

3. Post-training Code LMs: Supervised Fine-Tuning

Post-training: Recipe for SFT

➤ Instruction Data

- ❑ Objective: Steer LLMs toward targeted behaviors (instruction following, reasoning).
- ❑ Approach: Leverage instruction synthesis, structured response generation, and systematic quality evaluation.

➤ Training Curriculum

- ❑ Data Selection: Data mixing and scaling.
- ❑ Training Strategy: Iterative and multi-stage SFT, heavy/lightweight SFT.

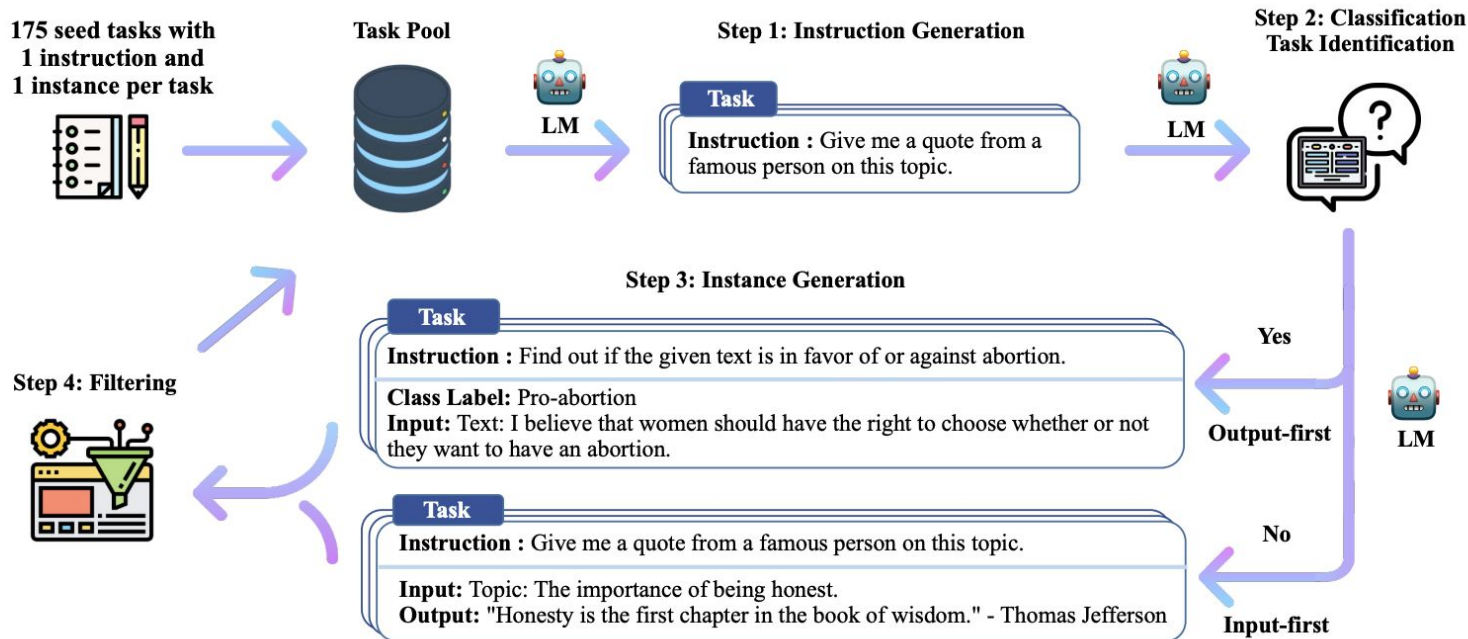
Recipe for SFT: Instruction Data

- Train code specific LLM (before)
 - Branching main pre-training run and continuing pre-training, followed by iterative SFT (Llama 3 approach)
 - Leverage to generate synthetic data, quality filtering

