# Graphical Abstract

**Retrieval Augmentation of Large Language Models for Lay Language Generation**
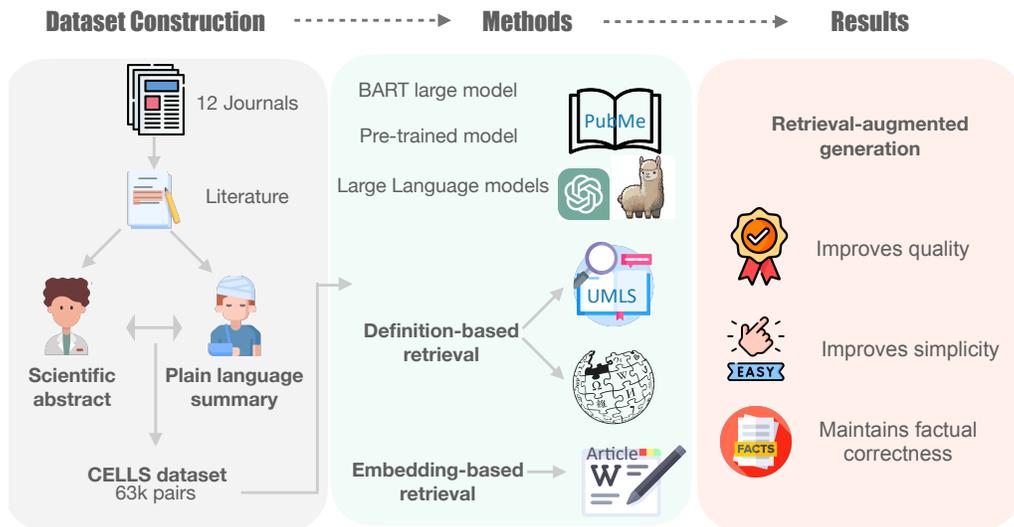
Yue Guo [1*], Wei Qiu [2*], Gondy Leroy[3], Sheng Wang [2], Trevor Cohen[1]
[1] Biomedical and Health Informatics, University of Washington
[2] Paul G. Allen School of Computer Science, University of Washington
[3] Management Information Systems, University of Arizona
[*] These authors contributed equally to this work.

# Highlights

**Retrieval Augmentation of Large Language Models for Lay Language Generation**

Yue Guo [1][*], Wei Qiu [2][*], Gondy Leroy[3], Sheng Wang [2], Trevor Cohen[1]
[1] Biomedical and Health Informatics, University of Washington
[2] Paul G. Allen School of Computer Science, University of Washington
[3] Management Information Systems, University of Arizona
[*] These authors contributed equally to this work.

- **Problem:** Automated lay summary generation can improve the accessibility of health information, but is challenging because of the need to provide background information absent in source documents.

- **What is already known:** Current models face constraints due to corpus size, topic diversity, and untested utility of external information retrieval.

- **What this paper adds:** We approach lay language generation by simplifying content and also generating background explanations, achieved through innovative Retrieval-Augmented Lay Language (RALL) methods. This paper also introduces CELLS, the largest (63k pairs) and most diverse (12 journals) parallel corpus for lay language generation, with a specialized subset to advance background explanation capabilities.

# Retrieval Augmentation of Large Language Models for Lay Language Generation

Yue Guo [1*], Wei Qiu [2*], Gondy Leroy[3], Sheng Wang [2], Trevor Cohen[1]

[1] Biomedical and Health Informatics, University of Washington

[2] Paul G. Allen School of Computer Science, University of Washington

[3] Management Information Systems, University of Arizona

[*] These authors contributed equally to this work.

## Abstract

The complex linguistic structures and specialized terminology of expert-authored content limit the accessibility of biomedical literature to the general public. Automated methods have the potential to render this literature more interpretable to readers with different educational backgrounds. Prior work has framed such lay language generation as a summarization or simplification task. However, adapting biomedical text for the lay public includes the additional and distinct task of *background explanation*: adding external content in the form of definitions, motivation, or examples to enhance comprehensibility. This task is especially challenging because the source document may not include the required background knowledge. Furthermore, background explanation capabilities have yet to be formally evaluated, and little is known about how best to enhance them. To address this problem, we introduce Retrieval-Augmented Lay Language (RALL) generation, which intuitively fits the need for external knowledge beyond that in expert-authored source documents. In addition, we introduce CELLS, the largest (63k pairs) and broadest-ranging (12 journals) parallel corpus for lay language generation. To evaluate RALL, we augmented state-of-the-art text generation models with information retrieval of either term definitions from the UMLS and Wikipedia, or embeddings of explanations from Wikipedia documents. Of these, embedding-based RALL models improved summary quality and simplicity while maintaining factual correctness, suggesting that Wikipedia is a helpful source for background explanation in this context. We also evaluated the ability of both open-soured Large Language Model (Llama 2) and closed-sourced Large Language Model (GPT-4) in background expla-

nation, with and without retrieval augmentation. Results indicate that these LLMs can generate simplified content, but that the summary quality is not ideal. Taken together, this work presents the first comprehensive study of background explanation for lay language generation, paving the path for disseminating scientific knowledge to a broader audience. Our code and data are publicly available at: https://github.com/LinguisticAnomalies/pls_retrieval.

## 1. Introduction

The COVID-19 pandemic underscored the difficulties the general public faces when attempting to use scientific information to guide their health-related decisions Soroya et al. (2021); Bin Naeem and Kamel Boulos (2021). Though widely *available* in scientific papers and preprints, the information required to guide health-related decision making is often not *accessible*: medical jargon Korsch et al. (1968), scientific writing styles Kurtzman and Greene (2016), and insufficient scientific background Crossley et al. (2014) make this information opaque to non-experts. Consequently, there is a pressing need to deliver scientific knowledge in lay language, which has motivated research on automated generation of lay language summaries.

Prior work has framed lay language generation as a summarization or simplification task Guo et al. (2021); Devaraj et al. (2021). However, adapting biomedical text for the lay public includes the distinct task of *background explanation*: adding external content in the form of definitions, history, or examples to enhance comprehensibility. Cognitive studies of text comprehension suggest that providing missing background information effectively improves reader comprehension, especially when readers lack the prerequisite domain knowledge to fill this in themselves McNamara et al. (1996).

Text simplification, which modifies content to improve readability while retaining its key points, has been widely studied Jonnalagadda et al. (2009); Qenam et al. (2017). However, generating background information is especially challenging, because the source document may not include the required background knowledge. Furthermore, background explanation capabilities have yet to be formally evaluated, and little is known about how best to enhance them. Retrieval augmentation methods, which use information retrieval to identify additional content to inform text generation, present an

2

intuitive fit for the need to acquire external knowledge. In the current work, we explore methods for Retrieval-Augmented Lay Language (RALL) generation, augmenting state-of-the-art text generation models with information retrieval of either term definitions from the UMLS Bodenreider (2004) and Wikipedia, or embeddings of explanations from Wikipedia documents Lewis et al. (2020). Our findings indicate that RALL models improve summary quality and simplicity while maintaining factual correctness, suggesting that general knowledge from Wikipedia in particular is a good source for background explanation. With Large Language Models (LLMs) becoming increasingly accessible, we also tested the ability of two LLMs for background explanation: we prompted both open-source Llama 2 Touvron et al. (2023) and closed-source GPT-4 OpenAI (2023) with and without external knowledge from Wikipedia. Results indicate that these LLMs improve simplicity but do not preserve the summary quality.

Abstractive summarization methods require source/summary pairs, with the summaries written in plain language. The limited size and topical breadth of publicly-available paired corpora constrain the scope of applicability of models trained for this task and limit the generalizability of published evaluations. Therefore, a further contribution of this work is the Corpus for Enhancement of Lay Language Synthesis (CELLS): 62,886 pairs of scientific abstracts with corresponding lay language summaries (Table 1). Summaries are written by abstract authors or other domain experts, assuring the quality of our dataset. CELLS is larger and more diverse than prior datasets Guo et al. (2021); Devaraj et al. (2021), aggregating papers from 12 journals (Table 2) spanning various biomedical domains. From CELLS, we derived a set of specialized paired corpora: 233,916 algorithm-aligned sentence pairs for *simplification* and 47,157 scientific/lay-language pairs emphasizing novel content that is absent from scientific abstracts for *background explanation* to support our research on background augmentation.

## 2. Related Work

### 2.1. Lay language summary generation

Text summarization and simplification are two important aspects of lay language generation. Text summarization is a widely-studied research topic Cohan et al. (2018); Cachola et al. (2020); Devaraj et al. (2022). It has been a focus of research attention in the biomedical domain Mishra et al. (2014); Bui et al. (2016); Givchi et al. (2022), with applications including

**Abstract:** Clinical reports of Zika Virus (ZIKV) RNA detection in breast milk have been described, but evidence conflicts as to whether this RNA represents infectious virus...

**Summary:** *Only 4 years have passed since the Zika virus outbreak in Brazil, and much remains to be understood about the transmission and health consequences of Zika infection.* To date, some case reports have detected Zika virus RNA in the breast milk of infected mothers, but the presence of a virus' RNA does not mean that intact virus is present...

Table 1: Example abstract/summary pair from CELLS. The lay language summary is written by the abstract's authors. There are two challenges in lay language summary generation: generating background explanations (*italicized*) and simplifying the original abstract (underlined).

summarization of radiology reports Cai et al. (2021); Zhang et al. (2018, 2020), biomedical literature Wang et al. (2021); Plaza (2014); Cai et al. (2022) and medical dialogue Chintagunta et al. (2021); Joshi et al. (2020). Several of the biomedical text simplification datasets and methods have also been reported in the literature Jonnalagadda et al. (2009); Li et al. (2020); Cao et al. (2020); Lu et al. (2023). However, these were designed for sentence-level text simplification, rather than translation of paragraphs and longer documents into interpretable lay language. Compared to other paragraph-level lay language generation efforts Guo et al. (2021); Devaraj et al. (2021), the current work is the first to focus on background explanation generation.

*2.2. Lay language summarization datasets*

Previous endeavors towards developing datasets for automated conversion of scientific text into lay language have been limited in scale and scope. The CL-SciSumm 2020 shared task series Chandrasekaran et al. (2020) provided a training dataset encompassing 572 articles and corresponding author-constructed lay summaries, collated from a diverse array of scientific journals published by Elsevier. Guo et al. Guo et al. (2021) and Devaraj et al. Devaraj et al. (2021) introduce datasets of ∼5k scientific abstract and lay language summary pairs drawn from systematic reviews in the Cochrane Library. Goldsack et al. Goldsack et al. (2022) present ∼30k biomedical literature abstract pairs from PLOS and eLife. Luo et al. Luo et al. (2022) developed a dataset from ∼28k biomedical abstract pairs from PLOS. Attal et al. Attal et al. (2023) describe the PLABA dataset, encompassing 750

pairs of abstracts, each set featuring a sentence-aligned adaptation generated by human authors. The dataset developed for our study differs from these prior efforts in that: 1) CELLS is a large (∼63K) abstract-level dataset which includes different article types besides systematic reviews; and 2) we address the need for background explanation in lay language generation, deriving a specialized subset emphasizing content that is absent from the abstracts.

## 2.3. Lay language summary generation methodologies

Present text summarization methodologies predominantly fall into two categories: extractive and abstractive Das and Martins (2007). Extractive summarization involves ranking and selecting critical elements of the original text and combining them to form a condensed version Erkan and Radev (2004); Cheng and Lapata (2016). In contrast, abstractive summarization introduces novel words and phrases absent from the original text Gupta and Gupta (2019). The necessity to provide pertinent background, explain terminology, and apply straightforward sentence structures makes lay language summarization intrinsically an abstractive task Guo et al. (2021). The emergence of Transformer-based approaches such as BART, T5, and PEGA-SUS has significantly advanced this field Zhang et al. (2021); Yadav et al. (2022). BART, especially when pre-trained on domain-specific data, has demonstrated strong performance in the simplification of biomedical review articles Guo et al. (2021); Goldsack et al. (2022) and the summarization of randomized controlled trials Wallace et al. (2021). We employ BART as the benchmark model in the current work, including a variant with additional PubMed-specific pre-training. In addition, newer work has indicated that auto-regressive LLMs can outperform other Transformer models in lay language generation tasks Goldsack et al. (2023). Therefore, we also evaluated the performance of two such LLMs, including Llama 2 Touvron et al. (2023) and GPT-4 OpenAI (2023), on our dataset.

