



# Materialize

WHITEPAPER

## Flink vs Materialize: Comprehensive Enterprise Platform Analysis

What companies do varies incredibly widely. When it comes to their data, though, what companies *need* is remarkably similar:

- **Consistency:** Both operational and business data must be trustworthy and consistent.
- **Agility:** Teams need to minimize the time it takes to build and modify real-time data pipelines (and the transformations that happen within them).
- **Data democratization:** Most companies need to create and compose data assets in a self-service way and apply them for everything from tracking fast-moving inventory to providing fresh context for AI/ML.
- **Few resource constraints:** Companies aren't in the data business. They need data tools that don't require dedicated technical expertise.

Engineering teams we've talked to tell us they are spending too much time and money on the complicated, undifferentiated work of building, configuring, and maintaining a set of services in and around integrating, transforming, and delivering fresh data. One common way we see people solve this problem is by adopting stream processors like Apache Flink, which still require external support services like storage, orchestration, and some sort of serving layer.

### Materialize vs Flink: Key decision factors

- **Development velocity:** Enterprise customers deploying Materialize in production report approximately 50% faster deployment cycles than their previous Flink installs.
- **Cost efficiency:** These same customers say that deploying Materialize is 45-50% of the cost they experienced deploying a stream processor like Flink
- **Data consistency:** Materialize provides strict serializability vs. Flink's eventual consistency
- **Team accessibility:** Standard Postgres SQL-only (Materialize) vs. JVM programming and the FlinkSQL dialect (Flink)
- **Consistency:** Guaranteed global consistency through strict serializability (Materialize) vs. eventual consistency through "exactly-once state semantics" (Flink)
- **Composability for higher order data products:** Materialize's strict consistency model makes Materialize views safe to compose, cache, and expose to applications or agents vs Flink's fundamental timing inconsistencies break system composability.

## Introduction

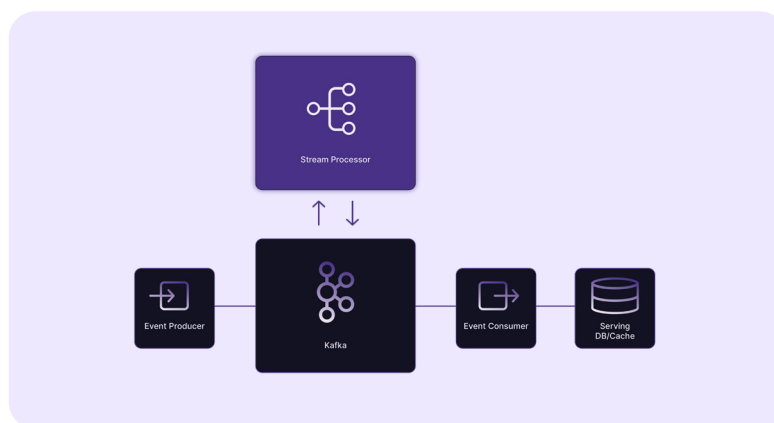
Apache Flink is a distributed stream processing framework for stateful computations over unbounded and bounded data streams. Flink specializes in stream processing, but it relies on other datastores for serving real-time data, which may not be optimized for that purpose.

Materialize is a SQL-based real-time data integration and transformation platform that combines cloud storage, incremental view maintenance, and a built-in, Postgres-compatible serving layer. This design enables streaming data processing within a single framework, reducing complexity and cost.

