

Reward Score Matching

Unifying reward-based fine-tuning for flow and diffusion models

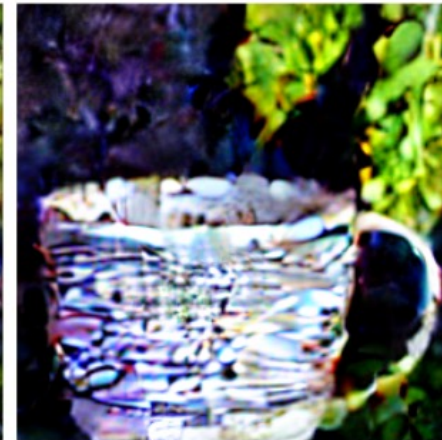
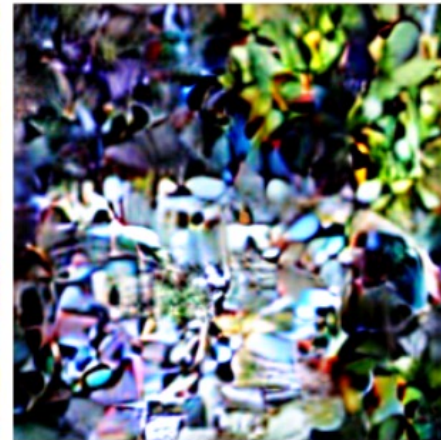
Jeongjae Lee^{*}, Jinho Chang^{*}, Jeongsol Kim[†], Jong Chul Ye[†]



ICML 2026 SPIGM Workshop

Flow-based models

$p_1(\mathbf{x}_1)$



$p_0(\mathbf{x}_0)$



$t = 1$

$t = 0$

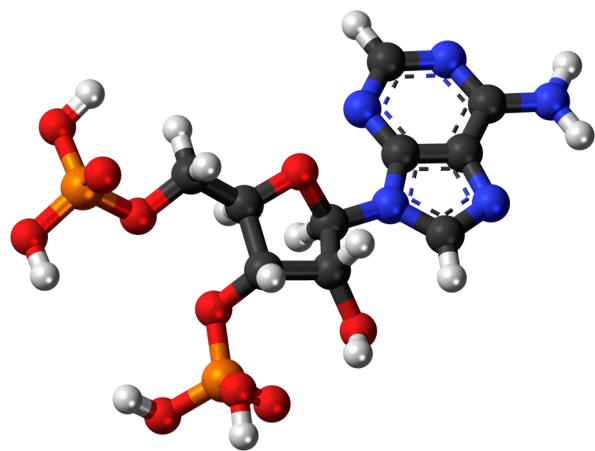
Reward model

Introduces an additional objective/constraint:

$$r : \mathbb{R}^d \rightarrow \mathbb{R}, \mathbf{x}_0 \mapsto r(\mathbf{x}_0)$$

QED

Drug likeliness
of molecule



Aesthetic Score

Aesthetic quality
of image



Reward alignment (fine-tuning)

Given: Pretrained generative model p^{ref}

Obtain: Finetuned generative model p^{\star}

$$p^{\star}(\mathbf{x}_0) = \arg \max_q \mathbb{E}_{\mathbf{x}_0 \sim q} [r(\mathbf{x}_0)]$$

Reward alignment (fine-tuning)

Given: Pretrained generative model p^{ref}

Obtain: Finetuned generative model p^{\star}

$$\begin{aligned} p^{\star}(\mathbf{x}_0) &= \arg \max_q \mathbb{E}_{\mathbf{x}_0 \sim q} [r(\mathbf{x}_0)] - \alpha \mathcal{D}_{\text{KL}}(q \| p^{\text{ref}}) \\ &= \frac{p^{\text{ref}}(\mathbf{x}_0)}{Z} \exp\left(\frac{r(\mathbf{x}_0)}{\alpha}\right) \end{aligned}$$

Reward-based fine-tuning → Reward-tilted distribution

Fragmented literature

Differentiable reward

Non-differentiable reward

Soft RL

SQDF

PPO

GRPO

GRPO
variants

Optimal
control

Adjoint
Matching

VGG-Flow

GFlowNet

Residual ∇ -DB

DAG

RWR

DiffusionNFT

AWM

One problem, Many methods.

RL fine-tuning

$$p^*(\mathbf{x}_t) = \frac{p^{\text{ref}}(\mathbf{x}_t)}{Z} \exp\left(\frac{V_t(\mathbf{x}_t)}{\alpha}\right)$$

RL fine-tuning

$$p^*(\mathbf{x}_t) = \frac{p^{\text{ref}}(\mathbf{x}_t)}{Z} \exp\left(\frac{V_t(\mathbf{x}_t)}{\alpha}\right)$$

$$s^*(\mathbf{x}_t) = s^{\text{ref}}(\mathbf{x}_t) + \frac{1}{\alpha} \nabla_{\mathbf{x}_t} V_t(\mathbf{x}_t)$$

