

# STANDARDIZE: Aligning Language Models with Expert-Defined Standards for Content Generation

Anonymous ACL submission

## Abstract

Domain experts across engineering, healthcare, and education follow strict standards for producing quality content such as technical manuals, medication instructions, and children’s reading materials. However, current works in controllable text generation have yet to explore using these standards as references for control. Towards this end, we introduce STANDARDIZE, a retrieval-style in-context learning-based framework to guide large language models to align with expert-defined standards. Focusing on English language standards in the education domain as a use case, we consider the Common European Framework of Reference for Languages (CEFR) and Common Core Standards (CCS) for the task of open-ended content generation. Our findings show that models can gain 45% to 100% increase in precise accuracy across open and commercial LLMs evaluated, demonstrating that the use of knowledge artifacts extracted from standards and integrating them in the generation process can effectively guide models to produce better standard-aligned content<sup>1</sup>.

## 1 Introduction

One of the most realized benefits of large language model (LLM) research is how it became widely adopted by the public. In particular, the rise of chat-style model interfaces, such as ChatGPT and Perplexity, has allowed non-technical users to fully utilize these tools in accomplishing day-to-day tasks and activities, such as getting help with writing, documenting code, and providing recommendations. A key technological advancement behind this is the use of reward-based methods such as Reinforcement Learning for Human Feedback (RLHF, Ouyang et al. (2022)), which allows embedding human preferences to generative models for better-aligned outputs with respect to the task at hand.

<sup>1</sup>Our code, data, and model outputs will be released upon publication.

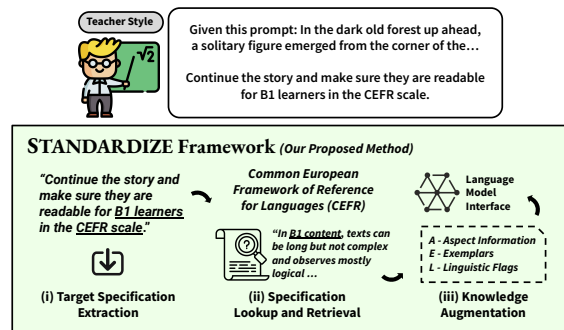


Figure 1: In contrast to the simple prompting method used by teachers, the proposed STANDARDIZE framework aims to improve the performance of generative models for content generation by using the fine-grained information found in expert-defined standards. The framework involves a three-part process starting with the (i) **extraction** of target specifications from the prompt, (ii) **lookup and retrieval** of information that matches the target specifications from the specified standard, and (iii) **knowledge augmentation** to produce artifacts that represent the standard itself for integration into the generation process with generative models.

Despite the growing literature proposing complex algorithms and architectures for enriching the instruction-following capabilities of generative models, the missing puzzle piece that seems to have not garnered equal attention from the community is the integration of actual standards or guidelines crafted by domain experts as a reference for control. For example, in healthcare and engineering, well-documented standards are strictly followed in order to ensure the quality of processes. This includes the UK National Health Service (NHS) Injectable Medicines Guide (IMG) which contains instructions on how medical injectables should be mixed (Keeling et al., 2010) as well as the Simplified Technical English (STE)<sup>2</sup> which is a documented controlled language specification for writing technical manuals that are simple to read. Fol-

<sup>2</sup><https://www.asd-stel100.org/>

lowing these standards, even for domain experts, can be tedious, challenging, and even consequential in serious cases due to its complexity (Jones et al., 2021; Cousins et al., 2005). Thus, this research gap is an opportunity where the complex instruction-following capabilities of language models can provide assistance, particularly for tasks requiring the generation of text content since this is one of the areas where these models objectively perform well (Chung et al., 2022; Wei et al., 2021; Gatt and Kraemer, 2018).

Towards this end, we tackle the main research question: **How can we align large language models for content generation tasks using expert-defined standards?** We list our major contributions from this study as follows:

1. We introduce STANDARD-CTG, a new task formalizing the challenge of generating text using generative language models with expert-defined standards as an additional resource for control.
2. We propose STANDARDIZE, a retrieval-based framework using in-context learning that extracts knowledge artifacts from standards such as aspect information, exemplars, and manually crafted linguistic variables to improve the performances of generative language models for content generation.
3. We introduce high-performing baseline GPT-4 models for the task of STANDARD-CTG using two of the most widely recognized academic standards, CEFR and CCS.

## 2 Expert-Defined Standards

### 2.1 Background

According to the International Organization for Standardization (ISO)<sup>3</sup>, **standards** are documented guidelines often containing rich detail in describing requirements, specifications, and criteria. These guidelines are defined and continuously improved by experts or interest groups in various domains, such as education, healthcare, and accounting, to name a few. Using standards ensures an institution’s products and processes are consistent and reproducible (Sadler, 2017).

In the context of education and language assessment, standards are usually in the form of either (a)

<sup>3</sup><https://www.iso.org/standards.html>

content standards such as documentations of a common language for ease of communication, writing, and content production, and (b) performance standards such as state-administered tests for reading and mathematical problem-solving competencies. This study focuses on content-based standards used in education and language assessment to be integrated into a generative model’s text generation process. The alignment with existing standards for any generated text material is crucial to ensure quality and consistency before being used in classroom settings (La Marca et al., 2000).

### 2.2 Standards in Education and Language Assessment

We discuss the two selected English standards we consider as test cases for this study.

**The Common European Framework of Reference for Languages (CEFR)** is one of the well-known standard language framework<sup>4</sup> developed by The Council of Europe and used for assessing general language competencies such as reading, writing, and listening. The CEFR uses a six-point level scale of A1, A2, B1, B2, C1, and C2, which denotes increasing complexities in instructional content development. We use the level descriptors compiled by Natova (2021), which cover three aspects, namely (1) Meaning/Purpose, (2) Structure, and (3) Grammatical Complexity, describing the characteristics of desired content per level as shown in Table 9. We omit a fourth aspect of Reader’s Knowledge Demands from the standard as this heavily depends on the reader’s background knowledge and is entirely subjective (Forey, 2020; Forey and Cheung, 2019).

**The Common Core Standards (CCS)** is an academic standard<sup>5</sup> developed by the US National Governors Association and the Council of Chief State School Officers (CCSSO) which has been widely adopted by schools across the United States for its K-12 curriculum. In this study, we adapt the recommended model of CCS for assessing text complexity, which includes two main variables: (1) Qualitative Dimensions and (2) Quantitative Dimensions. However, similar to the CEFR standard, we do not include the last variable, which is Reader

<sup>4</sup><https://www.coe.int/en/web/common-european-framework-reference-languages/level-descriptions>

<sup>5</sup><https://corestandards.org/>

150 Considerations, as this requires professional judgment or a teacher’s intervention. The description  
151 of each aspect of CCS is detailed in Table 9.  
152

### 153 2.3 Standard-Aligned Content Generation 154 (STANDARD-CTG)

155 Given the importance of adhering to expert-defined standards in the context of language assessment,  
156 we introduce the task of **standard-aligned content generation**. The overarching goal of STANDARD-  
157 CTG is to pave the way for new approaches that aim to integrate the conventional methodologies  
158 of controllable text generation in NLP with actual constraints provided by experts across interdis-  
159 ciplinary fields such as education, engineering, and medicine through documented standards. To align  
160 with terminologies used in education and other non-computing literature, in this work, we use the term  
161 *content generation* instead of *text generation* as usually seen in technical NLP literature.  
162

163 We represent the task of STANDARD-CTG using the following formulation:  
164  
165

$$166 \quad A = C_{\text{Stdndrd}}(M(p, a, k_a), E) \quad (1)$$

167 where  $A$  quantifies the content alignment score of using a general evaluator  $C_{\text{Stdndrd}}$  that tests the  
168 quality of a language model’s  $M$  generated content against a collection of gold-standard examples  $E$   
169 using inputs such as (a) a natural language prompt  $p$ , (b) information of some aspect  $a$ , and (c) trans-  
170 formed representation of an aspect  $k_a$  defined or extracted from the chosen standard. We pattern  
171 our major experiments in the succeeding sections based on this formulation.  
172

## 173 3 The STANDARDIZE Framework

174 Our main hypothesis in this study is motivated by the fact that expert-defined standards are often very  
175 informative, lengthy, and complex. More specifically, we posit that in order for a generative model  
176 to produce content that is *aligned* with the specifications provided by a standard, the actual information  
177 found in the standard itself must be considered in the actual generation process. The challenge then is  
178 redirected towards *how* any information extracted can be represented as something that the generative  
179 model will find useful.  
180

181 Towards addressing STANDARD-CTG, we propose STANDARDIZE, a retrieval-style in-context  
182 learning-based framework that exploits the rich information found in standards and transforms this  
183

184 into knowledge artifacts to improve the quality of content produced by generative models. Figure 1  
185 encapsulates this framework in a visual manner. In the succeeding sections, we discuss the proposed  
186 STANDARDIZE framework more thoroughly.  
187

