

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19

Bridging the Gap in Health Literacy: Harnessing the Power of Large Language Models to Generate Plain Language Summaries from Biomedical Texts

Carolina Salazar-Lara^{1,2,3*}, Andrés Felipe Arias Russi^{1,4}, Rubén Manrique^{1,4}

¹School of Engineering, Universidad de Los Andes, Bogotá, Colombia

²Department of Biomedical Engineering, Universidad de Los Andes, Bogotá, Colombia

³Department of Industrial Engineering, Analytics, Universidad de Los Andes, Bogotá, Colombia

⁴Department of Computer and Systems Engineering, Universidad de Los Andes, Bogotá, Colombia

*Corresponding author

E-mail: c.salazar499@uniandes.edu.co (CS)

20 **Abstract**

21 Health literacy is essential for individuals to navigate the healthcare system and make informed decisions about their health.
22 Low health literacy levels have been associated with negative health outcomes, particularly among older populations and those
23 financially restricted or with lower educational attainment. Plain language summaries (PLS) are an effective tool to bridge the
24 gap in health literacy by simplifying content found in biomedical and clinical documents, in turn, allowing the general audience
25 to truly understand health-related documentation. However, translating biomedical texts to PLS is time-consuming and
26 challenging, for which they are rarely accessible by those who need them. We assessed the performance of Natural Language
27 Processing (NLP) for systematizing plain language identification and Large Language Models (LLMs), Generative Pre-trained
28 Transformer (GPT) 3.5 and GPT 4, for automating PLS generation from biomedical texts. The classification model achieved high
29 precision (97·2%) in identifying if a text is written in plain language. GPT 4, a state-of-the-art LLM, successfully generated PLS
30 that were semantically equivalent to those generated by domain experts and which were rated high in accuracy, readability,
31 completeness, and usefulness. Our findings demonstrate the value of using LLMs and NLP to translate biomedical texts into
32 plain language summaries, and their potential to be used as a supporting tool for healthcare stakeholders to empower patients
33 and the general audience to understand healthcare information and make informed healthcare decisions.

34 **Keywords**

35 *Health literacy, Plain Language, Large Language Models, Analytics, Natural Language Processing, Natural Language*
36 *Generation, Biomedicine*

37 **Introduction**

38 Health literacy refers to an individual's capacity to access, understand, and use health information [1]. It empowers patients and
39 their families to navigate healthcare systems, comprehend and act upon a diagnosis or medical instruction, adhere to medication
40 regimens, and make informed decisions, otherwise considered daunting, regarding participation in clinical trials, treatment
41 options, or medical procedures [2-4]. Low health literacy levels have been consistently associated with higher mortality rates,
42 increased instances of preventable hospitalizations, and poor treatment adherence [3]. Paradoxically, while health literacy is
43 crucial for positive health outcomes, the 2015 European Health Literacy Survey revealed that almost half of the respondents
44 had inadequate health literacy, particularly among older populations, those who are financially restricted, or who have lower
45 educational attainment [5-6].

46 With the growing expectation for individuals to participate in healthcare decisions, enhancing health literacy becomes a
47 significant attribute in improving public health and reducing health disparities [1, 7-8].

48 Improving health literacy in the population extend beyond actions taken to increase individual health literacy levels. In line with
49 the General Data Protection Regulation (GDPR) principle of transparency, stakeholders such as healthcare providers,
50 policymakers, and pharmaceutical companies should strategize to improve their organizational health literacy (OHL) by ensuring
51 the clarity and comprehensibility of health documentation [9-10].

52 One strategy to do so is by simplifying clinical and scientific research language into lay-friendly summaries, known as plain
53 language summaries (PLS).

54 There are different techniques and guidelines that can be used to translate complex scientific and biomedical concepts into PLS,
55 for example, eliminating the use of technical jargon, replacing passive voice by active, or using short sentences and paragraphs
56 [6, 11]. However, authoring a PLS can be time consuming and challenging, particularly in areas like clinical settings which
57 typically involve documents with technical and domain-specific vocabulary.

58 With the advancement of technology, new methods have been developed to automate the simplification of biomedical texts. In
59 2022, a review by Oldov et al. analyzed 32 tools or methods using either rule-based approach or Natural Language Processing
60 (NLP) and concluded that NLP methods offer more promising outputs but were limited by scarcity of training data, resulting in
61 continued reliance on rule-based methods [12]. Large Language Models (LLMs) with their immense data training potential and
62 text generation capabilities, present a promising solution to tackle this challenge and automate the generation of PLS from
63 technical documents.

64 With the objective of bridging the gap in health literacy by facilitating the translation of biomedical texts to comprehensible
65 summaries designed for patients, our study demonstrates the potential of NLP to develop a classification system to identify if a
66 text is written in plain language, and LLMs to automate the generation of accurate, complete, and comprehensible PLS.

67 **Materials and Methods**

68 Our methodology, outlined in Figure 1, consisted of 3 main steps: collecting and processing of sample texts in technical and
69 plain language, conducting a quantitative analysis of the texts to generate a plain language classification model and a qualitative
70 analysis to generate the prompts for the LLMs, and using the LLMs to generate PLS and test them.

71 **Data Collection and Processing**

72 We collected biomedical texts, both in technical and plain language (see the data sources in *Supplementary Table 1*) and
73 assembled them into a dataset of 14,441 texts. This “*main dataset*” was then divided into training and testing sets, consisting of
74 4,596 plain and 6,721 technical texts for training, and 1,149 plain and 1,975 technical texts for testing.

75 We enlarged each dataset by treating each paragraph of a minimum of 250 words as a distinct unit, while excluding texts with
76 fewer than 250 words. As a result, our "*augmented dataset*" had 61,354 texts, divided into 16,731 plain and 31,740 technical for
77 training, and 5,090 plain and 7,793 technical for testing.

78 **Analysis of Plain Language**

79 We conducted qualitative and quantitative analysis of the texts to identify unique linguistic traits and variables that classify a text
80 as plain language.

81 **Qualitative Analysis**

82 Driven by the varying and broad-scope guidance on creating high-quality PLS [13], we analyzed a subset of our plain texts and
83 created a 'criteria checklist' (see *Supplementary Table 2*) with the linguistic attributes most commonly present in plain texts. Key
84 resources used in this process were guides and reviews, such as: Your Guide to CLEAR WRITING by CDC [11], Federal Plain
85 Language Guidelines [14], Health Literacy Universal Precautions Toolkit by Agency for Healthcare Research and Quality
86 (AHRQ) [15], Just Plain Clear Glossary by United Health Group [16], EU 536/2014 Summary of Clinical Results for Laypersons
87 [17], and results presented by Stoll et al, in their systematic review of theory, guidelines, and empirical research on PLS [13].
88 We used the resultant checklist to complement the qualitative findings described in the next section and aid in developing the
89 prompt detailed in the section LLM Prompt for Plain Language Summary Generation.

1. Data Collection and Processing

Technical Documents	Plain Language Documents
U.S National Library of Medicine (NIH) ClinicalTrials.gov	Pfizer Results Plain Language Summaries Trial Summaries by Citeline Regulatory
Cochrane Library by Wiley	

2. Analysis of Plain Language

- Quantitative analysis** of readability metrics and language variables to create **Plain Text Classification Model**.
- Qualitative analysis** of subset of documents to create a **criteria checklist of plain texts linguistic attributes and LLM prompt for PLS generation**.

