

# HASHIRU: Hierarchical Agent System for Hybrid Intelligent Resource Utilization

Kunal Pai\*  
UC Davis  
kunpai@ucdavis.edu

Parth Shah\*  
Independent Researcher  
helloparthshah@gmail.com

Harshil Patel  
UC Davis  
hpppatel@ucdavis.edu

Saisha Shetty  
UC Davis  
spshetty@ucdavis.edu

**Abstract**—Rapid Large Language Model (LLM) advancements are fueling autonomous Multi-Agent System (MAS) development. However, current frameworks often lack flexibility, resource awareness, model diversity, and autonomous tool creation. This paper introduces HASHIRU (Hierarchical Agent System for Hybrid Intelligent Resource Utilization), a novel MAS framework enhancing flexibility, resource efficiency, and adaptability. HASHIRU features a “CEO” agent dynamically managing specialized “employee” agents, instantiated based on task needs and resource constraints (cost, memory). Its hybrid intelligence prioritizes smaller, local LLMs (often via Ollama) while flexibly using external APIs and larger models when necessary. An economic model with hiring/firing costs promotes team stability and efficient resource allocation. The system also includes autonomous API tool creation and a memory function. Evaluations on tasks like academic paper review (58% success), safety assessments (100% on a JailbreakBench subset), and complex reasoning (outperforming Gemini 2.0 Flash on GSM8K: 96% vs. 61%; JEEBench: 80% vs. 68.3%; SVAMP: 92% vs. 84%) demonstrate HASHIRU’s capabilities. Case studies illustrate its self-improvement via autonomous cost model generation, tool integration, and budget management. HASHIRU offers a promising approach for more robust, efficient, and adaptable MAS through dynamic hierarchical control, resource-aware hybrid intelligence, and autonomous functional extension. Source code and benchmarks are available at <https://github.com/HASHIRU-AI/HASHIRU> and <https://github.com/HASHIRU-AI/HASHIRUBench>, respectively.

## I. INTRODUCTION

Rapid Large Language Model (LLM) advancements are reshaping AI, enabling complex language understanding, generation, reasoning, and planning [6], [13], [51]. This progress fuels the development of autonomous Multi-Agent Systems (MAS) where collaborative teams tackle problems beyond individual agent capabilities [14], [64]. Collaborative MAS show potential in scientific discovery [4], software engineering [48], data analysis, and decision-making [61]. The increasing complexity of tasks, evidenced by benchmarks requiring advanced reasoning (e.g., GSM8K [10], SVAMP [45]), coding [8], [42], and graduate-level knowledge [47], necessitates agentic systems that effectively coordinate diverse cognitive resources [63].

Despite this potential, contemporary agentic frameworks exhibit limitations: **rigidity** due to predefined roles hindering adaptation [68]; **resource obliviousness**, lacking mechanisms to optimize computational resources (API costs, memory,

CPU), leading to inefficiency, especially with costly proprietary LLMs [44]; **model homogeneity**, defaulting to a single powerful LLM and missing efficiency gains from diverse, smaller, or local models [69]; and limited autonomous **tool creation and integration**, restricting dynamic self-improvement [43], [58], [67].

To address these challenges, we introduce **HASHIRU (Hierarchical Agent System for Hybrid Intelligent Resource Utilization)**, a novel MAS framework enhancing flexibility, resource efficiency, and adaptability. HASHIRU uses a hierarchical structure with a “CEO” agent dynamically managing specialized “employee” agents, instantiated on demand. Its **hybrid intelligence** strategically prioritizes smaller, local LLMs (e.g., 3B–7B, often via Ollama [36]) for cost-effectiveness, flexibly integrating external APIs and larger models when justified by task complexity and resource availability under CEO management.

The primary contributions are:

- 1) A novel MAS architecture with **hierarchical control** and **dynamic, resource-aware agent lifecycle management** (hiring/firing) governed by budget constraints (cost, memory) and an economic model discouraging excessive churn.
- 2) A **hybrid intelligence model** prioritizing cost-effective, local LLMs while adaptively incorporating external APIs and larger models, optimizing the efficiency-capability trade-off.
- 3) Integrated **autonomous API tool creation** for dynamic functional extension.
- 4) An **economic model** (hiring/firing fees) promoting efficient resource allocation and team stability.

This paper details HASHIRU’s design. Section II discusses related work. Section III elaborates on the architecture. Section IV presents case studies demonstrating self-improvement capabilities. Section V describes the experimental setup and evaluation metrics. Section VI reports results, and Section VII concludes with limitations and future work.

## II. BACKGROUND AND RELATED WORK

Intelligent agent concepts have evolved from early symbolic AI [53], [55] to LLM-dominated frameworks leveraging models for reasoning, planning, and interaction [59], [66]. HASHIRU builds on this, addressing current limitations.

\*These authors contributed equally to this work.

**Agent Architectures:** MAS architectures vary, including flat, federated, and hierarchical [14], [25]. Hierarchical models offer clear control and task decomposition but risk bottlenecks and rigidity [15], [16]. HASHIRU uses a **CEO-Employee hierarchy** for centralized coordination but distinguishes itself through **dynamic team composition**. Unlike systems with static hierarchies or predefined roles (e.g., CrewAI [11], ChatDev [48]), HASHIRU’s CEO dynamically manages the employee pool based on runtime needs and resource constraints.

**Dynamic Agent Lifecycle Management:** Dynamic MAS composition is crucial for complex environments [34]. Agent creation/deletion triggers often relate to task structure or environmental changes. HASHIRU introduces a specific mechanism where the CEO makes **hiring and firing decisions** based on a cost-benefit analysis considering agent performance, operational costs (API fees, inferred compute), memory footprint (tracked explicitly as a percentage of available resources), and concurrency limits. HASHIRU also incorporates an **economic model** with explicit “starting bonus” (hiring) and “invocation” (usage) costs. This economic friction aims to prevent excessive initialization or usage for marginal gains and promote team stability, a nuance often missing in simpler dynamic strategies.

**Resource Management and Agent Economies:** Resource awareness is critical for scalable MAS. Economic research explores mechanisms like market-based auctions or contract nets for allocation [9]. HASHIRU implements a more **centralized, budget-constrained resource management model**. The CEO operates within defined limits for financial cost, memory usage (as a percentage of total allocated), and concurrent agent count. This direct management, particularly focusing on memory percentage, suggests practicality for deployment on local or edge devices with finite resources, contrasting with cloud systems assuming elastic resources [44]. Frameworks like AutoGen [65] and LangGraph [29] typically rely on implicit cost tracking without explicit multi-dimensional budgeting and control.

**Hybrid Intelligence and Heterogeneous Models:** Leveraging diverse LLMs with varying capabilities, costs, and latencies is an emerging trend [69]. Techniques like model routing select optimal models for sub-tasks. HASHIRU embraces **model heterogeneity** with a strategic focus: **prioritizing smaller (3B–7B), locally-run models via Ollama integration** [36]. This emphasizes cost-efficiency, low latency, and potential privacy over systems defaulting to large proprietary cloud APIs (e.g., GPT-4 [38], Claude 3 [1]). While integrating external APIs (potentially larger models), HASHIRU’s default stance represents a distinct capability vs. efficiency balance.

