

High-Fidelity Contextual Generation via Prompt-Controlled Large-Context RAG

Havishya Vally Sunkara¹, Divya Sree Jampani², Manasa Tamma³, Visweswara Sujith Kumar Grandhi⁴,
Doddapaneni V Subba Rao⁵

¹²³⁴*Students, Department of Computer Science and Engineering SRK Institute of Technology, Vijayawada,
NTR District, Andhra Pradesh, India*

⁵*Associate Professor, Department of Computer Science and Engineering SRK Institute of Technology,
Vijayawada, NTR District, Andhra Pradesh, India*

Abstract—Deploying Large Language Models (LLMs) effectively in knowledge-intensive domains necessitates generating responses that are not only accurate but also demonstrably grounded in specific contextual documents. Conventional fine-tuning struggles to guarantee such inference-time faithfulness and requires constant retraining for dynamic data, while standard Retrieval-Augmented Generation (RAG) is often hampered by context fragmentation. This paper introduces Prompt-Controlled Large-Context RAG (PCLC-RAG), a methodology designed to unlock high-fidelity, contextually-aware generation by synergizing the vast capacity of modern large context window LLMs (e.g., $\geq 1\text{M}$ tokens) with precise inference-time prompt engineering via Structured Prompt Architectures (SPAs). PCLC-RAG utilizes direct document ingestion governed by the SPA, offering a practical alternative to fine-tuning, particularly for commercial applications with evolving knowledge bases, as long as documents fit the context window. Initial validation confirmed its capacity for exceptionally high qualitative accuracy (estimated $\sim 95\%$) without fine-tuning. We present PCLC-RAG as a highly effective paradigm for tasks demanding deep contextual understanding and verifiable accuracy, offering advantages in control, adaptability, and reduced maintenance overhead compared to fine-tuning. Scalability challenges and mitigation strategies are discussed, positioning PCLC-RAG as a viable architecture for demanding real-world applications.

Index Terms—Large Language Models (LLMs), Retrieval-Augmented Generation (RAG), Prompt Engineering, Structured Prompt Architecture (SPA), Large Context Window, Contextual Grounding, Contextual Synthesis, Verifiable AI, Zero-Shot Learning, Faithful Generation, Fine-tuning Alternatives, Commercial Applications, Dynamic Data.

I. INTRODUCTION

Large Language Models (LLMs) [1]- [3] hold immense potential, yet ensuring their outputs are factually accurate and strictly grounded in authoritative knowledge sources remains a critical barrier to deployment in high-stakes applications [5]. Fine-tuning [9], [10], while useful for specialization, offers weak guarantees against generating content inconsistent with specific documents provided only at inference time [6]. Furthermore, maintaining specialized models via fine-tuning becomes impractical in many commercial settings where knowledge bases (e.g., policies, manuals, regulations) are constantly updated. Retrieval-Augmented Generation (RAG) [13] improves grounding, but conventional snippet-based retrieval [14] can fragment context and impede holistic reasoning [15]. We introduce Prompt-Controlled Large-Context RAG (PCLC-RAG), leveraging recent LLMs with expansive context windows ($\geq 1\text{M}$ tokens) [18], [19]. PCLC-RAG combines direct, full-document ingestion (where context window allows) with sophisticated inference-time control via a Structured Prompt Architecture (SPA). This approach provides a powerful alternative to fine-tuning, allowing for high adaptability to dynamic document sets prevalent in commercial environments. The SPA guides the LLM's reasoning, enforces strict grounding to explicit text, and crucially, enables controlled contextual synthesis.

Initial validation via a prototype [?], which employed Google Gemini 1.5 Flash, demonstrated PCLC-RAG's potential. This prototype achieved

exceptionally high qualitative accuracy (est. ~95%) on document-based QA entirely without fine-tuning. This strongly suggests the viability of achieving high fidelity through inference-time control over large contexts using appropriately capable foundation models, making it an attractive option for commercial deployment where data evolves rapidly.

This paper formalizes PCLC-RAG as a powerful methodology. We detail its architecture (Section III), analyze its synthesis capability (Section IV), discuss validation insights (Section V), outline ideal applications, including commercial scenarios (Section VI), examine trade-offs (Section VII), and conclude (Section VIII).