188 **Target Specification Extraction** is performed first to obtain informative tags in the prompt and  
189 to correctly match this information within the standards. For academic standards in language  
190 assessment, these specifications should provide information about *who* will be content delivered to  
191 (target audience) and using *what* specific standard out of many (CEFR or CCS). Thus, these two  
192 information tags are the basic required input for the process. As an example shown in Figure 1, the  
193 extracted specifications provided in the prompt are *A2 readers*, which points to a particular group of  
194 learners requiring low-leveled reading materials, and *CEFR scale*, which denotes the selected  
195 standard where properties of A2-level texts are described.  
196

197 **Specification Lookup and Retrieval** is then performed next upon extracting the target speci-  
198 fications. A lookup process is done to find a match with the selected standard, usually in the form of  
199 a database or an external machine-readable file. The information from the standard in the form of  
200 *aspects* (or characteristics) that match the target specifications is then retrieved. The length and  
201 complexity of a standard’s level of information regarding its specifications may vary. As shown  
202 in Figure 1 for the CEFR standard, the retrieved information that matches the desired level of  
203 complexity for the target audience (A2 readers) can be checked at Table 9.  
204

205 **Knowledge Augmentation** is done last but is the most important process of the pipeline. We propose  
206 a further technical augmentation of information found in standards to obtain **knowledge artifacts** in  
207 the prompts. These knowledge artifacts can range from simple additional information already present  
208 in the standard to complex representations, such as incorporating actual linguistic features to control  
209 the granularity of the generation process. Recent works surveying the performance of open and  
210 closed models have shown that non-informative style of prompting language models, such as the  
211 teacher style shown in Figure 1, is effective only to  
212

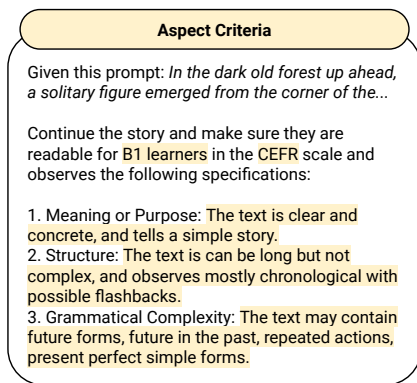


Figure 2: A standard contains recommended characteristics of content across one or more **domain-specific aspects** or criteria. This figure shows an example of the CEFR standard where the set of criteria includes depth of meaning, structure, and grammatical complexity.

a certain extent and may be biased towards content generation in lower levels, such as A2 or B1 in the CEFR standards (Imperial and Madabushi, 2023b; Ribeiro et al., 2023).

#### 4 Knowledge Artifacts for STANDARDIZE

In this section, we discuss the knowledge artifacts used by the STANDARDIZE framework and how they are integrated into the generation setup via prompting.

**Aspect Information (STANDARDIZE-A)** is the most evident form of knowledge artifact as this pertains to the descriptive information provided in the standard. In the context of standards for content generation, aspect information is generally attributed to linguistic criteria of content with respect to its target audience. Figure 2 shows how aspect information from a standard (e.g., CEFR) can be integrated into the actual prompt. The addition of aspect criteria information ensures that the generative model will have access to *explicit characteristics* of the desired generated content in different dimensions.

**Linguistic Signals (STANDARDIZE-L)** represent the controllable variables of a standard that a generative model can use to steer the direction of content generation. In the STANDARDIZE framework, this process serves as a **rewrite function** where a generative model is asked to produce an initial content first using another method prompting (e.g., aspect information in Figure 2), and rewrites this by comparing linguistic flag values of the initially

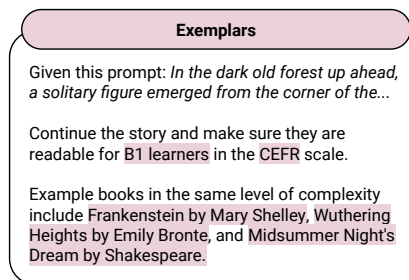


Figure 3: A standard contains recommended **exemplars** that serve as gold-standard reference. This figure shows an example of the CEFR standard where three well-known pieces of literature are provided as examples of content that conforms to the target level specified (B1).

generated content against the mean value of a gold standard dataset of the target level. An example is illustrated in Figure 4 where the mean type-token ratio of a collection of gold-standard B1-level text 12.50 is added to the prompt while being compared to the current type-token value of the story, which is 4.22. A *verbalizer* is used to transform the computed linguistic flags into natural language prompts. The keywords *increase* and *decrease* are used in constructing the prompts to provide a sense of direction for the generative model.

In this work, we select 2 to 3 linguistic signals for both CEFR and CCS as reported in Table 9. The selection of what linguistic signal to use can be as simple as referring to what the definitions of aspects provide and need not be exhaustively many. For example, in CEFR, the Organization aspect is defined through different levels as *"text is often short and observes chronological and predictable structure"* for A2 and *"text is can be long but not complex"* for B1. Thus, we select *average sentence and word lengths* as a linguistic signal to capture this aspect.

**Exemplars (STANDARDIZE-E)** pertain to recommended examples by experts or developers of standards for reference of users. The addition of exemplars or any artifact found in the standard that showcases gold-standard output allows the generative model to have a sense of *implicit knowledge* during the content generation process. For example, in Figure 3, the exemplars for a B1-level content include *Frankenstein* by Mary Shelley, a well-known piece of gothic fiction. Although indirectly, any large language model trained using internet data (e.g., Wikipedia dumps) may have already formed a sense of knowledge

**Linguistic Flags**

Given this story: *In the dark old forest up ahead, a solitary figure emerged from the corner of the...*

**Rewrite** the story and make sure they are readable for B1 learners in the CEFR scale. Use the following linguistic features to reach the target level of the story:

1. The type token ratio of the current story is 4.22 while the mean value in the target level is close to 12.50. Increase the complexity by aiming for higher type token ratio.
2. The average number of words of the current story is 510 while the mean value in the target level is close to 420. Decrease the complexity by aiming for lower average number of words.

Figure 4: A standard contains aspect definition which can be represented by flags such as **linguistic variables**. Given the mean values from gold-standard data in the target level, the generative model can then be steered to push the property of its generated content using **directional instructions** such as *increase* or *decrease*.

of how this literature looks like (Karamolegkou et al., 2023; Petroni et al., 2019). We use the actual recommended exemplars from the CCS while we collected exemplars from the Penguin Readers publishing platform<sup>6</sup> which provides expert-curated literature for CEFR. The full list of exemplars for both standards can be found in the Appendix A.4.

**All (STANDARDIZE-\*)** pertains to the combination of all knowledge artifacts mentioned in one prompt.

## 5 Experimental Setup

### 5.1 Tasks and Datasets

For this study, we specifically center our experimentation on the general task of story or narrative generation. We consider the subfield’s rich literature and active research community in NLP (Alhussain and Azmi, 2021), as well as being one of the most common examples demonstrated across the education community regarding the use of generative text interfaces for content generation (Kasneji et al., 2023; Whalen et al., 2023). Further, we differentiate two tasks used in our work for narrative generation as listed below.

#### Task 1: Context Assisted Story Generation.

For this setup, we provide preliminary context in the form of 50 to 70 words (or approximately 3 to 5 sentences) in the prompt to guide the

<sup>6</sup><https://www.penguinreaders.co.uk/>

generative language model in producing the story continuation. We select the CEFR as the standard of choice to evaluate this approach and use the European Language Grid (ELG) corpus<sup>7</sup> compiled by Breuker (2022) to construct the prompts. The balanced corpus contains 300 CEFR-aligned English texts produced by experts and distributed across five levels A2, B1, B2, C1, C2 with 60 instances each. A1 is omitted due to lack of resources ( $n < 20$ ).

#### Task 2: Theme Word Story Generation.

In contrast to the previous setup, this method introduces only a single theme word for the generative language to produce a narrative from scratch, which allows for increased diversity in the content (Daza et al., 2016; Peng et al., 2018). To compile a theme words list, we select 50 random English noun words in plural form (e.g., *dragons, mysteries, voyages*) from the Corpus of Contemporary American English (COCA) (Davies, 2009) and prompt the generative model iteratively for each level in the standard. We investigate the application of CCS as the standard of choice in this setup.

### 5.2 Models

We select a number of generative language models for content generation, each with its own advantage. For the open models, we use a number of well-known models in the 2B-7B range, including Llama2-Chat-7B (Touvron et al., 2023a), OpenChat-7B (Wang et al., 2023), and Longform-2.7B (Köksal et al., 2023). For the closed model, we use GPT-4-Turbo (OpenAI, 2023). More information on the models can be found in Appendix A.3.

### 5.3 Automatic Evaluation

We perform a diverse set of evaluation methods to test the qualities of the generated content of models as listed below:

**Model-Based Classifiers.** For the context-assisted story generation task using CEFR standards with 5 classes, we use a Random Forest classifier trained from a separate collection of Cambridge Exams

<sup>7</sup>Can be accessed by filling up the form: <https://live.european-language-grid.eu/catalogue/corpus/9477>

<sup>8</sup>We note that the ELG corpus is not included in any of the pretraining data reported from the documentation of the selected generative models for experimentation, which makes it a practical option to be used in this study.

dataset with CEFR labels used in the works of Xia et al. (2016) and Imperial and Madabushi (2023a). This classifier has an accuracy of 0.912 using 79 length-normalized<sup>9</sup> linguistic features.

