# Exploration of Python and Scala in Databricks Advantages and Challenges

## Hari Prasad Bomma

Data Engineer, USA
haribomma2007@gmail.com

**Abstract**

**In the ever evolving realm of big data processing and analysis, the choice of programming language has become a critical factor in determining the efficiency and effectiveness of data driven solutions. Two prominent contenders in this area are Python and Scala. Each boasting a unique set of features and capabilities that cater to the diverse needs of data scientists and engineers. This article explores the performance and scalability of large scale data processing tasks in Databricks using Python and Scala. It delves into the specific advantages and challenges associated with each language and provide practical insights and recommendations for data engineers and scientists.**

**Keywords: Databricks, Python, Spark Clusters, Machine Language, Apache Spark, Notebook, Delta Lake, Data Lake, Pyspark, JVM**

**Introduction:**

Databricks, a cloud based data processing and analytics platform. Founded in 2013 by the creators of Apache Spark, Databricks has emerged as a game changer in the realm of big data and machine learning. It rapidly evolved, becoming a leading provider of a unified analytics platform that simplifies the management and processing of large scale data.

Databricks is built upon the foundation of Apache Spark, a powerful open source distributed computing framework that excels in processing and analyzing vast amounts of data [2]. Spark's ability to perform tasks in memory, rather than relying on disk based storage, has made it a preferred choice for a wide range of data intensive applications, from real time stream processing to advanced machine learning [3].

The evolution of Databricks has been tied to the growth and adoption of Apache Spark. As Spark's popularity has increased, Databricks has leveraged its expertise and insights to create a comprehensive platform that integrates Spark with a variety of other tools and technologies.

One of the key features that set Databricks apart is its ability to handle different data types and sources, from structured data in databases to unstructured data from various sources. This versatility, combined with the power of Spark, has made Databricks a popular choice for organizations looking to extract valuable insights from their data. The Databricks platform provides a collaborative and user friendly environment for data engineers, data scientists, and analysts to work together in a unified workspace. By offering a range of capabilities, including interactive notebooks for data exploration,

automated model deployment pipelines, and scalable data processing infrastructure, Databricks empowers teams to accelerate their data driven initiatives and unlock new business opportunities [2] [3].

The company has expanded its offerings to include features like Delta Lake, a data lake management system, and MLflow, an open source platform for machine learning model management [3]. These advancements have solidified Databricks' position as a comprehensive data analytics platform, catering to the growing needs of organizations grappling with the challenges and opportunities presented by the era of big data and artificial intelligence. [4] [3][5]

## Language Overview:

### Python:

Python is a high level, interpreted programming language introduced by Guido van Rossum in 1991, renowned for its emphasis on code readability and simplicity. It's clear and intuitive syntax, paired with an extensive array of libraries (such as Pandas, NumPy, and Scikit learn), makes it a popular choice among both beginners and experienced developers. Python is an interpreted language, which facilitates quick testing and debugging, and its dynamic typing allows for rapid and flexible coding. Within the Apache Spark ecosystem, Python is primarily utilized through PySpark, which enables the creation of Spark applications in Python. PySpark's blend of simplicity and efficiency makes it ideal for data scientists and engineers to prototype and test data pipelines and machine learning models, particularly in interactive environments like Jupyter Notebooks within Databricks.

Key features: Python's high level nature allows programmers to focus more on problem solving rather than low level mechanical constructs, making it an appealing choice for financial analysts and quantitative programmers. Python's ecosystem is home to a vast array of libraries and frameworks, such as Scikit learn, Pandas, and NumPy, which provide powerful tools for data analysis, machine learning, and scientific computing. Python's ability to interoperate with C and C++ programming languages allows for the integration of new Python based programs with existing code investments, enabling efficient and seamless development [8] [1].

### Advantages of Python in Databricks

Python is a popular choice for data scientists and engineers due to its simplicity, extensive ecosystem of libraries, and ease of use. Within the Databricks environment, Python offers several advantages. Python is a widely adopted language with a large user base, providing access to a vast array of libraries and tools, such as pandas and scikit learn, which are well integrated with Databricks. The simplicity and readability of Python code make it an attractive option for data teams, as it enables rapid development and prototyping of data driven solutions.

### Scala:

Scala, is a high level, statically typed programming language introduced by Martin Odersky in 2003. Combining object oriented and functional programming paradigms, Scala is praised for its conciseness, scalability, and robust performance. Running on the Java Virtual Machine (JVM), it is fully compatible

with Java libraries and frameworks. Scala's support for functional programming, with its static typing and concurrency model, enhances code safety, performance, and parallel processing capabilities. Scala holds a special place in the Apache Spark ecosystem, because Spark was originally written in Scala. Its performance benefits from being a compiled language and its seamless integration with core Spark APIs. The functional programming style of Scala aligns well with Spark's data processing model, enabling efficient and concise data transformations, thereby leveraging Scala's inherent power and flexibility.

