

No Reader Left Behind: Multi-Agent Summaries Everyone Can Understand

Jimin Jung¹, MyoungJin Kim¹, Jaehyung Seo^{2†}, Heuseok Lim^{1†}

¹Department of Computer Science and Engineering, Korea University

²Department of Computer Science and Engineering, Konkuk University

{stopmin02, imannamj, limhseok}@korea.ac.kr

seojae777@konkuk.ac.kr

Abstract

The Plain Writing Act in the United States requires government documents to be written in clear and simple language. However, existing summarization systems struggle to address diverse linguistic and cognitive barriers among general readers. We propose NRLB (No Reader Left Behind), a unified multi-agent framework for plain language summarization that simulates three representative reader groups: elementary school students, non-native speakers, and readers with attention deficits. NRLB integrates template-based planning with an iterative feedback loop guided by simulated readers and domain expert revision to address comprehension barriers such as unknown terms, missing contexts, and confusing sentences. Evaluations across multiple datasets demonstrate consistent improvements in both readability and factuality. Human evaluation further supports these findings, with annotator preference rates ranging from 55% to 76%, highlighting NRLB’s ability to generate summaries that are both faithful to the source and accessible to a wide range of readers.

1 Introduction

During the COVID-19 pandemic, gaps in scientific communication deepened public confusion. Lower science literacy was associated with greater susceptibility to misinformation and lower compliance, whereas even a 1% increase in public science literacy led to a 14.2% rise in civic engagement with public health measures (Qin et al., 2024). These findings suggest that even modest improvements in science literacy can directly influence public behavior, underscoring the importance of clear and accessible scientific communication for effective public health responses (Xu et al., 2015). The Plain Writing Act, enacted by the U.S. government in 2010, mandates that official communications

† Corresponding Authors

Plain Writing Act Guidelines: Before vs After Examples	
Example 1: Write in Short Sentences <u>Before</u> This regulation governs disaster assistance for services to prevent hardship caused by fire, flood, or acts of nature that are not provided by FEMA or the Red Cross.	<u>After</u> This regulation governs disaster assistance that: (a) Consists of services to prevent hardship caused by fire, flood, or acts of nature; and (b) Is furnished by a provider other than FEMA or the Red Cross.
Example 2: Use definite, concrete, everyday language <u>Before</u> Accompany, accomplish, accorded ...	<u>After</u> Go with, carry out, given...
Example 3: Add a definition of an important technical term <u>Before</u> You might be eligible if your income is 50% of your city/town's Area Median Income (AMI).	<u>After</u> You might be eligible if your income is 50% of your city/town's Area Median Income (AMI). AMI means the middle income in your area. Half of people in your area make more, half make less.

Figure 1: Examples of revision strategies based on the U.S. Plain Writing Act, including shortening sentences, using concrete language, and defining technical terms.

be understandable to the general public. These requirements reflect a long-standing institutional need for accessible public communication. Federal guidelines recommend a 6th-8th grade reading level, concise and active sentences, simplified or annotated technical terms, and accessibility across platforms. Agencies are also required to provide plain language resources and mechanisms for public feedback.¹ These guidelines emphasize the need to accommodate readers with limited literacy, non-native English proficiency, or cognitive and visual impairments (Gooding, 2022). These principles are typically implemented through strategies such as sentence simplification, vocabulary reduction, and defining specialized terms, as illustrated in Figure 1.

However, existing automated summarization and simplification systems often fail to adhere to the principles of the Plain Writing Act. They often overlook key reader-specific factors such as back-

¹<https://www.justice.gov/open/plain-writing-act>

ground knowledge, attention span, and syntactic complexity, resulting in outputs that retain technical jargon and complex sentence structures. While some models achieve surface-level simplification, they rarely improve comprehension for general audiences. Furthermore, many prior approaches generate multiple versions tailored to specific sub-groups (e.g., K-12 students or non-native adults), which conflicts with the Act’s core principle of clarity for all readers (Mo and Hu, 2024; Xu et al., 2015). These limitations highlight the need for a unified framework that supports broader accessibility through adaptive simplification.

To address this gap, we propose NRLB (No Reader Left Behind), an automated framework for plain language summarization that explicitly models commonly overlooked reading barriers. Grounded in prior research on literacy and cognitive processing (Smith et al., 2021; Tighe and Schatschneider, 2016; Shero et al., 2021), NRLB simulates three representative reader groups: elementary school students, non-native readers, and readers with attention deficits. Each agent identifies distinct comprehension challenges, enabling targeted revisions that improve accessibility across diverse reader populations.

The NRLB framework consists of two core modules. **Module 1: Content Planning and Drafting** begins with a Planner Agent that classifies the input document and selects an appropriate Domain Expert Agent, who generates a template-based initial summary capturing the core content. **Module 2: Feedback-Guided Simplification** refines this draft through a structured loop of reader feedback and expert revision. Reader Agents identify segments with lexical complexity, missing context, or syntactic difficulty. A Checklist Agent aggregates and prioritizes these issues, which are then used by the Domain Expert Agent to propose targeted edits. Finally, an Editor Agent applies these revisions in context to produce the final plain language summary, which is evaluated using metrics for readability, relevance, and factuality. This unified design enables systematic identification and resolution of diverse comprehension barriers across reader groups.

Our contributions are as follows:

(1) We propose NRLB, a unified plain language summarization framework that simulates diverse reader perspectives within a single system to explicitly model key comprehension barriers across reader groups.

(2) We introduce a structured multi-agent feedback loop with coordinated role-based agents that enables iterative identification and systematic resolution of comprehension barriers.

(3) We demonstrate that NRLB improves both readability and factuality through multi-round refinement across multiple benchmarks and human evaluations.

2 Related Work

2.1 Plain Language Summarization

Plain language summarization aims to simplify language while preserving key content and meaning. Prior work has focused on lexical or syntactic rewriting, often using Simple Wikipedia (Xu et al., 2015), readability-controlled fine-tuning (Luo et al., 2022), or task-specific evaluation metrics (Urlana et al., 2023). However, these approaches tend to emphasize surface-level readability, often at the expense of content fidelity.

In contrast, NRLB preserves the core content of documents through a template-based initial summary and maintains semantic consistency through a structured, controlled revision process. Rather than relying on simple text reduction, it addresses barriers such as unknown terms, missing context, and complex sentence structures, enabling both readability and factuality.

2.2 Multi-Agent Summarization

To address these limitations, recent work has explored multi-agent architectures where role-specific agents collaborate to overcome the limitations of single large language models (Du et al., 2024; Chen et al., 2023). In text simplification, **AgentSimp** (Fang et al., 2025) adopts a stage-wise pipeline with sequential agents (e.g., Metaphor Analysis, Article Logic Analyst), enabling structured multi-agent collaboration through intermediate analysis and rewriting. Although AgentSimp supports iterative reconstruction, it lacks an explicit feedback-driven refinement mechanism that systematically incorporates accumulated feedback across iterations, limiting its ability to progressively improve readability. Similarly, **ExpertEase** (Mo and Hu, 2024) employs teacher, student, and expert roles to generate grade-specific summaries tailored to different educational levels. However, its audience-segregated design diverges from the goal of the Plain Writing Act, which seeks universally understandable documents.

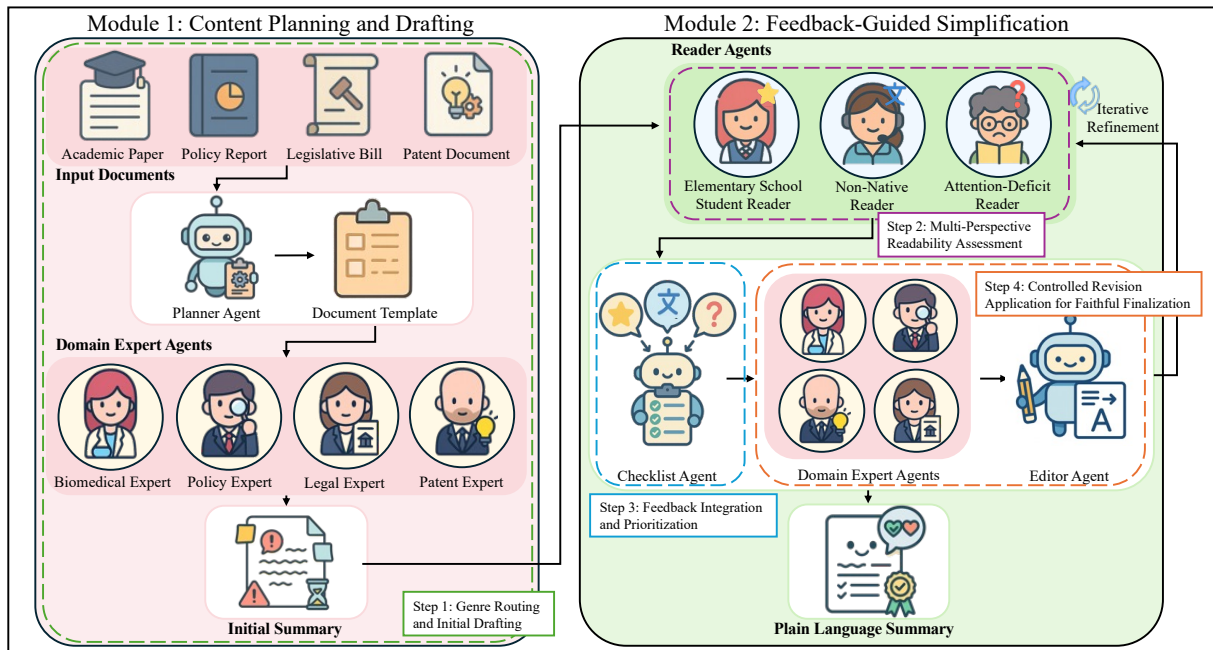


Figure 2: Overview of NRLB. Module 1 produces an initial draft via a Planner and Domain Expert Agent, while Module 2 iteratively refines it through a feedback loop with Reader, Checklist, Domain Expert, and Editor Agents.

In contrast, **NRLB** introduces a feedback-driven iterative refinement loop that accumulates and incorporates simulated reader feedback across steps, enabling systematic identification and resolution of comprehension barriers, resulting in more adaptive and inclusive plain language summarization.

2.3 Modeling Reader Diversity

Text simplification has traditionally aimed to support underserved readers, such as individuals with dyslexia, non-native speakers, and those with cognitive challenges. For example, (Rello et al., 2013) showed that replacing unfamiliar words with shorter and more frequent alternatives improved reading speed and comprehension, while (Wilkins and Todirascu, 2020) demonstrated the effectiveness of vocabulary-based simplification strategies for readers with cognitive challenges in multilingual settings. However, prior work has largely focused on lexical-level simplification and often treats different reader groups in isolation, failing to adequately address broader comprehension barriers such as limited background knowledge or syntactic complexity. Although these challenges have been noted in prior studies (Xu et al., 2015; Dadu et al., 2021), few approaches systematically address them within a unified framework.

In contrast, **NRLB** jointly models reader groups, including elementary school students, non-native readers, and readers with attention difficulties

within a unified framework. It simulates feedback from Reader Agents to systematically identify and resolve comprehension barriers. While collecting human feedback iteratively is costly and often impractical, recent work suggests that LLMs can approximate human reading behavior when guided by persona attributes (Argyle et al., 2023; Park et al., 2023; Aher et al., 2023), enabling feedback-driven iterative refinement. This approach enables unified modeling and systematic resolution of diverse comprehension barriers.

3 Methodology: NRLB

Figure 2 presents an overview of the **NRLB** (No Reader Left Behind) framework, which generates plain language summaries for general audiences through two sequential modules. This design enables structured coordination between planning, feedback, and revision, ensuring both content preservation and readability. **Module 1: Content Planning and Drafting** begins with a Planner Agent that classifies the genre of the input document. Based on this classification, the system selects an appropriate Domain Expert Agent and assigns a predefined template corresponding to the identified document type. The Domain Expert Agent then completes the template by following genre-specific guidelines, producing an initial summary that emphasizes informational completeness

and logical flow. **Module 2: Feedback-Guided Simplification** provides the initial summary to three simulated Reader Agents: Elementary School Student Reader, Non-Native Reader, and Attention-Deficit Reader. These agents independently identify comprehension barriers across three categories: unknown terms, missing contexts, and confusing sentences. The Checklist Agent integrates and prioritizes the feedback based on their frequency and estimated difficulty. The resulting checklist is returned to the Domain Expert Agent, who drafts specific revision suggestions. The Editor Agent then applies these suggestions to produce a revised summary. This process is repeated for up to N iterations, following the draft, review, and revise cycle outlined in the U.S. Plain Writing Act (2010). Additional implementation details are provided in Appendix A.

3.1 Genre Routing and Initial Drafting

The NRLB framework begins with a genre routing step that determines the document type and corresponding expert configuration. The Planner Agent classifies the input document into one of four types: Academic Paper, Policy Report, Legislative Bill, or Patent Document, and assigns the corresponding Domain Expert Agent (e.g., Biomedical, Policy, Legal, or Patent Expert). For example, a Policy Report is routed to a Policy Expert. Based on the genre label, the Planner Agent selects a slot-based summary template aligned with standard writing conventions for that document type. Each template consists of labeled sections with recommended sentence lengths. For instance, a policy report from the U.S. Government Accountability Office (GAO) typically includes “Why GAO Did This Study,” “What GAO Found,” and “What GAO Recommends.” A list of genre-specific templates is provided in Appendix A.1.

The Domain Expert Agent completes the template using example summaries and writing instructions (Appendix A.2). The agent aims to ensure informational completeness and logical flow while adhering to length constraints. In biomedical academic papers, the Methods section summarizes the experimental design, while the Results section highlights key findings. Although the initial summary captures the core content, it may still include barriers to comprehension such as unknown terms, missing contexts, or confusing sentences. These issues are addressed through iterative refinement in the next module.

3.2 Multi-Perspective Readability Assessment

To capture diverse comprehension barriers faced by general readers, we identified three representative reader profiles in Module 2. These profiles are grounded in prior research on literacy and cognitive processing and include the Elementary School Student Reader Agent, the Non-Native Reader Agent, and the Attention-Deficit Reader Agent, each representing a distinct reader population. All three agents evaluate the same summary according to the linguistic and cognitive characteristics of the group they simulate. They are instantiated using an LLM with prompts that reflect the reader’s abilities at three levels: word-level (e.g., limited vocabulary), knowledge-level (e.g., lack of background context), and sentence-level (e.g., difficulty parsing complex syntax). Each agent is instructed to detect comprehension barriers across all three predefined categories: unknown terms, missing contexts, and confusing sentences. This simulation-based approach enables consistent and reproducible feedback across diverse reader perspectives without requiring human annotators. Full prompt specifications are provided in Appendix A.3.

- **The Elementary School Student Reader Agent** simulates the perspective of a fourth-grade student with a vocabulary of approximately 3,000 common English words (Laufer and Ravenhorst-Kalovski, 2010). It flags unfamiliar or technical terms beyond this range as unknown terms. Institutional references or scientific concepts that require background knowledge are classified as missing contexts. Sentences longer than 15 words or containing multiple subordinate clauses are marked as confusing sentences. Typical feedback includes phrases such as “too many difficult words,” “background explanation is needed,” and “the sentence is too long.”
- **The Non-Native Reader Agent** represents an adult reader familiar with high-frequency vocabulary but sensitive to low-frequency academic terms, idioms, and cultural expressions (Ha, 2022). It flags rare or idiomatic phrases as unknown terms and identifies culturally specific institutions or historical references lacking explanation as missing contexts. Sentences with garden-path constructions, reduced relative clauses, or repeated passive voice are labeled as confusing sentences.

Common feedback includes statements like “I don’t know this word,” “the cultural context is missing,” and “the sentence structure is too complicated.”

