal_Info

Name	Age
Alice	25
Jack	31
Bob	26

0
1
2

Result of Filter

Name	Age
Alice	25
Bob	26

0
1

Result of Order By

Name	Age
Bob	26
Alice	25

1

Lineage of Filter

iid	oid
0	0
2	1



Lineage of Order By

iid	oid
1	0
0	1

DuckDB

Personal_Info

Name	Age
Alice	25
Jack	31
Bob	26

Age < 30

Age
25
31
26

0
2

Smoked Duck

Personal_Info

Name	Age
Alice	25
Jack	31
Bob	26

Filter

Age < 30

DuckDB

0
2

move

Operator Lineage

pinned operator lineage in memory

Background:

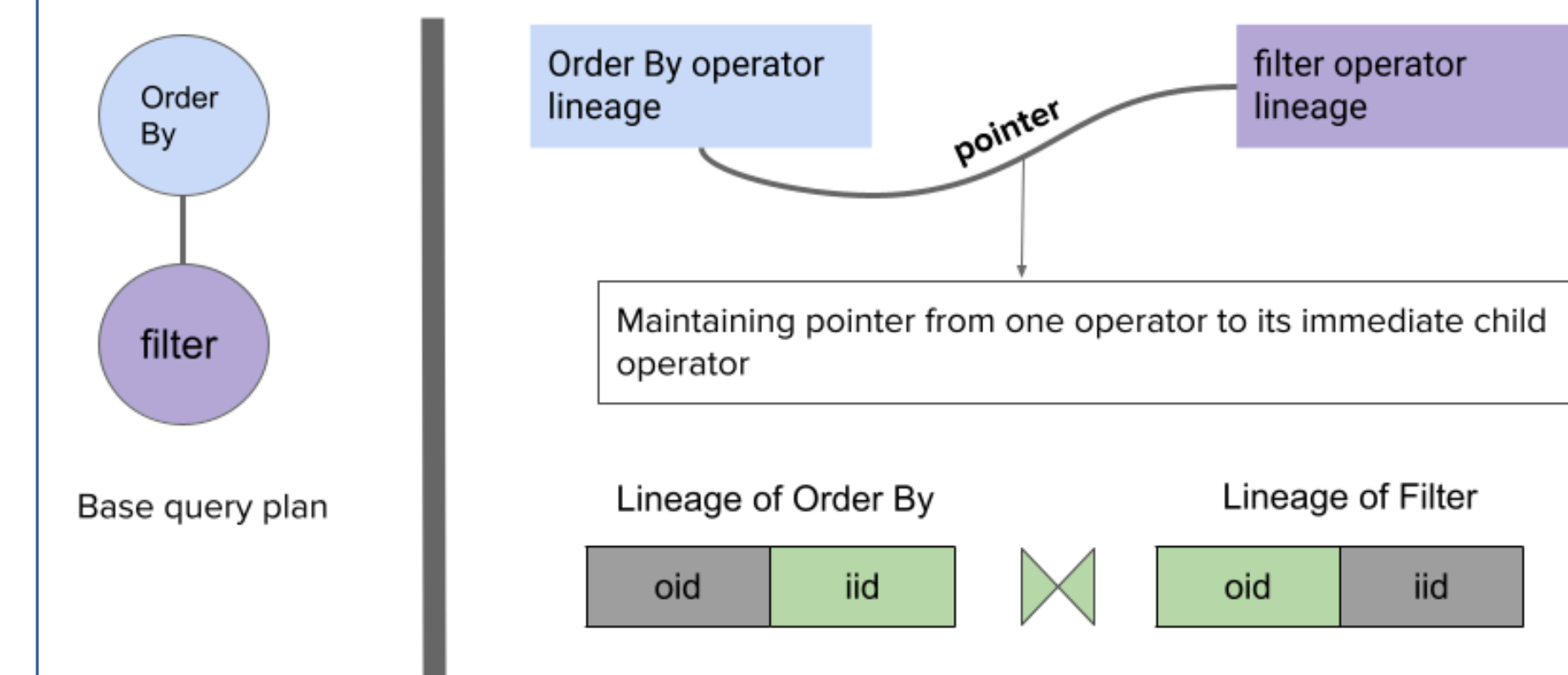
DuckDB is a vectorized execution engine that supports late materialization. It generates data structures during operator execution that encode lineage. Smoked Duck captures lineage by pinning these data structures into memory.

