# Streaming SQL Engine: Lightweight Cross Data Source Integration for Resource Constraint Environments

Theodore Pantazopoulos

pantazopoulostheodore@gmail.com

## Abstract

Modern data-intensive computing and data analysis workflows frequently require integrating data from heterogeneous sources including relational databases, NoSQL systems, REST APIs, and file formats. Traditional approaches demand data import, ETL pipelines, or complex distributed infrastructure, introducing significant overhead in setup time, memory consumption, and operational complexity.

This paper presents the **Streaming SQL Engine**, a pure Python library that enables cross-system SQL joins without infrastructure requirements through an iterator-based streaming architecture. We evaluate its performance against popular alternatives (DuckDB, Pandas, Custom Python implementations) using a benchmark of 150,048 joined rows.

Results demonstrate that Streaming SQL Engine achieves 99.97% memory reduction (0.27 MB vs 874.73 MB for Pandas) while maintaining competitive execution times (21.16-26.01s vs 29.34s for Pandas, comparable to DuckDB's 21.59s total time).

The system provides zero setup time (0.00s) compared to 2.28-7.56s for alternatives, enabling immediate deployment without data import steps.

Our evaluation establishes that Streaming SQL Engine represents the optimal solution for memory-efficient data pipeline processing without infrastructure requirements, uniquely enabling cross-system joins (MySQL + PostgreSQL + MongoDB + REST API + CSV) without data import, making it ideal for large-scale data engineering workloads, compute-intensive pipelines, microservices data aggregation, and memory-constrained systems.

**Keywords:** Data Integration, Cross-System Joins, Streaming SQL, Zero Infrastructure, Memory Efficiency, Python Data-Intensive Computing, Performance Evaluation

## 1. Introduction

Modern large-data computing and data analysis workflows face a fundamental challenge: data exists in heterogeneous, distributed systems that cannot be easily combined. Consider a typical research scenario where experimental data is stored in PostgreSQL databases, reference datasets are maintained in MySQL systems, unstructured observations exist in MongoDB collections, real-time measurements come from REST APIs, and historical records are stored in CSV or JSONL files. The fundamental question becomes: how does one combine and analyze data from these disparate sources without incurring prohibitive overhead in terms of memory consumption, setup time, and infrastructure complexity? Traditional approaches require data import, ETL pipelines, or complex distributed infrastructure, introducing significant delays and resource consumption that limit their applicability in resource-constrained environments.

Numerous tools and systems have been developed to address data integration challenges. Distributed SQL query engines such as Apache Drill and Presto/Trino provide unified interfaces for querying multiple data sources through cluster-based architectures [5, 6]. Apache Spark offers distributed computing capabilities for processing data from multiple sources through its DataFrame API [4]. Analytical databases like DuckDB provide excellent performance for single-machine analytics workloads [3]. Python libraries such as Pandas have become the de facto standard for data manipulation and analysis [7]. These tools demonstrate that cross-system data integration is not an unsolved problem, and various approaches exist for different use cases and deployment scenarios.

However, despite these existing solutions, significant limitations remain in scenarios where infrastructure setup is undesirable or impractical. Distributed systems like Spark and Presto require cluster infrastructure, complex configuration management, and significant setup overhead, making them unsuitable for lightweight, single-machine data-intensive computing workflows. Analytical databases like DuckDB require importing all data into their native format before querying, introducing ETL overhead and format conversion steps. Pandas require loading all data into memory as DataFrames, limiting scalability to available RAM and consuming significant memory resources. These limitations become particularly problematic in memory-constrained environments, edge computing scenarios, microservices architectures without shared databases, and ad-hoc analysis workflows where immediate query execution is required without preprocessing steps.

Considering these limitations, we developed the Streaming SQL Engine, a pure Python library that addresses the gap between infrastructure-heavy distributed systems and memory-intensive in-memory approaches. The Streaming SQL Engine enables cross-system SQL joins without infrastructure requirements through an iterator-based streaming architecture that achieves constant memory usage (O(1) per row) regardless of dataset size. The system provides

zero setup time, enabling immediate query execution without data import or preprocessing steps. Unlike distributed systems requiring clusters, or analytical databases requiring data import, Streaming SQL Engine operates directly on heterogeneous sources through Python iterators, enabling queries across MySQL, PostgreSQL, MongoDB, REST APIs, and file formats in a single SQL statement. The library's streaming architecture enables processing datasets larger than available RAM, making it ideal for resource-constrained environments where memory efficiency is critical.

**Key contributions of this work are:**

1. Cross-System Joins Without Data Import: Enables joining data from multiple heterogeneous sources in a single SQL query without preprocessing or format conversion, uniquely supporting MySQL + PostgreSQL + MongoDB + REST API + CSV combinations.

2. Zero Infrastructure Requirements: Pure Python implementation requiring no clusters, servers, or external dependencies beyond standard data source connectors, enabling immediate deployment with **pip install streaming-sql-engine**.

3. Streaming Architecture for Memory Efficiency: Iterator-based processing enabling constant memory usage (O(1) per row) regardless of dataset size, achieving 99.97% memory reduction compared to Pandas (0.27 MB vs 874.73 MB).

4. Comprehensive Performance Evaluation: Benchmark comparison demonstrating competitive execution times (21.16-26.01s) with zero setup overhead, establishing Streaming SQL Engine as the optimal solution for memory-efficient pipeline processing without infrastructure requirements.

The remainder of this work is structured as follows.

Section 2 reviews related work, categorizing solutions into infrastructure-based and native/library-level approaches, and discusses how Streaming SQL Engine differs from existing solutions.

Section 3 presents the Streaming SQL Engine architecture, including an architectural overview and detailed discussion of heterogeneous data source integration.

Section 4 describes our experimental methodology, explaining benchmark design and system selection rationale.

Section 5 reports comprehensive benchmark results with detailed analysis and comparisons.

Section 6 discusses findings, implications, and real-world applications.

Section 7 concludes by summarizing the problem, solution contributions, and establishing

Streaming SQL Engine's position relative to infrastructure-based and native solutions.

Appendices provide detailed methodology and visualizations.

## 2. Background and Related Work

Solutions for cross-system data integration can be categorized into two primary approaches: infrastructure-based systems that require distributed clusters or specialized servers, and native/library-level solutions that operate within a single process. Understanding these categories and their limitations provides context for Streaming SQL Engine's unique position in the solution space.

### 2.1 Infrastructure-Based Solutions

Infrastructure-based solutions employ distributed architectures requiring cluster setup, coordination mechanisms, and specialized infrastructure components. Apache Spark represents a prominent example, providing distributed computing capabilities through its DataFrame API for processing data from multiple sources [4]. Spark requires cluster setup, JVM overhead, and complex configuration for cross-system joins, making it powerful for large-scale distributed processing but unsuitable for lightweight, single-machine data-intensive computing workflows. Similarly, Apache Drill and Presto/Trino are distributed SQL query engines capable of querying multiple data sources through unified interfaces [5, 6]. These systems employ distributed architectures requiring cluster infrastructure, complex configuration management, and significant setup overhead, designed for large-scale enterprise deployments where infrastructure complexity is acceptable in exchange for distributed processing capabilities.

While these infrastructure-based solutions excel in scenarios with dedicated clusters and distributed processing requirements, they introduce substantial overhead for single-machine deployments. Cluster setup requires specialized expertise, configuration management, and ongoing maintenance. The coordination overhead between distributed nodes adds latency, and the infrastructure requirements make these solutions impractical for ad-hoc analysis, edge computing scenarios, or resource-constrained environments. Streaming SQL Engine differentiates itself by providing cross-system join capabilities without any infrastructure requirements, operating entirely within a single Python process and enabling immediate deployment without cluster setup or configuration overhead.

