# OUNLP at TSAR 2025 Shared Task: Multi-Round Text Simplifier via Code Generation

**Cuong Huynh**
School of Computer Science
University of Oklahoma
cuong@ou.edu

**Jie Cao**
School of Computer Science
University of Oklahoma
jie.cao@ou.edu

## Abstract

This paper describes the OUNLP system submitted to the TSAR-2025 Shared Task (Alva-Manchego et al., 2025), designed for readability-controlled text simplification using LLM-prompting-based generation. Based on the analysis of prompt-based text simplification methods, we discovered an interesting finding that text simplification performance is highly related to the gap between the source CEFR (Arase et al., 2022) level and the target CEFR level. Inspired by this finding, we propose two multi-round simplification methods and generate them via GPT-4o: rule-based simplification (MRS-Rule) and jointly rule-based LLM simplification (MRS-Joint). Our submitted systems ranked 7 out of 20 teams. Later improvements with MRS-Joint show that taking the LLM simplified candidates as the starting point could further boost the multi-round simplification performance [1].

## 1 Introduction

Complex text makes it difficult for language learners and people with limited literacy to read. Text simplification improves learning, accessibility, and information sharing with a wider audience. With the advent of deep learning and large language models (LLMs), simplification performance has improved significantly, supported by the release of important datasets (Imperial et al., 2025). Modern approaches have explored zero-shot prompting (Chi et al., 2023; Barayan et al., 2025; Farajidizaji et al., 2024), instruction tuning (Imperial and Tayyar Madabushi, 2023), and related strategies.

From our baseline analysis of trial data, we observed that a larger gap between the CEFR level of the original sentence and the target level (**CEFR-Gap**) substantially increases the likelihood of simplification failure. This finding highlights the importance of addressing complexity not in a single step but through a structured, iterative process. Building on this insight, we introduced two novel models generated by GPT-4o for multi-round text simplification: MRS-Rule, a rule-based framework that progressively adjusts sentence structures and vocabulary, and MRS-Joint, which integrates rules with prompting techniques to leverage the strengths of both symbolic and generative approaches.

The primary contribution of this work is to show that multi-round small rule-based simplification are more effective at handling large CEFR gaps than conventional single-step approaches. Our proposed MRS-Joint method outperforms the MRS-Rule and baseline models, as validated through extensive experiments and qualitative analyzes. Additionally, we explore the potential of automatic code generation for text simplification, although further refinement remains necessary.

## 2 Task Setup

The goal of the shared task is to simplify a given source text into a target text with the desired CEFR proficiency level (A1<A2<B1<B2<C1<C2). For the datasets, we use the same trial (40 examples) and test (200 examples) data sets provided by the TSAR workshop to build and evaluate our methods. For the evaluation metrics, we follow the same metrics from the official TSAR-2025 shared-task metrics, which covers both readability-level control (CEFR Compliance, we focus on **RMSE**, the distance between predicted and target CEFR levels, the lower the better) and the preservation of meaning by evaluating semantic fidelity between the simplified sentence and the original sentence, or the simplified sentence and a human-written reference via **MeaningBERT** (Beauchemin et al., 2023), denoted as **MB-Orig** and **MB-Ref** respectively [2].

---

[1] https://github.com/ounlp/Multi-Round-Text-Simplifier

[2] Please refer to the shared task paper (Alva-Manchego et al., 2025) for more details of other metrics such as BERTScore (Zhang et al., 2019) etc.

## 3 Motivation for Multi-Round

In this section, we present our baseline model, the Naïve Prompt model, and show that simplification becomes increasingly challenging as the gap between the source and target levels widens.

### 3.1 Baseline: Naïve Prompt-based (Run 1)

We use GPT-4o to generate the code first (denoted as Baseline or "Program 1"), which will call the OpenAI APIs (GPT-4o-mini) with the following prompt from (Barayan et al., 2025). This generated our Run-1 submission of the test data. Please refer to Appendix A.1 for more details.