RL fine-tuning as score matching

$$p^*(\mathbf{x}_t) = \frac{p^{\text{ref}}(\mathbf{x}_t)}{Z} \exp\left(\frac{V_t(\mathbf{x}_t)}{\alpha}\right)$$

$$s^*(\mathbf{x}_t) = s^{\text{ref}}(\mathbf{x}_t) + \frac{1}{\alpha} \nabla_{\mathbf{x}_t} V_t(\mathbf{x}_t)$$

$$\mathcal{L}_{\text{RSM}}(\theta) = \mathbb{E}_{t_i, \mathbf{x}_{t_i}, \epsilon} \left[C_1(t_i) \underbrace{\left(\|\mathbf{s}_{t_i}^\theta - (s_{t_i}^{\text{ref}} + \gamma(t_i) \hat{\Psi}_{t_i})\|^2 \right)}_{\text{Guidance}} + C_2(t_i) \underbrace{\|\mathbf{s}_{t_i}^\theta - \mathbf{s}_{t_i}^{\theta_{\text{old}}}\|^2}_{\text{Regularization}} \right]$$

Many methods = Score matching against value-tilted score

Flow-GRPO

```
ratio = torch.exp(log_prob - sample["log_probs"][:, j])
```

Flow-GRPO (Liu et. al., NeurIPS 2025)

$$p(\mathbf{x}) = (\text{const}) \exp\left(-\frac{\|\mathbf{x} - \boldsymbol{\mu}\|^2}{(\text{const})}\right)$$

$$\log p(\mathbf{x}) = (\text{const}) + (\text{const})\|\mathbf{x} - \boldsymbol{\mu}\|^2$$

Flow-GRPO

```
ratio = torch.exp(log_prob - sample["log_probs"][:, j])
```

Flow-GRPO (Liu et. al., NeurIPS 2025)

$1 + \log \widehat{\rho}_t(\theta)$ is a better estimator of $\rho_t(\theta)$ than $\widehat{\rho}_t(\theta)$

Avoid exponentiation \Rightarrow Much lower variance

PCPO (Lee and Ye, ICLR 2026)

Flow-GRPO is also L_2 regression

```
ratio = torch.exp(log_prob - sample["log_probs"][:, j])
```

Flow-GRPO (Liu et. al., NeurIPS 2025)

$1 + \log \widehat{\rho}_t(\theta)$ is a better estimator of $\rho_t(\theta)$ than $\widehat{\rho}_t(\theta)$

Avoid exponentiation \Rightarrow Much lower variance

PCPO (Lee and Ye, ICLR 2026)

$$\mathcal{L}_{\text{RSM}}(\theta) = \mathbb{E}_{t_i, \mathbf{x}_{t_i}, \epsilon} \left[C_1(t_i) \underbrace{\left(\|\mathbf{s}_{t_i}^\theta - (\mathbf{s}_{t_i}^{\text{ref}} + \gamma(t_i) \hat{\Psi}_{t_i})\|^2 \right)}_{\text{Guidance}} + C_2(t_i) \underbrace{\|\mathbf{s}_{t_i}^\theta - \mathbf{s}_{t_i}^{\theta_{\text{old}}}\|^2}_{\text{Regularization}} \right]$$

Flow-GRPO is also L_2 regression against the value-tilted score.

Three design knobs

$$\mathcal{L}_{\text{RSM}}(\theta) = \mathbb{E}_{t_i, \mathbf{x}_{t_i}, \epsilon} \left[C_1(t_i) \underbrace{\left(\|\mathbf{s}_{t_i}^\theta - (\mathbf{s}_{t_i}^{\text{ref}} + \gamma(t_i) \hat{\Psi}_{t_i})\|^2 \right)}_{\text{Guidance}} + C_2(t_i) \underbrace{\|\mathbf{s}_{t_i}^\theta - \mathbf{s}_{t_i}^{\theta_{\text{old}}}\|^2}_{\text{Regularization}} \right]$$

$\hat{\Psi}_{t_i}$

value guidance estimator

$C_1(t_i), \gamma(t_i)$

guidance strength

$C_2(t_i), \text{clipping}$

regularization

These three knobs explain differences among methods.

Knob 1: Estimator design

$$\hat{\Psi}_{t_i}^{\text{LA},1} = \frac{\sigma_{t_i}^2}{\alpha\Omega(t_i)} \frac{1}{K_i} \sum_{k=1}^{K_i} \nabla_{\mathbf{x}_{t_{i-1}}^{(k)}} r(\hat{\mathbf{x}}_{0|t_j}^{(k)}) \approx \frac{\sigma_{t_i}^2}{\alpha\Omega(t_i)} \mathbb{E}_{\mathbf{x}_{t_{i-1}:t_j}} \left[\nabla_{\mathbf{x}_{t_{i-1}}} r(\hat{\mathbf{x}}_{0|t_j}) \right]$$

$$\hat{\Psi}_{t_i}^{\text{LA},0} = \frac{\sigma_{t_i}}{\alpha\Omega(t_i)} \frac{1}{K_i} \sum_{k=1}^{K_i} r(\hat{\mathbf{x}}_{0|t_j}^{(k)}) \epsilon_{t_i}^{(k)} \approx \frac{\sigma_{t_i}}{\alpha\Omega(t_i)} \mathbb{E}_{\mathbf{x}_{t_{i-1}:t_j}} \left[r(\hat{\mathbf{x}}_{0|t_j}) \epsilon_{t_i} \right]$$

