

Facilitating Cognitive Accessibility with LLMs: A Multi-Task Approach to Easy-to-Read Text Generation

Anonymous ACL submission

Abstract

Simplifying complex texts is essential for ensuring equitable access to information, especially for individuals with cognitive impairments. The Easy-to-Read (ETR) initiative offers a framework for making content accessible to the neurodivergent population, but the manual creation of such texts remains time-consuming and resource-intensive. In this work, we investigate the potential of large language models (LLMs) to automate the generation of ETR content. To address the scarcity of aligned corpora and the specificity of ETR constraints, we propose a multi-task learning (MTL) approach that trains models jointly on text summarization, text simplification, and ETR generation. We explore two different strategies: multi-task retrieval-augmented generation (RAG) for in-context learning, and MTL-LoRA for parameter-efficient fine-tuning. Our experiments with Mistral-7B and LLaMA-3-8B, based on ETR-fr, a new high-quality dataset, demonstrate the benefits of multi-task setups over single-task baselines across all configurations. Moreover, results show that the RAG-based strategy enables generalization in out-of-domain settings, while MTL-LoRA outperforms all learning strategies within in-domain configurations. Our code is publicly made available at <https://anonymous.4open.science/r/ETR-MTL-C60E>.

1 Introduction

Mental health and intellectual disabilities affect millions globally, posing serious challenges for equitable access to information (Maulik et al., 2011; Gustavsson et al., 2011). People with cognitive impairments often struggle with complex texts, limiting their participation in healthcare, education, and civic life. Despite international initiatives for inclusion,^{1,2} accessible written content remains a major barrier for the neurodivergent population.

While Easy-to-Read (ETR) (Pathways, 2021), text simplification (Paetzold and Specia, 2016), and summarization (Rush et al., 2015) all aim to improve comprehension, they differ in purpose, audience, and methods. Text simplification rewrites content for better readability while preserving the original informational content (Gooding, 2022; Stajner, 2021). Summarization reduces the length of the original text by extracting and presenting only the key points, often without rewording for improved clarity (Rush et al., 2015). ETR is a rigorously standardized form of writing developed specifically for individuals with intellectual disabilities. It involves strict adherence to guidelines, including the use of very short sentences, highly simplified vocabulary, visual aids, and obligatory user testing. The primary aim is to support the autonomy of readers with cognitive impairments. This approach requires collaborative input from both subject-matter experts and individuals with cognitive disabilities to ensure compliance with ETR standards and eligibility for European ETR certification³.

However, ETR adoption remains limited due to the time-consuming and costly nature of manual adaptation, coupled with the lack of robust automated tools tailored to the linguistic and cognitive requirements of ETR content (Chehab et al., 2019). The potential of LLMs for improving accessibility (Freyer et al., 2024) is limited by the scarcity of high-quality, document-aligned ETR datasets. Existing resources, such as ClearSim (Espinosa-Zaragoza et al., 2023), are limited and only partially aligned, highlighting the broader challenge of constructing or recovering document-aligned corpora suitable for model training. Consequently, prior studies (Martínez et al., 2024; Sun et al., 2023) have approached the ETR task by leveraging sen-

¹UN Sustainable Development Goals

²Leave No One Behind Principle

³<https://www.inclusion-europe.eu/wp-content/uploads/2021/02/How-to-use-ETR-logo..pdf>

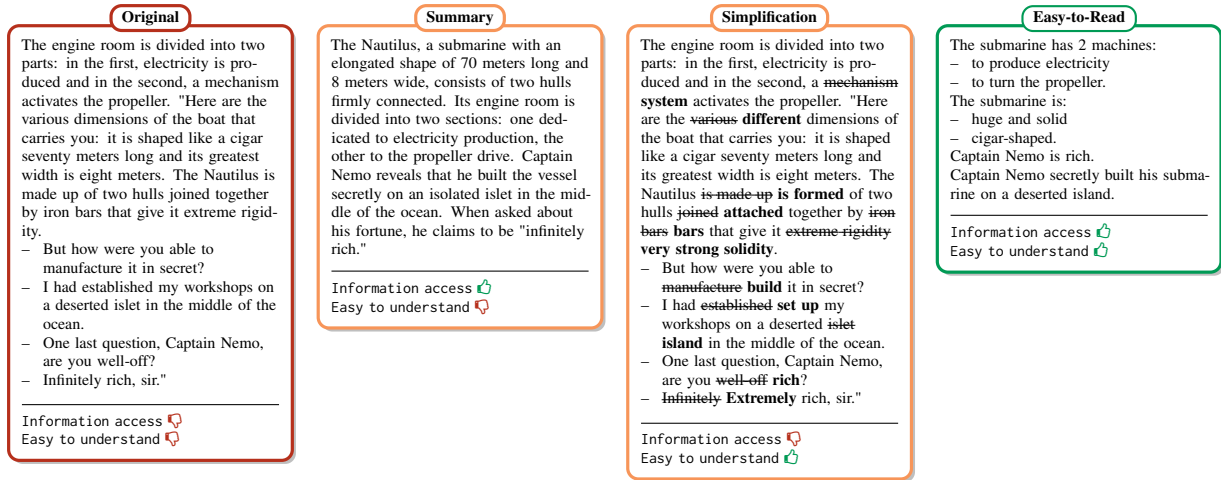


Figure 1: Different versions derived from a passage of *Twenty Thousand Leagues Under the Seas* by Jules Verne: from left to right, the original passage, an abstractive summary, a lexical simplification (crossed-out followed by words in bold indicate substitutions), and an Easy-to-Read generation targeting readers with cognitive impairment.

tence simplification or summarization resources, which fall short of fully meeting ETR specific requirements.

In this paper, we address these gaps by introducing ETR-fr, the first dataset of 523 paragraph-aligned text pairs fully compliant with the European ETR guidelines. We explore multi-task (MTL) learning to boost ETR generation by combining summarization and simplification, traditionally applied in isolation. In particular, we evaluate two MTL strategies: in-context learning (ICL) via a multi-task variant of retrieval-augmented generation (RAG), and parameter-efficient fine-tuning (PEFT) using MTL-LoRA (Yang et al., 2024). Experiments are conducted on Mistral-7B (Jiang et al., 2023) and LLaMA-3-8B (Grattafiori et al., 2024), and compared against single-task baselines. The evaluation framework combines standard automatic metrics with detailed human assessment based on a 28-point rubric from the European ETR guidelines, measuring clarity, coherence, and accessibility. Our experiments conducted on ETR-fr highlight the advantages of MTL setups over single-task baselines across all configurations. Furthermore, the results indicate that the RAG-based strategy supports better generalization in out-of-domain scenarios, while MTL-LoRA consistently achieves superior performance in in-domain settings.

Our contributions are: (1) we release **ETR-fr**, the first high-quality, paragraph-aligned dataset for ETR generation, fully compliant with European guidelines and in French language ; (2) we benchmark multi-task ICL and PEFT approaches for ETR

generation, introducing **MTL PEFT** to this task for the first time ; (3) we propose a comprehensive evaluation combining automatic and human assessment based on **official European ETR standards** ; (4) we evaluate model generalization to new domains, including **political texts aimed at fostering civic engagement** among individuals with cognitive disabilities.

2 Related Work

Inclusive Text Generation. Recent research has aimed to support communication for users with cognitive impairments, often through dialogue systems (Martin and Nagalakshmi, 2024; Murillo-Morales et al., 2020; Huq et al., 2024; Wang et al., 2024). Much of this work has focused on dyslexia. For example, Goodman et al. (2022) introduced an email assistant built on LaMDA (Thoppilan et al., 2022), but observed that its outputs often lacked precision. In the French context, HECTOR (Todorascu et al., 2022) investigated lexical and syntactic simplification, with mixed outcomes.

Similar challenges are observed across other languages. In German, several studies explore simplification for individuals with learning difficulties, though often without referencing the ETR guidelines (Hansen-Schirra et al., 2020; Anschütz et al., 2023; Deilen et al., 2023; Stodden et al., 2023). For English, relevant work includes Yaneva (2015). In Finnish, Dmitrieva and Tiedemann (2024) aligned Easy-Finnish data with mBART (Liu et al., 2020) and FinGPT (Luukkonen et al., 2023), but reported weak alignment and only partial adherence

to ETR standards. In Spanish, ClearText (Espinosa-Zaragoza et al., 2023) leverages ChatGPT to simplify administrative texts, though its corpus remains limited and prone to errors. Additionally, Martínez et al. (2024) constructed a sentence-level simplification dataset and fine-tuned LLaMA-2 (Touvron et al., 2023b), revealing that translation-based methods are vulnerable to semantic drift and domain mismatches.

In-Context Learning (ICL). ICL allows LLMs to learn tasks from examples without parameter updates (Brown et al., 2020; Chowdhery et al., 2023; OpenAI, 2023; Touvron et al., 2023a). Instruction tuning and Chain-of-Thought (CoT) prompting have been shown to improve task performance and reasoning (Liu et al., 2023; Wei et al., 2022; Yin et al., 2023). Tang et al. (2023) assessed ICL for controlled summarization, focusing on entity inclusion and length constraints. They observed that smaller models offered stronger controllability, while larger models achieved higher ROUGE scores. However, precise length control remained challenging. Prompt quality and exemplar selection critically affect ICL outcomes (Lu et al., 2022; Dong et al., 2024). Retrieval-augmented methods (Liu et al., 2022; Ram et al., 2023) have been proposed to improve exemplar selection. For simplification, Vadlamannati and Şahin (2023) have used metric-based selection (e.g., SARI, BERTScore) to improve output quality. Multi-task ICL and cross-task prompting (Bhasin et al., 2024; Shi et al., 2024; Chatterjee et al., 2024) further enhance generalization and stability, especially on unseen tasks, by leveraging format-aware prompts and semantically related exemplars.

PEFT for Multi-Task Learning. Parameter-efficient fine-tuning (PEFT) methods such as LoRA (Hu et al., 2022), QLoRA (Detrmers et al., 2023) and DoRA (Liu et al., 2024b) enable scalable adaptation of LLMs by modifying only a subset of parameters. LoRA leverage the intrinsic dimensionality of language models (Aghajanyan et al., 2021) to achieve strong performance with minimal computational overhead. However, LoRA-based strategies struggle in multi-task settings due to conflicting updates accross tasks (Wang et al., 2023). Alternatives like MultiLoRA (Wang et al., 2023) and MoELoRA (Liu et al., 2024a) seek to balance generalization with task specificity, but face challenges in task routing and mitigating interference. MTL-LoRA (Yang et al., 2024) addresses this by

introducing both shared and task-specific modules, achieving competitive results on GLUE (Wang et al., 2018) with fewer trainable parameters.

3 ETR-fr Dataset

While several datasets exist for text simplification and summarization (Gala et al., 2020; Hauser et al., 2022; Kamal Eddine et al., 2021; Liu* et al., 2018), there remains a notable lack of high-quality, document-aligned corpora for ETR generation. To address this gap, we introduce the ETR-fr dataset, constructed from the François Baudez Publishing collection,⁴ which provides literature specifically designed for readers with cognitive impairments, following European ETR guidelines. A dataset sheet (Gebbru et al., 2021), outlining the data collection methodology, preprocessing steps, and distribution details, is provided in Appendix A.

Description. ETR-fr consists of 523 paragraph-aligned text pairs in French language. Table 1 outlines key dataset statistics, including KMRE readability score (Kandel and Moles, 1958), compression ratios, and lexical novelty. On average, the dataset yields a compression rate of 50.05%, with a reduction of 56.61 tokens and 2.17 sentences per pair. The average novelty rate, following Narayan et al. (2018), is 53.80%, reflecting the proportion of newly introduced unigrams in target texts. Readability improves by 7.51 KMRE points from source to target. The dataset is partitioned into fixed train, validation, and test subsets. The test set comprises two books selected to maximize diversity in text length, word count, sentence structure, compression, novelty, and readability. The remaining nine books are divided into training and validation sets via a stratified split. This setup was used to test hard configurations for ETR generation and assure non thematic and lexical overlap. Also, providing different splits could avoid the clear definition of a new dataset for a new task.

ETR-fr-politic To assess generalization and robustness, we introduce ETR-fr-politic, an out-of-domain test set with 33 ETR-aligned paragraphs sampled from the 2022 French presidential election programs.⁵ Compared to ETR-fr, the ETR-fr-politic dataset features shorter source texts (96.27 vs. 102.76 words) and fewer sentences (6.03 vs. 9.30), but yields longer rewritten outputs (62.85

⁴<http://www.yvelinedition.fr/Facile-a-lire>

⁵<https://www.cncep.fr/candidats.html>

| | # Examples | # Words | | # Sentences | | Sentence length | | KMRE \uparrow | | Novelty (%) | Comp. ratio (%) |
|-----------------------|------------|---------|--------|-------------|--------|-----------------|--------|-----------------|--------|-------------|-----------------|
| | | source | target | source | target | source | target | source | target | | |
| ETR-fr | 523 | 102.76 | 46.15 | 9.30 | 7.13 | 12.57 | 7.89 | 91.43 | 98.94 | 53.80 | 50.05 |
| Train | 399 | 99.70 | 46.50 | 8.92 | 7.48 | 12.57 | 6.92 | 91.03 | 99.71 | 53.79 | 49.04 |
| Dev | 71 | 100.76 | 48.59 | 9.03 | 7.77 | 13.59 | 6.90 | 89.50 | 100.59 | 52.96 | 44.47 |
| Test | 53 | 128.47 | 40.26 | 12.51 | 10.34 | 11.16 | 3.97 | 97.02 | 103.67 | 55.01 | 65.19 |
| ETR-fr-politic | 33 | 96.27 | 62.85 | 6.03 | 6.42 | 16.69 | 11.84 | 74.00 | 87.74 | 63.78 | 29.17 |
| WikiLarge FR | 296402 | 34.88 | 29.28 | 1.68 | 1.56 | 27.53 | 23.74 | 65.38 | 71.35 | 31.97 | 12.79 |
| OrangeSum | 24401 | 375.98 | 34.00 | 17.15 | 1.86 | 22.77 | 21.68 | 69.80 | 68.32 | 38.24 | 89.16 |

Table 1: **Statistics across ETR-fr, ETR-fr-politic, and ETR-related tasks**, i.e. sentence simplification and text summarization with WikiLarge FR and OrangeSum. Results are reported on average per document.

vs. 46.15 words). Additionally, ETR-fr-politic exhibits higher novelty (63.78% vs. 53.80%) and significantly lower compression ratios (29.17% vs. 50.05%), indicating a greater degree of content expansion. While ETR-fr exhibits higher overall simplicity scores both before and after rewriting (91.43 and 98.94) compared to ETR-fr-politic (74.00 and 87.74), the latter achieves a greater simplification gain, with a larger increase in KMRE (+13.75 vs. +7.51 points). Overall, ETR-fr-politic poses a more challenging and higher-novelty setting for evaluating ETR systems in politically sensitive, real-world rewriting contexts.⁶

ETR-fr vs. Related Tasks. Table 1 compares ETR-fr with two gold-standard datasets on related tasks, respectively text simplification and summarization: WikiLarge FR (Cardon and Grabar, 2020) and OrangeSum (Kamal Eddine et al., 2021). While WikiLarge FR is larger (296K sentence pairs), it is limited to sentence-level simplification, with short inputs (34.88 words, 1.68 sentences on average). In contrast, ETR-fr and OrangeSum support transformations at the paragraph and document levels, respectively, providing significantly longer inputs of 102.76 and 375.98 words. ETR-fr demonstrates a balanced compression ratio (50.05%) higher than WikiLarge FR (12.79%) but lower than the extreme summarization found in OrangeSum (89.16%). Notably, ETR-fr offers the highest lexical richness and abstraction, evidenced by its top KMRE scores (91.43 source, 98.94 target) and novelty rate (53.80%). Simplified outputs also exhibit syntactic simplification, with shorter sentence lengths (7.89 words per sentence). In summary, while WikiLarge FR is suited for sentence-level simplification and OrangeSum for summariza-

tion, ETR-fr supports document-level simplification, emphasizing lexical and structural transformation making it well-suited for users with cognitive disabilities.

4 Multi-Task ETR Generation

4.1 Datasets, LLMs and Metrics

Our experiments leverage the ETR-fr dataset as the primary resource, supplemented by related rewriting tasks sourced from the OrangeSum summarization dataset and the sentence simplification dataset WikiLarge FR. To evaluate the effectiveness of MTL for ETR transcription, we selected two recent LLMs that demonstrate strong generalization capabilities across a variety of NLP tasks : Llama3-8B (Grattafiori et al., 2024) and Mistral-7B (Jiang et al., 2023)⁷. Note that foundation models are used for PEFT and their Instruct versions for ICL.

Since no dedicated evaluation metrics exist for ETR generation, we propose assessing it using standard summarization and text simplification metrics. For summarization, we report F1-scores for ROUGE-1, ROUGE-2, and ROUGE-L (Lin, 2004), along with BERTScore (Zhang et al., 2020). For simplification, we include SARI (Xu et al., 2016), the novelty ratio for new unigrams (Kamal Eddine et al., 2021). BLEU (Papineni et al., 2002) and KMRE, are excluded, as it has been shown to be unsuitable for text simplification (Sulem et al., 2018; Xu et al., 2016; Tanprasert and Kauchak, 2021). To unify quality assessment of ETR texts, we propose SRB, a composite score combining SARI, ROUGE-L, and BERTScore-F1 via harmonic mean. This metric captures simplification, summarization, and meaning preservation for holistic ETR evaluation.