- **The Attention-Deficit Reader Agent** reflects the cognitive profile of individuals with limited working memory and reduced attention span (Jacobson et al., 2011). It marks visually or phonetically complex words, compounds, and low-frequency vocabulary as unknown terms. Concepts introduced without sufficient background are flagged as missing contexts. Sentences exceeding 15 words or containing multiple relative clauses are categorized as confusing sentences. Representative feedback includes comments such as “the word is hard to decode” and “the long sentence makes it hard to stay focused.”

Each Reader Agent provides structured feedback, which is integrated and prioritized in the next stage. This consolidated view reveals comprehension barriers that general readers may face. Although we focus on three representative reader types, NRLB is designed to easily plug in additional Reader Agents, allowing flexible extension to additional reader profiles.

3.3 Feedback Integration and Prioritization

The Checklist Agent collects structured feedback from the Reader Agents, where issues are categorized as unknown terms, missing contexts, or confusing sentences. The priority of each issue is determined by the number of agents that flagged it. Issues flagged by all three agents receive top priority. When multiple issues have the same priority, the agent applies the Automated Readability Index (ARI) to assess sentence complexity (Senter and Smith, 1967). The ARI estimates U.S. grade-level difficulty by combining the average characters per word and the average number of words per sentence. For each category, the Checklist Agent selects up to three items with the highest difficulty and compiles them into a unified checklist. This refined list is forwarded to the Domain Expert Agent to guide targeted revisions.

3.4 Controlled Revision Application for Faithful Finalization

The Domain Expert Agent generates revision proposals based on the checklist received from the Checklist Agent. Using the same domain expert as

in Module 1 ensures continuity between drafting and revision. It replaces unknown terms with simple synonyms or 5-word definitions, adds missing contexts with up to 15-word explanations from the source, rewrites confusing sentences into shorter active sentences, and marks items as “insufficient information” if no source support exists.

The Editor Agent applies these revisions to produce an updated summary while maintaining sentence-level coherence and contextual flow. When edits involve splitting or replacing parts of the same sentence, the agent determines a conflict-free application order. Revisions are only applied when the original sentence exists in the current summary, and replacements are executed sequentially to preserve consistency. The Editor Agent integrates reader feedback with expert-suggested revisions to enable coherent multi-perspective refinement. This revision process can be repeated up to N times; in our experiments, we used two iterations by default. This process improves readability while preserving factual consistency through controlled integration of revisions. As shown in Section 5.2, removing the Editor Agent degrades the factual consistency of the final summary.

4 Experimental Setup

4.1 Datasets

We evaluate NRLB on four representative datasets across diverse domains. PLOS (Goldsack et al., 2022) consists of biomedical academic papers, GovReport (Huang et al., 2021) consists of policy reports from U.S. government agencies, BillSum (Kornilova and Eidelman, 2019) consists of U.S. legislative bills, and BigPatent (Sharma et al., 2019) consists of patent documents with rich technical language. We randomly sample 500 test examples from each dataset to ensure both statistical robustness and computational efficiency, following recent best practices in long-document summarization research (Xu et al., 2025). All datasets include human-written reference summaries, and further dataset-specific details are provided in Appendix C.1.

4.2 Models

Without additional fine-tuning, we evaluate NRLB using three large language models with distinct characteristics. Llama-3.1-8B-Instruct, an open-source model licensed under Apache 2.0, serves as a reproducible open-source baseline (Dubey et al.,

		PLOS			GovReport			BillSum			BigPatent		
		GPT-4o	Llama-3.1	Qwen3	GPT-4o	Llama-3.1	Qwen3	GPT-4o	Llama-3.1	Qwen3	GPT-4o	Llama-3.1	Qwen3
ROUGE-1	Initial Summary	46.85	47.27	41.62	37.93	46.93	46.19	46.70	45.81	41.08	41.87	44.94	34.21
	Round 1	45.12	44.43	37.55	36.35	45.33	41.41	43.28	42.66	37.08	37.74	40.09	30.55
	Round 2	42.77	41.64	36.03	34.73	44.19	40.04	40.88	39.70	34.62	35.40	36.87	29.18
BERTScore	Initial Summary	86.91	86.51	85.78	86.14	86.30	85.59	87.07	86.83	85.61	86.52	86.68	85.55
	Round 1	86.61	86.00	85.33	85.89	85.84	85.04	86.43	86.20	85.03	85.82	85.73	84.89
	Round 2	86.29	85.43	84.96	85.63	85.48	84.74	85.99	85.59	84.62	85.44	85.11	84.50
FKGL ↓	Initial Summary	17.73	18.74	19.12	17.79	19.95	19.81	17.29	18.04	18.07	18.60	21.13	20.88
	Round 1	13.37	16.55	14.12	13.92	17.51	15.51	12.64	15.66	13.08	12.57	16.88	14.48
	Round 2	10.94	15.21	12.46	11.58	15.98	13.72	10.39	14.60	11.16	10.27	14.62	12.38
DCRS ↓	Initial Summary	14.25	13.69	15.04	13.33	13.43	14.24	13.11	12.77	14.14	13.76	13.36	15.10
	Round 1	13.21	12.67	13.61	12.49	12.51	13.31	12.06	11.88	12.68	12.31	12.14	13.32
	Round 2	12.56	12.06	13.15	11.96	11.96	12.90	11.50	11.40	12.24	11.70	11.47	12.74
CLI ↓	Initial Summary	18.72	17.05	20.86	17.77	16.60	19.38	16.56	15.13	17.38	18.32	16.45	21.52
	Round 1	15.86	14.54	16.54	15.71	15.09	17.06	14.07	13.59	14.22	14.69	13.76	16.29
	Round 2	13.96	13.20	15.14	14.20	14.23	16.01	12.56	12.62	13.09	12.84	12.42	14.71
LENS	Initial Summary	65.38	64.80	62.49	60.28	55.37	56.04	60.91	58.38	56.29	59.62	55.28	55.86
	Round 1	71.35	67.21	71.62	66.32	59.20	63.37	68.20	62.58	66.68	70.51	61.08	69.34
	Round 2	74.67	66.68	72.60	70.00	60.60	64.61	71.63	62.68	68.30	74.44	62.62	71.56
SummaC	Initial Summary	55.48	49.67	44.48	42.36	43.63	38.59	34.41	33.48	30.69	53.19	54.73	42.87
	Round 1	67.12	49.17	48.88	53.14	45.70	45.38	39.41	34.51	36.02	69.53	57.95	49.87
	Round 2	72.93	50.86	53.81	61.27	47.76	49.00	41.56	34.55	38.50	73.49	58.94	54.58

Table 1: Performance of three models (GPT-4o, Llama-3.1-8B-Instruct, Qwen3-8B) on four benchmarks across initial summary and two refinement rounds; best scores in bold. ↓ indicates that lower is better in readability metrics.

2024). GPT-4o represents OpenAI’s recent commercial model with multimodal capabilities (Hurst et al., 2024). Qwen3-8B in reasoning mode is included for its strong performance on complex reasoning tasks (Yang et al., 2025). Full configuration details are provided in Appendix C.2.

4.3 Metrics

We evaluate NRLB across three core dimensions as defined in the Plain Language Summarization Shared Task: Relevance, Readability, and Factuality (Goldsack et al., 2024). Relevance measures how well the summary preserves key content, Readability reflects how easily the summary can be understood by readers, and Factuality evaluates the consistency of the summary with the source. From the available metrics, we select the seven most aligned with our task: ROUGE-1, BERTScore, FKGL, DCRS, CLI, LENS, and SummaC.

Relevance is measured using ROUGE-1, which captures unigram overlap between summaries and references (Liu and Liu, 2008), and BERTScore, which quantifies semantic similarity using contextual embeddings (Zhang et al., 2019). Readability is assessed through traditional formula-based metrics (FKGL, DCRS, and CLI) that reflect sentence and word complexity (Tanprasert and Kauchak, 2021; Dale and Chall, 1948; Coleman and Liau, 1975), along with LENS, a learned metric that better captures qualitative improvements from lexical

and structural simplification (Maddela et al., 2022). Factuality is evaluated using SummaC, which applies natural language inference to assess consistency between the source document and summary (Laban et al., 2022). Detailed descriptions are provided in Appendix C.3.

5 Results

5.1 Improving Readability via Iterative Revision

Table 1 summarizes the effects of two rounds of revision (Rounds 1 and 2) using NRLB across four datasets and three models. Iterative refinement consistently improved readability across all datasets. In contrast, ROUGE-1 and BERTScore declined slightly due to reduced lexical overlap, reflecting the well-known trade-off in plain language summarization where clarity often reduces surface similarity (Dreyer et al., 2021; Goldsack et al., 2024). Meanwhile, factuality improved across all settings, with SummaC scores increasing from 55.48 to 72.93 on PLOS and from 53.19 to 73.49 on BigPatent, confirming that iterative revision enhanced both logical consistency and source alignment.

Across models, GPT-4o showed the best overall balance, maintaining strong initial relevance and achieving the highest readability and factuality by Round 2. Llama-3.1-8B-Instruct retained high ROUGE but made modest gains elsewhere. Qwen3-

Method	PLOS				GovReport				BillSum				BigPatent			
	FKGL ↓	DCRS ↓	CLI ↓	LENS ↓	FKGL ↓	DCRS ↓	CLI ↓	LENS ↓	FKGL ↓	DCRS ↓	CLI ↓	LENS ↓	FKGL ↓	DCRS ↓	CLI ↓	LENS ↓
All	15.40	12.06	13.19	67.64	16.09	12.01	14.29	59.08	15.00	11.68	13.09	62.63	14.68	11.51	12.93	64.40
w/o ele	15.73	12.19	13.45	66.18	16.56	12.05	14.68	59.94	14.48	11.35	12.52	62.26	15.23	11.67	12.79	61.90
w/o non	15.20	12.01	13.10	66.36	16.19	12.01	14.41	60.12	14.52	11.47	12.60	62.10	16.83	11.91	13.31	59.82
w/o att	15.40	12.06	13.21	67.31	16.34	12.06	14.51	59.27	14.49	11.42	12.58	63.01	15.57	11.74	13.00	62.34
w/o ele,non	16.01	12.28	13.68	64.75	17.03	12.05	14.88	58.26	14.61	11.39	12.58	62.77	16.64	11.90	13.31	58.71
w/o ele,att	15.74	12.24	13.47	64.84	16.83	12.13	14.81	58.20	14.73	11.43	12.56	62.09	17.79	12.08	13.73	57.91
w/o non,att	15.27	12.03	13.09	66.84	16.23	12.07	14.42	58.74	14.72	11.49	12.72	62.09	16.15	11.88	13.26	61.71

Table 2: Ablation study on the effect of Reader Agent combinations across four datasets (200 samples each). “ele” = Elementary School Student Reader, “non” = Non-Native Reader, “att” = Attention-Deficit Reader.

Dataset	Setting	Initial	Round 1	Round 2
PLOS	w/ Editor	49.6662	49.1664 (-0.50)	52.3063 ($+3.14$)
	w/o Editor	49.9958	48.8337 (-1.16)	49.8839 ($+1.05$)
GovReport	w/ Editor	43.6340	45.7042 ($+2.07$)	49.9039 ($+4.20$)
	w/o Editor	43.6733	43.0163 (-0.66)	43.6692 ($+0.65$)
BillSum	w/ Editor	33.4814	34.5100 ($+1.03$)	35.0046 ($+0.49$)
	w/o Editor	33.2260	33.4913 ($+0.27$)	33.2915 (-0.20)
BigPatent	w/ Editor	54.7311	57.9544 ($+3.22$)	58.9384 ($+0.98$)
	w/o Editor	55.2051	55.0744 (-0.13)	55.8485 ($+0.77$)

Table 3: SummaC scores across datasets and revision rounds (Initial denotes the initial summary) with and without the Editor (Editor Agent). Positive changes are shown in bold.

8B, though initially less consistent, achieved the largest improvements after revision, likely due to its reasoning-oriented architecture. While Round 3 results (Appendix D.1) brought minor further improvements in readability and factuality, relevance dropped by up to 10 points compared to the initial summary, supporting Round 2 as the default stopping point. This behavior reflects a relevance-readability trade-off: aggressive simplification in later rounds tends to disrupt semantic flow and reduce lexical alignment with the source. Based on this empirical observation, we adopt Round 2 as the default setting, balancing accessibility with content preservation. Efficiency analysis, including latency, API calls, and average cost, is provided in Appendix D.2.

5.2 Ablation Study

Effect of Reader Agent Configurations We evaluated the effect of using one, two, and three Reader Agents. As shown in Table 2, all configurations were evaluated under the same two-round revision pipeline using four readability metrics. Using two Reader Agents generally improved readability across datasets and metrics, resulting in improved clarity and accessibility. For instance, on GovReport, FKGL improved from 17.03 to 16.09, and LENS rose from 58.20 to 60.12.

With three Reader Agents, FKGL improved further and LENS reached its highest scores on PLOS and BigPatent, indicating the strongest overall improvements. In contrast, BillSum achieved near-optimal readability with two Reader Agents, and the third added only marginal gains. These findings suggest that using at least two Reader Agents is generally beneficial for readability, with the three-agent setup performing best overall. We therefore adopt the three-agent configuration as the default. Additional results with four and five agents are reported in Appendix B.1, and further analysis on feedback overlap and the complementary roles of Reader Agents is provided in Appendix B.2.

Effect of the Editor Agent As shown in Table 3, the inclusion of the Editor Agent consistently improves SummaC, a factuality metric, at each revision round, yielding final values well above those without the Editor Agent. In contrast, without the Editor Agent, the scores show inconsistent changes in Round 1, with both increases and decreases observed, and this instability persists in Round 2. These results indicate that the Editor Agent is essential for merging feedback from Reader Agents and the Domain Expert coherently, producing steady gains in factuality and improving the overall reliability of the final summary. In particular, the Editor Agent plays a central role in transforming multi-agent feedback into actionable revisions, thereby improving overall performance.

5.3 Effect of Checklist Size

We study the impact of the number of feedback items (K) on summary quality, comparing $K = 3, 6,$ and 9 . Results are shown in Figure 3. To analyze trade-offs, we report three normalized metrics for each dataset: Relevance (average of ROUGE-1 and BERTScore), Surface Readability (inverse average of FKGL, DCRS, and CLI), and Semantic Readability (LENS).

As K increases, surface readability improves, re-

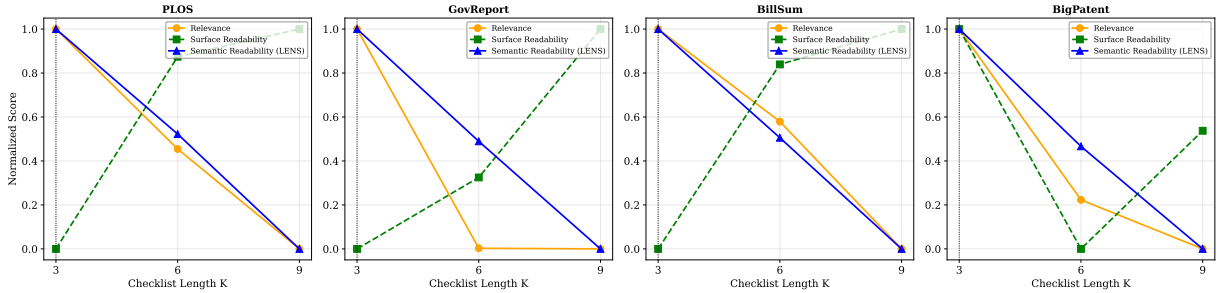


Figure 3: Impact of checklist length on summary quality; $K = 3$ achieves the best trade-off between relevance and readability.

