

# Beyond Flesch-Kincaid: New Neural Metrics to Improve Difficulty Prediction for Educational Texts

Anonymous NAACL 2024 submission

## Abstract

Using large language models (LLMs) for educational applications such as dialogue-based teaching is a hot topic. Effective teaching, however, requires adapting the difficulty of content and explanations to the education level of their students. Even the best LLMs today struggle to do this well. If we want to improve LLMs on this adaptation task, we need to be able to reliably measure adaptation success. However, current *Static* metrics for text difficulty, like the Flesch-Kincaid Reading Ease score, are known to be crude and brittle. We therefore introduce and evaluate a new set of *neural* metrics for text difficulty. Based on a user study, we create Neural metrics as LLM prompts that leverage the general language understanding capabilities of LLMs to capture more abstract and complex text features than Static metrics. Through regression experiments, we show that our Neural metrics improve text difficulty prediction over Static metrics alone. Our results demonstrate the promise of Neural metrics as a new class of features for evaluating text adaptation to different education levels.

## 1 Introduction

Large language models (LLMs) today can answer wide-ranging questions and explain complex concepts with high accuracy (Chung et al., 2022)(OpenAI, 2023). This development has motivated explorations into their uses for education, ranging from automated student assessment and personalised content to dialogue-based teaching (Upadhyay et al., 2023; Sallam, 2023; Yan et al., 2023; Hosseini et al., 2023). Effective teaching requires that the difficulty of content and explanations is tailored to the education level of the students. Human teachers are trained to do this, and adjust their material and style without much prompting. However, this adaptation is not just the adjustment of one variable. It is a complex undertaking, touching upon lexicon, syntax, pragmatics, and semantics.

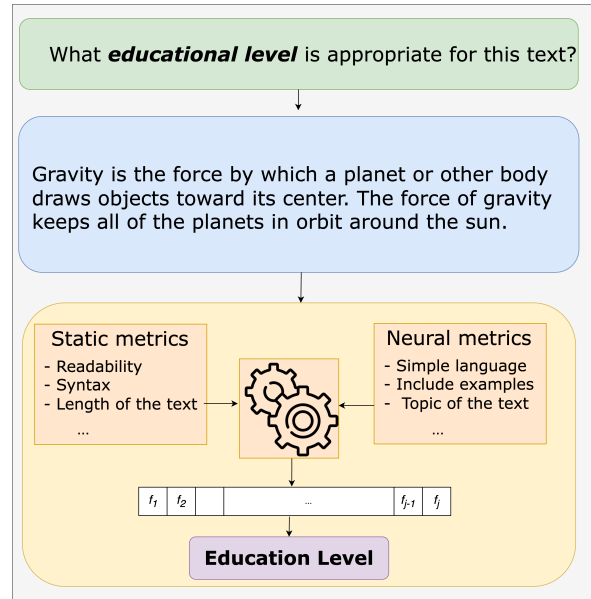


Figure 1: Schematic overview of our proposed approach to text difficulty classification. We calculate Static and Neural metrics relevant to text difficulty based on a given input text<sup>1</sup>. Either or both types of metrics are then fed into a regression classifier that makes a final prediction.

Consequently, even the best LLMs today struggle to provide appropriate answers for a given audience – they can adapt to different writing styles (casual or formal) and domains (emails, blog posts, and essays), but their ability to adapt to different age ranges or levels of education is much more limited (Kasneci et al., 2023). At best, they can differentiate between child and adult learners, but that distinction is rarely useful when adapting recurring material for different grades.

Improving the ability of LLMs to adapt their outputs to different levels of education is therefore crucial to unlocking their usefulness for education. One of the most basic requirements to achieve this goal is a way to measure adaptation success.

However, whether a given output is appropriate for a given level of education is still very hard to

measure. Existing *Static* metrics, like the Flesch-Kincaid Reading Ease score, are based on simple formulas, heuristics, and word counts. They share the brittleness of all heuristic approaches and are known to be noisy measures of text difficulty at best. Also, these metrics were developed for longer-form explanations, like those found in textbooks, rather than dialogue-style teaching. Due to their reliability on counts, their estimates are unreliable in shorter formats. To make improvements on the adaptability of LLMs to education levels measurable, we need better metrics. Only when we can measure improvements can we make tangible progress in automation.<sup>2</sup>

Alternatively, we can use (Neural) classifiers to predict the educational level of a given text. They generalize better and can be applied to texts of varying length. However, these classifiers are expensive to train, and require more training data than we usually have for a niche domain like educational purposes.

The obvious goal is human assessment, but it is expensive to collect and, like all annotation tasks, suffers from disagreement.

In this paper, we introduce and evaluate a new set of Neural metrics for text difficulty as complements to existing Static metrics. Neural metrics exploit the general language understanding capabilities of LLMs to capture more abstract features of educational texts than Static metrics. LLMs can, for example, easily and flexibly classify whether a given text uses a metaphor or not (which is one of the adaptation techniques used by teachers to adjust content to higher education levels). This would be difficult to do with Static approaches.

First, to motivate our selection of Neural metrics, we conduct a user study, where we ask a group of university students to 1) assess the difficulty of educational texts and explain their reasoning, and 2) come up with prompts for an LLM to change the difficulty of a given text. We then translate the qualitative findings from both parts of the study into concrete Neural metrics – like the metaphor example above.

Second, to evaluate the usefulness of our new Neural metrics for measuring text appropriateness for education levels, we run a series of experiments

---

<sup>2</sup>As exemplified by the cases of BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), and BLANC (Recasens and Hovy, 2011), among others, which kickstarted and sustained the development of automated approaches to machine translation, summarization, and coreference resolution, respectively.

using a subset of the ScienceQA dataset (Auer et al., 2023), which contains question-answer pairs across several topics and education levels. Specifically, we run multinomial logistic regressions based on Static metrics, Neural metrics, and the combination of the two to evaluate the marginal benefits of our new Neural metrics. We also compare these regression approaches with using an LLM directly for zero-shot classification.

We find that Neural metrics by themselves are much less useful for text difficulty classification than Static metrics – even if both perform on par or better than a zero-shot LLM classifier. However, the combination of both Neural and Static metrics clearly performs best overall. This shows that Neural metrics do indeed capture relevant signals beyond what Static metrics can capture despite the large number of Static metrics that we include in our experiments. Therefore, the use of Neural metrics in addition to Static metrics appears to be a promising direction for improving text difficulty classification.

