

Flow-GRPO: Training Flow Matching Models via Online RL

Jie Liu*, Gongye Liu*, Jiajun Liang, Yangguang Li,
Jiaheng Liu, Xintao Wang, Pengfei Wan, Di Zhang, Wanli Ouyang✉



- 1. Existing post-training algorithms for FM models**
- 2. Flow-GRPO: Training FM models via online RL**
- 3. Experiments**

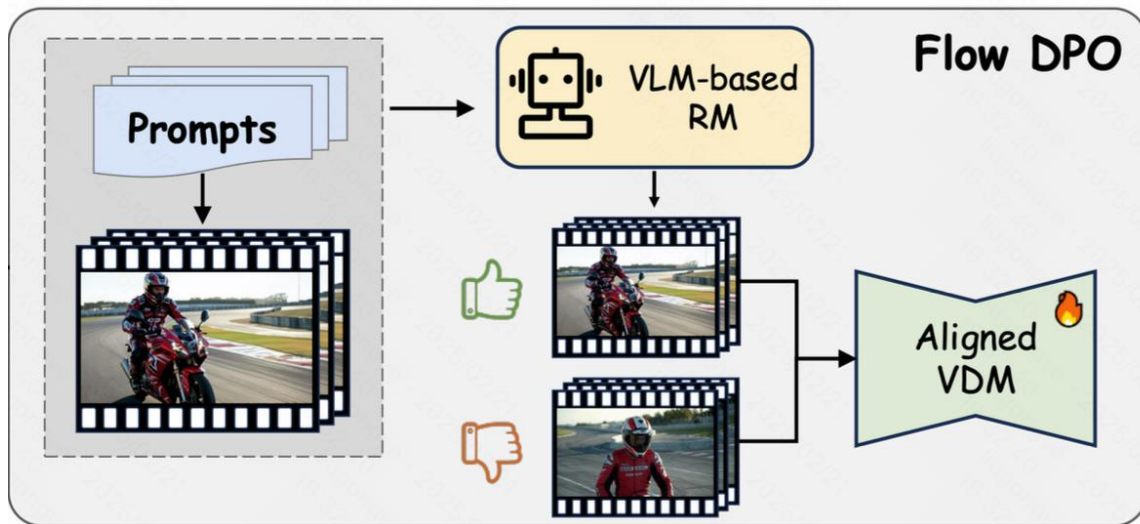
Post-Training Algorithms

- Evolution in LLM
 - 2024: DPO / Online-DPO dominated alignment methods
 - 2025: GRPO / PPO replaced DPO across LLMs
 - Offering better stability, on-policy updates, and higher reward efficiency.