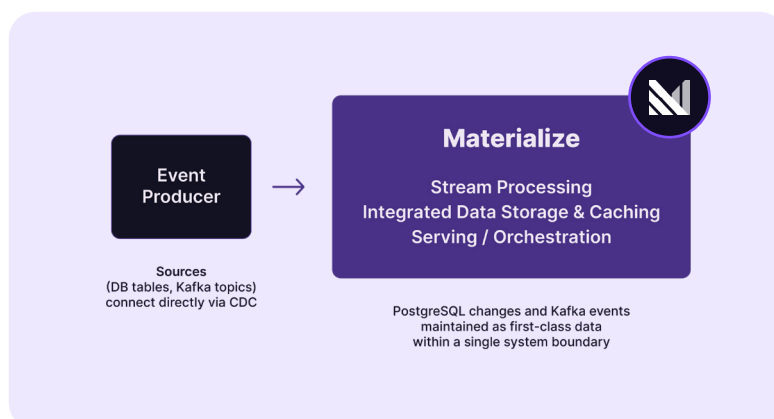
## Not quite an apples-to-apples comparison

Before we dig into the gritty details of each platform, it's important to understand that when comparing Flink and Materialize we are comparing a stream processor and a real-time data integration platform:

- **Stream processors** handle the core data transformation, but they also rely on Kafka for intermediate storage, and some type of serving layer (i.e. a separate database, or a Redis-like tool). In practice, stream processors work in a context of multiple supporting services:



- **Materialize** is one tier up in the software abstraction stack. Materialize combines the compute, orchestration, and load balancing into one platform, so users see a unified real-time data integration and transformation platform that presents as a SQL database.



The functional difference between Flink and Materialize is fundamental:

**Flink** delivers raw stream processing. It functions as the compute layer, and requires external services to materialize data views: a message bus or CDC connector (e.g., Kafka, Debezium) as the input layer, and an output system (e.g., Kafka, Elasticsearch, a database) to persist results.

**Materialize's** database abstraction means Materialize doesn't deliver raw stream processing. At the heart of Materialize, [Timely Dataflow](#), is mature and stable open source software originally developed at Microsoft Research starting more than a decade ago. Materialize's platform builds the database abstractions and distributed architecture around this core stream processing library.

**For the user, Materialize presents as a Postgres wire-compatible data streaming platform with incrementally — and continually — updated materialized views. This means you can store raw data, transform, and serve all within a single system.**

## Architecture and system integration

Materialize hides the complexity of a stream processor behind a database abstraction — a fundamental shift in real-time data processing — because we think that this is the best abstraction to serve the broad majority of operational data use cases. Use cases that your engineering teams are likely grappling with right now, because so far they are not served well by traditional databases (wrong computation paradigm). Stream processors like Flink represent a step in the right direction, but they create complexity by requiring a system of services.

### Apache Flink: Capabilities and architectural constraints

Apache Flink operates as a distributed stream processing engine with a clear architectural boundary: it processes data in motion without persistence or native query capabilities.

Understanding Flink's design reveals both its strengths and the infrastructure requirements it imposes. A production Flink deployment follows a three-tier architecture.

- 1. Input Layer:** Data ingestion occurs through message buses like Kafka or change data capture systems such as Debezium. These systems provide the streaming data interface, handling backpressure, partitioning, and initial fault tolerance.
- 2. Compute Layer:** Flink job clusters execute the core stream processing logic. The engine excels at complex event processing, windowing operations, and stateful computations with exactly-once processing guarantees. However, Flink only maintains any transient state required for computation and immediately forwards results downstream.
- 3. Output Layer:** Processed results must be materialized in external systems — additional Kafka topics, Elasticsearch clusters, or traditional databases. This externalization is mandatory since Flink provides no native storage or query interface.



This multi-tier architecture creates several operational implications:

- **Chaining jobs is complex.** Multi-stage processing pipelines require external coordination through intermediate storage systems. Each stage boundary introduces network I/O, serialization overhead, and additional failure points. Pipeline orchestration becomes a distributed systems problem requiring careful attention to backpressure propagation and failure recovery across system boundaries.
- **There's no native support for querying.** Consuming processed data from Flink for business intelligence requires additional infrastructure investment. Real-time dashboards, operational APIs, and ad-hoc analytics require separate query-serving systems that read from Flink's outputs. This typically involves data warehouses, search indices, or purpose-built serving layers.
- **Operational boundaries are unclear.** System availability depends on the coordinated operation of multiple independent components. Each system in the pipeline maintains its own durability guarantees, backup procedures, and failure modes. Operational teams must understand and manage the combined reliability characteristics of the entire stack.