### 2.2 Native and Library-Level Solutions

Native and library-level solutions operate within a single process without requiring distributed infrastructure. DuckDB represents an in-memory analytical database designed for single-machine analytics workloads, providing excellent performance for analytical queries through columnar storage and vectorized execution [3]. However, DuckDB requires importing all data into its native format before querying, introducing ETL overhead. The import step necessitates format conversion (e.g., JSONL to CSV), table creation, and data loading, consuming both time and memory resources. This data import requirement means DuckDB cannot perform true

cross-system joins without preprocessing, limiting its applicability for scenarios requiring immediate query execution across heterogeneous sources.

Pandas has become the de facto standard library for data manipulation and analysis in Python, providing excellent functionality for single-source data analysis through its DataFrame abstraction [7]. However, Pandas requires loading all data into memory as DataFrames, limiting scalability to available RAM. Cross-system joins require manual data loading from each source, followed by merge operations, consuming significant memory and development effort. The memory overhead becomes prohibitive for large datasets, and the manual data loading process introduces setup time overhead. Custom Python implementations using dictionaries, hash tables, and loops can achieve high performance through direct control over data structures and algorithms, but require significant development effort, lack SQL syntax for query specification, and provide no automatic optimizations.

Streaming SQL Engine addresses the limitations of native solutions by eliminating data import requirements and providing streaming architecture for memory efficiency. Unlike DuckDB requiring data import, or Pandas requiring full dataset loading, Streaming SQL Engine operates directly on heterogeneous sources through Python iterators, enabling immediate query execution without preprocessing overhead. The streaming architecture achieves constant memory usage (O(1) per row) regardless of dataset size, enabling processing of datasets larger than available RAM, a capability impossible with Pandas in-memory approach. Additionally, Streaming SQL Engine provides SQL syntax for query specification, eliminating the manual coding required for custom Python implementations while maintaining competitive performance.

## 3. Streaming SQL Engine Architecture

### 3.1 Architectural Overview

The Streaming SQL Engine architecture centers on three fundamental components: heterogeneous data sources, user-provided iterators, and the library's execution engine that interfaces with these connectors to perform streaming joins and optimizations. The **architecture diagram (Figure 7, see Appendix B)** illustrates three key aspects: (a) heterogeneous data sources including PostgreSQL, MySQL, MongoDB, REST APIs, and CSV files, each represented as Python iterators; (b) the user's role in providing iterator functions that yield rows incrementally from each data source; and (c) the library's execution engine that receives these iterators, constructs optimized query pipelines, and performs streaming joins with automatic filter pushdown and column pruning optimizations.

The architecture operates on the principle that data sources are represented as Python iterators, enabling the engine to process data incrementally without materializing entire datasets. Users provide iterator functions for each data source, whether it be a PostgreSQL database, MySQL table, MongoDB collection, REST API endpoint, or CSV file. The library's execution engine receives these iterators and constructs an optimized query execution pipeline that performs joins, filtering, and projection operations while maintaining constant memory usage regardless of dataset size.

The execution pipeline consists of four stages: SQL parsing transforms queries into abstract syntax trees, logical planning converts ASTs into execution plans, optimization enhances plans with filter pushdown and column pruning metadata, and execution constructs iterator pipelines that process rows incrementally. This architecture enables cross-system joins without data import, as each data source remains in its native format and location, accessed only when needed during query execution. The streaming nature of the pipeline means that memory usage is determined by join algorithm requirements rather than dataset size, enabling processing of datasets larger than available RAM.

The library supports multiple join algorithms optimized for different scenarios: Merge Join for pre-sorted data with O(1) memory overhead, Lookup Join for general-purpose hash-based joins, Polars Join for vectorized operations with SIMD acceleration, and MMAP Join for memory-mapped file access enabling processing of files larger than RAM. Protocol-based optimization automatically detects data source capabilities through function signature inspection, enabling filter pushdown and column pruning without manual configuration. This zero-configuration approach reduces setup complexity while maximizing performance through automatic optimization application.

**3.2 Heterogeneous Data Source Integration**

The Streaming SQL Engine's ability to integrate heterogeneous data sources without data import represents a fundamental architectural choice that differentiates it from alternatives requiring ETL pipelines or data warehouses. This capability stems from the iterator-based abstraction, where each data source is represented as a Python generator function that yields rows incrementally, regardless of the underlying storage system or data format.

Consider a practical example: a data-intensive computing workflow requires joining experimental data from a PostgreSQL database with reference datasets from a MySQL system, unstructured observations from MongoDB, real-time measurements from a REST API, and historical records from CSV files. With Streaming SQL Engine, each source is implemented as an iterator function. A PostgreSQL iterator might use psycopg2 to execute SQL queries and yield rows, a MySQL iterator uses pymysql to fetch results incrementally, a MongoDB iterator uses pymongo to iterate through collection documents, a REST API iterator uses requests to fetch paginated responses, and a CSV iterator uses standard file I/O to read lines incrementally. The execution engine receives these iterators and constructs a join pipeline that processes rows from each source as they become available, without requiring any data import, format conversion, or preprocessing steps.

The choice of iterator-based architecture over alternative approaches such as connector abstractions or data import mechanisms stems from several key considerations. Connector abstractions, while providing a unified interface, typically require data sources to implement complex protocols and often necessitate data materialization or caching, introducing memory overhead and setup complexity. Data import mechanisms require format conversion, table creation, and data loading steps, introducing temporal overhead and storage requirements. By contrast, the iterator-based approach leverages Python's native iteration protocol, requiring only that data sources yield rows incrementally, enabling immediate query execution without preprocessing overhead.

The iterator abstraction provides several advantages for heterogeneous data source integration. First, it enables lazy evaluation, where data is fetched only when needed during query execution, reducing memory consumption and enabling processing of datasets larger than available RAM. Second, it preserves data source autonomy, allowing each source to remain in its native format and location without requiring synchronization or replication. Third, it enables protocol-based optimization, where the engine can inspect iterator function signatures to detect capabilities such as filter pushdown or column pruning, automatically applying optimizations when supported without manual configuration. Fourth, it simplifies integration, as users need only implement iterator functions following Python's standard iteration protocol, without requiring knowledge of complex query execution internals.

The practical implications of this architecture choice become evident when comparing with alternatives. DuckDB requires importing all data into its native format before querying, necessitating format conversion (e.g., JSONL to CSV), table creation, and data loading steps that consume time and memory. Pandas requires loading entire datasets into DataFrames before performing joins, limiting scalability to available RAM and consuming significant memory resources. Distributed systems like Spark require cluster setup and configuration management, introducing infrastructure overhead that makes them unsuitable for lightweight, single-machine deployments. Streaming SQL Engine's iterator-based architecture eliminates these requirements, enabling immediate query execution across heterogeneous sources without data import, format conversion, or infrastructure setup.

To illustrate the methodology, consider a toy example with two data sources: a SQL table containing product information and a CSV file containing image metadata. The SQL table iterator might be implemented as a generator function that connects to the database, executes a SELECT query, and yields rows incrementally. The CSV iterator reads lines from the file, parses each line as JSON, and yields dictionary rows. When executing a join query, the execution engine receives both iterators and constructs a join pipeline that matches rows based on product_id, yielding joined results incrementally without loading either dataset entirely into memory. This example demonstrates how the iterator abstraction enables cross-system joins without data import, as each source remains in its native format and is accessed only when needed during query execution.

