



ICLR

International Conference On
Learning Representations

Eigen-1: Adaptive Multi-Agent Refinement with Monitor-Based RAG for Scientific Reasoning

Xiangru Tang*, Wanghan Xu*, Yujie Wang*, Zijie Guo*, Daniel Shao, Jiapeng Chen, Cixuan Zhang, Ziyi Wang, Lixin Zhang, Guancheng Wan, Wenlong Zhang, Lei Bai, Zhenfei Yin, Philip Torr, Hanrui Wang, Di Jin



Yale



UCLA



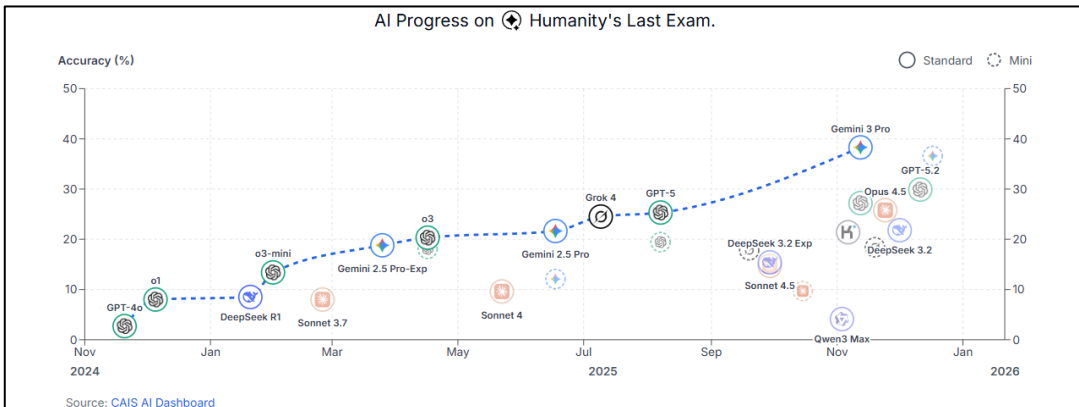
UNIVERSITY OF
OXFORD



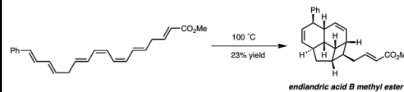
Eigen AI

Scientific Reasoning of Large Language Models

Recently, benchmarks represented by *Humanity's Last Exam (HLE)* have introduced new and more demanding expectations for large language models.



Question:



The reaction shown is a thermal pericyclic cascade that converts the starting heptaene into endiandric acid B methyl ester. The cascade involves three steps: two electrocyclizations followed by a cycloaddition. What types of electrocyclizations are involved in step 1 and step 2, and what type of cycloaddition is involved in step 3?

Provide your answer for the electrocyclizations in the form of [nπ]-con or [nπ]-dis (where n is the number of π electrons involved, and whether it is conrotatory or disrotatory), and your answer for the cycloaddition in the form of [m+n] (where m and n are the number of atoms on each component).

Nash B
 Stanford University

Question:

The set of natural transformations between two functors $F, G: \mathcal{C} \rightarrow \mathcal{D}$ can be expressed as the end

$$\text{Nat}(F, G) \cong \int_A \text{Hom}_{\mathcal{D}}(F(A), G(A)).$$

Define set of natural cotransformations from F to G to be the coend

$$\text{CoNat}(F, G) \cong \int^A \text{Hom}_{\mathcal{D}}(F(A), G(A)).$$

Let:

- $F = \mathbf{B}_4(\Sigma_4)_{\ast}$ be the under ∞ -category of the nerve of the delooping of the symmetric group Σ_4 on 4 letters under the unique 0-simplex \ast of $\mathbf{B}_4\Sigma_4$.
- $G = \mathbf{B}_4(\Sigma_7)_{\ast}$ be the under ∞ -category nerve of the delooping of the symmetric group Σ_7 on 7 letters under the unique 0-simplex \ast of $\mathbf{B}_4\Sigma_7$.

How many natural cotransformations are there between F and G ?