## 2.4. Retrieval-augmented text generation

Background explanation helps laypeople understand biomedical concepts Srikanth and Li (2020). Furthermore, experiments in cognitive psychology have shown that providing explanatory content improves the recall of readers with limited domain knowledge Britton and Gülgöz (1991); McNamara et al. (1996). Information retrieval methods present an intuitive approach to identify content to inform background explanations, with established utility for clinical question answering Simpson et al. (2014); Roberts et al. (2015); Luo

et al. (2022), biomedical text summarization Alambo et al. (2022); Mishra et al. (2014); Plaza (2014) and clinical outcome prediction Naik et al. (2021). There are two main categories of information retrieval methods that have been used to augment the generation of natural language text. *Definition-based* retrieval methods identify terms that exist in predefined lexicons, and use their definitions to inform text generation Alambo et al. (2022); Moradi and Ghadiri (2018). *Embedding-based* retrieval methods retrieve documents with similar low-dimensional representations, instead of depending upon lexical overlap between terms Deerwester et al. (1990); Cao and Xiong (2018); Guu et al. (2020); Karpukhin et al. (2020); Lewis et al. (2020). Retrieval augmentation has been shown to improve the performance of question answering systems Lewis et al. (2020), and reduce the frequency of so-called "hallucinations" (statements without grounding in training data) in text generated by language models Shuster et al. (2021). However, these approaches have not been explored for lay language generation, despite their intuitive fit to the subtask of background explanation in particular. In the current work, we explore both definition- and embedding-based retrieval approaches and evaluate the utility of external information from the UMLS and Wikipedia for this important subtask.

## 3. Materials and methods

### 3.1. The CELLS Dataset

We present CELLS, the largest dataset of parallel scientific abstracts and expert-authored lay language summaries (LLSs) developed to date (Section 3.1.1), offering unique opportunities to study the performance of lay language generation models. To facilitate research on key LLS generation subtasks, we have also derived subsets for simplification and background explanation (Section 3.1.2).

### 3.1.1. Data compilation

To develop CELLS, we manually reviewed biomedical journals and identified 19 with a LLS section (see Appendix Table A.2). We collected scientific abstracts (*source*) and their aligned LLSs (*target*) from these journals. We excluded abstracts where LLSs are not associated with a full-length paper (i.e., LLS in a separate section for the journal's website or social media feed) that required extensive human inspection. After further excluding non-biomedical topics, we obtained 75,205 pairs of abstracts and LLSs. To ensure

6

|                                      |         | Length |     |
|--------------------------------------|---------|--------|-----|
| Journal                              | Num.    | Src    | Tgt |
| PNAS                                 | 25,647  | 227    | 124 |
| PLOS Genetics                        | 8,030   | 256    | 192 |
| PLOS Pathogens                       | 7,345   | 260    | 193 |
| PLOS Neglected Tropical Diseases     | 7,185   | 315    | 198 |
| PLOS Computational Biology           | 7,072   | 253    | 188 |
| Cochrane                             | 5,377   | 624    | 334 |
| PLOS Biology                         | 2,149   | 243    | 212 |
| Health Technology Assessment         | 557     | 645    | 318 |
| Health Services and Delivery Response | 510    | 623    | 316 |
| Public Health Research               | 93      | 624    | 331 |
| Programme Grants for Applied Research | 78     | 722    | 311 |
| Efficacy and Mechanism Evaluation    | 70      | 659    | 341 |

Table 2: Journals included in CELLS. The average length (token level) of lay language summaries (Tgt) is shorter than that of scientific abstracts (Src).

data quality, we identified outliers using source-target lexical similarity and length. As a result, we excluded pairs from eLife, Annals of the Rheumatic Diseases, and Reproductive Health. This left a set of 62,886 source-target pairs from 12 journals.

### 3.1.2. Dataset applications

Using CELLS, we developed three evaluation tasks:

*Lay language generation.* For this task, we used the full-length scientific abstract and LLS pairs in CELLS for abstract-level lay language generation. As mentioned in Section 1, this task requires paragraph-level simplification, summarization, and background explanation to produce understandable summaries for laypeople. The following tasks focus on two of these challenges: abstract simplification and background explanation generation.

*Simplification.* Paragraph-level simplification fits the lay language generation task, but simplification is difficult to isolate because of the frequent insertion of background explanations. To focus on sentence-level simplification as a subtask, we developed a Greedy Paired Sentence Search (GPSS) algorithm (Algorithm 1) to align sentences from the abstracts and LLSs. The underlying idea is to identify matched source and target sentences based on lexical

overlap and sentence sequence. An example is provided in Figure 1. After applying GPSS, each source and target sentence was labeled as "matched" or "unmatched", resulting in a large set of 233,916 matching abstract- and LLS-derived sentence pairs for simplification.

---

**Algorithm 1** Greedy paired sentences search (GPSS) algorithm

---

   **Input:** SRC – the list of sentences in the source abstract, TGT – the list of sentences in the target abstract      **Output:**   P – the set including the indices of the paired source and target sentences

 1: **function** GPSS(src_start, src_end, tgt_start, tgt_end, score)
 2:     **if** src_start > src_end  **or** tgt_start > tgt_end **then**
 3:         **return** ∅
 4:     src_max,tgt_max ← argmax$_{i,j}$(score[i,j], src_start ≤ i < src_end, tgt_start ≤ j < tgt_end)
 5:      pairs ← {(src_max, tgt_max)}
 6:      pairs ← pairs ∪ GPSS(src_start, src_max-1, tgt_start, tgt_max-1, score)
 7:      pairs ← pairs ∪ GPSS(src_max+1, src_end, tgt_max+1, tgt_end, score)
 8:      **return** pairs
 9: $N_{src}$ ← the number of sentences in SRC
10: $N_{tgt}$ ← the number of sentences in TGT
11: **for** i ← 0 to $N_{src}$-1 **do**
12:     **for** j ← 0 to $N_{tgt}$-1 **do**
13:         S[i,j] ← ROUGE-L(TGT[j], SRC[i])
14: P ← GPSS(0, $N_{src}$, 0, $N_{tgt}$, S)

---

*Background explanation.* As mentioned in Section 1, adding explanations is a common strategy to enhance comprehension in 'Background' section. To support this research, we derived a large-scale paragraph-level dataset that emphasizes the insertion of novel content, i.e., background explanations. Focusing on explanation requires a reliable approach to extract sentences containing additional content in the LLS 'Background' section. Human annotation is reliable but costly, and exhaustive annotation of the 62,886 pairs in CELLS is infeasible. Therefore, we obtained paired source/target subparagraphs with the aforementioned GPSS algorithm. After applying GPSS,
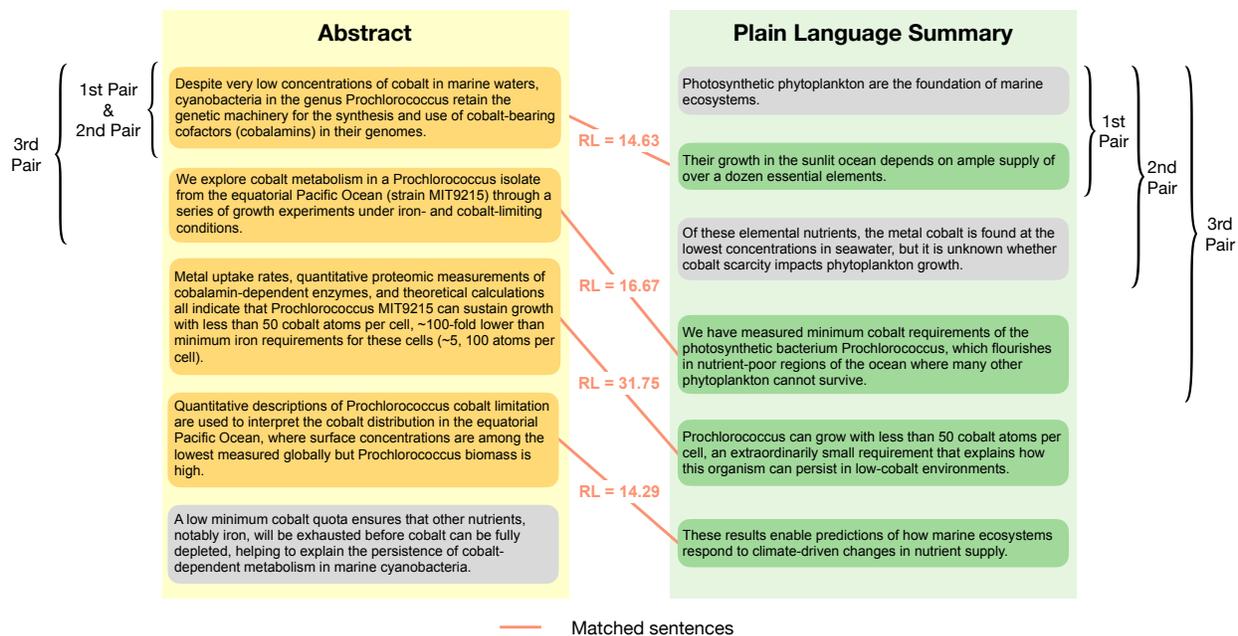
8

Figure 1: An example application of the GPSS algorithm. RL indicates the F1 score from ROUGE-L between the sentences in the abstract and plain language summary. For the background explanation subset, we combined unaligned target sentences (grey blocks) with proximal aligned sentences (green blocks). The example presented illustrates the generation of three paired examples ("pair") for the background explanation subset. All three pairs include the initial explanatory content that precedes the first matched sentence (RL = 14.63), as well as the sentence in the lay language summary that matches it. The second pair also includes the explanatory content after this matched sentences, and the third pair adds the following matched sentence also (i.e. the second sentence in the source abstract, and the lay language summary sentence that aligns with it). These combinations allow for the possibility that added content may relate to the preceding, or the subsequent sentence.

9

we considered the unaligned ("unmatched") sentences as putative *explanations*. We targeted the Background section, but section headers are unavailable for most abstracts. We therefore conducted human annotation to examine the utility of different empirically-defined boundaries, and the presence of external information. Fifty randomly selected abstracts from CELLS were annotated by two annotators: one medical student and one graduate student without medical training but with good familiarity with the dataset. Cohen's Kappa among the annotators is 0.74, indicating substantial agreement Artstein and Poesio (2008). Informed by the results of the annotation process, we selected the 2nd pair (refer to Figure 1) to demarcate the 'background' section, given its superior integration of background and external information. Further details can be found in Appendix A.1. Overall, we extracted 47,157 source/target pairs for background explanation [1].