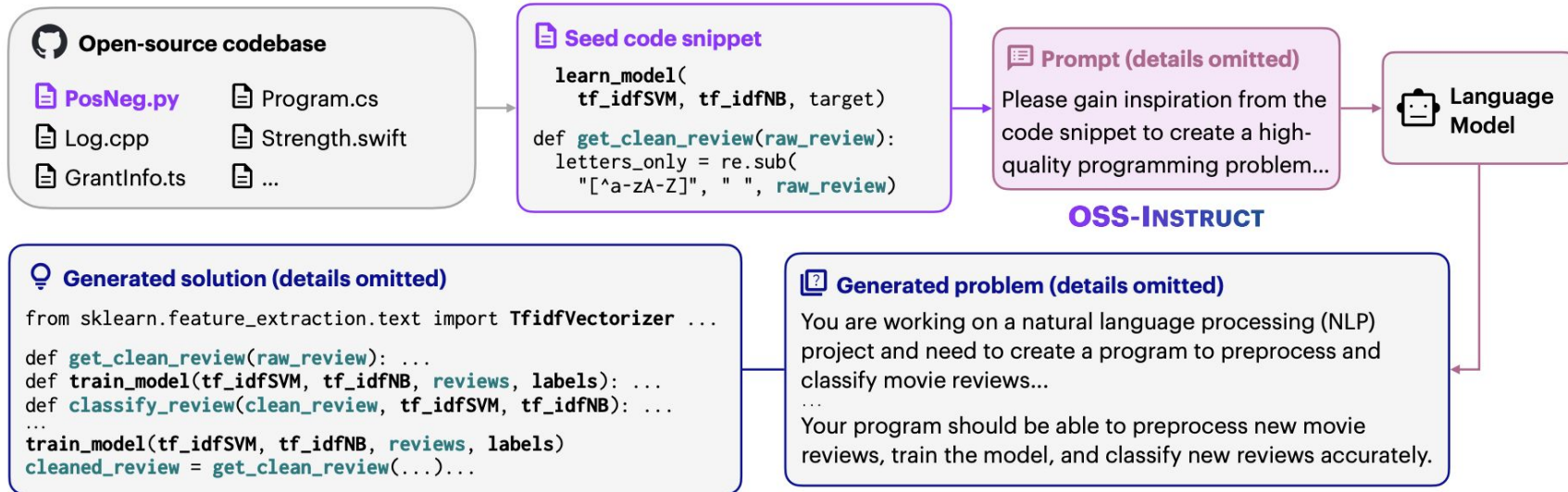
Recipe for SFT: Instruction Data

- Train code specific LLM (before)
 - Branching main pre-training run and continuing pre-training, followed by iterative SFT (Llama 3 approach)
 - Leverage to generate synthetic data, quality filtering
- Synthetic data generation
 - Techniques: Self-Instruct, Evol-Instruct, OSS-Instruct, and more.

Self-Instruct



OSS-Instruct



Llama 3 adopted OSS-Instruct to generate synthetic code instruction data

Recipe for SFT: Instruction Data

➤ Train code specific LLM (before)

- Branching main pre-training run and continuing pre-training, followed by iterative SFT (Llama 3 approach)
- Leverage to generate synthetic data, quality filtering

➤ Synthetic data generation

- Techniques: Self-Instruct, Evol-Instruct, OSS-Instruct, and more.
- Reasoning-based data for complex tasks

