# A Vision in Explainable AI (XAI)

**Gurpreet Singh[1], Brahmleen Kaur[1], Satinder Kaur[1*], Satveer Kour[1], Mehakdeep Kaur[1], Kumari Sarita[2]**

Department of Computer Engineering and Technology, Guru Nanak Dev University, Amritsar[1]

Department of Computer Science, Guru Nanak Dev University, Amritsar[2]

**Abstract:** Artificial Intelligence (AI) has transformed industries and everyday life with its ability to automate complex tasks and make predictions based on large datasets. However, one of the biggest challenges with AI, particularly with advanced models such as deep learning, is the lack of transparency. These models, often referred to as "black boxes," provide predictions and decisions, but the reasoning behind them is not immediately clear to users. This lack of interpretability has led to the development of Explainable AI (XAI), a set of techniques that aim to make AI systems more transparent, understandable, and trustworthy. XAI is crucial for building confidence in AI, especially in high-stakes areas like healthcare, finance, law, and autonomous vehicles. The aim of this work is to provide a comprehensive guide that delves into the components, methods and techniques, future scope and applications of XAI. It concludes by providing a detailed understanding about XAI, how it enhances AI models by making them more interpretable and accountable. So, it provides a new vision to researchers how they can justify decision making and results with AI.

**Keywords:** LIME, SHAP, Post-Hoc Explainability, Intrinsic Explainability

## 1. INTRODUCTION

Artificial Intelligence (AI) has transformed industries and everyday life with its ability to automate complex tasks and make predictions based on large datasets. However, one of the biggest challenges with AI, particularly with advanced models such as deep learning, is the lack of transparency [1]. These models, often referred to as "black boxes," provide predictions and decisions, but the reasoning behind them is not immediately clear to users. This lack of interpretability has led to the development of Explainable AI (XAI), a set of techniques that aim to make AI systems more transparent, understandable, and trustworthy [2]. XAI is crucial for building confidence in AI, especially in high-stakes areas like healthcare, finance, law, and autonomous vehicles [3]. The following sections delve into the components, methods and techniques, future scope and applications of XAI, thus, providing a detailed understanding of how it enhances AI models by making them more interpretable and accountable.

## 2. COMPONENTS OF EXPLAINABLE AI (XAI)

XAI can be broken down into several fundamental components, including the distinctions between interpretability and explainability, as well as global and local explanations [4]. Each of these components plays a vital role in providing clarity about AI's decision-making process, fostering trust and accountability in the models. Here's an expanded and detailed table covering the components of Explainable AI (XAI):

| Component | Definition | Purpose | Key Characteristics | Examples & Techniques | Use Cases |
|---|---|---|---|---|---|
| **Interpretability[5]** | The degree to which a human can understand or make sense of an AI model's internal workings. | Helps users grasp the model's decision-making process by simplifying | - Focuses on transparency of model structure.- Easier in simple models like decision trees.- Directly shows | - Decision Trees- Linear Regression- Feature Importance (e.g., weights in linear models) | - Finance (credit scoring models)- Healthcare (interpreting lab test results)- Autonomous systems |

