# Nemotron3 Nano

## Open, Efficient Mixture-of-Experts Hybrid Mamba-Transformer Model

### Focus on Data Mixture & Categories

# Executive Summary & Key Achievements

## Nemotron 3 Nano 30B-A3B Core Achievements

- Pretrained on 25 Trillion Tokens (3T New Unique Tokens)
- Better Accuracy, <1/2 Activated Parameters
- 3.3x Higher Inference Throughput
- Supports 1M Context Length
- Enhanced Agentic, Reasoning, & Chat Abilities

## Performance Comparison

Legend: Accuracy | Throughput | Context Length

Nemotron 3 Nano 30B-A3B: Accuracy 67%, Throughput 29.1, Context Length 93%

GPT-OSS-20B: Accuracy 26%, Throughput 30%, Context Length 17%

Qwen3-30B-A3B-Thinking-2507: Accuracy 25%, Throughput 17%, Context Length 35%

Y-axis: 0, 10, 20, 30, 40, 50

# Model Architecture Overview

MoE
MoE
MoE
Mamba-2
Mamba-2
Attention
Attention
Attention
MoE

52 Total Layers

**31.6B**
Total Parameters

**3.2B**
Activated (3.6B with embeddings)

**Granular MoE**

6 Activated

Shared Experts

Shared Experts

128 Experts, 6 Activated + 2 Shared

# Architecture Specifications & Hyperparameters

**Model Dimension: 2688**
- 52 Layers
- 32 Q-heads
- 2 KV-heads

## Mamba
- 128 Mamba State Dimension
- 8 Mamba Groups
- 64 Mamba Heads

## Granular MoE
- 128 Total Routable Experts
- 6 Activated Experts
- 2 Shared Experts

# Pretraining Strategy Overview

## Two-Phase Pretraining Approach

**Phase 1 (Diverse)**

23.5T Tokens
Diverse Data, Promoting Diversity

**Phase 2 (High-Quality)**

1.5T Tokens
High-Quality, Targeted Data

## Learning Rate Schedule

Warmup 8.4B

Stable at 10⁻³

Decay

Tokens

$10^{-3}$

$10^{-3}$

8.4B

Learning Rate

5T

$10^{-5}$

## Data Mixtures

**Phase 1 Mixture**

- Web Crawl 55%
- Code 10%
- Code
- Code 10%
- Code
- Code 10%
- Academic 6%
- Other

**Phase 1 Mixture**

- Web Crawl 55%
- Code 10%
- Code 10%
- Web 7%
- Academic 8%
- Academic 6%
- Academic 6%
- Other 1%

**Phase 2 Mixture**

- Wikipedia 30%
- Wikipedia 30%
- Wikipedia 7%
- Journals 20%
- Journals 20%
- Books 20%
- Books 25%
- Books 25%

# Pretraining Data: Nemotron-Pretraining-Code-v2

## GitHub Sourced Data

Cut-off: Apr 15, 2025.

Multi-stage filtering & deduplication.

High-quality code refinement.

## Synthetic Generation

Qwen3 32B generated.

Q&A pairs, Student-Teacher dialogue (Python), Code-review dialogue (Python/C++).

## Code Rewriting

Style-Guided (SGCR), Self-Contained Optimization (SCOR).

Python to C++ transpilation, improves downstream C++.

## Quality Assessment

Pylint-based analysis.

Automated quality scoring.

Ensures code adherence to best practices.

# Data Mixture & Categories

## Web Crawl & Quality Groups

- Crawl-Medium (15%) — **15%**
- Medium-High (12%)
- Syn-Medium-High (10%)
- Crawl-High (8%)
- Syn-Crawl-High (5%)

Phase 1 & Phase 2 Mixtures (Quality Prioritization)

## Structured & Specialized Data

- Academic Text (6%)
- Math (5%) — **5%**
- Wikipedia (4%)
- Code — **6%** (6%)
- Nemotron-CC-Code (4%)
- Crawl++ (OpenWebText, BigScience, Reddit) (5%) — **5%**

## Multilingual & Synthetic SFT

- Multilingual (19 languages) (6%)
- General-SFT (5%)
- STEM-SFT (5%)
- Code-SFT (4%)

