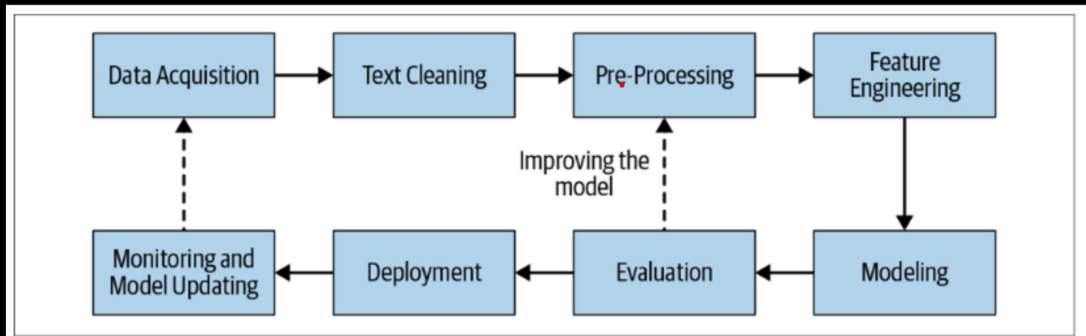# NLP pipeline



① Data Acquisition

② Text cleaning

③ pre-processing

④ feature engineering

⑤ Modeling

⑥ evaluation

⑦ Deployment

⑧ Monitoring and Model updating

## Data Acquisition

① This is the initial stage where raw textual data is collected
The data may come from various sources Web scraping
APIs
data repositories
user input

② The goal is to gather a large amount of diverse and relevant
text data that can be used for NLP task

## Text cleaning

① Once data is acquired It contains noise — HTML tags, punctuation
stop words

② Text cleaning ⟶ lowering all text
⟶ Removal special characters, extra spaces, number
⟶ Handling missing or reduced

# Pre-processing

After cleaning, the data undergoes pre-processing to prepare it for modeling

1. Tokenization :- text ~ words or sub words
2. Stemming or lemmation: Reducing words to their base or root form   acting → act
3. Stop word Removal : eliminating common words, that are (the, and, in) not informative
4. POS → none, verb → structure of sentence

# Feature Engineering

In this stage, the cleaned and pre-processed text is transformed into features that the model can understand

1. TF-IDF — text ⌒ vector
2. Word embedding → Word2vec, BERT
3. N-grams

# Modeling

[ ML model
  Deep learning ]

This is where machine learning or deep learning models are trained using the features

# Evaluation