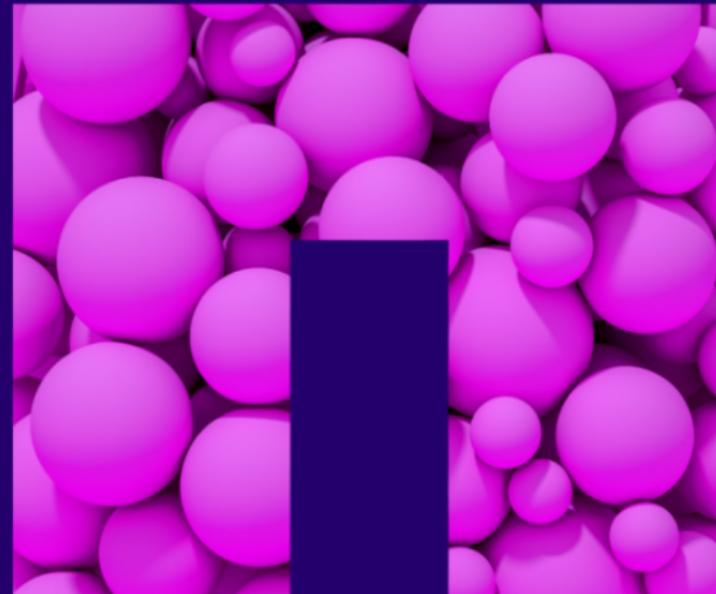


PostgreSQL and DuckDB

Supercharging
ad-hoc Data Analysis
and ETL

Josef Machytka <josef.machytka@netapp.com>
NetApp Open Source Services
2025-01-09 - PostgreSQL MeetUp for All



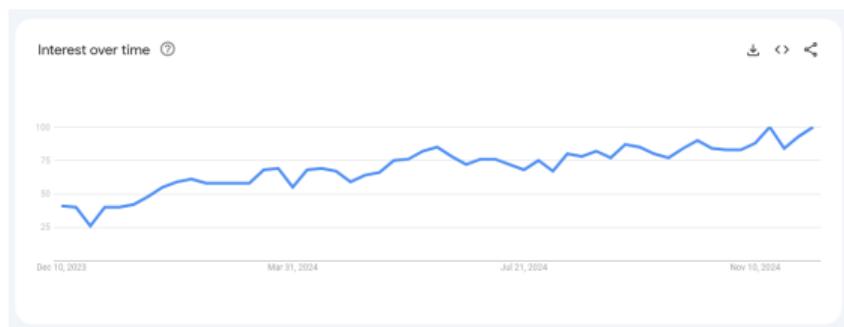
Josef Machytka

- Professional Service Consultant - PostgreSQL specialist at NetApp Open Source Services / Credativ
- 30+ years of experience with different databases.
- PostgreSQL (12y), BigQuery (7y), Oracle (15y), MySQL (12y), Elasticsearch (5y), MS SQL (5y).
- 10+ years of experience with Data Ingestion pipelines, Data Analysis, Data Lake and Data Warehouse
- 2 years of practical experience with different LLMs / AI including their architecture and principles.
- From Czechia, living now 11 years in Berlin.

- LinkedIn: linkedin.com/in/josef-machytka
- ResearchGate.com: researchgate.net/profile/Josef-Machytka
- Academia.edu: netapp.academia.edu/JosefMachytka
- Medium.com: medium.com/@josef.machytka

Hype around DuckDB

- Significant hype around DuckDB in recent months



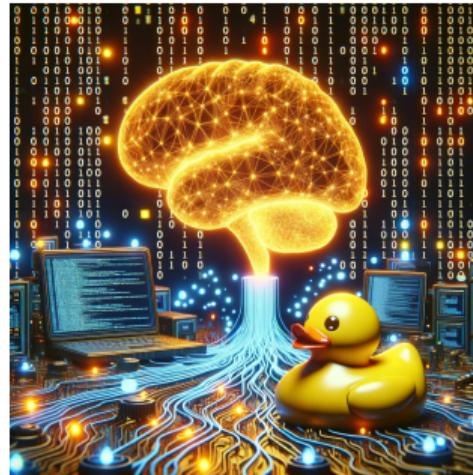
(Search for "duckdb" by Google trends)

How DuckDB handles data not fitting into memory?	600	354
<small>4 min read · Nov 13, 2024 · View story</small>	<small>Views</small>	<small>Reads</small>
Extending DuckDB ETL Capabilities with Python	532	315
<small>3 min read · Nov 24, 2024 · View story</small>	<small>Views</small>	<small>Reads</small>
Easy Cross-Database Selects with DuckDB	358	175
<small>2 min read · Nov 26, 2024 · View story</small>	<small>Views</small>	<small>Reads</small>
DuckDB Database File as a New Standard for Sharing ...	326	235
<small>4 min read · Dec 30, 2024 · View story</small>	<small>Views</small>	<small>Reads</small>
DuckDB as a Rudimentary Data Migration Tool	241	153
<small>3 min read · Nov 30, 2024 · View story</small>	<small>Views</small>	<small>Reads</small>
Using DuckDB as an Intelligent ETL tool for PostgreSQL	203	128
<small>4 min read · Nov 2, 2024 · View story</small>	<small>Views</small>	<small>Reads</small>

