

# The Specifics of Development in Highly Specialized Technical Domains: The Experience of Using Scala and Spark for ML Tasks

Danylo Liakhovetskyi  
Middle Java Backend Engineer at Agile Engine  
Pensacola, FL, USA

\*\*\*\*\*

## Abstract:

This article is dedicated to the use of the Scala programming language and the Apache Spark framework in specialized domains related to machine learning tasks. The purpose of the work is to analyze the advantages of these technologies in processing large volumes of data and ensuring high performance. The methodology includes an overview of the characteristics of specialized domains, an analysis of Scala and Spark's capabilities, genetic information analysis, and the processing of streaming data in the Internet of Things. The ways of integrating Scala and Spark with external libraries are discussed, expanding their application range. The research on these approaches is based on the analysis of publicly available articles on the internet, allowing for a comprehensive exploration of the topic. The results show that using Scala and Spark contributes to solving data processing tasks and creating adaptive machine learning models. Importantly, these technologies ensure scalability and enable parallel processing of information. The conclusions highlight that using Scala and Spark accelerates development, increases solution accuracy, and ensures the reliability of the obtained results. This work will be useful for developers, machine learning specialists, and researchers working on the application of technologies in specialized domains. The results confirm that the use of Scala and Spark enhances the effectiveness of machine learning solutions, improving performance.

**Keywords**— Specialized Domains, Scala, Apache Spark, machine learning, big data, performance, forecasting.

\*\*\*\*\*

## Introduction

Modern technologies require the processing of large volumes of data and the creation of intelligent systems, which is increasingly necessary in specialized technical fields. Sectors such as bioinformatics, industrial monitoring, and real-time information analysis impose strict requirements on tools that ensure flexibility and performance. Machine learning tasks in these areas are associated with the need to account for data complexity, high computation speed, and scalability. In such conditions, the choice of technologies becomes a crucial element for achieving the desired results.

The Scala programming language and the Apache Spark framework have long established themselves as effective tools for working with big data and developing algorithms. Scala offers capabilities for data manipulation and integration of computational approaches, while Spark provides distributed processing. The joint use of these tools opens up possibilities for developing solutions that meet specific requirements.

The topic remains relevant as there is an ongoing need to create solutions for big data processing and effective operation in multitasking and real-time environments. The advantages of Scala and Spark are recognized, but their application in analysis and forecasting tasks requires further scientific work to identify the potential and limitations of these technologies.

The goal of this work is to analyze the capabilities of the Scala and Apache Spark frameworks in processing large volumes of data and ensuring high performance.

## **Materials and Methods**

The article by Erraissi A., Azouazi M., Belangour A., Banane M. [1] discusses event prediction methods using machine learning models. The paper analyzes algorithms for handling data, including processing data streams and making predictions based on historical data. An important aspect is the selection of algorithms and features that improve the accuracy of predictions.

The article by Radulovic N., Boulegane D., Bifet A. [2] describes the Scalar-a platform, designed for machine learning competitions using streaming data. The focus is on algorithms that process information, with an emphasis on developing systems with minimal latency and high throughput when analyzing dynamically incoming data.

In the work by Yang J. [3], the use of the Spark platform to build scalable systems capable of processing large volumes of information is discussed. The focus is on designing systems in the context of the Internet of Things, where it is essential to process data coming from numerous devices in real time. Architectural solutions, algorithm optimization, and computational task distribution that affect system efficiency are also examined.

The article by Yu B. et al. [4] presents the Sparker method, aimed at simplifying machine learning models when using the Spark platform. The primary task is to reduce computational costs when training large models, which is crucial for processing big data. The work proposes methods for optimizing the process, thereby improving its efficiency.

The article by Mu C., Zheng J., Chen C. [5] discusses advancements in the development of chips for machine learning that provide computational efficiency compared to traditional neural networks. Recent developments presented at the International Conference on Semiconductor Technologies are discussed, and future directions that could impact performance in the coming years are predicted.

Thus, the discussed sources focus on the efficiency of computing systems, scalability, and the application of machine learning in specific fields such as real-time processing and epidemiology. Contradictions arise in different approaches to data processing and the need to reduce computational costs. Little attention is given to the integration of platforms, which would enhance the flexibility of models in applications. Furthermore, the issues of model resilience in the face of changing data and emergency situations, such as pandemics, are insufficiently addressed.