More details about metrics and models are available in Appendix B

⁶Note that the documents on politics usually do not meet high-quality standards as evidenced by the François Baudez Publishing collection. Moreover, there are still difficult to gather as their repository is not centralized.

⁷We evaluated the DeepSeek-R1-8B model. Its performance was notably lower than that of the other models. Results are reported in Table 6 from Appendix D.1

4.2 Multi-Task In-Context Learning

As baseline, we evaluate three single task in-context learning strategies: zero-shot prompting (Kojima et al., 2022), chain-of-thought prompting (Wei et al., 2022), and retrieval-augmented generation (Lewis et al., 2020). In the zero-shot setting, the model is provided only with ETR task-specific instruction, without any examples, serving as a baseline to assess the model’s ability to generalize purely from the prompt. To enhance reasoning in more complex tasks, we incorporate CoT prompting, which explicitly elicits intermediate reasoning steps in the prompt. For a fair and reproducible evaluation, we use consistent instruction-based prompt templates across all models, as detailed in Appendix C.

Multi Task RAG. To enable few-shot multi-task ICL, we implement a multi-task RAG. Demonstrations from multiple tasks are retrieved and incorporated into the prompt. We explore three sequencing strategies for organizing demonstrations within the prompt context, which are listed as follows.

Random Ordering: Examples from all 3 tasks are interleaved in a fully randomized manner (e.g., $t_1, t_3, t_3, t_2, t_1, t_1, t_3, t_2, t_2$), serving as a baseline to assess robustness to prompt structure.

Task-Grouped Ordering: Examples are grouped by task, presenting all demonstrations from one task before moving to the next one (e.g., $t_1, t_1, t_1, t_2, t_2, t_2, t_3, t_3, t_3$). This structure emphasizes intra-task consistency.

Task-Interleaved Ordering: Examples alternate across tasks at each shot level, maintaining a round-robin pattern (e.g., $t_1, t_2, t_3, t_1, t_2, t_3, t_1, t_2, t_3$). This configuration aims to balance exposure across tasks within the prompt.

The impact of the number of shots per task and example orderings is shown in Appendix C (Figure 4 and Figure 5). Note that to encode examples into dense vector representations, we use the jina-embeddings-v3 (Sturua et al., 2024) model, and for distance computation, we employ the L2 distance metric.

4.3 Multi-Task PEFT

LoRA. As baseline, we implement LoRA (Hu et al., 2022). LoRA approximates full fine-tuning by decomposing weight matrices into low-rank

components. To reduce dimensionality, the weight matrix $\mathbf{W}_0 \in \mathbb{R}^{d \times k}$ is approximated by the product of two lower-rank matrices: $\mathbf{B} \in \mathbb{R}^{d \times r}$ and $\mathbf{A} \in \mathbb{R}^{r \times k}$, with $r \ll \min(d, k)$. This low-rank update preserves the backbone while enabling efficient adaptation, such that $h = \mathbf{W}_0 x + \frac{\alpha}{r} \mathbf{B} \mathbf{A} x$. LoRA can be applied to each linear layer in the Transformer architecture, such as $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V, \mathbf{W}_O$ matrices projections in the attention layers.

MTL-LoRA. Yang et al. (2024) introduce MTL-LoRA. Given task input x_t , MTL-LoRA first applies a shared standard LoRA down-projection via matrix \mathbf{A} . To retain task-specific information, it inserts a task-specific low-rank matrix $\Lambda_t \in \mathbb{R}^{r \times r}$ between the down- and up-projections, transforming $\mathbf{A}x_t$. Instead of a single shared up-projection, MTL-LoRA uses n matrices $\mathbf{B}^i \in \mathbb{R}^{d \times r}$ to support diverse knowledge-sharing strategies. Outputs are combined via a weighted average, where weights $w_t \in \mathbb{R}^{n \times 1}$ are learned per task as in Equation 1.

$$h_t = \mathbf{W}x_t + \sum_{i=1}^n \frac{\exp(w_t^i/\tau) \mathbf{B}^i}{\sum_{j=1}^n \exp(w_t^j/\tau)} \Lambda_t \mathbf{A}x_t \quad (1)$$

Here, τ controls the softness of the weighting. Each Λ_t is initialized as a diagonal identity matrix to ensure $\Delta \mathbf{W}_t = 0$ at start.

MTL Loss for ETR Generation. The model is trained to generate outputs conditioned on instructions. Given an instruction sequence $I = i_1, i_2, \dots, i_m$ and a corresponding completion sequence $C = c_1, c_2, \dots, c_n$, where I may contain special prompt tokens (e.g., $\langle \text{Input} \rangle$ and $\langle \text{Output} \rangle$), the full input is represented as $x = i_1, \dots, i_m, c_1, \dots, c_n$. The model is trained to autoregressively predict each token in C conditioned on all preceding tokens in I and C as defined in Equation 2.

$$P(C|I) = \prod_{j=1}^n P(c_j | i_1, \dots, i_m, c_1, \dots, c_{j-1}) \quad (2)$$

Based on the findings from Huerta-Enochian and Ko (2024), the objective is to minimize the negative log-likelihood of the completion sequence given the instruction as defined in Equation 3.

$$\mathcal{L} = - \sum_{j=1}^n \log P(c_j | i_1, \dots, i_m, c_1, \dots, c_{j-1}) \quad (3)$$

To account for imbalance across different instruction-following tasks, we apply a task-specific weighting scheme during training. Let N_t be the number of training examples for task t , and let $N = \sum_t N_t$ be the total number of training examples across all tasks. Each task’s contribution to the overall loss is scaled by a factor $w_t = \frac{N_t}{N}$, such that the final loss is redefined in Equation 4.

$$\mathcal{L}_{MTL} = \sum_{t=1}^T w_t \times \mathcal{L}_t \quad (4)$$

5 Results

The best models are selected based on the highest SRB score on the ETR-fr validation set, following a grid search hyperparameter tuning strategy.⁸ To complement this analysis, all models are run five times with different seeds, and detailed average results can be found in Appendix D.

5.1 In-Domain Quantitative Results

ICL Performance. As shown in Table 2, ICL models evidence steady improvements when transitioning from zero-shot and CoT prompting to RAG-based prompting. For LLaMA-3-8B, RAG achieves the best results with ETR-fr only inputs (e.g., 33.43/12.99/24.38 ROUGE-1/2/L and 42.16 SARI), outperforming zero-shot by a large margin. Adding related tasks does not consistently improve performance under ICL, and in some cases, leads to reduced novelty and compression ratio.

Impact of Fine-Tuning. PEFT significantly outperforms ICL methods. The best overall performance is achieved by LLaMA-3-8B with MTL-LoRA fine-tuned on ETR-fr and WikiLarge FR, obtaining highest scores across SARI (44.67), BERTScore-F1 (74.05), SRB (39.60), and compression ratio (56.11), while maintaining strong novelty (33.05).

LLM Comparison. Across both prompting and fine-tuning paradigms, LLaMA-3-8B outperforms Mistral-7B in most metrics. For instance, with LoRA fine-tuning on ETR-fr, LLaMA-3-8B achieves higher ROUGE-L (25.04 vs. 24.02), SARI (42.15 vs. 42.09), and SRB (38.77 vs. 37.98). This suggests that the architectural or scale advantages of LLaMA-3-8B translate effectively into more efficient capabilities.

⁸Hyperparameter tuning is detailed in Appendix B.

Combination of Tasks. Incorporating auxiliary tasks such as text summarization and simplification can provide complementary supervision, as seen in PEFT strategies. However, they do not yield performance gains in the ICL setting. Notably, MTL-LoRA with ETR-fr and WikiLarge FR for LLaMA-3-8B achieves the highest SARI and compression ratio, suggesting the relevance of sentence simplification data to the ETR generation task. However, inclusion of all three tasks does not universally yield the best results, and in some cases introduces performance regressions in BERTScore and novelty. This implies that careful curation of task mixtures is essential to avoid dilution or conflict between training objectives. Overall, these results highlight that while RAG improves performance in ICL, parameter-efficient fine-tuning (particularly MTL-LoRA) remains the most effective approach for high-quality in-domain ETR-fr.

5.2 Out-of-Domain Quantitative Results

ICL Performance. As shown in Table 3, among prompting strategies, RAG consistently outperforms zero-shot and CoT in all major content preservation metrics (ROUGE-1/2/L, BERTScore-F1) and the composite SRB score. On LLaMA-3-8B, using RAG with all three tasks (E,O,W) achieves the highest overall SRB score (41.52) and the best ROUGE-L (28.43), indicating its strong generalization and content fidelity. Moreover, it yields the highest SARI (42.63) and BERTScore-F1 (73.39), showcasing a balanced ability to simplify while preserving semantics. Interestingly, zero-shot exhibits extremely poor compression ratios, especially on Mistral-7B (-309.24), suggesting potential prompt misalignment or excessive hallucination. However, it achieves the highest novelty score (55.37) on LLaMA-3-8B, implying that despite poor content fidelity, more diverse lexical outputs are generated.

Impact of Fine-Tuning. While PEFT strategies generally lag behind RAG in terms of SRB and BERTScore, they offer stable and interpretable performance, with notably better compression ratios than zero-shot, CoT and most RAG-based strategies. The best PEFT model in terms of SRB, LLaMA-3-8B+LoRA trained solely on ETR-fr, achieves a relatively low compression ratio (6.38), indicating only moderate summarization. However, this comes at the expense of lower ROUGE, SARI, and BERTScore metrics compared to RAG-based

| | Method | Task | R-1 ↑ | R-2 ↑ | R-L ↑ | SARI ↑ | BERT-F1 ↑ | SRB ↑ | Comp. ratio | Novelty |
|--------------------------------|-----------|-------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| In-Context Learning | | | | | | | | | | |
| Mistral-7B | Zero-Shot | E | 23.92 | 7.09 | 16.28 | 37.07 | 69.75 | 29.20 | −64.14 | 35.70 |
| | CoT | E | 23.58 | 7.22 | 16.17 | 37.39 | 68.80 | 29.10 | −60.53 | <u>36.09</u> |
| | RAG | E | 32.14 | 10.47 | 22.72 | 40.05 | 72.41 | 36.24 | 44.32 | 26.55 |
| | | E,O | 31.12 | 9.58 | 21.92 | 39.54 | 71.29 | 35.32 | 48.45 | 26.61 |
| | | E,W | 30.29 | 9.69 | 21.29 | 38.69 | 71.59 | 34.56 | 33.80 | 23.01 |
| | | E,O,W | 29.84 | 9.57 | 21.58 | 39.53 | 71.06 | 35.01 | 46.42 | 25.85 |
| LlaMA-3-8B | Zero-Shot | E | 24.94 | 8.23 | 17.37 | 38.59 | 70.29 | 30.70 | −21.56 | 38.73 |
| | CoT | E | 27.57 | 8.96 | 18.72 | 38.26 | 71.02 | 32.04 | 7.80 | 31.10 |
| | RAG | E | 33.43 | 12.99 | 24.38 | 42.16 | 72.58 | 38.21 | 46.18 | 27.14 |
| | | E,O | 31.10 | 10.87 | 22.37 | 39.94 | 71.27 | 35.81 | 39.22 | 24.29 |
| | | E,W | 33.03 | 11.62 | 23.28 | 40.59 | 72.14 | 36.83 | 41.89 | 25.26 |
| | | E,O,W | 29.35 | 9.97 | 20.54 | 39.03 | 70.84 | 33.93 | 25.94 | 23.69 |
| Paramter-Efficient Fine-Tuning | | | | | | | | | | |
| Mistral-7B | LoRA | E | 32.47 | 12.40 | 24.02 | 42.09 | 73.56 | 37.98 | 44.42 | 18.35 |
| | MTL-LoRA | E,O | 32.67 | 12.74 | 24.33 | 41.95 | 73.52 | 38.20 | 53.48 | 24.17 |
| | | E,W | 32.62 | 12.92 | 24.28 | 42.53 | <u>73.90</u> | 38.35 | <u>53.62</u> | 24.99 |
| | | E,O,W | 33.65 | 12.83 | 24.93 | 42.25 | 73.62 | 38.77 | 48.93 | 23.38 |
| LlaMA-3-8B | LoRA | E | 31.76 | 13.17 | 25.04 | 42.15 | 72.93 | 38.77 | 50.66 | 18.87 |
| | MTL-LoRA | E,O | <u>33.44</u> | 13.22 | 24.24 | 43.04 | 73.86 | 38.45 | 51.36 | 23.06 |
| | | E,W | 32.54 | <u>13.56</u> | <u>25.08</u> | 44.67 | 74.05 | <u>39.60</u> | 56.11 | 33.05 |
| | | E,O,W | 32.78 | 13.64 | 25.67 | <u>43.53</u> | 73.28 | 39.69 | 53.24 | 24.39 |

Table 2: **Performance comparison, on ETR-fr test set**, across ICL methods and PEFT strategies on three tasks: ETR-fr (E), OrangeSum (O) and WikiLarge FR (W). Best results are in **bold**, second-best are underlined.

| | Method | Task | R-1 ↑ | R-2 ↑ | R-L ↑ | SARI ↑ | BERT-F1 ↑ | SRB ↑ | Comp. ratio | Novelty |
|--------------------------------|-----------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| In-Context Learning | | | | | | | | | | |
| Mistral-7B | Zero-Shot | E | 28.36 | 11.02 | 19.29 | 39.87 | 68.10 | 32.75 | −309.24 | 48.37 |
| | CoT | E | 29.78 | 11.22 | 19.90 | 39.62 | 69.40 | 33.37 | −261.30 | <u>50.85</u> |
| | RAG | E | 39.22 | 15.28 | 28.12 | 41.33 | 73.15 | <u>40.86</u> | 11.03 | 25.49 |
| | | E,O | 37.87 | 14.59 | 26.43 | 39.51 | 72.08 | <u>38.96</u> | 14.37 | 18.41 |
| | | E,W | 39.77 | 15.55 | 27.74 | 40.32 | 72.47 | 40.19 | 10.80 | 17.81 |
| E,O,W | | 39.12 | 15.97 | <u>28.26</u> | 40.74 | 72.87 | 40.73 | 14.63 | 18.33 | |
| LlaMA-3-8B | Zero-Shot | E | 29.60 | 10.84 | 18.83 | 40.55 | 68.68 | 32.50 | −180.74 | 55.37 |
| | CoT | E | 31.68 | 11.46 | 20.14 | 40.80 | 69.87 | 33.91 | −83.36 | 45.41 |
| | RAG | E | 37.48 | 13.98 | 26.94 | 41.05 | 73.18 | 39.92 | 11.37 | 41.63 |
| | | E,O | 40.53 | 15.15 | 27.47 | 41.14 | 72.75 | 40.29 | −12.56 | 31.01 |
| | | E,W | 39.72 | <u>16.02</u> | 26.83 | <u>41.99</u> | <u>73.32</u> | 40.15 | 13.75 | 35.70 |
| E,O,W | | <u>40.12</u> | 16.55 | 28.43 | 42.63 | 73.39 | 41.52 | −4.79 | 30.08 | |
| Paramter-Efficient Fine-Tuning | | | | | | | | | | |
| Mistral-7B | LoRA | E | 35.13 | 12.23 | 25.93 | 38.04 | 70.28 | 37.94 | 21.55 | 11.79 |
| | MTL-LoRA | E,O | 29.36 | 11.02 | 21.87 | 38.68 | 69.22 | 34.87 | 36.68 | 40.29 |
| | | E,W | 34.32 | 12.56 | 24.85 | 38.72 | 70.54 | 37.38 | <u>22.51</u> | 19.10 |
| | | E,O,W | 36.45 | 13.22 | 26.21 | 38.39 | 70.97 | 38.32 | 18.33 | 10.55 |
| LlaMA-3-8B | LoRA | E | 35.53 | 13.83 | 26.94 | 39.90 | 71.30 | 39.37 | 6.38 | 16.13 |
| | MTL-LoRA | E,O | 32.77 | 12.20 | 24.23 | 38.84 | 69.74 | 36.88 | 18.26 | 19.30 |
| | | E,W | 37.46 | 13.74 | 27.06 | 38.26 | 71.30 | 38.90 | 8.45 | 6.44 |
| | | E,O,W | 36.48 | 13.69 | 25.90 | 36.19 | 70.97 | 37.35 | 8.68 | 2.06 |

Table 3: **Performance comparison, on ETR-fr-politic test set**, across ICL methods and PEFT strategies on three tasks: ETR-fr (E), OrangeSum (O) and WikiLarge FR (W). Best results are in **bold**, second-best are underlined.

approaches. Additionally, MTL-LoRA configurations do not demonstrate performance improve-

ments over single-task LoRA in out-of-domain (OOD) settings, particularly on LLaMA-3-8B, suggesting a tendency toward overspecialization on the target task of ETR derived from children’s books.

Combination of Tasks. Prompting or training with multiple datasets (E,O,W) can improve OOD generalization. LLaMA-3-8B+RAG and Mistral-7B+RAG show substantial gains across all metrics compared to single-task prompting, confirming the benefits of multi-domain exposure in OOD settings. This situation is mitigated for the PEFT strategy, where performance improvement is backbone-dependent. While Mistral-7B+MTL-LoRA steadily benefits from concurrent learning achieving best results in terms of SRB with its (E,O,W) configuration, overall best results with LLaMA-3-8B are obtained with single task setting.

