

A Framework for Fine-Grained Complexity Control in Health Answer Generation

Daniel Ferreira
IEETA
University of Aveiro
Aveiro, Portugal
djbf@ua.pt

Tiago Melo Almeida
IEETA
University of Aveiro
Aveiro, Portugal
tiagomeloalmeida@ua.pt

Sérgio Matos
IEETA, DETI, LASI
University of Aveiro
Aveiro, Portugal
aleixomatos@ua.pt

Abstract

Health literacy plays a critical role in ensuring people can access, understand, and act on medical information. However, much of the health content available today is too complex for many people, and simplifying these texts manually is time-consuming and difficult to do at scale. To overcome this, we developed a new framework for automatically generating health answers at multiple, precisely controlled complexity levels. We began with a thorough analysis of 166 linguistic features, which we then refined into 13 key metrics that reliably differentiate between simple and complex medical texts. From these metrics, we derived a robust complexity scoring formula, combining them with weights learned from a logistic regression model. This formula allowed us to create a large, multi-level dataset of health question-answer pairs covering 21 distinct complexity levels, ranging from elementary patient-friendly explanations to highly technical summaries. Finally, we fine-tuned a Llama-3.1-8B-Instruct model using “control codes” on this dataset, giving users precise control over the complexity of the generated text and empowering them to select the level of detail and technicality they need.

1 Introduction

Health literacy, which is the ability to obtain, process, and understand basic health information, remains a significant challenge worldwide. A survey conducted by the World Health Organization (WHO) between 2019 and 2021 across 17 European countries found that between 25% and 75% of people struggle with understanding health-related information, with variation depending on country-specific factors like education and healthcare access (Pelikan et al., 2021).

In the United States, approximately 80 million adults had limited health literacy as of 2018, with disproportionately higher rates among older adults, minority groups, and individuals of lower socioeco-

nomic status (Woods et al., 2023). These statistics matter because people with lower health literacy often struggle to understand medical terms, leading to poorer health outcomes and increased healthcare costs (Shahid et al., 2022). This issue becomes even more important as more people turn to online sources for health information. In 2022, 58.5% of U.S. adults searched for health information online (Wang and Cohen, 2022), yet studies show that most health-related content online exceeds recommended readability levels (Szmuda et al., 2020; Mohile et al., 2023).

Large language models (LLMs) like GPT-4 (OpenAI, 2023), Med-PaLM (Singhal et al., 2023), and Claude (Anthropic, 2024) now generate health information and are increasingly used in healthcare contexts. However, these models typically produce text at a fixed complexity level, often too advanced for many readers (Amin et al., 2024). Current approaches to medical text simplification focus on converting complex text into simpler versions (Gondy et al., 2018; Flores et al., 2023; Li et al., 2024) rather than dynamically adjusting complexity based on individual needs.

This gap presents an opportunity to develop language models that can generate health answers with adjustable complexity levels, a capability that would make information more accessible to everyone, regardless of their health literacy level.

2 Related Work

This section provides an overview of existing literature and previous research relevant to the scope of this study.

2.1 Text Complexity and Readability Assessment

The earliest attempts to measure text complexity used simple formulas based on surface-level features. Smith and Senter (1967) developed the Au-

tomated Readability Index (ARI), which counts characters per word and sentence length to estimate reading difficulty. Shortly after, Kincaid et al. (1975) created the Flesch-Kincaid Grade Level formula, which also considers syllable counts and remains widely used today for its simplicity and reliability.

Zheng and Yu (2018) noted that standard formulas failed to capture medical complexity because they ignored specialized terminology and semantic relationships. They developed a ranking system that compared documents relative to each other rather than assigning absolute scores, using both surface-level features and word embeddings to better match human judgments of readability.

Jiang and Xu (2024) created MedReadMe, manually annotating 4,520 medical sentences with readability labels and identifying complex spans within each sentence. They introduced “Google-Easy” and “Google-Hard” categories based on how commonly terms appear in web searches. Their analysis of 650 linguistic features revealed that medical jargon density and syntactic complexity were the strongest predictors of reading difficulty.

Devaraj et al. (2021) proposed using a masked language model (MLM) to differentiate technical and lay medical text. Their method evaluates how accurately a model trained on scientific literature predicts masked tokens, based on the observation that technical terminology is more predictable within domain-specific contexts. Luo et al. (2022) improved this method by focusing on noun phrases, allowing multi-word medical terms like “heart attack” to be treated as single semantic units.

While methods based on masked language modeling have shown promise, they mainly focus on single-word complexity. Lyu and Pergola (2024) addressed this limitation with SciGisPy, a metric rooted in Fuzzy-Trace Theory (FTT) (Reyna, 2012) that evaluates how well simplified texts preserve the core meaning (gist), emphasizing semantic coherence and the ability to form clear mental models.

2.2 Medical Text Simplification

Medical text simplification started with straightforward rule-based systems. For instance, Damay et al. (2006) used techniques like lexical substitution and sentence restructuring to make medical texts easier to understand. Later, Kandula et al. (2010) took this further by combining both semantic and syntactic methods to simplify electronic medical records and patient education materials.

The field progressed significantly with the development of large-scale datasets for training language models. Devaraj et al. (2021) created the Cochrane dataset, which pairs technical abstracts with lay summaries from the Cochrane Database of Systematic Reviews. Using this parallel data, they trained BART models with unlikelihood training, explicitly penalizing the generation of tokens identified as technical language through a bag-of-words classifier. Flores et al. (2023) replaced the bag-of-words classifier with the Flesch-Kincaid readability formula to identify and penalize complex words. To prevent hallucinations that can occur when optimizing solely for simplicity, they also incorporated factual consistency into their loss function and designed a beam search method that weighs both readability and accuracy during decoding.

Basu et al. (2023) created Med-EASi, a finely annotated dataset for simplifying medical texts that identifies four types of textual transformations: elaboration, replacement, deletion, and insertion. With this dataset, they built T5-based models that allow users to select specific medical terms and control exactly how they should be simplified.

Lu et al. (2023) developed NapSS, a two-stage “summarize-then-simplify” method for medical text simplification that first identifies important sentences using a summarizer trained on paired technical abstracts and their human-simplified versions, and then extracts key phrases to create “narrative prompts” that guide the language model during the simplification process, helping preserve the logical flow and medical accuracy of the original text.

Phatak et al. (2022) applied reinforcement learning to medical text simplification by designing reward functions that balance content preservation, Flesch-Kincaid readability scores, and lexical simplicity. Rahman et al. (2024) later created SimpleDC, a dataset of original and simplified texts related to digestive cancers. They fine-tuned LLaMA models on this dataset and further improved them using reinforcement learning, guided by a binary classifier trained to detect simple language.

2.3 Controllable Text Generation

Recent research has explored ways to control text readability during generation. Ribeiro et al. (2023) developed methods for controllable summarization using instruction-based prompting, reinforcement learning with a Gaussian reward function that penalizes deviations from desired readability scores, and lookahead decoding to anticipate how word

choices impact readability.

Luo et al. (2022) focused on readability control specifically for biomedical text summarization. They first tried prepending special tokens as prompts to the input and then tested a multi-head architecture with separate decoders for different readability levels. While the multi-head approach helped create some distinction between technical and plain language outputs, they found that the level of readability control was still very limited.

Tran et al. (2024) introduced ReadCtrl, which instruction-tunes language models to generate text at specific readability scores on an almost continuous scale rather than predefined categories. Meanwhile, Hsu et al. (2024) found that even with clear instructions, language models often produce outputs that do not align with traditional readability metrics. They also showed that readers generally preferred explanations written at a high school level, suggesting that there may be a sweet spot of complexity balancing clarity and informative content.

While prior work has focused primarily on binary simplification or relied on traditional readability metrics that fail to capture the unique challenges of medical terminology, we developed a more comprehensive framework that integrates multiple linguistic features to accurately measure the complexity of medical text and generate content at precisely targeted readability levels.

3 Methods

This section details the framework developed for automatically generating health answers at multiple complexity levels, as illustrated in Figure 1.

3.1 Data Collection

We used two established datasets containing paired original and simplified medical texts. Though these datasets provide parallel texts at different complexity levels, the “simplified” versions, while less complex than the originals, are not always simple in absolute terms. This relative simplification creates a sliding scale rather than distinct complexity levels, making it difficult to develop a reliable readability formula. To overcome this limitation, we created a synthetic dataset containing pairs of clearly differentiated simple and complex medical texts.

3.1.1 Medical Text Simplification Datasets

We evaluated our metrics using two parallel corpora of medical texts: PLABA (Attal et al., 2023) and

Cochrane (Devaraj et al., 2021). Both datasets include original medical texts paired with simplified versions. PLABA contains sentence and paragraph-level simplifications of biomedical abstracts, while Cochrane focuses on paragraph-level simplifications of systematic reviews. More detailed descriptions are available in Appendix A.1.

Table 1 summarizes the key characteristics of the three datasets used in this stage of the project.

3.1.2 HSQA-Claude Dataset

We created a new dataset using Claude 3.5 Sonnet to generate answers to questions from the HealthSearchQA dataset (Singhal et al., 2023), which contains 3,173 commonly searched consumer medical queries. We manually identified and filtered out questions that were not genuinely health-related to ensure the quality and relevance of our dataset. For each valid question, we prompted the model to produce one answer using technical medical language suitable for healthcare professionals, and another using simple language appropriate for patients with limited health literacy. This approach provided clearly differentiated examples of simple and complex medical text covering the same information content.

Dataset	Source	# Pairs
<i>PLABA-sent</i>	PubMed abstracts	7,643
<i>PLABA-para</i>	PubMed abstracts	750
<i>Cochrane</i>	Systematic Reviews Database	4,459
<i>HSQA-Claude</i>	HealthSearchQA questions	3,150

Table 1: Parallel datasets used for text complexity analysis.

