

Sample Reward Soups: Query-Efficient Multi-Reward Guidance for Text-to-Image Diffusion Models

Yinghua Yao¹, Yuangang Pan^{1†}, Guoji Fu², Ivor Tsang^{1,3}

¹Agency for Science, Technology and Research, Singapore

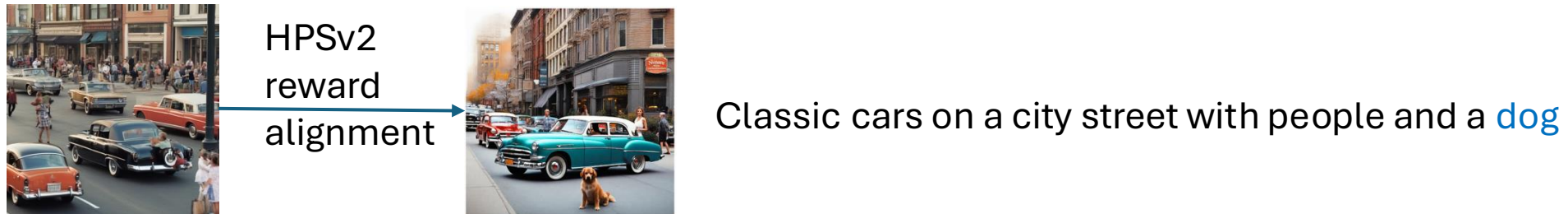
²National University of Singapore, Singapore

³Nanyang Technological University, Singapore

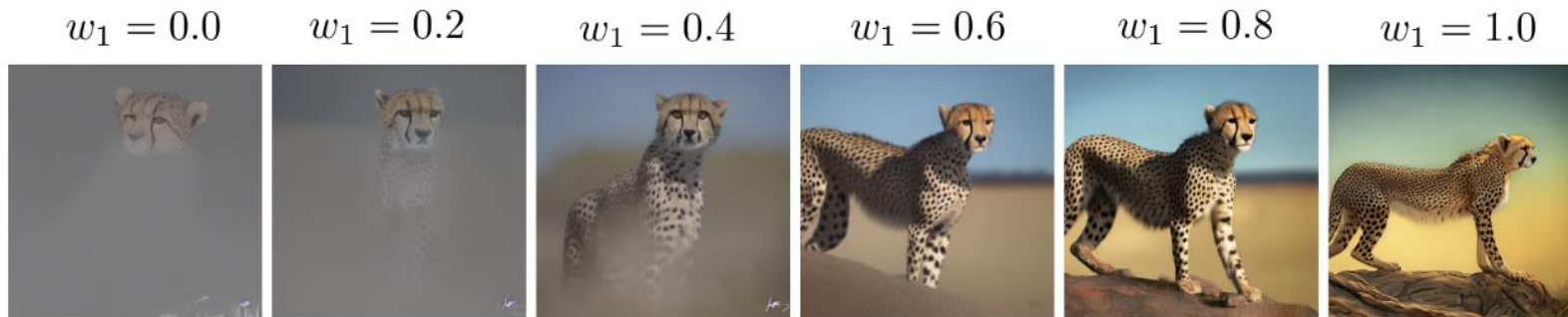


Alignment for Text-to-Image Models

- Alignment of Text-to-Image (T2I) models (diffusion/flow models)