5.3 Human Evaluation

Manual evaluation is essential for assessing ETR text quality and compliance with European guidelines, which include 57 weighted questions covering clarity, simplicity, and accessibility,⁹ to ensure content is understandable and appropriate for the target audience. We validated our approach through human evaluation with 10 native French speakers, 7 NLP researchers and 3 linguists, all volunteers, who assessed outputs from the ETR-fr and ETR-politic test sets.¹⁰ We evaluated outputs generated by two model configurations: (1) Llama-3-8B+RAG augmented with ETR-fr (E) and WikiLarge FR (W), and (2) Llama-3-8B+MTL-LoRA trained on ETR-fr, OrangeSum (O), and WikiLarge FR, alongside their respective single-task variants. These models were chosen as the best performing ones, respectively for ICL and PEFT, for in-domain settings. The evaluation was performed on 6 source documents (3 from ETR-fr and 3 from ETR-fr-politic test sets). Each annotator reviewed 24 outputs, resulting in 60 samples per model and a total of 240 different samples evaluated. The assessment prioritized the most critical ETR guideline criteria, including information selection, sentence construction, word choice, and illustrations, covering 28 detailed questions (see Table 10 in Appendix). Additionally, we assessed general text generation quality metrics such as Fluency, Gram-

mar/Spelling, Relevance, Textual Coherence, and Overall Perceived Quality, through additional five questions. ETR criteria were rated on a binary scale (respected, not respected, not applicable), whereas human judgments used a 5-point Likert scale (1–5).

In-domain Results. Figures 2 presents the human evaluation results.¹¹ On ETR-fr, all methods perform well with respect to the European ETR guidelines. LoRA achieves the highest overall validation rate of 0.91, particularly excelling in word choice and sentence construction. MTL-LoRA+(E,O,W) shows the best results for sentence construction, while RAG+(E,W) outperforms other models in information selection. In terms of text generation quality, RAG leads with an overall score of 4.24, driven by strong performance in fluency, grammar, and coherence. While MTL-LoRA+(E,O,W) and LoRA are competitive across individual criteria, with MTL-LoRA+(E,O,W) scoring best on 3 out of 4 dimensions, their overall quality scores are comparable (3.95). Although automatic metrics indicate improved performance in multi-task settings, human evaluation results are more mixed, revealing no clear advantage for single- versus multi-task strategies, except in the Illustrations dimension.

Out-of-domain Results Overall performance declines on the more challenging ETR-fr-politic, yet RAG+(E,W) remains the most robust across both ETR criteria and text quality evaluations, underscoring the value of the multi-task setting. Specifically, RAG+(E,W), trained on a broader mix of tasks combining ETR and sentence simplification, achieves a total validation rate of 0.80 for ETR guidelines and an overall quality score of 3.76. In contrast, MTL-LoRA+(E,O,W) exhibits the sharpest drop in quality (2.62), indicating difficulties in managing politically nuanced content, although it still outperforms the single-task configuration in 3 out of 5 evaluation dimensions. Furthermore, in terms of European ETR compliance, MTL-LoRA+(E,O,W) struggles to generalize in out-of-domain settings, showing improvement only in the Illustrations criterion.

6 Conclusion

In this paper, we introduced ETR-fr, the first dataset fully compliant with the European ETR guidelines targeting neurodivergent populations, and explored

⁹https://www.unapei.org/wp-content/uploads/2020/01/liste_verification-falc-score_v2020-01-14-1.xlsx

¹⁰All evaluators received training and were blind to model development to prevent bias.

¹¹Overall scores are provided in a table in Appendix D.2.

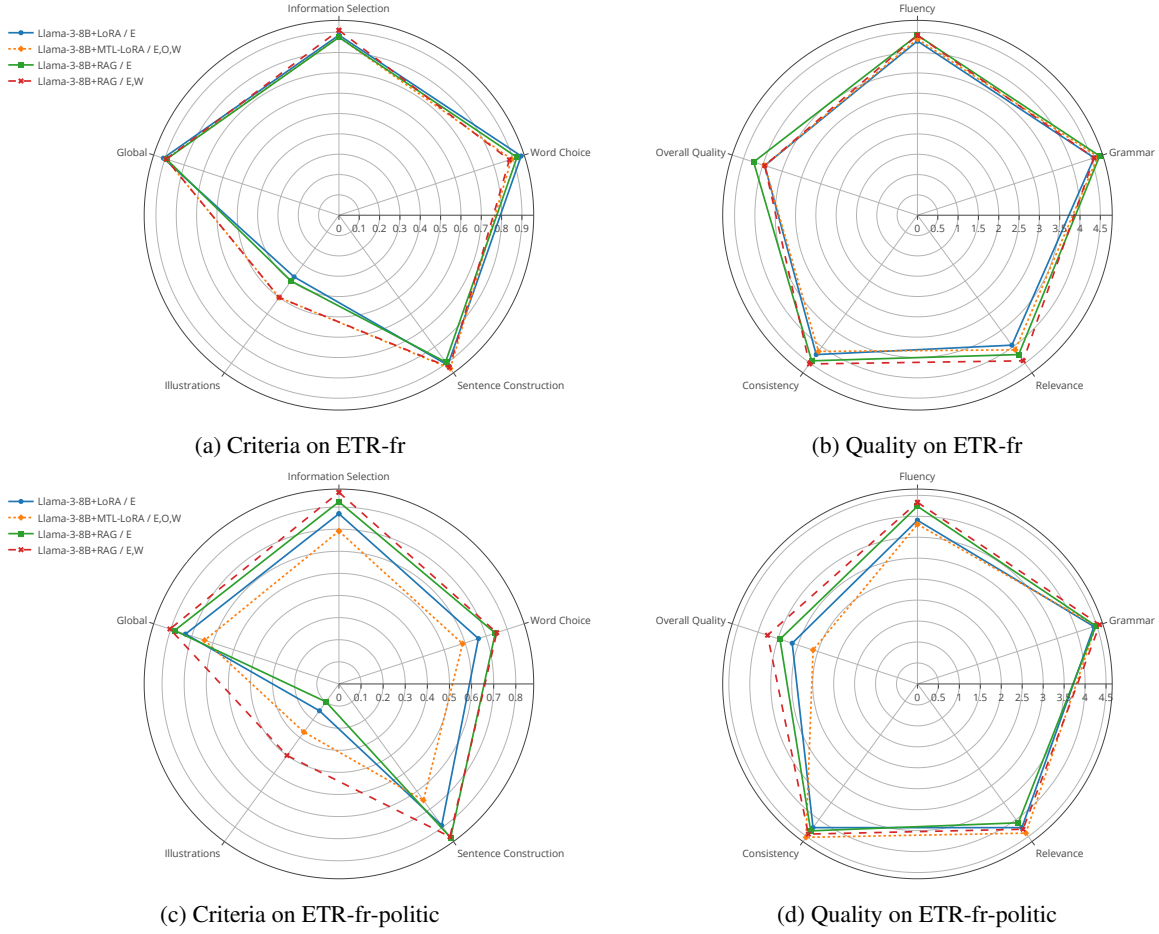


Figure 2: **Human evaluation of generation quality on ETR-fr and ETR-fr-politic** using their optimal ICL and MTL configurations. Subfigures (a) and (c) show average scores based on the ETR guideline criteria. Subfigures (b) and (d) present average human ratings for text generation quality.

multi-task learning to improve ETR generation with LLMs. Our experiments show that multi-task setups, particularly RAG for ICL and MTL-LoRA for PEFT, consistently improve performance in both in-domain and OOD settings according to automatic metrics. While human evaluation reveals more nuanced outcomes, it nonetheless confirms the benefits of multi-task learning across a broad range of ETR criteria and text quality dimensions.

7 Limitations

The development of ETR generation models introduces important constraints and considerations that reflect the complexity of cognitive accessibility and language model behavior.

Misalignment with deployment contexts. While our evaluation combines automatic and human assessments, it does not simulate usage in real-world settings such as assistive reading tools or educational platforms. Thus, the practical utility of outputs for neurodivergent users remains untested.

Absence of direct end-user feedback. Human evaluation was conducted by proxy annotators, which limits insights into subjective usability, emotional response, and real-world accessibility, central concerns in ETR adoption.

No explicit modeling of cognitive load. Though our models optimize for readability and fluency,

they do not account for cognitive effort. Even simplified outputs may challenge users when processing abstract or ambiguous content.

ETR guidelines as a fixed supervision target.

We use European ETR guidelines as a normative framework. While they offer structure, rigid adherence may exclude culturally specific or individualized accessibility strategies, limiting generalization.

Simplification-centric task framing. Our formulation treats ETR as summarization and simplification. However, this may overlook strategies unique to ETR, such as intentional redundancy, explicit inference resolution, and narrative scaffolding, often crucial for accessibility.

Susceptibility to hallucinations. As with most generative models, hallucinations and factual drift remain concerns, especially with RAG-based systems. This is particularly risky for audiences who may interpret outputs literally or depend on high textual reliability.

8 Impact and Ethical Considerations

Social and Ethical Challenge. Identifying limitations is essential for transparency and inclusive design. ETR generation impacts neurodivergent readers and intersects with accessibility, language rights, and communicative equity. As such, simplification systems must be evaluated not only on linguistic performance but on their potential to oversimplify or marginalize. By clarifying the limitations of our work, we aim to support responsible development and deployment. Acknowledging these boundaries also helps position ETR generation as a socio-technical task, one that demands sensitivity to both linguistic quality and lived experience.

Risks of Oversimplification. Simplified language is not neutral, it involves choices about what meaning is retained or lost. In some cases, simplification may erase nuance, flatten perspective, or reinforce harmful stereotypes. This tension is particularly acute for readers who engage with language differently.

Toward Responsible Design. Mitigating risks requires human-in-the-loop systems, participatory evaluation involving end users, and adaptation strategies that go beyond surface-level clarity. ETR guidelines should be viewed as a starting point, not a universal solution.

Positioning ETR as a Research Problem. ETR remains underexplored in NLP. By introducing aligned data, task-specific metrics, and a critical lens on modeling assumptions, we aim to establish it as a standalone task, one that demands linguistic sensitivity, practical design, and participatory validation.

References

- Armen Aghajanyan, Sonal Gupta, and Luke Zettlemoyer. 2021. [Intrinsic dimensionality explains the effectiveness of language model fine-tuning](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 7319–7328, Online. Association for Computational Linguistics.
- Miriam Anschütz, Joshua Oehms, Thomas Wimmer, Bartłomiej Jezierski, and Georg Groh. 2023. [Language models for German text simplification: Overcoming parallel data scarcity through style-specific pre-training](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 1147–1158, Toronto, Canada. Association for Computational Linguistics.
- Harmon Bhasin, Timothy Ossowski, Yiqiao Zhong, and Junjie Hu. 2024. [How does multi-task training affect transformer in-context capabilities? investigations with function classes](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, pages 169–187, Mexico City, Mexico. Association for Computational Linguistics.
- C. H. Björnsson. 1983. [Readability of Newspapers in 11 Languages](#). *Reading Research Quarterly*, 18(4):480–497.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, and Arvind Neelakantan. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Rémi Cardon and Natalia Grabar. 2020. [French Biomedical Text Simplification: When Small and Precise Helps](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 710–716, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Anwoy Chatterjee, Eshaan Tanwar, Subhabrata Dutta, and Tanmoy Chakraborty. 2024. [Language models can exploit cross-task in-context learning for data-scarce novel tasks](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11568–11587, Bangkok, Thailand. Association for Computational Linguistics.
- Nael Chehab, Hadmut Holken, and Mathilde Malgrange. 2019. [Simples - etude recueil des besoins falc](#). Technical report, SYSTRAN and EPNAK and EPHE and CHArt-LUTIN.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, and Paul Barham. 2023. [Palm: Scaling language modeling with pathways](#). *J. Mach. Learn. Res.*, 24:240:1–240:113.
- Silvana Deilen, Sergio Hernández Garrido, Ekaterina Lapshinova-Koltunski, and Christiane Maaß. 2023. [Using ChatGPT as a CAT tool in easy language translation](#). In *Proceedings of the Second Workshop on Text Simplification, Accessibility and Readability*, pages 1–10, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. [Qlora: Efficient finetuning of quantized llms](#). *Advances in neural information processing systems*, 36:10088–10115.
- Anna Dmitrieva and Jörg Tiedemann. 2024. [Towards Automatic Finnish Text Simplification](#). In *Proceedings of the Workshop on DeTermIt! Evaluating Text Difficulty in a Multilingual Context @ LREC-COLING 2024*, pages 39–50, Torino, Italia. ELRA and ICCL.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Baobao Chang, Xu Sun, Lei Li, and Zhifang Sui. 2024. [A survey on in-context learning](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 1107–1128, Miami, Florida, USA. Association for Computational Linguistics.
- Isabel Espinosa-Zaragoza, José Abreu-Salas, Paloma Moreda, and Manuel Palomar. 2023. [Automatic Text Simplification for People with Cognitive Disabilities: Resource Creation within the ClearText Project](#). In *Proceedings of the Second Workshop on Text Simplification, Accessibility and Readability*, pages 68–77, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.
- Nils Freyer, Hendrik Kempt, and Lars Klöser. 2024. [Easy-read and large language models: on the ethical dimensions of llm-based text simplification](#). *Ethics and Information Technology*, 26(3).
- Núria Gala, Anaïs Tack, Ludivine Javourey-Drevet, Thomas François, and Johannes C. Ziegler. 2020. [Alector: A Parallel Corpus of Simplified French Texts with Alignments of Misreadings by Poor and Dyslexic Readers](#). In *12th Language Resources and Evaluation Conference*, pages 1353–1361, Marseille, France. European Language Resources Association.
- Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and Kate Crawford. 2021. [Datasheets for datasets](#). *Commun. ACM*, 64(12):86–92.
- Sian Gooding. 2022. [On the ethical considerations of text simplification](#). In *Ninth Workshop on Speech and Language Processing for Assistive Technologies (SLPAT-2022)*, pages 50–57, Dublin, Ireland. Association for Computational Linguistics.

