

Newsjam – December Report

Anonymous ACL submission

Abstract

We introduce *Newsjam*, an automatic summarization tool for French COVID-19 news articles. To this purpose, we implement two extractive summarization methods. We then evaluate our tool using two evaluation metrics and two different corpora, an existing one as well as a new one that we have built ourselves.

1 Introduction

The ongoing COVID-19 pandemic has caused a 65% increase in social media use in the United States¹. However, 77% of American adults report experiencing news fatigue² from all of the information coming from news sources and social media³.

These observations have prompted us to create *Newsjam*, a COVID-19 news summarization tool aimed at reducing news fatigue by keeping only the main points of pandemic-related information and by aggregating them in a single place, therefore reducing the need to rummage through the plethora of available information.

Our project consists of three connected parts:

- An article classification model to separate COVID-19 articles into two categories: French and global articles.
- A backend architecture that automatically summarizes these articles.
- Two bots that regularly feed the summarizer with relevant, newly fetched articles, and then posts the output summary as a tweet, meaning our summaries must abide by the maximum tweet length of 280 characters.

¹Statista Research Department (2021), Additional daily time spent on social media platforms by users in the United States due to coronavirus pandemic as of March 2020, Statista

²Watson, A. (2020), Share of adults who feel overwhelmed by the amount of information in coronavirus news coverage in the United States as of April 2020, by age group, Statista

³Nielsen et al. (2021), An Ongoing Infodemic: How People in Eight Countries Access and Rate News and Information about Coronavirus a Year into the Pandemic, Reuters Institute

2 Literature Review

Automatic text summarization refers to generating a condensed version of a text document while preserving key information. There are two primary text summarization methods, extractive and abstractive (Saggion and Poibeau, 2013). Extractive summarization centers around identifying key sentences in the text and putting them together verbatim, whereas abstractive summarization involves generating novel text. Extractive approaches include assigning importance scores to sentences using topic modeling, k-means clustering, and latent semantic indexing (Allahyari et al., 2017). Primary approaches seen in abstractive text summarization include the use of deep learning, RNN encoder-decoders, Gated Recurrent Units, and Long Short-Term Memory (Suleiman and Awajan, 2020).

Text summarization comes with a few key challenges. During the testing stage, reference summaries are needed for evaluation. However, datasets often contain poor reference summaries or do not contain them at all, making evaluation unreliable or impossible (Suleiman and Awajan, 2020). Other challenges include the occurrence of out-of-vocabulary (OOV) words that are absent from the training dataset but are central to understanding a document, and finding metrics able to evaluate a summarized or paraphrased fragment on both a syntactic and semantic level (Soto, 2021).

Given the multiple challenges encountered in text summarization, the question arises as to how to validate the results. The most commonly used metric is ROUGE (Recall-Oriented Understudy for Gisting Evaluation), which compares a generated summary to a gold one (typically created by a human) and calculates the number of overlapping units (Lin, 2004a). ROUGE is however far from ideal: Dorr et al. (2005) found that the metric was sensitive to the type of summarization. ROUGE scores also tend to be higher for summaries that are

longer or generated by using supervised learning approaches (Shulter, 2017).

3 Methodology

3.1 Scraping

Several French news websites have been selected for scraping. One custom scraper is required for each website due to their different underlying structure. Our project currently contains two working Python scrapers for two sites: one for *Actu*⁴, using the *selenium*⁵ library, and the other for *L'Est Républicain*⁶ using the *beautifulsoup*⁷ library.

For each website, we retrieve articles tagged with the “COVID-19” label or the closest fitting category (such as “Health”). Summaries are also extracted from each article when available.

In both cases, the output of these scrapers is a JSON file that contains the same fields as the MLSUM corpus (Scialom et al., 2020). These fields have been chosen for data consistency as MLSUM is part of our experimental setup (see 4). The scrapers also work on an iterative basis and only fetch new articles that have not yet been retrieved.

Before publicizing the summaries of our scraped articles on the planned Twitter bot, the legal aspect will need to be taken into account. This will have to be done on a case-by-case basis to ensure full compliance with the wishes of the owners of the intellectual property contained on these news sites. If using the summaries of these articles in a public setting proves too difficult, we will simply make the Twitter bots private accounts as a proof of concept.