SQL Execution

Contribution:

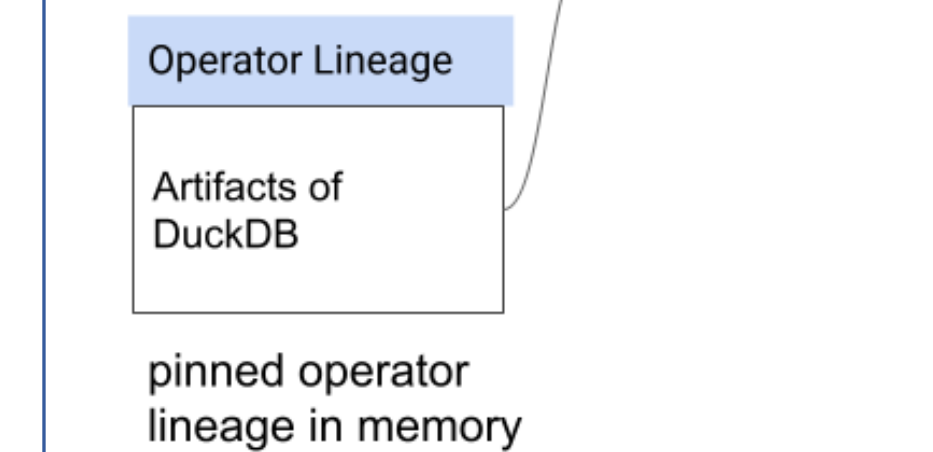
The two expensive costs mentioned are resolved using changes proposed by Smoked Duck. The changes are integrated into DuckDB's execution engine as mentioned below. Thus by leveraging DuckDB's vectorized execution engine we compute fine grained provenance at interactive speeds.

lineage_query(base query, oid)