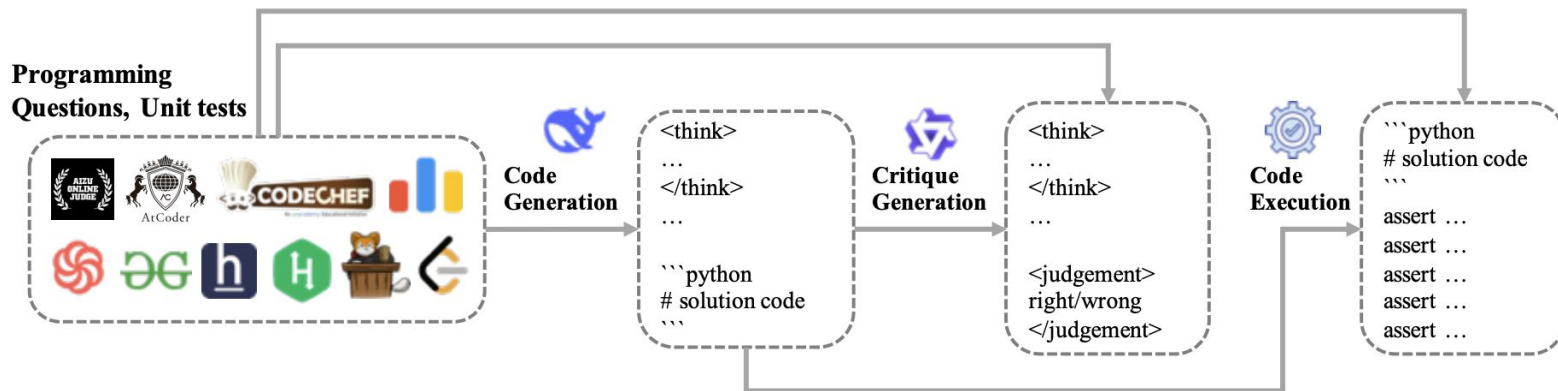
Reasoning vs. Non-reasoning Mode

Also known as “Thinking” vs “Non-Thinking” Mode

Thinking Mode	Non-Thinking Mode
<pre>< im_start >user {query} /think< im_end > < im_start >assistant <think> {thinking-content} </think> {response}< im_end ></pre>	<pre>< im_start >user {query} /no_think< im_end > < im_start >assistant <think> </think> {response}< im_end ></pre>

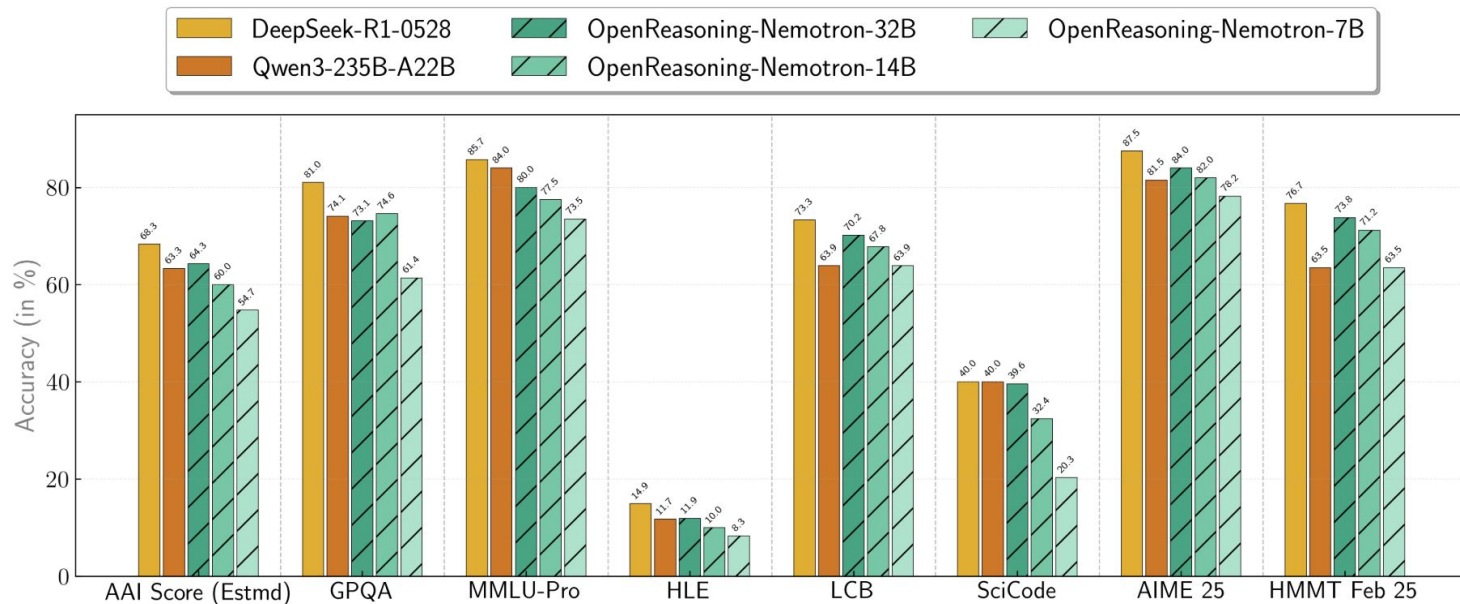
Reasoning-based Data for Competitive Coding

Strong-to-Weak distillation (off-policy)



