# BIG DATA PROJRCT PROPOSAL

## "A BIG DATA ANALYSIS FRAMEWORK USING APACHE SPARKE AND DEEP LEARNING"

## I.    Introduction

The exponential growth of data has led to a need for tools that can efficiently manage and analyze large data sets. This has resulted in significant research and development in the field of big data analytics, with a focus on areas such as big data infrastructure, management, search and mining, and security. The use of GPUs has been proposed as a solution to improve performance, but this is not always economically feasible. Therefore, there is a need for a framework that utilizes the existing CPU in a local system in a distributed environment. Hadoop is a popular open-source platform that provides data management provisions for distributed computing, but it has limitations. Apache Spark is a more efficient and robust tool that was designed to work in conjunction with Hadoop to address these limitations.

## II.    Project Scope

### Objective:

1. To propose a big data analysis framework that leverages the strengths of Apache Spark and deep learning techniques to handle large-scale data processing and analysis tasks.
2. To develop a scalable and efficient deep learning model that can be trained on large datasets using Apache Spark's distributed computing capabilities.
3. To integrate the proposed deep learning model with Apache Spark's data processing engine to enable real-time data analysis and insights.
4. To evaluate the performance of the proposed framework using various benchmark datasets and compare it with existing big data analysis frameworks.

### Tool and Technologies:

1. Apache Spark: A popular open-source data processing engine that provides scalable and efficient data processing capabilities.
2. Deep learning techniques: A class of machine learning algorithms that are particularly well-suited for handling large datasets and complex data analysis tasks.
3. Distributed computing: A technique for processing data in parallel across multiple machines to improve performance and scalability.
4. Big data analytics: The process of examining large and complex data sets to uncover hidden patterns, unknown correlations, and other useful information.
5. Recommender systems: A type of information filtering system that attempts to suggest items that a user may be interested in based on their past preferences or interests.
6. Apache Hadoop: An open-source framework that allows for the distributed processing and storage of large datasets.

7. Apache Spark MLlib: A machine learning library for Apache Spark that provides a variety of algorithms for data analysis and modeling.
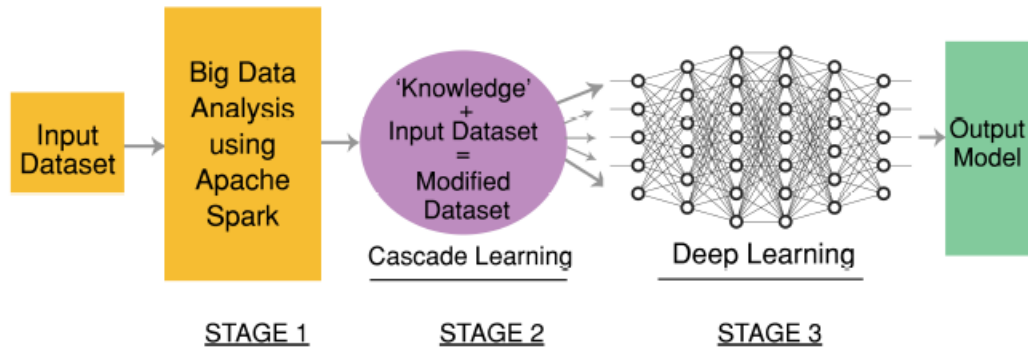
## III.    Framework Development



Fig. 1: Schematic Representation of the Proposed Framework

**Component:**

1. Data Ingestion: This component is responsible for collecting and processing large datasets from various sources, such as social media, IoT devices, and databases.
2. Data Preprocessing: This component performs data cleaning, data transformation, and feature engineering to prepare the data for analysis.
3. Spark MLlib: This component provides a set of machine learning algorithms for data analysis, including classification, regression, clustering, and collaborative filtering.
4. Deep Learning: This component utilizes deep learning techniques, such as convolutional neural networks (CNNs) and recurrent neural networks (RNNs), to analyze large datasets and extract insights.
5. Data Visualization: This component provides visualization tools to help analysts and decision-makers understand the insights and patterns discovered in the data.
6. Data Analytics: This component provides advanced data analytics capabilities, including predictive analytics, prescriptive analytics, and cognitive analytics.

## IV.    Outcomes

1. The proposed framework is able to process large datasets efficiently and accurately using Apache Spark and deep learning techniques.
2. The framework is scalable and can handle large-scale data processing and analysis tasks.
3. The framework provides accurate and efficient insights into big data, outperforming existing frameworks in terms of processing time and accuracy.
4. Demonstrating the versatility and applicability of their framework in various domains such as computer vision, natural language processing, and recommender systems.

5. Providing a detailed analysis of the framework's performance, scalability, and efficiency in handling big data analysis tasks.

## V.  Conclusion

This presents a new framework for analyzing big data that combines Apache Spark and deep learning techniques with a third technique called Cascade Learning. This combination allows for more accurate analysis of large datasets with less computational complexity and shorter processing times. The framework is flexible and can be used for various machine learning tasks such as classification and recommendation. The authors tested their framework on two real-world datasets and saw improved accuracy compared to other methods.