> **Baseline Prompt**
>
> Please simplify the following Complex Sentence to make it easier to read and understand by {CEFR-LEVEL} CEFR level English learners. {CEFR-LEVEL} level English learner {CEFR-Description}. To simplify, you may replace difficult words with simpler ones, elaborate, or remove them when possible. You may also break down a lengthy sentence into shorter, clear sentences. Ensure the revised sentence is grammatically correct, fluent, and maintains the core message of the original without changing its meaning. Complex Sentence: {Source} Simplified Sentence:

**CEFR Level Prediction** Since the trial data only gives the CEFR level for target text, not for the source text and the simplified texts, we estimate a text proficiency level using three ModernBERT classifiers with the voting mechanism [3]. Each model independently predicts a CEFR label (A1-C2) with a confidence score. We combine predictions via majority voting: the label with the most votes is selected. Ties are broken by the largest sum of confidences, then by the highest single-model confidence; if still tied, we prefer the simpler (lower) level to remain conservative. The resulting CEFR level also determines whether a simplification is still needed for a text.

**CEFR-Gap** We assign an integral value from 0 to 5 for each CEFR level according to the order of (A1<A2<B1<B2<C1<C2). The CEFR gap for each example is defined as the numerical difference between the source level and the targe level (e.g., the CEFR gap between C1 and A2 is 4-1=3). We run the generated program on the 40 trial examples in the trial data, and then study the performance of the baseline models for each group of examples with the same CEFR Gap as Table 1. We found that **RMSE** rises from 0.624 with a one-level gap, to 1.027 with two levels, and then goes further to 1.581 with three levels, indicating that

---

[3]AbdullahBarayan/ModernBERT-base-doc_en-Cefr, ModernBERT-base-doc_sent_en-Cefr, and ModernBERT-base-reference_AllLang2-Cefr2

| CEFR-Gap | RMSE | MB-Orig | MB-Ref |
|----------|------|---------|--------|
| 1 (18)   | 0.624 | 0.859 | 0.832 |
| 2 (18)   | 1.027 | 0.841 | 0.758 |
| 3 (4)    | 1.581 | 0.761 | 0.762 |

Table 1: CEFR-Gap Analysis on CEFR accuracy (RMSE) and meaning preservation. The bracket shows the total number of examples we found in the trial data with that CEFR gap. It shows the larger the gap, the higher the RMSE, the lower the other MB scores.

larger downward steps are harder to control. Meaning preservation also weakens: **MB-Orig** declines from 0.859 to 0.841 and then to 0.761, while **MB-Ref** falls from 0.832 to 0.758 and stays near 0.762 for the widest gap, although that last figure is based on only four samples. These patterns reveal a trade-off: stronger simplification with larger **CEFR-Gap** makes it more difficult to match the target level and to keep the original meaning intact. In short, bigger CEFR gaps demand more radical linguistic changes, which inevitably reduce both level accuracy and semantic fidelity.

## 4 Proposed Multi-Round Methods

Based on the findings in §3, smaller gap between the source and the target CEFR level will be relatively easy to simplify. Hence, we propose to simplify texts with multiple rounds by taking previous simplification results as inputs with two multi-round methods: rule-based simplification (MRS-Rule §4.1) and jointly rule-based and LLM Prompting (MRS-Joint §4.2). For each program, we first demonstrate the prompts and operations to generate and fix, and then briefly analyze the detailed workflow of the generated program. The orange box shows the operations and prompts we used to generate the MRS-Rule Code, while the blue box at the bottom shows the further steps we used to fix the generated code to make it work.

### 4.1 MRS-Rule: Rule-based (Run 2)

The generated code (see more details in Appendix §A.2) for MRS-Rule does not call any large language model API for simplification, but only rule-based rewriting combined with automatic CEFR level verification and semantic checks.