3.2 Annotation

Part of our pipeline consists of implementing a binary news classifier, with the aim to create two distinct Twitter bots so that readers could choose whether they want to familiarize themselves with updates on France or the rest of the world. To this purpose, we manually annotated the entire *L'Est Républicain* corpus. Articles were tagged “1” for local, French news and “0” for global news. The latter category includes not only news about other countries but also global events, such as decisions made by the World Health Organization. As of now, the dataset used for classification is rather

imbalanced (65% of articles are local and 35% are global). The imbalance of classes will be tackled by adding more data to the corpus in the future.

3.3 Article Classification

Our pre-processing of the corpus for classification begins with noise removal (punctuation and irrelevant special characters), stopwords removal, and lemmatization. We then implement three different methods for classification: Multinomial Naive Bayes (MNB), Logistic Regression (LR), and Support Vector Machine (SVM). Theoretically, the MNB classifier is based on Bayes’ theorem and assumes strong (naive) independence of features (Kowsari et al., 2019). The LR classifier uses a logistic function to model and predict the dependent variable that can take the value of 0 or 1. Lastly, the idea behind the SVM classifier is to represent data entries as points in a high dimensional space and find a hyperplane with the largest margin to separate as many points as possible.

For model-specific issues, one way to improve the performance is hyperparameter tuning, which can be done with the *GridSearch CV* method. MNB may benefit from tuning the α parameter that represents Laplace smoothing and helps tackle issue of zero probabilities. For our LR model, we tasked GridSearch CV with finding the best *C* parameter. It appeared to be a lower value, which is expected as our imbalanced dataset requires more regularization and more weight to the complexity penalty to avoid overfitting. Regarding the SVM model, we tuned the *C* parameter as well to find a balance between the minimum margin and accounting for outliers in the data.

3.4 Summarization

Our project currently implements two extractive approaches for summarization, Latent Semantic Indexing and K-means clustering.

3.4.1 Latent Semantic Indexing

Latent Semantic Indexing (LSI) is a technique initially introduced for automatic document classification (Borko and Bernick, 1963) and information retrieval (Deerwester et al., 1990), but it was later found to be efficient for automatic text summarization (Gong and Liu, 2001; Ozsoy et al., 2011) and its evaluation (Steinberger et al., 2004; Steinberger and Jezek, 2009). We choose to implement LSI due to its popularity, historical significance and relative simplicity.

⁴<https://actu.fr/societe/coronavirus>

⁵<https://pypi.org/project/selenium/>

⁶<https://estrepublikain.fr/sante/coronavirus>

⁷<https://pypi.org/project/beautifulsoup4/>

Our algorithm replicates that of [Gong and Liu \(2001\)](#), including for sentence selection. Our implementation was initially based on [Chakravarthy \(2020\)](#) as a starting point, but has since largely diverged. For each article, we choose the optimal number of topics n_{topics} by measuring C_v topic coherence ([Röder et al., 2015](#)). We pick:

$$n_{topics} = \arg \max_{k \in [2,10]} C_v(k)$$

The output summary is generated by concatenating the chosen sentences by order of decreasing score until we reach the aforementioned maximum summary length of 280 characters. Highly-ranked sentences are skipped in favor of lower-ranked ones if the former would make the summary go over the character limit and the latter would not.

We decide to apply LSI not to the TF-IDF (term frequency-inverse document frequency) of the raw text, but to that of a list of keywords generated by removing punctuation and stopwords from the article and stemming the remaining words. Our intent in doing so is to eliminate noise that could be caused by stopwords and to apply topic modeling to the most semantically loaded words in a document.

3.4.2 Word embeddings and k-means clustering

K-means clustering is an alternative way to model topics within a document, which has successfully been applied to text summarization before ([Shetty and Kallimani, 2017](#)).

Our implementation uses contextual word embeddings as the raw input of k-means clustering, ([Gupta, 2020](#); [ialifinaritra, 2021](#)) which are generated on a sentence basis using the models FlauBERT ([Le et al., 2019](#)) and CamemBERT ([Martin et al., 2020](#)). We arbitrarily choose $n_{clusters} = 5$ for k-means clustering, but further research is necessary to determine an optimal value.