$\mathbb{E}_{\mathbf{x}}[\nabla r(\mathbf{x})]$ vs $\mathbb{E}_{\mathbf{x},\epsilon}[r(\mathbf{x})\epsilon]$

first-order vs zeroth-order

$\mathbb{E}_{\mathbf{x}} \Rightarrow \frac{1}{K_i} \sum_{k=1}^{K_i} \mathbf{x}^{(k)} \dots$

number of particles

$\mathbb{E}_{\mathbf{x}_0}[r(\mathbf{x}_0) \dots]$ vs $\mathbb{E}_{\mathbf{x}_{t_j}}[r(\hat{\mathbf{x}}_{0|t_j}) \dots]$

rollout depth

Knob 1: Estimator design (Part a)

$$\hat{\Psi}_{t_i}^{LA,1} \approx (\text{coefficient}) \mathbb{E}_{\mathbf{x}} \left[\nabla_{\mathbf{x}} r(\mathbf{x}) \right]$$
$$\hat{\Psi}_{t_i}^{LA,0} \approx (\text{coefficient}) \mathbb{E}_{\mathbf{x}} \left[r(\mathbf{x}) \epsilon_{t_i} \right]$$

$\mathbb{E}_{\mathbf{x}}[\nabla r(\mathbf{x})]$ vs $\mathbb{E}_{\mathbf{x},\epsilon}[r(\mathbf{x})\epsilon]$

first-order vs zeroth-order

Knob 1: Estimator design (Part **b**)

$$\hat{\Psi}_{t_i}^{\text{LA},1} = (\text{coefficient}) \frac{1}{K_i} \sum_{k=1}^{K_i} \nabla_{\mathbf{x}^{(k)}} r(\hat{\mathbf{x}}^{(k)}) \approx (\text{coefficient}) \mathbb{E}_{\mathbf{x}} \left[\nabla_{\mathbf{x}} r(\mathbf{x}) \right]$$

$$\hat{\Psi}_{t_i}^{\text{LA},0} = (\text{coefficient}) \frac{1}{K_i} \sum_{k=1}^{K_i} r(\hat{\mathbf{x}}^{(k)}) \boldsymbol{\epsilon}^{(k)} \approx (\text{coefficient}) \mathbb{E}_{\mathbf{x}} \left[r(\mathbf{x}) \boldsymbol{\epsilon}_{t_i} \right]$$

$\mathbb{E}_{\mathbf{x}}[\nabla r(\mathbf{x})]$ vs $\mathbb{E}_{\mathbf{x},\boldsymbol{\epsilon}}[r(\mathbf{x})\boldsymbol{\epsilon}]$

first-order vs zeroth-order

$$\mathbb{E}_{\mathbf{x}} \Rightarrow \frac{1}{K_i} \sum_{k=1}^{K_i} \mathbf{x}^{(k)} \dots$$

number of particles

Knob 1: Estimator design (Part c)

$$\hat{\Psi}_{t_i}^{\text{LA},1} = (\text{coefficient}) \frac{1}{K_i} \sum_{k=1}^{K_i} \nabla_{\mathbf{x}_{t_{i-1}}^{(k)}} r(\hat{\mathbf{x}}_{0|t_j}^{(k)}) \approx (\text{coefficient}) \mathbb{E}_{\mathbf{x}_{t_{i-1}:t_j}} \left[\nabla_{\mathbf{x}_{t_{i-1}}} r(\hat{\mathbf{x}}_{0|t_j}) \right]$$

$$\hat{\Psi}_{t_i}^{\text{LA},0} = (\text{coefficient}) \frac{1}{K_i} \sum_{k=1}^{K_i} r(\hat{\mathbf{x}}_{0|t_j}^{(k)}) \epsilon_{t_i} \approx (\text{coefficient}) \mathbb{E}_{\mathbf{x}_{t_{i-1}:t_j}} \left[r(\hat{\mathbf{x}}_{0|t_j}) \epsilon_{t_i} \right]$$

$\mathbb{E}_{\mathbf{x}}[\nabla r(\mathbf{x})]$ vs $\mathbb{E}_{\mathbf{x},\epsilon}[r(\mathbf{x})\epsilon]$

first-order vs zeroth-order

$\mathbb{E}_{\mathbf{x}} \Rightarrow \frac{1}{K_i} \sum_{k=1}^{K_i} \mathbf{x}^{(k)} \dots$

number of particles

$\mathbb{E}_{\mathbf{x}_0}[r(\mathbf{x}_0) \dots]$ vs $\mathbb{E}_{\mathbf{x}_{t_j}}[r(\hat{\mathbf{x}}_{0|t_j}) \dots]$

