# Great Learning

# AIML

## CAPSTONE PROJECT

# NATURAL
# LANGUAGE
# PROCESSING

## MACHINE TRANSLATION ASSIGNMENT

# *PROBLEM STATEMENT*

- **DOMAIN:** MACHINE TRANSLATION.

- **CONTEXT:**

Machine Translation is the automated translation of source material into another language without human intervention. The database comes from ACL2014 Ninth workshop on Statistical Machine Translation. This workshop mainly focusses on language translation between European language pairs. The idea behind the workshop is to provide the ability for two parties to communicate and exchange the ideas from different countries.

- **DATA DESCRIPTION:**

The database is basically sentences in German/English of various events. Three datasets are obtained from Statistical Machine Translation workshop. Either the dataset can be downloaded from the link or can be used from the shared files. Three datasets are,
  - Europarl v7
  - Common Crawl corpus
  - News Commentary

Link to download the dataset:  https://statmt.org/wmt14/translation-task.html

- **PROJECT OBJECTIVE:**

Design a Machine Translation model that can be used to translate sentences from German language to English language or vice-versa.

- **PROJECT TASK:** [ Score: 100 points]

  1. **Milestone 1:** [ Score: 40 points ]
     ‣ **Input**: Context and Dataset
     ‣ **Process**:
         ‣ Step 1: Import and merge all the three datasets. [ 5 points ]
         ‣ Step 2: Data cleansing [ 7 points ]
         ‣ Step 3: NLP pre processing - Dataset suitable to be used for AIML model learning [ 8 points ]
         ‣ Step 4: Design, train and test simple RNN & LSTM model [ 10 points ]
         ‣ Step 5: Interim report [ 10 points ]
     ‣ **Submission**: Interim report, Jupyter Notebook with all the steps in Milestone-1

  2. **Milestone 2:** [ Score: 60 points ]
     ‣ **Input**: Preprocessed output from Milestone-1
     ‣ **Process**:
         ‣ Step 1: Design, train and test RNN & LSTM model with embeddings [ 7 points ]
         ‣ Step 2: Design, train and test bidirectional RNN & LSTM model [ 8 points ]
         ‣ Step 3: Design, train and test Encoder-Decoder RNN & LSTM model **(Optional-If interested can try, but marks will not be reduced if not attempted)**
         ‣ Step 4: Choose the best performing model and pickle it. [ 5 points ]
         ‣ Step 5: Final Report [40 Points]
     ‣ **Submission**: Final report, Jupyter Notebook with all the steps in Milestone-1 and Milestone-2

  3. **Milestone 3:** [ Optional ]
     ‣ **Process**:
         ‣ Step 1: Design a clickable UI based Translation interface.
              Hint: Input - Sentence in one language(German/English), Output - Translated sentence in other language(English/German)
     ‣ **Submission**: Final report, Jupyter Notebook with the addition of clickable UI based interface

‣ Hints:
   ‣ Please refer to the research papers to understand how to Machine Translation: https://statmt.org/wmt14/papers.html
   ‣ To make GUI as a desk app you can use TKINTER library.
   ‣ To make web service GUI you can use FLASK or DJANGO library.
   ‣ Reference: https://www.mygreatlearning.com/academy/learn-for-free/courses/machine-translation

## POINTS TO REMEMBER

1. A maximum of 100 points will be awarded for this project

2. Project to be submitted within 6 weeks of date of release. Late submission will be accepted under genuine situation. Score will be given as per the below formula:

   If the current score is greater than 40 then the final score will be capped at 40.

   Else the current score will be awarded.

3. Any form of plagiarism is strictly prohibited. No score will be awarded in this case.

# HAPPY LEARNING!

Great Learning