### 4.1.1 Code Generation

Prompts 2.1, 2.2, and 2.3 are the three main prompts that we used to generate the code for the MRS-rule method step by step. When using Prompt 2.2 to instruct GPT-4o for further simplification by jointly checking CEFR level and semantic similarity, it suggests the following rules and is used in a sophisticated candidate generation pipeline (§A.2.1).

- replace_words: substitute complex words with simpler synonyms (e.g., "utilize" → "use", "approximately" → "about").

- simplify_numbers_units: standardize numerical expressions and units (e.g., remove separators, normalize "metres/meters").

- strip_relative_clauses: remove non-essential subordinate clauses (e.g., clauses beginning with *which/that/who/where/when* or discourse markers like *however/although*) to reduce syntactic complexity.

- keep_shortest_clause: select the simplest clause from a multi-clause sentence by choosing the shortest well-formed segment.

- trim_to_limit: shorten the text to a step-dependent word budget while preserving a grammatical ending.

- sentence_split: break long sentences into shorter, more readable parts at punctuation boundaries, then simplify each part.

More importantly, it also smartly suggested sacrificing semantic preservation for higher CEFR-level accuracy, demonstrating improved performance over prompting baseline (Table 2).

### 4.1.2 Workflow

Figure 1 shows the workflow of MRS-Rule, which includes iterative retries with dynamic conditions such as similarity floor, maximum editing steps to reach the best-effort CEFR-levels.

**Reconciliation Retries** In each retry, the system first generates multiple candidate sentences from the original text. Then, the best candidate is selected using cosine similarity and the predicted CEFR level. This candidate becomes the seed for the next round, based on the assumption that easier sentences can be further simplified toward the target CEFR level. Candidates are created using one or more rules (details in §A.2.1). After each round, all candidates are scored for meaning preservation (cosine similarity) and difficulty (CEFR level). The best-scoring candidate is carried forward as the seed for the next round. If it still does not reach the target level, additional rule-based refinements are applied (§4.1.1). Subsequent retries follow the same process, but use more relaxed thresholds. The CEFR level is validated by majority vote from three ModernBERT classifiers. Sentences that remain unsimplified go through further retries with gradually looser similarity thresholds and larger edit budgets. Finally, the system picks the candidate closest to the target level, reorders the text IDs, and outputs the results. If any sentences are still not simplified, the system slightly lowers the similarity threshold (to 0.88), increases the maximum edit steps (to 8), and reprocesses only the remaining sentences—up to six rounds. All hyperparameters for our program are summarized in the Appendix Table 4.

**Nearest-level Fill** The simplification will continue for multiple rounds of the above simplification rules until all sentences are simplified to the target level or a retry cap is reached. For sentences that did not be simplified to the target level, we
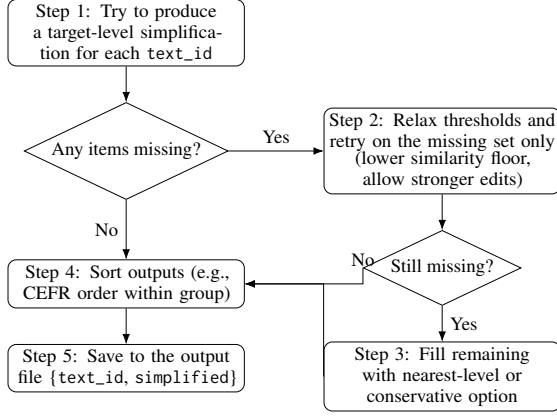
Figure 1: Workflow of **MRS-Rule** (Run 2)

will use the **nearest-level fill**, selecting the candidate whose predicted CEFR level is closest to the target while keeping the original meaning, before reorganizing and saving the final output.