Synthetic SFT-style Datasets

Comprehensive data mixture prioritizing quality and diversity across 15 categories, including optimized web crawl, structured data, and targeted synthetic datasets for enhanced performance.
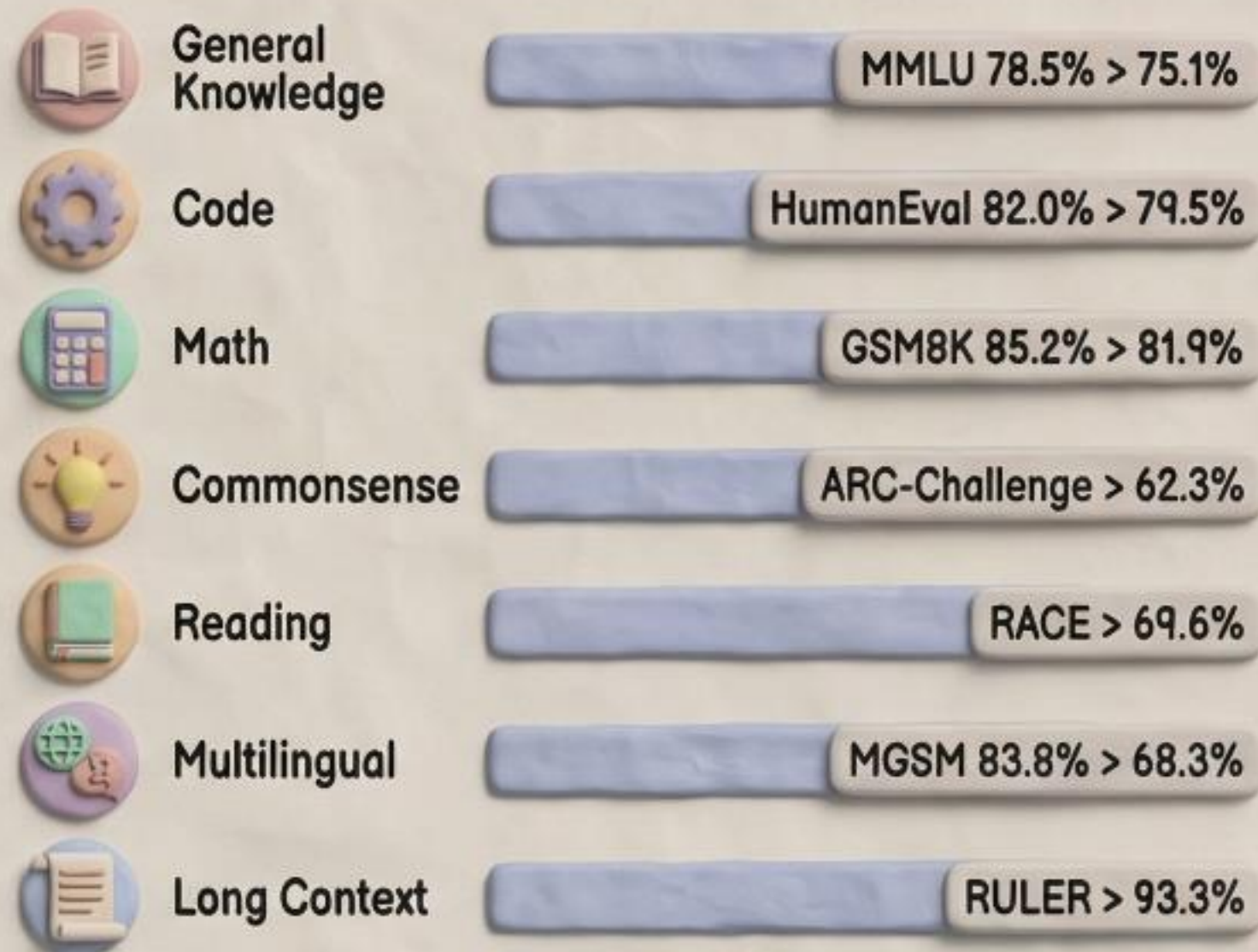
# Base Model Evaluation Results

## Comprehensive Benchmark Comparison: Nemotron 3 Nano vs. Qwen3

### Nemotron 3 Nano 30B-A3B Base

| Category | Benchmark |
|---|---|
| General Knowledge | MMLU 78.5% > 75.1% |
| Code | HumanEval 82.0% > 79.5% |
| Math | GSM8K 85.2% > 81.9% |
| Commonsense | ARC-Challenge > 72.0% |
| Reading | RACE > 74.7% |
| Multilingual | MGSM > 79.7% |
| Long Context | RULER > 77.9% |

### Qwen3-30B-A3B-Base

| Category | Benchmark |
|---|---|
| General Knowledge | MMLU 78.5% > 75.1% |
| Code | HumanEval 82.0% > 79.5% |
| Math | GSM8K 85.2% > 81.9% |
| Commonsense | ARC-Challenge > 62.3% |
| Reading | RACE > 69.6% |
| Multilingual | MGSM 83.8% > 68.3% |
| Long Context | RULER > 93.3% |

Highlight: Nemotron 3 Nano demonstrates superior performance in most categories, leading in General, Code, and Math benchmarks.

