
BIG DATA ANALYSIS IN FOOTBALL

Khushi Talaviya*¹, Rushi Kamble*², Satvik Nayak*³, Saniddhya Dubey*⁴,

Dr. Hemant Gianey*⁵

*^{1,2,3,4}Dept. Of Computer Engineering Narsee Monjee Institute Of Management Studies Shirpur, India.

*⁵Dept. Of Computer Engineering Narsee Monjee Institute Of Management Studies Shirpur, India.

DOI : <https://www.doi.org/10.56726/IRJMETS45613>

ABSTRACT

The world of football is undergoing a data-driven transformation, and the demand for real-time analytics is greater than ever. This project explores the realm of real-time football analytics using Apache Spark, with the aim of enhancing match analysis, engaging fans, and preventing injuries. The research begins by setting clear objectives and systematically collecting data from various sources, including player tracking systems, social media platforms, and in-game statistics. The heart of the project lies in Apache Spark, a real-time data processing framework that efficiently handles vast volumes of data generated during live matches. This technology allows for data transformation, feature extraction, and the application of machine learning algorithms to provide real-time insights. User-friendly dashboards and interactive elements elevate the fan experience, offering live commentary and sentiment analysis during matches. The accuracy and effectiveness of the real-time analytics are rigorously assessed through comparisons with actual match outcomes and user feedback. By sharing research findings and methodologies, this project contributes to the advancement of real-time football analytics. The potential impact is profound, revolutionizing how football matches are analyzed, enhancing fan engagement, and offering invaluable insights to decision-makers in the world of football. The project underscores the power of Big Data technologies and a passion for the beautiful game in the ever-evolving landscape of football analytics.

Keywords: Big-Data, Scouting, Apache Spark, XG Boost-Classifer.

I. INTRODUCTION

The world of professional sports is constantly changing, especially football. In the past, success was largely dependent on the physical abilities of athletes. However, today, data and analytics play a crucial role. Big data has revolutionized the way football clubs, coaches, players, and analysts operate. It provides them with a treasure trove of insights that can be used to gain a competitive edge.

This project will explore the fascinating intersection of big data and football, showcasing how Apache Hadoop and Pig, two powerful components of the Apache ecosystem, are revolutionizing the way we understand, strategize, and optimize the beautiful game.

Football has always been a complex sport that requires strategy, skill, and teamwork. With the advent of digital technology and the proliferation of data-generating sources, the amount and variety of data available to football organizations has grown exponentially. This data includes player statistics, performance metrics, injury analysis, and even fan engagement data. When properly harnessed and analyzed, this data can provide a wealth of information.

To tackle the immense challenge of processing and analyzing this deluge of data, football clubs and organizations have turned to Apache Hadoop and Pig. These open-source, distributed computing frameworks are capable of handling massive datasets and executing complex data processing tasks in a scalable and efficient manner. With the help of these tools, football clubs can unearth hidden patterns, predict player performance, and devise winning strategies that were previously unimaginable.

This project aims to shed light on the applications of Apache Hadoop and Pig in the football domain. We will explore the various ways in which these technologies enable data-driven decision-making in areas like player scouting, game analysis, injury prevention, and fan engagement. By the end of this journey, you will have a comprehensive understanding of how big data and Apache ecosystem tools have come together to redefine the future of football, offering clubs and enthusiasts a new realm of possibilities to explore and conquer in the world's most popular sport.

In simpler terms:

Apache Hadoop and Pig are two powerful tools that football clubs and organizations can use to analyze massive amounts of data. This data can be used to make better decisions about players, games, injuries, and fans.

Here are some examples of how Apache Hadoop and Pig can be used in football:

- **Player scouting:** Hadoop and Pig can be used to analyze player statistics, performance metrics, and injury data to identify potential talent and assess the risk of injuries.
- **Game analysis:** Hadoop and Pig can be used to analyze game data to identify trends, patterns, and weaknesses. This information can be used to develop winning strategies and improve player performance.
- **Injury prevention:** Hadoop and Pig can be used to analyze injury data to identify risk factors and develop prevention strategies.
- **Fan engagement:** Hadoop and Pig can be used to analyze fan data to understand their needs and interests. This information can be used to develop targeted marketing campaigns and improve the overall fan experience.