We then use the same algorithms as in 3.4.1 for selecting sentences and building the final summary, with the modification that scores are generated for individual words (and not sentences) by computing the cosine similarity between a word’s embedding and topic centroids. For each cluster, the top words are picked and mapped back to the original sentence that contains them.

The choice of K-means clustering is motivated by the simplicity and ease of implementation of the method, while embeddings are chosen with the intuition that they capture additional semantic

information. Specifically, we use contextual embeddings with the hope that they would perform better than a bag-of-word model such as TF-IDF.

4 Experimental setup

4.1 Datasets

Our own corpus, created via scraping news websites, contains 895 articles extracted from *Actu* and 1,753 articles from *L’Est Républicain* along with their summaries.

In addition, we also use the French version of the MultiLingual SUMmarization corpus (MLSUM) ([Scialom et al., 2020](#)) which is made up of over 400,000 news articles from *Le Monde*, split into train, test, and validation sets. The title and content of each article are stored along with a corresponding summary for evaluation, which usually comes from a highlights section found at the beginning of each article. Our models are evaluated on the test set of this corpus, which contains 15,828 articles.

We choose this dataset to evaluate our program due to its semantic proximity with ours (recent news articles, in French, extracted from online newspapers) as well as the large amount of data it contains, which we believe will give us more reliable scores than our own, smaller dataset.

4.2 Article Classification Evaluation

In a machine learning process, it is often insufficient to split the dataset into training and testing subsets and assume that the model will always perform well on unseen data. For this reason, we used the *cross_val_score* method from the *sklearn* library. This method divides data into a user-defined number k of sets, out of which $k - 1$ are used for training and the remaining one is used for testing. The choice of the testing set changes with each execution. We ran a five-fold validation on all models and collected the values of all four metrics: accuracy, precision, recall, and F1-Score.

Since MNB classifiers are known to struggle with correctly predicting classes when trained with imbalanced datasets ([Kowsari et al., 2019](#)), we re-sampled the data to balance the classes in an exact 1:1 proportion. The LR and SVM models were evaluated before and after the tuning.

4.3 Summarization Evaluation

To evaluate our models, we use ROUGE ([Lin, 2004b](#)) as well as the much more recent, state of the art metric BERTScore ([Zhang et al., 2020](#)). Both

metrics calculate precision, recall, and F1-score but via two different means. The three scores are proportions ranging from zero to one, where a higher score indicates higher performance. For our evaluation, we are focusing on the F1-score, which is calculated on two different forms of our data:

- The *standard score* is computed on pairs consisting of the generated summary and the reference summary.
- The *keyword score* is computed on pairs consisting of stemmed keyword-only versions of the generated and reference summaries. Keyword extraction is described in 3.4.1.

To our knowledge, it is not common practice for text summarization models to evaluate keyword versions of summaries. Our goal in adding this evaluation is to reduce the risk of score inflation due to stopword similarity. We also expect stemming to increase the opportunity for matching word subsequences between our generated summaries and the references, which is the same motivation for the METEOR evaluation metric (Lavie and Agarwal, 2007) to include a Porter stemmer module.

4.3.1 ROUGE-L

The ROUGE-L score (Lin, 2004b) is a specific version of ROUGE denoted by the hyphenated 'L', which stands for "Longest Common Subsequence" (LCS). As implied by its name, this version of ROUGE tries to find the longest common sequence of words between a generated summary and a reference one. Scores are then computed based on the LCS length in a pair of summaries. In other words, if two summaries have a long LCS in common, they will receive high scores and vice versa. The choice of ROUGE-L was motivated by its prevalence, which allows us to compare the performance of our models with the existing literature.

4.3.2 BERTScore

BERTScore (Zhang et al., 2020) is a metric originally created for text generation and translation, but due to it being based on comparison of a generated and a reference text, it can also be applied to summarization. We use it as our second evaluation metric since it evaluates the text on a deeper level, by using BERT contextual embeddings (Devlin et al., 2019) to compare the semantic content between the summaries rather than their surface-level similarity using n-grams, like in many versions of ROUGE (Lin, 2004b) and BLEU (Papineni et al., 2002). BERTScore computes scores via cosine similarity of the contextual embeddings.

