# TnCSimplifier - Summarizes Terms and Conditions/ Privacy policies

**Devansh Goswami**
devansh21460@iiitd.ac.in

**Ritisha Singh**
Ritisha21089@iiitd.ac.in
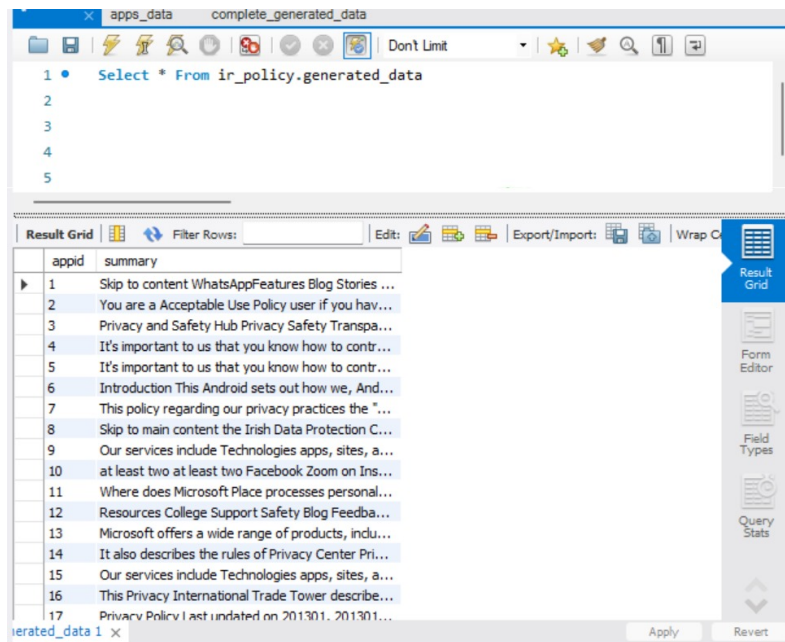
**Samanyu Kamra**
samanyu21487@iiitd.ac.in

**Shriya Verma**
shriya21490@iiitd.ac.in

## 1   Updated Baseline Results

We evaluated our dataset using five different LLM models: Pegasus, BART, and three variants of T-5. After assessing the cosine similarity scores, we determined that Pegasus outperformed the other models for summarizing privacy policies and terms and conditions.

### 1.1   Updation of database:

Following the selection of the Pegasus model, we utilized it to condense the privacy policy texts contained within our database. These condensed summaries were then incorporated into a new table named 'generateddata'. This table serves as a valuable resource for future use, allowing users to access pre-generated summaries without the need for the model to regenerate the output each time.



Figure 1: Generated Summaries using Pegasus

## 2   Data Pre-processing

Our preprocessing function aims to prepare text for input into the PEGASUS summarization model, with a focus on handling named entities and facilitating intelligent chunking for simplifying terms and conditions.

The preprocessing includes:

1. Removal of URLs, HTML tags, and decoding of HTML entities.
2. Named entities are temporarily replaced with placeholders to prevent them from being split during subsequent processing steps.
3. All Special characters, excluding common punctuation marks, and the placeholders for entities are removed.
4. The text is split into sentences and rejoined to ensure proper sentence boundaries.

This preprocessing technique is optimal for inputting terms and conditions into PEGASUS because the removal of extra elements and simplified structure of the text allows the model to focus more effectively on summarizing the essential content of the terms and conditions text.

## 3   Pegasus Model

Pegasus is a specific model that gives generated abstracts based on the article which means the model reads the article text and writes a suitable headline. This abstractive text summarization is one of the most challenging tasks,involving understanding long passages, information compression, and language generation.

The Pegasus model is impressive symbolic reasoning, as researchers concluded based on performance and evaluation results. Some variations show not just the human-level of natural language understanding, even beating the human performance in generating summarized output.

### 3.1   Architecture:

The current practice for this task would be to train a language model by predicting the masked out token at the end of the sequence. The task is known as self-supervised autoregressive language modeling. Next, fine-tune it on a dataset of labeled summaries and have the generalized output of an article or a document that hasn't seen on the training test, and contains the most information from the source documents. The aim of fine-tuning this kind of architecture is to gain the best results on downstream NLP tasks such as yelp-review, natural language inference, and question-answering.
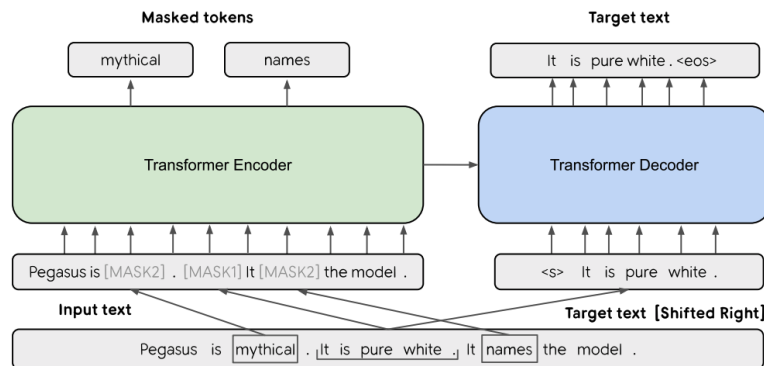


Figure 2: Pegasus Model Architecture

In our code we used the pre-processed tokens as inputs to generate the summary using the Pegasus pre-trained model.

# 4  Mid Term - 1 Results

In the baseline section, we had curated the data that we used to train and test our summarization models.In addition to this, We tested a number of summarization models and assessed their performance based on their cosine similarity scores, leading us to the selection of Pegasus.

In this segment of the project, We utilized the pretrained pegasus model to generate summaries of the privacy policies we had available in our database.



Figure 3: Updated Database Schema

In addition to this, We developed a basic web application prototype that takes input text and returns summary back to user.



Figure 4: Website Prototype