

## AutoCode: LLMs as Problem Setters for Competitive Programming

Shang Zhou <sup>1,\*</sup>, Zihan Zheng <sup>2,\*</sup>, Kaiyuan Liu <sup>3,\*</sup>, Zeyu Shen <sup>4,\*</sup>, Zerui Cheng <sup>4,\*</sup>, Zexing Chen <sup>1</sup>, Hansen He <sup>5</sup>, Jianzhu Yao <sup>4</sup>, Huanzhi Mao <sup>7</sup>, Qiuyang Mang <sup>7</sup>, Tianfu Fu <sup>6</sup>, Beichen Li <sup>8</sup>, Dongruixuan Li <sup>9</sup>, Wenhao Chai <sup>4,†</sup>, Zhuang Liu <sup>4,†</sup>, Aleksandra Korolova <sup>4,†</sup>, Peter Henderson <sup>4,†</sup>, Natasha Jaques <sup>3,†</sup>, Pramod Viswanath <sup>4,10,†</sup>, Saining Xie <sup>2,†</sup>, Jingbo Shang <sup>1,†</sup>

<sup>1</sup> UC San Diego   <sup>2</sup> New York University   <sup>3</sup> University of Washington   <sup>4</sup> Princeton University   <sup>5</sup> Canyon Crest Academy   <sup>6</sup> OpenAI   <sup>7</sup> UC Berkeley   <sup>8</sup> MIT   <sup>9</sup> University of Waterloo   <sup>10</sup> Sentient Labs

## Motivation

Problem setting requires a **deeper understanding of algorithms** than solving. It includes all challenges of solving and even more. Current benchmarks suffer from:

- **High FPR:** Incorrect solutions pass weak test suites.
- **High FNR:** Valid solutions crash on bad inputs.

This **distorts evaluation** and pollutes RL training data. AutoCode solves these problems to build better datasets.

# AutoCode Framework

A closed loop **Validator, Generator, and Checker** framework. It mirrors how human experts create contest problems.

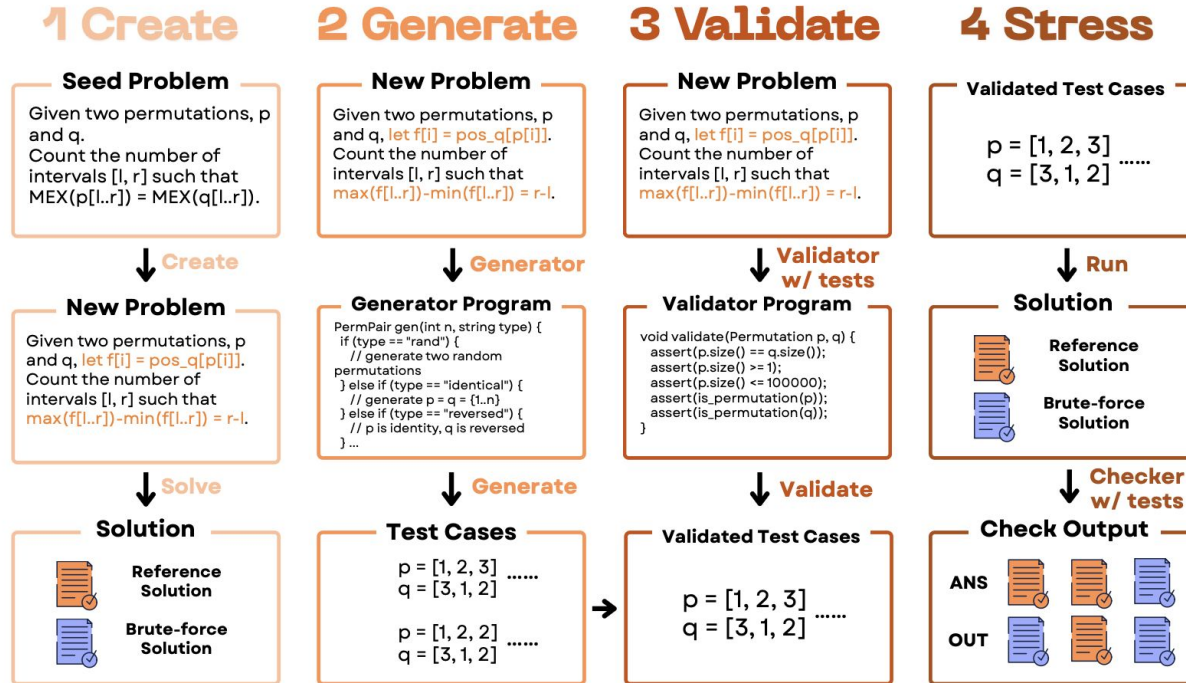


Figure 1. AutoCode pipeline: (1) Create from seed, (2) Generate test cases, (3) Validate inputs, (4) Stress test via dual verification.

