

# DeepSeek-V3.2: Pushing the Frontier of Open Large Language Models



Authors: DeepSeek-AI



# Motivation and Research Problem



Long Sequences &  
Vanilla Attention



Closed-source fortress  
operates smoothly

Computational Resources: Post-Training



- The paper introduces **DeepSeek-V3.2**, aimed at enhancing AI models' efficiency and performance in reasoning tasks.



- A significant problem is the **inefficiency in processing long sequences** due to reliance on vanilla attention mechanisms.



- Generalization and instruction-following deficits have been noted, along with declining performance in comparison to **closed-source models**. These limitations are critical as they hinder scalability, practical application, and widen the **performance gap**.



- Additionally, open-source models face limitations in **computational resources during post-training**, affecting their ability to handle complex tasks effectively.

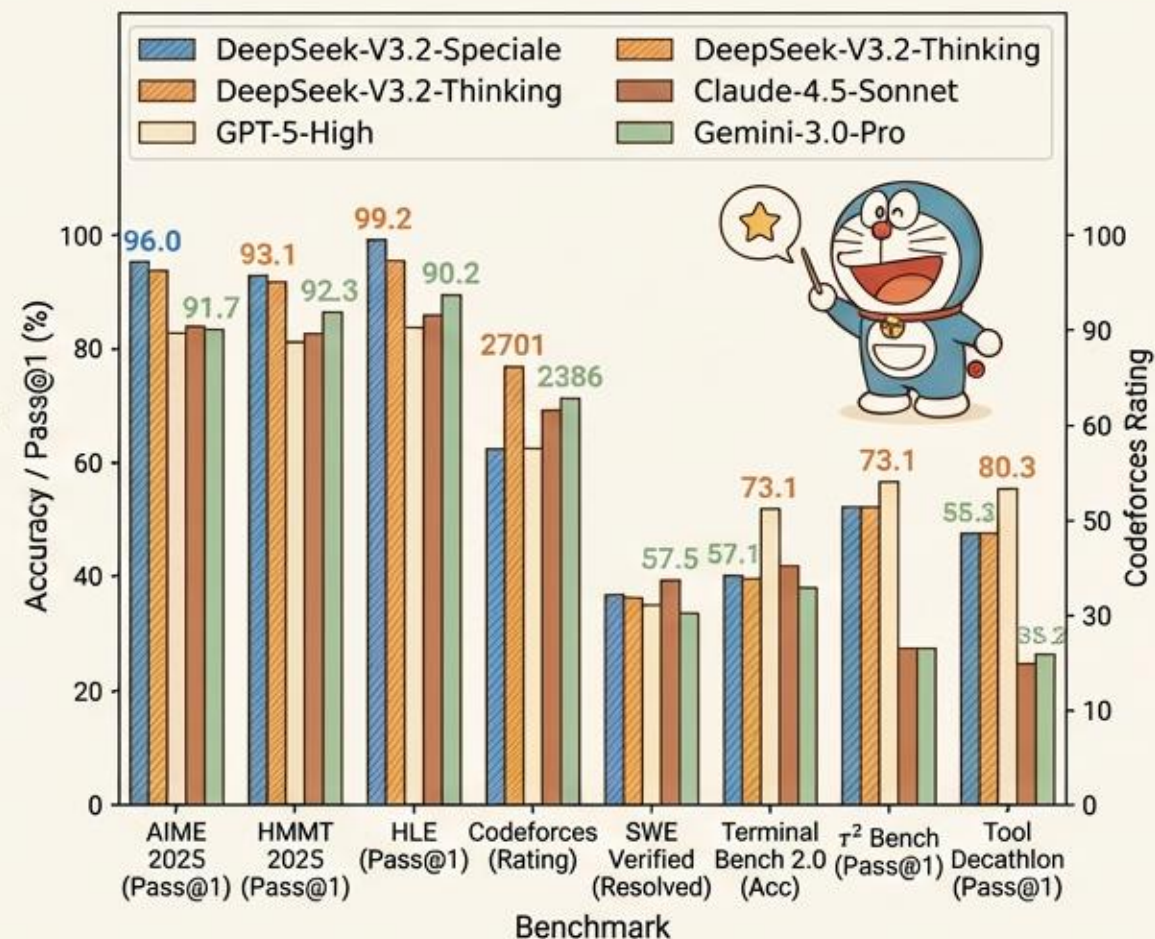
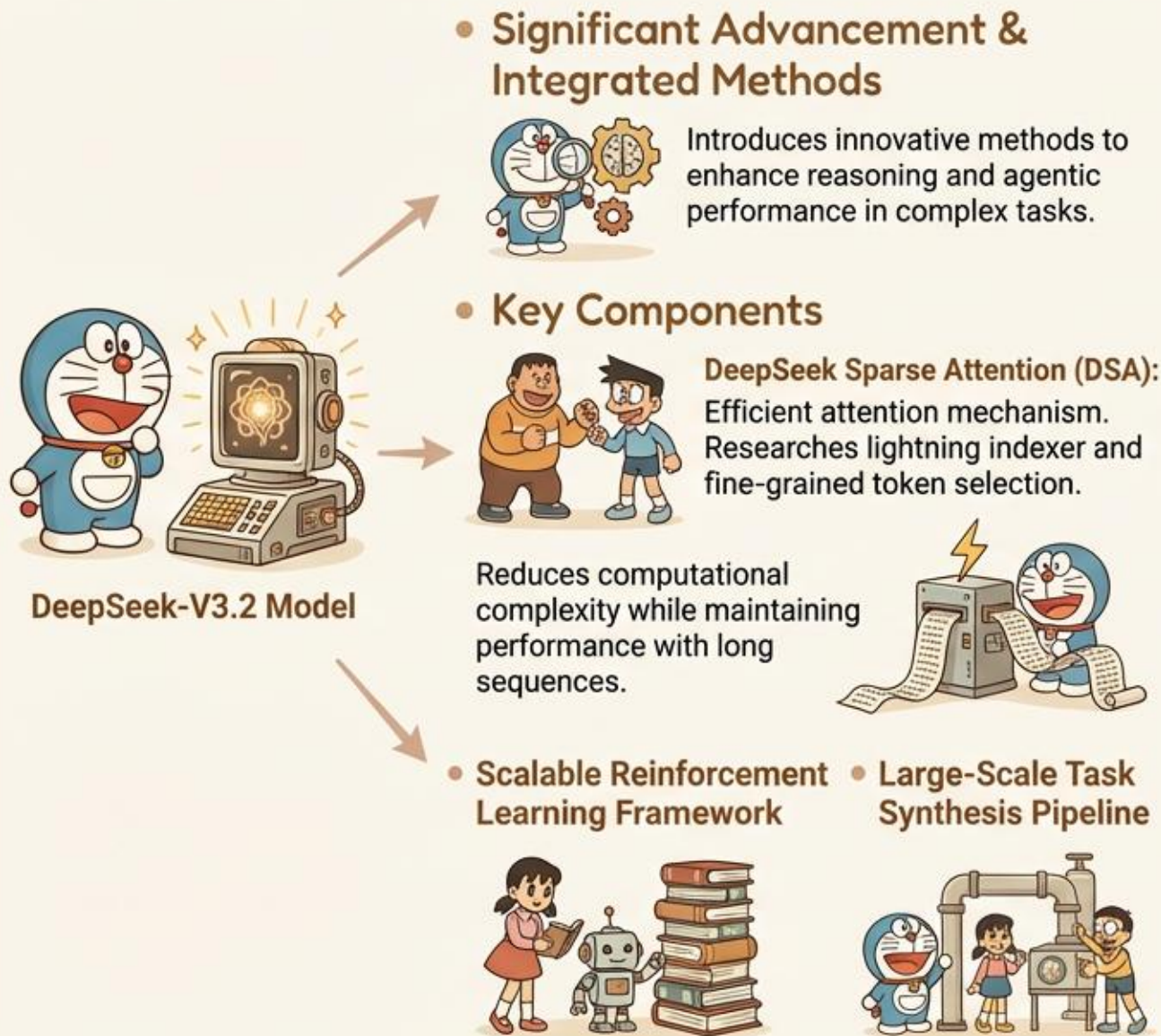


- The research is motivated by these gaps and seeks to develop a model that addresses these inefficiencies.





