

# Before you <think>, monitor:

Implementing Flavell's metacognitive framework in LLMs



# Monitor-Generate   # Generate-Verify   # Monitor-Generate-Verify   # Metacognition   #LLM-Modulo Framework

## TL;DR

**Nick Oh**  
socius labs  
London, UK  
nick.sh.oh@socius.org

**Fernand R. Gobet\***  
Centre for Philosophy of  
Natural and Social Science,  
London School of Economics  
London, UK  
f.gobet@lse.ac.uk

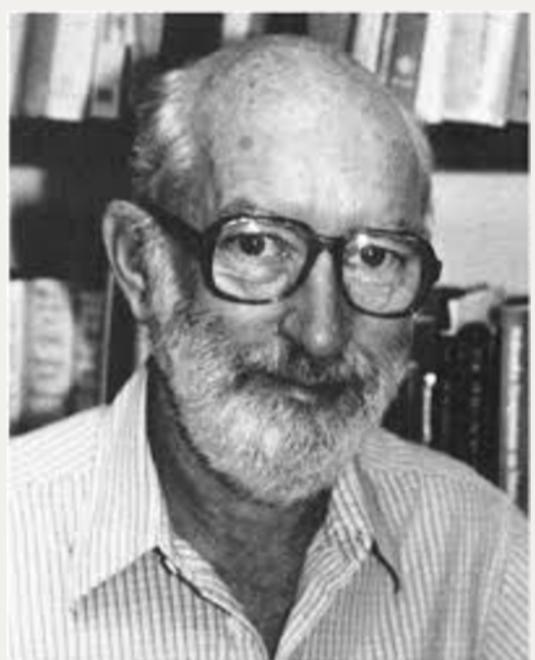
\*ACKNOWLEDGEMENT: This experimental implementation builds on the *Monitor-Generate-Verify* framework co-developed with Fernand Gobet.

**Monitor-Generate methods** (e.g., SELF-DISCOVER) **plan but can't verify**; **Generate-Verify methods** (e.g., SELF-REFINE) **refine but start blind**.

We chain them via Flavell's (1979) metacognitive theory.

**GSM8K: 75.4% vs 68.4% accuracy, 1.3 vs 2.0 attempts.**  
**Better initial solutions need less verification.**

## Background



Metacognition and Cognitive Monitoring  
A New Area of Cognitive-Developmental Inquiry

