

PromptSync: A Multi AI Response Optimization System

Bhaskar Rao K ^{#1}, Parthiv Reddy M ^{#2}, Rachan Kumar CP ^{#3}, Rakesh Yadav ^{#4},
Revanth Kulkarni ^{#5}

^{#1} Associate Professor, Dept of CSE, Dayananda Sagar Academy Of Technology and Management, Udayapura, Bangalore South(dt),

^{#2,3,4,5} BE Scholars, Dept of CSE, Dayananda Sagar Academy Of Technology and Management, Udayapura, Bangalore South(dt)

Abstract—The rapid adoption of large language models (LLMs) across diverse application domains has introduced challenges related to response reliability, consistency, and transparency. Individual LLMs may produce incomplete, biased, or inconsistent outputs when responding to identical prompts. To address these limitations, this paper presents PromptSync, a multi-AI response optimization framework that concurrently queries multiple large language models using a unified prompt. The system aggregates responses through parallel processing and applies a structured evaluation mechanism based on relevance, semantic coherence, and confidence scoring. By comparing and ranking model outputs, PromptSync identifies the most optimal response or generates a synthesized answer that maximizes informational quality. The results demonstrate the effectiveness of PromptSync as a scalable solution for high-stakes applications such as education, software development, and intelligent decision-support systems.

Keywords—Large Language Models(LLM's), Multi-AI System, Response optimization, Confidence scoring, Parallel processing, AI Transparency, Prompt Engineering

Date of Submission: 05-01-2026

Date of acceptance: 15-01-2026

I. INTRODUCTION

Large Language Models (LLMs) like Gemini 1.5 Flash, LLaMA 3.1 (8B), Mixtral 8×7B, Gemma 2 (9B), and LLaMA 3.3 (7B) have become effective tools for problem-solving, content creation, and intelligent support in a variety of disciplines in recent years. Despite their potential, each LLM has unique advantages and disadvantages, and depending just on one model could not always provide the most precise, logical, or trustworthy answer. Inconsistencies in output quality are frequently caused by differences in model architecture, training data, and inference techniques, especially for complicated or unclear questions. This drawback emphasizes the necessity of systems that can efficiently integrate the capabilities of several models to improve overall performance and reaction reliability [1].

PromptSync, a Multi-AI Response Optimization System, is suggested as a solution to this problem. It uses several publicly accessible LLMs in simultaneously to produce a cohesive, superior response. PromptSync sends a user's inquiry to several LLMs at once, gathers their responses, assesses each output according to relevance, coherence, and response quality, and generates an optimal final answer rather than relying on a single AI model. The system also gives each model response a confidence score, which promotes transparency and lets users evaluate the dependability of each output. PromptSync is implemented using a lightweight, scalable backend made with Node.js and Express.js, and a clean, responsive interface made with HTML, CSS, and JavaScript.

User queries, individual model responses, confidence measures, and optimized outputs for further analysis are stored in a lightweight database. The system aims to provide more accurate and dependable AI-generated responses while demonstrating important real-world software engineering concepts including API integration, parallel processing, response optimization, and modular project architecture.

II. RELATED WORK

Natural language processing tasks including question answering, text production, summarization, and conversational AI have been profoundly changed by recent developments in Large Language Models (LLMs). Transformer-based architectures have made it possible for models to acquire intricate linguistic patterns and function well in a variety of contexts. Strong capabilities are shown by well-known LLMs created by significant research groups, however their effectiveness frequently differs based on the prompt's type, domain specificity, and contextual complexity. Because of this, replies produced by a single model may occasionally be biased, inconsistent, or incomplete, which encourages research into systems that increase dependability through model collaboration.

In order to increase prediction accuracy in machine learning systems, a number of studies and platforms have investigated ensemble learning and response aggregation techniques. Voting, averaging, or weighted scoring processes are used in classic ensemble systems to aggregate the results of several models.