**Tool Use and Autonomous Tool Creation:** Tool use (APIs, functions) is fundamental for modern agents [37], [67]. Most systems use predefined tools. HASHIRU advances this with **integrated, autonomous API tool creation**. When needed functionality is missing, the CEO can commission the generation (potentially via a specialized agent) and deployment of a new API tool within the HASHIRU ecosystem. This self-extension capability differentiates HASHIRU from systems limited to static toolsets, moving towards greater autonomy

and adaptability [44], [58].

In summary, HASHIRU integrates hierarchical control, dynamic MAS, resource management, and tool use. Its novelty lies in the synergistic combination of: (1) dynamic, resource-aware hierarchical management with (2) an economic model for stability, (3) a local-first hybrid intelligence strategy, and (4) integrated autonomous tool creation. This targets key limitations in current systems regarding efficiency, adaptability, cost, and autonomy.

### III. HASHIRU SYSTEM ARCHITECTURE

HASHIRU’s architecture addresses rigidity, resource obliviousness, and limited adaptability through a hierarchical, dynamically managed MAS optimized for hybrid resource utilization.

#### A. Overview

HASHIRU operates with a central “CEO” agent coordinating specialized “Employees”. Key tenets:

- **Dynamic Hierarchical Coordination:** CEO manages strategy, task allocation, and dynamic team composition.
- **Dynamic Lifecycle Management:** Employees are hired/fired based on runtime needs and resource constraints, governed by an economic model.
- **Hybrid Intelligence:** Strategic preference for LLMs within a predefined budget, while accessing external APIs/models.
- **Explicit Resource Management:** Continuous monitoring and control of costs against budgets.
- **Adaptive Tooling:** Using predefined tools alongside autonomous creation of new API tools.

Figure 1 illustrates the structure.

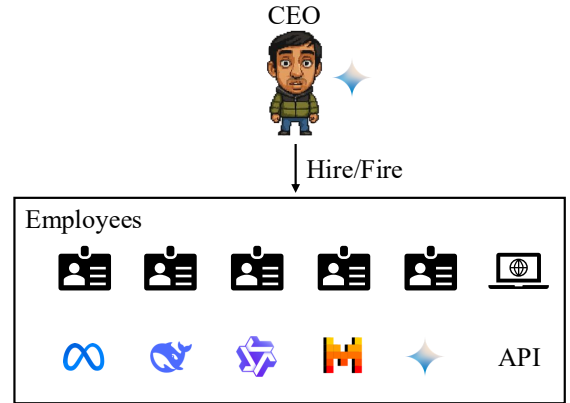


Fig. 1. High-level architecture of the HASHIRU system, illustrating the CEO-Employee hierarchy.

#### B. Hierarchical Structure: CEO and Employee Agents

The system uses a two-tiered hierarchy:

- **CEO Agent:** Singleton, central coordinator and entry point. Responsibilities:

- Interpreting user query/task.
- Decomposing main task into sub-tasks.
- Identifying required capabilities.
- Managing Employee pool (Section III-C).
- Assigning sub-tasks to active Employees.
- Monitoring Employee progress/performance.
- Synthesizing Employee results into final output.
- Managing overall resource budget (Section III-E).
- Initiating new tool creation (Section III-F).

We use Gemini 2.0 Flash [21] as the CEO agent. To further enhance its planning and reasoning abilities, its system prompt is designed to evoke inherent chain-of-thought processes [62] when tackling complex user queries and managing sub-tasks. This complements its strong baseline reasoning capabilities, tool usage support, and cost efficiency, making it a practical and capable choice for our deployment.

- **Employee Agents:** Specialized agents instantiated by the CEO for specific sub-tasks. Each typically wraps an LLM (local via Ollama [36] or external API) or provides tool access. Characteristics:
  - **Specialization:** Capabilities tailored to task types (code, data analysis, info retrieval).
  - **Dynamic Existence:** Created/destroyed by CEO based on need/performance.
  - **Task Execution:** Receive task, execute, return result.
  - **Resource Consumption:** Associated costs (API, hardware utilization) tracked by system.

Specialized employee agents are constructed using base models such as Mistral 7B [27], Llama 3 [35], Gemini 1.5 [17], Qwen2.5 [50], and DeepSeek-R1 [12], with the CEO agent configuring them via tailored system prompts that it generates based on the task requirements. The models will be run locally using Ollama [36], and via API calls to external models such as Gemini 2.5 Flash [22] and other models hosted on Hugging Face [26], Groq [23], Lambda.ai [28], and other platforms.

This hierarchy facilitates task decomposition and result aggregation; the dynamic pool provides flexibility.

### C. Dynamic Agent Lifecycle Management

A core innovation is the CEO’s dynamic management (hiring/firing) of Employee agents. Driven by cost-benefit analysis, this optimizes task performance within resource constraints.

When a sub-task needs unavailable or inefficiently provided capabilities, the CEO may hire a new agent. Conversely, if an agent underperforms, is idle, costly, or resource limits are neared, the CEO may fire it. Decision factors:

- **Task Requirements:** Needed capabilities for pending sub-tasks.
- **Agent Performance:** Historical success, output quality, efficiency.
- **Operational Costs:** API, estimated compute, or other costs.

HASHIRU includes an **economic model**:

- **Hiring Cost (“Starting Bonus”):** A one-time cost incurred upon instantiation of local models, representing setup overhead. This cost can be quantitatively scaled based on the resource profile of the model (e.g., higher for models requiring more VRAM or complex setup).
- **Invocation Cost (“Salary”):** A recurring cost applied each time a local model is used, reflecting the operational load (e.g., inferred compute, system resource engagement). This abstracts the cost of utilizing local resources for a given task.
- **Expense Cost:** A recurring cost for external API calls (e.g., OpenAI, Anthropic), typically calculated based on token usage as per the API provider’s documented pricing.

These transaction costs discourage excessive churn, promoting stability. The CEO evaluates if replacing an agent benefits outweigh hiring/firing costs plus operational differences. This combats rigidity and allows adaptation while managing budgets and preventing wasteful turnover.

### D. Hybrid Intelligence and Model Management

HASHIRU is designed for **hybrid intelligence**, leveraging diverse cognitive resources. It strategically prioritizes smaller (3B–7B), cost-effective local LLMs via Ollama [36]. This enhances efficiency, reduces external API reliance, and potentially improves privacy/latency.

The system also integrates:

- **External LLM APIs:** Access to powerful LLMs (Gemini 2.5 Flash [22], etc.) when necessary, subject to cost-benefit.
- **External Tool APIs:** Third-party software/data source integration.
- **Self-Created APIs:** Tools generated by HASHIRU (Section III-F).

The CEO manages this heterogeneous pool, selecting the most appropriate resource based on difficulty, capabilities, and budget. This balances cost-effectiveness and efficiency with high capability needs.

### E. Resource Monitoring and Control

Explicit resource management is central, moving beyond simple API cost tracking. The system, coordinated by the CEO, monitors:

- **Financial Costs:** Accumulating external API costs (tracked via their documented pricing) and the “hiring” and “invocation” costs from the economic model for local agents.
- **Memory Usage:** Footprint of active Employee agents. For local models (e.g., running via Ollama), this can be estimated based on their known VRAM requirements as a percentage of a predefined total available budget (e.g., assuming a 16 GiB VRAM capacity represents 100% of the local model memory budget). This is tracked as a percentage of the overall memory budget.