#### THE UPSHOT:

Flink provides sophisticated stream processing primitives, but assembling a robust, composable, and queryable real-time data system with Flink can be complex and cumbersome because it requires managing multiple moving parts.

## The real-time data integration approach: Materialize

Materialize consolidates traditional stream processing architecture by integrating ingestion, computation, storage, and serving all within a unified system boundary.

The platform's separated storage, compute, and serving planes all operate and scale independently, so you can ingest arbitrary volumes of data to elastic storage (for example: S3), then spin up unlimited numbers of compute instances to read from, transform, and write this data back — and you can serve results to as many concurrent connections as you like.

- **Unified data model:** Sources — Kafka topics, PostgreSQL tables, MySQL databases — connect directly through change data capture protocols. But unlike Flink's ephemeral processing model, this data is persisted durably within Materialize itself. PostgreSQL changes and Kafka events are maintained as first-class data within a single system boundary — eliminating the overhead and complexity of coordinating external services for data availability.



- **Incremental materialization:** In Materialize, the system maintains materialized views, live reports that update continuously as source data changes. This differs fundamentally from traditional batch materialization because Materialize updates the current view with each new event, but with much greater efficiency — only processing any incremental changes rather than doing a full recomputation.
- **[Virtual time](#) ensures global consistency:** Virtual time is a technique for distributed systems where events are timestamped prescriptively rather than descriptively. [Materialize uses virtual time](#) to efficiently compute, maintain, and return the specific correct answers at specific virtual times by recording, transforming, and reporting continually evolving, explicitly timestamped histories of collections of data. These explicit histories eliminate any potential concurrency by promptly and unambiguously communicating the exact contents of a collection at each of an ever-growing set of times.
- **Native query interface:** Materialized views support direct SQL querying through standard database protocols (PostgreSQL wire protocol compatibility). This enables immediate consumption by existing business intelligence tools, application frameworks, and analytical workloads without additional serving infrastructure.

## Key architectural advantages

Materialize's real-time model consolidates multiple operational concerns within a single system boundary, creating distinct advantages and constraints:

- **Operational simplification:** Materialize's single system boundary means that platform teams manage just one system instead of needing to coordinate external state stores or intermediate topics.
- **Incremental computation with durability:** Inputs are ingested once and retained as needed for recomputation or recovery.
- **Direct interactivity:** Queries can be executed directly against maintained views. This enables use cases like operational dashboards, AI feature stores, or real-time APIs with no additional infrastructure.

### THE UPSHOT:

Materialize's real-time data integration platform **minimizes architectural sprawl, reduces operational cost, and allows teams to focus on business logic** instead of pipeline plumbing.

## Consistency & composability

### Flink's consistency model: Guarantees vs. observable behavior

Flink's "exactly-once state semantics" represents one of the most frequently cited — and misunderstood — guarantees in stream processing. The core limitation lies in when outputs become visible and how they align across jobs.

- **Internal consistency mechanics:** Flink achieves globally consistent state snapshots within individual jobs through its distributed snapshot algorithm, derived from the Chandy-Lamport protocol. This ensures that internal state across all parallel operators reflects a coherent point in the data stream at checkpoint time. The algorithm is sophisticated and works reliably for Flink's internal bookkeeping.
- **Compromised output visibility:** The challenge emerges when considering how and when processed results become observable downstream. Flink sinks emit data as it flows through the processing topology, operating independently of the snapshot boundaries that govern internal consistency. And multiple Flink applications that process the same inputs may do so at different times and produce results in differing orders. This fundamental asynchrony creates a disconnect between internal state consistency and external output visibility.
- **Transaction boundary loss:** Consider a common scenario where upstream database transactions contain multiple related changes—perhaps updating both an account balance and transaction log simultaneously. When Flink processes CDC events from this database, it cannot guarantee that these related changes appear together in downstream systems. The processing topology may emit one change before the other, breaking the original transactional relationship.

