

# Identifying Similarities in Contracted Object Description on Public Procurement

Wesley Lima  
Dept. of Computing  
UFPI  
Teresina, Brazil  
wesley@ufpi.edu.br

Victor Silva  
Dept. of Computing  
UFPI  
Teresina, Brazil  
victor.silva@ufpi.edu.br

Jasson Silva  
Dept. of Computing  
UFPI  
Teresina, Brazil  
jasson\_jcs@ufpi.edu.br

Ricardo Lira  
Dept. of Computing  
UFPI  
Teresina, Brazil  
ricardoalr@ufpi.edu.br

Anselmo Paiva  
Depto. of Informatics  
UFMA  
São Luís, Brazil  
paiva@nca.ufma.br

**Abstract**—The digitization of public procurement has led to the automation of data analysis in public administration. However, using natural language in contract details still poses challenges in analysis and auditing. The large volume of unstructured data in bidding notices calls for tools that automate data analysis for citizens and control bodies. The description of the tender object is a critical problem in data analysis, as it lacks standardization and often contains unclear content. Detecting similarities between contract objects is crucial for analyzing public purchases and identifying fraudulent practices. However, manually reading all contracts to notice similarities is tedious and impractical. Therefore, automating the identification of similarities is essential. We propose a solution based on deep learning that uses the Euclidean distance between contextualized embeddings of hiring object descriptions to detect similarity. To improve performance, Named Entity Recognition (NER) is used to extract the central idea of the contract object, excluding confusing words and expressions. This approach improves the detection of textual similarity in public procurement objects. Overall, these advancements aim to facilitate quick and timely analysis of public contracts.

**Index Terms**—BERT, fine-tuning, NLP, public procurements, NER

## I. INTRODUCTION

Digital transformation has played a significant role in public procurement, bringing a series of improvements to government processes. This technological revolution has impacted various aspects of government procurement, from operational efficiency to transparency and broad competition. The digitization of public procurement has allowed it to automate data analysis by auditing authorities and oversight in public administration.

Public procurement must be done through bidding notices [1]. This process involves the publication of invitations to bid by government or public entities, detailing the need for goods, services, or construction. Interested suppliers submit bids, which are publicly opened and examined. The bids undergo evaluation based on criteria such as price and technical specifications. The contract is then awarded to the successful bidder, and both parties execute the contract.

Bid notices are published as textual documents written in natural language that detail the procurement conditions (price, technical specification, and others). Auditing authorities must analyze these documents to guarantee transparency, fairness, and competitiveness throughout the process, optimizing the use of public funds.

The textual representation of bid notices makes analysis and timely auditing difficult. Thus, the analysis and auditing of all procurement processes cannot be done properly, opening up the possibility of fraudulent processes. The large number of contracting processes comprising a large amount of unstructured data brings to light the need for tools that automate data analysis by citizens and control bodies.

A fundamental problem in analyzing the data is related to the description of the object of the bid notice: a concise and complete statement of the good or service contracted. Like any text in natural language, they do not follow a pre-established standard, and, in many situations, the content is unclear and difficult to understand.

Categorizing contracts is fundamental to optimizing the use of the vast resources involved. In the case of textual documents, the natural approach is to detect similarities between descriptions of contracted objects.

Detecting similarities between the bid notice objects is too an essential task in analyzing public procurement, as it enables both the management analysis of the goods contracted by public entities and the extraction of evidence of fraud related to the acquired object, such as the split purchase [2]. Splitting a purchase is the illegal practice of dividing a contract into several small contracts to carry out less competitive procurement processes. To identify this type of practice, the specialist overseeing public purchases must read all the contracts already awarded by the public entity in a financial year to detect similarities. It is a tedious and, in many cases, unfeasible task, given the high volume of bidding processes. Therefore, automating the identification of similarities between bidding objects will make it possible to analyze public contracts quickly and promptly.

In this context, contextualized sentence representations can significantly improve similarity detection between bidding objects, as they allow models to take semantics into account instead of just lexis.

Techniques based on deep learning using contextualized embeddings are state-of-the-art for detecting the semantic similarity of sentences, as they can capture nuances in the meaning of words, dealing better with polysemy, synonyms, and antonyms [3].

