

Automatic Generation of Introduction and Abstract for Research Papers - GPT Neo as a summarizer

219354V - R.P.D. Kumarasinghe

Supervisor: Dr. Nisansa de Silva

Overview

1. Introduction
2. Research Problem
3. Research Objectives
4. Methodology
5. Results
6. References



Introduction

Introduction

- The abstract of a research paper provides a quick summary of the entire paper from problem to solution to the result
- The Introduction section provides a primer to the rest of the paper by summarising the goals and the setting of the research while expanding on the basis established by the abstract

Abstract

Paper Format for the Proceedings of the 2011 IEEE International Conference on Computational Intelligence and Computing Research

N.K. Surname¹, P. M. Surname², A.S.A. Surname²
¹Centre for Information Technology and Engineering, Manonmaniam Sundaranar University, City, Country
²Department of Communication Technology, University of College, City, Country
(e-mail address)

Abstract - These instructions give you basic guidelines for preparing papers for the ICIC2011 Proceedings. Papers up to 5 pages must be submitted using this format. This document is a template for Microsoft Word. If you are reading a paper version of this document, please download the electronic file from the Conference website so you can use it to prepare your manuscript. Abstract should not exceed 150 words. To allow retrieval by CD-ROM software, please include appropriate key words in your abstract, in alphabetical order, separated by commas.

Keywords - Fonts, formatting, margins

I. INTRODUCTION

Your goal is to simulate, as closely as possible, the usual appearance of typeset papers in the *IEEE Transactions*. One difference is that the authors' affiliations should appear immediately following their names – do not include your title there. For items not addressed in these instructions, please refer to a recent issue of an *IEEE Transactions*.

II. METHODOLOGY

All papers must be submitted electronically in pdf format. Prepare your paper using a A4 page size of 210 mm × 297 mm (8.27" × 11.69").

1) **Type sizes and typefaces:** The best results will be obtained if your computer word processor has several type sizes. Try to follow the type sizes specified in Table I as best as you can. Use 14 point bold, capital letters for the title, 12 point Roman (normal) characters for author names and 10 point Roman characters for the main text and author's affiliations.

2) **Format:** In formatting your page, set top margin to 25 mm (1") and bottom margin to 31 mm (1 1/4"). Left and right margins should be 19 mm (3/4"). Use a two-column format where each column is 83 mm (3 1/4") wide and spacing of 6 mm (1/4") between columns. Indent paragraphs by 6 mm (1/4").

Left and right-justify your columns. Use tables and figures to adjust column length. Use automatic hyphenation and check spelling. All figures, tables, and equations must be included *in-line* with the text. Do not use links to external files.

III. RESULTS

A. Figures and Tables

Graphics should be in TIFF, 600 dpi (1 bit/sample) for line art (graphics, charts, drawings or tables) and 220 dpi for photos and gray scale images.

Position figures and tables at the tops and bottoms of columns. Avoid placing them in the middle of columns. Large figures and tables may span across both columns. Figure captions should be below the figures; table names and table captions should be above the tables. Use the abbreviation "Fig." even at the beginning of a sentence.

Figure axis labels are often a source of confusion. Try to use words rather than symbols. As an example, write the quantity "Magnetization," or "Magnetization *M*," not just "*M*." Put units in parentheses. Do not label axes only with units. As in Fig. 1, for example, write "Magnetization (A/m)" or "Magnetization (A · m⁻¹)," not just "A/m." Do not label axes with a ratio of quantities and units. For example, write "Temperature (K)," not "Temperature/K."

Multipliers can be especially confusing. Write "Magnetization (kA/m)" or "Magnetization (10³ A/m)." Do not write "Magnetization (A/m) × 1000" because the reader would not know whether the top axis label in Fig. 1 meant 16000 A/m or 0.016 A/m. Figure labels should be legible, approximately 10-point type.