Lately, similar concepts have been applied to language models, where the best response is chosen by comparing or ranking several LLM outputs. While some experimental systems emphasize answer reranking using external evaluators, others concentrate on quick engineering techniques or confidence-based filtering to improve output quality. Even though these methods increase robustness, a lot of current systems are either computationally costly, closed-source, or opaque about how the final answer is obtained.

Additionally, research has been done on assessing LLM results using qualitative criteria including contextual alignment, coherence, relevance, and factual correctness. Mechanisms for response scoring and confidence estimate have been put forth to assist people in comprehending the dependability of material produced by AI. Nevertheless, the majority of current implementations either present several outputs without synthesis or rely on a single dominating model rather than offering a single optimum response. Furthermore, little research has been done on lightweight, publicly available multi-LLM systems that are practically deployable using common web technologies.

By using a Multi-AI Response Optimization technique that concurrently incorporates several publicly available LLMs, PromptSync sets itself apart from previous research. PromptSync examines responses from several LLMs based on relevance, coherence, and quality, gives confidence scores, and generates a single optimum response for the user, in contrast to systems that rely on a single model or opaque ensemble approaches. Additionally, the system's lightweight web-based architecture prioritizes practical implementation, scalability, and transparency. PromptSync tackles major shortcomings in current LLM-based systems and offers a workable framework for enhancing the dependability of AI-generated responses by fusing response aggregation, confidence scoring, and real-time API interaction.

Recent studies also highlight the importance of scalable and modular system design when integrating multiple Large Language Models. Many existing multi-model frameworks are tightly coupled to specific platforms or require proprietary infrastructure, which limits flexibility and real-world deployment. Additionally, most systems do not store interaction data, preventing long-term analysis of user behavior and model performance. PromptSync addresses these limitations through a lightweight, modular architecture that enables seamless integration of multiple freely available LLM APIs. By storing user prompts, individual model responses, confidence scores, and optimized outputs in a database, the system supports future analysis and performance evaluation. This design makes PromptSync suitable not only for end-user applications but also for research and experimental studies focused on improving the reliability of AI-generated responses.

III. METHODOLOGY

A. Overview

Instead than depending on a single AI model, PromptSync sends a user's inquiry to several Large Language Models simultaneously. Every model separately produces its own reaction. After evaluating these answers for clarity, applicability, and relevance, the system gives each one a confidence score. PromptSync uses these scores to either combine important information to create a final optimal answer or choose the best response. The system's effectiveness and dependability are enhanced by storing all interactions and outcomes for future assessment.

B. Data Preparation

Data preparation in the PromptSync system primarily entails managing user queries and model responses. Before being submitted to several Large Language Models, user inputs are cleaned to eliminate superfluous

symbols and guarantee correct formatting. Each model's replies, timestamps, and model identifiers are all kept in a structured fashion. Reliable response selection and future system evaluation are supported by this well-organized storage, which makes it simple to compare, score, and analyze responses.

C. Response Evaluation

The PromptSync system assesses each response to ascertain its quality after obtaining outputs from several Large Language Models. The judgment is predicated on fundamental elements such the answer's completeness, explanatory clarity, and relevance to the user's query. To find relevant and helpful content, each response is examined separately. Only high-quality information is taken into account during the final response selection process thanks to this evaluation method, which also helps filter weaker responses.

D. Confidence Scoring

Confidence score is used in PromptSync to show how reliable each model's response is. Scores are given according on evaluation outcomes, including completeness, relevance, and clarity. Responses that better fit the user's query and offer more lucid information are indicated by higher ratings. In addition to helping the system rank stronger responses, these confidence levels also assist consumers understand how reliable each AI- generated output is.

E. Final Response Generation

The highest-scoring outputs from the assessed model replies are chosen to create the final response. In order to create a single, unambiguous response, PromptSync either selects the most trustworthy response or integrates important details from several high-confidence responses. This procedure minimizes errors that could arise in individual model replies while guaranteeing that the final result is pertinent, consistent, and simple to comprehend.