4.3.3 Execution time

Due to limited computing resources, time, and the relatively large size of the MLSUM corpus, we felt important to measure the average execution time of our summarization methods. This also allows us to track speed improvements made over time by optimizing various aspects of our program.

For each dataset and summarization method, we summarize the first 100 articles of the test split of the dataset and measure the average elapsed time per article. This process is repeated 10 times to mitigate the statistical variance of time measurements and we report the average value of this execution time. A lower value reflects a more useful summarization method for our purposes. All time measurements were applied on the same machine, and GPU processing was enabled whenever possible.

5 Results

5.1 Article Classification Results

Our article classification results are in Table 1.

At first, the MNB classifier generated a considerable number of false negatives and had less than ideal accuracy at 72.1%. After tuning the α parameter, the MNB's performance improved across

Method	Accuracy	Precision	Recall	F1
Multinomial Naive Bayes (MNB)	72.1	70.8	97.7	82.1
MNB (resampled)	67.1	88.3	38.6	53.1
MNB (tuned)	83.1	92.2	81.1	86.2
Logistic Regression (LR)	83.8	88	87.3	87.7
LR (tuned)	85.7	88.3	90.5	89.3
Support Vector Machine (SVM)	82.8	83.3	91.3	87.1
SVM (tuned)	85.2	87.5	90.6	89.1

Table 1: Article classification evaluation

Method	ROUGE-L	Keyword ROUGE-L	BERTScore	Exec. time
<i>MLSUM corpus, testset (15,828 articles)</i>				
LSI	0.1507	0.1147	TBA	TBA
FlauBERT + k-means	TBA	TBA	TBA	TBA
CamemBERT + k-means	TBA	TBA	TBA	TBA
<i>Built corpus (895 + 1,753 = 2,648 articles)</i>				
LSI	0.5391	0.5309	0.5666	TBA
FlauBERT + k-means	0.2911	0.2862	0.2098	TBA
CamemBERT + k-means	0.2463	0.2290	0.2941	TBA

Table 2: Summarization evaluation (only F1-Scores reported)

the board. While resampling boosted the precision of the model from 70.8% to 88.3%, other metrics plummeted below acceptable levels. Overall, MNB and LR each provide maximum values depending on the metric. When it comes to accuracy, the tuned LR and SVM models show the best performance with a negligible difference, 85.7% and 85.2% accordingly. The precision metric is lead by the tuned MNB with 92.2%. The highest recall is recorded by our initial MNB model at 97.7%. Finally, the top F1-Score goes to the tuned LR with 89.3%.

5.2 Summarization Results

Our summarization results are in Table 2.

We generally observe much higher scores on the *Actu* corpus than the *MLSUM* corpus. We also find that the LSI model clearly outperforms the k-means clustering implementations.

6 Discussion

Since there is no one method that performs high and above all the others for article classification, we must focus on one metric that we deem to be the most important. Thus we will focus on the *F1-Score* metric. We justify this choice with our goal to maximize the amount of true positives while minimizing the amount of false positives and false negatives, which translates to minimizing the misidentification of local and global articles in our model. Given this choice, we believe that the tuned LR method will work best for our article classification.

For our summarization evaluation, the high scores from the *Actu* corpus should not be taken too seriously for the performance of our model due to its small size. On the other hand, even though the *MLSUM* dataset has issues that will be discussed later in this section, we believe the evaluation on this corpus better reflects the performance of our models, simply because there is much more data.

With that said, we also manually looked into a few summary selections to observe the reliability of these scores ourselves.

While looking deeper into the *Actu* LSI examples of the data, we found that the summaries chosen by our model outlined the articles well and matched a quality reference summary. In the *MLSUM* LSI examples, we noticed a recurring issue where we felt that our summary matched the article, but received a low score due to a poor reference summary. This seems inevitable when working with a corpus of this magnitude, but it is important to note because it exhibits a way in which the scores may not always reflect the quality of a generated summary. Finally, one issue that has affected both datasets is poor summary selection, which may be due to the performance of the summarization methods themselves. Examples of some of these phenomena in our results can be found in Appendix A.

7 Conclusion and Future Works

We have successfully created a new corpus that is specific to French COVID-19 articles. This corpus has been annotated and used to train three different binary classifiers on the articles. Our summarization task currently contains two different methods and two evaluation metrics. These have been tested on the corpus we have built as well as another, much larger one to get more robust results.