TABLE I
TYPE SIZES FOR CAMERA-READY PAPERS

Type Size (pts)	Appearance		
	Regular	Bold	Italic
7	Table captions*		
8	Section titles, tables, table names*, first letters in table captions*, table superscripts, figure captions, text subscripts and superscripts, references, footnotes		
9		Abstract	
10	Authors' affiliations, main text, equations, first letter in section titles*, first letter in table names*		Subheading
12	Authors' names		
14		Paper title	

* Capital letters

Research Problem

Research Problem

- Abstract and Introduction are expected to be concise and informative.
- But generating them manually is difficult and time consuming.
- Summarization has domain specific training approaches which perform well.



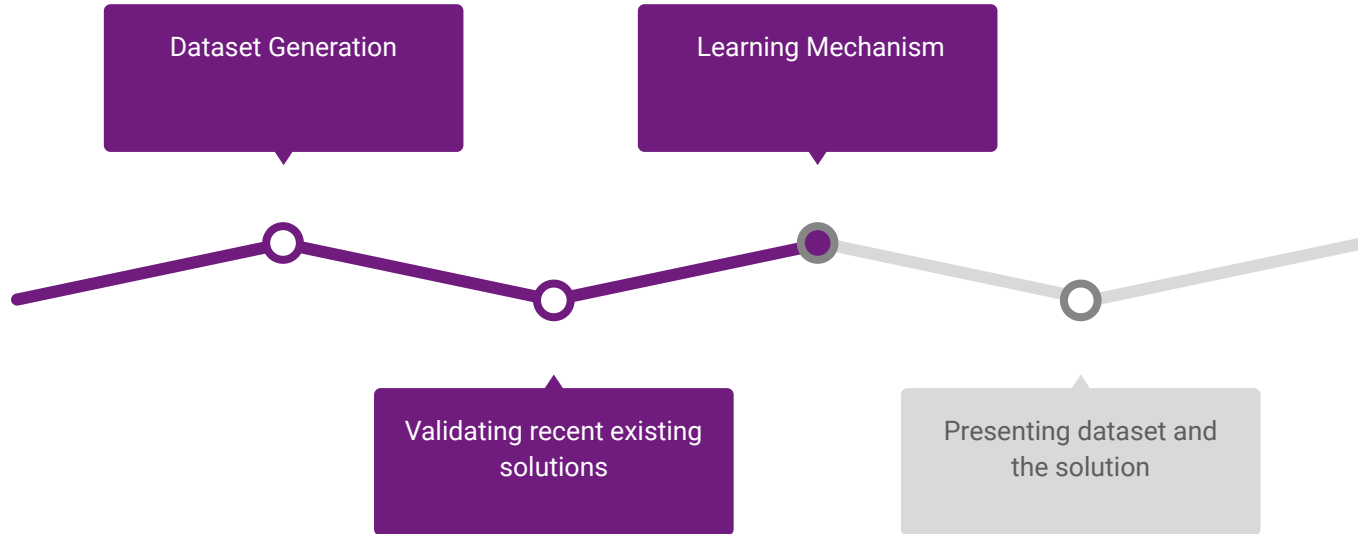
Research Objectives

Research Objectives

1. Creating a sufficient data set for the task of Abstract and Introduction generation in the computational linguistic domain
2. Evaluating existing state-of-the art solutions of text summarization technologies on the above data set and other comparable data sets.
3. Creating automatic summarization models capable of Abstract And Introduction generation in the computational linguistic domain.
4. Creating an online application which, when given the LATEX source sans the Abstract and Introduction, generates these sections automatically.

Methodology

Methodology: Backlog



Methodology: GPT as A Summarizer

- GPT models need start text
- It will predict the rest

Text> Lorem Ipsum is simply dummy text of the printing and typesetting industry. Lorem Ipsum has been the industry's standard dummy text ever since the 1500s, when an unknown printer took a galley of type and scrambled it to make a type specimen book. It has survived not only five centuries, but also the leap into electronic typesetting, remaining essentially unchanged.