| | | its logic and structure. | relationships between input and output. | | (interpretable control systems) |
|---|---|---|---|---|---|
| **Explainability[5]** | Providing understandable and coherent explanations for how and why an AI model made a specific decision. | Offers insights into the model's decision-making process, even if the model itself remains complex. | - Used when models are complex ("black box" models).- Post-hoc techniques explain predictions.- Aims to improve trust and accountability. | - LIME (Local Interpretable Model-agnostic Explanations) - SHAP (SHapley Additive exPlanations)- Attention Mechanisms (in deep learning) | - AI-driven hiring tools (explaining hiring decisions)- Healthcare (interpreting medical diagnoses)- E-commerce (personalized recommendations) |
| **Global Explanations[5-6]** | Provides an overarching understanding of how a model behaves as a whole. | Helps users understand general relationships between inputs and outputs across the entire dataset. | - Focuses on overall patterns in the model.- Identifies important features influencing predictions.- Used for debugging and policy compliance. | - Feature Importance Analysis- Surrogate Models (e.g., decision tree approximation of a complex model)- Decision Rules Extraction | - Regulatory compliance (ensuring fair AI decisions)- Business intelligence (understanding market trends)- Scientific research (analyzing AI-driven discoveries) |
| **Local Explanations[5-6]** | Focuses on individual predictions, explaining why a model made a particular decision for a specific input. | Ensures transparency and accountability in individual decisions, especially in critical applications like healthcare and finance. | - Provides instance-specific insights.- Helps users validate AI decisions.- Useful in sensitive applications (e.g., loan approvals, medical diagnoses). | - LIME (explains single predictions)- SHAP (assigns importance scores to features for an individual prediction)- Counterfactual Explanations (shows alternative inputs that would change the decision) | - Loan approvals (explaining rejected applications)- Fraud detection (understanding flagged transactions)- Medical diagnosis (explaining why a disease was predicted) |

## 3. METHODS AND TECHNIQUES IN EXPLAINABLE AI

As AI systems grow increasingly complex, ensuring that these systems' decisions are transparent and explainable is critical for both trust and accountability [7]. Several methods and techniques have been developed to enhance the

interpretability of AI models, and they can be categorized into two primary groups: intrinsic explainability and post-hoc explainability [8]. Both approaches play an essential role in making AI systems more accessible to users, stakeholders, and domain experts.

## 3.1 Intrinsic Explainability

Intrinsic explainability refers to the creation of AI models that are inherently interpretable by design. These models are built in a way that their internal workings are transparent, meaning that their decisions are easy to follow and understand without the need for additional tools or methods. Models that exhibit intrinsic explainability are typically simpler and less complex, making them more suitable for tasks where transparency is a priority.

**3.1.1 Linear Models**: Linear models [9], such as linear regression and logistic regression, are classic examples of intrinsically interpretable models. These models are widely used because of their simplicity and the clear relationship they establish between input features and output predictions. For example, in linear regression, the model's output is a weighted sum of the input features, where the coefficients of the features represent the strength of their influence on the prediction. The relationship between each input and output is explicit and can be easily described in a simple mathematical equation, making it straightforward for humans to understand how predictions are made. This simplicity enables transparency and makes it easier to interpret the model's decision-making process. Linear models are often favoured in scenarios where model transparency is critical, but they may not perform as well in complex tasks that require capturing non-linear relationships in the data.
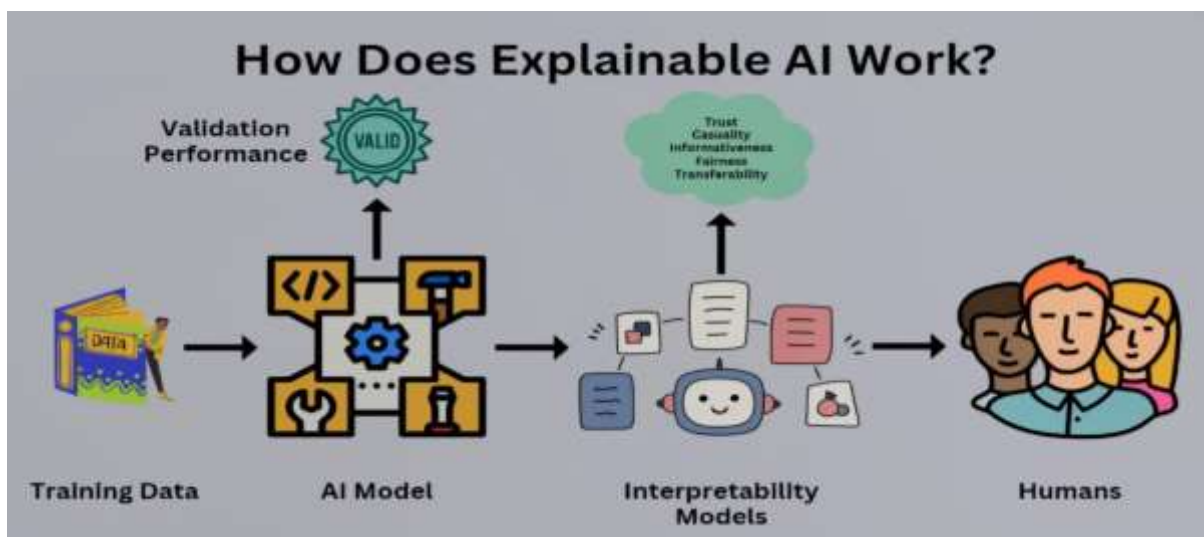