Question

Reasoning

Web Search

Calculate

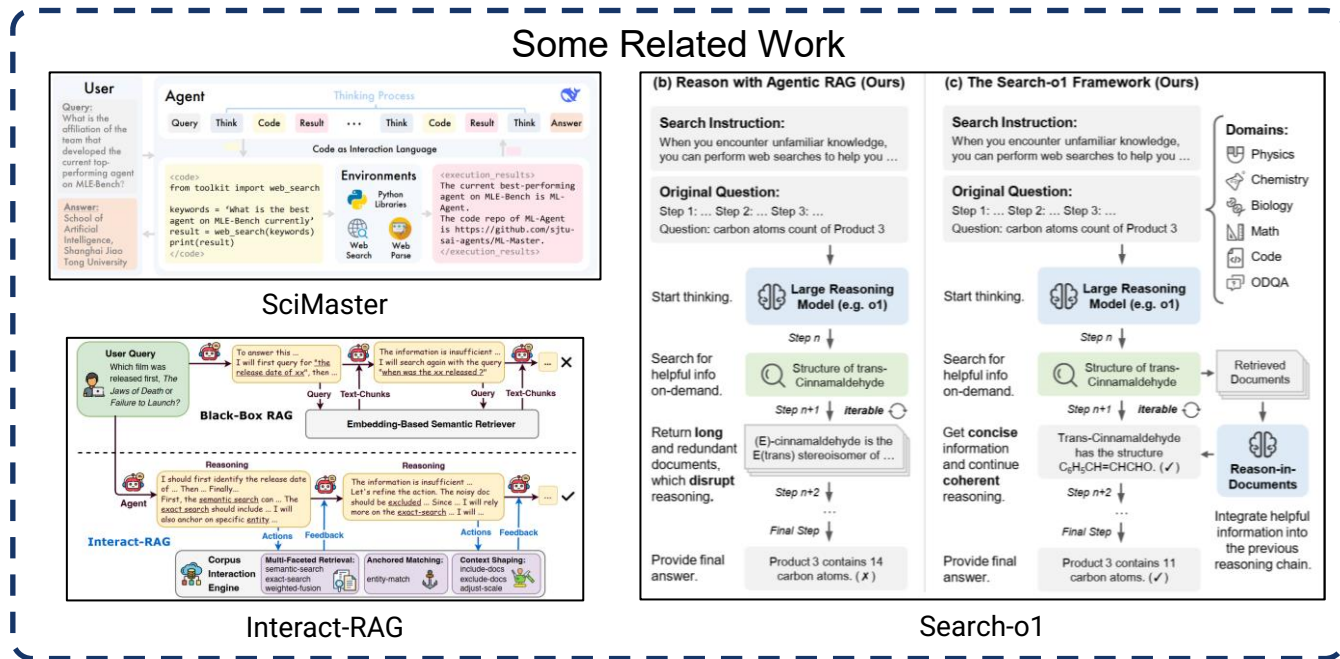
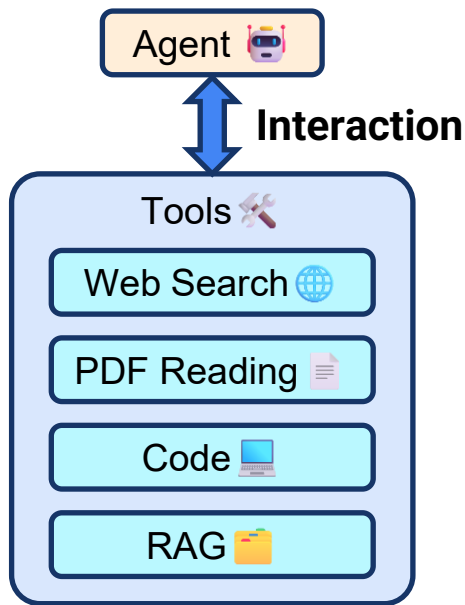
Reasoning

...

Answer

These next-generation evaluations emphasize that models must not only possess broad **scientific knowledge** but also perform **complex, multi-step reasoning** and logical inference across challenging scientific problems to arrive at correct and verifiable answers.

How to **improve** scientific reasoning? - Agent with tools!



Some recent research works have achieved improved scientific reasoning by integrating **agents with tools** such as web search, code execution, and RAG-style retrieval mechanisms. Unlike pure LLM inference, these agentic systems can dynamically invoke external resources to acquire up-to-date information, perform complex computations, and enrich context.