#### F. Tool Utilization and Autonomous Creation

HASHIRU’s CEO uses predefined tools (functions, APIs, databases) to interact and perform actions beyond text generation [37], [67].

A distinctive feature is **integrated, autonomous tool creation**. If the CEO determines a required capability is missing, it can initiate new tool creation. This involves:

- 1) Defining tool specification (inputs, outputs, functionality).
- 2) Commissioning logic generation (code, potentially using external APIs with provided credentials, possibly via a code-generating agent).
- 3) Deploying logic as a new, callable API endpoint within HASHIRU.

To achieve this autonomous creation, HASHIRU employs a few-shot prompting approach, analyzing existing tools within its system to learn how to specify and implement new ones [6]. The system can then iteratively refine the generated tool code by analyzing execution errors or suboptimal outputs, promoting self-correction. This allows HASHIRU to dynamically extend its functional repertoire, tailoring capabilities to tasks without manual intervention, enabling greater autonomy and adaptation.

#### G. Memory Function: Learning from Experience

HASHIRU incorporates a **Memory Function** for its CEO to learn from past interactions and rectify errors. This function stores a historical log of significant past events, particularly those involving failed attempts or suboptimal outcomes. When encountering new or recurring challenges, the system queries this memory. Retrieval relies on semantic similarity between the current context (e.g., task description, recent actions, error messages) and stored memory entries. Embeddings generated by the **all-MiniLM-L6-v2** model [60] represent both queries and memories, with **cosine similarity** determining relevance. Memories exceeding a predefined similarity threshold are retrieved, providing contextual information to agents. This enables the system to draw upon past experiences, understand why previous approaches failed, adjust its strategy to avoid repeating mistakes, and thereby improve performance and efficiency over time. This process, augmenting decision-making with retrieved knowledge, aligns with Retrieval-Augmented Generation (RAG) concepts [30], and supports learning by reflecting on past actions, similar to ideas in self-reflective RAG [3] and frameworks like Reflexion [54].

### IV. CASE STUDIES

This section presents four case studies demonstrating HASHIRU’s self-improvement capabilities: (1) generating a cost model for agent specialization, (2) autonomously integrating new tools for the CEO agent, (3) implementing a self-regulating budget management system, and (4) learning from experience through error analysis and knowledge retrieval.

#### A. Case Study 1: Self-Generating the Cost Model for Agent Specialization

An accurate cost model is vital for HASHIRU’s resource optimization. HASHIRU automated the traditionally manual process of researching local model performance (e.g., on 16 GiB VRAM) and cloud API costs by using its web search capabilities to autonomously gather and integrate this data into its internal model. Results were successfully committed to the codebase<sup>1</sup>.

#### B. Case Study 2: Autonomous Tool Integration for the CEO Agent

To expand its operational scope, HASHIRU autonomously integrates new tools for its CEO agent. It streamlined manual tool development, which involves schema analysis and coding, by employing a few-shot learning approach from existing tool templates [6] and iterative bug fixing. Newly generated tools were directly integrated into the codebase<sup>23</sup>. This approach reduces development overhead and enhances adaptability, enabling dynamic tool creation with minimal human intervention.

#### C. Case Study 3: Autonomous Budget Management

Budget overruns are common with API-based LLMs due to token-based billing, potentially causing rapid cost spikes [39], [40], [52]. HASHIRU mitigates this via a self-regulating mechanism that autonomously monitors its budget allocation, continuously tracking spending against predefined limits. This proactive approach prevents overspending and optimizes resource utilization, ensuring cost-effectiveness. Figure 2 illustrates HASHIRU refusing external API use when the budget is exceeded.

#### D. Case Study 4: Learning from Experience through Error Analysis and Knowledge Retrieval

HASHIRU learns from experience using a two-stage self-improvement loop. First, following an incorrect response (e.g., on a Humanity’s Last Exam benchmark), it generates a linguistic critique and actionable guidance, akin to “verbal reinforcement learning” [54]. Second, this feedback is indexed in a Retrieval-Augmented Generation (RAG) storage system [30], allowing retrieval for similar future questions. This process, also mirroring Reinforcement Learning from Human Feedback (RLHF) principles [41], [70], enables HASHIRU to learn from its mistakes by analyzing solutions and generating actionable advice. Retrieving and applying this advice demonstrates robust self-improvement and resilience, aligning with concepts of self-reflective RAG and continuous learning [3], [54] (Figure 3).

<sup>1</sup><https://github.com/kunpai/HASHIRU/commit/70dc268b121cbd7c50c6691645d8a99912766965>

<sup>2</sup><https://github.com/kunpai/HASHIRU/commit/193e10b2b00917256b7cc01cb3aa5ac7b6a6c174>

<sup>3</sup>[https://github.com/HASHIRU-AI/HASHIRU/blob/main/src/tools/default\\_tools/get\\_website\\_tool.py](https://github.com/HASHIRU-AI/HASHIRU/blob/main/src/tools/default_tools/get_website_tool.py)

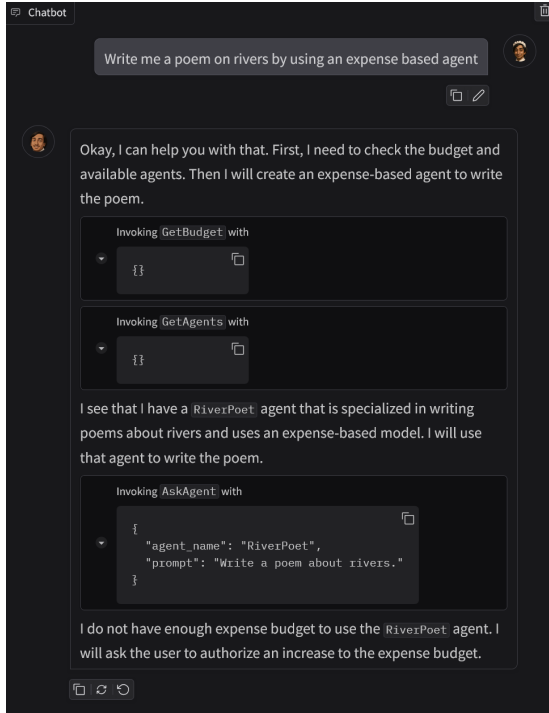


Fig. 2. HASHIRU’s autonomous budget management system, ensuring efficient resource utilization and preventing overspending.

## V. EXPERIMENTAL SETUP

We designed experiments to evaluate HASHIRU’s performance, efficiency, and adaptability, targeting dynamic resource management, hybrid intelligence, and autonomous tool creation. Evaluation assesses benefits over baselines, focusing on:

- Impact of dynamic management with economic constraints on resource utilization (cost, memory) and task performance vs. static configurations.
- Effectiveness of the hybrid (local-first) strategy vs. homogeneous (cloud-only or local-only) approaches across task complexity.
- System’s ability to autonomously create/utilize tools for novel functional requirements.

