

The background features a dark blue gradient with faint, white technical diagrams on the left side, including circular gauges with numerical scales (40, 150, 160, 170, 180, 190, 200, 210, 220, 230, 240, 250, 260) and arrows. At the bottom, there is a silhouette of a mountain range under a starry night sky.

TERMS OF SERVICE SUMMARIZATION

E. RUFFOLI, F. HUDEMA, T. BALDI

CLICK-TO-AGREE

Problem

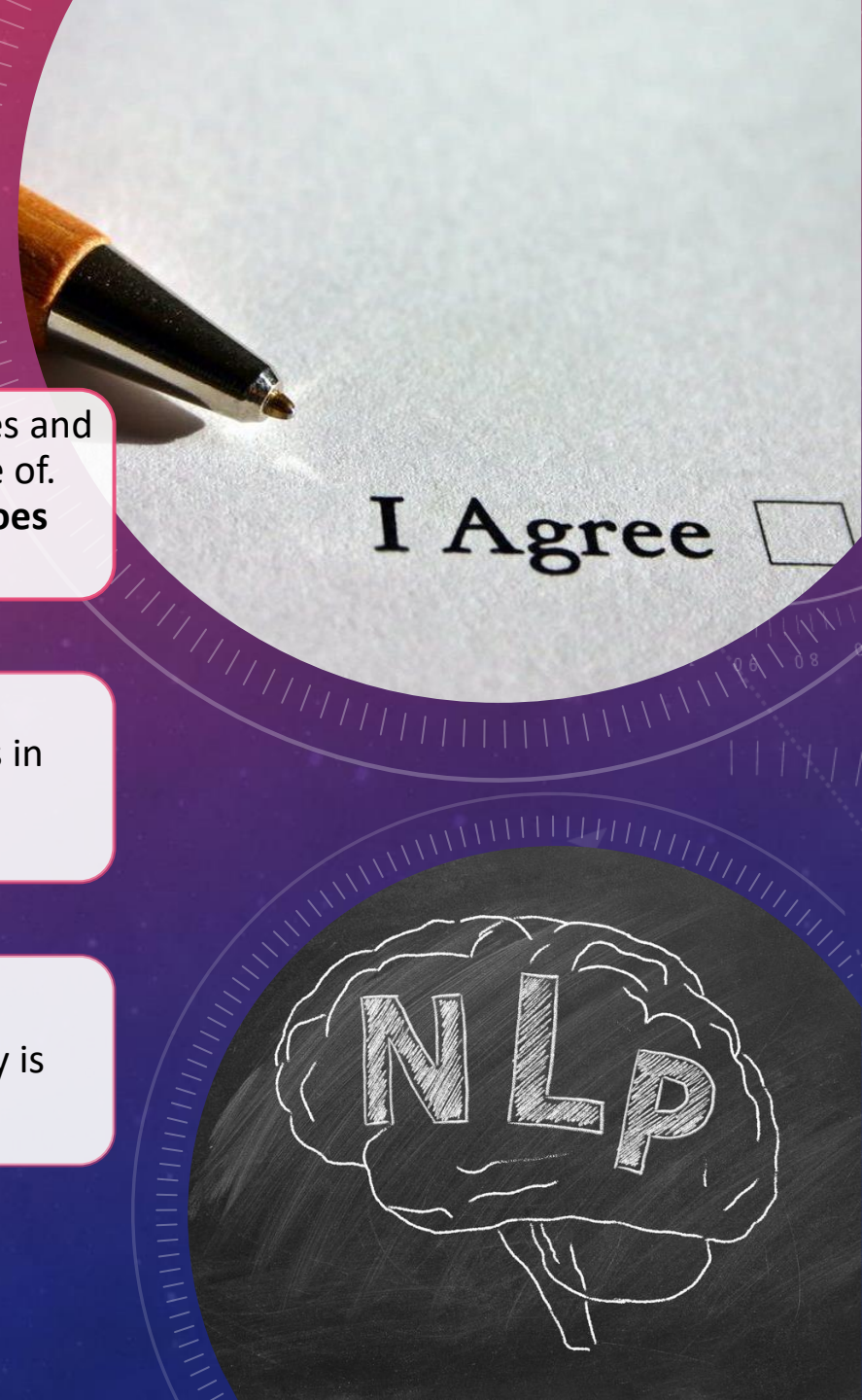
- Users give the right to keep, analyze and sell their data to web-based services and third parties by accepting services contracts about which they are not aware of. Increasingly often, **people click away their right to go to court if anything goes wrong.**

Solution

- Automated synthesis system that will extract key information from contracts in order to make them more understandable to the users.

