



PySpark: Basics, Supported Options, and Missing Features

Overview of PySpark

What is PySpark? PySpark is the Python API for the Apache Spark engine. It lets you use familiar Python syntax while leveraging Spark's distributed computing capabilities to perform **real-time, large-scale data processing** across clusters ¹. The API provides an interactive shell and integrates with the Spark runtime to combine Python's ease-of-use with Spark's speed and scalability ¹.

Core components and features

- **Spark SQL & DataFrames:** PySpark supports all of Spark's structured-data features such as Spark SQL and DataFrames. DataFrames are collections of data organized into named columns that enable you to read, write and transform data using both Python code and SQL expressions on the same execution engine ². Because Spark automatically optimizes DataFrame queries, they are the recommended high-level API; the low-level RDD API is available but lacks automatic query optimization and is considered more difficult to use ³.
- **Structured Streaming:** PySpark includes **Structured Streaming**, a scalable and fault-tolerant stream processing engine built on the Spark SQL engine. Streaming computations are expressed similarly to batch computations; Spark incrementally updates results as new data arrives ⁴.
- **Machine Learning (MLlib) and Pipelines:** MLlib provides a high-level API for building scalable machine-learning pipelines. Users can train and tune models and chain together transformers and estimators to create end-to-end ML workflows ⁵.
- **Spark Core and RDDs:** The Spark Core engine supplies low-level constructs like Resilient Distributed Datasets (RDDs) and in-memory processing ⁶. While RDDs give fine-grained control and support unstructured data, they require more code and don't benefit from query optimization; therefore DataFrames are preferred for most workloads ³.
- **Pandas API on Spark:** To ease the transition for pandas users, PySpark exposes a pandas API on Spark. This layer runs pandas operations on Spark clusters without changing most of the code ⁷. It allows a single codebase that works on small datasets with pandas and scales to larger datasets with Spark ⁸.
- **Spark Connect:** Newer versions offer Spark Connect, a client-server architecture that allows remote connectivity from any application. PySpark acts as the client, enabling Spark to be used as a service ⁹.

Options and Capabilities Supported in PySpark

PySpark exposes nearly all of Spark's capabilities through Python, giving users flexibility to build diverse data workloads:

Data processing APIs

- **DataFrames and Spark SQL** – support a wide range of data formats (Parquet, ORC, CSV, JSON), SQL queries, grouping and aggregation, window functions and user-defined functions (UDFs). Because Spark SQL and DataFrame operations run on the same optimized engine, you can mix SQL and Python code seamlessly ¹⁰.
- **RDD API** – provides low-level transformations (`map`, `filter`, `reduce`) and actions (`collect`, `count`). RDDs are fault-tolerant and flexible for unstructured data but require more verbose code and lack query optimization ³.
- **Pandas API on Spark** – allows pandas users to scale DataFrame operations using familiar pandas syntax ⁷.

Streaming and event processing

- **Structured Streaming** – processes continuous streams with the same API used for batch data. It automatically manages state and can output to files, databases or other sinks ⁴.
- **Legacy Spark Streaming (DStreams)** – an older streaming API still present for backward compatibility; however, Structured Streaming is recommended ¹¹.

Machine learning and analytics

- **Mlib** – includes algorithms for classification, regression, clustering, recommendation and dimensionality reduction ⁵.
- **Pipelines** – provide a declarative framework for building, maintaining and testing data pipelines, allowing users to focus on transformations rather than execution mechanics ¹².

Additional libraries

- **Graph analysis** – although Spark's GraphX library is implemented in Scala/Java and not directly available via PySpark, users can leverage **GraphFrames**, an extension that brings graph processing to the PySpark DataFrame API. GraphFrames enables PageRank, breadth-first search and motif finding on graphs stored in DataFrames.
- **External integrations and connectors** – PySpark integrates with Hive, HDFS, Amazon S3, Azure Data Lake, JDBC-compliant databases and other storage systems through built-in readers and writers. It can run on local mode, Yarn, Kubernetes or Mesos and can be embedded in notebooks or scheduled via Airflow or other orchestrators.

Limitations and Missing Features

Despite its rich feature set, PySpark has several limitations and areas where further development would improve usability and performance:

Python-JVM overhead

- **UDF performance penalties:** When you call a Python UDF, Spark treats your function as a black box. The Catalyst optimizer cannot inspect the logic, so optimizations such as predicate push-down and whole-stage code generation are disabled ¹³. Each row must cross the JVM-Python boundary, incurring serialization and deserialization overhead ¹⁴. In addition, the Python worker executes UDF code sequentially because of the Global Interpreter Lock (GIL), further reducing throughput ¹⁵. Memory management is separate from the JVM, so Python worker processes may run out of memory or throw exceptions that crash tasks ¹⁶.
- **Single-threaded UDF execution:** Because Python code executes in a single thread per task, custom UDFs cannot take full advantage of multi-core CPUs. Scala or Java UDFs are often faster because they run directly on the JVM and benefit from Spark's code generation.