*3.2. Human Validation of Dataset*

To ensure the robustness of the dataset used for background explanation and simplification, two expert annotators (same as above) assessed 250 paragraph pairs from the background explanation subset and 500 pairs from the simplification subset. The annotators were tasked with evaluating the pairs from the background explanation by: 1) confirming their presence within the background section; 2) identifying any external data not originally in the source; 3) classifying any external information as either definition, motivation, or example (detailed definitions of the categories can be found in Section 4.4); and 4) discerning whether the target and source information are aligned, where alignment is defined by the presence of common entities, i.e., at least one shared "triple" (A triple consists of three components: A subject, a predicate, and an object). Our annotators determined that 92.8% of pairs were situated in the background section and 62.8% included external information. Among the identified external information, 76.4% were motivations, 38.2% were definitions, and 10.8% were examples (these labels are not mutually exclusive). In addition, 87.2% of background explanation pairs were found to be in alignment. For the simplification subset, we focused on the alignment of paired sentences and found a 68.6% alignment. The chal-

---

[1] The term "background explanation" refers to specific sentences found within a Background section, but not every sentence in this section qualifies as a "background explanation." Instead, a background explanation is specifically defined as unmatched sentences serving explanatory purposes.
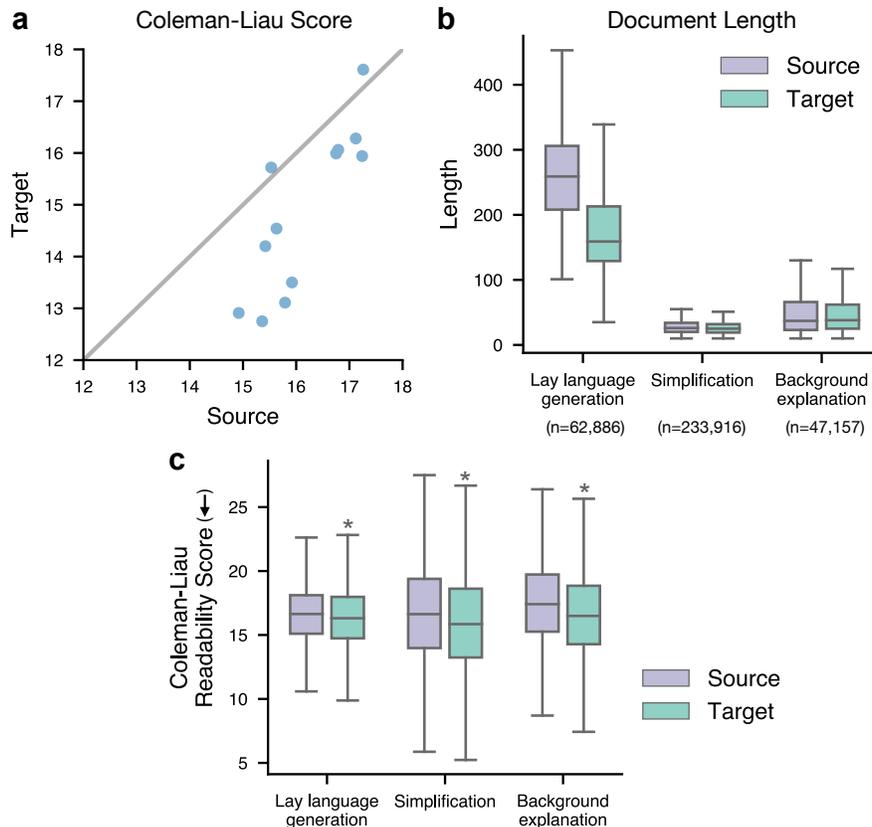
Figure 2: Dataset analysis. a, source and target Coleman-Liau readability scores for the 12 journals included in CELLS. Each dot represents one journal. Lower score indicates text is easier to read. b,c, Average length and Coleman-Liau readability score for source and target text for three tasks (i.e., lay language generation, simplification and background explanation). On average, target text is shorter and easier to read for all three tasks. "*" indicates that the score of the target significantly lower than that of the source with p-value < 0.05 (paired t-test).

lenge of sentence-level alignment in scientific summaries remains an active area of investigation Krishna et al. (2023), emphasizing the ongoing need for the advancement of alignment algorithms. Taking into account the inherent complexities of sentence alignment and external information detection, we consider the observed alignment and external rates to be within acceptable bounds for the purpose of the current work. However, these results also reinforce the need for further research on sentence alignment.

### 3.2.1. Dataset analysis

Dataset statistics are shown in Table 2. CELLS covers various topics including genetics, pathogens, neglected tropical diseases, computational biology, health services, and biomedical research. This diversity of topics and journals provides opportunities to study model generalizability. Background explanations notwithstanding, the average length of the source (professional language) is longer than the target (plain language summary) for each journal. The readability scores for each journal are shown in Figure 2a. The Coleman-Liau readability score indicates the estimated years of education required to understand a piece of text. Most of the average readability scores for the target are lower than those for the source, indicating that the target LLSs are generally easier to understand.

Figures 2b and 2c show lexical features of CELLS components for three tasks. On average, LLSs are shorter than corresponding scientific abstracts. Although the readability of both source and target texts is at the college level Karačić et al. (2019), the difference in readability between them is statistically significant (paired t-test), indicating LLSs are easier to understand than source text.

We randomly split the dataset into 45,280; 11,295; and 6,311 abstract/LLS pairs as the train, validation, and test sets respectively.

### 3.3. Methods

We investigated the performance of language models with intermediate pre-training (i.e. further pre-training on in-domain text) and retrieval-augmented lay language generation (RALL).

### 3.3.1. In-domain pre-training for simplification

Abstractive models are more applicable than extractive ones for our tasks since extractive summaries are written in the same professional language as their source documents. Therefore, we applied a state-of-the-art abstractive summarization model – *BART* Lewis et al. (2020) – to our tasks. BART uses hidden state representations of text sequences that are encoded bi-directionally (as is the case with BERT Devlin et al. (2018)) to inform a decoder model that predicts the next word in the sequence (as is the case with GPT-series and related models e.g. Brown et al. (2020)). During semi-supervised pre-training of the model, input sequences are perturbed with a range of transformations (for example, some tokens may be masked), and the model attempts to reconstruct the original sequence. As such, BART

has both the ability to encode hidden state representations that take an entire sequence into account, and a convenient mechanism to generate text in response to an input sequence. Once the model has been pre-trained on unlabeled text, it can be fine-tuned for particular tasks, such as summarization, by training it to generate target sequences in response to source sequences. We adopted a $\text{BART}_{\text{Large}}$ model that has been fine-tuned for general-domain text summarization on the widely-used CNN/DM summarization dataset (Cable News Network / Daily Mail Nallapati et al. (2016)) as our baseline. To be concise, we use *Vanilla* to denote the BART-Large-CNN model in the following text. Due to the complexity of our task and the relatively small size of our data, we employed intermediate pre-training - further semi-supervised pre-training on additional unlabeled in-domain text - to attempt to improve the performance of BART. Previous work shows that adaptive pre-training with domain-relevant unlabeled data can improve model performance Gururangan et al. (2020). Therefore, we further pre-trained the BART model on a corpus from the biomedical domain. We first perturbed 300K PubMed abstracts[2] by text substitution and sentence shuffling, and trained the BART model to reconstruct the original text. The pre-trained model was then further fine-tuned for the tasks of summarization, sentence simplification and background explanation, using our datasets.

### 3.3.2. Definition-based explanation retrieval

As the source documents in our dataset may omit required background knowledge, models should be able to retrieve relevant background knowledge from external sources. We evaluated two approaches to retrieving this knowledge.

The definition-based retrieval model uses a straightforward method to add explanations of terms to the text, by identifying definitions of terms that exist in a predefined lexicon. In our experiments, we used datasets derived from the UMLS Bodenreider (2004) and Wikipedia to augment the context of the source document. The UMLS includes medical term (entity $e$) and definition ($d$) pairs $\mathcal{D}_u = \{(e_i, d_i)\}$. For each source document $s$, we used Scispacy Neumann et al. (2019) to identify the expression of normalized UMLS terms $e_{u_1}, e_{u_2}, ..., e_{u_m}$ in $s$. Then we added corresponding UMLS term definitions $d_{u_1}, d_{u_2}, ..., d_{u_m}$ to $s$ to obtain $\hat{s}$. The Wikipedia dataset includes keyword

---

[2]https://www.kaggle.com/cvltmao/pmc-articles

($e$) and definition ($d$) pairs $\mathcal{D}_w = \{(e_i, d_i)\}$. For each source document $s$, we applied KeyBERT[3], which uses BERT embeddings and cosine similarity to find the sub-phrases in a document most similar to the document itself, to identify three keywords $(e_{w_1}, e_{w_2}, e_{w_3}) = \text{KeyBERT}(s)$. We obtained the definitions of those keywords, $d_{w_1} d_{w_2}, d_{w_3}$, from the Wikipedia dataset and added them to the end of the source document, $s$, to obtain $\hat{s}$. Lastly, we fine-tuned the BART model using $\hat{s}$.

### 3.3.3. Embedding-based explanation retrieval

For the embedding-based retrieval method, we adopted a state-of-the-art dense retrieval model to augment the source with related documents from an external set. Specifically, we applied the retrieval-augmented generation (RAG) model Lewis et al. (2020) using another Wikipedia-derived dataset $Z_w$ of 21M 100-word documents $z$. The retrieval component $p_\phi(z|s) \propto \exp(\mathbf{d}(\text{z})^\text{T} \mathbf{q}(\text{s}))$ is based on the Dense Passage Retriever (DPR) Karpukhin et al. (2020), where $\mathbf{d}(z) = \text{BERT}_d(z)$ and $\mathbf{q}(s) = \text{BERT}_q(s)$ are the representations of the documents (in our case, the source document to be summarized and the Wikipedia documents that provide candidates for retrieval) produced by two $\text{BERT}_\text{Base}$ encoders. The DPR model retrieves the top $k$ documents $z$ with the highest prior probability $p_\phi(z|s)$ using the Maximum Inner Product Search method Johnson et al. (2019). After applying DPR, we concatenated the source $s$ and the retrieved content $z$ as the input. We used the RAG-Sequence model whose generator produces the output sequence probabilities for each concatenated document:

$$p(t|s) \approx \sum_{z \in top\text{-}k(p_\phi(\cdot|s))} p_\phi(z|s) p_\theta(t|s, z).$$

$p_\theta(t|s, z)$ is the generator, and we used BART for this purpose. As can be seen from the formula, both the content of the retrieved documents ($z$) and their probabilities of retrieval ($p_\phi(z|s)$) inform the generated text. During training, we fixed $\text{BERT}_d(\cdot)$, and only fine-tuned $\text{BERT}_q(\cdot)$ and the BART generator.

### 3.4. LLMs

To evaluate the performance of LLMs in generating background explanations or plain language summaries, we utilized Llama 2 Touvron et al. (2023)

---

[3]https://github.com/MaartenGr/KeyBERT

and GPT-4 OpenAI (2023). We explored two prompts: 1) `"Summarize in plain language: input"`, and 2) `"Summarize in plain language, providing necessary explanations: input"`. To further assess the impact of the retrieval-augmented approach on LLMs, we established two settings for input: the source alone and the source combined with Wikipedia definitions as identified using KeyBERT.

### 3.5. Experiments

### 3.5.1. Experimental setup

All experiments except LLMs were run using a single NVIDIA Tesla V-100 GPU. Models were developed using PyTorch Paszke et al. (2019). We used the Fairseq[4] BART implementation, and the HuggingFace Transformers Library Wolf et al. (2019) to implement the RAG model. For RAG, we retrieved the top 5 documents for each input. The maximum length of generated texts was set to 700 for paragraph-level lay language generation and 150 for background explanation and sentence-level simplification. Other hyper-parameters were set to their default values.

We used the `Llama-2-70B-chat`[5] model for Llama 2 [6] and GPT-4 [7] for the GPT model. The generation process was configured with a maximum length of 150 tokens. All other parameters were set to their default values.