By using Apache Hadoop and Pig, football clubs and organizations can gain a competitive edge in the world's most popular sport.

II. LITERATURE REVIEW

Fen Gong[1] et al. talk about how due to the The success or failure of any team sport hinges on the extent of collaborative effort within the team, and to gauge this, it is essential to rely on valid metrics. In the context of football, a highly team-oriented sport, there exists a multitude of metrics, making it crucial to identify those that exert a significant influence on the final outcome. This study employs correlation analysis to sift through an extensive array of data metrics, singling out the most effective performance indicators. Subsequently, it employs the entropy weight method to assign weightage to each of these performance indicators.[1]

The introduction of Pass Times Matrices provides a means to characterize the level of cooperation among different players in the entire team. Consequently, an objective, equitable, and all-encompassing evaluation model for team performance is established. To validate the model, it is applied to the data from Everton's 17-18 Premier League season, analyzing all 38 games in detail. The study conducts a correlation test between the outcomes of each match and the evaluation results, affirming the accuracy of the evaluation model.[1]

Through the process of correlation analysis, this research discerns meaningful performance indicators from a wealth of data metrics. The entropy weight method is employed to determine the significance of each performance metric, with a specific focus on those indicators positively associated with team cooperation. Consequently, a comprehensive evaluation score formula for team collaboration is formulated. Simultaneously, Pass Times Matrices are leveraged to represent the degree of cooperation among team members.[1]

Correlation analysis highlights a positive relationship between the comprehensive team cooperation score and match results. Improved cooperation within a football team is evidently linked to better game results. Furthermore, it is observed that the comprehensive evaluation of teamwork is positively correlated with specific indicators. Consequently, a targeted coaching approach focusing on these specific metrics can enhance the team's overall cooperation, thereby improving match results and achieving the desired competitive goals.[1]

Yao Wu[2] et al. elaborate that they established a social network based on player positions and passing sequences, drawing upon existing literature. Various metrics were employed to comprehensively evaluate the significance of different positions on the field.[2]

Our findings revealed that, within the context of football passing, the attacking midfielder (AMC) held the most pivotal position, closely followed by the central defending midfielder (MC). Through two-sample difference tests, it became evident that winning teams consistently exhibited superior performance in certain positions, particularly on the left wing, central forward, and both left and right defenders.[2]

To assess the influence of player positions on the overall network and test the robustness of passing networks, we systematically removed n positions (where n ranged from 1 to 10) from a team and evaluated network efficiency, primarily focusing on the positions on the left side.[2]

Leveraging social network analysis techniques, we probed the organizational efficiency of football passing

networks, employing various measures. A comparison across different playing positions revealed that, during the football passing process, the AMC played the most critical role, followed by the MC. Further, when comparing winning and losing teams in various matches, the analysis demonstrated that winning teams typically excelled in LFW (left forward), FW (forward), DL (defender on the left), and DR (defender on the right). Notably, LFW stood out as a key contributor to winning teams, emphasizing their efficient use of the field's width. Additionally, winning teams displayed superior defensive performance, with the roles of DL and DR proving more vital than DC (central defender) in thwarting opponents.[2]

It is essential to recognize that passing patterns evolve during specific game moments, such as counterattacks. By scrutinizing the characteristics of passing networks in different time periods and situations (e.g., offensive and defensive phases), we can gain dynamic insights into a team's performance, delving deeper into their distinctive attributes.[2]

Historically, existing literature primarily focused on general player positions, running performance, and passing traits. However, football matches are characterized by dynamic complexity, and analyzing individual positions in isolation fails to capture the complete picture. Moreover, players' positional characteristics influence their movement and passing requirements, emphasizing the need for a more nuanced approach. In light of this, we propose a future analysis of local network features. We contend that the proportion of core local network activity within the team's overall network during a game holds greater significance for a team's performance, as exemplified by Barcelona's MSN trio, comprising Messi, Suarez, and Neymar.[2]