Fig.1. Working of Explainable AI

**3.1.2 Decision Trees**: Decision trees [10] are another example of models with high intrinsic explainability. A decision tree is a flowchart-like structure where each node in the tree represents a decision based on a particular feature or condition. At each branching point, the model splits the data according to the most relevant feature, with each split leading to a decision or outcome at the leaf nodes. The decision-making process can be easily followed by tracing the path taken from the root node to the leaf node. This structure allows for clear reasoning behind each decision made by the model. Because the model's structure is visual and intuitive, it is easier for humans to understand the logic behind the model's predictions. Decision trees are particularly useful in applications where the rationale behind each decision needs to be explicit and understandable, such as in healthcare or financial services [11]. However, they can sometimes be prone to overfitting, especially in the case of deep trees.

**3.1.3 Rule-Based Models**: Rule-based models are another form of intrinsically interpretable AI. These models make decisions based on a set of predefined rules, often expressed in simple "if-then" statements. For example, in medical diagnostic systems, an AI model might use a series of rules to diagnose a disease based on observed symptoms [12]. If a patient presents with a high fever and a cough, the model may use a rule such as "If fever > 100°F and cough = yes, then diagnose as flu." These rules are explicitly defined, making the model's decision process easy to explain to non-experts. Rule-based models are transparent and allow users to directly understand the reasoning behind a model's decision. However, they can be limited in their ability to handle more complex, nuanced data, and they may require significant manual effort to define the rules.

### 3.2 Post-Hoc Explainability

While intrinsically interpretable models are valuable, many state-of-the-art machine learning models—such as deep learning networks, random forests, and support vector machines—are considered "black-box" models due to their complexity and lack of transparency. These models often provide high accuracy but are difficult to interpret directly. As a result, post-hoc explainability techniques [13] have been developed to add interpretability to these black-box models after they have been trained. Post-hoc methods help to generate human-readable explanations of why a model made a particular decision or prediction, even if the internal workings of the model are not easily understandable.

**3.2.1 LIME (Local Interpretable Model-agnostic Explanations)**: LIME is one of the most widely used post-hoc explanation techniques [14]. The central idea behind LIME is to generate local explanations for individual predictions made by complex, black-box models. LIME works by approximating the complex model with a simpler, interpretable model, such as a linear regression, for a specific instance or input. To do this, LIME perturbs the input data slightly and observes how the model's predictions change, generating a set of new instances that are similar to the original data point. The simpler model is then trained on these perturbed instances, allowing LIME to approximate the decision-making process of the complex model in that particular region of the input space [13]. The result is an interpretable explanation of which features were most influential in the decision for that specific instance. LIME is useful for understanding individual predictions, which is crucial for applications like credit scoring or medical diagnosis, where users need to know why a specific decision was made for them.