Model Prompt Anatomy

Context or Rationale	<i>Using the following clinical trial protocol text as input, generate a Protocol Plain Language Summary (PPLS) compliant with GDPR and understandable by any patient, regardless of their health literacy.</i>	Technical Version	An echocardiogram will be used to diagnose cardiac failure.
Output: Structure and Format	<i>The generated PPLS should be presented in a logical order, using the following headings: [1, 2, 3, ..., n]</i>	Plain Language Version	An echocardiogram (a non-invasive test that uses sound waves to create images of the heart) will be used to identify if your heart is not working right.
Content: Examples and rules	<i>In the section [] : Answer "What are the goals of the study?" Specify the main and secondary objectives of the trial and how they will be measured (eg. the main trial endpoint is the percent change in the number....).</i>		
Restrictions and Considerations	<i>The generated PPLS must follow these plain language guidelines: [...] The AI model should not invent information or add content that is not present in the input protocol.</i>		

**This is an excerpt, the complete prompt can be found in supplementary tables*

3. Testing LLMs

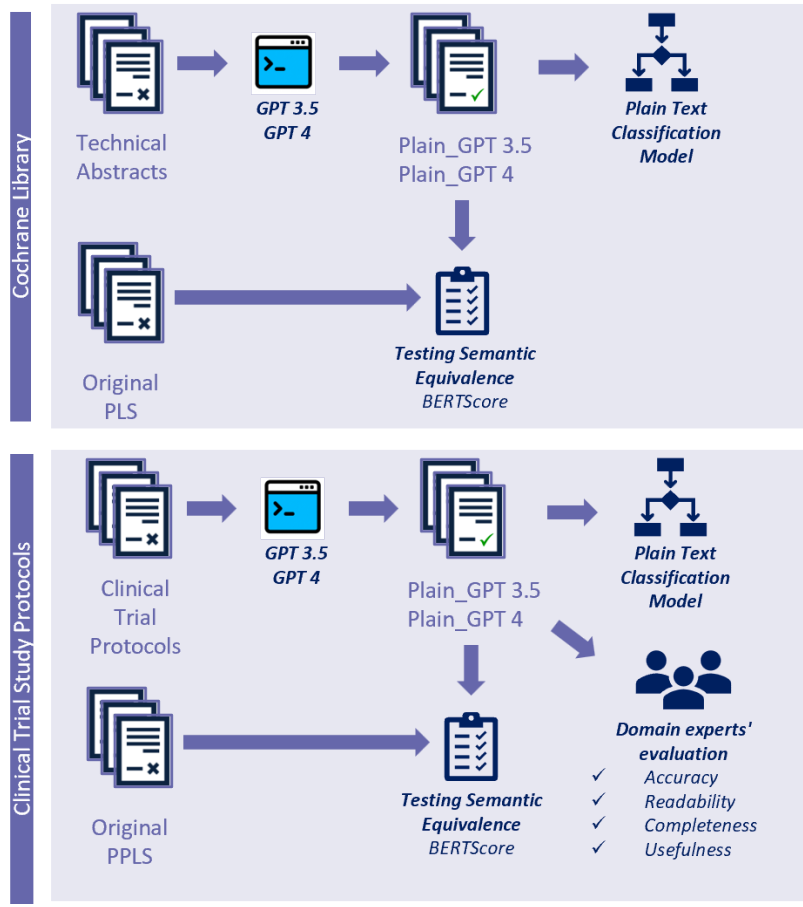


Figure 1. Methodology. Our methodology involved three steps: 1) collection and processing of biomedical texts (technical documents and plain language documents) into datasets for training and testing, 2) quantitative analysis of the texts to create a plain language classification model, and qualitative analysis to identify linguistic traits in plain texts to guide the engineering of a prompt that could translate biomedical text into Plain Language Summaries (PLS) using Language Learning Models (LLMs; and 3) testing the effectiveness of the LLMs in generating PLS quantitatively with our classification model and with semantic equivalence (BERTScore) and qualitatively with domain experts' evaluation.

90 Quantitative Analysis

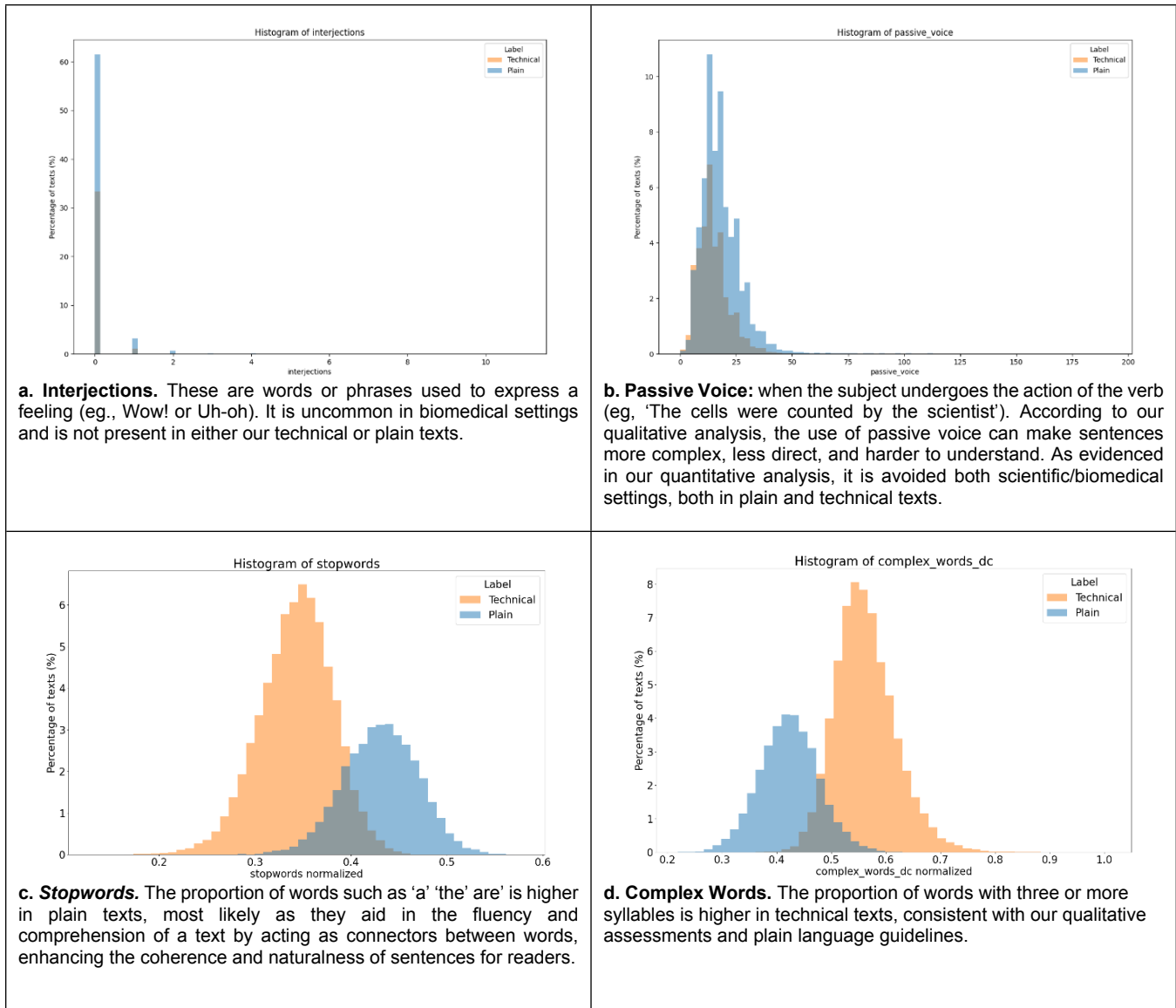
91 We computed readability metrics and language variables for each text in the augmented dataset using the Readability library
92 [18] and SpaCy [19], respectively. This resulted in 64 variables presenting each text's readability and linguistic traits (see
93 *Supplementary Table 3*).