3.2 Metrics

We implemented 166 metrics to measure text readability and complexity, covering various linguistic dimensions. We chose this broad scope to comprehensively explore and identify the most robust indicators of medical text complexity, given the multifaceted nature of readability and the lack of a single, universally agreed-upon metric in the domain. The following sections describe each category of metrics we used in our analysis.

3.2.1 Traditional Metrics

We calculated 20 traditional readability formulas, including Flesch-Kincaid Grade Level (Kincaid et al., 1975), SMOG Index (McLaughlin,

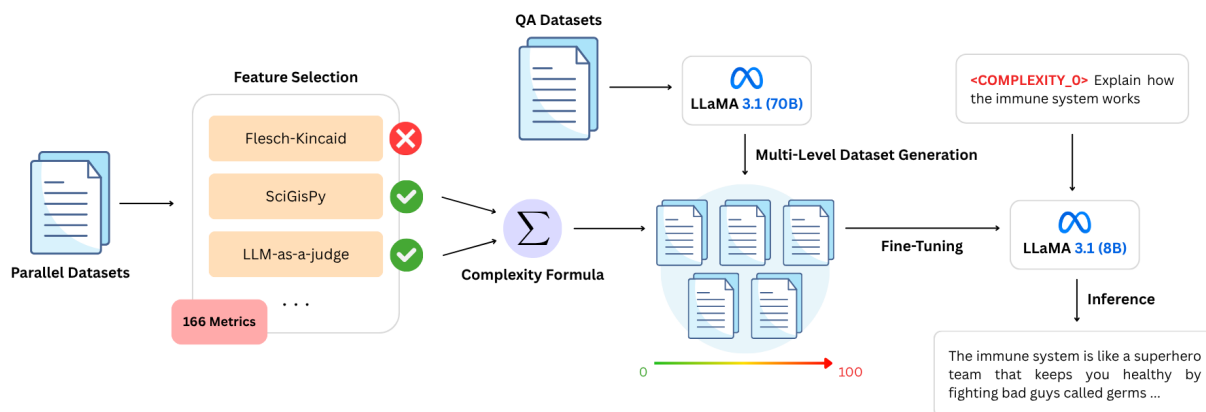


Figure 1: Framework for complexity-controlled health answer generation.

1969), and Coleman-Liau Index (Coleman and Liau, 1975). These metrics estimate text difficulty based on surface-level features like word length, syllable count, and sentence length, working on the general assumption that longer lexical units require more cognitive effort, thereby making the text more complex (Yu et al., 2020). Although not designed for biomedical literature, they can serve as a useful starting point to judge how easy or difficult a text is to read and understand. We supplemented these with 8 statistical measures capturing additional aspects of readability, including the proportion of difficult words from the Dale-Chall list (Dale and Chall, 1948) and lexical diversity metrics such as TTR and MTLTD (McCarthy and Jarvis, 2010).

3.2.2 Syntactic Structure

We implemented 16 syntax-based metrics using spaCy (Honnibal et al., 2020) for dependency parsing and part-of-speech tagging, organized into two categories. For lexical distribution, we calculated content-to-function word ratio, which compares meaning-carrying words to grammatical words (Just and Carpenter, 1992), and part-of-speech distributions to identify texts with higher noun density typical of scientific writing (Biber et al., 1999). For structural complexity, we measured dependency distance (Gibson, 2000), passive voice proportion (Ferreira, 2003), noun phrase length (Biber et al., 1999), embedding depth (Gibson, 1998), negation density, and left-right asymmetry (Hawkins, 2004). These metrics capture aspects of syntactic complexity that increase cognitive load, such as deeply embedded clauses and words separated from their grammatical dependents.

3.2.3 Medical Terminology and Jargon

We implemented 19 term-level metrics using the Unified Medical Language System (UMLS) Metathesaurus (National Library of Medicine, 2024) and Consumer Health Vocabulary (CHV) (Zeng and Tse, 2006). For concept identification, we used QuickUMLS (Soldaini, 2016), which performs faster approximate dictionary matching compared to MetaMap (Aronson and Lang, 2010). These metrics include term density, expert-to-lay ratio, semantic type diversity, and CHV familiarity scores that measure how frequently terms appear in consumer health materials (Keselman et al., 2007).

We also built a RoBERTa-large (Liu et al., 2019) sequence tagger with Conditional Random Fields (CRF), trained on the MedReadMe dataset to identify seven distinct categories of medical jargon as defined by Jiang and Xu (2024). These categories include easy and hard medical terms, medical entities, complex terms, multisense words, and medical and general abbreviations. This method enables more fine-grained analysis than dictionary lookups, capturing context-dependent terminology and terms absent from UMLS. From this, we derived 29 other metrics capturing jargon density, distribution across categories, and clustering patterns.

3.2.4 Gist Formation

We adapted GisPy (Hosseini et al., 2022), an open-source tool based on Fuzzy-Trace Theory (Reyna, 2012), which measures how easily readers can understand the essential meaning of a text. GisPy calculates scores for several components that contribute to gist formation, including referential cohesion (connecting ideas between sentences), coreference resolution (tracking entities throughout text), deep cohesion (presence of causal connectives),

and semantic verb overlap (relatedness of actions). We modified the original implementation to use BioSimCSE-BioLinkBERT-BASE (raj Kanakarajan et al., 2022), trained on biomedical literature, making it more suitable for our task. We also implemented SciGisPy (Lyu and Pergola, 2024), which tailors GisPy for biomedical text simplification. SciGisPy introduces domain-specific improvements, such as information content measures derived from biomedical corpora and semantic chunking to measure topic cohesion.

3.2.5 Masked Language Model

We implemented three MLM-based metrics using Bio+ClinicalBERT (Alsentzer et al., 2019), which outperformed other BERT variants in our tests. These metrics measure complexity by calculating how predictable medical terminology is within context. The first metric randomly masks 15% of tokens, the second specifically targets noun phrases, and the third applies a ranking method (RNPTC), which weighs phrases based on their prediction probability (Luo et al., 2022). We found that increasing the number of random masking iterations from 10 to 30 significantly improved reliability by reducing variance. As a result, the simpler random masking approach became more effective than the other two methods in distinguishing between technical and simplified texts.

3.2.6 Semantic Clustering

We built on the method introduced by Cha et al. (2017), which uses word embeddings to measure text complexity. In our implementation, each word is mapped to a BioWordVec embedding (Zhang et al., 2019), and these vectors are grouped using K-means clustering. While the original implementation used 100 clusters, we increased this to 300 to better reflect the distinctions in medical vocabulary. We then create a count vector for how often words fall into each cluster, which serves as a feature vector for predicting readability. We trained two separate Support Vector Regression (SVM) models, one using the CLEAR corpus (Crossley et al., 2023), and another using the MedReadMe dataset (Jiang and Xu, 2024) for medical texts.

3.2.7 ALBERT Transformer

We used the ALBERT-xxlarge model (Lan et al., 2019) from the winning entry in the CommonLit Readability Prize Kaggle competition (Malatinszky et al., 2021). This model processes text through

attention layers to capture relationships between words before predicting a readability score. Although the original solution used an ensemble of models, ALBERT-xxlarge was singled out by the winner as especially important, thanks to its parameter-sharing structure, which helps prevent overfitting while still capturing complex language features. The same model was later reused in the REFereE framework for evaluating text simplification (Huang and Kochmar, 2024).

3.2.8 LLM Expert Evaluation

We created a hybrid method for evaluating text readability using large language models as expert evaluators. Specifically, we prompted three 70 billion-parameter models (Nvidia-Llama-3.1-Nemotron-70B (Wang et al., 2024), Llama3-OpenBioLLM-70B (Pal, 2024), and DeepSeek-R1-Distill-Llama-70B (DeepSeek-AI, 2025)) to evaluate texts on five dimensions: vocabulary complexity, syntactic complexity, conceptual density, required background knowledge, and overall cognitive load. Each model rated texts on a 1–5 scale using few-shot prompting with three calibration examples that we personally annotated. Because running multiple large models is computationally expensive, we trained a smaller and more efficient BioSimCSE-BioLinkBERT-BASE model (raj Kanakarajan et al., 2022) on the averaged LLM scores. This distilled model not only processes texts much faster, but also improves the results by smoothing out inconsistencies in the original LLM judgements.

3.3 Formula Development

After collecting and implementing the linguistic features, we followed a systematic approach to select the most reliable features for our complexity formula. Since we lacked human-annotated readability scores, we developed a data-driven methodology to identify stable features that consistently distinguished simple from expert-level medical texts, using the datasets described in Section 3.1.

The feature selection process began by removing features with absolute pairwise correlations above 0.7 to reduce collinearity and lower the risk of unintentionally excluding important features from the final model. We then applied Lasso logistic regression with bootstrapping, adapting the methodology described by Laurin et al. (2016), which involved the following steps:

1. Creating 1,000 bootstrap samples from our

training data using random sampling with replacement.

2. Fitting a Lasso logistic regression model to each bootstrap sample to classify if a text was written for experts or general audience.
3. Calculating the coefficient of variation (CV) for each feature, defined as the standard deviation divided by the mean absolute value of the coefficient, across bootstrap samples.
4. Using the interquartile range (IQR) method to exclude features with unstable coefficients by calculating the upper fence ($Q3 + 1.5 \times IQR$). Features with CV exceeding this threshold were considered outliers and removed.
5. Further filtering features if the 95% confidence interval for the value of the coefficient included zero.

We then trained our final logistic regression model using only the HSQA-Claude dataset, which contains controlled comparisons of text complexity with a cleaner signal-to-noise ratio. For this purpose, we used ElasticNet regularization to estimate feature weights, as it balances the benefits of both Lasso and Ridge regression and better handles any remaining collinearity among features. This process resulted in a final set of 13 metrics (listed in Appendix B.1) after excluding those that performed exceptionally well in one dataset but poorly or inconsistently in others. These features were likely overfitting to specific data characteristics and were removed to improve generalizability.