| | | | |
|-----|---|--|-----|
| 786 | Steven M. Goodman, Erin Buehler, Patrick Clary, Andy | Mathew Huerta-Enochian and Seung Yong Ko. 2024. | 844 |
| 787 | Coenen, Aaron Donsbach, Tiffanie N. Horne, Michal | Instruction fine-tuning: Does prompt loss matter? | 845 |
| 788 | Lahav, Robert MacDonald, Rain Breaw Michaels, | In <i>Proceedings of the 2024 Conference on Empiri-</i> | 846 |
| 789 | Ajit Narayanan, Mahima Pushkarna, Joel Riley, Alex | <i>cal Methods in Natural Language Processing</i> , pages | 847 |
| 790 | Santana, Lei Shi, Rachel Sweeney, Phil Weaver, Ann | 22771–22795, Miami, Florida, USA. Association for | 848 |
| 791 | Yuan, and Meredith Ringel Morris. 2022. LaMPost: | Computational Linguistics. | 849 |
| 792 | Design and Evaluation of an AI-assisted Email Writ- | | |
| 793 | ing Prototype for Adults with Dyslexia . In <i>Proceed-</i> | Syed Mahmudul Huq, Rytis Maskeliūnas, and Robertas | 850 |
| 794 | <i>ings of the 24th International ACM SIGACCESS Con-</i> | Damaševičius. 2024. Dialogue agents for artificial | 851 |
| 795 | <i>ference on Computers and Accessibility</i> , ASSETS | intelligence-based conversational systems for cogni- | 852 |
| 796 | '22, pages 1–18, New York, NY, USA. Association | tively disabled: a systematic review . <i>Disability and</i> | 853 |
| 797 | for Computing Machinery. | <i>Rehabilitation: Assistive Technology</i> , 19(3):1059– | 854 |
| | | 1078. | 855 |
| 798 | Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, | Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men- | 856 |
| 799 | Abhinav Pandey, Abhishek Kadian, Ahmad Al- | sch, Chris Bamford, Devendra Singh Chaplot, Diego | 857 |
| 800 | Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, | de las Casas, Florian Bressand, Gianna Lengyel, Guil- | 858 |
| 801 | Alex Vaughan, et al. 2024. The llama 3 herd of mod- | laume Lample, Lucile Saulnier, L  lio Renard Lavaud, | 859 |
| 802 | els . | Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, | 860 |
| | | Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, | 861 |
| 803 | Anders Gustavsson, Mikael Svensson, Frank Jacobi, | and William El Sayed. 2023. Mistral 7b . <i>Preprint</i> , | 862 |
| 804 | Christer Allgulander, Jordi Alonso, Ettore Beghi, | arXiv:2310.06825. | 863 |
| 805 | Richard Dodel, Mattias Ekman, Carlo Faravelli, | | |
| 806 | Laura Fratiglioni, Brenda Gannon, David Hilton | Moussa Kamal Eddine, Antoine Tixier, and Michalis | 864 |
| 807 | Jones, Poul Jennum, Alben Jordanova, Linus J  n- | Vazirgiannis. 2021. BARThez: a Skilled Pretrained | 865 |
| 808 | son, Korinna Karampampa, Martin Knapp, Gisela | French Sequence-to-Sequence Model . In <i>Proceed-</i> | 866 |
| 809 | Kobelt, Tobias Kurth, Roselind Lieb, Mattias Linde, | <i>ings of the 2021 Conference on Empirical Methods</i> | 867 |
| 810 | Christina Ljungcrantz, Andreas Maercker, Beatrice | <i>in Natural Language Processing</i> , pages 9369–9390, | 868 |
| 811 | Melin, Massimo Moscarelli, Amir Musayev, Fiona | Online and Punta Cana, Dominican Republic. Asso- | 869 |
| 812 | Norwood, Martin Preisig, Maura Pugliatti, Juergen | ciation for Computational Linguistics. | 870 |
| 813 | Rehm, Luis Salvador-Carulla, Brigitte Schlehofer, | | |
| 814 | Roland Simon, Hans-Christoph Steinhausen, Lars Ja- | Liliane Kandel and Abraham Moles. 1958. Application | 871 |
| 815 | cob Stovner, Jean-Michel Vallat, Peter Van den | de l'indice de Flesch    la langue fran  aise. <i>Cahiers</i> | 872 |
| 816 | Bergh, Jim van Os, Pieter Vos, Weili Xu, Hans-Ulrich | <i>Etudes de Radio-T  l  vision</i> , 19. | 873 |
| 817 | Wittchen, Bengt J  nsson, and Jes Olesen. 2011. Cost | | |
| 818 | of disorders of the brain in europe 2010 . <i>European</i> | J. P. Kincaid and And Others. 1975. Derivation of New | 874 |
| 819 | <i>Neuropsychopharmacology</i> , 21(10):718–779. | Readability Formulas (Automated Readability Index, | 875 |
| | | Fog Count and Flesch Reading Ease Formula) for | 876 |
| 820 | Lasse Hansen, Ludvig Renbo Olsen, and Kenneth | Navy Enlisted Personnel. Technical report, National | 877 |
| 821 | Enevoldsen. 2023. Textdescriptives: A python pack- | Technical Information Service, Springfield, Virginia | 878 |
| 822 | age for calculating a large variety of metrics from | 22151 (AD-A006 655/5GA, MF \$2. ERIC Number: | 879 |
| 823 | text . <i>Journal of Open Source Software</i> , 8(84):5153. | ED108134. | 880 |
| | | | |
| 824 | Silvia Hansen-Schirra, Walter Bisang, Arne Nagels, | Takeshi Kojima, Shixiang (Shane) Gu, Machel Reid, Yu- | 881 |
| 825 | Silke Gutermuth, Julia Fuchs, Liv Borghardt, Silvana | taka Matsuo, and Yusuke Iwasawa. 2022. Large lan- | 882 |
| 826 | Deilen, Anne-Kathrin Gros, Laura Schiffl, and Jo- | guage models are zero-shot reasoners . In <i>Advances in</i> | 883 |
| 827 | hanna Sommer. 2020. Intralingual Translation into | <i>Neural Information Processing Systems</i> , volume 35, | 884 |
| 828 | Easy Language - or how to reduce cognitive process- | pages 22199–22213. Curran Associates, Inc. | 885 |
| 829 | ing costs , page 197–225. Easy - Plain - Accessible. | | |
| 830 | Frank & Timme. | Ariel N. Lee, Cole J. Hunter, and Nataniel Ruiz. 2023. | 886 |
| | | Platypus: Quick, Cheap, and Powerful Refinement | 887 |
| 831 | Renate Hauser, Jannis Vamvas, Sarah Ebling, and Mar- | of LLMs . | 888 |
| 832 | tin Volk. 2022. A multilingual simplified language | | |
| 833 | news corpus . In <i>2nd Workshop on Tools and Re-</i> | Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio | 889 |
| 834 | <i>sources to Empower People with READING Difficul-</i> | Petroni, Vladimir Karpukhin, Naman Goyal, Hein- | 890 |
| 835 | <i>ties (READI) within the 13th Language Resources</i> | rich K  ttler, Mike Lewis, Wen-tau Yih, Tim Rock- | 891 |
| 836 | <i>and Evaluation Conference (LREC)</i> , pages 25–30, | t  schel, Sebastian Riedel, and Douwe Kiela. 2020. | 892 |
| 837 | Marseille, France. European Language Resources | Retrieval-augmented generation for knowledge- | 893 |
| 838 | Association. | intensive nlp tasks . In <i>Advances in Neural Infor-</i> | 894 |
| | | <i>mation Processing Systems</i> , volume 33, pages 9459– | 895 |
| 839 | Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen- | 9474. Curran Associates, Inc. | 896 |
| 840 | Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu | | |
| 841 | Chen. 2022. LoRA: Low-rank adaptation of large | Chin-Yew Lin. 2004. ROUGE: A Package for Auto- | 897 |
| 842 | language models . In <i>International conference on</i> | matic Evaluation of Summaries . In <i>Text Summariza-</i> | 898 |
| 843 | learning representations . | <i>tion Branches Out</i> , pages 74–81, Barcelona, Spain. | 899 |
| | | Association for Computational Linguistics. | 900 |

| | | | |
|-----|--|--|------|
| 901 | Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, | Teven Scao, Thomas Wolf, Osmo Suominen, Samuli | 958 |
| 902 | Lawrence Carin, and Weizhu Chen. 2022. What | Sairanen, Mikko Merioksa, Jyrki Heinonen, Aija | 959 |
| 903 | makes good in-context examples for gpt-3? In <i>Pro-</i> | Vahtola, Samuel Antao, and Sampo Pyysalo. 2023. | 960 |
| 904 | <i>ceedings of Deep Learning Inside Out: The 3rd Work-</i> | FinGPT: Large Generative Models for a Small Lan- | 961 |
| 905 | <i>shop on Knowledge Extraction and Integration for</i> | guage . In <i>Proceedings of the 2023 Conference on</i> | 962 |
| 906 | <i>Deep Learning Architectures, DeeLIO@ACL 2022,</i> | <i>Empirical Methods in Natural Language Processing,</i> | 963 |
| 907 | <i>Dublin, Ireland and Online, May 27, 2022, pages</i> | pages 2710–2726, Singapore. Association for Com- | 964 |
| 908 | 100–114. Association for Computational Linguistics. | putational Linguistics. | 965 |
| 909 | Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, | Lara J Martin and Malathy Nagalakshmi. 2024. Bridg- | 966 |
| 910 | Hiroaki Hayashi, and Graham Neubig. 2023. Pre- | ing the social & technical divide in augmentative | 967 |
| 911 | train, prompt, and predict: A systematic survey of | and alternative communication (aac) applications for | 968 |
| 912 | prompting methods in natural language processing. | autistic adults . <i>arXiv preprint arXiv:2404.17730</i> . | 969 |
| 913 | <i>ACM Comput. Surv.</i> , 55(9). | | |
| 914 | Peter J. Liu*, Mohammad Saleh*, Etienne Pot, Ben | Paloma Martínez, Alberto Ramos, and Lourdes Moreno. | 970 |
| 915 | Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam | 2024. Exploring large language models to generate | 971 |
| 916 | Shazeer. 2018. Generating wikipedia by summariz- | Easy to Read content . <i>Frontiers in Computer Science,</i> | 972 |
| 917 | ing long sequences . In <i>International Conference on</i> | 6. Publisher: Frontiers. | 973 |
| 918 | <i>Learning Representations</i> . | | |
| 919 | Qidong Liu, Xian Wu, Xiangyu Zhao, Yuanshao Zhu, | Pallab K. Maulik, Maya N. Mascarenhas, Colin D. | 974 |
| 920 | Derong Xu, Feng Tian, and Yefeng Zheng. 2024a. | Mathers, Tarun Dua, and Shekhar Saxena. 2011. | 975 |
| 921 | When MOE Meets LLMs: Parameter Efficient Fine- | Prevalence of intellectual disability: A meta-analysis | 976 |
| 922 | tuning for Multi-task Medical Applications . In <i>Pro-</i> | of population-based studies . <i>Research in Develop-</i> | 977 |
| 923 | <i>ceedings of the 47th International ACM SIGIR Con-</i> | <i>mental Disabilities</i> , 32(2):419–436. | 978 |
| 924 | <i>ference on Research and Development in Information</i> | | |
| 925 | <i>Retrieval, SIGIR '24, page 1104–1114, New York,</i> | Tomas Murillo-Morales, Peter Heumader, and Klaus | 979 |
| 926 | NY, USA. Association for Computing Machinery. | Miesenberger. 2020. Automatic Assistance to Cogni- | 980 |
| 927 | Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo | tive Disabled Web Users via Reinforcement Learning | 981 |
| 928 | Molchanov, Yu-Chiang Frank Wang, Kwang-Ting | on the Browser . In <i>Computers Helping People with</i> | 982 |
| 929 | Cheng, and Min-Hung Chen. 2024b. DoRA: Weight- | <i>Special Needs</i> , pages 61–72, Cham. Springer Interna- | 983 |
| 930 | decomposed low-rank adaptation . In <i>Proceedings of</i> | <i>tional Publishing</i> . | 984 |
| 931 | <i>the 41st International Conference on Machine Learn-</i> | | |
| 932 | <i>ing</i> , volume 235 of <i>Proceedings of Machine Learning</i> | Shashi Narayan, Shay B. Cohen, and Mirella Lapata. | 985 |
| 933 | <i>Research</i> , pages 32100–32121. PMLR. | 2018. Don't Give Me the Details, Just the Sum- | 986 |
| 934 | Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey | mary! Topic-Aware Convolutional Neural Networks | 987 |
| 935 | Edunov, Marjan Ghazvininejad, Mike Lewis, and | for Extreme Summarization . In <i>Proceedings of the</i> | 988 |
| 936 | Luke Zettlemoyer. 2020. Multilingual Denoising | <i>2018 Conference on Empirical Methods in Natural</i> | 989 |
| 937 | Pre-training for Neural Machine Translation . <i>Trans-</i> | <i>Language Processing</i> , pages 1797–1807, Brussels, | 990 |
| 938 | <i>actions of the Association for Computational Linguis-</i> | Belgium. Association for Computational Linguistics. | 991 |
| 939 | <i>tics</i> , 8:726–742. Place: Cambridge, MA Publisher: | | |
| 940 | MIT Press. | OpenAI. 2023. GPT-4 technical report . <i>CoRR</i> , | 992 |
| 941 | Ilya Loshchilov and Frank Hutter. 2019. Decoupled | abs/2303.08774. | 993 |
| 942 | weight decay regularization . In <i>7th International</i> | | |
| 943 | <i>Conference on Learning Representations, ICLR 2019,</i> | Gustavo Paetzold and Lucia Specia. 2016. Unsuper- | 994 |
| 944 | <i>New Orleans, LA, USA, May 6-9, 2019</i> . OpenRe- | vised lexical simplification for non-native speakers . | 995 |
| 945 | view.net. | In <i>Proceedings of the AAAI Conference on Artificial</i> | 996 |
| 946 | Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, | <i>Intelligence</i> , volume 30. | 997 |
| 947 | and Pontus Stenetorp. 2022. Fantastically ordered | | |
| 948 | prompts and where to find them: Overcoming few- | Kishore Papineni, Salim Roukos, Todd Ward, and Wei- | 998 |
| 949 | shot prompt order sensitivity . In <i>Proceedings of the</i> | Jing Zhu. 2002. Bleu: a Method for Automatic Eval- | 999 |
| 950 | <i>60th Annual Meeting of the Association for Comput-</i> | uation of Machine Translation . In <i>Proceedings of</i> | 1000 |
| 951 | <i>ational Linguistics (Volume 1: Long Papers)</i> , pages | <i>the 40th Annual Meeting of the Association for Com-</i> | 1001 |
| 952 | 8086–8098, Dublin, Ireland. Association for Compu- | <i>putational Linguistics</i> , pages 311–318, Philadelphia, | 1002 |
| 953 | tational Linguistics. | Pennsylvania, USA. Association for Computational | 1003 |
| 954 | Risto Luukkainen, Ville Komulainen, Jouni Luoma, | Linguistics. | 1004 |
| 955 | Anni Eskelinen, Jenna Kanerva, Hanna-Mari Kupari, | | |
| 956 | Filip Ginter, Veronika Laippala, Niklas Muennighoff, | Pathways. 2021. Information for all: European stan- | 1005 |
| 957 | Aleksandra Piktus, Thomas Wang, Nouamane Tazi, | dards for making information easy to read and under- | 1006 |
| | | stand . | 1007 |
| | | Clifton Poth, Hannah Sterz, Indraneil Paul, Sukannya | 1008 |
| | | Purkayastha, Leon Engländer, Timo Imhof, Ivan | 1009 |
| | | Vulić, Sebastian Ruder, Iryna Gurevych, and Jonas | 1010 |
| | | Pfeiffer. 2023. Adapters: A Unified Library for | 1011 |
| | | Parameter-Efficient and Modular Transfer Learning . | 1012 |
| | | In <i>Proceedings of the 2023 Conference on Empirical</i> | 1013 |

| | | |
|------|--|------|
| 1014 | <i>Methods in Natural Language Processing: System Demonstrations</i> , pages 149–160, Singapore. Association for Computational Linguistics. | 1070 |
| 1015 | | 1071 |
| 1016 | | 1072 |
| 1017 | Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. In-context retrieval-augmented language models . <i>Transactions of the Association for Computational Linguistics</i> , 11:1316–1331. | 1073 |
| 1018 | | 1074 |
| 1019 | | 1075 |
| 1020 | | 1076 |
| 1021 | | 1077 |
| 1022 | Alexander M. Rush, Sumit Chopra, and Jason Weston. 2015. A Neural Attention Model for Abstractive Sentence Summarization . In <i>Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing</i> , pages 379–389, Lisbon, Portugal. Association for Computational Linguistics. | 1078 |
| 1023 | | 1079 |
| 1024 | | 1080 |
| 1025 | | 1081 |
| 1026 | | 1082 |
| 1027 | | 1083 |
| 1028 | Abigail See, Peter J. Liu, and Christopher D. Manning. 2017. Get to the point: Summarization with pointer-generator networks . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1073–1083, Vancouver, Canada. Association for Computational Linguistics. | 1084 |
| 1029 | | 1085 |
| 1030 | | 1086 |
| 1031 | | 1087 |
| 1032 | | 1088 |
| 1033 | | 1089 |
| 1034 | | 1090 |
| 1035 | Zhenmei Shi, Junyi Wei, Zhuoyan Xu, and Yingyu Liang. 2024. Why larger language models do in-context learning differently? In <i>Proceedings of the 41st International Conference on Machine Learning</i> , volume 235 of <i>Proceedings of Machine Learning Research</i> , pages 44991–45013. PMLR. | 1091 |
| 1036 | | 1092 |
| 1037 | | 1093 |
| 1038 | | 1094 |
| 1039 | | 1095 |
| 1040 | | 1096 |
| 1041 | Sanja Stajner. 2021. Automatic text simplification for social good: Progress and challenges . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 2637–2652, Online. Association for Computational Linguistics. | 1097 |
| 1042 | | 1098 |
| 1043 | | 1099 |
| 1044 | | 1100 |
| 1045 | | 1101 |
| 1046 | Regina Stodden, Omar Momen, and Laura Kallmeyer. 2023. DEplain: A German parallel corpus with intralingual translations into plain language for sentence and document simplification . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 16441–16463, Toronto, Canada. Association for Computational Linguistics. | 1102 |
| 1047 | | 1103 |
| 1048 | | 1104 |
| 1049 | | 1105 |
| 1050 | | 1106 |
| 1051 | | 1107 |
| 1052 | | 1108 |
| 1053 | | 1109 |
| 1054 | Saba Sturua, Isabelle Mohr, Mohammad Kalim Akram, Michael Günther, Bo Wang, Markus Krimmel, Feng Wang, Georgios Mastrapas, Andreas Koukounas, Nan Wang, and Han Xiao. 2024. jina-embeddings-v3: Multilingual embeddings with task lora . <i>Preprint</i> , arXiv:2409.10173. | 1110 |
| 1055 | | 1111 |
| 1056 | | 1112 |
| 1057 | | 1113 |
| 1058 | | 1114 |
| 1059 | | 1115 |
| 1060 | Elior Sulem, Omri Abend, and Ari Rappoport. 2018. Simple and Effective Text Simplification Using Semantic and Neural Methods . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 162–173, Melbourne, Australia. Association for Computational Linguistics. | 1116 |
| 1061 | | 1117 |
| 1062 | | 1118 |
| 1063 | | 1119 |
| 1064 | | 1120 |
| 1065 | | 1121 |
| 1066 | | 1122 |
| 1067 | Renliang Sun, Zhixian Yang, and Xiaojun Wan. 2023. Exploiting summarization data to help text simplification . In <i>Proceedings of the 17th Conference of the</i> | 1123 |
| 1068 | | 1124 |
| 1069 | | 1125 |
| | | 1126 |
| | | 1127 |
| | | 1128 |
| | <i>European Chapter of the Association for Computational Linguistics</i> , pages 39–51, Dubrovnik, Croatia. Association for Computational Linguistics. | |
| | | |
| | Yuting Tang, Ratish Puduppully, Zhengyuan Liu, and Nancy Chen. 2023. In-context learning of large language models for controlled dialogue summarization: A holistic benchmark and empirical analysis . In <i>Proceedings of the 4th New Frontiers in Summarization Workshop</i> , pages 56–67, Singapore. Association for Computational Linguistics. | |
| | | |
| | Teerapaun Tanprasert and David Kauchak. 2021. Flesch-kincaid is not a text simplification evaluation metric . In <i>Proceedings of the 1st Workshop on Natural Language Generation, Evaluation, and Metrics (GEM 2021)</i> , pages 1–14, Online. Association for Computational Linguistics. | |
| | | |
| | Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, YaGuang Li, Hongrae Lee, Huaixiu Steven Zheng, Amin Ghafouri, Marcelo Menegali, Yanping Huang, Maxim Krikun, Dmitry Lepikhin, James Qin, Dehao Chen, Yuanzhong Xu, Zhifeng Chen, Adam Roberts, Maarten Bosma, Vincent Zhao, Yanqi Zhou, Chung-Ching Chang, Igor Krivokon, Will Rusch, Marc Pickett, Pranesh Srinivasan, Laichee Man, Kathleen Meier-Hellstern, Meredith Ringel Morris, Tulsee Doshi, Renelito Delos Santos, Toju Duke, Johnny Soraker, Ben Zevenbergen, Vinodkumar Prabhakaran, Mark Diaz, Ben Hutchinson, Kristen Olson, Alejandra Molina, Erin Hoffman-John, Josh Lee, Lora Aroyo, Ravi Rajakumar, Alena Butryna, Matthew Lamm, Viktoriya Kuzmina, Joe Fenton, Aaron Cohen, Rachel Bernstein, Ray Kurzweil, Blaise Aguerre-Arcas, Claire Cui, Marian Croak, Ed Chi, and Quoc Le. 2022. Lamda: Language models for dialog applications . <i>arXiv preprint arXiv:2201.08239</i> . | |
| | | |
| | Amalia Todirascu, Rodrigo Wilkens, Eva Rolin, Thomas François, Delphine Bernhard, and Núria Gala. 2022. HECTOR: A Hybrid TEXT Simplification Tool for Raw Texts in French . In <i>Proceedings of the Thirteenth Language Resources and Evaluation Conference</i> , pages 4620–4630, Marseille, France. European Language Resources Association. | |
| | | |
| | Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. Llama: Open and efficient foundation language models . <i>Preprint</i> , arXiv:2302.13971. | |
| | | |
| | Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan | |

| | | | |
|------|---|---|------|
| 1129 | Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, | Wei Xu, Courtney Napoles, Ellie Pavlick, Quanze Chen, | 1188 |
| 1130 | Isabel Kloumann, Artem Korenev, Punit Singh Koura, | and Chris Callison-Burch. 2016. Optimizing Sta- | 1189 |
| 1131 | Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di- | tistical Machine Translation for Text Simplification . | 1190 |
| 1132 | ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar- | <i>Transactions of the Association for Computational</i> | 1191 |
| 1133 | tiniet, Todor Mihaylov, Pushkar Mishra, Igor Moly- | <i>Linguistics</i> , 4:401–415. Place: Cambridge, MA Pub- | 1192 |
| 1134 | bog, Yixin Nie, Andrew Poulton, Jeremy Reizen- | lisher: MIT Press. | 1193 |
| 1135 | stein, Rashi Rungta, Kalyan Saladi, Alan Schel- | | |
| 1136 | ten, Ruan Silva, Eric Michael Smith, Ranjan Sub- | Victoria Yaneva. 2015. Easy-read documents as a gold | 1194 |
| 1137 | ramanian, Xiaoqing Ellen Tan, Binh Tang, Ross | standard for evaluation of text simplification out- | 1195 |
| 1138 | Taylor, Adina Williams, Jian Xiang Kuan, Puxin | put . In <i>Proceedings of the Student Research Work-</i> | 1196 |
| 1139 | Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, An- | <i>shop</i> , pages 30–36, Hissar, Bulgaria. INCOMA Ltd. | 1197 |
| 1140 | gela Fan, Melanie Kambadur, Sharan Narang, Aure- | Shoumen, BULGARIA. | 1198 |
| 1141 | lien Rodriguez, Robert Stojnic, Sergey Edunov, and | | |
| 1142 | Thomas Scialom. 2023b. Llama 2: Open Founda- | Yaming Yang, Dilxat Muhtar, Yelong Shen, Yuefeng | 1199 |
| 1143 | tion and Fine-Tuned Chat Models . <i>arXiv preprint</i> . | Zhan, Jianfeng Liu, Yujing Wang, Hao Sun, Denvy | 1200 |
| 1144 | ArXiv:2307.09288 [cs]. | Deng, Feng Sun, Qi Zhang, Weizhu Chen, and Yun- | 1201 |
| | | hai Tong. 2024. Mtl-lora: Low-rank adaptation for | 1202 |
| 1145 | Subhadra Vadlamannati and Gözde Şahin. 2023. Metric- | multi-task learning . <i>CoRR</i> , abs/2410.09437. | 1203 |
| 1146 | based in-context learning: A case study in text sim- | | |
| 1147 | plification . In <i>Proceedings of the 16th International</i> | Fan Yin, Jesse Vig, Philippe Laban, Shafiq Joty, Caim- | 1204 |
| 1148 | <i>Natural Language Generation Conference</i> , pages | ing Xiong, and Chien-Sheng Wu. 2023. Did you read | 1205 |
| 1149 | 253–268, Prague, Czechia. Association for Computa- | the instructions? rethinking the effectiveness of task | 1206 |
| 1150 | tional Linguistics. | definitions in instruction learning . In <i>Proceedings</i> | 1207 |
| | | <i>of the 61st Annual Meeting of the Association for</i> | 1208 |
| 1151 | Ashish Vaswani, Samy Bengio, Eugene Brevdo, Fran- | <i>Computational Linguistics (Volume 1: Long Papers)</i> , | 1209 |
| 1152 | cois Chollet, Aidan N. Gomez, Stephan Gouws, Llion | pages 3063–3079, Toronto, Canada. Association for | 1210 |
| 1153 | Jones, Łukasz Kaiser, Nal Kalchbrenner, Niki Parmar, | Computational Linguistics. | 1211 |
| 1154 | Ryan Sepassi, Noam Shazeer, and Jakob Uszkoreit. | | |
| 1155 | 2018. Tensor2tensor for neural machine translation . | Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. | 1212 |
| 1156 | <i>CoRR</i> , abs/1803.07416. | Weinberger, and Yoav Artzi. 2020. Bertscore: Evalu- | 1213 |
| | | ating text generation with BERT . In <i>8th International</i> | 1214 |
| 1157 | Alex Wang, Amanpreet Singh, Julian Michael, Felix | <i>Conference on Learning Representations, ICLR 2020,</i> | 1215 |
| 1158 | Hill, Omer Levy, and Samuel Bowman. 2018. GLUE: | <i>Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenRe- | 1216 |
| 1159 | A multi-task benchmark and analysis platform for nat- | <i>view.net</i> . | 1217 |
| 1160 | ural language understanding . In <i>Proceedings of the</i> | | |
| 1161 | <i>2018 EMNLP Workshop BlackboxNLP: Analyzing</i> | | |
| 1162 | <i>and Interpreting Neural Networks for NLP</i> , pages | | |
| 1163 | 353–355, Brussels, Belgium. Association for Com- | | |
| 1164 | putational Linguistics. | | |
| 1165 | Xi Wang, Procheta Sen, Ruizhe Li, and Emine Yil- | | |
| 1166 | maz. 2024. Simulated Task Oriented Dialogues | | |
| 1167 | for Developing Versatile Conversational Agents . In | | |
| 1168 | <i>Advances in Information Retrieval</i> , pages 157–172, | | |
| 1169 | Cham. Springer Nature Switzerland. | | |
| 1170 | Yiming Wang, Yu Lin, Xiaodong Zeng, and Guan- | | |
| 1171 | nan Zhang. 2023. Multilora: Democratizing lora | | |
| 1172 | for better multi-task learning . <i>arXiv preprint</i> | | |
| 1173 | <i>arXiv:2311.11501</i> . | | |
| 1174 | Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten | | |
| 1175 | Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, | | |
| 1176 | and Denny Zhou. 2022. Chain-of-thought prompt- | | |
| 1177 | ing elicits reasoning in large language models . In | | |
| 1178 | <i>Advances in Neural Information Processing Systems</i> , | | |
| 1179 | volume 35, pages 24824–24837. Curran Associates, | | |
| 1180 | Inc. | | |
| 1181 | Sander Wubben, Antal van den Bosch, and Emiel Krah- | | |
| 1182 | mer. 2012. Sentence Simplification by Monolingual | | |
| 1183 | Machine Translation . In <i>Proceedings of the 50th An-</i> | | |
| 1184 | <i>annual Meeting of the Association for Computational</i> | | |
| 1185 | <i>Linguistics (Volume 1: Long Papers)</i> , pages 1015– | | |
| 1186 | 1024, Jeju Island, Korea. Association for Computa- | | |
| 1187 | tional Linguistics. | | |

A ETR-fr Dataset Sheet

The dataset description follows the recommendations and template proposed by [Gebru et al. \(2021\)](#).

Motivation

For what purpose was the dataset created? Was there a specific task in mind? Was there a specific gap that needed to be filled? Please provide a description.

The ETR-fr dataset was created to address the lack of high-quality, document-aligned corpora suitable for generating Easy-to-Read (ETR) text. It supports the task of generating cognitively accessible texts for individuals with cognitive impairments by providing paragraph-aligned text pairs that follow the European ETR guidelines. This dataset enables the training and evaluation of automatic systems for ETR generation in French, targeting the linguistic and cognitive accessibility requirements typically overlooked by existing simplification or summarization.

Who created this dataset (e.g., which team, research group) and on behalf of which entity (e.g., company, institution, organization)?

The dataset was constructed by the authors of the this paper on ETR-fr.

Composition

What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)? Are there multiple types of instances (e.g., movies, users, and ratings; people and interactions between them; nodes and edges)? Please provide a description.

Each instance in the ETR-fr dataset consists of a pair of paragraph-aligned French texts: a source text and its corresponding Easy-to-Read (ETR) version. These are designed to support document-level simplification, emphasizing both lexical and structural transformation.

How many instances are there in total (of each type, if appropriate)?

The dataset contains 523 paragraph-aligned text pairs. Additionally, an out-of-domain subset, ETR-fr-politic, includes 33 paragraph pairs from 2022 French presidential election programs.

What data does each instance consist of? “Raw” data (e.g., unprocessed text or images) or features? In either case, please provide a description.

Each instance consists of “raw” French text paragraphs: a complex source text and its corresponding simplified (ETR) version. These are aligned at the paragraph level and include natural language text only.

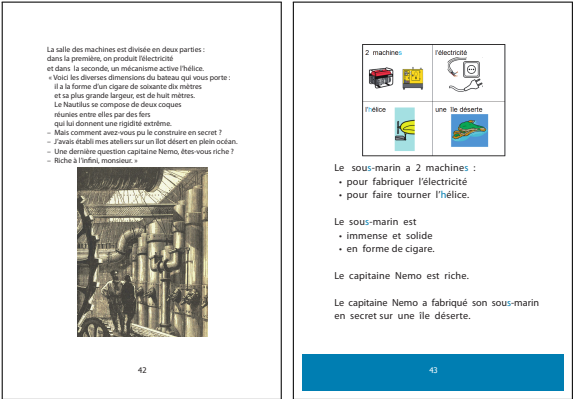


Figure 3: Extract of the ETR book *Twenty Thousand Leagues Under the Seas* by Jules Verne from François Baudez Publishing. **Left page** is the original text with an illustration. **Right page** is the ETR transcription with the main information plus its captioned vignettes.

Is there a label or target associated with each instance? If so, please provide a description.

Yes. The target is the simplified (ETR-compliant) version of the source paragraph, forming a supervised text-to-text pair for generation tasks.

Is any information missing from individual instances? If so, please provide a description, explaining why this information is missing (e.g., because it was unavailable). This does not include intentionally removed information, but might include, e.g., redacted text.

The pictograms present with the original texts have not been extracted.

Are relationships between individual instances made explicit (e.g., users’ movie ratings, social network links)? If so, please describe how these relationships are made explicit.

No such relationships exist or are made explicit in this dataset.

Are there recommended data splits (e.g., training, development/validation, testing)? If so, please provide a description of these splits, explaining the rationale behind them.

Yes. The dataset is divided into training (399 pairs), validation (71 pairs), and test (53 pairs) subsets. The test set comprises two distinct books chosen to ensure diversity in linguistic features such as text

length, structure, and readability. The remaining books were split into training and validation sets using a stratified approach to minimize thematic and lexical overlap. Additionally, the ETR-fr-politic test set (33 pairs) was introduced to assess model generalization on out-of-domain content not seen during training.

Are there any errors, sources of noise, or redundancies in the dataset? If so, please provide a description.

No specific mention of noise or redundancy issues is made in the source document.

Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g., websites, tweets, other datasets)? If it links to or relies on external resources, a) are there guarantees that they will exist, and remain constant, over time; b) are there official archival versions of the complete dataset (i.e., including the external resources as they existed at the time the dataset was created); c) are there any restrictions (e.g., licenses, fees) associated with any of the external resources that might apply to a future user?

The dataset is self-contained, it does not rely on external resources.

Does the dataset contain data that might be considered confidential (e.g., data that is protected by legal privilege or by doctor-patient confidentiality, data that includes the content of individuals non-public communications)? If so, please provide a description.

No. All texts are from published sources and are intended for public consumption.

Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety? If so, please describe why.

No such content is reported or expected in the dataset.

Does the dataset relate to people? If not, you may skip the remaining questions in this section.

No. The dataset is composed of literary and political texts and does not contain personal information.

Collection Process

How was the data associated with each instance acquired? Was the data directly observable (e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly inferred/derived from other data (e.g., part-of-speech tags, model-based guesses for age

or language)? If data was reported by subjects or indirectly inferred/derived from other data, was the data validated/verified? If so, please describe how.

The data are directly observable from published ETR books. Each ETR version is produced by a pair of trained transcribers working collaboratively, in accordance with the European Easy-to-Read guidelines (Pathways, 2021), to obtain official ETR certification.

What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or sensor, manual human curation, software program, software API)? How were these mechanisms or procedures validated?

To collect the data from ETR books, we first obtained the PDF versions and manually curated them to identify pairs of pages containing the original text and its corresponding ETR version. The textual content was then extracted using the Python library pypdfium2¹².

If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic, probabilistic with specific sampling probabilities)?

The dataset is not sampled from a larger set; it includes the complete collection of available aligned texts selected for the study.

Who was involved in the data collection process (e.g., students, crowdworkers, contractors) and how were they compensated (e.g., how much were crowdworkers paid)?

Unknown for the manual book transcriptions. The data collection was carried out by the main author of this paper as part of their research work.

Over what timeframe was the data collected? Does this timeframe match the creation timeframe of the data associated with the instances (e.g., recent crawl of old news articles)? If not, please describe the timeframe in which the data associated with the instances was created.

The exact creation dates of the original books are unknown. However, the dataset itself was constructed between May 2023 and June 2023.

Were any ethical review processes conducted (e.g., by an institutional review board)? If so, please provide a description of these review processes, including the outcomes, as well as a link or other access point to any supporting documentation.

No ethical review.

¹²<https://github.com/pypdfium2-team/pypdfium2>

| | | | |
|------|--|---|------|
| 1392 | Does the dataset relate to people? If not, you may | | |
| 1393 | skip the remaining questions in this section. | | |
| 1394 | No. | Are there tasks for which the dataset should not be used? If so, please provide a description. | 1438 |
| | | No. | 1439 |
| 1395 | | | 1440 |
| 1396 | | Any other comments? | 1441 |
| 1397 | | | |
| 1398 | | | |
| 1399 | Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing of missing values)? | Distribution | 1442 |
| 1400 | If so, please provide a description. If not, you may skip the remainder of the questions in this section. | | 1443 |
| 1401 | Manual cleaning was performed to remove chapter titles from the original texts, as these were not present in the corresponding ETR versions. | Will the dataset be distributed to third parties outside of the entity (e.g., company, institution, organization) on behalf of which the dataset was created? If so, please provide a description. | 1444 |
| 1402 | | | 1445 |
| 1403 | | | 1446 |
| 1404 | | | 1447 |
| 1405 | | | 1448 |
| 1406 | Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support unanticipated future uses)? If so, please provide a link or other access point to the “raw” data. | How will the dataset will be distributed (e.g., tarball on website, API, GitHub) Does the dataset have a digital object identifier (DOI)? | 1449 |
| 1407 | Yes. The raw data is provided alongside the cleaned version. | The dataset will be available on GitHub repository. | 1450 |
| 1408 | | | 1451 |
| 1409 | | | 1452 |
| 1410 | | | 1453 |
| 1411 | | When will the dataset be distributed? | 1454 |
| 1412 | Is the software used to preprocess/clean/label the instances available? If so, please provide a link or other access point. | The dataset will be released pending agreement from the ETR books publisher. | 1455 |
| 1413 | | | 1456 |
| 1414 | | Will the dataset be distributed under a copyright or other intellectual property (IP) license, and/or under applicable terms of use (ToU)? If so, please describe this license and/or ToU, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms or ToU, as well as any fees associated with these restrictions. | 1457 |
| 1415 | <ul style="list-style-type: none">• pypdfium2: https://github.com/pypdfium2-team/pypdfium2 | The dataset will be released under a custom license, subject to approval from the ETR books publisher. Redistribution and use will be permitted for research purposes only, with appropriate citation. No commercial use will be allowed without explicit permission. | 1460 |
| 1416 | | | 1461 |
| 1417 | <ul style="list-style-type: none">• cleantext: https://pypi.org/project/cleantext/ | | 1462 |
| 1418 | | | 1463 |
| 1419 | | | 1464 |
| 1420 | | | 1465 |
| 1421 | Has the dataset been used for any tasks already? If so, please provide a description. | | 1466 |
| 1422 | No. | Have any third parties imposed IP-based or other restrictions on the data associated with the instances? If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms, as well as any fees associated with these restrictions. | 1467 |
| 1423 | | No. | 1468 |
| 1424 | What (other) tasks could the dataset be used for? | | 1469 |
| 1425 | This dataset could also be used for text classification and style transfer. | Do any export controls or other regulatory restrictions apply to the dataset or to individual instances? If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any supporting documentation. | 1470 |
| 1426 | | No restrictions. | 1471 |
| 1427 | Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses? For example, is there anything that a future user might need to know to avoid uses that could result in unfair treatment of individuals or groups (e.g., stereotyping, quality of service issues) or other undesirable harms (e.g., financial harms, legal risks) If so, please provide a description. Is there anything a future user could do to mitigate these undesirable harms? | | 1472 |
| 1428 | No. | | 1473 |
| 1429 | | | 1474 |
| 1430 | | | 1475 |
| 1431 | | | 1476 |
| 1432 | | | 1477 |
| 1433 | | | 1478 |
| 1434 | | | 1479 |
| 1435 | | | 1480 |
| 1436 | | | 1481 |
| 1437 | | | 1482 |
| | | Maintenance | 1483 |
| | | | 1484 |

Who will be supporting/hosting/maintaining the dataset?