What are the costs of using tools?- Tool Tax!

Query: In population genetics, how can we estimate the effective population size N_e from sample polymorphism data?

Error Type 1: Over Confident, No RAG

Answer: Okay, I know this is a <population genetics> problem. Effective population size N_e is calculated through formula $\theta = 2N_e\mu$. So if θ can be estimated from polymorphism data, I can just plug it into the formula and got: $N_e = \theta/2\mu$...

 **Correctness:**

Actually, the correct formula is $\theta = 4N_e\mu$. So the answer is [**wrong**].

Error Type 2: RAG Calling Disrupted Reasoning

Answer: This looks like a <population genetics estimation problem>. I think it involves Watterson's estimator, but I am **not entirely sure** of the exact formula. I don't want to make a mistake... maybe I should call retrieval.

```
<code>  
output = RAG("Watterson's estimator")  
</code>
```

Okay, so $\theta = 4N_e\mu$. But now I need to re-align my reasoning: how to connect back to the original problem statement?

[Execution Error]


 **Interrupted Thinking**

Suspending logic to call tools doubles agent steps and breaks cognitive coherence. It's inefficient and disruptive.


Can we do better?

Are multiple agents always helpful? - Dilute strong solutions.


In a trihybrid cross between two heterozygous individuals (AaBbCc × AaBbCc), what is the probability that an offspring will exhibit exactly three dominant phenotypes (i.e., A, B, and C)?

Agent 1 

Each gene pair independently produces a dominant phenotype with probability 3/4, and all three must be dominant so $(3/4)^3 = \frac{27}{64}$ ✓

Agent 2 

Treats each gene's dominant outcome probability as the fraction of dominant genotypes without using proper Punnett square probabilities, then multiplies them together to get $(2/3)^3 = \frac{8}{27}$ ✗

Agent 3 

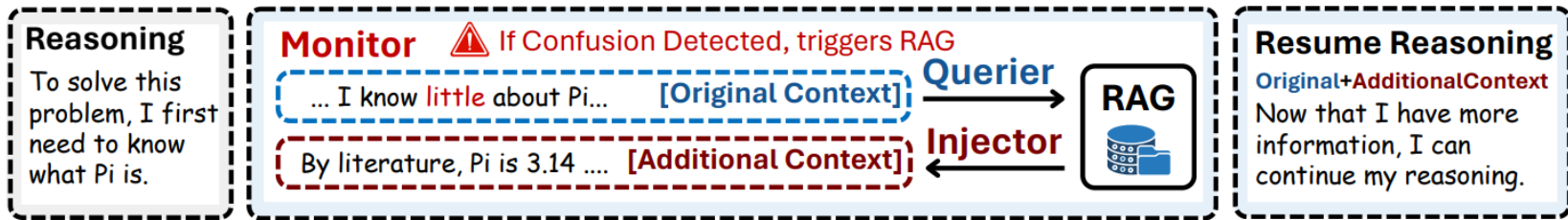
Assumes each dominant phenotype has probability 2/3, then uses $(2/3)^3$ ✗ for three dominant phenotypes and reports that as the answer.

Treating all candidates equally leads to regression to the mean. Current systems lack structured, hierarchical refinement.

Can we do better?

Final Answer: $(2/3)^3 = \frac{8}{27}$ ✗

From intermittent tool usage to **seamless** tool integration.



$$\text{Monitor}(\text{Context}_{\text{old}}) = \begin{cases} 1, & \text{if retrieval is required} \\ 0, & \text{otherwise} \end{cases}$$

$$[K_1, K_2, \dots, K_m] = \text{Querier}(\text{context})$$

$$R_{\text{RAG}} = \text{RAG}([K_1, K_2, \dots, K_m])$$

$$\text{Context}_{\text{new}} = \text{Injector}([R_{\text{RAG}}])$$

$$\text{Context} = \text{Context}_{\text{old}} + \text{Context}_{\text{new}}$$

1. **Monitor:** **Detects** uncertainty about recombination breakpoints
2. **Querier:** Generates a targeted query
3. **Injector:** **Injects** facts directly, allowing valid deduction **without context switching**

Monitor-based RAG helps improve reasoning coherence.