II. RELATED WORK

PCLC-RAG synthesizes concepts from: LLMs and Fine-tuning [1], [4], [9]- [12]. PCLC-RAG offers an alternative focused on inference-time control vs. weight adaptation, reducing the need for constant retraining on dynamic data sources common in commercial use. RAG [13]- [16]. PCLC-RAG differs significantly by leveraging large context for *direct, potentially holistic ingestion*, shifting complexity from retrieval heuristics to prompt-based control (SPA) operating over the full provided context. Prompt Engineering & Instruction Following [2], [21]- [23]. PCLC-RAG employs advanced SPAs as the core logic mechanism, extending beyond simple instruction following. Large Context Models [18], [19]. PCLC-RAG is enabled by this architectural shift; research explores effective utilization [24]. Grounding & Verifiable AI [5]- [8]. PCLC-RAG contributes an inference-time, prompt-controlled grounding strategy designed for high fidelity.

III. PCLC-RAG: METHODOLOGY AND ARCHITECTURE

PCLC-RAG targets NLP tasks demanding demonstrable faithfulness to dynamic source documents. See Fig. 1. Its architecture comprises four key components:

A. High-Capacity Foundation LLM (M)

The methodology requires a state-of-the-art foundation model possessing specific capabilities:

- **Expansive Context Window:** A context capacity (C) large enough to ingest the full text of relevant documents for the target task (ideally $\geq 1M$ tokens

for complex scenarios). This large window is key to enabling direct ingestion as an alternative to retrieval or continuous fine-tuning.

- **Strong Reasoning Abilities:** Capable of understanding relationships, drawing logical inferences, and synthesizing information based on the provided context.
- **High Instruction Fidelity:** Accurately and reliably follows complex, multi-part instructions and constraints embedded within the SPA.
- **Multimodal Support (Optional):** For applications involving non-textual data within documents.

Models like Google's Gemini series [18] or Anthropic's Claude 3 family [19] exemplify the class of models suitable for PCLC-RAG.

B. Direct Context Ingestion Module

Simply concatenates the full text of all relevant source documents $D = \{d_i\}$ into a single context string D_{concat} , potentially with clear separators, ensuring the total prompt length remains below the model's context limit C . This preserves the holistic context, avoiding fragmentation inherent in snippet retrieval.

C. Structured Prompt Architecture (SPA) (T)

The core control mechanism. The SPA is a task-specific, engineered prompt template incorporating detailed instructions and constraints. Key design principles include:

- 1) **Operational Constraints:** Defining the LLM's role, scope of allowed knowledge (strictly D), and task objective.
- 2) **Grounding & Relevance Mandates:** Explicit rules requiring outputs to be directly supported by D , including specific rules governing how synthesis can occur (e.g., combining evidence) while forbidding external facts.
- 3) **Task Execution Logic:** Step-by-step reasoning guidance if needed (e.g., find relevant passages, synthesize, format output).
- 4) **Output Specification:** Defining the required format, length, and structure of the response R .
- 5) **Behavioral Guardrails:** Prohibitions against speculation, bias, or generating harmful content.
- 6) **Contingency Handling:** Instructions for cases where an answer cannot be found or synthesized solely from D .
- 7) **Dynamic Parameter Integration:** Placeholders for

inserting the User Query (Q) and the concatenated documents (D_{concat}).

Effective SPA engineering is crucial for PCLC-RAG's success and requires iterative refinement.

D. Generation Process

The final prompt P is constructed by instantiating the SPA Template T with the specific Query Q and the ingested documents D_{concat} . The LLM M processes this comprehensive prompt P , generating the response R strictly according to the rules embedded within the SPA, operating directly over the provided context D . See Algorithm 1.

Algorithm 1 PCLC-RAG Inference Workflow

Require: Doc set $D = \{d_i\}$, Query Q , SPA Template T , LLM M , Max Context C .
Ensure: Grounded/Synthesized Response R or Error E .

```

0:  $D_{\text{concat}} \leftarrow \text{ConcatenateWithSeparators}(D)$ 
0:  $P \leftarrow \text{InstantiateTemplate}(T, Q, D_{\text{concat}})$ 
0:  $\text{prompt\_len} \leftarrow \text{TokenCount}(P, M_{\text{tokenizer}})$ 
0: if  $\text{prompt\_len} > C$  then
0:   return  $E_{\text{contextLimit}}$  {Requires Hybrid approach or error}
0: end if
0:  $R \leftarrow M.\text{generate}(P, \text{GenerationParams})$  {Guided by SPA rules over D}

