# 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   Introduction

In today's digital age, the proliferation of mobile apps and online services has made it increasingly challenging for users to navigate through the lengthy and complex terms and conditions, privacy policies, and legal agreements associated with them. Often, users tend to skip reading these documents due to their tedious and time-consuming nature. However, this lack of awareness can expose users to potential risks, as cyber criminals may exploit their personal information without their informed consent.

The *Tnc-Simplifier* app aims to address this issue by providing a solution that simplifies the process of understanding and evaluating the policies of apps and websites. By presenting policy information in a summarized and user-friendly manner, our app empowers users with the essential insights they need to make informed decisions about granting access to their personal information. This project is motivated by a commitment to protect users and enhance digital transparency in an increasingly complex online landscape.

## 2   Problem Statement

*Tnc-Simplifier* aims to develop a solution that simplifies the process of reading and understanding privacy policies. Users tend to skip privacy policies primarily due to their length and the legal jargon within them that the average user can not understand. Our web-application aims to provide a summarized version of the policies, highlighting the crucial permissions the user is granting to various applications. Primary motivation for the development of such a web application is the increased number of cyber crimes in today's world. We aim to protect users and minimize the exploitation of data. The main challenge in developing such an application will be differentiating between irrelevant and essential portions of a privacy policy. In order to address this issue, We will use keywords in order to segregate essential parts from non-essential parts. Input to our summarization system can be given in text/image format. A chat-bot will be in place to further solve targeted questions about the privacy policies being queried. Moreover, Users will be able to translate the outputs of our algorithm into any language of their choice, increasing accessibility and user understanding. We will also allow voice-enabled queries which will process, and answer questions pertaining to the terms and conditions of the application being queried, increasing the accessibility of our system. We are working to ensure a safer digital experience for users of *Tnc-Simplifier*.

## 3   Literature Survey

The literature survey section provides an overview of the existing research and applications related to simplifying terms and conditions and privacy policies. We review relevant studies, applications, and

technologies that address similar challenges in the field of digital privacy and user understanding of legal agreements.

There are two main types of text summarization, namely extractive and abstractive. Traditional extractive summarization techniques work by extracting phrases one at a time and modelling the relationships between them. A popular approach is the TextRank algorithm, which uses a graph-based model. As for abstractive techniques, Google's Pegasus [4] is one of the most popular models. For sequence-to-sequence learning, it employs a transformer-based encoder-decoder paradigm. In order to help individuals and companies better understand the potential risks and consequences of the agreements they are entering into, AI and natural language processing are taking the fields of privacy policies and legal text summarization by storm. The evolution in summarizing legal documents through NLP and machine learning techniques is underscored by several key studies, each contributing to the field's development. Initially, the application of algorithms like Multinomial Naive Bayes [1] offered a pathway to distill complex legal texts into more accessible summaries, setting a precedent for innovation. Building on this, LexRank [2] utilized eigenvector centrality to enhance extractive summarization, demonstrating the effectiveness of graph-based methods over traditional approaches by more accurately determining sentence importance.

Adding to the complexity, a study [5] introduced a CNN-based risk classifier to prioritize high-risk sections in privacy policies, marking a significant step towards tailored summarization focused on user concerns. This approach was further refined by integrating RoBERTa with sequence-to-sequence models, creating a hybrid system that extracts and paraphrases legal terms. Google's Pegasus has been fine tuned on legal corpus to create PEGASUSCourtOp/legalPEGASUS [3], this showcases the potential of legal domain specefic adaptations of existing systems.

Several existing platforms and services like Lexcheck and Lawgeex currently use AI and natural language processing to identify potentially risky clauses or language in contracts and policies, which is important for legal contracts and policies. These projects help those without legal expertise by highlighting key risks, preventing legal issues, and protecting against reputational harm. Our web-application, *Tnc-Simplifier*, takes this up a notch by increasing accessibility with our multilingual features and utilization of optical character recognition to increase the amount and types of queries that can be addressed.

## 4    Baseline Results

In the baseline section, we have curated the data that we will use to train and test our summarization models. In addition to this, We have tested a number of summarization models and assessed their performance.

### 4.1    Dataset Curation:

From various applications of different categories, like, texting, online-shopping, social media, health, etc. Data was collected and stored in the form of tables. This data is stored in our back-end and will be used for future processing.

