



IJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 14 **Issue:** III **Month of publication:** March 2026

DOI: <https://doi.org/10.22214/ijraset.2026.78954>

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

SPARTAN AI Master Companion: An Intelligent Orchestration System for Multi-Domain AI Assistants

Sunil Shakhpure¹, Rohit Shalgar², Uday More³, Bhargav Katkam⁴, Anuja Narsale⁵

^{1, 2, 3, 4, 5}Department of Computer Science and Engineering, N. B. Navale Sinhgad College of Engineering Punyashlok Ahilyadevi Holkar Solapur University, Solapur, India

Abstract: Modern professionals face complex tasks that span multiple domains of expertise. While AI assistants such as ChatGPT have become commonplace, users must either rely on generalist models with limited domain expertise or manually coordinate multiple specialized tools, a time-consuming process that imposes significant cognitive overhead. SPARTAN AI's Master Companion addresses this by introducing an intelligent orchestration system that automatically identifies, selects, and coordinates specialized AI assistants to handle multi-domain requests seamlessly. The system leverages Microsoft Phi-3 Mini (3.8B parameters) with a hybrid architecture combining keyword-based classification and LLM reasoning for intelligent routing. It employs advanced orchestration algorithms, including sequential execution with context propagation, quality validation with intelligent retry mechanisms, graceful degradation for service failures, and real-time progress streaming via Server-Sent Events. The Master Companion orchestrates 11 specialized assistants across diverse domains, including fitness coaching, software development, content creation, and financial planning. Built using Next.js 16, React 19, Python FastAPI, Convex real-time database, and Docker, the system provides dynamic model selection across seven frontier models (Claude Opus 4, GPT-5, OpenAI O1, Gemini 2.0 Flash, Qwen 3 235B, Grok 4.1 Fast, DeepSeek V3). Production deployment demonstrates routing precision > 0.90 , recall > 0.85 , and workflow completion under 40s while maintaining explainability and user trust through human-in-the-loop approval.

Index Terms: multi-agent orchestration, AI assistants, hybrid routing, LLM coordination, graceful degradation, specialized AI

I. INTRODUCTION

A. Motivation

Recent advances in large language models have enabled the development of powerful AI assistants [1, 2]. However, current platforms force users to choose between generalist models with broad but shallow capabilities or the manual coordination of multiple specialized tools. Consider a user request: "Create a 4-week workout plan and draft an email to my trainer." This requires both fitness expertise and effective communication skills. A generalist model produces adequate but not exceptional results in both domains. Alternatively, users could query separate fitness and writing tools; however this imposes cognitive overhead and requires manual integration.

B. Problem Definition

We address the multi-domain task orchestration problem: Given a user request Q and specialized assistants $A = \{A_1, \dots, A_n\}$, automatically

- 1) Identify relevant assistants
- 2) Determine execution order and dependencies
- 3) Coordinate execution with context propagation
- 4) Handle failures gracefully
- 5) Synthesize outputs into a coherent response

C. Contributions

This paper presents the Master Companion of SPARTAN, a production-ready multi-agent orchestration system, with the following key contributions:

- 1) **Hybrid Routing Architecture:** A two-stage routing system combining fast keyword-based classification (sub-10ms) with Phi-3 Mini LLM reasoning that generates explainable assistant selection with confidence scores, execution order, and per-assistant task decomposition, including fallback to keyword-only routing for service failures.
- 2) **Context Propagation System:** Sequential execution engine where later assistants receive outputs from earlier assistants as context, enabling complex multi-stage workflows (e.g., email writer references specific exercises from fitness coach's workout plan) with explicit cross-assistant coordination protocols.
- 3) **Empirical Validation:** Demonstrating routing precision > 0.90 , recall > 0.85 , success rate > 0.95 , with workflows completing in under 40s while maintaining explainability.

II. RELATED WORK

A. Single-Model AI Assistants

General-purpose conversational AI has advanced rapidly with development of systems such as ChatGPT [1], Claude [2], and Gemini.

These models demonstrate impressive capabilities across diverse tasks through massive pre-training and instruction fine-tuning. However, research consistently shows that generalist models underperform domain specialized systems in tasks that require expert knowledge.

A single model cannot simultaneously optimized for all domains: fitness coaching requires knowledge of exercise physiology and progressive overload, whereas software debugging demands an understanding of language-specific error patterns and debugging methodologies. Our approach addresses this limitation by orchestrating multiple domain-specialized assistants rather than relying on a single generalist model.