```

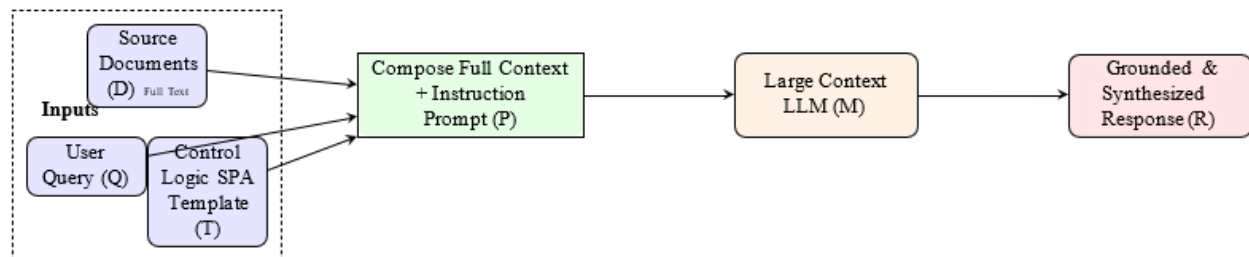


Fig. 1. Conceptual architecture of PCLC-RAG. Inputs (Documents D , Query Q , SPA Template T) form Prompt P . Large Context LLM (M) processes P guided by SPA to generate Response R (grounded extraction + controlled synthesis).

V. INITIAL VALIDATION AND HIGH ACCURACY OBSERVATIONS

The principles of PCLC-RAG were validated through the development and testing of a prototype question-answering system [20]. For this prototype, Google's Gemini 1.5 Flash model was utilized, as its capabilities and context window were sufficient for

0: return $R = 0$

IV. CONTROLLED CONTEXTUAL SYNTHESIS

A key capability enabled by PCLC-RAG is *controlled contextual synthesis*, distinct from ungrounded hallucination. While standard RAG often focuses on extracting snippets, PCLC-RAG, through its SPA, permits the LLM to leverage its powerful reasoning engine over the complete context D to:

- Synthesize answers by integrating evidence from multiple locations within D .
- Draw inferences and conclusions logically derived **only** from the information present in D .
- Utilize its pre-trained knowledge strategically to interpret, structure, and elaborate on information found within D 's scope, crucially **without** introducing external facts*.

This addresses complex queries requiring more than simple extraction, providing richer, more useful responses while remaining verifiably context-bound. The SPA acts as the control mechanism, explicitly defining the boundaries of permissible synthesis and forbidding the introduction of unrelated external information, thus guiding the LLM towards relevant and faithful synthesis.

the specific document types targeted (student-related records). It is important to note that while this specific model was used for validation, the PCLC-RAG methodology itself is applicable to any foundation model meeting the requirements outlined in Section III.A. This validation, while qualitative, yielded compelling results demonstrating the approach's potential across suitable LLMs.

A. Implementation

The prototype implemented direct context ingestion for relevant documents and utilized an iteratively refined SPA designed for high-fidelity QA. Generation was performed in a zero-shot manner, relying entirely on the SPA and the chosen model's inherent capabilities, without any task-specific fine-tuning.

B. Striking Qualitative Findings

Manual evaluation across a diverse set of document-based questions revealed remarkable performance:

- **Exceptional Accuracy & Faithfulness (Est. ~95% Qualitative):** The system exhibited outstanding accuracy for both direct extraction and, notably, for valid contextual synthesis. Responses showed extremely high fidelity to the scope and content of the source documents. This high accuracy level, achieved without fine-tuning, underscores the power of SPA-guided reasoning over large contexts.
- **Effective Controlled Synthesis:** The system successfully generated relevant and coherent synthesized answers requiring interpretation or integration of information from *D*, achieving this without fabricating information beyond the document's scope.
- **Robust Hallucination Mitigation:** The constraints embedded within the SPA proved highly effective in minimizing the generation of outputs untethered from the provided documents *D*. Instances of hallucination were significantly reduced compared to less constrained prompting methods.
- **Reliable SPA Adherence:** The LLM demonstrated strong adherence to the complex instructions and constraints specified in the SPA, indicating the feasibility of fine-grained control via prompt engineering in large-context models.