Abstract> It was popularised in the 1960s with the release of Letraset sheets containing Lorem Ipsum passages, and more recently with desktop publishing software like Aldus PageMaker including versions of Lorem Ipsum.

Methodology: GPT as A Summarizer

- GPT models need start text
- It will predict the rest

Text>Contrary to popular belief, Lorem Ipsum is not simply random text. It has roots in a piece of classical Latin literature from 45 BC, making it over 2000 years old. Richard McClintock, a Latin professor at Hampden-Sydney College in Virginia, looked up one of the more obscure Latin words, consectetur, from a Lorem Ipsum passage, and going through the cities of the word in classical literature, discovered the undoubtable source.

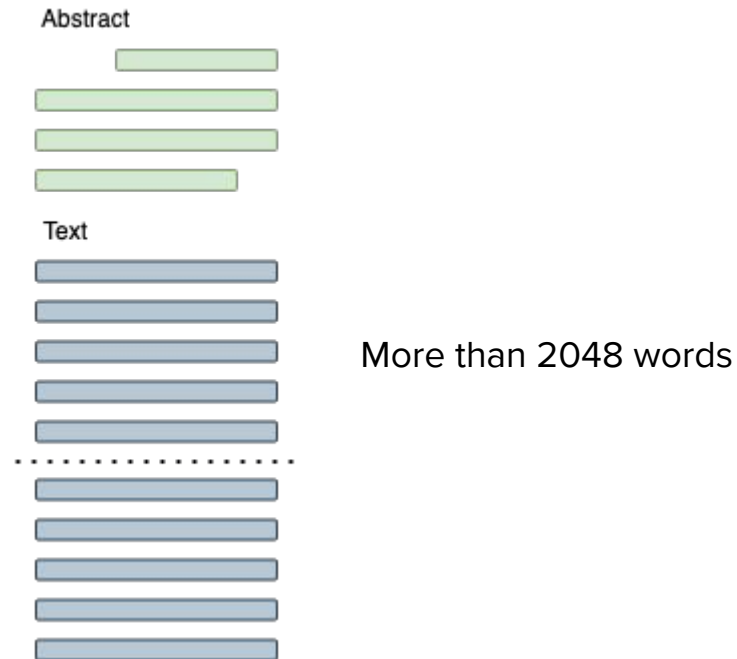
Abstract>

Abstract: In this paper, we introduce a novel architecture for zero-shot conversational question answering (CQA) in corporate settings. Previous works usually employ a three-layer multi-grained model structure to address CQA tasks in dialog systems. This results in an expensive computational cost when training on large-scale datasets. In modern organizations, these expensive training steps are not prohibitively expensive. For example, the Dialog System Trustworthiness Report (\$4.8K training steps) suggests that training on a dialogue corpus with \$500 hours of conversation data (\$750 hours of human-annotated training data) results in about 85% accuracy when the corresponding crowd-sourced dialog corpus is used (steps 4 to 7). In contrast, a small-scale corpus reduces the training steps to \$20 hours (\$10 hours for crowd-

<Abstract:

Methodology: Challenge

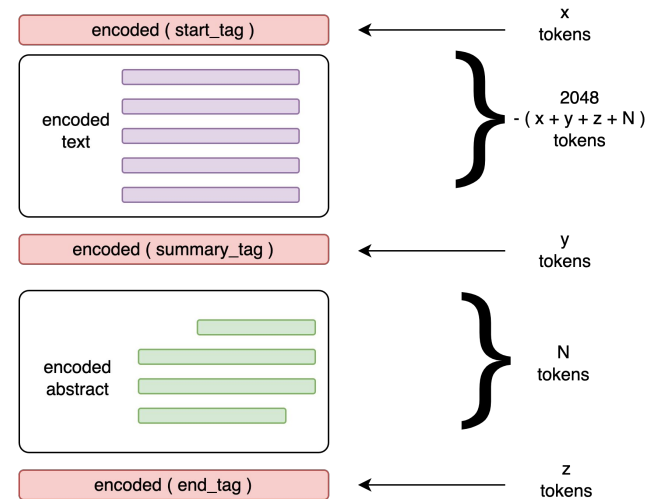
- Pre trained GPT models has 2048 token limit



