

Benchmark (Metric)		DeepSeek-V3.1-Terminus	DeepSeek-V3.2-Exp
General	MMLU-Pro (EM)	85.0	85.0
	GPQA-Diamond (Pass@1)	80.7	79.9
	Humanity’s Last Exam (Pass@1)	21.7	19.8
Search Agent	BrowseComp (Acc.)	38.5	40.1
	BrowseComp_zh (Acc.)	45.0	47.9
	SimpleQA (Acc.)	96.8	97.1
Code	LiveCodeBench (2408-2505) (Pass@1)	74.9	74.1
	Codeforces-Div1 (Rating)	2046	2121
	Aider-Polyglot (Acc.)	76.1	74.5
Code Agent	SWE Verified (Agent mode)	68.4	67.8
	SWE-bench Multilingual (Agent mode)	57.8	57.9
	Terminal-bench (Terminus 1 framework)	36.7	37.7
Math	AIME 2025 (Pass@1)	88.4	89.3
	HMMT 2025 (Pass@1)	86.1	83.6

Table 1 | Evaluations of DeepSeek-V3.1-Terminus and DeepSeek-V3.2-Exp. Overall, DeepSeek-V3.2-Exp does not show substantial performance degradation compared with DeepSeek-V3.1-Terminus. The performance of DeepSeek-V3.2-Exp on GPQA, HLE, and HMMT 2025 is lower than that of DeepSeek-V3.1-Terminus because DeepSeek-V3.2-Exp generates fewer reasoning tokens. However, this performance gap closes when using intermediate checkpoints that produce a comparable number of tokens.

agent tasks, we employ rule-based outcome reward, length penalty, and language consistency reward. For general tasks, we employ a generative reward model where each prompt has its own rubrics for evaluation. Our reward design carefully balances two key trade-offs: (1) length versus accuracy and (2) language consistency versus accuracy.

3. Evaluations

Model Capabilities. We evaluate DeepSeek-V3.2-Exp on a suite of benchmarks, which focus on diverse capabilities, and compare it with DeepSeek-V3.1-Terminus in Table 1. While DeepSeek-V3.2-Exp significantly improves computational efficiency on long sequences, we do not observe substantial performance degradation compared with DeepSeek-V3.1-Terminus, on both short- and long-context tasks. In addition, we also compare the reinforcement learning training curves of DeepSeek-V3.2-Exp and DeepSeek-V3.1-Terminus, as shown in Figure 2. The performance of both models on BrowseComp and SWE Verified improves steadily throughout the training process, with closely aligned curves, which reflects the training stability of DSA.

Inference Costs. DSA reduces the core attention complexity of the main model from $O(L^2)$ to $O(Lk)$, where k ($\ll L$) is the number of selected tokens. Although the lightning indexer still has a complexity of $O(L^2)$, it requires much less computation compared with MLA in DeepSeek-V3.1-Terminus. Combined with our optimized implementation, DSA achieves a significant end-to-end speedup in long-context scenarios. Figure 3 presents how token costs of DeepSeek-V3.1-Terminus and DeepSeek-V3.2-Exp vary with the token position in the sequence. These costs are estimated from benchmarking the actual service deployed on H800 GPUs, at