(Stats of my DuckDB related articles on medium.com)

Table of contents

- Big Data vs Small Data
- Data Formats
- Meet DuckDB
- Connect DuckDB and PostgreSQL
- Data Analysis and ETL
- Rudimentary Migration Tool

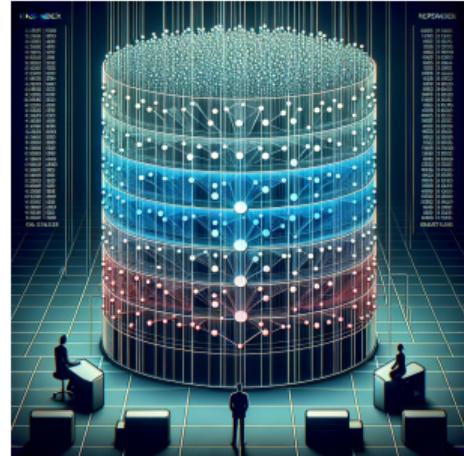


All AI images used in slides
created by the author
using DeepDreamGenerator

Big Data vs Small Data

Big Data

- Data sets too large or complex for traditional processing at that time
- Characterized by 5 V's:
- Volume: Large data size, low density of useful information
- Velocity: High speed of data in and out, crucial for real-time data
- Variety: Different data types - structured, unstructured, semi-structured
- Veracity: Data uncertainty, including noise, outliers, and errors
- Value: Potential to generate significant value from data



Big Data is (Mostly) Dead

- The boom of Big Data appears to be mostly over
- Companies often collected large data sets without clear usage plans
- Only a small percentage of companies truly need to store Big Data
- Some raw data can be deleted after a short period of time
- Companies mostly need aggregated results
- Useful data sets typically have only a few dozen GBs
- Hardware is advancing rapidly, instances can have TBs of memory
- Refer to the article [Big Data is Dead](#)



Small Data Manifesto

- DuckDB Team wrote [The Small Data Manifesto](#)
- Small Data can be more useful and accessible than Big Data
- Easier to understand, process, and share
- Sufficient for many use cases and decision-making processes
- Complex systems for Big Data lead to resource competition
- Modern notebooks are powerful but often underutilized



Data Formats

Relational Data are not Predominant Anymore

- 15-20 years ago, data was mainly stored in relational databases
- Data analysis was performed on classical database tables using SQL
- The digital world was "slower," allowing time to transform data
- Data was stored in relational tables for structured querying

- Today, real-time processing is the new norm
- Modern technologies utilize diverse data formats
- Data analysis must handle these diverse formats directly



Most Common Data File Formats

- CSV/TSV - Comma (Tab) Separated Values - row oriented, each line is a record
- JSON - JavaScript Object Notation - key-value pairs, supports nested structures
- Parquet - columnar storage format, optimized for efficient reading
- Avro - row oriented, schema-based, binary format, supports rich data structures
- ORC - Optimized Row Columnar - columnar storage format, optimized for read-heavy workloads



Meet DuckDB

Meet DuckDB

- Created at the National Research Institute for Mathematics and Computer Science, Amsterdam
- Non-profit DuckDB Foundation ensures long-term maintenance and development
- Open-source, column-oriented, in-memory relational database
- Single-node database, intended for embedding in applications, like SQLite
- Designed for heavy parallel analytical workloads
- Columnar-vectorized query processing engine

- API support for major languages: Go, Python, R, Java, C++, Rust
- Extremely portable, runs on all architectures, no dependences



DuckDB follows PostgreSQL conventions

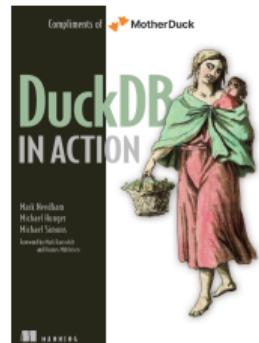
- SQL dialect closely follows PostgreSQL conventions
- Uses SQL parser from PostgreSQL
- New functionality can be added via custom extensions
- Uses Multi-Version Concurrency Control (MVCC) for transactions



DuckDB at a glance



[DuckDB Main Webpage - duckdb.org](https://duckdb.org)



[MotherDuck Webpage - motherduck.com](https://motherduck.com)

Single-node database

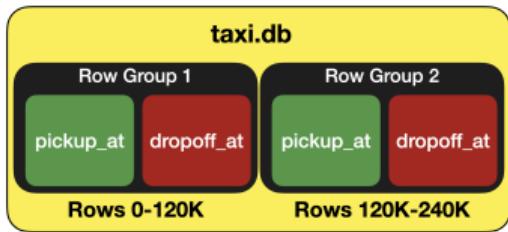
- Single-node database
- Not designed for distributed environments
- Lacks commands CREATE USER / ROLE, GRANT, REVOKE
- Designed mainly for datasets fitting into memory
- Successfully tested with datasets larger than memory



(Picture from Udemy course [DuckDB - The Ultimate Guide](#))

DuckDB Storage

- Single-file database format with columnar storage and lightweight compression
- Layout similar to Apache Parquet but columns split into fixed-size blocks
- Default block size 256 KB - default_block_size parameter (min 16 KB)
- Allows in-place ACID updates, deletes, and adding/dropping columns



(Image from the article [Lightweight Compression in DuckDB](#))

DuckDB Data Structures

- Supports all traditional database data types
- Includes multiple synonyms from older databases
- Python structures: Enums, Lists/Arrays, Maps, and Structs
- Built-in functions for operations with these objects
- Works directly with Pandas, NumPy, and Polars data frames
- Memory variables: SET VARIABLE, SELECT getvariable('name')
- Allows user-defined data types
- Implements UNION data type ("one of these types")

General-Purpose Data Types

The table below shows all the built-in general-purpose data types. The alternatives listed in the aliases column can be used to refer to these types as well, however, note that the aliases are not part of the SQL standard and hence might not be accepted by other database engines.