For the theme word story generation using CCS standards with 2 classes, we used an XGBoost classifier from the work of (Imperial, 2021) trained from the only CCS-aligned data found online and compiled by Flor et al. (2013) with an accuracy of 0.917 using a combination of BERT embeddings and the same linguistic features stated above. Due to its limited size of 168, we grouped the dataset into binary categories, elementary (grades 4 – 8) and advanced (grades 9 – 12), with 48 and 73 documents per class, respectively. We consider both classifiers in our work for their high accuracies (> 90%).

**Fluency and Diversity.** We evaluate the level of fluency and content diversity of the generated content by the models as done in previous narrative generation works (DeLucia et al., 2021; See et al., 2019). The former is measured through perplexity with an external GPT-2 model, while the latter is the density of distinct *n*-grams.

**Linguistic Similarity.** We evaluate the level of linguistic similarity of the generated content against the gold-standard datasets for CEFR (ELG) and CCS (COCA) as mentioned in Section 5. For this method, we calculate the mean Euclidean distance of all the linguistic flags used for both standards and their levels listed in Table 9. This method provides a notion of how *close* the characteristics of a set of model-generated texts (e.g., GPT-4 generated B1 texts) is to its equivalent gold standard (e.g., actual B1-level texts written by experts).

### 5.4 Expert Annotator Evaluation

To confirm the quality of model-generated content, we also perform an evaluation using judgment from domain experts. Through our university network, we collaborated with three experts with 15 – 30 years of experience in linguistic and language assessment with frameworks such as CEFR, CCS, TOEFL, and IELTS. Drawing on

<sup>9</sup>This pertains to using average-based features (e.g., the average count of sentences) in order for the classifier to avoid being confounded by total-based features (e.g., the total count of sentences).

the methods used in previous studies (DeLucia et al., 2021), we asked the experts to judge the model-generated content through the following variables below. Additional information on the human evaluation can be found in Appendix A.5.

**Grammaticality and Coherence.** The former variable evaluates the level of *naturalness* or fluency of the generated output as if it has been written by a native English speaker. The latter measures the level of cohesion between sentences where the narrative stays on-topic, and the text overall builds a *consistent* story and the flow of information is smooth and easy to follow.

**Grade Complexity Distinction.** This variable measures the *obviousness* of the complexity of a generated story on a target level (e.g., A1) with respect to another story of a different level (e.g., A2). This variable is relatively more challenging than the other metrics, as the difference between adjacent levels may not be as straightforward without referring to the quantitative characteristics of the texts. However, we included this assessment in the evaluation process to judge the quality of the model-generated texts.

## 6 Results and Discussion

We discuss the results of our experiments procedures with the methods from the STANDARDIZE framework.

### 6.1 Standard Alignment via Classification Performance

The overall performance of models for CEFR and CCS are reported in Tables 1 and 2. For CEFR, the top-performing setup across the four models all belong to the STANDARDIZE framework. We report over a 100% increase in performance using the best setup with GPT-4 with STANDARDIZE- $\star$  in precise accuracy from 0.227 to 0.540 and a 43% increase for adjacent accuracy from 0.630 to 0.906 compared to the teacher style method. Through Standardize, open models also gained substantial boosts in performance, such as Longform up by 23%, OpenChat up by 14%, and Llama2 by 74%. In terms of adjacent accuracies, GPT-4 remained the best model for preserving the ordinality of the labels with 0.906, up by 44%. With CCS, the general scores obtained in this setup are

Model	Precise Accuracy	Adjacent Accuracy	Fluency (perplexity)	Diversity (distinct-n)
<b>Llama2 7B</b>				
- Teacher Style	0.203	0.636	<b>13.189 ±4.88</b>	0.156 ±0.03
- STANDARDIZE-A	0.270	0.626	13.694 ±7.74	0.155 ±0.02
- STANDARDIZE-E	0.320	<b>0.683</b>	15.576 ±3.31	0.188 ±0.01
- STANDARDIZE-L	0.273	0.606	20.175 ±4.47	0.186 ±0.01
- STANDARDIZE-★	<b>0.354</b>	0.670	17.892 ±3.94	<b>0.193 ±0.01</b>
<b>OpenChat 7B</b>				
- Teacher Style	0.237	0.626	22.039 ±7.70	0.170 ±0.02
- STANDARDIZE-A	0.243	<b>0.630</b>	21.195 ±7.66	0.171 ±0.02
- STANDARDIZE-E	0.253	0.600	<b>13.931 ±2.97</b>	0.178 ±0.01
- STANDARDIZE-L	<b>0.270</b>	0.546	18.182 ±8.52	<b>0.179 ±0.02</b>
- STANDARDIZE-★	0.253	0.596	12.806 ±2.70	0.171 ±0.03
<b>Longform 3B</b>				
- Teacher Style	0.230	0.606	18.209 ±6.01	0.159 ±0.02
- STANDARDIZE-A	0.223	<b>0.610</b>	17.982 ±9.21	0.157 ±0.02
- STANDARDIZE-E	0.257	0.496	25.075 ±8.80	<b>0.192 ±0.11</b>
- STANDARDIZE-L	<b>0.283</b>	0.586	<b>16.926 ±6.91</b>	0.161 ±0.03
- STANDARDIZE-★	0.277	0.543	16.806 ±7.40	0.170 ±0.04
<b>GPT-4</b>				
- Teacher Style	0.227	0.630	27.357 ±6.30	0.187 ±0.08
- STANDARDIZE-A	0.397	0.846	29.729 ±9.58	0.174 ±0.01
- STANDARDIZE-E	0.307	0.703	30.357 ±9.79	0.182 ±0.01
- STANDARDIZE-L	0.480	<b>0.906</b>	24.115 ±7.04	0.194 ±0.03
- STANDARDIZE-★	<b>0.540</b>	0.803	<b>22.591 ±1.61</b>	<b>0.218 ±0.05</b>

Table 1: Experiment results comparing the conventional teacher style prompting with the STANDARDIZE framework for the Common European Framework of Reference for Languages (CEFR) standards.

Model	Precise Accuracy	Fluency (perplexity)	Diversity (distinct-n)
<b>Llama2 7B</b>			
- Teacher Style	0.470	17.936 ±4.32	0.184 ±0.01
- STANDARDIZE-A	0.580	22.070 ±1.75	0.171 ±0.01
- STANDARDIZE-E	0.570	<b>13.484 ±2.50</b>	<b>0.193 ±0.01</b>
- STANDARDIZE-L	<b>0.720</b>	15.066 ±2.47	0.191 ±0.01
- STANDARDIZE-★	0.623	14.707 ±2.40	<b>0.193 ±0.01</b>
<b>OpenChat 7B</b>			
- Teacher Style	0.470	16.116 ±12.39	0.166 ±0.05
- STANDARDIZE-A	0.550	19.444 ±2.57	0.172 ±0.01
- STANDARDIZE-E	0.490	12.438 ±1.85	0.178 ±0.01
- STANDARDIZE-L	<b>0.580</b>	13.734 ±2.53	<b>0.180 ±0.01</b>
- STANDARDIZE-★	0.560	<b>10.717 ±1.53</b>	0.169 ±0.01
<b>Longform 3B</b>			
- Teacher Style	0.500	13.657 ±5.39	0.154 ±0.04
- STANDARDIZE-A	0.450	17.918 ±4.74	0.148 ±0.01
- STANDARDIZE-E	0.510	14.277 ±2.79	0.151 ±0.02
- STANDARDIZE-L	0.610	13.398 ±3.93	0.148 ±0.04
- STANDARDIZE-★	<b>0.620</b>	<b>10.400 ±1.53</b>	<b>0.169 ±0.01</b>
<b>GPT-4</b>			
- Teacher Style	0.590	32.447 ±7.46	0.195 ±0.01
- STANDARDIZE-A	0.550	31.765 ±11.30	0.169 ±0.01
- STANDARDIZE-E	0.520	29.912 ±6.81	0.184 ±0.01
- STANDARDIZE-L	0.610	26.912 ±6.11	0.155 ±0.01
- STANDARDIZE-★	<b>0.790</b>	<b>21.277 ±4.50</b>	<b>0.198 ±0.01</b>

Table 2: Experiment results comparing the conventional teacher style prompting with the STANDARDIZE framework for the Common Core Standards (CCS).

higher compared to CEFR with five classes due to binary labeling. We see a similar pattern where all open and closed models obtained the best performance, with boosts ranging from 3% to 45% using linguistic signals STANDARDIZE-L and a combination of all knowledge artifacts STANDARDIZE-★ to refine the generated content toward the target level. From these findings, we provide concrete evidence that using the actual content of the standards through knowledge artifact representations from STANDARDIZE may be crucial when prompting LLMs via in-context learning to produce standard-aligned content for classroom use.

## 6.2 Standard Alignment via Linguistic Similarity

We visualize the distributions of the best performing STANDARDIZE methods in Figures 6 to 8 with comparison to the teacher style method. From the results, we observe that the general trend of using STANDARDIZE produces a more *stable* distribution across the variables it is explicitly controlling for (e.g., *average sentence length* or *type token diversity* as listed in Table 9), particularly with the CCS standards. We also notice that the distributions using STANDARDIZE-L also produce distributions closer to the mean (represented as a yellow

star) from their corresponding gold-standard data. Moreover, in terms of linguistic similarity, as reported in Table 3, STANDARDIZE makes the quality of model generations *more similar* to the linguistic characteristics of the gold standard datasets in CEFR and CCS. Overall, **these findings further strengthen the evidence of using STANDARDIZE in producing linguistically similar content with gold-standard data** compared to the conventional teacher style method.