# Framework Overview



**\*\*Figure 1(focus: architecture): Benchmark of DeepSeek-V3.2 and its counterparts.**

For HMMT 2025, we report the February competition, consistent with the baselines. For HLE, we report the text-only subset.



# Methodology Details



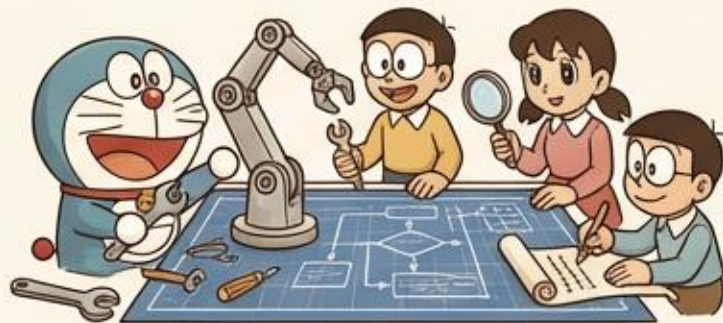
- **DeepSeek Sparse Attention (DSA)** is a critical component, significantly reducing computational complexity. The **lightning indexer** computes relevance scores to determine precedence among tokens, utilizing **top-k** selections.

Table 2 (focus: impact of DSA)

Benchmark	Current Accuracy	Improved Accuracy
MMLU-Pro	87.5	90.1



- The scalable reinforcement learning framework provides robust reinforcement learning protocols and expands computational resources during post-training to boost model performance.
- Lastly, the large-scale agentic task synthesis pipeline integrates reasoning within tool-use scenarios, generating training data to improve the model's generalization and instruction-following capabilities.