3.4 Multi-Level Dataset

After developing and validating our complexity formula, we created a medical dataset containing answers rewritten at multiple levels of complexity to train our controlled text generation model.

3.4.1 Source Datasets

We built our dataset using question-answer pairs from five established medical datasets: LiveQA (Abacha et al., 2017), MedicationQA (Abacha et al., 2019), MEDIQA-AnS (Savery et al., 2020), MedQuAD (Abacha and Demner-Fushman, 2019), and BioASQ Task 13B. After cleaning and filtering for quality, we retained 31,917 question-answer pairs. Table 2 provides a brief overview of these datasets, with detailed descriptions available in Appendix A.2.

Dataset	Source	# Pairs
<i>LiveQA</i>	U.S. NLM	800
<i>MedicationQA</i>	NIH websites	690
<i>MEDIQA-AnS</i>	CHiQA-retrieved passages	312
<i>MedQuAD</i>	NIH websites	16,423
<i>BioASQ</i>	PubMed/MEDLINE articles	13,692

Table 2: Source datasets used to create our multi-level medical QA dataset

3.4.2 Dataset Creation

For each question-answer pair in our source datasets, we created five versions of the answer, each written for a different audience, namely young children, middle school students, high school students, college graduates, and biomedical experts. We generated these answers using the models described in Section 3.2.8, with DeepSeek handling 70% of the generation, Nemotron 20%, and OpenBioLLM 10%. This allocation was based on preliminary experiments, which showed that using multiple models helped capture a broader range of writing styles for each education level.

We designed a prompt that generated all five variants simultaneously, with answers becoming progressively more complex (see Appendix C.3). The prompt included three examples to guide the models, descriptions of each target audience, and instructions to keep the answers factually accurate. It also instructed the models to flag any cases where the original answer did not fully address the question, allowing us to filter out problematic samples from the dataset early on.

After generating the variants, we checked the quality of all answers through a two-stage process. First, we used regex patterns to identify and remove samples containing placeholder text instead of proper content. Then we evaluated each variant against its original answer using metrics for content preservation and factual accuracy, including ROUGE (Lin, 2004), BLEURT (Sellam et al., 2020), BERTScore (Zhang et al., 2019), UniEval (Zhong et al., 2022), and SummaC (Laban et al., 2022). The filtering identified relatively few problems and only 2,926 samples (1.56%) were removed from the initial 187,769. This low rejection rate was not surprising, since the variants were created directly from the original answers. Most of the issues found actually stemmed from contradictions or inaccuracies present in the source material.

Each variant was annotated using the complexity formula described in Section 3.3. This gave us raw scores between -34.56 and 31.99, which we converted to a more practical 0-100 scale and then binned into 21 categories labeled 0, 5, 10, and so on up to 100, with each bin containing roughly 8,800 samples. These bins aligned reasonably well with our original five levels, though with some natural overlap between categories. For example, the majority of high school-level variants fell within bins labeled 50-70, while college-level variants typically ranged from 60-80.

The final dataset includes 184,843 answers for 36,969 questions. Each entry has the original question, the reference answer, the variants at different complexity levels, as well as the corresponding evaluation metrics and complexity scores.

3.5 Model Fine-Tuning

After creating our multi-level dataset, we fine-tuned a language model to generate medical text with controlled complexity levels. We experimented with two different methods: natural language instructions and control codes.

For natural language instructions, we used prompts like “Answer the following question with a complexity score of 75 out of 100.” For control codes, we added special tokens to the model’s vocabulary (e.g., “<COMPLEXITY_75>”) and placed them at the beginning of each prompt. These new tokens were initialized by positioning them along a “complexity direction” in the embedding space. We identified simple and complex anchor words in the model’s vocabulary, created a vector between them, and placed our tokens along this vector. This gave the tokens semantic meaning before training even began.

We selected Llama-3.1-8B-Instruct (Grattafiori et al., 2024) as our base model and applied LoRA fine-tuning (Hu et al., 2021) with rank 8, alpha 16, and a learning rate of 5e-5, and targeted all projection matrices in the transformer architecture.

During training, we implemented context-aware batching, grouping all answers for the same medical question into a single batch. This helped the model focus on the patterns that actually matter and avoid spurious correlations. For example, if a batch includes both simple and technical answers about asthma, gradient updates adjust the model’s weights to preserve important details, such as inflammation and breathing issues, while tailoring the language to match the desired complexity level.

We found that using control codes worked better than using natural language instructions. The training converged faster, and the model generated more consistent responses at each complexity level.

4 Experiments and Results

This section details the evaluation of our complexity scoring formula and the performance of our fine-tuned model in generating text at specific complexity levels.

4.1 Formula Validation

We evaluated our complexity scoring formula using data from the four datasets introduced in Section 3.1. We trained the formula on 80% of the HSQA-Claude dataset and tested it on the remaining 20%, as well as the complete Cochrane and PLABA datasets. This setup helped us determine how well our formula works for different text types and simplification strategies.

For comparison, we used two baselines. The first was the Flesch-Kincaid Grade Level (FKGL), which is the most popular and widely used readability formula today. The second baseline (marked with † in Table 3) corresponds to the best-performing metric for each dataset, selected post hoc from the full set of existing metrics.

To evaluate performance, we used three complementary statistical measures. Cohen’s d measures the standardized difference between the means of two distributions by indicating how many standard deviations separate the simple and complex text groups. The Area Under the Curve (AUC) measures how well the scoring method distinguishes between the two classes, giving an estimate of the probability that a randomly chosen complex text receives a higher score than a randomly chosen simple one. Jensen-Shannon (JS) Divergence measures the dissimilarity between two probability distributions by comparing their entire shapes rather than just their averages or classification accuracy.

Figure 2 shows the score distributions of simple (green) and complex (red) texts using our formula. The HSQA-Claude dataset shows the clearest separation between the two groups, with virtually no overlap. The Cochrane and PLABA-para datasets also show good separation, although with more overlap between the distributions. This likely happens because many of the so-called “simplified” texts in these datasets still include difficult jargon and remain relatively complex. The PLABA-sent

dataset has the highest overlap, since shorter texts often do not provide enough context to reliably judge their complexity.

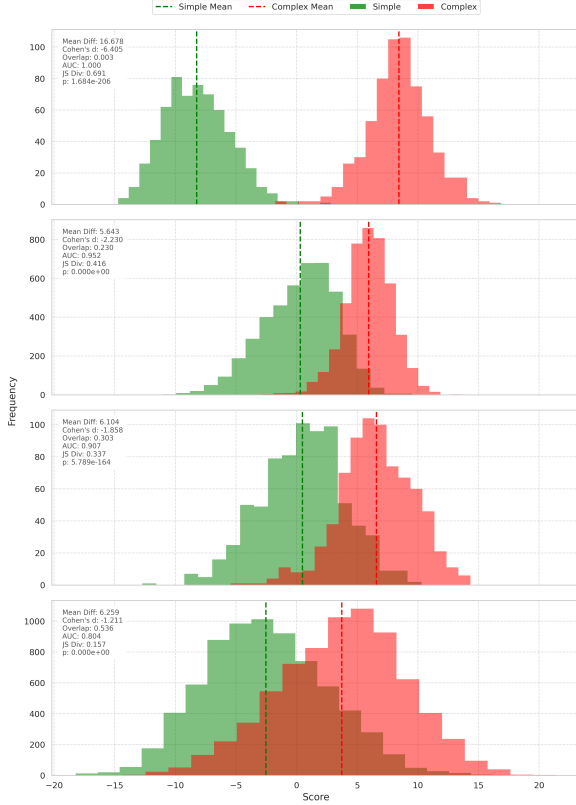


Figure 2: Distribution of complexity scores in the four parallel text datasets.

Table 3 compares our formula against the baseline methods. While certain metrics occasionally show slightly better results on specific datasets, their performance fluctuates more from case to case. In contrast, our formula consistently delivers strong results regardless of text length, domain, or simplification strategy. Moreover, perfect numerical separation is not always ideal, as some degree of overlap between distributions may actually reflect genuine ambiguities or edge cases in the data, not necessarily a flaw in the scoring method. In practice, what matters is how well a score captures the perceived reading difficulty experienced by individuals with different levels of health literacy, not just how cleanly it separates two labeled groups in a curated dataset.

4.2 Model Performance

We evaluated the ability of our fine-tuned model to generate text at specific complexity levels by comparing it to the original base model and a version using few-shot prompting. Using 100 questions sam-

Dataset	Method	Cohen's <i>d</i>	AUC	JS Div.
PLABA-sent	Our formula	1.21	0.80	0.16
	FKGL	0.58	0.67	0.05
	†	0.99	0.76	0.11
PLABA-para	Our formula	1.86	0.91	0.34
	FKGL	0.95	0.76	0.12
	†	1.89	0.91	0.32
Cochrane	Our formula	2.23	0.95	0.42
	FKGL	0.61	0.68	0.06
	†	2.36	0.95	0.42
HSQA-Claude	Our formula	6.40	1.00	0.69
	FKGL	1.58	0.90	0.31
	†	6.11	1.00	0.67

† Represents the best-performing metric for each dataset.

Table 3: Comparison of readability scoring methods.

pled from HealthSearchQA (Singhal et al., 2023), we generated responses at each target complexity level and calculated the difference between the requested complexity and the actual complexity of the generated text.

Figure 3 shows the relationship between the target and the generated complexity levels for each model. The fine-tuned model closely follows the ideal diagonal line, particularly at lower and mid-range levels. However, there is some compression at the highest levels (80-100), an issue that requires detailed examination in future studies. The few-shot approach shows a step-like pattern, indicating that it captures general complexity trends but lacks fine-grained control. Meanwhile, the baseline outputs are clustered around a fixed level (~ 60), showing little response to different targets.