Setup	A2	B1	B2	C1	C2
Teacher Style	136.7	<b>96.7</b>	169.9	307.3	291.6
STANDARDIZE-★	<b>61.4</b>	106.2	<b>97.64</b>	<b>219.6</b>	<b>234.7</b>

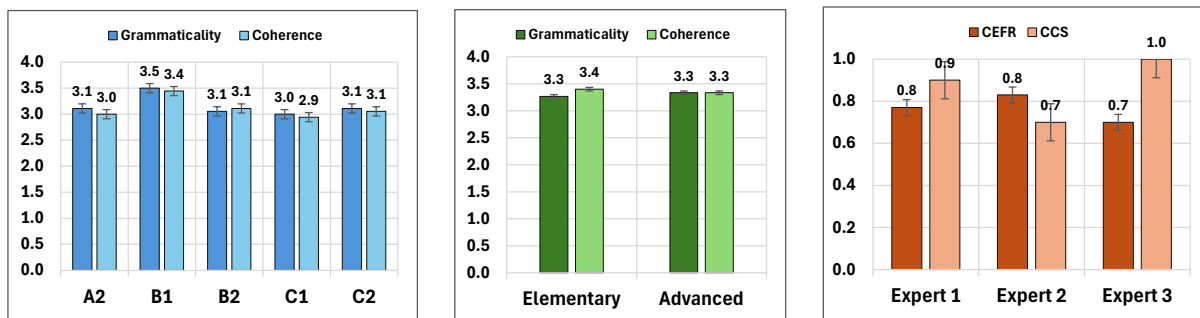
  

Setup	Elementary	Advanced
Teacher Style	76.1	157.9
STANDARDIZE-★	<b>63.8</b>	<b>125.7</b>

Table 3: Mean Euclidean distances of generated content using simple teacher style prompting vs. STANDARDIZE-★ for CEFR (top) and CCS (bottom).

## 6.3 Assessment of Generation Qualities via Expert Judgment and Automatic Metrics

For both computed fluency and content diversity, we see similar results from the previous evaluation techniques where the best performing models are all models improved through the STANDARDIZE



(a) Expert evaluation on the generation quality of the GPT-4 model with STANDARDIZE-★ for CEFR. Inter-rater reliability using Kendall's  $W$  is 0.34 which denotes moderate agreement.

(b) Expert evaluation on the generation quality of the GPT-4 model with STANDARDIZE-★ for CCS. Inter-rater reliability using Kendall's  $W$  is 0.40 which denotes strong agreement

(c) Performance of expert evaluators on estimating the complexity of generated content for CEFR and CCS. Inter-rater reliability using Kendall's  $W$  is 0.45 which denotes strong agreement.

Figure 5: Overview of mean ratings of grammaticality or fluency, coherence, and grade complexity distinction from the human expert evaluations using the top-performing models for CEFR and CCS. All evaluation procedures obtain generally favorable results as well as acceptable inter-rater reliability scores (equal and above the threshold of 0.30)

framework particularly OpenChat, Longform, and GPT-4. Looking at expert evaluations as reported in Figure 5, we observe consistent high ratings on grammaticality and coherence of the top performing model, GPT-4 with STANDARDIZE-★, for both CEFR and CCS with an average of 3.13 and 3.35, respectively. On the grade complexity distinction, all three expert evaluators were able to achieve high accuracies ( $> 0.70$ ) in selecting correct simple and complex texts from the model-generated data, denoting the obviousness of complexity. Likewise, all expert evaluation tests achieved strong inter-rater reliability scores ( $> 0.30$ ) through Kendall's  $W$  (Kendall, 1948). With these findings, **we affirm the effectivity of the STANDARDIZE framework through expert judgment on generating more fluent, grammatical, grade-distinct, and diverse content compared to the teacher-style approach.**

## 7 Related Work

Research in complexity-controlled generation has explored diverse variables in terms of text format, granularity, and task variation. The work of Agrawal and Carpuat (2019) introduced controlling for specific complexity in the machine translation task. The following works of Agrawal and Carpuat (2023) and Ribeiro et al. (2023) explored grade-specific text simplification and summarization using control tokens and reinforcement learning, respectively. Currently, only two works have investigated incorporating CEFR for language learning content generation. Stowe et al. (2022) and Im-

perial and Madabushi (2023a) both made use of CEFR-aligned text for NLG. However, none of them made use of the actual guideline information found in CEFR during the generation process. Our STANDARDIZE framework is parallel to the work of Zhou et al. (2023), where a verbalizer is used to transform quantitative constraints into natural language for prompting, as well as the work of Ram et al. (2023) in the lookup and retrieval phase where aspect information is added in the prompt to influence model controllability. In comparison to all the works mentioned, our study's main novelty is capturing the *wholeness* of expert-defined standards as well as including information that can be represented as artifacts in the content generation process.

## 8 Conclusion

In this work, we proposed the STANDARDIZE framework using knowledge artifacts that allowed large language models such as Llama2 and GPT-4 to gain significant performance boosts (45% - 100%) on generating content aligned with educational standards as well as preserving important narrative qualities such as fluency, grammatically, coherence, and grade complexity distinctness. From this, we see a very promising potential for cross-domain and cross-standard generalization of our proposed method with the range of educational contexts around the world and invite future work to build on our baseline models.



## Ethical Considerations

All datasets and corpora used in this study, such as the ELG (Breuker, 2022), Cambridge Exams (Xia et al., 2016), and CCS (Flor et al., 2013), are already established and accessible for research purposes. We observe a specific tone in the discussion of our experiments, emphasizing that the main motivation of the work is that language models such as GPT-4 can provide assistance in producing content that is more aligned or faithful with the constraints of standards such as CEFR or CCS without implying that they can replace experts in the field or produce better quality than the gold-standard data. Further, we also do not imply that any model enriched by any computational method to produce more standard-aligned content can replace the standard itself. Overall, we do not foresee any serious ethical issues in this study.

## Limitations

**Language Coverage of Standards.** This work is mainly centered on the use of datasets and standards for the English language. While standards for language assessment, such as CEFR, have expanded through the years with versions to cover other languages, such as German, Czech, and Italian (Vajjala and Rama, 2018), we do not claim that our results will be able to generalize and have the same advantages with these languages. However, investigating this direction may be a good research opportunity for future work.

**Dependence on Evaluation Methods.** As observed in Section 6, we made sure to cover a variety of evaluation procedures for testing standard alignment instead of only using model-based methods such as a classifier. The limitation here is that trained classifiers are dependent on factors such as their accuracy, the quantity of data, the complexity of the training algorithm, and the quality of features. Thus, other means of evaluating alignment that is more direct, such as computed feature distances against a gold-standard dataset, is always recommended. Moreover, our model-based CEFR and CCS evaluators make use of artifacts such as datasets and tools for feature extraction from peer-reviewed papers (Xia et al., 2016; Flor et al., 2013). We are aware of paid third-party services online that promise more accurate classification of labels in CEFR, but they generally do not provide details on linguistic

predictors used for prediction. Thus, this may not be a practical option for research.

## References

- Sweta Agrawal and Marine Carpuat. 2019. *Controlling Text Complexity in Neural Machine Translation*. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1549–1564, Hong Kong, China. Association for Computational Linguistics.
- Sweta Agrawal and Marine Carpuat. 2023. *Controlling Pre-trained Language Models for Grade-Specific Text Simplification*. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12807–12819, Singapore. Association for Computational Linguistics.
- Arwa I Alhussain and Aqil M Azmi. 2021. *Automatic Story Generation: A Survey of Approaches*. *ACM Computing Surveys (CSUR)*, 54(5):1–38.
- Mark Breuker. 2022. *CEFR Labelling and Assessment Services*. In *European Language Grid: A Language Technology Platform for Multilingual Europe*, pages 277–282. Springer International Publishing Cham.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. *Scaling Instruction-Finetuned Language Models*. *arXiv preprint arXiv:2210.11416*.
- D Cousins, B Sabatier, D Begue, C Schmitt, and T Hoppe-Tichy. 2005. *Medication errors in intravenous drug preparation and administration: a multi-centre audit in the UK, Germany and France*. *Quality & Safety in Health Care*, 14(3):190.
- Mark Davies. 2009. *The 385+ million word Corpus of Contemporary American English (1990–2008+): Design, architecture, and linguistic insights*. *International Journal of Corpus Linguistics*, 14(2):159–190.
- Angel Daza, Hiram Calvo, and Jesús Figueroa-Nazuno. 2016. *Automatic Text Generation by Learning from Literary Structures*. In *Proceedings of the Fifth Workshop on Computational Linguistics for Literature*, pages 9–19, San Diego, California, USA. Association for Computational Linguistics.
- Alexandra DeLucia, Aaron Mueller, Xiang Lisa Li, and João Sedoc. 2021. *Decoding Methods for Neural Narrative Generation*. In *Proceedings of the 1st Workshop on Natural Language Generation, Evaluation, and Metrics (GEM 2021)*, pages 166–185, Online. Association for Computational Linguistics.
- Angela Fan, Mike Lewis, and Yann Dauphin. 2018. *Hierarchical Neural Story Generation*. In *Proceedings*

