# QUESTION 6:

## TOKENIZERS:

**Indic -BERT tokenizer** is specifically trained for languages of the Indian subcontinent, including Hindi. It is fine-tuned to better understand the composition and structure of these languages, which should lead to more accurate tokenization, especially for tasks like word grouping in Hindi. **mBERT tokenizer** is trained on a multilingual corpus, including Hindi. But it is not specifically optimized for Hindi, instead trained on a diverse range of languages and can handle multilingual text reasonably well. So, it should be able to capture word groups effectively in Hindi. **BPE (Byte Pair Encoding) tokenizer** might not be able to make right word grouping because it iteratively replaces the most frequent pair of consecutive characters with a new symbol, it can handle only some level of morphological differences but for Hindi morphological differences is more complex. **Unigram tokenizer** splits text into individual words based on whitespace or punctuation. So, it may work decently in most languages.

| | UNIGRAM -1000 | BPE- 1000 | BPE- 2000 | mBERT - 1000 | mBERT - 2000 | Indic-BER T -1000 | Indic-BER T -2000 | WHITE SPACE |
|---|---|---|---|---|---|---|---|---|
| **PRECISION** | 0.02117061 021170610 3 | 0.0298102 98102981 03 | 0.0298102 981029810 3 | 0.02053140 096618357 6 | 0.0205314 009661835 76 | 0.00991735 537190082 7 | 0.0099173 553719008 27 | 0.0573122 52964426 88 |
| **RECALL** | 0.09239130 434782608 | 0.1195652 17391304 35 | 0.1195652 173913043 5 | 0.09239130 434782608 | 0.0923913 043478260 8 | 0.03260869 565217391 | 0.0326086 956521739 1 | 0.1576086 95652173 92 |
| **F-SCORE** | 0.03444782 168186423 5 | 0.0477223 42733188 726 | 0.0477223 427331887 26 | 0.03359683 794466403 | 0.0335968 379446640 3 | 0.01520912 547528517 3 | 0.0152091 254752851 73 | 0.0840579 71014492 76 |

## COMPARISON:

**1-** According to the values of precision obtained from question-5 i am getting the best precision for white space tokenizer. This is because white space tokenizer is separating every single word which is matching with all **punctuation word groups** and all **conjunction word groups** as these are single words. But other models are breaking the words even which makes it difficult to form the right word group. Indic-Bert is supposed to responde best but actually it's giving syllables instead of words so right match is not possible.

Recall gives the true positives out of all actual positive instances in the data, and is also best for white space tokenizer. And hence, its F-score is also highest.

**2-** Also we can observe that the metric values for the same models (e.g. BPE) on different vocab sizes (1000 and 2000) are the same. From this we interpret that the token formed by them is very much the same. So, vocab size is not helping much in improving the right word group division for the sentences.