How

- The service we offer is a software that uses algorithms to summarize textual documents. This specific field of Artificial Intelligence is called NLP and today is certainly one of the branches on which more research is carried out.



TEXT SUMMARIZATION

Extractive

Alice and Bob took the train to **visit the zoo**. They **saw** a baby giraffe, a lion and **a flock of colorful tropical birds**.

- **Summary:**

Alice and Bob visit the zoo saw a flock of birds

Abstractive

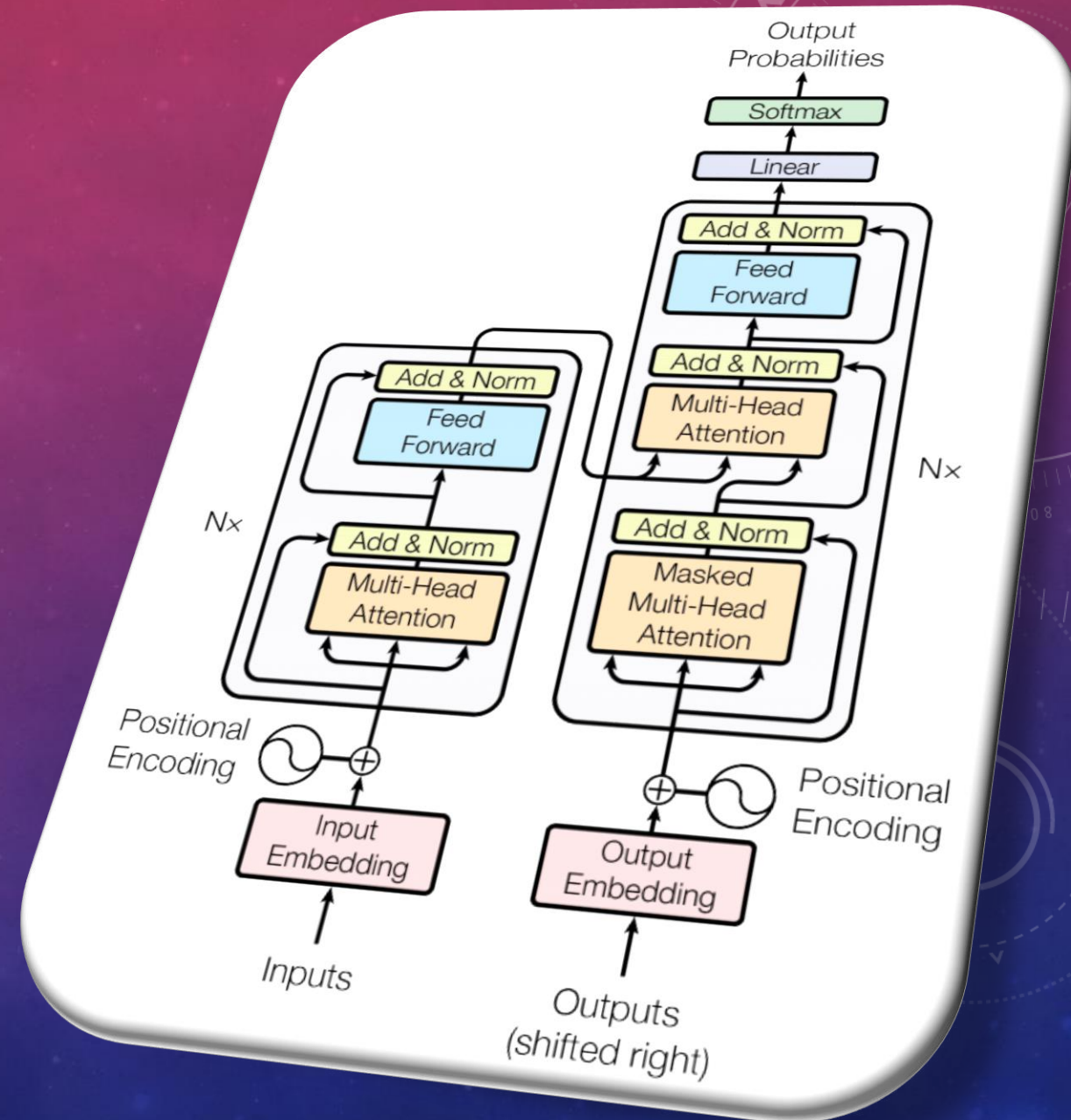
Alice and Bob took the train to visit the zoo. They saw a baby giraffe, a lion and a flock of colorful tropical birds