C. Limitations and Implications

While these initial findings are highly encouraging, they are based on qualitative assessment within a single task/LLM configuration (Gemini 1.5 Flash on specific documents) and lack comparison against rigorous quantitative baselines. The estimated 95% accuracy requires formal validation across various models and tasks. However, the sheer effectiveness

observed strongly motivates the formalization of PCLC-RAG presented in this paper and highlights its significant potential as a paradigm for building highly reliable and verifiable AI systems adaptable to dynamic knowledge sources, using a range of capable foundation models.

III. STRENGTHS AND IDEAL USE CASES

PCLC-RAG offers distinct advantages for applications where deep understanding, context-bound accuracy, and verifiable reasoning are paramount.

A. Core Strengths

Deep Contextual Understanding: Processes full documents (within context limits), enabling holistic reasoning.

High Accuracy & Faithfulness: SPA control promotes strict grounding and enables accurate synthesis.

Verifiable Grounding: Outputs are more easily traced back to source documents via SPA constraints.

• **Adaptability to Dynamic Data:** Easily adapts to new or updated documents at inference time without retraining.

- **Practical Fine-tuning Alternative:** Offers a strong alternative to continuous fine-tuning, especially in commercial settings with frequently updated knowledge bases (policies, manuals, etc.), reducing maintenance overhead and the need for large labeled datasets for grounding, provided documents fit the context window.
- **Reduced Labeled Data Needs (for Grounding):** Achieves high accuracy via prompting/SPA design, minimizing reliance on large labeled datasets specifically for fine-tuning grounding behaviors.

B. Ideal Application Domains

PCLC-RAG excels in scenarios demanding high reliability, verifiable grounding, and adaptability to changing information:

- **Commercial Knowledge Management:** Querying internal documentation, policies, product specs, or support knowledge bases that evolve over time.
- **Knowledge-Intensive QA:** Interacting with dense manuals, research papers, legal documents, or financial reports where accuracy and grounding are critical.
- **Compliance and Regulation:** Verifying procedures or outputs against potentially changing internal

- policies or external regulations.
- Contextual Summarization: Generating summaries guaranteed to reflect only the source material, useful for reports or briefings.
- Automated Reporting: Creating structured reports from logs, transcripts, or data feeds based on defined templates and source data.
- Domain-Specific Assistants: Providing grounded assistance in legal, medical, or financial domains based on authoritative, potentially updated, sources.

Its suitability for commercial use is enhanced by its ability to adapt to new information without the development lifecycle costs associated with fine-tuning for every data update.

VI.DISCUSSION: TRADE-OFFS, SCALABILITY, AND ENGINEERING

While PCLC-RAG offers significant advantages, its implementation involves key trade-offs and engineering considerations.

A. Core Advantages Revisited

PCLC-RAG's primary benefits stem from its unique architecture: enhanced verifiability through SPA control over full context, superior adaptability to dynamic knowledge sources compared to fine-tuning (providing a commercially viable alternative for evolving document sets, assuming context limits are met), potential for high accuracy on grounded tasks without labeled data overhead, and granular control over output characteristics, including enabling controlled synthesis.

B. Limitations and Engineering Considerations

Effective implementation requires addressing several factors:

- 1) Context Window Limits: Finite C is the primary constraint. PCLC-RAG is most directly applicable when the relevant document set fits within the model's context window. For larger corpora, hybrid strategies (e.g., SPA-controlled RAG over pre-filtered chunks, hierarchical processing) become necessary, adding complexity.
- 2) Long-Context Fidelity: Potential for degraded attention [24] requires careful SPA structuring or using models optimized for long context processing.