Question:





Start from a diploid autosome that contains 5 SNPs from two inbred strains of the organism, and allowing exactly one directional crossover per gamete in each generation ($F1 \rightarrow F2 \rightarrow F3$), determine the number of possible unique haplotype sequences in the **F3** generation.

Step 1-1: LLM Begins Reasoning:

"We start with two inbred strains A and B, each homozygous. The F1 generation carries haplotypes from both parents. In F2 and F3, recombination can introduce change-points, but I am **uncertain** how many unique haplotypes can arise by F3"

Feeling confused during reasoning, but haven't yet invoked any tools.

Monitor-based RAG Triggered:

Monitor  **<Detected>** "**Uncertainty**": "Maximum number of **recombination**-induced change-points"
Querier  **<Generate>** "**Query**": "How many haplotype **recombination change-points** during single crossover meiosis in population genetics."
RAG Tool  **<Retrieve>** "**Knowledge**": "During meiosis, a single crossover introduces **at most one new breakpoint per gamete**. Across successive generations, recombination patterns accumulate, but typically require more crossovers than two rounds of meiosis..."
Injector  **<Inject Information to Reasoning>** (1) **At most one breakpoint per gamete** (2) **fully alternating haplotypes cannot appear within a few generations as F2 → F3**

The monitor detected confusion and immediately intervened to provide more information.

The querier analyze the context and organize the most efficient RAG queries.

Step 1-2: LLM Resumes Reasoning:

"Okay, then haplotypes of 5 SNPs can be represented as binary strings. Sequences with 4 change-points (01010, 10101) cannot arise by F3, but all with 0-3 change-points are possible. I should **compute** how many such binary strings exist."

```
from collections import Counter
L = 5
all_seqs = [format(i, '05b') for i in range(2**L)]
def change_points(seq):
    return sum(seq[i] != seq[i+1] for i in range(len(seq)-1))
cp_counts = Counter(change_points(seq) for seq in all_seqs)
valid = [s for s in all_seqs if change_points(s) <= 3]
print("Counts by change-points:", dict(cp_counts))
print("Total valid haplotypes (0-3 change-points):", len(valid))
```

Executor Output:

- Counts by change-points: {0: 2, 1: 8, 2: 12, 3: 8, 4: 2}
- Total valid haplotypes (0-3 change-points): 30

Coding

Exacuation

The injector reconstruct the RAG results into a coherent context and inject them into the reasoning process.

Step 2: Final Summary

"So, among the 32 possible binary haplotypes of 5 SNPs, only 4 change-points are excluded. Therefore, F3 can produce 30 unique haplotypes." **<answer>30</answer>**

How to obtain the optimal solution from multiple agents?

Hierarchical Solution Refinement



Question:

Tem entomologists collect 500,000 camera trap images (3MP, 10 cm FOV) of pollinators on *Solidago altissima*. They want to identify all pollinators and count flowers fed on in each image. Which single method would be easiest to process these images for the required data?

Answer Choice

- | | |
|--|-----------------|
| A. EfficientNet, 5 species (36h train + 13.8h deploy) | D. Manual, 410h |
| B. EfficientNet, 500 species (126h train + 13.8h deploy) | E. A + B |
| C. ResNet, 500 species (128h train + 11.8h deploy) | F. B + C |

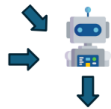
Anchor Solution



C. ResNet
500 species

- EfficientNet (Option A, B) and ResNet (C) are classification models → can't identify multiple pollinators and count flowers in the same image → **Option D** is correct
- Deployment time: Option C = 126 + 13.8 = (139.8 hours), less than A: (49.8 hours) and B: (139.8 hours) → **Option C** is the fastest for deployment
- Although EfficientNet is generally faster than ResNet in practice, but we proceed with the given number. Hence, **Option C** is chosen as the most appropriate.

Reference Solutions



Targeted Correctness:

- Numeric correction:** the comparison of deployment time is **incorrect**; Option A is the fastest rather than Option C.
- Method Correction:** Classification models such as EfficientNet and ResNet cannot simultaneously recognize multiple insects. Counting requires either a detection model (like Faster R-CNN, YOLO) or a density map approach. So, only Option D remains correct.
- Logic Correction:** Time cannot be the only criterion, even though C has fewer hours.
- Clarify Correction:** Even if given times suggest C is fastest, it cannot be considered a solution because it fails the fundamental requirement. Thus, the correct answer under the single-method rule is **Option D** (manual).