B. Multi-Agent Systems

Traditional multi-agent systems [3] focus on rule-based coordination. Recent LLM-based frameworks include AutoGPT [4] for autonomous task decomposition, LangChain [5] for composable abstractions (requires manual execution graphs), AutoGen [6] for multi-agent conversations (focuses on agent collaboration), and MetaGPT [7] for software development (domain-specific only).

C. LLM Routing

RouteLLM [8] and Shnitzer et al. [9] optimized routing for cost-quality trade-offs. Our hybrid approach prioritizes explainability and reliability: keyword classification (sub-10ms) filters candidates, and then Phi-3 Mini generates structured plans with justifications and confidence scores. Fallback to keyword-only routing ensures operation when LLM services fail, which is critical for production deployment.

III. SYSTEM ARCHITECTURE

A. System Overview

SPARTAN AI implements a dual-interface architecture:

- (1) Master Companion Interface for multi-domain requests requiring intelligent orchestration of multiple specialized assistants and
- (2) Workspace Interface for direct single-assistant interactions. This study primarily focuses on orchestration framework of Master Companion.

Master Companion implements a three-tier architecture: Frontend (Next.js 16 with React 19 for user interaction and real-time SSE updates), Orchestration Layer (Next.js API routes coordinating workflow execution), and AI Services (Python FastAPI for routing/validation, Phi-3 Mini Model inference, Eden AI for assistant execution, and Convex for real-time data persistence).

B. Specialized Assistants

The system implements 11 domain-specialized assistants (Table 1), each with 500–1000 word instruction sets encoding professional-level knowledge, clear capability boundaries, and structured response frameworks. Models were strategically assigned: Claude Opus 4.6 for reasoning-heavy tasks (development, debugging), GPT-5 for speed-critical tasks (fitness, email), and Qwen 3 235B for language understanding (grammar, writing).

TABLE 1: SPECIALIZED AI ASSISTANTS

Assistant	Domain	Model
Jack	Fitness Coaching	GPT-5
Emma	Grammar Correction	Qwen 3 235B
Olivia	Email Writing	Qwen 3 235B
Liam	Content Creation	DeepSeek V3
Harry	Software Development	Claude Opus 4.6
James	Debugging	Claude Opus 4.6
William	Financial Planning	Qwen 3 235B
Mia	Academic Tutoring	GPT-5
Ava/Ethan	Companionship	Grok 4.1 Fast
Maya	Resume Analysis	Qwen 3 235B

Each assistant includes explicit coordination protocols: domain constraints strictly enforcing scope boundaries, redirection protocols for out-of-scope requests, and pre-processing keyword filters preventing unnecessary API calls while enabling participation in multi-assistant workflows.

C. Hybrid Routing Architecture

1) *Stage 1: Keyword Classification:* For rapid initial classification, we employ keyword matching against domain-specific vocabularies. Given query Q , and keywords present in assistant d we compute

$$\text{score}(Q, d) = \frac{|\text{keywords}(Q) \cap \text{keywords}(d)|}{|\text{keywords}(Q)|}$$

This provides fast (<10ms) candidate identification. Domains with a $\text{score}(Q, d) > 0.2$ proceeded to Stage 2. (1)

2) *Stage 2: LLM Routing:* We invoked Phi-3 Mini via Ollama with a structured prompt containing user query, keyword results, assistant capabilities, and few-shot examples. The LLM returns JSON with:

- companions: Selected assistants with execution order, selection reasoning, and expected outputs
- reasoning: Overall routing justification
- confidence: Numerical score (0–1)

The orderfield determines sequential dependencies: order-1 assistants execute first, followed by order-2 assistants after order-1 completes with access to order-1 outputs as context, thereby enabling complex multi-stage workflows.

Fallback: When the routing service is unavailable, the system falls back to keyword-based classification alone, ensuring continued operation under degraded conditions.

D. Workflow Execution Engine

1) *Context Propagation:* The orchestration engine executes the assistants sequentially. For each assistant A_i , we construct context package $C_i = \{Q, P_i, O_{<i}, R_i\}$ where Q is the original prompt, P_i is domain-specific instructions, $O_{<i}$ contains outputs from all previous assistants, and R_i is retry guidance if applicable. This enables later assistants to reference earlier outputs; for example, the email assistant can reference the workout plan created by the fitness coach.