Name	Aliases	Description
BIGINT	INT8, LONG	signed eight-byte integer
BIT	BITSTRING	string of 1s and 0s
BLOB	BYTEA, BINARY, VARBINARY	variable-length binary data
BOOLEAN	BOOL, LOGICAL	logical boolean (true/false)
DATE		calendar date (year, month, day)
DECIMAL(prec, scale)	NUMERIC(prec, scale)	fixed-precision number with the given width (precision) and scale, defaults to prec = 18 and scale = 3
DOUBLE	FLOAT8,	double precision floating-point number (8 bytes)
FLOAT	FLOAT4, REAL	single precision floating-point number (4 bytes)

(Documentation [DuckDB Data Types](#))

DuckDB CLI

- Starting "duckdb" command or container, we run both engine and command line client
- CLI has "[dot commands](#)" - starting with "." - similar to psql "\" commands
- .help - List all available dot commands
- .help CMD - Show help for specific command CMD
- .excel - Display the output of next command in spreadsheet
- .mode MODE - Set the output mode to MODE (box, line, table, csv, json, markdown)
- .system CMD ARGS... - Run CMD ARGS... in a system shell
- .tables TABLE - List names of tables matching LIKE pattern TABLE
- .timer on|off - Turn SQL timer on or off
- .exit or .quit - Exit the DuckDB shell

DuckDB Extensions

```
D select * from duckdb_extensions;
```

extension_name varchar	loaded boolean	installed boolean	install_path varchar	description varchar	aliases varchar[]	extension_version varchar	install_mode varchar	installed_from varchar
arrow	false	false		A zero-copy data integration between Apache Arrow and DuckDB	[]			
autocomplete	true	true	(BUILT-IN)	Adds support for autocomplete in the shell	[]		STATICALLY_LINKED	
aws	false	false		Provides features that depend on the AWS SDK	[]			
azure	false	false		Adds a filesystem abstraction for Azure blob storage to DuckDB	[]			
delta	false	false		Adds support for Delta Lake	[]			
excel	false	false		Adds support for Excel-like format strings	[]			
fts	true	true	(BUILT-IN)	Adds support for Full-Text Search Indexes	[]	v1.1.3	STATICALLY_LINKED	
httpfs	false	false		Adds support for reading and writing files over a HTTP(S) connection	[http, https, s3]			
iceberg	false	false		Adds support for Apache Iceberg	[]			
icu	true	true	(BUILT-IN)	Adds support for time zones and collations using the ICU library	[]	v1.1.3	STATICALLY_LINKED	
inet	false	false		Adds support for IP-related data types and functions	[]			
jemalloc	true	true	(BUILT-IN)	Overwrites system allocator with JEMalloc	[]	v1.1.3	STATICALLY_LINKED	
json	true	true	(BUILT-IN)	Adds support for JSON operations	[]	v1.1.3	STATICALLY_LINKED	
motherduck	false	false		Enables motherduck integration with the system	[md]			
mysql_scanner	false	false		Adds support for connecting to a MySQL database	[mysql]			
parquet	true	true	(BUILT-IN)	Adds support for reading and writing parquet files	[]	v1.1.3	STATICALLY_LINKED	
postgres_scanner	false	true	/home/josef/.duckdb/extensions/v1.1.3/linux_amd64-	Adds support for connecting to a Postgres database	[postgres]	03eae7	REPOSITORY	core
shell	true	true		Adds CLI-specific support and functionalities	[]		STATICALLY_LINKED	
spatial	false	false		Geospatial extension that adds support for working with spatial data and functions	[]			
sqlite_scanner	false	false		Adds support for reading and writing SQLite database files	[sqlite, sqlite3]			
substrait	false	false		Adds support for the Substrait integration	[]			
tpcds	false	false		Adds TPC-DS data generation and query support	[]			
tpch	true	true	(BUILT-IN)	Adds TPC-H data generation and query support	[]	v1.1.3	STATICALLY_LINKED	
vss	false	false		Adds indexing support to accelerate Vector Similarity Search	[]			

24 rows

9 columns

Implemented JOINS

- Usual Joins: LEFT JOIN, RIGHT JOIN, FULL JOIN, CROSS JOIN (Cartesian product)
- LATERAL JOIN - join with a subquery that can reference columns from the left table
- NATURAL JOIN - i.e. join done using attributes with the same names
- SEMI JOIN - returns rows from the left table that have at least one matching row in the right table
- ANTI JOIN - returns rows from the left table that have no matching rows in the right table
- POSITIONAL JOIN - always FULL OUTER JOIN, used for ordered data, joins on row number
- ASOF JOIN - used for time series data, joins on the nearest timestamp