Nobuyoshi Hirotsu[3] et al. talk about how this paper presents a statistical model of a football match that proves valuable for gaining insights into team characteristics. The authors evaluate these characteristics through the utilization of maximum likelihood estimators, focusing on factors like home advantage, offensive and defensive strength, and their interplay. In their analysis, they consider not only the goals scored but also ball possession. The model is constructed using data from the 1999-2000 season of the English Premier League. The research illustrates the team characteristics, providing a means to visualize offensive and defensive capabilities, home game preferences, and relative performance against specific opponent teams.[3]

Liu Tianbiao[4] et al. talks about how in early football game research employing data mining techniques, the typical approach involved incorporating all fundamental game actions into the statistics and deriving descriptive outcomes. However, a more refined approach, using modified data structures and updated algorithms, has enabled the generation of diagnostic results presented through trend networks. To obtain these results, the researchers initiated the process by cleansing the collected data according to the adjusted data structure. Subsequently, data processing was carried out utilizing a data mining program. The focus of this study was the analysis of offensive actions by the German team, aiming to identify not only general playing patterns but also critical offensive combinations associated with goal-scoring opportunities.[4]

The research delves into the performance of the German women's team, employing the modified Apriori algorithm. During the first half of the game, the team exhibited a strong performance, with particular efficiency displayed by Player No. 4 on the left side. However, the second half did not match the prowess of the first. Nevertheless, Player No. 2 remained active, particularly on the right side.[4]

The findings of this study align with prior research, emphasizing the efficacy of the data mining method and the adapted data structure for game description and simulation. The descriptive results obtained here complement previous research [22]. This innovative approach excels in extracting useful, typical descriptive results from a sea of similar data, thus enabling coaches to make informed decisions without being overwhelmed by an abundance of information.[4]

This novel method not only reveals the general playing pattern but also identifies key combinations and players who positively influence scoring opportunities. The use of a tendency network facilitates an easy and insightful comparison between general and effective playing patterns.[4].

Yiannakos[5] et al. highlights that the objective of this study was to examine the goal scoring patterns in elite-level soccer matches. The study sample comprised 32 games from the European Championship (Euro 2004). Data analysis was conducted using cross-tabulation and chi-square methods.

The findings revealed that a higher proportion of goals

were scored in the second half (57.4%) compared to the first half (42.6%). Concerning the type of offense, organized attacks accounted for the highest frequency (44.1%), followed by goals resulting from set plays (35.6%) and counter-attacks (20.3%). When assessing the actions leading to goals, long passes were the most frequent (34.1%). A specific focus was placed on dead-ball situations, with corners and free kicks demonstrating a higher frequency during the games.[5]

In terms of scoring attempts, the study reported the following percentages: 44.4% in the penalty area, 32.2% in the goal area, and 20.4% outside the penalty area. The results underline the importance of training for dead-ball situations. Additionally, attention should be given to player fatigue toward the end of a game, which can lead to opponent team goal scoring, necessitating effective training to address this issue.[5]

In summary, this study offers two key insights. Firstly, it emphasizes the necessity to enhance physical conditioning, particularly players' aerobic fitness and stamina, to meet the heightened demands of late-game situations, where goal-scoring opportunities are more prevalent. Secondly, it suggests a focus on game tactics and training configurations, with an emphasis on set plays, corner-kicks, and free-kicks, which account for a significant portion of goals scored. It also advocates training players in effective marking and unpredictable execution of dead-ball situations to increase goal-scoring chances.[5]

Furthermore, dedicating additional training time to team tactics to achieve synchronized team movement and specific training to counter long passes is recommended. This study contributes valuable insights into contemporary football, highlighting the challenges players face in maintaining performance until the end of matches and the inclination of coaches toward organized offensive strategies in elite European soccer tournaments.[5].

