

# IBM StreamSets

Build, deploy and manage real-time streaming data pipelines with an intuitive visual oriented design



## Highlights

Streaming data integration coupled with extensive connectivity

Drag-and-drop visual pipeline designer

Data security and privacy with IBM StreamSets

Full-Featured SDK for Python

In today's highly-competitive environment, organizations seeking a strategic advantage need to leverage all the information at their disposal. The pressure to respond to market volatility, customer demands and digital transformation initiatives necessitates the access to continuous, real-time data across the enterprise. To meet these needs, companies must streamline data integration, helping to enable innovation while maintaining control over data flows. This requires implementing resilient, adaptive streaming pipelines that can handle high-velocity data and unexpected changes.

Given the unique needs of enterprises and due to specific use cases, the IBM Data Fabric architecture is entirely fit-for-purpose. Customers can choose from a set of seamlessly integrated data integration products that fit their needs, whether they be for artificial intelligence (AI), business intelligence and analytics, or other industry-specific requirements.

IBM's Data Integration portfolio includes industry-leading tools such as IBM DataStage for moving and transforming mission-critical data with ETL/ELT processing, and IBM StreamSets is a strategic addition to this portfolio. StreamSets enables customers to build real-time streaming data pipelines, complementing the existing capabilities of IBM DataStage. With IBM Databand, the data observability solution for data pipeline monitoring and issue remediation underpinning the entire portfolio, StreamSets is a seamless and comprehensive solution for designing, deploying and managing data pipelines across all data sources and integration patterns.

With StreamSets organizations get an innovative visual-oriented approach to building real-time data pipelines to capture and stream data in real-time, regardless of its structure or complexity, allowing for faster responses to changing business conditions, more informed decisions and greater innovation. Workloads can run in the same physical location where data resides as StreamSets offers hybrid cloud support, integrating data across multiple cloud platforms and on-premises systems.



# IBM defines real-time data

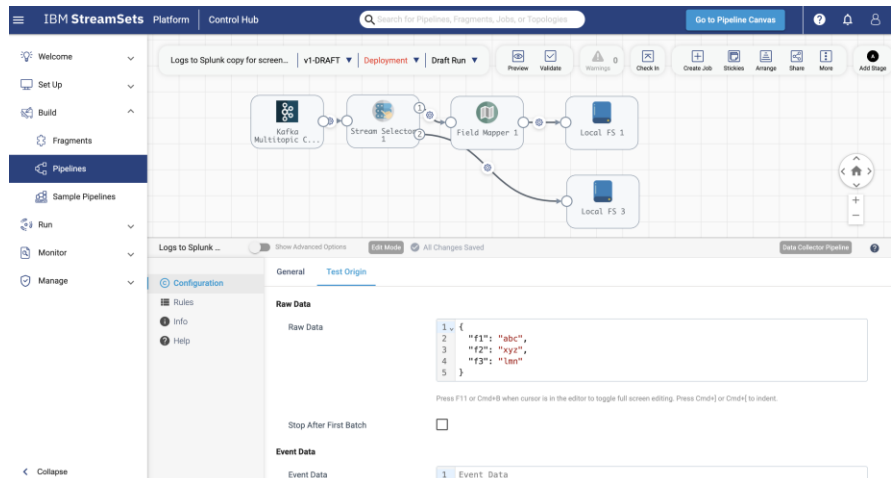
integration as the ability to ingest, process, and write data as soon as it's available, instead of on an intermittent or scheduled basis.

Use cases that benefit from real-time data integration are those for which the extraction of insights with minimal delay within **seconds** provides business value.

## Real-time data integration coupled with extensive connectivity

Real-time data streaming presents several challenges. One primary challenge is ensuring low latency, as delays can negatively impact user experience or decision-making processes. Scalability is another issue, as systems must handle varying loads efficiently without performance degradation. Data consistency and reliability are crucial, especially when dealing with distributed systems where data might arrive out of order or get lost. Managing the infrastructure costs associated with high-throughput and low-latency requirements can be significant. Additionally, ensuring data security and security during transmission and processing is essential to protect against breaches and comply with regulations. Real-time data streaming systems also need robust error handling and fault-tolerance mechanisms to maintain service continuity in the face of failures or network issues.