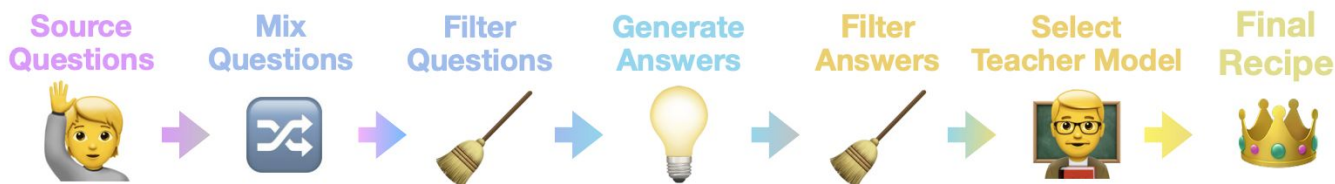
Overview of the OpenCodeReasoning development pipeline.

Distillation from DeepSeek-R1-0528



Distill(Qwen2.5-* -Instruct) => OpenReasoning-Nemotron-*

Reasoning-based Data for Competitive Coding



 Open Thoughts

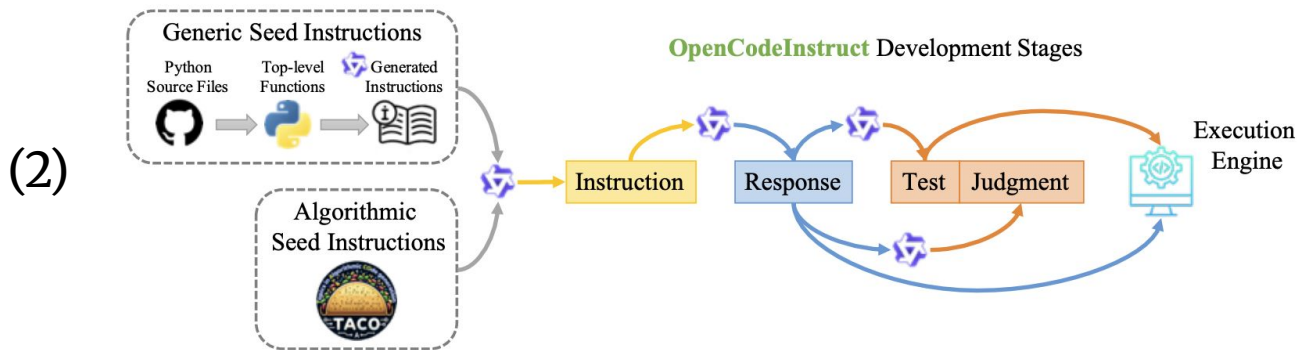
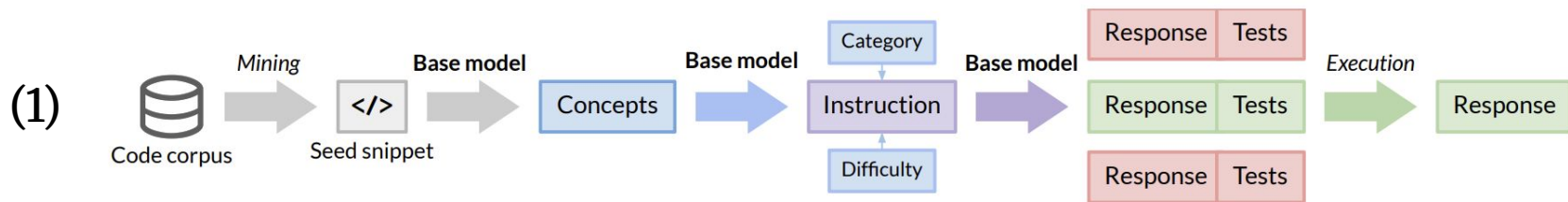
DATA RECIPES FOR REASONING MODELS

Recipe for SFT: Instruction Data

- Train code specific LLM (before)
 - Branching main pre-training run and continuing pre-training, followed by iterative SFT (Llama 3 approach)
 - Leverage to generate synthetic data, quality filtering
- **Synthetic data generation**
 - Techniques: Self-Instruct, Evol-Instruct, OSS-Instruct, and more.
 - Reasoning-based data for complex tasks
 - **Systematic quality assessment**