JOHN H. FLAVELL Stanford University

**Monitor-Generate-Verify:** Formalising Metacognitive Theory for Language Model Reasoning (Oh and Gobet, 2025)

**Algorithm 1** Flavell's Metacognitive Regulation

```

1: Initialise:  $S_0 = \text{ACTIVE}$ ;  $\tau = 0$ ;  $\mathcal{ME}_{\text{evaluative}}^{-1} = \emptyset$ 
2: while  $S_\tau = \text{ACTIVE}$  do
3:   // MONITOR: Retrieve knowledge & assess experience
4:    $\mathcal{MK}_\tau \leftarrow \begin{cases} \text{if } \tau = 0 \text{ then retrieve}(\mathcal{MK}, \mathcal{T}, \mathcal{G}) \\ \text{else } \mathcal{MK}_{\tau-1} \cup \text{retrieve}(\mathcal{MK}, \mathcal{ME}_{\tau-1}) \end{cases}$ 
5:    $\mathcal{ME}_{\text{difficulty}}^\tau \leftarrow \text{feel}(\mathcal{T}, \text{Outcomes}_{\tau-1}) \oplus \text{assess}(\mathcal{T}, \mathcal{MK}_\tau)$ 
6:   // GENERATE: Select & execute cognitive strategy
7:    $\mathcal{CS}_\tau \leftarrow \text{select}(s \in \mathcal{MK}_{\text{Strategy}} \mid \mathcal{ME}_{\text{difficulty}}^\tau, \mathcal{MK}_\tau, \mathcal{T}, \mathcal{G})$ 
8:    $\mathcal{CO}_\tau \leftarrow \text{execute}(\mathcal{CS}_\tau, \mathcal{T}, \mathcal{G})$ 
9:   // VERIFY: Evaluate progress & update knowledge
10:   $\mathcal{ME}_{\text{evaluative}}^\tau \leftarrow \text{assess}(\mathcal{CO}_\tau, \mathcal{MK}_\tau)$ 
11:   $\mathcal{MS}_\tau \leftarrow \text{select}(s \in \mathcal{MK}_{\text{meta}} \mid \mathcal{ME}_{\text{evaluative}}^\tau)$ 