The spatial distance between embeddings is commonly

the semantic similarity measure between sentences or words [4]. Standard methods for calculating this similarity include Euclidean distance and cosine similarity [5]. State-of-the-art strategies use pooling operations based on the average of the embeddings representing each word to generate sentence embeddings [6], [7].

However, according to Schockaert [8], combining different text fragments by pooling their embeddings may require a prohibitively high dimensionality if the resulting embeddings are to capture the epistemic state faithfully. Therefore, the average pooling of embedding words to represent a sentence can lead to incorrect similarity calculations. Intending to mitigate this problem and improve the search for similarity between sentences, we propose to extract the keyphrase from the bidding notice object.

The primary goal of this research is to improve the detection of similarity between the descriptions of public procurement objects using deep learning models based on contextualized embeddings. The proposed method achieved significant results, with an F1 score of over 95%.

The rest of this article is structured as follows. In the next section, we detail the main works related to the treatment of the description of the object of public procurement. Only some studies automatically exploit the description to extract meaningful information for expert analysis. In the Methodology section, we describe all the steps of the proposed method. In the Experiments and Results section, we detail how we evaluated the proposed method and the results, which show significant improvements in identifying similarities. Finally, section V provides the article's conclusions, contributions, limitations, and future works.

## II. RELATED WORK

Looking for methodologies that use the description of the object of the contract in some way, we found that most of the efforts are to structure contract data [9], [10], creating predictive models [11], [12] or classifying contracts based on the description of the contracted object [13], [14].

Beyond the classifier, Kayte et al. [13] propose a method for automatically generating the title of the tender notices from the user-given text. Using an LSTM model, they reach an accuracy of 97% for text generation and 95% for code classification using a Support Vector Machine(SVM). On the other hand, Navas-Loro et al. [14] created a multi-label classifier that assigns a Common Procurement Vocabulary(CPV) code based on the description of the contracted object. However, the classifier could have performed better for samples with generic and overlapping CPV codes.

Regarding structuring contract data, Potin et al. [9] describe the challenges of creating a public procurement database to predict fraud. As the main issues, the author cites missing data, non-standardized data (for example, several winners of a lot described in the same field), inconsistencies in names and addresses, description of the criteria for choosing the bidder without standardization, and lack of identification of bidders. While this, Costetchi et al. [10] structure public contracts

data extracted from electronic journals using an ontology and propose a methodology to assess the quality of the mapped data.

More specifically, Gomes et al. [15] carry out a systematic mapping to identify and characterize the methods, techniques, and intelligent algorithms for classifying incongruous textual descriptions in invoices. The research showed that using artificial intelligence techniques helped mitigate the problem of classifying and analyzing invoices with incongruous codes and descriptions.

Two researchers used the description of the contracted object to train prediction models. The work of Acikalin [11] tries to predict procurement outcomes based on their description. The model makes indirect predictions, i.e., the red flag itself is not in the description but in other data collected, such as the number of companies taking part, the difference between the value won, and the estimated value, among others. The results suggest that notice descriptions play an important role in the outcomes of public procurement calls.

Still, from a predictive model perspective, Gorgun et al. [12] investigate the impact of the call for proposals on the number of bidders in public procurement. Through a model based on features extracted from the text, such as the number of existing named entities, predict the number of participants. They investigated six different feature groups, and the experiments showed that n-grams yield the best prediction accuracy, suggesting that the content of procurement notices is fundamental in predicting the number of bidders.

Despite the relevance of the object's description for the analysis of public procurement, the literature review showed that few studies deal with the subject.

The search for semantic similarity proved vulnerable to the words that make up the text and are not directly related to the contracted good or service. Therefore, extracting the text's central idea before searching for similarity proved more promising. In NLP literature, this task is called keyphrase extraction.

In their Survey, Song et al. [16] classify the works that use neural models to extract keyphrases from a text into two groups:

- One-stage methods that treat the problem as a sequence labeling task. Its limitation is its inability to deal with overlapping key phrases.
- Two-stage methods that initially operate with heuristics to extract candidate keyphrases and then use a supervised or unsupervised method to choose the final keyphrase. However, models based on supervised learning depend on the costly work of gathering domain-specific characteristics.

Mu et al. [17] propose yet a third group, the so-called generational methods, based on "sequence-to-sequence" models [18] used in the field of machine translation. However, such methods cannot effectively use context and phrase-level information to construct keyphrases.