### 4.1.1    Database:

The data has been divided into 3 tables, One stores the primary data about the application; which includes - An application ID, Type of the application, Name of the application, the Privacy policy text and whether the application is paid or not. Another table stores the permissions that an application may ask for and how it will use the access. The third table stores the features and which applications have those features.
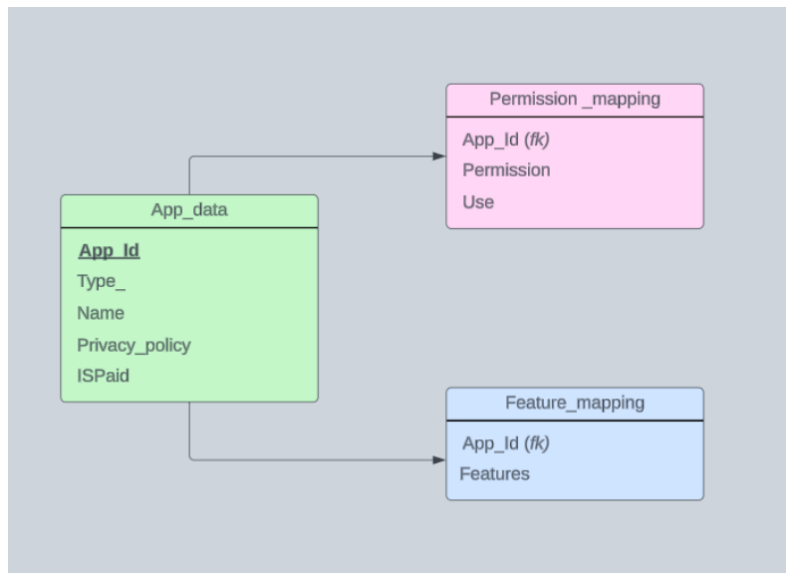
Figure 1: Structure of the database



Figure 2: Before Summarization



Figure 3: After Summarization

## 4.2 Summarization Models:

5 different LLMs models were run, of which 3 were T-5 variants, BART, and Pegasus. We compared the accuracies of these models on the basis of cosine similarity. It can be concluded that Pegasus performed better than the other LLMs.

```
For T-5
Model: t5-small
Summary: this privacy policy applies to you if you are a Weixin user. you
Difflib Similarity: 0.0007284150144544854
Cosine Similarity: 0.7054513692855835

Model: t5-base
Summary: we respect your concerns about privacy and appreciate your trust
Difflib Similarity: 0.0015009551532793597
Cosine Similarity: 0.636496365070343

Model: t5-large
Summary: WeChat respects your concerns about privacy and appreciates your
Difflib Similarity: 0.0010449320794148381
Cosine Similarity: 0.9254676103591919

For New Models
Model: facebook/bart-large-cnn
Summary: This Privacy Policy only applies to you if you are a WeChat user.
Difflib Similarity: 0.0002729692227201383
Cosine Similarity: 0.6110114455223083

Model: google/pegasus-large
Summary: This Privacy Policy only applies to you if you are a WeChat user,
Difflib Similarity: 0.006655950736908832
Cosine Similarity: 0.9366061687469482
```

Figure 4: Results obtained for one of the iterations

# 5    References

The reference section includes a comprehensive list of all the sources, research papers, articles, books, and online resources cited throughout the research paper. Proper citation is essential to give credit to prior work and to allow readers to access the sources for further study and verification.

1. Perera, Thenuka, and Theja Perera. "Barrister-Processing and Summarization of Terms Conditions / Privacy Policies," April 2, 2021. https://doi.org/10.1109/i2ct51068.2021.9418090.

2. LexRank: Graph-based Lexical Centrality as Salience in Text Summarization. (n.d.). https://www.cs.cmu.edu/afs/cs/project/jair/pub/volume22/erkan04a-html/erkan04a.html

3. Aashish Ghimire, Raj Shrestha, and John Edwards. "Too Legal; Didn't Read (TLDR): Summarization of Court Opinions," May 12, 2023. https://doi.org/10.1109/ietc57902.2023.10152119.

4. Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter J. Liu. 2020. PEGASUS: pre-training with extracted gap-sentences for abstractive summarization. In Proceedings of the 37th International Conference on Machine Learning (ICML'20), Vol. 119. JMLR.org, Article 1051, 11328–11339.

5. Moniba Keymanesh, Micha Elsner, and Srinivasan Sarthasarathy. "Toward Domain-Guided Controllable Summarization of Privacy Policies.," January 1, 2020, 18–24.

6. Dalal, Sarthak, Amit Singhal, and Brejesh Lall. "LexRank and PEGASUS Transformer for Summarization of Legal Documents." Lecture Notes in Electrical Engineering, January 1, 2023, 569–77. https://doi.org/10.1007/978-981-99-0085-5$_4$6.$https://huggingface.co/docs/transformers/en/model_doc/bart$