- **Summary:**

Alice and Bob **visited** the zoo **and** saw **animals** and birds

TRANSFORMERS

- Introduced by Google in December 2017 in the paper «Attention is all you need».
- Transformers turned out to be a gamechanger for natural language processing.
- They are based on the concept of «Attention»: a mechanism that prioritizes some parts of the input data. In this way the network can devote more focus to the important parts of the data.



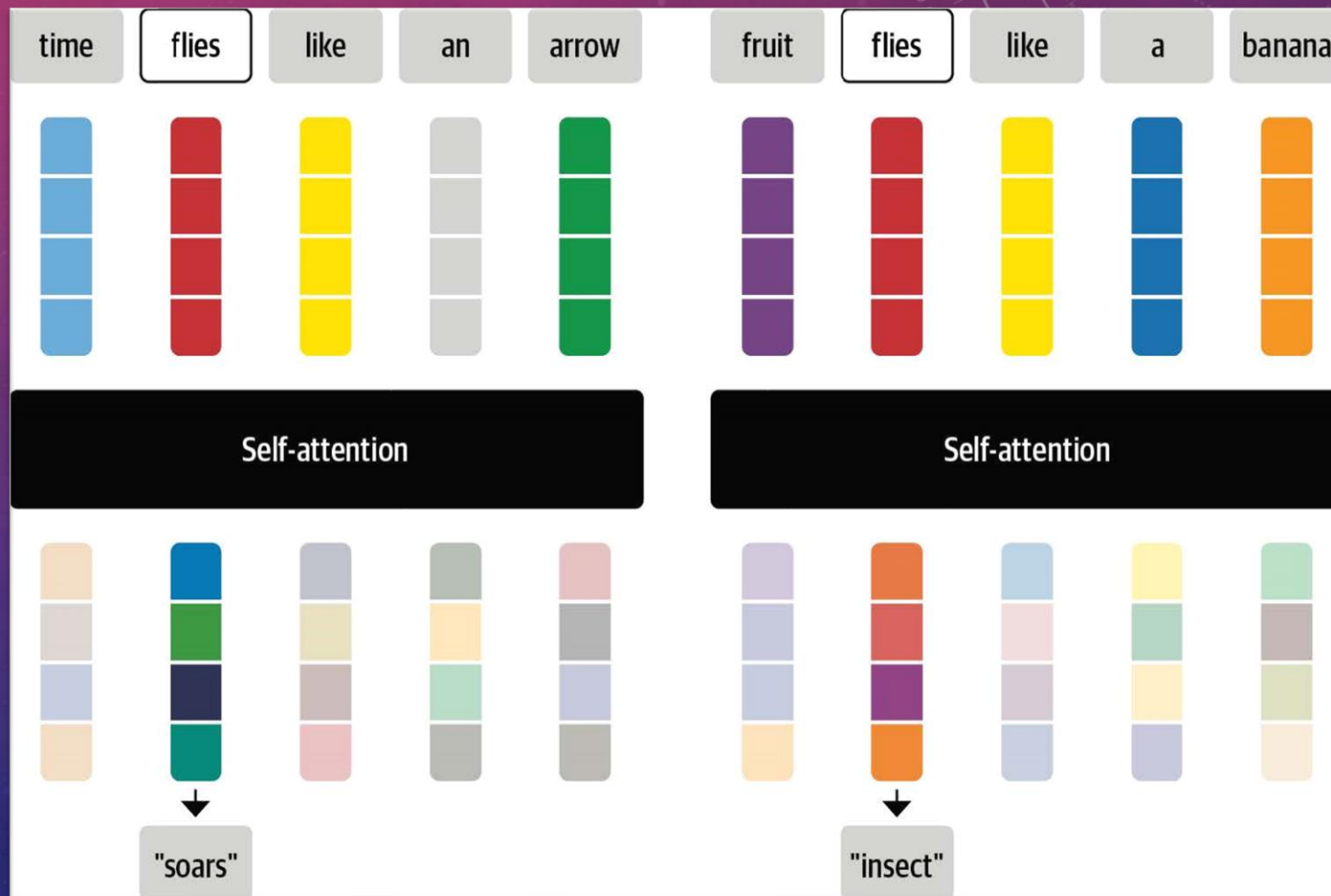
ATTENTION

For text sequences, the elements are token embeddings where each token is mapped to a vector of fixed dimension.

$$x'_i = \sum_{j=1}^n w_{ji} x_j$$

Attention updates raw token embeddings (x) into contextualized embeddings (x') to create representations that incorporate information from the whole sequence.

Attention can generate two different representations for the word “flies” based on the context.



SCRAPING

Collection of contracts from *tosdr.org*, through a Python script that exploits the API provided by the site.

Terms of services and privacy policies of major Internet services and companies.