94 We analyzed the language variables in our dataset to identify their potential to classify a text as technical or plain. We used
95 statistical hypothesis test for each of the variables of the *main dataset*. For each variable, we created a random sample of size
96 n from the plain texts ($X_1, X_2 \dots X_n \sim P_X$) and a random sample of size n from the technical texts ($Y_1, Y_2 \dots Y_n \sim Q_Y$), and tested if our
97 data supported either of the following hypotheses:

- 98 • *Null Hypothesis*, $H_0: P = Q$, the distributions of the proportion of the variable of interest for both samples (text and
99 technical) are the same.
- 100 • *Alternative Hypothesis*, $H_0: P \neq Q$, the distributions of the proportion of the variable of interest for both samples (text
101 and technical) are different.

102 We evaluated the null hypothesis by comparing our 2 distributions using non-parametric tests: Wilcoxon, Kolmogorov-Smirnov
103 (KS), and Mann–Whitney U. Given the multiple hypothesis tests, one for each variable, we adjusted the significance levels to
104 control the probability of Type I errors by using the Bonferroni correction to lower the alpha value by dividing the desired
105 significance level $\alpha = 0.05$ by the total number of tests $m = 64^{18}$.

106 Figure 2 illustrate examples of the comparison of the distributions of some of the variables in technical and plain texts. Out of
107 the 64 variables, only 'Interjections' and 'Passive Voice' did not provide sufficient evidence to reject the null hypothesis (p -value
108 > 0.0008). The other 62 variables were significantly distinct between the types of text and were included in our classification
109 model.



110 **Figure 2.** Comparison of the distribution of a sample of readability metrics or language variables between plain and technical texts.

111 Plain Texts Classification Model

112 We used the *augmented dataset - train* and the 62 distinct variables between text types (Section Quantitative Analysis), to build
113 the classification model. We used Gradient Boosting (GB) and Random Forest (RF) machine learning models.

114 LLM Prompt for Plain Language Summary Generation

115 Our objective was to design a prompt for LLMs capable of translating biomedical technical documents into PLS.

116 Beginning with a clinical trial protocol from ClinicalTrials.Gov (see data sources in *Supplementary Table 1*), we used a simple
117 initial prompt: 'Using the following clinical trial protocol text as input, create a plain language summary'. We tested this prompt

118 using both GPT3.5 and GPT4, analyzed the generated output, and iteratively refined the prompt by adding details and
119 instructions.

120 We aimed to produce a PLS that met the following qualitative criteria: (1) **Accuracy**: the content is clinically and scientifically
121 accurate. (2) **Readability**: the content is comprehensible by a lay person, as defined by the plain language criteria checklist
122 (*Supplementary Error! Reference source not found.*). (3) **Completeness**: the content adheres with the expectations of a
123 Protocol Plain Language Summary (PPLS) as specified by EU CTR No 536/2014 [16]. (4) **Usefulness**: the generated PLS can
124 be used as a first version to draft the study PPLS.

125 Our final prompt, provided in *Supplementary Table 4*, was designed specifically to generate a PLS of a clinical trial protocol. It
126 includes the following elements:

- 127 • **Context**: a clear rationale on why a PLS is needed for the given clinical trial protocol.
- 128 • **Output**: the desired structure and format for the generated summary, including the specific sections of the output.
- 129 • **Content**: the expected content within each section, with examples and rules to guide the generation process.
- 130 • **Restrictions**: limitations of the output (e.g., word count limitations, the inclusion of only the information provided in the
131 original protocol, and adherence to the criteria checklist for plain language as set out in *Supplementary Error!*
132 *Reference source not found.*).

133 After finalizing the prompt for generating a PPLS, we used the same approach to create a prompt to generate Cochrane Reviews
134 PLS (see the description of this data source in *Supplementary Table 1*, and the prompt in *Supplementary Table 5*).

135 We used our prompts with GPT 3.5 and GPT 4 to translate technical biomedical texts, Cochrane Reviews and Study Protocols,
136 into their respective PLS: Cochrane PLS and Protocol PLS. We quantitatively tested the generated PLS for plainness and
137 semantic equivalence. For PPLS, we also performed a qualitative assessment of the outputs by three experts in Clinical Trial
138 Operations and Regulatory Medical Writing, who rated each GPT 3.5 and GPT 4 text on a 5-point Likert Scale (1-Strongly
139 Disagree to 5-Strongly Agree). They evaluated the texts for accuracy, readability, completeness, and usefulness as defined in
140 the section: *LLM Prompt for Plain Language Summary Generation*.

141 Results

142 Plain Text Classification Model

143 The classification models accurately distinguished whether an input text was plain or technical. The Gradient Boosting model
144 showed slightly superior results with a precision rate of 97.2% (See Table 1).

145 **Table 1.** Comparison of tested classification models in terms of F1 Score or predictive performance, Accuracy, Recall, and Precision.

Model	F1 Score (Predictive Performance)	Accuracy	Recall	Precision
Random Forest	0.971	0.980	0.973	0.969
Gradient Boosting	0.975	0.982	0.977	0.972

146

147 **LLM Prompt for Plain Language Summary Generation**

148 **Cochrane Reviews: Plain Language Summaries**

149 We randomly selected a sample of 600 Cochrane texts from the *main dataset*: 300 technical abstracts and the corresponding
150 300 plain summaries. We then used our prompt in both GPT 3.5 and GPT 4 to generate the plain language summary from the
151 technical abstracts resulting in 300 Plain-GPT 3.5 and 300 Plain-GPT 4 summaries.

152 We tested the LLM-generated texts with our best model, Gradient Boosting, for plain language classification, and BERTScore
153 to test semantic equivalence against the original Cochrane plain summaries.

154 Our model classified 96% of GPT 3.5 texts and 99.6% of GPT 4 texts as plain. Hence, our prompt is effective in generating PLSs
155 that meet quantitative plain language requirements as defined in our classification model, with GPT 4 showing higher adherence.

156 The semantic equivalence score, BERTScore, confirmed both GPT 3.5 and GPT 4 successfully retained the original message.
157 However, GPT 4 produced plain summaries that outperformed GPT 3.5 in all parameters (Precision, Recall, and F1-Score) with
158 a significant difference (p -value < 0.05) (Table 2).

159 **Table 2.** Semantic equivalence score (BERT) between the GPT-generated plain summaries from Cochrane technical abstract vs. original
160 Cochrane PLS.

Semantic Equivalence	Plain_GPT 3.5	Plain_GPT4	p -value
Precision	0.790 ± 0.010	0.791 ± 0.015	0.027
Recall	0.772 ± 0.017	0.773 ± 0.016	0.003
F1-Score	0.780 ± 0.015	0.782 ± 0.014	0.001

161

162 **Protocol Plain Language Summaries**

163 We randomly selected a sample of nine clinical trial protocols from ClinicalTrials.Gov. Given that their corresponding PPLS
164 were not yet publicly published, we used Trial Summaries by Citeline Regulatory to find the corresponding Results Plain
165 Language Summaries (RPLS) and extracted four sections that are equivalent in a PPLS: 'Why is this study needed?':
166 Background and hypothesis of the trial (*Rationale*), 'Who will take part in this study?' (*Population*), 'How is this study designed?'
167 (*Trial Design*), and 'What treatments are being given during the study?' (*Interventions*).

168 Quantitative Analysis

169 We used our prompt specific for PPLS with both GPT 3.5 and GPT 4 to generate the plain language summary from the technical
170 protocols. We used our Gradient Boosting model to verify if LLM-generated texts were plain and BERTScore to check semantic
171 equivalence to the content on the RPLS. All LLM-generated PPLS were classified as plain, and BERTScore confirmed semantic
172 agreement with the content in the RPLS (Table 3). Consistent with Cochrane results, GPT 4 produced PPLS with higher semantic
173 equivalence than GPT 3.5 (no statistical analysis due to small sample size).