Methodology: Challenge

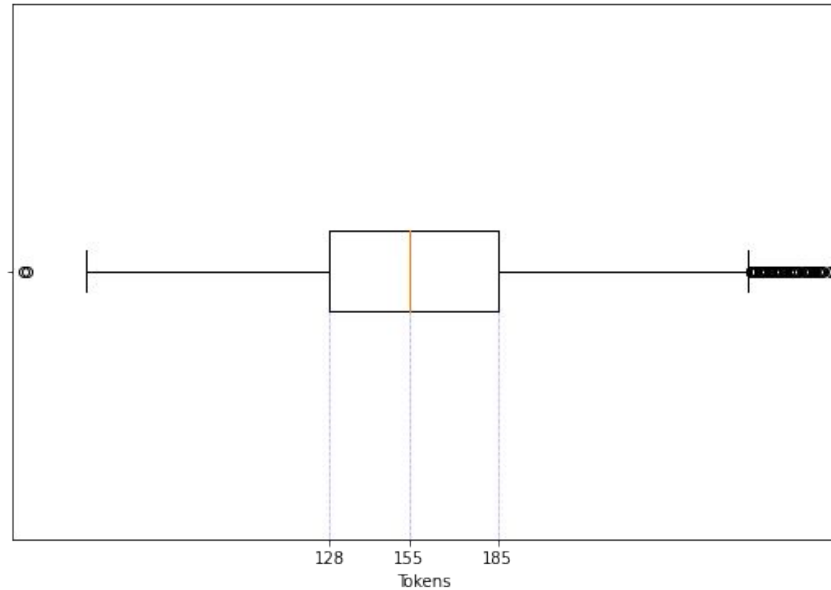
- We should have a fixed size for text length