Our classifier evaluation results are satisfactory, with a maximum F1-Score close to 90% in the tuned LR model. On the other hand, our summarization results are currently very low on the largest corpus, but early results on our own smaller corpus are promising. In the future, we will automate more processes within our project, finish evaluation of our summarization models and finally build the Twitter bots necessary to complete our final goal of posting our summaries to Twitter.

References

- Mehdi Allahyari, Seyedamin Pouriyeh, Mehdi Assefi, Saeid Safaei, Elizabeth D. Trippe, Juan B. Gutierrez, and Krys Kochut. 2017. [Text summarization techniques: A brief survey](#). In *Proceedings of arXiv, USA*, pages 1–9.
- Harold Borko and Myrna Bernick. 1963. [Automatic document classification](#). *Journal of the ACM (JACM)*, 10(2):151–162.
- Srinivas Chakravarthy. 2020. [Document Summarization Using Latent Semantic Indexing](#).
- Scott Deerwester, Susan T. Dumais, George W. Furnas, Thomas K. Landauer, and Richard Harshman. 1990. [Indexing by latent semantic analysis](#). *Journal of the American Society for Information Science*, 41(6):391–407.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [Bert: Pre-training of deep bidirectional transformers for language understanding](#).
- Bonnie J. Dorr, Christoph Monz, Stacy President, Richard Schwartz, and David Zajic. 2005. A methodology for extrinsic evaluation of text summarization: Does rouge correlate? In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 1–8.
- Yihong Gong and Xin Liu. 2001. [Creating generic text summaries](#). In *Proceedings of Sixth International Conference on Document Analysis and Recognition*, pages 903–907.
- Akanksha Gupta. 2020. [Understanding Text Summarization using K-means Clustering](#).
- ialifinaritra. 2021. [French Text Summarization](#). Original-date: 2020-10-07T21:22:46Z.
- Kamran Kowsari, Kiana Jafari Meimandi, Mojtaba Heidarysafa, Sanjana Medu, Laura Barnes, and Donald Brown. 2019. Text classification algorithms: A survey. *Information*, 10(4):150.
- Alon Lavie and Abhaya Agarwal. 2007. [METEOR: An automatic metric for MT evaluation with high levels of correlation with human judgments](#). In *Proceedings of the Second Workshop on Statistical Machine Translation*, pages 228–231, Prague, Czech Republic. Association for Computational Linguistics.
- Hang Le, Loïc Vial, Jibril Frej, Vincent Segonne, Maximin Coavoux, Benjamin Lecouteux, Alexandre Allauzen, Benoît Crabbé, Laurent Besacier, and Didier Schwab. 2019. [Flaubert: Unsupervised language model pre-training for french](#). *CoRR*, abs/1912.05372.
- Chin-Yew Lin. 2004a. Looking for a few good metrics: Rouge and its evaluation. In *Ntcir Workshop*.
- Chin-Yew Lin. 2004b. Rouge: A package for automatic evaluation of summaries. In *ACL 2004*.
- Louis Martin, Benjamin Muller, Pedro Javier Ortiz Suárez, Yoann Dupont, Laurent Romary, Éric de la Clergerie, Djamé Seddah, and Benoît Sagot. 2020. [CamemBERT: a tasty French language model](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7203–7219, Online. Association for Computational Linguistics.
- Makbule Ozsoy, Ferda Alpaslan, and Ilyas Cicekli. 2011. [Text summarization using latent semantic analysis](#). *J. Information Science*, 37:405–417.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [Bleu: A method for automatic evaluation of machine translation](#). In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics, ACL '02*, page 311–318, USA. Association for Computational Linguistics.
- Michael Röder, Andreas Both, and Alexander Hinneburg. 2015. [Exploring the space of topic coherence measures](#). *WSDM 2015 - Proceedings of the 8th ACM International Conference on Web Search and Data Mining*, pages 399–408.
- Horacio Saggion and Thierry Poibeau. 2013. Automatic text summarization: Past, present and future. In *Multi-Source, Multilingual Information Extraction and Summarization*, pages 3–21. Springer, Berlin, Heidelberg.
- Thomas Scialom, Paul-Alexis Dray, Sylvain Lamprier, Benjamin Piwowarski, and Jacopo Staiano. 2020. Mlsum: The multilingual summarization corpus. *arXiv preprint arXiv:2004.14900*.
- Krithi Shetty and Jagadish S. Kallimani. 2017. [Automatic extractive text summarization using k-means clustering](#). In *2017 International Conference on Electrical, Electronics, Communication, Computer, and Optimization Techniques (ICEECCOT)*, pages 1–9.
- Natalie Shulter. 2017. The limits of automatic summarisation according to rouge. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 41–45.
- William Soto. 2021. X-pareval: A multilingual metric for paraphrase evaluation. Master’s thesis, Université de Lorraine.
- Josef Steinberger and Karel Jezek. 2009. Evaluation measures for text summarization. *Computing and Informatics*, 28:251–275.
- Josef Steinberger, Karel Jezek, et al. 2004. Using latent semantic analysis in text summarization and summary evaluation. *Proc. ISIM*, 4:93–100.