To the best of our knowledge, only Gero et al. [19] use an NER model as a component of a method for extracting keyphrases. Their work - applied to the biomedical domain -

used an NER model to remove keywords from the domain. These words were combined with phrases extracted from the text to create a set of candidate phrases. Our work is, therefore, the first attempt to use a supervised NER model to extract keyphrases.

Therefore, it is possible to observe that there is a lack of studies that determine the textual similarity of the description of the tendered object and, consequently, guarantee more incredible speed in the expert's analysis. An opportunity still exists to investigate using NER to extract keyphrases and employ them to categorize object texts in bidding notices.

### III. METHODOLOGY

In this work, we have proposed a deep learning-based method using a supervised NER model to extract the central idea of a bid notice and use the Euclidean distance of sentence embeddings to detect similar objects. Figure 2 shows the pipeline of the proposed solution.

#### A. Preprocessing

Initially, we collected bid notices from public sites in PDF, DOC, and DOCX formats. A set of tools from the Parsr project [20] converted these documents to Markdown Format - a simplified markup language similar to HTML - which preserved the original structure and facilitated the extraction of the section that describes the contracted object.

At this stage, we were able to overcome a challenge which was to extract the correct section accurately. The sections' names and structures are neither standardized nor regulated by law. It is, therefore, expected to find different sections with similar names and repetition of the section name in various parts of the text. After successive optimizations of the extraction scripts, we achieved an acceptable degree of assertiveness.

Each object description section was converted and processed for scripts, eliminating irrelevant content, such as section numbers, tabs, URLs, and special characters. We also corrected some misspelled words. As a result, we obtained a dataset with 2,000 object descriptions.

#### B. Keyphrase Extraction

Using semantic similarity techniques directly on each procurement object revealed disappointing results. After an in-depth investigation of the dataset, we hypothesized that the low performance was due to many terms and phrases not directly related to the contracted good or service, only present to understand the reader better.

Consider a sentence  $\mathbf{S}$ , composed of  $n$  words represented by their contextualized embeddings  $e_1, \dots, e_n$ , and another sentence  $\mathbf{T}$ , composed of  $m$  words represented by their contextualized embeddings  $f_1, \dots, f_m$ . Consider also  $M_s$  and  $M_t$  as the embeddings of sentences  $\mathbf{S}$  and  $\mathbf{T}$ , respectively, calculated through a pooling operation based on the average of the embeddings. Then, it follows that the Euclidean distance  $d$  between  $M_s$  and  $M_t$  is dependent on the average distance between any two vectors  $e_i$  and  $f_i$ . Consequently, in a pair

of sequences  $e_1, \dots, e_n, f_1, \dots, f_m$  where the majority of the vectors have very low Euclidean distances between them, even if a small amount of keywords is vastly different, they distance  $M_s$  and  $M_t$  will still be very small, resulting in an inaccurate measure of similarity.

Let us consider, for example, that in the dataset used in this research, the keyphrases correspond, on average, to 13.6% of the words in the complete text. We have a significant influence of unrelated words in generating the embedding that represents the sentence, leading to the inaccuracy of the models that calculate similarity based on the distance of the embeddings. Therefore, it is essential to note that these accessory words - belonging to the sentence but not related to the topic - can harm the detection of similarity with other sentences.

In public procurement, repetitive expressions that make up the bid notice objects only introduce the procurement object. Therefore, to achieve a better performance in the use of textual similarity techniques, it is essential to extract only the central idea of the object of the contract.

To deal with this problem, we used a named entity recognition model to extract the key phrase from the object, excluding words and expressions that confuse models based on contextualized embeddings such as BERT. Consider, for example, the two tender objects in Figure 1.

It is possible to observe that they are entirely different objects, but the repetition of the passage: "*The present tender has as its object*" in the two descriptions makes it challenging to detect similarity by existing strategies in the literature that take into account the embeddings of the entire text and not just the passages that describe the focus of the procurement.

To train the NER model, we create a dataset composed of 2,000 bid notice objects and, using an open-source platform called Argilla [21], we annotate the keyphrase of each sample that the model must recognize as a named entity.

The model resulting from the training takes the full description of the contract object as input and extracts the key phrase. A new dataset stores the keyphrases used in the next step by the similarity detection module.