Connect DuckDB and PostgreSQL

Connect PostgreSQL to DuckDB

- PostgreSQL uses Foreign Data Wrappers (FDW) for external data sources
- Access DuckDB database files using `duckdb_fdw`
- Version 1.13 supports basic DML operations
- DuckDB as an analytical and ETL engine: extensions `pg_duckdb` and `pg_analytics`
- Run DuckDB queries and ETL tasks inside PostgreSQL
- Suitable for Data Lakehouse architecture - Apache Iceberg & Delta Lake integration
- Extensions are in early versions (0.2.x), not ready for production



Connect DuckDB to PostgreSQL

- Use the ATTACH command to connect to a remote PostgreSQL database
- Login credentials can be included in the ATTACH command or stored as a SECRET
- PostgreSQL tables accessible with local synonym, and remote schema nad table names
- Perform SQL operations on PostgreSQL tables directly from DuckDB
- Database can be attached also as READ_ONLY

```
INSTALL POSTGRES; LOAD POSTGRES;

-- direct attach with credentials
ATTACH 'dbname=postgres user=postgres password=postgres host=localhost port=5436' as pgdb (TYPE POSTGRES);

-- create secret for PostgreSQL connection
CREATE SECRET local_pg (TYPE POSTGRES, HOST 'localhost', PORT 5436, DATABASE postgres, USER 'postgres', PASSWORD 'postgres');

ATTACH '' AS pgdb (TYPE POSTGRES, SECRET local_pg);

-- work with PostgreSQL tables
SELECT * FROM pgdb.schema.table_name;

CREATE TABLE pgdb.schema.new_table AS SELECT * FROM local_duckdb_table;

CREATE TABLE local_duckdb_table AS SELECT * FROM pgdb.schema.table_name;

-- detach PostgreSQL database
DETACH pgdb;
```

Cross Databases Joins

- DuckDB can join tables from different databases of various versions easily.
- It can also join tables with external data files seamlessly.
- PostgreSQL can achieve the same, but it is more complex and requires more steps.

```
D ATTACH 'host=localhost port=5433 user=postgres password=postgres dbname=test' AS pg13 (TYPE POSTGRES, SCHEMA 'public');
D ATTACH 'host=localhost port=5434 user=postgres password=postgres dbname=postgres' AS pg14 (TYPE POSTGRES, SCHEMA 'public');
D ATTACH 'host=localhost port=5435 user=postgres password=postgres dbname=orders' AS pg15 (TYPE POSTGRES, SCHEMA 'public');
D SELECT
    u.username,
    u.email,
    o.order_date,
    o.total_amount,
    p.product_name,
    od.quantity,
    p.price,
    (od.quantity * p.price) AS total_price
FROM
    pg13.users u
JOIN
    pg15.orders o ON u.user_id = o.user_id
JOIN
    pg15.order_details od ON o.order_id = od.order_id
JOIN
    pg14.products p ON od.product_id = p.product_id
ORDER BY
    u.username, o.order_date;
```

username varchar	email varchar	order_date date	total_amount decimal(10,2)	product_name varchar	quantity int32	price decimal(10,2)	total_price decimal(18,2)
Alice	alice@example.com	2024-11-20	150.00	Mouse	2	20.00	40.00
Alice	alice@example.com	2024-11-20	150.00	Keyboard	1	50.00	50.00
Bob	bob@example.com	2024-11-21	250.00	Monitor	1	200.00	200.00
Charlie	charlie@example.com	2024-11-22	300.00	Laptop	1	1000.00	1000.00

Data Analysis and ETL

Testing Data

- Testing tables: `probe_measurements` and `user_data`
- Data generated using Python: 10 GB + 100 GB CSV files for both tables
- `probe_measurements`: PG 15 GB, 143M rows / PG 144 GB, 1.4B rows
- `user_data`: PG 11 GB, 40M rows / PG 111 GB, 400M rows
- Tables fitting into memory are easy to process remotely
- Large tables should be copied to DuckDB for efficient processing

```
CREATE TABLE probe_measurements (
    measurement_id BIGSERIAL PRIMARY KEY,
    probe_id INTEGER NOT NULL,
    measurement_time TIMESTAMP NOT NULL,
    temperature NUMERIC(5,2) NOT NULL,
    pressure NUMERIC(6,2) NOT NULL,
    humidity NUMERIC(5,2) NOT NULL,
    voltage NUMERIC(5,2) NOT NULL,
    current NUMERIC(5,2) NOT NULL,
    resistance NUMERIC(10,2) NOT NULL,
    sample_rate INTEGER NOT NULL);
```

```
CREATE TABLE user_data (
    user_id BIGSERIAL PRIMARY KEY,
    first_name VARCHAR(50) NOT NULL,
    last_name VARCHAR(50) NOT NULL,
    email VARCHAR(100) NOT NULL UNIQUE,
    signup_date TIMESTAMP NOT NULL,
    last_login TIMESTAMP,
    is_active BOOLEAN NOT NULL DEFAULT TRUE,
    account_balance NUMERIC(12,2) NOT NULL DEFAULT 0.00,
    country_code CHAR(2) NOT NULL,
    favorite_number INTEGER,
    profile_text TEXT,
    checksum UUID NOT NULL);
```