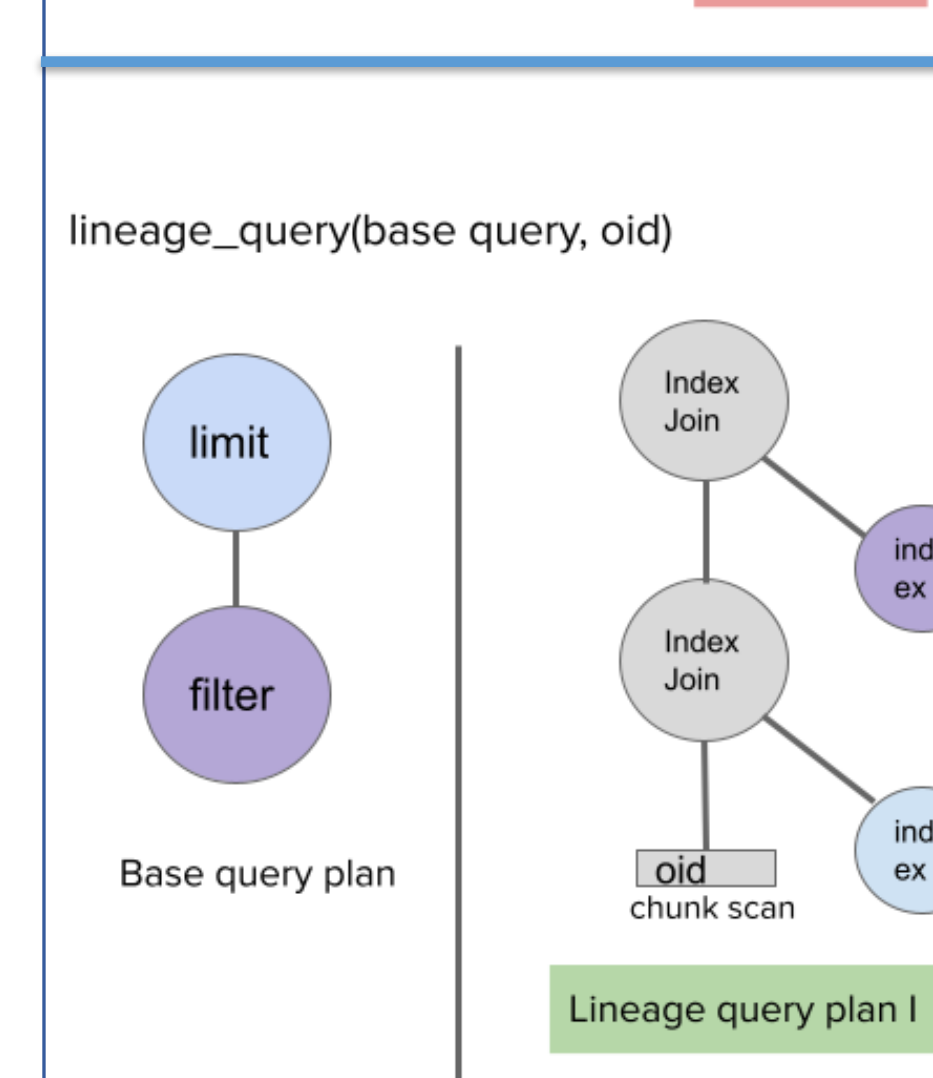


To avoid costs of relational table representation join the data structures pinned into memory using a pointer. Embed the pointer into operator state of DuckDB.

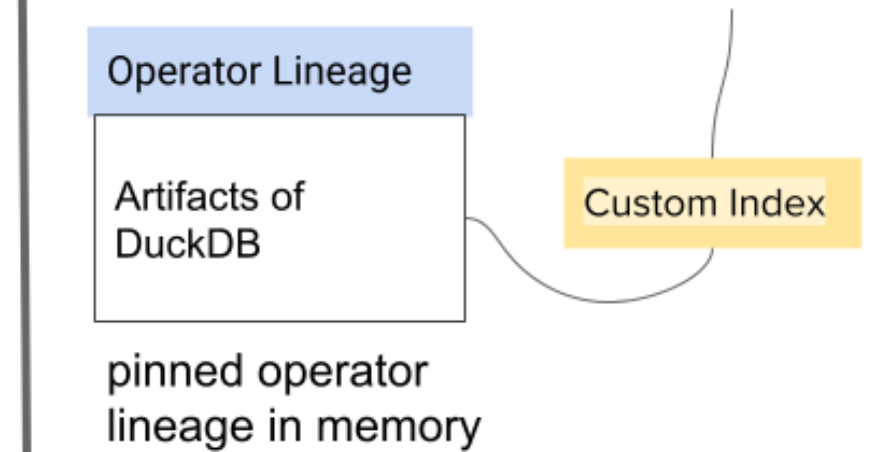
lineage_query(oid)



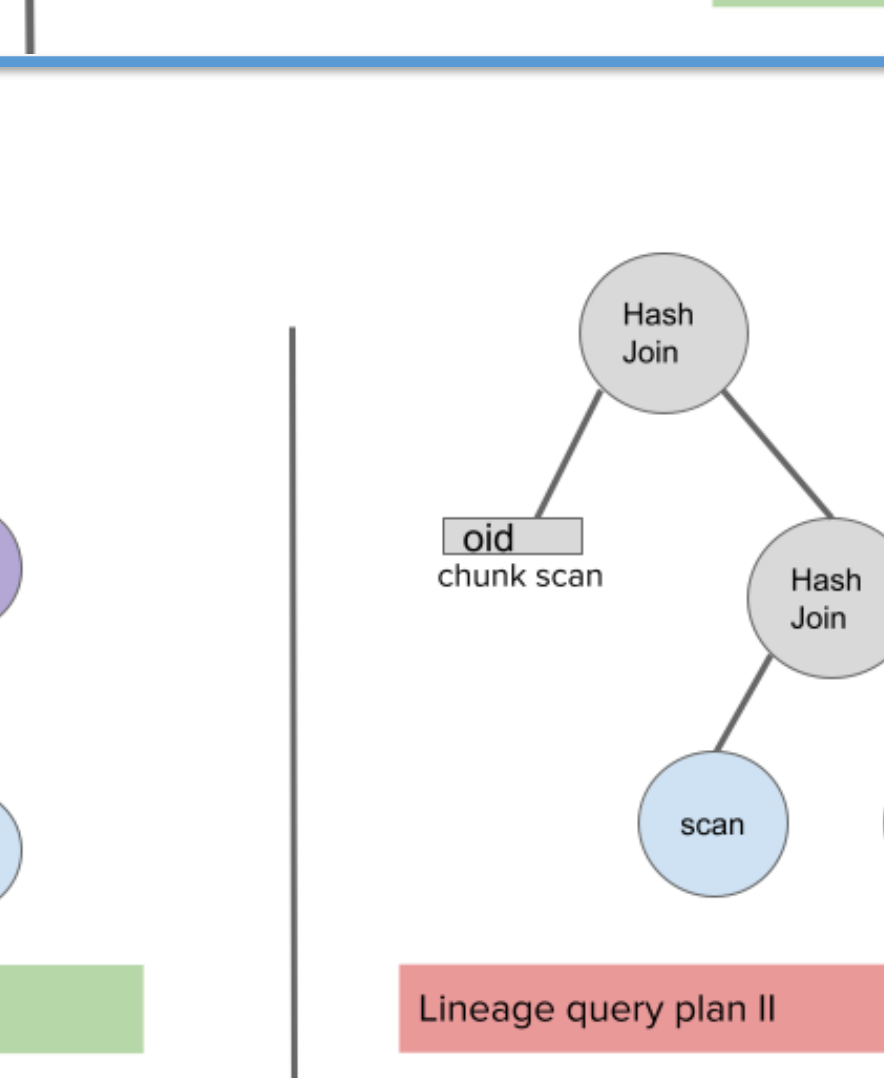
lineage_query(base query, oid)