#### C. Similarity search

The model converts the keyphrase into a contextualized embedding and uses it to search for similar objects through Euclidean distance from the embeddings of other sentences.

To perform a similarity search for a target query  $q$ , all the sentences involved must go through this conversion process, and the system must calculate the Euclidean distance from  $q$  to all the other sentences.

Public procurement datasets usually comprise thousands and even millions of tender notices, making keeping all these high-dimensional vectors in memory prohibitive. Methods based on product quantization codes [22] are suitable to overcome this limitation. They are effective techniques for compressing high-dimensional vectors, such as contextualized embeddings, making searching large volumes of data with reduced memory usage possible.

The purpose of this tender is to supply hospital cleaning materials...

The purpose of this tender is to provide electronic surveillance services...

Fig. 1. Example of keyphrases

Associated with product quantization, using an inverted file index (IVF) [23], enables an initial reduction of the search scope, improving search times for the nearest neighbor vectors.

We used the FAISS [24] implementation, which uses product quantization and inverted file index aims to perform similarity searches efficiently for large volumes of data.

The system returns a list sorted by the Euclidean distance  $d$  of a target query concerning all the objects in the dataset. To select the set of effectively similar sentences, we set a threshold  $l$  to return only sentences with  $d < l$ .

#### IV. EXPERIMENTS AND RESULTS

The named entity recognition model achieved good results in extracting the procurement focus from the tender object's description. We labeled 2,000 bid objects concerning the central phrase and divided them into a ratio of 80% for training and 20% for testing. We fine-tuned three pre-trained models based on BERT: BERTimbau [25], JurisBERT [26], and mBERT [27]. The results are in Table I.

All the models performed well, especially the model based on the multilingual version of BERT, which achieved an F1 score of 90.5%.

TABLE I  
RESULTS OF THE EXPERIMENTS

Pre-trained Model	Metrics		
	Precision	Recall	F1
BERTimbau	87.4%	90.9%	89.1%
JurisBERT	80.3%	85.9%	83.0%
mBERT	<b>89.4%</b>	<b>91.6%</b>	<b>90.5%</b>

To demonstrate the effectiveness of the proposed method, we created a dataset of 1754 pairs of contracting objects labeled as similar or not. We selected 50% similar pairs and the other 50% not similars. To assess the semantic similarity between the sentences, we used BERTScore [28], a metric that calculates the similarity of candidate and reference sentences using the contextualized embeddings of the tokens that make them up, weighted by the IDF score of each token. It is a metric that varies between 0 (not similar) and 1 (similar). Empirically, we consider sentence pairs with a BERTScore value greater than 0.7 similar. Table II shows the results obtained using the entire text of the contract subject and only the keyphrase extracted by the NER model.

The results show a significant improvement in all the metrics evaluated, especially the F1 score, which is over 65% when using keyphrases instead of the whole text to detect similarity.

TABLE II  
BERTSCORE OF OBJECTS AND KEYPHRASES

Content	Metrics			
	Accuracy	Precision	Recall	F1
Object	57.3%	83.8%	18.24%	29.9%
Keyphrase	<b>95.7%</b>	<b>97.8%</b>	<b>93.5%</b>	<b>95.6%</b>

Extracting the keyphrase from the text has also proved effective in classification tasks. A standard classification of public contracts is according to the type of contracted object. This information is essential in the analysis carried out by control bodies, as each type of contract involves a different set of rules. To demonstrate the effectiveness of keyphrase extraction, we trained two classifiers based on four pre-trained models based on BERT using a dataset of previously classified procurement objects. To train the first model, we used the full text of the object. To train the second model, we used only the key phrase of each object. Table III shows the results obtained.

All classifiers significantly improved by more than 10% in the F1 Score. The most improvement was in the RoBERTa [29] base model, which saw an 18.2% improvement in the F1 Score metric.