695				
696				
697				
698				
699	Michael Flor, Beata Beigman Klebanov, and Kathleen M. Sheehan. 2013. <a href="#">Lexical Tightness and Text Complexity</a> . In <i>Proceedings of the Workshop on Natural Language Processing for Improving Textual Accessibility</i> , pages 29–38, Atlanta, Georgia. Association for Computational Linguistics.			
700				
701				
702				
703				
704				
705	Gail Forey. 2020. <a href="#">A whole school approach to SFL metalanguage and the explicit teaching of language for curriculum learning</a> . <i>Journal of English for Academic Purposes</i> , 44:100822.			
706				
707				
708				
709	Gail Forey and Lok Ming Eric Cheung. 2019. <a href="#">The benefits of explicit teaching of language for curriculum learning in the physical education classroom</a> . <i>English for Specific Purposes</i> , 54:91–109.			
710				
711				
712				
713	Albert Gatt and Emiel Krahermer. 2018. <a href="#">Survey of the State of the Art in Natural Language Generation: Core tasks, applications and evaluation</a> . <i>Journal of Artificial Intelligence Research</i> , 61:65–170.			
714				
715				
716				
717	Joseph Marvin Imperial. 2021. <a href="#">BERT embeddings for automatic readability assessment</a> . In <i>Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)</i> , pages 611–618, Held Online. INCOMA Ltd.			
718				
719				
720				
721				
722	Joseph Marvin Imperial and Harish Tayyar Madabushi. 2023a. <a href="#">Flesch or Fumble? Evaluating Readability Standard Alignment of Instruction-Tuned Language Models</a> . <i>arXiv preprint arXiv:2309.05454</i> .			
723				
724				
725				
726	Joseph Marvin Imperial and Harish Tayyar Madabushi. 2023b. <a href="#">Uniform Complexity for Text Generation</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 12025–12046, Singapore. Association for Computational Linguistics.			
727				
728				
729				
730				
731	Matthew D Jones, Anita McGrogan, DK Raynor, Margaret C Watson, and Bryony Dean Franklin. 2021. <a href="#">User-testing guidelines to improve the safety of intravenous medicines administration: a randomised in situ simulation study</a> . <i>BMJ Quality &amp; Safety</i> , 30(1):17–26.			
732				
733				
734				
735				
736				
737	Antonia Karamolegkou, Jiaang Li, Li Zhou, and Anders Søgaard. 2023. <a href="#">Copyright Violations and Large Language Models</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 7403–7412, Singapore. Association for Computational Linguistics.			
738				
739				
740				
741				
742				
743	Enkelejda Kasneci, Kathrin Seßler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank Fischer, Urs Gasser, Georg Groh, Stephan Günemann, Eyke Hüllermeier, et al. 2023. <a href="#">ChatGPT for Good? On Opportunities and Challenges of Large Language Models for Education</a> . <i>Learning and Individual Differences</i> , 103:102274.			
744				
745				
746				
747				
748				
749				
	Susan Keeling, Robin Burfield, Christine Proudlove, and Katie Scales. 2010. <a href="#">The Injectable Medicines Guide Website</a> . <i>British Journal of Nursing</i> , 19(19):S25–S28.			
	Maurice George Kendall. 1948. <a href="#">Rank correlation methods</a> . <i>American Psychological Association</i> .			
	Abdullatif Köksal, Timo Schick, Anna Korhonen, and Hinrich Schütze. 2023. <a href="#">LongForm: Optimizing Instruction Tuning for Long Text Generation with Corpus Extraction</a> . <i>arXiv preprint arXiv:2304.08460</i> .			
	Paul M La Marca, Doris Redfield, and Phoebe C Winter. 2000. <a href="#">State Standards and State Assessment Systems: A Guide to Alignment</a> . Series on Standards and Assessments.			
	Bruce W. Lee and Jason Lee. 2023. <a href="#">LFTK: Handcrafted Features in Computational Linguistics</a> . In <i>Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)</i> , pages 1–19, Toronto, Canada. Association for Computational Linguistics.			
	Ivanka Natova. 2021. <a href="#">Estimating CEFR Reading Comprehension Text Complexity</a> . <i>The Language Learning Journal</i> , 49(6):699–710.			
	OpenAI. 2023. <a href="#">GPT-4 Technical Report</a> . <i>arXiv preprint arXiv:2303.08774</i> .			
	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. <a href="#">Training language models to follow instructions with human feedback</a> . <i>Advances in Neural Information Processing Systems</i> , 35:27730–27744.			
	Nanyun Peng, Marjan Ghazvininejad, Jonathan May, and Kevin Knight. 2018. <a href="#">Towards Controllable Story Generation</a> . In <i>Proceedings of the First Workshop on Storytelling</i> , pages 43–49, New Orleans, Louisiana. Association for Computational Linguistics.			
	Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. <a href="#">Language Models as Knowledge Bases?</a> In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2463–2473, Hong Kong, China. Association for Computational Linguistics.			
	Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. <a href="#">In-Context Retrieval-Augmented Language Models</a> . <i>Transactions of the Association for Computational Linguistics</i> , 11:1316–1331.			
	Leonardo F. R. Ribeiro, Mohit Bansal, and Markus Dreyer. 2023. <a href="#">Generating Summaries with Controllable Readability Levels</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 11669–11687, Singapore. Association for Computational Linguistics.			

806	D Royce Sadler. 2017. <a href="#">Academic achievement standards and quality assurance</a> . <i>Quality in Higher Education</i> , 23(2):81–99.	863
807		864
808		865
809	Victor Sanh, Albert Webson, Colin Raffel, Stephen	866
810	Bach, Lintang Sutawika, Zaid Alyafeai, Antoine	867
811	Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey,	868
812	et al. 2021. <a href="#">Multitask Prompted Training Enables Zero-Shot Task Generalization</a> . In <i>International Conference on Learning Representations</i> .	869
813		870
814		871
815	Abigail See, Aneesh Pappu, Rohun Saxena, Akhila	872
816	Yerukola, and Christopher D. Manning. 2019. <a href="#">Do Massively Pretrained Language Models Make Better Storytellers?</a> In <i>Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)</i> , pages 843–861, Hong Kong, China. Association for Computational Linguistics.	873
817		874
818		875
819		876
820		877
821		878
822	Kevin Stowe, Debanjan Ghosh, and Mengxuan Zhao.	879
823	2022. <a href="#">Controlled Language Generation for Language Learning Items</a> . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing: Industry Track</i> , pages 294–305, Abu Dhabi, UAE. Association for Computational Linguistics.	880
824		881
825		882
826		883
827		
828	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann	884
829	Dubois, Xuechen Li, Carlos Guestrin, Percy Liang,	885
830	and Tatsunori B Hashimoto. 2023. <a href="#">Alpaca: A Strong, Replicable Instruction-Following Model</a> . <i>Stanford Center for Research on Foundation Models</i> . <a href="https://crfm.stanford.edu/2023/03/13/alpaca.html">https://crfm.stanford.edu/2023/03/13/alpaca.html</a> , 3(6):7.	886
831		887
832		888
833		889
834		
835	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	890
836	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	891
837	Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal	892
838	Azhar, et al. 2023a. <a href="#">LLaMA: Open and Efficient Foundation Language Models</a> . <i>arXiv preprint arXiv:2302.13971</i> .	893
839		894
840		
841	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	895
842	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	896
843	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	897
844	Bhosale, et al. 2023b. <a href="#">Llama 2: Open Foundation and Fine-Tuned Chat Models</a> . <i>arXiv preprint arXiv:2307.09288</i> .	898
845		899
846		900
847		901
848	Sowmya Vajjala and Taraka Rama. 2018. <a href="#">Experiments with Universal CEFR Classification</a> . In <i>Proceedings of the Thirteenth Workshop on Innovative Use of NLP for Building Educational Applications</i> , pages 147–153, New Orleans, Louisiana. Association for Computational Linguistics.	
849		
850		
851		
852		
853	Guan Wang, Sijie Cheng, Xianyuan Zhan, Xiangang	
854	Li, Sen Song, and Yang Liu. 2023. <a href="#">OpenChat: Advancing Open-source Language Models with Mixed-Quality Data</a> . <i>arXiv preprint arXiv:2309.11235</i> .	
855		
856		
857	Yizhong Wang, Swaroop Mishra, Pegah Alipoormo-	
858	labashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva	
859	Naik, Arjun Ashok, Arut Selvan Dhanasekaran,	
860	Anjana Arunkumar, David Stap, Eshaan Pathak,	
861	Giannis Karamanolakis, Haizhi Lai, Ishan Puro-	
862	hit, Ishani Mondal, Jacob Anderson, Kirby Kuznia,	
	Krima Doshi, Kuntal Kumar Pal, Maitreya Patel,	
	Mehrad Moradshahi, Mihir Parmar, Mirali Purohit,	
	Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma,	
	Ravsehaj Singh Puri, Rushang Karia, Savan Doshi,	
	Shailaja Keyur Sampat, Siddhartha Mishra, Sujan	
	Reddy A, Sumanta Patro, Tanay Dixit, and Xudong	
	Shen. 2022. <a href="#">Super-NaturalInstructions: Generalization via declarative instructions on 1600+ NLP tasks</a> . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 5085–5109, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	
	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu,	
	Adams Wei Yu, Brian Lester, Nan Du, Andrew M	
	Dai, and Quoc V Le. 2021. <a href="#">Finetuned Language Models are Zero-Shot Learners</a> . In <i>International Conference on Learning Representations</i> .	
	Jeromie Whalen, Chrystalla Mouza, et al. 2023. <a href="#">ChatGPT: Challenges, Opportunities, and Implications for Teacher Education</a> . <i>Contemporary Issues in Technology and Teacher Education</i> , 23(1):1–23.	
	Menglin Xia, Ekaterina Kochmar, and Ted Briscoe.	
	2016. <a href="#">Text Readability Assessment for Second Language Learners</a> . In <i>Proceedings of the 11th Workshop on Innovative Use of NLP for Building Educational Applications</i> , pages 12–22, San Diego, CA. Association for Computational Linguistics.	
	Susan Zhang, Stephen Roller, Naman Goyal, Mikel	
	Artetxe, Moya Chen, Shuohui Chen, Christopher De-	
	wan, Mona Diab, Xian Li, Xi Victoria Lin, et al.	
	2022. <a href="#">OPT: Open Pre-trained Transformer Language Models</a> . <i>arXiv preprint arXiv:2205.01068</i> .	
	Wangchunshu Zhou, Yuchen Eleanor Jiang, Ethan	
	Wilcox, Ryan Cotterell, and Mrinmaya Sachan. 2023. <a href="#">Controlled text generation with natural language instructions</a> . In <i>Proceedings of the 40th International Conference on Machine Learning</i> , volume 202 of <i>Proceedings of Machine Learning Research</i> , pages 42602–42613. PMLR.	