IBM StreamSets offers a comprehensive data integration solution, connecting with various streaming technologies like Apache Pulsar and Apache Kafka. StreamSets also excels in Change Data Capture (CDC), capturing real-time data changes from sources like Oracle, PostgreSQL, and MongoDB, and delivering them to downstream systems. The platform enhances these capabilities with smart data pipelines, allowing users to schedule and monitor jobs, which define the pipeline and execution engine for data processing. Users can perform transformations securely and efficiently within these pipelines. Furthermore, StreamSets introduces the concept of pipeline fragments, which are reusable sets of connected stages that can be incorporated into multiple pipelines, ensuring consistent processing logic across different data flows and simplifying pipeline design and maintenance.



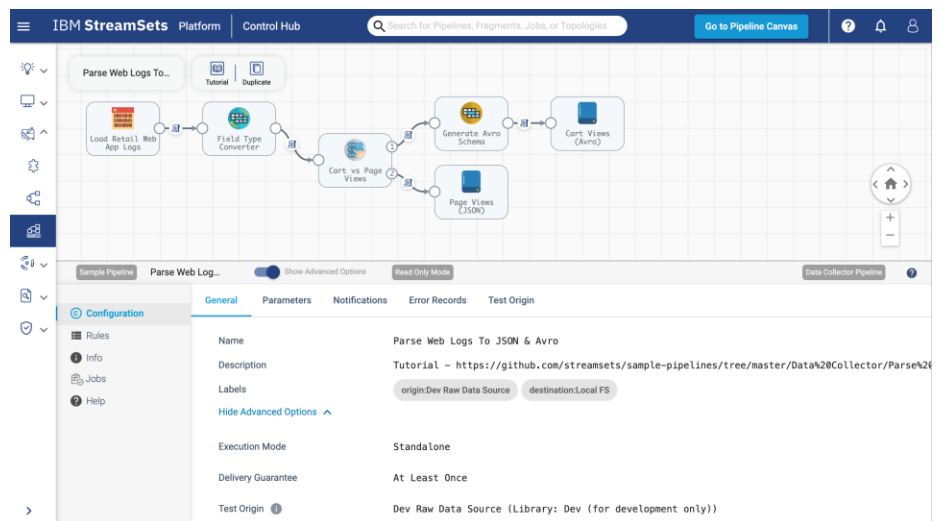
# 87%

organizations require data to be ingested and analyzed within 1 day or faster<sup>1</sup>

## Drag-and-drop visual pipeline designer

IBM StreamSets provides a comprehensive web-based user interface (UI) in which users can design, execute and monitor their pipelines. At its core is the pipeline canvas, a workspace dedicated to these tasks. Users can view and adjust the properties of a pipeline or its selected stages, with panels that can be resized, minimized, or maximized for convenience. The Preview panel shows data entering and exiting specific stages or groups of stages, along with stage properties and preview configurations, while the Monitor panel delivers real-time metrics and statistics for running pipelines. The UI also allows users to view a list of all available pipelines and related information, with filtering options such as Running Pipelines.

A graphical data integration methodology offers numerous business benefits, including ease of use and enhanced productivity, as it simplifies complex processes and makes them accessible to users with varying technical expertise. The intuitive drag-and-drop interfaces and visual workflows streamline the creation and management of data pipelines, speeding up development time and allowing for quicker adaptation to changing requirements. These tools also improve collaboration by facilitating better communication among team members and stakeholders, reducing errors through visual representations that make it easier to identify and correct issues. Additionally, GUI-based tools offer scalability to handle growing data volumes and complexity, enhanced monitoring and management capabilities for real-time visibility into data flows, and cost efficiency by lowering development and maintenance costs. They provide flexibility and integration with a wide range of data sources and existing systems, support compliance and governance needs, and empower users with self-service capabilities, reducing dependency on IT departments and fostering a culture of data-driven decision-making.