174 **Table 3.** Semantic equivalence score (BERT) between the GPT-generated PPLS from clinical trial protocols vs. the original content written for
175 the PLS.

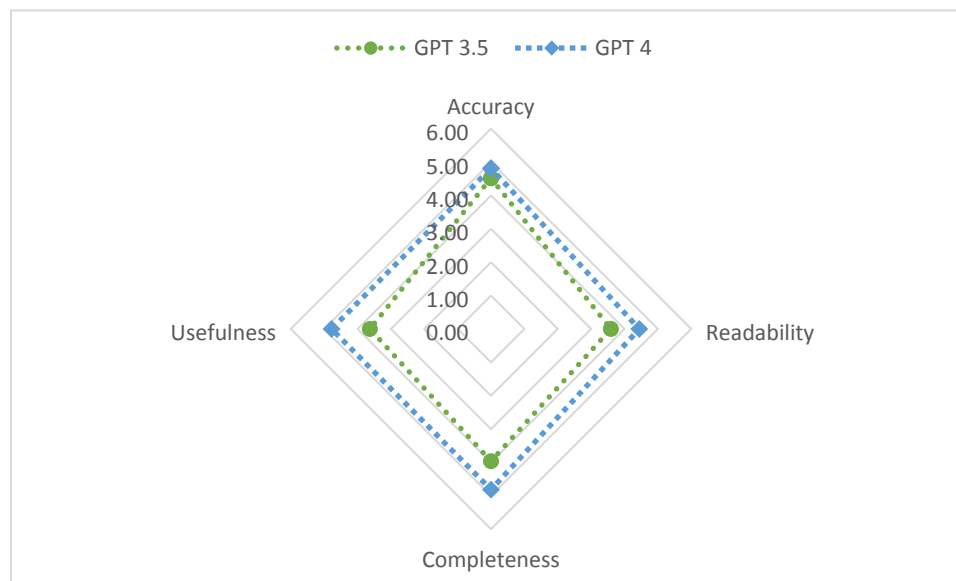
Semantic Equivalence	PPLS_GPT 3.5	PPLS_GPT4
Precision	0.8040 ± 0.0068	0.8073 ± 0.0208
Recall	0.7940 ± 0.0138	0.7975 ± 0.0129
F1-Score	0.7989 ± 0.0076	0.8023 ± 0.0109

176

177 Qualitative Analysis

178 Ratings by 3 domain experts who evaluated each LLM-generated text, demonstrated that GPT 4 outperformed GPT 3.5 in all
179 four criteria: accuracy, readability, completeness and usefulness, as indicated by an average score of 4.71 for GPT 4 texts as
180 compared to 3.93 for GPT 3.5 (see Figure 3 and Table 4).

181



182 **Figure 3.** Radar diagram comparing the qualitative assessment of the LLM-generated texts in 4 criteria: Accuracy, Readability, Completeness,
183 and Usefulness.

184 **Table 4.** Ratings for GPT 3.5 and GPT 4 plain summaries in 4 criteria: Accuracy, Readability, Completeness, and Usefulness.

	Accuracy	Readability	Completeness	Usefulness	Overall Score
GPT 3.5	4.52	3.59	3.96	3.63	3.93
GPT 4	4.81	4.44	4.81	4.78	4.71

185

186 In terms of accuracy, both GPT 3.5 and GPT 4 received high scores. Reviewers noted that both language models exhibited
187 scientific accuracy and relied exclusively on the input text (study protocol). Notably, even when the content in the original RPLS
188 contained inconsistencies (e.g. incorrect age limit or indication), both language models generated accurate PLS. This finding
189 suggests that language models can be used to automatically generate a first draft of a PLS while minimizing data inaccuracies
190 resulting from human error.

191 Regarding readability, both GPT 3.5 and GPT 4 generated texts that were likely to be understood by a lay audience. This
192 observation aligned with the results obtained through the classification model. However, GPT 3.5 occasionally employed
193 complicated medical jargon (e.g., 'chronic', 'randomized', 'double-blind') and longer words and sentences (e.g., 'approximately
194 640 adults' vs 'about 640 adults'). Similarly, GPT 4, despite its outstanding performance, occasionally preferred passive voice
195 over active voice, compromising clarity and concise writing. This highlights the importance of quality control by a healthcare
196 professional who should verify the content and style of the automatically generated PLS draft.

197 Completeness, which assessed the compliance of PPLS content and structure with EU CTR No 536/2014 guidelines, revealed
198 inconsistencies in the outputs generated by GPT 3.5. These inconsistencies manifested as the creation of new, unrequested
199 sections and summaries, with significant variation among the nine generated PLS. Conversely, GPT 4 consistently generated
200 PLSs that adhered to the specified format and content expectations, and complied with the guidelines, showing a remarkable
201 value in automating the time-consuming task of guaranteeing the content to be standardized and aligned with industry specific
202 and rigorous guidelines.

203 The usefulness ratings, indicating the suitability of the generated PLSs as draft versions, correlated with the findings in other
204 criteria. GPT 3.5 received moderate scores in generating draft PLS, while GPT 4 scored 4.78, indicating that the generated PLS
205 were highly suitable as draft versions of the PLS.

206 **Discussion**

207 In this study, we used NLP and LLMs to improve health literacy by generating PLS from biomedical texts. Our two-part strategy
208 involved creating a classification model for identifying if a text was written in plain language, and using LLMs (specifically GPT
209 3.5 and GPT 4) for the automated generation of the PLS.

210 The classification model achieved over 97% accuracy, indicating its effectiveness in distinguishing between the text types:
211 technical and plain. This is a very useful stand-alone strategy which could support authoring teams in identifying if their texts
212 targeted for patients or the general audience are compliant with plain language guidelines.

213 The LLMs exhibited outstanding performance in generating PLS, with GPT 4 outperforming GPT 3.5 in creating content that was
214 both plain and semantically similar. In a qualitative review by domain experts, GPT 4 also surpassed GPT 3.5 by generating
215 high-quality drafts of PLS. These drafts were scientifically accurate, compliant with plain language requirements, and met
216 expectations in content and structure. These results underlines the value of LLMs in supporting healthcare stakeholders to
217 streamline the generation of plain documents, and with that, promote equitable access to biomedical information, engagement
218 of the lay audience in health-related decision making, and improved health outcomes.

219 Our study highlights the importance of using well-designed, structured, and domain-specific prompts to guarantee the creation
220 of high-quality, easily comprehensible PLS. This is particularly vital when accuracy in biomedical facts is essential. This requires
221 the collection of feedback from stakeholders who are experts in the domain or field of interest. Such feedback would help to fine-
222 tune the prompts and guarantee that the output fullfils the purposes of different document types. Our study exemplified this with
223 various document types (e.g., Cochrane reviews, PPLS), some of which adhere to strict industry standards.

224 While the findings of our study are promising, they also underscore opportunities for further research to fully harness the potential
225 of NLP and LLMs in this context. Future studies could involve direct audience feedback in evaluating the understandability of
226 PLS. This would ensure that the generated content aligns with the comprehension levels of the intended audience, such as
227 patients in clinical settings, and would provide cues for ways in which the they could improve their interaction with biomedical
228 content, improving adherence to treatment plans or educating them about a disease or diagnosis. Additionally, depending on
229 the intended use and field of interest, refining the models could potentially account for specific linguistic nuances, exploring
230 advanced techniques like Retrieval Augmented Generation (RAG) could enhance factual accuracy, and expanding the dataset
231 to include a wider range of texts and languages could enhance the generalizability of the classification model and applicability
232 of the LLMs. Different interesting opportunities to leverage NLP and LLMs to serve society by simplifying what would otherwise
233 be daunting.

234 In conclusion, by leveraging the capabilities of NLP and LLMs, we have taken a significant step towards bridging the gap between
235 complicated biomedical texts and comprehensible summaries designed for the general audience. This framework paves the way
236 for prospective innovations in the field of health literacy, which, in turn, holds the potential to enhance health outcomes and
237 foster health equity.