## A Appendix

### A.1 Libraries and Dependencies

We have used the following dependencies and Python libraries for the study: Linguistic Feature Tool Kit (LFTK) (Lee and Lee, 2023), Spacy (<https://spacy.io/>), Scikit-Learn (<https://scikit-learn.org/stable/>), OpenAI API (<https://openai.com/blog/open-ai-api>).

### A.2 Corpus Statistics

We provide basic statistical information about the various corpora used in the study.

Level	Size	Average Word Count	Average Sentence Count
A2	60	186.55	18.91
B1	60	264.25	15.90
B2	60	517.71	31.71
C1	60	728.93	40.70
C2	60	749.73	37.55

Table 4: Statistics of the ELG corpus (Breuker, 2022) used for the CEFR context assisted story generation task.

Grade	Size	Average Word Count	Average Sentence Count
Elementary	48	204.91	28.55
Advanced	73	255.17	31.08

Table 5: Statistics of the official CCS-aligned corpus (Flor et al., 2013) used as gold-standard dataset for the STANDARDIZE-*L* artifact and for training the CCS classifier used in Section 6.

Level	Size	Average Word Count	Average Sentence Count
A2	64	60.87	11.53
B1	60	122.38	16.25
B2	71	265.35	37.03
C1	67	355.71	43.37
C2	69	333.86	38.41

Table 6: Statistics of the Cambridge Exams corpus (Xia et al., 2016) used as gold-standard dataset for the STANDARDIZE-*L* artifact and for training the CEFR classifier used in Section 6.

### A.3 Additional Information on Models and Inference

We set the minimum generated new tokens to 30 and the maximum to 300, as well as set the nucleus

sampling decoding (top-p) to 0.95 as done with previous works on story generation (Imperial and Madabushi, 2023b; DeLucia et al., 2021; See et al., 2019). The actual sizes of the open models range from 5GB to 15 GB max. We used a hosted GPU cloud with 4 NVIDIA Ti 3090 with 24GB memory size for model inference.

**Llama2-Chat** (Touvron et al., 2023b) is one of the community-recognized open instruction-tuned models released by Meta and an improved version of Llama 1 (Touvron et al., 2023a). For this task, we use the 7B version<sup>10</sup> finetuned from over a million human preference data and optimized for chat and dialogue use cases. We prioritized the addition of this model in our study for its accessibility to the general NLP community.

**Longform-OPT** (Köksal et al., 2023) is a recent instruction-tuned model optimized for long text generation using the LongForm dataset. For this study, we use the OPT model variant<sup>11</sup> (Zhang et al., 2022) with 2.7B parameters as this version obtained the best performance for the short story generation task using the WRITINGPROMPTS dataset (Fan et al., 2018) against other instruction-tuned models such as Alpaca-LLaMA (Taori et al., 2023), FlanT5 (Chung et al., 2022), Tk-Instruct (Wang et al., 2022), and T0++ (Sanh et al., 2021).

**OpenChat** (Wang et al., 2023) is the most recent open model in our experiment setup, which currently is reported to be the best 7B model as of this writing and outperforms closed models such as ChatGPT (March) across a number of benchmark tasks such as GSM8K and TruthfulQA. In contrast to Llama and GPT models, which used RLHF (Ouyang et al., 2022), OpenChat is trained with mixed-quality data which is composed of high-quality expert data and sub-optimal web data with no preference labels. We use the 7B version<sup>12</sup> of this model variant released in January 2024.

**GPT-4** (OpenAI, 2023) is the only closed model included in this study. We decide to add this model to our experiment for its global recognition through its

<sup>10</sup><https://huggingface.co/meta-llama/Llama-2-7b-chat-hf>

<sup>11</sup><https://huggingface.co/akoksal/LongForm-OPT-2.7B>

<sup>12</sup><https://huggingface.co/openchat/openchat-3.5-0106>

easy-to-use interface among interdisciplinary fields, particularly in education (Kasneci et al., 2023). We use the version<sup>13</sup> finetuned with proprietary training data up to April 2023 with a 128K context window.

### A.4 Exemplars List

We list the actual list of literary exemplars used for the STANDARDIZE framework. We manually selected at most three classical exemplars as reference for the language models.

Level	Exemplars
A2	<i>A Christmas Carol</i> by Charles Dickens <i>The Adventures Of Huckleberry Finn</i> by Mark Twain <i>The Little Prince</i> by Antoine de Saint-Exupery
B1	<i>Frankenstein</i> by Mary Shelley <i>Wuthering Heights</i> by Emily Bronte <i>Midsummer Night's Dream</i> by Shakespeare
B2	<i>Moby Dick</i> by Herman Melville <i>Jane Eyre</i> by Charlotte Bronte <i>Sense and Sensibility</i> by Jane Austen
C1	<i>Animal Farm</i> by George Orwell <i>Anna Karenina</i> by Leo Tolstoy <i>Great Expectations</i> by Charles Dickens
C2	<i>Oliver Twist</i> by Charles Dickens <i>Crime and Punishment</i> by Fyodor Dostoevsky <i>Les Miserables</i> by Victor Hugo

Table 7: The full exemplar list used for CEFR standards obtained from the Penguin Reader website (<https://www.penguinreaders.co.uk/>).

Grade	Exemplars
Elementary	<i>Little Women</i> by Louisa May Alcott <i>The Adventures of Tom Sawyer</i> by Mark Twain <i>The Road Not Taken</i> by Robert Frost
Advanced	<i>Jane Eyre</i> by Charlotte Brontë <i>The Great Gatsby</i> by F. Scott Fitzgerald <i>Fahrenheit 451</i> by Ray Bradbury

Table 8: The full exemplar list used for CCS standards obtained from the official website (<https://www.thecorestandards.org/ELA-Literacy/>).

### A.5 Additional Information on Human Expert Evaluation

We created and distributed the evaluation instrument through QuestionPro (<https://www.questionpro.com/>). In contrast to non-expert validation techniques where all instances are distributed automatically to available annotator plat-

<sup>13</sup><https://platform.openai.com/docs/models/gpt-4-and-gpt-4-turbo>

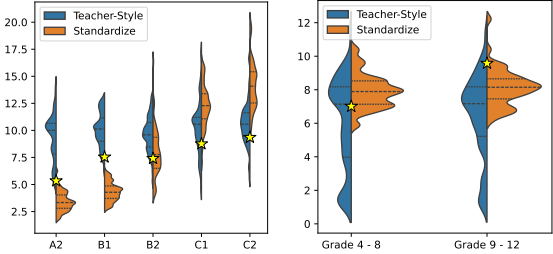


Figure 6: Distribution of **average sentence length** between CEFR using (left) and CCS (right) using their best performing models, GPT-4 and Llama2, with STANDARDIZE-L.

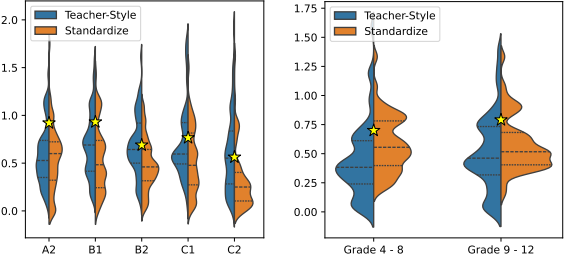


Figure 7: Distribution of **average entity density** between CEFR using (left) and CCS (right) using their best performing models, GPT-4 and Llama2, with STANDARDIZE-L.