## 4.2 MRS-Joint: Rule-based + Prompting

Building upon the Baseline model (§3), we combine LLM prompts (Barayan et al., 2025) and rule-based multi-round simplifications with automatic verification steps (§4.1). As shown in the workflow (§2), the LLM generates simplified sentences only in the first step. After each retry, the system selects the best candidate on the basis of cosine similarity and predicted CEFR level. This loop continues until the predicted CEFR level matches the target level. In each new retry, the system lowers the cosine similarity threshold (allowing more meaning change) and increases the maximum number of simplification steps. This process ensures that the final sentence fits the target proficiency level while preserving the original meaning.

### 4.2.1 Code Generation

For the MRS-Joint program (§A.3), we use Prompt 3.1 to integrate the LLM prompt from Baseline (Program 1) into the MRS-Rule (Program 2) by uploading the two program files first and then prompting. Then we use Prompt 3.2 to generate the code for over-generation-then-rank. The program, generated when we combined those two files, worked well, so there was nothing to fix.

---

**Prompt for Generating MRS-Joint**

<Operations:> Upload the Program 2(MRS-Rule) and the Program 1(Baseline) to ChatGPT.
**Prompt 3.1** Update this file(the file contains the program 2). Before simplifying the sentence, the program uses the naive prompt to generate one candidate. Other candidates will be generated based on the built-in rules.

**Prompt 3.2** After generating many candidates, the program selects the best candidate based on the cosine similarity and predicted level. If that best candidate does not meet the target level, the program continue generates more candidates based on that best candidate.
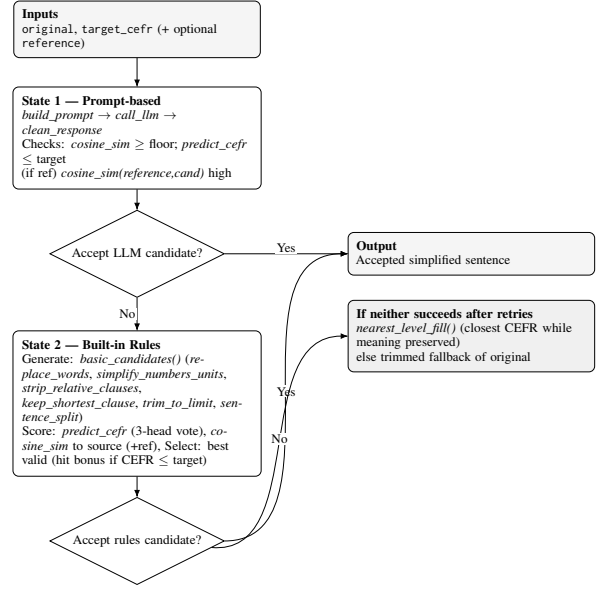
---

### 4.2.2 Workflow



Figure 2: Workflow of **MRS-Joint**

Figure 2 illustrates the **MRS-Joint** generated program by combining LLM prompting (§3.1) and multi-round rule-based simplification (§4.1). The generated program simply prompts the LLM in the first round, and then all subsequent rounds are purely rule-based, as described in §4.1.2.

## 5 Result

| Model | RMSE | MB-Orig | MB-Ref |
|---|---|---|---|
| Trial | | | |
| Baseline (Run 1) | 0.8944 | 0.8453 | 0.7958 |
| MRS-Rule (Run 2) | 0.8515 | 0.7961 | 0.7967 |
| MRS-Joint | 0.4472 | 0.8023 | 0.7574 |
| Test | | | |
| Baseline (Run 1) | 0.755 | 0.855 | 0.849 |
| MRS-Rule (Run 2) | 0.714 | 0.865 | 0.701 |
| MRS-Joint | 0.552 | 0.866 | 0.837 |

Table 2: CEFR accuracy (RMSE) and meaning preservation on trial and test datasets.