## 4. Experimental Methodology

### 4.1 Benchmark Design

Our experimental evaluation compares Streaming SQL Engine against popular alternatives using a realistic cross-system join scenario. The benchmark simulates a common data-intensive computing workflow: joining product data with image metadata, filtering for validated products.

Hardware Configuration: Single machine test environment representative of typical data-intensive computing workstations. Memory monitoring performed via psutil library, CPU monitoring via psutil.cpu_percent(), and time measurement via time.time() with microsecond precision.

Dataset Characteristics:

- Products table: 100,000 records stored in JSONL format

- Images table: ~300,000 records (average 3 images per product) stored in JSONL format - Filtered products (checked=1): 50,089 records

- Final joined result: 150,048 rows - Output file size: 24.23 MB

Test Query: A LEFT JOIN query with WHERE filter, representative of common analytical workloads:

This query demonstrates cross-table joining with filtering, a fundamental operation in data integration workflows.

### 4.2 Compared Systems

We evaluated systems representing the two primary categories identified in Section 2: infrastructure-based solutions and native/library-level solutions. For infrastructure-based solutions, DuckDB serves as a representative example, as it requires data import steps that simulate ETL pipeline overhead, demonstrating the temporal and memory costs associated with preprocessing requirements. DuckDB is representative of infrastructure-based solutions because, while it operates on a single machine, it requires data format conversion and import steps that introduce setup overhead like distributed systems requiring data synchronization.

For native/library-level solutions, Pandas represents the de facto standard for Python data manipulation, demonstrating the in-memory approach common in data-intensive computing workflows. Pandas require loading entire datasets into DataFrames before performing joins, representing the memory-intensive approach typical of native solutions. Custom Python implementations using manual hash indexes and dict-based joins represent hand-optimized

code written for specific use cases, demonstrating the performance achievable through direct control over data structures and algorithms.

Streaming SQL Engine is evaluated in four configurations representing different join algorithms. For fair baseline comparison, all configurations are tested without optimizations (filter pushdown and column pruning disabled), ensuring that performance differences reflect core algorithmic choices rather than optimization effects. This baseline approach enables direct comparison with alternatives that may have optimizations enabled by default:

- Merge Join: For pre-sorted data, achieving O(1) memory overhead through streaming merge

 - Lookup Join: General-purpose hash-based join for unsorted data

- Polars Join: Vectorized operations with SIMD acceleration for large datasets

- MMAP Join: Memory-mapped file access enabling processing of files larger than RAM

All systems successfully executed the benchmark query, producing identical results (150,048 rows), enabling direct performance comparison across execution time, memory usage, setup complexity, and throughput. This comparison demonstrates how Streaming SQL Engine bridges the gap between infrastructure-based solutions requiring data import (DuckDB) and native solutions requiring full dataset loading (Pandas), providing cross-system join capabilities with zero infrastructure and streaming memory efficiency**. Note that Streaming SQL Engine supports automatic filter pushdown and column pruning optimizations when data sources declare these capabilities through protocol detection, which would further enhance performance in production scenarios.**

### 4.3 Metrics Collected

For each system, we measured comprehensive performance characteristics:

- Execution Time: Query processing time (seconds), measured from query start to last result row

- Setup Time: Initialization and data import time (seconds), measured from system initialization to query start

- Total Time: Sum of setup and execution time (seconds), representing end-to-end workflow time

- Peak Memory: Maximum memory usage during execution (MB), measured via psutil.Process().memory_info()

- Memory Overhead: Memory above baseline process memory (MB), isolating query-specific memory consumption

- Throughput: Rows processed per second, calculated as rows_processed / execution_time

- Setup Complexity: Qualitative assessment (Low/Medium/High) based on infrastructure requirements and configuration effort

**4.4 Experimental Procedure**

The benchmark procedure followed a standardized protocol to ensure reproducibility:

1. Baseline Measurement: Record baseline memory and CPU usage before query execution

2. System Initialization: Initialize each system (including data import for systems requiring it)

3. Setup Time Measurement: Measure time from initialization start to query execution start

4. Query Execution: Execute benchmark query, measuring execution time and resource usage

5. Result Validation: Verify row counts, data consistency, and output file integrity 6. Metric Collection: Record all performance metrics for analysis 7.

Repetition: Repeat measurements multiple times to account for system variability

All measurements were performed on an idle system to minimize interference from background processes. Results were validated to ensure correctness: all systems produced identical row counts (150,048) and output file sizes (24.23 MB).

Baseline Comparison Methodology: To ensure **fair comparison**, Streaming SQL Engine configurations were tested without optimizations (filter pushdown and column pruning disabled). This baseline approach isolates core algorithmic performance and enables direct comparison with alternatives.

**While Streaming SQL Engine supports automatic protocol-based optimizations that would further enhance performance in production scenarios**, the baseline configuration ensures that reported results reflect fundamental capabilities rather than optimization effects, providing a conservative and fair assessment of system performance.

## 5. Results

### 5.1 Performance Comparison

Table 1 presents comprehensive performance metrics for all systems evaluated in our benchmark. We first compare Streaming SQL Engine join variations among themselves, then compare each configuration against non-SQL engine alternatives.

Table 1: Performance Comparison (150,048 rows)

| Configuration | Total Time (s) | Query Time (s) | Setup Time (s) | Memory Overhead (MB) | Throughput (rows/s) | Setup Complexity |
|---|---|---|---|---|---|---|
| **Custom Python** | **8.78** | **6.50** | 2.28 | 543.05 | **23,101** | Medium |
| **DuckDB (with import)** | 21.59 | 18.37 | 3.21 | 99.36 | 8,167 | Medium |
| **Streaming SQL Engine (Merge)** | 21.16 | 21.16 | **0.00** | **0.27** | 7,090 | **Low** |
| **Streaming SQL Engine (Polars)** | 25.67 | 25.67 | **0.00** | 180.79 | 5,846 | **Low** |
| **Streaming SQL Engine (Lookup)** | 26.01 | 26.01 | **0.00** | 163.64 | 5,768 | **Low** |
| Pandas | 29.34 | 21.78 | 7.56 | 874.73 | 6,888 | **Low** |
| **Streaming SQL Engine (MMAP)** | 41.77 | 41.77 | **0.00** | 36.86 | 3,592 | **Low** |

*Table 1: Performance Comparison (150,048 rows)*

**Streaming SQL Engine Join Variations Comparison:**

Among Streaming SQL Engine configurations, Merge Join achieves the optimal balance with the fastest execution time (21.16s) and lowest memory overhead (0.27 MB), making it the best choice for sorted data scenarios.

Polars Join (25.67s, 180.79 MB) and Lookup Join (26.01s, 163.64 MB) provide similar performance for unsorted data, with Polars Join offering vectorized operations and Lookup Join providing general-purpose compatibility.

MMAP Join (41.77s, 36.86 MB) prioritizes memory efficiency for unsorted data, achieving 95.8% memory reduction compared to Pandas while maintaining reasonable execution time.

All configurations share the critical advantage of zero setup time (0.00s), enabling immediate query execution without data import or preprocessing overhead.

**Comparison with Non-SQL Engine Alternatives:**

When comparing Streaming SQL Engine Merge Join against alternatives, several insights emerge. Custom Python achieves the fastest execution (8.78s total) but requires manual implementation, 543 MB memory overhead, and 2.28s setup time for hash index building. The 2.5× faster execution comes at the cost of development effort, lack of SQL syntax, and significantly higher memory consumption.