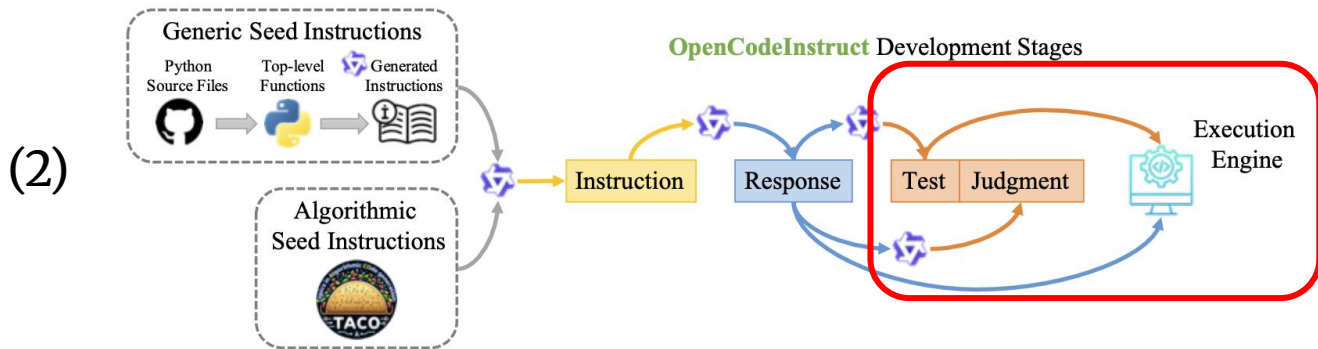
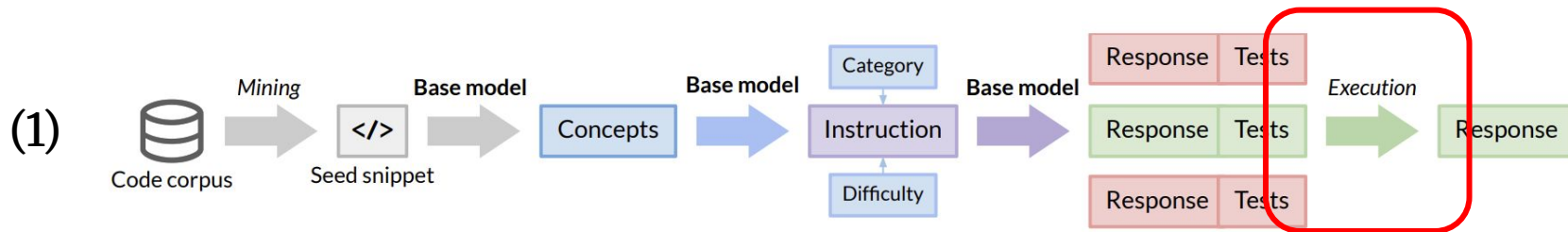
Quality Assessment

Based on execution feedback or LLM judgements



Quality Assessment

Based on execution feedback or LLM judgements



Quality Assessment

Based on execution feedback or LLM judgements

Selection strategy	Data size	Execution pass rate	Pass@1
Random selection (all)	27.7k	24.1%	61.6
Random selection (subset)	15.6k	24.2%	61.6
Failures only	15.6k	0%	57.9
Passes only	15.6k	100.0%	65.2

Table: Pass@1 on HumanEval+ with different response selection strategies from SelfCodeAlign.

Post-training: Recipe for SFT

➤ Instruction Data

- ❑ Objective: Steer LLMs toward targeted behaviors (instruction following, reasoning).
- ❑ Approach: Leverage instruction synthesis, structured response generation, and systematic quality evaluation.

➤ Training Curriculum

- ❑ Data Selection: Data mixing and scaling.
- ❑ Training Strategy: Iterative and multi-stage SFT, heavy/lightweight SFT.

Training Curriculum: Data Mixing for SFT

- Mixing code data with other sources

Dataset	% of examples	Avg. # turns	Avg. # tokens	Avg. # tokens in context	Avg. # tokens in final response
General English	52.66%	6.3	974.0	656.7	317.1
Code	14.89%	2.7	753.3	378.8	374.5
Multilingual	3.01%	2.7	520.5	230.8	289.7
Exam-like	8.14%	2.3	297.8	124.4	173.4
Reasoning and tools	21.19%	3.1	661.6	359.8	301.9
Long context	0.11%	6.7	38,135.6	37,395.2	740.5
Total	100%	4.7	846.1	535.7	310.4

Table: Statistics of SFT data used for Llama 3 post-training.