The dataset will be maintained by the primary author of the paper.

How can the owner/curator/manager of the dataset be contacted (e.g., email address)?

By submitting an issue on the dataset’s GitHub repository.

Is there an erratum? If so, please provide a link or other access point.

Yes, errata can be reported and tracked via GitHub issues.

Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)?

If so, please describe how often, by whom, and how updates will be communicated to users (e.g., mailing list, GitHub)?

Yes, updates will be handled by the repository maintainer on GitHub. Users can receive update notifications by subscribing to the repository.

If the dataset relates to people, are there applicable limits on the retention of the data associated with the instances (e.g., were individuals in question told that their data would be retained for a fixed period of time and then deleted)? If so, please describe these limits and explain how they will be enforced.

This dataset does not contain or pertain to any personal data.

Will older versions of the dataset continue to be supported/hosted/maintained? If so, please describe how. If not, please describe how its obsolescence will be communicated to users.

Yes, previous versions will remain available in the “Releases” section of the GitHub repository.

If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so? If so, please provide a description. Will these contributions be validated/verified? If so, please describe how. If not, why not? Is there a process for communicating/distributing these contributions to other users? If so, please provide a description.

Yes, contributors may open a GitHub issue and submit a pull request. They should mention the maintainer and clearly describe their proposed changes, which will then be reviewed and validated before being merged.

B Implementation Details

B.1 Multi-Task Methods

MTL-LoRA LLMs are trained for 6 epochs maximum, using the AdamW optimizer (Loshchilov and Hutter, 2019) with the following parameters: $\epsilon = 10^{-9}$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, and a weight decay of $\lambda = 0.01$. A linear learning rate scheduler with a 10% warm-up ratio is employed. The training batch size is fixed at 4, with 4 steps gradient accumulation and training tasks are randomly sampled. The learning rate is chosen from the set $\{1 \cdot 10^{-5}, 2 \cdot 10^{-5}, 5 \cdot 10^{-5}, 1 \cdot 10^{-4}\}$, and hyperparameter selection is performed to maximize SRB \uparrow . According to experimental findings, LoRA and MTL-LoRA hyperparameters are set to $r = 128$ and $attn_matrices = W_{QKVO}$. Moreover, we chose $\alpha = r$ to keep a 1:1 ratio so as not to overpower the backbone (Lee et al., 2023). For MTL-LoRA configuration, sharpness of the weight distribution is fixed at 0.5 and the optimal n up-projections is selected among $\{1, 2, 3\}$. We rely on the implementation provided by Adapters library (Poth et al., 2023) for all PEFT methods. Best hyperparameters for PEFT methods are in Table 4.

MTL-RAG To facilitate few-shot multi-task learning within the in-context learning framework, we develop a multi-task extension of Retrieval-Augmented Generation (RAG). Our approach retrieves demonstrations from various tasks and integrates them into the prompt. We conduct experiments using 1, 2, and 3 examples per task, analyzing how the ordering of tasks and examples within the prompt influences performance. We investigate three strategies for sequencing demonstrations in the prompt as mentioned in Section 4.2: random, grouped and interleaved orderings.

The optimal hyperparameters for in-context learning are summarized in Table 5.

B.2 Models

We utilize the following instruct models for In-Context Learning (ICL):

- Llama-3.1-8B-Instruct

- Mistral-7B-Instruct-v0.3

For experiments involving Parameter-Efficient Fine-Tuning (PEFT), we employ the following base models:

- Llama-3.1-8B

| | | | Batch size | lr | Acc. steps | Epochs | $\alpha = \tau$ | Attn. matrices | n up proj. | τ |
|------------|----------|-------|------------|-------------------|------------|--------|-----------------|----------------|------------|--------|
| LlaMA-3-8B | LoRA | E | 4 | $1 \cdot 10^{-4}$ | 4 | 6 | 128 | W_{QKVO} | - | - |
| | MTL-LoRA | E,O,W | 4 | $1 \cdot 10^{-4}$ | 4 | 6 | 128 | W_{QKVO} | 3 | 0.5 |
| | | E,O | 4 | $1 \cdot 10^{-4}$ | 4 | 6 | 128 | W_{QKVO} | 3 | 0.5 |
| | | E,W | 4 | $1 \cdot 10^{-4}$ | 4 | 6 | 128 | W_{QKVO} | 3 | 0.5 |
| Mistral-7B | LoRA | E | 4 | $1 \cdot 10^{-4}$ | 4 | 6 | 128 | W_{QKVO} | - | - |
| | MTL-LoRA | E,O,W | 4 | $1 \cdot 10^{-4}$ | 4 | 6 | 128 | W_{QKVO} | 3 | 0.5 |
| | | E,O | 4 | $5 \cdot 10^{-5}$ | 4 | 6 | 128 | W_{QKVO} | 3 | 0.5 |
| | | E,W | 4 | $1 \cdot 10^{-4}$ | 4 | 6 | 128 | W_{QKVO} | 3 | 0.5 |

Table 4: PEFT hyperparameter configurations chosen based on SRB performance on the ETR-fr validation set. E, O, and W refer to ETR-fr, OrangeSum, and WikiLarge FR, respectively.

| | | | k | Ordering |
|------------|-----------|-------|---|-------------|
| Mistral-7B | Zero-Shot | E | - | - |
| | CoT | E | - | - |
| | RAG | E | 7 | Random |
| | | E,O | 3 | Random |
| | | E,W | 3 | Random |
| | | E,O,W | 3 | Interleaved |
| LlaMA-3-8B | Zero-Shot | E | - | - |
| | CoT | E | - | - |
| | RAG | E | 9 | Random |
| | | E,O | 3 | Random |
| | | E,W | 3 | Random |
| | | E,O,W | 2 | Random |

Table 5: ICL hyperparameter configurations selected based on SRB performance on the ETR-fr validation set. Here, E denotes ETR-fr, O denotes OrangeSum, and W denotes WikiLarge FR.

- Mistral-7B-v0.3
- DeepSeek-R1-Distill-Llama-8B

B.3 Metrics

Text Descriptive Statistics To calculate the descriptive statistics presented in Table 1, such as word count, sentence length, compression ratio, KMRE, and others, we employ the `TextDescriptives` (Hansen et al., 2023) and `textacy` Python libraries, both of which use the `fr_core_news_md-3.8.0` model from `SpaCy`.

ROUGE (Recall-Oriented Understudy for Gisting Evaluation) (Lin, 2004) is a widely used metric for assessing the quality of automatically generated summaries by measuring n-gram and sequence overlap with reference texts. Specifically, we report the F1-scores for ROUGE-1 (ROUGE-1), ROUGE-2 (ROUGE-2), and ROUGE-L (ROUGE-L), which capture overlap of unigrams, bigrams, and longest

common subsequences, respectively. The F1-score represents the harmonic mean of precision and recall. For evaluation, we use Hugging Face’s interface to Google’s official implementation.

BERTScore (Zhang et al., 2020) is based on the contextual word representations generated by BERT-like encoders. Unlike traditional metrics like BLEU or ROUGE, which rely on exact lexical matches, BERTScore uses embeddings to capture finer semantic similarities, offering more flexibility with respect to context and greater robustness to word reordering and synonyms. For each word in the generated text, BERTScore finds the most similar word in the reference text using cosine similarities of their representations. The goal of this step is to align the words in the generated text with those in the reference text. These similarity scores for the aligned word pairs are then aggregated to obtain recall, precision, and F1-score. For reproducibility, we use the Hugging Face’s wrapper coupled with bert-base-multilingual-cased model.

SARI (Sentence-level Accuracy Rating for Text Simplification) (Xu et al., 2016) is commonly used to evaluate sentence and text simplification. Unlike other metrics like BLEU or ROUGE, which focus primarily on lexical similarity to reference texts, SARI takes into account three key aspects of simplification: content preservation (keep), information addition (add), and information deletion (del). For each word or n-gram generated, SARI evaluates whether the word should be kept, added, or deleted by comparing it with its source and the ground truth. The mathematical expression of SARI is the average of the F1-score of these three measures.

$$\text{SARI} = \frac{F1_{\text{keep}} + F1_{\text{add}} + F1_{\text{del}}}{3}$$

For evaluation, we use Hugging Face’s in-

terface, which is adapted from TensorFlow’s tensor2tensor implementation (Vaswani et al., 2018).

KMRE (Kandel-Moles Reading Ease) (Kandel and Moles, 1958) is the French adaptation of the Flesch-Kincaid Reading Ease (FKRE) (Kincaid and Others, 1975), originally designed for English. It measures the complexity of French texts based on sentence length and word length without the need for comparison with a reference text:

$$\text{KMRE} = 207 - 1.015 \left(\frac{\# \text{words}}{\# \text{sentences}} \right) - 73.6 \left(\frac{\# \text{syllables}}{\# \text{words}} \right)$$

KMRE, like the FKRE, is theoretically bounded between 0 and 100. However, it can exceed 100 in rare cases, particularly when the text contains very short sentences and simple, monosyllabic words. This is often the case in ETR documents, which are specifically designed for ease of reading. Moreover, Wubben et al. (2012) advises not to use this metric alone, as it does not account for grammar quality or meaning preservation. This is why we pair it with BERTScore, ROUGE, and SARI, and we do not monitor it for hyperparameter tuning.

LIX (läsbarhetsindex) (Björnsson, 1983) is a readability measure to indicate the difficulty of reading a text. It is calculated using two factors: the average sentence length and the percentage of long words (more than 6 letters) :

$$\text{LIX} = \left(\frac{\# \text{words}}{\# \text{periods}} \right) + 100 \left(\frac{\# \text{long words}}{\# \text{words}} \right)$$

LIX scores typically range from 20 ("very easy") to 60 ("very difficult"). The formula is considered objective and quick to compute compared to other readability measures.

SRB is proposed to measure the quality of a ETR text by aggregating metrics related to ETR transcription characteristics, *i.e.* simplification, summarization, and meaning preservation. To do this, we compute the harmonic mean of SARI, ROUGE-L, and BERTScore-F1:

$$\text{SRB} = \frac{3}{\frac{1}{\text{SARI}} + \frac{1}{\text{R-L}} + \frac{1}{\text{BERTScore-F1}}}$$

Novelty is used to evaluate abstractiveness, measured by the percentage of n-grams in the generated text that do not appear in the source document (See

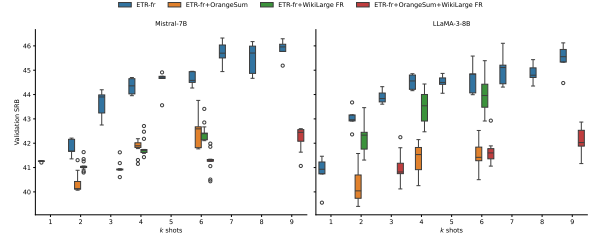


Figure 4: SRB performance score of Mistral-7B and LLaMA-3-8B on the ETR-fr validation set with varying number of in-context examples ($k = 1-9$) and task combinations.

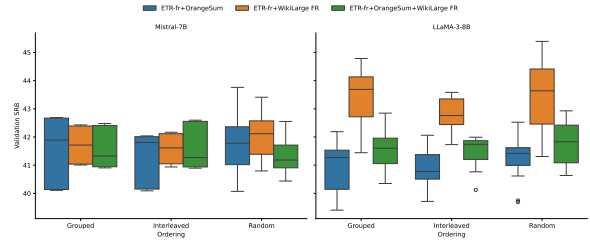


Figure 5: SRB performance of Mistral-7B and LLaMA-3-8B on the ETR-fr validation set under different example ordering strategies and task combination configurations.

et al., 2017; Kamal Eddine et al., 2021). We report only novel 1-grams, excluding stop words (commonly used words in a language).

Compression ratio is the proportion of the document that has been removed. A higher compression ratio indicates more reduction, meaning the summary is more condensed compared to the original document.

$$\text{Comp. Ratio} = 1 - \frac{\# \text{words in ETR}}{\# \text{words in source}}$$

C In-Context Learning Hyperparameters Effects

Figure 6 illustrates examples of prompts used for zero-shot (Fig. 6a), chain-of-thought (Fig. 6c) and few-shot (Fig. 6b).

C.1 Impact of the Number of Shots on ETR-fr Performance

Figure 4 presents the performance of LLaMA-3-8B and Mistral-7B on the French text simplification benchmark (ETR-fr) across varying numbers of in-context learning (ICL) examples ($k = 1$ to 9) and under different training configurations.

LLaMA-3-8B Performance. For the LLaMA-3-8B model, performance generally increases with larger k values. The basic task ETR-fr alone yields steadily rising median SRB scores, from 40.93 at $k = 1$ to 45.96 at $k = 9$. The incorporation of auxiliary datasets (OrangeSum and WikiLarge FR) leads to varied results. For instance, combining ETR-fr with WikiLarge FR at $k = 2$ raises the median from 42.96 to 42.33, while the three-dataset combination at $k = 6$ has a lower median of 41.60 compared to 44.84 for ETR-fr alone. This suggests diminishing returns or even negative interference when too many tasks are combined.

Mistral-7B Performance. The Mistral-7B model demonstrates a similar trend of improved performance with increasing k values for the ETR-fr task. Median SRB rise from 41.26 at $k = 1$ to 45.96 at $k = 9$. However, Mistral exhibits less variation across configurations. The inclusion of OrangeSum and WikiLarge FR improves SRB modestly, and the three-dataset combination remains slightly below the single-task performance. For example, at $k = 6$, ETR-fr alone achieves a median of 44.58, whereas the triple combination achieves only 41.28.

Comparative Insights. When comparing both models, LLaMA-3-8B tends to show greater gains from dataset combinations than Mistral-7B, although it also experiences more variance. For both models, the highest performances are obtained when using ETR-fr alone at higher k values, indicating that overloading the prompt context with multiple tasks may dilute performance. Moreover, the higher maximum SRB for LLaMA across configurations (e.g., up to 46.12) suggest it may have a higher performance ceiling, but with more fluctuation.

C.2 Conclusion

In summary, increasing the in-context learning size (k) generally improves model performance. Task combination has mixed effects: beneficial in some configurations but detrimental in others, especially when too many tasks are combined. LLaMA-3-8B appears more sensitive to these changes than Mistral-7B, highlighting important considerations for prompt engineering.

C.2.1 Impact of the Tasks Ordering on ETR-fr Performance

Figure 5 presents the impact of task ordering on model performance under different multi-task training configurations. For both models, three types of example ordering are compared: *grouped*, *interleaved*, and *random*. Each ordering is evaluated with different training task combinations, such as ETR-fr+OrangeSum, ETR-fr+WikiLarge FR, and ETR-fr+OrangeSum+WikiLarge FR.

LLaMA-3-8B Performance. For LLaMA-3-8B, performance consistently improves when WikiLarge FR data is added to the training set. The configuration using only ETR-fr+WikiLarge FR yields the highest SRB scores across all ordering methods, particularly under the random strategy, which achieves the highest maximum score (45.39). Overall, grouped and random orderings tend to result in higher median and upper-quartile SRB compared to interleaved ordering, indicating that the sequential arrangement of examples plays a role in performance.

Mistral-7B Performance. For Mistral-7B, the impact of training set composition is similarly positive, with improvements observed upon including WikiLarge FR. However, the differences among the three ordering strategies are more subtle. grouped and interleaved yield very similar statistics, with slight advantages in median SRB depending on the training data. The highest maximum score for Mistral-7B (43.76) occurs under the random strategy with the ETR-fr+OrangeSum dataset, although this configuration does not have the most consistent results across runs.

Comparative Insights. Comparing the two models, LLaMA-3-8B generally outperforms Mistral-7B in terms of median and maximum SRB, particularly when trained with ETR-fr and WikiLarge FR. Mistral-7B demonstrates more stable performance with narrower score ranges but slightly lower central tendency metrics. These results suggest that while both models benefit from enriched prompts, LLaMA-3-8B exhibits greater potential for high-end performance when paired with appropriate example ordering and task combinations.

D Complementary Evaluation Results

D.1 Quantitative Results

The average performances of various methods on the ETR-fr and ETR-fr-politic test sets is presented in tables 6a and 6b, respectively. These results compare In-Context Learning (ICL) techniques, such as Zero-shot, Chain-of-Thought (CoT), and Retrieval-Augmented Generation (RAG), against Parameter-Efficient Fine-Tuning (PEFT) methods including LoRA and MTL-LoRA. Evaluations are conducted across different LLM models (Mistral-7B, LLaMA-3-8B and DeepSeek-R1-8B) and task combinations (E: ETR-fr, O: OrangeSum, W: WikiLarge FR). Metrics such as ROUGE (R-1, R-2, R-L), SARI, BERTScore-F1, SRB, Compression Ratio, and Novelty are used to provide a comprehensive performance overview.

The experimental results clearly highlight the performance benefits of both retrieval augmentation and fine-tuning approaches, particularly under multi task settings.

In-Context Learning (ICL) Zero-Shot and CoT settings generally underperform across all metrics compared to RAG and PEFT. While CoT shows a slight improvement in novelty and informativeness over Zero-Shot, gains are marginal. RAG consistently improves performance over basic prompting, especially on the main ETR-fr test set. For both Mistral-7B and LLaMA-3-8B, RAG with task combinations (E, E+O, E+W, E+O+W) achieves substantial boosts in ROUGE and SARI scores. Notably, RAG yields the highest performance in most individual metrics under the ICL category.