Direct Work with External Data Files

- DuckDB implements *table functions* API for external data sources
- API allows querying data files as tables in the FROM clause
- Direct access: `SELECT * FROM 'mydata.parquet';`
- Function-based access: `SELECT * FROM read_parquet('mydata.parq', params...);`
- Automatically analyzes the structure of data files
- Supports CSV, Parquet, JSON, Excel files, and HTTP requests (cloud storage)
- Supports Apache Iceberg and Delta Lake external tables

Direct Import of Data Files into PostgreSQL with DuckDB

- Create PostgreSQL table directly from external data file using DuckDB
- Comparable speed to direct import into PostgreSQL
- No need for manual PostgreSQL table creation
- DuckDB automatically recognizes structure

```
ATTACH 'dbname=postgres user=postgres password=postgres host=localhost port=5436' as pgdb
(TYPE POSTGRES, SCHEMA 'public');

-- table is directly created in PostgreSQL
CREATE TABLE pgdb.probe_measurements AS SELECT * FROM 'probe_measurements_100gb.csv';

-- insert of data into existing PostgreSQL table
INSERT INTO pgdb.probe_measurements SELECT * FROM 'probe_measurements_100gb.csv';
```

Import of legacy DBF file into PostgreSQL

- Combining DuckDB with Pandas/Polars to read DBF file and write to PostgreSQL

```
import pandas as pd
import duckdb
from dbfread import DBF

table = DBF('dbase_sample_data.dbf', load=True)

dataframe = pd.DataFrame(iter(table))

conn = duckdb.connect(database=':memory:')

conn.execute("""
    ATTACH 'dbname=duckdb_test user=postgres password=postgres host=localhost port=5432'
    AS pg (TYPE POSTGRES, SCHEMA 'public')")
""")

conn.execute("CREATE TABLE IF NOT EXISTS pg.dbase_sample_data AS SELECT * FROM dataframe")

conn.close()
```

Conversion of data formats

- DuckDB can convert CSV, JSON, and other formats to Parquet
- Parquet format is de facto standard for Data Lakes/ Lakehouses
- Apache Iceberg, Delta Lake, Apache Hudi are all built on top of Parquet
- Can be used also to export data into other formats

```
-- export to Parquet
COPY 'user_data.csv' TO 'user_data.parquet' (format parquet);

COPY (SELECT col1, col2, sum(col3) AS col3_sum FROM datafram GROUP BY ALL)
      TO 'legacy_dbf_data.parquet' (format parquet);

COPY pgdb.probe_measurements TO 'pgdb_probe_measurements.parquet' (format parquet);

INSTALL avro FROM community; LOAD avro;
COPY (SELECT * FROM read_avro('some_avro_file.avro')) TO 'avro_data.parquet' (format parquet);

-- export to CSV, JSON
COPY pgdb.user_data TO 'user_data.csv' (format csv);

COPY pgdb.user_data TO 'user_data.json' (format json);
```

Data Analysis Commands

- DuckDB has commands: SUMMARIZE, HISTOGRAM, PIVOT
- SUMMARIZE - statistics for all columns of table/data file/data frame
- HISTOGRAM - creates a histogram of a column
- PIVOT - pivots a table based on a column
- Advanced aggregations and window functions



SUMMARIZE command

```
D summarize probe measurements;
```

```
100% [REDACTED]
```

column_name	column_type	min	max	approx_unique	avg	std	q25	q50	q75	count	null_percentage
	varchar	varchar	varchar	int64	varchar	varchar	varchar	varchar	varchar	int64	decimal(9,2)
measurement_id	BIGINT	1	1413016000	1527453031	706508000.5	487902584.12896883	353257459	706307700	1059392748	1413016000	0.00
probe_id	BIGINT	1	1800	1232	500.5851361102776	288.6756694833399	251	500	750	1413016000	0.00
measurement_time	TIMESTAMP	2024-12-18 23:39:37	2024-12-30 16:20:48	1109862	-2.7343037870909387e-05	28.86754646060938	2024-12-21 22:39:00.23257	2024-12-24 20:02:28.896443	2024-12-27 17:27:23.33628	1413016000	0.00
temperature	DOUBLE	-50.0	50.0	9552	-2.7343037870909387e-05	28.86754646060938	-24.996478102538034	-0.0011950089665148972	24.99023243889359	1413016000	0.00
pressure	DOUBLE	900.0	1100.0	18127	999.9984155742529	57.73468372202699	949.9449746301059	1000.00098645416814	1049.977315786229	1413016000	0.00
humidity	DOUBLE	0.0	100.0	7761	50.00012248850552	28.86750987892046	25.007864495794472	50.00065217644044	75.02619651710863	1413016000	0.00
voltage	DOUBLE	0.0	5.0	555	2.500065617607927	1.4433790933627602	1.250503211815084	2.5001711260409593	3.7460747999880537	1413016000	0.00
current	DOUBLE	0.0	10.0	1177	5.000136832894723	2.88667209380417	2.49937980548707	5.001974925863886	7.49916284729135	1413016000	0.00
resistance	DOUBLE	0.0	10000.0	1812389	5000.003452122564	2886.755667209653	2502.1300879043476	5000.173776842859	7500.563493510334	1413016000	0.00
sample_rate	BIGINT	1	1000	1232	500.5025258864726	288.67384154007715	251	500	751	1413016000	0.00