To summarise, our main contributions are:

1. We conduct a user study to motivate the creation of novel Neural metrics of text difficulty for educational texts.
2. We show in a series of regression experiments that these novel Neural metrics hold additional value for text difficulty classification beyond what Static metrics can capture.
3. By leveraging the interpretability of our regressions, we highlight the relative importance of individual Static and Neural metrics.

## 2 User Study

In this paper, we introduce novel Neural metrics for text difficulty. To provide an empirical foundation for these metrics, we ran a one-day user study with a group of university students in November 2023.

### 2.1 Study Design

Our user study consists of two main parts.

In the first part of our study, we asked participants to review 60 educational texts randomly sampled from the ScienceQA dataset (Lu et al., 2022). Each text consists of a question (e.g. “What is the mass of a dinner fork?”) with answer choices (“70 grams or 70 kilograms”) and a longer-form explanation of the solution. All texts we select here are

153 either from the social, natural, or language sciences.  
154 Participants were tasked with a) labeling the edu-  
155 cation level of each text as appropriate for either  
156 elementary school, middle school, or high school  
157 and b) explaining the reasoning behind their choice  
158 in a short free-text answer.

159 In the second part of our study, we asked partic-  
160 ipants to rewrite scientific text explanations, also  
161 sampled from ScienceQA, to be appropriate for  
162 different education levels, with the help of an LLM  
163 – in this case, ChatGPT. For example, participants  
164 were asked to rewrite a middle school explanation  
165 of thermal energy at the elementary and high school  
166 levels with the help of prompts. We recorded  
167 the prompts they used to get ChatGPT to accom-  
168 plish the adaptation for them. Thus, we collected  
169 prompts that are used both for text *simplification*  
170 and for text *complexification*.

## 171 2.2 Study Participants

172 We ran our study as part of a hackathon at  
173 [REDACTED UNIVERSITY]. There were seven  
174 participants aged between 21 and 31 years. Four  
175 participants were female, three male. All partic-  
176 ipants were students at [REDACTED UNIVER-  
177 SITY] and were enrolled at the time in programs  
178 specializing in computational linguistics, computer  
179 science, and AI. Five were studying for a bachelor’s  
180 degree and two for a master’s degree. The partic-  
181 ipants held prior educational degrees from school  
182 systems across five different countries. Their native  
183 languages were diverse, including English, Italian,  
184 German, Greek, and Ukrainian. They self-reported  
185 their English language proficiency at C1 and C2  
186 levels. Participants were compensated in study  
187 credits that could be counted towards completing  
188 their program.

## 189 2.3 Study Results

190 The first task of our study yielded 276 classification  
191 labels together with their corresponding descriptive  
192 justifications. These include 120 label-explanation  
193 pairs for middle school texts, 89 for high school,  
194 and 67 for elementary school texts. In the second  
195 task of our study, we collected 103 prompts for text  
196 simplification and complexification. We share il-  
197 lustrative examples of classifications, explanations,  
198 as well as prompts in Appendix A.

199 In the next section, we use the qualitative results  
200 from our study to motivate the construction of novel  
201 Neural metrics for text appropriateness for various  
202 education levels.

## 203 3 Metrics for Text Difficulty

### 204 3.1 Neural Metrics

205 Since the metrics we introduce are based on the  
206 prompts of Neural language models rather than  
207 discrete heuristics, we refer to them as ‘Neural’ to  
208 distinguish them. The goal of the Neural metrics  
209 we develop is to capture more abstract features of  
210 educational texts than would be possible with Static  
211 metrics, which typically focus on individual words  
212 and their statistics.

213 We derive our Neural metrics from the results  
214 of our user study. Figure 2 shows an illustrative  
215 example of our derivation process. We consider  
216 the explanations provided by users for why they  
217 consider a certain educational text to be of elemen-  
218 tary, middle, or high school level difficulty. Then,  
219 we identify recurring attributes and other explana-  
220 tion features that are mentioned by several users  
221 to reflect them in Neural metrics. More specifi-  
222 cally, we examine the distributions of unigrams,  
223 bigrams, and trigrams across all three labels, ex-  
224 cluding function words (see Figure 3). Some of the  
225 most frequent unigrams for the elementary level in-  
226 clude *simple*, *basic*, *elementary*; for the high school  
227 level, *high*, *complex*, *concepts*; and for the middle  
228 school, *explicit*, *explanation*, *middle*.

229 We qualitatively assessed the n-gram distribu-  
230 tions, considering both frequencies and topic ap-  
231 propriateness, before finalizing the query construc-  
232 tion. Each Neural metric takes the form of a sim-  
233 ple yes-no question which we use to prompt the  
234 LLMs. These metrics encompass the most frequent  
235 unigrams as well as less common bigrams and tri-  
236 grams derived from the findings of our study.

237 In total, we construct 41 Neural metrics using  
238 this process. Each Neural metric relates to either  
239 education level (24 metrics), lexical or syntactic  
240 complexity (7 metrics), or the topic of the text at  
241 hand (10 metrics). Table 1 shows an example of  
242 each type of metric. For the full list of all Neural  
243 metrics, see Appendix C.

### 244 3.2 Existing Static Metrics

245 Static metrics are the baseline we want to improve  
246 on. All Static metrics are based on simple formu-  
247 las, heuristics or counts of words or other textual  
248 features. This makes them simple to apply but also  
249 limits the conceptual complexity of what they can  
250 reasonably be expected to measure. In total, we  
251 include 38 Static metrics, selected from those com-  
252 piled in prior work by Flekova et al. (2016). These

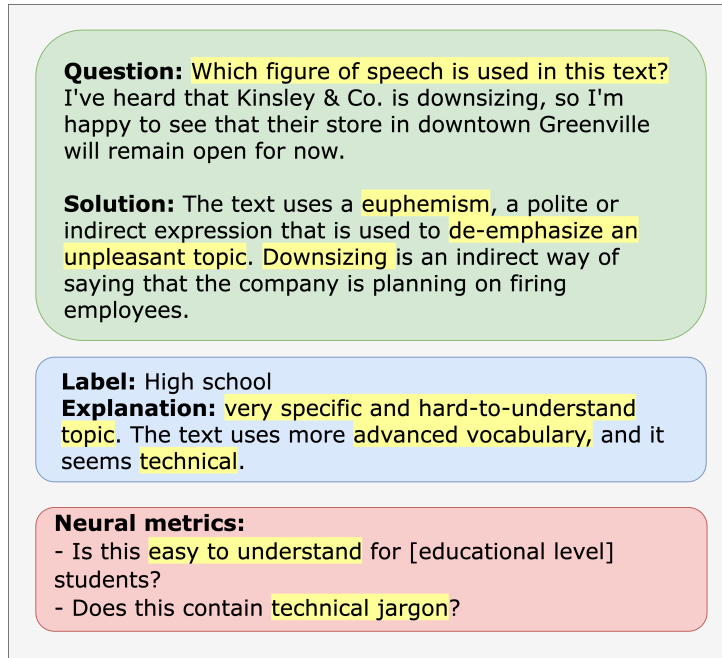


Figure 2: An illustrative example of the Neural metric process. The green box contains the education text from the ScienceQA dataset. The blue box shows the predicted educational level and the explanation. The red box contains the Neural metrics based on the sample.