**Can an online RL method like GRPO can be
applied to Flow Matching Models?**

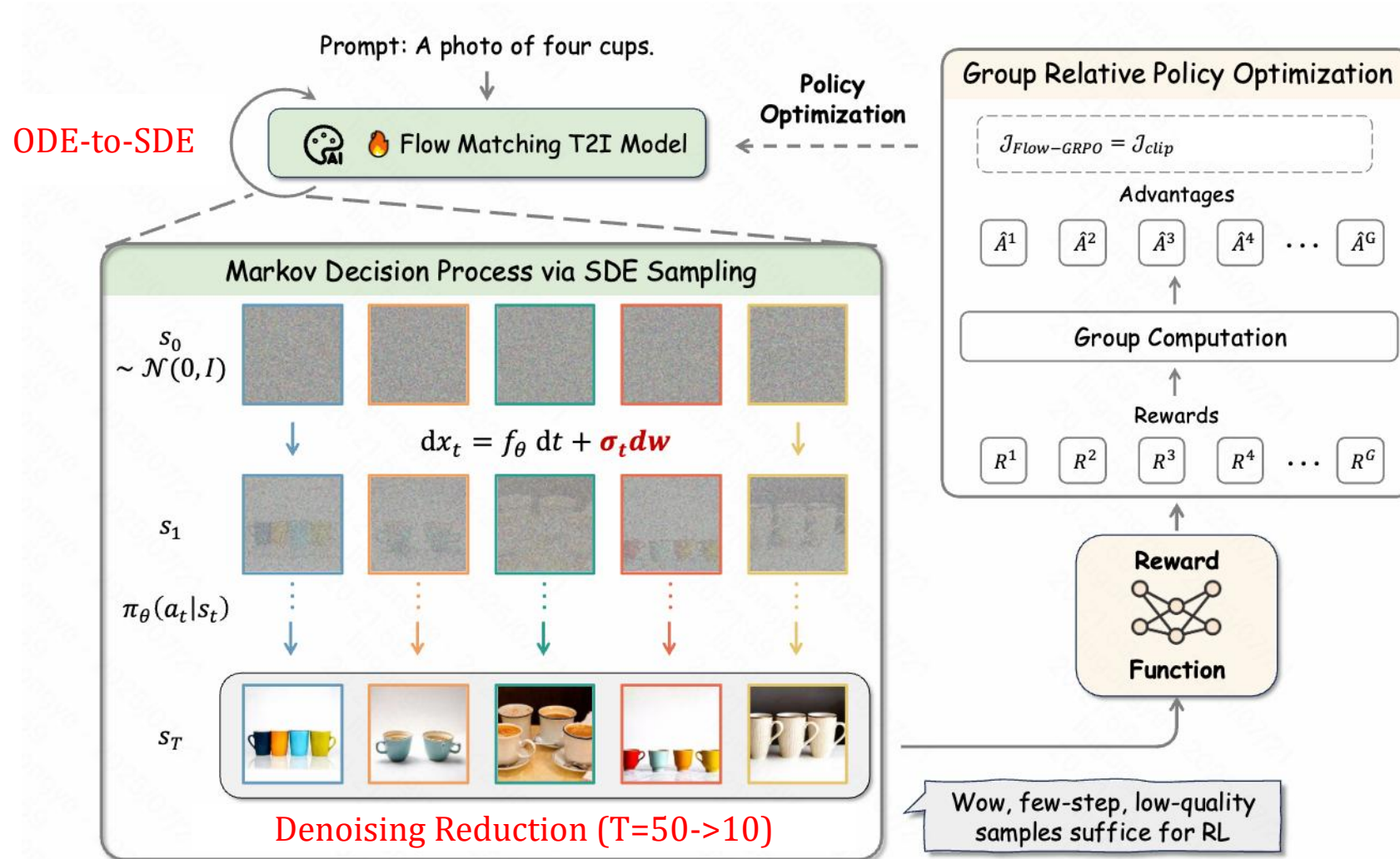
Existing Post-Training Algorithms: DPO

- Limitations of Flow-DPO
 - DPO: Learns directly from preference pairs → low data efficiency, limited performance gain.
 - GRPO: Trains a reward model to relabel online generated new samples → more stable training with steadily improving rewards.



$$-\mathbb{E} \left[\log \sigma \left(-\frac{\beta_t}{2} \left(\|\mathbf{v}^w - \mathbf{v}_\theta(\mathbf{x}_t^w, t)\|^2 - \|\mathbf{v}^w - \mathbf{v}_{\text{ref}}(\mathbf{x}_t^w, t)\|^2 \right. \right. \right. \\ \left. \left. \left. - (\|\mathbf{v}^l - \mathbf{v}_\theta(\mathbf{x}_t^l, t)\|^2 - \|\mathbf{v}^l - \mathbf{v}_{\text{ref}}(\mathbf{x}_t^l, t)\|^2) \right) \right) \right],$$

Post-Training Algorithms: GRPO



Challenges & Solutions

Challenge1: GRPO require a statistic sampling, but FM is trained for deterministic sampling.

Solution1: ODE-to-SDE

$$d\mathbf{x}_t = \mathbf{v}_t dt,$$

$$-\nabla \cdot [f_{\text{SDE}} p_t(\mathbf{x})] + \frac{1}{2} \nabla^2 [\sigma_t^2 p_t(\mathbf{x})] = -\nabla \cdot [\mathbf{v}_t(\mathbf{x}_t, t) p_t(\mathbf{x})]$$

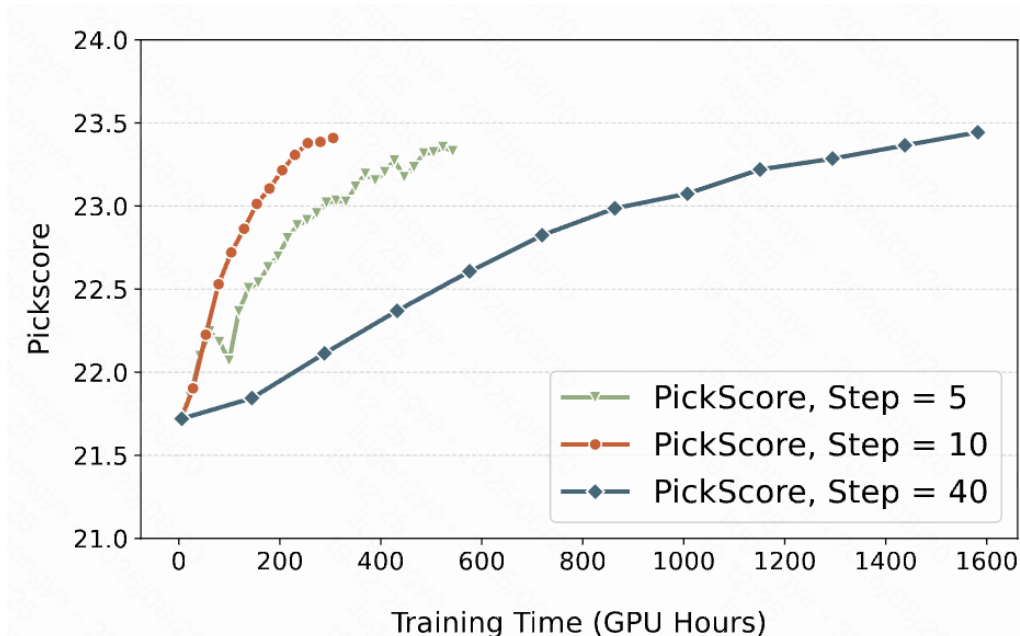
with equal marginal distribution

$$\mathbf{x}_{t+\Delta t} = \mathbf{x}_t + \left[\mathbf{v}_\theta(\mathbf{x}_t, t) + \frac{\sigma_t^2}{2t} (\mathbf{x}_t + (1-t)\mathbf{v}_\theta(\mathbf{x}_t, t)) \right] \Delta t + \sigma_t \sqrt{\Delta t} \epsilon$$