Parameter-Efficient Fine-Tuning (PEFT) PEFT models consistently outperform in-context learning (ICL) methods across all evaluation metrics. Both the LoRA and MTL-LoRA setups yield notable gains in fluency, simplicity, and informativeness. Among them, LLaMA-3-8B-MTL-LoRA achieves the best overall performance, excelling in metrics such as SARI, BERTScore-F1, and compression ratio, highlighting its effectiveness in producing simplified text that remains semantically faithful. The Multi-task LoRA (E+W) variant records the highest scores for SARI (44.67), BERTScore (74.05), and compression ratio (56.11), suggesting a well-balanced approach that preserves meaning while substantially reducing text length. Additionally, we report results for the DeepSeek-R1-8B model; however, its

performance is consistently lower than other LLM configurations, regardless of the fine-tuning strategy applied.

Out-of-Domain (ETR-fr-politic) Performance

On the political subset, the performance gap between ICL and PEFT narrows slightly; however, PEFT models continue to demonstrate a clear advantage. Among the ICL methods, RAG-based approaches retain their relative lead, particularly when augmented with additional context (E+W and E+O+W), indicating stronger generalization capabilities. Notably, the zero-shot LLaMA-3-8B model achieves the highest novelty score (55.73), which could signal greater output diversity, though it might also suggest reduced fidelity. Similar to previous findings, DeepSeek-R1-8B consistently underperforms compared to other LLM configurations, regardless of the fine-tuning method used.

D.2 Human Evaluation

We conduct a comprehensive human evaluation on two datasets, ETR-fr and ETR-fr-politic, assessing the generated explanations along dimensions guided by the ETR framework and general language quality metrics. Results are reported in Tables 7 and 8.

Explanation Criteria (ETR dimensions). On ETR-fr, all methods exhibit strong performance across information selection, word selection, and sentence construction (scores >0.88), with the LoRA method slightly outperforming others in word selection (0.94) and overall global quality (0.91). Illustration quality, however, remains a consistent weakness across methods, with high variance indicating instability or inconsistent strategy for visual grounding.

For the more challenging ETR-fr-politic, overall scores decrease across all explanation criteria. Notably, RAG with joint training on E and W achieves the best global score (0.80), outperforming LoRA and MTL-LoRA. While RAG maintains high scores in information selection and sentence construction, illustration scores remain low across the board, underscoring the difficulty of generating coherent examples or analogies in politically sensitive domains.

General Language Quality. As shown in Table 8, RAG again performs competitively on both datasets. On ETR-fr, it achieves the highest ratings in grammar and coherence (both > 4.4), with

| | Method | Task | R-1 \uparrow | R-2 \uparrow | R-L \uparrow | SARI \uparrow | BERT-F1 \uparrow | SRB \uparrow | Comp. ratio | Novelty |
|---------------------|-----------|------------------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|
| In Context Learning | | | | | | | | | | |
| Mistral-7B | Zero-Shot | E | 23.96 \pm 0.04 | 7.08 \pm 0.01 | 16.25 \pm 0.03 | 37.07 \pm 0.00 | 69.75 \pm 0.00 | 29.17 \pm 0.03 | -64.14 \pm 0.00 | 35.70 \pm 0.00 |
| | CoT | E | 23.53 \pm 0.06 | 7.23 \pm 0.01 | 16.20 \pm 0.04 | 37.39 \pm 0.00 | 68.80 \pm 0.00 | 29.12 \pm 0.05 | -60.53 \pm 0.00 | <u>36.09</u> \pm 0.00 |
| | RAG | E | <u>31.91</u> \pm 0.66 | <u>10.77</u> \pm 0.65 | <u>22.54</u> \pm 0.75 | <u>40.14</u> \pm 0.57 | <u>72.17</u> \pm 0.30 | <u>36.08</u> \pm 0.80 | 45.23 \pm 1.17 | 27.27 \pm 0.58 |
| | | E,O | 30.36 \pm 0.47 | 9.61 \pm 0.34 | 21.80 \pm 0.30 | 39.49 \pm 0.12 | 71.07 \pm 0.18 | 35.19 \pm 0.29 | <u>47.99</u> \pm 1.91 | 26.80 \pm 0.84 |
| | | E,W | 30.46 \pm 0.48 | 9.93 \pm 0.17 | 21.72 \pm 0.34 | 38.76 \pm 0.43 | 71.57 \pm 0.14 | 34.96 \pm 0.34 | 35.08 \pm 2.13 | 23.32 \pm 0.31 |
| | E,O,W | 29.85 \pm 0.04 | 9.58 \pm 0.03 | 21.53 \pm 0.05 | 39.53 \pm 0.00 | 71.06 \pm 0.00 | 34.98 \pm 0.05 | 46.42 \pm 0.00 | 25.83 \pm 0.00 | |
| LlaMA-3-8B | Zero-Shot | E | 24.90 \pm 0.20 | 8.16 \pm 0.25 | 17.10 \pm 0.38 | 38.48 \pm 0.38 | 70.15 \pm 0.17 | 30.38 \pm 0.48 | -22.52 \pm 2.47 | 39.13 \pm 0.92 |
| | CoT | E | 27.23 \pm 0.91 | 8.81 \pm 0.21 | 18.34 \pm 0.57 | 38.15 \pm 0.23 | 70.79 \pm 0.52 | 31.62 \pm 0.65 | 7.59 \pm 4.82 | 30.33 \pm 1.75 |
| | RAG | E | <u>33.05</u> \pm 0.72 | <u>12.23</u> \pm 0.44 | <u>23.77</u> \pm 0.68 | <u>41.66</u> \pm 0.45 | <u>72.59</u> \pm 0.38 | <u>37.57</u> \pm 0.70 | <u>43.36</u> \pm 2.62 | 27.06 \pm 0.29 |
| | | E,O | 30.77 \pm 0.35 | 10.85 \pm 0.31 | 22.10 \pm 0.35 | 39.84 \pm 0.22 | 71.13 \pm 0.17 | 35.54 \pm 0.32 | 24.36 \pm 30.13 | 25.02 \pm 1.84 |
| | | E,W | 32.14 \pm 0.56 | 11.70 \pm 0.34 | 23.11 \pm 0.19 | 40.49 \pm 0.32 | 71.88 \pm 0.18 | 36.64 \pm 0.24 | 42.30 \pm 1.59 | 26.70 \pm 0.92 |
| | E,O,W | 30.53 \pm 0.74 | 10.67 \pm 0.45 | 21.65 \pm 0.71 | 39.24 \pm 0.20 | 71.21 \pm 0.26 | 35.00 \pm 0.67 | 31.18 \pm 4.94 | 24.08 \pm 1.37 | |
| PEFT | | | | | | | | | | |
| Mistral-7B | LoRA | E | 32.45 \pm 0.03 | 12.38 \pm 0.02 | 23.99 \pm 0.05 | 42.09 \pm 0.00 | 73.56 \pm 0.00 | 37.95 \pm 0.04 | 44.42 \pm 0.00 | 18.35 \pm 0.00 |
| | MTL-LoRA | E,O | 32.62 \pm 0.04 | 12.73 \pm 0.01 | 24.29 \pm 0.04 | 41.95 \pm 0.00 | 73.52 \pm 0.00 | 38.16 \pm 0.03 | 53.48 \pm 0.00 | 24.17 \pm 0.00 |
| | | E,W | 32.68 \pm 0.05 | <u>12.91</u> \pm 0.01 | 24.25 \pm 0.03 | <u>42.53</u> \pm 0.00 | <u>73.90</u> \pm 0.00 | 38.33 \pm 0.03 | <u>53.62</u> \pm 0.00 | <u>24.99</u> \pm 0.00 |
| | | E,O,W | 33.60 \pm 0.05 | 12.81 \pm 0.05 | <u>24.89</u> \pm 0.04 | 42.25 \pm 0.00 | <u>73.62</u> \pm 0.00 | <u>38.74</u> \pm 0.03 | 48.93 \pm 0.00 | 23.38 \pm 0.00 |
| LlaMA-3-8B | LoRA | E | 31.80 \pm 0.03 | 13.16 \pm 0.09 | 24.92 \pm 0.18 | 42.15 \pm 0.01 | 72.84 \pm 0.17 | 38.67 \pm 0.17 | 50.50 \pm 0.28 | 18.37 \pm 0.88 |
| | MTL-LoRA | E,O | <u>33.38</u> \pm 0.06 | 13.16 \pm 0.05 | 24.20 \pm 0.04 | 43.06 \pm 0.01 | 73.88 \pm 0.01 | 38.42 \pm 0.03 | 50.90 \pm 0.40 | 23.25 \pm 0.17 |
| | | E,W | 32.54 \pm 0.05 | 13.50 \pm 0.06 | 25.01 \pm 0.06 | 44.67 \pm 0.00 | 74.05 \pm 0.00 | 39.54 \pm 0.05 | <u>56.11</u> \pm 0.00 | <u>33.05</u> \pm 0.00 |
| | | E,O,W | 32.78 \pm 0.02 | 13.67 \pm 0.03 | 25.55 \pm 0.16 | 43.58 \pm 0.10 | 73.33 \pm 0.09 | 39.62 \pm 0.09 | 52.66 \pm 1.00 | 24.27 \pm 0.21 |
| DeepSeek-R1-8B | LoRA | E | 20.45 \pm 0.65 | 7.72 \pm 0.29 | 15.40 \pm 0.13 | 41.29 \pm 0.04 | 66.02 \pm 0.26 | 28.76 \pm 0.16 | -4.61 \pm 3.83 | 21.86 \pm 0.29 |
| | MTL-LoRA | E,O | 23.70 \pm 0.32 | 8.86 \pm 0.04 | 18.18 \pm 0.33 | 42.91 \pm 0.06 | 66.72 \pm 0.24 | 32.15 \pm 0.37 | <u>8.57</u> \pm 1.08 | 27.92 \pm 0.86 |
| | | E,W | <u>25.38</u> \pm 0.11 | <u>9.35</u> \pm 0.05 | <u>18.52</u> \pm 0.07 | <u>43.06</u> \pm 0.03 | <u>68.08</u> \pm 0.14 | <u>32.64</u> \pm 0.08 | -0.52 \pm 2.52 | <u>36.16</u> \pm 0.30 |
| | | E,O,W | 22.70 \pm 0.10 | 7.93 \pm 0.01 | 16.59 \pm 0.02 | 42.94 \pm 0.00 | 67.18 \pm 0.00 | 30.47 \pm 0.02 | -9.35 \pm 0.00 | 29.50 \pm 0.00 |

(a) Performance on ETR-fr test set.

| | Method | Task | R-1 \uparrow | R-2 \uparrow | R-L \uparrow | SARI \uparrow | BERT-F1 \uparrow | SRB \uparrow | Comp. ratio | Novelty |
|---------------------|-----------|------------------|------------------|------------------|------------------|------------------|--------------------|------------------|--------------------|------------------|
| In Context Learning | | | | | | | | | | |
| Mistral-7B | Zero-Shot | E | 28.42 \pm 0.12 | 10.98 \pm 0.07 | 19.31 \pm 0.03 | 39.87 \pm 0.00 | 68.10 \pm 0.00 | 32.77 \pm 0.03 | -309.24 \pm 0.00 | 48.37 \pm 0.00 |
| | CoT | E | 29.80 \pm 0.03 | 11.21 \pm 0.05 | 19.88 \pm 0.08 | 39.62 \pm 0.00 | 69.40 \pm 0.00 | 33.35 \pm 0.07 | -261.30 \pm 0.00 | 50.85 \pm 0.00 |
| | RAG | E | 40.19 \pm 0.63 | 16.07 \pm 0.60 | 28.25 \pm 0.31 | 41.40 \pm 0.46 | 73.01 \pm 0.34 | 40.96 \pm 0.35 | 9.00 \pm 3.96 | 23.21 \pm 2.39 |
| | | E,O | 37.49 \pm 0.61 | 14.50 \pm 0.35 | 26.38 \pm 0.69 | 39.46 \pm 0.35 | 72.27 \pm 0.26 | 38.92 \pm 0.58 | 14.26 \pm 2.65 | 17.57 \pm 1.61 |
| | | E,W | 39.65 \pm 0.19 | 15.36 \pm 0.35 | 27.85 \pm 0.38 | 40.08 \pm 0.36 | 72.35 \pm 0.29 | 40.17 \pm 0.23 | 8.72 \pm 1.73 | 17.47 \pm 1.68 |
| E,O,W | | 39.14 \pm 0.04 | 15.96 \pm 0.09 | 28.40 \pm 0.11 | 40.74 \pm 0.00 | 72.87 \pm 0.00 | 40.82 \pm 0.07 | 14.63 \pm 0.00 | 18.33 \pm 0.00 | |
| LlaMA-3-8B | Zero-Shot | E | 29.10 \pm 0.40 | 10.68 \pm 0.35 | 18.70 \pm 0.41 | 40.68 \pm 0.48 | 68.65 \pm 0.11 | 32.39 \pm 0.51 | -178.23 \pm 7.77 | 55.73 \pm 1.07 |
| | CoT | E | 31.15 \pm 0.99 | 10.47 \pm 0.81 | 19.54 \pm 0.65 | 39.80 \pm 0.63 | 69.66 \pm 0.43 | 33.09 \pm 0.74 | -70.57 \pm 8.09 | 47.80 \pm 1.71 |
| | RAG | E | 37.68 \pm 0.53 | 14.46 \pm 0.65 | 26.09 \pm 0.60 | 42.05 \pm 0.90 | 73.01 \pm 0.20 | 39.57 \pm 0.41 | 1.47 \pm 6.45 | 41.78 \pm 0.86 |
| | | E,O | 37.43 \pm 2.11 | 14.28 \pm 0.89 | 25.92 \pm 1.42 | 40.95 \pm 0.90 | 72.41 \pm 0.61 | 39.05 \pm 1.37 | -7.72 \pm 14.32 | 31.85 \pm 1.69 |
| | | E,W | 39.99 \pm 1.10 | 16.27 \pm 0.61 | 27.84 \pm 1.10 | 42.41 \pm 0.43 | 73.83 \pm 0.47 | 41.06 \pm 0.96 | 13.46 \pm 2.37 | 36.72 \pm 2.01 |
| | | E,O,W | 38.33 \pm 1.46 | 15.12 \pm 1.08 | 26.89 \pm 1.10 | 41.08 \pm 0.94 | 72.86 \pm 0.51 | 39.86 \pm 1.13 | 6.34 \pm 7.54 | 29.92 \pm 0.48 |
| | | PEFT | | | | | | | | |
| Mistral-7B | LoRA | E | 35.10 \pm 0.04 | 12.28 \pm 0.04 | 25.97 \pm 0.03 | 38.04 \pm 0.00 | 70.28 \pm 0.00 | 37.96 \pm 0.02 | 21.55 \pm 0.00 | 11.79 \pm 0.00 |
| | MTL-LoRA | E,O | 29.29 \pm 0.07 | 11.02 \pm 0.01 | 21.90 \pm 0.04 | 38.68 \pm 0.00 | 69.22 \pm 0.00 | 34.90 \pm 0.03 | 36.68 \pm 0.00 | 40.29 \pm 0.00 |
| | | E,W | 34.32 \pm 0.06 | 12.60 \pm 0.07 | 24.87 \pm 0.11 | 38.72 \pm 0.00 | 70.54 \pm 0.00 | 37.40 \pm 0.09 | 22.51 \pm 0.00 | 19.10 \pm 0.00 |
| | | E,O,W | 36.34 \pm 0.10 | 13.24 \pm 0.02 | 26.29 \pm 0.08 | 38.39 \pm 0.00 | 70.97 \pm 0.00 | 38.37 \pm 0.06 | 18.33 \pm 0.00 | 10.55 \pm 0.00 |
| LlaMA-3-8B | LoRA | E | 34.65 \pm 1.43 | 13.34 \pm 0.85 | 26.40 \pm 0.95 | 39.70 \pm 0.35 | 70.73 \pm 0.99 | 38.85 \pm 0.90 | 4.67 \pm 2.97 | 16.19 \pm 0.11 |
| | MTL-LoRA | E,O | 32.17 \pm 0.52 | 11.94 \pm 0.23 | 23.98 \pm 0.22 | 39.35 \pm 0.44 | 69.49 \pm 0.21 | 36.81 \pm 0.06 | 17.14 \pm 0.98 | 20.01 \pm 0.62 |
| | | E,W | 37.58 \pm 0.12 | 13.68 \pm 0.05 | 27.02 \pm 0.03 | 38.26 \pm 0.00 | 71.30 \pm 0.00 | 38.88 \pm 0.02 | 8.45 \pm 0.00 | 6.44 \pm 0.00 |
| | | E,O,W | 36.38 \pm 0.22 | 13.72 \pm 0.07 | 25.75 \pm 0.23 | 36.19 \pm 0.00 | 70.94 \pm 0.04 | 37.24 \pm 0.17 | 8.76 \pm 0.13 | 2.04 \pm 0.05 |
| DeepSeek-R1-8B | LoRA | E | 23.89 \pm 0.27 | 7.57 \pm 0.30 | 18.48 \pm 0.30 | 39.34 \pm 0.32 | 63.60 \pm 0.24 | 31.49 \pm 0.34 | -50.45 \pm 2.83 | 24.56 \pm 1.15 |
| | MTL-LoRA | E,O | 26.81 \pm 1.84 | 8.41 \pm 0.40 | 19.60 \pm 1.15 | 39.03 \pm 0.16 | 65.02 \pm 0.42 | 32.58 \pm 1.08 | -38.60 \pm 0.76 | 25.44 \pm 0.10 |
| | | E,W | 26.53 \pm 0.79 | 9.77 \pm 0.86 | 18.97 \pm 0.73 | 39.47 \pm 0.49 | 65.37 \pm 0.23 | 32.14 \pm 0.83 | -49.42 \pm 0.94 | 21.95 \pm 0.85 |
| | | E,O,W | 29.83 \pm 0.04 | 11.18 \pm 0.04 | 21.13 \pm 0.07 | 36.58 \pm 0.00 | 67.35 \pm 0.00 | 33.51 \pm 0.06 | -46.02 \pm 0.00 | 4.32 \pm 0.00 |

strong fluency and relevance. MTL-LoRA slightly improves grammaticality, but this does not translate to gains in perceived overall quality.