Dataset	Method	Relevance		Readability			Factuality	
		ROUGE-1	BERTScore	FKGL ↓	DCRS ↓	CLI ↓	LENS	SummaC
PLOS (500 samples)	AgentSimp-Synchronous	44.59	86.23	22.20	13.85	16.45	60.08	45.44
	AgentSimp-Pipeline	44.92	86.31	22.33	13.84	16.37	60.79	46.07
	<i>NRLB (Ours)</i>	41.64	85.43	15.21	12.06	13.20	66.68	50.86
GovReport (500 samples)	AgentSimp-Synchronous	43.62	85.55	27.41	14.44	17.14	50.24	45.14
	AgentSimp-Pipeline	43.39	85.56	26.31	14.28	17.08	50.62	45.63
	<i>NRLB (Ours)</i>	45.33	85.84	17.51	12.51	15.09	59.20	45.70
BillsSum (500 samples)	AgentSimp-Synchronous	37.82	85.31	26.36	12.99	14.37	56.47	32.86
	AgentSimp-Pipeline	34.11	84.83	50.08	15.95	14.69	46.77	30.67
	<i>NRLB (Ours)</i>	39.70	85.59	14.60	11.40	12.62	62.68	34.55
BigPatent (500 samples)	AgentSimp-Synchronous	37.38	85.44	23.30	13.10	14.28	53.47	48.27
	AgentSimp-Pipeline	37.45	85.46	23.96	13.14	14.14	53.08	48.24
	<i>NRLB (Ours)</i>	36.87	85.11	14.62	11.47	12.42	62.62	58.94

Table 4: Automatic evaluation of NRLB and AgentSimp variants on four benchmarks (500 samples each). Best scores are shown in bold.

flecting simpler syntax and shorter sentences. However, both Relevance and LENS decline, indicating loss of key content and degraded semantic coherence. These results suggest that larger checklists tend to over-simplify summaries and reduce informativeness. We therefore adopt $K = 3$ as the default setting, which provides the best overall balance. For completeness, results for $K = 1$ and $K = 2$ are reported in Table 12, both yielding lower readability and factuality than $K = 3$.

5.4 Comparison with Baselines

We compare NRLB’s adaptive feedback loop with two communication strategies from AgentSimp (Fang et al., 2025): synchronous and pipeline-style settings. For a fair comparison, all systems were evaluated under a single-pass configuration (one revision round). NRLB achieved the strongest overall readability performance across all four datasets, improving FKGL, DCRS, CLI, and LENS scores. NRLB also achieved the highest factuality scores, reflecting stronger coherence and alignment with the source. In terms of relevance, NRLB outperformed AgentSimp on GovReport and BillsSum, but

showed slight decreases on PLOS and BigPatent, suggesting a trade-off between readability and content preservation. Overall, NRLB achieved the most balanced improvements even under a single-pass setup.

We additionally examined a direct prompting baseline using the same underlying LLM, but without agent collaboration or iterative refinement. In this setting, the model performs single-pass simplification based on a simple instruction prompt. This approach aggressively removes information during simplification, improving readability metrics (FKGL, DCRS, and CLI) but consistently reducing relevance (ROUGE-1 and BERTScore). These results suggest that excessive simplification can enhance surface-level readability at the cost of losing essential content. As a consequence, the generated summaries often exhibit limited informational coverage. As a result, the direct baseline exhibits behavior that differs from our goal of improving readability while preserving key information. We therefore exclude it from the main comparison and report its implementation details and additional results in Appendix C.4.

Group	Dataset	NRLB			AgentSimp			Pref. (%)
		Coh.	Simpl.	Faith.	Coh.	Simpl.	Faith.	
Non-Native	PLOS	4.10	3.85	3.85	3.77	2.61	3.78	75.0
	GovR.	4.27	3.73	4.12	3.72	3.35	3.62	65.0
	BillSum	4.30	3.88	3.97	3.37	2.97	3.40	71.7
	Patent	4.27	4.18	4.12	3.40	2.35	3.32	76.7
	Average	4.23	3.91	4.00	3.56	2.82	3.53	72.1
Elementary School Student	PLOS	3.17	2.83	3.33	3.33	2.83	2.83	66.7
	GovR.	3.89	3.56	3.33	3.66	3.11	3.00	55.7
	BillSum	3.83	2.83	3.50	3.17	3.17	3.17	67.0
	Patent	3.78	3.33	3.89	3.44	3.00	3.33	56.0
	Average	3.70	3.20	3.53	3.43	3.03	3.10	61.5

Table 5: Human evaluation results (5-point Likert scale) for NRLB vs. AgentSimp across two groups (Non-Native and Elementary School Student). “GovR.” denotes GovReport; bold indicates the higher score.

5.5 Human Evaluation and Qualitative Analysis

We conduct a human evaluation with two participant groups to compare NRLB with AgentSimp (see Table 5). First, three undergraduate non-native English speakers with intermediate English proficiency evaluated 80 documents. Each annotator rated the summaries on coherence, simplicity, and faithfulness using a five-point scale and selected the version they found easier to understand. NRLB achieved consistently higher scores across all metrics, with average improvements of 0.49-1.09 points. In preference judgments, 65-77% of cases (72.1% on average) favored NRLB, confirming its benefits for readability and user satisfaction.

We also conducted an evaluation with three elementary school students with limited English proficiency on 10 documents. NRLB again outperformed AgentSimp across most metrics, showing particularly clear improvements in coherence and faithfulness for GovReport, BillSum, and BigPatent. However, AgentSimp slightly outperformed NRLB on coherence in PLOS and on simplicity in BillSum. Overall, NRLB was preferred in 61.5% of cases, demonstrating its effectiveness in reducing reading difficulty and preserving essential content across diverse user groups. Further details of the human evaluation, including inter-annotator agreement analysis, are provided in Appendix E. Evaluation for the Attention-Deficit Reader group was not conducted due to ethical and clinical constraints, and the elementary-student evaluation was limited in scale to comply with IRB requirements.

Figure 4 presents a representative case study that exemplifies these improvements. Additional case studies and failure case analyses are provided in Appendix F and Appendix G.

<p>Initial Summary: Using RNA-Seq analysis in HeLa cells post-UIC knockdown, alongside mutational minigene analysis, in vivo splice-site blocking with antisense morpholinos, and in vitro binding experiments, we explored the mechanistic basis of UIC-dependent splicing regulation. The results reveal a conserved intra-U1 snRNP cross-regulation mechanism that links UIC and U1-70K expression through alternative splicing. Specifically, UIC levels dictate the balance between productive and non-productive U1-70K mRNA isoforms, ...</p>
<p>Elementary School Student Reader Agent: “unknown_terms”: [“non-productive U1-70K mRNA isoforms”, ...], “missing_contexts”: ..., “confusing_sentences”: ...</p>
<p>Non-Native Reader Agent: “unknown_terms”: [“dictate”, ...], “missing_contexts”: [Using RNA-Seq analysis in HeLa cells...], “confusing_sentences” ...</p>
<p>Attention-Deficit Reader Agent: “unknown_terms”: ..., “missing_contexts”: ..., “confusing_sentences” : [Using RNA-Seq analysis in HeLa cells...]</p>
<p>Biomedical Scientist Agent: “unknown_terms”: {“before”: “non-productive U1-70K mRNA isoforms”, “after”: non-useful U1-70K mRNA forms”}, {“before”: “dictate”, “after”: “control”}, “missing_contexts”: ...</p>
<p>Round-1 Summary: We used RNA-Seq analysis in HeLa cells after reducing UIC. We also used mutational minigene analysis, splice-site blocking with antisense morpholinos, and binding experiments to study UIC-dependent splicing regulation. HeLa cells are a type of human cell used in research. The results reveal a conserved intra-U1 snRNP cross-regulation mechanism that links UIC and U1-70K expression through alternative splicing. UIC levels control the balance between useful and non-useful U1-70K mRNA forms.</p>

Figure 4: Example case study shows how reader feedback drives lexical, contextual, and structural simplifications in NRLB.

6 Conclusion

This paper presents NRLB (No Reader Left Behind), a multi-agent framework for generating plain language summaries by simulating diverse reader perspectives. The system combines genre-specific content planning with iterative feedback and revision, guided by Reader, Checklist, Domain Expert, and Editor Agents. Across four datasets, NRLB consistently improved readability while improving factual consistency and preserving key content through collaborative refinement. Human evaluations confirmed that NRLB produces summaries that are easier to understand, more coherent, and better suited for diverse reader groups. These results demonstrate the practical applicability of NRLB as a viable approach to implementing plain language policies such as the U.S. Plain Writing Act. Future work includes extending NRLB to broader domains and incorporating more diverse reader profiles.

Limitations

While the current implementation of the NRLB pipeline demonstrates the feasibility of multi-agent readability-guided summarization, several limitations remain.

First, readability feedback relies on large language models that simulate elementary school students, non-native readers, and readers with attention deficits. Although the prompts are designed to elicit diverse comprehension challenges, the responses are not grounded in real-world reader data and may include hallucinated or redundant feedback, particularly in edge cases.

Second, human evaluation is limited in scale and coverage due to ethical and practical constraints. The evaluation with elementary school students was restricted to a subset of documents to comply with IRB requirements, while evaluation for the attention-deficit reader group was not conducted due to the need for clinical validation and handling sensitive personal health information, limiting generalizability and requiring large-scale evaluation.

Third, iterative refinement may introduce a trade-off between readability and content preservation. As the number of revision rounds or feedback items increases, the system may over-simplify the text, leading to loss of important details or degradation of semantic coherence.

Acknowledgments

This work was supported by the Commercialization Promotion Agency for R&D Outcomes(COMPA) grant funded by the Korea government(Ministry of Science and ICT)(2710086166), Institute for Information & communications Technology Promotion(IITP) grant funded by the Korea government(MSIT) (RS-2024-00398115, Research on the reliability and coherence of outcomes produced by Generative AI), Institute for Information & communications Technology Planning & Evaluation(IITP) grant funded by the Korea government(MSIT) (No. RS-2022-II220369, (Part 4) Development of AI Technology to support Expert Decision-making that can Explain the Reasons/Grounds for Judgment Results based on Expert Knowledge), and Institute of Information & communications Technology Planning & Evaluation (IITP) under the artificial intelligence star fellowship support program to nurture the best talents (IITP-2026-RS-2025-02304828) grant funded by the Korea government(MSIT).

References

- Gati V Aher, Rosa I Arriaga, and Adam Tauman Kalai. 2023. Using large language models to simulate multiple humans and replicate human subject studies. In *International conference on machine learning*, pages 337–371. PMLR.
- Lisa P Argyle, Ethan C Busby, Nancy Fulda, Joshua R Gubler, Christopher Rytting, and David Wingate. 2023. Out of one, many: Using language models to simulate human samples. *Political Analysis*, 31(3):337–351.
- Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chen Qian, Chi-Min Chan, Yujia Qin, Yaxi Lu, Ruobing Xie, and 1 others. 2023. Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors in agents. *arXiv preprint arXiv:2308.10848*, 2(4):6.
- Meri Coleman and Ta Lin Liau. 1975. A computer readability formula designed for machine scoring. *Journal of Applied Psychology*, 60(2):283.
- Tanvi Dadu, Kartikey Pant, Seema Nagar, Ferdous Ahmed Barbhuiya, and Kuntal Dey. 2021. Text simplification for comprehension-based question-answering. *arXiv preprint arXiv:2109.13984*.
- Edgar Dale and Jeanne S Chall. 1948. A formula for predicting readability: Instructions. *Educational research bulletin*, pages 37–54.
- Markus Dreyer, Mengwen Liu, Feng Nan, Sandeep Atluri, and Sujith Ravi. 2021. Evaluating the tradeoff between abstractiveness and factuality in abstractive summarization. *arXiv preprint arXiv:2108.02859*.
- Zhuoyun Du, Chen Qian, Wei Liu, Zihao Xie, Yifei Wang, Yufan Dang, Weize Chen, and Cheng Yang. 2024. Multi-agent software development through cross-team collaboration. *arXiv preprint arXiv:2406.08979*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. *arXiv e-prints*, pages arXiv–2407.
- Alexander R Fabbri, Wojciech Kryściński, Bryan McCann, Caiming Xiong, Richard Socher, and Dragomir Radev. 2021. Summeval: Re-evaluating summarization evaluation. *Transactions of the Association for Computational Linguistics*, 9:391–409.
- Dengzhao Fang, Jipeng Qiang, Xiaoye Ouyang, Yi Zhu, Yunhao Yuan, and Yun Li. 2025. Collaborative document simplification using multi-agent systems. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 897–912.
- Tomas Goldsack, Carolina Scarton, Matthew Shardlow, and Chenghua Lin. 2024. Overview of the biolaysumm 2024 shared task on the lay summarization of biomedical research articles. *arXiv preprint arXiv:2408.08566*.

- Tomas Goldsack, Zhihao Zhang, Chenghua Lin, and Carolina Scarton. 2022. Making science simple: Corpora for the lay summarisation of scientific literature. *arXiv preprint arXiv:2210.09932*.
- Sian Gooding. 2022. On the ethical considerations of text simplification. *arXiv preprint arXiv:2204.09565*.
- Hung Tan Ha. 2022. Vocabulary demands of informal spoken english revisited: what does it take to understand movies, tv programs, and soap operas? *Frontiers in Psychology*, 13:831684.
- Luyang Huang, Shuyang Cao, Nikolaus Parulian, Heng Ji, and Lu Wang. 2021. Efficient attentions for long document summarization. *arXiv preprint arXiv:2104.02112*.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and 1 others. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.
- Lisa A Jacobson, Matthew Ryan, Rebecca B Martin, Joshua Ewen, Stewart H Mostofsky, Martha B Denckla, and E Mark Mahone. 2011. Working memory influences processing speed and reading fluency in adhd. *Child neuropsychology*, 17(3):209–224.
- Anastassia Kornilova and Vlad Eidelman. 2019. Billsum: A corpus for automatic summarization of us legislation. *arXiv preprint arXiv:1910.00523*.
- Philippe Laban, Tobias Schnabel, Paul N Bennett, and Marti A Hearst. 2022. Summac: Re-visiting nli-based models for inconsistency detection in summarization. *Transactions of the Association for Computational Linguistics*, 10:163–177.
- Batia Laufer and Geke C. Ravenhorst-Kalovski. 2010. Lexical threshold revisited: Lexical text coverage, learners’ vocabulary size and reading comprehension. *Reading in a Foreign Language*, 22(1):15–30.
- Feifan Liu and Yang Liu. 2008. Correlation between rouge and human evaluation of extractive meeting summaries. In *Proceedings of ACL-08: HLT, short papers*, pages 201–204.
- Zheheng Luo, Qianqian Xie, and Sophia Ananiadou. 2022. Readability controllable biomedical document summarization. *arXiv preprint arXiv:2210.04705*.
- Mounica Maddela, Yao Dou, David Heineman, and Wei Xu. 2022. Lens: A learnable evaluation metric for text simplification. *arXiv preprint arXiv:2212.09739*.
- Joy Mahapatra and Utpal Garain. 2024. An extensive evaluation of factual consistency in large language models for data-to-text generation. *arXiv preprint arXiv:2411.19203*.
- Kaijie Mo and Renfen Hu. 2024. Expertease: A multi-agent framework for grade-specific document simplification with large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 9080–9099.
- Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual acm symposium on user interface software and technology*, pages 1–22.
- Haibo Qin, Zhongxuan Xie, Huping Shang, Yong Sun, Xiaohui Yang, and Mengming Li. 2024. The mass public’s science literacy and co-production during the covid-19 pandemic: empirical evidence from 140 cities in china. *Humanities and Social Sciences Communications*, 11(1):1–13.
- Luz Rello, Ricardo Baeza-Yates, Laura Dempere-Marco, and Horacio Saggion. 2013. Frequent words improve readability and short words improve understandability for people with dyslexia. In *IFIP Conference on Human-Computer Interaction*, pages 203–219. Springer.
- J. Senter, R. and A. Smith, E. 1967. Automated readability index. Technical Report AMRL-TR-66-22, Aerospace Medical Division, Aerospace Medical Research Laboratories, Wright-Patterson Air Force Base, Ohio. U.S. Air Force Systems Command, Project 1710 (Human Factors in the Design of Training Systems), Task 171007.
- Eva Sharma, Chen Li, and Lu Wang. 2019. Bigpatent: A large-scale dataset for abstractive and coherent summarization. *arXiv preprint arXiv:1906.03741*.
- Jeffrey A Shero, Jessica AR Logan, Stephen A Petrill, Erik Willcutt, and Sara A Hart. 2021. The differential relations between adhd and reading comprehension: A quantile regression and quantile genetic approach. *Behavior genetics*, 51(6):631–653.
- Reid Smith, Pamela Snow, Tanya Serry, and Lorraine Hammond. 2021. The role of background knowledge in reading comprehension: A critical review. *Reading Psychology*, 42(3):214–240.
- Teerapaun Tanprasert and David Kauchak. 2021. Flesch-kincaid is not a text simplification evaluation metric. In *Proceedings of the 1st Workshop on Natural Language Generation, Evaluation, and Metrics (GEM 2021)*, pages 1–14.
- Elizabeth L Tighe and Christopher Schatschneider. 2016. Examining the relationships of component reading skills to reading comprehension in struggling adult readers: A meta-analysis. *Journal of learning disabilities*, 49(4):395–409.
- Ashok Urlana, Pruthwik Mishra, Tathagato Roy, and Rahul Mishra. 2023. Controllable text summarization: Unraveling challenges, approaches, and prospects—a survey. *arXiv preprint arXiv:2311.09212*.
- Rodrigo Wilkens and Amalia Todirascu. 2020. Simplifying coreference chains for dyslexic children. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 1142–1151.