Challenges & Solutions

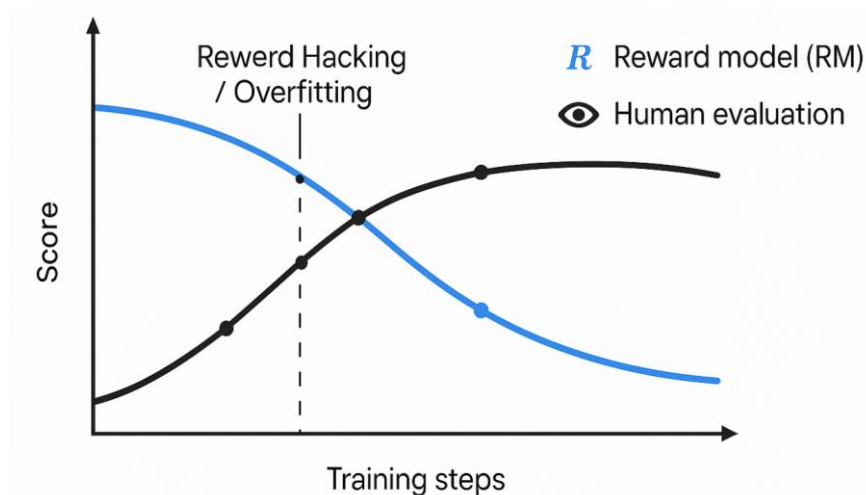
Challenge2: GRPO requires multiple samplings, and each sampling takes about 50 steps -> Too slow! 😞

Solution2: Denoising Reduction. Just 10-step training are sufficient to boost performance for 40-step inference -> Fast! 😊

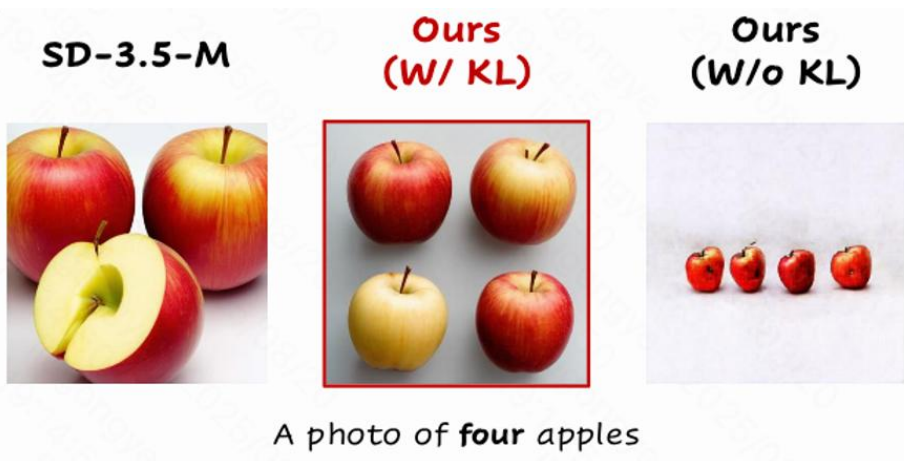


Challenges & Solutions






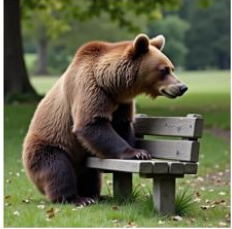


















Challenge3: Reward Hacking problem-> Overfitting to the imperfect rm



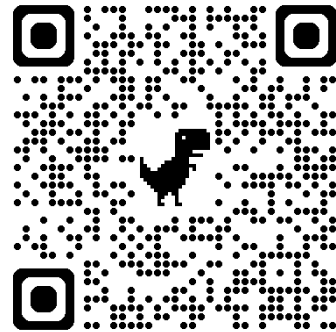
Solution3: KL divergence can effectively prevent reward hacking.



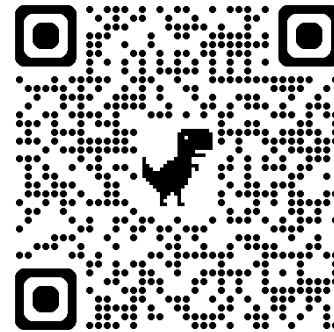
Experiments

	A photo of four giraffes	A photo of a white sandwich	A photo of a red dog	A photo of a brown giraffe and a white stop sign	A photo of a red orange and a purple broccoli	A photo of a bench left of a bear
FLUX.1 Dev						
GPT-4o						
SD-3.5-M						
SD-3.5-M + Flow-GRPO						

Thanks



paper



code