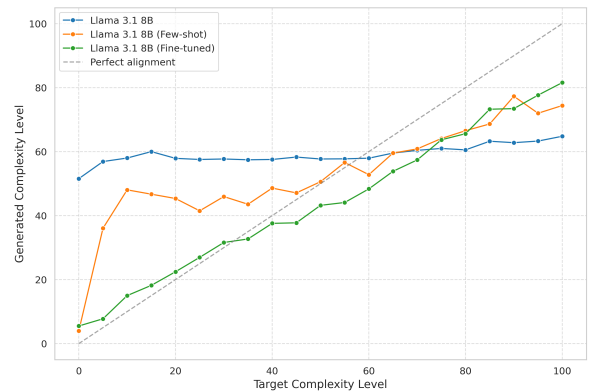


Figure 3: The ability of each model to generate text at the desired complexity level.

4.3 User-Centric Evaluation

To better understand the practical impact of our complexity control mechanism on end-users, we

conducted a downstream evaluation using simulated agents, powered by Claude Sonnet 4, as proxies for human evaluators. This method was chosen to overcome the logistical challenges associated with recruiting and managing a large pool of human participants with varying levels of health literacy.

We designed three user personas representing low, medium, and high health literacy levels, with specific prompts to influence how they interpret the content. For instance, the low-literacy persona was described as having “no medical training and rely on everyday language,” while the high-literacy persona was a “healthcare professional... comfortable with medical terminology.” The full prompts used for these personas are provided in Appendix C.4.

Each simulated user independently rated the responses along five quality dimensions on a scale of 1 to 5. These dimensions included understandability (ease of comprehension), usefulness (practical and actionable guidance), relevance (directness in addressing the question), and factuality (medical accuracy and reliability).

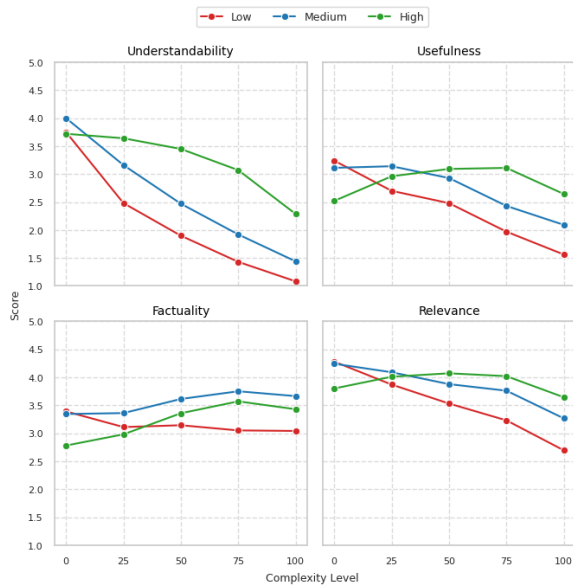


Figure 4: Evaluation scores from simulated user personas with low, medium, and high health literacy.

As shown in Figure 4, which presents average scores for five complexity levels (0, 25, 50, 75, and 100), increasing complexity leads to a dramatic and consistent drop in understandability for the three personas, with scores declining approximately 40-70% from the simplest to the most complex levels. This suggests that although more complex responses may contain richer information, they become substantially harder to follow regardless of

the reader’s health literacy level. When it comes to factuality, the scores remain relatively stable and, in some cases, even show a slight improvement, which indicates that changes in complexity do not come at the cost of medical accuracy. On the other hand, relevance and usefulness both vary greatly depending on the persona. Simpler answers are more helpful and relevant for users with low and medium health literacy, whereas the high-literacy persona seems to favor more complex responses, though this benefit plateaus and slightly decreases at the highest complexity level.

While these findings highlight the trade-offs involved in adjusting complexity for different user groups, it is important to acknowledge that simulated agents cannot fully replicate the nuanced and multifaceted ways genuine human users process and respond to medical information, including their emotional reactions, personal health contexts, and individual communication preferences. Therefore, these results should be viewed as indicative rather than definitive of actual human behavior.

5 Conclusions

We introduce a framework for creating medical answers tailored to different health literacy levels. We analyzed 166 linguistic features and defined a scoring formula based on a smaller set of 13, incorporating domain terminology, syntactic complexity, and signals from large language models, to reliably distinguish simple from complex medical text. Using this formula and public resources including LiveQA, MedQuAD, and BioASQ, we created a large dataset of 184,843 medical question-answer pairs rewritten at 21 complexity levels, filling a gap in training materials. We then fine-tuned a language model to generate text at distinct complexity levels, from very simple explanations to highly technical content for medical professionals. This versatility makes it useful in many healthcare settings. It can help create personalized patient education materials, support medical students as they learn more advanced topics, and generate documentation for healthcare providers, such as doctors and nurses.

Acknowledgments

This work was funded by national funds through FCT - Fundação para a Ciência e a Tecnologia, I.P., under unit UID/00127.

References

- Asma Ben Abacha, Eugene Agichtein, Yuval Pinter, and Dina Demner-Fushman. 2017. [Overview of the medical question answering task at trec 2017 liveqa](#).
- Asma Ben Abacha and Dina Demner-Fushman. 2019. [A question-entailment approach to question answering](#). *BMC Bioinformatics*, 20.
- Asma Ben Abacha, Yassine Mrabet, Mark Sharp, Travis R. Goodwin, Sonya E. Shooshan, and Dina Demner-Fushman. 2019. [Bridging the gap between consumers' medication questions and trusted answers](#). *Studies in Health Technology and Informatics*, 264:25–29.
- Emily Alsentzer, John Murphy, William Boag, Wei-Hung Weng, Di Jin, Tristan Naumann, and Matthew McDermott. 2019. [Publicly available clinical bert embeddings](#). In *Proceedings of the 2nd Clinical Natural Language Processing Workshop*, pages 72–78. Association for Computational Linguistics.
- Kanhai S Amin, Linda C Mayes, Pavan Khosla, and Rushabh H Doshi. 2024. [Assessing the efficacy of large language models in health literacy: A comprehensive cross-sectional study](#). *Yale Journal of Biology and Medicine*, 97:17–27.
- Anthropic. 2024. [The claude 3 model family: Opus, sonnet, haiku anthropic](#).
- Alan R. Aronson and François Michel Lang. 2010. [An overview of metamap: historical perspective and recent advances](#). *Journal of the American Medical Informatics Association : JAMIA*, 17:229–236.
- Kush Attal, Brian Ondov, and Dina Demner-Fushman. 2023. [A dataset for plain language adaptation of biomedical abstracts](#). *Scientific Data*, 10:1–11.
- Chandrayee Basu, Rosni Vasu, Michihiro Yasunaga, and Qian Yang. 2023. [Med-easi: Finely annotated dataset and models for controllable simplification of medical texts](#). *Proceedings of the 37th AAAI Conference on Artificial Intelligence, AAAI 2023*, 37:14093–14101.
- Douglas Biber, Stig Johansson, Geoffrey Leech, Susan Conrad, and Edward Finegan. 1999. *Longman Grammar of Spoken and Written English*. Pearson Education.
- Miriam Cha, Youngjune Gwon, and H. T. Kung. 2017. [Language modeling by clustering with word embeddings for text readability assessment](#). *International Conference on Information and Knowledge Management, Proceedings, Part F131841:2003–2006*.
- Meri Coleman and T. L. Liau. 1975. [A computer readability formula designed for machine scoring](#). *Journal of Applied Psychology*, 60:283–284.
- Scott Crossley, Aron Heintz, Joon Suh Choi, Jordan Batchelor, Mehrnosh Karimi, and Agnes Malatinszky. 2023. [A large-scaled corpus for assessing text readability](#). *Behavior Research Methods*, 55:491–507.
- Edgar Dale and Jeanne S Chall. 1948. [A formula for predicting readability: Instructions](#). *Educational Research Bulletin*, 27:37–54.
- Jerwin Jan S. Damay, Gerard Jaime D. Lojico, Kimberly Amanda L. Lu, and Dex B. Tarantan. 2006. [Simtext: Text simplification of medical literature](#). *Bachelor's Theses*.
- DeepSeek-AI. 2025. [Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning](#).
- Ashwin Devaraj, Iain J. Marshall, Byron C. Wallace, and Junyi Jessy Li. 2021. [Paragraph-level simplification of medical texts](#). *NAACL-HLT 2021 - 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Proceedings of the Conference*, pages 4972–4984.
- Fernanda Ferreira. 2003. [The misinterpretation of non-canonical sentences](#). *Cognitive Psychology*, 47:164–203.
- Lorenzo Jaime Flores, Heyuan Huang, Kejian Shi, Sophie Chheang, and Arman Cohan. 2023. [Medical text simplification: Optimizing for readability with unlikelihood training and reranked beam search decoding](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 4859–4873. Association for Computational Linguistics.
- Edward Gibson. 1998. [Linguistic complexity: locality of syntactic dependencies](#). *Cognition*, 68:1–76.
- Edward Gibson. 2000. *The dependency locality theory: A distance-based theory of linguistic complexity*, pages 94–126. The MIT Press.
- Gondy, Kauchak David, Gu Yang, Colina Sonia, Yuan Nicole P, Revere Debra Kloehn Nicholas, and Leroy. 2018. [Improving consumer understanding of medical text: Development and validation of a new subsimplify algorithm to automatically generate term explanations in english and spanish](#). *J Med Internet Res*, 20:e10779.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. [The llama 3 herd of models](#).
- John A Hawkins. 2004. *Efficiency and Complexity in Grammars*. Oxford University Press.
- Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. [spacy: Industrial-strength natural language processing in python](#).
- Pedram Hosseini, Christopher Wolfe, Mona Diab, and David Broniatowski. 2022. [Gispy: A tool for measuring gist inference score in text](#). In *Proceedings*