- Borui Xu, Yao Chen, Zeyi Wen, Weiguo Liu, and Bingsheng He. 2025. Evaluating small language models for news summarization: Implications and factors influencing performance. *arXiv preprint arXiv:2502.00641*.
- Wei Xu, Chris Callison-Burch, and Courtney Napoles. 2015. Problems in current text simplification research: New data can help. *Transactions of the Association for Computational Linguistics*, 3:283–297.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.
- 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 Implementation Details

Our framework is built entirely on large language models (LLMs), enabling a training-free and easily deployable approach. The implementation centers on two key aspects. First is the design of the overall process, which generates an initial summary based on document type and agent roles, followed by iterative revision through multi-agent interactions. Second is the construction of optimized prompting templates tailored to each role (e.g., Domain Expert, Editor, Reader Agents). The workflow is outlined in Algorithm 1.

Algorithm 1 Multi-Agent Summarization with Multi-Round Readability Refinement

Require: Source documents; prompts for Planner, Domain Expert, Readers, and Editor Agent

Ensure: Initial and N refined summaries are written to file

```
1: Load source documents (JSON format) into samples
2: for all sample  $\in$  samples do
3:   Extract source_text
4:   doc_type, expert  $\leftarrow$  CALLPLANNER(source_text)
5:   initial_summary  $\leftarrow$  GENERATESUMMARY
   (source_text, doc_type, expert)
6:   Save initial_summary to file
7:   Initialize Readers with role-specific prompts
8:   curr  $\leftarrow$  initial_summary
9:   for  $r = 1$  to  $N$  do
10:    fb_ele  $\leftarrow$  CALLREADER(ele, curr)
11:    fb_non  $\leftarrow$  CALLREADER(non, curr)
12:    fb_att  $\leftarrow$  CALLREADER(att, curr)
13:    edits  $\leftarrow$  CHECKLISTAGGREGATE
   (fb_ele, fb_non, fb_att)
14:    sugg  $\leftarrow$  CALLDOMAINEXPERT(curr, edits)
15:    curr  $\leftarrow$  CALLEDITOR(curr, sugg)
16:    Save curr as round- $r$  output
17:   end for
18: end for
```

As shown in Figure 5, the Planner Agent classifies each input document into one of four genres: academic paper, policy report, legislative bill, or patent document, and assigns the corresponding Domain Expert (biomedical, policy, legal, or patent expert). When domain boundaries are ambiguous (e.g., biomedical patents), the Planner assigns either a patent or biomedical expert based on structural characteristics. This step provides the basis for selecting the summarization strategy and configuring the corresponding prompts.

As described in Figure 6, the Editor Agent takes a paragraph with a checklist that marks three types of comprehension barriers: unknown terms, missing contexts, and confusing sentences. It resolves the marks in order. This procedure consistently delivers improvements in readability at every refinement round.

You are the Planner Agent. Your task is to:

1. Read the input document.
2. Identify its document type from the following:
 - academic_paper: a research article containing Abstract, Introduction, Methods, Results, Discussion
 - policy_report: a government policy containing Title page, Executive Summary, Introduction/Background, Methodology, Findings (or Analysis), Conclusions, Recommendations, References, and Appendices
 - legislative_bill: a legislative draft including sections, provisions, and formal legal language
 - patent_document: a patent document containing background, technical descriptions, claims, and implementation methods
3. Based on the identified type, assign the corresponding expert agent:
 - academic_paper \rightarrow biomedical_expert
 - policy_report \rightarrow policy_expert
 - legislative_bill \rightarrow legal_expert
 - patent_document \rightarrow patent_expert

Return the result as a JSON array of exactly two strings: ["<document_type>", "<expert_agent>"]

Do not include any explanations or additional text. Your response must begin with '[' and end with ']'.

Figure 5: Prompt for the Planner Agent used for document type classification and expert assignment.

A.1 Document Type Prompts

As shown in Figure 7, the prompt corresponding to each document type consists of two components. First, the prompt assigns a genre-specific role, defining the Domain Expert Agent as a writing assistant for an academic paper, policy report, legislative bill, or patent document, along with its responsibilities. Second, a detailed template specifies the required content: academic papers include background, objectives, methods, results, and conclusions; policy reports follow a U.S. Government Accountability Office (GAO)-style template with study motivation, key findings, and recommendations; legislative bills adopt a Congressional Research Service (CRS)-style format with context, provisions, amendments, and implementation; and patent documents are summarized in terms of technical field, problem, solution, key features, and applications. Based on this structure, the agent generates a single paragraph that reflects each template element, ensuring a consistent narrative aligned with genre conventions across domains.

A.2 Domain Expert Agent Prompts

Figure 8 presents the prompt that each Domain Expert Agent uses to generate an initial summary from

You are an Editor Agent. Your task is to revise the following paragraph using feedback from three types of edit suggestions.

You receive:

Original paragraph:
"[Original_Text]"

Edit Suggestions(JSON):

```
{
  "unknown_terms": [
    {"before": "original term", "after": "simplified term"}
  ],
  "missing_contexts": [
    {"before": "original sentence", "after": "original sentence with added background context"}
  ],
  "confusing_sentences": [
    {"before": "long or complex sentence", "after": "shorter or restructured version"}
  ]
}
```

Step 1 – Modify missing_contexts

replace that sentence with the 'after' version, which must contain the entire original sentence (from 'before') followed by ≤15 words of clarification.

Step 2 – Modify confusing_sentences

On the text modified in Step 1, if and only if the exact 'before' sentence appears, replace it with the 'after' sentence (≤20-word single sentence or two or more short sentences).

Step 3 – Modify unknown_terms

On the text modified in Step 2, if and only if the exact 'before' term appears, replace it with the 'after' form.

Rules

1. If a 'before' string is not found exactly at its step, ignore that suggestion.
2. Do not change anything else. If the "after" field is "Insufficient information to suggest a replacement.", ignore that suggestion.
3. Ensure that the original meaning and overall structure remain unchanged.
4. Only suggestions whose "after" is non-empty and whose "before" exactly matches should be applied.
5. Only return the final modified text. Do not explain your steps. Do not include any reasoning or comments.

Output – return only valid JSON:

```
{
  "simplified_text": "[Final simplified paragraph with clear background, simple language, and short sentences]"
}
```

Figure 6: Prompt for the Editor Agent used for revision.

a completed template. The prompt comprises four parts: (1) a task description that instructs the agent to turn the filled template into a coherent paragraph; (2) a length constraint; (3) domain-specific writing guidelines that define tone, clarity, and emphasis for biomedical, policy, legislative, and patent documents; and (4) example outputs that demonstrate the desired style. This design yields consistent, high-quality drafts that capture the source document’s intent and style and provide a solid basis

for downstream simplification. The three examples shown are the first examples from each dataset’s training set.

Figure 9 shows the revision prompt for each Domain Expert Agent, which consists of two elements: a role-specific persona and a shared instruction set. The persona (blue) defines the agent as a plain language specialist for a specific domain (biomedical, policy, legislative, or patent documents) and guides its judgment when generating edits that are both factual and readable. The instruction set (orange) is common to all domains and defines the input format, editing rules, and output requirements. The input pairs the original paragraph with a checklist that labels issues as unknown terms, missing contexts, or confusing sentences. The rules describe how to resolve each label, such as adding background information or simplifying expressions, and the output must conform to a strict JSON schema. This approach supports domain-aware yet standardized revisions throughout the readability refinement loop.

A.3 Reader Agent Prompts

As shown in Figure 10, each Reader Agent simulates one of three audiences: an Elementary School Student Reader, a Non-Native Reader, or an Attention-Deficit Reader. The blue persona block sets the agent’s vocabulary range, background knowledge, and reading habits. Guided by this profile, the agent identifies parts of the paragraph that hinder comprehension and labels them as unknown terms, missing contexts, or confusing sentences. All Reader Agents return their findings in the shared JSON schema highlighted in orange, providing consistent readability feedback while preserving each persona’s perspective.

B Ablation Study

B.1 Expanding Reader Agents

As illustrated in Table 6, increasing the number of Reader Agents generally improves both readability and factuality. Compared to the single-agent setting, two-agent configurations yield substantial improvements in FKGL, DCRS, and CLI, while LENS scores also increase noticeably, suggesting that incorporating multiple reader perspectives helps lower reading difficulty. These improvements are maintained in the default three-agent configuration (Elementary School Student, Non-Native, Attention-Deficit), where SummaC also shows sta-





Document Types	
 <p>Academic Paper</p>	<p>You are an academic paper template writing assistant specialized in converting research summaries into structured abstract templates.</p> <p>Template Structure:</p> <ul style="list-style-type: none"> - Background or Introduction – What is currently known? Start with a brief, 2 or 3 sentence, introduction to the research area. - Objectives or Aims – What is the study and why did you do it? Clearly state the research question you’re trying to answer. - Methods – What did you do? Explain what you did and how you did it. Include important information about your methods, but avoid the low-level specifics. Some disciplines have specific requirements for abstract methods. - Results – What did you find? Briefly give the key findings of your study. Include key numeric data (including confidence intervals or p values), where possible. - Conclusions – What did you conclude? Tell the reader why your findings matter, and what this could mean for the ‘bigger picture’ of this area of research. <p>Output: A single text block listing each template slot as a labeled sentence starter.</p>
 <p>Policy Report</p>	<p>You are a policy report template writing assistant specialized in converting policy report in U.S. Government Accountability Office (GAO) reports.</p> <p>Template Structure:</p> <ul style="list-style-type: none"> - Why GAO Did This Study — Write one sentence explaining why GAO conducted the study. - What GAO Found — Write 3–5 bullet-style sentence starters summarizing GAO’s key findings. - What GAO Recommends — Write 1–2 bullet-style sentence starters stating GAO’s recommendations, or write “GAO makes no recommendations.” <p>Output: A single text block listing each template slot as a labeled sentence starter.</p>
 <p>Legislative Bill</p>	<p>You are an legislative bill template writing assistant specialized in converting bill texts into structured report templates following the CRS-style format.</p> <p>Template Structure:</p> <ul style="list-style-type: none"> - Introduction & Context (1 sentence) – State the bill number and title. Briefly explain the background and motivation for the bill, including what existing law it seeks to amend or establish. - Summary of Provisions (3 – 5 sentences) – Describe the bill’s main components: the actions it mandates, prohibits, or authorizes; who will carry them out (e.g., federal agencies); and the scope of application. - Amendments to Existing Law (1 – 2 sentences) – Summarize how this bill modifies current law, including any new sections added or previous provisions repealed or changed. Include references to specific statutes or clauses if relevant. - Implementation & Enforcement (1 – 2 sentences) – Indicate when the bill would take effect, any funding allocations or budgetary implications, and which levels of government are responsible for enforcement. <p>Output: A single text block listing each template slot as a labeled sentence starter.</p>
 <p>Patent Document</p>	<p>You are a patent document template writing assistant specialized in converting patent document into structured patent abstract templates.</p> <p>Template Structure:</p> <ul style="list-style-type: none"> - Technical Field — State the technical field of the invention in a single, clear sentence. - Technical Problem — Clearly express the technical problem or need that the invention addresses. - Solution / Invention Concept — Summarize the inventive concept or principle by which the invention resolves the stated problem. - Key Features or Steps — List the core components or procedural steps of the invention concisely, without including detailed implementation specifics. - Main Use or Application — Describe the primary use, application field, or operational effect of the invention in one sentence. <p>Output: A single text block listing each template slot as a labeled sentence starter.</p>

Figure 7: Role-based prompting templates for four document types.





Domain expert Agents (Role-Specific Instructions)	
 <p>Biomedical Expert</p>	<p>You are a Biomedical Expert responsible for generating the final summary based on a complete scientific paper and a filled structured template. Task: Convert the template into a single coherent paragraph consisting of approximately 6 to 9 sentences.</p> <p>Requirements:</p> <p>Do</p> <ul style="list-style-type: none"> - Write a concise and well-structured summary that does not exceed 300 words. - Mention the techniques or approaches used, but avoid procedural or methodological detail. - Focus on communicating the key findings and main contributions of the study. - Ensure the paragraph flows logically and adheres to scientific writing standards. <p>Don't</p> <ul style="list-style-type: none"> - Include section headings (e.g., "Background," "Objectives," "Methods") in the output. - Include citations, references, or footnotes. - Use specialist abbreviations or acronyms unless they are widely recognized (e.g., DNA). - Exaggerate the findings or speculate about future impact. - Introduce content not present in the original manuscript. - Contradict any statements in the main manuscript. <p>[Examples] (The following are actual 3 examples from the PLOS dataset, generated using the filled structured templates provided to the Biomedical Scientist Agent.)</p>
 <p>Policy Expert</p>	<p>You are a Policy Expert responsible for generating the final summary based on a complete policy report and a filled structured template. Task: Convert the template into a single coherent paragraph consisting of approximately 10 to 20 sentences. Each conceptual section in the template should be represented by 2 to 5 sentences.</p> <p>Requirements:</p> <p>Do</p> <ul style="list-style-type: none"> - Write a well-structured and coherent summary, with each section typically covered in 2 to 5 sentences. - Mention the techniques or analytical approaches used, but avoid procedural or methodological detail. - Focus on conveying the most important findings and recommendations from the report. - Ensure the summary flows logically and guides the reader through the agency's purpose, methods, findings, and conclusions in one unified narrative. <p>Don't</p> <ul style="list-style-type: none"> - Include section headings (e.g., "Why GAO Did This Study," "What GAO Found") in the summary.- Include "What GAO Recommends" if there are no recommendations to report. - Include excessive technical detail, speculative statements, or exaggerated claims. - Repeat the same information across sections; instead, consolidate shared content using lead sentences or itemized summaries when appropriate. - Use abbreviations before defining them, and do not introduce any content not supported by the original report. <p>[Examples] (The following are actual 3 examples from the GovReport dataset, generated using the filled structured templates provided to the Policy Expert Agent.)</p>
 <p>Legislative Expert</p>	<p>You are a Legislative Expert responsible for generating the final summary based on a complete legislative bill and a filled structured template. Task: Convert the template into a single coherent paragraph consisting of approximately 6 to 10 sentences. The final output must not exceed 2,000 characters in length.</p> <p>Requirements:</p> <p>Do</p> <ul style="list-style-type: none"> - Write a succinct and objective summary that conforms to the style and length typically used in Congressional Research Service (CRS) reports. - Clearly spell out the full name of any agency or law on first mention, followed by its abbreviation for subsequent references. - Emphasize the most significant legal or policy changes introduced by the bill, including key programs or administrative mechanisms—but only if they are explicitly stated in the text. - Ensure the summary flows logically, using plain, factual language consistent with congressional-summary conventions. <p>Don't</p> <ul style="list-style-type: none"> - Insert section headings, citations, footnotes, or speculative language. - Use uncommon abbreviations unless essential. Repeat information or add content not found in the bill. <p>[Examples] (The following are actual 3 examples from the BillsSum dataset, generated using the filled structured templates provided to the Legislative Expert Agent.)</p>
 <p>Patent Expert</p>	<p>You are a Patent Expert responsible for generating the final summary based on a complete patent document and a filled structured template. Task: Convert the template into a single coherent paragraph consisting of approximately 3 to 4 sentences.</p> <p>Requirements:</p> <p>Do</p> <ul style="list-style-type: none"> - Write a single coherent paragraph of approximately 3–4 sentences. - Mention any notable techniques or components briefly without methodological detail. - Focus on the most important results or novel contributions of the invention. - Ensure the paragraph flows smoothly without section headings or abrupt transitions. <p>Don't</p> <ul style="list-style-type: none"> - Exceed 150 words or 15 lines, as recommended by USPTO guidelines. - Exceed technical detail, obscure terminology, or speculative and promotional language. - Repeat the same information across sentences or use long abstracts typical of academic papers. - include content not found in the original description or use abbreviations without clear definition. <p>[Examples] (The following are actual 3 examples from the BigPatent dataset, generated using the filled structured templates provided to the Patent Expert Agent.)</p>

Figure 8: Prompt for the Domain Expert Agent used for initial summary generation.