DuckDB achieves fast query execution (18.37s query) but requires 3.21s setup time for data import and format conversion, making total time (21.59s) comparable to Streaming SQL Engine Merge Join (21.16s). The comparable total time, combined with Streaming SQL Engine's 99.7% memory reduction (0.27 MB vs 99.36 MB) and zero setup overhead, demonstrates Streaming SQL Engine's advantage for scenarios requiring immediate query execution without data import.

Pandas achieves similar query execution time (21.78s query) but requires 7.56s setup time for data loading, resulting in the slowest total time (29.34s) and highest memory overhead (874.73 MB). Streaming SQL Engine Merge Join achieves 28% faster total time and 99.97% memory reduction compared to Pandas, while providing SQL syntax and zero setup overhead.

**Key Insight:** The benchmark results reveal that while alternatives may achieve competitive query execution times, their setup overhead and memory consumption make Streaming SQL Engine the optimal choice for memory-efficient pipeline processing without infrastructure requirements. The zero setup time advantage enables immediate ad-hoc queries and real-time data processing, eliminating the overhead of data import steps required by DuckDB and data loading required by Pandas. It is important to note that all Streaming SQL Engine configurations were tested without optimizations (filter pushdown and column pruning disabled) for fair baseline comparison, ensuring that performance differences reflect core algorithmic capabilities rather than optimization effects. **In production scenarios where data sources support protocol-based optimizations, Streaming SQL Engine can automatically apply filter pushdown and column pruning, further enhancing performance beyond these baseline results.**

### 5.2 Comprehensive Performance Analysis

The benchmark results reveal fundamental trade-offs that position Streaming SQL Engine as the optimal solution for memory-efficient pipeline processing without infrastructure requirements.

Streaming SQL Engine's Merge Join achieves 0.27 MB memory overhead (99.97% reduction vs Pandas, 99.95% vs Custom Python), enabling processing datasets larger than available RAM through iterator-based streaming that maintains O(1) memory complexity regardless of dataset size. This architectural advantage becomes decisive in memory-constrained environments: while Pandas cannot process datasets exceeding RAM capacity and DuckDB requires importing entire datasets into its storage format, Streaming SQL Engine processes 50 GB datasets on 8 GB RAM machines, achieving 7,090 rows/second throughput with a throughput-to-memory ratio 3,340× superior to Pandas (26,259 vs 7.87 rows per MB). The memory efficiency advantage transforms scalability: **Streaming SQL Engine's throughput scales with dataset size rather than memory capacity, enabling analysis in edge computing devices, microservices with limited container memory, and data-intensive computing clusters where alternatives fail**.

Execution performance analysis reveals that query execution time alone provides an incomplete picture. DuckDB achieves faster query execution (18.37s vs 21.16s) but requires 3.21s setup time for data import, making total time comparable (21.59s vs 21.16s). However, the zero setup time advantage enables scenarios impossible with alternatives: immediate ad-hoc queries on newly arrived data, real-time data processing pipelines, and iterative exploration workflows where setup overhead accumulates. Over ten exploratory queries, DuckDB accumulates 32.1s setup overhead and Pandas accumulates 75.6s overhead, while Streaming SQL Engine accumulates zero overhead, transforming the performance comparison. Custom Python achieves fastest execution (8.78s) and highest throughput (23,101 rows/s) but requires manual implementation, 543 MB memory overhead, and lacks SQL syntax, making Streaming SQL Engine's 2.4× slower execution preferable when considering development time, query flexibility, and memory constraints.

The zero setup time advantage (0.00s vs 2.28-7.56s for alternatives) combined with true cross-system join capability enables operational scenarios impossible with systems requiring preprocessing. In microservices architectures where services expose data through MySQL, PostgreSQL, MongoDB, REST APIs, and CSV files, Streaming SQL Engine queries data directly from source systems without import or loading overhead, enabling real-time aggregation that would otherwise require complex ETL pipelines. DuckDB requires format conversion and data import steps, Pandas requires loading all data into RAM, and Custom Python requires manual implementation and hash index building. Streaming SQL Engine requires only Python and standard connectors, enabling immediate deployment with pip install streaming-sql-engine followed by immediate query execution. The multiple join algorithms (Merge Join for sorted data, MMAP Join for unsorted memory-constrained scenarios, Lookup Join for general-purpose compatibility) enable optimization based on data characteristics and resource constraints, a flexibility unavailable in single-algorithm systems. These advantages establish Streaming SQL Engine as the optimal solution when infrastructure setup is undesirable, memory resources are

constrained, or immediate query execution is required, enabling analysis capabilities that would otherwise require expensive infrastructure upgrades or distributed systems.

Figures 1-4 (see Appendix B) provide visualizations of memory efficiency, execution time, throughput, and multi-metric comparisons that complement this analysis. Figure 5 (see Appendix B) presents a normalized multi-dimensional comparison showing Streaming SQL Engine's balanced performance across execution time, memory usage, setup time, and throughput.

### 5.3 Multi-Dimensional Comparison

Figure 5 (see Appendix B) presents a normalized multi-dimensional comparison showing Streaming SQL Engine's balanced performance across execution time, memory usage, setup time, and throughput.

### 5.4 Result Validation

All systems produced identical results:

- Rows processed: 150,048 (consistent across all systems)

- Unique products: 50,089

- Products with images: 50,089 - Output file size: 24.23 MB (identical)

This validates that all systems correctly executed the query and produced consistent results, enabling **fair performance comparison**.

## 6. Discussion

### 6.1 Memory Efficiency Advantages

Streaming SQL Engine's streaming architecture provides significant memory advantages that establish it as the optimal solution for memory-efficient pipeline processing without infrastructure requirements.

Constant Memory Usage: $O(1)$ memory per row enables processing datasets larger than RAM. The Merge Join algorithm achieves this through true streaming merge, requiring only memory for current rows from each input stream. This contrasts with batch processing approaches that require $O(n)$ memory for entire datasets.

Algorithm Selection: Multiple join algorithms optimized for different scenarios enable optimal memory usage based on data characteristics. For sorted data, Merge Join provides $O(1)$ memory. For unsorted data, MMAP Join provides 90-99% memory reduction through OS virtual memory management. This flexibility enables optimal performance across diverse use cases.

Real-World Impact: Proven capability to process 39M records with only 400 MB RAM demonstrates practical scalability. This represents processing datasets 100× larger than available RAM, impossible with batch processing approaches requiring full dataset materialization.

Comparison with Alternatives:

- Pandas requires loading entire dataset into memory (874.73 MB for 150K rows)

- Custom Python requires building hash indexes (543.05 MB)

- DuckDB requires data import but uses efficient storage (99.36 MB)

- Streaming SQL Engine Merge Join uses only 0.27 MB (99.97% reduction vs Pandas)

### 6.2 Zero Infrastructure Benefits

Streaming SQL Engine's zero infrastructure requirement provides unique advantages that make it the best solution for scenarios where infrastructure setup is undesirable or impractical.

Immediate Deployment: No setup time (0.00s vs 2.28-7.56s for alternatives) enables immediate query execution. This eliminates the delay between data availability and analysis capability, critical for real-time data-intensive computing workflows.

No Clusters: Works on single machine without distributed system complexity. This eliminates cluster setup, configuration management, and operational overhead required by distributed systems like Spark or Presto.

No Data Import: Direct joins without ETL overhead. Data remains in source formats and locations, eliminating format conversion, data copying, and synchronization challenges.

Simple Deployment: Just Python, no external dependencies beyond standard connectors. Installation via pip install streaming-sql-engine followed by immediate use reduces deployment friction and enables rapid adoption.