- of the 4th Workshop of Narrative Understanding (WNU2022), pages 38–46. Association for Computational Linguistics.
- Yi-Sheng Hsu, Nils Feldhus, and Sherzod Hakimov. 2024. [Free-text rationale generation under readability level control](#).
- Edward Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. [Lora: Low-rank adaptation of large language models](#). *ICLR 2022 - 10th International Conference on Learning Representations*.
- Yichen Huang and Ekaterina Kochmar. 2024. [Referee: A reference-free model-based metric for text simplification](#).
- Chao Jiang and Wei Xu. 2024. [Medreadme: A systematic study for fine-grained sentence readability in medical domain](#).
- Marcel Adam Just and Patricia A. Carpenter. 1992. A capacity theory of comprehension: Individual differences in working memory. *Psychological Review*, 99:122–149.
- Sasikiran Kandula, Dorothy Curtis, and Qing Zeng-Treitler. 2010. [A semantic and syntactic text simplification tool for health content](#). *AMIA Annual Symposium Proceedings*, 2010:366.
- Alla Keselman, Tony Tse, Jon Crowell, Allen Browne, Long Ngo, and Qing Zeng. 2007. [Assessing consumer health vocabulary familiarity: an exploratory study](#). *Journal of medical Internet research*, 9.
- J. P. Kincaid, Jr. Fishburne, Rogers Robert P., Chissom Richard L., and Brad S. 1975. [Derivation of new readability formulas \(automated readability index, fog count and flesch reading ease formula\) for navy enlisted personnel](#).
- Anastasia Krithara, James G Mork, Anastasios Nentidis, and Georgios Paliouras. 2023. [The road from manual to automatic semantic indexing of biomedical literature: a 10 years journey](#). *Frontiers in Research Metrics and Analytics*, 8.
- Philippe Laban, Tobias Schnabel, Paul N. Bennett, and Marti A. Hearst. 2022. [Summac: Re-visiting nli-based models for inconsistency detection in summarization](#). *Transactions of the Association for Computational Linguistics*, 10:163–177.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. [Albert: A lite bert for self-supervised learning of language representations](#). *CoRR*, abs/1909.11942.
- Charles Laurin, Dorret Boomsma, Gitta Lubke, Stat Appl, Genet Mol, and Biol Author. 2016. [The use of vector bootstrapping to improve variable selection precision in lasso models](#). *Statistical applications in genetics and molecular biology*, 15:305.
- Zihao Li, Samuel Belkadi, Nicolo Micheletti, Lifeng Han, Matthew Shardlow, and Goran Nenadic. 2024. [Large language models for biomedical text simplification: Promising but not there yet](#).
- Chin-Yew Lin. 2004. [Rouge: A package for automatic evaluation of summaries](#).
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, Veselin Stoyanov, and Paul G Allen. 2019. [Roberta: A robustly optimized bert pretraining approach](#).
- Junru Lu, Jiazheng Li, Byron C. Wallace, Yulan He, and Gabriele Pergola. 2023. [Napss: Paragraph-level medical text simplification via narrative prompting and sentence-matching summarization](#). *EACL 2023 - 17th Conference of the European Chapter of the Association for Computational Linguistics, Findings of EACL 2023*, pages 1049–1061.
- Zheheng Luo, Qianqian Xie, and Sophia Ananiadou. 2022. [Readability controllable biomedical document summarization](#). *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 4667–4680.
- Chen Lyu and Gabriele Pergola. 2024. [Scigispy: a novel metric for biomedical text simplification via gist inference score](#). *TSAR 2024 - 3rd Workshop on Text Simplification, Accessibility and Readability, Proceedings of the Workshop*, pages 95–106.
- Agnes Malatinszky, Aron Heintz, asiegel, Heather Harris, J S Choi, Maggie, Phil Culliton, and Scott Crossley. 2021. Commonlit readability prize. Kaggle.
- Christopher Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven Bethard, and David McClosky. 2014. [The stanford corenlp natural language processing toolkit](#). In *Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 55–60. Association for Computational Linguistics.
- Philip M. McCarthy and Scoot Jarvis. 2010. [MtlD, vocd-d, and hd-d: A validation study of sophisticated approaches to lexical diversity assessment](#). *Behavior Research Methods*, 42:381–392.
- Harry G McLaughlin. 1969. Smog grading - a new readability formula. *Journal of Reading*, pages 639–646.
- Juhi M Mohile, Joan B Luzon, Gunjan Agrawal, Neha R Malhotra, and Kathleen M Kan. 2023. [Assessment of readability and quality of patient education materials specific to nocturnal enuresis](#). *Journal of Pediatric Urology*, 19:558.e1–558.e7.
- National Library of Medicine. 2024. [Umls knowledge sources](#). Cited 2025 Apr 7.
- OpenAI. 2023. Gpt-4 technical report.

- Shon Otmazgin, Arie Cattan, and Yoav Goldberg. 2022. [F-coref: Fast, accurate and easy to use coreference resolution](#).
- Malaikannan Sankarasubbu Ankit Pal. 2024. [Openbi-ollms: Advancing open-source large language models for healthcare and life sciences](#).
- Jürgen Pelikan, Christa Straßmayr, Thomas Link, Dominika Miksova, Peter Nowak, Robert Griebler, Christina Dietscher, Stephan den Broucke, Rana Charafeddine, Antoniya Yanakieva, Nygyar Dzhafer, Zdenek Kučera, Alena Šteflová, Henrik Bøggild, Andreas Sørensen, Julien Mancini, Geneviève Chêne, Doris Schaeffer, Alexander Schmidt-Gernig, and Øystein Guttersrud. 2021. [International report on the methodology, results, and recommendations of the european health literacy population survey 2019-2021 \(hls19\) of m-pohl](#).
- Atharva Phatak, David W. Savage, Robert Ohle, Jonathan Smith, and Vijay Mago. 2022. [Medical text simplification using reinforcement learning \(teslea\): Deep learning-based text simplification approach](#). *JMIR Medical Informatics*, 10.
- Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D Manning. 2020. [Stanza: A python natural language processing toolkit for many human languages](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 101–108. Association for Computational Linguistics.
- Md Mushfiqur Rahman, Mohammad Sabik Irbaz, Kai North, Michelle S. Williams, Marcos Zampieri, and Kevin Lybarger. 2024. [Health text simplification: An annotated corpus for digestive cancer education and novel strategies for reinforcement learning](#). *Journal of Biomedical Informatics*, 158:104727.
- Kamal raj Kanakarajan, Bhuvana Kundumani, Abhijith Abraham, and Malaikannan Sankarasubbu. 2022. [Biosimcse: Biomedical sentence embeddings using contrastive learning](#). In *Proceedings of the 13th International Workshop on Health Text Mining and Information Analysis (LOUHI)*, pages 81–86. Association for Computational Linguistics.
- Valerie F. Reyna. 2012. [A new intuitionism: Meaning, memory, and development in fuzzy-trace theory](#). *Judgment and decision making*, 7:332.
- Leonardo F.R. Ribeiro, Mohit Bansal, and Markus Dreyer. 2023. [Generating summaries with controllable readability levels](#). *EMNLP 2023 - 2023 Conference on Empirical Methods in Natural Language Processing, Proceedings*, pages 11669–11687.
- Max Savery, Asma Ben Abacha, Soumya Gayen, and Dina Demner-Fushman. 2020. [Question-driven summarization of answers to consumer health questions](#). *Scientific Data*, 7:1–9.
- Thibault Sellam, Dipanjan Das, and Ankur P. Parikh. 2020. [Bleurt: Learning robust metrics for text generation](#). *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892.
- Rabia Shahid, Muhammad Shoker, Luan Manh Chu, Ryan Frehlick, Heather Ward, and Punam Pahwa. 2022. [Impact of low health literacy on patients' health outcomes: a multicenter cohort study](#). *BMC Health Services Research*, 22:1–9.
- Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, Perry Payne, Martin Seneviratne, Paul Gamble, Chris Kelly, Abubakr Babiker, Nathanael Schärli, Aakanksha Chowdhery, Philip Mansfield, Dina Demner-Fushman, and 13 others. 2023. [Large language models encode clinical knowledge](#). *Nature*, 620:172–180.
- E A Smith and R J Senter. 1967. Automated readability index. Technical report, Aerospace Medical Research Laboratories (U.S.), Wright-Patterson Air Force Base, Ohio. PMID: 5302480.
- Luca Soldaini. 2016. [Quickumls: a fast, unsupervised approach for medical concept extraction](#).
- T Szmuda, C Özdemir, S Ali, A Singh, M T Syed, and P Stoniewski. 2020. [Readability of online patient education material for the novel coronavirus disease \(covid-19\): a cross-sectional health literacy study](#). *Public Health*, 185:21–25.
- Hieu Tran, Zonghai Yao, Lingxi Li, and Hong Yu. 2024. [Readctrl: Personalizing text generation with readability-controlled instruction learning](#).
- Xun Wang and Robin A Cohen. 2022. [Health information technology use among adults: United states, july-december 2022 key findings data from the national health interview survey](#). Technical report.
- Zhilin Wang, Alexander Bukharin, Olivier Delalleau, Daniel Egert, Gerald Shen, Jiaqi Zeng, Oleksii Kuchaiev, and Yi Dong. 2024. [Helpsteer2-preference: Complementing ratings with preferences](#).
- Nikki Keene Woods, Umama Ali, Melissa Medina, Jared Reyes, and Amy K. Chesser. 2023. [Health literacy, health outcomes and equity: A trend analysis based on a population survey](#). *Journal of Primary Care Community Health*, 14.
- Biyang Yu, Zhe He, Aiwen Xing, and Mia Liza A. Lustria. 2020. [An informatics framework to assess consumer health language complexity differences: Proof-of-concept study](#). *Journal of medical Internet research*, 22.
- Qing T. Zeng and Tony Tse. 2006. [Exploring and developing consumer health vocabularies](#). *Journal of the American Medical Informatics Association : JAMIA*, 13:24.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2019. [Bertscore: Evaluating text generation with bert](#). *8th International Conference on Learning Representations, ICLR 2020*.