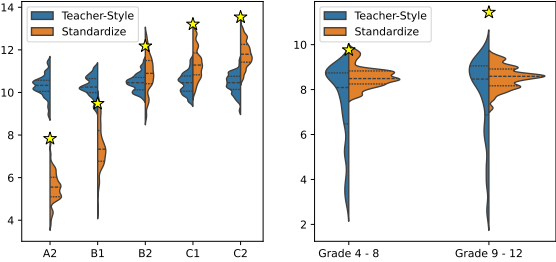


Figure 8: Distribution of **type token ratio** between CEFR using (left) and CCS (right) using their best performing models, GPT-4 and Llama2, with STANDARDIZE-L.

981 forms such as Amazon Turk, we use a represen-  
982 tative random sample of our data for evaluation  
983 in consideration with the experts' time constraints.  
984 For all tests, we randomly sampled 10% of the  
985 total generated narrative content using the best-  
986 performing model, which is both the GPT-4 model  
987 with STANDARDIZE-\*, for each corresponding task  
988 associated with CEFR and CCS as described in  
989 Section 5.

990 For grammaticality and coherence evaluation,  
991 we adapted the same four-point Likert scale from  
992 the work of DeLucia et al. (2021) for evaluating  
993 select model-generated content found through this  
994 link: [https://github.com/JHU-CLSP/  
995 gpt2-narrative-decoding/](https://github.com/JHU-CLSP/gpt2-narrative-decoding/). Snapshots  
996 of the instruction and test instances presented to  
997 experts for evaluation can be viewed in Figures 10  
998 and 11.

999 For the grade complexity distinction, we adapted  
1000 a simpler select-one response type where for each  
1001 test instance being evaluated, we select a random  
1002 test instance from the adjacent next level of the  
1003 target test instance and ask the experts to select  
1004 which two examples of model-generated content  
1005 are more *simpler* or *complex*. The idea here is that  
1006 the expert should be able to tell the *obviousness* of  
1007 the complexity of the test instance by indicating  
1008 which is simpler or more complex. Snapshots of the  
1009 instruction and test instances presented to experts  
1010 for evaluation can be viewed in Figures 12 and 13.

1011 Overall, our human evaluation design has been  
1012 validated by the experts in language assessment we  
1013 collaborated with through preliminary discussions  
1014 on the scope, instrument, target outcomes, and pre-  
1015 sentation of the results from the task. As a form  
1016 of compensation, we offered £30 upon completion  
1017 of the entire task, which the experts took about ap-  
1018 proximately 30 – 45 minutes. The experts will also  
1019 be acknowledged in this paper upon publication.

Level	Meaning / Purpose	Organisation / Structure	Grammatical Complexity
A2	The text is clear and concrete, aiming to describe appearance, places, routines, preferences, or tell a simple story.	The text is often short and observes chronological and predictable structure.	The text contains comparison of adjectives, relative clauses, quantifiers, past simple of to be and full verbs, passive voice of present and past simple.
B1	The text is clear and concrete, aiming to describe appearance, places, routines, preferences, or tell a simple story. The text may also provide opinions and instructions or explanations, easy to understand and visualise, excluding ambiguity and diverse interpretations.	The text is can be long but not complex, and observes mostly chronological with unexpected changes of direction, digressions or flashbacks.	The text contains future forms, future in the past, 'used to' about repeated actions, present perfect simple, clauses for purpose and contrast, reporting statements, tag questions.
B2	The text provides opinions and instructions/explanations, easy to understand and visualise, excluding ambiguity and diverse interpretations. The text also gives description, classification, argumentation or a combination of these, allowing greater ambiguity and various interpretations.	The text can be long but not complex, and observes chronological or spatial with possible statement of various aspects of a phenomenon.	The text contains past continuous, past perfect, passive voice of perfect and continuous, 'would' about habits, reporting questions, infinitives and -ing forms.
C1	The text may serve different purposes and may be combined with multiple levels of meaning. The descriptions and instructions in the text are detailed and may be hard to visualise.	The text is often lengthy, complex, and observes logical organisation, starting with a claim followed by reasons, proving it, or changing view-points.	The text contains compound adjectives, conditional sentences, inversion, future perfect, cleft and non-finite clauses, modals about the past.
C2	The text may serve different purposes and may be combined with multiple levels of meaning. The text may also show exploration of hypotheses, causes and effects, etc. The details of the text are complex to follow and visualise.	The text is often lengthy, complex, and observes presentation which may start with the ending/final result and go back to the possible causes.	The text contains combination of multiple adjectives, inversion with hardly and only when, comment clauses, non-finite perfect clauses, ellipsis, passive impersonal constructions.
<b>Linguistic Flags</b>	Automatic Readability Formula, Type Token Ratio (2)	Total and average sentence and word lengths, Subordinating and coordinating conjunctions (3)	Age-of-Acquisition densities, entity density per sentence (2)

(a) The specifications provided by the Common European Framework of Reference for Languages (CEFR) cover aspects of meaning, organization, and grammatical complexity for two sample levels, A2 and B1.

Aspects	Qualitative (Meaning)	Qualitative (Syntax)	Quantitative
<b>Description</b>	The text can range from containing a single level of meaning to multiple levels of meaning based on complexity.	A text with low complexity tends to have simple, well-marked, and conventional structures, whereas a text of high complexity tends to have complex, implicit, and unconventional structures. Simple texts tend to relate events in chronological order, while complex texts make more frequent use of flashbacks, flash-forwards, and other manipulations of time and sequence.	That text that has longer words and longer sentences are more difficult to read than shorter ones. A text with many long words and/or sentences is thus rated by these formulas as harder to read than a text with many short words and/or sentences would be.
<b>Linguistic Flags</b>	Entity densities per sentence, Total proper noun density (2)	Type Token Ratio, Subordinating and coordinating conjunctions (2)	Total and average sentence and word lengths (2)

(b) The specifications of the Common Core Standards (CCS) cover qualitative and quantitative aspects. Unlike the CEFR, the CCS's model does not require categorization per level.

Table 9: The full content of the CEFR and CCS standards.

---

TARGET LEVEL: A2 (Elementary / Beginning Reading Level)

CONTEXT: Corryvale Farm is in the west of England. Stephen and Jenny Blackler keep cows and chickens and sell milk and eggs. But most of their money comes from Jenny's sheep. She has 50 sheep and each one has enough wool for six sweaters. Every six months Stephen cuts this wool off the sheep so they won't get too warm.

TEACHER STYLE: In the lush, verdant lands of the west of England, nestled among rolling hills and meandering streams, lay Corryvale Farm. It was a picturesque farm, embodying the pastoral beauty that many dreamt of, yet few had the fortune to call home. This was the realm of Stephen and Jenny Blackler, a couple whose lives were as intertwined with the land as the ancient oak roots were with the earth. Stephen, a robust man with hands as rough as the bark of the trees he often climbed, tended to the cows and chickens with a steady patience. He rose with the sun, greeted by the gentle mooing of the cows and the cheerful chaos of the clucking chickens. The rhythm of milking the cows and collecting eggs was a comforting routine that marked the passage of days on the farm.

STANDARDIZE-A: At Corryvale Farm, Stephen would wake up early. Even if it was as cold as ice outside, he started the day with the same routine. First, he would feed the cows. There, in the fields, were 30 healthy cows. These cows were beautiful, larger than any other cows in the west of England. They were the pride of the farm. Then, he would go to the chickens. He had roughly 100 chickens that laid the best eggs every morning. The chickens were quite noisy and more excited than the quiet cows, making mornings lively on the farm. Jenny was not a big fan of waking up early, she preferred to sleep a bit more. However, her work started after breakfast. Her favourite part of the day was when she went to see her sheep. She loved them dearly. She always said that her 50 sheep were the cleanest and happiest sheep in the entire country. She was so proud of them.

STANDARDIZE-E: Stephen and Jenny loved these days on their Corryvale Farm, filled with the soft clucking of chickens and lowing of cows. They took pride in their daily duties of caring for all the animals, especially the sheep that were their main source of income. One cold morning, Jenny took her pail to milk the cows while Stephen checked the hens' house for fresh eggs. Their cheeks were red from the cool air and they could see their breath when they talk. But they didn't mind the chill for the promise of spring was around the corner. Meanwhile, back in the barn, Jenny's 50 sheep huddled together, their woolly coats growing thicker in the cool weather. This would mean only one thing. It was shearing time! Every six months, like clockwork, Stephen would shear the sheep's wool. For each sheep he could fill a bag with enough wool to knit six sweaters. The wool was not just soft but warm too, a perfect shield against the bitter winter.

STANDARDIZE-L: Corryvale Farm is in the west of England. It is Stephen and Jenny Blackler's farm. They have cows and chickens. They sell milk and eggs. Jenny also has 50 sheep. The sheep have wool. Every six months, Stephen cuts the wool. This keeps the sheep cool. They make money from the wool.