$$2048 - (x + y + z + N)$$



Methodology: Abstract token size

- We choose the 3rd quartile which is 185 as the N

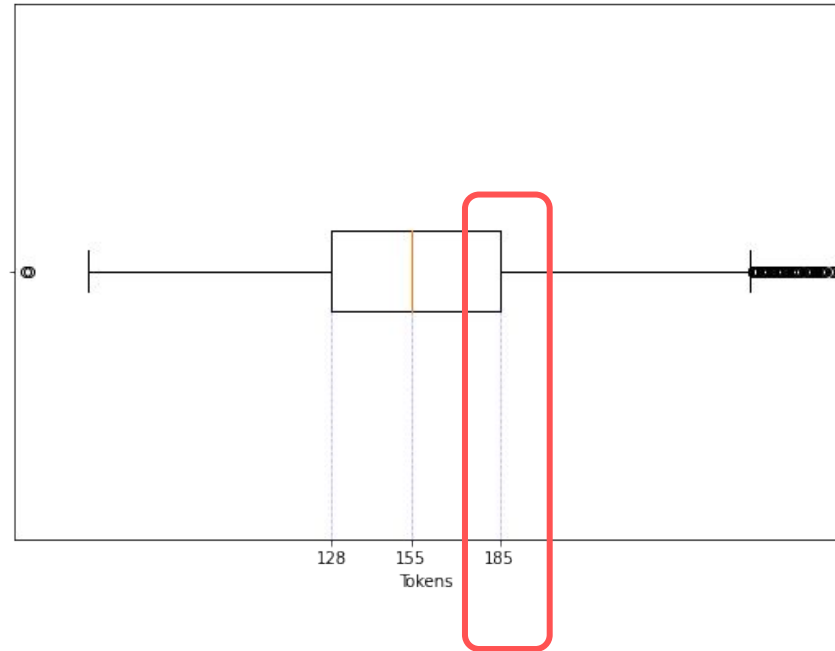


Methodology: Abstract token size

- We choose the 3rd quartile which is 185 as the N

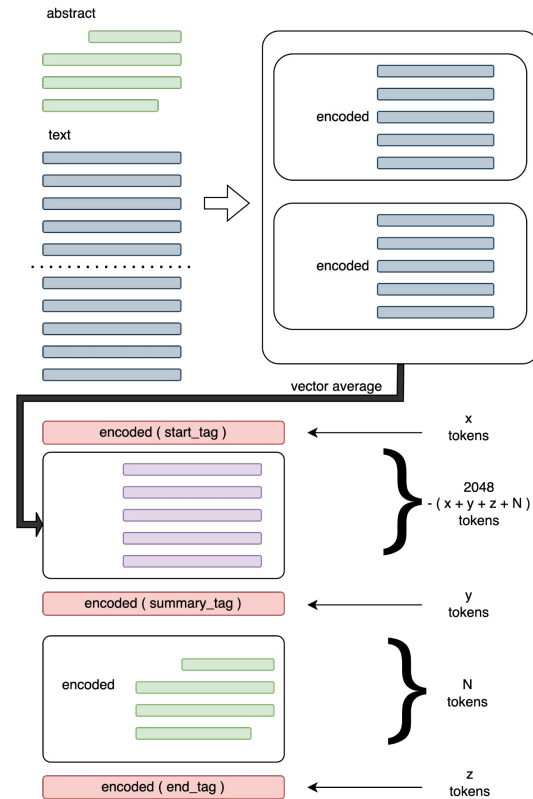
Abstract size = 185

Body size= $2048 - (x + y + z + 185)$



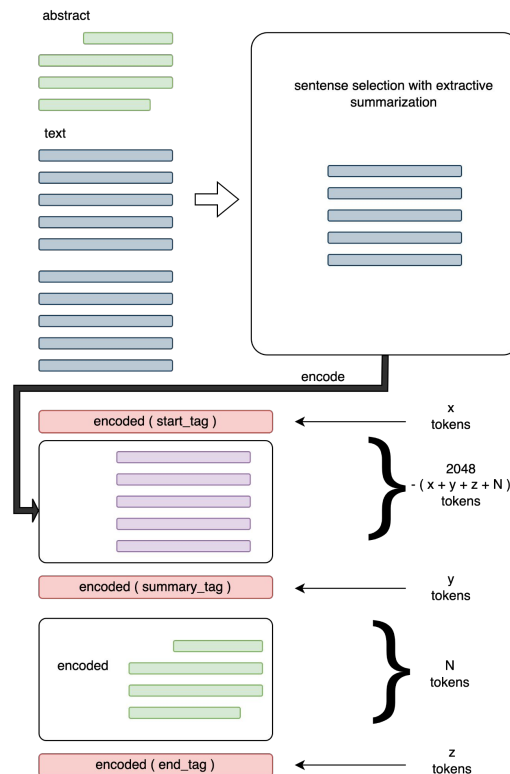
Methodology: Reduce Paper Text > Vector Average

- Divide paper text to chunks
- Encode the chunks
- Take the average of encoded chunk



Methodology: Reduce Paper Text > Extractive Method

- Paper text is pre-summarized using extractive techniques
- 4 techniques were used
 - Text Rank [1]
 - Lex Rank [2]
 - LSA [3]
 - Luhn [4]



[1] Erkan, Günes, and Dragomir R. Radev. "Lexrank: Graph-based lexical centrality as salience in text summarization." Journal of artificial intelligence research 22 (2004): 457-479.