# SFT Data Categories & Composition

## SFT Data Blend Visualization

| Competition Math 57% | Competition Code 30% | Conversational Tool Use 17% | Long Context 18% | Formal Proofs 50% | General Chat 21% | Instruction Following 20% | Safety, SWE, Science, GenSelect, CUDA |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | | | Terminal Use 10% | Multilingual 10% | | | |

### Competition Math
GPT-OSS 120B,
Tool-integrated traces

### Competition Code
DeepSeek-R1
responses

### Conversational Tool Use
Synthetic multi-turn trajectories
with LM judge filtering

### Formal Proofs
580k theorems autoformalized
to 550k Lean 4, 920k proof
traces, 300k examples

### Multilingual
5 target languages

### Long Context
128k mean, 256k max
synthetic data

### Formal Proofs
580k theorem
autoformalize data

### Terminal Use
Terminal Bench
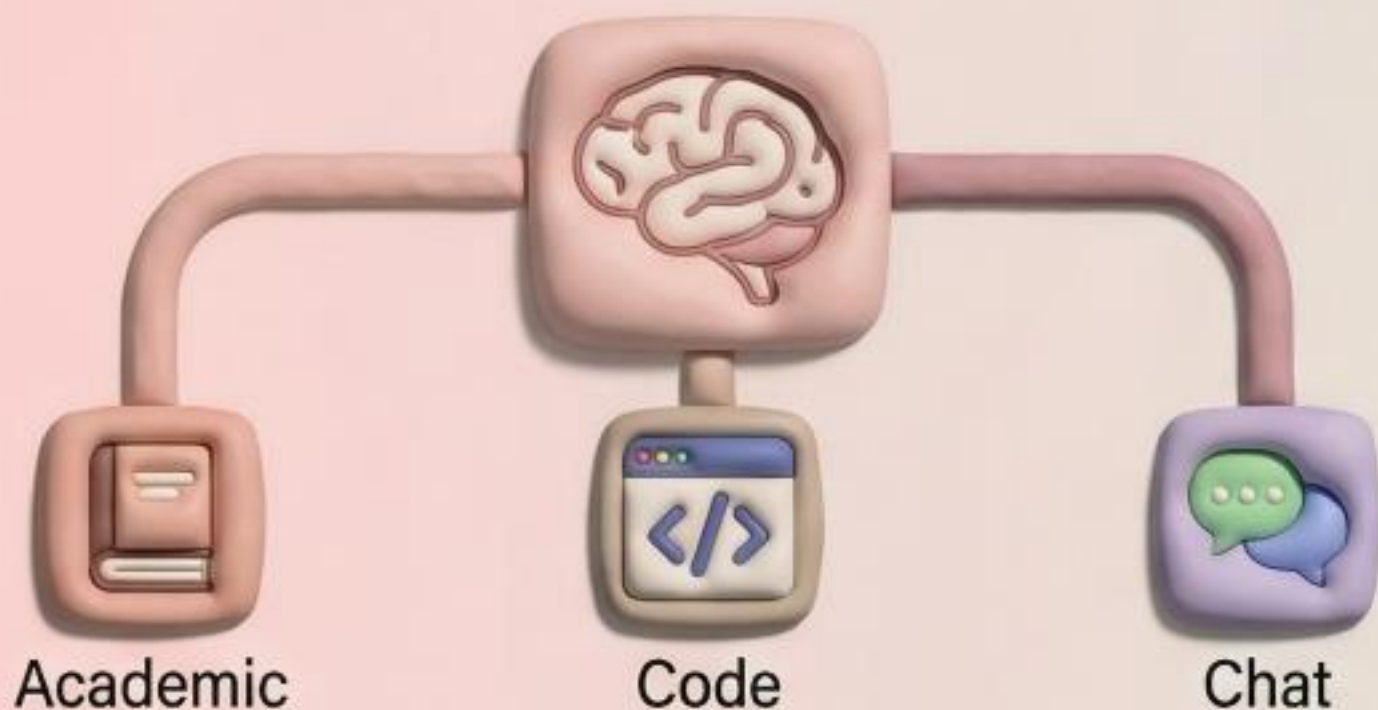verifiable tasks

### General Chat & IF
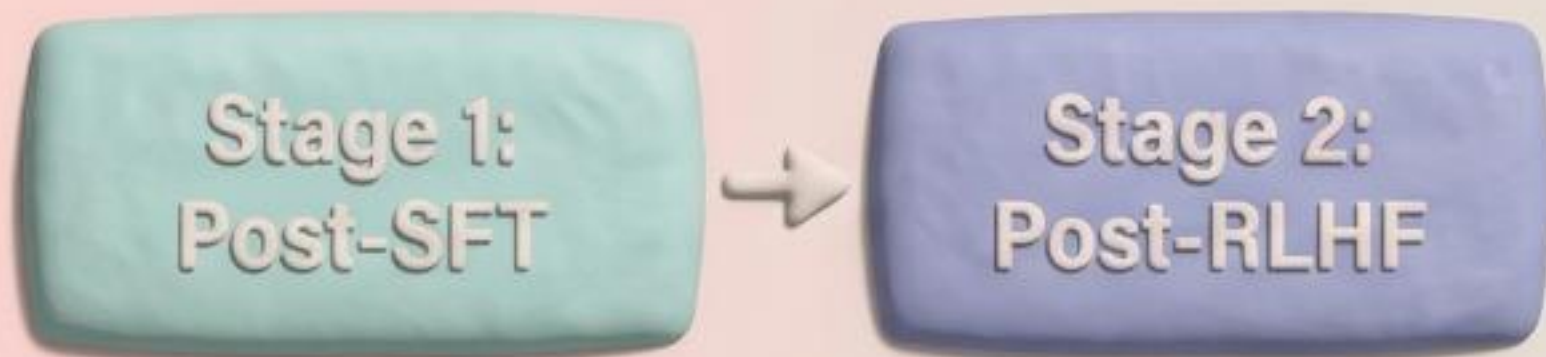LMSYS, WildChat,
IFeval, IFBench

### Safety, Science & CUDA
GitHub issues, Physics,
Chemistry, Biology,
21k PyTorch-CUDA pairs

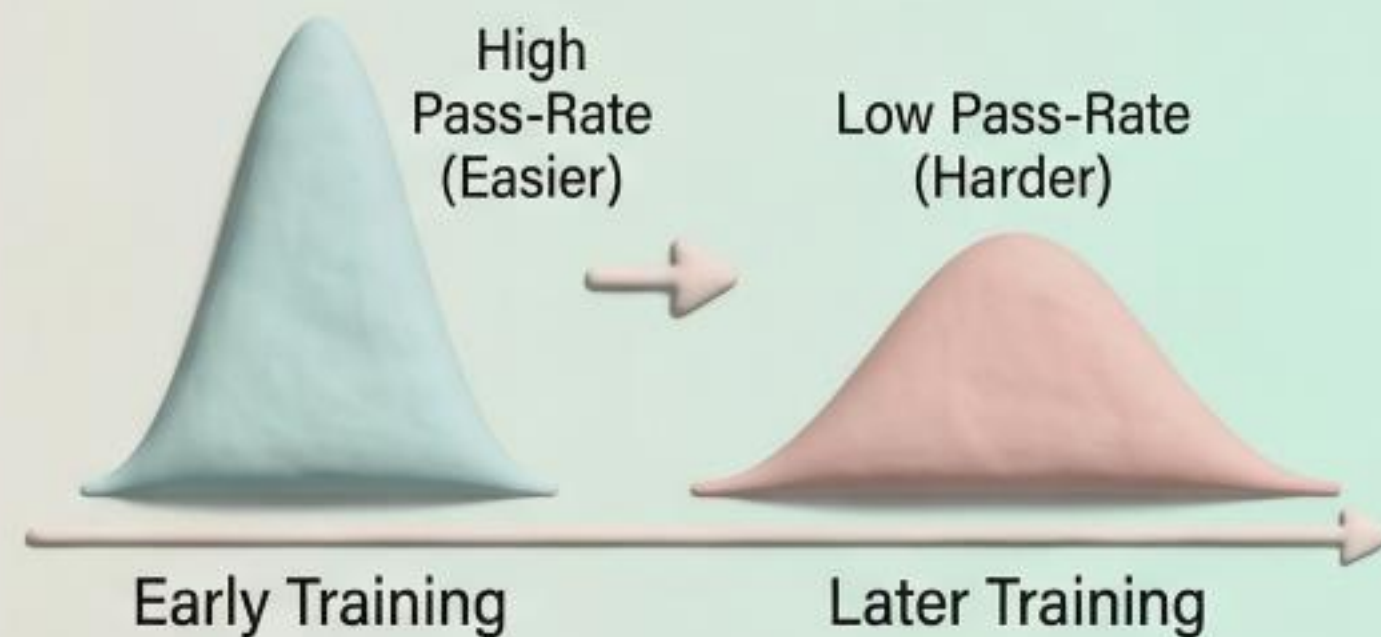# Multi-Environment RLVR: Training Strategy

## Unified RLVR Model

Academic

Code

Chat

## Two-Stage Training Pipeline

Stage 1: Post-SFT → Stage 2: Post-RLHF

Simultaneous training across all environments for stable gains.

