

# DuckDB 在数据处理中的应用

叶家炜

Jan 16, 2026

副标题：从入门到高级应用，探索 DuckDB 如何简化大数据处理

作者：技术博客作者 / 发布日期：2024-10-01 / 标签：DuckDB、数据处理、SQL、嵌入式数据库、大数据  
想象一下，你作为数据分析师，手握一台普通笔记本电脑，需要处理数 GB 甚至 TB 级别的 Parquet 文件。传统方案如 Pandas 往往因内存爆炸而崩溃，Spark 则需要复杂的集群部署和数小时的等待。这时，DuckDB 横空出世，它是一个开源的嵌入式列式 SQL OLAP 数据库，专为分析型查询而生，无需服务器、无需配置，直接在你的进程中运行，就能以惊人的速度执行复杂查询。DuckDB 的核心在于其向量化查询引擎和零拷贝机制，能够按 SIMD 指令批量处理数据，比传统行式数据库快上数十倍。它支持多种数据格式如 Parquet、CSV 和 Arrow，直接查询文件而无需 ETL 预处理。这篇文章将带你从基础入手，深入探索 DuckDB 在数据处理中的核心优势和实际应用场景。我们针对数据分析师、数据工程师以及 Python 或 R 用户，逐步展示如何用 DuckDB 简化本地数据探索、大规模 ETL 和实时分析。接下来，我们从基础知识开始，一步步揭开它的革命性实践。

## 1 DuckDB 基础知识

DuckDB 的架构设计独树一帜，它采用嵌入式模式，直接在宿主进程中运行，无需独立的服务器进程，这意味着零部署成本，特别适合笔记本电脑或容器环境。其列式存储结合向量化执行引擎，只读取查询所需的列，并利用 SIMD 指令（如 AVX-512）批量处理向量数据，这让它在 OLAP 工作负载下比 Pandas 快 10 到 100 倍。同时，DuckDB 原生支持 Parquet、CSV、JSON 和 Apache Arrow 等格式，你可以直接用 SQL 查询海量文件，而无需先加载到内存。此外，它扩展了标准 SQL，内置窗口函数、CTE 和 JSON 操作符，完美符合 ANSI SQL 标准并针对分析优化。

安装 DuckDB 极其简单。在 Python 环境中，只需运行 `pip install duckdb` 即可集成到你的 Jupyter Notebook 或脚本中。对于 CLI 用户，官网提供预编译二进制文件，一键下载即可使用。让我们来看一个快速上手示例，假设你有一个名为 `sales.csv` 的本地文件，包含订单数据。我们用 DuckDB 查询其月度总销售额，并与 Pandas 对比性能。

```
1 import duckdb
2 import pandas as pd
3 import time
4
5 # DuckDB 查询
6 con = duckdb.connect()
7 start = time.time()
8 result = con.execute("""
9     SELECT
10        month,
11        SUM(sales) AS total_sales
12    FROM sales
13   GROUP BY month
14 ORDER BY month
15""")
```

```

9   SELECT DATE_TRUNC('month', order_date) AS month,
10      SUM(amount) AS total_sales
11  FROM 'sales.csv'
12  GROUP BY 1
13  ORDER BY 1
14  """).fetchdf()
15 duckdb_time = time.time() - start
16 print(result)
17 print(f"DuckDB 时间: {duckdb_time:.2f}s")

```

这段代码首先导入 DuckDB 和 Pandas，并创建一个内存数据库连接 con。DATE\_TRUNC('month', order\_date) 是 SQL 标准函数，用于截取日期到月级别；SUM(amount) 计算总销售额，按月分组并排序。关键是 FROM 'sales.csv'，DuckDB 直接扫描 CSV 文件而无需加载全表，这避免了 Pandas 的内存峰值。执行 fetchdf() 将结果转为 Pandas DataFrame，便于后续可视化。假设文件为 1GB，该查询通常在 1 秒内完成，而 Pandas 版本 (pd.read\_csv + groupby) 可能需 10 秒以上，且内存占用高出数倍。这展示了 DuckDB 的零拷贝优势：数据在列式格式下直接向量化处理，无需序列化。