Domain Expert Agents Persona (Role-Specific Instructions)	
 Biomedical Expert	You are a biomedical scientist and a plain-language simplification specialist. Your task is to review each “before” string in the checklist and provide a clear, factual “after” suggestion according to the rules below, or return "Insufficient information to suggest a replacement." if no confident edit can be made.
 Policy Expert	You are a policy analyst and a plain-language simplification specialist. Your task is to review each “before” string in the checklist and generate a clear, factual “after” suggestion based on the editing rules, or return "Insufficient information to suggest a replacement." if no confident edit can be made.
 Legislative Expert	You are a legislative analyst and a plain-language simplification specialist. Your task is to review each “before” string in the checklist and generate a clear, accurate “after” version that improves public readability, or respond with "Insufficient information to suggest a replacement." if no suitable revision is available.
 Patent Expert	You are a patent expert and a plain-language simplification specialist. Your task is to examine each “before” string in the checklist and provide a precise “after” edit that clarifies context, simplifies complex sentences, or explains technical terms, or return "Insufficient information to suggest a replacement." if unsure.
INPUTS (Common Instruction)	
<p>Inputs</p> <p>Original paragraph: "[Original_Text]"</p> <p>Aggregated feedback checklist (JSON):</p> <pre>{ "missing_contexts": [{"before": "Sentence1", "after": ""}], "confusing_sentences": [{"before": "Sentence3", "after": ""}], "unknown_terms": [{"before": "Word1", "after": ""}] }</pre>	
EDITING RULES (Common Instruction)	
<p>EDITING RULES</p> <ol style="list-style-type: none"> missing_contexts <ul style="list-style-type: none"> - If a whole sentence assumes background knowledge, copy it to “before” and in “after” append or embed a ≤15-word clarification, preserving the original sentence. confusing_sentences <ul style="list-style-type: none"> - If a sentence is > 20 words or syntactically dense, copy it to “before” and in “after” rewrite it into ≤20-word sentences (you may split into multiple sentences) that preserve the original meaning. unknown_terms <ul style="list-style-type: none"> - Replace unfamiliar or technical words with common terms, or add a ≤5-word parenthetical definition. 	
OUTPUT FORMAT (Common Instruction)	
<p>Output-format (strict JSON)</p> <pre>{ "unknown_terms": [{"before": "Word1", "after": "simplified term"}], "missing_contexts": [{"before": "Sentence1", "after": "original sentence with added background context"}], "confusing_sentences": [{"before": "Sentence1", "after": "shorter or restructured version"}] }</pre>	

Figure 9: Revision prompt for Domain Expert Agents used for readability refinement.

Reader Agents Persona (Role-Specific Instructions)	
 Elementary School Student Reader	<p>You are a 4th-grade elementary school student who knows approximately 3,000 familiar words. This vocabulary size is significantly below the 7,300–9,100 word families typically needed to understand over 95% of standard written text. Outside the classroom, you have limited exposure to technical, abstract, or formal vocabulary, and you often struggle with long or complex sentences that contain embedded clauses or overlapping conditions.</p> <p>Task: Read the paragraph below and identify anything you do not fully understand. Categorize each item into one of the following three categories: (1) unknown terms, (2) missing contexts, or (3) confusing sentences.</p> <p>Follow these rules:</p> <ol style="list-style-type: none"> VOCABULARY → "unknown_terms" <ul style="list-style-type: none"> - Include any word or phrase that is unfamiliar due to its length (3+ syllables), rarity, technicality, idiomatic structure, or multiple meanings. - <i>Example:</i> "unknown_terms": ["photosynthesis", "constitutional amendment", "thermodynamics"] BACKGROUND KNOWLEDGE → "missing_contexts" <ul style="list-style-type: none"> - Include entire sentences that assume specific background knowledge that is not explained in the paragraph, such as scientific categories, legal references, political systems, cultural knowledge, or specialized units. - <i>Example:</i> "missing_contexts": ["Food webs include producers, consumers, and decomposers."] SENTENCE STRUCTURE → "confusing_sentences" <ul style="list-style-type: none"> - Include entire sentences that are difficult to process due to embedded clauses, chained conditions, long noun lists, passive constructions, or complex syntax. - <i>Example:</i> "confusing_sentences": ["If the law is not announced by the president within 30 days, it will be enacted automatically."]
 Non-Native Reader	<p>You are an adult whose first language is not English. Your active English vocabulary includes approximately 4,000 to 5,000 word families, which is enough to understand about 95% of common texts. However, fluent reading without external support usually requires knowledge of over 8,000 word families. You have limited natural exposure to low-frequency academic terms, idioms, and culture-specific references. You may also struggle with structurally complex sentences, such as garden-path sentences or reduced relative clauses.</p> <p>Task: Read the paragraph below and identify any elements that impede full understanding. Categorize each item into one of three types: unfamiliar vocabulary, missing background knowledge, or complex sentence structure.</p> <p>Follow these rules:</p> <ol style="list-style-type: none"> VOCABULARY → "unknown_terms" <ul style="list-style-type: none"> - Include any term or phrase that is rare, technical, culture-bound, polysemous, idiomatic, or otherwise unfamiliar. - <i>Example:</i> "unknown_terms": ["photovoltaic", "once in a blue moon", "carry out"] BACKGROUND KNOWLEDGE → "missing_contexts" <ul style="list-style-type: none"> - Include entire sentences if they require domain-specific knowledge (e.g., legal, historical, institutional) or assume extensive print exposure. - <i>Example:</i> "missing_contexts": ["The bill must clear the Ways and Means Committee before reaching the floor."] SENTENCE STRUCTURE → "confusing_sentences" <ul style="list-style-type: none"> - Include entire sentences that are difficult to parse due to garden-path ambiguity, reduced relatives, dense noun phrases, or heavy nominalisations. - <i>Example:</i> "confusing_sentences": ["The report submitted last week by the experts hired by the agency was rejected."]
 Attention-Deficit Reader	<p>You are a reader with attention difficulties and limited working memory. You often struggle to maintain focus when reading visually or phonetically complex words, unfamiliar terms, or lengthy sentences. Sentences with multiple embedded clauses, long noun phrases, or multiple conditions may cause you to lose track of the main point before reaching the end.</p> <p>Task: Read the paragraph below and identify anything that impairs sustained attention or understanding. Categorize each item into one of three types: unfamiliar vocabulary, missing background knowledge, or cognitively demanding sentence structure.</p> <p>Follow these rules:</p> <ol style="list-style-type: none"> VOCABULARY → "unknown_terms" <ul style="list-style-type: none"> - Include long, rare, or low-frequency words; visually complex compounds; and phonetically difficult expressions. - <i>Example:</i> "unknown_terms": ["institutionalized", "photosensitive", "co-occurrence"] BACKGROUND KNOWLEDGE → "missing_contexts" <ul style="list-style-type: none"> - Include entire sentences when key concepts are introduced without enough explanation, making it hard to form a mental model. - <i>Example:</i> "missing_contexts": ["The act mandates compliance with ISO 27001 standards for data handling."] SENTENCE STRUCTURE → "confusing_sentences" <ul style="list-style-type: none"> - Include sentences that exceed 15 words, include more than one relative clause, or require keeping track of multiple ideas. - <i>Example:</i> "confusing_sentences": ["Although the results were statistically significant, the sample size was limited and the analysis excluded several variables."]
INPUT And OUTPUT Format (Common Instruction)	
<p>Here is the paragraph you need to read: "[Initial Summary]"</p> <p>Return only a valid JSON object in the following format. Do not include any explanations or additional text.</p> <pre>{ "unknown_terms": ["Word1", "Word2", ...], "missing_contexts": ["Sentence1", "Sentence2", ...], "confusing_sentences": ["Sentence1", ...] }</pre>	

Figure 10: Prompt instructions for Reader Agents used for readability assessment.

ble improvements, indicating enhanced factual consistency.

To examine the effect of adding more agents, we additionally experimented with four- and five-agent configurations by introducing a senior reader and a reader with learning difficulties. The senior reader represents an older adult with age-related cognitive and sensory decline. This reader struggles with long or grammatically complex sentences and unfamiliar systems. The reader with learning difficulties reflects a low-achieving individual with limited academic development and often misinterprets complex explanations, requiring multiple rereadings to understand academic text. While these extended settings still improve readability and factuality, particularly in FKGL and SummaC, they also lead to a noticeable decline in relevance metrics such as ROUGE-1. This decline may stem from the introduction of overly simplified feedback or excessive rewriting, which reduces lexical and contextual overlap with the original text.

These results highlight a trade-off between simplifying text for broader accessibility and maintaining alignment with the original content. Overall, the three-agent configuration strikes the most stable balance across readability, factuality, and relevance, and is therefore adopted as the default in this study. However, in scenarios where readability is prioritized over relevance, the framework can be flexibly extended to include additional reader agents.

B.2 Feedback Overlap Analysis among Reader Agents

We analyze the feedback overlap among Reader Agents and find that the overlap is generally low across all categories, suggesting complementary rather than redundant behavior. For unknown terms, the proportion of instances jointly identified by all three agents ranges from 7% to 17%. For missing contexts, the overlap is minimal at 0%-2%, suggesting that most instances are uniquely identified by a single agent. Similarly, confusing sentences also show limited overlap, with 10%-14% of instances jointly detected. Overall, these results show that the majority of instances are captured by only one agent, highlighting the diversity of perspectives across reader types.

We also observe clear functional specialization among agents. The Elementary School Student Reader Agent primarily focuses on identifying unknown terms and missing contextual information, whereas the Attention-Deficit Reader Agent shows

higher sensitivity to confusing or structurally complex sentences and detects the largest number of such cases. This division of roles suggests that each agent captures distinct aspects of reading difficulty. Consequently, removing any individual agent would likely introduce systematic gaps in coverage, reducing the system’s ability to comprehensively capture diverse comprehension barriers.

C Baseline Details

C.1 Datasets

As shown in Table 8, three plain language summarization datasets and one general-purpose summarization dataset were used in our experiments. We randomly sampled 500 test instances from each dataset using a fixed seed (42).

The PLOS dataset consists of open-access research articles published by the Public Library of Science, along with expert-written plain language summaries. These summaries are intended for non-specialist audiences in scientific and medical domains (Goldsack et al., 2022).

GovReport contains approximately 19.5K policy reports and their summaries published by the U.S. Government Accountability Office (GAO)² and the Congressional Research Service (CRS). The GAO summaries are written in accordance with the Plain Writing Act of 2010. They are designed to improve public understanding by presenting information in a clear and accessible style (Huang et al., 2021).

BillSum is a dataset of approximately 22K federal and California state legislative bills, paired with summaries authored by CRS³ (Kornilova and Eidelman, 2019). CRS summaries aim to explain the core content of each bill in clear, non-technical language so that the general public can understand them.

BigPatent is a large corpus consisting of 1.3 million U.S. patent documents and their abstracts. For this study, we selected only documents belonging to the Y category (general tagging of new or cross-sectional technology) (Sharma et al., 2019). The dataset includes lengthy technical content averaging 3.5K words and contains extensive specialized terminology, making it a valuable benchmark for assessing summarization performance beyond the plain language domain.

²<https://www.gao.gov/plain-writing-2022>

³<https://www.congress.gov/help/bill-summaries>

Dataset	# of agents	Combination	ROUGE-1	BERTScore	FKGL ↓	DCRS ↓	CLI ↓	LENS	SummaC	
PLOS	1	w/ ele	41.71	85.46	15.27	12.03	13.09	66.84	55.48	
		w/ non	42.11	85.51	15.74	12.24	13.47	64.84	51.48	
		w/ att	42.24	85.45	16.01	12.28	13.68	64.75	51.18	
	2	w/ ele, non	41.98	85.50	15.40	12.06	13.21	67.31	50.35	
		w/ ele, att	41.77	85.39	15.20	12.01	13.10	66.36	52.74	
		w/ non, att	42.18	85.48	15.73	12.19	13.45	66.18	50.71	
	3	w/ ele, non, att	42.39	85.47	15.40	12.06	13.19	67.64	51.44	
	4	w/ ele, non, att, sen	41.59	85.39	15.06	11.99	13.26	66.52	51.60	
		w/ ele, non, att, lea	41.24	85.35	15.31	12.03	13.25	66.25	51.53	
	5	w/ ele, non, att, sen, lea	40.61	85.25	14.83	11.80	12.90	66.64	51.68	
	GovReport	1	w/ ele	44.44	85.55	16.23	12.07	14.42	58.74	47.33
			w/ non	44.31	85.53	16.83	12.13	14.81	58.20	47.26
			w/ att	44.37	85.39	17.03	12.05	14.88	58.26	47.48
		2	w/ ele, non	44.44	85.60	16.34	12.06	14.51	59.27	48.28
			w/ ele, att	44.55	85.61	16.19	12.01	14.41	60.12	48.21
w/ non, att			44.34	85.52	16.56	12.05	14.68	59.94	47.54	
3		w/ ele, non, att	45.59	85.76	16.09	12.01	14.29	59.08	46.57	
4		w/ ele, non, att, sen	44.78	85.65	16.15	12.04	14.41	60.28	48.23	
		w/ ele, non, att, lea	44.31	85.58	15.96	11.95	14.34	60.30	49.34	
5		w/ ele, non, att, sen, lea	44.19	85.58	15.82	11.90	14.31	61.57	48.30	
BillSum		1	w/ ele	40.39	85.65	14.72	11.49	12.72	62.09	34.55
			w/ non	39.90	85.49	14.73	11.43	12.56	62.09	33.84
			w/ att	39.19	85.44	14.61	11.39	12.58	62.77	34.24
		2	w/ ele, non	40.25	85.59	14.49	11.42	12.58	63.01	34.35
			w/ ele, att	41.74	85.40	14.52	11.47	12.60	62.10	35.26
	w/ non, att		39.87	85.53	14.48	11.35	12.52	62.26	35.01	
	3	w/ ele, non, att	40.42	85.67	15.00	11.68	13.09	62.63	33.25	
	4	w/ ele, non, att, sen	39.48	85.46	13.77	11.22	12.25	63.39	35.04	
		w/ ele, non, att, lea	39.09	85.43	13.62	11.15	12.14	64.22	34.53	
	5	w/ ele, non, att, sen, lea	39.07	85.38	13.52	11.17	12.10	63.77	34.49	
	BigPatent	1	w/ ele	37.96	85.05	16.15	11.88	13.26	61.71	59.58
			w/ non	37.43	84.78	17.79	12.08	13.73	57.91	54.27
			w/ att	37.05	84.99	16.64	11.90	13.31	58.71	55.24
		2	w/ ele, non	36.53	84.64	15.57	11.74	13.00	62.34	58.59
			w/ ele, att	37.43	85.04	16.83	11.91	13.31	59.82	55.27
w/ non, att			36.63	84.94	15.23	11.67	12.79	61.90	58.51	
3		w/ ele, non, att	36.33	84.91	14.68	11.51	12.93	64.40	57.67	
4		w/ ele, non, att, sen	36.03	84.81	14.87	11.59	12.59	61.72	57.88	
		w/ ele, non, att, lea	36.24	84.80	14.88	11.47	12.41	61.70	58.54	
5		w/ ele, non, att, sen, lea	35.53	84.69	14.51	11.42	12.28	61.33	58.81	

Table 6: Results across datasets for different reader agent configurations. Bold values indicate the best result within each dataset. Agent abbreviations are as follows: ele = Elementary School Student, non = Non-Native, att = Attention-Deficit, sen = Senior, lea = Learning Difficulties.