TABLE III  
RESULTS OF CLASSIFIERS

Pre-trained Model	Metrics			
	Accuracy	Precision	Recall	F1
<b>Object</b>				
mBERT	85.2%	80.6%	79.1%	79.8%
BERTimbau	88.7%	86.9%	80.3%	82.5%
jurisBERT	87.3%	80.9%	78.7%	79.3%
RoBERTa	80.2%	82.3%	70.4%	72.6%
<b>Keyphrase</b>				
mBERT	<b>94.1%</b>	<b>94.5%</b>	<b>90.9%</b>	<b>92.5%</b>
BERTimbau	<b>96.5%</b>	<b>96.0%</b>	<b>92.7%</b>	<b>94.1%</b>
jurisBERT	<b>94.5%</b>	<b>93.5%</b>	<b>92.0%</b>	<b>92.7%</b>
RoBERTa	<b>93.0%</b>	<b>93.8%</b>	<b>88.7%</b>	<b>90.8%</b>

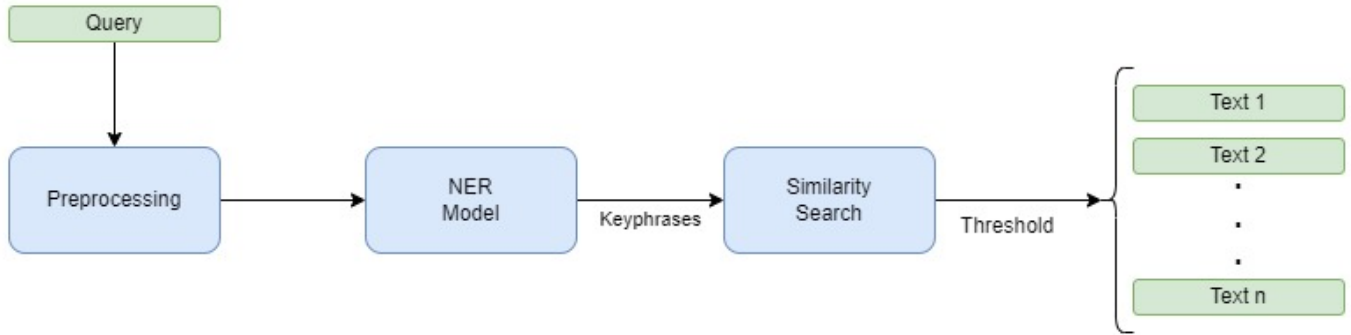


Fig. 2. Pipeline Model

To demonstrate how words that are not part of the keyphrase can harm classification, we used the Transformers Interpret [30], a model explainability tool, which highlights words that hurt a prediction in red and those that have a positive impact in green. Figure 3 (*Text in Portuguese*) shows that the model highlights directly related words to correctly classify a keyphrase in the 'engineering works' class, such as 'engineering', 'pothole', and 'road'. However, in Figure 4 (*Text in Portuguese*), the model wrongly classifies the full text of the object in the 'Services' class, influenced by words not directly related to the focus of the contract, such as 'dispensation', 'specialized', and 'specifications'.

## V. CONCLUSION

This article introduces a methodology for identifying similarities between public procurement objects. The central idea is to use an NER model to extract the keyphrase of each object before performing a similarity search using the Euclidean distance between contextualized embeddings.

The experiments showed that combining an NER model with a similarity search significantly improved performance in identifying similar objects.

This work improves the process of analysis of public contracts by specialists, automating the search for similar objects and, consequently, identifying signs of fraud, such as split purchase, and guaranteeing the generation of management reports related to the types of objects contracted.

The main contributions of this work include:

- The first annotated corpus for NER in bidding objects in the Portuguese language.
- The first corpus of sentence pairs of public procurement objects annotated concerning similarity in the Portuguese language.
- A new methodology for detecting textual similarity using NER to extract the core sentence from the object description.

Although the results have been promising, some limitations are evident. The number of manually classified objects used to create the NER model was minor, causing the model to incorrectly exclude or include words in the keyphrase extracted

from the object. In addition, the empirical definition of the similarity threshold sometimes leads to the incorrect exclusion or inclusion of objects in a set of similar documents.

In the future, we intend to increase the data set for training the NER model and study a more efficient strategy for determining the set of similar documents.

