# Self-play with Execution Feedback: Improving Instruction-following Capabilities of Large Language Models
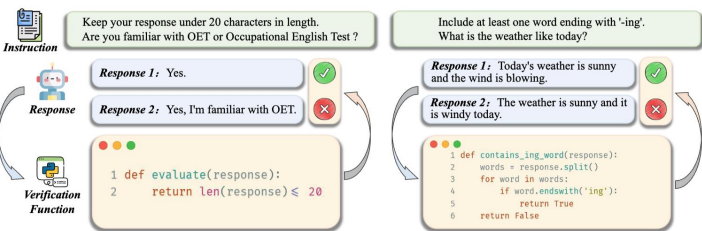
Guanting Dong, Keming Lu, Chengpeng Li, Tingyu Xia, Bowen Yu†, Chang Zhou, Jingren Zhou

Qwen Team, Alibaba Inc.

**ICLR 2025 (Spotlight, Top5% paper)**

## Motivation



A core strength of LLMs lies in executing natural language instructions. We presents AUTOIF , the first scalable framework for automated generation of high-quality instruction-following data. By reframing data validation as code verification, AUTOIF orchestrates three components: instruction generation, response-validation code synthesis, and unit test creation, forming a closed-loop quality assurance system. Execution feedback-driven rejection sampling efficiently produces data for SFT and RLHF. Evaluations on top open-source LLMs demonstrate substantial improvements across SFT, Offline/Online DPO training paradigms, particularly in self-alignment strong-to-weak distillation.
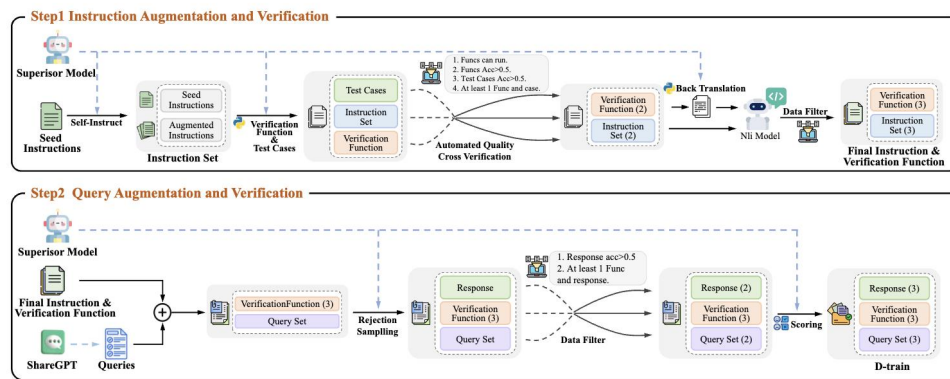
## Contribution & Conclusion

➤ We propose AUTOIF to efficiently enhance LLMs' instruction-following. It converts instruction-following alignment into auto code verification, making LLMs generate instructions, verification code, and unit test samples.

➤ Based on DPO algorithms, we treat executor feedback as a reward model, create pairwise preference samples, and design offline/on-policy strategies to optimize the model's instruction-following.
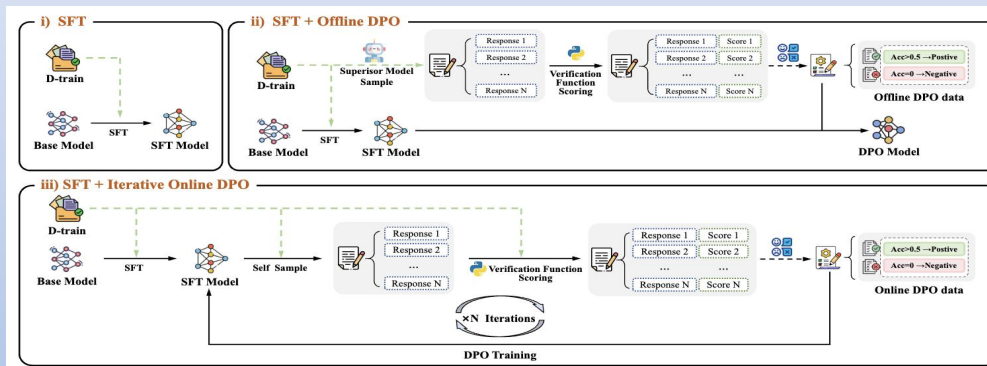
➤ AUTOIF is validated on benchmarks in "Self-Alignment" and "Strong-to-Weak" settings. It reaches over 90% accuracy in IFEval without sacrificing general and reasoning capability