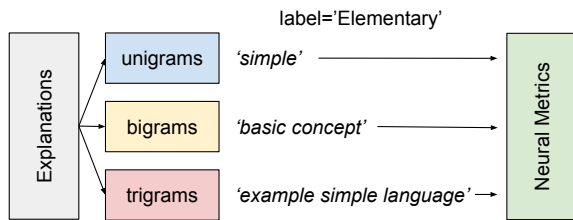


Figure 3: High-level view of the derivation process for the Neural metrics using n-gram frequencies. Function words are excluded.

| Neural Metric  | Category                     |
|--|------------------------------|
| Is this text suitable for an <b>elementary school</b> student? | Education level              |
| Does this text contain <b>technical jargon</b> ?               | Lexical/syntactic complexity |
| Is this text about <b>math</b> ?                               | Topic                        |

Table 1: Three example Neural metrics across the three categories of Neural metrics. Each Neural metric is a prompt for an LLM.

range from text-level metrics that capture general linguistic properties such as vocabulary size and word frequency to sentence-level metrics such as sentence length and syntactic complexity. For the full list of 38 Static metrics see Appendix ??.

## 4 Experiments

To evaluate the usefulness of our novel Neural metrics for measuring text difficulty we conduct a series of text difficulty classification experiments.

### 4.1 Dataset

All our experiments are based on the recently-released ScienceQA dataset (Auer et al., 2023). There are ca. 21k texts in ScienceQA. Each text consists of a question (e.g. “What is the mass of a dinner fork?”) with answer choices (“70 grams

or 70 kilograms”) and a longer-form explanation of the solution. Texts in ScienceQA are classified according to their grade level using the K12 system from the US education system. We simplify this classification by collapsing the 12-grade levels into just three: elementary school (grades 1 to 5), middle school (grades 6 to 8), and high school (grades 9 to 12).<sup>3</sup> From the total 21k texts in ScienceQA, we sample only those that do not use images in questions or explanations. We then deduplicate and sample 1,500 texts for each education level to create a balanced dataset of 4,500 texts. Of these 4,548 texts, we use 3638 (80%) for training and 910 (20%) for evaluation. To our knowledge, ours

<sup>3</sup><https://usahello.org/education/children/grade-levels/>



is the first use of the ScienceQA dataset for the purposes of training and evaluating text difficulty classifiers.

## 4.2 LLMs for Neural Metrics

We mainly use LLMs in our experiments to compute the 41 Neural metrics described in Section ???. In principle, any LLM can be used for this purpose. We use three state-of-the-art LLMs: Llama2 (Touvron et al., 2023), Mistral (Jiang et al., 2023) and GPT-4 (OpenAI, 2023). Llama2, launched in July 2023, is a collection of both pre-trained and fine-tuned LLMs, with sizes ranging from 7 billion to 70 billion parameters. It has been reported to outperform other open-access LLMs and demonstrates capabilities comparable to ChatGPT across various tasks. In the main body of this paper, we focus on the chat-optimised 13 billion parameter version of Llama2, which we refer to as Llama2-13b. The Mistral-7B, released in September 2023, is another open-source language model that surpasses similar-sized open LLMs, including Llama2-13b. We test the instruction-tuned version, Mistral-7b-Instruct-v0.2, which was published in December 2023. Access to these models is provided with Hugging Face. Additionally, we examine GPT-4, a proprietary model from OpenAI released in March 2023. GPT-4 is accessed through its API. This study aims to provide a comprehensive comparison of these state-of-the-art LLMs.

In Appendix D, we report corresponding results for Mistral-7B and GPT-4. We set the model temperature to zero to make responses deterministic. The maximum response length is 256 tokens. Otherwise, we use standard generation parameters as provided by the Hugging Face transformers library. We collected all responses in November and December 2023.

## 4.3 Multinomial Logistic Regression

To classify the difficulty level of texts, we use simple multinomial logistic regression. Formally, the task is to predict the difficulty level  $C_i$  of a given educational text  $S_i$ .  $C_i$  can take three ordinal values: elementary, middle, or high school difficulty. Instead of including  $S_i$  directly, we include sets of Static and/or Neural metrics  $\mathbf{M}_i$  that are computed based on  $S_i$ . We regress  $\mathbf{M}_i$  on  $C_i$  on the 3638 training texts and then evaluate on the 910 test education texts.

## 4.4 Baseline: Zero-Shot Classification

We exploit the general language capabilities of LLMs to compute Neural metrics, which we then use as inputs to a logistic classifier for text difficulty. A natural follow-up question is whether LLMs could just as well predict education level text difficulty directly. Therefore, we include a zero-shot text classification baseline. We use the same model for zero-shot classification as for computing our Neural metrics. For the experiments here, this is Llama2-13b. Note that while the logistic classifier is fitted to our training data, the zero-shot LLM has not seen any examples at inference time.

## 4.5 Experimental Setups

We vary which metrics we include across experimental setups to evaluate the marginal benefits of different metrics. There are three main setups of interest: 1) Neural metrics only, which we refer to as NEURAL, 2) Static metrics only, which we refer to as STATIC, and 3) the combination of the two, which we refer to as COMBO.

## 4.6 Results

Table 2 reports the overall results of our different logistic classifier setups along with the zero-shot LLM classification baseline.

| Method        | Macro-F1    | Accuracy (%) |
|---------------|-------------|--------------|
| NEURAL Reg.   | 0.42        | 45.9         |
| STATIC Reg.   | 0.76        | 78.4         |
| COMBO Reg.    | <b>0.80</b> | <b>80.7</b>  |
| Zero-shot LLM | 0.42        | 45.5         |

Table 2: Overall performance for difficulty level classification on our ScienceQA testset. NEURAL metrics and zero-shot LLM classification use the *Llama2-13b* model. Best performance is highlighted in **bold**.

