

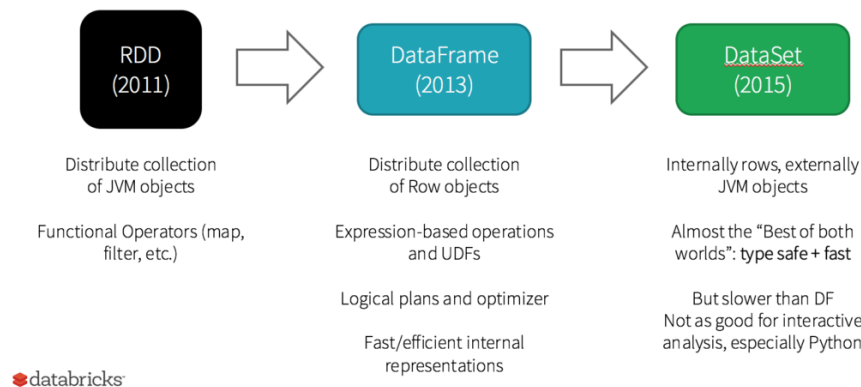


(1)

Resilient Distributed Dataset (RDD)

[Back to glossary \(/glossary/\)](#)

History of Spark APIs



RDD was the primary user-facing API in Spark since its inception. At the core, an RDD is an immutable distributed collection of elements of your data, partitioned across nodes in your cluster that can be operated in parallel with a low-level API that offers transformations and actions.

5 Reasons on When to use RDDs

1. You want low-level transformation and actions and control on your dataset;
2. Your data is unstructured, such as media streams or streams of text;
3. You want to manipulate your data with functional programming constructs than domain specific expressions;
4. You don't care about imposing a schema, such as columnar format while processing or accessing data attributes by name or column; and
5. You can forgo some optimization and performance benefits available with DataFrames and Datasets for structured and semi-structured data.

What happens to RDDs in Apache Spark 2.0?

Are RDDs being relegated as second class citizens? Are they being deprecated? The answer is a resounding NO! What's more is you can seamlessly move between DataFrame or Dataset and RDDs at will—by simple API method calls—and DataFrames and Datasets are built on top of RDDs.

A Tale of Three Apache Spark APIs: RDDs, DataFrames, an...



Additional Resources

- **A Tale of Three Apache Spark APIs: RDDs vs DataFrames and Datasets** (<https://databricks.com/blog/2016/07/14/a-tale-of-three-apache-spark-apis-rdds-dataframes-and-datasets.html>)
- **Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing** Research Paper (<https://databricks.com/research/resilient-distributed-datasets-a-fault-tolerant-abstraction-for-in-memory-cluster-computing-nsdi-2012>)

[Back to glossary \(/glossary/\)](/glossary/)

TRY DATABRICKS FOR FREE (/TRY-DATABRICKS?ITM_DATA=GLOSSARY-FOOTERCTA-TRIAL)

Product
(</product/data-lakehouse>)

Platform Overview
(</product/data-lakehouse>)

Learn & Support
(</spark/about>)

Documentation
(</documentation>)

Glossary (</glossary>)

Solutions (</solutions>)

By Industries
(</solutions#by-industry>)

By Role
(</solutions#by-role>)

Company
(<https://databricks.com/company/about-us>)

About Us
(<https://databricks.com/company/about-us>)

Pricing
(/product/pricing)

Open Source Tech
(/product/open-source)

Try Databricks (/try-databricks?itm_data=SiteWide-Footer-Trial)

Demo
(/discover/demos)

Training & Certification
(https://databricks.com/learn/training/home)

Help Center
(https://help.databricks.com/s/)

Legal (/legal)

Online Community
(https://community.databricks.com/s/)

Professional Services
(https://databricks.com/professional-services)

Careers at Databricks
(https://databricks.com/company/careers)

Diversity and Inclusion
(/company/diversity)

Newsroom
(https://databricks.com/newsroom)

Company Blog
(https://databricks.com/blog/category/careers)

Contact Us
(/company/contact)



See Careers at Databricks (/company/careers)

(https://databricks.com)

 Worldwide

(https://www.linkedin.com/company/databricks)

(https://www.facebook.com/pages/Databricks/560203)

(https://twitter.com/databricks)

(https://databricks.com/feed)

(https://www.glassdoor.com/Overview/Working-at-Databricks-EI_IE954734.11,21.htm)

(https://www.youtube.com/c/Databricks)

Databricks Inc. 160 Spear Street, 13th Floor San Francisco, CA 94105 1-866-330-0121

© Databricks 2022. All rights reserved. Apache, Apache Spark, Spark and the Spark logo are trademarks of the **Apache Software Foundation**. (https://www.apache.org/)

Privacy Policy (/privacypolicy) | Terms of Use (/terms-of-use) | Modern Slavery Statement (/wp-content/uploads/2021/10/DS-Databricks-Modern-Slavery-Policy-Statement-FY22.pdf)