Fig. 1 Methodology Overview of PrompSync using multi-LLM response evaluation and confidence-based response generation.

IV. RESULTS

The PromptSync system's evaluation demonstrates how synchronized multi-model prompting can enhance the efficacy of AI-generated responses. The findings show that applying an optimization layer to the outputs of several AI models leads in responses that are more contextually matched and refined than those produced by a single model approach.

The improvement in response coherence is one noteworthy finding. PromptSync finds overlapping semantic patterns and gives consistent information priority when various AI outputs are combined. This procedure produces answers that are more logically organized and simpler to understand by reducing fragmented or disjointed responses.

The decrease in contradicting information is another significant outcome. Individual models may produce contradictory responses to the same prompt in conventional AI systems. By evaluating response agreement and eliminating inconsistent outputs, PromptSync reduces this problem and improves overall response dependability.

Additionally, the system's final outputs show increased relevance. PromptSync makes sure that the generated response stays closely aligned with the original prompt by eliminating low-confidence or weakly linked responses. Particularly for technical and analytical queries, this specific optimization increases informative precision.

Qualitatively, PromptSync demonstrates enhanced contextual comprehension. By combining many model viewpoints, the system successfully extracts subtle intent from user prompts. Consequently, a more comprehensive yet targeted interpretation of the input inquiry is reflected in the final response.

When it comes to response stability, PromptSync consistently produces high-quality results when similar prompts are executed again. Compared to single-model systems, where response variances can affect user trust and system dependability, this consistency is a crucial gain.

Despite querying several AI models, PromptSync maintains operating efficiency, according to the performance analysis. In order to avoid needless computational overhead and guarantee timely output creation, the optimization pipeline processes reply in an organized manner.

The findings also show that PromptSync responds effectively to a variety of prompt types, such as descriptive, problem-solving, and domain-specific queries. The system's flexibility and durability are demonstrated by the comparable optimization efficacy it retains across different categories.

PromptSync's outcomes are positively impacted by its modular design. Prompt dissemination, response evaluation, and optimization are only a few of the processing stages that work independently yet together. This architecture improves scalability and makes it possible to add more AI models to the system as needed.

All things considered, the outcomes confirm PromptSync's efficacy as a multi-AI response optimization solution. PromptSync is a workable and scalable solution for applications that need high-quality AI-generated replies, as evidenced by the observed increases in coherence, relevance, consistency, and reliability