### 3.5.2. Evaluation

*Automated evaluation.* We first evaluated generation quality using ROUGE-L Lin (2004), BERTScore Zhang* et al. (2020), BLEU Papineni et al. (2002), and METEOR Banerjee and Lavie (2005)[8] to compare generated text to professionally-authored plain language target text. ROUGE-L depends on $n$-gram overlap, while BERTScore uses the similarity between embeddings and as such may be less sensitive to differences in word choice between human-authored and automatically-generated LLS. BLEU computes $n$-gram preci-

---

[4]https://github.com/pytorch/fairseq

[5]Model: https://ai.meta.com/llama/

[6]Implementation: The model was quantized to 4 bits using OPTQ as implemented in https://github.com/PanQiWei/AutoGPTQ, and hosted on a local server using https://github.com/turboderp/exllama for inference

[7]Implementation: https://openai.com

[8]Implementation: Fabbri et al. (2020) BERTScore hash code: `bert-base-uncased_L8_no-idf_version = 0.3.12(hug_trans=4.27.3)`

sion of generated text against target texts, including a brevity penalty. ME-TEOR employs a relaxed matching criterion based on the F-measure, and addresses the exact match restrictions and recall consideration of BLEU.[9] The Coleman-Liau readability score Coleman and Liau (1975) assesses the ease with which a reader can understand a passage, and word familiarity Leroy and Kauchak (2014) measures the inverse document frequency of uni-grams in text using frequency counts from Wikipedia. *Lower* Coleman-Liau score and word familiarity indicate that text is *easier* to read.[10]

To directly evaluate how representative of an LLS the generated text is, we trained a RoBERTa Liu et al. (2019) model to classify the source of sentences from the original abstracts and LLSs. Specifically, we used the paired source-target sentences from the GPSS algorithm with a sentence length between 10 and 150 words. The input to the RoBERTa model is a sentence and the label is 0 for a source sentence (from a scientific abstract) and 1 for a target sentence (from a LLS). The RoBERTa model achieved an AUROC of 0.83 and an F1 score of 0.74 on the held-out test set, demonstrating that there are detectable and generalizable differences in language use by the intended audience. As the model is trained to output a higher prediction for a target sentence (i.e., a LLS sentence), we used the prediction of the model to evaluate how "plain" the input text is. We refer to the predicted probability of this model as the "Plainness Score". A *higher* Plainness Score indicates that the text is more representative of an LLS.

*Human evaluation.* We set up our human evaluation similarly to Guo et al. (2021), providing pairs of source and summary text to the human evaluators, where the summary is either expert-written or generated by one of our two best-performing BART models and two GPT-4 models (evaluators were not informed which summaries were human-authored). We asked human evalua-tors to rate the summary for grammatical correctness, meaning preservation, understandability, factual correctness, and the relevance of external informa-tion, each on a 1-4 point Likert scale (1-very poor, 4-very good). Questions can be found in Appendix A.2. The study was considered exempt upon institutional IRB review. Twelve evaluators were recruited using an institu-tional NLP interest group channel. Each of them annotated four examples

---

[9]Please see BLEU and METEOR scores in Appendix.

[10]The familiarity measure is derived from *inverse* document frequency which is higher for rare terms, so *lower* familiarity scores indicate the use of *more familiar* words.

from the test set for background explanation. Three evaluators reviewed each example. All the evaluators have at least an undergraduate degree, lack specialized biomedical training, and half speak English at home. The Krippendorff's alpha coefficient Krippendorff (1970) was 0.40 for the four background explanation texts among all evaluators. Krippendorff's alpha coefficient measures multiple inter-rater agreements in ordinal data, and values range from 0 to 1. Considering the subjectivity of the task and the multiple choices per question, we considered this level of agreement to be acceptable.

## 4. Results

We experimented with five models using BART: the base (*Vanilla*) model, *Vanilla* further pre-trained on PubMed abstracts (*PubMed pre-trained*), *Vanilla* with UMLS (*UMLS definition-based retrieval*) and Wikipedia (*Wiki definition-based retrieval*) definition-based retrieval, and *Vanilla* with embedding-based retrieval using Wikipedia (*Wiki embedding-based retrieval*). Additionally, we experimented with three prompts using LLMs for background explanation: "Summarize in plain language: input" (*summary*), "Summarize in plain language, providing necessary explanations: input (*explain*), and "Summarize in plain language: input with Wiki definition-based retrieval (*wiki*).

### 4.1. RALL improves generation performance

We first evaluated the text generation performance of our models (Figure 3), using ROUGE-L and BERTScore to compare generated LLS text to the corresponding human-authored lay language text (the target) for a given abstract or sentence (the source). Due to the input length limitation of the BERT (512 tokens) and BART (1024 tokens) models, we did not perform retrieval-augmented generation for the abstract-level lay language generation task (Figure 3, 1st panel). However, for this task, pre-training on unlabeled data was not helpful.

We next compared text generation performance on the sentence-level simplification task (Figure 3, 2nd panel). Results indicate that the PubMed pre-trained model achieved better performance than the Vanilla model, suggesting that pre-training on domain-specific unlabeled data is helpful for sentence-level simplification, which aligns with the results from our prior work Guo et al. (2021). It can also be observed that the models with information retrieval from Wikipedia (Wiki definition- and embedding-based
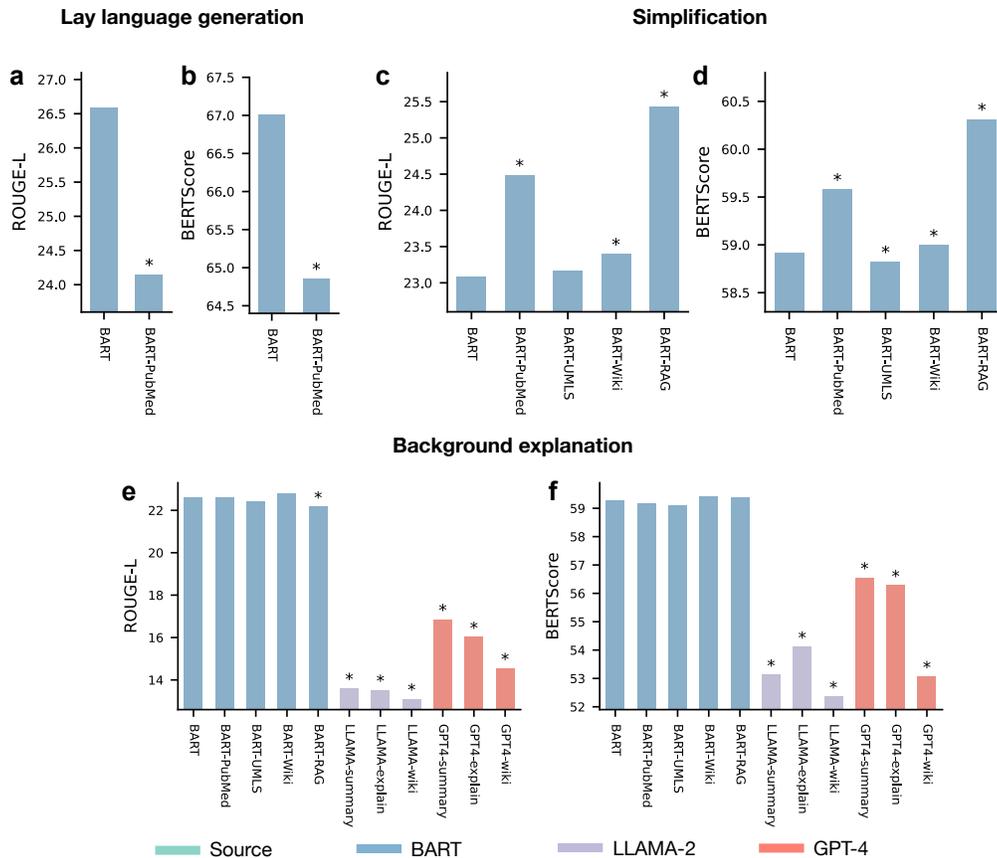
**Lay language generation**

**Simplification**

**Background explanation**

Figure 3: Models' performance in text generation. We used the F1 score of ROUGE-L and BERTScore to evaluate the generation quality of models on lay language generation, simplification, and background explanation tasks. P-values obtained through the t-test are employed to evaluate the performance of various models compared to the Vanilla model (BART). * indicates statistical significance with Bonferroni-Holm correction for multiple hypothesis testing Holm (1979).

retrieval) achieve higher ROUGE-L scores than the Vanilla model, suggesting that retrieving this external information may also be helpful for text simplification tasks. One reason for this may be that Wikipedia articles target a broader audience than the intended audience of specialized biomedical literature, and are therefore written to be more accessible.

For background explanation (Figure 3, 3rd panel), the PubMed pretrained model only shows marginal improvements, suggesting that pre-training on domain-specific unlabeled data is helpful but insufficient for this challeng-

18

ing task, perhaps because the authors of biomedical literature assume expert knowledge on the part of their readers and therefore seldom include the explanatory content that a non-expert reader might require. Furthermore, the BART models with retrieval from the Wikipedia dataset (Wiki definition- and embedding-based retrieval) achieved higher BERTScore, establishing the benefits of information retrieval techniques for background explanation with BART. ROUGE-L results show a smaller advantage for Wiki definition- based RALL generation, and unlike with BERTScore this advantage is not apparent for embedding-based RALL methods. With auto-regressive LLMs ROUGE and BERTscore results are remarkably consistent: GPT-4 outper- forms Llama 2 with all three prompts; however, both models are notably out- performed by BART-based architectures. Our results suggest that prompting exclusively for summarization yields superior outcomes compared to combin- ing summarization with explanation. The weakest performance is observed when using the Wiki definition-based retrieval source combined with sum- marization prompting. This disparity could be attributed to our reliance on zero-shot LLMs, whereas BART benefits from fine-tuning. This suggests there may be avenues for improvement, such as exploring few-shot learn- ing approaches within LLMs for background explanation, or using low-rank adaptation techniques to further improve auto-regressive LLM performance. To offer a comprehensive view of performance, BLEU and METEOR scores from the test set are also presented in Appendix Figure A.1. The observed BLEU patterns are consistent with those for ROUGE and BERTScore, while METEOR highlights advantages for 'explain' LLM queries and BART-RAG.

In acknowledgement of the limitations of our GPSS algorithm, where the background explanation doesn't always pair paragraphs with external content and the simplification subset doesn't always produce aligned pairs, we provide results from the 'gold' annotated dataset. These are available in Appendix Figure A.2 and Appendix Figure A.3. This 'gold' background explanation subset features pairs in alignment with external content, while the 'gold' simplification subset showcases aligned pairs. The results derived from these 'gold' validated subsets are consistent with those from the test set, indicating that the observed patterns are not attributable to errors in algorithmic alignment.

Overall, our results suggest that both biomedical domain pre-training and information retrieval are helpful for background explanation and sentence- level simplification. Furthermore, the information retrieval models based on BART using Wikipedia produced text that was most similar to human-

authored lay language. This indicates that general domain information written for a broader audience (e.g., Wikipedia) is a good resource for background explanation generation using BART.

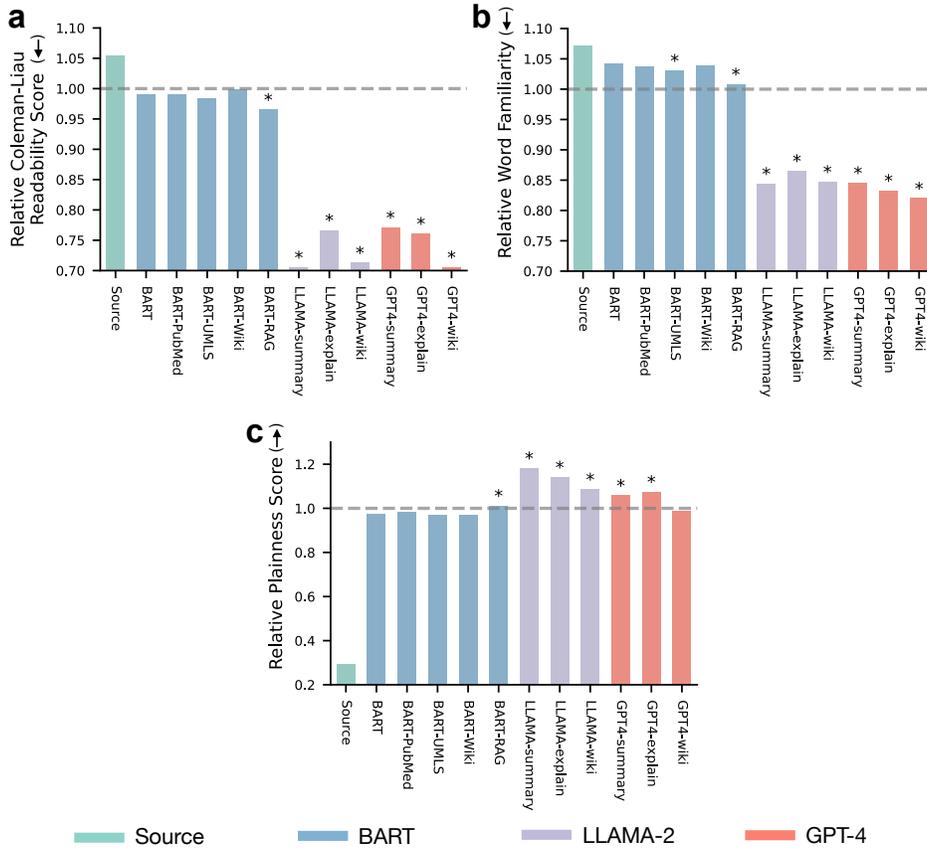## 4.2. RALL improves text interpretability



Figure 4: Readability, familiarity and plainness of the background explanation subset, relative to professionally-authored lay language text. (a) Relative Coleman-Liau readability score, (b) word familiarity, and (c) plainness score of the source and models' generated text. The relative score is calculated by dividing by the score of the target text. A lower readability score and word familiarity indicate that the text is easier to read (values below the dashed line are lower than those from professionally-authored plain language). A higher Plainness Score indicates that the text is more representative of an LLS. P-values obtained through the t-test are employed to evaluate the performance of various models compared to the Vanilla model (BART). * indicates statistical significance with Bonferroni-Holm correction for multiple hypothesis testing Holm (1979).