Operationally, this results in eventual consistency — outputs appear when convenient, not when correct, as researcher [Jamii Brandon](#) succinctly puts it. The implications extend beyond individual jobs: downstream consumers cannot safely compare or join outputs from separate Flink jobs, and they can't infer or determine consistency boundaries through observation alone.

**Composability breakdown:** These fundamental timing inconsistencies break system composability. Building higher-order data products means combining multiple Flink job outputs — creating systems that quickly become difficult to reason about. Teams often discover correctness issues only after deploying complex multi-job pipelines, leading to expensive architectural remediation or acceptance of data inconsistency.

While Flink does support batch / adhoc queries, it is really not built to easily go back and forth — i.e., you can't create a streaming view and then easily query it to test results. This is because Flink doesn't include storage so you always have to connect to external services, even if you are just chaining Flink jobs. Also, Flink does not have drivers. You aren't going to query Flink from a programming language — you are going to sink results into some third-party system and then read the data using that system's integrations.



## Materialize's consistency architecture: Strict serializability

Materialize handles consistent with a fundamentally different architectural approach: strict serializability across all system components with a global logical timeline:

- **Global timeline model:** Every batch of input data receives a logical timestamp within a unified global timeline. This timeline serves as the consistency foundation for all downstream processing, ensuring that computations across different views or queries reflect identical points in logical time.
- **Synchronized view updates:** All materialized views update atomically at single logical timestamps. When source data changes, related views progress together through the timeline, eliminating the asynchrony that creates consistency gaps in traditional stream processing. Applications reading across multiple views observe a coherent snapshot of the data state.
- **Transactional integrity preservation:** When ingesting data through change data capture, Materialize maintains the transactional boundaries from source systems. Multiple changes originating from the same database transaction become visible simultaneously downstream, preserving the semantic relationships that upstream applications depend upon.
- **Correctness before visibility:** Results become observable only after achieving correctness at the designated logical timestamp. This approach prevents the partial result visibility that can lead to incorrect business decisions or downstream processing errors.

**Fundamental composability:** Materialize's strict consistency model makes Materialize views safe to compose, cache, and expose to applications or agents — and it aligns with the requirements of modern data products that must behave predictably and support correctness guarantees.

Flink, on the other hand, offers eventual consistency that fundamentally undermines composability. Flink-based systems typically require additional application logic to handle consistency gaps, not to mention operational reasoning that is more complex by an order of magnitude: Flink deployments require that you understand — and handle — the interaction between processing topology, checkpoint timing, and sink behavior.

### THE UPSHOT:

Materialize guarantees strict serializability: **global consistency across the entire database**. At any given moment, you can trust that no matter how you are slicing your data up, you can always see every view is going to show results as of the same inputs, as of the same result, and the numbers always tie out. And you can reason across these, you can join them, you can compose them, and the numbers will still make sense.





## Developer experience & data accessibility

### Flink's Programming Model: Power with Complexity

Flink's computational model represents a significant departure from traditional batch processing and relational database paradigms. While this approach unlocks sophisticated stream processing capabilities, it introduces substantial learning curves and operational overhead for development teams.