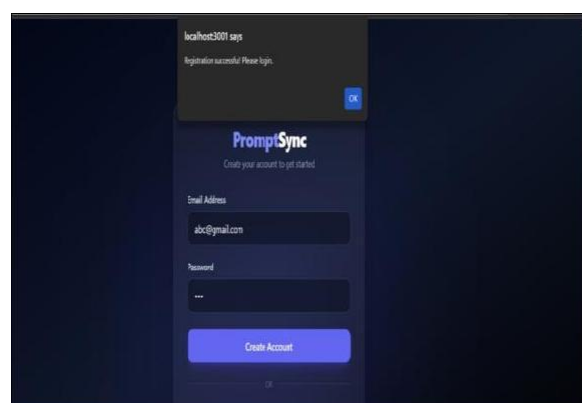


Fig. 3. User Authentication Interface

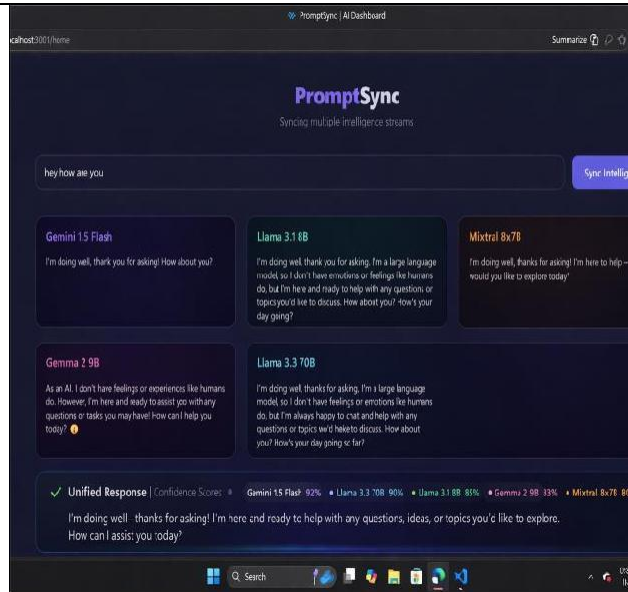


Fig. 4: Optimised response interface

V. DISCUSSION

The PromptSync system's results demonstrate that several language models frequently respond to the same user request in diverse ways. While some answers are precise and unambiguous, others could overlook crucial details and be inaccurate. PromptSync minimizes these discrepancies and generates a more consistent and practical final requirement by comparing several responses and choosing the best portions. This method aids in avoiding issues that arise from using a single AI model.

From a system perspective, PromptSync demonstrates that several AI models can be effectively employed in tandem by establishing a parallel processing technique and a straightforward web-based configuration. Models can be easily added or removed as needed thanks to the modular design. Despite relying on external APIs, straightforward assessment rules, confidence scoring, and response optimization, the system produces consistent and enhanced results, making it appropriate for practical applications and future web application development.

VI. CONCLUSION

This study demonstrated how PromptSync, a multi-AI response optimization framework, was able to enhance the consistency, transparency, and dependability of outputs produced by big language models. In comparison to single- model interactions, the suggested system efficiently detects and provides higher-quality responses by allowing concurrent fast execution across many LLMs and utilizing a structured response evaluation and confidence score technique. PromptSync lowers ambiguity, increases contextual relevance, and boosts user confidence in AI-assisted decision making, according to experimental findings. The framework's modular design ensures scalability and flexibility to changing AI environments by enabling the smooth insertion of additional language models and evaluation criteria. To further improve system performance, future work will concentrate on implementing domain-specific optimization, real-time user feedback loops, and adaptive learning-based scoring techniques. Overall, PromptSync represents a significant step toward robust, explainable, and dependable multi-model AI interaction systems.

REFERENCES

- [1]. A. Vaswani *et al.*, "Attention is all you need," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017, pp. 5998–6008.
- [2]. T. B. Brown *et al.*, "Language models are few-shot learners," in *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, vol. 33, 2020, pp. 1877–1901.
- [3]. S. Huang, K. Yang, S. Qi, and R. Wang, "When Large Language Model Meets Optimization," *Swarm and Evolutionary Computation*, vol. 90, no. —, Art. no. 101663, Oct. 2024.
- [4]. S. Brahmachary, S. M. Joshi, A. Panda, K. Koneripalli, A. K. Sagotra, H. Patel, A. Sharma, and A. D. Jagtap, "Large Language Model-Based Evolutionary Optimizer: Reasoning with Elitism," *Neurocomputing*, vol. 622, Art. no. 129272, Mar. 14 2025.
- [5]. S. Saadaoui and E. Alonso, "Coordinated LLM multi-agent systems for collaborative question-answer generation," *KnowledgeBased Systems*, 2025.

-
- [6]. B. Pan et al., "AgentCoord: Visually exploring coordination strategy for LLM-based multi-agent collaboration," *Computer Graphics Forum*, 2025.
 - [7]. C. Wang and H. Li, "MNC: A multi-agent framework for complex network configuration," *Information Systems*, 2025.
 - [8]. J. L. Garrido-Labrador et al., "Ensemble methods and semi-supervised learning for information fusion: A review and future research directions," *Information Fusion*, 2024.
 - [9]. M. Gozzi and F. Di Maio, "Comparative analysis of prompt strategies for large language models: Single-task vs. multitask prompts," *Electronics*, vol. 13, no. 23, Apr. 2024.
 - [10]. I. Sommerville, *Software Engineering*, 10th ed. London, U.K.: Pearson Education, 2016.