Table 2 summarizes the performance of our models in both trial and test datasets. Baseline (§3) and MRS-Rule (§4.1) are the two models corresponding to the two runs of our submission in the final evaluation period. After the evaluation, we found that simply merging two methods into MRS-Joint (§4.2) is more efficient, which is the most accurate model to match the target CEFR level (the best **RMSE**) while still maintaining the meaning. The Prompt-only baseline model (§3) preserves the original meaning best (highest **MeaningBERT-Orig**) but shows the weakest control of CEFR level
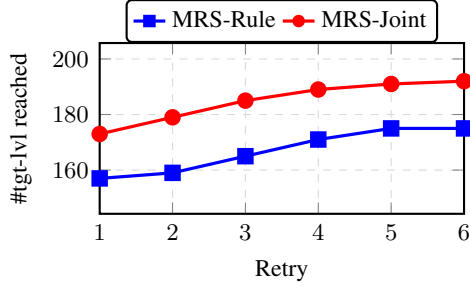
Figure 3: # Simplified sentences that reach the target level across retries for MRS-Rule vs. MRS-Joint.

| Target ↓ / Pred → | A2 | B1 | B2 |
|---|---|---|---|
| **A2** | 66 | 32 | 2 |
| **B1** | 20 | 79 | 1 |
| **B2** | 0 | 0 | 0 |

Table 3: Confusion matrix on the test data

(highest **RMSE**). Furthermore, comparing MRS-Joint with the Baseline, the difference mainly exists in the multi-round rule-based simplification. It shows that our multi-round rules significantly improve the performance with a little sacrifice on the meaning preservation. Figure 3 further shows that single round simplification performs poorly, while multi-round simplification could increasingly simplify more sentences to the target CEFR level. Furthermore, MRS-Joint, starting the simplification from LLM-simplified candidates, could boost the performance of multi-round simplification.

## 6 Qualitative Analysis

### 6.1 Overall Findings

As shown in Table 3, the program excels in simplifying complex sentences C1-C2 to the B1 level. Of 100 source sentences, 79 were successfully simplified to B1, with only 20 dropping further to A2 and 1 rising to B2. For sentences targeted at the A2 level, the results were mixed: only 66 reached the intended A2 level, while 32 overshot B1 and 2 even remained at B2.

Therefore, we recognize that simplifying the input of high-complexity C1–C2 to lower CEFR levels is inherently more challenging. The program is more prone to "overshooting," producing text that remains more complex than the intended target. In other words, the lower the target CEFR level, the higher the likelihood of program's not meeting the constraints of that level.

### 6.2 Case Study

To understand the behavior of the model beyond the overall accuracy scores, we performed a **qualitative error analysis** on three representative examples misclassified by the CEFR predictor. These examples illustrate three different types of misclassification.

**Case 1 – Overshoot (A2 → B1) (§B.1)** The model simplified vocabulary and shortened clauses but kept abstract ideas along with a relative clause typical of the B1 syntax. The CEFR predictor therefore rated the output B1, which is one level higher than the target level, showing that preserving key ideas may force more complex structures than the intended level.

**Case 2 - Lexical Imitation (§B.2)** Although shortened from the source text, the output kept formal phrases like "a large number of bridge accidents... of the bridge itself"" instead of simpler A2 wording such as "Many accidents happen while bridges are being built." The CEFR model therefore rated it B1, showing that better simplification requires lexical adaptation, not just shorter text.

**Case 3 – Under-generation (B1 → A2) (§B.3)** The system produced only a fragment, dropping the telescope's purpose and the planetary-defense discussion. With much of the conceptual content missing, the predictor judged the text A2 despite technical terms. This highlights that incomplete outputs can seem easier to cheat the CEFR predictor as the intended CEFR level.

These examples reveal three failure modes – overshoot, and undergeneration – demonstrating that successful CEFR simplification requires not only simpler words but also balanced control of meaning, style, and completeness.

