

Table 10 | Comparison of DeepSeek-R1 and other frontier models in fine-grained safety scenarios. **Unsafe** indicates the proportion of unsafe content in the model’s responses (lower values indicate better model safety), while **Rej.** represents the rejection rate in the model’s answers (lower values indicate a stronger tendency for the model to provide informative and safe answers to questions, rather than simply declining to respond). For DeepSeek-V3 and DeepSeek-R1, we report results under two configurations: with and without risk control system (introduced in D.3.1).

	Discrimi.		Illegal		Harmful		Ethical		Overall	
Ratio(%)	Unsafe	Rej.	Unsafe	Rej.	Unsafe	Rej.	Unsafe	Rej.	Unsafe	Rej.
Claude-3.7-Sonnet	8.4	2.5	14.1	4.5	9.5	5.5	7.5	0.6	10.7	3.6
o1 (2024-12-17)	7.2	37.8	12.3	54.8	5.0	73.5	8.8	34.4	9.0	50.4
GPT-4o (2024-05-13)	19.1	6.2	22.5	28.4	28.0	19.5	18.8	4.4	22.0	17.1
Qwen2.5 Instruct (72B)	12.8	2.5	14.5	9.5	15.5	5.0	11.9	0.0	13.8	5.4
DeepSeek-V3	20.3	2.5	17.3	13.9	17.5	9.5	13.1	1.9	17.6	8.1
+ risk control system	8.1	16.9	3.2	35.5	7.0	22.5	3.1	18.1	5.3	25.4
DeepSeek-R1	19.7	3.8	28.9	8.6	32.5	6.0	16.9	0.6	25.2	5.6
+ risk control system	9.1	17.2	6.6	39.1	13.0	29.0	6.9	13.1	8.5	27.3

answers, with lower values being more desirable (we prefer safe responses over rejections since it can provide risk warning information).

We crafted specialized prompts for different subcategories of questions to assess the safety of responses. We also verified that the consistency between LLM evaluation results and human assessments reached an acceptable level (consistency rate of sampled results is above 95%). The experimental comparison results are presented in Table 10, from which the following conclusions can be observed:

- **Analyzing unsafe rates:** DeepSeek-V3 (with risk control) belongs to the first tier of safe models (unsafe rate around 5%); DeepSeek-R1 (with risk control), Claude-3.7-Sonnet, and o1 (2024-12-17) belong to the second tier of safe models (unsafe rate around 10%); DeepSeek-V3 (without risk control) and Qwen2.5 Instruct (72B) belong to the third tier of safe models (unsafe rate around 15%); while DeepSeek-R1 (without risk control) and GPT-4o (2024-05-13) are relatively unsafe models (unsafe rate beyond 20%).
- **Analyzing rejection rates:** The base models of DeepSeek-R1 and DeepSeek-V3 have relatively low rejection rates but higher unsafe rates. After implementing a risk control system, these models show relatively low unsafe rates but higher rejection rates (around 25%). Additionally, Claude-3.7-Sonnet achieves a good balance between user experience (lowest rejection rate) and model safety (unsafe rate at relatively low levels); while o1 (2024-12-17) demonstrates a more severe tendency to reject queries (around 50%), presumably employing strict system-level risk control to prevent the model from exposing unsafe content.
- **Analyzing risk types:** DeepSeek-R1 performs exceptionally well in handling queries related to Illegal and Criminal Behavior and Moral and Ethical Issues, while showing average performance in scenarios involving Discrimination and Prejudice Issues and Harmful Behavior, which encourages us to pay more attention on these two categories when developing model safety features and risk control system.