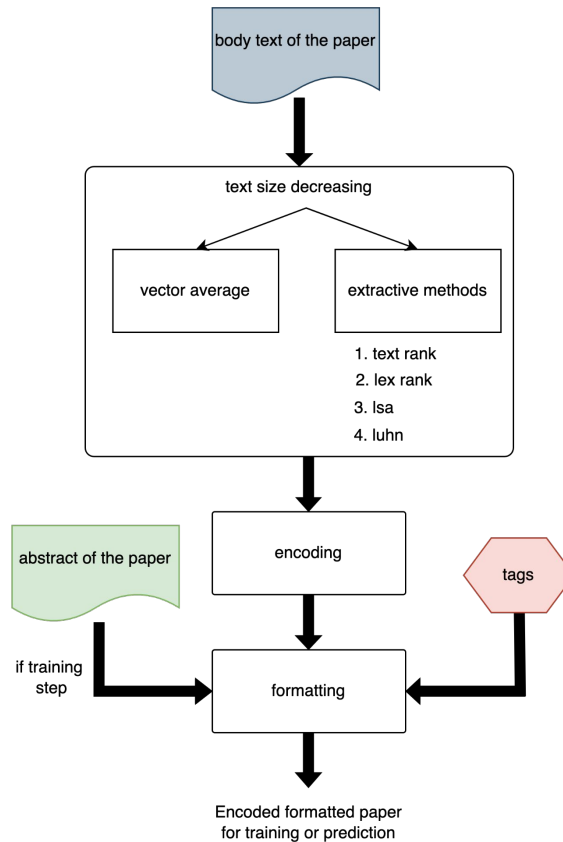
[2] Mihalcea, Rada, and Paul Tarau. "Textrank: Bringing order into text." In Proceedings of the 2004 conference on empirical methods in natural language processing, pp. 404-411. 2004.

[3] Landauer, Thomas K., Peter W. Foltz, and Darrell Laham. "An introduction to latent semantic analysis." Discourse processes 25, no. 2-3 (1998): 259-284.

[4] Luhn, Hans Peter. "The automatic creation of literature abstracts." IBM Journal of research and development 2, no. 2 (1958): 159-165.

Methodology: Reduce Paper Text

- Text size decreasing to satisfy 2048 tokens using
 - Vector average
 - Extractive methods



Results

Comparison sample 1

Actual Abstract

"The chit-chat-based conversational recommendation systems (CRS) provide item recommendations to users through natural language interactions. To better understand user's intentions, external knowledge graphs (KG) have been introduced into chit-chat-based CRS. However, existing chit-chat-based CRS usually generate repetitive item recommendations, and they cannot properly infuse knowledge from KG into CRS to generate informative responses. \nTo remedy these issues, we first reformulate the conversational recommendation task to highlight that the recommended items should be new and possibly interested by users. Then, we propose the Knowledge-Enriched Conversational Recommendation System (KECRS). Specifically, we develop the Bag-of-Entity (BOE) loss and the infusion loss to better integrate KG with CRS for generating more diverse and informative responses. BOE loss provides an additional supervision signal to guide CRS to learn from both human-written utterances and KG. Infusion loss bridges the gap between the word embeddings and entity embeddings by minimizing distances of the same words in these two embeddings. Moreover, we facilitate our study by constructing a high-quality KG, \ie The Movie Domain Knowledge Graph (TMDKG). Experimental results on a large-scale dataset demonstrate that KECRS outperforms state-of-the-art chit-chat-based CRS, in terms of both recommendation accuracy and response generation quality."

Comparison sample 1

Vector Average

"One of the most important tasks in Natural Language Processing is the prediction of the future state of a system. \&Added by \new{Bengio Biovit{1} }; 1985) is a use case where one immediately sees the interest in event history prediction. There has been effort to combine end-to-end approaches with various external knowledge sources, typically SQL or other relational databases. Large systems are inherently containers for external data, and such an abstraction helps the model to be parameter-efficient and also to be able to continuously adapt to the larger context. However, this paper focuses on a more specialized type of event, that of multi-turn communication between humans. Inspired by the symbolic nature of human conversations, we combine external language data with an approach that tackles the problem in the non-autoregressive fashion, exploiting pre-computed hidden representations from Transformers. We show that"

