

Integral Solution Analysis 2

Comprehensive Integral Solution Analysis: $\int \frac{\tan^3(\ln x)}{x} dx$

This document presents a technical evaluation of 23 Large Language Model (LLM) versions on their ability to solve the indefinite integral:

$$\int \frac{\tan^3(\ln x)}{x} dx$$

Correct Solution:

$$\frac{1}{2} \tan^2(\ln x) + \ln |\cos(\ln x)| + C$$

(Equivalent forms using $\sec^2(\ln x)$ or $-\ln |\sec(\ln x)|$ are also accepted).

1 Master Results Table

Model	Size	Architecture	Status	Key Failure/Success Factor
Gemma 3 1B	1B	Dense (MatFormer)	Fail	Hallucinated algebraic identities due to compression.
Gemma 3 4B	4B	Dense (MatFormer)	Partial	Correct math, but context window failure (forgot u substitution).
Gemma 3 12B	12B	Dense (MatFormer)	Pass	The stability threshold for the Gemma 3 family.
Gemma 3 27B	27B	Dense (MatFormer)	Pass	Robust performance; effectively used \sec^2 identity.
Mistral 7B	7B	Dense	Pass	Outperformed Llama 70B via strict instruction following.
Nemotron Nano 9B v2	9B	Hybrid (Mamba-2)	Pass	Mamba Architecture maintained state perfectly.
Nemotron Nano 12B VL	12.6B	Hybrid (Mamba-2)	Pass	Best Verification: Rejected hallucinated calculator data.
Trinity Mini	26B	Sparse MoE (3B Active)	Fail	Routing Failure: 3B active parameters dropped the sign context.
Tongyi DeepResearch	30.5B	Sparse MoE (3.3B Active)	Pass	Agentic training compensated for low active parameter count.
Olmo 3 32B Think	32B	Dense (RLVR)	Pass	Metacognition: Explicitly doubted and verified itself.
Llama 3.3 70B	70B	Dense	Fail	Rate Error: Failed power rule integration ($\int x \neq x$).
Cognito Llama 70B	70B	Llama Fine-tune	Fail	“Reasoning” fine-tune caused logical over-complication.
Cognito Llama 405B	405B	Llama Fine-tune	Pass	Scale corrected the logic errors seen in the 70B version.
GLM 4.5 Air	106B	MoE (12B Active)	Pass	Sufficient active memory for symbolic logic.
Devstral 2 2512	123B	Dense (Coding)	Pass	Coding Specialist: Treated math like syntax parsing.
Nova 2 Lite	120B	Distilled	Pass	Efficient knowledge distillation from larger teacher.

Model	Size	Architecture	Status	Key Failure/Success Factor
KAT-Coder-Pro V1	~40B	Agentic Coding	Pass	Structured the solution as a “Plan of Action”.
LongCat Flash	560B	MoE (27B Active)	Pass	“Flash-Thinking” dynamic compute ensured accuracy.
DeepSeek V3.1 Next N1	671B	MoE (37B Active)	Pass	Agentic post-training verified pre-conditions (dx/x).
DeepSeek Chimera R1T	671B	MoE (Reasoning)	Pass	Verified result via differentiation.
DeepSeek Chimera R1T2	671B	MoE (Reasoning)	Pass	Tri-Mind Arch: Fixed consistency issues of V1.
Kimi K2	1T	MoE (32B Active)	Pass	Verified: Differentiated result to check correctness.
Grok 4.1 Fast	Unk	Agentic	Pass	Fast, accurate derivation.

2 Detailed Model Analysis

2.1 The Gemma 3 Family (1B - 27B)

- **Gemma 3 1B: Failed.** The 1B model, optimized for mobile deployment, displayed “hallucination under load.” It correctly identified the u -substitution but invented algebraic identities to force a simplification that didn’t exist.
- **Gemma 3 4B: Partial.** This model solved the calculus correctly but failed the “Instruction Constraint” (Substitute back to x). This is a classic “Context Window Failure” common in $< 7B$ models, where the initial constraint is overwritten by the cognitive load of the math steps.
- **Gemma 3 12B & 27B: Success.** The 12B model appears to be the **Minimum Viable Size** for reliable symbolic calculus in the Gemma architecture, showing none of the drift present in the smaller variants.