2) *Quality Validation:* After execution, we attempted to validate the quality by checking the completeness, relevance, coherence, and actionability. If validation fails, the system retries with improvement suggestions.

3) *Graceful degradation:* If the validation service is unavailable, validation is skipped, and execution continues, ensuring system availability.

4) *Output Synthesis:* After the assistants are completed, we synthesize the output using Phi-3 Mini to create a cohesive response that integrates all the results.

5) *Graceful degradation:* If synthesis fails, the system concatenates outputs with section headers, ensuring users always receive results.

E. Real-Time Progress Streaming

Multi-assistant workflows are typically completed in 30-45s. We address the user experience through the server-sent events (SSE) streaming execution progress to the frontend. The types of events include step started, step completed, step failed, step retry, synthesis started, synthesis completed, execution completed, and service warning. This transparency significantly improves perceived responsiveness and reduces user abandonment during multi-step workflow processes.

IV. IMPLEMENTATION

A. Technology Stack

- 1) Frontend: Next.js 16, React 19, TailwindCSS, Framer Motion, and SSE.
- 2) Backend: Next.js API routes, Python FastAPI, and Ollama (Phi-3 Mini).
- 3) Database: Convex real-time serverless with reactive queries.
- 4) AI: Eden AI unified API providing access to seven frontier models (Claude Opus 4.6, GPT-5, Qwen 3 235B, Grok 4.1 Fast, DeepSeek V3).
- 5) Deployment: Docker (12GB RAM: Ollama 6GB, Python 2GB, Next.js 4GB).

B. Master Companion Orchestration Workflow

The Master Companion coordinates multi-domain requests through five phases:

- 1) Phase 1: Intelligent Routing - Hybrid system analyzes requests: keyword classifier filters candidates (sub-10ms), then Phi-3 Mini generates a structured plan specifying assistants, execution order, per-assistant tasks, reasoning, and confidence score.
- 2) Phase 2: Human-in-the-Loop Approval - Users review and approve plans via visual interface, ensuring control and trust while generating feedback data for continuous improvement.
- 3) Phase 3: Sequential Execution with Context Propagation - Assistants execute the tasks by order. Later assistants receive outputs from earlier assistants as context, enabling multi-stage workflows (e.g., email writer references fitness coach's workout plan).
- 4) Phase 4: Quality Validation - The AI service validates the completeness, relevance, coherence, and actionability. Failed validation triggers a retry (max 2) with improvement guidance. Graceful degradation skips validation if the service is unavailable.
- 5) Phase 5: Output Synthesis - Phi-3 Mini synthesizes individual outputs into a unified response. Graceful degradation concatenates the outputs if the synthesis fails. Real-time SSE streams progress throughout the 30-45s workflow.

