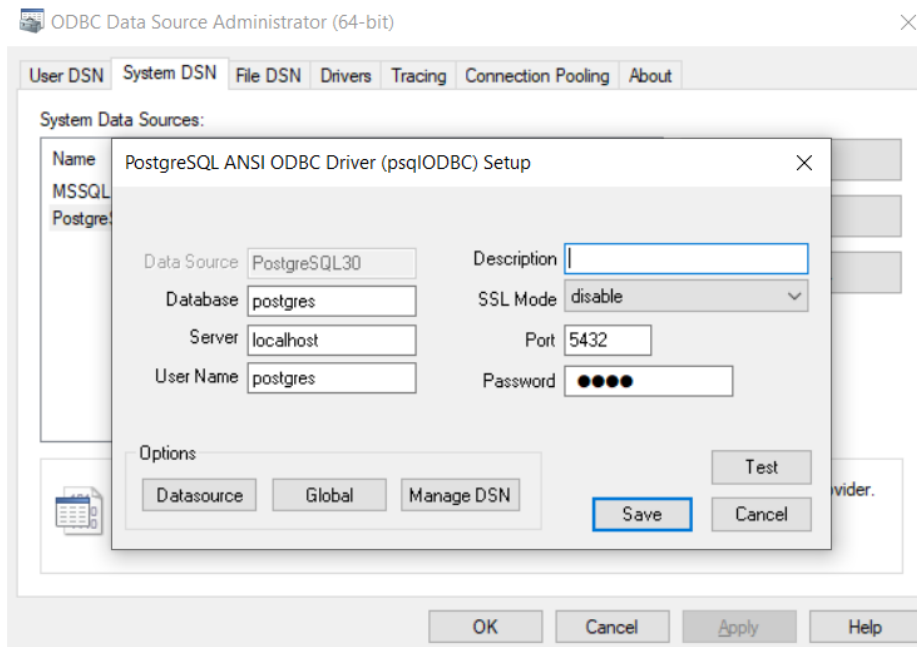SETTING UP THE LINKED SERVER:

I am using a heterogeneous distributed database system of one instance of MSSQL and one instance of PostgreSQL.

We begin by adding the drivers to the System DSN in the ODBC gateway.

PostgreSQL-



MSSQL-

Now in MSSQL under 'Linked Servers', we create POSTGRESQL. We specify the Data Source as PostgreSQL30 (from ODBC). The security settings is set to the fourth option to login with password.

We create a table called LIKES in PostgreSQL and insert 4 rows into the table.

*CREATE TABLE LIKES ( postId INT NOT NULL, likedUserId INT NOT NULL, PRIMARY KEY(postId, likedUserId);*

QUERYING DISTRIBUTED DATABASES:

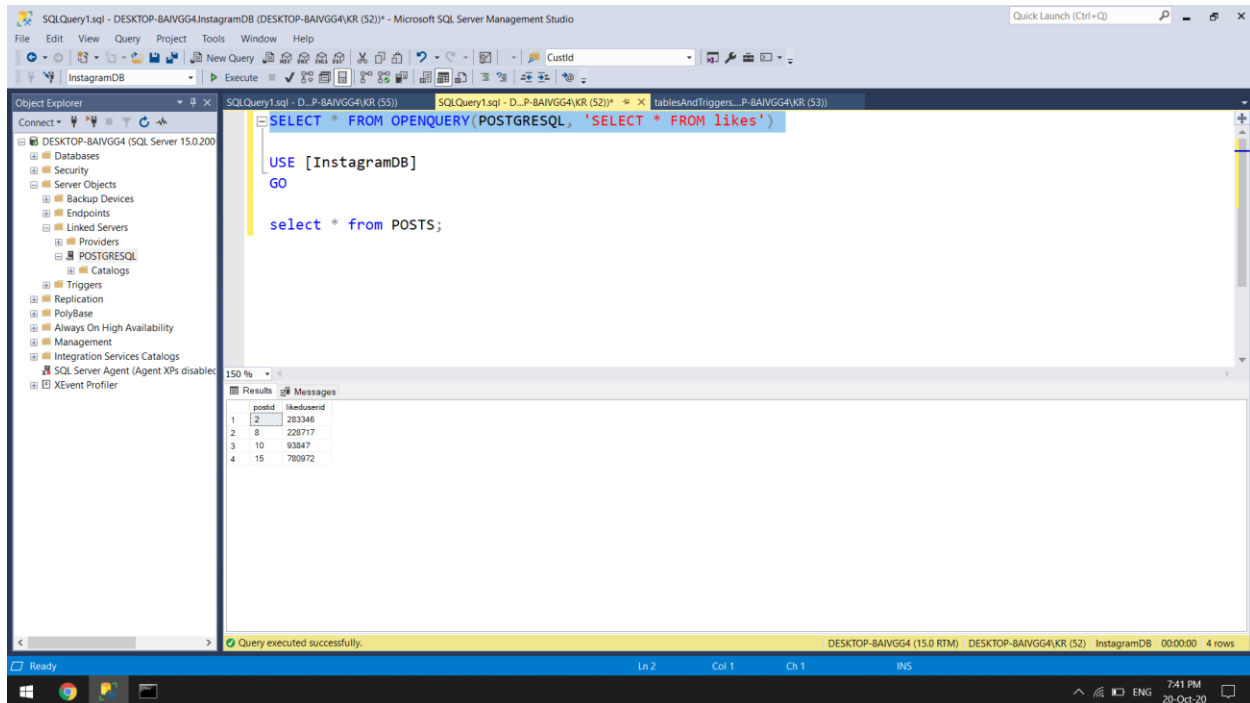We can perform a query on the table LIKES in PostgreSQL from MSSQL using the POSTGRESQL method.



