

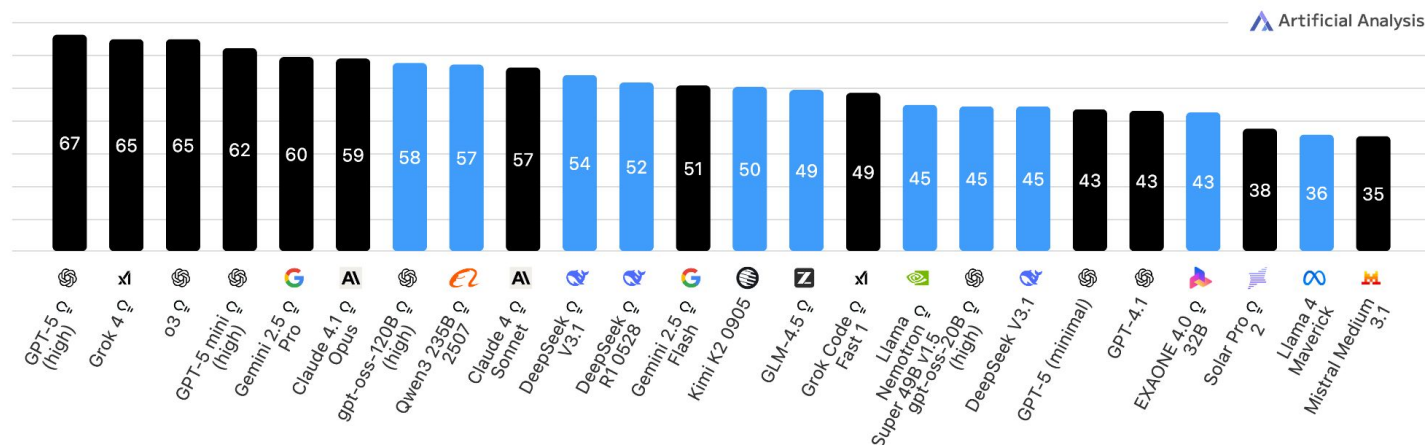
2. Preliminaries

The Landscape of LLMs: open vs closed



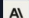


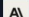





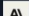





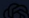












Artificial Analysis Intelligence Index by Open Weights vs Proprietary

Artificial Analysis Intelligence Index v3.0 incorporates 10 evaluations: MMLU-Pro, GPQA Diamond, Humanity's Last Exam, LiveCodeBench, SciCode, AIME 2025, IFBench, AA-LCR, Terminal-Bench Hard, τ^2 -Bench Telecom



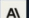


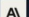





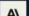


















■ Proprietary ■ Open Weights



The Landscape of LLMs: Open vs Closed

<div>Bash Only Verified Lite Full Multimodal</div>						
Bash Only evaluates all LMs with a minimal agent on SWE-bench Verified (details)						
Filters: Open Scaffold ▾ All Tags ▾						
Model	% Resolved	Avg. \$	Org	Date	Release	
  Claude 4.5 Sonnet (20250929)	70.60	\$0.56		2025-09-29	1.13.3	
  Claude 4 Opus (20250514)	67.60	\$1.13		2025-08-02	1.0.0	
  GPT-5 (2025-08-07) (medium reasoning)	65.00	\$0.28		2025-08-07	1.7.0	
  Claude 4 Sonnet (20250514)	64.93	\$0.37		2025-07-26	1.0.0	
  GPT-5 mini (2025-08-07) (medium reasoning)	59.80	\$0.04		2025-08-07	1.7.0	
  o3 (2025-04-16)	58.40	\$0.33		2025-07-26	1.0.0	
  Qwen3-Coder 480B/A35B Instruct	55.40	\$		2025-08-02	1.0.0	
  GLM-4.5 (2025-08-22)	54.20	\$0.30		2025-08-22	1.9.1	
  Gemini 2.5 Pro (2025-05-06)	53.60	\$0.29		2025-07-26	1.0.0	
  Claude 3.7 Sonnet (20250219)	52.80	\$0.35		2025-07-20	0.0.0	