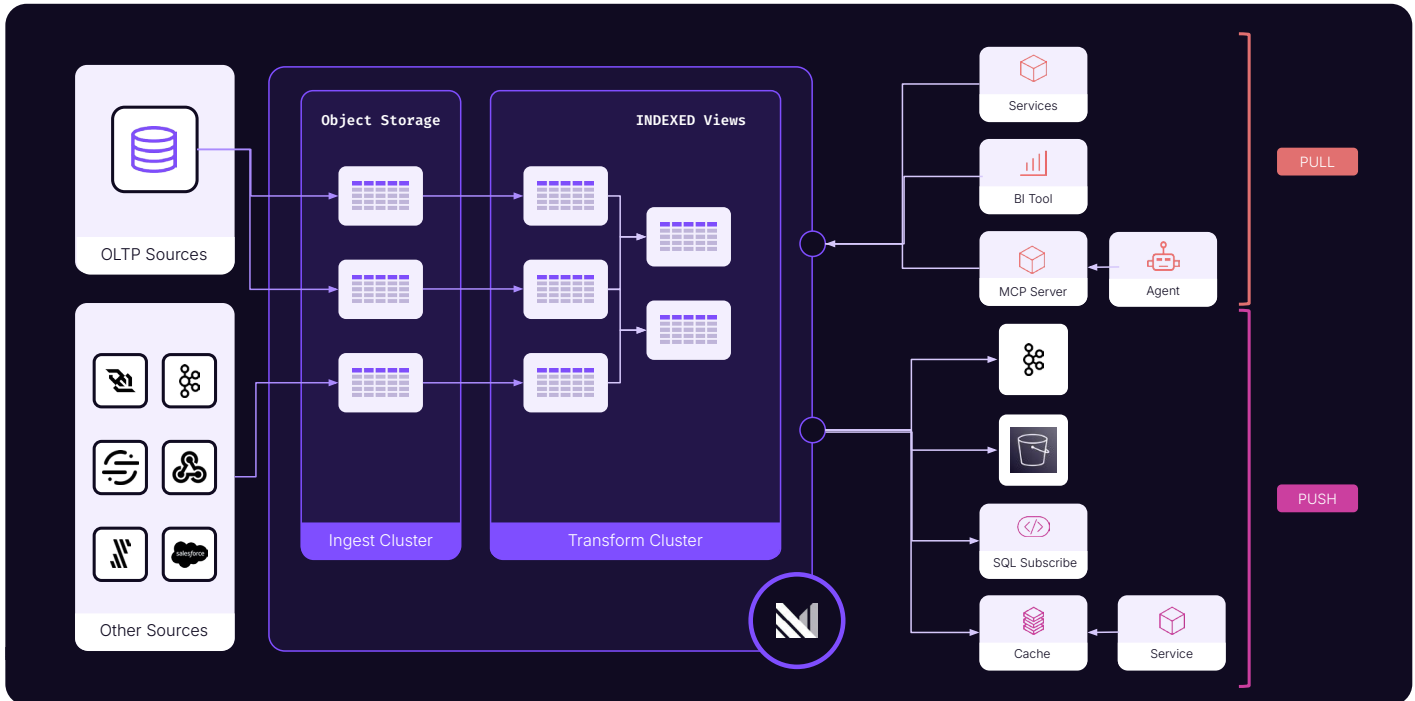
- **Custom programming:** Flink development centers on imperative programming through specialized APIs in Java, Scala, or Python. Developers must explicitly define transformation logic, manage operator chaining, and handle data serialization. This approach provides fine-grained control over processing behavior but requires teams to master streaming-specific programming patterns that differ fundamentally from familiar batch-oriented development.
- **Bespoke SQL dialect:** Flink relies on its own customized version of SQL, so developers must specifically learn FlinkSQL in order to run queries on streaming (and batch) datasets.
- **Temporal complexity:** Producing correct results in Flink requires deep understanding of event time semantics, watermark generation, and out-of-order data handling. Developers must reason about when computations should trigger, how late data affects results, and how to balance latency against completeness. These concepts have no direct analogues in traditional application development, creating significant knowledge barriers for teams transitioning to stream processing.
- **Complex configuration for state and fault tolerance:** Checkpointing behavior, state backend selection, and recovery strategies require explicit configuration and ongoing tuning. Developers must understand how state size affects performance, how checkpoint intervals impact latency, and how to optimize for their specific workload characteristics. These operational concerns become part of the development process rather than being abstracted away by the platform.
- **Testing and Debugging Challenges:** Flink's asynchronous processing model creates non-deterministic output timing that complicates testing strategies. Traditional unit testing approaches often prove insufficient, requiring developers to create complex test harnesses that account for timing variations. Debugging production issues becomes particularly challenging due to limited runtime introspection capabilities and the distributed nature of stream processing execution.