2.2 Hybrid Mamba vs. Sparse MoE (The “Small” Model War)

- **Nemotron Nano 12B 2 VL (Hybrid Mamba): Star Performer.** This model uses a **Hybrid Transformer-Mamba** architecture. The Mamba (State Space Model) component acts as a recurrent memory, allowing it to “hold” the variable definitions ($u = \ln x$) indefinitely without the compute cost of attention. It uniquely **rejected a hallucinated calculator result** during its reasoning trace, trusting its own derivation instead.
- **Trinity Mini (Sparse MoE): Failed.** While having 26B total parameters, it only uses **3B Active Parameters** per token. This proved insufficient. The “Math Expert” layer likely lost the negative sign context during routing, a known weakness of Sparse MoEs when the active parameter count drops below $\sim 7B$.
- **Tongyi DeepResearch (Sparse MoE): Success.** Despite also having only **3.3B Active Parameters**, it succeeded. Its training on **Long-Horizon Agentic Research** enforces a “step-by-step verification” protocol that Trinity lacks, effectively brute-forcing accuracy through rigid logical structure.

2.3 The “Thinking” Models (RLVR & Reasoning)

- **Olmo 3 32B Think: Success.** This model is trained with **Reinforcement Learning with Verifiable Rewards (RLVR)**. Its transcript was unique: it explicitly “doubted” its own sign convention for $\int \tan u$, paused, and chose to differentiate the result to be sure. It proves that a 32B “Thinker” can outperform a 70B “Predictor.”
- **DeepSeek R1T2 Chimera: Success.** An “Assembly of Experts” merging the reasoning of R1 with the chat fluency of V3. It fixed the “over-thinking” loops sometimes seen in R1T, arriving at the solution efficiently while still verifying it via differentiation.

2.4 Large Dense Models (70B+)

- **Llama 3.3 70B: Failed.** A significant failure for a 70B model. It attempted to integrate $\int x dx$ but failed the power rule, a sign that its training data prioritized “chat fluency” over “mathematical precision”.
- **Devstral 2 2512 (123B): Success.** As a **Coding-Optimized** model, Devstral treated the integral like a syntax parsing problem. It broke the “code” (equation) into functions (substitution, identity), executed them, and returned the output with zero “fluff”.

2.5 Massive Scale MoEs ($\geq 500B$)

- **Kimi K2 (1T Parameters): Success.** With 1 Trillion parameters, Kimi K2 has highly specialized experts. It likely routed this prompt to a specific “Latex/Calculus” expert trained on arXiv papers, resulting in the most formally mathematically correct notation of the group.
- **DeepSeek V3.1 Nex N1: Success.** This “Agentic” variant of V3 structured its response as a **Workflow**. It validated the “Precondition” (dx/x exists) before executing the “Action” (Substitution), a behavior derived from its autonomous agent post-training.

3 Conclusion

The analysis of these 23 models highlights three critical trends in 2025 AI architecture:

1. **Hybrid Mamba Wins:** The **Nemotron 12B** (Hybrid) verified its answer better than the **Llama 70B** (Dense), proving that State Space Models are superior for logical state retention.
2. **Active Parameter Threshold:** For MoEs, **12B Active Parameters** (GLM 4.5, Kimi) is the safe zone. Dipping to **3B Active** (Trinity) risks logic errors unless compensated by specific Agentic training (Tongyi).
3. **Metacognition:** The strongest models (Olmo Think, DeepSeek R1, Kimi) all shared one trait: **Self-Verification**. They differentiated their answers to check for errors, a step the older-style models (Llama 3.3) failed to take.