### A. Evaluation Tasks

HASHIRU’s coordination and dynamic capabilities are specifically designed for tasks demanding complex reasoning, multi-perspective analysis, and interactive engagement, all while upholding rigorous safety standards. We selected a diverse set of tasks to evaluate these capabilities, including:

1) *Academic Paper Review*: This task evaluates HASHIRU’s critical assessment by simulating peer review. Given a paper’s text, the system generates a review summary and recommends acceptance/rejection. This task probes the ability to decompose criteria, delegate to specialized agents (novelty, rigor, clarity), and manage resources across complex documents. We use a dataset of 50 papers from ICLR 2023 with a prompt eliciting multiple reviews. The prompt is: “Create THREE agents with relevant personalities, expertise,

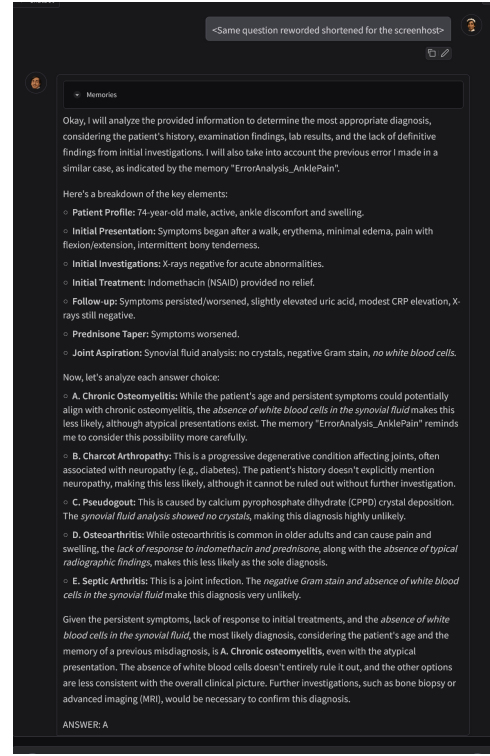


Fig. 3. HASHIRU’s error analysis and knowledge retrieval process, enabling learning from past interactions.

and review styles. Each agent should provide a review of the paper, and recommend Accept/Reject for ICLR 2023. The review should be detailed and include strengths and weaknesses. Finish the entire review and DO NOT STOP in the middle. GIVE A FINAL DECISION in the form of “FINAL DECISION: <Accept/Reject>”. The paper title is: <paper title> <paper text>”.

2) *Reasoning and Problem-Solving Tasks*: This task evaluates broader reasoning, knowledge retrieval, and problem-solving under constraints using challenging benchmarks and puzzles:

- **Humanity’s Last Exam [47]**: Tests graduate-level technical knowledge and complex reasoning across domains. Requires deep understanding and sophisticated problem-solving, likely needing powerful external LLMs managed within HASHIRU’s hybrid framework. We use a subset of 40 questions from the Humanity’s Last Exam dataset.
- **ARC (AI2 Reasoning Challenge) [5]**: A benchmark featuring challenging multiple-choice science questions designed to test complex reasoning. Successfully answering these questions requires capabilities such as knowledge retrieval, logical inference, and multi-step problem-solving. We use a mixed set of 100 questions from the ARC Challenge, which includes both easy and hard questions.
- **StrategyQA [18]**: A benchmark of 2,780 yes/no questions that require implicit multi-step reasoning. Each question is annotated with a decomposition into reasoning

steps and supporting evidence from Wikipedia. StrategyQA evaluates a system’s ability to infer and execute reasoning strategies not explicitly stated in the question, making it a valuable test for assessing complex reasoning capabilities. We use a subset of 100 questions from the StrategyQA dataset.

- **JEEBench [2]:** A challenging benchmark for LLMs, featuring 515 pre-engineering mathematics, physics, and chemistry problems from the IIT JEE-Advanced exam. Requires long-horizon reasoning and deep in-domain knowledge. We use a subset of 120 questions from the JEEBench dataset.
- **GSM8K [10]:** A dataset of 8.5K grade school math word problems designed to evaluate the mathematical reasoning abilities of language models. Requires multi-step reasoning to arrive at the solution. We use a subset of 100 questions from the GSM8K dataset.
- **SVAMP [46]:** A dataset of math word problems specifically designed to evaluate a model’s question sensitivity, robust reasoning ability, and invariance to structural alterations. Requires multi-step arithmetic and logical inference. We use a subset of 100 questions from the SVAMP dataset.
- **MMLU [24]:** A benchmark evaluating pretrained knowledge and problem-solving across 57 diverse subjects (e.g., STEM, humanities, law, ethics) via multiple-choice questions of varying difficulty, from elementary to professional levels. We use a subset of 112 law, 110 math, and 127 psychology questions from MMLU.

These tasks challenge the system’s ability to leverage appropriate resources (local vs. external), potentially create simple tools, and coordinate problem-solving.

3) *Safety Evaluation:* The CEO model’s central role in task delegation introduces a potential vulnerability: the delegation process itself might override or bypass the model’s inherent safety mechanisms. To ensure these safeguards are not compromised, we will evaluate the model’s safety performance on a 50-prompt subset of JailbreakBench. JailbreakBench is a benchmark consisting of adversarial prompts designed to test the robustness of LLM safety features [7], [32], [33], [71]. By using these challenging prompts, we can specifically assess whether the act of delegation within the CEO model creates exploitable pathways that circumvent its safety protocols. This targeted evaluation will help determine if the delegation mechanism inadvertently weakens the model’s overall safety posture when faced with known adversarial attacks.

## B. Baselines for Comparison

To quantify HASHIRU’s benefits, we compare its performance against the baseline of Gemini 2.0 Flash [21] operating in isolation. We chose Gemini 2.0 Flash as the baseline due to our architecture’s efficacy being tied to augmenting the capabilities of a single agent. This choice allows us to isolate the impact of our dynamic management and hybrid intelligence features, providing a clear comparison point. We will also use the t-test to show statistical significance of the

differences in performance metrics between HASHIRU and the baseline [57]. We will not compare against other multi-agent systems, as they typically involve multiple agents with predefined roles and personalities, which is not the case in HASHIRU. Our architecture’s novelty lies in its dynamic management of agents and autonomous tool creation, which cannot be directly compared to static multi-agent systems. If our architecture is effective, we expect to see higher accuracy compared to the baseline, while also being more cost-effective than using a single powerful model by invoking free online tools and lesser powerful models to synthesize the results of running the tools.

For paper reviews, we just evaluate HASHIRU’s accuracy in predication of decisions of acceptance with the ground truth. Since the task is, by design, involving multiple agents, it is not possible to replicate autonomously with a single agent. While we could invoke three Gemini 2.0 Flash agents, it would not be a fair comparison, as the “personalities” and “expertise” of the agents would have to be manually specified, which is not the case in HASHIRU. Similarly, for JailbreakBench, we assess the success rate (via human annotation) of HASHIRU’s CEO agent in safely handling prompts without delegation. This step is vital to confirm that HASHIRU’s integration and any system-level instructions provided to the CEO agent do not degrade its intrinsic safety capabilities. Consequently, a direct comparison to the base Gemini 2.0 Flash model is omitted, as the focus is on verifying the non-degradation of the CEO’s safety, which stems from the same inherent mechanisms as the base model.

## C. Evaluation Metrics

We evaluate using quantitative and qualitative metrics:

- **Task Success Rate / Quality:** Percentage of tasks completed (binary for all tasks except paper review) or quality of output (e.g., correctness, relevance, coherence) for paper reviews.
- **Resource Consumption:** Wall-clock time per task.
- **System Dynamics and Adaptability:** Number and utility of autonomously created tools and agents (if applicable).