The Landscape of LLMs: Open vs Closed

<div>Bash Only Verified Lite Full Multimodal</div>						
Bash Only evaluates all LMs with a minimal agent on SWE-bench Verified (details)						
Filters: Open Scaffold ▾ All Tags ▾						
Model		% Resolved	Avg. \$	Org	Date	Release
  Claude 4.5 Sonnet (20250929)		70.60	\$0.56		2025-09-29	1.13.3
  Claude 4 Opus (20250514)		67.60	\$1.13		2025-08-02	1.0.0
  GPT-5 (2025-08-07) (medium reasoning)		65.00	\$0.28		2025-08-07	1.7.0
  Claude 4 Sonnet (20250514)		64.93	\$0.37		2025-07-26	1.0.0
  GPT-5 mini (2025-08-07) (medium reasoning)		59.80	\$0.04		2025-08-07	1.7.0
  o3 (2025-04-16)		58.40	\$0.33		2025-07-26	1.0.0
  Qwen3-Coder 480B/A35B Instruct		55.40	\$		2025-08-02	1.0.0
  GLM-4.5 (2025-08-22)		54.20	\$0.30		2025-08-22	1.9.1
  Gemini 2.5 Pro (2025-05-06)		53.60	\$0.29		2025-07-26	1.0.0
  Claude 3.7 Sonnet (20250219)		52.80	\$0.35		2025-07-20	0.0.0

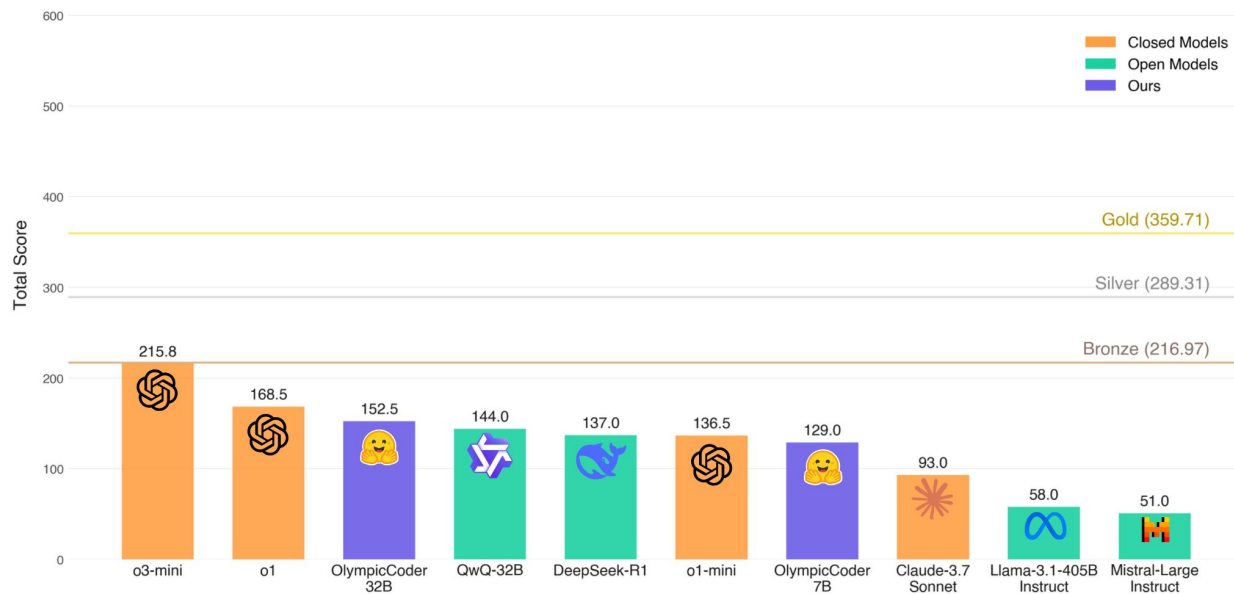
Open
models!



The Landscape of LLMs: Open vs Closed

International Olympiad in Informatics (2024)

Performance of 10 selected models across 50 submissions



What Are Code LMs?