In the political domain, quality metrics decline, consistent with the ETR scores. RAG trained on E and W maintains robust fluency and coherence, achieving the best overall quality score (3.76). In contrast, MTL-LoRA’s performance degrades notably in global quality (2.62), despite competitive scores in coherence and relevance, suggesting potential trade-offs introduced by multitask learning in more nuanced domains.

Summary. These results highlight RAG’s robustness across both explanation and linguistic quality metrics, particularly when trained jointly on E and W. The consistent underperformance in illustration generation across all models indicates a need for future work on grounded or multimodal explanation strategies, especially in high-stakes domains like politics.

D.3 Comparison of Ground Truth and Generated ETR Outputs

The table 9 presents a detailed comparison of different model configurations (Mistral-7B and LLaMA-8B), training methods (RAG, LoRA, MTL-LoRA), and task combinations (E: Explanation, O: Observation, W: Writing). Metrics include the average number of words and sentences, sentence length, KMRE (higher is better), novelty, and compression ratio.

Overall, models trained with MTL-LoRA tend to generate more concise outputs while maintaining strong performance in terms of KMRE. For instance, LLaMA-8B + MTL-LoRA (E,W) achieves the highest KMRE score (102.98) and the highest novelty (33.05), indicating its ability to produce informative and diverse content.

RAG-based methods generally generate longer texts, with higher sentence lengths (up to 11.07 words on average for LLaMA-8B + RAG (E,O,W)), but often at the expense of novelty. This suggests that RAG relies more heavily on retrieved content, which may reduce the originality of generated text.

Compared to the ground truth, the generated texts generally contain more words and exhibit equal or greater sentence lengths. Notably, the MTL-LoRA configurations achieve higher compression ratios, highlighting their ability to effectively condense information. While no method fully replicates the characteristics of the test set,

defined by its notably short sentences and high compression. LLaMA-8B MTL-LoRA trained on Wikilarge (W) and ETR-fr (E) yields outputs that most closely resemble the test set in terms of both compression and sentence structure.

E Human Evaluation Questions

Table 10 presents a comprehensive set of human evaluation questions based on the ETR European guidelines, organized into four key categories: Information Choice, Sentence Construction and Word Choice, Illustrations, and Overall Quality. Each category includes multiple criteria designed to assess the clarity, structure, and accessibility of information provided in a text. For example, the Information Choice section evaluates whether essential information is prioritized, logically ordered, and clearly grouped. Sentence Construction and Word Choice emphasizes linguistic simplicity, clarity, and consistency, discouraging complex vocabulary, metaphors, or abbreviations unless adequately explained. The Illustrations section assesses the use of relatable examples to clarify abstract ideas, while the Quality section covers fluency, grammar, factual correctness, coherence, and other aspects of textual integrity. These criteria serve as a structured framework to ensure texts are understandable, reader-friendly, and fit for purpose.

| | Method | Task | Informations | Words | Sentences | Illustrations | Global |
|------------|----------------|-------|-----------------|-----------------|-----------------|-----------------|-----------------|
| ETR-fr | | | | | | | |
| LlaMA-3-8B | LoRA | E | 0.89 \pm 0.08 | 0.94 \pm 0.04 | 0.91 \pm 0.05 | 0.38 \pm 0.40 | 0.91 \pm 0.04 |
| | MTL-LoRA | E,O,W | 0.88 \pm 0.06 | 0.89 \pm 0.07 | 0.93 \pm 0.04 | 0.50 \pm 0.65 | 0.89 \pm 0.04 |
| | RAG | E | 0.88 \pm 0.07 | 0.92 \pm 0.05 | 0.89 \pm 0.04 | 0.40 \pm 0.52 | 0.89 \pm 0.04 |
| | | E,W | 0.91 \pm 0.05 | 0.88 \pm 0.07 | 0.92 \pm 0.04 | 0.50 \pm 0.44 | 0.89 \pm 0.04 |
| | ETR-fr-politic | | | | | | |
| | LoRA | E | 0.77 \pm 0.14 | 0.66 \pm 0.11 | 0.79 \pm 0.11 | 0.15 \pm 0.24 | 0.73 \pm 0.08 |
| LlaMA-3-8B | MTL-LoRA | E,O,W | 0.69 \pm 0.13 | 0.59 \pm 0.11 | 0.65 \pm 0.12 | 0.27 \pm 0.27 | 0.64 \pm 0.08 |
| | RAG | E | 0.82 \pm 0.09 | 0.74 \pm 0.10 | 0.86 \pm 0.07 | 0.10 \pm 0.23 | 0.78 \pm 0.05 |
| | | E,W | 0.87 \pm 0.06 | 0.75 \pm 0.09 | 0.85 \pm 0.08 | 0.40 \pm 0.37 | 0.80 \pm 0.06 |

Table 7: **Human evaluation of generations based on ETR guideline criteria**, comparing various methods on the ETR-fr and ETR-fr-politic test sets using their optimal ICL and MTL configurations. Each method is evaluated along four explanation dimensions: Informations (information selection), Words (lexical choice), Sentences (sentence construction), Illustrations, and Global representing the overall quality score. Training tasks are abbreviated as E (ETR-fr), O (OrangeSum), and W (WikiLarge FR). Reported scores are means with 95% confidence intervals.

| | Method | Task | Fluency | Grammar | Relevance | Coherence | Overall Quality |
|------------|----------------|-------|-----------------|-----------------|-----------------|-----------------|-----------------|
| ETR-fr | | | | | | | |
| LlaMA-3-8B | LoRA | E | 4.29 \pm 0.26 | 4.57 \pm 0.23 | 3.95 \pm 0.39 | 4.24 \pm 0.32 | 3.95 \pm 0.37 |
| | MTL-LoRA | E,O,W | 4.33 \pm 0.33 | 4.67 \pm 0.22 | 4.10 \pm 0.38 | 4.14 \pm 0.39 | 3.95 \pm 0.44 |
| | RAG | E | 4.43 \pm 0.27 | 4.71 \pm 0.21 | 4.24 \pm 0.38 | 4.43 \pm 0.34 | 4.24 \pm 0.35 |
| | | E,W | 4.43 \pm 0.23 | 4.57 \pm 0.23 | 4.43 \pm 0.34 | 4.52 \pm 0.27 | 3.95 \pm 0.34 |
| | ETR-fr-politic | | | | | | |
| | LoRA | E | 3.90 \pm 0.52 | 4.43 \pm 0.42 | 4.24 \pm 0.43 | 4.24 \pm 0.45 | 3.14 \pm 0.62 |
| LlaMA-3-8B | MTL-LoRA | E,O,W | 3.81 \pm 0.45 | 4.48 \pm 0.34 | 4.40 \pm 0.38 | 4.52 \pm 0.23 | 2.62 \pm 0.55 |
| | RAG | E | 4.24 \pm 0.38 | 4.48 \pm 0.34 | 4.10 \pm 0.35 | 4.33 \pm 0.30 | 3.45 \pm 0.44 |
| | | E,W | 4.33 \pm 0.33 | 4.57 \pm 0.23 | 4.29 \pm 0.29 | 4.43 \pm 0.27 | 3.76 \pm 0.40 |

Table 8: **Human ratings of fluency, grammar, relevance, coherence, and overall quality** for different methods evaluated on the ETR-fr and ETR-fr-politic test sets, using their optimal ICL and MTL configurations. Training tasks are abbreviated as E (ETR-fr), O (OrangeSum), and W (WikiLarge FR). Scores are reported as means with 95% confidence intervals.

| | Method | Tasks | # Words | # Sentences | Sentence length | KMRE \uparrow | Novelty | Comp. ratio |
|------------|---------------------|-----------------|---------|-------------|-----------------|-----------------|---------|-------------|
| | <i>Ground Truth</i> | <i>Test Set</i> | 40.26 | 8.91 | 4.64 | 102.99 | 55.01 | 65.19 |
| Mistral-7B | RAG | E | 66.38 | 7.70 | 8.76 | 99.77 | 26.55 | 44.32 |
| | | E,O | 60.91 | 6.13 | 10.05 | 97.21 | 26.61 | 48.45 |
| | | E,W | 80.74 | 7.83 | 10.67 | 97.37 | 23.01 | 33.80 |
| | | E,O,W | 62.45 | 6.15 | 10.25 | 97.62 | 25.85 | 46.42 |
| LlaMA-8B | RAG | E | 63.72 | 7.87 | 8.38 | 101.70 | 27.14 | 46.18 |
| | | E,O | 74.19 | 7.57 | 9.92 | 97.45 | 24.29 | 39.22 |
| | | E,W | 69.72 | 7.64 | 9.49 | 100.34 | 25.26 | 41.89 |
| | | E,O,W | 87.17 | 8.40 | 11.07 | 97.48 | 23.69 | 25.94 |
| Mistral-7B | LoRA | E | 65.55 | 9.26 | 7.73 | 101.20 | 18.35 | 44.42 |
| | MTL-LoRA | E,O | 56.75 | 8.25 | 7.38 | 102.61 | 24.17 | 53.48 |
| | | E,W | 54.08 | 9.28 | 6.46 | 104.23 | 24.99 | 53.62 |
| | | E,O,W | 60.08 | 8.81 | 7.23 | 101.80 | 23.38 | 48.93 |
| LlaMA-8B | LoRA | E | 56.96 | 8.64 | 7.62 | 100.93 | 18.87 | 50.66 |
| | MTL-LoRA | E,O | 60.08 | 9.87 | 7.00 | 100.84 | 23.06 | 51.36 |
| | | E,W | 50.09 | 9.19 | 6.50 | 102.98 | 33.05 | 56.11 |
| | | E,O,W | 54.06 | 8.77 | 7.42 | 101.35 | 24.39 | 53.24 |

Table 9: Comparison of different model configurations (Mistral-7B and LlaMA-8B) and training methods (RAG, LoRA, MTL-LoRA) across various task combinations (E: ETR-fr, O: OrangeSum, W: WikiLarge FR). The metrics include word count, sentence count, average sentence length, KMRE (higher is better), novelty, and compression ratio. Ground truth statistics from the test set are also provided for reference.

Rewrite this text by following the principles of clarity and accessibility below:

- Provide only essential information. Avoid information overload.
- Present the information in a logical and easy-to-follow order.
- Highlight the main message right from the start.
- Group related information together.
- Repeat important information if it helps understanding.
- Use short and simple sentences.
- Choose easy-to-understand words.
- Clearly explain difficult words, and repeat the explanation if needed.
- Use language appropriate for the intended audience.
- Use the same word to refer to the same thing throughout the text.
- Avoid abstract ideas, metaphors, and complex comparisons.
- Don't use foreign or obscure words without explanation.
- Avoid contractions and texting-style language.
- Speak directly to the reader in a clear and accessible way.
- Ensure that pronouns are always clear and unambiguous.
- Prefer positive phrasing over negative.
- Use the active voice as much as possible.
- Choose simple punctuation.
- Use bullet points or numbers for lists, not commas.
- Write numbers as digits (e.g., 1, 2, 3), not in words.
- Explain acronyms the first time they appear.
- Don't use unexplained abbreviations.
- Write dates out in full for better clarity.
- Limit use of percentages or large numbers, and explain them simply.
- Don't use unnecessary special characters.
- Use concrete examples to explain complex ideas.
- Prefer examples from everyday life.

###Input: <input_text>

###Output:

(a) Zero Shot Prompt

Rewrite this text by following the principles of clarity and accessibility below:

- Provide only essential information. Avoid information overload.
- Present the information in a logical and easy-to-follow order.
- Highlight the main message right from the start.
- Group related information together.
- Repeat important information if it helps understanding.
- Use short and simple sentences.
- Choose easy-to-understand words.
- Clearly explain difficult words, and repeat the explanation if needed.
- Use language appropriate for the intended audience.
- Use the same word to refer to the same thing throughout the text.
- Avoid abstract ideas, metaphors, and complex comparisons.
- Don't use foreign or obscure words without explanation.
- Avoid contractions and texting-style language.
- Speak directly to the reader in a clear and accessible way.
- Ensure that pronouns are always clear and unambiguous.
- Prefer positive phrasing over negative.
- Use the active voice as much as possible.
- Choose simple punctuation.
- Use bullet points or numbers for lists, not commas.
- Write numbers as digits (e.g., 1, 2, 3), not in words.
- Explain acronyms the first time they appear.
- Don't use unexplained abbreviations.
- Write dates out in full for better clarity.
- Limit use of percentages or large numbers, and explain them simply.
- Don't use unnecessary special characters.
- Use concrete examples to explain complex ideas.
- Prefer examples from everyday life.

###Example 1

Task: <task_name>

Input: <example_input>

Output: <example_output>

...

Complete the following example:

Task: ETR

Input: <input_text>

Output:

(b) Few Shot Prompt


```
1. Analyze the text to identify what can be simplified or clarified.
2. Briefly note the points that need improvement (syntax, vocabulary, structure...).
3. Rewrite the text by applying the following guidelines:
- Provide only essential information. Avoid information overload.
- Present the information in a logical and easy-to-follow order.
- Highlight the main message right from the start.
- Group related information together.
- Repeat important information if it helps understanding.
- Use short and simple sentences.
- Choose easy-to-understand words.
- Clearly explain difficult words, and repeat the explanation if needed.
- Use language appropriate for the intended audience.
- Use the same word to refer to the same thing throughout the text.
- Avoid abstract ideas, metaphors, and complex comparisons.
- Don't use foreign or obscure words without explanation.
- Avoid contractions and texting-style language.
- Speak directly to the reader in a clear and accessible way.
- Ensure that pronouns are always clear and unambiguous.
- Prefer positive phrasing over negative.
- Use the active voice as much as possible.
- Choose simple punctuation.
- Use bullet points or numbers for lists, not commas.
- Write numbers as digits (e.g., 1, 2, 3), not in words.
- Explain acronyms the first time they appear.
- Don't use unexplained abbreviations.
- Write dates out in full for better clarity.
- Limit use of percentages or large numbers, and explain them simply.
- Don't use unnecessary special characters.
- Use concrete examples to explain complex ideas.
- Prefer examples from everyday life.
Start by reasoning step by step, then finish by providing the final version.
###Input: <input_text>
###Output:
```

(c) Chain of Thought Prompt

Figure 6: Zero Shot, Chain of Thought and Few Shot Prompts

| Information Choice | Code | Description |
|---------------------------------------|-------|---|
| Information Choice | CI3 | Providing too much information can create confusion. Only important information should be given. Is this criterion met? |
| | CI4 | Are the pieces of information placed in an order that is easy to follow and understand? |
| | CI5 | Is the main information easy to find? |
| | CI6 | Are pieces of information about the same topic grouped together? |
| | CI8 | Are important pieces of information repeated? |
| Sentence construction and word choice | CPM1 | Are the sentences short? |
| | CPM2 | Are the words easy to understand? |
| | CPM3 | Are difficult words clearly explained when you use them? |
| | CPM4 | Are difficult words explained more than once? |
| | CPM5 | Is the language used the most suitable for the people who will use the information? |
| | CPM6 | Is the same word used throughout the document to describe the same thing? |
| | CPM7 | Difficult and abstract ideas like metaphors should not be used. Is this criterion met? |
| | CPM8 | Uncommon words in a foreign language should not be used. Is this criterion met? |
| | CPM9 | Contracted words, like text messaging slang, should not be used. Is this criterion met? |
| | CPM10 | Does the author address directly the people for whom the information is intended? |
| | CPM11 | Can you easily identify to whom or what the pronouns correspond? |
| | CPM12 | Are positive sentences rather than negative ones used whenever possible? |
| | CPM13 | Is the active voice used instead of the passive voice whenever possible? |
| | CPM14 | Is the punctuation simple? |
| | CPM15 | Are bullets or numbers used instead of lists of words separated by commas? |
| | CPM16 | Are numbers written in digits (1, 2, 3) rather than words? |
| | CPM17 | Acronyms should be avoided or explained when used. Is this criterion met? |
| | CPM18 | Abbreviations should not be used. Is this criterion met? |
| | CPM19 | Are dates written out in full? |
| | CPM20 | The use of percentages or large numbers should be limited and always explained. Is this criterion met? |
| | CPM21 | Special characters should not be used. Is this criterion met? |
| Illustrations | I1 | Are there examples to illustrate complex ideas? |
| | I2 | Are examples, as much as possible, drawn from everyday life? |
| Quality | CA1 | Language fluency |
| | CA2 | Grammar / Spelling |
| | CA3 | Factual accuracy |
| | CA4 | Textual coherence |
| | CA5 | Presence of copies from the original text? |
| | CA6 | Presence of chaotic repetitions? |
| | CA7 | Presence of hallucinations? |
| | CA8 | Overall perceived quality |

Table 10: Evaluation criteria, extracted from ETR European guidelines, for information clarity, sentence construction, illustrations, and quality.