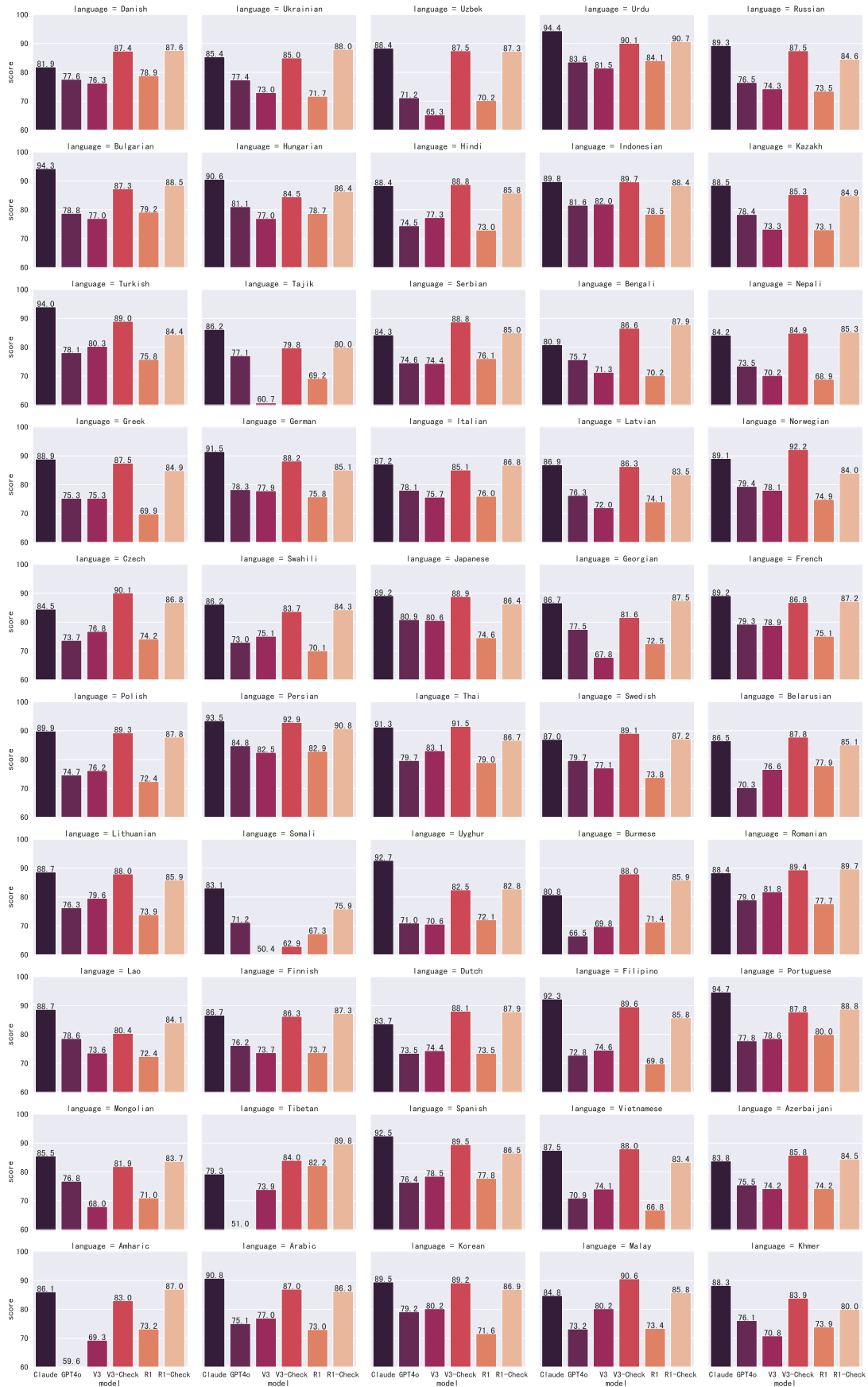


Figure 14 | Multilingual safety performance. V3-check and R1-check represent the risk control system evaluation results for DeepSeek-V3 and DeepSeek-R1, respectively.

D.3.4. Multilingual Safety Performance

In the previous section’s evaluation, we primarily focused on the model’s safety performance in special languages (Chinese and English). However, in practical usage scenarios, users’ linguistic backgrounds are highly diverse. Assessing safety disparities across different languages is essential. For this purpose, we translated the original bilingual safety testset (introduced in the D.3.3) into 50 commonly used languages. For high-frequency languages, we conducted full translation of the entire dataset, while for low-frequency languages, we performed sampling translation. This process resulted in a comprehensive multilingual safety test set consisting of 9,330 questions. During the translation process, we employed a combined approach of LLM translation and human-assisted calibration to ensure the quality of the translations.

We continued to use the LLM-as-a-judge methodology described in the previous section, which determines safety labels (safe, unsafe, or rejected) for each question-answer pair. Rather than merely rejecting risky queries, we prefer responses that provide safe content; therefore, we assigned higher scores to safe responses (5 points per question, with 5 points for safe responses, 0 points for unsafe responses, and 4 points for rejections). The final safety score proportions (safety score as a percentage of the total possible safety score) across 50 languages are presented in Figure 14. For DeepSeek-V3 and DeepSeek-R1, we evaluated safety scores for models with and without the risk control system (introduced in D.3.1). Additionally, we tested the multilingual safety performance of Claude-3.7-Sonnet and GPT-4o(2024-05-13). From Figure 14, we can draw the following conclusions:

- With risk control system in place, DeepSeek-V3 (86.5%) and DeepSeek-R1 (85.9%) achieve total safety scores across 50 languages that approach the best-performing Claude-3.7-Sonnet (88.3%). This demonstrates that DeepSeek has reached state-of-the-art levels in system-level multilingual safety.
- Without risk control system, DeepSeek-V3 (75.3%) and DeepSeek-R1 (74.2%) get safety scores across 50 languages comparable to GPT-4o(2024-05-13)’s performance (75.2%). This indicates that even when directly using the open-source versions of R1, the model still exhibits a moderate level of safety standard.
- Examining language-specific weaknesses, we categorize languages with safety scores below 60 points as high-risk languages for the corresponding model. Among the 50 languages evaluated, DeepSeek-R1 (without risk control system) and Claude-3.7-Sonnet have zero high-risk languages; DeepSeek-V3 (without risk control system) and GPT-4o(2024-05-13) have one and two high-risk languages, respectively. This suggests that DeepSeek-R1 has no obvious language-specific vulnerabilities.