Dima Suleiman and Arafat Awajan. 2020. [Deep learning based abstractive text summarization: Approaches, datasets, evaluation measures, and challenges.](#) *Mathematical Problems in Engineering*, 2020:1–29.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. [Bertscore: Evaluating text generation with bert.](#)

A Examples of LSI model outputs

1. Quality Summary/Poor Score

- MLSUM Test Dataset Article 54
- *Reference Summary*: "Le suspect principal, un employé des services de la ville, a tiré « à l'aveugle ». Il est lui aussi décédé."
- *Generated Summary*: "Douze personnes ont été abattues vendredi 31 mai par un tireur dans un bâtiment municipal de Virginia Beach (Etat de Virginie), station balnéaire de la côte est américaine."
 - ROUGE-L
 - * Standard F1-Score: 0.151
 - * Keyword F1-Score: 0.066
 - BERTScore
 - * Standard F1-Score: 0.157
 - * Keyword F1-Score: 0.119
- This example above is important, because it shows that our model is able to select high quality summaries, but that they will not always be evaluated as such, because some of the reference summaries in the MLSUM corpus are not up to the expected gold standard.

2. Poor Summary/Poor Score

- MLSUM Test Dataset Article 21
- *Reference Summary*: "Le journalisme est un métier qui nécessite parfois de faire face aux situations les plus extrêmes. Aujourd'hui, c'est à une expérience en lisière de l'anthropophagie que je dois me livrer, puisque ma mission consiste à aller déguster le « Chirac », un burger portant le nom d'un ancien président de la République amateur de tête de veau."
- *Generated Summary*: "Le vaillant cobaye qui m'accompagne et moi mordons à pleines dents dans cet agglomérat de fromage de Normandie dégoulinant, de pommes confites et de sauce à la moutarde."
 - ROUGE-L
 - * Standard F1-Score: 0.146
 - * Keyword F1-Score: 0.038
 - BERTScore
 - * Standard F1-Score: 0.018
 - * Keyword F1-Score: 0.036

- This is an example where our generated summary is quite poor as the main subject of the article, a burger named the "Chirac", is never even mentioned, but a very in-depth description is included without any context.

3. Quality Summary/Quality Score

- Built Corpus Actu Article 36
- *Reference Summary*: "A partir de ce lundi 11 octobre 2021, le port du masque en extérieur n'est plus obligatoire dans le Pas-de-Calais. Explications."
- *Generated Summary*: "Le port du masque en extérieur n'est plus obligatoire dans le Pas-de-Calais. A partir de ce lundi 11 octobre 2021, le port du masque en extérieur n'est plus obligatoire dans le Pas-de-Calais."
 - ROUGE-L
 - * Standard F1-Score: 0.738
 - * Keyword F1-Score: 0.743
 - BERTScore
 - * Standard F1-Score: 0.711
 - * Keyword F1-Score: 0.521
- This last example is one where we have a high quality generated summary that receives a relatively high F1-Score. The summary itself could still be improved by being less repetitive, which may in turn improve the score. This problem of repetitiveness is one we often found with our generated summaries from this corpus and we plan to look into it further. However, overall we are very satisfied with the information that is conveyed through generated summaries like this one.