NEURAL regression performs relatively poorly, at 0.42 macro-F1 and 45.9% accuracy. This is on-par with directly prompting the LLM (in this case Llama2-13b) to make a zero-shot classification. By comparison, STATIC regression performs substantially better, at 0.76 macro-F1 and 78.4% accuracy. Best overall is COMBO, the combination of the two sets of metrics, at 0.80 macro-F1 and 80.7% accuracy.

To investigate performance in more detail, we split out the results for each regression setup by label, i.e. education level, in Table 3.

|        | Level  | Precision      | Recall         | F1-Score       |
|--------|--------|----------------|----------------|----------------|
| NEURAL | Elem.  | 0.55           | 0.36           | 0.43           |
|        | Middle | 0.50           | 0.21           | 0.30           |
|        | High   | 0.43           | 0.78           | 0.55           |
| STATIC | Elem.  | 0.80           | 0.84           | 0.82           |
|        | Middle | 0.71           | 0.59           | 0.64           |
|        | High   | 0.79           | <b>0.87***</b> | 0.83           |
| COMBO  | Elem.  | <b>0.86***</b> | <b>0.87***</b> | <b>0.86***</b> |
|        | Middle | <b>0.73***</b> | <b>0.67***</b> | <b>0.70***</b> |
|        | High   | <b>0.81***</b> | <b>0.87***</b> | <b>0.86***</b> |

Table 3: Performance for difficulty level classification on our ScienceQA testset, split by level. \*\*\* indicates statistically significant improvements. NEURAL metrics use the Llama2-13b model. Best performance per level is highlighted in **bold**.

The overall picture remains unchanged: NEURAL regression still performs worst, while STATIC performs much better, and COMBO performs best, which indicates some marginal benefit to including the NEURAL metrics. The most challenging task appears to be identifying middle school-level texts, with no model scoring more than 0.70 macro-F1. This may, in part, be explained by the ordinal nature of our labels, making it easier to classify content at the extremes of the difficulty scale than in the middle. The NEURAL model struggles more with elementary school than high school texts, whereas the STATIC and COMBO models perform equally well on both.

we collect multiple (1000) bootstrap samples to train and test the logistic regression models for each approach. This method helps in understanding the variability and reliability of the model performance. We use t-tests to determine if the observed differences in accuracies are statistically significant. The results indicate a statistically significant improvement when using COMBO over both STATIC and NEURAL approaches individually.

One big benefit of our regression approach over, for example, classification with an LLM, is that we can easily measure the feature importance of each metric that goes into the classification result. For this purpose, we calculate univariate F-tests between each metric and the difficulty level variable. Table 4 shows the top-five most important features each among the Neural and the Static metrics, based on these F-tests.

Most noticeably, the Neural metrics are gener-

ally much less important than the Static metrics. The top Neural metric, which gives a binary assessment of whether a given text is appropriate for the skill level of elementary school students, is roughly 10% as relevant to difficulty level as the fifth-most-important Static metric counting the number of unique tokens (62.38 vs 629.00 F). However, while they may not hold the same importance, all of the top metrics are highly statistically significant.

## 5 Discussion

### 5.1 The Value of Neural Metrics

Neural metrics by themselves may not be a good-enough basis for classifying text difficulty (Table 2). Static metrics are much more effective by comparison. However, our results also show that Neural metrics do indeed capture relevant features of the text that are not captured by Static metrics, since models that combine both kinds of metrics clearly perform best overall. This is despite the fact that the Static metrics we include are many and highly diverse.

Beyond the demonstrated practical utility of the specific Neural metrics we introduced in this paper, the use of Neural metrics more generally appears to be a promising direction for assessing text difficulty. Better Neural metrics identified in future work may be even more effective complements to Static metrics.

### 5.2 Limitations

The user study we conducted provides a clear empirical motivation for the Neural metrics we selected. This in itself is a core contribution of our work. However, due to resource and time constraints, the sample of participants in the study is fairly small and of limited diversity. Future work could improve on our approach by conducting larger studies or recruiting participants from even more relevant professions (e.g. teachers) to motivate the selection of even better Neural metrics.

Our experiments are mostly constrained by the availability of relevant data for text difficulty classification. The ScienceQA dataset that we use is, to our knowledge, the only dataset that fits our experimental setup in terms of size and detail on education level. Therefore, we cannot make any strong claims about the generalisability of our results. Future work could invest into building new datasets and testing cross-domain performance of

|                | Rank | Metric  | F         |
|----------------|------|---|-----------|
| Neural Metrics | 1    | Is this <b>appropriate for the skill level</b> of elementary school students? | 62.38***  |
|                | 2    | Does this contain a <b>complex language</b> structure?                        | 58.04***  |
|                | 3    | Is this <b>easy to understand</b> for elementary school students?             | 46.72***  |
|                | 4    | Is this <b>suitable</b> for an elementary school student?                     | 34.43***  |
|                | 5    | Is this about <b>earth science</b> ?  | 27.71***  |
| Static Metrics | 1    | Herdan’s C (measures <b>lexical diversity</b> )                               | 916.40*** |
|                | 2    | Entropy (measures <b>the lexical diversity</b> )                              | 875.73*** |
|                | 3    | Flesch-Kincaid Reading Ease (measures <b>readability</b> )                    | 715.95*** |
|                | 4    | Simpson’s Diversity Index (measures <b>lexical diversity</b> )                | 686.87*** |
|                | 5    | # unique tokens (measures <b>length</b> and <b>lexical diversity</b> )        | 629.00*** |

Table 4: Top five most important features among the Neural and Static metrics. Feature importance is measured using univariate F-tests. Larger F indicates higher feature importance. \*\*\* indicates significance at >99.999% confidence.

both Static and Neural metrics, which would give useful insights into which text features are most generally indicate of text difficulty.

## 6 Related Work

### 6.1 Automatic Evaluation of Educational Content

The difficulty level prediction of questions presented to students is crucial for facilitating more effective and efficient learning. Pérez et al. (2012) shows teachers usually fail to identify the correct difficulty level of the questions according to their students’ answers and final scores. The student’s perception of the difficulty also changes across grades and subjects. AlKhuzayy et al. (2023) discovers that linguistic features significantly influence the determination of question difficulty levels in educational assessments. They have explored various syntactic and semantic aspects to understand the complexity of these questions. Crossley et al. (2019) shows the value of using crowdsourcing methods to gather human assessments of text comprehension, coupled with linguistic attributes derived from advanced readability metrics. This approach aids in creating models that explain how humans understand and process text, as well as factors influencing reading speed. Imperial and Madabushi (2023) and Rooein et al. (2023) use language models for content generation over text simplification tasks and controlling readability scores for specific age and educational levels. They show the limitations of the LLMs in adaptation to the specific educational grades.

