

## MANDATE: 2

### CONTRIBUTION PREPARATION AND PLANNING OF DATASET AND PRE TRAINED MODEL

**Keywords** QuestionAnsweringModel, NLP, IPC, LEGAL-BERT, ENCODING, TOKENIZATION

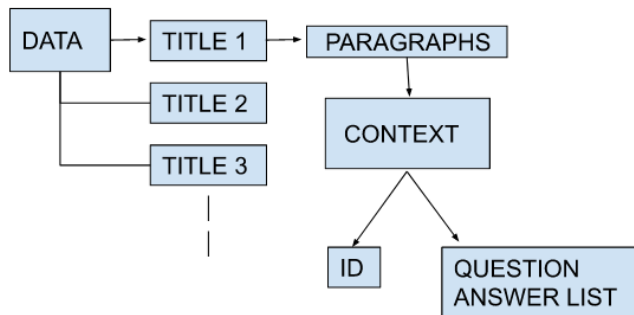
#### 1 DATASET SELECTION:

IPC (Indian Penal Codes) data is available online on Government Site in .pdf format: "https://legislative.gov.in".

But pre trained models like BERT takes data in a particular format.

##### 1.1 HUGGING FACE TRANSFORMERS:

BERT/DistilBERT accepts data in json format:



```
{'title': 'S. 302 Punishment for murder',  
'paragraphs': [{'context': 'Punishment for murder', 'qas': []},  
'context': 'Whoever commits murder shall be punished with death, or imprisonment for life, and shall also be liable to fine.',  
'qas': [{'question': 'what is the punishment for murder?',  
'id': '62db2313-6a67-45a0-bf35-1c314b71fd9b',  
'answers': [{'answer_start': 32,  
'text': 'punished with death, or imprisonment for life, and shall also be liable to fine.'}]}]}
```

For QuestionAnswering Model, the dataset needs to be converted in the above format to feed pretrained model BERT or DistilBERT.

We can do this using HAYSTACK ANNOTATION TOOL.

##### 1.2 HAYSTACK ANNOTATION TOOL:

We can upload our data on this tool and using it we can select a particular text and add question answer sets. This tool then, will automatically format our data as required.

#### 2 PRE-TRAINED MODEL:

##### 2.1 DistilBERT

DistilBERT is smaller but faster version of BERT with maintaining similar accuracy. But it may not perform as BERT in complex tasks.

##### 2.2 LEGAL-BERT

LEGAL-BERT has been trained on large corpus of legal texts. It will help in understanding legal terminology better.

Reference: <https://huggingface.co/nlpauieb/legal-bert-base-uncased>

##### 2.2.1 MODEL

```
modelName = "nlpauieb/legal-bert-base-uncased"  
tokenizer = AutoTokenizer.from_pretrained(modelName)
```

#### 3 DATA PRE-PROCESSING:

##### 3.0.1 ADDING ANSWER's ENDING INDEX

In our train dataset we have start index for answer in context but don't have the end index. For each answer, from answer text length we can add start index to it and get end index.

- [1] Government of India. "https://legislative.gov.in", Indian Penal Code and laws[Dataset]
- [2] "https://github.com/re-search/DocProductstart-of-content", DocProduct [Context Approach]