rollout depth

Knob 1: Estimator design

$$\hat{\Psi}_{t_i}^{\text{LA},1} = \frac{\sigma_{t_i}^2}{\alpha\Omega(t_i)} \frac{1}{K_i} \sum_{k=1}^{K_i} \nabla_{\mathbf{x}_{t_{i-1}}^{(k)}} r(\hat{\mathbf{x}}_{0|t_j}^{(k)}) \approx \frac{\sigma_{t_i}^2}{\alpha\Omega(t_i)} \mathbb{E}_{\mathbf{x}_{t_{i-1}:t_j}} \left[\nabla_{\mathbf{x}_{t_{i-1}}} r(\hat{\mathbf{x}}_{0|t_j}) \right]$$
$$\hat{\Psi}_{t_i}^{\text{LA},0} = \frac{\sigma_{t_i}}{\alpha\Omega(t_i)} \frac{1}{K_i} \sum_{k=1}^{K_i} r(\hat{\mathbf{x}}_{0|t_j}^{(k)}) \boldsymbol{\epsilon}_{t_i}^{(k)} \approx \frac{\sigma_{t_i}}{\alpha\Omega(t_i)} \mathbb{E}_{\mathbf{x}_{t_{i-1}:t_j}} \left[r(\hat{\mathbf{x}}_{0|t_j}) \boldsymbol{\epsilon}_{t_i} \right]$$

$$\mathbb{E}_{\mathbf{x}} [\nabla r(\mathbf{x})] \quad \text{vs} \quad \mathbb{E}_{\mathbf{x}, \epsilon} [r(\mathbf{x}) \epsilon]$$

First-order vs Zeroth-order

$$\mathbb{E}_{\mathbf{x}} \Rightarrow \frac{1}{K_i} \sum_{k=1}^{K_i} \mathbf{x}^{(k)} \dots$$

More particles = More compute, Less variance

$$\mathbb{E}_{\mathbf{x}_0} [r(\mathbf{x}_0) \dots] \quad \text{vs} \quad \mathbb{E}_{\mathbf{x}_{t_j}} [r(\hat{\mathbf{x}}_{0|t_j}) \dots]$$

More rollout depth = More compute, Less bias

Three design knobs

$$\mathcal{L}_{\text{RSM}}(\theta) = \mathbb{E}_{t_i, \mathbf{x}_{t_i}, \epsilon} \left[\underbrace{C_1(t_i) \left(\|\mathbf{s}_{t_i}^\theta - (\mathbf{s}_{t_i}^{\text{ref}} + \gamma(t_i) \hat{\Psi}_{t_i})\|^2 \right)}_{\text{Guidance}} + \underbrace{C_2(t_i) \|\mathbf{s}_{t_i}^\theta - \mathbf{s}_{t_i}^{\theta_{\text{old}}}\|^2}_{\text{Regularization}} \right]$$