Francisco Henriques[6] talks about how the research in question delves into the intersection of football clubs, Big Data, and competitive advantage. In the modern sports landscape, the ability to gain a competitive edge is paramount for football clubs, and technology has played a pivotal role in this pursuit. The study examines how the utilization of Big Data can benefit football clubs in two key areas: injury prevention and player decision-making.[6]

For injury prevention, the research found that the ability to closely manage and regulate players' workloads throughout the season can significantly reduce the occurrence of injuries. The study notes a decline in overall injury rates in elite European football clubs between 2001 and 2017, attributing this improvement to advanced training methods informed by data analysis. This underscores the vital role of Big Data in enhancing football clubs' injury prevention strategies.[6]

Regarding player decision-making, the study highlights the potential of Big Data to improve a player's ability to make effective decisions during matches. Access to data on game moments and performance indicators simplifies and accelerates the evaluation of players' decision-making capabilities. The research also notes the use of sport-specific software to train players' decision-making skills, emphasizing that Big Data has the potential to enhance this crucial aspect of the game.[6]

The study's ultimate aim is to explore how Big Data applications in injury prevention and decision-making can influence sporting results. While the study acknowledges that success in football is influenced by various factors beyond data analysis, it underscores the potential of Big Data to tip the balance in favor of teams, even in a sport known for its low-scoring nature. The research advocates for the continuous evolution of technology and data analysis to provide football clubs with the tools they need to maintain a competitive advantage on the field. In summary, the study underscores the growing importance of Big Data in football and its potential to revolutionize injury prevention and player performance, ultimately contributing to a competitive edge in this highly competitive sport.[6].

Ashish Chouhan[7] et al. talk about how in this paper, Shotifier, a binary classification pipeline, is introduced, leveraging the concept of hybrid parallelism. The primary focus of the Shotifier pipeline is on forwards and strikers in the realm of football. It utilizes match statistics, spatial variables like the initial and final ball positions, as well as tactical elements such as events and subevents to classify whether a shot results in a goal (shot conversion). Given the low-scoring nature of football, the utmost emphasis is on securing the winning goal. However, due to the game's unpredictability, coaches must consider numerous factors like team formation, weather conditions, lineup, and past results when making strategic decisions. Identifying the pivotal

factors that significantly influence a player's ability to convert a shot into a goal is a complex challenge. This paper focuses its research on forwards and strikers to predict the influential factors on shot conversion.[7] The Shotifier pipeline adopts a hybrid parallelism approach, consisting of two key steps. First, it employs Kernel Density Estimator (KDE) to visually depict concentrated zones on the football pitch during specific match events and to identify zones with maximum ball activities based on the ball's starting and ending positions. Second, it employs historical match statistical data to determine the factors with a substantial impact on shot conversion, using machine learning models such as Support Vector Machine (SVM), XGBoost, Random Forest (RF), and Multi-layer Perceptron (MLP). The pipeline employs k-fold cross-validation to assess model effectiveness and reports a precision of 70.39% for the XGBoost classifier, outperforming other models considered in the study.[7]

In summary, this paper introduces Shotifier, a binary classification pipeline that integrates various factors to predict shot conversion in football, with a primary focus on forwards and strikers. It employs a hybrid parallelism approach and machine learning models to enhance the accuracy of shot conversion predictions, offering valuable insights for coaches and managers in strategizing their formations and tactics.[7]

Matthias Kempe [8] et al. talk about in the realm of sports analytics, there has been a notable evolution, particularly in soccer analytics, owing to the increased availability of substantial tracking data. The study at hand focuses on evaluating passing performance in soccer, aiming to support the notion that tactical behavior in team sports can be effectively analyzed exclusively using tracking data. To substantiate this claim, the research explores the correlation between changes in spatiotemporal variables related to passing and essential performance indicators. The findings clearly demonstrate that spatiotemporal variables possess the capability to predict pass accuracy and key performance indicators at the individual player level, thereby confirming the initial hypothesis.[8]

Additionally, the study introduces a straightforward composite performance indicator for the assessment of passes and players, all based on tracking data. In summary, the results of this study offer a viable approach for real-time assessment of tactical behavior and present a novel method for scouting and evaluating players in soccer and team sports in general.[8]

In essence, this paper furnishes empirical evidence that tactical performance in soccer can be thoroughly assessed using tracking data, thus eliminating the need for human observations. This discovery is significant for the field, showcasing the advantages of intelligent data scouting. Moreover, this approach marks a pioneering development in real-time game analysis, affording coaches more reliable data for in-game decision-making. Furthermore, it opens up new avenues for player evaluation and scouting, potentially granting clubs a substantial edge in the competitive transfer market. From the perspective of data science, this study underscores the potential of data science techniques in exploring intricate human behaviors.[8]