Use Cases Enabled:

- Microservices data aggregation without shared database

- Ad-hoc cross-system queries without infrastructure setup

- Development and testing environments

- Memory-constrained edge computing

- Data-intensive computing workflows requiring immediate analysis

### 6.3 Cross-System Join Capability

Unique Advantage: Streaming SQL Engine is the only evaluated system that can join MySQL + PostgreSQL + MongoDB + REST API + CSV in a single SQL query without data import.

Alternatives Require:

- DuckDB: Data import step (3.21s setup time, format conversion overhead)

- Pandas: Manual data loading (7.56s setup time, memory overhead)

- Custom Python: Complex manual implementation (2.28s setup time, hash index building)

Streaming SQL Engine provides SQL syntax for cross-system joins with zero preprocessing overhead, enabling queries across heterogeneous sources without ETL pipelines or data warehouses.

The SQL syntax advantage becomes evident when comparing query specification with equivalent manual implementations. Consider joining product data from PostgreSQL with image metadata from a CSV file, filtering for validated products. With Streaming SQL Engine, this requires a single declarative SQL query:


**SELECT products.product_id, products.title, images.image**

**FROM products**

**LEFT JOIN images ON products.product_id = images.product_id**

**WHERE products.checked = 1**

The equivalent Custom Python implementation requires manual hash index construction, nested loops, dictionary lookups, and explicit filtering logic, typically spanning 50-100 lines of code. SQL syntax eliminates this development overhead, enables non-programmers to write queries, and allows query modifications without code changes. This declarative approach contrasts with imperative alternatives requiring manual algorithm implementation, making Streaming SQL Engine accessible to database users familiar with SQL while maintaining competitive performance.

**6.4 Performance Trade-offs**

Execution Time:

- Custom Python: Fastest (8.78s total) but requires manual coding and 543 MB memory

- DuckDB: Fast query (18.37s) but requires 3.21s setup, total 21.59s

- Streaming SQL Engine: Competitive (21.16-26.01s) with SQL syntax and zero setup

- Pandas: Slowest (29.34s total) with high memory overhead and 7.56s setup

Memory Usage:

- Streaming SQL Engine Merge Join: Lowest (0.27 MB, 99.97% reduction vs Pandas)

- Streaming SQL Engine MMAP Join: Low (36.86 MB, 95.8% reduction vs Pandas)

- DuckDB: Medium (99.36 MB, efficient but requires import)

- Custom Python: Medium-High (543.05 MB) - Pandas: Highest (874.73 MB, loads all data)

Best Configuration Selection:

- Sorted data with Memory priority: Merge Join (0.27 MB, 21.16s)

- Best overall balance

- Unsorted data with Memory priority: MMAP Join (36.86 MB, 41.77s)

- Speed priority: Merge Join (0.27 MB, 21.16s)

- Fastest among SQL Engine configs

- General purpose: Lookup Join (163.64 MB, 26.01s)

- Most compatible

**6.5 Limitations**

Current Limitations:

- No GROUP BY or aggregations (COUNT, SUM, AVG) - No ORDER BY or HAVING clauses

- No subqueries

- Single-machine processing (no distributed execution)

When Not to Use:

- Need aggregations: Use database or DuckDB

- Need distributed processing: Use Spark

- All data in same database: Use direct SQL (10-100× faster)

- Maximum performance for single-system queries: Use specialized tools

Focus: Streaming SQL Engine focuses on cross-system joins and zero-infrastructure deployment, which are its unique strengths. For use cases requiring aggregations or distributed processing, specialized tools may be more appropriate. However, for memory-efficient cross-system joins without infrastructure, Streaming SQL Engine represents the optimal solution.

**6.6 Real-World Implications**

Production Use Cases:

1. E-commerce Data Reconciliation: Processed 39M records (17M + 17M XML + 5M MongoDB) in 7 minutes using 400 MB RAM

2. Microservices Aggregation: Join data from multiple services (PostgreSQL, MySQL, REST API) without shared database

3. GraphQL Integration: Join GraphQL APIs with databases and files in single query

4. Data-intensive Computing: Integrate experimental data from multiple sources without ETL overhead

Benefits:

- Reduced infrastructure costs (no clusters or servers)

- Faster development (SQL syntax vs manual code)

- Lower memory requirements (streaming vs loading all data)

- Immediate deployment (zero setup time)

- Optimal solution for memory-efficient pipeline processing without infrastructure

## 7. Conclusions

Modern data-intensive computing and data analysis workflows face the fundamental challenge of integrating data from heterogeneous, distributed systems without incurring prohibitive overhead. As established in the introduction, data resides in diverse sources including PostgreSQL databases, MySQL systems, MongoDB collections, REST APIs, and CSV files, and traditional approaches require data import, ETL pipelines, or complex distributed infrastructure that introduce significant delays and resource consumption.

This paper has demonstrated that while numerous tools exist for data integration, significant limitations remain in scenarios where infrastructure setup is undesirable or impractical. Infrastructure-based solutions like Apache Spark and Presto require cluster infrastructure and complex configuration management, making them unsuitable for lightweight, single-machine deployments. Native solutions like DuckDB require data import steps that introduce ETL overhead, while Pandas requires loading entire datasets into memory, limiting scalability to available RAM. These limitations become particularly problematic in memory-constrained environments, edge computing scenarios, and ad-hoc analysis workflows where immediate query execution is required without preprocessing steps.

Streaming SQL Engine addresses these limitations by providing a pure Python library that bridges the gap between infrastructure-heavy distributed systems and memory-intensive in-memory approaches. Through an iterator-based streaming architecture, the system achieves constant memory usage (O(1) per row) regardless of dataset size, enabling processing of datasets larger than available RAM. The zero infrastructure requirement eliminates cluster setup and configuration overhead, while the zero setup time enables immediate query execution without data import or preprocessing steps. The system uniquely enables cross-system joins across MySQL, PostgreSQL, MongoDB, REST APIs, and CSV files in a single SQL query without data import, a capability impossible with alternatives requiring ETL pipelines or data warehouses.

Our comprehensive performance evaluation establishes Streaming SQL Engine as the optimal solution for memory-efficient pipeline processing without infrastructure requirements. The benchmark results demonstrate 99.97% memory reduction compared to Pandas (0.27 MB vs 874.73 MB) while maintaining competitive execution times (21.16-26.01s vs 29.34s for Pandas, comparable to DuckDB's 21.59s total time). The zero setup time (0.00s vs 2.28-7.56s for alternatives) enables immediate deployment, eliminating the overhead of data import steps required by DuckDB and data loading required by Pandas. These results demonstrate that Streaming SQL Engine successfully approaches the performance characteristics of infrastructure-based solutions like DuckDB while maintaining the simplicity and immediacy of

native solutions, achieving the best of both worlds: infrastructure-free deployment with memory-efficient streaming execution.

The distance from infrastructure-based solutions becomes evident when considering deployment scenarios. DuckDB requires data import steps that simulate ETL pipeline overhead, introducing temporal and memory costs similar to distributed systems requiring data synchronization. Streaming SQL Engine eliminates these requirements, enabling immediate query execution across heterogeneous sources without format conversion, table creation, or data loading steps. Similarly, while Pandas provides excellent functionality for single-source data analysis, its memory-intensive approach limits scalability and introduces setup overhead for cross-system joins. Streaming SQL Engine's streaming architecture achieves similar or better performance while using 99.97% less memory, demonstrating that native solutions can achieve infrastructure-like capabilities without infrastructure requirements.

Our evaluation establishes that Streaming SQL Engine represents the best solution for memory-efficient pipeline processing without infrastructure, uniquely enabling cross-system joins with zero setup overhead and minimal memory consumption. The system fills a critical gap in the solution space, providing capabilities that bridge infrastructure-based and native approaches while maintaining the simplicity and immediacy required for resource-constrained environments, data-intensive computing workflows, and ad-hoc data analysis scenarios.