Summaries of the main components of the contracts written by hand.

Uber Privacy Notice:

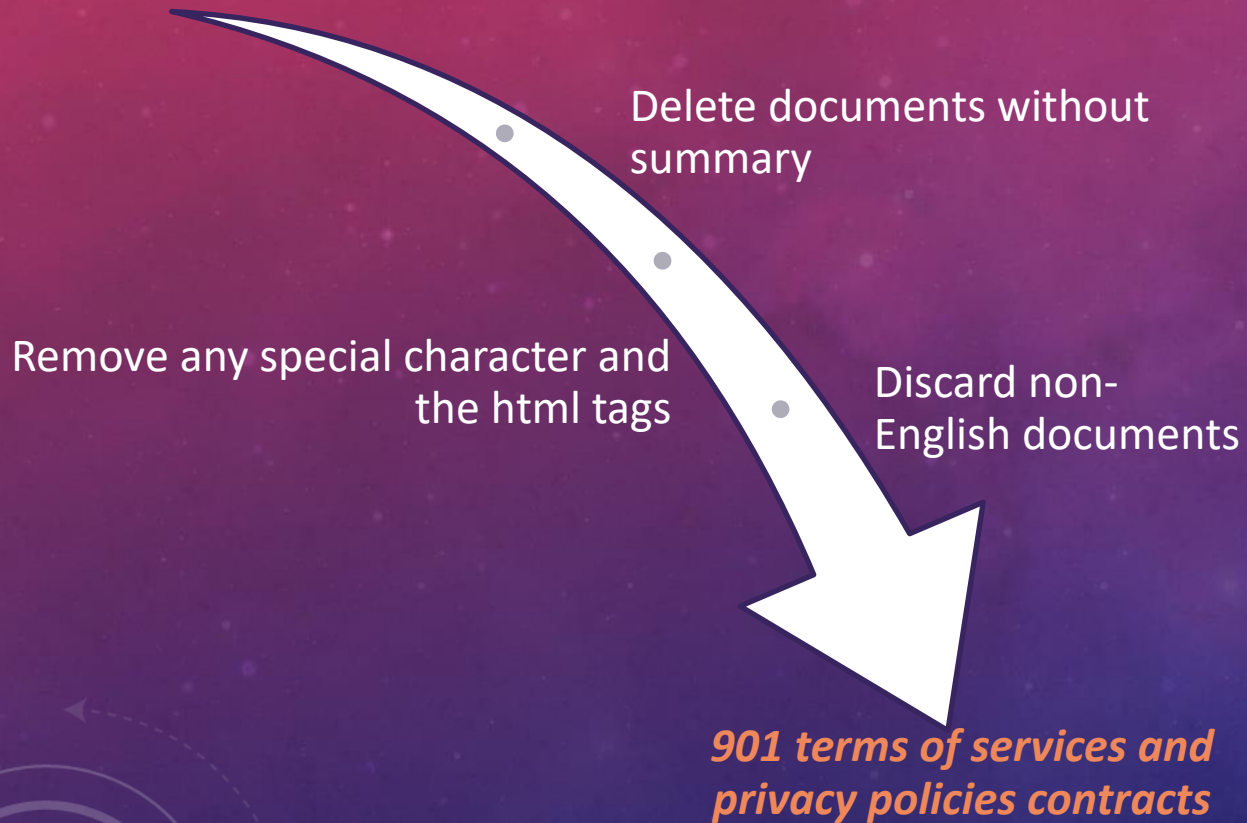
[...]We collect data when users create or update their Uber accounts. This may include their name, email, phone number, login name and password, address, profile picture, payment or banking information (including related payment verification information), driver's license and other government identification documents (which may indicate document numbers as well as birth date, gender, and photo). We may use the photos submitted by users to verify their identities, such as through facial recognition technologies. How we use personal data Location data (driver and delivery person): Uber collects this data when the Uber app is running in the foreground (app open and on-screen) or background (app open but not on-screen) of their mobile device. Location data (riders and delivery recipients). We collect precise or approximate location data from riders' and delivery recipients' mobile devices if they enable us to do so.[...]