Samah Alouf[9] et al. highlights how in the realm of sentiment analysis research, there has been a notable utilization of online soccer content generated by both fans and teams. This study introduces a specialized approach for comprehending the sentiment expressed in soccer-related conversations. To achieve this, the researchers have created a soccer-specific lexicon, which plays a crucial role in developing a sentiment model trained on soccer datasets. The outcomes of this research highlight the efficacy of this approach in accurately discerning fans' emotions during soccer events. The study's core contribution lies in elucidating the process of automatically crafting a domain-specific lexicon from a soccer dataset. The results underscore the effectiveness of the proposed approach in recognizing sentiments expressed within soccer conversations, encompassing pre-match, in-game, and post-match discussions. Furthermore, an interesting revelation in the research is the connection between fan activities and the occurrence of goal-scoring events. It was observed that goal occurrences often resulted in a shift in sentiment, with positive expressions directed toward the scoring team and negative sentiments aimed at the opposing team.[9]

It's essential to note that the sentiment analysis in this work was conducted exclusively on English tweets. As a future enhancement, the researchers contemplate incorporating multilingual tweets to offer a more comprehensive representation of sentiment. Additionally, there's a plan to refine the soccer lexicon by integrating domain-specific expert knowledge to further enhance the accuracy and relevance of sentiment analysis in the soccer domain.[9]

III. METHODOLOGY

To embark on this research endeavor, it is imperative to commence with a clear articulation of your research goals. Define the specific aspects of football analytics that you intend to enhance through real-time data. Whether your focus is on real-time match analysis, engaging fans, preventing injuries, or a combination of these objectives, the establishment of well-defined objectives provides a robust foundation for guiding your research efforts. At the core of real-time analytics lies data collection. It is vital to identify and establish reliable data sources that can provide real-time information during football matches. These sources encompass player tracking systems, social media platforms, in-game statistics, and other repositories that contribute relevant data. In parallel, specify the critical data points essential for your analysis, including real-time player positions, dynamic event data (such as goals, fouls, passes), and fan sentiment data, among other information pertinent to your research objectives.

To ensure the timely and uninterrupted acquisition of real-time data, it is crucial to develop and implement data collection pipelines equipped to continuously stream and process data during live matches. Managing the substantial volumes of data generated during these matches is a fundamental consideration. Implement a robust Big Data storage solution, such as the Hadoop Distributed File System (HDFS), to efficiently store and organize this wealth of real-time data. Prioritize data reliability and integrity by incorporating redundancy mechanisms and backup strategies into your data management approach. This is crucial to prevent data loss and ensure data availability. Efficiently processing real-time data is at the heart of your research. Integrate real-time data processing frameworks like Apache Kafka or Apache Flink to ingest, process, and analyze data as it arrives. These frameworks play a pivotal role in ensuring that the data is handled promptly and efficiently. Develop and employ data transformation and feature extraction algorithms to prepare raw real-time data for analysis. These algorithms should be designed to extract actionable insights from the continuously updated data streams.

Select the most appropriate machine learning and data analysis algorithms based on your research objectives. These algorithms may include techniques for sentiment analysis, predictive analytics, and tactical insights, depending on the scope of your research. Ensure that the selected algorithms are capable of operating effectively in a real-time context where timely results are of utmost importance to decision-makers and fans. The training of machine learning models is a crucial aspect of your research. Use historical football data to train your chosen machine learning models, serving as a baseline for real-time predictions and insights. These models should be designed to continuously adapt and learn as new real-time data becomes available. Develop user-friendly dashboards and visualization tools to present real-time insights to coaches, analysts, and fans during matches. The user interface should be intuitive, allowing users to interact with and comprehend the real-time data effectively. Implement interactive elements, such as live commentary or displays of fan sentiment, to engage the audience and enhance the overall viewing experience. Assess the accuracy and effectiveness of your real-time analytics by comparing the predictions and insights generated with the actual outcomes of football matches. Conduct user surveys and feedback collection to evaluate the impact of your real-time analytics on decision-making by coaches and the level of engagement experienced by fans. Given the real-time nature of the data, privacy and security considerations are paramount. Address data security and privacy concerns, especially when dealing with real-time data that may include personal information or sensitive data. Ensure strict compliance with relevant data protection regulations to protect user privacy. Continuously monitor the performance of your real-time data processing infrastructure and make necessary optimizations to accommodate increased data volumes or improve response times. Scalability is crucial to ensure that the infrastructure can effectively handle the dynamic nature of real-time data streams. Consider future enhancements and expansion of your research, such as multi-language support for sentiment analysis and the integration of domain-specific expert knowledge to enhance the accuracy and relevance of real-time insights.