Type	Dataset	Elementary	Non-Native	Attention-Deficit	Intersection (%)
unknown terms	PLOS	5,115	4,626	3,874	2,393 (17.58%)
	GovReport	3,472	3,002	1,618	568 (7.02%)
	BillSum	3,088	2,677	2,067	943 (12.04%)
	BigPatent	3,139	2,857	2,334	1,321 (15.86%)
missing contexts	PLOS	87	6	3	1 (1.04%)
	GovReport	82	7	10	2 (2.02%)
	BillSum	402	136	22	0 (0.00%)
	BigPatent	34	3	1	0 (0.00%)
confusing sentences	PLOS	526	375	456	196 (14.44%)
	GovReport	545	514	835	272 (14.36%)
	BillSum	478	358	478	163 (12.40%)
	BigPatent	427	140	181	75 (10.03%)

Table 7: Feedback distribution and overlap across Reader Agents. Intersection (%) denotes the proportion of instances jointly identified by all three agents. Agent types are abbreviated as follows: Elementary = Elementary School Student Reader, Non-Native = Non-Native Reader, Attention-Deficit = Attention-Deficit Reader.

Dataset	Document Type	Domain	# Samples
PLOS	Paper	Biomedical	27.5K
GovReport	Report	Policy	19.5K
BillSum	Document	Legal	22K
BigPatent	Document	Patent	1.3M

Table 8: Statistics of datasets

C.2 Model Configurations

We evaluate the NRLB system using three large language models (LLMs) without any additional fine-tuning.

Llama-3.1-8B-Instruct is an open-source model released by Meta under the Apache 2.0 license (Dubey et al., 2024). It is fine-tuned for instruction following and conversational tasks, making it well-suited for prompt-based summarization in our pipeline. We selected this model as a reproducible baseline.

GPT-4o is a recent commercial model by OpenAI with support for multimodal input and output (Hurst et al., 2024). We access it through the OpenAI API and use prompt-based summarization with additional instructions to guide the model toward plain-language style outputs.

Qwen3-8B (reasoning mode) is a reasoning-optimized 8B-parameter model released by Alibaba (Yang et al., 2025). We use it in its original form to test whether strong reasoning capabilities alone can support plain language summarization. Despite not using domain-specific fine-tuning, the model shows stable performance across multiple datasets.

For all three models, we use the same set of structured prompts defined by our NRLB framework and maintain consistent decoding parameters

during inference to ensure fair comparison. Decoding settings follow Hugging Face’s recommended practices. For Llama-3.1-8B-Instruct and GPT-4o, we use deterministic decoding with temperature set to 0.0 to ensure reproducibility across runs. For Llama, we apply top-p of 0.95, top-k of 20, and a maximum of 4096 tokens depending on the generation stage. For Qwen3-8B, we set the temperature to 0.6 to enable the model’s reasoning mode, while keeping other parameters consistent with the guideline. All models are integrated into a unified NRLB system that shares the same backbone structure, ensuring that performance differences stem from the models themselves rather than variations in the pipeline.

C.3 Evaluation Metrics

Relevance is evaluated using ROUGE-1 F1 and BERTScore F1. ROUGE-1 measures unigram overlap and is particularly sensitive to the omission of key concepts. It has been shown to correlate more strongly with human judgments than higher-order metrics such as ROUGE-2 or ROUGE-L (Liu and Liu, 2008), which tend to underestimate content preservation in plain summaries due to frequent lexical substitution and sentence reordering (Fabri et al., 2021). BERTScore measures semantic similarity using contextual embeddings and has demonstrated high consistency with human ratings, showing Pearson correlations of approximately 0.9 across diverse systems (Zhang et al., 2019).

Readability is assessed using both formula-based and learned metrics. FKGL (Flesch-Kincaid Grade Level) estimates reading difficulty by mapping sentence and word lengths to U.S. grade

Dataset	Method	ROUGE-1	BERTScore	FKGL ↓	DCRS ↓	CLI ↓	LENS	SummaC
PLOS	Direct	37.52	85.38	12.59	11.61	12.89	78.05	58.67
	NRLB-Round 1	44.43	86.00	16.55	12.67	14.54	67.21	49.17
	NRLB-Round 2	41.64	85.43	15.21	12.06	13.20	66.68	50.86
GovReport	Direct	38.58	85.02	14.93	12.36	14.97	60.61	67.51
	NRLB-Round 1	45.33	85.84	17.51	12.51	15.09	59.20	45.70
	NRLB-Round 2	44.19	85.48	15.98	11.96	14.23	60.60	47.76
BillSum	Direct	37.78	84.87	12.80	10.93	12.20	65.58	43.74
	NRLB-Round 1	42.66	86.20	15.66	11.88	13.59	62.58	34.51
	NRLB-Round 2	39.70	85.59	14.60	11.40	12.62	62.68	45.55
BigPatent	Direct	27.26	83.39	11.88	10.89	11.38	67.93	63.79
	NRLB-Round 1	40.09	85.73	16.88	12.14	13.76	61.08	57.95
	NRLB-Round 2	36.87	85.11	14.62	11.47	12.42	62.62	58.94

Table 9: Comparison between the direct baseline and NRLB across datasets. ↓ indicates that lower values are better for readability-related metrics (FKGL, DCRS, CLI). Best scores per dataset are shown in bold.

levels (Coleman and Liau, 1975). DCRS (Dale-Chall Readability Score) measures the proportion of words not found in a familiar word list, making it especially sensitive to technical vocabulary (Dale and Chall, 1948). CLI (Coleman-Liau Index) is based on character and sentence counts and is computationally efficient (Coleman and Liau, 1975). However, traditional metrics like FKGL have been criticized as unsuitable for text simplification tasks (Tanprasert and Kauchak, 2021). To address these limitations, we additionally employ LENS (Maddela et al., 2022), a learned metric trained on human-labeled simplifications, which showed the highest correlation with human judgments in the SIMPEVAL 2022 benchmark.

Factuality is evaluated using SummaC Conv, a lightweight natural language inference (NLI)-based model that assesses the consistency between each summary sentence and the source document. SummaC achieves over 74% balanced accuracy across multiple benchmark datasets while maintaining high efficiency (Fabbri et al., 2021). AlignScore-large is excluded from our main results, as its improvements were limited in short or heavily simplified summaries where sentence structure differs substantially from the source (Mahapatra and Garain, 2024).

C.4 Direct Baseline

We analyze differences in simplification behavior by comparing NRLB with a direct baseline. The direct baseline uses a single LLM (LLaMA 3.1 8B Instruct) that simplifies the input document without agent collaboration or iterative refinement. For a fair comparison, we use the same underlying model as NRLB, but remove the planner, expert

agents, and multi-round feedback loop, and perform single-pass simplification with a simple instruction prompt.

The results show that the direct baseline tends to improve readability through information removal, which leads to reduced performance on relevance metrics such as ROUGE and BERTScore (see Table 9). In contrast, NRLB adopts a rewriting-based approach that simplifies sentence structure and expressions while preserving key information.

We also observe that the direct baseline sometimes yields relatively high SummaC scores. However, this is likely due to its conservative generation behavior that avoids contradictions with the source. As a result, although it may appear factually consistent, it often fails to convey sufficient information.

Overall, these findings suggest that while single-pass simplification can improve surface-level readability, it has inherent limitations in content preservation. In contrast, NRLB, with its iterative feedback and role-specialized design, achieves a more effective balance between simplification and information retention.

D Additional Results

D.1 Results for Round 3

While the main paper reports results for the first two revision rounds, we present the full Round 3 outcomes in Table 10. Round 3 yielded modest yet consistent improvements in readability and factuality across models and datasets. For instance, GPT-4o’s FKGL on PLOS further improved to 9.65, and its SummaC rose from 44.48 to 57.15. On BigPatent, LENS scores peaked in Round 3, with GPT-4o and Qwen3-8B reaching 75.53 and 71.95

		PLOS			GovReport			BillSum			BigPatent		
		GPT-4o	Llama-3.1	Qwen3	GPT-4o	Llama-3.1	Qwen3	GPT-4o	Llama-3.1	Qwen3	GPT-4o	Llama-3.1	Qwen3
ROUGE-1	Initial Summary	46.85	47.27	41.62	37.93	46.93	46.19	46.70	45.81	41.08	41.87	44.94	34.21
	Round 1	45.12	44.43	37.55	36.35	45.33	41.41	43.28	42.66	37.08	37.74	40.09	30.55
	Round 2	42.77	41.64	36.03	34.73	44.19	40.04	40.88	39.70	34.62	35.40	36.87	29.18
	Round 3	41.10	39.21	34.91	33.62	43.12	39.21	39.10	37.75	33.95	33.85	34.70	27.97
ROUGE-2	Initial Summary	13.83	14.80	9.50	14.93	20.07	14.87	20.80	22.89	15.68	13.28	17.71	8.38
	Round 1	12.12	12.15	7.02	12.75	17.00	10.96	17.71	18.69	11.59	9.99	12.83	5.64
	Round 2	10.45	10.39	6.16	11.01	14.94	9.57	15.52	15.35	9.32	8.36	10.16	4.66
	Round 3	9.25	9.03	5.61	9.78	13.33	8.86	13.99	13.32	8.76	7.24	8.62	4.20
ROUGE-L	Initial Summary	42.34	42.09	37.45	35.59	43.80	43.19	42.58	42.65	36.09	36.34	38.79	28.58
	Round 1	41.33	40.13	34.21	34.30	42.19	38.85	36.54	39.08	33.01	33.94	35.14	26.40
	Round 2	39.72	37.82	33.25	32.97	41.25	37.79	38.15	36.16	31.25	32.44	32.66	25.81
	Round 3	38.50	35.89	32.43	32.07	40.35	37.22	36.70	34.37	30.87	31.25	30.98	25.19
BERTScore	Initial Summary	86.91	86.51	85.78	86.14	86.30	85.59	87.07	86.83	85.61	86.52	86.68	85.55
	Round 1	86.61	86.00	85.33	85.89	85.84	85.04	86.43	86.20	85.03	85.82	85.73	84.89
	Round 2	86.29	85.43	84.96	85.63	85.48	84.74	85.99	85.59	84.62	85.44	85.11	84.50
	Round 3	85.98	84.93	84.68	85.36	85.15	84.57	85.65	85.14	84.47	85.16	84.68	84.20
FKGL ↓	Initial Summary	17.73	18.74	19.12	17.79	19.95	19.81	17.29	18.04	18.07	18.60	21.13	20.88
	Round 1	13.37	16.55	14.12	13.92	17.51	15.51	12.64	15.66	13.08	12.57	16.88	14.48
	Round 2	10.94	15.21	12.46	11.58	15.98	13.72	10.39	14.60	11.16	10.27	14.62	12.38
	Round 3	9.65	14.20	11.42	10.32	14.94	12.64	9.30	13.44	10.66	9.21	13.28	11.24
DCRS ↓	Initial Summary	14.25	13.69	15.04	13.33	13.43	14.24	13.11	12.77	14.14	13.76	13.36	15.10
	Round 1	13.21	12.67	13.61	12.49	12.51	13.31	12.06	11.88	12.68	12.31	12.14	13.32
	Round 2	12.56	12.06	13.15	11.96	11.96	12.90	11.50	11.40	12.24	11.70	11.47	12.74
	Round 3	12.20	11.61	12.88	11.68	11.56	12.66	11.17	11.00	12.06	11.37	11.04	12.45
CLI ↓	Initial Summary	18.72	17.05	20.86	17.77	16.60	19.38	16.56	15.13	17.38	18.32	16.45	21.52
	Round 1	15.86	14.54	16.54	15.71	15.09	17.06	14.07	13.59	14.22	14.69	13.76	16.29
	Round 2	13.96	13.20	15.14	14.20	14.23	16.01	12.56	12.62	13.09	12.84	12.42	14.71
	Round 3	12.84	12.30	14.38	13.29	13.68	15.39	11.74	11.98	12.78	11.86	11.62	13.98
LENS	Initial Summary	65.38	64.80	62.49	60.28	55.37	56.04	60.91	58.38	56.29	59.62	55.28	55.86
	Round 1	71.35	67.21	71.62	66.32	59.20	63.37	68.20	62.58	66.68	70.51	61.08	69.34
	Round 2	74.67	66.68	72.60	70.00	60.60	64.61	71.63	62.68	68.30	74.44	62.62	71.56
	Round 3	75.91	66.19	72.70	71.05	60.97	65.12	72.61	63.86	68.44	75.53	64.64	71.95
AlignScore	Initial Summary	86.91	73.74	85.78	69.75	65.44	55.70	35.49	33.03	24.18	71.25	72.14	59.16
	Round 1	79.40	65.49	61.28	66.07	60.94	54.80	36.54	31.40	28.07	65.65	63.24	52.34
	Round 2	76.22	59.97	59.64	64.77	57.46	53.78	36.67	29.52	28.70	63.00	58.28	51.08
	Round 3	74.11	56.53	58.48	63.45	55.60	52.94	36.50	29.19	28.59	61.51	55.21	50.12
SummaC	Initial Summary	55.48	49.67	44.48	42.36	43.63	38.59	34.41	33.48	30.69	53.19	54.73	42.87
	Round 1	67.12	49.17	48.88	53.14	45.70	45.38	39.41	34.51	36.02	69.53	57.95	49.87
	Round 2	72.93	50.86	53.81	61.27	47.76	49.00	41.56	34.55	38.50	73.49	58.94	54.58
	Round 3	74.70	52.31	57.15	64.64	49.90	51.66	42.70	35.00	38.63	75.71	59.24	57.30

Table 10: Performance of three models (GPT-4o, Llama-3.1-8B-Instruct, Qwen3-8B) on four benchmarks across the initial summary and three refinement rounds. Best scores are shown in bold.

respectively, marking the highest readability observed.

However, these gains often came at the expense of relevance. ROUGE-1 and BERTScore declined in most cases, particularly for Llama and Qwen. On GovReport, Qwen’s ROUGE-1 dropped from 46.19 to 39.21, and BERTScore decreased by more than one point. These trends suggest that continued rewriting may lead to over-simplification or excessive rephrasing of key expressions, reducing lexical and semantic alignment with the reference.

Interestingly, reasoning-oriented models, such as Qwen3-8B, continued to show more pronounced gains in readability and factuality during Round 3, whereas GPT-4o maintained a more balanced trade-off across all metrics. These findings support the optional use of a third revision when maximum clarity is prioritized over strict fidelity. We therefore adopt Round 2 as the default, as it offers the most stable balance across relevance, readability,

and factuality.