**3.2.2 SHAP (SHapley Additive exPlanations)**: SHAP [14] is another popular technique for post-hoc explanation, and it is based on cooperative game theory. SHAP values provide a way to understand the contribution of each feature to a model's prediction by treating the model as a cooperative game where each feature "contributes" to the final outcome. The technique calculates the Shapley values for each feature, which measure the average contribution of each feature across all possible combinations of features. SHAP values provide a unified and consistent way to explain the importance of features, both locally for individual predictions and globally for the overall behaviour of the model.
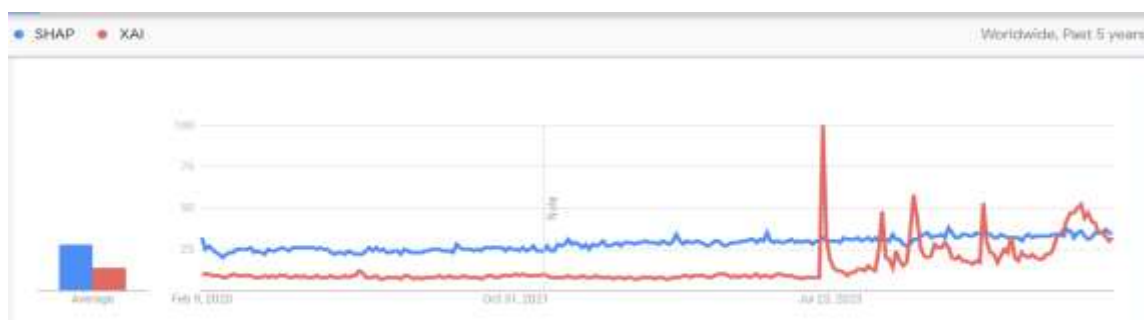


Fig. 2. SHAP in Google Trends

SHAP is particularly useful for complex models like gradient-boosted trees and deep neural networks, as it can provide both feature importance scores and localized explanations for specific predictions. The primary advantage of SHAP over other post-hoc methods is its theoretical foundation in game theory, which ensures that the feature importance scores are fair and consistent [15].

**3.2.3 Feature Importance**: In many machine learning models, particularly tree-based models like random forests or gradient boosting machines, it is possible to compute feature importance to understand which features had the greatest influence on the model's predictions [16]. Feature importance techniques assess how much each feature contributes to the reduction in error or impurity (for decision trees) across all the data points in the dataset. Methods like permutation feature importance involve randomly permuting the values of each feature and measuring the impact on the model's performance, providing a way to quantify the importance of each feature. Gradient-based methods, such as integrated gradients, are another technique used to calculate feature importance in neural networks. By identifying the most influential features, feature importance methods provide valuable insights into how the model arrives at its predictions, helping users to understand the decision-making process better.

**3.2.4 Surrogate Models**: Surrogate models [17] are interpretable models used to approximate the behaviour of more complex models. After training a black-box model, a simpler surrogate model, such as a decision tree or linear regression, is trained on the same input-output data that the complex model has used. The surrogate model serves as an interpretable approximation of the black-box model's decision-making process. Surrogate models can help to explain the overall behaviour of a complex model by providing a simpler, more transparent model that can be understood by humans.

However, surrogate models are not always a perfect representation of the original model, as they may not capture all the intricacies of the black-box model's behaviour. Despite this, surrogate models are still widely used as a tool for making black-box models more interpretable.

## 3.3 Visualization Techniques

Visualization plays a crucial role in Explainable AI (XAI) by helping researchers and users interpret complex machine learning models [18]. These techniques provide insights into how a model makes decisions, making AI systems more transparent and trustworthy. Below are three widely used visualization methods: Attention Maps, Feature Importance Scores, and Activation Maximization.

**3.3.1 Attention Maps:** Attention maps are a powerful tool used to visualize which parts of the input data an AI model focuses on when making predictions. They are widely used in computer vision and natural language processing (NLP) models, particularly in deep learning architectures like transformers and convolutional neural networks (CNNs).

In computer vision, attention maps highlight the regions of an image that contribute the most to a model's classification decision. For example, in an image recognition task where the model classifies a dog, the attention map may show that the network is primarily focusing on the dog's face or body rather than irrelevant background elements. Grad-CAM (Gradient-weighted Class Activation Mapping) is a commonly used technique for generating such visual explanations.