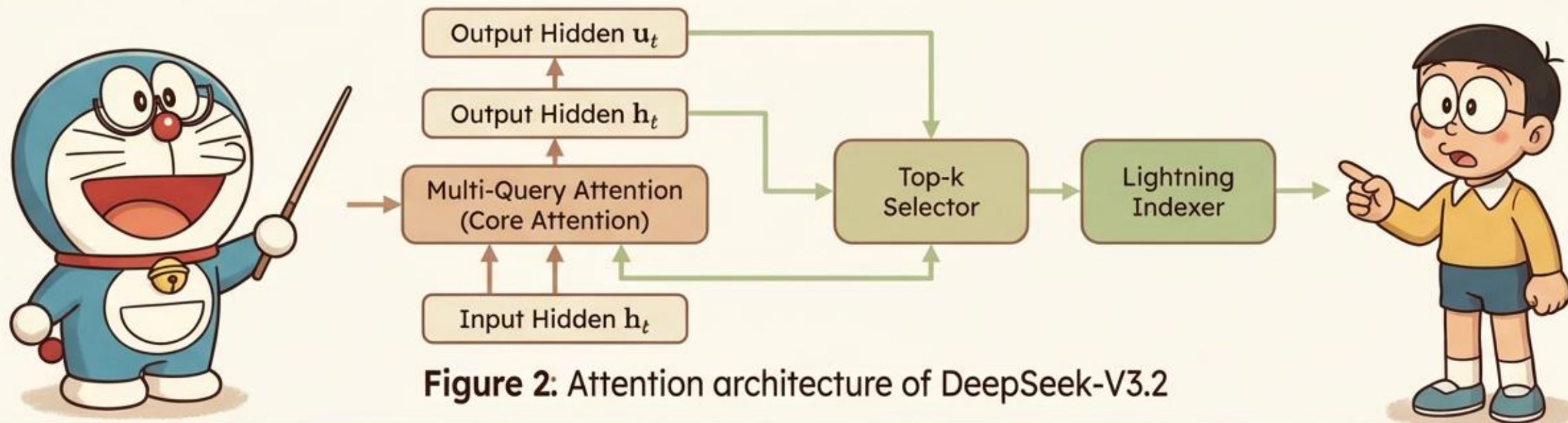




# Mathematical Formulations



DeepSeek-V3.2 employs specific mathematical formulations to optimize its performance. The Sparse Attention Output Equation is given by  $\mathbf{u}_t = \text{Attn}(\mathbf{h}_t, \{\mathbf{c}_s \mid I_{t,s} \in \text{Top-}k(I_{t,:})\})$ , where  $\mathbf{u}_t$  is the attention output vector for the  $t^{\text{th}}$  query token. The lightning indexer calculates index scores with  $I_{t,s} = \sum_{j=1}^{H^I} w_{t,j}^I \cdot \text{ReLU}(\mathbf{q}_{t,j}^I \cdot \mathbf{k}_s^I)$ , determining relevance among tokens.





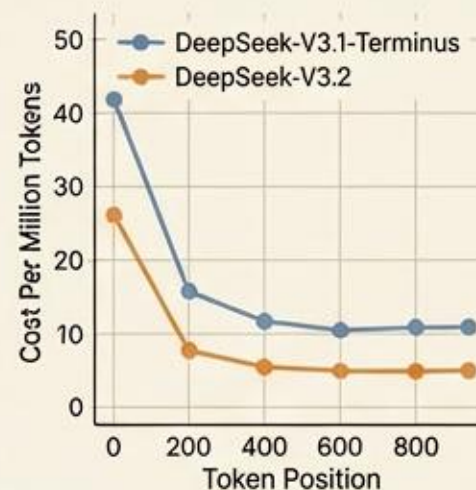
# Results and Experiments



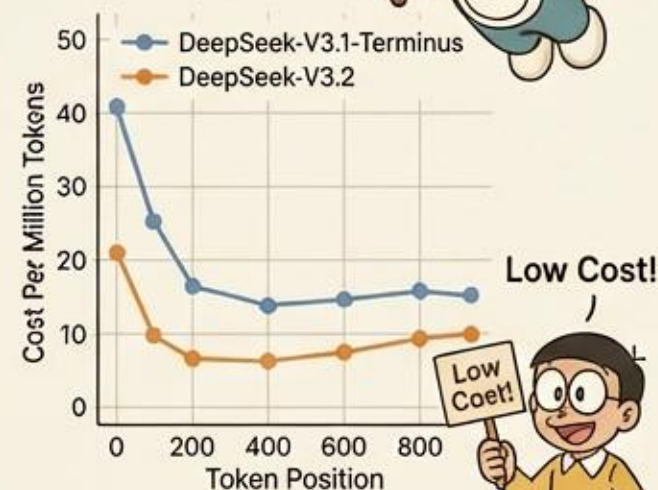
- **DeepSeek-V3.2's** performance was rigorously tested across various benchmarks.
- **Competitive results** in key coding and mathematical competitions, indicating a reduction in the performance gap with proprietary models.
- **Evaluated on datasets:** LiveCodeBench and Codeforces, achieving high accuracy and code resolution scores.
- In **AIME 2025**, it reached a pass rate of **93.1%**.
- **Ablation studies** demonstrated that incorporating the DSA and scalable RL framework significantly improved task efficiency and expanded capability.

**Table 3** (focus: comparison):

Benchmark	DeepSeek-V3.2	Baseline
HMMT Feb 2025	92.5%	88.3%



(a) Prefilling



(b) Decoding

**Figure 3** (focus: performance comparisons): Inference costs of DeepSeek-V3.1-Terminus and DeepSeek-V3.2 on H800 clusters.



# Conclusions

- DeepSeek-V3.2 addresses significant limitations in open-source LLM performance, applying innovative elements like DSA and reinforced RL frameworks to improve reasoning capabilities.
- It successfully reduces computational inefficiencies and enhances task performance.
- The main contributions lie in closing performance gaps with proprietary models, increasing computational efficiency, and demonstrating superiority in complex problem-solving scenarios.
- Key findings include high pass rates in mathematical competitions and notable performance improvements over traditional models.

**Table 4** (focus: competition success):

Competition	Pass Rate
IMO 2025	<b>35/42</b>