- Multi-reward alignment: aesthetics, compressibility, etc



$w_1 f_1 + w_2 f_2$: aesthetic reward f_1 and compressibility reward f_2

Alignment Techniques

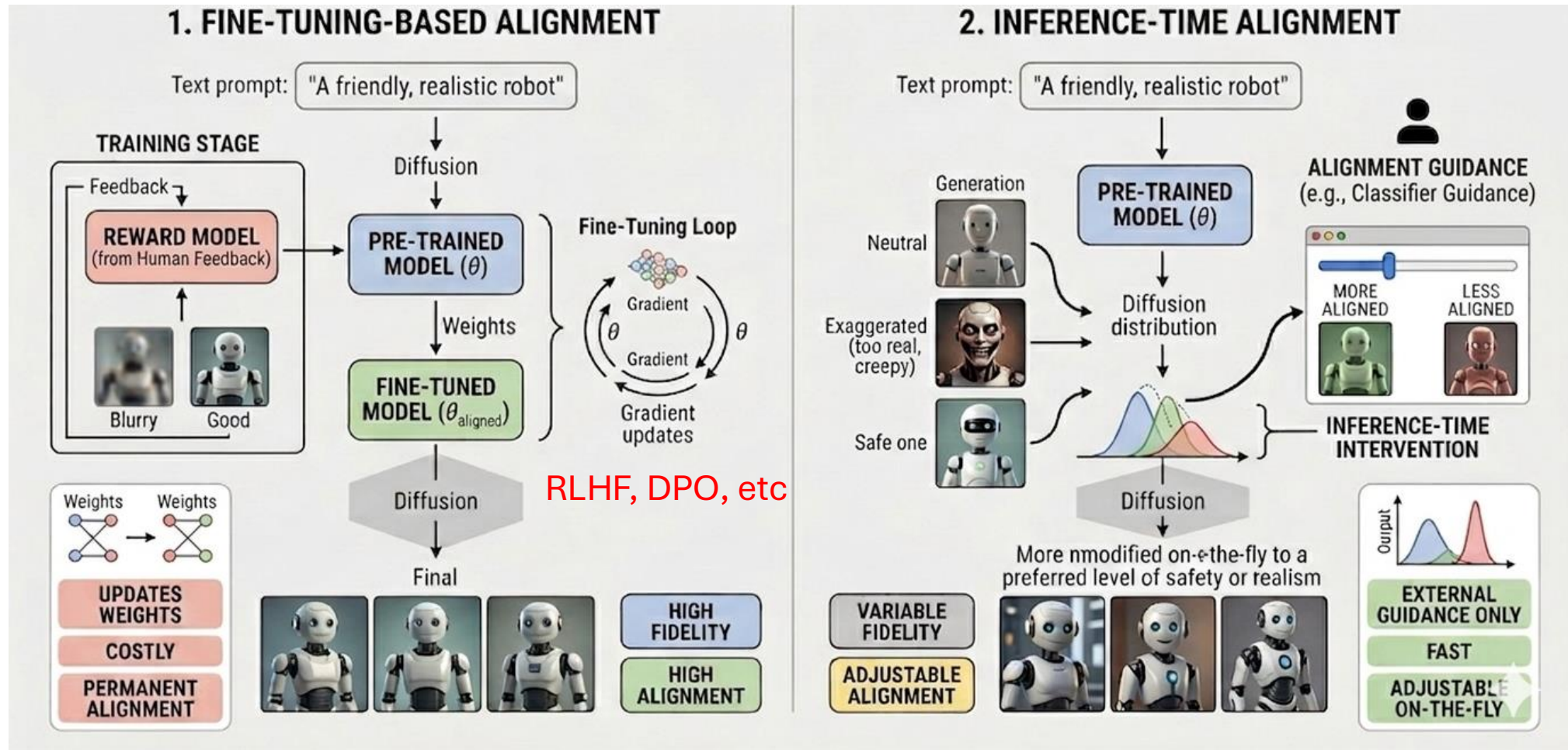
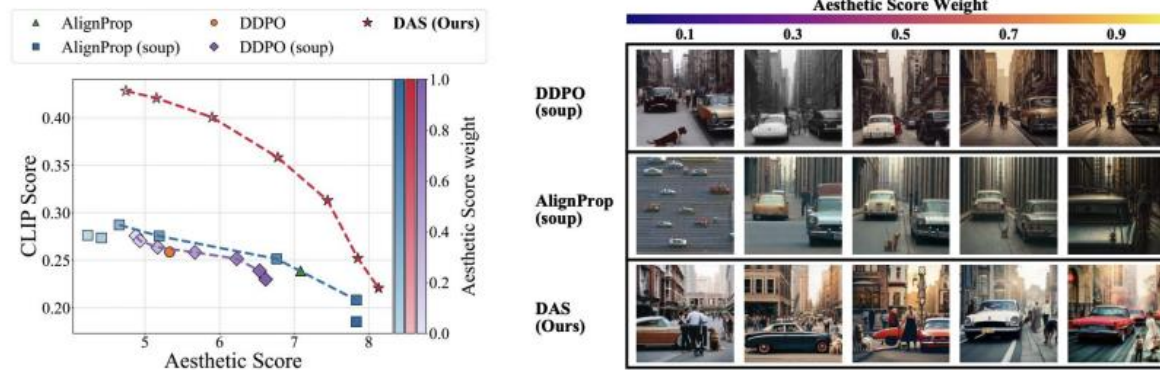


Figure is created by Gemini.