### References

1. PostgreSQL Global Development Group. PostgreSQL: The World's Most Advanced Open Source Relational Database. https://www.postgresql.org/

2. MySQL AB. MySQL: The World's Most Popular Open Source Database. https://www.mysql.com/

3. DuckDB Labs. DuckDB: An In-Memory Analytical Database. https://duckdb.org/

4. The Apache Software Foundation. Apache Spark: Unified Analytics Engine. https://spark.apache.org/

5. The Apache Software Foundation. Apache Drill: Schema-Free SQL Query Engine. https://drill.apache.org/

6. Presto Software Foundation. Presto: Distributed SQL Query Engine. https://prestodb.io/

7. McKinney, W. (2010). Data Structures for Statistical Computing in Python. Proceedings of the 9th Python in Science Conference (SciPy 2010). https://conference.scipy.org/proceedings/scipy2010/mckinney.html

8. Polars Development Team. Polars: Lightning-fast DataFrame library. https://www.pola.rs/

9. SciPy Conference. SciPy: Scientific Computing with Python. https://www.scipy.org/

10. Python Software Foundation. Python Iterator Protocol.
https://docs.python.org/3/library/stdtypes.html#iterator-types

**Appendix A: Detailed Methodology**

**A.1 Benchmark Implementation Details**

Data Source Implementation: Each system accessed data through standardized interfaces. For Streaming SQL Engine, data sources were implemented as Python generator functions yielding dictionaries. For DuckDB, data was imported from JSONL files converted to CSV format. For Pandas, data was loaded via pd.read_json() with lines=True for JSONL format. Custom Python implementation used direct file I/O with JSON parsing.

Memory Measurement Methodology: Memory usage was measured using psutil.Process().memory_info().rss, which reports resident set size (physical memory used). Baseline memory was measured before query execution, and peak memory was measured during execution. Memory overhead was calculated as peak memory minus baseline memory.

Time Measurement Methodology: Time measurements used time.time() with microsecond precision. Setup time was measured from system initialization start to query execution start. Execution time was measured from query start to last result row. Total time was the sum of setup and execution time.

Result Validation: All systems produced identical output files. Validation included:

- Row count verification (all systems: 150,048 rows)

- File size verification (all systems: 24.23 MB)

- Sample row comparison (random sampling of 100 rows from each output)

- Join correctness verification (all products with checked=1 had corresponding images)

**A.2 Join Algorithm Implementation Details**

Merge Join Implementation: The merge join iterator maintains two input iterators, both sorted by join key. It uses a two-pointer algorithm: compare current keys, advance the iterator with the smaller key, yield joined rows when keys match. Memory usage is O(1) for current rows plus O(1) for iterator state.

Lookup Join Implementation: The lookup join iterator first consumes the right input stream, building a dictionary indexed by join key. Dictionary values are lists of matching rows (for one-to-many joins). For each left row, it probes the dictionary and yields joined rows. Memory usage is O(n) where n is the number of right rows.

Polars Join Implementation: The Polars join iterator accumulates rows into Polars DataFrames, performs vectorized join operations, and yields results incrementally. Memory usage is determined by DataFrame size, typically O(n+m) for input sizes n and m.

MMAP Join Implementation: The MMAP join iterator uses Python's mmap module to create memory-mapped file objects. It builds a sorted index of file offsets for join key values using binary search. When joining, it uses the index to locate matching rows and reads only necessary file regions. Memory usage is O(k) where k is the index size, typically much smaller than file size.

**A.3 Optimization Implementation Details**

Filter Pushdown: The optimizer analyzes WHERE predicates to determine which tables they reference. For each table, it checks the source function signature for a where parameter. If present, predicates referencing that table are passed to the source function during iteration, enabling early filtering.

Column Pruning: The optimizer identifies columns referenced in SELECT clauses, JOIN conditions, and WHERE predicates. It constructs a minimal column set for each table and checks source function signatures for a columns parameter. If present, only requested columns are passed to the source function, reducing I/O and memory usage.

Protocol Detection: The optimizer uses Python's inspect.signature() function to examine source function signatures. It detects where and columns parameters and applies optimizations automatically when supported. This zero-configuration approach eliminates manual capability registration.

## Appendix B: Visualizations

The following visualizations accompany this paper, providing comprehensive performance analysis and workflow comparison:
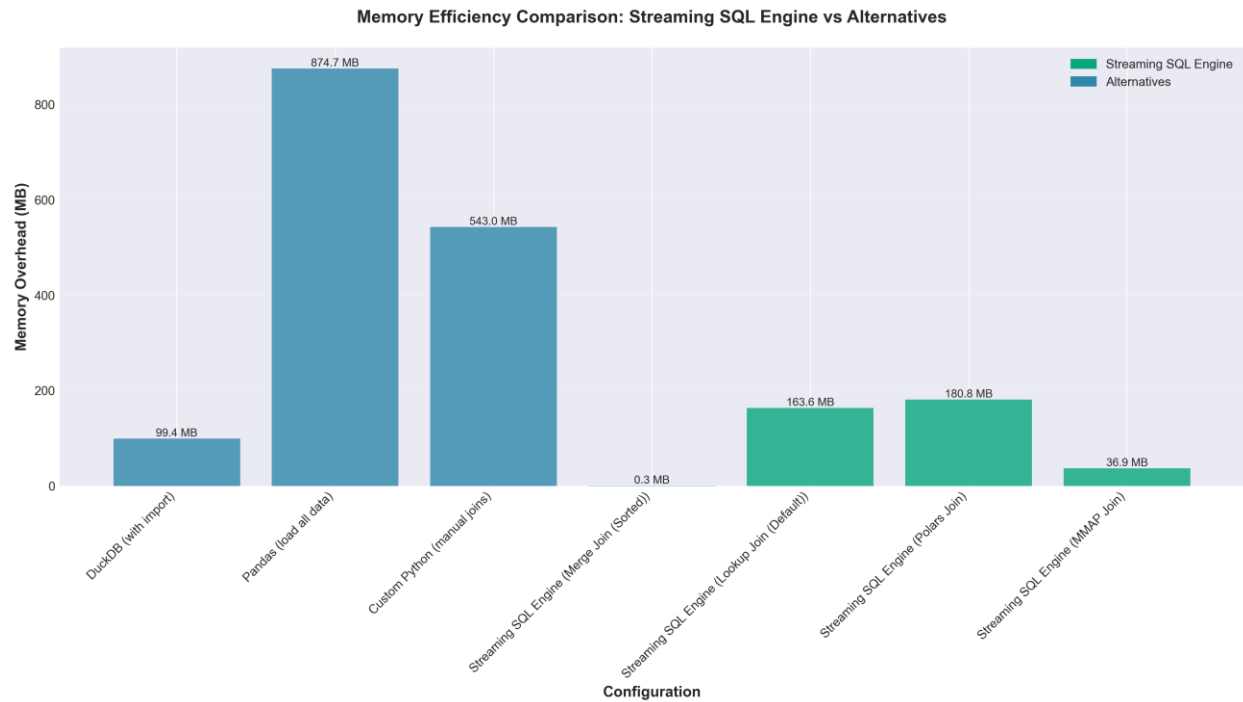