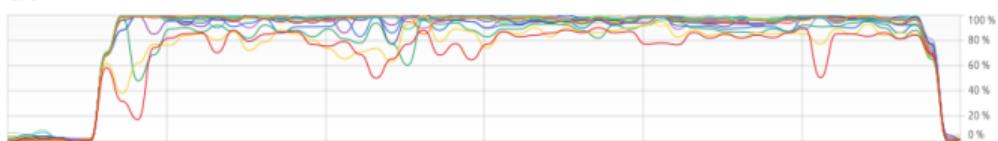
10 rows

12 columns

Run Time (s): real 257.893 user 4876.925018 sys 11.167547

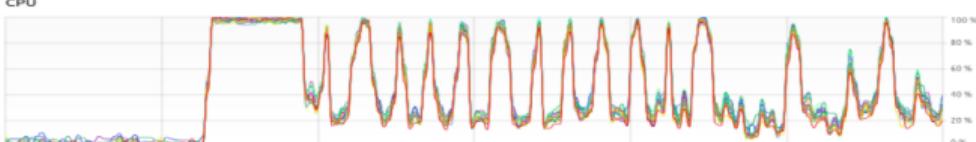
Local DuckDB table:

 CPU



Remote PostgreSQL table:

 CPU



HISTOGRAM command

```
select
    key_value['key'] as bucket_max,
    key_value['value'] as data_points
from (
    select unnest(map_entries(histogram(round(humidity/10)*10))) as key_value
    from probe_measurements where humidity is not null);
```

bucket_max double	data_points uint64
0.0	70589516
10.0	141280714
20.0	141311460
30.0	141304167
40.0	141291430
50.0	141299513
60.0	141312375
70.0	141298756
80.0	141312497
90.0	141287019
100.0	70728553

11 rows 2 columns

Run Time (s): real 4.541 user 89.779258 sys 0.040516

PIVOT command

- Can be very memory intensive, depends on the number of columns and rows

PIVOT (

```
SELECT probe_id, measurement_time, temperature
FROM probe_measurements
WHERE measurement_time BETWEEN '2024-12-01' AND '2024-12-30')
measurement_time::date
    AVG(temperature) AS avg_temperature GROUP BY probe_id;
```

Run Time (s): real: 16.

Direct work with data file

- Functions like SUMMARIZE, HISTOGRAM, PIVOT can run on external data files
- DuckDB uses parallel readers to process large data files efficiently
- Enables real-time data analysis without the need for data import

```
D summarize from 'probe_measurements_100gb.csv';
100%
```

column_name	column_type	min	max	approx_unique	avg	std	q25	q50	q75	count	null_percentage
varchar	varchar	varchar	varchar	int64	varchar	varchar	varchar	varchar	varchar	int64	decimal(9,2)
measurement_id	BIGINT	1	14130160000	1527453031	706508000..5	407902584..1289552	353460881	706262751	1059619539	1413016000	0.00
probe_id	BIGINT	1	1000	1232	500..5051361102776	288..67566940331545	250	501	750	1413016000	0.00
measurement_time	TIMESTAMP	2024-12-18 23:39:37	2024-12-30 16:20:48	1109862	2024-12-21 22:32:59..337439	2024-12-24 20:00:33..598275	2024-12-27 17:28:59..862436	2024-12-27 17:28:59..862436	1413016000	0.00	
temperature	DOUBLE	-50..0	50..0	9552	-2..7343037870649e-05	28..86754646608487	-24..97861983122555	0..0010965803637129083	24..99057013871463	1413016000	0.00
pressure	DOUBLE	900..0	1100..0	18127	999..99841557480876	57..7346837228424	950..0124387507778	1000..15833945617	1049..9635827052468	1413016000	0.00
humidity	DOUBLE	0..0	100..0	7761	50..00012488505006	28..86750987920785	24..9760938314825	49..987955528933654	75..00055998729647	1413016000	0.00
voltage	DOUBLE	0..0	5..0	555	2..500065617607933	1..443799833626929	1..2510427633778147	2..4994305489855044	3..749225082297946	1413016000	0.00
current	DOUBLE	0..0	10..0	1177	5..000136832894714	2..88667293042673	2..5014135758725295	4..9994664650448	7..499699661017365	1413016000	0.00
resistance	DOUBLE	0..0	10000..0	1812389	5000..009352122519	2866..7558672095927	2501..0407878575295	4999..7076553837607	7498..754915867427	1413016000	0.00
sample_rate	BIGINT	1	1000	1232	500..502528864726	288..6738415400651	250	501	750	1413016000	0.00

10 rows 12 columns

Run Time (s): real 306.288 user 5008.840000 sys 47.870835

```
D PIVOT (
  SELECT probe_id, measurement_time, temperature
  FROM 'probe_measurements_100gb.csv'
  WHERE measurement_time BETWEEN '2024-12-01' AND '2024-12-30')
ON measurement_time::date
USING AVG(temperature) AS avg_temperature GROUP BY probe_id;
100%
```

probe_id	2024-12-18_avg_tem_double	2024-12-19_avg_tem_double	2024-12-20_avg_tem_double	2024-12-21_avg_tem_double
17	4..781382978723404	0..0659882476780733	0..16312828213207278	-0..07717688773880788
915	4..841684210526315	-0..009986179062481	-0..008431021601162	0..1795473186896114
143	-1..957603305785129	0..145583758525312	0..07717094975725977	-0..12720375882521211
9	1..5455056179775282	0..09855607778879723	0..03765412237862668	0..04697267741617673
693	-1..71313131313128	-0..013763761349344	0..097529871959382	0..13564543425474802