We next evaluated the interpretability of the generated text using the data from the background explanation subset (Figure 4). Existing text interpretability metrics consistently return better scores for the models' outputs than for the source text. Of note, the Coleman-Liau readability scores of the models' outputs are even lower than those of the target text (Figure 4a.) This indicates that our datasets help the BART model to generate more straightforward and readily interpretable text. We also found that the retrieval-augmented BART models performed well in this interpretability evaluation, suggesting that the UMLS and Wikipedia datasets may be easier to understand than professional-language abstracts. Overall, the embedding-based RALL model, which used Wikipedia as a source, had the best readability, familiarity, and plainness scores. These results further support the utility of retrieval augmentation for lay language generation, suggesting it can benefit the style of generated text, as well as its content. For LLMs, the model outputs consistently score well across all metrics. Outputs generated with the Wiki-based definition retrieval source show improved results in the readability score and relative word familiarity compared to those without it. However, this advantage doesn't extend to the plainness score.

### 4.3. Human evaluation

Figure 5 shows the human rating scores across four pairs of target and generated texts. Evaluators generally rated generated background explanations higher than those from the expert-generated LLS. It is interesting to note that the Wiki definition-based retrieval BART model was judged to have the least relevant external information but the best understandability, according to raters. Exploring the tradeoff between the amount of external information and understandability, and jointly optimizing them presents a challenging direction for future work. These results confirm the effectiveness of our dataset for improving automated models' ability to generate LLS with relevant external information added. For GPT-4, when prompted with summarization and explanation, yields the highest scores in understandability, meaning preservation, and information correctness. However, it falls short in incorporating relevant external information. This suggests that GPT-4's explanatory outputs should be meticulously vetted to prevent potential misalignment with the topic. In contrast, the embedding-based Wiki approach (BART-RAG) excels in meaning preservation, maintaining the accuracy of key information, and integrating relevant external information. However, its outputs can be challenging to comprehend. This raises the potential of syn-
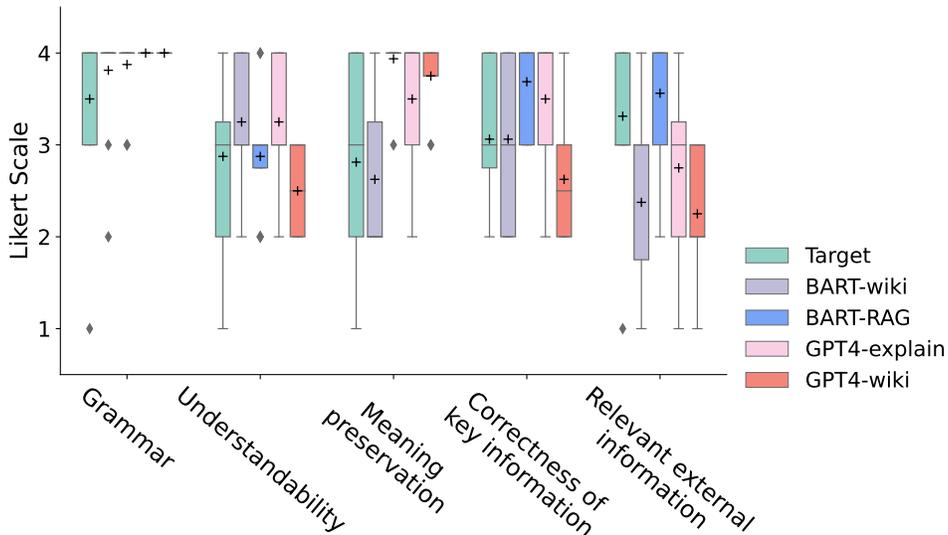
Figure 5: Human evaluation results for four generated texts from the background explanation task. Each text was assigned to four raters. For the Likert scale, 1 is very poor, and 4 is very good. "+" indicates the mean.

ergizing the embedding-based method with LLMs to achieve the ideal lay language summary.

## 4.4. Selected examples

To define the scope of background explanation, we identified three types of explanation in the dataset, as shown in Table 3. We did not aim to enumerate all possible categories. Rather, our goal was to provide some initial insights into explanation phenomena.

The most common explanation type we observed is a *definition*, including "common" medical words, technical terms, and abbreviations, to avoid misunderstanding. *Motivation*, including prevalence, risk factors, history, etc., sustains readers' interest and establishes whether the topic under discussion meets their information needs. Providing a *concrete example* allows readers to link an otherwise obscure concept with a more familiar one. For example, connecting the increasing temperature in the ocean with coral reefs makes it easy for the reader to understand the importance of the study.

We present background explanations for the two best BART models and two GPT-4 models in Table 3. This provides evidence that the retrieval-augmented model can generate both term definitions and motivations for

22

the main topic. However, the generated external content may not be aligned with the target (e.g. Marburg virus does not cause Ebola virus disease), highlighting the importance of improving the relevance and correctness of generated abstractive summaries as an area for future research. In addition, the models appear unable to produce illustrative examples. This ability goes beyond retrieving evidence, and appears difficult to learn.

## 5. Discussion

We have two key observations from the model development and evaluation. Regarding BART models, RALL variants outperform the vanilla model, though these improvements are larger with the sentence simplification subset than when background explanation is emphasized. Background explanation generation is challenging and requires considering both the knowledge sources from which information is drawn, and the understandability of generated text. Examples suggest that the models can add term definitions and relevant epidemiological data but fail to provide illustrative examples for related concepts. These abilities may be beyond current models' capabilities and fall outside the scope of the resources used in our study for information retrieval. Therefore, generating high-quality explanations may require acquiring other knowledge resources, or decomposing the background explanation task into more granular subtasks. Another key observation is that human ratings are essential to assessment of background explanation task performance. Although we included automated evaluation metrics for generation quality and text simplicity, they cannot capture explanation quality. Rater evaluations of the external content for existence, relevance, and correctness of background explanations suggest additional advantages for RALL models that are opaque to automated evaluation methods.

To the best of our knowledge, CELLS is the largest lay language generation dataset developed to date, and the derived dataset for background explanation serves as the first explanation generation benchmark. We envision these data broadly applying to biomedically-related applications, and other NLP methods. On the biomedical applications side, we provide a benchmark to develop new NLP tools to generate LLSs for scientific literature. With the assistance of such NLP tools, researchers can write more understandable text, allowing healthcare consumers to interpret and apply the information it contains to guide their health-related decision-making. From an NLP methodological perspective, these datasets offer an excellent oppor-

23

|  | **Definition** | **Motivation** | **Example** |
|---|---|---|---|
| **Source** | Marburg virus (MARV) disease is lethal, with fatality rates up to 90%… | Clinical reports of Zika Virus (ZIKV) RNA detection in breast milk have been described, but evidence conflicts as to whether this RNA represents infectious virus… | …Very little is known about how environmental changes such as increasing temperatu… affect disease dynamics in the ocean, especially at large spatial scales… |
| **Target** | Marburg virus ( MARV) is a member of the Filoviridae family that causes severe hemorrhagic fevers with high fatality rates in humans… | Only 4 years have passed since the Zika virus outbreak in Brazil, and much remains to be understood about the transmission and health consequences of Zika infection. To date, some case reports have detected Zika virus RNA in the breast milk… | …This hypothesis is supporte… by local observations--for example, that some coral diseases become more prevalent in the summertime--but it has never been tested a… large spatial scales or over relatively long periods…. |
| **BART-Wiki** | Marburg virus ( MARV) is a highly lethal pathogen that is transmitted by the bite of infected mosquitoes… | Zika Virus ( ZIKV) is a mosquito-borne virus that has been circulating in tropical and subtropical regions of the world since 2015, and it has been linked to an outbreak of microcephaly and Guillain-Barre syndrome ( GBS) in Brazil in 2015. ZIKV RNA detection in breast milk… | …Disease dynamics in the ocean are influenced by environmental changes such… increasing temperature, but little is known about how thes… changes affect disease dynamics over large spatial scales… |
| **BART-RAG** | Marburg virus ( MARV) is the causative agent of Ebola virus disease, which is fatal in up to 90% of cases… | Zika virus ( ZIKV) is a mosquito-borne flavivirus that has caused a worldwide health crisis since its discovery in the Americas in 2015. Clinical reports of Zika virus RNA detection in human breast milk… | …Disease dynamics in the ocean are affected by environmental changes such… temperature, but very little is known about how these changes affect disease dynamics at large scales… |
| **GPT4-Wiki** | The Marburg virus is a deadly disease with a high fatality rate. | The Zika virus, which gained worldwide attention due to an outbreak in the Americas from 2015 to 2016, affects a protein-coding gene in a species called Drosophila melanogaster… | We don't know much about how changes in the environment, like rising temperatures, impact the spread and development of diseases in the ocean, particularly on a large scale. |
| **GPT4-explain** | The virus called Marburg Virus Disease is extremely dangerous and can kill up to 90% of the people it infects. | The Zika virus has gained worldwide attention over the last five years mainly because of its comeback in the Americas from 2015 to 2016… | We don't know a lot about ho… changes in the environment, like rising temperatures, impa… the spread and behavior of diseases in the ocean. This is… especially true when we're looking at larger areas. |

Table 3: Typical types of background explanation from scientific abstracts to lay language summaries, and the corresponding generated text using two best-performing models. Our retrieval-augmented models can generate both term definitions and related epidemiological data for the main topic, but fail to link examples to related concepts.

tunity to develop and evaluate novel LLS generation techniques. CELLS can also support sentence-level and paragraph-level simplification research, and with additional annotation could provide a basis for open-question answering and informational retrieval tasks.

While we evaluated the correctness of key information in human evaluation, it remains difficult for non-experts to identify the factuality or external information relevance. An improved model with factuality enforcement could promote sequences with higher accuracy. Medical experts (i.e. medical students) could be recruited for evaluation of factual correctness. More abstracts and human raters are required to confirm the apparent appeal of LLS from retrieval-augmented text generation models. Furthermore, to improve the quality of the dataset for background explanation, larger-scale verification is needed. We also note that our strategies for adding entity-driven explanations are straightforward, and that we did not perform a hyperparameter search to optimize the relatively expensive dense retrieval procedure, on account of resource constraints.