- 3) SPA Engineering Complexity: Designing robust, effective SPAs is skill-intensive and iterative.
- 4) SPA Robustness/Completeness: Ensuring SPAs handle diverse inputs and edge cases requires careful design and testing.
- 5) Synthesis Control Boundary: Precisely defining permissible synthesis vs. hallucination in SPA logic remains challenging.
- 6) Evaluation Challenges: Quantifying faithfulness and synthesis quality requires specialized metrics and potentially significant human evaluation effort.
- 7) Implicit vs. Explicit Knowledge Trade-off: PCLC-RAG prioritizes explicit grounding on provided documents. It leverages the LLM's pre-trained abilities for reasoning *about* the context but captures less implicit domain knowledge than extensive fine-tuning might. This is often desirable for verifiability but might be a limitation if deep, unstated domain assumptions are required.
- 8) Input Document Quality Dependency: Performance relies on reasonably clean input; robust pre-processing or multimodal models may be needed.
- 9) Rigorous Benchmarking Needed: Requires careful comparison against strong baselines using appropriate context-bound metrics across various capable LLMs.

C. Scalability and Efficiency Considerations

Deploying PCLC-RAG at scale requires addressing the latency and cost inherent in processing large contexts:

- Latency Mitigation: Employ optimized inference infrastructure (GPUs/TPUs, serving frameworks like vLLM), use faster model variants (e.g., Gemini Flash, if sufficient for the task) where possible, implement Hybrid RAG (pre-filtering) to intelligently reduce effective context size for the main LLM, and utilize output streaming for perceived responsiveness.
- Throughput Enhancement: Leverage horizontal scaling with load balancing, dynamic batching in inference servers, asynchronous processing queues for non-interactive tasks, and autoscaling infrastructure.

- **Cost Management:** Combine Hybrid RAG (pre-filtering) with prompt token optimization (if possible without losing necessary context), response length constraints via SPA, strategic caching, appropriate model selection (balancing capability vs. cost), and efficient hardware utilization.

A multi-pronged strategy is typically required to make PCLC- RAG practical for high-load scenarios.

VIII.CONCLUSION AND FUTURE WORK

PCLC-RAG presents a compelling methodology for high- fidelity, contextually-grounded NLP by integrating large con- text LLMs with precise SPA control. It offers a practical and adaptable alternative to fine-tuning, particularly for commercial applications requiring verifiable accuracy against dynamic document sets, provided these documents fit within the model's context window. Enabling both strict grounding and controlled contextual synthesis, it addresses key limitations of prior approaches. Initial validation demonstrated exceptional qualitative accuracy without fine-tuning, highlighting its potential.

While PCLC-RAG introduces engineering challenges, particularly around SPA design, context limits, and evaluation, practical mitigation strategies exist. Key future work must focus on:

- **Rigorous Quantitative Benchmarking:** Formally evaluating PCLC-RAG across different LLMs and diverse tasks against relevant baselines using metrics designed for context-bound faithfulness and synthesis quality.
- **Advancing SPA Engineering:** Developing systematic methodologies and tools for designing and optimizing robust SPAs.
- **Strategies for Ultra-Large Contexts:** Investigating and refining hybrid or hierarchical approaches to extend PCLC-RAG principles beyond single-context limits.
- **Standardized Evaluation Protocols:** Creating and vali- dating protocols and datasets for assessing large-context grounding and controlled synthesis.

PCLC-RAG represents a significant step towards realizing more reliable, adaptable, and contextually intelligent LLM ap- plications suitable for demanding

real-world and commercial use cases.

REFERENCES

- [1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems (NIPS)*, 2017, vol. 30.
- [2] T. B. Brown et al., "Language models are few-shot learners," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2020, vol. 33, pp. 1877–1901.
- [3] H. Touvron et al., "Llama 2: Open foundation and fine-tuned chat models," *arXiv preprint arXiv:2307.09288*, 2023.
- [4] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proc. NAACL-HLT*, 2019, pp. 4171–4186.
- [5] Z. Ji, N. Lee, R. Frieske, T. Yu, D. Su, Y. Xu, E. Ishii, Y. J. Bang, A. Madotto, and P. Fung, "Survey of hallucination in natural language generation," *ACM Computing Surveys*, vol. 55, no. 12, pp. 1–38, 2023.
- [6] Y. Zhang, Y. Li, L. Cui, D. Cai, L. Liu, T. Huang, F. Wei, "A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions," *arXiv preprint arXiv:2311.05232*, 2023.
- [7] O. Evans, O. Stuhlmüller, N. D. Goodman, "Truthful AI: Developing and governing AI that does not lie," *arXiv preprint arXiv:2110.06674*, 2021.
- [8] T. Gao, A. H. Liu, D. Chen, P. C. Pathak, "Reducing Activation Regions Reduces Hallucinations," *arXiv preprint arXiv:2210.06453*, 2022.
- [9] J. Howard and S. Ruder, "Universal language model fine-tuning for text classification," in *Proc. ACL*, 2018, pp. 328–339.
- [10] C. Raffel et al., "Exploring the limits of transfer learning with a unified text-to-text transformer," *J. Mach. Learn. Res.*, vol. 21, no. 140, pp. 1–67, 2020.
- [11] E. J. Hu et al., "LoRA: Low-rank adaptation of large language models," in *Proc. ICLR*, 2022.
- [12] N. Houlsby et al., "Parameter-efficient transfer learning for NLP," in *Proc. ICML*, PMLR, 2019,