Rudimentary Migration Tool

Migration from MySQL to PostgreSQL

- Example: Migrating a simple MySQL table with special data types

```
-- Create the table
CREATE TABLE special_data_types (
    id INT AUTO_INCREMENT PRIMARY KEY,
    name VARCHAR(50) NOT NULL,
    status ENUM('active', 'inactive', 'pending') NOT NULL,
    permissions SET('read', 'write', 'execute') NOT NULL,
    small_number TINYINT NOT NULL,
    medium_number MEDIUMINT NOT NULL,
    description TEXT,
    data BLOB,
    created_at DATE NOT NULL
);

-- Insert 10 rows of data
INSERT INTO special_data_types (name, status, permissions, small_number, medium_number, description, data, created_at) VALUES
('Alice', 'active', 'read,write', 5, 1000, 'Alice description', 'Alice data', '2023-01-01'),
('Bob', 'inactive', 'read', 10, 2000, 'Bob description', 'Bob data', '2023-02-01'),
('Charlie', 'pending', 'write,execute', 15, 3000, 'Charlie description', 'Charlie data', '2023-03-01'),
('David', 'active', 'read,write,execute', 20, 4000, 'David description', 'David data', '2023-04-01'),
('Eve', 'inactive', 'execute', 25, 5000, 'Eve description', 'Eve data', '2023-05-01'),
('Frank', 'pending', 'read,write', 30, 6000, 'Frank description', 'Frank data', '2023-06-01'),
('Grace', 'active', 'read', 35, 7000, 'Grace description', 'Grace data', '2023-07-01'),
('Hank', 'inactive', 'write,execute', 40, 8000, 'Hank description', 'Hank data', '2023-08-01'),
('Ivy', 'pending', 'read,write,execute', 45, 9000, 'Ivy description', 'Ivy data', '2023-09-01'),
('Jack', 'active', 'execute', 50, 10000, 'Jack description', 'Jack data', '2023-10-01');
```

Migration from MySQL to PostgreSQL

- Steps for migration:
- ATTACH MySQL database
- ATTACH PostgreSQL database
- CREATE TABLE in PostgreSQL AS SELECT FROM MySQL

```
D ATTACH 'host=mysql_container port=3306 user=root password=root database=duckdb_test' AS mysql_duckdb_test (TYPE MYSQL);
D SELECT * FROM mysql_duckdb_test.special_data_types;
```

id int32	name varchar	status varchar	permissions varchar	small_number int8	medium_number int32	description varchar	data blob	created_at date
1	Alice	active	read,write	5	1000	Alice description	Alice data	2023-01-01
2	Bob	inactive	read	10	2000	Bob description	Bob data	2023-02-01
3	Charlie	pending	write,execute	15	3000	Charlie description	Charlie data	2023-03-01
4	David	active	read,write,execute	20	4000	David description	David data	2023-04-01
5	Eve	inactive	execute	25	5000	Eve description	Eve data	2023-05-01
6	Frank	pending	read,write	30	6000	Frank description	Frank data	2023-06-01
7	Grace	active	read	35	7000	Grace description	Grace data	2023-07-01
8	Hank	inactive	write,execute	40	8000	Hank description	Hank data	2023-08-01
9	Ivy	pending	read,write,execute	45	9000	Ivy description	Ivy data	2023-09-01
10	Jack	active	execute	50	10000	Jack description	Jack data	2023-10-01

10 rows 9 columns

```
D ATTACH 'host=postgres_container port=5432 user=postgres password=postgres dbname=duckdb_test' as pg_duckdb_test (TYPE POSTGRES, SCHEMA 'public');
D CREATE TABLE pg_duckdb_test.special_data_types as SELECT * FROM mysql_duckdb_test.special_data_types;
```

