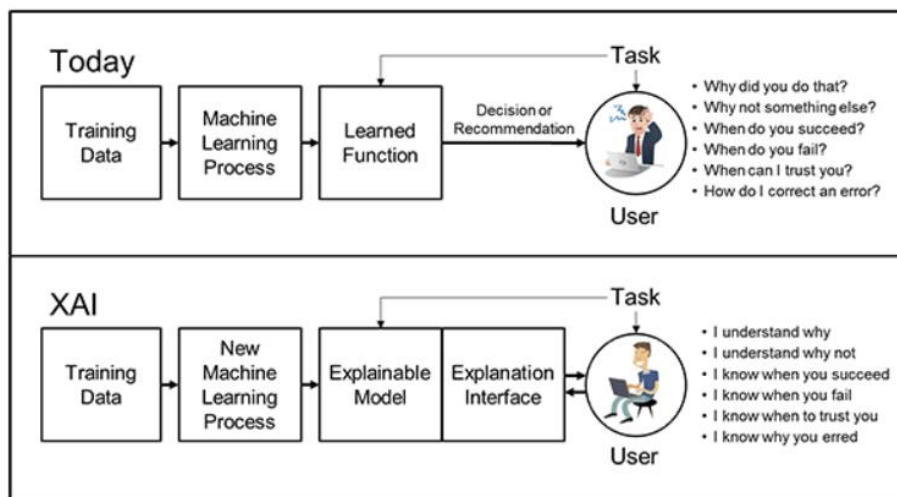


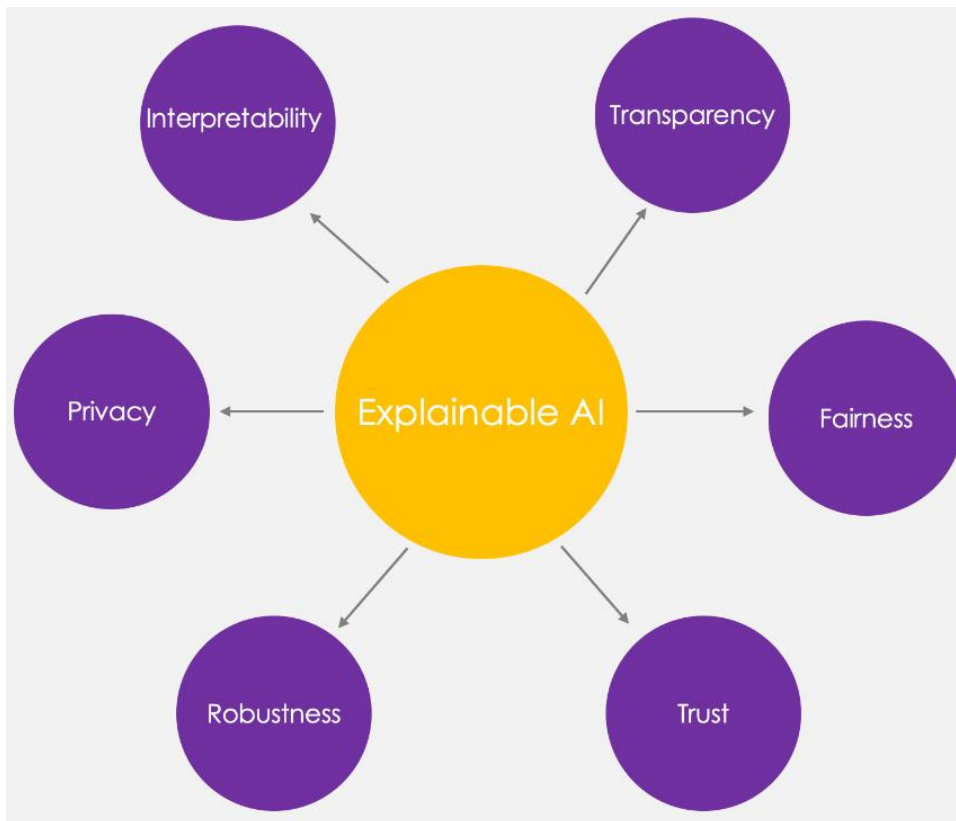
# Explainable AI

- ❖ Explainable Artificial Intelligence (XAI) encompasses a set of procedures and techniques.
- ❖ Its primary purpose is to facilitate the production of understandable and reliable results from machine learning algorithms, catering to human users.
- ❖ XAI serves as a pivotal element within the fairness, accountability, and transparency (FAT) paradigm of machine learning.
- ❖ It addresses the need for transparency and interpretability in AI systems.
- ❖ XAI is frequently discussed in the context of deep learning, which often involves complex and intricate neural network architectures.
- ❖ Its application in deep learning aims to make these sophisticated models more interpretable.
- ❖ Organizations deploying AI seek to establish trust in the technology they utilize.
- ❖ XAI plays a vital role in this process, allowing users to comprehend the behavior of AI models.



## Why Explainable AI is needed?

- **Transparency.** Ensuring stakeholders understand the models' decision-making process.
- **Fairness.** Ensuring that the models' decisions are fair for everyone, including people in protected groups (race, religion, gender, disability, ethnicity).
- **Trust.** Assessing the confidence level of human users using the AI system.
- **Robustness.** Being resilient to changes in input data or model parameters, maintaining consistent and reliable performance even when faced with uncertainty or unexpected situations.
- **Privacy.** Guaranteeing the protection of sensitive user information.
- **Interpretability.** Providing human-understandable explanations for their predictions and outcomes.



## **Benefits of explainable AI (XAI)**

### **Improved Decision-making:**

- Explainable AI enhances decision-making by providing insights into influential factors in model predictions.