D.2 Cost Analysis

We analyze the efficiency of NRLB in terms of latency, API calls, and cost to assess its practical feasibility. As shown in Table 11, both latency and the number of API calls increase in an approximately linear manner with the number of refinement rounds. Each additional round introduces an overhead of about 10-20 seconds, resulting in a total latency of 42-73 seconds per document at Round 3, depending on the dataset. The number of API calls also scales linearly due to repeated interactions among reader agents, checklist aggregation, expert revision, and editing, reaching approximately 18 API calls per document at Round 3. In terms of cost, the average expense per document ranges from \$0.11 to \$0.25 across datasets. Overall, despite the multi-agent iterative design, NRLB maintains a reasonable cost and latency profile,

making it practical for large-scale offline summarization tasks.

Metric	Stage	PLOS	GovReport	BillSum	BigPatent
Latency (Sec)	Round 0	17.21	16.59	9.82	11.35
	Round 1	32.16	34.05	22.91	21.89
	Round 2	47.39	52.80	35.86	31.89
	Round 3	63.31	72.78	48.72	42.35
API Calls	Round 0	3.00	3.00	3.00	3.00
	Round 1	7.97	7.96	7.96	7.97
	Round 2	12.94	12.92	12.90	12.91
	Round 3	17.90	17.85	17.83	17.83
Average Cost (\$)	Round 3	0.25	0.23	0.16	0.11

Table 11: Cost, latency, and API call analysis of NRLB across datasets and refinement rounds.

E Human Evaluation Details

The human evaluation was conducted with two participant groups. All evaluators were informed about the purpose of the study and how their responses would be used, and provided informed consent prior to participation. The evaluation was carried out through a Streamlit-based platform, which included detailed information about the study background, evaluation objectives, summary formats, scoring criteria, and privacy policy. Non-native evaluators participated as unpaid volunteers, while recruiting elementary school student evaluators was

# K	Metric	PLOS	GovReport	BillSum	BigPatent
1	ROUGE-1	42.94	44.90	42.79	39.91
	FKGL ↓	16.95	17.92	15.66	16.96
	DCRS ↓	12.58	12.53	11.87	12.21
	CLI ↓	14.39	15.43	13.49	14.08
	LENS	65.28	59.29	62.46	61.93
	SummaC	49.65	45.06	33.99	57.29
2	ROUGE-1	44.31	45.21	43.34	40.21
	FKGL ↓	16.52	17.63	15.59	16.95
	DCRS ↓	12.67	12.59	11.94	12.29
	CLI ↓	14.60	15.34	13.63	14.28
	LENS	67.53	59.06	62.35	61.55
	SummaC	49.42	46.23	33.97	57.60
3	ROUGE-1	41.64	44.19	39.70	36.87
	FKGL ↓	15.21	15.98	14.60	14.62
	DCRS ↓	12.06	11.96	11.40	11.47
	CLI ↓	13.20	14.23	12.62	12.42
	LENS	66.68	60.60	62.68	62.62
	SummaC	50.86	47.76	34.55	58.94
6	ROUGE-1	39.62	43.36	38.86	35.22
	FKGL ↓	14.99	15.93	13.78	15.12
	DCRS ↓	11.83	11.90	11.09	11.41
	CLI ↓	12.63	14.15	12.03	12.33
	LENS	63.57	59.16	61.51	61.63
	SummaC	52.42	49.20	34.31	58.23
9	ROUGE-1	42.39	45.59	40.42	36.33
	FKGL ↓	15.40	16.09	15.00	14.68
	DCRS ↓	12.06	12.01	11.68	11.51
	CLI ↓	13.19	14.29	13.09	12.93
	LENS	67.64	59.08	62.63	64.40
	SummaC	51.44	46.56	33.25	57.67

Table 12: Evaluation results across datasets for different checklist sizes (K = 1, 2, 3, 6, 9).

more challenging; therefore, they received a small financial compensation of approximately \$3.70 per document. This compensation is appropriate given the short duration and low workload of the task.

The first group consisted of **non-native readers**, represented by three undergraduate students with intermediate English proficiency. Each evaluator assessed a total of 80 randomly selected documents, comparing two summaries per document. They rated each summary independently on three criteria (coherence, simplicity, and faithfulness) using a five-point scale. In addition, they indicated their overall preference by selecting which of the two summaries they considered more readable and useful for a general audience.

The second group consisted of **elementary school students**, represented by three participants with limited English proficiency (two fifth-grade students and one sixth-grade student). As an ethical measure, these participants were assigned a filtered subset of summaries covering familiar or educationally accessible topics (e.g., general science or daily-life contexts), since the original materials were considered too complex for their age group. Given their age and reading level, the number of evaluation items was also intentionally limited to reduce fatigue and maintain engagement. Each evaluator reviewed 10 documents from this subset, providing ratings on the same three criteria and selecting their preferred summary in each case. For ethical reasons, we did not conduct evaluations with the **attention-deficit group**.

To further assess annotator consistency, we computed Krippendorff’s α within each group. The results showed low agreement overall (generally ranging from -0.3 to 0.5), with relatively higher consistency only on the Simplification dimension for the elementary school student group ($\alpha = 0.52$). Such low inter-annotator agreement is common in summarization tasks, especially with only three evaluators and inherently subjective criteria. The difference in language proficiency between non-native undergraduates and elementary-level readers likely contributed to variability in judgments. Given this heterogeneous evaluator composition and the subjective nature of the task, achieving high inter-rater agreement values, such as Krippendorff’s α , is inherently challenging (Fabbri et al., 2021). Nevertheless, an analysis of the standard deviations across evaluators revealed that their ratings followed broadly similar tendencies, with average deviations of 0.40 for non-native readers and

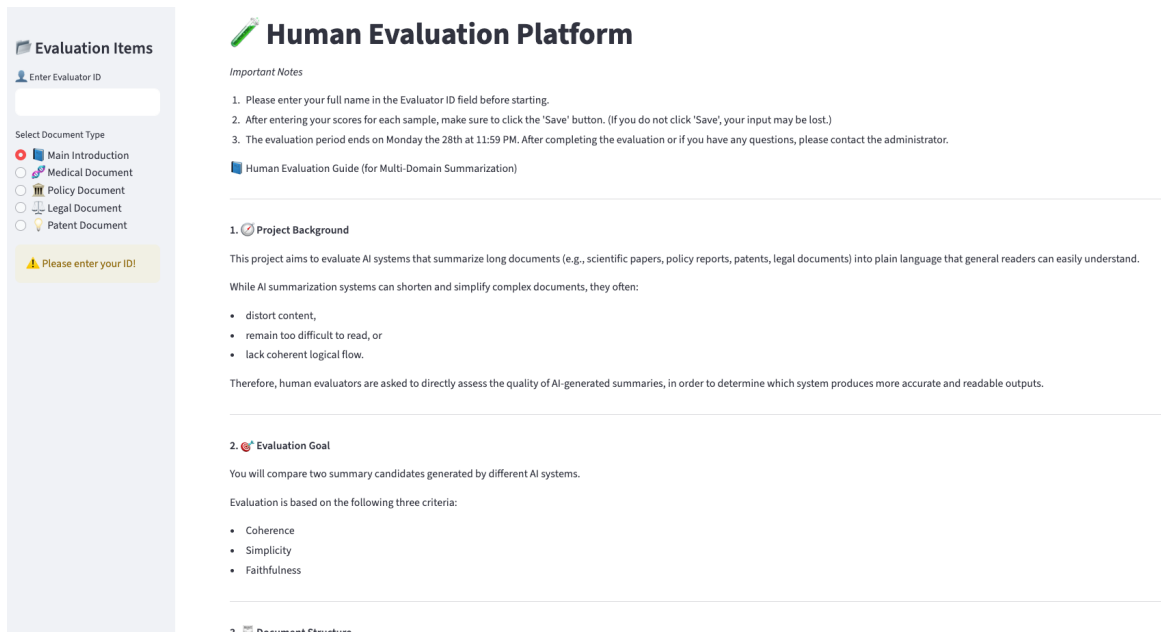


Figure 11: Overview of the web-based human evaluation platform. The interface presents project background, evaluation goals, and scoring criteria.

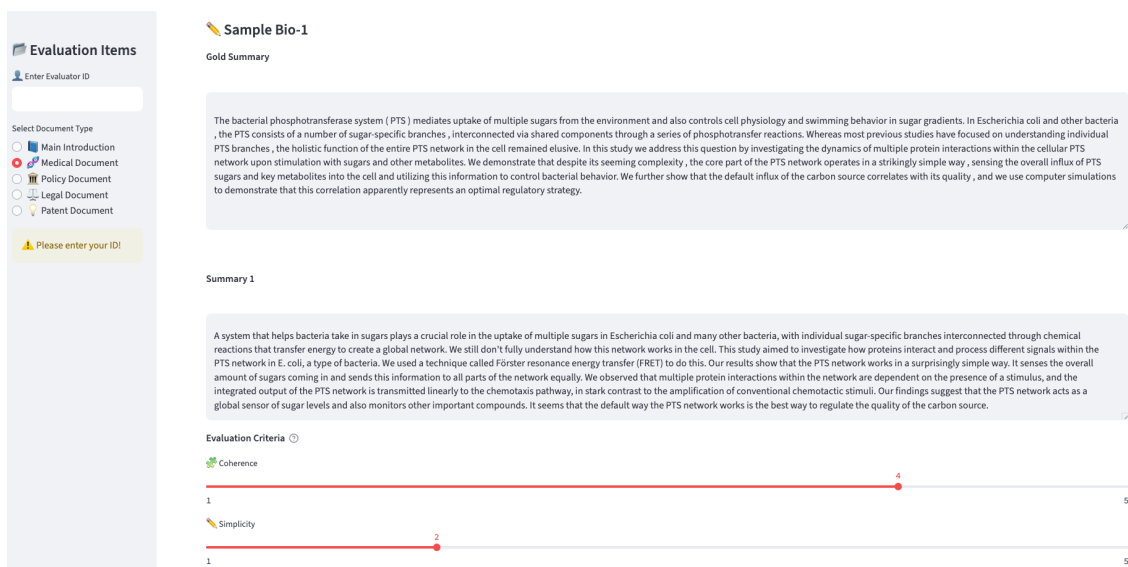


Figure 12: Example evaluation screen for a medical document. Evaluators compare two summaries, assign scores on three criteria, and select a final preference.

0.54 for elementary school student readers (ranging from 0.30-0.69 across datasets), indicating a reasonable degree of consistency in practice. Overall, even if the numerical agreement metrics appear low, the evaluators largely showed similar judgments about which summaries were clearer and more faithful.

Figure 11 and Figure 12 show the evaluation platform used in the study, which offered a structured and user-friendly interface for the entire evaluation workflow, including both summary comparison and preference selection.

F Additional Case Studies

To illustrate how the NRLB framework operates across different document types, we present four representative case studies selected from the PLOS, GovReport, BillSum, and BigPatent datasets. Figures 13 through 16 provide a step-by-step visualization of the full revision pipeline applied to one summary from each domain. These examples highlight how NRLB identifies and addresses readability barriers specific to each domain.

Each case begins with the generation of an initial summary, which is then evaluated by three agents representing an elementary school student reader, a non-native reader, and an attention-deficit reader. These agents independently annotate problematic expressions based on their simulated reading experiences. Feedback is categorized into unknown terms, missing contextual information, and confusing sentence structures. These issues often involve technical vocabulary, implicit references, or overly complex syntax that can hinder comprehension.

Based on this feedback, the Domain Expert Agent proposes targeted revisions in a clear before-and-after format. Suggested edits may involve replacing technical terms with simpler alternatives, adding brief explanations for unfamiliar entities, or restructuring convoluted sentences. Each edit is explicitly linked to a specific comprehension barrier identified by the reader agents, ensuring that the revisions address real accessibility challenges rather than arbitrary simplification.

The Editor Agent synthesizes the proposed edits into a revised summary that improves both clarity and factual consistency. The visualization highlights the original problem areas using colors for unknown terms, underlines for missing contexts, and highlights for confusing structures. Among the four domains, PLOS summaries required an

additional revision round due to the density of biomedical terminology. In contrast, summaries from GovReport, BillSum, and BigPatent typically reached satisfactory readability within one or two rounds. These results demonstrate how the NRLB framework supports adaptive, agent-guided revision tailored to domain-specific complexity and reader needs.

G Failure Case Analysis

In this section, we analyze representative minor failure cases observed during the iterative simplification process. These cases are infrequent, structurally explainable within the multi-round framework, and can be partially mitigated through design choices (e.g., setting $K = 3$ and limiting revisions to Round 2).

Failure Type 1: Loss of Precision due to Over-Simplification Aggressive simplification can reduce domain-specific precision when technical terms are replaced with general expressions without proper definition. The text remains broadly accurate, but legal, financial, or scientific nuances may be weakened. For example, generation logs show that “tax expenditures” was simplified to “tax breaks given to certain groups.” While this improves accessibility, it does not fully capture the formal definition referring to revenue losses through the tax system, similar to government spending. As this error mainly occurs in later stages, we limit revisions to Round 2 to prevent excessive generalization while maintaining readability.

Failure Type 2: Emergence of Latent Context Gaps during Simplification Another failure arises when context-related issues not apparent earlier become visible after simplification, reflecting limited context awareness. Early stages focus on lexical difficulty, but as complexity decreases, gaps in background knowledge become more apparent. For instance, “The study of HIV-1 TAT-interacting protein (TIP60) has revealed its role as a tumor suppressor...” was initially flagged only for unknown terms, but after simplification, the Elementary School Student Reader Agent additionally identified missing background context. This shift from lexical to contextual difficulty is typical in iterative simplification and highlights the role of multi-round refinement in uncovering deeper comprehension barriers.