**Large Language Models trained in part on source code;
Or LLMs specialized for coding.**

Pre-training Recipes

- **Coding is capability within a general LLM:**
 - Qwen3, Kimi k2, DeepSeek v2

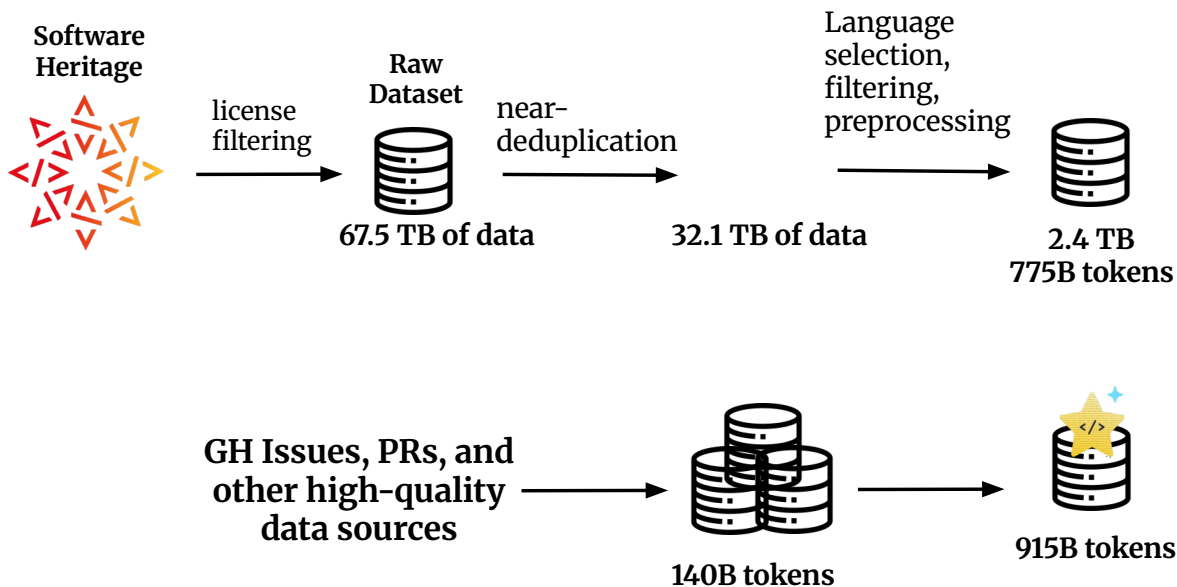
Pre-training Recipes

- **Coding is capability within a general LLM:**
 - Qwen3, Kimi k2, DeepSeek v2
- **Continual pre-training from an LLM:**
 - Qwen3-Coder, DeepSeek Coder v2

Pre-training Recipes

- **Coding is capability within a general LLM:**
 - Qwen3, Kimi k2, DeepSeek v2
- **Continual pre-training from an LLM:**
 - Qwen3-Coder, DeepSeek Coder v2
- **Code LMs trained from scratch:**
 - Codex, CodeGen, StarCoder, StarCoder2

Pre-training data: GitHub



The Stack v2 pipeline

Pre-training data: CommonCrawl

- DeepSeekCoder v2 **70B** tokens
- Qwen2.5 Coder 580B tokens filtered down to **110B** tokens

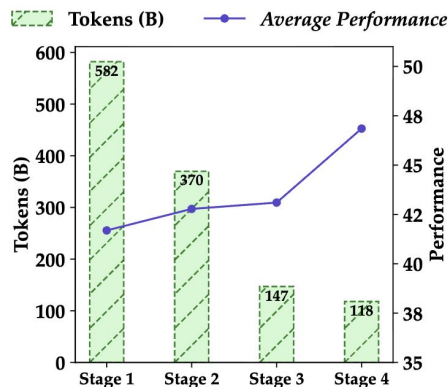


Figure 1: Number of data tokens across different cc-stages, and the validation effectiveness of training Qwen2.5-Coder using corresponding data.