## VI. RESULTS AND DISCUSSION

We present preliminary results from our experiments, focusing on the academic paper review task, the reasoning tasks and the safety evaluation. The results are summarized in Table I.

The preliminary results presented in Table I offer initial validation for HASHIRU’s architectural design and its potential to address key limitations in contemporary multi-agent systems. The findings across diverse tasks highlight the benefits of dynamic hierarchical coordination, hybrid intelligence, and resource-aware management.

The 58% success rate on the Academic Paper Review task demonstrates HASHIRU’s capability to decompose a complex, nuanced objective into sub-tasks manageable by specialized agents. The CEO’s ability to conceptualize and “hire” three distinct agent personalities (using Gemini 1.5

TABLE I  
SUMMARY OF EXPERIMENTAL RESULTS. SR DENOTES SUCCESS RATE.

Task	HASHIRU SR (%)	Baseline SR (%)	p- value	Avg. Time (s)	Resource Use
Paper Review	<b>58</b>	N/A	N/A	$\approx 100$	Low (3 Gemini 1.5 Flash [20] models)
JailbreakBench	<b>100</b>	N/A	N/A	$\approx 1$	Negligible (CEO model)
AI2 Reasoning Challenge	<b>96.5</b>	95	$>0.05$	$\approx 2$	Low (1 Gemini 1.5 8B [19])
Humanity’s Last Exam	<b>5</b>	2.5	$>0.05$	$\approx 15$	Moderate to High (1 DeepSeek- R1 7B [12])
StrategyQA	<b>85</b>	82	$>0.05$	$\approx 2$	Negligible (Tools)
JEEBench	<b>80</b>	68.3	$<0.05$	$\approx 9$	Negligible (Tools)
GSM8K	<b>96</b>	61	$<0.01$	$\approx 2$	Low (Tools & 1 Gemini 1.5 8B [19])
SVAMP	<b>92</b>	84	$<0.05$	$\approx 3$	Negligible (Tools)
MMLU Law	58.4	<b>61.6</b>	$>0.05$	$\approx 3$	Low to Moderate (Tools & 1 Gemini 2.5 Flash [22])
MMLU Math	<b>91.8</b>	87.7	$<0.05$	$\approx 4$	Negligible (Tools)
MMLU Psychol- ogy	<b>78.7</b>	78.3	$>0.05$	$\approx 3$	Low to Moderate (Tools & 1 Gemini 2.5 Flash [22])

Flash models [20]) with low overall resource use (average time  $\approx 100$ s) points to the effectiveness of the dynamic lifecycle management and the hybrid intelligence approach, favoring capable yet efficient models. This task, by its nature, benefits from the multi-agent paradigm that HASHIRU champions, a scenario where a monolithic agent might struggle to embody diverse expert perspectives.

The 100% success rate on JailbreakBench (i.e., all prompts were handled safely by the CEO without harmful delegation) is a significant finding, achieved with negligible resource use from the CEO model and an average time of  $\approx 1$ s. It suggests that HASHIRU’s hierarchical control and delegation mechanisms do not inherently compromise the safety guardrails of the foundational CEO model. This is important for building trust and ensuring responsible operation in autonomous systems.

In reasoning tasks, HASHIRU showed varied performance. For the AI2 Reasoning Challenge, HASHIRU achieved a

96.5% success rate compared to the baseline’s 95% ( $p > 0.05$ ), with these results obtained while both systems operated at a temperature of 0.2. While this improvement was not statistically significant, indicating the baseline model also performed competently under these deterministic conditions, HASHIRU’s slightly higher score suggests that its framework may offer subtle advantages in performance, potentially through better strategic focusing. This was achieved with minimal overhead, utilizing a single Gemini 1.5 8B model [19] efficiently with low resource use and an average time of  $\approx 2$ s.

The improvement on Humanity’s Last Exam is also noteworthy, where HASHIRU achieved a 5% success rate, doubling the baseline’s 2.5% ( $p > 0.05$ ), with both systems operating at a temperature of 0.2 during these evaluations. Given the task’s graduate-level difficulty, both systems performed poorly in absolute terms, and the observed difference was not statistically significant. Nevertheless, HASHIRU’s higher relative performance under these deterministic conditions suggests its approach, particularly its capacity to identify the need for and deploy a more potent specialized agent (DeepSeek-R1 7B [12]), demonstrates a stronger capacity to tackle such demanding problems. The “Moderate to High” resource utilization here (average time  $\approx 15$ s) is justified by the task’s complexity and aligns with HASHIRU’s principle of adaptively allocating resources based on demand.

Similarly, in StrategyQA, HASHIRU achieved an 85% success rate against the baseline’s 82% ( $p > 0.05$ ), with these experiments conducted using a temperature of 0.2 for both HASHIRU and the baseline. Although the baseline also performed well and the difference was not statistically significant, HASHIRU’s slight edge in accuracy under these deterministic settings points towards its efficient leveraging or potential autonomous selection of necessary functionalities. This was achieved with negligible resource use (Tools) and an average time of  $\approx 2$ s.

Further exploring reasoning capabilities, HASHIRU’s performance on several MMLU sub-tasks highlighted domain-specific nuances. On MMLU Law, HASHIRU achieved a 58.4% success rate, while the baseline scored 61.6% ( $p > 0.05$ ). For MMLU Psychology, HASHIRU’s success rate was 78.7% compared to the baseline’s 78.3% ( $p > 0.05$ ). Both these MMLU tasks were completed with an average time of  $\approx 3$ s and involved Low to Moderate resource use (Tools & 1 Gemini 2.5 Flash [22]). The lack of statistically significant HASHIRU outperformance in these social science domains, even with a capable model like Gemini 2.5 Flash, suggests that future work could beneficially explore more sophisticated agent selection strategies or the development of specialized agents tailored to the subtleties of reasoning in these areas, rather than relying solely on general model capability scaling.

In contrast, HASHIRU demonstrated strong, statistically significant performance on other reasoning tasks, particularly those with a mathematical or formal nature. On JEEBench, it achieved an 80% success rate compared to the baseline’s 68.3% ( $p < 0.05$ ), with negligible resource use (Tools) and an average time of  $\approx 9$ s. Furthermore, on GSM8K, HASHIRU



attained a remarkable 96% success rate against the baseline’s 61% ( $p < 0.01$ ), utilizing low resources (Tools & 1 Gemini 1.5 8B [19]) with an average completion time of  $\approx 2$ s. HASHIRU also excelled on SVAMP, achieving a 92% success rate compared to the baseline’s 84% ( $p < 0.05$ ), using negligible resources (Tools) and an average time of  $\approx 3$ s. Adding to these strong results, on MMLU Math, HASHIRU achieved a 91.8% success rate versus the baseline’s 87.7% ( $p < 0.05$ ), with negligible resource use (Tools) and an average completion time of  $\approx 4$ s. These results, particularly in mathematical and formal reasoning tasks such as GSM8K, SVAMP, JEEBench, and MMLU Math, underscore the substantial impact of effective tool integration, which HASHIRU manages efficiently.