$$\hat{\Psi}_{t_i}$$

value guidance estimator



$$C_1(t_i), \gamma(t_i)$$

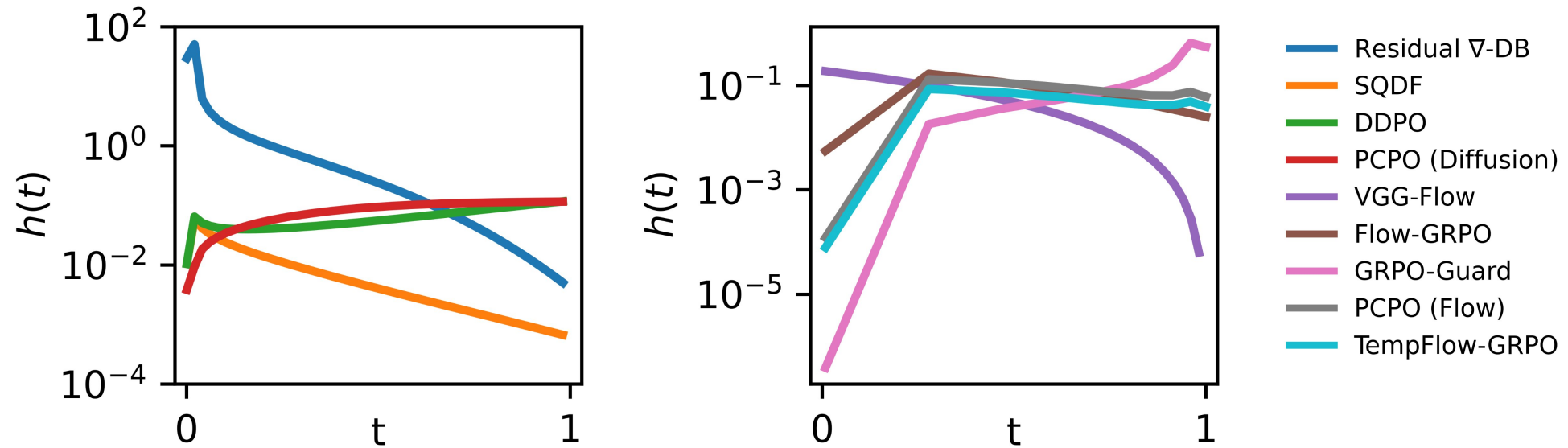
guidance strength

$$C_2(t_i), \text{clipping}$$

regularization

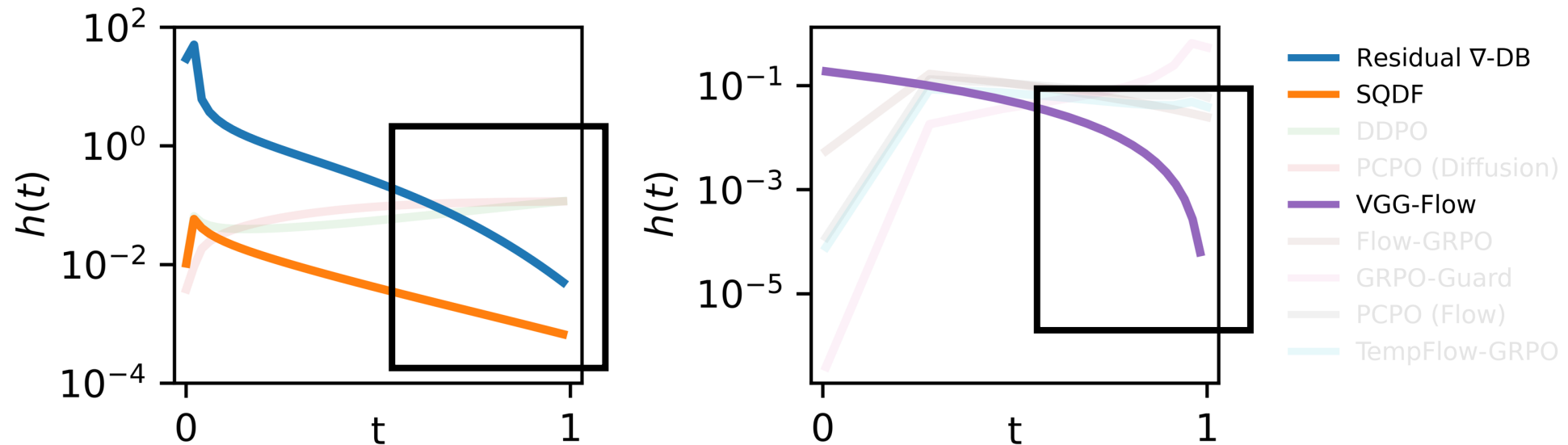
These three knobs explain differences among methods.

Knob 2: Guidance strength $h(t)$



$$h(t_i) := \boxed{C_1(t_i)\gamma(t_i)} \frac{\delta(t_i)\sigma_{t_i}^2}{\alpha\Omega(t_i)}$$

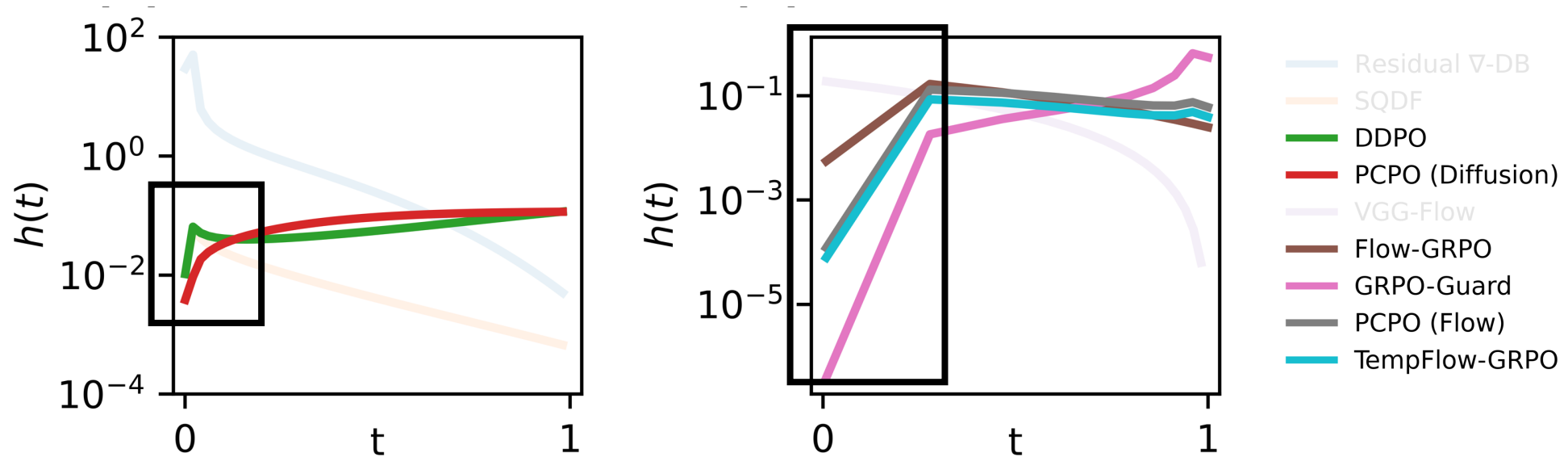
Knob 2: Guidance strength $h(t)$



$$h(t_i) := \boxed{C_1(t_i)\gamma(t_i)} \frac{\delta(t_i)\sigma_{t_i}^2}{\alpha\Omega(t_i)}$$

Existing methods suppress guidance at unreliable timesteps.

Knob 2: Guidance strength $h(t)$



$$h(t_i) := \boxed{C_1(t_i)\gamma(t_i)} \frac{\delta(t_i)\sigma_{t_i}^2}{\alpha\Omega(t_i)}$$

Existing methods suppress guidance at unreliable timesteps.