In NLP, attention maps are used in transformer models like BERT and GPT to show which words in a sentence are most relevant to a given prediction. For instance, in a sentiment analysis task, an attention map might highlight negative words such as "terrible" or "disappointing" when classifying a review as negative. By visualizing these attention scores, users can better understand how the model interprets textual data.

**3.3.2 Feature Importance Scores:** Feature importance scores help identify which input variables (features) have the most significant influence on a model's predictions [19]. This technique is especially useful for interpreting complex machine learning models like decision trees, random forests, and deep neural networks.

There are several ways to compute feature importance. In tree-based models like Random Forests and XGBoost, importance scores are often derived from how frequently a feature is used for splitting and how much it reduces impurity in decision nodes. Higher scores indicate that a feature contributes more significantly to predictions.

For black-box models like deep neural networks, techniques like SHapley Additive exPlanations (SHAP) and Local Interpretable Model-agnostic Explanations (LIME) provide feature importance insights by approximating the model's decision boundary. These scores allow data scientists to verify whether the model is making logical decisions. For example, in a credit approval system, a high feature importance score for "income level" over "age" may indicate the model is prioritizing financial stability rather than age-related bias.

**3.3.3 Activation Maximization:** Activation Maximization [20] is a technique used to understand what patterns in input data cause neurons in deep neural networks to fire strongly. This helps in visualizing what a model has learned at different layers, providing insights into how it processes information.

Activation Maximization generates synthetic images that maximize the response of specific neurons. For example, if a neuron in a CNN is responsible for detecting "cat-like" features, the activation maximization process might generate a visual pattern resembling a cat's fur texture or facial features. This technique is widely used in a visualization method developed by Google that enhances neural network interpretations by amplifying the patterns it recognizes.

In, activation maximization can help understand how individual words or phrases activate different layers of a neural network. By studying these activations, researchers can diagnose potential biases or identify whether the model is relying on meaningful linguistic structures.

## 4. APPLICATIONS OF EXPLAINABLE AI

Explainable AI is being applied in a wide range of fields where transparency and accountability are paramount. Some key applications include:

### 4.1 Healthcare

In healthcare, AI models are increasingly used for medical diagnosis, drug discovery, and treatment recommendations [16]. It is essential that healthcare professionals understand how AI systems make these decisions in order to verify the accuracy and reliability of their recommendations. XAI techniques, such as LIME and SHAP, are being used to explain predictions made by models in medical imaging, clinical decision support systems, and more.

### 4.2 Finance

AI is widely used in finance for tasks like fraud detection, credit scoring, and risk management [17]. Since these decisions can significantly affect people's financial well-being, it is crucial that they are explainable. Financial institutions are using

XAI to ensure that AI systems are making fair, unbiased decisions and that customers can understand the reasoning behind their credit approvals or loan rejections.

## 4.3 Autonomous Vehicles

Autonomous vehicles rely on AI models to make real-time decisions related to navigation, safety, and traffic management [18]. Understanding how these models make decisions is critical, particularly in the case of accidents or safety violations. XAI can help make the decision-making process of autonomous vehicles more transparent and accountable.

## 4.4 Legal and Criminal Justice Systems

AI is being used in the legal field for tasks like predicting recidivism risk and assisting in legal research. In these contexts, explainability is essential to ensure that decisions are fair and free from biases. XAI can help explain how AI systems arrive at decisions in areas like sentencing, parole, and bail.

Hence, Explainable AI (XAI) is a rapidly evolving field that addresses the transparency challenges posed by complex machine learning models [19]. With its focus on improving the interpretability and accountability of AI systems, XAI is helping build trust and confidence in AI technologies, particularly in critical domains like healthcare, finance, and autonomous systems. While challenges such as balancing performance with interpretability remain, the advancements in XAI techniques like LIME, SHAP, and surrogate models are paving the way for more understandable AI systems. As AI continues to shape the future, ensuring that its decisions are explainable, ethical, and fair will be crucial for ensuring its responsible deployment across various industries.