Evaluating lay language generation inherently poses challenges due to the multifaceted nature of the task, including aspects such as incorporating background explanations and omitting technical terms. While the ROUGE (Lin, 2004) and BLEU (Papineni et al., 2002) metrics are conventionally applied to evaluate lay language generation, their applicability is constrained due to inherent limitations associated with their reliance on lexical overlap, and the need for high-quality reference summaries. Moreover, these metrics are not adept at detecting hallucinations, a critical consideration, especially in the healthcare domain where the accuracy of lay language plays a pivotal role in informing health decisions (Wallace et al., 2021; Pagnoni et al., 2021). A recent investigation indicated that ROUGE, BLEU, METEOR, and BERTScore face challenges in capturing text simplification precisely Guo et al. (2023). Additionally, Mac et al. found that automated readability scores frequently display inconsistency and lack accuracy Mac et al. (2022). While human evaluations provide comprehensive feedback, they are resource-intensive, making them challenging to scale to extensive datasets. Therefore, a metric tailored specifically for lay language generation is much needed, and one should exercise caution when interpreting results using existing measures.

The GPSS algorithm searches for external content by calculating the lexical similarity between the source and target sentences using the ROUGE score. However, this may fail to recognize alignment at the semantic level

25

when meanings but not terms overlap. To address this challenge, end-to-end approaches that can learn embedding-derived similarities from the source and target and classify the external content accordingly may be worth exploring. Since language models pre-trained in the medical domain have achieved state-of-the-art performance on several biomedical NLP tasks, exploring the benefits of these models is an important direction for future work. Regarding lay language generation, one remaining challenge that is a possible direction for future work involves directly applying retrieval-augmented methods to full-length abstracts instead of the background sections. Also, it would be intriguing yet challenging to generate LLS for different education levels. One potential solution may be incorporating a reward function that responds to readability, interpretability, or plainness metrics.

LLMs show promise in the realm of lay language generation. While the outputs from LLMs may not align closely with the target, the produced text is notably easy to comprehend. This ability to simplify addresses a key challenge in the existing lay language generation datasets, which typically offer only a single target. This suggests the potential to develop a pipeline that first broadens the source and then tailors the content for varying levels of readability. Moreover, the less-than-optimal performance of LLMs underscores the potential of exploring few-shot learning further. Finally, incorporating Wiki-definitions could be problematic for LLMs operating in a zero-shot learning mode. For those wishing to employ retrieval-augmented methods with LLMs, a more judicious selection of external resources or a thorough vetting of the incorporated resources is imperative.

## 6. Conclusion

To improve the interpretability of lay language text generated by neural language models, we applied state-of-the-art text generation models augmented with retrieval components and achieved promising quality and readability scores as compared with reference lay language summaries generated by human experts. Results from human evaluation support the benefits of retrieval-augmented lay language generation for the generation of background explanations in particular. The new dataset and results provide a foundation for advancement in the challenging area of automated background explanation generation and lay language generation, with the potential to mediate clearer communication of biomedical science for better informed health-related decision making.

## 7. Acknowledgments

## References

S. H. Soroya, A. Farooq, K. Mahmood, J. Isoaho, S.-e. Zara, From information seeking to information avoidance: Understanding the health information behavior during a global health crisis, Information processing & management 58 (2021) 102440.

S. Bin Naeem, M. N. Kamel Boulos, Covid-19 misinformation online and health literacy: a brief overview, International journal of environmental research and public health 18 (2021) 8091.

B. M. Korsch, E. K. Gozzi, V. Francis, Gaps in doctor-patient communication: I. doctor-patient interaction and patient satisfaction, Pediatrics 42 (1968) 855–871.

E. T. Kurtzman, J. Greene, Effective presentation of health care performance information for consumer decision making: a systematic review, Patient education and counseling 99 (2016) 36–43.

S. A. Crossley, H. S. Yang, D. S. McNamara, What's so simple about simplified texts? a computational and psycholinguistic investigation of text comprehension and text processing., Reading in a Foreign Language 26 (2014) 92–113.

Y. Guo, W. Qiu, Y. Wang, T. Cohen, Automated lay language summarization of biomedical scientific reviews, in: Proceedings of the AAAI Conference on Artificial Intelligence, volume 35, 2021, pp. 160–168.

A. Devaraj, I. Marshall, B. C. Wallace, J. J. Li, Paragraph-level simplification of medical texts, in: Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, 2021, pp. 4972–4984.

D. S. McNamara, E. Kintsch, N. B. Songer, W. Kintsch, Are good texts always better? interactions of text coherence, background knowledge, and levels of understanding in learning from text, Cognition and instruction 14 (1996) 1–43.

S. Jonnalagadda, L. Tari, J. Hakenberg, C. Baral, G. Gonzalez, Towards effective sentence simplification for automatic processing of biomedical text,

in: Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics, Companion Volume: Short Papers, Association for Computational Linguistics, Boulder, Colorado, 2009, pp. 177–180. URL: https://aclanthology.org/N09-2045.

B. Qenam, T. Y. Kim, M. J. Carroll, M. Hogarth, et al., Text simplification using consumer health vocabulary to generate patient-centered radiology reporting: translation and evaluation, Journal of medical Internet research 19 (2017) e8536.

O. Bodenreider, The unified medical language system (umls): integrating biomedical terminology, Nucleic acids research 32 (2004) D267–D270.

P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, T. Rocktäschel, et al., Retrieval-augmented generation for knowledge-intensive nlp tasks, Advances in Neural Information Processing Systems 33 (2020) 9459–9474.

H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, et al., Llama 2: Open foundation and fine-tuned chat models, arXiv preprint arXiv:2307.09288 (2023).

OpenAI, Gpt-4 technical report, 2023. arXiv:2303.08774.

A. Cohan, F. Dernoncourt, D. S. Kim, T. Bui, S. Kim, W. Chang, N. Goharian, A discourse-aware attention model for abstractive summarization of long documents, in: Proceedings of NAACL-HLT, 2018, pp. 615–621.

I. Cachola, K. Lo, A. Cohan, D. Weld, Tldr: Extreme summarization of scientific documents, Findings of EMNLP (2020).

A. Devaraj, W. Sheffield, B. C. Wallace, J. J. Li, Evaluating factuality in text simplification, in: Proceedings of the conference. Association for Computational Linguistics. Meeting, volume 2022, NIH Public Access, 2022, p. 7331.

R. Mishra, J. Bian, M. Fiszman, C. R. Weir, S. Jonnalagadda, J. Mostafa, G. Del Fiol, Text summarization in the biomedical domain: a systematic

review of recent research, Journal of biomedical informatics 52 (2014) 457–467.

D. D. A. Bui, G. Del Fiol, J. F. Hurdle, S. Jonnalagadda, Extractive text summarization system to aid data extraction from full text in systematic review development, Journal of biomedical informatics 64 (2016) 265–272.

A. Givchi, R. Ramezani, A. Baraani, Graph-based abstractive biomedical text summarization, Journal of Biomedical Informatics (2022) 104099.

X. Cai, S. Liu, J. Han, L. Yang, Z. Liu, T. Liu, Chestxraybert: A pretrained language model for chest radiology report summarization, IEEE Transactions on Multimedia (2021).

Y. Zhang, D. Y. Ding, T. Qian, C. D. Manning, C. P. Langlotz, Learning to summarize radiology findings, in: Proceedings of the Ninth International Workshop on Health Text Mining and Information Analysis, 2018, pp. 204–213.

Y. Zhang, D. Merck, E. Tsai, C. D. Manning, C. Langlotz, Optimizing the factual correctness of a summary: A study of summarizing radiology reports, in: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, 2020, pp. 5108–5120.

M. Wang, M. Wang, F. Yu, Y. Yang, J. Walker, J. Mostafa, A systematic review of automatic text summarization for biomedical literature and ehrs, Journal of the American Medical Informatics Association (2021).

L. Plaza, Comparing different knowledge sources for the automatic summarization of biomedical literature, Journal of biomedical informatics 52 (2014) 319–328.

X. Cai, S. Liu, L. Yang, Y. Lu, J. Zhao, D. Shen, T. Liu, Covidsum: A linguistically enriched scibert-based summarization model for covid-19 scientific papers, Journal of Biomedical Informatics 127 (2022) 103999.

B. Chintagunta, N. Katariya, X. Amatriain, A. Kannan, Medically aware gpt-3 as a data generator for medical dialogue summarization, in: Proceedings of the Second Workshop on Natural Language Processing for Medical Conversations, 2021, pp. 66–76.

A. Joshi, N. Katariya, X. Amatriain, A. Kannan, Dr. summarize: Global summarization of medical dialogue by exploiting local structures., in: Findings of the Association for Computational Linguistics: EMNLP 2020, 2020, pp. 3755–3763.

S. Jonnalagadda, L. Tari, J. Hakenberg, C. Baral, G. Gonzalez, Towards effective sentence simplification for automatic processing of biomedical text, in: Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics, Companion Volume: Short Papers, 2009, pp. 177–180.

J. Li, C. Lester, X. Zhao, Y. Ding, Y. Jiang, V. V. Vydiswaran, Pharmmt: A neural machine translation approach to simplify prescription directions, in: Findings of the Association for Computational Linguistics: EMNLP 2020, 2020, pp. 2785–2796.

Y. Cao, R. Shui, L. Pan, M.-Y. Kan, Z. Liu, T.-S. Chua, Expertise style transfer: A new task towards better communication between experts and laymen, in: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, 2020, pp. 1061–1071.

J. Lu, J. Li, B. C. Wallace, Y. He, G. Pergola, Napss: Paragraph-level medical text simplification via narrative prompting and sentence-matching summarization, arXiv preprint arXiv:2302.05574 (2023).

M. K. Chandrasekaran, G. Feigenblat, E. Hovy, A. Ravichander, M. Shmueli-Scheuer, A. de Waard, Overview and insights from the shared tasks at scholarly document processing 2020: Cl-scisumm, laysumm and long-summ, in: Proceedings of the First Workshop on Scholarly Document Processing, 2020, pp. 214–224.

T. Goldsack, Z. Zhang, C. Lin, C. Scarton, Making science simple: corpora for the lay summarisation of scientific literature, arXiv preprint arXiv:2210.09932 (2022).

Z. Luo, Q. Xie, S. Ananiadou, Readability controllable biomedical document summarization, arXiv preprint arXiv:2210.04705 (2022).

K. Attal, B. Ondov, D. Demner-Fushman, A dataset for plain language adaptation of biomedical abstracts, Scientific Data 10 (2023) 8.

D. Das, A. Martins, A Survey on Automatic Text Summarization, Technical Report, Carnegie Mellon University, 2007. `https://www.cs.cmu.edu/~afm/Home_files/Das_Martins_survey_summarization.pdf`.

G. Erkan, D. R. Radev, Lexrank: Graph-based centrality as salience in text summarization, Journal of Artificial Intelligence Research (2004).

J. Cheng, M. Lapata, Neural summarization by extracting sentences and words, arXiv preprint arXiv:1603.07252 (2016).

S. Gupta, S. K. Gupta, Abstractive summarization: An overview of the state of the art, Expert Systems with Applications 121 (2019) 49–65.

L. Zhang, R. Negrinho, A. Ghosh, V. Jagannathan, H. R. Hassanzadeh, T. Schaaf, M. R. Gormley, Leveraging pretrained models for automatic summarization of doctor-patient conversations, arXiv preprint arXiv:2109.12174 (2021).

D. Yadav, J. Desai, A. K. Yadav, Automatic text summarization methods: A comprehensive review, arXiv preprint arXiv:2204.01849 (2022).

B. C. Wallace, S. Saha, F. Soboczenski, I. J. Marshall, Generating (factual?) narrative summaries of rcts: Experiments with neural multi-document summarization, AMIA Summits on Translational Science Proceedings 2021 (2021) 605.

T. Goldsack, Z. Luo, Q. Xie, C. Scarton, M. Shardlow, S. Ananiadou, C. Lin, Overview of the biolaysumm 2023 shared task on lay summarization of biomedical research articles, in: Proceedings of the 22st Workshop on Biomedical Language Processing, Toronto, Canada. Association for Computational Linguistics, 2023.