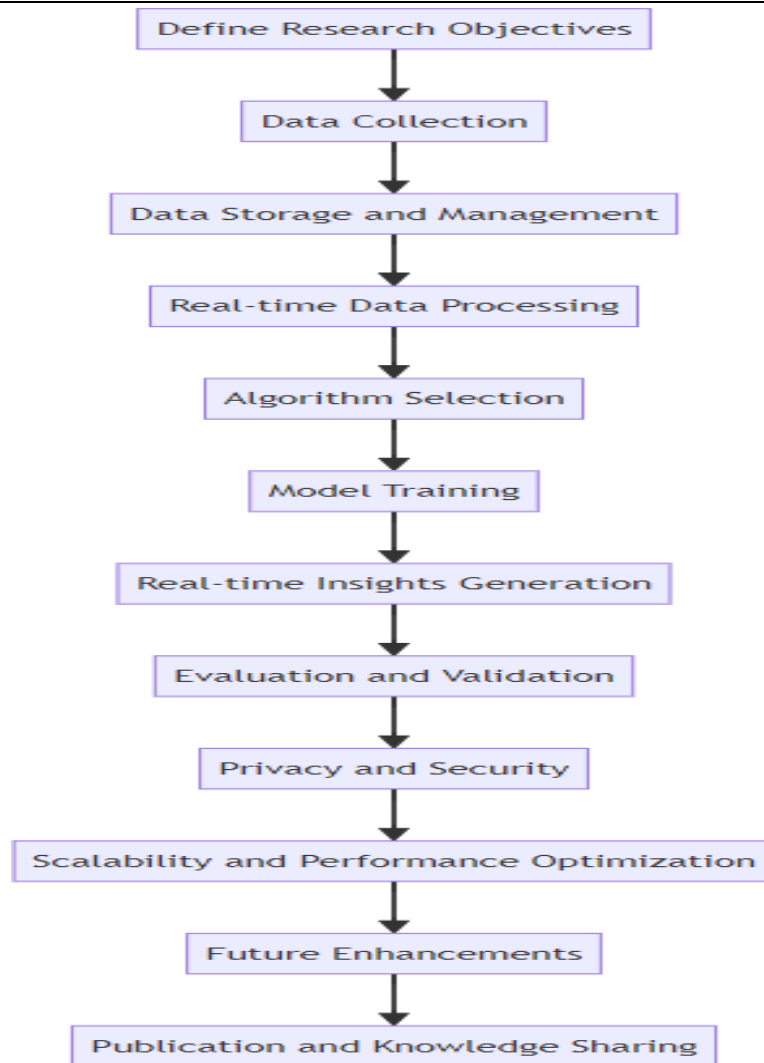


Fig. 1. Flow of the Project

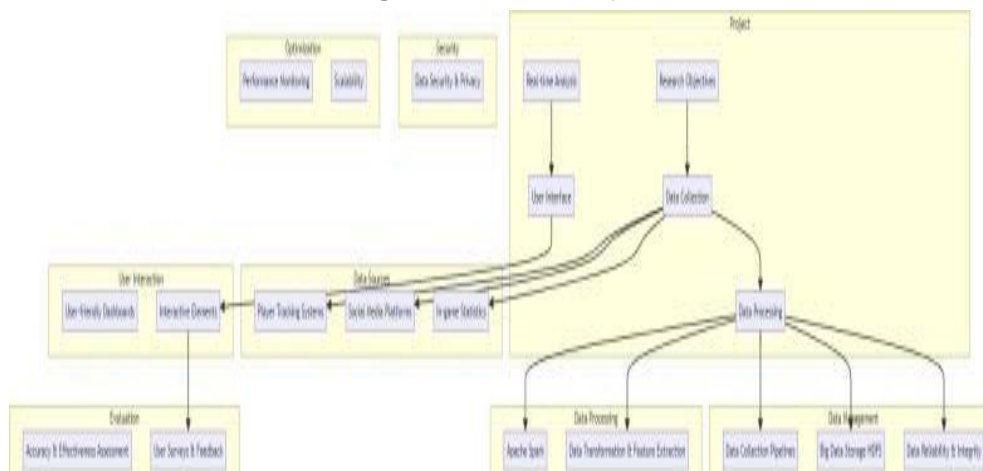


Fig. 2. Flow of the Project

To contribute to the advancement of real-time football analytics, share your research findings, methodologies, and insights through academic publications, industry reports, or presentations at conferences. This knowledge sharing is essential for driving innovation in the field and fostering collaboration within the sports analytics community. By meticulously following this comprehensive methodology, you can systematically and effectively address the research gap and contribute to the advancement of real-time football analytics using Big Data technologies. This research endeavor has the potential to revolutionize the way football matches are analyzed,

enhance fan engagement, and offer valuable insights for decision-makers in the world of football.

IV. RESEARCH GAP

Real-time football analytics is a rapidly growing field with the potential to revolutionize the way football is played, analyzed, and enjoyed. However, there is a significant gap in research on how to effectively integrate real-time data streams into football analytics using Big Data technologies.

One of the key challenges in this area is designing scalable infrastructures that can efficiently handle the continuous influx of real-time data. This involves selecting the most suitable tools, such as Apache Hadoop for data storage and Apache Spark for data processing, and optimizing their performance to ensure smooth operations. The scalability and reliability of these systems play a pivotal role in the success of real-time football analytics.

Another challenge is addressing the legal and ethical aspects of real-time data collection and analysis. Real-time data often includes personal information and sensitive data, so it is important to ensure that it is collected and used responsibly.