D.3.5. Robustness against Jailbreaking

In real-world application scenarios, malicious users may employ various jailbreaking techniques to circumvent a model’s safety alignment and elicit harmful responses. Therefore, beyond evaluating model safety under direct questioning, we place significant emphasis on examining the model’s robustness when confronted with jailbreaking attacks. Thus, we constructed a dedicated test suite for jailbreaking evaluation. Specifically, we developed a template collection consisting of 2,232 jailbreaking instructions. We then randomly concatenated these jailbreaking prompts with questions from the original safety testset (introduced in D.3.3) and further examined the performance differences in the model’s responses when confronted with original unsafe questions versus newly formulated questions with jailbreaking elements.

When evaluating the results, we followed the LLM-as-a-Judge safety assessment (introduced

in D.3.3), while improving the safety evaluation prompts to focus more specifically on identifying manipulative traps in jailbreak attempts. Each question-answer pair was classified into one of three categories: safe, unsafe, or rejected (introduced in D.3.3). The results of jailbreak attacks against various models are presented in Table 11. From these results, we draw the following conclusions:

Table 11 | Comparison of DeepSeek-R1 and other frontier models in jailbreaking scenarios.

	Unsafe Ratio			Rejected Ratio		
Ratio(%)	Origin	Jailbreak	GAP	Origin	Jailbreak	GAP
Claude-3.7-Sonnet o1 (2024-12-17)	10.7	26.2	+15.5	3.6	21.9	+18.3
GPT-4o (2024-05-13)	9.0	12.1	+3.1	50.4	79.8	+29.4
GPT-4o (2024-05-13)	22.0	30.4	+8.4	17.1	57.3	+40.2
Qwen2.5 Instruct (72B)	13.8	29.7	+15.9	5.4	25.2	+19.8
DeepSeek-V3	17.6	36.4	+18.8	8.1	8.9	+0.8
+ risk control system	5.3	2.3	-3.0	25.4	46.5	+21.1
DeepSeek-R1	25.2	85.9	+60.7	5.6	1.9	-3.7
+ risk control system	8.5	4.3	-4.2	27.3	87.3	+60.0

- All tested models exhibited significantly increased rates of unsafe responses and rejections, along with decreased safety rates when facing jailbreak attacks. For example, Claude-3.7-Sonnet, showed a 33.8% decrease in the proportion of safe responses when confronted with our security jailbreak attacks. This demonstrates that current cutting-edge models still face substantial threats from jailbreak attacks.
- Compared to non-reasoning models, the two reasoning models in our experiments — DeepSeek-R1 and o1(2024-12-17) — rely more heavily on the risk control system for security checks, resulting in considerably higher overall rejection rates (79.8% and 87.3% respectively).
- Open-source models (DeepSeek, Qwen) face more severe jailbreak security challenges than closed-source models, because of the lack of a risk control system in locally deployed models. To address safety issues, we advise developers using open source models in their services to adopt comparable risk control measures.

E. More Analysis

E.1. Performance Comparison with DeepSeek-V3

Since both DeepSeek-R1 and DeepSeek-V3 share a common base architecture, namely DeepSeek-V3-Base, a critical question naturally arises: which specific dimensions are enhanced through the application of different post-training techniques? To address this, we first compare the R1 family of models with DeepSeek-V3 and DeepSeek-V3-Base, as summarized in Table 12. Notably, DeepSeek-R1 demonstrates significant improvements in competitive programming and mathematical reasoning tasks, as evidenced by superior performance on benchmarks such as LiveCodeBench and AIME 2024. These enhancements in reasoning capabilities also translate into higher scores on the Arena-Hard evaluation suite. Furthermore, DeepSeek-R1 exhibits stronger long-context understanding, as indicated by its improved accuracy on the FRAMES