The methodology of this work involves reviewing the characteristics of domains and analyzing the capabilities of Scala and Spark, along with research into the processing of streaming data in the Internet of Things.

## **Results and Discussion**

In some fields, the required task involves working with data that differs in structural heterogeneity, making it necessary to use tools that can ensure the preprocessing and cleaning of data for further analysis. Consequently, an essential element is the ability to solve multiple tasks simultaneously while maintaining the required performance. For example, systems working with streaming data from sensors or video feeds must process information quickly, distributing resources as needed, which is crucial in the presence of time constraints that can impact the outcome.

Working in specialized areas involves developing models that account for physical or biological characteristics. Unlike standard machine learning tasks that focus on improving accuracy, it is necessary to integrate domain-specific knowledge into the algorithm. This requires the use of hybrid approaches that combine traditional machine learning methods with simulations, which ultimately allows solving tasks where mathematical machine learning methods alone are insufficient, and simulations help improve model quality by considering the specific features and complexities of a given domain. In the application of Scala and Spark for such machine learning tasks, hybrid methods can ensure an effective data processing workflow, accelerate computations, and facilitate the integration of different approaches [1, 5]. The development process in highly specialized technical domains is illustrated in Figure 1.

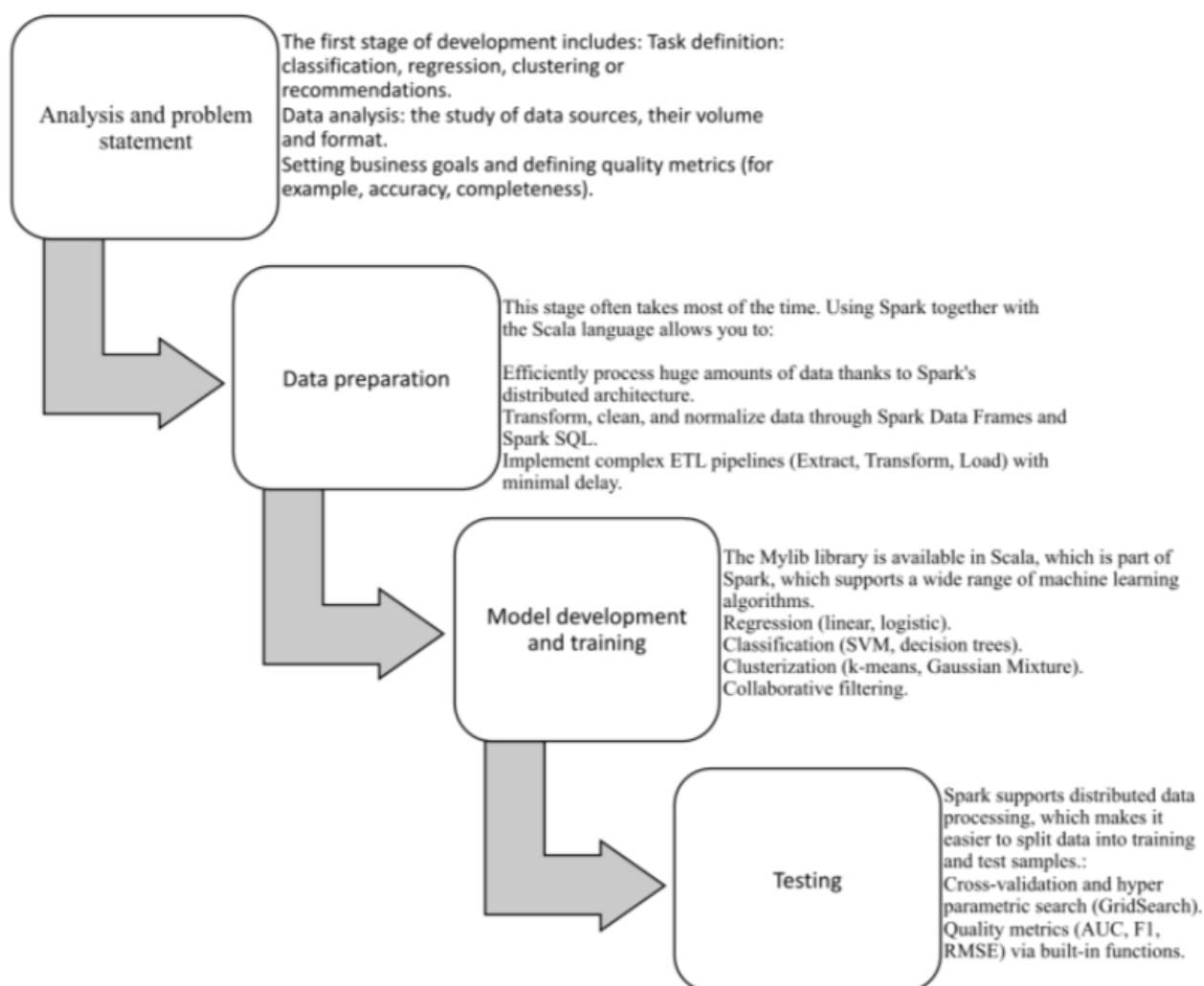


Fig. 1. The development process in highly specialized technical domains (compiled by the author)

The Scala programming language combines functional and object-oriented approaches, allowing for the development of systems for processing large data sets. Additionally, its use with the JVM enhances performance when handling such data, which is crucial when there is a high load.

The key features of Scala lie in the use of principles such as immutability of data and higher-order functions. This approach reduces the number of errors caused by state changes, improves the organization of parallel computations, ultimately enhancing system stability and increasing the predictability of processes, thus lowering the likelihood of failures.

For tasks related to time series analysis and streaming data processing, Scala provides tools for developing algorithms that meet the requirements [1, 2, 5]. Below, Table 1 compares the approaches of Scala and Apache Spark.

**Table 1.** Comparison of programming approaches [2, 5].

Criterion	Scala	Spark	Future Trends and Impact