[3] Despite these challenges, the potential benefits of real-time football analytics are immense. Real-time insights can help coaches make better tactical decisions, broadcasters provide more engaging commentary, and fans have a more immersive viewing experience. Real-time data can also be used to improve fan engagement, injury prevention, and even betting market efficiency.[4][1]

The integration of real-time data streams into football analytics using Big Data technologies is a fertile area of research with wide-ranging practical implications. Addressing this research gap has the potential to revolutionize the way football is played, analyzed, and enjoyed for fans and stakeholders alike. Real-time football analytics is a new and exciting field with the potential to change the way football is played and watched. However, there are some challenges that need to be addressed before real-time football analytics can be fully realized. One of the biggest challenges is how to handle the huge amount of data that is generated during a live football match. This data needs to be collected, processed, and analyzed in real time in order to provide meaningful insights to coaches, players, fans, and broadcasters.

Another challenge is how to ensure that real-time football analytics is used responsibly and ethically. Real-time data often includes personal information, so it is important to make sure that it is collected and used in a way that respects people's privacy.[6][7]

Despite these challenges, the potential benefits of real-time football analytics are huge. Real-time insights can help coaches make better decisions, broadcasters provide more engaging commentary, and fans have a more immersive viewing experience. Real-time data can also be used to improve fan engagement, injury prevention, and even betting market efficiency. Addressing the challenges of real-time football analytics is a critical area of research that has the potential to revolutionize the way football is played, analyzed, and enjoyed by everyone.