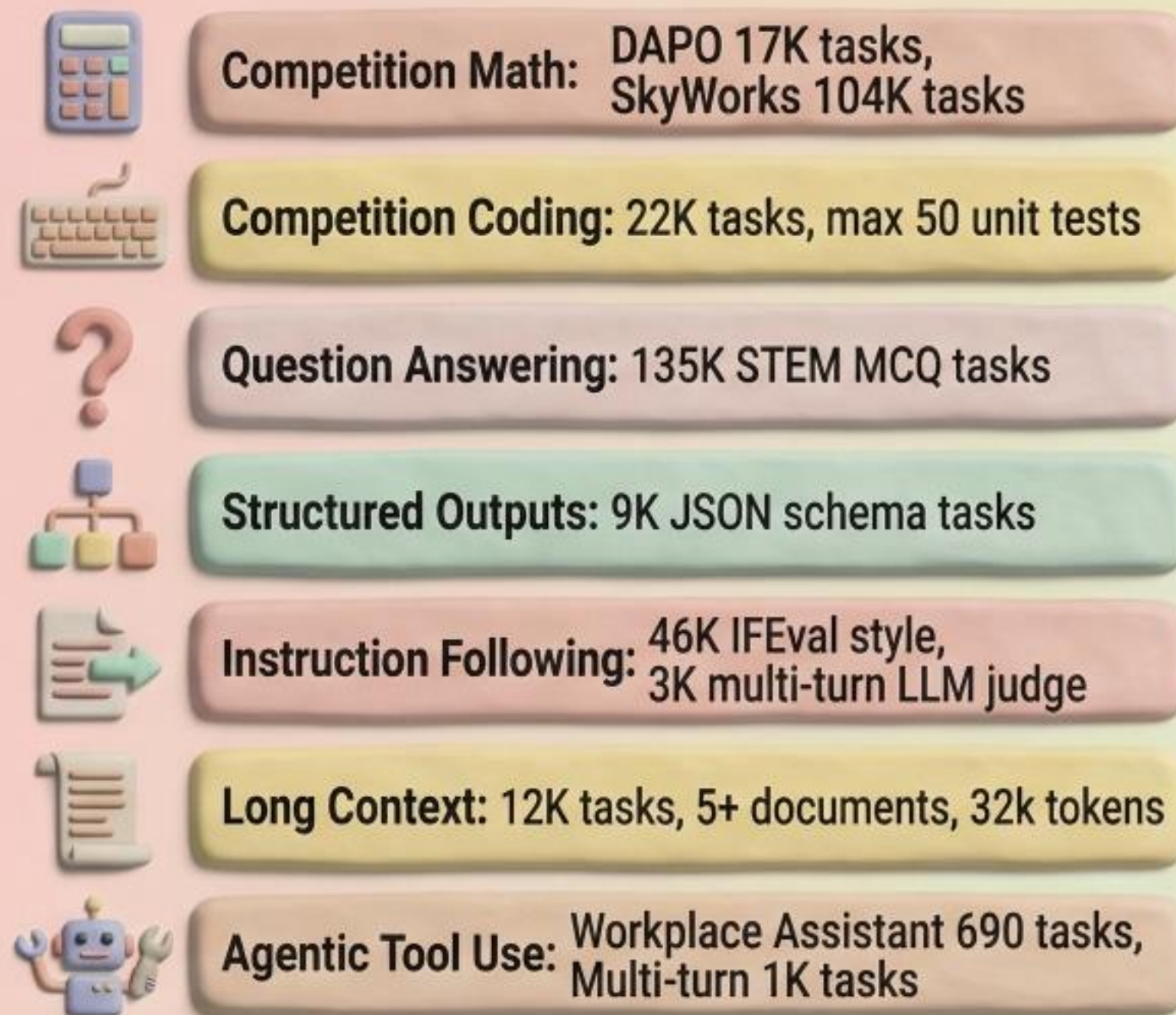## Curriculum Training with Gaussian Sampling

High Pass-Rate (Easier) → Low Pass-Rate (Harder)

Early Training

Later Training

## Batch-wise Pass Rate Evolution

Batch 1K        Batch 5K        Batch 10K

❌ Contrast: Single-environment training causes unrecoverable degradation.