We try to execute a join operation on attributes from table POSTS which is on MSSQL and LIKES which is on PostgreSQL. See the output below.

select P.postUserId,
        P.caption,
        P.likes,
        L.likedUserId
from POSTS P,
    (select * from openquery(POSTGRESQL, 'SELECT * FROM likes')) L
where P.postId = L.postId;

See the execution plan of the query- a clustered index scan is being performed on POSTS whereas a Remote Scan is being performed on LIKES, following which the projected attributes of both are joined.



In the remote scan, we can notice two things-

- The estimated I/O cost is 0 because there is no way for MSSQL to be able to estimate the cost of scan on a table in postgresql.
- Both attributes of LIKES table are seen in the output list, because we use postId in the where condition and likedUserId in the selection list.



Now we try another join operation with the only difference that the selList and where condition, both, include postId from LIKES. So we do not need the likedUserId attribute from LIKES anymore. Now in the execution plan, in the remote scan on LIKES, the output list contains only postId.

```
select P.postUserId,
       P.caption,
       P.likes,
       L.postId
from POSTS P,
    (select * from openquery(POSTGRESQL, 'SELECT * FROM likes')) L
where P.postId = L.postId;
```

However, the query compiler is not reliable in such a scenario for being able to choose the most optimum execution plan. So the developer/user who is querying the distributed database must be aware of the schema and locations of the attributes involved in a distributed database query and try to write the most optimal form of the query.

Screenshot 1 (7:52 PM):

SQL query visible:
```sql
where P.postId = L.postId;

select P.postUserId,
       P.caption,
       P.likes,
       L.postId
from POSTS P,
     (select * from openquery(POSTGRESQL, 'SELECT * FROM likes')) L
where P.postId = L.postId;
```

Execution plan tooltip — Remote Scan:
Scan rows in a table stored in a database or file other than the current database server.

| Physical Operation | Remote Scan |
| Logical Operation | Remote Scan |
| Actual Execution Mode | Row |
| Estimated Execution Mode | Row |
| Actual Number of Rows for All Executions | 4 |
| Actual Number of Batches | 0 |
| Estimated Operator Cost | 3.36333 (87%) |
| Estimated I/O Cost | 0 |
| Estimated CPU Cost | 3.36333 |
| Estimated Subtree Cost | 3.36333 |
| Number of Executions | 1 |
| Estimated Number of Executions | 1 |
| Estimated Number of Rows Per Execution | 10000 |
| Estimated Row Size | 11 B |
| Actual Rebinds | 0 |
| Actual Rewinds | 0 |
| Node ID | 1 |

Output List
[MSDASQL].postid
Remote Source
PostgreSQL30
Remote Object
SELECT * FROM likes

Query 1: Query cost (relative to the bat
select P.postUserId, P.caption, P.likes,
Hash Match (Inner Join) Cost: 5 % 0.012s 4 of 10000 (0%)
SELECT Cost: 0 %
Clustered Inde [POSTS].[PK_P

Query executed successfully.

DESKTOP-8AIVGG4 (15.0 RTM)  DESKTOP-8AIVGG4\KR (52)  InstagramDB  00:00:00  4 rows

---



Screenshot 2 (7:53 PM):

SQL query visible:
```sql
where P.postId = L.postId;

select P.postUserId,
       P.caption,
       P.likes,
       L.postId
from POSTS P,
     (select * from openquery(POSTGRESQL, 'SELECT * FROM likes')) L
where P.po
```

Execution plan tooltip — Hash Match:
Use each row from the top input to build a hash table, and each row from the bottom input to probe into the hash table, outputting all matching rows.

| Physical Operation | Hash Match |
| Logical Operation | Inner Join |
| Actual Execution Mode | Row |
| Estimated Execution Mode | Row |
| Actual Number of Rows for All Executions | 4 |
| Actual Number of Batches | 0 |
| Estimated Operator Cost | 0.194875 (5%) |
| Estimated I/O Cost | 0 |
| Estimated CPU Cost | 0.194866 |
| Estimated Subtree Cost | 3.85138 |
| Number of Executions | 1 |
| Estimated Number of Executions | 1 |
| Estimated Number of Rows Per Execution | 10000 |
| Estimated Row Size | 73 B |
| Actual Rebinds | 0 |
| Actual Rewinds | 0 |
| Node ID | 0 |

Output List
[InstagramDB].[dbo].[POSTS].postUserId, [InstagramDB].[dbo].
[POSTS].caption, [InstagramDB].[dbo].[POSTS].likes,
[MSDASQL].postid
Hash Keys Probe
[InstagramDB].[dbo].[POSTS].postid

Query 1: Query cost
select P.postUserId
Hash Ma (Inner J Cost: 0.012 4 10000
SELECT Cost: 0 %

Query executed successfully

DESKTOP-8AIVGG4 (15.0 RTM)  DESKTOP-8AIVGG4\KR (52)  InstagramDB  00:00:00  4 rows