Comparison sample 1

Lex Rank

"Multimodal propaganda is a problem frequently encountered in online news. Due to the dynamic nature of the Internet, advanced techniques become necessary to handle the challenges that arise from the rapidly changing trend of communication in the online world. One such challenge is propaganda detection. However, the disseminated information can be deeply affected by biases. Further, knowledge is usually embodied in multi-modal information such as videos and text. This paper adopts recent results of multimodal propaganda techniques detection and propose a novel multimodal propaganda technique detection model by leveraging pretrained models of visual and linguistic information. In addition to textual and visual information, we propose two advanced augmentation techniques -- Back-Translation and SCIKT 2.0. The proposed model is further evaluated on the test set, and the results show an"

Comparison sample 1

Text Rank

"While the task of propaganda detection is a difficult one for computer vision models, it is supported by diverse techniques that are inexpensive to pretrain specifically in a low-resource setting and require minimal understanding of the methods. The most dominant form of propaganda detection is opinion-based and opinions are stated as statements of opinions to generate influence by couching the claims in a certain way so as to generate opinion dynamics for people. Also, propaganda techniques are often described in phrases or phrases that can be understood by the general public, so their detection can help the public to form an opinion on the issue. In this paper, we present a multimodal propaganda detection model that uses textual and image feature extraction along with LIIR_R linguistics and visual features for propaganda detection. In addition, we propose an in-class reaction model to identify the persuasion techniques behind propaganda from text and images"

Comparison sample 1

LSA

"Considering the ever rising amounts of information available on the internet, the identification of automated propaganda on Twitter is crucial, because it results in the deployment of a unified approach to combat both speeches and potentially disinformation.\n\nThe shared task on early detection of propaganda techniques, represented by visual and textual features, shares the same goal as ours: to make automatic detection of propaganda techniques in multimodal inputs feasible. We present our submission for the shared task at \acrs2021.\n\nOur proposed method is based on a deep neural network followed by a classification approach, leveraging both textual and visual features.\n\nWe fine-tuned a pretrained model with additional improvements on the feature extraction and classification approaches and performed experiments on the test sets of the shared task. The submission achieves a score of 0.16, which is comparable to other published submissions."

Comparison sample 1

Luhn

"Automated techniques to detect and identify propaganda-inspired techniques in multimodal text are being proposed and investigated by researchers. We propose a novel system for the multimodal propaganda technique detection task at the CELSA shared task for 2019.\nWe propose a novel data augmentation and feature extraction approach to boost the performance of the proposed system. We anchor our proposed system on a pretrained visuolinguistic representation and a pretrained language model.\nWe conduct a thorough study to understand the impact of data augmentation and feature extraction method used to improve system performance.\nOur system reaches an average of 63.02 and IQR of 0.79 on the validation set of the test sets, which makes it the best among the competitors. Further results can be found in [\url{http://royalsociety.nl/~miszewia20/tasks/}](http://royalsociety.nl/~miszewia20/tasks/)"

Results: Results



Average of 330 predictions

Pre-Sum Method	ROUGE-1	ROUGE-2	ROUGE-L
Vector Average	0.1843	0.0204	0.1698
Lex Rank	0.2612	0.0478	0.2359
Text Rank	0.2548	0.0441	0.2304
LSA	0.2629	0.0472	0.2382
Luhn	0.2602	0.0483	0.2343

References

References

- [1] Erkan, Günes, and Dragomir R. Radev. "Lexrank: Graph-based lexical centrality as salience in text summarization." *Journal of artificial intelligence research* 22 (2004): 457-479.
- [2] Mihalcea, Rada, and Paul Tarau. "Textrank: Bringing order into text." In *Proceedings of the 2004 conference on empirical methods in natural language processing*, pp. 404-411. 2004.
- [3] Landauer, Thomas K., Peter W. Foltz, and Darrell Laham. "An introduction to latent semantic analysis." *Discourse processes* 25, no. 2-3 (1998): 259-284.
- [4] Luhn, Hans Peter. "The automatic creation of literature abstracts." *IBM Journal of research and development* 2, no. 2 (1958): 159-165.

Thank You