#### THE UPSHOT:

These characteristics make Flink potentially powerful for teams with dedicated streaming expertise. **But they also make Flink inaccessible:** the majority of teams out there don't have the significant training and operational resources needed to deploy Flink in a real-time data streaming environment.

## Materialize's SQL-centric approach: Familiar paradigms for Stream Processing

Materialize addresses the accessibility challenge by building on familiar relational database concepts while providing streaming capabilities, all behind a standard (and, for many, familiar) SQL interface.



- **Declarative programming model:** In Materialize, all transformation logic expresses through standard SQL syntax familiar to any developer with database experience. Teams define materialized views using `CREATE VIEW` statements, aggregate data with `GROUP BY` clauses, and join streams using familiar SQL `JOIN` syntax. This eliminates the need to learn specialized streaming programming models while maintaining the expressiveness required for complex real-time processing.
- **Internal incrementalization:** Materialize handles time semantics, incremental computation, and consistency management internally. Developers focus on defining what results they want rather than how the system should process changing data. Because Materialize is built on Timely Dataflow, the platform automatically manages consistency across complex view hierarchies without exposing concepts like watermarks or out of order concerns to application developers.

- **Deterministic query behavior:** Given identical logical timestamps, queries return identical results regardless of execution timing or resource allocation. This deterministic behavior simplifies testing, debugging, and reasoning about system correctness. Developers can validate logic using familiar SQL tools and techniques rather than building specialized streaming test frameworks.
- **Broad integration:** Standard PostgreSQL wire protocol compatibility enables immediate integration with existing development tools, business intelligence platforms, and application frameworks. Teams can use dbt for transformation logic management, connect Tableau for visualization, or integrate with web applications using standard database drivers. This compatibility eliminates the custom integration work typically required when adopting new stream processing platforms.

#### THE UPSHOT:

Materialize's stream-processor-wrapped-in-a-database approach offers a **drastically lower barrier to entry and faster time-to-value**. Teams can ship real-time features without hiring streaming specialists or re-architecting their systems.

## Enterprise implications

Flink is a flexible, general-purpose stream processor well-suited to complex event processing or bespoke transformation pipelines. However, it provides no consistent query interface, no composability guarantees, and no support for direct consumption of results in real-time applications.

Materialize, by contrast, offers strict serializability and composable views. It integrates ingestion, storage, and compute within a single system, using a standard SQL interface to offer broad accessibility for existing data teams.

These fundamentally divergent approaches to data stream processing have real-world consequences for enterprise adoption:

- **Skills & resources:** Materialize can be used by anyone with SQL competency, while Flink typically requires engineers who are very well versed with JVM platforms and the nuances of streaming and distributed systems.
- **Consistency guarantee:** Materialize provides strict serializability, ensuring that no matter how you slice your data, every view shows results as of the same inputs. Flink offers eventually-consistent data where you do not know if the output is correct at any given moment.

- **Cost structure:** Materialize is typically around half the cost of a stream processor like Flink or Spark when factoring in development, operational, and maintenance costs.
- **Business value:** One of Materialize's biggest values is enabling customers to access transactions in real time via a BI tool, while Flink requires additional systems for query access.

## Total cost of ownership

**Compared with Flink's standalone stream processor, Materialize offers lower overall TCO.**

Companies need fewer resources to work with Materialize's integrated data streaming and transformation platform. Because deploying Materialize is drastically less complicated, they can also work with less experienced engineers — it's much easier for organizations to spin up.

Teams that have moved to Materialize further report savings related to denormalization and solving cache invalidation. Some organizations have been able to reduce their data warehouse spend by offloading expensive queries to run through Materialize. Users also tell us that they are able to reduce secondary costs associated with pipeline development, updates, and maintenance. In terms of dollars and cents, Materialize typically lands at half the cost of a stream processor like Flink in terms of operating costs.