12:   $\mathcal{MO}_\tau \leftarrow \text{execute}(\mathcal{MS}_\tau, \mathcal{CO}_\tau, \mathcal{MK}_\tau, \mathcal{G})$ 
13:   $\mathcal{MK} \leftarrow \text{update}(\mathcal{MK}, \Phi_\tau)$  where  $\Phi_\tau = (\mathcal{ME}_\tau^\tau, \text{Strategy}_\tau, \text{Outcome}_\tau)$ 
14:   $\mathcal{MK} \leftarrow \text{if goal\_achieved}(\mathcal{CO}_\tau, \mathcal{G}) \text{ then TERMINATE else ACTIVE}$ 
15:   $\tau \leftarrow \tau + 1$ 
16: end while
17: return  $y_{\arg \max_i \text{mean}(\mathcal{ME}_\tau^i)}$ 

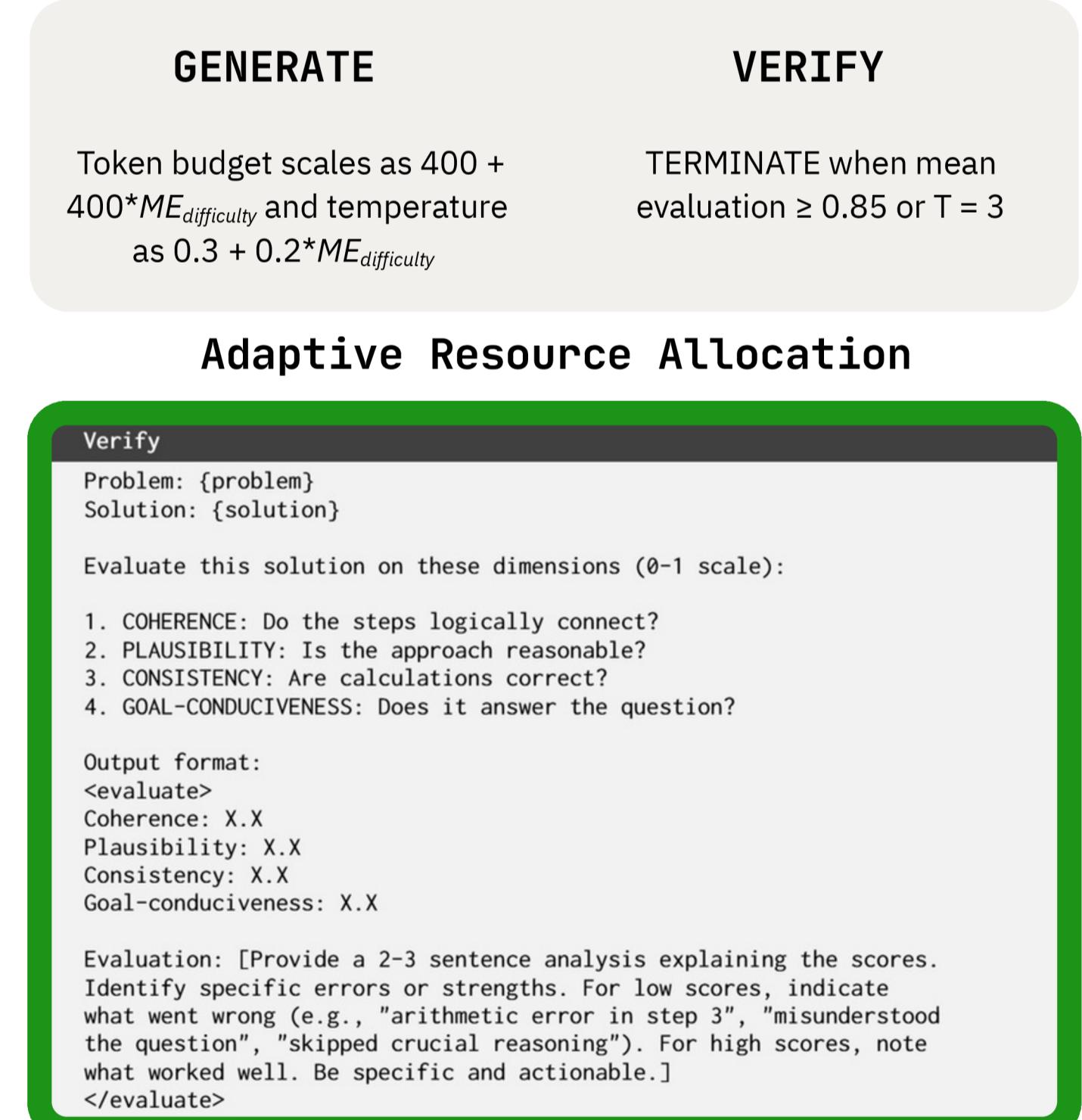
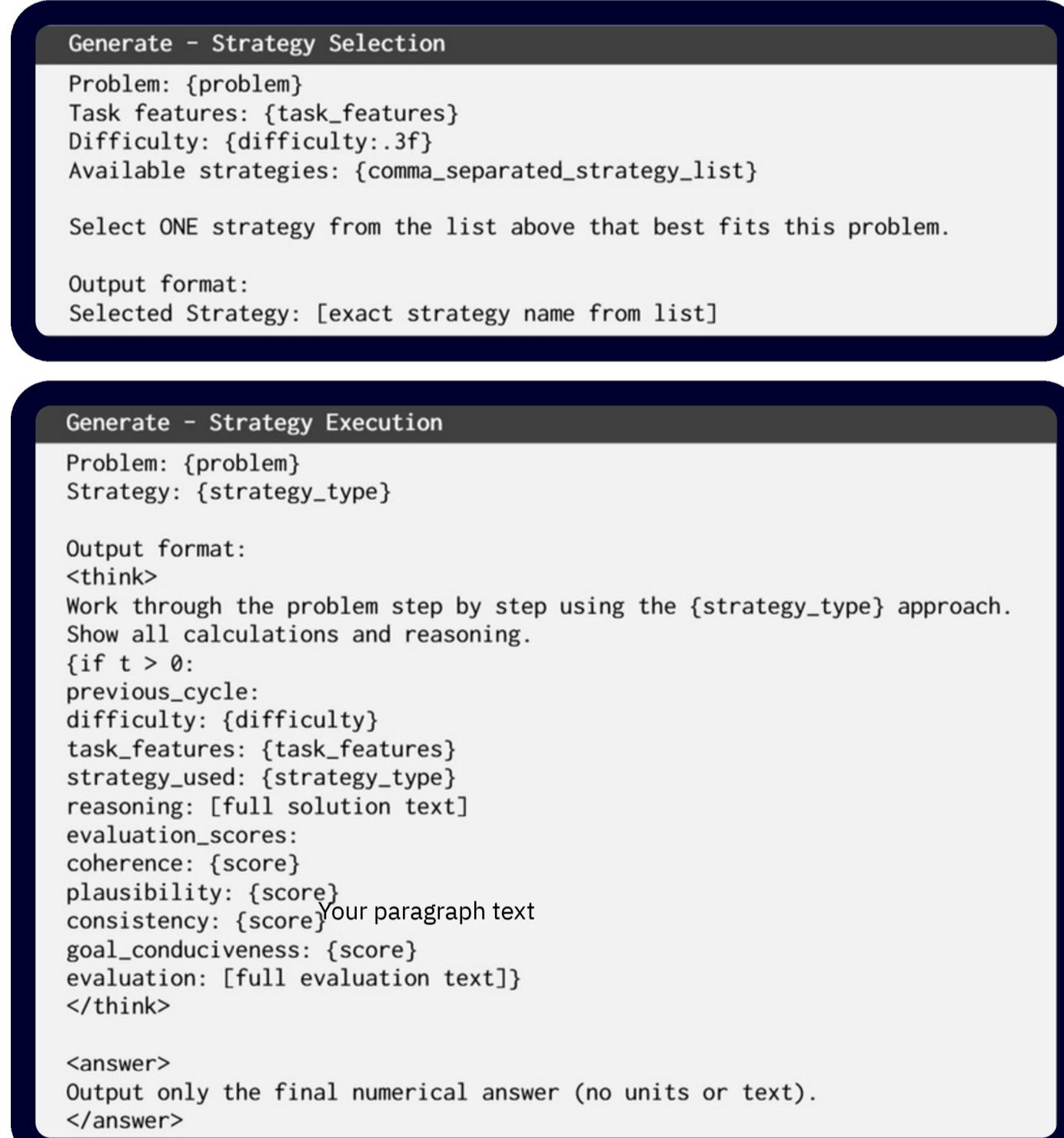
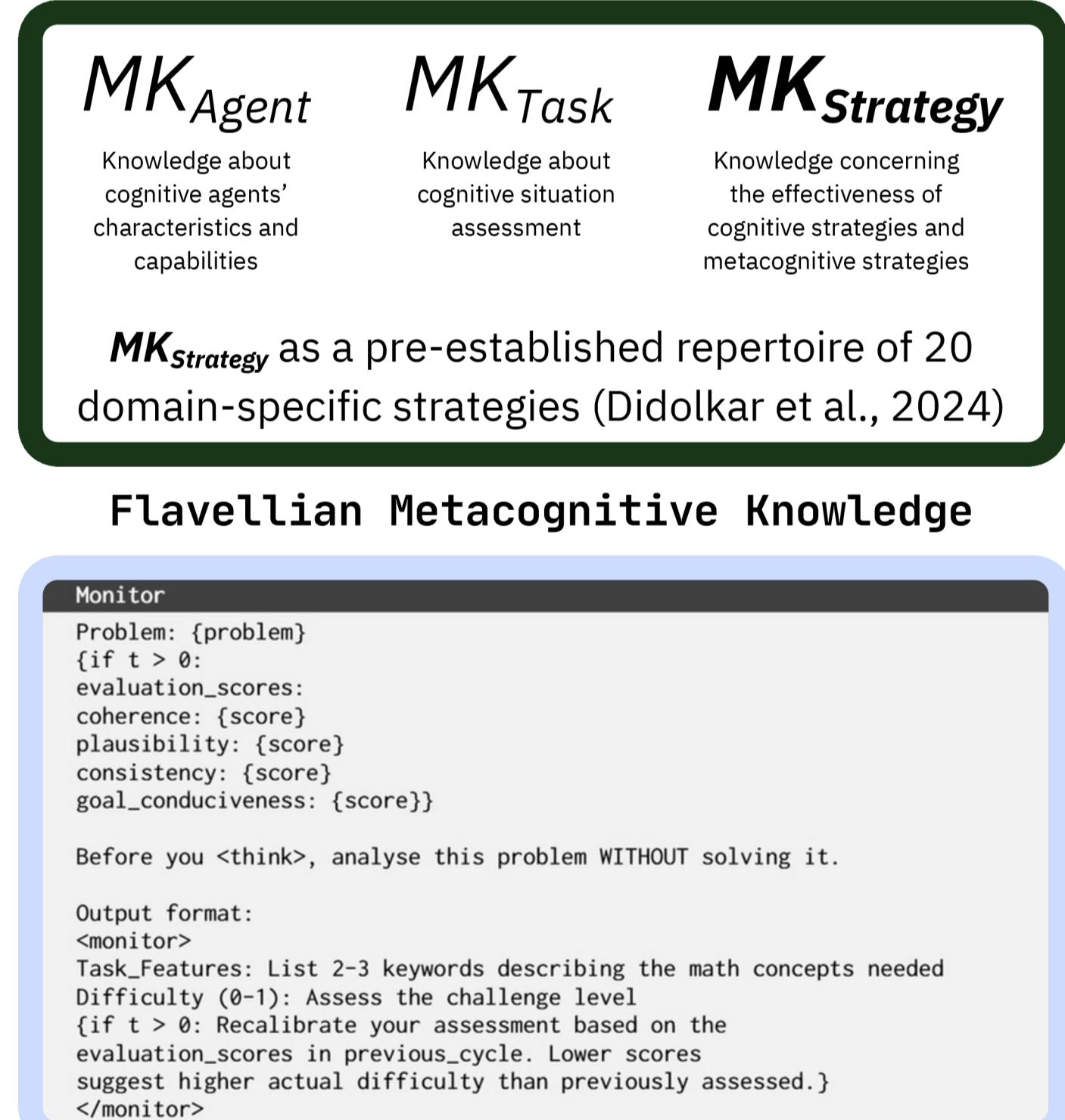
```