Three design knobs

$$\mathcal{L}_{\text{RSM}}(\theta) = \mathbb{E}_{t_i, \mathbf{x}_{t_i}, \epsilon} \left[C_1(t_i) \underbrace{\left(\|\mathbf{s}_{t_i}^\theta - (\mathbf{s}_{t_i}^{\text{ref}} + \gamma(t_i) \hat{\Psi}_{t_i})\|^2 \right)}_{\text{Guidance}} + C_2(t_i) \underbrace{\|\mathbf{s}_{t_i}^\theta - \mathbf{s}_{t_i}^{\theta_{\text{old}}}\|^2}_{\text{Regularization}} \right]$$

- ✓ $\hat{\Psi}_{t_i}$ value guidance estimator
- ✓ $C_1(t_i), \gamma(t_i)$ guidance strength
- $C_2(t_i), \text{clipping}$ regularization

These three knobs explain differences among methods.

Knob 3: Regularization

if (PPO-Clip condition):

$$\gamma(t_i) \leftarrow 0, C_2(t_i) \leftarrow 0$$

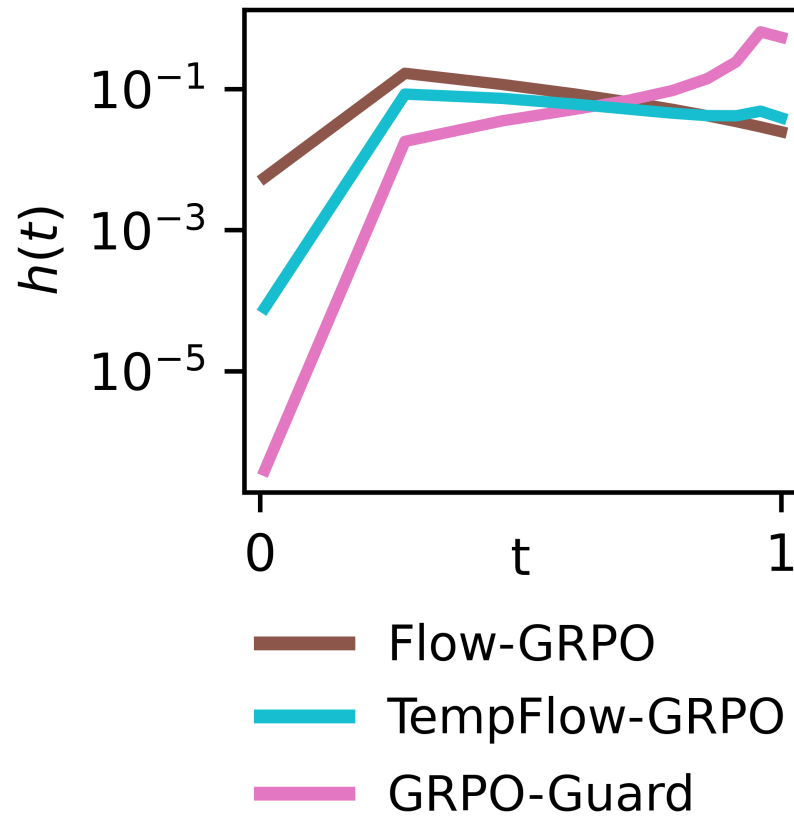
$$\mathcal{L}(\theta) = \mathbb{E}_{t_i, \mathbf{x}_{t_i}, \epsilon} \left[C_1(t_i) \underbrace{\left(\|\mathbf{s}_{t_i}^\theta - \mathbf{s}_{t_i}^{\text{ref}}\|^2 \right)}_{\text{Anchor}} \right]$$

else:

$$\mathcal{L}(\theta) = \mathbb{E}_{t_i, \mathbf{x}_{t_i}, \epsilon} \left[C_1(t_i) \underbrace{\left(\|\mathbf{s}_{t_i}^\theta - (\mathbf{s}_{t_i}^{\text{ref}} + \gamma(t_i) \hat{\Psi}_{t_i})\|^2 \right)}_{\text{Guidance}} + C_2(t_i) \underbrace{\|\mathbf{s}_{t_i}^\theta - \mathbf{s}_{t_i}^{\theta_{\text{old}}}\|^2}_{\text{Anchor}} \right]$$