# RLVR Algorithm & Performance

## Algorithm & Training Pipeline

- Synchronous GRPO with masked importance sampling
- 128 Prompts/Step, 16 Gens/Prompt, Batch Size 2048 — 16
- On-Policy Updates, MoE Router Weights Frozen, Aux-loss-free Load Balancing
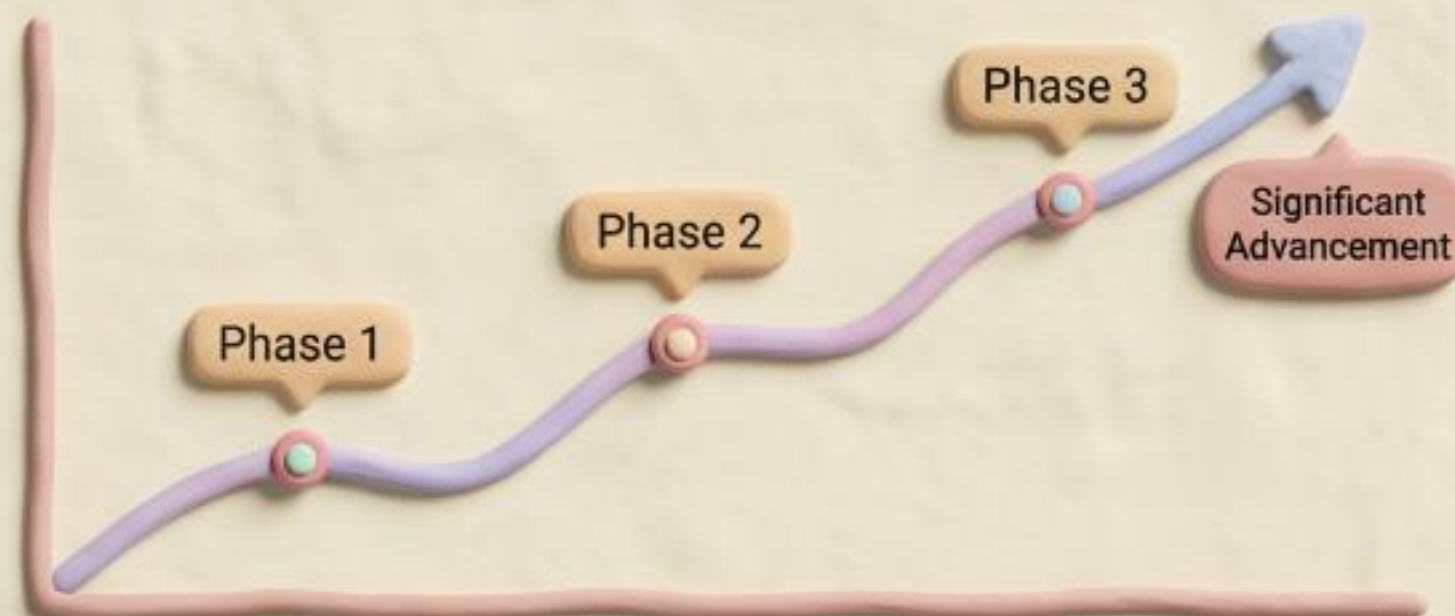- Maximum Generation Length 49K, Overlong Filtering Enabled