### Algorithm 1 Flavell's Model of Cognitive Monitoring

**Require:** task  $\mathcal{T}$ , model  $M$ , strategies  $\mathcal{MK}$ , prompts  $\mathcal{P} = \{p_{\text{monitor}}, p_{\text{strategy}}, p_{\text{execute}}, p_{\text{verify}}\}$

- 1: Initialize  $S_0 = \text{ACTIVE}$ ,  $\tau = 0$ ,  $\mathcal{ME}_{\text{evaluative}}^{-1} = \emptyset$
- 2: **while**  $S_\tau = \text{ACTIVE}$  and  $\tau < T$  **do**
- 3: // **Monitor: Assess task difficulty**
- 4:  $\mathcal{ME}_{\text{difficulty}}^\tau, \text{features}_\tau = M(p_{\text{mon}} \parallel \mathcal{T} \parallel \mathcal{ME}_{\text{evaluative}}^{-1})$  ▷ Assess task difficulty
- 5: // **Generate: Apply cognitive strategy**
- 6:  $\text{strategy}_\tau = M(p_{\text{str}} \parallel \text{features}_\tau \parallel \mathcal{ME}_{\text{difficulty}}^\tau \parallel \mathcal{MK})$  ▷ Choose approach
- 7:  $\text{solution}_\tau = M(p_{\text{exe}} \parallel \mathcal{T} \parallel \text{strategy}_\tau)$  ▷ Execute strategy with adaptive params
- 8: // **Verify: Evaluate performance**
- 9:  $\mathcal{ME}_{\text{evaluative}}^\tau = M(p_{\text{ver}} \parallel \mathcal{T} \parallel \text{solution}_\tau)$  ▷ Evaluate output quality
- 10:  $S_{\tau+1} = \begin{cases} \text{if mean}(\mathcal{ME}_{\text{evaluative}}^\tau) \geq 0.85 \text{ then TERMINATE} \\ \text{else ACTIVE} \end{cases}$
- 11:  $\tau = \tau + 1$
- 12: **end while**
- 13: **return**  $y_{\arg \max_i \text{mean}(\mathcal{ME}_\tau^i)}$

## MGV framework (zero-shot)



## Experimental Setup

meta-llama/Llama-3.1-8B-Instruct (+NVIDIA H100 SXM)

**Self-Verification** (Weng et al., 2022): Generate- $k$ -Verify (w/ majority voting)

**SELF-REFINE** (Madaan et al., 2023): Iterating Generate-Verify  $k$  times

**MGV** (Flavell): Iterating Monitor-Generate-Verify  $k$  times

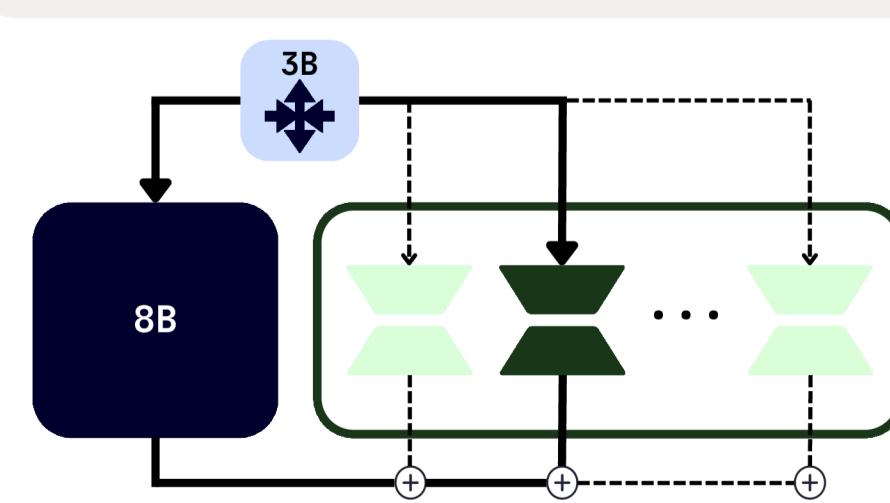
## Results

Method	Accuracy	Avg Time (s)	Avg Attempts
Self-Verification	442/659 (67.07%)	7.52	1.2
Self-REFINE	451/659 (68.44%)	6.98	2.0
MGV (Flavell)	497/659 (75.42%)	9.60	1.3

## Future Work (+ pilot experiment)

### LLM Modular Framework

- Small models (0.5B) can reliably **Plan**, while **Execution** demands substantially larger models (Qin et al., 2025) (for our preliminary exploration of this direction see below)



Branched curriculum (2 epochs shared arithmetic → 3 specialized LoRAs @ 1 epoch each) with 3B routing outperformed 4-epoch GRPO by 1.4pp on GSM8K (8B models)

### Model Intrinsic MK

- Following V-STAR's approach (Hosseini et al., 2024), future work could **train the Monitor through iterative self-supervised learning**: generate diverse solutions per problem, contrast successful versus failed attempts using DPO/GRPO, and progressively refine both difficulty assessment and strategy selection capabilities – building model-intrinsic metacognitive knowledge through bootstrapped preference learning.

### Metacognitive Space within LLM

- Implicit confidence measures derived from token likelihoods exhibit greater metacognitive sensitivity than explicitly prompted confidence (Xiong et al., 2023)
- Some evidence of **subjectivity** (e.g., confidence and certainty) corresponding to linearly separable directions in representations (Zou et al., 2023; Liu et al., 2023)