Project Title: **Fake News Detection Using NLP**

## Project Steps:

**Phase 1**: Problem Definition and Design Thinking

## Problem Definition:

The problem is to develop a fake news detection model using a Kaggle dataset. The goal is to distinguish between genuine and fake news articles based on their titles and text. This project involves using natural language processing (NLP) techniques to preprocess the text data, building a machine learning model for classification, and evaluating the model's performance. NLP techniques can be leveraged to analyze the text content of news articles, social media posts, and other textual sources to identify misleading or false information. Fake news detection using NLP is a multifaceted problem that requires a combination of data collection, preprocessing, machine learning modeling, and ongoing monitoring to effectively combat the spread of misinformation.

## Design Thinking:

1. **Data Source:**

   Choose the fake news dataset available on Kaggle, containing articles titles and text, along with their labels (genuine or fake).This dataset should represent various topics and sources to ensure the model's generalizability.

2. **Data Preprocessing:**

   Clean and preprocess the textual data to prepare it for analysis. This may involve tasks such as tokenization, lowercasing, removing stop words, and stemming/lemmatization.

3. **Feature Extraction:**

   Convert the text data into numerical features that can be used as input for machine learning algorithms. Common techniques include:

   - TF-IDF (Term Frequency-Inverse Document Frequency): Capturing the importance of words in a document relative to a corpus.

- Word Embeddings: Representing words as dense vectors using models like Word2Vec, GloVe, or BERT embeddings.

- N-grams: Capturing sequences of words to capture context.

## 4. Model Selection:

Choose an appropriate machine learning or deep learning model for fake news detection. Common choices include logistic regression, Naive Bayes, Support Vector Machines (SVM), recurrent neural networks (RNNs), or transformers like BERT.

## 5. Model Training:

Train the selected model using the preprocessed data, using labeled examples to learn the patterns associated with real and fake news.The model learns to distinguish between real and fake news based on the provided features.

## 6. Evaluation:

Evaluate the model's performance using metrics like accuracy, precision, recall, F1-score, and ROC-AUC. It's essential to consider the class imbalance issue, as fake news might be a minority class.