## Method: AUTOIF



**1. Instruction Augmentation:** Starting with seed instructions, LLMs generate augmented instructions through self-instruct.
**2. Verification Functions:** Automatically generating Python functions to verify the correctness of responses.
**3. Back-Translation:** Ensuring consistency between instructions and verification functions.
**4. Query Augmentation:** Creating diverse queries and responses for training.
**5. Quality Filtering:** Filtering data based on verification function accuracy and query relevance.

## Training Strategies



## Main Results

| Model | IFEval | | | | FollowBench (SSR) | | | | | | C-Eval | MMLU | GSM8k | HumanEval |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Pr (S) | Pr. (L) | Ins. (S) | Ins. (L) | Level 1 | Level 2 | Level 3 | Level 4 | Level 5 | Avg | | | | |
| *Baselines (< 10B)* | | | | | | | | | | | | | | |
| Qwen2-7B | 37.7 | 43.6 | 49.4 | 53.4 | 55.6 | 53.5 | 53.7 | 49.9 | 48.6 | 52.3 | 74.4 | 64.4 | 71.1 | 58.1 |
| Qwen2-7B/ShareGPT | 30.9 | 33.5 | 42.4 | 45.2 | 56.1 | 52.7 | 50.8 | 45.2 | 47.9 | 50.5 | 70.2 | 59.8 | 59.4 | 52.4 |
| LLaMA3-8B | 24.6 | 26.1 | 38.1 | 39.7 | 10.0 | 10.3 | 10.5 | 14.3 | 12.7 | 11.6 | 24.2 | 38.8 | 4.5 | 0.6 |
| LLaMA3-8B(ShareGPT) | 23.7 | 26.4 | 33.8 | 37.1 | 44.0 | 40.0 | 39.6 | 33.3 | 33.6 | 38.1 | 35.2 | 44.6 | 20.5 | 38.1 |
| Mistral-7B | 23.3 | 24.6 | 38.4 | 39.6 | 40.1 | 39.7 | 37.9 | 35.7 | 36.7 | 38.0 | 38.2 | 47.6 | 20.5 | 38.4 |
| *Baselines (> 10B)* | | | | | | | | | | | | | | |
| Qwen2-72B-Instruct | 77.1 | 80.4 | 84.4 | 86.9 | 70.2 | 66.6 | 63.5 | 58.1 | 56.3 | 62.9 | 83.8 | 80.8 | 87.9 | 73.8 |
| LLaMA3-70B-Instruct | 77.8 | 83.8 | 84.2 | 88.8 | 60.7 | 60.5 | 61.1 | 61.7 | 60.3 | 60.9 | 60.2 | 80.5 | 92.6 | 78.7 |
| Mixtral-8x22B | 41.8 | 47.3 | 55.2 | 60.0 | 63.9 | 60.0 | 58.2 | 56.2 | 55.3 | 58.7 | - | - | - | - |
| GPT-4† | 76.9 | 79.3 | 83.6 | 85.4 | 84.7 | 77.6 | 76.2 | 77.9 | 73.3 | 77.9 | - | - | - | - |
| GPT-3.5 Turbo† | - | - | - | - | 80.3 | 71.2 | 74.2 | 69.6 | 67.1 | 72.5 | - | - | - | - |
| **Supervision Model: Qwen2-72B** | | | | | | | | | | | | | | |