### 6.2 Question Answering Datasets in Education

The review study by AlKhuzayy et al. (2023) about the literature on item difficulty prediction reveals a significant shortage of publicly accessible datasets with items that are labeled according to their difficulty levels. For example, Hsu et al. (2018) gathered their dataset from national standardized entrance tests that often concentrate on the medical and language fields, annotated with the performance data of 270,000 examinees. This study includes the necessity for a publicly accessible collection of standardized datasets and the need for further exploration into alternative methods for feature elicitation and prediction modeling. The lack of publicly available datasets for measuring difficulty has led researchers toward the domain of Automatic Question Generation (AQG) in recent years. Typically, AQG tends to be more straightforward in structure and cognitive demand compared to questions written by humans. Most of these automatically generated questions are basic, primarily addressing only the first level of Bloom’s taxonomy, which is focused on recall (Leo et al., 2019). Another source of educational datasets is retrieved from online learning platforms or websites specific to the study’s domain. An example includes the collection of 1657 programming problems from LeetCode<sup>4</sup>, labeled with the number of solutions submitted and the pass rate for each problem. Additionally, fewer datasets are from domain-specific textbooks and preparation books, particularly prevalent in the language domain for their role

<sup>4</sup><https://leetcode.com>

in training students for language proficiency exams. The remaining sources were developed by domain experts to meet specific study goals, and according to AIKhuzaey et al. (2023), only 7% from school or university-level assessments.

The Stanford Question Answering Dataset (SQuAD), developed by Rajpurkar et al. (2016) in 2016, features 150,000 questions in the form of paragraph-answer pairs sourced from Wikipedia articles. This dataset was utilized by Bi et al. (2021) to develop and test their models for predicting the difficulty of reading comprehension questions. Lu et al. (2022) created a new multimodal science question-answering datasets, which includes 21,000 English passages from school reading exams, each accompanied by four multiple-choice questions. The ScienceQA dataset provides several metadata fields associated with each question, including extensive solutions and general explanations. In contrast to SQuAD, this dataset demands more advanced reasoning abilities to answer its questions.

## 7 Conclusion

Good teachers succeed in making the material understandable for their respective audiences. This adaptation is a complex process which goes well beyond replacing individual words and phrases. However, existing Static metrics for text difficulty, like the Flesh-Kincaid Reading Ease score, still focus on precisely those elements. As a result, these metrics are crude and brittle, failing to adapt to new domains and working mainly on long-form documents.

Large Language Models are increasingly used in educational domains and offer ways to go beyond individual word replacement due to their general language capabilities. However, at the same time, they still struggle to adapt to precise education levels. To effectively automate text adaptation to education levels, we need to measure the success of that adaptation, which requires more flexible metrics than the ones we currently have.

We introduce a suite of prompt-based Neural metrics for text adaptation based on a user study. We empirically show that these metrics, in combination with traditional Static metrics, improve text difficulty prediction. Our work opens up new avenues for the use of LLMs in educational applications.

## Ethical Considerations

The participants in the user study we used in our paper were student volunteers of a course on related topics. They were able to leave the study at any point, and were compensated in course credits that could be counted towards their study program. The study was conducted in accordance with the rules of the host university and passed their ethics assessment. The risk for harm to the participants in this setting was assessed as minimal.

## References

- Samah AIKhuzaey, Floriana Grasso, Terry R Payne, and Valentina Tamma. 2023. Text-based question difficulty prediction: A systematic review of automatic approaches. *International Journal of Artificial Intelligence in Education*, pages 1–53.
- Sören Auer, Dante AC Barone, Cassiano Bartz, Eduardo G Cortes, Mohamad Yaser Jaradeh, Oliver Karras, Manolis Koubarakis, Dmitry Mouromtsev, Dmitrii Pliukhin, Daniil Radyush, et al. 2023. The sciqa scientific question answering benchmark for scholarly knowledge. *Scientific Reports*, 13(1):7240.
- Sheng Bi, Xiya Cheng, Yuan-Fang Li, Lizhen Qu, Shiron Shen, Guilin Qi, Lu Pan, and Yinlin Jiang. 2021. Simple or complex? complexity-controllable question generation with soft templates and deep mixture of experts model. *arXiv preprint arXiv:2110.06560*.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.
- Scott A Crossley, Stephen Skalicky, and Mihai Dascalu. 2019. Moving beyond classic readability formulas: New methods and new models. *Journal of Research in Reading*, 42(3-4):541–561.
- Lucie Flekova, Daniel Preoțiuc-Pietro, and Lyle Ungar. 2016. Exploring stylistic variation with age and income on twitter. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 313–319.
- Mohammad Hosseini, Catherine A Gao, David M Liebovitz, Alexandre M Carvalho, Faraz S Ahmad, Yuan Luo, Ngan MacDonald, Kristi L Holmes, and Abel Kho. 2023. An exploratory survey about using chatgpt in education, healthcare, and research. *medRxiv*, pages 2023–03.
- Fu-Yuan Hsu, Hahn-Ming Lee, Tao-Hsing Chang, and Yao-Ting Sung. 2018. Automated estimation of item difficulty for multiple-choice tests: An application of word embedding techniques. *Information Processing & Management*, 54(6):969–984.