238 References

- 239 1. Nielsen-Bohlman L, Panzer AM, Kindig DA. Health Literacy: A Prescription to End Confusion: National Academies
240 Press 2004;2, What Is Health Literacy? <https://www.ncbi.nlm.nih.gov/books/NBK216035/>
241
- 242 2. Berkman ND, Sheridan SL, Donahue KE, Halpern DJ, Crotty K. Low health literacy and health outcomes: an updated
243 systematic review. *Annals of Internal Medicine* 2011;155(2):97-107 [https://www.acpjournals.org/doi/10.7326/0003-](https://www.acpjournals.org/doi/10.7326/0003-4819-155-2-201107190-00005?url_ver=Z39.88-2003&rfr_id=ori:rid:crossref.org&rfr_dat=cr_pub%20%20pubmed)
244 [4819-155-2-201107190-00005?url_ver=Z39.88-2003&rfr_id=ori:rid:crossref.org&rfr_dat=cr_pub%20%20pubmed](https://www.acpjournals.org/doi/10.7326/0003-4819-155-2-201107190-00005?url_ver=Z39.88-2003&rfr_id=ori:rid:crossref.org&rfr_dat=cr_pub%20%20pubmed)
245
- 246 3. Berkman ND, Sheridan SL, E Donahue K, et al. Health literacy interventions and outcomes: an updated systematic
247 review. *Evidence Report/ Technology Assessment* 2011;199:1-941
248 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4781058/>
249
- 250 4. Miller TA. Health literacy and adherence to medical treatment in chronic and acute illness: A meta-analysis. *Patient*
251 *Education and Counseling* 2016;99(7):1079-1086 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4912447/>
252
- 253 5. Sørensen K, Røthlin F, Pelikan JM, et al. Health literacy in Europe: comparative results of the European health literacy
254 survey (HLS-EU). *European Journal of Public Health* 2015;25(6):1053-1058
255 <https://academic.oup.com/eurpub/article/25/6/1053/2467145?login=true>
256
- 257 6. Bahador B, Baedorf Kassis S, Gawrylewski H, et al. Promoting equity in understanding: A cross-organizational plain
258 language glossary for clinical research. *Medical Writing* 2020;29(4):10-15 [https://journal.emwa.org/writing-for-](https://journal.emwa.org/writing-for-patients/promoting-equity-in-understanding-a-cross-organisational-plain-language-glossary-for-clinical-research/)
259 [patients/promoting-equity-in-understanding-a-cross-organisational-plain-language-glossary-for-clinical-research/](https://journal.emwa.org/writing-for-patients/promoting-equity-in-understanding-a-cross-organisational-plain-language-glossary-for-clinical-research/)
260
- 261 7. Stormacq C, Van den Broucke S, Wosinski J. Does health literacy mediate the relationship between socioeconomic
262 status and health disparities? *Integrative review. Health Promotion International* 2019;34(5):e1-e17
263 <https://academic.oup.com/heapro/article-abstract/34/5/e1/5068634?redirectedFrom=fulltext&login=true>
264
- 265 8. Schillinger D. Social Determinants, Health Literacy, and Disparities: Intersections and Controversies. *HLRP: Health*
266 *Literacy Research and Practice* 2021;5(3):233-243 [https://journals.healio.com/doi/full/10.3928/24748307-20210712-](https://journals.healio.com/doi/full/10.3928/24748307-20210712-01?rfr_dat=cr_pub++0pubmed&url_ver=Z39.88-2003&rfr_id=ori%3Arid%3Acrossref.org)
267 [01?rfr_dat=cr_pub++0pubmed&url_ver=Z39.88-2003&rfr_id=ori%3Arid%3Acrossref.org](https://journals.healio.com/doi/full/10.3928/24748307-20210712-01?rfr_dat=cr_pub++0pubmed&url_ver=Z39.88-2003&rfr_id=ori%3Arid%3Acrossref.org)

268

269 9. GDPR. General Data Protection Regulation (GDPR) - The principle of Transparency. Accessed December 22, 2023.

270 [https://health.ec.europa.eu/latest-updates/updated-document-questions-answers-clinical-trials-regulation-eu-no-](https://health.ec.europa.eu/latest-updates/updated-document-questions-answers-clinical-trials-regulation-eu-no-5362014-2023-09-29_en)

271 [5362014-2023-09-29_en](https://health.ec.europa.eu/latest-updates/updated-document-questions-answers-clinical-trials-regulation-eu-no-5362014-2023-09-29_en)

272

273 10. Trezona A, Rowlands G, Nutbeam D. Progress in Implementing National Policies and Strategies for Health Literacy-

274 What Have We Learned so Far? International Journal of Environmental Research and Public Health 2018;15(7):1554

275 <https://www.mdpi.com/1660-4601/15/7/1554>

276

277 11. Centers for Disease Control and Prevention. Your Guide to CLEAR WRITING. May 9, 2022. Accessed November 15,

278 2023. <https://www.cdc.gov/nceh/clearwriting/docs/clear-writing-guide-508.pdf>

279

280 12. B. Ondov, K. Attal and D. Demner-Fushman, "A survey of automated methods for biomedical text simplification,"

281 Journal of the American Medical Informatics Association, vol. 29, no. 11, pp. 1976-1988, 2022.

282

283 13. Stoll M, Kerwer M, Lie K, Chasiotis A. Plain language summaries: A systematic review of theory, guidelines, and

284 empirical research. PLoS ONE 2022; 17(6): e0268789

285 journals.plos.org/plosone/article?id=10.1371/journal.pone.0268789

286

287 14. The Plain Language Action and Information Network. (2011). Federal Plain Language Guidelines. pp. 1–14. Accessed

288 November 20, 2023. <https://www.plainlanguage.gov/media/FederalPLGuidelines.pdf>

289

290 15. Brach C, ed. AHRQ Health Literacy Universal Precautions Toolkit, 3rd Edition. Rockville, MD. Agency for Healthcare

291 Research and Quality. AHRQ Publication No. 23-0075. Accessed November 20, 2023. [https://www.ahrq.gov/health-](https://www.ahrq.gov/health-literacy/improve/precautions/index.html)

292 [literacy/improve/precautions/index.html](https://www.ahrq.gov/health-literacy/improve/precautions/index.html)

293

294 16. United Health Group. Just Plain Clear Glossary. Accessed December 5, 2023. <https://www.justplainclear.com/en>

295

296 17. European Union. Q&A: Clinical Trial Regulation (EU) No 536/2014 2023. Accessed December 26, 2023.