Summary:

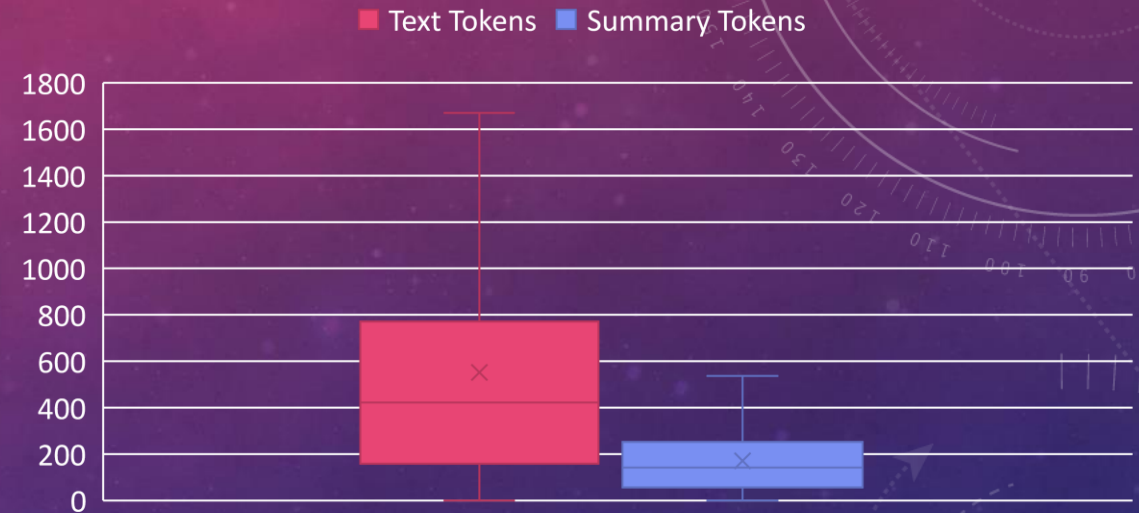
Many different types of personal data are collected. Your biometric data is collected. Information is provided about how your personal data is used. This service receives your precise location through GPS coordinates.

DATASET PREPROCESSING AND ANALYSIS

Raw dataset



Tokens Distribution



Since summarization models have a **maximum context length** in terms of tokens, so we need to analyze the **size of the input and output** documents.

HOW TO FINE TUNE

Transformers model

- Transformer models (BART, T5, etc.) have been trained as language models. This means they have been trained on large amounts of raw text in a self-supervised fashion.

Fine-tuning

- Adaptation of the model to the summarization task. For our application we select pretrained models fine-tuned on English language that has been adapted to the summarization task using datasets as the *cnn_dailymail*.

Training a Fine-tuned model

- Perform additional training with a dataset specific to the task. Fine-tuning in this step requires way less data and the amount of time and resources needed to get good results are much lower.



The Hugging Face ecosystem consists of two main parts: a family of **library** and the **Hub**. The library supports the major deep learning frameworks including PyTorch, and provides a standardized interface to a wide range of transformer models.



Google Colab is a hosted **Jupyter notebook** service that requires no setup to use, while providing access free of charge to **computing resources** including GPUs.