lineage_query(oid)



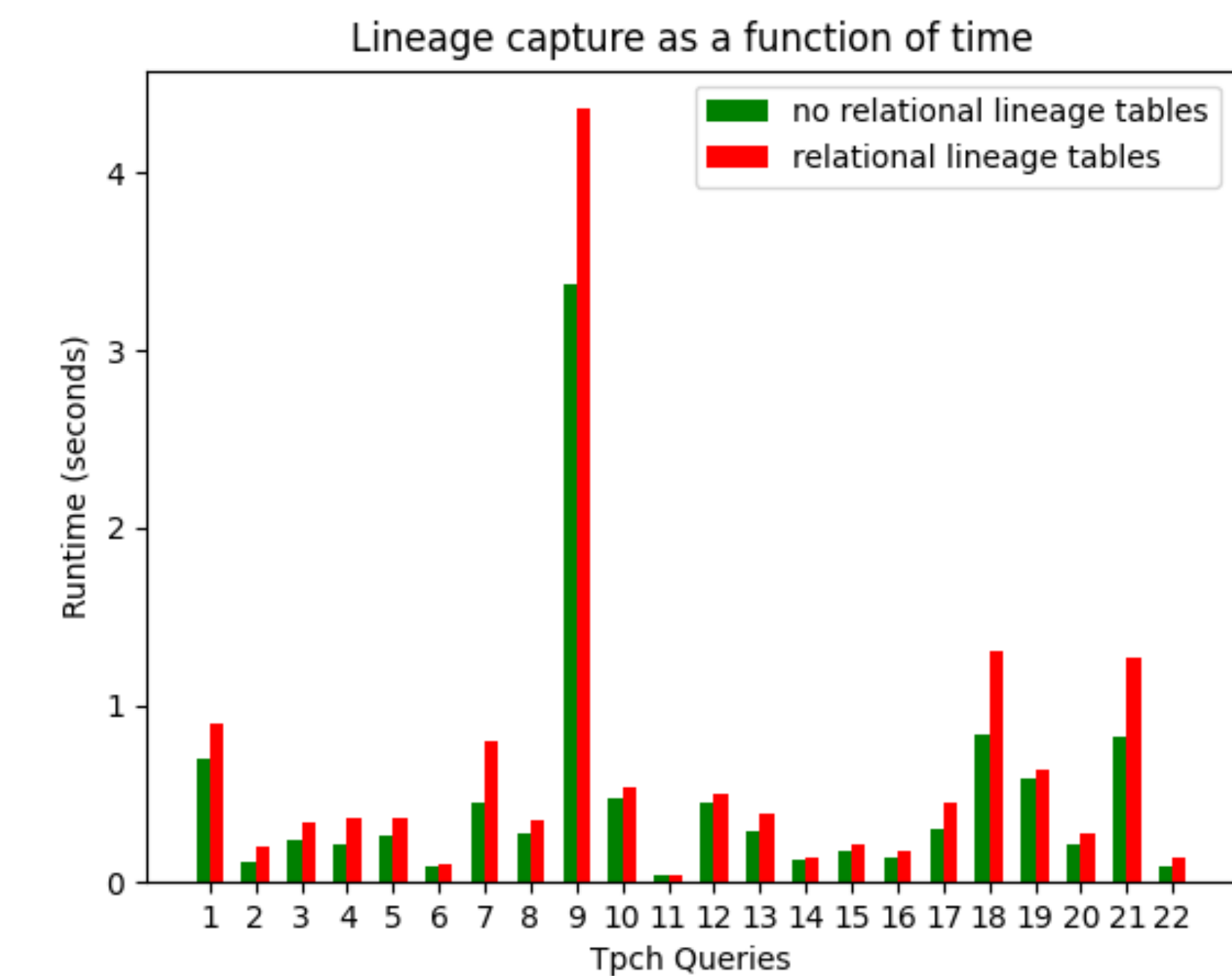
lineage_query(base query, oid)