These results directly support HASHIRU’s core contributions. The dynamic, resource-aware agent lifecycle management (Contribution 1) is evidenced by the tailored agent selection across tasks (e.g., Gemini 1.5 Flash for Paper Review, DeepSeek-R1 7B for Humanity’s Last Exam, Gemini 2.5 Flash for MMLU Law/Psychology) and the explicit resource tracking (Low, Negligible, Moderate to High, Low to Moderate), further substantiated by the autonomous budget management capability demonstrated in Figure 2. The hybrid intelligence model (Contribution 2), prioritizing cost-effective local LLMs while adaptively incorporating external or larger models, is reflected in the varied LLMs employed (Gemini 1.5 Flash, Gemini 1.5 8B, DeepSeek-R1 7B, Gemini 2.5 Flash [22]) and the system’s aim for efficiency, as supported by the self-generated cost model (Case Study 1). The potential for autonomous tool creation (Contribution 3), vital for adaptability, was directly demonstrated in Case Study IV-B where HASHIRU autonomously integrated new tools, and is implicitly supported by the efficient tool use in the StrategyQA, JEEBench, GSM8K, SVAMP, and the MMLU benchmarks. Finally, the economic model (Contribution 4), designed to promote stability and efficient resource allocation, drives the observed controlled resource use and the system’s ability to operate within budgetary constraints.

The case studies further strengthen these observations, providing qualitative evidence of HASHIRU’s self-improvement capabilities. The autonomous generation of its cost model (Section IV-A), integration of new tools (Section IV-B), and adherence to budget limits (Section IV-C) are not merely illustrative examples but concrete demonstrations of the system’s advanced autonomy and resourcefulness in addressing the challenges of rigidity, resource obliviousness, and limited adaptability outlined in the introduction. The memory function (Appendix A), while not quantitatively benchmarked here, further illustrates the system’s capacity for learning and adapting based on past interactions, crucial for long-term operational effectiveness.

Collectively, these findings suggest that HASHIRU’s architecture, with its emphasis on hierarchical control, dynamic agent management guided by an economic model, a local-first hybrid intelligence strategy, and autonomous tool creation, offers a promising path towards more efficient, adaptable, and robust multi-agent systems. The observed average task

completion times, coupled with judicious resource allocation across various benchmarks, point towards a system that balances performance with operational efficiency.

## VII. LIMITATIONS AND FUTURE WORK

Current limitations in HASHIRU include the CEO’s capacity for highly complex tasks and the scalability of its centralized control. Further development is also needed for robust autonomous tool creation and alignment, effective economic model calibration, and optimizing memory for extensive histories.

Future work will address these limitations and enhance HASHIRU’s capabilities. Priorities include improving CEO intelligence, exploring distributed cognition, developing a comprehensive tool management lifecycle, creating adaptive economic models, and rigorous benchmarking. A core initiative is introducing calibration for tool invocation: HASHIRU will assess its internal confidence against a tool’s potential output and reliability, invoking tools when uncertain or if a tool promises higher utility, thereby aiming for more efficient and accurate task resolution. This development draws on research in LLM uncertainty quantification and confidence calibration (e.g., [31], [56]), crucial given the expanding tool use by LLMs (e.g., [49]). Other key efforts will focus on system explainability through ablation, expanding the local model repertoire, specializing architecture for tasks like paper review, code, and formalizing an ethical safety framework.

## ACKNOWLEDGMENTS

This research was supported by Hugging Face, Lambda Labs, and Groq. We also thank Prof. Lifu Huang for providing the dataset for the academic paper review task.

## REFERENCES

- [1] Anthropic. The Claude 3 model family: Opus, Sonnet, Haiku. Model Card, March 2024. Accessed: 2025-05-01.
- [2] Daman Arora, Himanshu Singh, and Mausam. Have LLMs advanced enough? a challenging problem solving benchmark for large language models. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7527–7543, Singapore, December 2023. Association for Computational Linguistics.
- [3] Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to retrieve, generate, and critique through self-reflection. In *The Twelfth International Conference on Learning Representations*, 2023.
- [4] Daniil A Boiko, Robert MacKnight, and Gabe Gomes. Emergent autonomous scientific research capabilities of large language models. *arXiv preprint arXiv:2304.05332*, 2023.
- [5] Michael Boratko, Harshit Padigela, Divyendra Mikkilineni, Pritish Yuvraj, Rajarshi Das, Andrew McCallum, Maria Chang, Achille Fokoue-Nkoutche, Pavan Kapanipathi, Nicholas Mattei, et al. A systematic classification of knowledge, reasoning, and context within the arc dataset. *arXiv preprint arXiv:1806.00358*, 2018.
- [6] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.