TRANSFORMERS PIPELINE

Raw text

«We don't sell your personal data to advertisers, and we don't share information that directly identifies you (such as your name, email address or other contact information) with advertisers unless you give us specific permission. Instead, advertisers can tell us things such as the kind of audience that they want to see their ads, and we show those ads to people who may be interested. We provide advertisers with reports about the performance of their ads that help them understand how people are interacting with their content. See Section 2 below to learn more.»

Tokenizer

Input Ids

- Splitting the input into words, subwords, or symbols (like punctuation) that are called tokens
 - Mapping each token to an integer
- «[201,2023, 2607, ...]»

Model

Summary

«Your personal data is used for advertising.»

EVALUATION METRICS: ROUGE-N

"the brown fox jumps"

reference text

"the fox"

$$\frac{\text{match}(\text{gram}_2)}{\text{count}(\text{gram}_2)} = \frac{1}{1} = 100\% \text{ recall}$$

"the hello a cat dog fox jumps"

["the", "fox", "jumps"]

1.0 recall
0.43 precision

$$2 * \frac{0.43 * 1.0}{0.43 + 1.0} = 0.6$$

60% f1 score

Model

reference text

"the fox jumps" → ["the", "fox", "jumps"]

"the hello a cat dog fox jumps"

$$\frac{\text{count}_{\text{match}}(\text{gram}_n)}{\text{count}(\text{gram}_n)} = \frac{3}{7} = 43\% \text{ precision}$$

EVALUATION METRICS: ROUGE-L

"the hello a cat dog fox jumps"

longest common
subsequence length = 2

"the fox jumps"

reference text

$\text{LCS}(\text{gram}_n)$

$\text{count}(\text{gram}_n)$

$$\frac{2}{3} = 66\% \text{ recall}$$



$$2 * \frac{0.29 * 0.66}{0.29 + 0.66} = 0.6$$

40% f1 score

"the hello a cat dog fox jumps"

longest common
subsequence length = 2

"the fox jumps"