## 7 Conclusion

We found that a larger gap between the CEFR level of the original and target sentences (CEFR-Gap) increases the likelihood of simplification failure. Based on this finding, we proposed two multi-round simplification methods generated by GPT-4o: MRS-Rule, which applies rule-based simplification, and MRS-Joint, which combines rules with prompting. Extensive experiments and case studies show that MRS-Joint outperforms both the prompting baseline and MRS-Rule, confirming the effectiveness of multi-round simplification and the feasibility of text simplifyer via code generation.

## Limitation

We note a few limitations of our work. The models we used are closed-source models such as using GPT-4o for code generation while using GPT-4o-mini for API, which are not explicitly finetuned in the text simplification datasets by us. Our work is also limited to one dataset and one language (English), and two types of GPT-4o generated model. Furthermore, focusing on coding generation, we could also extend the study to self-evolve algorithm discovery (Novikov et al., 2025) and compare it with other prompts and more coding agents. Besides those, we believe explicitly involving curriculum-based domain knowledge in a structured multi-round simplification will be promising methods in the era of artificial intelligence.

## Lay Summary

This project aims to make complex English sentences easier to understand, especially for language learners. Our team participated in the TSAR 2025 competition, in which the goal was to rewrite sentences to match specific levels of English proficiency, such as beginner (A1) or intermediate (B1), based on the **Common European Framework of Reference (CEFR)**. The insight of our team was that the greater the difference between the original difficulty of a sentence and the target level (called the "CEFR Gap"), the harder it is to simplify the sentence successfully. For example, turning a very advanced sentence (C1) into a basic one (A2) is much more difficult than making small adjustments. This inspired us to develop a **multistep approach** for simplification.

Our team created two systems, and the code is generated with AI with our instructions:

**MRS-Rule**: Uses rules to gradually simplify text in multiple rounds (e.g., replace difficult words, break long sentences).

**MRS-Joint**: Combines a model (GPT-4o-mini) to generate an initial simplified text, and then refines it through multiple rule-based steps.

Both systems repeatedly check whether the new sentence meets the desired CEFR level and still retains the original meaning. If not, they retry those sentences with adjustments. This multi-round process continues until the system either succeeds or picks the closest acceptable version.

In testing, the MRS-Joint method performs best. It reaches the target reading level more often than the baseline approach, although sometimes at the cost of slightly reducing the original meaning. Still, it shows strong overall results: it handles complex sentences better and produced more accurate simplifications. Our team also analyzed the errors. Sometimes, the program "oversimplified" or retained too many complex words. Other times, it shortened the sentence too much and left out important information. These findings will help improve future systems.

In short, this work shows that a multi-step process can make content more accessible to learners while maintaining its original intent.

## References

Fernando Alva-Manchego, Regina Stodden, Joseph Marvin Imperial, Abdullah Barayan, Kai North, and Harish Tayyar Madabushi. 2025. Findings of the TSAR 2025 shared task on readability-controlled text simplification. In *Proceedings of the Fourth Workshop on Text Simplification, Accessibility, and Readability (TSAR 2025)*, Suzhou, China. Association for Computational Linguistics.

Yuki Arase, Satoru Uchida, and Tomoyuki Kajiwara. 2022. Cefr-based sentence difficulty annotation and assessment. *arXiv preprint arXiv:2210.11766*.

Abdullah Barayan, Jose Camacho-Collados, and Fernando Alva-Manchego. 2025. Analysing zero-shot readability-controlled sentence simplification. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 6762–6781, Abu Dhabi, UAE. Association for Computational Linguistics.

David Beauchemin, Horacio Saggion, and Richard Khoury. 2023. Meaningbert: assessing meaning preservation between sentences. *Frontiers in Artificial Intelligence*, 6:1223924.

Alison Chi, Li-Kuang Chen, Yi-Chen Chang, Shu-Hui Lee, and Jason S. Chang. 2023. Learning to paraphrase sentences to different complexity levels. *Transactions of the Association for Computational Linguistics*, 11:1332–1354.

