

# DuckDB

DuckDB runs SQL right on your DataFrames and files with no server, here from a first query to joins, exports, and EXPLAIN.

## Connect & Query







One import, one connection, SQL straight away.

PURPOSE	COMMAND	RESULT
Run SQL with zero setup.	<code>duckdb.sql("SELECT 42 AS answer")</code>	 SQL → DuckDB → answer:42 in-memory by default
Open an in-memory connection.	<code>con = duckdb.connect()</code>	 con → DuckDB RAM dies on exit
Open or create a file database.	<code>con = duckdb.connect("my.duckdb")</code>	 DuckDB → my.duckdb persists
Run SQL on that connection.	<code>con.sql("SELECT * FROM range(3)")</code>	 DuckDB → Relation 0 1 2
Execute a statement (DDL, params).	<code>con.execute("INSERT INTO t VALUES (?)", [1])</code>	 ? → DuckDB no Relation back
Close when done.	<code>con.close()</code>	 con → DuckDB ✓ flushed

import duckdb, then duckdb.sql(...) runs against a default in-memory db; connect() for your own, connect("f.duckdb") to persist

## Query Files Directly







Point SQL at Parquet, CSV, or a glob of files with no load step.

PURPOSE	COMMAND	RESULT
Read one Parquet file in SQL.	<code>duckdb.sql("SELECT * FROM 'sales.parquet'")</code>	 parquet → DuckDB → no load
Read a folder of files (glob).	<code>FROM 'data/*.parquet'</code>	 * → Relation auto-union
CSV with auto type detection.	<code>FROM read_csv_auto('data.csv')</code>	 csv → INT, VARCHAR, DATE types sniffed
Read CSV via the Python helper.	<code>duckdb.read_csv("data.csv")</code>	 csv → Relation Python wrapper
Read CSV from a URL.	<code>FROM 'https://.../d.csv'</code>	 globe → DuckDB https over wire
Peek at the first rows.	<code>FROM 'sales.parquet' LIMIT 5</code>	 LIMIT 5 first 5 rows

a path inside the SQL reads the file directly; a glob auto-unions a folder, read\_csv\_auto sniffs types, https:// streams remotely

## SQL on a DataFrame


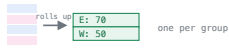



Query a pandas or Polars DataFrame by its Python variable name, no copy.

PURPOSE	COMMAND	RESULT
Query a pandas frame by name.	<code>duckdb.sql("SELECT * FROM sales_df")</code>	
Query a Polars frame by name.	<code>duckdb.sql("SELECT region FROM pl_df")</code>	
Aggregate the frame in SQL.	<code>SELECT region, sum(amt) FROM sales_df GROUP BY region</code>	
Register a frame under a name.	<code>con.register("orders", sales_df)</code>	
Send the result back to pandas.	<code>duckdb.sql("SELECT * FROM sales_df").df()</code>	
Mix a frame and a file.	<code>SELECT * FROM sales_df JOIN 'ref.parquet' USING (id)</code>	

replacement scans put a DataFrame's Python variable name straight in the FROM clause; .df() / .pl() hand the answer back

## Filter, Aggregate, Window







WHERE filters rows, GROUP BY rolls up, OVER computes without collapsing.

PURPOSE	COMMAND	RESULT
Filter rows with WHERE.	<code>SELECT * FROM t WHERE amt &gt; 100</code>	
Group and aggregate.	<code>region, sum(amt) AS total ... GROUP BY region</code>	
Order and limit (top N).	<code>ORDER BY amt DESC LIMIT 10</code>	
Running total with a window.	<code>sum(amt) OVER (ORDER BY ts) AS running</code>	
Partitioned window (per group).	<code>rank() OVER (PARTITION BY region ORDER BY amt DESC)</code>	
Condense the whole pipeline.	<code>WHERE amt &gt; 0 GROUP BY region HAVING avg &gt; 50</code>	

GROUP BY collapses many rows into one per group; a window OVER (PARTITION BY ... ORDER BY ...) computes without collapsing

## Joins







INNER keeps matches, LEFT keeps all left rows, ASOF matches the nearest earlier time.