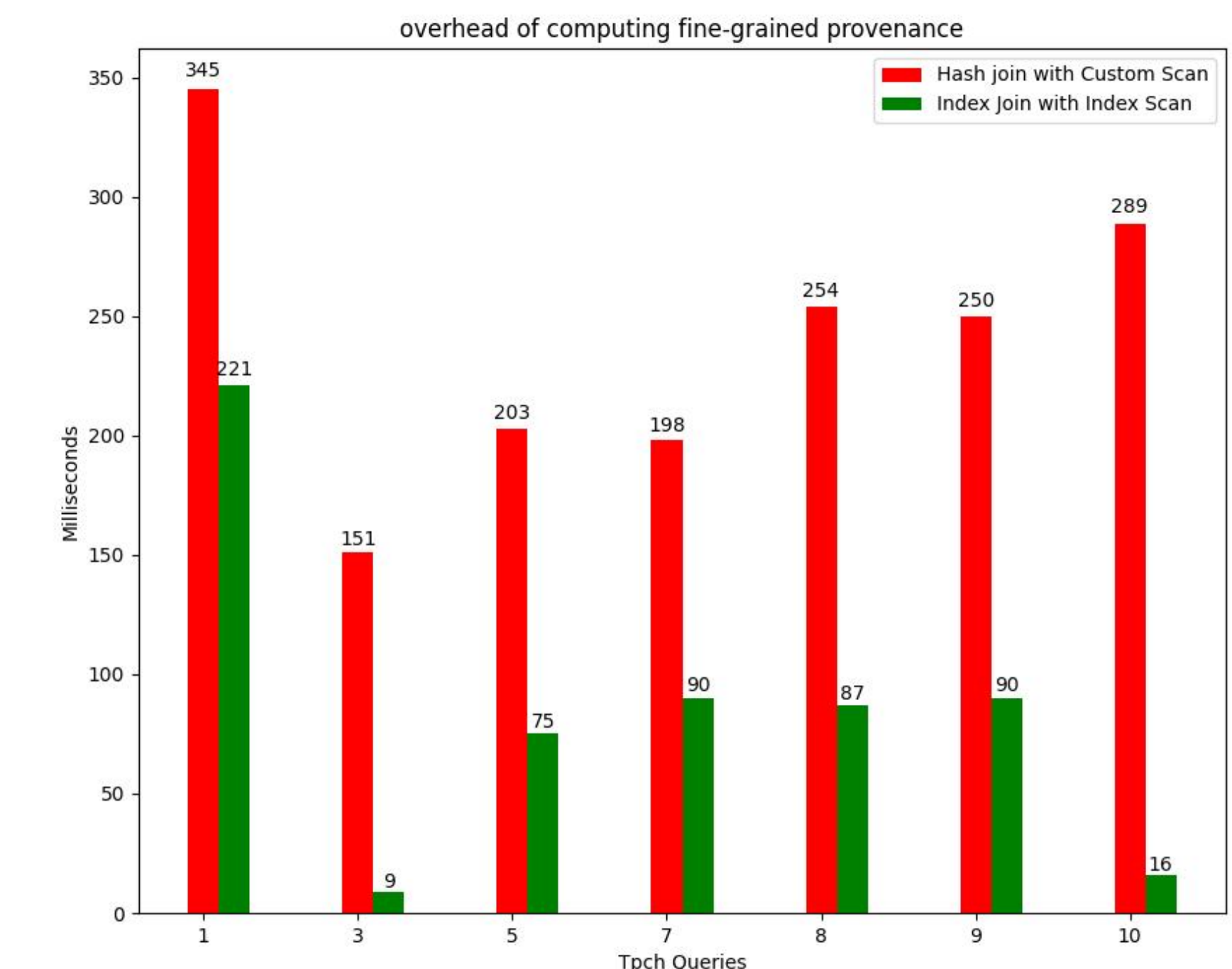
To avoid full scans use custom index and add these indexes to the generic index interface of DuckDB.

