# Ex1b: Hidden Markov Model (HMM) based Predictive Text System

**Learning Objective:**

Toimplementa predictivetextsystemusing Hidden Markov Model (HMM), enabling students to understand contextual probability, POS transitions, and next-word prediction using the Brown Corpus.

**Steps:**
1. Initialize the environment by installing and importing the Natural Language Toolkit (nltk). Download the Brown Corpus (for text data) and the Universal Tagset (for simplified Part-of-Speech tags).
2. Extract raw sentences from the Brown Corpus using brown.sents().
3. Tokenization: Convert all words to lowercase and flatten the nested list into a single list of tokens to verify the total volume of training data.
4. Load the tagged version of the corpus (brown.tagged_sents) which provides the "Hidden States" (POS tags) linked to the "Observations" (words).
5. Train the HMM Model. Use the HiddenMarkovModelTrainer to perform Supervised Learning. Calculate the internal parameters of the HMM: Initial State Probabilities, Transition Probabilities, Emission Probabilities
6. Build POS Transition Probabilities
7. Build Emission Probabilities
8. Define Prediction Function
9. Test the Model
10. Verify POS prediction on certain words.

**Program:**

```
import nltk
from nltk.corpus import brown
from collections import defaultdict, Counter

nltk.download('brown')
nltk.download('universal_tagset')

tagged_words = brown.tagged_words(tagset='universal')

transition_counts = defaultdict(Counter)
emission_counts = defaultdict(Counter)
tag_counts = Counter()

for i in range(len(tagged_words) - 1):
  word, tag = tagged_words[i]
  next_word, next_tag = tagged_words[i + 1]

  transition_counts[tag][next_tag] += 1
  emission_counts[tag][word.lower()] += 1
  tag_counts[tag] += 1

def predict_next_word_hmm(previous_word):
  previous_word = previous_word.lower()

  possible_tags = [
    tag for tag in emission_counts
    if previous_word in emission_counts[tag]
    ]
```

```
    if not possible_tags:
     return "No prediction available"
     current_tag = possible_tags[0]
     next_tag = transition_counts[current_tag].most_common(1)[0][0]
     predicted_word = emission_counts[next_tag].most_common(1)[0][0]
     return predicted_word

    word = "the"
    prediction = predict_next_word_hmm(word)

    print("Previous Word:", word)
    print("Predicted Next Word:", prediction)
```

**Output:**



**Learning Outcome:**
Upon completion of this experiment, HMM based model was built to predict the next word
with POS-based suggestions, and evaluated the effectiveness of context-driven predictive text
generation using the Brown Corpus.