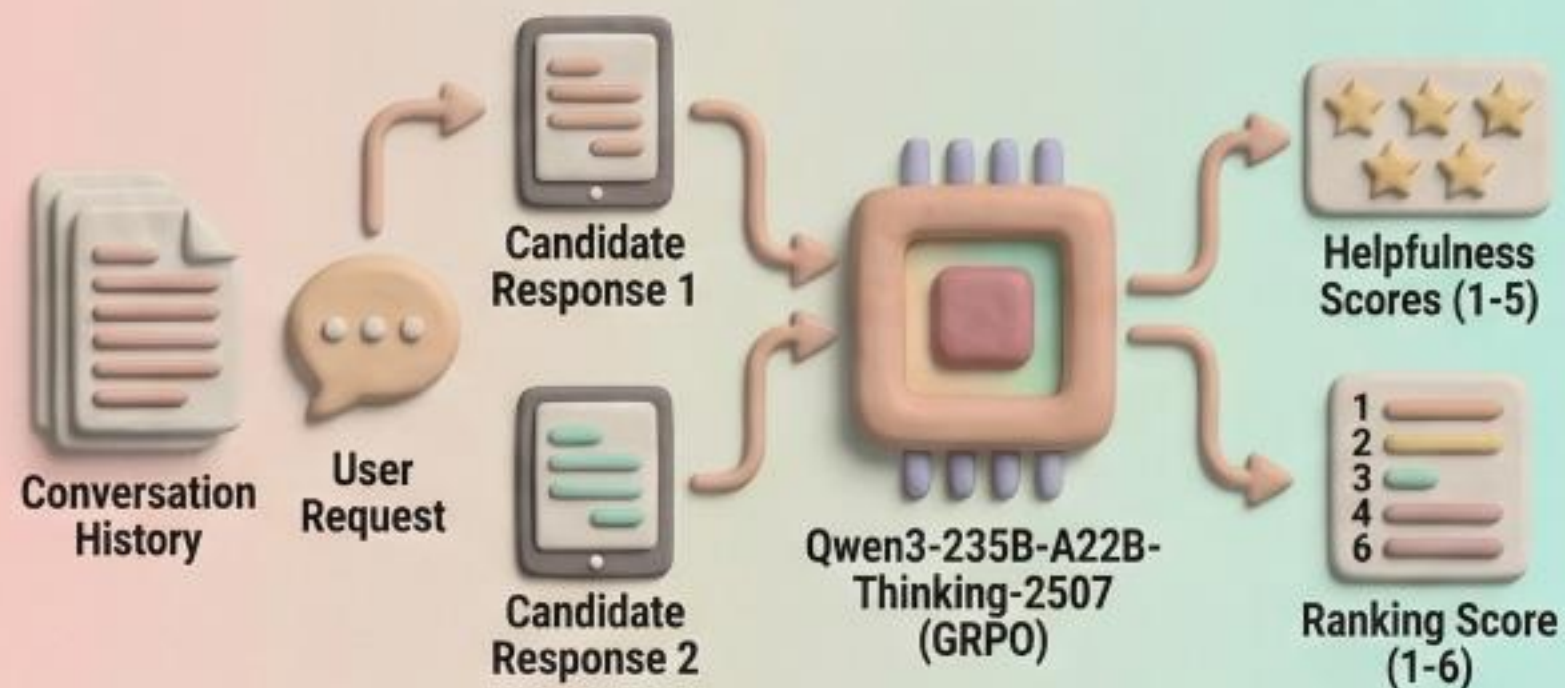## Performance Benchmarks

### RLVR vs. SFT Baseline

- GPQA: 20.6% (SFT)
- LiveCodeBench: 20.7% (SFT)
- AIME 2025: 31.3% (SFT)
- IFBench Prompt: 40.3% (SFT)

### RLVR Benchmark Performance Throughout Training

- Phase 1
- Phase 2
- Phase 3
- Significant Advancement

# Generative Reward Model Training

## GenRM Reasoning Flow & Scoring

Conversation History

User Request

Candidate Response 1

Candidate Response 2

Qwen3-235B-A22B-Thinking-2507 (GRPO)

Helpfulness Scores (1-5)

Ranking Score (1-6)

Reasons over inputs to generate helpfulness scores and a final ranking.

## Reward Function & Training Config

$$R = -C_1 \cdot I_{format} - |P_{h1} - G_{h1}| - |P_{h2} - G_{h2}| - C_2 \cdot |P_r - G_r|$$

$C_1 = 10$

$C_2 = 1$

128 prompts/batch

8 generations

One Gradient Step

Reward function guides training with specific batch and generation settings.
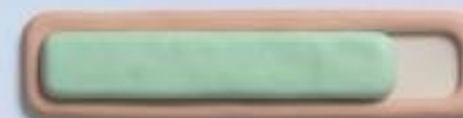
## Data Sources & Performance Gains