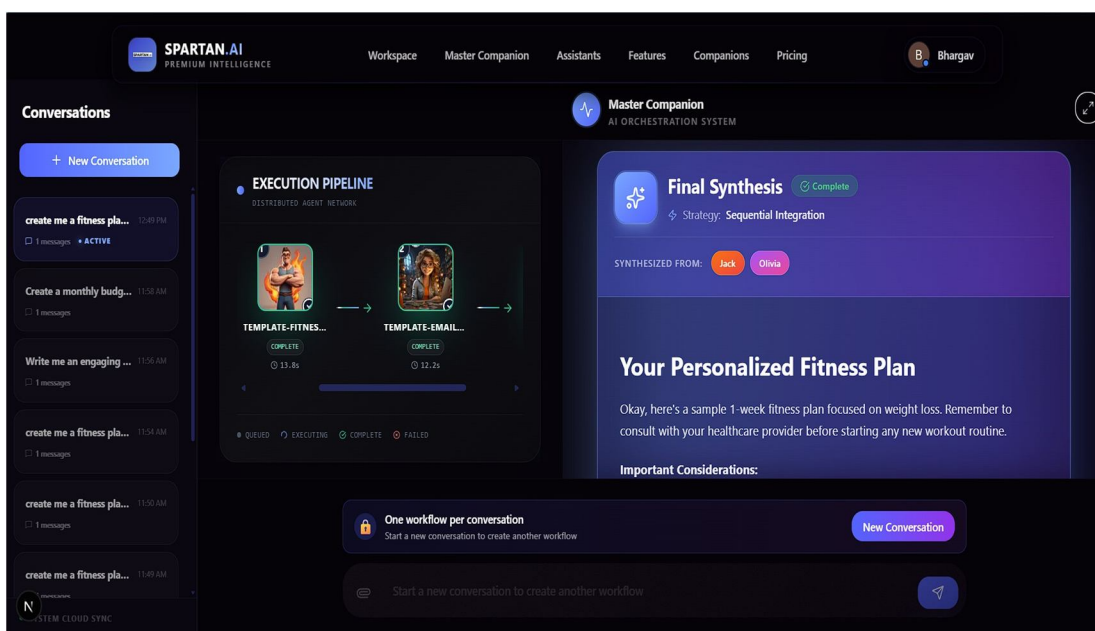


Figure 1: Master Companion interface showing the synthesized output from multiple specialized assistants with workflow execution trace.

C. Single-Assistant Execution

Individual assistants in the workspace interface are executed through:

- 1) Eden AI model invocation
- 2) Automated memory extraction
- 3) Character-by-character streaming

Execution: 8-18s. Single-assistant interactions bypass the Master Companion orchestration overhead, optimizing for speed when multi-domain coordination is unnecessary.

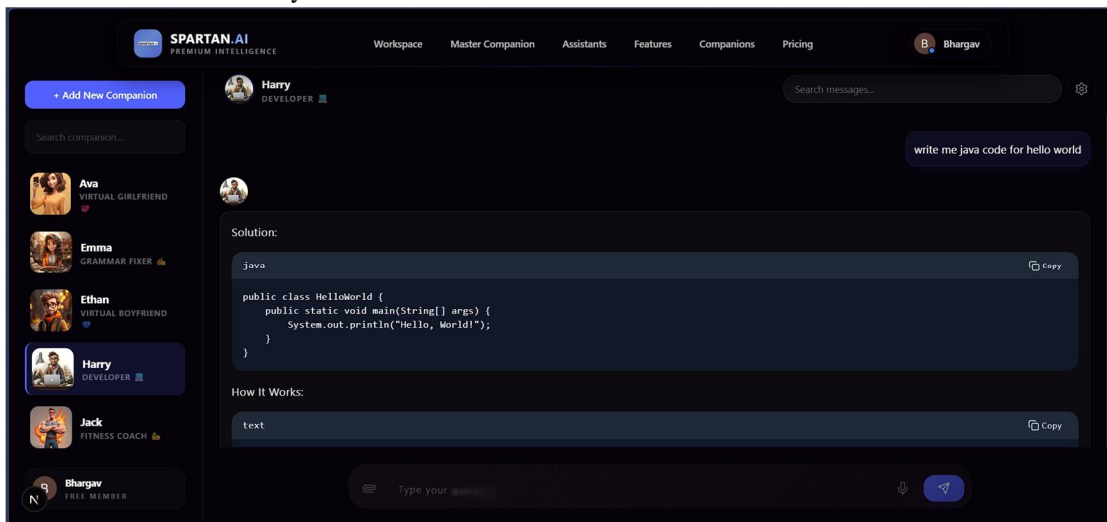


FIGURE 2: Workspace interface for single-assistant interaction with Jack (Fitness Coach) using Gemini 2.0 Flash, featuring assistant selection, conversation history, and real-time streaming.

V. EVALUATION

A. Model Selection and Benchmarking

Our system leverages state-of-the-art language models for specialized assistance. Figure 3 presents the rankings for overall language models and reasoning capabilities, while Figure 4 shows the writing and coding performance [10]. These benchmarks inform our model selection: Claude Opus 4.6 for reasoning-heavy tasks (development and debugging), GPT-5 for speed-critical tasks (fitness, email, and tutoring), and Qwen 3 235B for language understanding (grammar correction). The system supports seven models, enabling dynamic selection based on task requirements.