Bypass DuckDB's internal planner and build the optimized plan with the pushed down filter.

Results



Overhead of capturing fine-grained provenance



Overhead of querying fine-grained provenance

DuckDB users can compute fine-grained provenance

Create input relation and insert values

```

In [18]: import duckdb
conn = duckdb.connect()
conn.execute("CREATE TABLE personal_info(name VARCHAR, age INTEGER);")
conn.execute("INSERT INTO personal_info VALUES ('ALICE', 25), ('JACK', 31), ('BOB', 26);")
conn.execute("SELECT * FROM personal_info")
print(conn.fetchall())

Input relation
[('ALICE', 25), ('JACK', 31), ('BOB', 26)]

```

Tuples sorted in descending order of age < 30 : Output Relation

```

In [19]: conn.execute("PRAGMA trace_lineage='ON'")
conn.execute("SELECT * FROM personal_info WHERE age < 30 ORDER BY age DESC")
print(conn.fetchall())
print(conn.fetchall())

PragmatraceLineage 1
PostProcess time: 1.4e-05 sec
Output relation
[('BOB', 26), ('ALICE', 25)]
PostProcess time: 8e-06 sec

```

Compute lineage of the above query and the 0th tuple in the output relation

```

In [20]: conn.execute("PRAGMA trace_lineage='OFF'")
conn.execute("PRAGMA lineage_query('SELECT * FROM personal_info WHERE age < 30 ORDER BY age DESC', 'ALICE', 0)")
print("Lineage of 0th tuple in the output is ")
print(conn.fetchall())

Lineage of 0th tuple in the output is PragmatraceLineage 0
Time Taken : 0
[()]

```

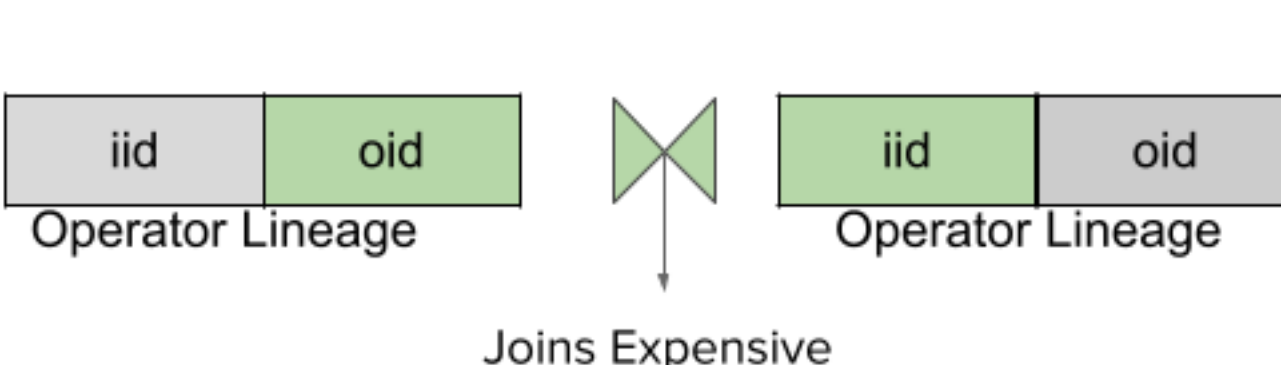
DuckDB users can now query fine-grained provenance

Acknowledgements. This work builds on Smoked Duck, a project in the WuLab. I would like to thank Charlie and Haneen for providing immense guidance throughout this project. Also would like to thank Prof. Eugene Wu and the members of WuLab for constant support for the completion and presentation of the work.

Expensive costs:

- 1) Formatting data structure pinned into memory to a relational table.
- 2) Joining lineage data structures to compute fine grained provenance which entails full scans over lineage data

lineage_query(base query, oid)



Operator Lineage
pinned operator lineage in memory

Expensive

oid	iid

Joins Expensive