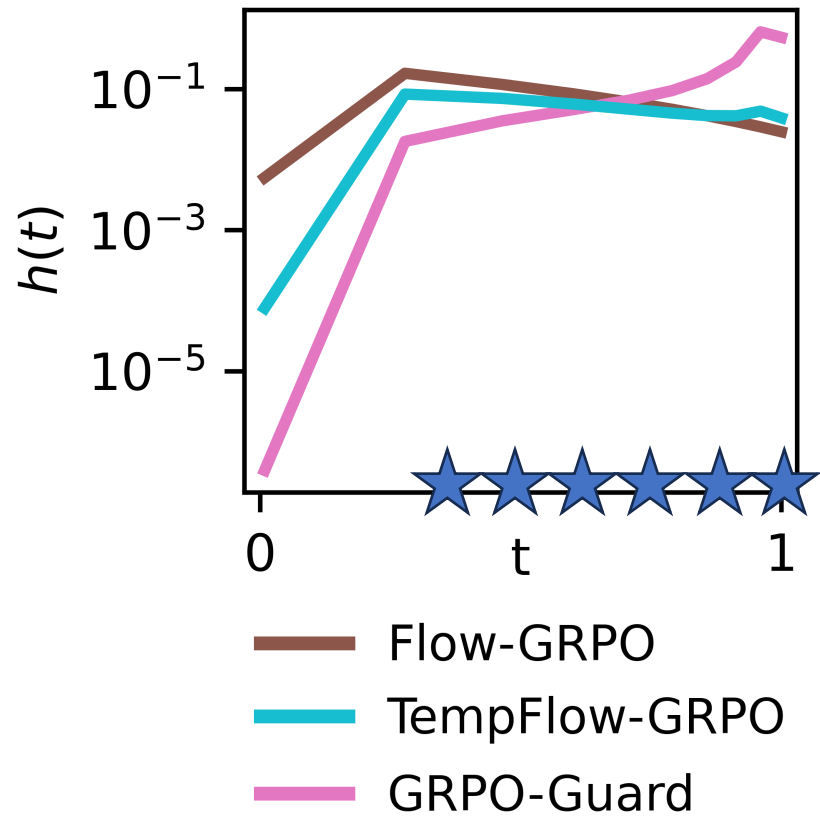
Off-policy updates are heuristically regularized.

Case Study: GRPO with GenEval



TempFlow-GRPO suppresses guidance at $t_i \approx 0$, but maintains $K_i = \text{const}$

Case Study: GRPO with GenEval

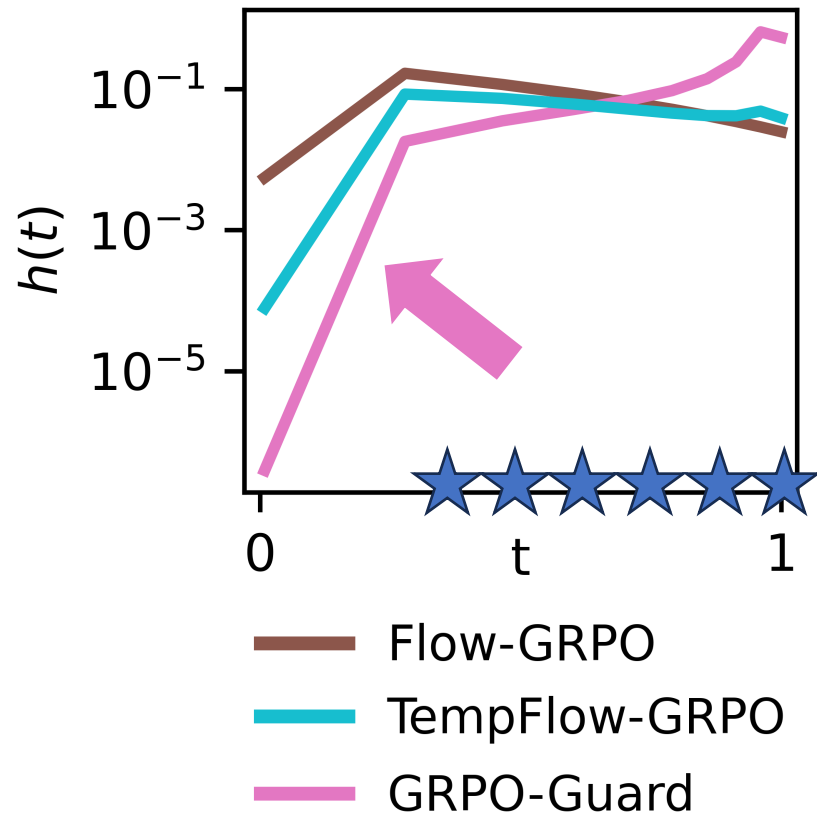


TempFlow-GRPO suppresses guidance at $t_i \approx 0$, but maintains $K_i = \text{const}$