# Test Case Generation Results

**Table 1. 7,538 problem benchmark.** 195,988 submissions, 50% correct / 50% incorrect. AutoCode uses o3.

Method	Consistency $\uparrow$	FPR $\downarrow$	FNR $\downarrow$
CodeContests	72.9%	7.7%	46.3%
CodeContests+	79.9%	8.6%	31.6%
TACO	80.7%	11.5%	26.9%
HardTests	81.0%	12.1%	25.8%
<b>AutoCode (Ours)</b>	<b>91.1%</b>	<b>3.7%</b>	<b>14.1%</b>

## 720 Problem Benchmark (Unfiltered Codeforces)

**98.7%**  
Consistency

**1.3%**  
FPR

**1.2%**  
FNR

# Ablation Study

---

Table 2. Ablation on 720 problems. GPT 5 High, 33 submissions per problem.

Configuration	Consist. $\uparrow$	FPR $\downarrow$	FNR $\downarrow$
w/o Exhaustive	98.4%	1.7%	1.3%
w/o Random/Extreme	98.4%	1.6%	1.3%
w/o TLE Inducing	98.6%	1.4%	1.3%
w/o Prompt Optimization	98.0%	1.8%	2.9%
<b>Full Framework</b>	<b>98.7%</b>	<b>1.3%</b>	<b>1.2%</b>

# Novel Problem Generation

AutoCode creates new variants from seed problems with reference solutions for verification.

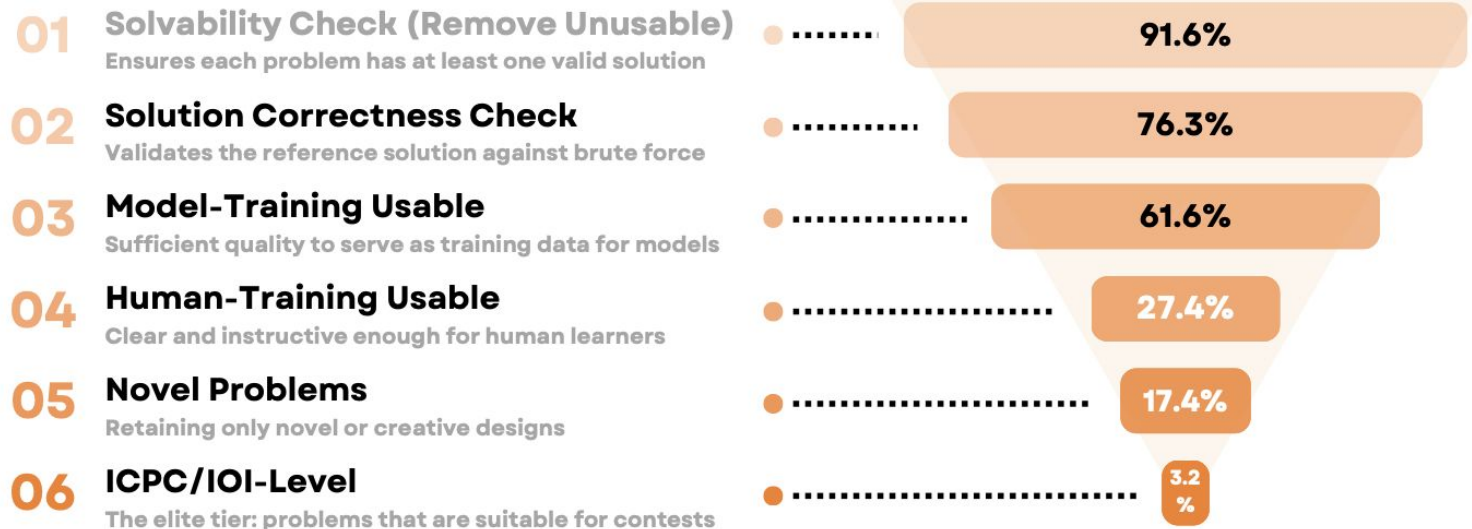


Figure 2. **Quality levels.** 61.6% suitable for training; 3.2% reach ICPC/IOI.

## Key Findings

- **Beyond limits:** LLMs create problems they cannot solve (4.2%).
- **Combining ideas:** LLMs mix existing frameworks over inventing new ones.
- **Elo gain:** Novel problems yield +498 Elo compared to +108.
- **Evaluation gap:** Human and LLM quality correlation is only  $r = 0.07$ .
- **Difficulty proxy:** Difficulty aligns well with human rating ( $r = 0.60$ ).

## Conclusion

- **98.7%** match with official judgments
- **50%** drop in both FPR and FNR
- **61.6%** of problems are good for training