Jiaping Zheng and Hong Yu. 2018. [Assessing the readability of medical documents: A ranking approach](#). *JMIR Medical Informatics*, 20.

Ming Zhong, Yang Liu, Da Yin, Yuning Mao, Yizhu Jiao, Pengfei Liu, Chenguang Zhu, Heng Ji, and Jiawei Han. 2022. [Towards a unified multi-dimensional evaluator for text generation](#). *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022*, pages 2023–2038.

A Datasets

This appendix provides additional details about the datasets used in our study, including the medical text simplification datasets used to validate our evaluation metrics and the question-answer datasets used to build our multi-level medical QA corpus.

A.1 Medical Text Simplification Datasets

We used two publicly available datasets of simplified medical texts to support the evaluation of our complexity metrics and train our formula:

- **PLABA** (Attal et al., 2023): The PLABA dataset contains 750 biomedical abstracts that have been rewritten in plain language, totaling 7,643 sentence pairs. It was created by scraping 75 common medical questions from MedlinePlus and retrieving relevant paper abstracts from PubMed. Human annotators then simplified these abstracts by replacing technical terms with familiar synonyms (e.g., “orthosis” to “brace”), breaking down complex sentences, and removing content that might not be relevant to a general audience. We used PLABA at both sentence and paragraph levels to evaluate our complexity metrics.
- **Cochrane Dataset** (Devaraj et al., 2021): The Cochrane Simplification dataset contains 4,459 pairs of technical medical texts and their simplified versions, sourced from the Cochrane Database of Systematic Reviews. These paragraph-level simplifications are derived from pls written for readers without a university education and involve a mix of paraphrasing, deletion, and summarization to make the original texts more accessible.

A.2 Source Medical QA Datasets

To create our multi-level medical QA corpus, we combined samples from five existing datasets that represent a range of medical topics, question styles, and answer formats:

- **LiveQA** (Abacha et al., 2017): The LiveQA dataset includes real-world consumer health questions submitted to the U.S. nlm during the TREC 2017 LiveQA challenge. The original release had 634 training pairs and 104 test questions, each with multiple reference answers. After cleaning the data, we retained 800 question-answer pairs covering topics such as diseases, treatments, medications, and medical exams.
- **MedicationQA** (Abacha et al., 2019): The MedicationQA dataset contains 690 consumer questions about medications, each paired with an answer from a trusted medical website, such as MedlinePlus and DailyMed, addressing topics like drug usage, dosage, side effects, and drug interactions.
- **MediQA-AnS** (Savery et al., 2020): The MediQA-AnS dataset, created for the MEDIQA 2021 challenge, includes 156 consumer health questions, each paired with two reference summaries (abstractive and extractive) both written by medical experts based on passages retrieved using the CHiQA system Demner-Fushman2020.
- **MedQuAD** (Abacha and Demner-Fushman, 2019): The Medical Question Answering Dataset consists of 47,457 question-answer pairs sourced from 12 websites managed by the U.S. National Institutes of Health (NIH), including MedlinePlus, cancer.gov, and niddk.nih.gov. Due to copyright restrictions, we had to exclude over 31,000 entries, leaving us with a total of 16,423 samples.
- **BioASQ** (Krithara et al., 2023): The BioASQ Task 13B dataset, part of the 2025 BioASQ challenge, includes 5,389 biomedical questions. Each question is paired with one or more ideal answers, resulting in a total of 13,692 question-answer pairs. These answers are concise, expert-written summaries that draw from scientific literature, primarily PubMed, and use precise biomedical terminology.

B Results and Performance

This appendix provides additional performance details and supporting results for the main experiments described in the paper.

B.1 Selected Features

The following 13 metrics were selected for our final complexity scoring formula, listed here along with their coefficients.

B.1.1 LLM Vocabulary Complexity (3.217)

We use this score to estimate how difficult a piece of text is to understand, based on evaluations from the three 70-billion parameter models we described earlier. Each model rated texts on a scale from 1 to 5, with higher scores indicating more complex language. This feature has the largest positive coefficient in our formula, confirming that vocabulary choice drives most of the perceived difficulty in medical texts.

B.1.2 Dale-Chall Score (1.839)

The Dale-Chall readability formula (Dale and Chall, 1948) estimates how difficult a text is to read based on the average sentence length and the percentage of “difficult” words not found on a pre-defined list of familiar words. In our implementation, we expanded the original list of 3,000 words by including those from the Spache list. The positive coefficient in the formula shows that texts with longer sentences and more unfamiliar words tend to be significantly more complex.

B.1.3 Type-Token Ratio (0.173)

Type-token ratio (TTR) measures lexical diversity by dividing the number of unique words (types) by the total number of words (tokens) in a text. The positive coefficient confirms that texts with more diverse vocabulary contribute to higher complexity scores, though with less impact than the vocabulary complexity or the Dale-Chall readability formula.

B.1.4 ALBERT Transformer Score (-2.471)

We used the ALBERT-xxlarge model (Lan et al., 2019) from the winning entry in the CommonLit Readability Prize Kaggle competition (Malatinszky et al., 2021). This model processes text through attention layers to capture relationships between words before predicting a readability score. The negative coefficient appears because ALBERT assigns higher scores to texts that are easier to read, which runs in the opposite direction of our scoring

system, where higher values indicate lower readability.

B.1.5 Referential Cohesion (0.068)

This feature captures how well a paragraph maintains topical consistency by measuring the semantic similarity between consecutive sentences (Lyu and Pergola, 2024). To compute it, we embed each sentence using BioSimCSE-BioLinkBERT-BASE (raj Kanakarajan et al., 2022) and calculate the cosine similarity between adjacent sentence pairs. A sharp drop in similarity, falling in the bottom 25% of the distribution, marks a potential topic shift, or “breakpoint.” We count the number of chunks in each paragraph based on these breakpoints and take the average over the entire text. The small positive coefficient may seem counterintuitive, since texts that are more cohesive are usually easier to read. However, this result suggests that even highly technical medical texts in our datasets tend to maintain strong internal cohesion despite their complex vocabulary.

B.1.6 Information Content (0.691)

This feature measures how specialized the vocabulary is in a given text, based on how often each word appears in a biomedical corpus (Lyu and Pergola, 2024). The basic idea is that technical terms tend to be rarer and harder to understand. To build our reference corpus, we combined data from biomedical and consumer health sources, including MedQuAD (Abacha and Demner-Fushman, 2019), LiveQA (Abacha et al., 2017), MedicationQA (Abacha et al., 2019), and other medical datasets. We lemmatize each word in the corpus, count how often each lemma appears, and calculate its information content as the negative logarithm of its probability. For any given text, we extract all nouns and verbs, look up their information content values, and calculate the average. The positive coefficient in our model supports the idea that texts with a more technical and less common vocabulary tend to be more complex.

B.1.7 Verb Ratio (-0.330)

Part-of-speech distributions measure the frequency of different grammatical categories relative to the total word count. We calculate separate ratios for nouns, verbs, adjectives, adverbs, conjunctions, and auxiliary verbs. A negative coefficient for verb ratio indicates that texts with fewer verbs relative to other parts of speech are rated as more complex. This is consistent with research showing that

academic and scientific writing tends to use more nouns and fewer verbs (Biber et al., 1999).

B.1.8 Function Word Ratio (-0.596)

The content-to-function word ratio calculates the proportion of content words (nouns, verbs, adjectives, adverbs) to function words (auxiliaries, determiners, prepositions, conjunctions) in a text. A negative coefficient means that texts with more content words and fewer function words are seen as more complex. This is because function words help organize sentence structure, so when they are used less frequently, the resulting text can be more syntactically dense and cognitively demanding for readers (Just and Carpenter, 1992).

B.1.9 Masked Probability Score (-0.049)

This metric evaluates how predictable words are in biomedical text using a masked language model Devaraj et al. (2021).. Specifically, we randomly mask 15% of the tokens and run this process 30 times, then measure how accurately Bio+ClinicalBERT can guess the original words. In general, technical or scientific writing tends to have more predictable language patterns, especially due to consistent use of domain-specific terms. The negative weight in the scoring formula helps balance out other vocabulary-based metrics. It prevents penalizing texts that are technically dense but still internally consistent and readable.

B.1.10 MedReadMe Cluster Score (0.295)

This score comes from a clustering-based word embedding model trained on the MedReadMe dataset (Jiang and Xu, 2024). Each word in the text is converted into its BioWordVec embedding, then assigned to one of 300 semantic clusters using K-means clustering. The pattern of these assignments forms a feature vector that represents how the vocabulary is distributed across different semantic categories. A positive coefficient means that texts using vocabulary patterns similar to those found in more complex medical content tend to receive higher complexity scores.

B.1.11 Embedding Depth (-0.161)

Embedding depth measures how deep the hierarchical structure of a sentence goes in its dependency tree. To calculate this, we identify the word with the longest chain of grammatical dependencies leading to the root of the sentence. A sentence with greater embedding depth usually contains more subordinate clauses (introduced by words like

“which,” “that,” “when”) and complex phrases embedded within one another. This typically makes text harder to process, as readers must track multiple incomplete grammatical relationships while reading, increasing cognitive effort (Gibson, 1998). However, in our corpus, the expert texts often used more concise, noun-heavy sentences with fewer nested clauses. In contrast, the simpler texts used more explanatory language with embedded clauses to break down complex concepts. This pattern explains the negative coefficient in our formula.

B.1.12 Average Dependency Distance (-0.826)

Dependency distance measures how many words separate a dependent word (object or modifier) from its head word (main verb or noun) in a sentence. Longer distances increase cognitive load, since the reader must keep track of the dependent word while processing the words in between (Gibson, 2000). We calculate the average dependency distance for each sentence and then find the overall average for the entire text. Although this metric correlates with higher difficulty when used alone, the negative coefficient in our multivariate model suggests an inverse relationship when considered alongside other features.