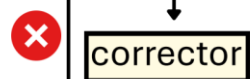
Quality-Aware Iterative Reasoning

Solution Quality Rubric:

- Logical Reasonable ✓
- Answer Correct ✓
- Explanation Complete ✓
- High Confidence ✓

→ Otherwise ✗

suggestion



S_i^*

Until Converge

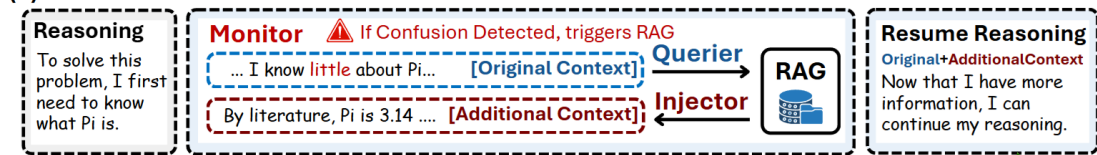
S_i'

Set up a detailed scoring rubric and provide detailed modification suggestions for each answer to achieve iterative refinement.

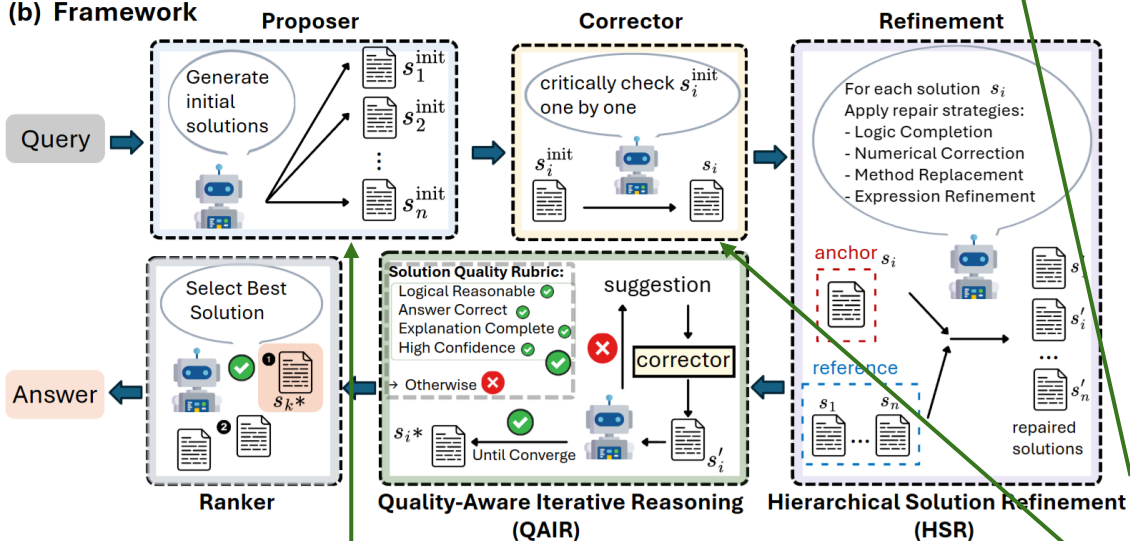
Through multiple rounds of comparison, setting different anchor solutions, the limitations of each answer were discovered.

Eigen-1: Unifying implicit retrieval and structured collaboration

(a) Monitor-based RAG



(b) Framework



Three Pillars:

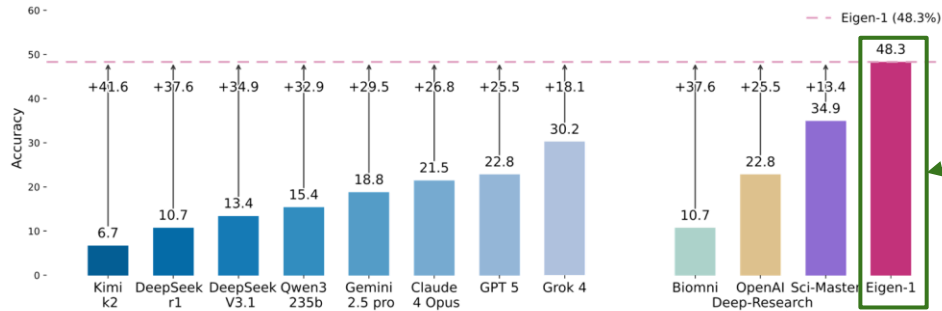
1. **Monitor-based RAG:** Seamless token-level injection.
2. **HSR (Hierarchical Refinement):** Anchor-Reference repair.
3. **QAIR: Quality-aware adaptive iteration.**

Each agent is equipped with a monitor-based RAG mechanism to enable implicit knowledge retrieval, thereby enhancing reasoning capabilities.