reference text

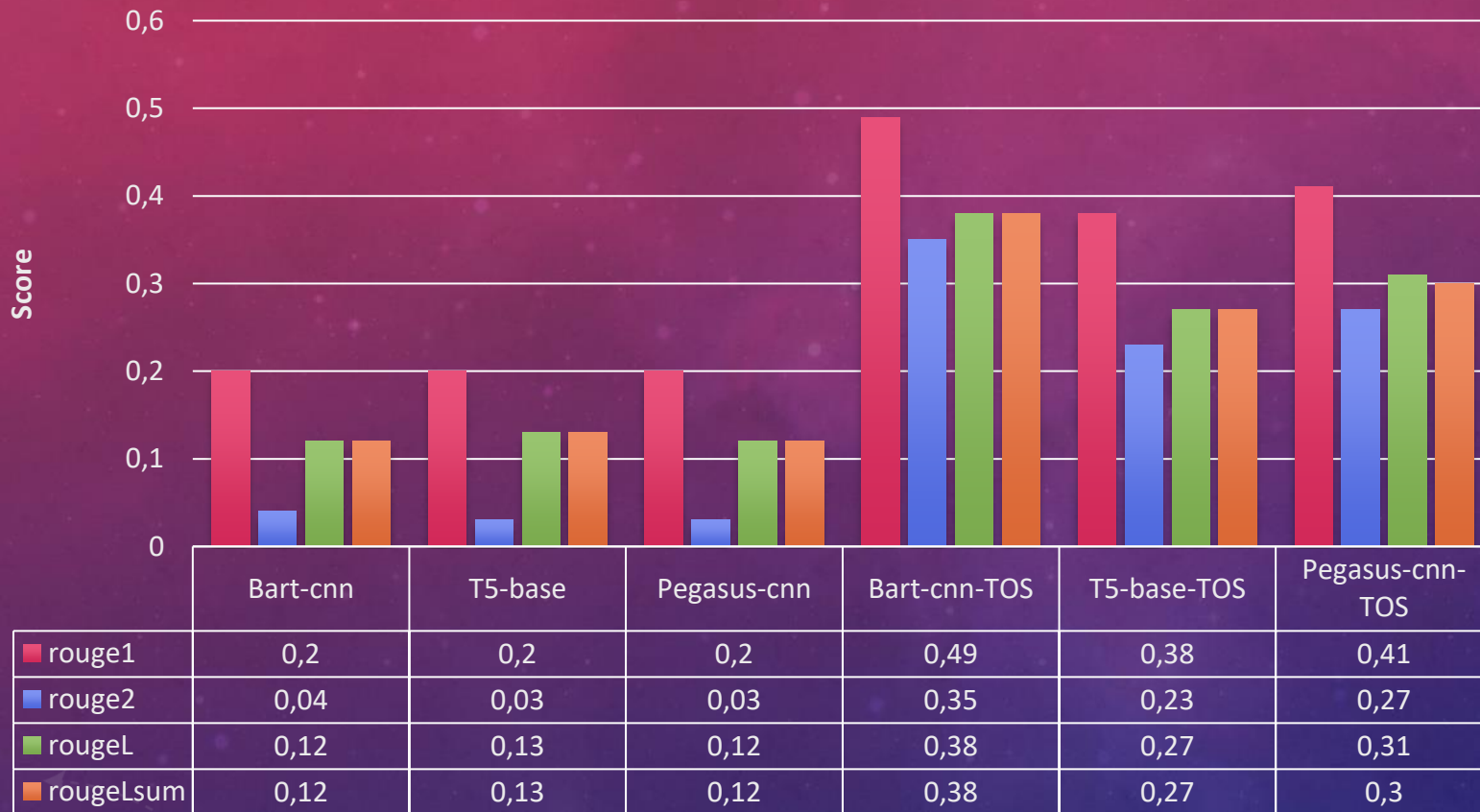
$\text{LCS}(\text{gram}_n)$

$\text{count}(\text{gram}_n)$

$$\frac{2}{7} = 29\% \text{ precision}$$



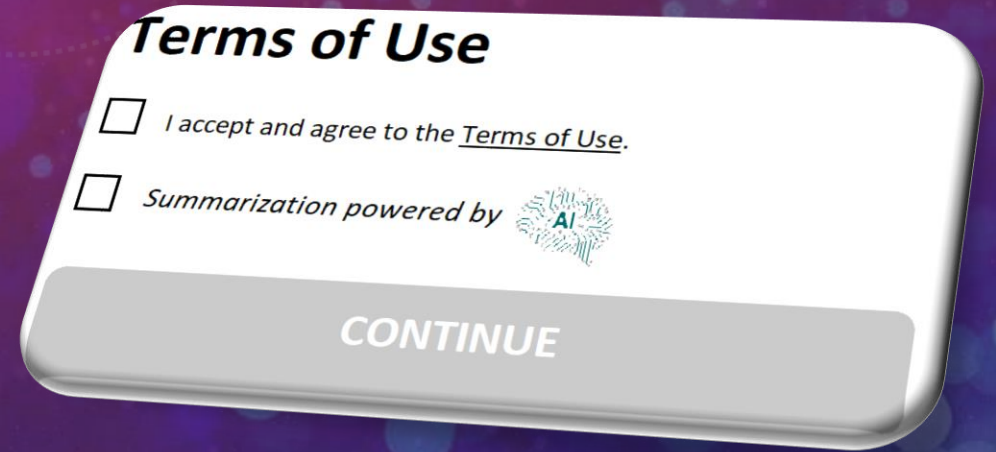
MODELS COMPARISON



- The limitation of these summarization models is that they have a maximum **input length** in terms of context tokens.
- Unfortunately, there is no single strategy to solve this problem, and it's an open and **active research question**.
- We adopted the solution to **split the text** in blocks smaller than the maximum number of tokens accepted as input by the models, being careful not to cut off sentences, and apply recursively the computation of summaries.

CONCLUSION

- Provide the user with a mechanism to summarize and extract key information from their contracts to be clear and transparent towards customers.
- The company can easily embed it in its web pages, in particular in contract boxes, in this way the user is forced to read the summarization in key points of the contract before continuing.
- Companies that operate in the digital industry and are configured as a hub of third-party services, in which users must accept new terms of service each time they purchase a service.



TRY IT!



[Terms Of Service Summarizer Demo](#)