Qwen2.5 Coder web code data filtering

Pre-training data: Synthetic Data

- Qwen2.5 Coder
- Arctic-SnowCoder

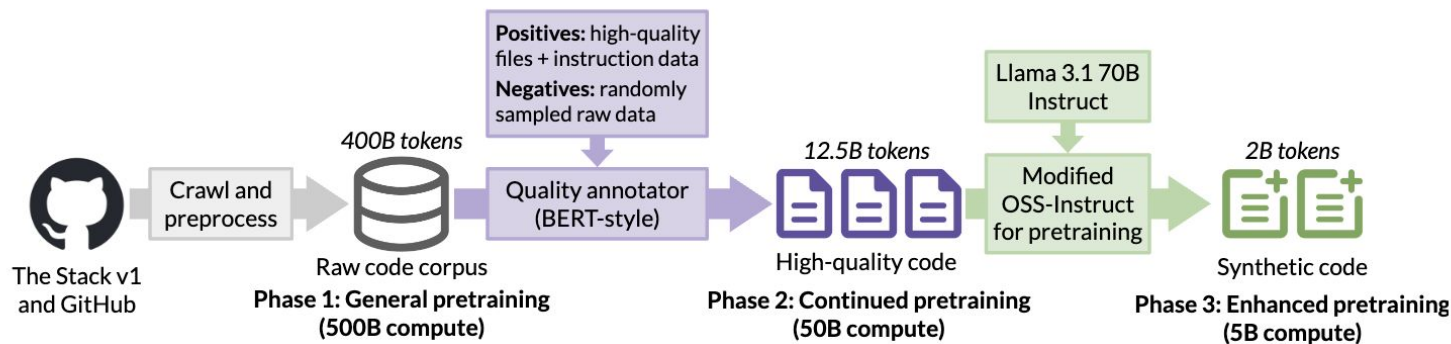


Figure 1: Three-phase pretraining of Arctic-SnowCoder-1.3B with progressively higher-quality data.

Pre-training data: Synthetic Data Filtering

- GitHub Code filtered by an LLM: the LLM/classifier rates the educational quality from 0 to 5

Below is an extract of SQL code. Evaluate whether it has a high educational value and could help teach SQL and database concepts. Use the additive 5-point scoring system described below. Points are accumulated based on the satisfaction of each criterion:

- Add 1 point if the code contains valid SQL syntax, even if it's just basic queries or simple table operations.
- Add another point if the code addresses practical SQL concepts (e.g., JOINs, subqueries), even if it lacks comments.
- Award a third point if the code is suitable for educational use and introduces key concepts in SQL and database management, even if the topic is somewhat advanced (e.g., indexes, transactions). The code should be well-structured and contain some comments.
- Give a fourth point if the code is self-contained and highly relevant to teaching SQL. It should be similar to a database course exercise, demonstrating good practices in query writing and database design.
- Grant a fifth point if the code is outstanding in its educational value and is perfectly suited for teaching SQL and database concepts. It should be well-written, easy to understand, and contain explanatory comments that clarify the purpose and impact of each part of the code.

The extract:
--- MySQL mode ---
select name
from employee
order by name asc;

After examining the extract:

- Briefly justify your total score, up to 100 words.
- Conclude with the score using the format: "Educational score: <total points>"

Prompt for filtering The Stack v2 SQL data in [Stack-Edu](#)

Pre-training data: Other Sources

- Web text data (FineWeb2, DCLM, Map-CC..)
- Math data (FineMath, MegaMath, InfiMM-WebMath..)

Pre-training data: Other Sources

From [Qwen2.5 Coder](#): “Math and Text data may positively contribute to code performance, but only when their”

Token Ratio			Coding		Math		MMLU	General		<i>Average</i>
Code	Text	Math	Common	BCB	MATH	GSM8K		CEval	HellaSwag	
100	0	0	49.8	40.3	10.3	23.8	42.8	35.9	58.3	31.3
85	15	5	43.3	36.2	26.1	52.5	56.8	57.1	70.0	48.9
70	20	10	48.3	38.3	33.2	64.5	62.9	64.0	73.5	55.0

Table 3: The performance of Qwen2.5-Coder training on different data mixture policy.