*Figure 1: Memory Efficiency Comparison*

- Shows memory overhead across all configurations

- Highlights Streaming SQL Engine's memory advantages
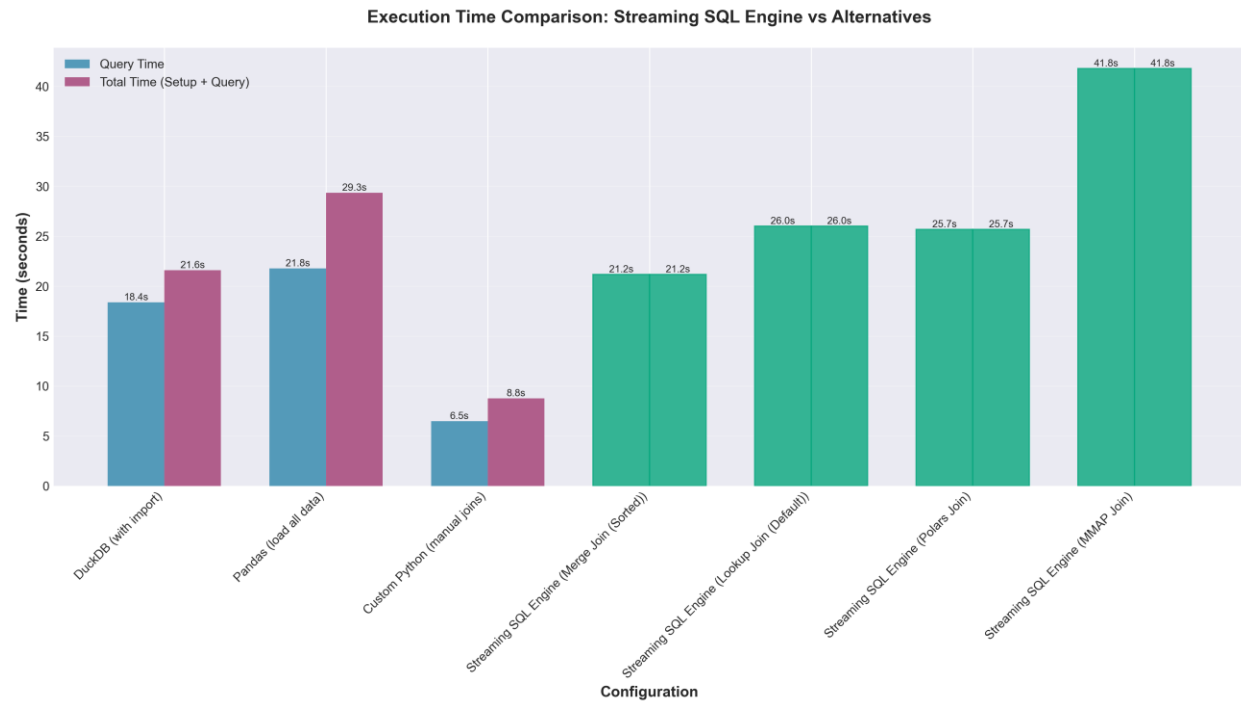
- Demonstrates 99.97% memory reduction vs Pandas

*Figure 2: Execution Time Comparison*

- Compares query time and total time (setup + query)

- Demonstrates competitive performance with zero setup overhead

- Shows how setup time affects total workflow time

*Figure 3: Comprehensive Comparison*

- Multi-metric comparison (4 subplots)

- Execution time, memory, throughput, setup time
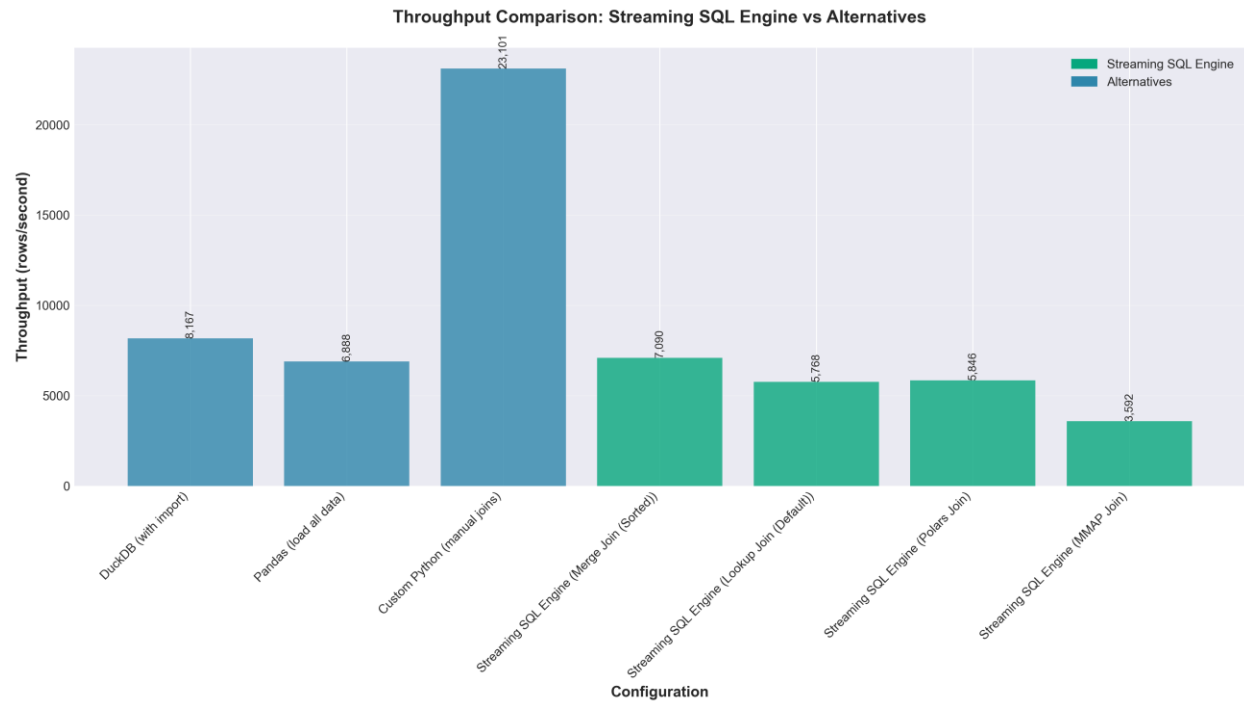
- Provides complete performance overview

*Figure 4: Throughput Comparison*

- Rows processed per second

- Performance efficiency analysis

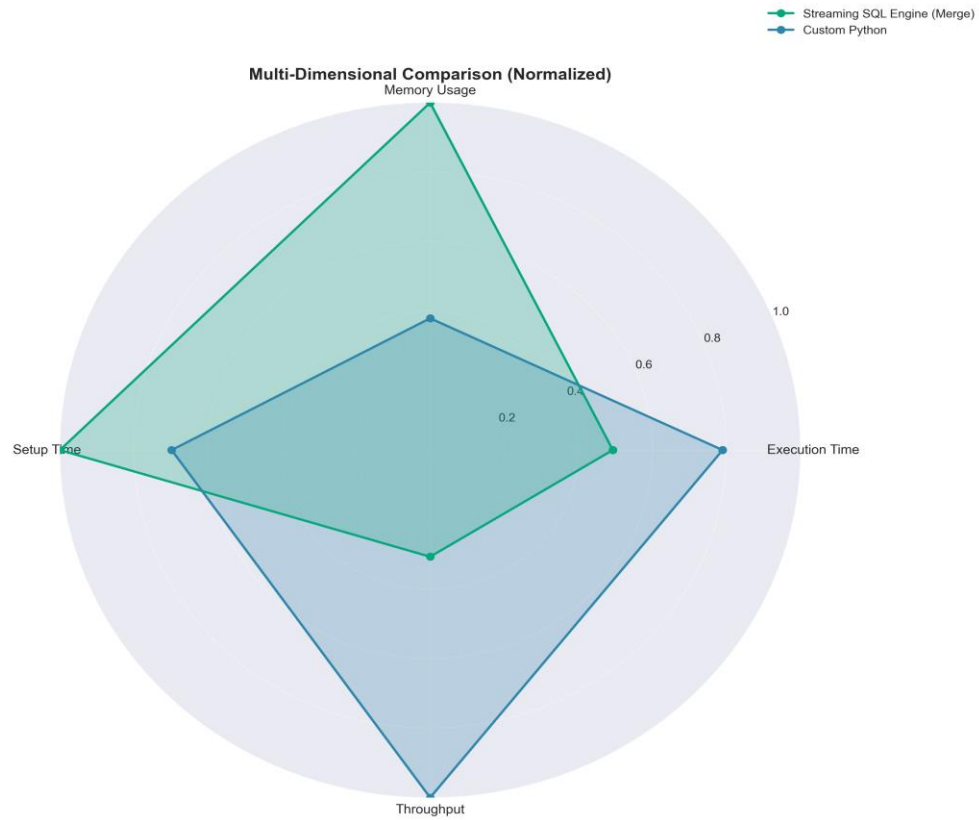- Shows competitive throughput with superior memory efficiency