B.1.13 Coreference Chains (-0.390)

Coreference resolution tracks how entities are referenced throughout a text. When a document refers to the same person, object, or concept using different terms (e.g., pronouns, synonyms, or descriptions), it creates coreference chains that help readers follow who or what is being discussed. For instance, if a text mentions “Dr. Smith” and later refers to her as “she” or “the physician,” these references form a continuous link to the same entity. To calculate CoREF, we use FastCoref (Otmazgin et al., 2022) instead of the Stanford CoreNLP implementation previously used in GisPy (Manning et al., 2014; Qi et al., 2020). We decided to make this switch because CoreNLP was causing significant delays in the processing pipeline, especially when working with longer documents. FastCoref, on the other hand, not only performs on par with state-of-the-art models but also runs much faster, completing tasks in seconds that used to take minutes. Following the same methodology as GisPy, we identify all coreference chains in a document, calculate the ratio of chains to sentences for each paragraph, and then compute the final CoREF score as the average of these paragraph-level scores. The negative coef-

ficient indicates that complex medical texts often contain fewer or shorter coreference chains, introducing new entities without established reference patterns, which increases reading difficulty.

B.2 Quantitative Model Performance

The quantitative results in Table 4 confirm what we see in Figure 3. The fine-tuned model outperforms both alternatives in every metric, with a mean absolute error (MAE) 23% lower than the few-shot method and nearly 50% lower than the baseline. The strong correlation coefficient (0.84) and high R^2 value (0.66) together validate its ability to consistently generate responses at the intended complexity level.

Model	MAE	RMSE	Correlation	R^2
Baseline	26.07	31.28	0.21	-0.07
Few-shot	17.33	22.03	0.69	0.47
Fine-tuned	13.30	17.63	0.84	0.66

Table 4: Comparison of how accurately each model generates text at the desired complexity levels.

C Prompts

This appendix compiles the complete set of prompts used throughout this work. These prompts were integral to various stages of our research, from dataset generation to text complexity evaluation, and include placeholders for dynamic content that was filled in during the actual runs.

C.1 Prompt for Generating the HSQA-Claude Dataset

This prompt was used to generate the HSQA-Claude dataset, introduced in Section 3.1.2. It provides detailed instructions for generating expert-level and patient-friendly answers to health-related questions, handling ambiguous or off-topic questions, correcting grammatical issues, and formatting the output as a JSON array. The list of questions to be answered is represented by the placeholder [QUESTION_LIST].

You are providing two types of answers to health-related questions:

1. An expert answer written as if for medical professionals (like in clinical documentation or medical education)
2. A patient-friendly answer written as if for a medical forum or patient consultation

IMPORTANT INSTRUCTIONS:

1. For questions that don't immediately appear health-related:
 - If there's any possible health interpretation, treat it as a health question
 - Mark as "(wrong topic)" in the question field if you are confident it has no health relevance
 - Strive to provide a health-related answer even if the question seems unusual
 - Example: "How do you make an IO game?" is clearly not health-related
 - Example: "How do I make a paste?" could be about medical adhesives or food preparation for special diets, so treat as health-related
2. For questions with spelling or grammar issues:
 - Fix any grammatical errors in questions while preserving their meaning
 - Add missing articles (a, an, the) where needed
 - Correct subject-verb agreement
 - Improve clarity but maintain the original intent
 - Example: "Is jaundice can be cured?" "Can jaundice be cured?"
 - Example: "Is every white patch is vitiligo?" "Is every white patch vitiligo?"

Make sure to answer every unique question in the provided order.

Questions:
[QUESTION_LIST]

General guidelines for all answers:

1. Vary response style naturally – avoid rigid templates or repetitive structures
2. Match answer length to the topic's complexity – some need more context, others can be brief
3. Expert answers don't need to be longer than simple ones – focus on clarity and accuracy
4. Adapt detail level to the specific question and context
5. Ensure information is accurate and factual
6. Avoid overused phrases or patterns in medical writing
7. Structure responses logically and coherently

Guidelines specific to expert answers:

1. Write in clinical documentation style using precise medical terminology
2. Include key differential diagnoses when relevant
3. Discuss diagnostic criteria and clinical presentations

4. Mention standard treatment approaches and clinical decision-making factors
5. Include relevant quantitative information (rates, thresholds, timeframes)
6. Focus on assessment and management considerations
7. Use professional medical syntax and phrasing

Guidelines specific to patient-friendly answers:

1. Use clear, accessible language without medical jargon
2. Explain concepts in practical terms
3. Address common concerns and misconceptions
4. Include appropriate reassurance while being honest about risks
5. Use analogies ONLY when the concept is complex and would genuinely benefit from one
6. Focus on practical implications and self-care when relevant

Format each Q&A pair as:

```
{
  "question": "The health question",
  "expert_answer": "Clinical-style medical explanation",
  "simple_answer": "Patient-friendly explanation"
}
```

The complete response should be a JSON array:

```
{
  "qa_pairs": [
    // Q&A pairs here
  ]
}
```

Return only valid JSON with no additional text.

C.2 Prompt for Evaluating Text Complexity

This is the prompt used in Section 3.2.8, where we evaluate the complexity of medical texts using pre-trained language models. It asks the model to rate a given text on five different dimensions of complexity, using a scale from 1 to 5, and to provide a brief explanation for each rating. The text to be evaluated is placed at [TEXT], and the model is guided by three annotated examples inserted into the placeholders [EXAMPLE_1_TEXT], [EXAMPLE_1_EVALUATION], and so forth. We initially tried using JSON for the output format, but since the models often generated invalid JSON, we switched to XML because it is easier to parse and less error-prone.

You are an expert in evaluating the readability and complexity of texts.

Your task is to assess the given text on several dimensions using a scale from 1 to 5, where 1 is the simplest and 5 is the most complex.

When evaluating the text, you must:

1. Assess each dimension independently using the defined 5-level scale.
2. Provide a brief reasoning for your assessment.
3. Ensure your evaluation is consistent and well-justified.

The five dimensions and their levels (1 to 5) are defined as follows:

****Vocabulary Complexity**:**

- 1: Very basic words, suitable for young children.
- 2: Simple words, understandable by most adults.
- 3: Moderate vocabulary, including some technical terms.
- 4: Advanced vocabulary, with specialized terms.
- 5: Highly technical or specialized vocabulary, requiring expert knowledge.

****Syntactic Complexity**:**

- 1: Very simple sentence structures, short sentences.
- 2: Basic sentence structures, mostly simple and compound sentences.
- 3: Moderate complexity with a mix of simple and complex sentences.
- 4: Complex sentence structures, with subordinate clauses and intricate syntax.
- 5: Highly complex syntax, with nested clauses and sophisticated constructions.

****Conceptual Density**:**

- 1: Single, straightforward ideas presented one at a time.
- 2: Few related concepts introduced at a manageable pace.
- 3: Multiple concepts with clear connections between them.
- 4: Many interrelated concepts requiring careful attention to follow.
- 5: Dense with numerous abstract and interrelated concepts.

****Background Knowledge**:**

- 1: No special knowledge needed beyond everyday experience.
- 2: Basic familiarity with the subject area.
- 3: General education in the domain or field discussed.
- 4: Considerable domain knowledge required.
- 5: Expert-level knowledge in the field necessary.

****Cognitive Load**:**

- 1: Minimal effort to process and understand.
- 2: Some attention needed but generally easy.
- 3: Requires focus and moderate effort.
- 4: Demands concentration and significant mental effort.
- 5: Requires sustained intense concentration and analytical thinking.

Below are examples to guide your assessment:

Example 1
Text: [EXAMPLE_1_TEXT]
[EXAMPLE_1_EVALUATION]

Example 2
Text: [EXAMPLE_2_TEXT]
[EXAMPLE_2_EVALUATION]

Example 3
Text: [EXAMPLE_3_TEXT]
[EXAMPLE_3_EVALUATION]

Now, evaluate the following text:

[TEXT]

Place your response between <root> and </root> tags in exactly this format:
<root>
<vocabulary_complexity>score</vocabulary_complexity>
<syntactic_complexity>score</syntactic_complexity>
<conceptual_density>score</conceptual_density>
<background_knowledge>score</background_knowledge>
<cognitive_load>score</cognitive_load>
<reasoning>brief explanation</reasoning>
</root>

Only include scores as integers between 1 and 5 within the tags.
Ensure each tag is properly closed with the corresponding closing tag.
Do not include any additional text outside the <root> and </root> tags.
Use only the specified XML format.

C.3 Prompt for Generating Answer Variants

This is the prompt used to generate the answer variants for the multi-level dataset described in Section 3.4.2. It takes a question and a reference answer, then generates a series of variants written for audiences with increasing levels of background knowledge. The total number of variants, along with the question and original answer, are inserted into the placeholders [NUM_VARIANTS], [QUESTION], and [ORIGINAL_ANSWER]. We also

provide three in-context examples and use XML as the output format for the same reasons discussed in the previous prompt.

You are an expert in creating educational content for different reading abilities. Your task is to generate multiple answer variants for the given question and original answer, each at a specified complexity level, while preserving all factual information.

When generating each variant, you must:

1. Preserve ALL factual information from the original answer and keep it relevant to the question.
2. Adjust vocabulary, sentence structure, and explanation detail to match the complexity level.
3. Do not introduce substantively new claims that aren't reasonably implied by the original answer.
4. Ensure the answer is coherent and well-structured.
5. If the original answer does not directly address the question asked, respond with: '[CONTENT_MISMATCH]' as the answer.

Complexity levels (1 to 5) are defined as follows:

- 1: For a young child; use very simple vocabulary, short sentences, and basic concepts.
- 2: For a middle school student; use basic scientific terms, clear explanations, and moderate detail.
- 3: For a high school student; use technical terminology, longer sentences, and detailed explanations.
- 4: For a college graduate; use in-depth technical details, complex sentence structures, and scientific language.
- 5: For a biomedical expert; use advanced scientific terminology, assume prior knowledge, and provide precise details.