Fig.3. Applications of XAI

## 5.    FUTURE SCOPE OF EXPLAINABLE AI (XAI)

Explainable AI (XAI) is set to revolutionize the future of artificial intelligence by making AI models more transparent, trustworthy, and ethical. As AI systems become more complex and integrated into critical sectors like healthcare, finance, law, and defence, the need for explainability grows exponentially [20]. In healthcare, for example, XAI can help doctors understand AI-driven diagnoses, ensuring that medical decisions are based on logical and justifiable reasons rather than black-box predictions. Similarly, in finance, XAI can enhance risk assessment models, allowing institutions to explain why a loan was approved or denied, thereby improving fairness and compliance with regulations.

Moreover, XAI will play a crucial role in legal and ethical AI governance. Governments and regulatory bodies are increasingly demanding transparency in AI decision-making, particularly in areas like hiring, criminal justice, and autonomous systems. Explainable AI can help organizations comply with policies such as the EU's General Data Protection Regulation (GDPR), which mandates that AI-driven decisions be interpretable. In addition, XAI can reduce biases in AI models by enabling developers to identify and rectify unfair patterns, leading to more equitable outcomes across different demographic groups.

Another significant future application of XAI is in autonomous systems, such as self-driving cars and robotic automation. These systems must make real-time decisions that impact human lives, and having an explainable framework will ensure that their actions can be audited and trusted. In cybersecurity, XAI can help analysts understand why certain threats are flagged, improving response strategies against cyber-attacks. As AI continues to evolve, explainability will be a key factor in gaining public trust, fostering responsible AI development, and ensuring that AI serves humanity in an ethical and transparent manner.

## 6. CONCLUSION

Explainable AI (XAI) is a critical advancement in the field of artificial intelligence, addressing the need for transparency, trust, and accountability in AI-driven decision-making. As AI systems become increasingly complex and influential across various industries, the importance of understanding their reasoning and ensuring fairness cannot be overstated. By distinguishing between interpretability and explainability, XAI enables users to comprehend both the overall behaviour of AI models and the specific factors influencing individual predictions. This capability is particularly vital in high-stakes applications such as healthcare, finance, autonomous vehicles, and the legal system, where decisions must be justifiable and free from bias.

The future of XAI is promising, with its potential to shape AI governance, regulatory compliance, and ethical AI development. As organizations and governments push for greater transparency in AI, explainability will become a standard requirement rather than an optional feature. Emerging techniques like LIME, SHAP, and surrogate models are paving the way for more interpretable AI systems, allowing stakeholders to validate AI-generated decisions. Moreover, XAI will play a crucial role in improving AI-driven automation, cybersecurity, and risk assessment, ensuring that AI technologies are not only powerful but also responsible and aligned with human values. By fostering greater trust and understanding, Explainable AI will help bridge the gap between human intuition and machine intelligence, ultimately leading to more ethical and effective AI applications in the future.

## REFERENCES

[1]. Holzinger, A., Langs, G., Denk, H., Zatloukal, K., & Müller, H. (2019). Causability and explainability of artificial intelligence in medicine. WIREs Data Mining and Knowledge Discovery, 9(4). https://doi.org/10.1002/widm.1312

[2]. Marrone, R., Cropley, D., & Medeiros, K. (2024). How does narrow AI impact human creativity? Creativity Research Journal,

[3]. 1–11. https://doi.org/10.1080/10400419.2024.2378264

[4]. Mirbabaie, M., Stieglitz, S., Brünker, F., Hofeditz, L., Ross, B., & Frick, N. R. J. (2020). Understanding collaboration with virtual assistants – the role of social identity and the extended self. Business &amp; Information Systems Engineering, 63(1), 21–37. https://doi.org/10.1007/s12599-020-00672-x