Challenges of Multi-Reward Alignment

- Weighted-Sum: $F(x, w_{1:M}) = \sum_{i=1}^M w_i f_i(x)$
- Reward combinations grows exponentially with the number of reward functions: $\{w_{1:M}^j\}_{j=1}^L, L \gg M$

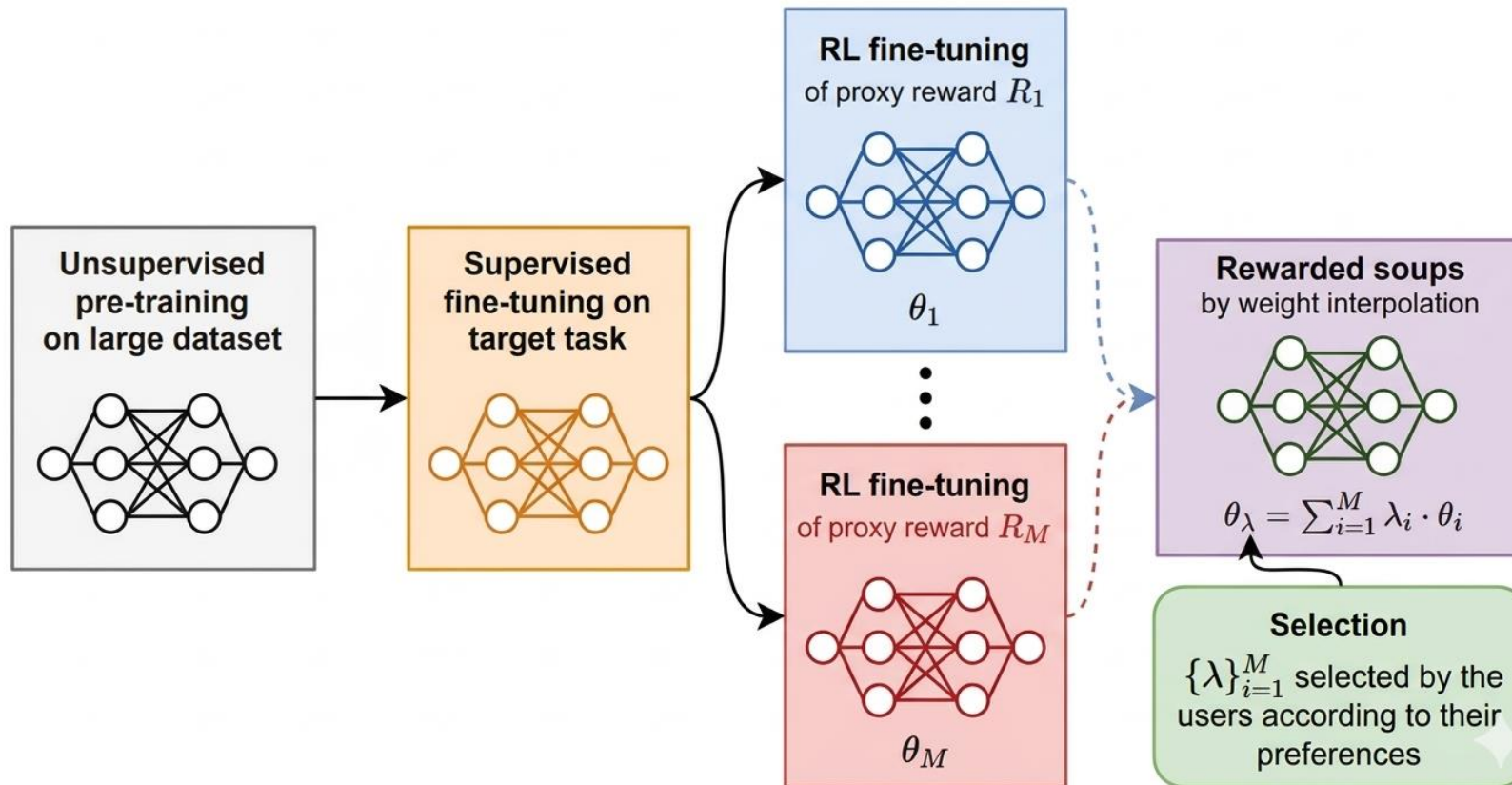
$$w \cdot \text{Aesthetic Score} + (1 - w) \cdot 20 \cdot \text{CLIPScore}$$



(a) Trade-off in multi-objective optimization.