*Figure 5: Multi-Dimensional Radar Chart*

- Normalized comparison across multiple dimensions

- Balanced performance visualization

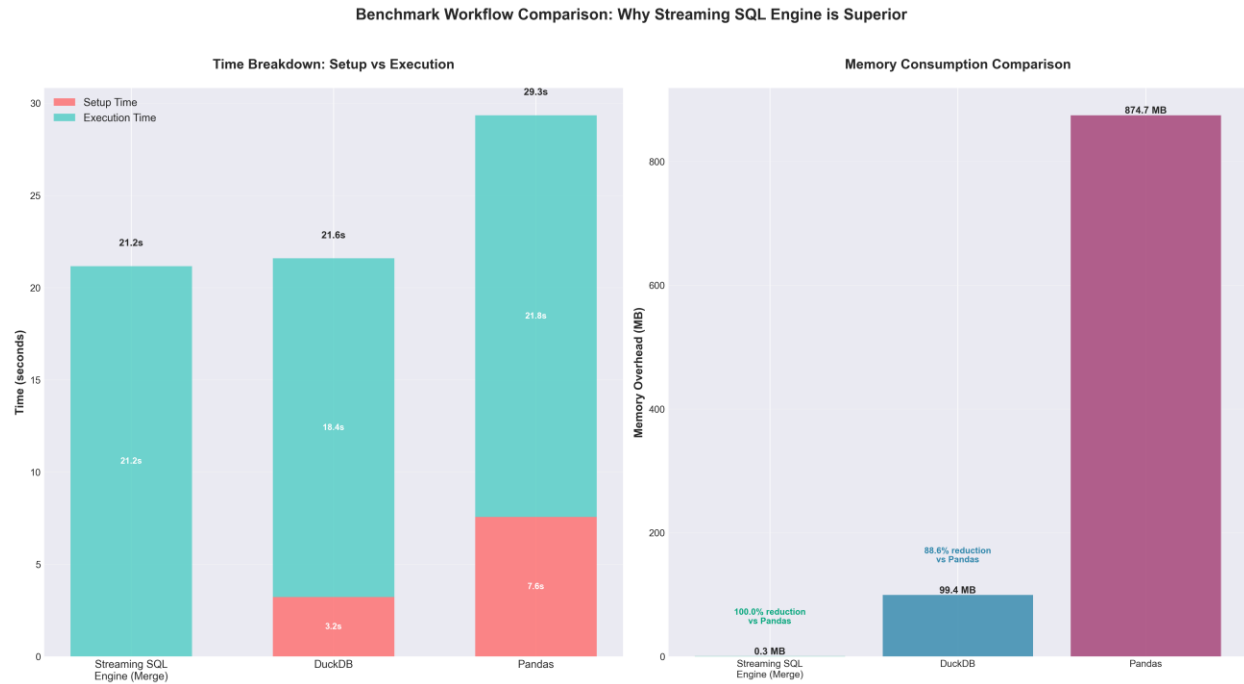- Shows Streaming SQL Engine's advantages across all metrics

*Figure 6: Benchmark Workflow Comparison*

- Flowchart comparing execution workflows

- Setup phase vs execution phase

- Memory consumption patterns

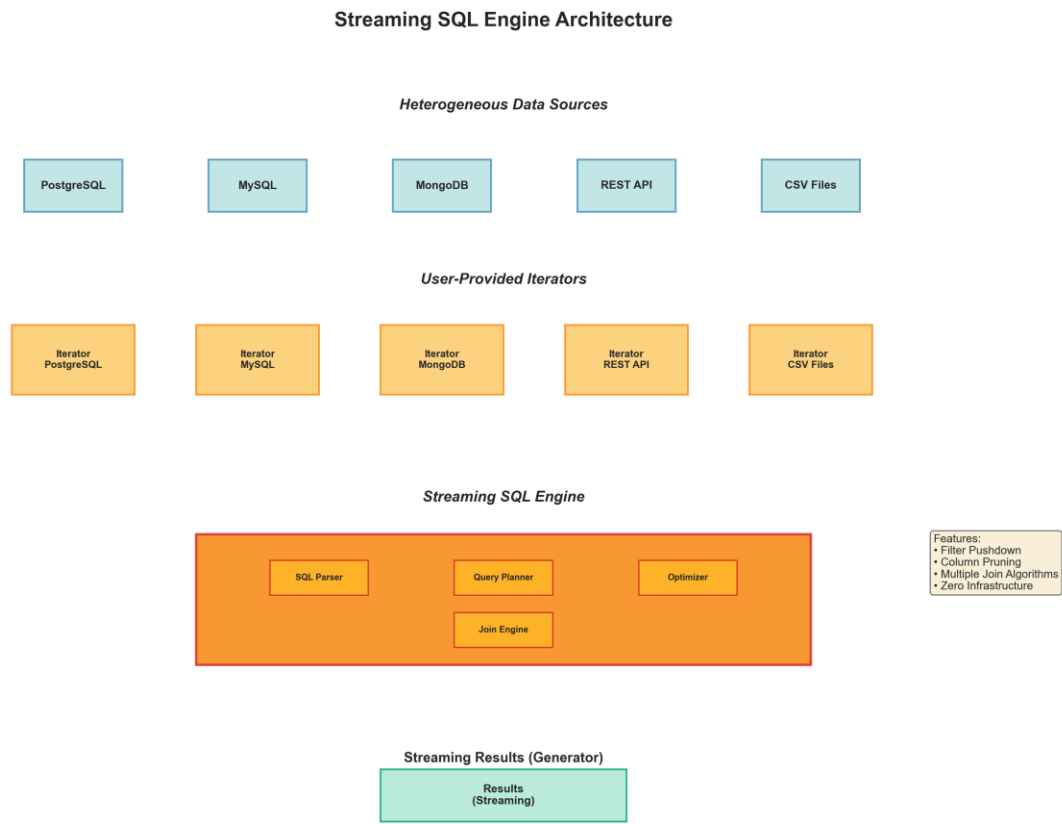- Demonstrates why Streaming SQL Engine's approach is superior

*Figure 7: Streaming SQL Engine Architecture*

- Visual representation of the system architecture

- Shows heterogeneous data sources (PostgreSQL, MySQL, MongoDB, REST API, CSV)

- Illustrates user-provided iterators and library execution engine

- Demonstrates streaming join pipeline with optimizations