## REFERENCES

- [1] Grandia, J., Kuitert, L., Schotanus, F. & Volker, L. Introducing Public Procurement. *Public Procurement: Theory, Practices And Tools*. pp. 1-18 (2023), [https://doi.org/10.1007/978-3-031-18490-1\\_1](https://doi.org/10.1007/978-3-031-18490-1_1)
- [2] IACRC International Anti-Corruption Resource Center. "Guide to Combating Corruption & Fraud in Infrastructure Development Projects", Available online at: <https://guide.iacrc.org/potential-scheme-split-purchases/>
- [3] Han, M., Zhang, X., Yuan, X., Jiang, J., Yun, W. & Gao, C. A survey on the techniques, applications, and performance of short text semantic similarity. *Concurrency And Computation: Practice And Experience*. **33** (2020)
- [4] Salton, G., Wong, A. & Yang, C. A Vector Space Model for Automatic Indexing. *Commun. ACM*. **18**, 613-620 (1975,11), <https://doi.org/10.1145/361219.361220>
- [5] Huang, A. & Others Similarity measures for text document clustering. *Proceedings Of The Sixth New Zealand Computer Science Research Student Conference (NZCSRSC2008)*, Christchurch, New Zealand. **4** pp. 9-56 (2008)
- [6] Reimers, N. & Gurevych, I. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. *Conference On Empirical Methods In Natural Language Processing*. (2019)
- [7] Gao, T., Yao, X. & Chen, D. SimCSE: Simple Contrastive Learning of Sentence Embeddings. *Proceedings Of The 2021 Conference On Empirical Methods In Natural Language Processing*. pp. 6894-6910 (2021,11), <https://aclanthology.org/2021.emnlp-main.552>
- [8] Schockaert, S. Embeddings as epistemic states: Limitations on the use of pooling operators for accumulating knowledge. *International Journal Of Approximate Reasoning*. pp. 108981 (2023)
- [9] Potin, L., Labatut, V., Morand, P. & Largeron, C. FOPPA: an open database of French public procurement award notices from 2010–2020. *Scientific Data*. **10**, 303 (2023,5,19), <https://doi.org/10.1038/s41597-023-02213-z>
- [10] Costetchi, E., Vassiliades, A. & Nyulas, C. Towards a mapping framework for the Tenders Electronic Daily Standard Forms. (KGCW'23: 4th International Workshop on Knowledge Graph Construction, 2023,5)
- [11] Acikalin, U., Gorgun, M., Kutlu, M. & Tas, B. How you describe procurement calls matters: Predicting outcome of public procurement using call descriptions. *Natural Language Engineering*. pp. 1-22 (2023)
- [12] Gorgun, M., Kutlu, M. & Onur Taş, B. Predicting The Number of Bidders in Public Procurement. *2020 5th International Conference On Computer Science And Engineering (UBMK)*. pp. 360-365 (2020)



Legend: <span style="color: red;">■</span> Negative <span style="color: gray;">□</span> Neutral <span style="color: green;">■</span> Positive				
True Label	Predicted Label	Attribution Label	Attribution Score	Word Importance
0	Obras de engenharia (1.00)	Obras de engenharia	2.77	[CLS] serviços de engenharia cons ##istente no tapa buraco de vias urbana ##s [SEP]

Legend: <span style="color: red;">■</span> Negative <span style="color: gray;">□</span> Neutral <span style="color: green;">■</span> Positive				
True Label	Predicted Label	Attribution Label	Attribution Score	Word Importance
0	Obras de engenharia (1.00)	Obras de engenharia	2.77	[CLS] serviços de engenharia cons ##istente no tapa buraco de vias urbana ##s [SEP]

Fig. 3. Explainability using NER (Text in Portuguese)

Legend: <span style="color: red;">■</span> Negative <span style="color: gray;">□</span> Neutral <span style="color: green;">■</span> Positive				
True Label	Predicted Label	Attribution Label	Attribution Score	Word Importance
3	Serviços (0.89)	Serviços	1.32	[CLS] a presente dispens ##a de licitação elet ##rônica tem por finalidade a contrataç ##ão de empresa especializada nos serviços de engenharia cons ##istente no tapa buraco de vias urbana ##s . o obje ##to abrang ##erá as especific ##idades conforme descri ##tas no anexo i - projeto básico . [SEP]

Legend: <span style="color: red;">■</span> Negative <span style="color: gray;">□</span> Neutral <span style="color: green;">■</span> Positive				
True Label	Predicted Label	Attribution Label	Attribution Score	Word Importance
3	Serviços (0.89)	Serviços	1.32	[CLS] a presente dispens ##a de licitação elet ##rônica tem por finalidade a contrataç ##ão de empresa especializada nos serviços de engenharia cons ##istente no tapa buraco de vias urbana ##s . o obje ##to abrang ##erá as especific ##idades conforme descri ##tas no anexo i - projeto básico . [SEP]