Migration from MySQL to PostgreSQL

- Result: PostgreSQL table after migration

```
duckdb_test=# \ds+ special_data_types
Table "public.special_data_types"
 Column |      Type       | Collation | Nullable | Default | Storage | Compression | Stats target | Description
-----+----------------+-----+-----+-----+-----+-----+-----+-----+
 id   | integer        |           |           |          | plain   |           |           |
 name | character varying |           |           |          | extended |           |           |
 status | character varying |           |           |          | extended |           |           |
 permissions | character varying |           |           |          | extended |           |           |
 small_number | smallint       |           |           |          | plain   |           |           |
 medium_number | integer        |           |           |          | plain   |           |           |
 description | character varying |           |           |          | extended |           |           |
 data | bytea          |           |           |          | extended |           |           |
 created_at | date           |           |           |          | plain   |           |           |
Access method: heap

duckdb_test=# select * from special_data_types ;
 id | name | status | permissions | small_number | medium_number | description | data | created_at
----+-----+-----+-----+-----+-----+-----+-----+-----+
 1 | Alice | active | read,write | 5 | 1000 | Alice description | \x416c6963652064617461 | 2023-01-01
 2 | Bob | inactive | read | 10 | 2000 | Bob description | \x426f622064617461 | 2023-02-01
 3 | Charlie | pending | write,execute | 15 | 3000 | Charlie description | \x436861726c69652064617461 | 2023-03-01
 4 | David | active | read,write,execute | 20 | 4000 | David description | \x44617669642064617461 | 2023-04-01
 5 | Eve | inactive | execute | 25 | 5000 | Eve description | \x4576652064617461 | 2023-05-01
 6 | Frank | pending | read,write | 30 | 6000 | Frank description | \x4672616e6b2064617461 | 2023-06-01
 7 | Grace | active | read | 35 | 7000 | Grace description | \x47726163652064617461 | 2023-07-01
 8 | Hank | inactive | write,execute | 40 | 8000 | Hank description | \x48616e6b2064617461 | 2023-08-01
 9 | Ivy | pending | read,write,execute | 45 | 9000 | Ivy description | \x4976792064617461 | 2023-09-01
10 | Jack | active | execute | 50 | 10000 | Jack description | \x4a61636b2064617461 | 2023-10-01
(10 rows)
```

Migration from MS SQL to PostgreSQL

- Combining DuckDB with Polars to read MS SQL table and write it to PostgreSQL

```
import pyodbc
import polars as pl
import duckdb

conn_str = "DRIVER={ODBC Driver 17 for SQL Server};SERVER=localhost,1433;DATABASE=master;UID=sa;PWD=..."
connection = pyodbc.connect(conn_str)

query = "SELECT * FROM dbo.small_mssql_table"
df = pl.read_database(query, connection)
connection.close()

duckdb_conn = duckdb.connect(database=':memory:')
duckdb_conn.execute("""
    ATTACH 'dbname=duckdb_test user=postgres password=postgres host=localhost port=5432'
    AS pg (TYPE POSTGRES, SCHEMA 'public')")
duckdb_conn.execute("CREATE TABLE pg.small_mssql_table AS SELECT * FROM df")
duckdb_conn.close()
```

Migration from MS SQL to PostgreSQL - Batch Processing

```
import pyodbc
import polars as pl
import duckdb

source_table_name = 'dbo.big_mssql_table'
target_table_name = 'pg.big_mssql_table'

batch_size = 1000
conn_str = "DRIVER={ODBC Driver 17 for SQL Server};SERVER=localhost,1433;DATABASE=master;UID=sa;PWD=..."
connection = pyodbc.connect(conn_str)

mssql_cursor = connection.cursor()
query = f"SELECT * FROM {source_table_name}"
mssql_cursor.execute(query)
rows = mssql_cursor.fetchmany(batch_size)

schema = {column[0]: pl.datatypes.DataType.from_python(column[1]) for column in mssql_cursor.description}
df = pl.DataFrame([dict(zip(schema.keys(), row)) for row in rows], schema=schema, orient="row")

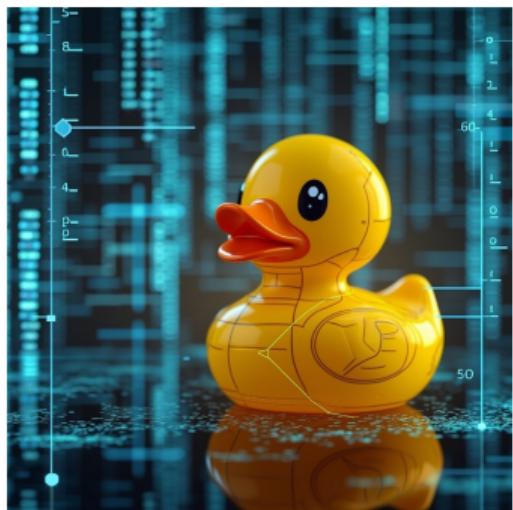
duckdb_conn = duckdb.connect(database=':memory:')
duckdb_conn.execute("""
    ATTACH 'dbname=duckdb_test user=postgres password=postgres host=localhost port=5432'
    AS pg (TYPE POSTGRES, SCHEMA 'public')"""
)
duckdb_conn.execute(f"CREATE TABLE {target_table_name} AS SELECT * FROM df")

while rows:
    rows = mssql_cursor.fetchmany(batch_size)
    if rows:
        df = pl.DataFrame([dict(zip(schema.keys(), row)) for row in rows], schema=schema, orient="row")
        duckdb_conn.execute(f"INSERT INTO {target_table_name} SELECT * FROM df")

mssql_cursor.close()
connection.close()
duckdb_conn.close()
```

Summary

- Powerful complementary tool for PostgreSQL
- Uses PostgreSQL SQL parser and follows its SQL conventions
- Directly works with Parquet, CSV, JSON, and other data formats
- Intelligent ETL tool for data imports and exports
- Simplifies data format conversions
- Enables easy cross-database joins
- Provides SUMMARIZE stats, HISTOGRAM, PIVOT on tables or data files
- Rudimentary migration tool for table snapshots
- Database file easily shared across different platforms
- Direct use of Python Pandas/Polars data frames extends capabilities



THANK YOU

- Questions?
- Josef Machytka <josef.machytka@netapp.com>

