**Please perform the following task:**

Duration: 1 day

1. Take the open-source dataset: https://huggingface.co/datasets/imdb
2. Load any model from huggingface and its tokenizer. Load another tokenizer from spacy for reference.
3. Tokenize the dataset using the two tokenizers and save the tokens from the two tokenizers separately in two lists.
4. Now, take top 1000 tokens based on the information value/entropy of the tokens, from the two lists. Now you have two lists, each of 1000 tokens.
5. Compare the two list of tokens (it is up to you how you want to compare, create plots, generate some metrics etc.). Based on the comparison which one do you think creates better tokens – does the result match your intuition? Which one would you use for a sentiment classification task? Why?