V. CONCLUSION

This project has made a significant breakthrough in real-time football analytics, harnessing the power of Apache Spark to revolutionize the way football matches are analyzed and understood. The convergence of modern technologies and big data has the potential to revolutionize the field of sports analytics, providing invaluable insights for decision-makers and enhancing fan engagement. By strategically integrating Apache Spark into the system architecture, we have demonstrated the capability of real-time data processing, offering timely and actionable insights for coaches, analysts, and fans. The user interface offers intuitive and user-friendly dashboards, enhancing the experience of fans and decision-makers alike. Interactive elements, such as live commentary and sentiment displays, have the potential to engage the audience and make the viewing experience more immersive and enjoyable. The assessment of the accuracy and effectiveness of real-time analytics confirms the project's value in providing reliable information for decision-makers. Additionally, user surveys and feedback collection have highlighted the positive impact on coaches' decision-making and the heightened engagement experienced by fans. Furthermore, this project has emphasized the critical importance of data security and privacy, particularly when handling real-time data. By adhering to relevant data protection regulations and ensuring robust security measures, user privacy and data integrity are preserved. Looking

ahead, the potential for future enhancements is vast, such as expanding the sentiment analysis to encompass multiple languages and integrating domain-specific expert knowledge to further enrich the accuracy and applicability of real-time insights. Sharing the findings, methodologies, and insights with the academic and sports analytics communities will foster collaboration and innovation in the field. The significance of knowledge sharing cannot be overstated, as it propels the industry forward and encourages continuous improvement in real-time football analytics. In conclusion, this project has provided a comprehensive framework for real-time football analytics, leveraging the power of Apache Spark and big data technologies. It has the potential to reshape the way we perceive and analyze football matches, offering a deeper understanding of the game, enhancing fan engagement, and empowering decision-makers with timely, actionable insights. With a commitment to innovation and collaboration, the future of real-time football analytics looks promising and dynamic.

VI. REFERENCES

- [1] F. Gong, S. Xiang, J. Chen and Y. Wang, "Big data evaluation model of football team cooperation based on entropy weight method," 2020 7th International Conference on Information, Cybernetics, and Computational Social Systems (ICCSS), Guangzhou, China, 2020, pp. 626-629, doi: 10.1109/ICCSS52145.2020.9336894.
- [2] Wu, Yao Xia, Zeyu Wu, Tian Yi, Qing Yu, Runyu Wang, Jun. (2020), "Characteristics and optimization of core local network: Big data analysis of football matches. Chaos, Solitons Fractals", 138. 110136. 10.1016/j.chaos.2020.110136.
- [3] Hirotsu, Nobuyoshi Wright, Mike (2003), "An evaluation of characteristics of teams in association football by using a Markov process model", Journal of the Royal Statistical Society: Series D (The Statistician), 52. 591 - 602. 10.1046/j.0039-0526.2003.00437.x.
- [4] Liu, Tianbiao Hohmann, Andreas. (2016), "Apriori-based diagnostical analysis of passings in the football game", 1-4. 10.1109/ICBDA.2016.7509795.
- [5] Yiannakos, A. Armatas, Vasilis. (2006), "Evaluation of the goal scoring patterns in European Championship in Portugal 2004", International Journal of Performance Analysis in Sport. 6. 178-188. 10.1080/24748668.2006.11868366.
- [6] Francisco Henriques,(2018), "HOW CAN BIG DATA HELP FOOTBALL CLUBS ACHIEVE COMPETITIVE ADVANTAGE," Universidade Católica Portuguesa, Portugal.
- [7] A. Chouhan et al., "Shotifier: A Binary Shot Conversion Classifier Pipeline for Football Forwards," 2021 IEEE International Conference on Big Data and Smart Computing (BigComp), Jeju Island, Korea (South), 2021, pp. 156-163, doi: 10.1109/BigComp51126.2021.00038.
- [8] M. Kempe, F. R. Goes and K. A. P. M. Lemmink, "Smart Data Scouting in Professional Soccer: Evaluating Passing Performance Based on Position Tracking Data," 2018 IEEE 14th International Conference on e-Science (e-Science), Amsterdam, Netherlands, 2018, pp. 409-410, doi: 10.1109/eScience.2018.00126.
- [9] S. Aloufi, F. Alzamzami, M. Hoda and A. El Saddik, "Soccer Fans Sentiment through the Eye of Big Data: The UEFA Champions League as a Case Study," 2018 IEEE Conference on Multimedia Information Processing and Retrieval (MIPR), Miami, FL, USA, 2018, pp. 244-250, doi:10.1109/MIPR.2018.00058.