Training Curriculum: Data Mixing for SFT

➤ Coarse-to-fine SFT

- Coarse: Large-scale lower-quality/diverse data
- Fine: Small-scale high-quality data

Stage	Data Source	# Examples
Stage1	RealUser-Instruct	0.7 M
	Large-scale Instruct	2.3 M
	Infinity-Instruct	1.0 M
Stage2	McEval-Instruct	36 K
	Evol-Instruct	111 K
	Verified-Instruct	110 K
	Package-Instruct	110 K

Table: Detailed SFT data statistics for **OpenCoder** post-training.

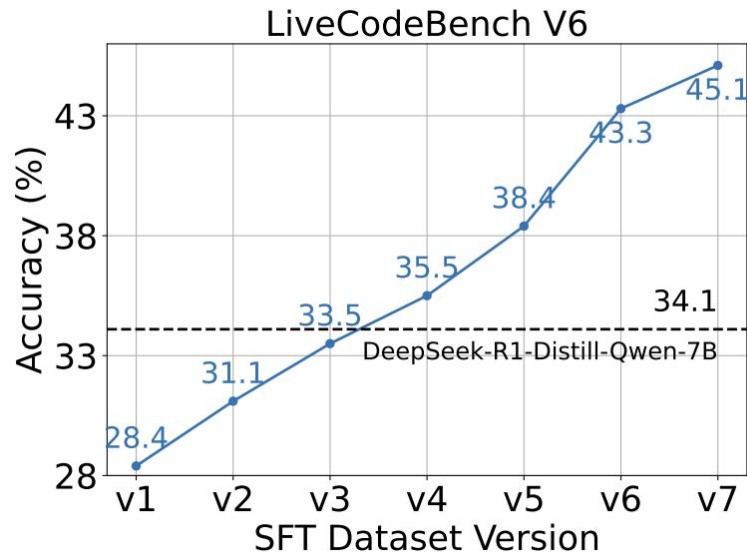
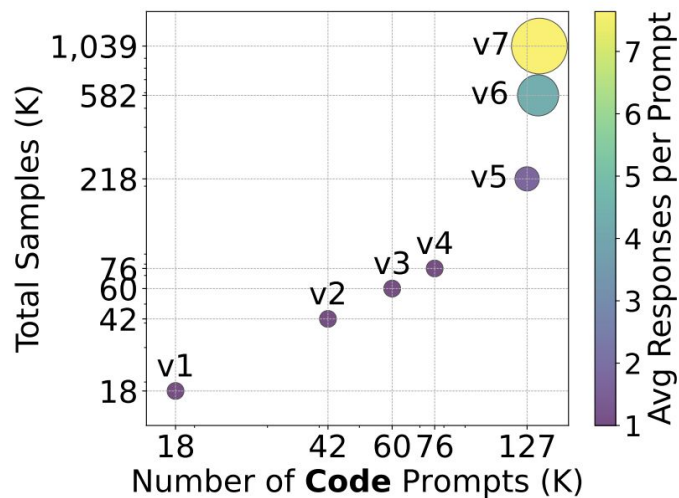
Training Curriculum: Data Mixing for SFT

➤ Reasoning vs. Non-reasoning data

- ❑ In general, **reasoning-based samples > non-reasoning samples**
- ❑ DeepSeek-R1 => 600k reasoning, 200k non-reasoning SFT samples
- ❑ Llama-Nemotron => 900k reasoning, 9M non-reasoning SFT code samples
- ❑ Nemotron-H => 5:1 and 1: 1 ratio of reasoning to non-reasoning samples for stage 1 and 2

Training Curriculum: Data Scaling

Supervised Fine-tuning of Qwen2.5-7B



Training Curriculum: SFT Strategy

- Iterative SFT (before)
 - Model improves and generates better synthetic data for subsequent iterations
- Multi-stage SFT (now)
 - Coarse-to-fine SFT: initial stages leverage large-scale, lower-quality data, while later stages refine the model using smaller but higher-quality datasets
- Lightweight vs. heavy SFT
 - SFT can over-constrain the model, restricting exploration during the online RL stage (lightweight SFT is suggested in Llama 4 post-training recipe)

Training Curriculum: SFT Strategy