- Identifying key factors helps prioritize actions and strategies for achieving desired outcomes.

### **Increased Trust and Acceptance:**

- Explainable AI builds trust and acceptance in machine learning models.
- Overcoming the opacity of traditional models, it accelerates adoption by making models more understandable and reliable.

### **Reduced Risks and Liabilities:**

- Explainable AI mitigates risks and liabilities associated with machine learning models.
- It establishes a framework for addressing regulatory and ethical considerations, minimizing potential negative impacts.

## **Facilitates Regulatory Compliance**

- Explainable AI aligns with regulatory requirements, providing a transparent framework.
- This ensures compliance with regulations and ethical standards, reducing legal and regulatory risks.

## **How does Explainable AI work?**

### **Machine Learning Model:**

- Core component representing the underlying algorithms for predictions.
- Utilizes various machine learning techniques (supervised, unsupervised, or reinforcement learning).
- Applicable in diverse domains such as medical imaging, natural language processing, and computer vision.

### **Explanation Algorithm:**

- Key component providing insights into factors influencing model predictions.
- Employs explainable AI approaches like feature importance, attribution, and visualization.
- Offers valuable insights into the decision-making process of the machine learning model.

Interface:

### **Component responsible for presenting insights to human users.**

- Can take the form of web applications, mobile apps, or visualizations.
- Offers a user-friendly and intuitive interface for accessing and interacting with explainable AI-generated insights.