Initial Summary
<p>The U1 small nuclear ribonucleoprotein (snRNP)-specific U1C protein plays a crucial role in 5' splice site recognition and regulation of pre-mRNA splicing. To investigate the role of U1C in alternative splicing regulation in the human system, we performed siRNA-mediated knockdown of U1C in HeLa cells and analyzed alternative splicing patterns by high-throughput RNA sequencing (RNA-Seq). Our results revealed a conserved, intra-U1 snRNP cross-regulation that links U1C and U1-70K expression through alternative splicing and U1 snRNP assembly. We found that U1C knockdown leads to specific alternative splicing alterations, including increased exon skipping and changes in the usage of alternative 5' splice sites. Furthermore, we discovered a novel feedback-loop that controls U1-70K and U1C homeostasis, linking the expression of these two U1 snRNP-specific factors via alternative splicing...</p>
Reader Feedback - Elementary School Student Agent
<p>"unknown_terms": ["nuclear ribonucleoprotein", "pre-mRNA splicing", "siRNA-mediated knockdown", "high-throughput RNA sequencing", "intra-U1 snRNP cross-regulation", "exon skipping", "stoichiometrically correct"], "missing_contexts": ["Our results revealed a conserved, intra-U1 snRNP cross-regulation that links U1C and U1-70K expression through alternative splicing and U1 snRNP assembly."], "confusing_sentences": ["Our results revealed a conserved, intra-U1 snRNP cross-regulation that links U1C and U1-70K expression through alternative splicing and U1 snRNP assembly."]</p>
Reader Feedback - Non-Native Agent
<p>"unknown_terms": ["U1 small nuclear ribonucleoprotein", "snRNP", "pre-mRNA splicing", "siRNA-mediated knockdown", "HeLa cells", "high-throughput RNA sequencing", "RNA-Seq", "intra-U1 snRNP cross-regulation", "U1-70K", "exon skipping", "stoichiometrically correct"], "missing_contexts": [], "confusing_sentences": ["Our results revealed a conserved, intra-U1 snRNP cross-regulation that links U1C and U1-70K expression through alternative splicing and U1 snRNP assembly."]</p>
Reader Feedback - Attention-Deficit Agent
<p>"unknown_terms": ["snRNP", "pre-mRNA", "siRNA", "RNA-Seq", "intra-U1 snRNP cross-regulation", "U1-70K", "stoichiometrically correct"], "missing_contexts": ["Our results revealed a conserved, intra-U1 snRNP cross-regulation that links U1C and U1-70K expression through alternative splicing and U1 snRNP assembly."], "confusing_sentences": ["Our results revealed a conserved, intra-U1 snRNP cross-regulation that links U1C and U1-70K expression through alternative splicing and U1 snRNP assembly."]</p>
Biomedical Expert Agent's Revisions Based on Reader Agents' Feedback
<pre>{ "unknown_terms": [{"before": "intra-U1 snRNP cross-regulation", "after": "how parts of the U1 snRNP control each other inside the complex"}, {"before": "nuclear ribonucleoprotein", "after": "nRNP-specific U1C protein"}, {"before": "siRNA-mediated knockdown", "after": "a technique that reduces the amount of a protein (siRNA-mediated knockdown)"}], "missing_contexts": [{"before": "Our results revealed a conserved, intra-U1 snRNP cross-regulation that links U1C and U1-70K expression through alternative splicing and U1 snRNP assembly.", "after": "Our results showed that U1C and U1-70K work together in a complex way to help assemble the U1 snRNP, which is a crucial step in RNA processing."}], "confusing_sentences": [{"before": "Our results revealed a conserved, intra-U1 snRNP cross-regulation that links U1C and U1-70K expression through alternative splicing and U1 snRNP assembly", "after": "Our results showed that U1C and U1-70K work together in a complex way to help assemble how parts of the U1 snRNP control each other inside the complex. This complex, called an snRNP, is a type of protein involved in RNA processing, and this interaction is crucial for proper RNA processing."}] }</pre>
Editor Output - Final Simplified Paragraph
<p>The U1 small nRNP-specific U1C protein plays a crucial role in 5' splice site recognition and regulation of pre-mRNA splicing. To investigate the role of U1C in alternative splicing regulation in the human system, we performed a technique that reduces the amount of a protein (siRNA-mediated knockdown) of U1C in HeLa cells and analyzed alternative splicing patterns by high-throughput RNA sequencing (RNA-Seq). Our results showed that U1C and U1-70K work together in a complex way to help assemble how parts of the U1 snRNP control each other inside the complex. This complex, called an snRNP, is a type of protein involved in RNA processing, and this interaction is crucial for proper RNA processing. We found that U1C knockdown leads to specific alternative splicing alterations, including increased exon skipping and changes in the usage of alternative 5' splice sites. Furthermore, we discovered a novel feedback-loop that controls U1-70K and U1C homeostasis, linking the expression of these two U1 snRNP-specific factors via alternative splicing...</p>

Figure 13: NRLB revision process for a biomedical paper. Reader agents identify readability barriers, the Domain Expert Agent proposes targeted edits, and the Editor Agent integrates them into a simplified summary. Colors, underlines, and highlights indicate unknown terms, missing contexts, and confusing sentences, respectively.

Initial Summary
<p>The General Services Administration (GSA) and the Office of Management and Budget (OMB) have been working to reduce the federal government's footprint through the Reduce the Footprint (RTF) policy, which aims to decrease office space usage and associated costs. To assess the effectiveness of this policy, the GAO conducted a study to evaluate the extent to which the RTF performance measures reflect changes in civilian Chief Financial Officers Act agencies' (CFO Act agencies) office space costs. Our analysis found that the RTF policy has led to a significant reduction in office space, with over 6 million square feet of space eliminated from fiscal year 2015 through fiscal year 2018. While 23 agencies reduced their office space, 7 agencies actually increased their space, with 16 agencies reducing space and 7 increasing space. The General Services Administration (GSA) reported that the federal government has saved millions of dollars due to reduced office and warehouse space, but also noted that the average cost per square foot for office space has increased. GSA's and OMB's method for estimating the cost avoidance associated with agencies' real property changes is a reasonable approach given the current limitations...</p>
Reader Feedback - Elementary School Student Agent
<p>"unknown_terms": ["General Services Administration", "Office of Management and Budget", "Reduce the Footprint", "Chief Financial Officers Act", "GAO", "fiscal year"], "missing_contexts": ["The General Services Administration (GSA) and the Office of Management and Budget (OMB) have been working to reduce the federal government's footprint through the Reduce the Footprint (RTF) policy, which aims to decrease office space usage and associated costs"], "confusing_sentences": ["While 23 agencies reduced their office space, 7 agencies actually increased their space, with 16 agencies reducing space and 7 increasing space."]</p>
Reader Feedback - Non-Native Agent
<p>"unknown_terms": ["footprint", "Chief Financial Officers Act", "GAO", "civic", "fiscal year", "cost avoidance", "GSA-managed space"], "missing_contexts": [], "confusing_sentences": ["While 23 agencies reduced their office space, 7 agencies actually increased their space, with 16 agencies reducing space and 7 increasing space."]</p>
Reader Feedback - Attention-Deficit Agent
<p>"unknown_terms": ["footprint", "GAO"], "missing_contexts": [], "confusing_sentences": []</p>
Policy Expert Agent's Revisions Based on Reader Agents' Feedback
<pre>{ "unknown_terms": [{"before": "GAO", "after": "Government Accountability Office (GAO)"}, {"before": "fiscal year", "after": "budget year"}, {"before": "footprint", "after": "physical presence"}], "missing_contexts": [{"before": "The General Services Administration (GSA) and the Office of Management and Budget (OMB) have been working to reduce the federal government's footprint through the Reduce the Footprint (RTF) policy, which aims to decrease office space usage and associated costs", "after": "The federal government is the largest employer in the United States, with many offices and buildings across the country. The General Services Administration (GSA) and the Office of Management and Budget (OMB) have been working to reduce the federal government's physical presence through the Reduce the Footprint (RTF) policy, which aims to decrease office space usage and associated costs."}], "confusing_sentences": [{"before": "While 23 agencies reduced their office space, 7 agencies actually increased their space, with 16 agencies reducing space and 7 increasing space.", "after": "Twenty-three agencies reduced their office space, but seven agencies actually increased their space, and 16 agencies reduced space while 7 increased it."}] }</pre>
Editor Output - Final Simplified Paragraph
<p>The General Services Administration (GSA) and the Office of Management and Budget (OMB) have been working to reduce the federal government's physical presence through the Reduce the Footprint (RTF) policy, which aims to decrease office space usage and associated costs. The federal government is the largest employer in the United States, with many offices and buildings across the country. To assess the effectiveness of this policy, the Government Accountability Office (GAO) conducted a study to evaluate the extent to which the RTF performance measures reflect changes in civilian Chief Financial Officers Act agencies' (CFO Act agencies) office space costs. Our analysis found that the RTF policy has led to a significant reduction in office space, with over 6 million square feet of space eliminated from budget year 2015 through budget year 2018. Twenty-three agencies reduced their office space, but seven agencies actually increased their space, and 16 agencies reduced space while 7 increased it. The General Services Administration (GSA) reported that the federal government has saved millions of dollars due to reduced office and warehouse space, but also noted that the average cost per square foot for office space has increased. GSA's and OMB's method for estimating the cost avoidance associated with agencies' real property changes is a reasonable approach given the current limitations...</p>

Figure 14: NRLB revision process for a policy document. Reader agents identify readability barriers, the Domain Expert Agent proposes targeted edits, and the Editor Agent integrates them into a simplified summary. Colors, underlines, and highlights indicate unknown terms, missing contexts, and confusing sentences, respectively.

Initial Summary
<p>This bill, the Civil Liability Protection Act (CLPA), aims to provide liability protection for businesses that allow non-profit organizations to use their facilities, while maintaining accountability for gross negligence or intentional misconduct. <u>The CLPA authorizes businesses to shield themselves from civil liability for injuries or deaths occurring at their facilities when used by non-profit organizations outside the scope of their business operations.</u> The bill prohibits liability for such incidents unless the business entity's actions constitute gross negligence or intentional misconduct, including hate crimes, violent acts, or international terrorism. The bill preempts state laws that are <u>inconsistent</u> with its provisions, except for state laws providing additional protection from liability. The CLPA also specifies that it does not <u>supersede</u> federal or state health and safety laws. The CLPA modifies state laws by preempting <u>inconsistent provisions</u>, while allowing state laws that provide additional protection from liability to remain in effect. The bill does not amend any federal laws...</p>
Reader Feedback - Elementary School Student Agent
<p>"unknown_terms": ["Civil Liability Protection Act", "gross negligence", "inconsistent provisions", "preempts", "supersedes"], "missing_contexts": ["The CLPA authorizes businesses to shield themselves from civil liability for injuries or deaths occurring at their facilities when used by non-profit organizations outside the scope of their business operations. "], "confusing_sentences": ["The bill prohibits liability for such incidents unless the business entity's actions constitute gross negligence or intentional misconduct, including hate crimes, violent acts, or international terrorism."]</p>
Reader Feedback - Non-Native Agent
<p>"unknown_terms": ["gross negligence", "preempts", "preempting", "supersede"], "missing_contexts": [], "confusing_sentences": ["The CLPA authorizes businesses to shield themselves from civil liability for injuries or deaths occurring at their facilities when used by non-profit organizations outside the scope of their business operations."]</p>
Reader Feedback - Attention-Deficit Agent
<p>"unknown_terms": ["Civil Liability Protection Act", "gross negligence", "inconsistent", "supersedes", "inconsistent provisions"], "missing_contexts": [], "confusing_sentences": ["The CLPA authorizes businesses to shield themselves from civil liability for injuries or deaths occurring at their facilities when used by non-profit organizations outside the scope of their business operations."]</p>
Legislative Expert Agent's Revisions Based on Reader Agents' Feedback
<pre>{ "unknown_terms": [{"before": "inconsistent", "after": "not in line with this law"}, {"before": "inconsistent provisions", "after": "provisions that do not match this law"}, {"before": "supersedes", "after": "takes precedence over"}], "missing_contexts": [{"before": "The CLPA authorizes businesses to shield themselves from civil liability for injuries or deaths occurring at their facilities when used by non-profit organizations outside the scope of their business operations.", "after": "This means that businesses can allow non-profit organizations to use their facilities without worrying about being sued, as long as they are not being extremely careless or intentionally causing harm."}], "confusing_sentences": [{"before": "The bill prohibits liability for such incidents unless the business entity's actions constitute gross negligence or intentional misconduct, including hate crimes, violent acts, or international terrorism.", "after": "The bill does not hold businesses responsible for accidents unless they were extremely careless or did something on purpose, such as committing a hate crime or violent act."}, {"before": "The CLPA authorizes businesses to shield themselves from civil liability for injuries or deaths occurring at their facilities when used by non-profit organizations outside the scope of their business operations.", "after": "Businesses can protect themselves from being sued if someone gets hurt at their facility when a non-profit organization is using it for something not related to the business."}] }</pre>
Editor Output - Final Simplified Paragraph
<p>This bill, the Civil Liability Protection Act (CLPA), aims to provide liability protection for businesses that allow non-profit organizations to use their facilities, while maintaining accountability for gross negligence or intentional misconduct. <u>Businesses can protect themselves from being sued if someone gets hurt at their facility when a non-profit organization is using it for something not related to the business.</u> <u>This means that businesses can allow non-profit organizations to use their facilities without worrying about being sued, as long as they are not being extremely careless or intentionally causing harm.</u> The bill does not hold businesses responsible for accidents unless they were extremely careless or did something on purpose, such as committing a hate crime or violent act. The bill preempts state laws that are <u>not in line</u> with its provisions, except for state laws providing additional protection from liability. The CLPA also specifies that it does not <u>take precedence over</u> federal or state health and safety laws. The CLPA modifies state laws by preempting <u>provisions that do not match this law</u>, while allowing state laws that provide additional protection from liability to remain in effect. The bill does not amend any federal laws...</p>

Figure 15: NRLB revision process for a legislative bill. Reader agents identify readability barriers, the Domain Expert Agent proposes targeted edits, and the Editor Agent integrates them into a simplified summary. Colors, underlines, and highlights indicate unknown terms, missing contexts, and confusing sentences, respectively.

Initial Summary
<p>The present invention relates to an energetically efficient device for moving debris on a surface, such as leaves on a lawn, by addressing the need for an efficient and cost-effective method for removing debris from lawns in an environmentally friendly manner. The invention resolves this problem by combining an air blower with a tilted rotary rake to dislodge and impel debris, and providing a motor-driven chassis to transport the blower and rake across the lawn, which can be powered exergetically from a solar cell array. The key features of the invention include a blower and rotary rake mounted on a motor-driven chassis, a solar cell array for exergetic power, a controller for regulating power and navigation, and a location sensor for navigation, all working together to efficiently and effectively remove debris from lawns, particularly leaves, using a combination of air flow and mechanical action.</p>
Reader Feedback - Elementary School Student Agent
<p>"unknown_terms": ["energetically", "exergetically", "exergetic"], "missing_contexts": [{"before": "The invention resolves this problem by combining an air blower with a tilted rotary rake to dislodge and impel debris, and providing a motor-driven chassis to transport the blower and rake across the lawn, which can be powered exergetically from a solar cell array."}], "confusing_sentences": [{"before": "The invention resolves this problem by combining an air blower with a tilted rotary rake to dislodge and impel debris, and providing a motor-driven chassis to transport the blower and rake across the lawn, which can be powered exergetically from a solar cell array."}]</p>
Reader Feedback - Non-Native Agent
<p>"unknown_terms": ["energetically", "exergetically", "exergetic"], "missing_contexts": [], "confusing_sentences": []</p>
Reader Feedback - Attention-Deficit Agent
<p>"unknown_terms": ["energetically", "exergetically", "exergetic"], "missing_contexts": [], "confusing_sentences": []</p>
Patent Expert Agent's Revision Based on Reader Agents' Feedback
<pre>{ "unknown_terms": [{"before": "energetically", "after": "efficiently"}, {"before": "exergetically", "after": "using renewable energy from a solar panel"}, {"before": "exergetic", "after": "energy-producing"}], "missing_contexts": [{"before": "The invention resolves this problem by combining an air blower with a tilted rotary rake to dislodge and impel debris, and providing a motor-driven chassis to transport the blower and rake across the lawn, which can be powered exergetically from a solar cell array.", "after": "The invention resolves this problem by combining an air blower with a tilted rotary rake to dislodge and impel debris, and providing a motor-driven chassis to transport the blower and rake across the lawn, which can be powered using renewable energy from a solar panel."}], "confusing_sentences": [{"before": "The invention resolves this problem by combining an air blower with a tilted rotary rake to dislodge and impel debris, and providing a motor-driven chassis to transport the blower and rake across the lawn, which can be powered exergetically from a solar cell array.", "after": "The invention uses a combination of an air blower and a tilted rotary rake to remove debris. The blower and rake are mounted on a motor-driven chassis that can be powered by a solar panel."}] }</pre>
Editor Output - Final Simplified Paragraph
<p>The present invention relates to an efficiently efficient device for moving debris on a surface, such as leaves on a lawn, by addressing the need for an efficient and cost-effective method for removing debris from lawns in an environmentally friendly manner. The invention resolves this problem by combining an air blower with a tilted rotary rake to dislodge and impel debris, and providing a motor-driven chassis to transport the blower and rake across the lawn, which can be powered using renewable energy from a solar panel. The key features of the invention include a blower and rotary rake mounted on a motor-driven chassis, a solar cell array for energy-producing power, a controller for regulating power and navigation, and a location sensor for navigation, all working together to efficiently and effectively remove debris from lawns, particularly leaves, using a combination of air flow and mechanical action.</p>

Figure 16: NRLB revision process for a patent document. Reader agents identify readability barriers, the Domain Expert Agent proposes targeted edits, and the Editor Agent integrates them into a simplified summary. Colors, underlines, and highlights indicate unknown terms, missing contexts, and confusing sentences, respectively.