# Evaluating RAG, Reflection Tuning, and Chain of Thought in LLM Translations: Individual and Combined Effects

Ajaykumar Premkumar Nair
B00968276
ajay.nair@dal.ca

Manish Shankar Jadhav
B00969328
mn649712@dal.ca

October 29, 2024

## Abstract

This research investigates the effectiveness of three advanced techniques - Retrieval-Augmented Generation (RAG), Reflection Tuning, and Chain of Thought - in enhancing Large Language Model (LLM) translation capabilities, specifically focusing on conversational French-to-English translation. The study particularly emphasizes challenging aspects of casual French conversation. While these techniques have shown promise individually in various natural language processing tasks, their specific impact on translation quality, both independently and in combination, remains unexplored, especially in preserving the cultural and contextual nuances of informal communication. Our study is to evaluate these techniques with everyday conversational French-to-English translation tasks. The evaluation framework combines automated metrics (BLEU, METEOR, chrF, BERTScore, and COMET) with human assessment criteria including adequacy, fluency, and pragmatic accuracy to provide a comprehensive quality assessment. Special attention will be paid to the preservation of humor and cultural context in the translations. We will assess each technique in isolation to establish individual performance metrics, examine paired combinations to identify potential effects, and analyze the impact of implementing all three techniques simultaneously. The research aims to quantify improvements in translation accuracy and quality against a baseline of standard LLM translations, particularly focusing on how well each approach handles the inherent ambiguity and cultural specificity of conversational French [1]. By comprehensively evaluating these techniques in the context of conversational French-to-English translation, this study seeks to contribute valuable insights to the field of machine translation and provide practical recommendations for optimizing LLM translation systems. The findings will help inform the development of more effective translation architectures and demonstrate whether the computational overhead of implementing multiple enhancement techniques is justified by corresponding improvements in translation quality.

# 1 Problem Statement

This research addresses the critical challenge of accurately translating conversational French to English using Large Language Models (LLMs), with a specific focus on handling informal expressions, cultural nuances, and contextual understanding [2] [3]. Current LLM translations often fail to capture the subtle intricacies of everyday conversation and are not culturaly aware for idiomatic expressions, leading to translations that appear artificial or lose their intended meaning [4]. This problem is significant because informal communication represents a substantial portion of real-world language use, and inaccurate translations can lead to miscommunication and cultural disconnects in increasingly global digital interactions. Additionally, as LLMs become more integrated into communication tools, the need for natural-sounding conversational translations becomes crucial for effective cross-cultural dialogue. The scope of our project encompasses the evaluation of three specific enhancement techniques: Retrieval-Augmented Generation (RAG), Reflection Tuning, and Chain of Thought. We will assess these techniques both independently and in combination, using a comprehensive evaluation framework that includes both automated metrics (BLEU, METEOR, chrF) and human assessment criteria, specifically testing their effectiveness on a curated dataset of conversational French texts that include informal expressions, slang, and cultural references. Our goal is to quantifiably determine which technique or combination of techniques most effectively improves the quality of conversational French-to-English translations while considering the practical implications of implementation complexity and computational requirements .

# 2 Possible Approaches

Our research methodology employs a systematic approach to evaluate the effectiveness of RAG for NLP tasks[5], Reflection Tuning [6], and Chain of Thought [7] in improving French-to-English conversational translations:

- **Phase 1: Individual Technique Evaluation**

    - Establish baseline performance using standard LLM translations.
    - Implement and test each technique independently:
        * **RAG**: Using curated French-English conversational pairs and idioms.
        * **Reflection Tuning**: Implementing self-evaluation and refinement steps.
        * **Chain of Thought**: Applying step-by-step reasoning for complex expressions.
    - Evaluate using automated metrics (BLEU, METEOR, chrF) and human assessment.

- **Phase 2: Combined Analysis**

    - Test technique pairs to identify synergies:
        * RAG + Reflection Tuning
        * RAG + Chain of Thought
        * Reflection Tuning + Chain of Thought
    - Implement all three techniques simultaneously.
    - Analyze potential trade-offs between translation quality and computational overhead.

# 3  Project Plan

## 3.1  Timeline and Major Milestones

1. **Phase 1: Foundation**

   - November 2: Establish baseline system
     - Set up initial metrics
     - Document baseline performance

2. **Phase 2: Core Implementation**

   - November 5: Implement RAG (Retrieval-Augmented Generation)
   - November 9: Integrate Reflection mechanism
   - November 12: Implement Chain of Thought reasoning

3. **Phase 3: Testing**

   - November 13-17: Individual component testing
     - Test RAG performance
     - Evaluate Reflection mechanism
     - Assess Chain of Thought accuracy
   - November 18-20: Combined testing
     - Combining all three methods and testing

4. **Phase 4: Analysis and Documentation**

   - November 21-30
     - Analyze performance of all the runs
     - Compile final results and prepare documentation

# References

[1] P. Chen, J. Tang, and A. Birch, "Evaluating the translation performance of large language models based on euas-20," *arXiv preprint arXiv:2408.03119*, 2024.

[2] H. Namukwaya, "Beyond translating french into english: Experiences of a non-native translator," *TranscUlturAl: A Journal of Translation and Cultural Studies*, vol. 5, p. 61, 03 2014.

[3] Z. Wang, *Study on the Importance of Cultural Context Analysis in Machine Translation*, 05 2013, vol. 212, pp. 29–35.

[4] P. Singh, M. Patidar, and L. Vig, "Translating across cultures: Llms for intralingual cultural adaptation," 2024. [Online]. Available: https://arxiv.org/abs/2406.14504

[5] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W. tau Yih, T. Rocktäschel, S. Riedel, and D. Kiela, "Retrieval-augmented generation for knowledge-intensive nlp tasks," 2021. [Online]. Available: https://arxiv.org/abs/2005.11401

[6] M. Li, L. Chen, J. Chen, S. He, H. Huang, J. Gu, and T. Zhou, "Reflection-tuning: Data recycling improves llm instruction-tuning," 2023. [Online]. Available: https://arxiv.org/abs/2310.11716

[7] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. Chi, Q. Le, and D. Zhou, "Chain-of-thought prompting elicits reasoning in large language models," 2023. [Online]. Available: https://arxiv.org/abs/2201.11903