297 https://health.ec.europa.eu/system/files/2023-09/regulation5362014_qa_en.pdf

298

- 299 18. The Python Package Index (PyPI) Readability 0.3.1. 2023. January 12, 2019. Accessed November 2023
300 <https://pypi.org/project/readability/>
301
- 302 19. SpaCy. 2016-2023. Accessed November 2023. <https://spacy.io/>
303
- 304 20. ClinicalTrials.gov by U.S National Library of Medicine (NIH). Accessed November 2023
305 <https://www.clinicaltrials.gov/about-site/about-ctg>
306
- 307 21. ClinicalTrials.gov API by U.S National Library of Medicine (NIH). Accessed November 2023
308 <https://classic.clinicaltrials.gov/api/gui>
309
- 310 22. The Python Package Index (PyPI) selenium 4.15.2, Python Software Foundation, 2023. Accessed December 2022
311 <https://pypi.org/project/selenium/>
312
- 313 23. The Python Package Index (PyPI) beautifulsoup 4 4.12.2, Python Software Foundation, 2023. Accessed December
314 2022 <https://pypi.org/project/beautifulsoup4/>
315
- 316 24. Pfizer Plain Language Study Results Summaries, 2023. Accessed September 2023
317 <https://www.pfizer.com/science/clinical-trials/plain-language-study-results-summaries>
318
- 319 25. Citeline Trial Summaries Citeline Regulatory, Pharma Intelligence UK Limited. Accessed September 2023
320 <https://www.trialssummaries.com/Home/LandingPage>
321