<b>Programming Language</b>	A functional and object-oriented programming language. Popular for high-performance computing.	Not a programming language, but provides an API for Scala, Python, and others. Used for distributed data processing.	Continued development of multi-language solutions, with Spark being used with various programming languages, while Scala will maintain its role in performance-critical tasks.
<b>Ease of Use</b>	Has a steep learning curve and requires knowledge.	Easy to use due to its API and integration with other languages.	Simplification of interfaces, and integration with available programming languages.
<b>Machine Learning Support</b>	Can be used to implement custom ML algorithms, but lacks the necessary library support.	Spark MLlib provides built-in algorithms for machine learning, including classification, regression, clustering, and data processing.	Growth in the use of neural networks, further integration of Spark with new frameworks like TensorFlow and PyTorch to enhance ML capabilities.
<b>Scalability</b>	Offers high performance on a single server, but scalability is limited by JVM capabilities and the need for task distribution.	Spark is optimized for distributed data processing and scales well on clusters.	Continued development of cloud solutions, hybrid computing, and improved scalability considering domain-specific needs.
<b>Integration with Other Tools</b>	Can integrate with JVM libraries and the Java ecosystem, but is limited compared to other universal frameworks.	Easily integrates with Hadoop, Hive, HBase, Cassandra, and supports Python.	Increased integration capabilities with cloud platforms and ML tools like Keras and TensorFlow.
<b>Parallelism Support</b>	Supports parallelism through multitasking in JVM but is less efficient when processing distributed data.	Spark provides the most suitable method for handling distributed data.	Growing interest in parallel data processing and optimization of algorithms for scalable architectures.
<b>Use in Specialized Domains</b>	Widely used in financial, telecommunications, and engineering applications, where high performance is required.	Used in domains that require the processing of large volumes of data.	With the development of technologies such as IoT and edge computing, Spark and Scala may be applied to solve problems in new fields.

Apache Spark is a platform for processing large data sets in distributed systems. It uses in-memory computation, which accelerates processing compared to methods based on disk storage. The framework includes the MLlib library, which contains algorithms for solving classification, regression, and clustering tasks in the development of recommendation systems.

Spark supports working with streaming data through the Spark Streaming module. This capability is used for applications that require data analysis. The integration of data processing with machine learning methods ensures data relevance and accuracy of results [3, 4]. Below, Table 2 outlines the features of Scala and Apache Spark integration.