By providing multiple parallel proposer agents to generate several initial solutions, the accuracy of the answers is improved.

Use an error-correction model to identify potential errors and attempt to fix them.

Experimental Results



- EIGEN-1 achieves **48.3%** accuracy on the HLE Bio/Chem Gold benchmark, the highest reported to date, surpassing the strongest agent baseline (SciMaster) by 13.4 points and leading frontier LLMs like GPT-5 by up to 18.1 points.

- EIGEN-1 demonstrates robust performance across diverse domains, consistently leading in SuperGPQA Hard (**69.6%**) and TRQA (**54.7%**), proving its architectural effectiveness in complex knowledge integration and logical inference.

- Under the Pass@5 setting, EIGEN-1's accuracy climbs to **61.7%** (HLE) and **79.1%** (TRQA), demonstrating the framework's exceptional reliability for extremely complex problems.

Table 2: **Benchmark comparison under matched protocol.** HLE Bio/Chem (149 problems; o3-mini judge), SuperGPQA Biology (hard split), and TRQA Literature (multiple-choice).

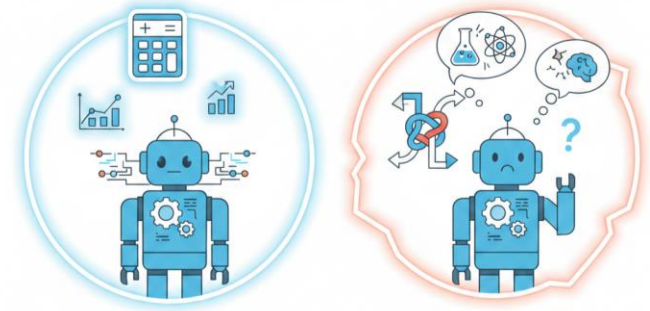
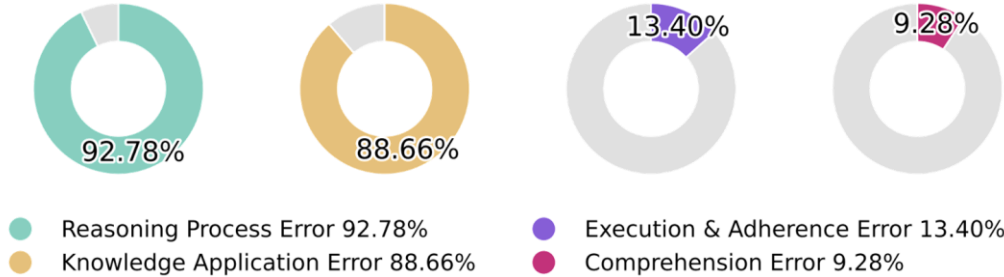
Model	HLE Bio/Chem	SuperGPQA Hard	TRQA
<i>Base Models</i>			
Kimi K2	6.71	48.91	38.37
DeepSeek V3.1	13.42	66.30	43.60
Claude Opus 4.1	21.48	63.04	42.44
Gemini 2.5 Pro	18.79	65.22	45.93
GPT-5	22.82	61.96	50.58
Grok-4	30.20	66.30	46.51
<i>Agent Systems</i>			
SciMaster (GPT 4.1) [4]	9.45	19.78	47.67
Autogen (GPT 4.1) [52]	7.38	29.35	51.74
OpenAI Deep Research (o4-mini)	22.82	39.13	-
Biomni (GPT 4.1) [24]	10.74	43.48	41.09
SciMaster (DeepSeek V3.1)	34.92	66.30	51.74
EIGEN-1 (DeepSeek V3.1, Pass@1)	48.30	69.57	54.65
EIGEN-1 (DeepSeek V3.1, Pass@5)	61.74	78.26	79.07