N. Srikanth, J. J. Li, Elaborative simplification: Content addition and explanation generation in text simplification, arXiv preprint arXiv:2010.10035 (2020).

B. K. Britton, S. Gülgöz, Using kintsch's computational model to improve instructional text: Effects of repairing inference calls on recall and cognitive structures., Journal of educational Psychology 83 (1991) 329.

M. S. Simpson, E. M. Voorhees, W. Hersh, Overview of the trec 2014 clinical decision support track, Technical Report, LISTER HILL NATIONAL CENTER FOR BIOMEDICAL COMMUNICATIONS BETHESDA MD, 2014.

K. Roberts, M. S. Simpson, E. M. Voorhees, W. R. Hersh, Overview of the trec 2015 clinical decision support track., in: TREC, 2015.

M. Luo, A. Mitra, T. Gokhale, C. Baral, Improving biomedical information retrieval with neural retrievers, arXiv preprint arXiv:2201.07745 (2022).

A. Alambo, T. Banerjee, K. Thirunarayan, M. Raymer, Entity-driven fact-aware abstractive summarization of biomedical literature, arXiv preprint arXiv:2203.15959 (2022).

A. Naik, S. Parasa, S. Feldman, L. L. Wang, T. Hope, Literature-augmented clinical outcome prediction, arXiv preprint arXiv:2111.08374 (2021).

M. Moradi, N. Ghadiri, Different approaches for identifying important concepts in probabilistic biomedical text summarization, Artificial intelligence in medicine 84 (2018) 101–116.

S. Deerwester, S. T. Dumais, G. W. Furnas, T. K. Landauer, R. Harshman, Indexing by latent semantic analysis, Journal of the American society for information science 41 (1990) 391–407.

Q. Cao, D. Xiong, Encoding gated translation memory into neural machine translation, in: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, 2018, pp. 3042–3047.

K. Guu, K. Lee, Z. Tung, P. Pasupat, M. Chang, Retrieval augmented language model pre-training, in: International Conference on Machine Learning, PMLR, 2020, pp. 3929–3938.

V. Karpukhin, B. Oguz, S. Min, P. Lewis, L. Wu, S. Edunov, D. Chen, W.-t. Yih, Dense passage retrieval for open-domain question answering, in: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), 2020, pp. 6769–6781.

K. Shuster, S. Poff, M. Chen, D. Kiela, J. Weston, Retrieval augmentation reduces hallucination in conversation, in: Findings of the Association for

Computational Linguistics: EMNLP 2021, Association for Computational Linguistics, Punta Cana, Dominican Republic, 2021, pp. 3784–3803. URL: https://aclanthology.org/2021.findings-emnlp.320. doi:10.18653/v1/2021.findings-emnlp.320.

R. Artstein, M. Poesio, Inter-coder agreement for computational linguistics, Computational linguistics 34 (2008) 555–596.

K. Krishna, E. Bransom, B. Kuehl, M. Iyyer, P. Dasigi, A. Cohan, K. Lo, Longeval: Guidelines for human evaluation of faithfulness in long-form summarization, arXiv preprint arXiv:2301.13298 (2023).

J. Karačić, P. Dondio, I. Buljan, D. Hren, A. Marušić, Languages for different health information readers: multitrait-multimethod content analysis of cochrane systematic reviews textual summary formats, BMC medical research methodology 19 (2019) 1–9.

M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, L. Zettlemoyer, Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension, in: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, 2020, pp. 7871–7880.

J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, Bert: Pre-training of deep bidirectional transformers for language understanding, arXiv preprint arXiv:1810.04805 (2018).

T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al., Language models are few-shot learners, Advances in neural information processing systems 33 (2020) 1877–1901.

R. Nallapati, B. Zhou, C. dos Santos, C. Gulcehre, B. Xiang, Abstractive text summarization using sequence-to-sequence rnns and beyond, in: Proceedings of the 20th SIGNLL Conference on Computational Natural Language Learning, 2016, pp. 280–290.

S. Gururangan, A. Marasović, S. Swayamdipta, K. Lo, I. Beltagy, D. Downey, N. A. Smith, Don't stop pretraining: Adapt language models to domains and tasks, in: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, 2020, pp. 8342–8360.

M. Neumann, D. King, I. Beltagy, W. Ammar, Scispacy: Fast and robust models for biomedical natural language processing, in: Proceedings of the 18th BioNLP Workshop and Shared Task, 2019, pp. 319–327.

J. Johnson, M. Douze, H. Jégou, Billion-scale similarity search with gpus, IEEE Transactions on Big Data 7 (2019) 535–547.

A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al., Pytorch: An imperative style, high-performance deep learning library, Advances in neural information processing systems 32 (2019).

T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, et al., Huggingface's transformers: State-of-the-art natural language processing, arXiv preprint arXiv:1910.03771 (2019).

C.-Y. Lin, Rouge: A package for automatic evaluation of summaries, in: Text summarization branches out, 2004, pp. 74–81.

T. Zhang*, V. Kishore*, F. Wu*, K. Q. Weinberger, Y. Artzi, Bertscore: Evaluating text generation with bert, in: International Conference on Learning Representations, 2020. URL: https://openreview.net/forum?id=SkeHuCVFDr.

K. Papineni, S. Roukos, T. Ward, W.-J. Zhu, Bleu: a method for automatic evaluation of machine translation, in: Proceedings of the 40th annual meeting of the Association for Computational Linguistics, 2002, pp. 311–318.

S. Banerjee, A. Lavie, Meteor: An automatic metric for mt evaluation with improved correlation with human judgments, in: Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization, 2005, pp. 65–72.

A. R. Fabbri, W. Kryściński, B. McCann, C. Xiong, R. Socher, D. Radev, Summeval: Re-evaluating summarization evaluation, arXiv preprint arXiv:2007.12626 (2020).

M. Coleman, T. L. Liau, A computer readability formula designed for machine scoring., Journal of Applied Psychology 60 (1975) 283.

G. Leroy, D. Kauchak, The effect of word familiarity on actual and perceived text difficulty, Journal of the American Medical Informatics Association 21 (2014) e169–e172.

Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, V. Stoyanov, Roberta: A robustly optimized bert pre-training approach, arXiv preprint arXiv:1907.11692 (2019).

K. Krippendorff, Estimating the reliability, systematic error and random error of interval data, Educational and Psychological Measurement 30 (1970) 61–70.

S. Holm, A simple sequentially rejective multiple test procedure, Scandinavian journal of statistics (1979) 65–70.

A. Pagnoni, V. Balachandran, Y. Tsvetkov, Understanding factuality in abstractive summarization with frank: A benchmark for factuality metrics, arXiv preprint arXiv:2104.13346 (2021).

Y. Guo, T. August, G. Leroy, T. Cohen, L. L. Wang, Appls: A meta-evaluation testbed for plain language summarization, arXiv preprint arXiv:2305.14341 (2023).

O. Mac, J. Ayre, K. Bell, K. McCaffery, D. M. Muscat, Comparison of readability scores for written health information across formulas using automated vs manual measures, JAMA Network Open 5 (2022) e2246051–e2246051.

## Appendix A. Appendix

*Appendix A.1. Background explanation annotation*

We provided annotators with the original abstract/LLS pair and the content (both matched and unmatched) before the 1st, 2nd, and 3rd matched sentence pairs. To make sure we have matching content from the corresponding LLS, the 1st matched sentence was included to capture situations in which explanation is provided before the first matching sentences (e.g. an introductory sentence defining terms). Examples of matching strategies can be found in Figure 1. We asked annotators to identify whether the filtered content 1) is in the background section (to confirm our heuristics indeed identify these sections); 2) truly contains external content (to confirm that unaligned GPSS sentences represent content that is absent from the source document); and 3) is paired (to confirm aligned GPSS sentences represent the same content). The results are shown in Appendix Table A.1.

*Appendix A.2. Human evaluation questions*

The questions we included in the human evaluation questionnaire are as follows:
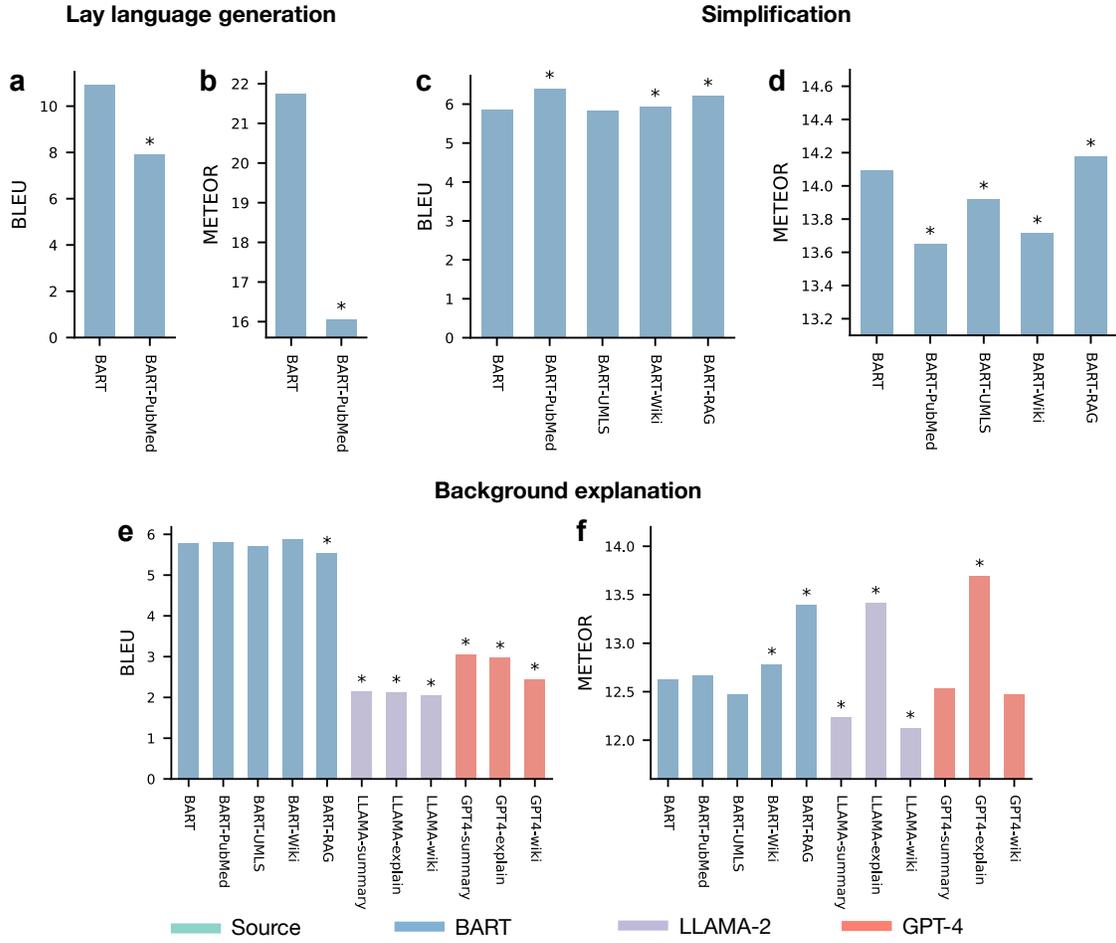
- Is the grammar of the plain text correct?

- Is the plain text easier to understand than the original text?

- Does the plain text provide all the important information from the original text?

- Is the information in the plain text correct compared to the original text?

- Does the plain text provide relevant additional information compared to the original text?