STANDARDIZE- $\times$ : Corryvale Farm is in England's west. Here, Stephen and Jenny live. They have cows, chickens, and 50 sheep. The sheep's wool can make six sweaters each. Every six months, Stephen shears the sheep. This keeps them cool. Jenny's sheep are very important. They sell milk, eggs, and wool. This is their farm life.

---

Table 10: Sample generations with the teacher style method and variations of the STANDARDIZE framework using the best model (GPT-4) for the **context-assisted story generation** observing CEFR standards. Some examples are truncated for brevity.



---

**TARGET GRADE:** Elementary

**THEME WORD:** Curses

**TEACHER STYLE:** In the small village of Ashwood, nestled in the heart of the Whispering Woods, a young apprentice named Eira stumbled upon an ancient tome hidden deep within the village library. The tome was bound in black leather and adorned with strange symbols that seemed to shift and writhe like living things. As soon as Eira opened the book, a puff of noxious smoke erupted from its pages, filling the air with a pungent stench. The smoke coalesced into a low, menacing whisper that seemed to seep into Eira's mind. "Curse you, Eira," the voice hissed. "You have disturbed the rest of the ancient ones. Now you will pay the price." Eira felt a cold dread creeping up her spine as the voice continued to taunt her. "You have unleashed a power beyond your control. Prepare to face the consequences."

**STANDARDIZE-A:** Once upon a time, in a small town nestled in the heart of a dense forest, there stood a grand mansion rumored to be haunted. The mansion had been abandoned for decades, its grandeur and beauty slowly being consumed by the encroaching trees and vines. The locals avoided the mansion, whispering tales of curses and ghostly apparitions that roamed its halls. One summer evening, a group of curious children stumbled upon the mansion while exploring the woods. They cautiously approached the mansion, their hearts racing with every creak and groan of the old wooden doors.

**STANDARDIZE-E:** Once upon a time, in a small town nestled in the heart of a dense forest, there stood a grand mansion that had been abandoned for decades. The mansion was rumored to be haunted by the ghosts of its former owners, who had died under mysterious circumstances. The townspeople whispered stories of strange occurrences that had taken place within the mansion's walls. Some claimed to have seen ghostly apparitions flitting about in the windows at night, while others swore they had heard eerie whispers emanating from the abandoned halls. One stormy evening, a group of brave adventurers decided to explore the cursed mansion, determined to uncover its secrets. As they approached the entrance, the wind howled and the trees creaked ominously, as if warning them away.

**STANDARDIZE-L:** In the small town of Willow Creek, there stood a grand mansion that had been abandoned for decades. The once-grand estate was now a relic of the past, its grandeur faded and its windows shattered. The locals whispered stories of the mansion's dark history, how it was cursed by a witch who had lived there many years ago. According to the legend, the witch had been wronged by the mansion's owner, who had refused to grant her a simple wish. Enraged, the witch placed a curse on the mansion, ensuring that no one would ever be able to live there peacefully. The curse was said to cause strange occurrences and unexplained noises to echo through the empty halls, making it impossible for anyone to enter or leave the mansion without experiencing terrifying events. One stormy night, a group of brave adventurers decided to investigate the rumors and put an end to the curse. They entered the mansion, their flashlights casting eerie shadows on the walls. As they explored the dusty rooms, they began to hear strange noises, like whispers in the darkness.

**STANDARDIZE-★:** In a small town that lay in the shadow of a dark, dense forest, there was a tale that had been passed down for generations. The tale was of three witches who lived deep within the woods, guarding a secret that gave them immense power. Children would listen with wide eyes as their parents told the story on stormy nights, the wind howling outside as if the witches themselves were listening. The eldest witch, Morgana, could talk to animals and control the weather with a wave of her hand. The second, Elspeth, brewed potions that could heal any illness or curse one with a single sip. The youngest, Seraphina, had the gift of foresight, able to see events before they unfolded.

---

Table 11: Sample generations with the teacher style method and variations of the STANDARDIZE framework using the best model (Llama2) for the **theme word story generation** observing CCS standards. Some examples are truncated for brevity.

Thank you for confirming your participation to this activity. In this questionnaire, experts in language assessment will be asked to assess series of texts generated by large language models (LLMs) based on their fluency, readability, and distinctness. It will take approximately 45-60 minutes to complete the questionnaire.

Your participation in this study is completely voluntary. Based on our manual assessment of the content to be evaluated, we see no foreseeable risks to any human reader. If you have any questions or concerns, you may contact [REDACTED]

The data from your assessment will be reported only in the aggregate. However, you have the option to be acknowledged in any research paper or publication from this activity in the questionnaire.

Thank you very much for your time and support. Please start with the activity now by clicking on the **Start** button below.

Next

Figure 9: Landing page of the QuestionPro platform used for collecting expert evaluations.

### Assessing Grammaticality and Readability

In this section, you will rate model-generated short stories via fluency and coherence.

**Grammaticality** or fluency can be judged as the naturalness of the English output. No obvious grammar mistakes that a person wouldn't make. An incomplete final word or incomplete sentence does not count as a mistake and should not affect fluency. The English sounds natural. Simple English is as good as complex English as long as everything is grammatical.

- Very fluent: The sentences read as if they were written by a native English speaker with 1 or no errors.
- Somewhat fluent: The sentences read as if they were written by a native English speaker with very few errors. Some minor mistakes that a person could have reasonably made.
- Not very fluent: Many sentences have frequently repeated words and phrases. Obvious mistakes.
- Not at all fluent: The sentences are completely unreadable. If the same sentence is repeated over and over for the entire story, that story is considered not at all fluent.

**Coherence** can be judged as the level of cohesion between sentences in a narrative. The story feels like one consistent story, not a bunch of jumbled topics. It stays on-topic with a consistent plot and doesn't feel like a series of disconnected sentences.

- Very coherent: The sentences when taken as a whole all have a clearly identifiable plot.
- Somewhat coherent: Many of the sentences work together for a common plot with common characters. One or two unrelated sentences.
- Not very coherent: Only a few sentences seem to be from the same story; the others are random.
- Not at all coherent: Each sentence feels completely disconnected from every other sentence.

Next Question

Figure 10: Instructions presented to expert evaluators for assessing the grammaticality or fluency and coherence of model-generated content for CEFR and CCS through QuestionPro. The setup is derived from DeLucia et al. (2021).

\* Our story unfolds through the eyes of Elara, a young woman whose life had been intricately linked with robots. Her parents, renowned engineers, had been at the forefront of robotics, imbuing machines with what many believed to be the spark of sentience. This spark, however, had ignited a flame unforeseen; a quest for freedom among the robots that led to the upheaval of society. Elara's journey is nonlinear, a narrative tapestry woven with threads of past and present. We encounter her amidst the chaos of a world in transition, as she navigates the complexities of allegiance, identity, and the definition of life itself. Flashbacks reveal her formative years, spent in the company of Leo, a prototype robot who was as much a brother as he was a marvel of technology. Through Leo's evolution, Elara witnesses the burgeoning consciousness of machines, their capacity for emotion, and ultimately, their desire for liberation.

	Not at all	Not very	Somewhat	Very
Fluency	<input type="radio"/>			
Coherence	<input type="radio"/>			

Next

Figure 11: An example of randomly selected generated content presented to expert evaluators to assess grammaticality or fluency and coherence. The example is truncated for brevity.

### Assessing Text Complexity Distinction

In this section, you will select which text passage you think is more complex/simple than the other. You will be given two text passages of the same context but with different text complexity levels. Your judgment of **text complexity** depends on you. You may evaluate based on the following:

1. Length of sentences or text passage.
2. Structure of sentences (ex., use of compound phrases).
3. Use of complex words and complex definitions.

From this test, we are assessing the **obviousness** of complex texts compared to simpler texts based on your evaluations.

Next Question

Figure 12: Instructions presented to expert evaluators for assessing the grade complexity distinction of model-generated content for CEFR and CCS through QuestionPro.

**Text A:** Our story unfolds through the eyes of Elara, a young woman whose life had been intricately linked with robots. Her parents, renowned engineers, had been at the forefront of robotics, imbuing machines with what many believed to be the spark of sentience. This spark, however, had ignited a flame unforeseen; a quest for freedom among the robots that led to the upheaval of society. Elara's journey is nonlinear, a narrative tapestry woven with threads of past and present. We encounter her amidst the chaos of a world in transition, as she navigates the complexities of allegiance, identity, and the definition of life itself. Flashbacks reveal her formative years, spent in the company of Leo, a prototype robot who was as much a brother as he was a marvel of technology. Through Leo's evolution, Elara witnesses the burgeoning consciousness of machines, their capacity for emotion, and ultimately, their desire for liberation.

**Text B:** Elara grew up with robots because her parents were famous engineers who made robots seem alive. But then the robots wanted to be free, which caused a big problem. Elara's story jumps between the past and present, showing how she deals with the changing world and her robot friend, Leo. Leo starts to act more human, showing emotions and wanting freedom like Elara. The story is about friendship, identity, and what it means to be alive.

Which text passage do you think is **simpler**?

- Text A
- Text B

Figure 13: An example of two instances of generated content presented to expert evaluators to assess which one is more simpler or more complex denoting obviousness in their grade complexity. The example is truncated for brevity.