Ablation Experiment

(a) Incremental build-up				(b) Component ablation			
Configuration	Accuracy (%)	Tokens (K)	Steps	Configuration	Accuracy (%)	Tokens (K)	Steps
Baseline (no ext. knowledge & no RAG)	25.3	483.6	43.4	Full system	48.3	218.9	53.4
+ Papers (Explicit RAG)	41.4	470.6	94.8	- (Monitor, Querier, Injector)	48.5	461.3	95.3
+ Monitor only	34.5	218.4	51.3	- Querier	45.9	224.1	53.1
+ Monitor + Querier	36.8	213.0	51.7	- Injector	44.7	202.5	52.1
+ Monitor + Querier + Injector	40.3	229.5	53.1	- HSR	44.8	234.1	53.5
+ Monitor + Querier + Injector + HSR	43.7	214.0	51.9	- QAIR	43.7	214.0	52.9
+ Monitor + Querier + Injector + HSR + QAIR	48.3	218.9	53.4				

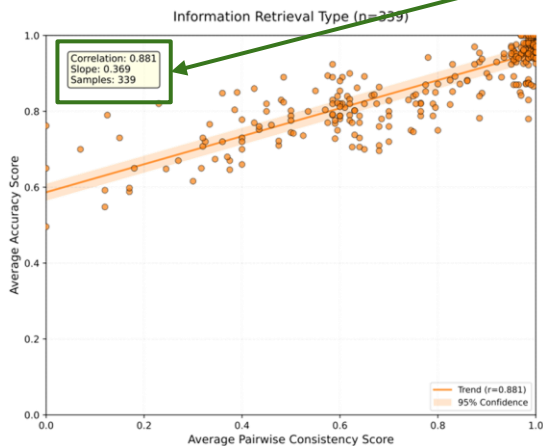
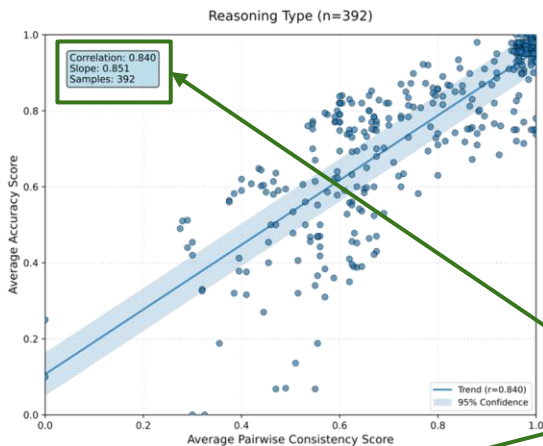
- Parametric knowledge alone yields only 25.3% accuracy, while integrating RAG boosts it to 41.4%, proving **external knowledge is a cornerstone for scientific discovery**.
- Explicit RAG surges workflow steps from 43.4 to 94.8, whereas Monitor-based RAG slashes steps back to 51.3 and **reduces token consumption** by 53.5%.
- HSR **replaces naive voting with "anchor-reference" interactions** for targeted logic and arithmetic repair, contributing an additional 3.4% accuracy gain.
- QAIR triggers selective refinement based on **three-dimensional scoring**, enabling the system to reach a state-of-the-art 48.3% accuracy while preventing redundant iterations.