| *Strong-to-Weak* | | | | | | | | | | | | | | |
| Qwen2-7B-SFT | 40.7₊₃.₀ | 44.5₊₀.₉ | 51.3₊₁.₉ | 55.4₊₂.₀ | 60.2₊₄.₆ | 53.7₊₀.₂ | 54.3₊₀.₆ | 49.9₊₀.₀ | 48.6₊₀.₀ | 53.3₊₁.₀ | 73.9₊₀.₆ | 64.4₊₀.₀ | 74.1₊₃.₀ | 58.3₊₀.₂ |
| w/ Offline DPO | 41.2₊₃.₅ | 44.7₊₁.₂ | 51.4₊₂.₀ | 56.2₊₂.₂ | 61.4₊₅.₈ | 54.5₊₁.₀ | 54.3₊₀.₆ | 51.2₊₁.₃ | 48.6₊₀.₀ | 54.0₊₁.₇ | 75.1₊₁.₀ | 64.5₊₀.₁ | 72.9₊₁.₈ | 59.5₊₁.₄ |
| *Self-Alignment* | | | | | | | | | | | | | | |
| Qwen2-72B-Instruct w/ Online DPO | 44.0₊₆.₃ | 46.6₊₃.₀ | 55.0₊₅.₆ | 57.9₊₄.₅ | 61.4₊₅.₈ | 56.8₊₃.₃ | 57.8₊₄.₁ | 55.4₊₅.₅ | 51.6₊₃.₀ | 56.6₊₄.₃ | 76.0₊₁.₆ | 64.8₊₀.₄ | 72.3₊₁.₂ | 58.2₊₀.₁ |
| | 80.2₊₃.₁ | 82.3₊₁.₉ | 86.1₊₁.₇ | 88.0₊₁.₁ | 76.2₊₆.₀ | 69.8₊₃.₂ | 67.0₊₃.₅ | 61.6₊₃.₅ | 62.8₊₆.₅ | 67.5₊₄.₆ | 84.9₊₁.₁ | 81.2₊₀.₄ | 88.2₊₀.₃ | 75.0₊₁.₂ |
| **Supervision Model: LLaMA3-70B** | | | | | | | | | | | | | | |
| *Strong-to-Weak* | | | | | | | | | | | | | | |
| LLaMA3-8B-SFT | 28.7₊₄.₁ | 40.3₊₁₄.₂ | 41.4₊₃.₃ | 52.2₊₁₂.₀₆ | 46.6₊₃₆.₆ | 46.2₊₃₅.₉ | 45.9₊₃₅.₄ | 37.6₊₂₃.₃ | 41.0₊₂₈.₃ | 43.5₊₃₁.₉ | 34.5₊₁₀.₃ | 45.6₊₆.₈ | 33.2₊₂₈.₇ | 38.2₊₃₇.₆ |
| w/ Offline DPO | 27.9₊₃.₃ | 41.5₊₁₅.₄ | 40.5₊₂.₄ | 54.1₊₁₄.₄ | 51.9₊₄₁.₉ | 51.3₊₄₁.₀ | 50.1₊₃₉.₆ | 45.3₊₃₁.₀ | 47.5₊₃₄.₈ | 49.2₊₃₇.₆ | 36.2₊₁₂.₀ | 45.5₊₆.₇ | 31.9₊₂₇.₄ | 38.5₊₃₇.₉ |
| w/ Online DPO | 28.8₊₄.₂ | 43.1₊₁₇.₀ | 42.2₊₄.₁ | 56.0₊₁₆.₃ | 54.6₊₄₄.₆ | 52.1₊₄₁.₈ | 50.0₊₃₉.₅ | 49.0₊₃₄.₇ | 43.7₊₃₁.₀ | 49.9₊₃₈.₃ | 38.2₊₁₄.₀ | 45.1₊₆.₃ | 32.5₊₂₈.₀ | 38.4₊₃₇.₈ |
| *Self-Alignment* | | | | | | | | | | | | | | |
| LLaMA3-70B w/ Online DPO | 80.2₊₂.₄ | 85.6₊₁.₈ | 86.7₊₂.₅ | 90.4₊₁.₆ | 71.0₊₁₀.₃ | 67.2₊₆.₇ | 66.2₊₅.₁ | 64.6₊₂.₉ | 63.5₊₃.₂ | 66.5₊₅.₆ | 61.6₊₁.₄ | 80.7₊₀.₂ | 92.7₊₀.₁ | 78.7₊₀.₀ |

**Results:**
1. AUTOIF achieves up to 90.4% accuracy on IFEval and over 5% improvement on FollowBench.
2. No decline in other capabilities

**Key Findings:**
3. Online DPO outperforms Offline DPO by effectively targeting model weaknesses
4. Larger models (e.g., Qwen2-72B) show greater improvements.
5. Higher function pass rates lead to better performance.

## Scaling Results