## 2 DuckDB 在数据处理中的核心应用场景

在本地数据探索与 ETL 场景中，DuckDB 闪耀光芒。数据分析师常在 Jupyter 中处理 GB 级 CSV 或 Parquet 文件，传统工具易卡顿。DuckDB 允许你用纯 SQL 进行聚合、JOIN 和窗口函数计算。以 TPC-H 基准数据集为例，假设有一个 10GB 的 orders.parquet 和 lineitem.parquet，我们计算供应商交付延迟统计。

```

1 result = con.execute("""
2     SELECT o.supplier_id,
3         AVG(DATE_PART('day', l.shipdate - l.receiptdate)) AS avg_delay
4     FROM 'orders.parquet' o
5     JOIN 'lineitem.parquet' l ON o.orderkey = l.orderkey
6     WHERE l.shipdate > l.receiptdate
7     GROUP BY 1
8     ORDER BY 2 DESC
9  """).fetchdf()

```

这里，DuckDB 的列式存储确保 JOIN 只涉及必要列，DATE\_PART('day', ...) 计算天数差，自动利用分区剪枝 (pruning) 跳过无关数据块。相比 Pandas 的 merge，内存使用降低 80%，查询时间从分钟级降至秒级。这种能力让 ETL 管道从繁琐脚本转为简洁 SQL。

DuckDB 与 Python/R 生态的无缝集成进一步放大其价值。通过 query().df() 或 pl.from\_arrow()，它可与 Polars 和 Pandas 互操作，甚至通过 Ibis 框架提供统一 SQL 接口。举例，从 S3 读取 Parquet 并结合 Polars 做特征工程：

```

1 import duckdb
2 import polars as pl

```

```

3 df = duckdb.query("""
5   SELECT user_id,
8     AVG(order_value) OVER (PARTITION BY region) AS avg_region_value
7   FROM 's3://bucket/sales.parquet'
9 """).pl() # 转为 Polars DataFrame
9 features = df.with_columns(pl.col("avg_region_value").rank("dense").alias("value_rank
    → "))

```

这段代码启用 HTTPFS 扩展 (DuckDB 内置)，直接访问 S3；窗口函数 AVG OVER 计算区域均值，Polars 接管后续排名特征生成。这种链式工作流让机器学习管道高效无比。

对于大规模数据处理，DuckDB 支持联邦查询和扩展。HTTPFS 允许查询云存储如 S3 或 GCS，Spatial 扩展处理地理数据。我们可以跨多个 Parquet 文件执行 UNION ALL 和 GROUP BY：

```

1 result = con.execute("""
2   SELECT region, SUM(revenue) AS total
3   FROM read_parquet(['s3://bucket/2023/*.parquet'])
4   GROUP BY 1
5 """).fetchdf()

```

`read_parquet` 自动并行扫描分区文件，`predicate pushdown` 将过滤条件推到存储层，极大提升效率。在实时场景，DuckDB 可集成 Kafka 或 Redis，例如流式日志管道中持续查询最新分区。

### 3 实际案例分析

让我们通过电商销售数据分析这个入门级案例，感受 DuckDB 的实战魅力。假设有一个 10GB 的 `orders.parquet`，包含用户订单记录。任务是计算月度 GMV、Top 用户和 RFM 模型 (Recency、Frequency、Monetary)。

```

1 gmv_query = """
2   SELECT DATE_TRUNC('month', order_date) AS month,
3         SUM(amount) AS gmv
4   FROM 'orders.parquet'
5   GROUP BY 1 ORDER BY 1
6 """
7 top_users = """
8   SELECT user_id, SUM(amount) AS total_spent,
9     NTILE(5) OVER (ORDER BY COUNT(*) DESC) AS rfm_f
10  FROM 'orders.parquet'
11  GROUP BY 1
12 """
13 con.execute(gmv_query).fetchdf()