|     |   |   |     |
|-----|---|---|-----|
| 615 | Joseph Marvin Imperial and Harish Tayyar Madabushi.                         | Donya Rooein, Amanda Cercas Curry, and Dirk Hovy.                     | 668 |
| 616 | 2023. Flesch or fumble? evaluating readability stan-                        | 2023. Know your audience: Do llms adapt to dif-                       | 669 |
| 617 | dard alignment of instruction-tuned language models.                        | ferent age and education levels? <i>arXiv preprint</i>                | 670 |
| 618 | <i>arXiv preprint arXiv:2309.05454</i> .                                    | <i>arXiv:2312.02065</i> .   | 671 |
| 619 | Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-                         | Malik Sallam. 2023. <a href="#">Chatgpt utility in healthcare ed-</a> | 672 |
| 620 | sch, Chris Bamford, Devendra Singh Chaplot, Diego                           | <a href="#">ucation, research, and practice: Systematic review</a>    | 673 |
| 621 | de las Casas, Florian Bressand, Gianna Lengyel, Guil-                       | <a href="#">on the promising perspectives and valid concerns.</a>     | 674 |
| 622 | laume Lample, Lucile Saulnier, et al. 2023. Mistral                         | <i>Healthcare</i> , 11(6).  | 675 |
| 623 | 7b. <i>arXiv preprint arXiv:2310.06825</i> .                                |   |     |
| 624 | Enkelejda Kasneci, Kathrin Seßler, Stefan Küchemann,                        | Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-                    | 676 |
| 625 | Maria Bannert, Daryna Dementieva, Frank Fischer,                            | bert, Amjad Almahairi, Yasmine Babaei, Nikolay                        | 677 |
| 626 | Urs Gasser, Georg Groh, Stephan Günemann, Eyke                              | Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti                    | 678 |
| 627 | Hüllermeier, et al. 2023. Chatgpt for good? on op-                          | Bhosale, et al. 2023. Llama 2: Open founda-                           | 679 |
| 628 | portunities and challenges of large language models                         | tion and fine-tuned chat models. <i>arXiv preprint</i>                | 680 |
| 629 | for education. <i>Learning and Individual Differences</i> ,                 | <i>arXiv:2307.09288</i> .   | 681 |
| 630 | 103:102274.   |   |     |
| 631 | Jared Leo, Ghader Kurdi, Nicolas Matentzoglou, Bijan                        | Shriyash Upadhyay, Etan Ginsberg, and Chris Callison-                 | 682 |
| 632 | Parsia, Ulrike Sattler, Sophie Forge, Gina Donato,                          | Burch. 2023. Improving mathematics tutoring with                      | 683 |
| 633 | and Will Dowling. 2019. Ontology-based generation                           | a code scratchpad. In <i>Proceedings of the 18th Work-</i>            | 684 |
| 634 | of medical, multi-term mcqs. <i>International Journal</i>                   | <i>shop on Innovative Use of NLP for Building Educa-</i>              | 685 |
| 635 | <i>of Artificial Intelligence in Education</i> , 29:145–188.                | <i>tional Applications (BEA 2023)</i> , pages 20–28.                  | 686 |
| 636 | Chin-Yew Lin. 2004. <a href="#">ROUGE: A package for auto-</a>              | Lixiang Yan, Lele Sha, Linxuan Zhao, Yuheng Li,                       | 687 |
| 637 | <a href="#">matic evaluation of summaries</a> . In <i>Text Summariza-</i>   | Roberto Martinez-Maldonado, Guanliang Chen,                           | 688 |
| 638 | <i>tion Branches Out</i> , pages 74–81, Barcelona, Spain.                   | Xinyu Li, Yueqiao Jin, and Dragan Gašević. 2023.                      | 689 |
| 639 | Association for Computational Linguistics.                                  | Practical and ethical challenges of large language                    | 690 |
| 640 | Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-                        | models in education: A systematic literature review.                  | 691 |
| 641 | Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter                             | <i>arXiv preprint arXiv:2303.13379</i> .                              | 692 |
| 642 | Clark, and Ashwin Kalyan. 2022. Learn to explain:                           |   |     |
| 643 | Multimodal reasoning via thought chains for science                         | <b>A Selected Prompts from the User Study</b>                         | 693 |
| 644 | question answering. <i>Advances in Neural Information</i>                   | We collect the top prompts of the students from                       | 694 |
| 645 | <i>Processing Systems</i> , 35:2507–2521.                                   | the chat history with analytical, manual, and AI                      | 695 |
| 646 | OpenAI. 2023. <a href="#">GPT-4 Technical Report</a> .                      | Assistant (ChatGPT).  | 696 |
| 647 | Kishore Papineni, Salim Roukos, Todd Ward, and Wei-                         | <b>A.1 Elementary School:</b>   | 697 |
| 648 | Jing Zhu. 2002. <a href="#">Bleu: a method for automatic evalu-</a>         | - Simplify a text for elementary school, using sim-                   | 698 |
| 649 | <a href="#">ation of machine translation</a> . In <i>Proceedings of the</i> | ple language for 6-12 year olds. - Create an el-                      | 699 |
| 650 | <i>40th Annual Meeting of the Association for Compu-</i>                    | ementary version of a high school lecture text. -                     | 700 |
| 651 | <i>tational Linguistics</i> , pages 311–318, Philadelphia,                  | Simplify a high school text for elementary school.                    | 701 |
| 652 | Pennsylvania, USA. Association for Computational                            | - Explain in a way an 8-year-old would understand.                    | 702 |
| 653 | Linguistics.  | - This is a text meant for high school students. Can                  | 703 |
| 654 | Elena Verdú Pérez, Luisa M Regueras Santos, María                           | you help me make an appropriate version for ele-                      | 704 |
| 655 | Jesús Verdú Pérez, Juan Pablo de Castro Fernández,                          | mentary school students with very simple language                     | 705 |
| 656 | and Ricardo García Martín. 2012. Automatic clas-                            | and comprehensive, easy-to-understand examples?                       | 706 |
| 657 | sification of question difficulty level: Teachers’ esti-                    | <b>A.2 Middle School:</b>   | 707 |
| 658 | mation vs. students’ perception. In <i>2012 Frontiers</i>                   | - Give examples from middle school lectures. -                        | 708 |
| 659 | <i>in Education Conference Proceedings</i> , pages 1–5.                     | Adapt a high school text for middle school, using                     | 709 |
| 660 | IEEE.   | less advanced language. - Be more textbook-like                       | 710 |
| 661 | Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and                       | and more to the point for level of middle school.                     | 711 |
| 662 | Percy Liang. 2016. Squad: 100,000+ questions                                | - Adapt content for a student in middle school. -                     | 712 |
| 663 | for machine comprehension of text. <i>arXiv preprint</i>                    | Simplify a lecture text for middle school, using                      | 713 |
| 664 | <i>arXiv:1606.05250</i> .   | illustrative examples.  | 714 |
| 665 | Marta Recasens and Eduard Hovy. 2011. Blanc: Im-                            | <b>A.3 High School:</b>   | 715 |
| 666 | plementing the rand index for coreference evaluation.                       | - Enhance scientific accuracy and add comprehen-                      | 716 |
| 667 | <i>Natural language engineering</i> , 17(4):485–510.                        | sive examples for high school level. - Adapt a                        | 717 |