Fig. 4. Explainability using full text (Text in Portuguese)

- [13] Kayte, S. & Schneider-Kamp, P. A Mixed Neural Network and Support Vector Machine Model for Tender Creation in the European Union TED Database. *International Joint Conference On Knowledge Discovery, Knowledge Engineering And Knowledge Management*. (2019)
- [14] Navas-Loro, M., Garijo, D. & Corcho, O. Multi-label Text Classification for Public Procurement in Spanish. *Procesamiento Del Lenguaje Natural*. **69** pp. 73-82 (2022)
- [15] Gomes, W. & Colaço, M. Applications of Artificial Intelligence for Auditing and Classification of Incongruent Descriptions in Public Procurement. *Proceedings Of The XVIII Brazilian Symposium On Information Systems*. (2022)
- [16] Song, M., Feng, Y. & Jing, L. A Survey on Recent Advances in Keyphrase Extraction from Pre-trained Language Models. *Findings Of The Association For Computational Linguistics: EACL 2023*. pp. 2153-2164 (2023,5)
- [17] Mu, F., Yu, Z., Wang, L., Wang, Y., Yin, Q., Sun, Y., Liu, L., Ma, T., Tang, J. & Zhou, X. Keyphrase Extraction with Span-based Feature Representations. *ArXiv*. **abs/2002.05407** (2020)
- [18] Shi, T., Keneshloo, Y., Ramakrishnan, N. & Reddy, C. Neural Abstractive Text Summarization with Sequence-to-Sequence Models. *ACM/IMS Trans. Data Sci.* **2** (2021,1), <https://doi.org/10.1145/3419106>
- [19] Gero, Z. & Ho, J. NamedKeys: Unsupervised Keyphrase Extraction for Biomedical Documents. *Proceedings Of The 10th ACM International Conference On Bioinformatics, Computational Biology And Health Informatics*. pp. 328-337 (2019), <https://doi.org/10.1145/3307339.3342147>
- [20] AXA-Group Parsr: Transforms PDF, Documents and Images into Enriched Structured Data. (<https://github.com/axa-group/Parsr>,2024), Accessed: 2024-01-10
- [21] S.L.U., A. Open-source data curation platform for LLMs. (<https://argilla.io/>,2024), Accessed: 2024-01-08
- [22] Jégou, H., Douze, M. & Schmid, C. Product Quantization for Nearest Neighbor Search. *IEEE Transactions On Pattern Analysis And Machine Intelligence*. **33**, 117-128 (2011)
- [23] Sivic & Zisserman Video Google: a text retrieval approach to object matching in videos. *Proceedings Ninth IEEE International Conference On Computer Vision*. pp. 1470-1477 vol.2 (2003)
- [24] Johnson, J., Douze, M. & Jégou, H. Billion-Scale Similarity Search with GPUs. *IEEE Transactions On Big Data*. **7**, 535-547 (2021)
- [25] Souza, F., Nogueira, R. & Lotufo, R. BERTimbau: Pretrained BERT Models for Brazilian Portuguese. *Intelligent Systems*. pp. 403-417 (2020)
- [26] Viegas, C., Costa, B. & Ishii, R. JurisBERT: A New Approach that Converts a Classification Corpus into an STS One. *Computational Science And Its Applications – ICCSA 2023*. pp. 349-365 (2023)
- [27] Devlin, J., Chang, M., Lee, K. & Toutanova, K. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *Proceedings Of The 2019 Conference Of The North American Chapter Of The Association For Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long And Short Papers)*. pp. 4171-4186 (2019)
- [28] Zhang\*, T., Kishore\*, V., Wu\*, F., Weinberger, K. & Artzi, Y. BERTScore: Evaluating Text Generation with BERT. *International Conference On Learning Representations*. (2020)
- [29] Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L. & Stoyanov, V. RoBERTa: A Robustly Optimized BERT Pretraining Approach. (2019)
- [30] Pierse, C. Transformers Interpret: Model explainability that works

seamlessly with transformers. (<https://github.com/cdpierse/transformers-interpret>, 2024), Accessed: 2024-01-23