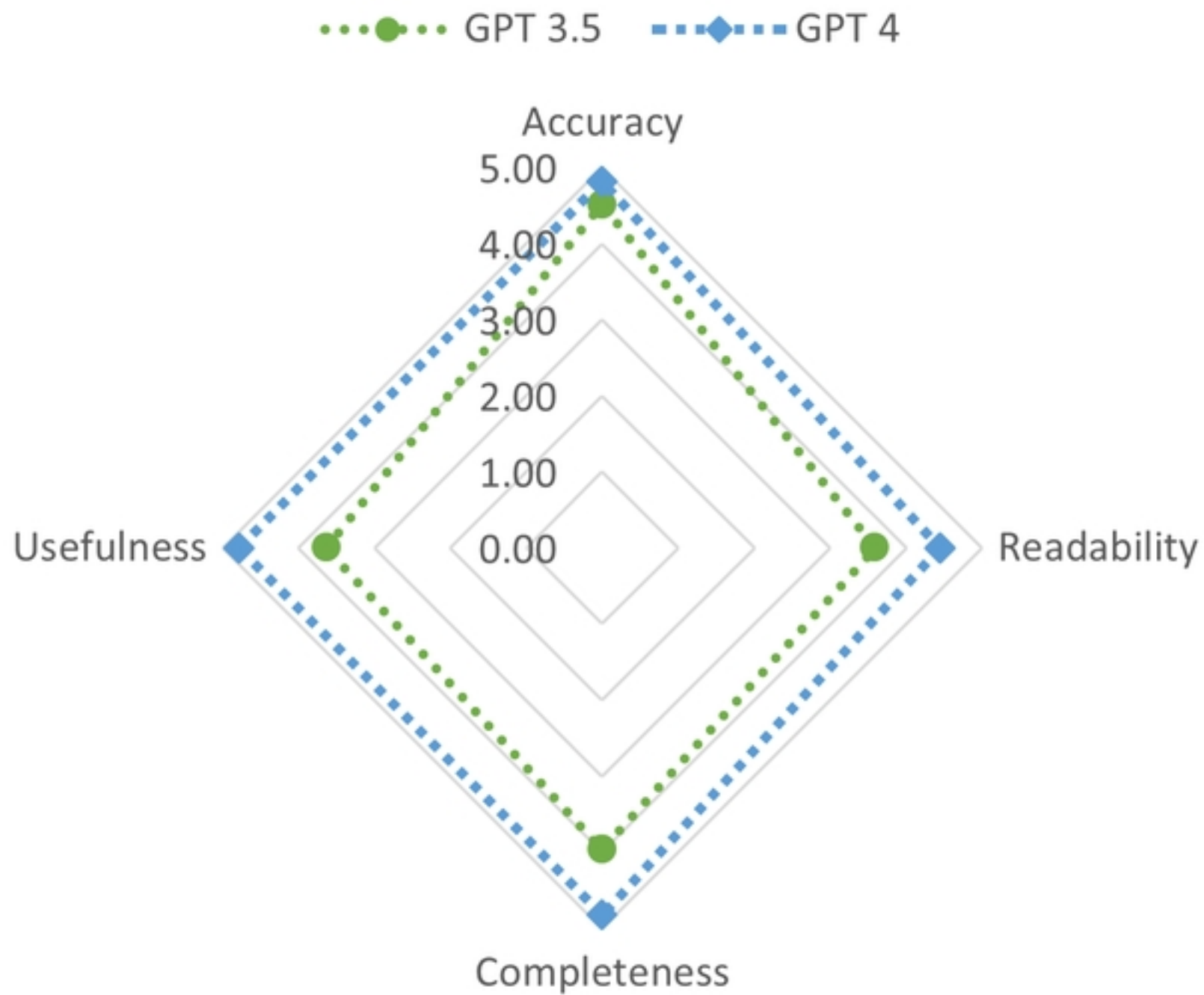


Fig3

Histogram of interjections

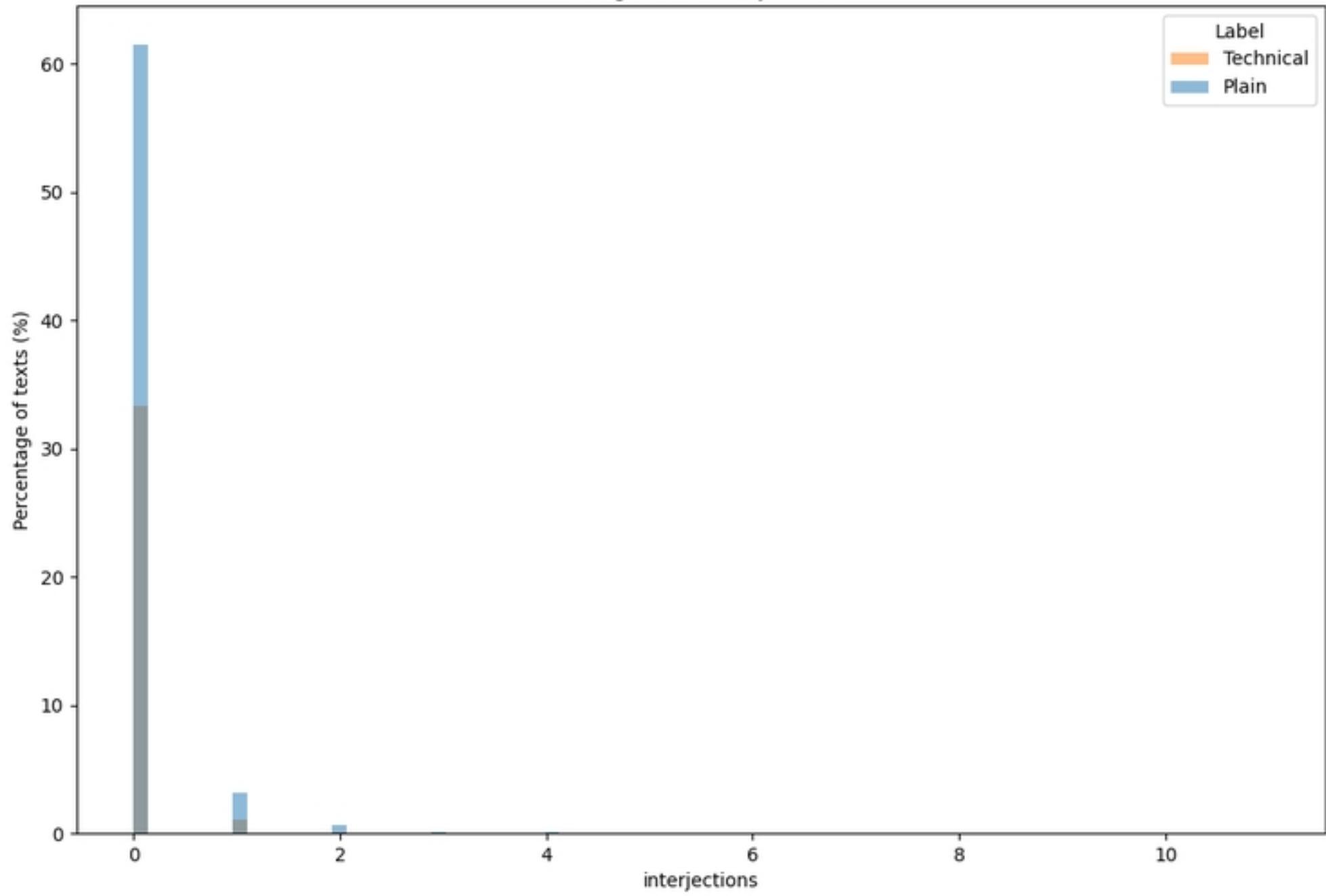


Fig2a

Histogram of passive_voice

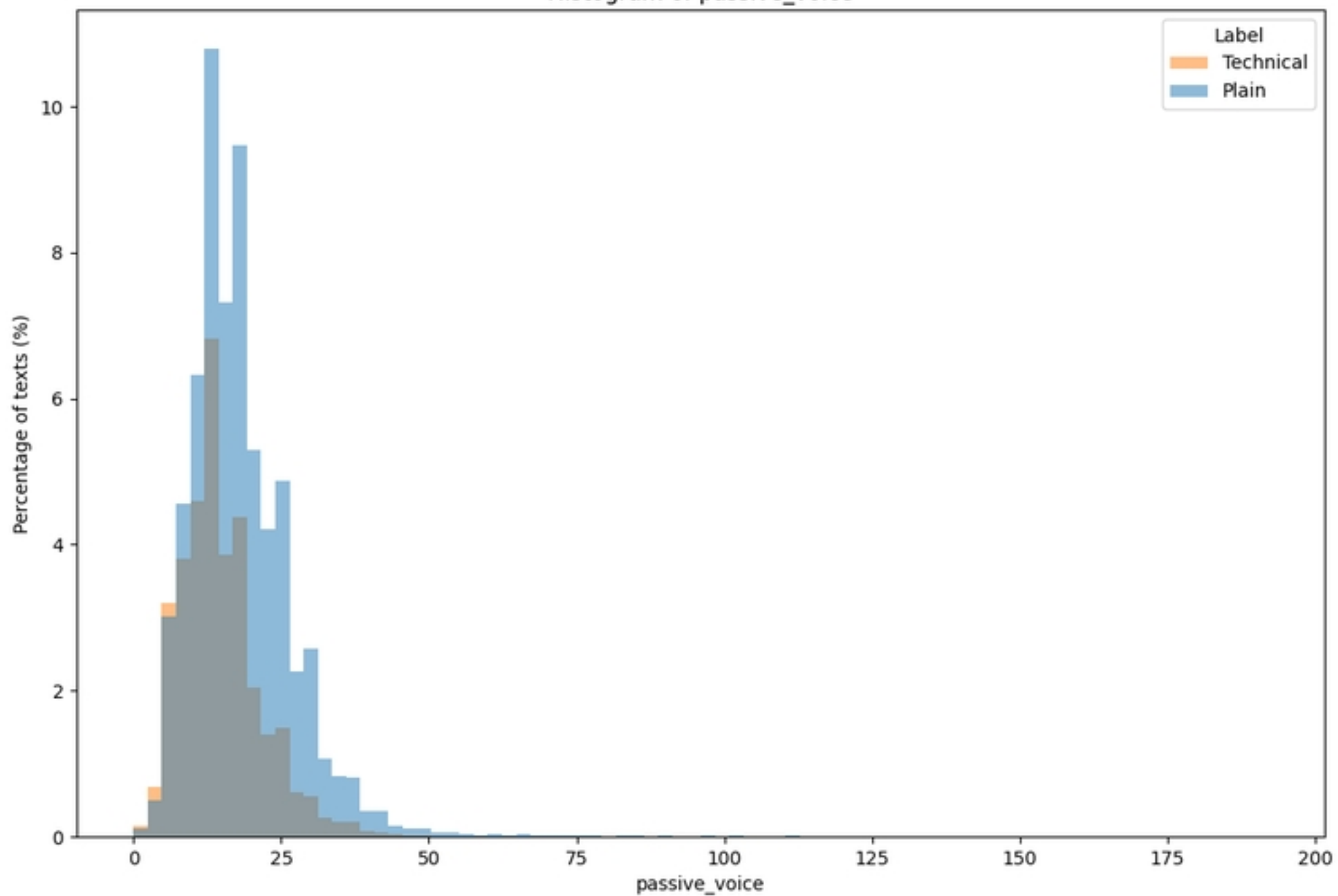


Fig2b

Histogram of stopwords

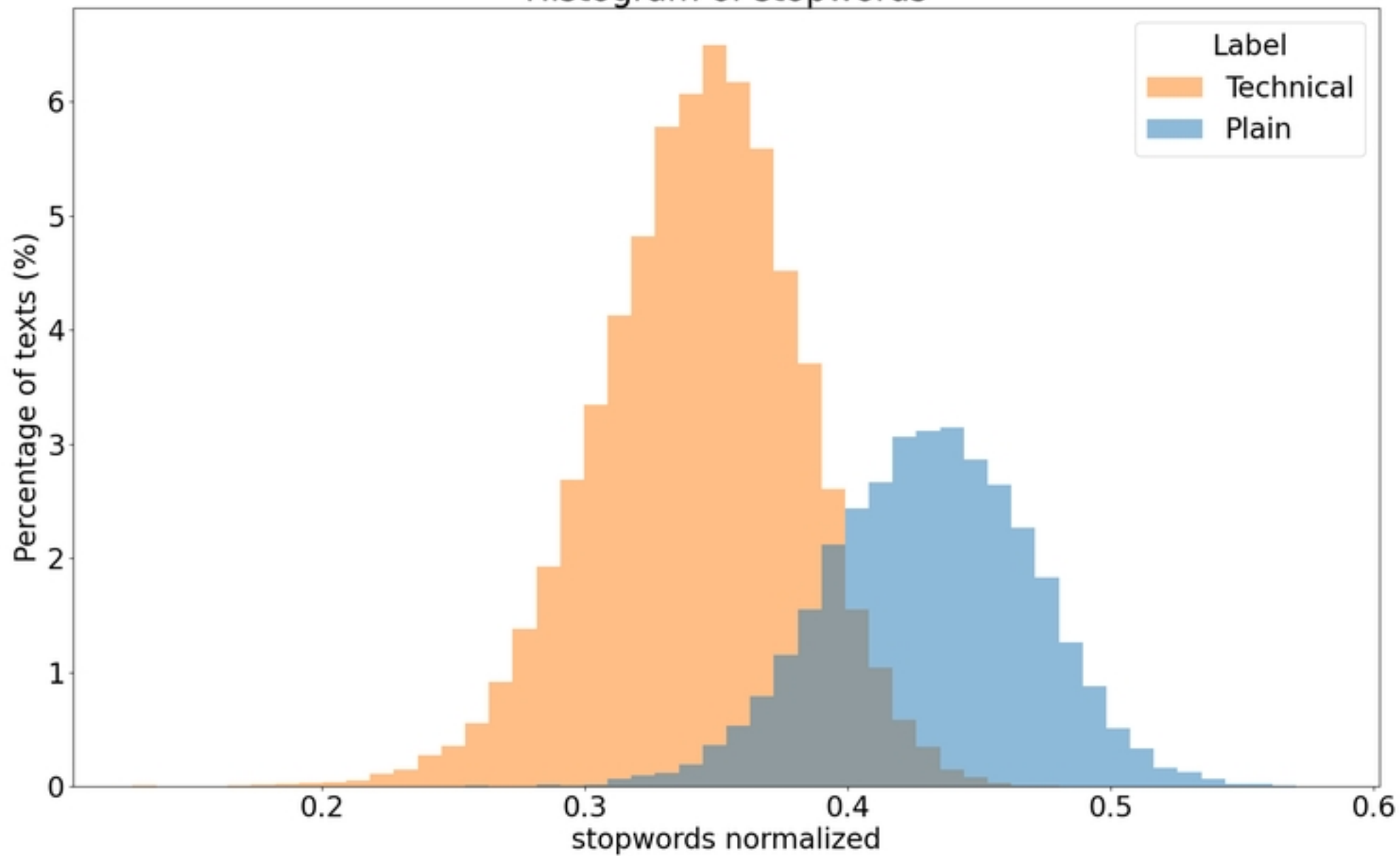


Fig2c

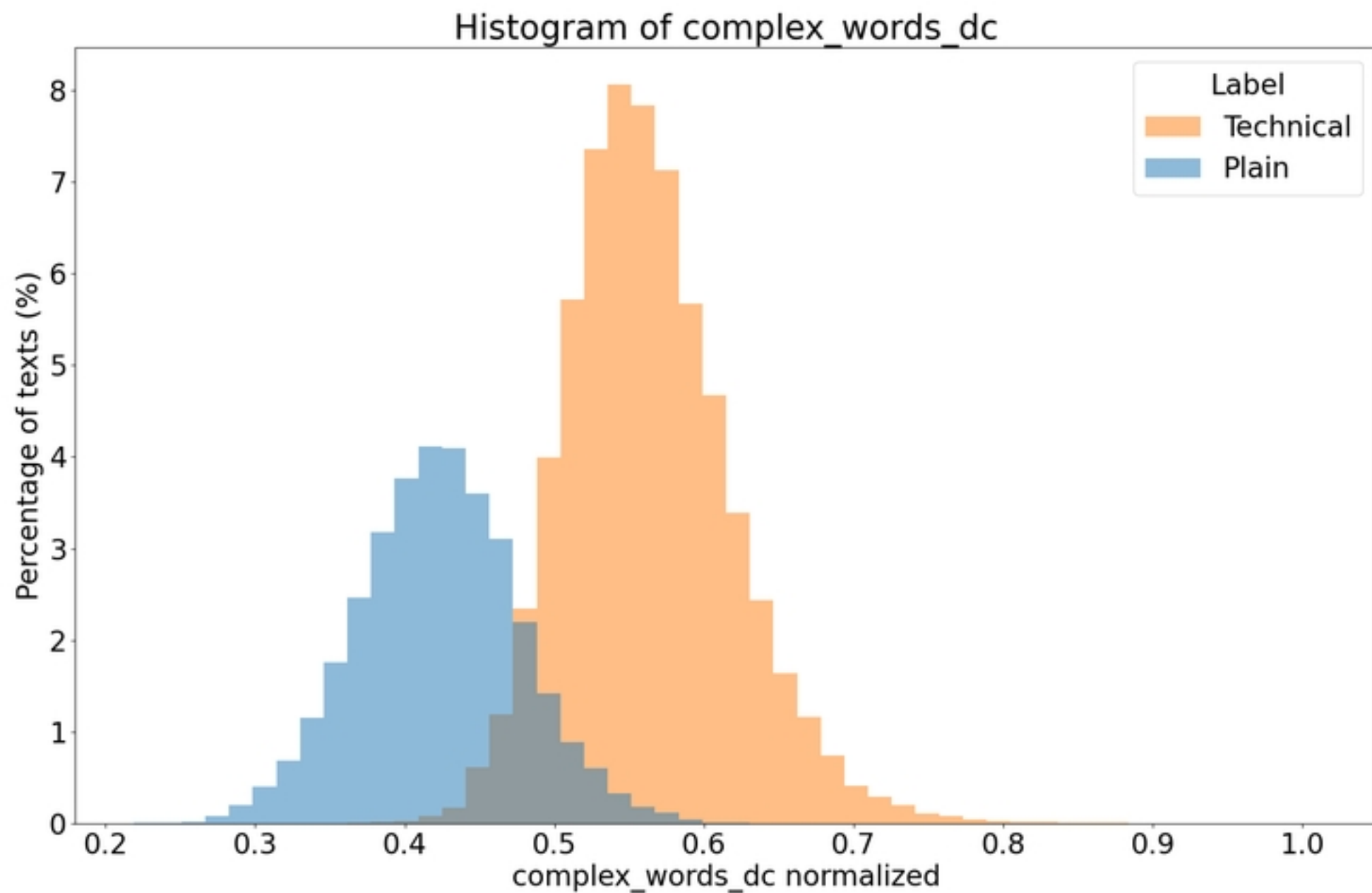
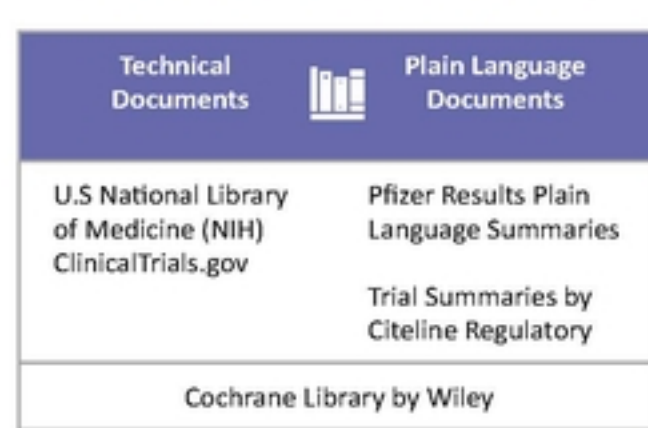


Fig2d

1. Data Collection and Processing

2. Analysis of Plain Language

3. Testing LLMs



Quantitative analysis of readability metrics and language variables to create Plain Text Classification Model.



Qualitative analysis of subset of documents to create a criteria checklist of plain texts linguistic attributes and LLM prompt for PLS generation.

Model Prompt Anatomy

Context or Rationale

Using the following clinical trial protocol text as input, generate a Protocol Plain Language Summary (PPLS) compliant with GDPR and understandable by any patient, regardless of their health literacy.

Output: Structure and Format