Keyfeaures: Scala's design principles make it particularly well suited for building scalable and concurrent systems, which is a crucial requirement in financial applications. Scala's support for functional programming constructs, such as immutable data structures and higher order functions, is well suited for building scalable and concurrent systems.

## Advantages of Scala in Databricks

Scala, has gained traction in the Databricks ecosystem due to its strong typing, functional programming capabilities, and compatibility with the widely adopted Apache Spark framework.

Scala's static typing and functional programming provides improved code reliability, scalability, and performance, in the context of large scale data processing.

Apache Spark, one of the most widely used open source processing engines for big data. It has rich language integrated APIs and a wide range of libraries, making Scala a natural choice for Databricks users.

Scala's ability to integrate with Java and its interoperability with Databricks' existing infrastructure and tools can be a huge advantage for organizations already invested in the Java ecosystem.

Apache Spark and its Scala based APIs have been instrumental in deploying big data solutions across a wide range of organizations, highlighting the platform's versatility and scalability.

## Literature Review:

The performance evaluation of data driven intelligent algorithms for big data ecosystems has been a subject of growing interest in the research community [6]. Researchers have highlighted the advantages of highly scalable machine learning architectures, such as Apache Spark, and their ability to handle massive datasets effectively. Additionally, the literature emphasizes the importance of comparing the performance of various frameworks, including Spark MLlib, scikit learn, and RapidMiner, to guide the selection of the most appropriate tools for big data analytics.[6]

Moreover, the literature outlines the challenges and requirements that have emerged in deploying Apache Spark to a wide range of organizations. These insights provide a valuable context for understanding the factors that influence the choice between Python and Scala in the Databricks environment.

A comparative analysis of SQL on Hadoop systems, including Spark SQL, found that each system employs different architectural approaches and optimization techniques, leading to varying performance characteristics [7].

**Challenges and Considerations:**

While both Python and Scala offer unique advantages in the Databricks environment, there are also challenges and considerations to be taken into account:

Code complexity and developer productivity: Scala's more rigid syntax and functional programming approach may present a slower learning curve for data teams which used to Python's simplicity.

Resource management: Effective resource allocation and management, such as executor memory configuration, can have a significant impact on the scalability and performance of big data systems in Databricks, regardless of the chosen language.

Language interoperability and integration: Transitioning between Python and Scala within Databricks, as well as integrating with other tools and frameworks, can introduce additional complexity for data teams.

Potential performance tradeoffs: While Scala's strong typing and functional programming can provide performance benefits, there may be scenarios where Python's high level abstractions and rich ecosystem of libraries offer comparable or better performance for specific use cases.

**Case study:** An ETL Pipeline Overview
**Extract**: Read a dataset.
**Transform**: Perform a simple transformation.
**Load**: Write the transformed data to a new location.

**Python (PySpark)**
*Step 1: Extract*
```
# Load necessary libraries
import pyspark.sql.functions as F
from pyspark.sql import SparkSession
# Initialize SparkSession
spark = SparkSession.builder.appName("Python ETL").getOrCreate()
# Read data
df = spark.read.csv("/path/to/input/data.csv", header=True, inferSchema=True)
```

*Step 2: Transform*
```
# Perform a transformation   e.g., adding a new column with prefix
transformed_df = df.withColumn("new_column", F.concat(F.lit("prefix_"), F.col("existing_column")))
```

*Step 3: Load*
```
# Write transformed data to a new location
transformed_df.write.mode("overwrite").csv("/path/to/output/data")
```

**Scala**

*Step 1: Extract*

```
// Load necessary libraries
import org.apache.spark.sql.SparkSession
import org.apache.spark.sql.functions.
// Initialize SparkSession
val spark = SparkSession.builder.appName("Scala ETL").getOrCreate()
// Read data
val df = spark.read.option("header", "true").csv("/path/to/input/data.csv")
```

*Step 2: Transform*
```
// Perform a transformation   e.g., adding a new column with prefix
val transformedDf = df.withColumn("new_column", concat(lit("prefix_"), col("existing_column")))
```

*Step 3: Load*
```
// Write transformed data to a new location
transformedDf.write.mode("overwrite").csv("/path/to/output/data")
```

**Comparison:**

**Ease of Development:**

> **Python**: Python's syntax is simpler and more readable, which can make the development process quicker for those familiar with the language.

> **Scala**: Scala's syntax, while more complex, can offer more power and flexibility, particularly advantageous for those already comfortable with functional programming paradigms.