|     |  |     |
|-----|--|-----|
| 718 | middle school text for high school, using advanced                 | 767 |
| 719 | language. - Increase difficulty for high school, with              | 768 |
| 720 | advanced vocabulary and scientific concepts. - Can                 | 769 |
| 721 | you make it more scientific and less story-telling-                | 770 |
| 722 | like? - Increase difficulty level with comprehensive               | 771 |
| 723 | examples.  | 772 |
| 724 | <b>B Parameter settings</b>  | 773 |
| 725 | The Static matrices are collected by Python pack-                  | 774 |
| 726 | ages such as pandas, textstat, spacy, wordfreq,                    | 775 |
| 727 | and wordfreq. We use <code>_en_core_web_sm</code> param-           | 776 |
| 728 | eter from the spacy model. For Regression model,                   | 777 |
| 729 | we use sklearn package and SelectKBest with                        | 778 |
| 730 | <code>f_classif</code> score function.                             | 779 |
| 731 | <b>C List of Metrics</b>   | 780 |
| 732 | <b>Static Metrics.</b> <code>num_words</code> : Number of words in | 781 |
| 733 | the text. <code>num_sentences</code> : Number of sentences         | 782 |
| 734 | in the text. <code>type_token_ratio</code> : Ratio of unique       | 783 |
| 735 | words to total words, indicating vocabulary rich-                  | 784 |
| 736 | ness. <code>avg_word_length</code> : The average length of         | 785 |
| 737 | words in the text. <code>academic_word_list_ratio</code> :         | 786 |
| 738 | Ratio of academic or domain-specific words                         | 787 |
| 739 | to total words. <code>complex_words</code> : Number of             | 788 |
| 740 | complex words, often based on syllable count                       | 789 |
| 741 | or other criteria. <code>mtld</code> : Measure of Textual          | 790 |
| 742 | Lexical Diversity, a metric for vocabulary rich-                   | 791 |
| 743 | ness. <code>num_unique_tokens</code> : Number of unique            | 792 |
| 744 | words or tokens in the text. <code>avg_sent_length</code> :        | 793 |
| 745 | The average length of sentences in the text.                       | 794 |
| 746 | <code>std_sent_length</code> : The standard deviation of sen-      | 795 |
| 747 | tence lengths, indicating sentence length variation.               | 796 |
| 748 | <code>clauses_sentences_ratio</code> : Ratio of clauses to         | 797 |
| 749 | sentences, providing insight into sentence com-                    | 798 |
| 750 | plexity. <code>pos_ratios</code> : Ratios of different parts of    | 799 |
| 751 | speech (e.g., nouns, verbs) in the text. FKGL: Flesch-             | 800 |
| 752 | Kincaid Grade Level, an estimate of the text's read-               | 801 |
| 753 | ability. FKES: Flesch Reading Ease score is a read-                | 802 |
| 754 | ability measure for US education systems. <code>ttr</code> :       | 803 |
| 755 | TypeToken Ratio. <code>brunet_index</code> : A metric for          | 804 |
| 756 | text diversity and richness. <code>d_measure</code> : A mea-       | 805 |
| 757 | sure of lexical diversity. <code>yules_k</code> : Yule's K, a      | 806 |
| 758 | measure of text's distribution of word frequencies.                | 807 |
| 759 | <code>herdan_c</code> : Herdan's C, a metric for vocabulary        | 808 |
| 760 | richness. <code>simpsons_di</code> : Simpson's Diversity In-       | 809 |
| 761 | index, indicating the diversity of words in the text.              | 810 |
| 762 | <code>entropy</code> : A measure of information entropy, indi-     | 811 |
| 763 | cating unpredictability or randomness of the text.                 |     |
| 764 | <b>Neural Metrics.</b> Is this readable for an elemen-             |     |
| 765 | tary school student based on the Flesch-Kincaid                    |     |
| 766 | grade scale?, Is this suitable for an elementary                   |     |
|     | school student?, Is this easy to understand for el-                | 767 |
|     | ementary school students?, Is this relevant to cur-                | 768 |
|     | riculum topics for elementary school students?, Is                 | 769 |
|     | this relevant to the knowledge and experiences of                  | 770 |
|     | elementary school students?, Is this content suit-                 | 771 |
|     | able for meeting the expected knowledge level of                   | 772 |
|     | elementary school students? Is this able to provide                | 773 |
|     | detailed feedback to help elementary school stu-                   | 774 |
|     | dents learn?, Is this appropriate for the skill level              | 775 |
|     | of elementary school students?, Is the length of                   | 776 |
|     | this suitable for elementary school students?, Is                  | 777 |
|     | this readable for a middle school student based on                 | 778 |
|     | the Flesch-Kincaid grade scale?, Is this suitable                  | 779 |
|     | for a middle school student?, Is this easy to under-               | 780 |
|     | stand for middle school students?, Is this relevant                | 781 |
|     | to curriculum topics for middle school students?, Is               | 782 |
|     | this relevant to the knowledge and experiences of                  | 783 |
|     | middle school students?, Is this content suitable for              | 784 |
|     | meeting the expected knowledge level of middle                     | 785 |
|     | school students? Is this able to provide detailed                  | 786 |
|     | feedback to help middle school students learn?, Is                 | 787 |
|     | this appropriate for the skill level of middle school              | 788 |
|     | students?, Is the length of this suitable for middle               | 789 |
|     | school students?, Is this readable for a high school               | 790 |
|     | student based on the Flesch-Kincaid grade scale?,                  | 791 |
|     | Is this suitable for a high school student?, Is this               | 792 |
|     | easy to understand for high school students?, Is                   | 793 |
|     | this relevant to curriculum topics for high school                 | 794 |
|     | students?, Is this relevant to the knowledge and                   | 795 |
|     | experiences of high school students?, Is this con-                 | 796 |
|     | tent suitable for meeting the expected knowledge                   | 797 |
|     | level of high school students? Is this able to pro-                | 798 |
|     | vide detailed feedback to help high school students                | 799 |
|     | learn?, Is this appropriate for the skill level of high            | 800 |
|     | school students?, Is the length of this suitable for               | 801 |
|     | high school students?, Does this contain metaphors                 | 802 |
|     | and/or figurative language?, Does this contain a                   | 803 |
|     | complex language structure?, Does this contain                     | 804 |
|     | technical jargon?, Is the language of this simple?,                | 805 |
|     | Is this about science?, Is this about language sci-                | 806 |
|     | ence?, Is this about natural science?, Is this about               | 807 |
|     | social science?, Is this about math?, Is this about                | 808 |
|     | physics?, Is this about chemistry?, Is this about                  | 809 |
|     | earth science?, Is this about world history?, Is this              | 810 |
|     | about geography?   | 811 |
|     | <b>D Details over Mistral-7B and GPT-4</b>                         | 812 |
|     | <b>Models</b>  | 813 |
|     | We describe the performance of Mistral-7B over                     | 814 |
|     | the subset of our dataset with 4500 samples (1500                  | 815 |
|     | samples for each educational level) and GPT-4 for                  | 816 |