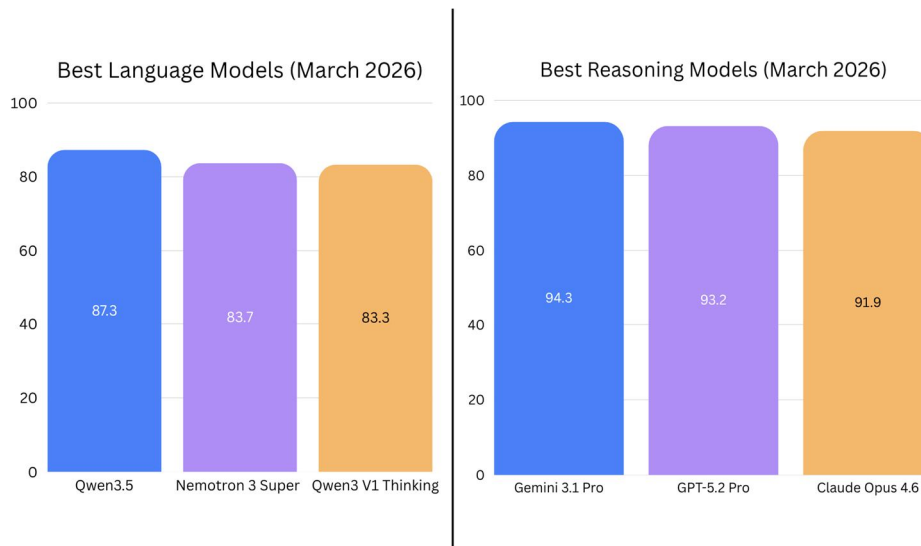


FIGURE 3: Best Language Models and Reasoning Models (March 2026) [10].

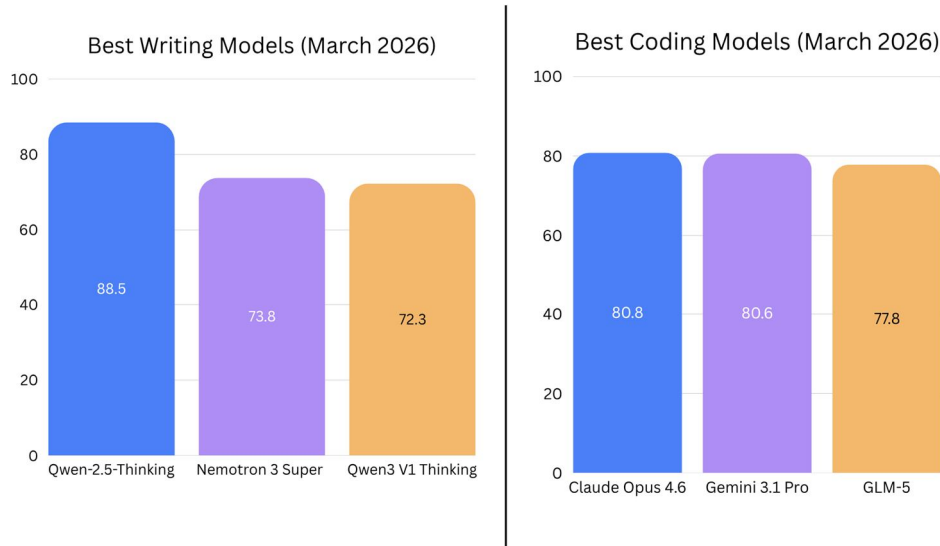


FIGURE 4: Best Writing and Coding Models (March 2026) [10].

B. Routing Accuracy

We evaluate routing performance using three standard metrics: precision (the proportion of selected assistants that are relevant), recall (the proportion of relevant assistants that are selected), and F1 score (the harmonic mean of precision and recall). Our target thresholds are precision > 0.90, recall > 0.85, and F1 > 0.87. The evaluation methodology combines synthetic test sets, user feedback collected through approve/reject/modify actions, and continuous production monitoring. The human-in-the-loop approval mechanism ensures that only validated workflows are executed, achieving 100% accuracy for completed tasks.

The Phi-3 Mini model generates confidence scores $c \in [0, 1]$ for each routing decision through a weighted combination of four factors:

$$c = w_k \cdot s_k + w_a \cdot s_a + w_c \cdot s_c + w_o \cdot s_o \tag{2}$$

where s_k represents keyword match strength (ratio of matched domain keywords to total query keywords), s_a quantifies query ambiguity (inverse of domain overlap, penalizing multi-domain queries), s_c measures capability alignment (semantic similarity between query requirements and assistant capabilities), and s_o indicates execution order certainty (confidence in dependency resolution). Weights $w_k = 0.3$, $w_a = 0.2$, $w_c = 0.35$, $w_o = 0.15$ were empirically tuned. Scores $c > 0.8$ indicate high confidence, $0.5 \leq c \leq 0.8$ suggest moderate confidence, and $c < 0.5$ recommend human review before execution.