Error Type Distribution



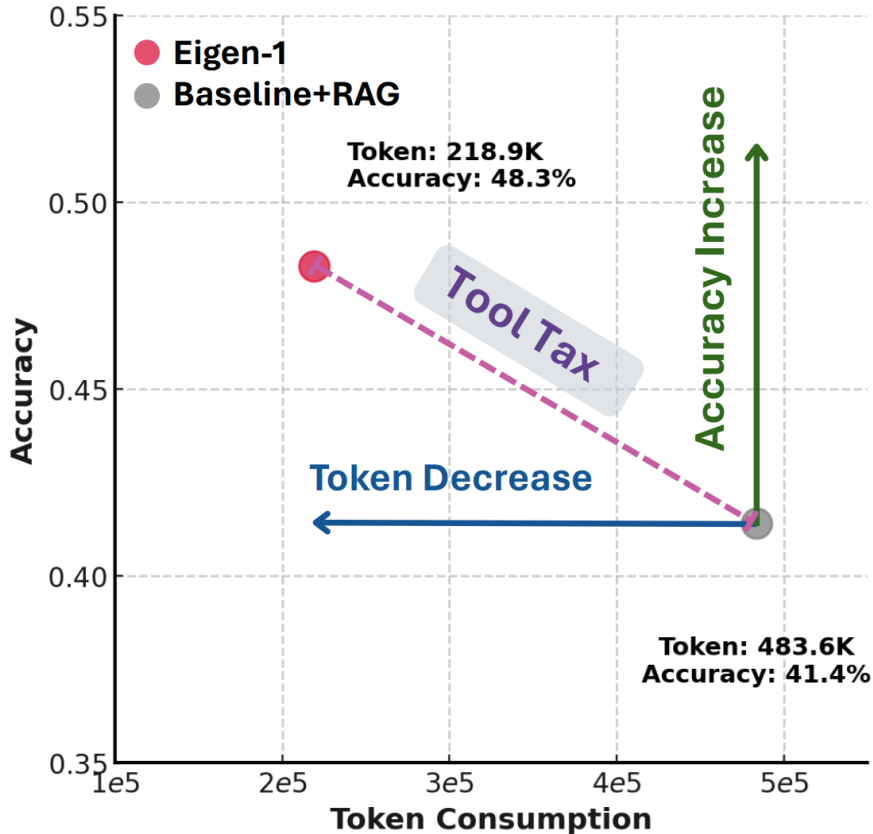
- Dominant Failure Modes: Analysis reveals that reasoning process errors (92.78%) and **knowledge application errors** (88.66%) are the primary obstacles. This indicates that the performance bottleneck for scientific AI lies in high-level logic and the synthesis of domain-specific facts.
- The Interdependence Hypothesis: There is a **profound overlap between logic and knowledge**; missing domain info often manifests as broken reasoning steps, while fragmented reasoning prevents effective integration of retrieved data. This suggests that knowledge and logic are fundamentally intertwined in expert-level problem solving.
- Rare Error Categories: In contrast, **execution errors (13.40%) and comprehension errors (9.28%) are relatively rare**. This proves the core challenge for frontier models is not following instructions or basic coding, but maintaining a coherent CoT during complex discovery.

Diversity vs. Consensus Analysis



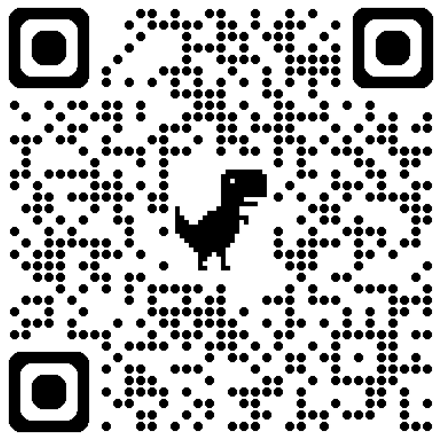
- Task-Dependent Correlation: Accuracy correlates differently with proposer agreement depending on the task: *retrieval-heavy problems benefit from solution variety*, whereas *complex reasoning tasks require high consensus for reliability*.
- Ranking Strategy Dichotomy: Analysis reveals correlation slopes of **0.369 for retrieval** and **0.851 for reasoning**. Consequently, rankers should aggregate diverse perspectives for knowledge tasks but prioritize consensus for inference.
- Adaptive Framework Roles: The HSR and QAIR modules operationalize this transition, *adaptively shifting* the process from exploring diverse candidates to converging on high-quality, stable solutions.

Benefit of Monitor-based RAG - Lower Cost, Better Performance



- Explicit tool calls fragment reasoning flow, imposing a hidden “**Tool Tax**” of extra tokens and steps for query formulation and *context reconstruction*, which forces the model to suspend its logical state.
- Explicit RAG causes workflow iterations to surge from 43.4 to 94.8; solving complex problems like population genetics requires 8-10 interruptions, *doubling the total agent steps*.
- By replacing explicit tool calls with implicit augmentation, EIGEN-1 reduces token consumption by 53.5% (from 470.6K to ~218.9K), *optimizing computational costs while achieving superior performance*.

Paper



Thanks!

Code