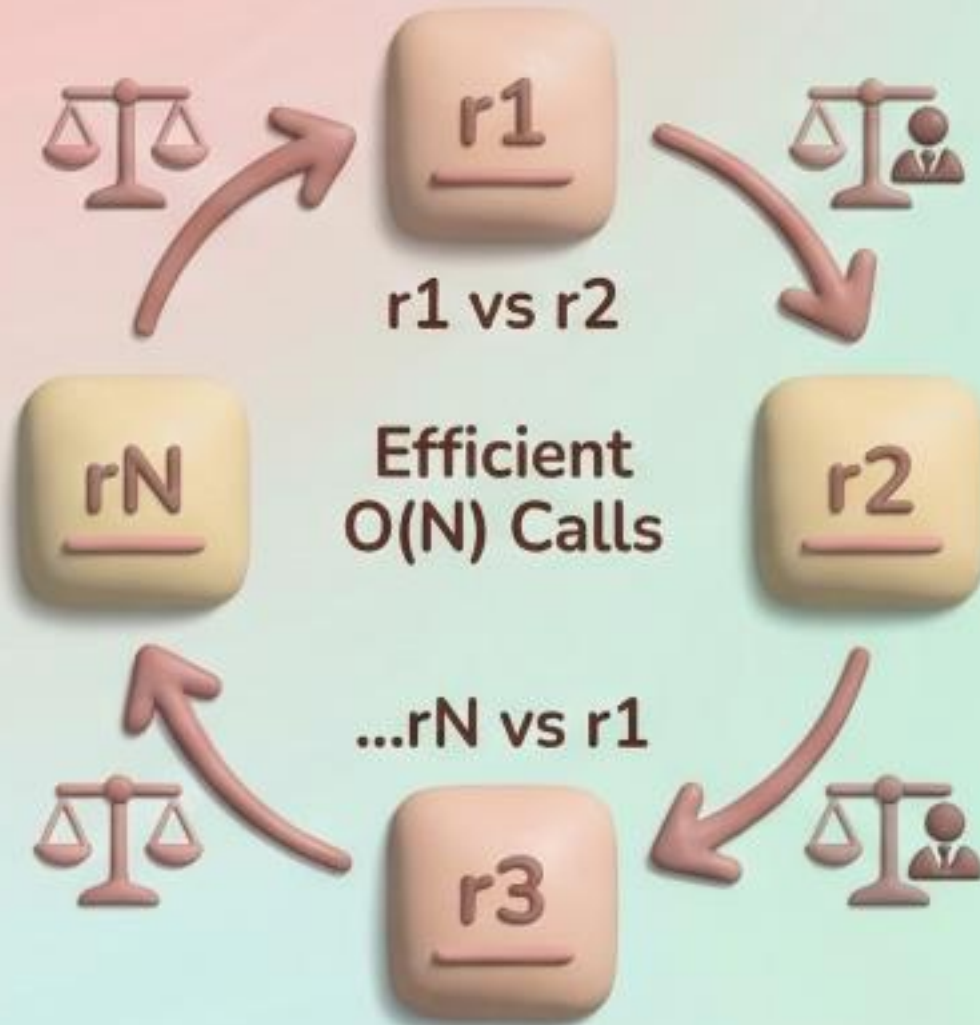HelpSteer3 Data + Synthetic Safety Blend

RM-Bench

JudgeBench

Internal-Val-Set

Data from HelpSteer3 and synthetic sources drive significant performance improvements across benchmarks.

# Post-Training Evaluation Results

**General Knowledge (MMLU-Pro)**

- Nemotron 3 Nano: 78.30
- Qwen3: 66.00
- GPT-OSS 20B: 57.8

**Reasoning**

AIME25
- 80.90
- 33.00
- 10.90

GPQA
- 75.00
- 34.00
- 10.90

LiveCodeBench
- 73.04
- 38.00
- 54.80

**Agentic**

TauBench V2 avg
- 48.00
- 22.00
- 5.00

SWE-Bench
- 38.76
- 10.00
- 12.05

**Chat & Instruction Following**

Arena-Hard-V2 Avg
- 67.65
- 48.55
- 51.00

IFBench
- 72.10
- 25.90
- 43.03

**Long Context**

RULER @ 1M
- 59.50
- N/A
- 34.00

**Multilingual**

MMLU-ProX
- 86.20
- 65.00
- 77.60

Legend: Nemotron 3 Nano, Qwen3, GPT-OSS 20B

# Quantization Results & Trade-offs

FP8  BF16

## 99%
**Median Accuracy Recovery vs BF16**

| Benchmark | % | FP8 | BF16 |
|---|---|---|---|
| MMLU-Pro | 78.30 | 78.30 vs 77.48 | |
| AIME25 no tools | 89.06 | 89.06 vs 87.71 | |
| AIME25 with tools | 99.17 | 99.17 vs 98.80 | |
| GPQA no tools | 73.04 | 73.04 vs 72.47 | |
| GPQA with tools | 75.00 | 75.00 vs 73.40 | |
| LiveCodeBench | 68.25 | 68.25 vs 67.62 | |
| TauBench average | 48.00 | 48.00 vs 44.79 | |
| IFBench | 35.85 | 35.85 vs 36.06 | |
| AA-LCR | 59.50 | 59.50 vs 59.63 | |
| MMLU-ProX | 78.10 | 78.10 vs 77.48 | |

## Throughput Improvement

**KV Cache FP8 Quantization** → **Larger Batch Sizes** → **Significant Throughput Gains**

## Ablation Study Visualization

**Visualization of Different Quantization Configurations & Results**

# Conclusion & Key Contributions

## 🧠 Nemotron 3 Nano Overview ⚙️

Open, Efficient MoE Hybrid Mamba-Transformer for agentic reasoning

Better/on-par accuracy with up to 3.3x higher inference throughput.

Supports 1M context length.

## Key Innovations: Architecture

Granular MoE architecture with shared experts.

Two-phase pretraining on 25T tokens (diverse then high-quality data).

## Key Innovations: Training & Optimization 📈

Multi-environment RLVR training simultaneously.

GenRM-based RLHF with length control.

Selective FP8 quantization preserving 99% accuracy.

## ☁️ Released Assets & Impact 📇

All model weights, training recipes, data, and code released on HuggingFace.

Significant advancement in efficient yet capable language models for agentic applications. 🤗