PURPOSE	COMMAND	RESULT
Inner join on a key.	<code>FROM a JOIN b ON a.id = b.id</code>	
Left join (keep all left rows).	<code>FROM a LEFT JOIN b ON a.id = b.id</code>	
Join on a shared column name.	<code>FROM a JOIN b USING (id)</code>	
ASOF join (nearest earlier row).	<code>ASOF JOIN prices p ON t.sym=p.sym AND t.ts&gt;=p.ts</code>	
Stack rows from two sources.	<code>SELECT * FROM a UNION ALL SELECT * FROM b</code>	
Anti join (rows with no match).	<code>WHERE id NOT IN (SELECT id FROM b)</code>	

ASOF JOIN matches each left row to the nearest earlier right row ( $t.ts \geq p.ts$ ), the standard tool for prices-per-trade

## Move Results Out







A Relation is lazy. One fetch method turns it into pandas, Polars, Arrow, or Python.

PURPOSE	COMMAND	RESULT
To a pandas DataFrame.	<code>duckdb.sql("...").df()</code>	
To a Polars DataFrame.	<code>duckdb.sql("...").pl()</code>	
To an Arrow table.	<code>duckdb.sql("...").arrow()</code>	
To plain Python objects.	<code>duckdb.sql("...").fetchall()</code>	
To NumPy column arrays.	<code>duckdb.sql("...").fetchnumpy()</code>	
Lazy until you fetch.	<code>rel = duckdb.sql("SELECT ...")</code>	

nothing runs until you fetch or print; `.df()` / `.pl()` / `.arrow()` / `.fetchall()` / `.fetchnumpy()` pick the output format you want

## Persist & Export



Materialize a query into a stored table, or write it straight to a file with COPY.

PURPOSE	COMMAND	RESULT
Create a table from a query.	<code>CREATE TABLE summary AS SELECT ... GROUP BY region</code>	
Replace an existing table.	<code>CREATE OR REPLACE TABLE summary AS ...</code>	
Append more rows.	<code>INSERT INTO summary SELECT * FROM more</code>	
Write a query to Parquet.	<code>COPY(SELECT * FROM t) TO 'out.parquet'</code>	
Write to CSV with a header.	<code>COPY t TO 'out.csv' (HEADER, DELIMITER ',')</code>	
Write via the Relation helper.	<code>duckdb.sql("..."). write_parquet("out.parquet")</code>	

CREATE [OR REPLACE] TABLE ... AS SELECT freezes a query into a stored table; COPY (q) TO 'file' writes straight to disk

## Inspect & EXPLAIN

See the schema, list what exists, and read how a query will run.

PURPOSE	COMMAND	RESULT								
Describe a table or query.	<code>DESCRIBE summary</code>	<table border="1"> <thead> <tr> <th>column_name</th> <th>column_type</th> </tr> </thead> <tbody> <tr> <td>region</td> <td>VARCHAR</td> </tr> <tr> <td>total</td> <td>HUGEINT</td> </tr> </tbody> </table>	column_name	column_type	region	VARCHAR	total	HUGEINT		
column_name	column_type									
region	VARCHAR									
total	HUGEINT									
List tables in the database.	<code>SHOW TABLES</code>	<ul style="list-style-type: none"> <li>prices · summary</li> <li>sales · trades</li> </ul>								
Show all columns across tables.	<code>DESCRIBE or duckdb_columns()</code>	<table border="1"> <thead> <tr> <th>table.column</th> <th>type</th> </tr> </thead> <tbody> <tr> <td>sales.amt</td> <td>INTEGER</td> </tr> <tr> <td>summary.total</td> <td>HUGEINT</td> </tr> </tbody> </table> catalog	table.column	type	sales.amt	INTEGER	summary.total	HUGEINT		
table.column	type									
sales.amt	INTEGER									
summary.total	HUGEINT									
Summarize a dataset's stats.	<code>SUMMARIZE t</code>	<table border="1"> <thead> <tr> <th>min</th> <th>max</th> <th>approx_unique</th> <th>null%</th> </tr> </thead> <tbody> <tr> <td>██████</td> <td>██████</td> <td>██████</td> <td>██████</td> </tr> </tbody> </table> quick profile per column	min	max	approx_unique	null%	██████	██████	██████	██████
min	max	approx_unique	null%							
██████	██████	██████	██████							
Read the query plan.	<code>EXPLAIN SELECT region, sum(amt) ... GROUP BY</code>									
Read the plan with real timings.	<code>EXPLAIN ANALYZE SELECT ...</code>									

DESCRIBE / SHOW TABLES / SUMMARIZE inspect schema and stats; EXPLAIN prints the operator tree, EXPLAIN ANALYZE adds timings