**Performance:**

> **Execution Speed**: Scala generally offers faster execution due to being a compiled language and its strong typing system. PySpark (Python) might have slightly slower performance, especially with large datasets.

> **Concurrency**: Scala excels in handling concurrency due to its robust multi threading capabilities, whereas Python's Global Interpreter Lock (GIL) can pose challenges for multi threaded tasks.

**Integration and Ecosystem:**

> **Library Support**: Python's extensive libraries for data science and machine learning give it an edge in those domains, whereas Scala's interoperability with Java libraries is beneficial for leveraging existing Java based solutions.

**Results:**

After conducting the comparative analysis of Python and Scala in the Databricks environment, the key findings are as follows:

- Python and Scala both offer distinct advantages in the Databricks ecosystem, with Python's simplicity and extensive ecosystem of libraries appealing to data analysts and programmers,

while Scala's strong typing, functional programming capabilities, and integration with Apache Spark make it a compelling choice for large scale data processing and machine learning tasks.

- o The performance characteristics of Python and Scala within Databricks can vary depending on the specific use case and workload.

- o The ability to seamlessly integrate Python and Scala within the Databricks platform allows for a flexible and collaborative workflow, leveraging the strengths of both languages.

- o Effective resource management, such as executor memory configuration, is critical for ensuring the scalability and performance of big data systems in Databricks, regardless of the chosen language.

- o While Python's simplicity and developer productivity may be advantageous in certain scenarios, Scala's static typing and functional programming approach can provide improved code reliability, scalability, and performance, particularly for large scale data processing and machine learning tasks.

**Conclusion:**

The choice between Python and Scala in the Databricks environment ultimately depends on the specific requirements, team expertise, and the nature of the data processing and machine learning tasks at hand. Python's simplicity, extensive ecosystem, and ease of use make it a compelling choice for rapid prototyping, exploratory data analysis, and machine learning tasks, particularly for teams with a strong background in the language. Scala's focus on performance, scalability, and compatibility with Apache Spark make it a suitable choice for large scale data processing, complex business logic, and mission critical applications. Scala's interoperability with Java libraries is beneficial for leveraging existing Java based solutions and is a great fit for organizations with Java ecosystems.

To effectively leverage the strengths of both languages, data teams within the Databricks environment should consider a hybrid approach, where Python and Scala are used in tandem, leveraging the strengths of each language for different aspects of the data pipeline. Ultimately, the decision should be guided by a thorough understanding of the project requirements, the team's skill set, and a careful evaluation of the tradeoffs between the two languages in the context of the Databricks platform.

**References:**

[1].    Fletcher, S., & Gardner, C. (2012). Welcome to Python (p. 1). https://doi.org/10.1002/9780470685006.ch1

[2].    Vavilapalli, V. K., Murthy, A. C., Douglas, C., Agarwal, S., Konar, M., Evans, R., Graves, T., Lowe, J., Shah, H., Seth, S., Saha, B., Curino, C., O'Malley, O., Radia, S., Reed, B., & Baldeschwieler, E. (2013). Apache Hadoop YARN (p. 1). https://doi.org/10.1145/2523616.2523633

[3].    Guller, M. (2015). Big Data Analytics with Spark. In Apress eBooks. https://doi.org/10.1007/978-1-4842-0964-6

[4].    Saddad, E., El-Bastawissy, A., Hoda, M., & Hazman, M. (2020). Lake Data Warehouse Architecture for Big Data Solutions. In International Journal of Advanced Computer Science and

Applications (Vol. 11, Issue 8). Science and Information Organization. https://doi.org/10.14569/ijacsa.2020.0110854

[5]. Sheikh, R. A., & Goje, N. S. (2021). Role of Big Data Analytics in Business Transformation (p. 231). https://doi.org/10.1002/9781119711148.ch13

[6]. Junaid, M., Ali, S., Siddiqui, I. F., Nam, C.-S., Qureshi, N. M. F., Kim, J., & Shin, D. R. (2022). Performance Evaluation of Data-driven Intelligent Algorithms for Big data Ecosystem. In Wireless Personal Communications (Vol. 126, Issue 3, p. 2403). Springer Science+Business Media. https://doi.org/10.1007/s11277-021-09362-7

[7]. Tapdiya, A., & Fabbri, D. (2018). A comparative analysis of state-of-the-art SQL-on-Hadoop systems for interactive analytics. In arXiv (Cornell University). Cornell University. https://doi.org/10.48550/arxiv.1804.00224

[8]. Sheppard, K. (2014). Introduction to Python for Econometrics, Statistics and Data Analysis. http://janroman.dhis.org/AFI/Python/Python_introduction.pdf