```

首先，GMV 查询使用 DATE\_TRUNC 分组求和，整个 10GB 文件在 3 秒内处理完，内存峰值仅 1.5GB。其次，RFM 计算中 NTILE(5) 将用户按频次分桶，ORDER BY COUNT(\*) DESC 确保 Top 用户优先。这比 Pandas groupby + quantile 简单高效，后续可直接用 Matplotlib 绘图：gmv\_df.plot(x='month', y='gmv')。转向中级案例：TB 级 Nginx 日志处理与异常检测。数据为 JSON 格式日志，我们检测 Top IP 和异常峰值。

```

1 anomaly_query = """
2     SELECT ip,
3         COUNT(*) AS requests,
4         AVG(request_time) OVER (ORDER BY log_time
5             ROWS BETWEEN 100 PRECEDING AND CURRENT ROW) AS rolling_avg
6
7     FROM read_json_auto('logs/*.json')
8     WHERE request_time > 1.0 -- 慢请求
9     GROUP BY ip
10    HAVING requests > (SELECT AVG(requests) * 3 FROM (SELECT COUNT(*) as requests FROM
11        → read_json_auto('logs/*.json') GROUP BY window(log_time, '1 hour'))))
12
13"""
14
15 result = con.execute(anomaly_query).fetchdf()

```

read\_json\_auto 自动推断 schema，窗口函数计算过去 100 条的滚动平均，HAVING 子句用自连接检测 3 σ 峰值。整个 TB 级扫描只需分钟级，对比 Dask 的延迟调度，DuckDB 单机更快、更易调试。结果导出 Arrow 格式 con.arrow(result) 给 scikit-learn 训练异常模型。

高级案例转向企业级 BI Dashboard。我们集成 Streamlit，实现多源联邦查询：本地数据库 + S3 Parquet。

```

1 import streamlit as st
2 con = duckdb.connect()
3 query = st.text_area("输入SQL", value="")
4     SELECT * FROM postgres_query('host=localhost dbname=prod', 'SELECT * FROM sales
5         → LIMIT 100')
6     UNION ALL
7     SELECT * FROM 's3://bucket/reports.parquet' WHERE date > '2024-01-01'
8 """
9 if st.button("执行"):
10     st.dataframe(con.execute(query).fetchdf())

```

postgres\_query 扩展扫描远程 Postgres，UNION ALL 融合云数据。优化中，用 CREATE MATERIALIZED VIEW 预计算视图，并设置 PRAGMA threads=8 启用多核。

## 4 高级技巧与最佳实践

性能优化是 DuckDB 的强项。通过 PRAGMA threads=16; PRAGMA memory\_limit='8GB'; 配置线程数和内存上限，确保资源高效利用。优先用 SQL 原生函数而非 UDF，避免解释器开销；依赖分区剪枝和谓词下推，如在 WHERE 中指定日期范围，自动跳过无关 Parquet 行。调试时，EXPLAIN ANALYZE SELECT ... 输出查询

计划树，展示向量化 JOIN 和哈希表大小。

DuckDB 不适合高并发 OLTP，转而推荐 Postgres；对于云需求，可用 MotherDuck 服务。对于监控，查询 profile 揭示瓶颈，如 I/O 绑定的扫描需优化分区。

## 5 与其他工具对比

Pandas 以灵活 API 著称，但在大规模数据上内存饥饿，而 DuckDB 在低内存大数据场景中胜出，提供 SQL 简洁性。Polars 凭借 Rust 实现速度飞快，DuckDB 则以熟悉 SQL 语法取胜，无需学习新 API。ClickHouse 擅长海量分布式数据，DuckDB 更适合本地嵌入式原型。Spark 的分布式能力强大，但单机快速迭代时 DuckDB 更敏捷简便。

## 6 结论与展望

DuckDB 以其零配置、高性能和普适集成，彻底革新了数据处理范式，从本地探索到联邦查询，它让复杂任务化为优雅 SQL。立即安装试用吧，GitHub 示例仓库 [github.com/example/duckdb-blog](https://github.com/example/duckdb-blog) 含所有代码。展望未来，DuckDB 1.0 将强化稳定性，WASM 支持浏览器分析，更多扩展如 ML 集成将至。DuckDB 不是取代工具，而是你数据旅程中的瑞士军刀。

参考资源：

官网：[duckdb.org](https://duckdb.org)

文档：<https://duckdb.org/docs/>

论文：DuckDB: RadixJoin + Vectorwise

你的数据处理痛点是什么？欢迎评论区分享！