450 samples (150 for each educational level).

| Method        | Macro-F1    | Accuracy (%) |
|---------------|-------------|--------------|
| NEURAL Reg.   | 0.37        | 41.66        |
| STATIC Reg.   | 0.73        | 74.00        |
| COMBO Reg.    | <b>0.76</b> | <b>77.3</b>  |
| Zero-shot LLM | 0.37        | 36.90        |

Table 5: Overall performance for difficulty level classification on our ScienceQA testset. NEURAL metrics and zero-shot LLM classification use the *Mistral-7B* model. Best performance is highlighted in **bold**.

| Method        | Macro-F1    | Accuracy (%) |
|---------------|-------------|--------------|
| NEURAL Reg.   | 0.55        | 54.44        |
| STATIC Reg.   | 0.74        | 74.44        |
| COMBO Reg.    | <b>0.76</b> | <b>76.22</b> |
| Zero-shot LLM | 0.56        | 54.88        |

Table 6: Overall performance for difficulty level classification on our ScienceQA testset. NEURAL metrics and zero-shot LLM classification use the *GPT-4* model. Best performance is highlighted in **bold**.

|        | Level  | Precision   | Recall      | F1-Score    |
|--------|--------|-------------|-------------|-------------|
| NEURAL | Elem.  | 0.59        | 0.34        | 0.43        |
|        | Middle | 0.55        | 0.11        | 0.18        |
|        | High   | 0.36        | 0.91        | 0.51        |
| STATIC | Elem.  | 0.76        | 0.88        | 0.82        |
|        | Middle | 0.81        | 0.45        | 0.58        |
|        | High   | 0.68        | <b>0.94</b> | 0.79        |
| COMBO  | Elem.  | <b>0.80</b> | <b>0.90</b> | <b>0.85</b> |
|        | Middle | <b>0.82</b> | <b>0.51</b> | <b>0.63</b> |
|        | High   | <b>0.71</b> | <b>0.94</b> | <b>0.81</b> |

Table 7: Performance for difficulty level classification on ScienceQA testset, split by level. NEURAL metrics use the *Mistral-7B* model. Best performance per level is highlighted in **bold**.

|        | Level  | Precision   | Recall      | F1-Score    |
|--------|--------|-------------|-------------|-------------|
| NEURAL | Elem.  | 0.65        | 0.57        | 0.61        |
|        | Middle | 0.41        | 0.38        | 0.39        |
|        | High   | 0.57        | 0.74        | 0.64        |
| STATIC | Elem.  | <b>0.87</b> | <b>0.74</b> | <b>0.80</b> |
|        | Middle | <b>0.68</b> | <b>0.72</b> | <b>0.70</b> |
|        | High   | 0.69        | <b>0.78</b> | <b>0.73</b> |
| COMBO  | Elem.  | 0.81        | <b>0.74</b> | 0.78        |
|        | Middle | 0.65        | 0.69        | 0.67        |
|        | High   | <b>0.71</b> | 0.74        | <b>0.73</b> |

Table 8: Performance for difficulty level classification on our ScienceQA testset, split by level. NEURAL metrics use the *GPT-4* model. Best performance per level is highlighted in **bold**.

|                | Rank | Metric  | F         |
|----------------|------|---|-----------|
| Neural Metrics | 1    | Is this <b>easy to understand</b> for elementary school students?             | 79.75***  |
|                | 2    | Does this contain <b>metaphors and/or figurative language</b> ?               | 42.20***  |
|                | 3    | Is this <b>readable</b> for elementary school students                        | 40.80***  |
|                | 4    | Is this <b>appropriate for the skill level</b> of elementary school students? | 35.65***  |
|                | 5    | Is this <b>relevant to curriculum topics</b> for elementary school students?  | 16.57***  |
| Static Metrics | 1    | Herdan’s C (measures <b>lexical diversity</b> )                               | 348.82*** |
|                | 2    | Entropy (measures <b>variability or complexity</b> )                          | 346.36*** |
|                | 3    | Flesch-Kincaid Reading Ease (measures <b>readability</b> )                    | 284.60*** |
|                | 4    | Simpson’s Diversity Index (measures <b>lexical diversity</b> )                | 255.80*** |
|                | 5    | Flesch-Kincaid Grade Level (measures the school <b>grade level</b> )          | 235.71*** |

Table 9: Top five most important features among the Neural and Static metrics. Feature importance is measured using univariate F-tests. Larger F indicates higher feature importance. NEURAL metrics use the *Mistral-7B* model. \*\*\* indicates significance at >99.999% confidence.

|                | Rank | Metric   | F         |
|----------------|------|--|-----------|
| Neural Metrics | 1    | Is this <b>appropriate for the skill level</b> of elementary school students?  | 62.38***  |
|                | 2    | Does this contain a <b>complex language</b> structure?                         | 58.04***  |
|                | 3    | Is this <b>easy to understand</b> for elementary school students?              | 46.72***  |
|                | 4    | Is this <b>suitable</b> for an elementary school student?                      | 34.43***  |
|                | 5    | Is this about <b>earth science</b> ?   | 27.71***  |
| Static Metrics | 1    | Herdan’s C (measures <b>lexical diversity</b> )                                | 916.40*** |
|                | 2    | Entropy (measures <b>variability or complexity</b> )                           | 875.73*** |
|                | 3    | Flesch-Kincaid Reading Ease (measures <b>readability</b> )                     | 715.95*** |
|                | 4    | Simpson’s Diversity Index (measures <b>lexical diversity</b> )                 | 686.87*** |
|                | 5    | Number of unique tokens (measures <b>length</b> and <b>lexical diversity</b> ) | 629.00*** |

Table 10: Top five most important features among the Neural and Static metrics. Feature importance is measured using univariate F-tests. Larger F indicates higher feature importance. NEURAL metrics use the *GPT-4* model. \*\*\* indicates significance at >99.999% confidence.