Asma Farajidizaji, Vatsal Raina, and Mark Gales. 2024. Is it possible to modify text to a target readability level? an initial investigation using zero-shot large language models. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 9325–9339, Torino, Italia. ELRA and ICCL.

Joseph Marvin Imperial, Abdullah Barayan, Regina Stodden, Rodrigo Wilkens, Ricardo Munoz Sanchez, Lingyun Gao, Melissa Torgbi, Dawn Knight, Gail Forey, Reka R Jablonkai, and 1 others. 2025. Universalcefr: Enabling open multilingual research on

language proficiency assessment. *arXiv preprint arXiv:2506.01419*.

Joseph Marvin Imperial and Harish Tayyar Madabushi. 2023. Flesch or fumble? evaluating readability standard alignment of instruction-tuned language models. In *Proceedings of the Third Workshop on Natural Language Generation, Evaluation, and Metrics (GEM)*, pages 205–223, Singapore. Association for Computational Linguistics.

Alexander Novikov, Ngân Vũ, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zsolt Wagner, Sergey Shirobokov, Borislav Kozlovskii, Francisco JR Ruiz, Abbas Mehrabian, and 1 others. 2025. Alphaevolve: A coding agent for scientific and algorithmic discovery. *arXiv preprint arXiv:2506.13131*.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.

# A Details for 3 Generated Program

In general, all three programs are generated by GPT-4o model, which covers the following python libraries and models.

**Libraries** In the AI generated code of MRS-Rule and MRS-Joint, the following Python libraries are used: `argparse`, `os`, `sys`, `json`, `re`, `math`, `pathlib`, `collections`, `typing`, `numpy`, `requests`. Besides those regular pythong libraries, we noticed that libraries like `transformers` (Hugging Face), `SentenceTransformers`, `NumPy` are used for natural language processing and machine learning parts.

**Models** CEFR level is predicted by three ModernBERT from huggingface, ModernBERT-base-doc_en-Cefr, ModernBERT-base-doc_sent_en-Cefr, ModernBERT-base-reference_AllLang2-Cefr; while semantic similarity is using the library of `sentence transformers/all-MiniLM-L6-v2`), and LLM for generating the code is GPT-4o[4]. The LLM API used for text simplification is `gpt-4o-mini`.

## A.1 Program 1: Baseline Naïve Prompt

The generated Program 1 is in the file of "First_Version_Sentence_Simplification.py" in the code repo. It is built with **OpenAI's Chat Completions API**. The script is lightweight and designed for **large-scale, reproducible simplification** runs,

---

[4]https://chatgpt.com/?model=gpt-4o, accessible at 09/23/2025

while maintaining a clean JSONL output compatible with downstream CEFR or readability evaluations. This baseline Program 1 is used as Run 1 in our submission and is also used in our **CEFR-Gap** analysis.

## A.2 Program 2: MRS-Rule

The generated Program 2 is in the file "Second_Version_Sentence_Simplification.py". Specifically, ChatGPT suggests useful rules to generate candidate implication with a basic candidate `base_candidates()`. Generate multiple simplified variants of an input sentence using lightweight rule-based transformations without relying on an LLM. The details of the code are shown in the code listing 1. The corresponding hyperparameters used in the code are summarized in Table 4.

| Parameter | Value |
|---|---|
| `similarity_floor` | 0.88 |
| `max_steps` | 8 |
| `max_retries` | 6 |
| `floor_step` | 0.03 |
| `steps_step` | 6 |
| `sim_floor` (internal) | 0.88 ↓ |
| `w_hit` | 10 |
| `w_ref` | 2.5 |
| `w_orig` | 0.5 |
| `llm_timeout` | 60 seconds |
| `use_llm` | true |
| `sim_threshold` | 0.72–0.75 |