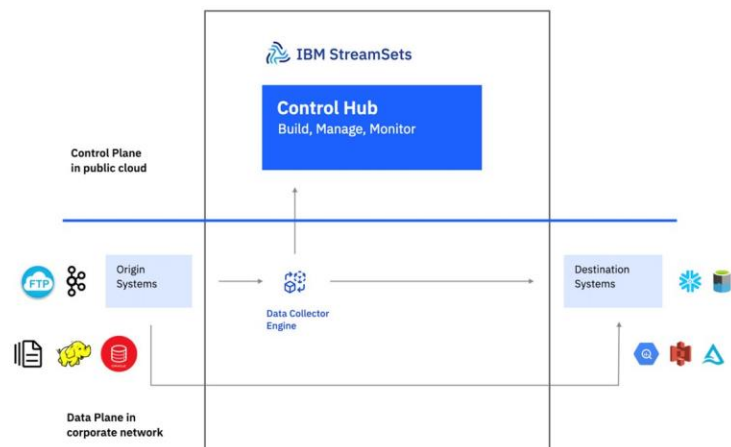
# StreamSets

provides powerful, flexible, and secure data integration solutions that cater to the diverse needs of modern enterprises

## Ensuring data security and privacy with IBM StreamSets

IBM StreamSets has a "secure by design" architecture in which a cloud-based control plane is completely separated from the customer-deployed data plane. The control plane, also known as Control Hub, is a cloud-native application that serves as the central management system for StreamSets' data integration platform. This powerful tool enables users to create, schedule, manage and monitor pipelines, engines and jobs from a single, unified interface. Designed specifically for cloud environments and managed by StreamSets, Control Hub offers a comprehensive solution for orchestrating data flows across an organization. In contrast, the Data Plane houses the actual ingestion and transformation engines, operating securely in either on-premise or cloud environments, providing flexibility to meet diverse infrastructure requirements. The StreamSets' service offering includes a unique feature that allows customers to snapshot pipeline data for configuration and debugging purposes. It's important to note that this feature is not enabled by default and is entirely under the customer's control. For organizations with stringent firewall policies, it offers the option to implement a direct connection between the user's browser and the data engines, further enhancing data security.

StreamSets' visibility into customer data is carefully delineated to balance functionality with privacy and security concerns. Organizations can access to three main categories of data: configuration data, usage data, and telemetry data. [Configuration data](#) encompasses information about the customer's assets, developed pipelines, and administrative settings such as user roles and permissions. This data is crucial for tailoring the instance to each customer's specific implementation. Given its sensitive nature, this data is securely maintained within the application. For customers using the SaaS version or DataOps platform, this data resides in a public cloud with robust security measures. Usage data, primarily in the form of [metering information](#), is collected for billing purposes and to help customers estimate costs for potential usage expansion. [Telemetry data](#) is gathered to support and improve StreamSets' services. This information helps the company prioritize feature development and enhancements based on customer usage patterns, such as identifying widely used integrations or stages. Additionally, telemetry data contributes to improving overall performance, efficiency, and uptime of stages, pipelines, and services. By collecting and analyzing this data.





### **For more information**

To learn more about IBM StreamSets, contact your IBM representative or IBM Business Partner, [give it a try](#) or [visit here to sign up for the latest updates](#).

1. Forrester, "The Future Of Data: Make It Fast", 2018
2. Jet Brains, "The State of Developer Ecosystem 2022", 2022

© Copyright IBM Corporation 2024  
IBM Corporation  
New Orchard Road  
Armonk, NY 10504

Produced in the  
United States of America  
August 2024

IBM, the IBM logo, is trademarks or registered trademarks of International Business Machines Corporation, in the United States and/or other countries. Other product and service names might be trademarks of IBM or other companies. A current list of IBM trademarks is available on [ibm.com/trademark](https://ibm.com/trademark).

This document is current as of the initial date of publication and may be changed by IBM at any time. Not all offerings are available in every country in which IBM operates.

THE INFORMATION IN THIS DOCUMENT IS PROVIDED "AS IS" WITHOUT ANY WARRANTY, EXPRESS OR IMPLIED, INCLUDING WITHOUT ANY WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND ANY WARRANTY OR CONDITION OF NON-INFRINGEMENT.

IBM products are warranted according to the terms and conditions of the agreements under which they are provided.