- pp. 2790–2799.
- [13] P. Lewis et al., "Retrieval-augmented generation for knowledge-intensive NLP tasks," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2020, vol. 33, pp. 9459–9474.
- [14] V. Karpukhin et al., "Dense passage retrieval for open-domain question answering," in *Proc. EMNLP*, 2020, pp. 6769–6781.
- [15] O. Ram et al., "In-context retrieval-augmented language models," *arXiv preprint arXiv:2302.00083*, 2023.
- [16] O. Khattab and M. Zaharia, "ColBERT: Efficient and effective passage search via contextualized late interaction over BERT," in *Proc. SIGIR*, 2020, pp. 39–48.
- [17] G. Izacard et al., "Unsupervised dense retrieval with contrastive learning," *arXiv preprint arXiv:2112.09118*, 2021.
- [18] Gemini Team, Google, "Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context," *Google DeepMind Blog*, Feb. 2024. [Online]. Available: <https://deepmind.google/technologies/gemini/gemini-1-5/>
- [19] Anthropic, "Introducing the Claude 3 models," *Anthropic Blog*, Mar. 2024. [Online]. Available: <https://www.anthropic.com/news/claude-3-family>
- [20] [Havishya], "A Web-Based Document QA Interaction Prototype developed for students using PCLC-RAG" <https://translite-demo.netlify.app>
- [21] L. Ouyang et al., "Training language models to follow instructions with human feedback," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2022, vol. 35, pp. 27730–27744.
- [22] J. Wei et al., "Chain-of-thought prompting elicits reasoning in large language models," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2022, vol. 35, pp. 24824–24837.
- [23] H. W. Chung et al., "Scaling instruction-finetuned language models," *arXiv preprint arXiv:2210.11416*, 2022.
- [24] N. F. Liu, K. Lin, J. Hewitt, A. Singh, P. Liang, and P. S. H. Dao, "Lost in the middle: How language models use long contexts," *arXiv preprint arXiv:2307.03172*, 2023.
- [25] O. Khattab et al., "DSPy: Compiling Declarative Language Model Calls into Self-Improving Pipelines," *arXiv preprint arXiv:2310.03714*, 2023.
- [26] Y. Zhou et al., "Large Language Models Are Human-Level Prompt Engineers," *arXiv preprint arXiv:2211.01910*, 2022.
- [27] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang, "SQUAD: 100,000+ questions for machine reading comprehension of text," in *Proc. EMNLP*, 2016, pp. 2383–2392.
- [28] T. Kwiatkowski et al., "Natural questions: a benchmark for question answering research," *Trans. Assoc. Comput. Linguist.*, vol. 7, pp. 453–466, 2019.
- [29] A. Kornilova and V. Eidelman, "BillSum: A Corpus for Automatic Summarization of US Legislation," in *Proc. 2nd Workshop on New Frontiers in Summarization (EMNLP)*, 2019, pp. 83–90.
- [30] P. Laban, T. D. Chowdhury, M. H. Shah, K. Kryścin'ski, and R. Socher, "SummaC: Revisiting NLI-based models for inconsistency detection in summarization," *Trans. Assoc. Comput. Linguist.*, vol. 10, pp. 163–177, 2022.
- [31] C.-Y. Lin, "ROUGE: A package for automatic evaluation of summaries," in *Proc. Workshop on Text Summarization Branches Out (ACL)*, 2004, pp. 74–81.
- [32] T. Zhang, V. Kishore, F. Wu, K. Q. Weinberger, and Y. Artzi, "BERTScore: Evaluating text generation with BERT," in *Proc. ICLR*, 2020.