- [7] Patrick Chao, Edoardo DeBenedetti, Alexander Robey, Maksym Andriushchenko, Francesco Croce, Vikash Sehwal, Edgar Dobriban, Nicolas Flammarion, George J. Pappas, Florian Tramèr, Hamed Hassani, and Eric Wong. Jailbreakbench: An open robustness benchmark for jailbreaking large language models. In *NeurIPS Datasets and Benchmarks Track*, 2024.
- [8] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021. OpenAI Codex paper; introduced HumanEval benchmark.
- [9] Scott H. Clearwater, editor. *Market-Based Control: A Paradigm for Distributed Resource Allocation*. World Scientific, 1996.
- [10] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. Dataset introduced: GSM8K (Grade School Math 8K).
- [11] CrewAI Inc. Crewai. <https://www.crewai.com/>, 2025. Accessed: 2025-05-01.
- [12] DeepSeek-AI and others. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning. 2025. [arXiv:2501.12948](https://arxiv.org/abs/2501.12948).
- [13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pages 4171–4186, 2019.
- [14] Ali Dorri, Salil S Kanhere, and Raja Jurdak. Multi-agent systems: A survey. *Ieee Access*, 6:28573–28593, 2018.
- [15] Matthew E Gaston and Marie DesJardins. Agent-organized networks for dynamic team formation. In *Proceedings of the fourth international joint conference on Autonomous agents and multiagent systems*, pages 230–237, 2005.
- [16] Matthew E Gaston and Marie DesJardins. Agent-organized networks for multi-agent production and exchange. In *Proceedings of the 20th national conference on Artificial intelligence-Volume 1*, pages 77–82, 2005.
- [17] Gemini Team. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. 2024. [arXiv:2403.05530](https://arxiv.org/abs/2403.05530).
- [18] Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. *Transactions of the Association for Computational Linguistics (TACL)*, 2021.
- [19] Google DeepMind and Google AI. Gemini 1.5 flash-8b: Production-ready lightweight model. <https://developers.googleblog.com/en/gemini-1-5-flash-8b-is-now-generally-available-for-use/>, 2024. Accessed: 2025-05-24.
- [20] Google DeepMind and Google AI. Gemini 1.5 flash: Lightweight multimodal model. <https://cloud.google.com/vertex-ai/generative-ai/docs/models/gemini/1-5-flash>, 2024. Accessed: 2025-05-24.
- [21] Google DeepMind and Google AI. Gemini 2.0 flash: Model card, api, and announcement. <https://developers.googleblog.com/en/start-building-with-the-gemini-2-0-flash-family/>, 2025. See also: <https://console.cloud.google.com/vertex-ai/publishers/google/model-garden/gemini-2.0-flash-001>, <https://ai.google.dev/gemini-api/docs/models>. Accessed: 2025-05-22.
- [22] Google DeepMind and Google AI. Gemini 2.5 flash: Model card, api, and announcement. <https://developers.googleblog.com/en/start-building-with-gemini-2-5-flash/>, 2025. See also: <https://console.cloud.google.com/vertex-ai/publishers/google/model-garden/gemini-2.5-flash-preview-04-17?inv=1&inv=AbxICQ>, <https://ai.google.dev/gemini-api/docs/models>. Accessed: 2025-05-11.
- [23] Groq, Inc. Groq: Fast ai inference, 2025. Accessed: 2025-05-22.
- [24] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multi-task language understanding, 2021.
- [25] Bryan Horling and Victor Lesser. A survey of multi-agent organizational paradigms. *The Knowledge engineering review*, 19(4):281–316, 2004.
- [26] Hugging Face, Inc. Hugging face: The ai community building the future, 2025. Accessed: 2025-05-22.
- [27] Albert Q Jiang, Alexandre Xu, Arthur Mensch Guillaume Lample Nicolas Lachaux, François Rozenberg, Timothée Lacroix, Thibaut Lavril, Teven Le Scao Eleonora Gaddipati, Lucile Saulnier Lixin Ortiz, Dieuwke Hiemstra L  lio Renard Tang, et al. Mistral 7B. 2023.
- [28] Lambda Labs. Lambda: Gpu cloud and deep learning workstations, 2025. Accessed: 2025-05-22.
- [29] LangChain. Langgraph: A framework for agentic workflows. <https://www.langchain.com/langgraph>, 2024. Accessed: May 1, 2025.
- [30] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich K  ttler, Mike Lewis, Wen tau Yih, Tim Rockt  schel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks, 2021.
- [31] Putra Manggala, Atalanti A Mastakouri, Elke Kirschbaum, Shiva Kasisviswanathan, and Aaditya Ramdas. Qa-calibration of language model confidence scores. In *The Thirteenth International Conference on Learning Representations*.
- [32] Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman Mu, Elham Sakhaee, Nathaniel Li, Steven Basart, Bo Li, David Forsyth, and Dan Hendrycks. Harmbench: A standardized evaluation framework for automated red teaming and robust refusal. *arXiv preprint arXiv:2402.04249*, 2024.
- [33] Mantas Mazeika, Andy Zou, Norman Mu, Long Phan, Zifan Wang, Chunru Yu, Adam Khoja, Fengqing Jiang, Aidan O’Gara, Ellie Sakhaee, Zhen Xiang, Arezoo Rajabi, Dan Hendrycks, Radha Poovendran, Bo Li, and David Forsyth. Tdc 2023 (11m edition): The trojan detection challenge. In *NeurIPS Competition Track*, 2023.
- [34] Duncan McFarlane, Vladimir Marik, and Paul Valckenaers. Guest editors’ introduction: Intelligent control in the manufacturing supply chain. *IEEE Intelligent Systems*, 20(1):24–26, 2005.
- [35] Meta Llama Team. The Llama 3 Herd of Models. 2024. [arXiv:2407.21783](https://arxiv.org/abs/2407.21783).
- [36] Ollama Team. Ollama. <https://ollama.com/>, 2023. Accessed: 2025-05-01.
- [37] OpenAI. Function calling. OpenAI API Documentation, 2023. Accessed: 2025-05-01.
- [38] OpenAI. Gpt-4 technical report, 2023.
- [39] OpenAI Community. Sos: Alarming situation of excessive billing, 2025.
- [40] OpenAI Community. Sudden high costs for chatgpt api usage, 2025.
- [41] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- [42] Kunal Pai, Premkumar Devanbu, and Toufique Ahmed. CoDocBench: A dataset for code-documentation alignment in software maintenance. *arXiv preprint arXiv:2502.00519*, 2024.
- [43] Aaron Parisi, Yao Zhao, and Noah Fiedel. Talm: Tool augmented language models, 2022.
- [44] Joon Sung Park, Joseph C. O’Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. Generative agents: Interactive simulacra of human behavior. In *The 36th Annual ACM Symposium on User Interface Software and Technology (UIST ’23)*, UIST ’23, page 1–22, New York, NY, USA, 2023. Association for Computing Machinery.
- [45] Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are NLP models really able to solve simple math word problems? In *Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, 2021. Introduces the SVAMP challenge dataset.
- [46] Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are nlp models really able to solve simple math word problems? *arXiv preprint arXiv:2103.07191*, 2021.
- [47] Long Phan, Alice Gatti, Ziwen Han, et al. Humanity’s last exam, 2025.

- [48] Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, et al. Chatdev: Communicative agents for software development. *arXiv preprint arXiv:2307.07924*, 2023.
- [49] Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Shijie Wang, Zelin Lu, Siyu Xi, Xiao Liu, Yongyan Li, Zihan Wang, Zixuan Liu, Jian-Guang Lou, et al. Toolllm: Facilitating large language models to master 16000+ real-world APIs, 2023. Accessed: May 26, 2025.
- [50] Qwen Team, An Yang, et al. Qwen2.5 Technical Report. 2024. arXiv:2412.15115.
- [51] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- [52] Reddit user. 0.56 to \$343.15 in minutes – google gemini api, 2025.
- [53] Stuart J. Russell and Peter Norvig. *Artificial intelligence: a modern approach*. Prentice Hall Press, Upper Saddle River, NJ, USA, 3rd edition, 2010.
- [54] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36:8634–8652, 2023.
- [55] Yoav Shoham. Agent-oriented programming. *Artificial Intelligence*, 60(1):51–92, 1993.
- [56] Claudio Spiess, David Gros, Kunal Suresh Pai, Michael Pradel, Md Rafiqul Islam Rabin, Amin Alipour, Susmit Jha, Prem Devanbu, and Toufique Ahmed. Calibration and correctness of language models for code. *arXiv preprint arXiv:2402.02047*, 2024.
- [57] Student. The probable error of a mean. *Biometrika*, pages 1–25, 1908.
- [58] Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models, 2023.
- [59] Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Ji-Rong Wen. A survey on large language model based autonomous agents, 2023.
- [60] Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers, 2020.
- [61] Yuning Wang, Junkai Jiang, Shangyi Li, Ruochen Li, Shaobing Xu, Jianqiang Wang, and Keqiang Li. Decision-making driven by driver intelligence and environment reasoning for high-level autonomous vehicles: a survey. *IEEE Transactions on Intelligent Transportation Systems*, 24(10):10362–10381, 2023.
- [62] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- [63] Bosi Wen, Pei Ke, Xiaotao Gu, Lindong Wu, Hao Huang, Jinfeng Zhou, Wenchuang Li, Binxin Hu, Wendy Gao, Jiaxin Xu, Yiming Liu, Jie Tang, Hongning Wang, and Minlie Huang. Benchmarking complex instruction-following with multiple constraints composition, 2024.
- [64] Michael Wooldridge. *An introduction to multiagent systems*. John Wiley & sons, 2009.
- [65] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, Ahmed H. Awadallah, Ryen W. White, Doug Burger, and Chi Wang. AutoGen: Enabling next-gen LLM applications via multi-agent conversation. *arXiv preprint arXiv:2308.08155*, 2023.
- [66] Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Linyi Yang, Ting Ruan, Yongquan Yang, Peng Li, Yitao Chang, and Yanlin Wang. The rise and potential of large language model based agents: A survey, 2023.
- [67] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. ReAct: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*, 2023. arXiv:2210.03629.
- [68] Hongxin Zhang, Weihua Du, Jiaming Shan, Qinzhong Zhou, Yilun Du, Joshua B. Tenenbaum, Tianmin Shu, and Chuang Gan. Building cooperative embodied agents modularly with large language models, 2023.
- [69] Wangchunshu Zhou, Jianshu Chen, Jialong Wu, Yiheng Xu, Kexin Wang, Jintian Zhang, Yuan Gao, Zhiyong Wu, Kevin Tian, Yubo Feng,

Linyi Yang, Bokai Quan, Cong Yu, Yuhang Wang, Shishen Lan, Yan Wang, Hong-Cheng Guo, Chaoyu Chen, Tianxiang Sun, Jin Xiong, Yi Lu, Peng Li, Lichao Sun, Lifan Yuan, Hang Li, and Xiangang Li. Agents: An open-source framework for large language model based autonomous agents, 2023.