The generated PPLS should be presented in a logical order, using the following headings: [1, 2, 3, ..., n]

Content: Examples and rules

In the section [] : Answer "What are the goals of the study?" Specify the main and secondary objectives of the trial and how they will be measured (eg. the main trial endpoint is the percent change in the number...).

Restrictions and Considerations

The generated PPLS must follow these plain language guidelines: [...] The AI model should not invent information or add content that is not present in the input protocol.

*This is an excerpt, the complete prompt can be found in supplementary tables

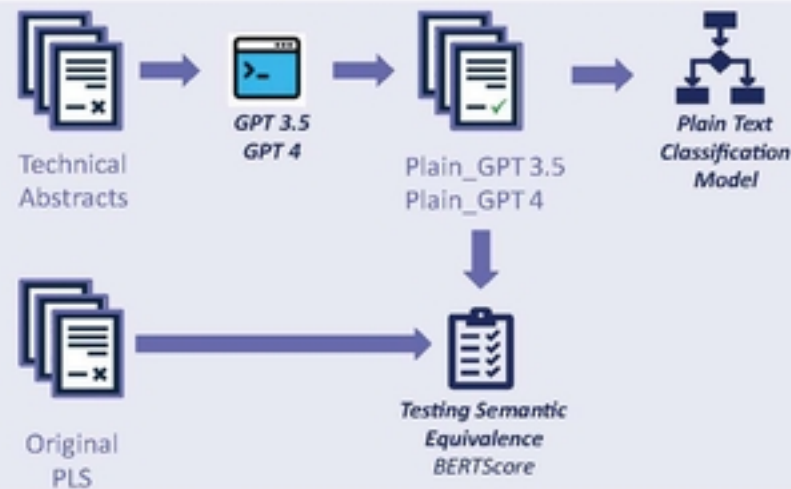
Technical Version

An echocardiogram will be used to diagnose cardiac failure.

Plain Language Version

An echocardiogram (a non-invasive test that uses sound waves to create images of the heart) will be used to identify if your heart is not working right.

Cochrane Library



Clinical Trial Study Protocols