**Table 2.** Features of Scala and Apache Spark integration [3, 4].

Category	Integration Features	Challenges	Ways to Minimize
<b>Integration with existing systems</b>	- Integrates with the Java ecosystem through JVM, as well as with Big Data tools such as Hadoop and Hive.	- Integration issues may arise with outdated or non-standard systems.	- Rewriting code is necessary.
<b>Big data processing</b>	- Spark can process distributed data through cluster computing. - Scala is effective for local processing and computations.	- Performance issues may occur when processing extremely large volumes of data if the infrastructure is poorly configured.	- Lack of experience in optimizing Spark parameters for unusual tasks.
<b>Scalability, optimization</b>	- Spark scales easily on clusters and supports horizontal scaling.	- Scala is limited in scalability without optimal environment configuration.	- Performance loss due to inefficient resource use.
<b>Integration with other tools</b>	- As previously mentioned, Scala integrates well with other libraries like Java and JVM. - Spark supports numerous other frameworks and programming languages.	- Configuration and consistency issues with different versions of libraries and frameworks.	- Version compatibility issues when integrating with external libraries.

The use of Scala with Apache Spark enables efficient management of computational processes and integration of business logic with machine learning algorithms. An example of this approach is the DataFrame API for data processing, where Scala code is compact, simplifying development, reducing the likelihood of errors, and increasing productivity.

Scala allows the development of algorithms, integration with Spark MLlib, and connection to third-party libraries, which broadens the ability to address specific tasks. This combination of tools solves problems that require customization, taking into account the specifics of the application domain.

It is also worth noting that using Scala and Spark allows the processing of large volumes of data while integrating mathematical models based on the physical characteristics of equipment. In such cases, machine learning algorithms must be connected to physical models, requiring knowledge in both areas. Distributed data processing with Spark accelerates model construction, allowing for quick responses to changes.

Internet of Things (IoT) systems generate data at high frequencies, necessitating the development of solutions for processing this data. Spark Streaming provides the necessary tools for building streaming applications capable of integrating data with machine learning models and making real-time decisions based on predictions.

Such tasks require not only data processing but also preprocessing—filtering, aggregation, and transformation into the appropriate format. Scala enables flexible organization of these processes, creating solutions tailored to the needs of specific applications [1, 3, 5].

Thus, the use of Scala and Apache Spark in solving forecasting tasks allows for handling large volumes of data. These tools support parallel processing and enable the creation of systems that make real-time decisions. In the context of increasing data volumes and the need for rapid predictions, their application is relevant in fields such as finance, trade, IoT, and industrial monitoring.

## **Conclusion**

The aspects of using the Scala programming language and the Apache Spark framework for solving machine learning tasks were discussed. The results demonstrated that these technologies provide capabilities for working with large volumes of data and creating models that account for the specific characteristics of individual domains.

Scala offers functionality and compatibility with Java, which promotes flexibility in algorithm development. Apache Spark provides tools for scalable data processing, supports stream processing, and in-memory computing, which is essential for tasks that require rapid response. The use of Scala and Spark, therefore, is an effective approach for solving machine learning problems in specialized fields. These technologies enhance the quality of data processing and the integration of domain knowledge into the development process.

## **References**

1. Erraissi A., Azouazi M., Belangour A., Banane M. Machine Learning model to predict the number of cases contaminated by COVID-19 // *International Journal of Computing and Digital Systems*. – 2020. – pp.1-24.
2. Radulovic N., Boulegane D., Bifet A. SCALAR - A Platform for Real-time Machine Learning Competitions on Data Streams // *Journal of Open Source Software*. – 2020. – Vol. 5 (56). – pp. 2676.
3. Yang J. Machine Learning Platform Design and Application Based on Spark // *International Conference on Internet of Things and Machine Learning (IoTML 2023)*. – SPIE, 2023. – Vol. 12937. – pp. 331-335.
4. Yu B. et al. Sparker: Efficient Reduction for More Scalable Machine Learning with Spark // *Proceedings of the 50th International Conference on Parallel Processing*. – 2021. – pp. 1-11.
5. Mu C., Zheng J., Chen C. Beyond convolutional neural networks computing: New trends on ISSCC 2023 machine learning chips // *J Semicond*. – 2023. – Vol. 44 (50), 050203. – pp. 1-4