Easy win

1. Concentrate $\{K_i\}$ on $\{t_i\}$ that matters

Case Study: GRPO with GenEval

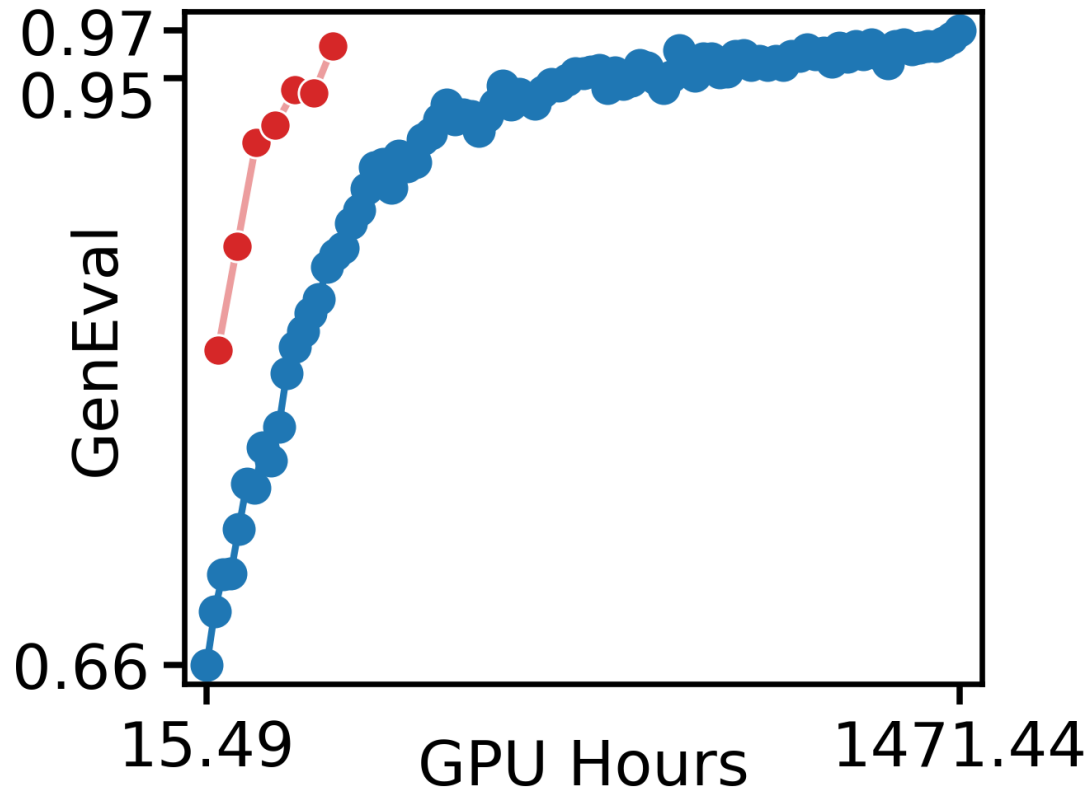


TempFlow-GRPO suppresses guidance at $t_i \approx 0$, but maintains $K_i = \text{const}$

Easy win

1. Concentrate $\{K_i\}$ on $\{t_i\}$ that matters
2. Double down on $h(t_i)$

Case Study: GRPO with GenEval



5x

wall-clock speedup
over **TempFlow-GRPO**

GenEval = 0.97

RSM \Rightarrow straightforward redesigns \Rightarrow improve efficiency

Reward-weighted Regression

“RWR is an approximation to Reward-Weighted MLE”

Diffusion RL tutorial (Uehara et. al., 2024)

$$\mathcal{L}_{\text{RSM}}(\theta) = \mathbb{E} \left[\underbrace{C_1(t_i) \|\mathbf{s}_{t_i}^\theta - \mathbf{s}_{t_i}^*\|^2}_{\text{Guidance}} + C_1(t_i)C_2(t_i) \underbrace{\|\mathbf{s}_{t_i}^\theta - \mathbf{s}_{t_i}^{\theta_{\text{old}}}\|^2}_{\text{Regularization}} \right] \quad \mathbf{s}_{t_i}^* \equiv \mathbf{s}_{t_i}^{\text{ref}} + \gamma(t_i) \hat{\Psi}_{t_i}$$
$$\mathcal{L}_{\text{NFT}}(\theta) \equiv \mathbb{E} \left[\underbrace{r(\mathbf{x}_0) \|\mathbf{v}_{t_i}^\theta - \mathbf{v}_{t_i}^{\text{fwd}}\|^2}_{\text{RWR}} + (\beta - r(\mathbf{x}_0)) \underbrace{\|\mathbf{v}_{t_i}^\theta - \mathbf{v}_{t_i}^{\theta_{\text{old}}}\|^2}_{\text{Regularization}} \right] \quad \mathbf{v}_{t_i}^{\text{fwd}} \equiv (\boldsymbol{\epsilon} - \mathbf{x}_0)$$

DiffusionNFT, AWM, RAM, ... are closely related.

Takeaways

$$\arg \max_q \mathbb{E}_{x_0 \sim q} [r(x_0)] - \alpha \mathcal{D}_{\text{KL}}(q \| p^{\text{ref}})$$

$$\mathcal{L}_{\text{RSM}}(\theta) = \mathbb{E}_{t_i, \mathbf{x}_{t_i}, \epsilon} \left[\underbrace{C_1(t_i) \left(\|\mathbf{s}_{t_i}^\theta - (\mathbf{s}_{t_i}^{\text{ref}} + \gamma(t_i) \hat{\Psi}_{t_i})\|^2 \right)}_{\text{Guidance}} + C_2(t_i) \underbrace{\|\mathbf{s}_{t_i}^\theta - \mathbf{s}_{t_i}^{\theta_{\text{old}}}\|^2}_{\text{Regularization}} \right]$$

Same problem, Same L_2 regression.

$\hat{\Psi}_{t_i}$

$C_1(t_i), \gamma(t_i)$

$C_2(t_i), \text{clipping}$

Understand three axes \Rightarrow Simple & effective redesigns.



{jaylee2000, jinhojsk515,
jeongsol, jong.ye}@kaist.ac.kr

Project Page