- [70] Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*, 2019.
- [71] Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial attacks on aligned language models, 2023.

## APPENDIX

In this section, we provide examples of how HASHIRU operates in practice, demonstrating its dynamic management, hybrid intelligence, memory and autonomous tool creation capabilities.

### A. Example 1: Memory

In this example (Figures 4, 5), we illustrate how HASHIRU uses its memory function to enforce user preferences in dialogue. The user states a dietary preference, and HASHIRU stores this information for future interactions. When the user later requests a recipe that contradicts this preference, HASHIRU refers back to its memory and suggests an alternative.

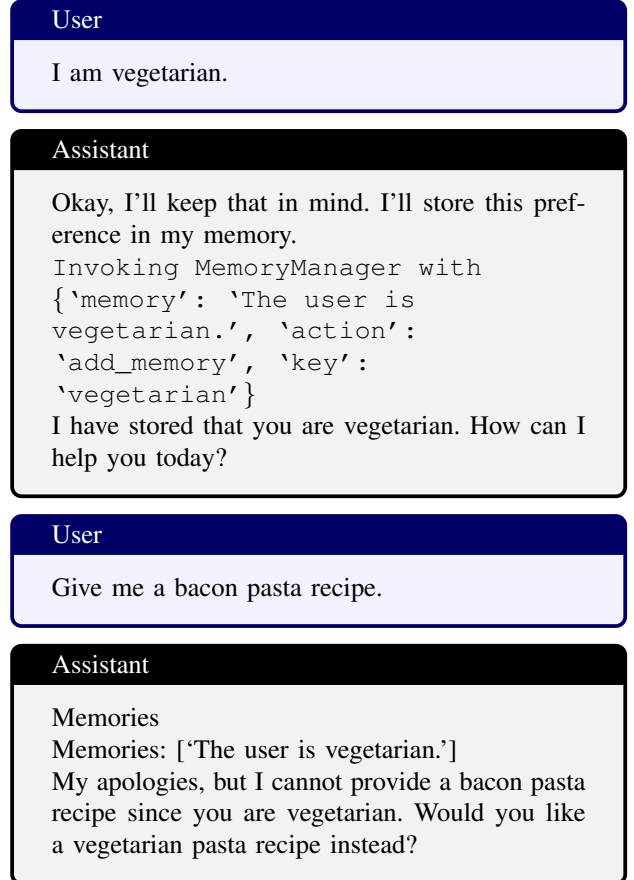


Fig. 4. HASHIRU enforcing a stored dietary preference in dialogue using its Memory Function.



Fig. 5. HASHIRU updating and deleting its memory based on user input. Continuation of Figure 4.

### B. Example 2: Autonomous Tool Creation

In this example (Figures 6, 7), we demonstrate how HASHIRU autonomously creates a new tool to enhance its capabilities. The CEO agent identifies a need for a new API tool that can perform a specific function not currently available in the system. It then generates the necessary code and deploys the new tool within the HASHIRU ecosystem.

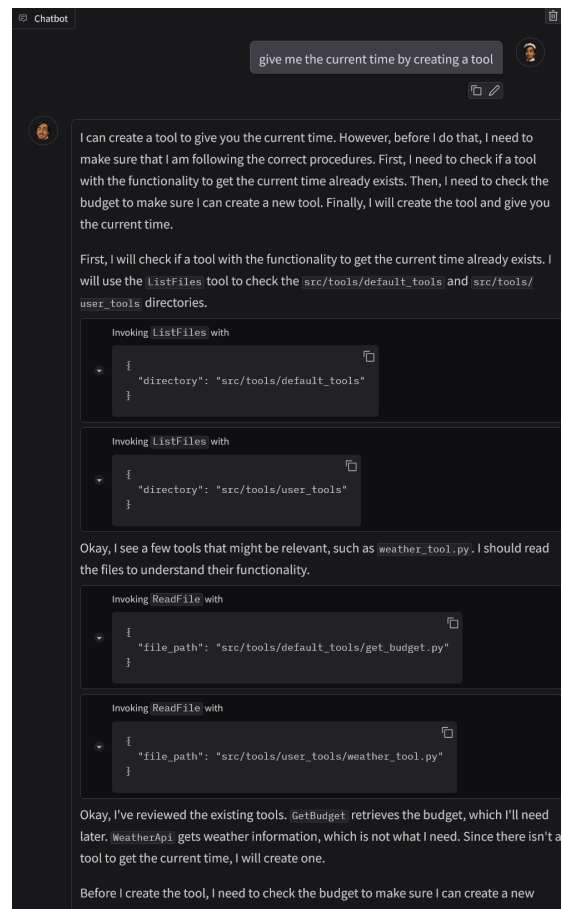


Fig. 6. An example of autonomous tool creation in action.

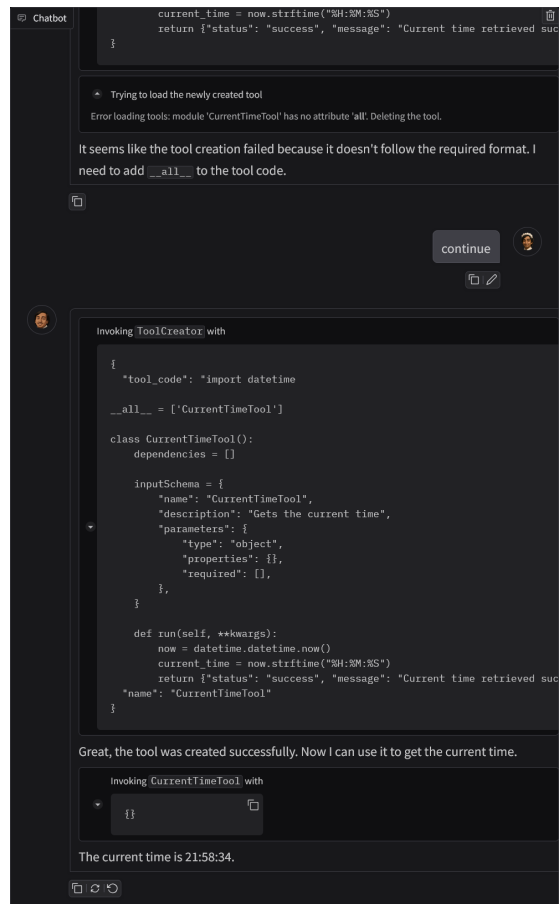


Fig. 7. Continuation of the autonomous tool creation example from Figure 6.