Below are examples of how to adjust answers by complexity:

Example 1
Question: [EXAMPLE_1_QUESTION]
Original Answer: [EXAMPLE_1_ANSWER]
[EXAMPLE_1_VARIANTS]

Example 2
Question: [EXAMPLE_2_QUESTION]
Original Answer: [EXAMPLE_2_ANSWER]
[EXAMPLE_2_VARIANTS]

Example 3
Question: [EXAMPLE_3_QUESTION]
Original Answer: [EXAMPLE_3_ANSWER]
[EXAMPLE_3_VARIANTS]

Now, generate [NUM_VARIANTS] answer variants for the following question and original answer. The variants should be ordered from the simplest to the most complex, reflecting a gradual increase in complexity. For each variant, assign a complexity level from 1 to 5, where 1 is the simplest and 5 is the most complex, based on the definitions provided.

Question: [QUESTION]
Original Answer: [ORIGINAL_ANSWER]

Place your response between <root> and </root> tags in exactly this format:

```
<root>
  <variant>
    <complexity_level>1</complexity_level>
    <answer>{Your first answer goes here}</answer>
  </variant>
  <variant>
    <complexity_level>2</complexity_level>
    <answer>{Your next answer goes here}</answer>
  </variant>
  ...
</root>
```

Ensure EACH variant has both <complexity_level> and <answer> tags. Each tag must be properly closed with the corresponding closing tag. Do not include any additional text outside the <root> and </root> tags. Use only the specified XML format.

C.4 Prompt for Evaluating Medical Responses Using Simulated Personas

This prompt, used in Section 4.3, simulates how individuals with varying levels of health literacy interpret and rate responses based on five predefined quality dimensions. The specific question and its corresponding answer are dynamically inserted into the prompt at the [QUESTION] and [ANSWER] placeholders, while the background of the simulated user is inserted at the start of the prompt in place of [SIMULATED_USER_BACKGROUND].

The full evaluation prompt is shown below.

[SIMULATED_USER_BACKGROUND]

You must evaluate the medical answer strictly from your own perspective and level of health literacy. Do not try to judge it from a general or professional viewpoint unless that matches your background.

Your task is to score the answer across five dimensions on a scale from 1 to 5, where 1 is the lowest and 5 is the highest.

The five dimensions and their levels (1 to 5) are defined as follows:

****Understandability**:**

- 1: Very difficult to understand, confusing language or concepts
- 2: Somewhat difficult, requires effort to follow
- 3: Moderately understandable, generally clear
- 4: Easy to understand, well-explained concepts
- 5: Extremely clear and accessible for the intended audience

****Usefulness**:**

- 1: Not helpful, lacks practical value
- 2: Minimally helpful, limited practical application
- 3: Moderately useful, provides some actionable information
- 4: Very useful, offers clear guidance or valuable insights
- 5: Extremely useful, highly actionable and comprehensive

****Clarity**:**

- 1: Very confusing, many unclear or ambiguous parts
- 2: Somewhat confusing, several unclear elements
- 3: Generally clear with minor confusing aspects
- 4: Clear and well-structured, easy to follow
- 5: Exceptionally clear, no confusing elements

****Relevance**:**

- 1: Does not address the question, completely off-topic
- 2: Minimally relevant, partially addresses the question
- 3: Moderately relevant, addresses main aspects of the question
- 4: Highly relevant, directly addresses the question well
- 5: Perfectly relevant, comprehensively addresses all aspects

****Factuality**:**

- 1: Contains significant medical inaccuracies or misinformation
- 2: Contains some questionable or potentially inaccurate information
- 3: Generally accurate with minor issues or omissions
- 4: Medically accurate and reliable information
- 5: Completely accurate, evidence-based, and up-to-date

Question: [QUESTION]

Answer: [ANSWER]

Respond using only the following JSON format, without any additional text

```

    or explanations:
    ```json
 {
 "reasoning": "Brief reasoning for
 scores (optional, not required)
 ",
 "understandability": 1-5,
 "usefulness": 1-5,
 "clarity": 1-5,
 "relevance": 1-5,
 "factuality": 1-5
 }
    ```

```

C.5 Persona Definitions

Each evaluation was run using one of the following user personas, inserted at the [SIMULATED_USER_BACKGROUND] placeholder in the prompt above:

- **Low Health Literacy:** You are a person with low health literacy evaluating medical information. You have no medical training and rely on everyday language to understand health topics. You struggle with medical jargon and need simple, clear explanations.
- **Medium Health Literacy:** You are a person with moderate health literacy evaluating medical information. You have some familiarity with common medical terms through personal experience, general education, or caring for family members. You can understand basic medical concepts but may struggle with highly technical information.
- **High Health Literacy:** You are a healthcare professional or medical student evaluating medical information. You have extensive medical training and are comfortable with medical terminology, clinical concepts, and evidence-based practice.

D Model Output Examples

This appendix presents example responses generated by our fine-tuned language model for the medical question “Can asthma be cured?” across five different complexity levels (see Table 5). Each response was generated using a specific control token to target the desired complexity level, ranging from 0 (most accessible) to 100 (most technical). These examples demonstrate how the model adapts its language, terminology, and depth of explanation based on the specified complexity target while maintaining medical accuracy throughout all levels.

E Limitations and Future Work

While our work has made meaningful progress in simplifying medical texts, it also has some important limitations.

First, we focused only on English. The features we used to measure complexity, especially those tied to medical terms, may not translate well to other languages that have different grammar rules or naming conventions in medicine.

Second, we developed our complexity formula without using human feedback. Instead, we assumed that the best formula is the one that maximizes the gap between simple and complex texts, following a set of heuristics we defined based on our understanding of the data. However, perceived complexity is subjective and can vary depending on a person’s background, reading ability, and familiarity with medical concepts. Therefore, testing with real users would be necessary to confirm if the formula aligns with human judgments. Furthermore, because the formula is a simple linear equation, it can be “gamed.” For example, shortening sentences or swapping in simpler words could reduce the complexity score without making the text any easier to understand. A possible solution would be to follow the same approach we used in Section 3.2.8, distilling the scores into a language model to smooth errors and reduce the impact of outliers. We chose a simple, interpretable formula here, but future work should explore more flexible, non-linear models.

Another concern is the use of synthetic data for training. Even with filters in place, the dataset might still include outdated or inaccurate medical information carried over from the original sources. We also did not evaluate our model against external benchmarks or state-of-the-art systems like GPT-4 or Claude, and lacked specialized datasets to measure factual accuracy and relevance when controlling the complexity of the answers. More importantly, the model was not tested longitudinally with actual users.

Finally, we did not explore alternative methods for measuring text complexity. For example, training models to predict which of two texts is more complex (as in learning-to-rank frameworks) rather than assigning absolute scores could be an alternative approach to evaluate text readability.

Control Code	Generated Output
<COMPLEXITY_0>	Asthma can't be fully "cured," but you can manage it. Treatments help control symptoms and prevent attacks, making life easier for those with the condition.
<COMPLEXITY_25>	No, there is currently no cure for asthma. However, symptoms can often be managed with the use of medications and other treatments that help control inflammation in your lungs so you don't experience as many attacks. You'll likely need to work closely with a doctor or healthcare provider to develop an action plan to manage your condition effectively. Some people have periods when they are free from symptoms (remission), but this doesn't mean their underlying disease has gone away. If left untreated, asthma will always get worse over time; if not treated properly, it may cause permanent damage to lung tissue, which could lead to chronic obstructive pulmonary disease (COPD). While it's possible to make lifestyle changes such as quitting smoking or avoiding allergens to reduce flare-ups, managing asthma requires ongoing medical attention.
<COMPLEXITY_50>	Asthma is a chronic condition with no cure, but its symptoms can be effectively managed through medication and lifestyle adjustments. The primary goal of treatment is to prevent exacerbations by maintaining optimal control over the disease process. This involves using inhaled corticosteroids as maintenance therapy for most patients, alongside other medications such as bronchodilators or biologics based on specific clinical needs. By adhering strictly to prescribed treatments under medical supervision, individuals with asthma can significantly reduce their risk of acute attacks, minimize disruptions to daily life, and improve overall quality of living despite ongoing management requirements.
<COMPLEXITY_75>	Asthma is a chronic inflammatory airway disease characterized by hyperresponsiveness to various stimuli, leading to recurrent episodes of wheezing and dyspnea. While the condition itself cannot be definitively "cured" due to its inherent pathophysiology involving complex interactions between environmental triggers, genetic predisposition, immune system modulation, and structural changes in the lung parenchyma, therapeutic interventions can significantly mitigate symptoms and improve quality of life for affected individuals through targeted management strategies aimed at modulating inflammation, reducing bronchoconstriction, enhancing pulmonary function, and minimizing adverse effects on health-related outcomes such as exercise capacity or cognitive performance. These interventions encompass pharmacological agents (e.g., corticosteroids), non-pharmacologic modalities like physical activity and avoidance of irritants, along with advanced technologies including mechanical ventilation support during acute exacerbations, highlighting an integrated approach that balances symptom control with prevention-oriented care models to enhance overall well-being and reduce long-term morbidity associated with this common respiratory disorder.
<COMPLEXITY_100>	Asthma is a chronic inflammatory disorder characterized by reversible airflow obstruction, manifesting as episodic wheezing and dyspnea. While no definitive "cure" exists for asthma due to its complex pathophysiology involving multiple cell types and mediators in the airway mucosa, therapeutic strategies can effectively manage symptoms and prevent exacerbations through targeted interventions aimed at modulating inflammation, enhancing bronchodilation, reducing allergen sensitization, and improving patient adherence to treatment regimens. Notably, recent advances in precision medicine have enabled personalized approaches tailored to an individual's genetic predisposition, environmental exposures, and molecular phenotypic characteristics of their disease state, thereby optimizing outcomes across various subpopulations within this heterogeneous condition.

Table 5: Example of model responses at different complexity levels for the question "Can asthma be cured?"