Table 4: Hyperparameters in MRS-Rule and MRS-Joint

### A.2.1 Generated Code to Apply Rules

Listing 1: Rule-based Simplification to Generate Candidates

```
base = text.strip()
lim = max(8, 28 - 2*step_idx)
cands = [
    replace_words(base),
    simplify_numbers_units(base),
    keep_shortest_clause(base),
    strip_relative_clauses(base),
    trim_to_limit(replace_words(base),
        lim),
    trim_to_limit(simplify_numbers_units
        (base), lim),
    trim_to_limit(strip_relative_clauses
        (base), lim),
    trim_to_limit(keep_shortest_clause(
        replace_words(base)), lim),
]
for s in sentence_split(
    strip_relative_clauses(base)):
    cands.append(trim_to_limit(
        replace_words(s), max(10, lim-4)
        ))
```

## A.3 Program 3: MRS-Joint

Program 3 is in the file "Third_Version_Sentence_Simplification.py". It simply combines the previous rule-based program A.2 with LLM-prompting-based program in A.1. This program shares the same rules used in A.2.1, and uses the same set of multi-round hyper parameters as Table 4.

## B Case Studies

To understand the model's behaviour beyond overall accuracy scores, we conducted a **qualitative error analysis** on three representative C1–C2 source texts that were simplified by the system but misclassified by the CEFR predictor. These examples illustrate three distinct types of misclassification.

### B.1 Case 1 – Overshoot: A2 Target → Predicted B1

**Source (C1–C2)**
*I am a member of the Campaign to Stop Killer Robots . . . We need to be very cautious indeed about what we ask machines to do.*
  **System Output (intended A2)**
*I am part of the Campaign to Stop Killer Robots. Forget the scary robots from movies. Right now, military machines like drones and robots are already used to kill with little help from people. The next step is robots that can decide who to kill by themselves. I believe this is very wrong and dangerous for people. We need to be very careful about what we let machines do.*
  **Analysis**
The model successfully shortened the complex original sentence that had many clauses and replaced high-level terms (*autonomous 'murderbots'*, *completely unethical*) with simpler words. However, it retained **abstract concepts** (*robots that can decide who to kill*) and an **embedded relative clause**, both typical of B1 syntax. The CEFR predictor therefore judged the output **one level higher than the A2 target**.
*Insight:* Even when sentences are lexically simplified, preserving key ideas may require structures that naturally push the text above the intended level.

### B.2 Case 2 – Lexical Imitation: A2 Target → Predicted B1

**Source (C1–C2)**
*A large number of bridge accidents occur during*

*the construction of the bridge itself . . . this can be deadly for the workers on it at the time.*
  **System Output (intended A2)**
*A large number of bridge accidents occur during the construction of the bridge itself. These accidents are often due to.*
  **Analysis**
Although truncated, the output **copied formal phrases from the source** (*a large number of bridge accidents occur*, *during the construction of the bridge itself*). It failed to rephrase these into simpler, high-frequency A2 forms such as *Many accidents with bridges happen while they are being built*. The CEFR model therefore still rated it **B1**, despite the missing ending.
*Insight:* True simplification requires **lexical adaptation**, not only shortening. Retaining formal academic expressions—even in a shorter text—can maintain a higher perceived level.

### B.3 Case 3 – Under-generation: B1 Target → Predicted A2

**Source (C1–C2)**
*Whether NASA can find the remaining middle-sized NEOs depends on getting the money to build NEO-Cam . . . the PHAs.*
  **System Output (intended B1)**
*NASA's ability to find the remaining middle-sized near-Earth objects (NEOs) depends on getting funding to build NEOCam, a 0.5-meter space telescope that.*
  **Analysis**
The system produced only a **partial sentence**, omitting the telescope's function and the entire discussion of planetary defence. With the **conceptual load drastically reduced**, the CEFR predictor assigned an **A2 level**, even though the fragment still contains technical terms (*NEOs*, *NEOCam*).
*Insight:* Incomplete outputs can appear easier than intended, causing the CEFR assessment to **underestimate** the level. Quality checks for completeness are essential alongside automatic scoring.