## **Explainable AI (XAI) Techniques**

### **LIME (Local Interpretable Model-agnostic Explanations):**

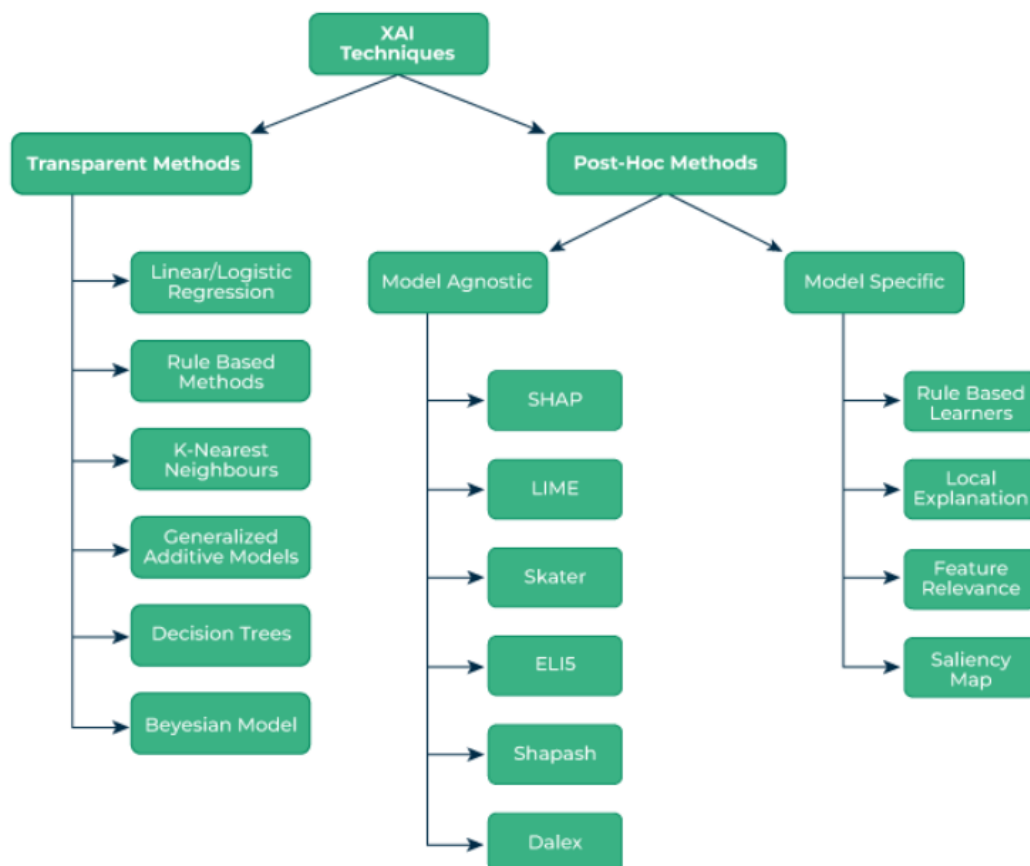
- Approach: LIME is a widely used XAI technique that creates a local model approximation to explain the interpretability of a global model.
- Implementation in Python: Utilize the 'lime' package, which offers tools and functions for generating and interpreting LIME explanations.
- Purpose: LIME provides insights into the factors most relevant in the local context of a specific prediction.

### SHAP (SHapley Additive exPlanations):

- Approach: SHAP leverages Shapley values from game theory to offer interpretable insights into the factors influencing a model's predictions.
- Implementation in Python: The 'shap' package provides a comprehensive set of tools and functions for implementing and interpreting SHAP explanations.
- Purpose: SHAP explanations reveal the contribution of each feature to the overall model prediction.

### ELI5 (Explain Like I'm 5):

- Approach: ELI5 simplifies complex model explanations using an intuitive language, making them accessible to non-experts.
- Implementation in Python: The 'eli5' package offers tools and functions for generating and interpreting ELI5 explanations in a user-friendly manner.
- Purpose: ELI5 provides interpretable insights into the relevant factors influencing model predictions in a straightforward manner.



## **Current Limitations of XAI:**

### **Computational Complexity:**

- XAI approaches can be computationally complex, demanding significant resources and processing power.
- Challenges real-time and large-scale applications, restricting deployment in these contexts.

### **Limited Scope and Domain-Specificity:**

- Many XAI methods are domain-specific and may not be universally applicable.
- Limited scope poses challenges in diverse machine learning models and applications.

### **Lack of Standardization and Interoperability:**

- Lack of standardization in the XAI field hinders interoperability.
- Different XAI approaches using varied metrics and algorithms make comparison and evaluation challenging.

## **Explainable AI Case Studies:**

### **Medical Imaging:**

- XAI used to diagnose diseases, visualizing factors influential in the diagnostic process.
- Identifies important features in cancer diagnosis, enhancing interpretability.

### **Natural Language Processing:**

- XAI applied to interpret and analyze text, revealing influential words in sentiment analysis.
- Provides insights into factors predictive of positive or negative sentiment.

### **Computer Vision:**

- XAI employed in computer vision for image recognition and classification.
- Identifies and visualizes crucial regions in images, aiding in object classification.

## **Companies Using Explainable AI:**

### **Google:**

Utilizes explainable AI in applications like medical imaging and computer vision.

Example: DALL-E model generates images from text, employing XAI for transparency.

### **Apple:**

Applies explainable AI in medical imaging, natural language processing, and computer vision.

Example: Core ML framework incorporates XAI to identify and mitigate biases.

### **Microsoft:**

Integrates explainable AI in applications such as medical imaging and computer vision.

Example: Explainable Boosting Machine algorithm provides insights and addresses biases.