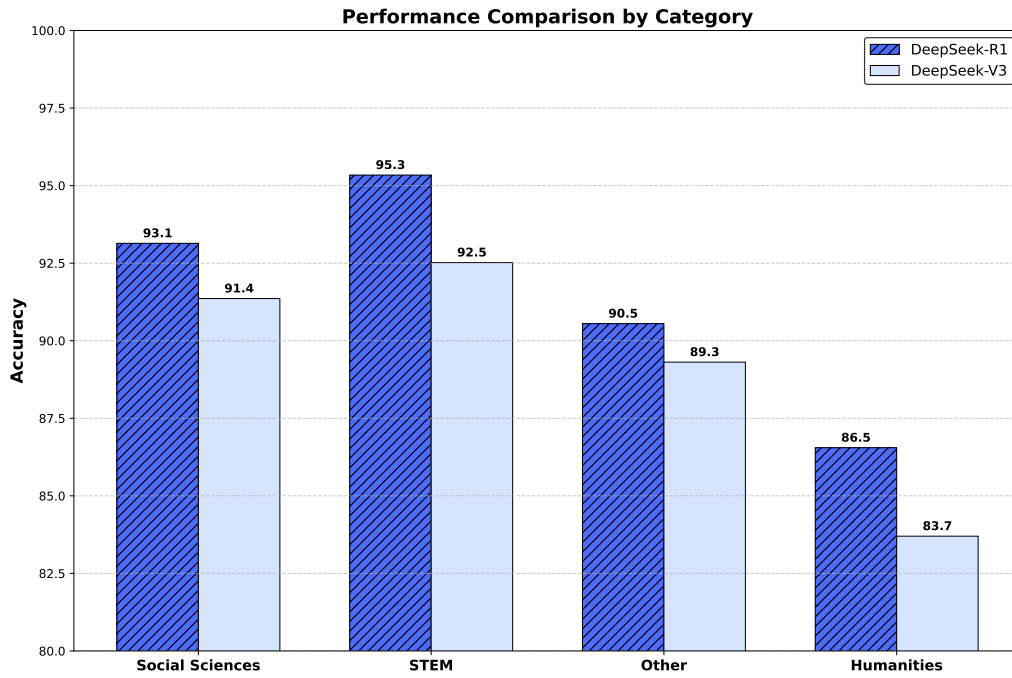


Figure 15 | The comparison of DeepSeek-V3 and DeepSeek-R1 across MMLU categories.

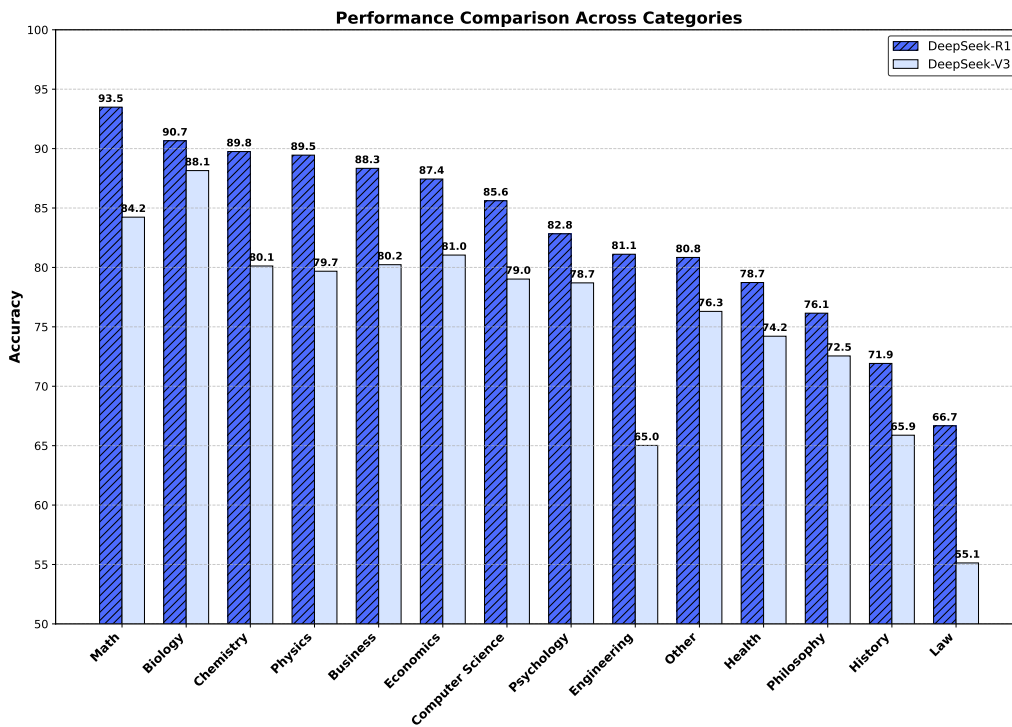


Figure 16 | The comparison of DeepSeek-V3 and DeepSeek-R1 across MMLU-Pro categories.