API and ecosystem gaps

- **No typed Dataset API:** The strongly typed Dataset API available in Scala and Java is not exposed in PySpark. Python developers must use DataFrames (untyped) or RDDs, losing compile-time type safety.
- **GraphX unavailability:** Spark's native GraphX library is Scala/Java only. Although the GraphFrames package fills some gaps, the full GraphX API—such as graph builders, graph-oriented optimizations and built-in algorithms—remains unavailable to Python users.
- **Complex debugging and performance tuning:** Distributed computing introduces complexity. Beginners face a steep learning curve around configuring executors, partitions, and memory to avoid shuffling and skew. Debugging distributed jobs requires examining Spark UI metrics and logs. Poor configurations can lead to bottlenecks or OOM errors. External articles note that PySpark requires careful performance tuning and dependency management, and that integrations with Python libraries like pandas or NumPy need explicit conversions that can be slow for large datasets ¹⁷.
- **Integration limitations:** While the pandas API on Spark helps, many Python libraries cannot run directly on workers. Converting large PySpark DataFrames to pandas via `.toPandas()` must collect all data to the driver and can exhaust memory ¹⁸.
- **Environment and dependency management:** Running PySpark often requires aligning Python versions, Spark versions, and Java dependencies. Inconsistent environments can cause runtime errors or poor performance ¹⁷.

Development and operational challenges

- **Steep learning curve:** Distributed data processing requires understanding partitions, shuffles and lazy evaluation. Beginners may struggle with mental models like lineage, caching and checkpointing ¹⁷.
- **Limited built-in data quality checks and lineage:** Unlike tools such as dbt, PySpark has no native support for data quality testing or lineage tracking. These capabilities must be built manually or integrated from other frameworks.
- **Limited declarative pipeline support:** Although PySpark includes Pipelines and job definitions, it lacks high-level declarative orchestration and version control found in modern orchestrators. Users often combine Spark with Airflow, Prefect or dbt to manage complex workflows.

Areas for Improvement

- **Arrow-based and vectorized UDFs:** Expanding support for pandas UDFs and Arrow-based execution can reduce serialization overhead and enable vectorized operations in Python. Further improvements in cross-language integration and asynchronous execution would minimize the cost of moving data between Python and JVM.
- **Typed Dataset and GraphX support:** Providing Python bindings for the Dataset API and GraphX would allow Python users to leverage compile-time type checks and advanced graph algorithms. Community projects like GraphFrames show potential, but official support would increase stability and features.
- **Better debugging and observability:** Integrating more robust debugging tools (similar to Spark UI but more Python-friendly) and easier ways to trace lineage would help developers troubleshoot performance issues and understand data flows.
- **Native data quality and testing tools:** Embedding built-in assertions, tests and expectations (e.g., Great Expectations-style) could help validate data transformations and catch issues early.
- **Simplified environment management:** Packaging Spark, Java and Python dependencies into unified installers or container images, plus better guidance on version compatibility, would reduce friction for new users.
- **Enhancements for the pandas API on Spark:** Improving parity with pandas functions and performance, adding more plotting and interactive capabilities, and expanding documentation will make it easier for pandas users to adopt PySpark.

Conclusion

PySpark empowers Python developers to harness the distributed computing power of Apache Spark. Its support for SQL, streaming, machine learning and the pandas API enables a wide range of data tasks. However, the Python-to-JVM boundary introduces overhead that can make UDFs slow ¹⁹. Missing typed

APIs and graph processing support, challenging debugging, and complex environment setup present areas where the ecosystem could evolve. Continued investment in cross-language optimization, richer Python APIs and better tooling will make PySpark more approachable and performant for the data community.

1 2 3 4 5 6 7 8 9 10 11 12 **PySpark Overview — PySpark 4.1.0 documentation**

<https://spark.apache.org/docs/latest/api/python/index.html>

13 14 15 16 19 **Why Your PySpark UDF Is Slowing Everything Down**

<https://www.canadiandataguy.com/p/why-your-pyspark-udf-is-slowing-everything>

17 18 **What is PySpark? Features, Benefits, and Getting Started**

<https://www.theknowledgeacademy.com/blog/what-is-pyspark/>