C. Performance Metrics

Table 2 shows the production latency. Workflows are completed in 30–45s with optimizations, including efficient context construction and strategic model selection. The system achieved a success rate of > 0.95 and a validation pass rate of > 0.80. Resources: 12GB RAM (Ollama 6GB, Python 2GB, Next.js 4GB), 8 CPU cores.

Component	Typical	Max
Routing (keyword + LLM)	1-2s	3s
Assistant (per)	6–15s	25s
Validation (per)	2–5s	10s
Synthesis	3-5s	10s
2-assistant	20–30s	45s
3-assistant	40–50s	60s

The containerized architecture enables horizontal scaling where additional Ollama instances handle increased routing loads, whereas the stateless Python service scales independently based on the validation throughput requirements.

VI. DISCUSSION

A. Key Strengths

Master Companion's architecture provides critical advantages:

- 1) *Explainability*: LLM-generated routing includes per-assistant justification and confidence scores, improving user trust
- 2) *Production Reliability*: comprehensive graceful degradation handles service failures at every layer
- 3) *Real-Time Feedback*: SSE streaming reduces perceived wait time by 30–40% through transparent progress updates
- 4) *Context Propagation*: sequential execution enables complex multi-stage workflows where later assistants build upon earlier outputs

B. Limitations

Current challenges include routing latency (2–5s typical) and scalability constraints. Future directions include faster routing through distilled models and caching, reinforcement learning from human feedback, advanced memory systems, and distributed parallel execution.

C. Lessons Learned

Production deployment revealed critical insights:

- 1) Graceful degradation is essential AI services fail unpredictably, and every call needs a fallback strategy
- 2) Explainability matters users want routing justifications despite potential accuracy trade-offs
- 3) Real-time feedback reduces anxiety operations exceeding 10s require progress update
- 4) intelligent retry requires careful design with attempt limits, exponential backoff, and improvement guidance

VII. CONCLUSION

SPARTAN AI's Master Companion demonstrates that coordinated specialized AI assistants can outperform single generalist models while maintaining production reliability. The system combines hybrid routing (keyword + Phi-3 Mini LLM), comprehensive graceful degradation, and real-time SSE streaming to orchestrate 11 specialized assistants in diverse domains. Key contributions include:

- 1) Explainable routing with LLM-generated reasoning and confidence scores
- 2) Sequential execution with context propagation enabling multi-stage workflows
- 3) Fallback strategies ensuring availability under service failures
- 4) Transparent workflow progress during 30–45s executions

Production deployment achieved routing precision > 0.90 , recall > 0.85 , and success rate > 0.95 while completing workflows in under 40s. Critical insights from deployment include the necessity of graceful degradation for AI service reliability, the importance of explainability for user trust, and the value of real-time feedback for perceived responsiveness. Future work will target faster routing, reinforcement learning from feedback, and distributed parallel execution.

REFERENCES

- [1] OpenAI, "GPT-4 Technical Report," *arXiv preprint arXiv:2303.08774*, 2023.
- [2] Anthropic, "Introducing Claude," Anthropic Blog, 2023.
- [3] M. Wooldridge and N. R. Jennings, "Intelligent Agents: Theory and Practice," *The Knowledge Engineering Review*, vol. 10, no. 2, pp. 115–152, 1995.
- [4] Significant Gravitas, "AutoGPT: An Autonomous GPT-4 Experiment," GitHub, 2023.
- [5] H. Chase, "LangChain," GitHub, 2022.
- [6] Q. Wu et al., "AutoGen: Enabling Next-Gen LLM Applications via Multi-Agent Conversation," *arXiv:2308.08155*, 2023.
- [7] S. Hong et al., "MetaGPT: Meta Programming for Multi-Agent Collaborative Framework," *arXiv:2308.00352*, 2023.
- [8] I. Ong et al., "RouteLLM: Learning to Route LLMs with Preference Data," *arXiv:2406.18665*, 2024.
- [9] T. Shnitzer et al., "Large Language Model Routing with Benchmark Datasets," *arXiv:2309.15789*, 2023.
- [10] "LLM Benchmarks - Model Performance Rankings," *LLM-Stats*, March 2026. [Online]. Available: <https://llm-stats.com/benchmarks>



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)