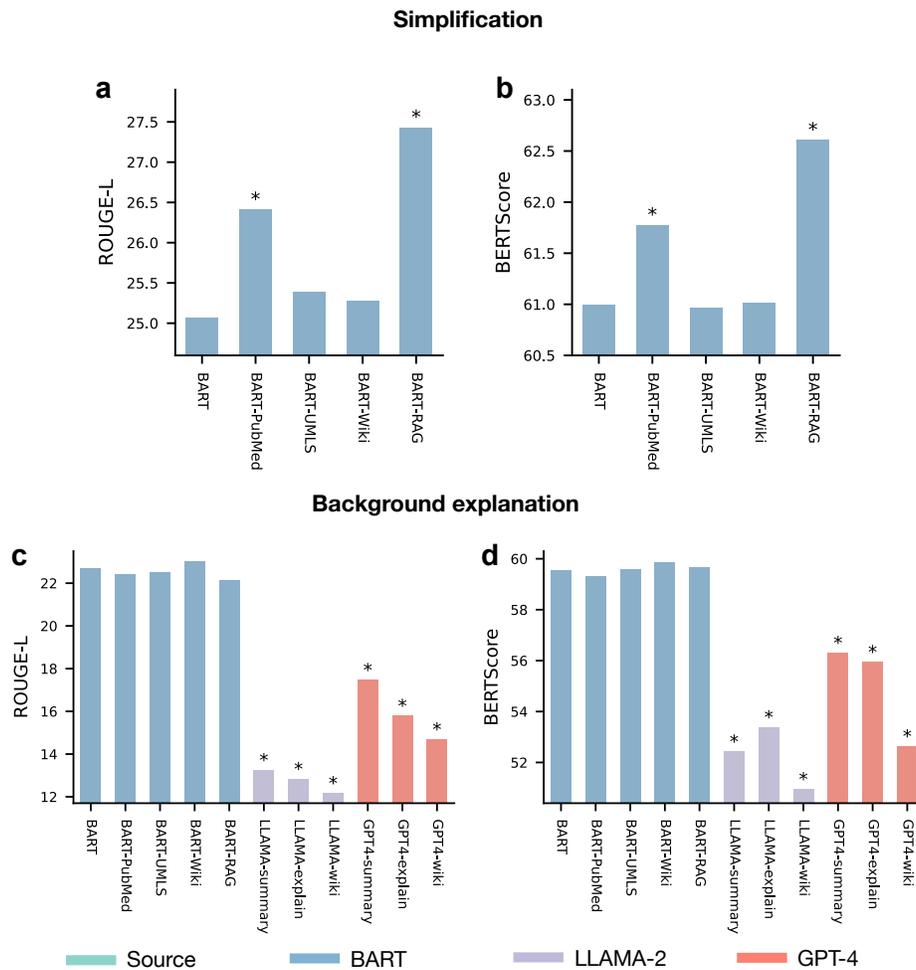|  | Annotator 1 | Annotator 2 |
|---|---|---|
| **1st Pair** | | |
| Background | 48 | 48 |
| External | 20 | 18 |
| Pair | 43 | 40 |
| **2nd Pair** | | |
| Background | 47 | 47 |
| External | 36 | 28 |
| Pair | 47 | 46 |
| **3rd Pair** | | |
| Background | 21 | 23 |
| External | 35 | 28 |
| Pair | 50 | 47 |

Appendix Table A.1: Background extraction annotation results. The columns show the number of the 50 annotated examples that each annotator labeled as containing content from the background section, including content external to the source, and being aligned with content from the source.

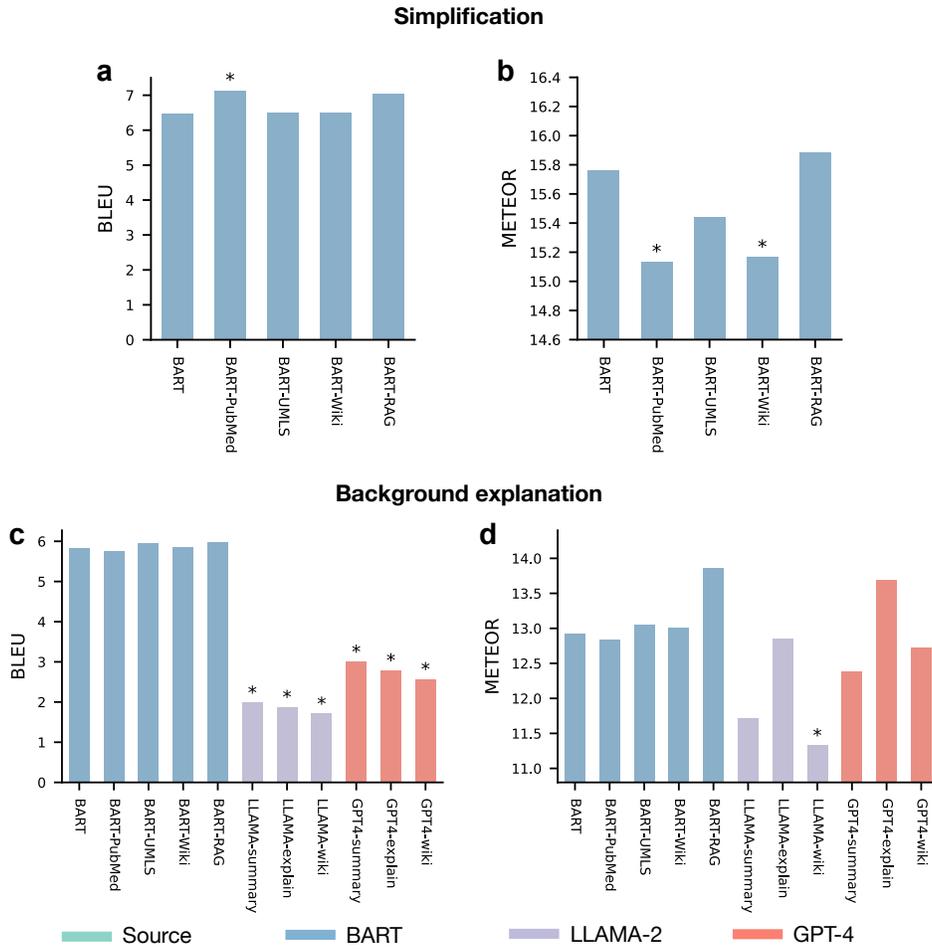| Publisher | Type | Name | Where are they displayed | Written by | Start from |
|---|---|---|---|---|---|
| Annals of the Rheumatic Diseases | Journal | Patient summary | Dedicated section of website | Authors | 2013 |
| Autism | Journal | Lay abstract | Dedicated section of website | Authors | 2011 |
| Autism Research | Journal | Lay abstract | Dedicated section of website | Authors | 2008 |
| British Journal of Dermatology | Journal | QAs | Issues searching page | Authors/editorial team | 2013 |
| Cochrane | Journal | Plain language summary | Within article | Authors/editorial team | 1997 |
| eLife | Blog | eLife digest | Section on blog | Editorial team | 2012 |
| European Urology | Journal | Patient summary | Within article | Authors | 2014 |
| NIHR Efficaacy and Mechanism Evaluation | Journal | Plain English summary | Within article | Authors | 2014 |
| NIHR Health Services and Delivery Response | Journal | Plain English summary | Within article | Authors | 2014 |
| NIHR Health Technology Assessment | Journal | Plain English summary | Within article | Authors | 2014 |
| NIHR Programme Grants for Applied Research | Journal | Plain English summary | Within article | Authors | 2015 |
| PLOS Biology | Journal | Author summary | Within article | Authors | 2007 |
| PLOS Computational Biology | Journal | Author summary | Within article | Authors | 2005 |
| PLOS Genetics | Journal | Author summary | Within article | Authors | 2005 |
| PLOS Medicine | Journal | Author summary | Within article | Authors | 2006 |
| PLOS Neglected Tropical Diseases | Journal | Author summary | Within article | Authors | 2007 |
| PLOS Pathogens | Journal | Author summary | Within article | Authors | 2005 |
| Proceedings of the National Academy of Sciences | Journal | Significance | Within article | Authors | 2013 |
| Reproductive Health | Journal | Plain English summary | Within article | Authors | 2016 |

Appendix Table A.2: Summary of journals with Plain Language Summary

**Lay language generation**

**Simplification**
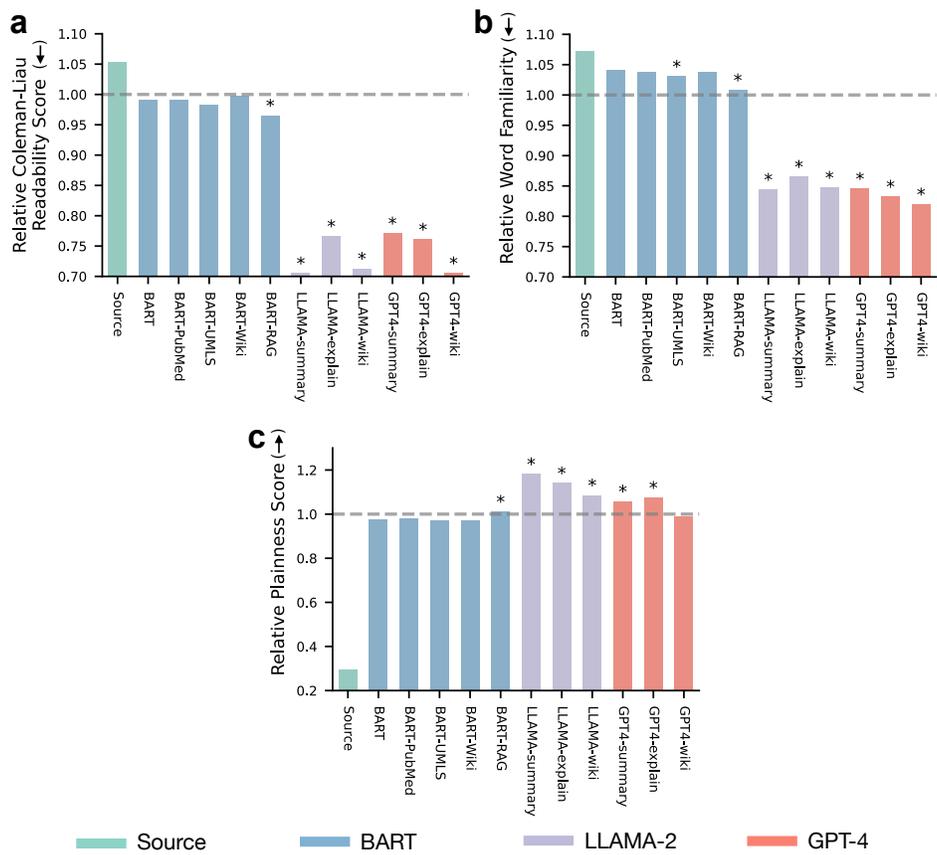


**Background explanation**



Appendix Figure A.1: Models' performance in text generation. We used the F1 score of BLEU and METEOR to evaluate the generation quality of models on lay language generation, simplification, and background explanation tasks. P-values obtained through the t-test are employed to evaluate the performance of various models compared to the Vanilla model (BART). A p-value less than 0.05 is indicated by (*).

**Simplification**



**Background explanation**



Appendix Figure A.2: Models' performance in text generation on the validated dataset. We used the F1 score of ROUGE-L and BERTScore to evaluate the generation quality of models on lay language generation, simplification, and background explanation tasks. P-values obtained through the t-test are employed to evaluate the performance of various models compared to the Vanilla model (BART). A p-value less than 0.05 is indicated by (*).

Appendix Figure A.3: Models' performance in text generation on the validated dataset. We used the F1 score of BLEU and METEOR to evaluate the generation quality of models on lay language generation, simplification, and background explanation tasks. P-values obtained through the t-test are employed to evaluate the performance of various models compared to the Vanilla model (BART). A p-value less than 0.05 is indicated by (*).

Appendix Figure A.4: Readibility, familiarity, and plainness of the validated background explanation subset. (a) Relative Coleman-Liau readability score, (b) word familiarity, and (c) plainness score of the source and models' generated text. The relative score is calculated by dividing by the score of the target text. A lower readability score and word familiarity indicate that the text is easier to read (values below the dashed line are lower than those from professionally-authored plain language). A higher Plainness Score indicates that the text is more representative of an LLS. P-values obtained through the t-test are employed to evaluate the performance of various models compared to the Vanilla model (BART). A p-value less than 0.05 is indicated by (*).