[5]. Ko, H., Lee, S., Park, Y., & Choi, A. (2022). A survey of recommendation systems: Recommendation models, techniques, and application fields. Electronics, 11(1), 141. https://doi.org/10.3390/electronics11010141

[6]. Szegedy, C. (2020). A promising path towards autoformalization and general artificial intelligence. In Lecture Notes in Computer Science (pp. 3–20). Springer International Publishing. https://doi.org/10.1007/978-3-030-53518-6_1

[7]. Poli, R. (2023). Super-Human and super-ai cognitive augmentation of human and human-ai teams assisted by brain computer interfaces. Proceedings of the Genetic and Evolutionary Computation Conference, 3–3. https://doi.org/10.1145/3583131.3603554

[8]. Hsu, F. (2022). Behind Deep Blue. https://doi.org/10.1515/9780691235141

[9]. Cheng, L., & Yu, T. (2019). A new generation of AI: A review and perspective on machine learning technologies applied to smart energy and electric power systems. International Journal of Energy Research, 43(6), 1928–1973. https://doi.org/10.1002/er.4333

[10]. Gupta, A., Anpalagan, A., Guan, L., & Khwaja, A. S. (2021). Deep learning for object detection and scene perception in self-driving cars: Survey, challenges, and open issues. Array, 10, 100057. https://doi.org/10.1016/j.array.2021.100057

[11]. Langley, C., Cirstea, B. I., Cuzzolin, F., & Sahakian, B. J. (2022). Theory of mind and preference learning at the interface of cognitive science, neuroscience, and AI: A review. Frontiers in Artificial Intelligence, 5. https://doi.org/10.3389/frai.2022.778852

[12]. Chowdhary, K. R. (2020). Natural language processing. In Fundamentals of Artificial Intelligence (pp. 603–649). Springer India. https://doi.org/10.1007/978-81-322-3972-7_19

[13]. Watchus, B. (2024). Towards self-aware AI: Embodiment, feedback loops, and the role of the insula in consciousness. MDPI AG. https://doi.org/10.20944/preprints202411.0661.v1

[14]. Dhruv, A. J., Patel, R., & Doshi, N. (2020). Python: The most advanced programming language for computer science applications. Proceedings of the International Conference on Culture Heritage, Education, Sustainable Tourism, and Innovation Technologies, 292–299. https://doi.org/10.5220/0010307902920299

[15]. Giorgi, F. M., Ceraolo, C., & Mercatelli, D. (2022). The R language: An engine for bioinformatics and data science. Life, 12(5), 648. https://doi.org/10.3390/life12050648

[16]. Gao, K., Mei, G., Piccialli, F., Cuomo, S., Tu, J., & Huo, Z. (2020). Julia language in machine learning: Algorithms, applications, and open issues. Computer Science Review, 37, 100254. https://doi.org/10.1016/j.cosrev.2020.100254

[17]. Zan, T., & Hu, Z. (2023). VoiceJava: A syntax-directed voice programming language for java. Electronics, 12(1), 250. https://doi.org/10.3390/electronics12010250

[18]. Kang, Y., Cai, Z., Tan, C.-W., Huang, Q., & Liu, H. (2020). Natural language processing (NLP) in management research: A literature review. Journal of Management Analytics, 7(2), 139–172. https://doi.org/10.1080/23270012.2020.1756939

[19]. Yogish, D., Manjunath, T. N., & Hegadi, R. S. (2019). Review on natural language processing trends and techniques using NLTK. In Communications in Computer and Information Science (pp. 589–606). Springer Singapore. https://doi.org/10.1007/978-981-13-9187-3_53

[20]. Ofoeda, J., Boateng, R., & Effah, J. (2019). Application programming interface (API) research. International Journal of Enterprise Information Systems, 15(3), 76–95. https://doi.org/10.4018/ijeis.2019070105