(b) Generated samples according to reward weights

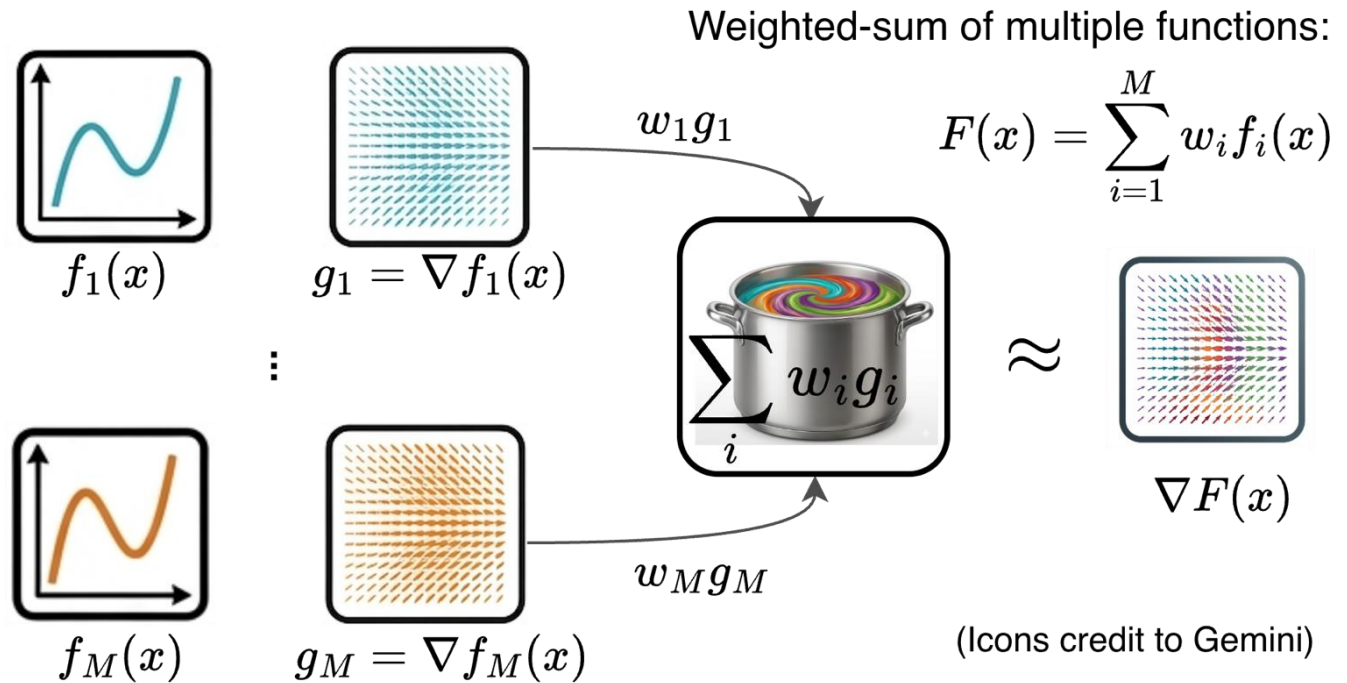
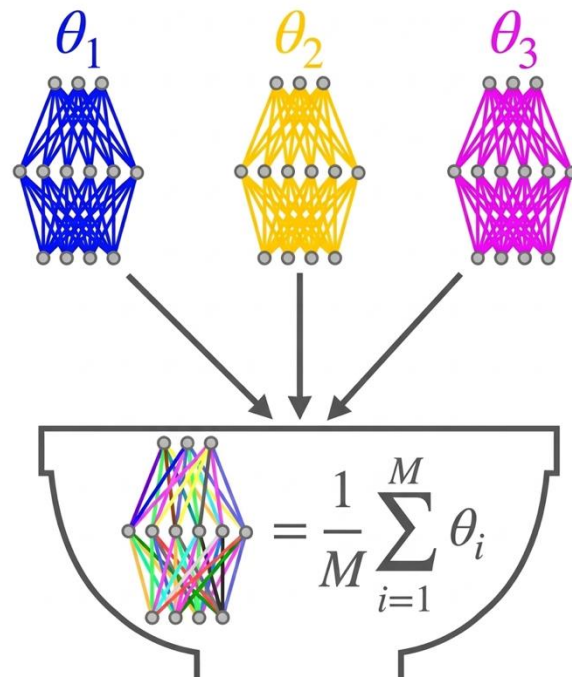
Rewarded Soups for Efficient Fine-tuning



NO efficient strategy for inference-time multi-reward alignment

Figure is from “Rewarded soups: towards Pareto-optimal alignment by interpolating weights fine-tuned on diverse rewards,” NeurIPS 2023. Small edits credit to Gemini.

Weight Soups \longrightarrow Gradient Soups



(Icons credit to Gemini)

First inference-time soup strategy!!!

- Reward guidance subject to the number of reward functions M , not the number of reward combinations L . $M \ll L$
- Gradient soups can approximate the gradient of weighted function under the black-box denoising alignment process

Black-box Reward Inference-time Alignment

- Optimize the denoising distribution at t step for high-reward generation

$$\max_{\boldsymbol{\mu}_{t-1}} \mathcal{F}(\boldsymbol{\mu}_{t-1}) = \mathbb{E}_{\mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{t-1}, \beta_t \mathbf{I})} [f(\mathbf{x}_{t-1})] \quad \boldsymbol{\mu}_{t-1} = \mu_{\theta}(\mathbf{x}_t, t)$$

- Reward-guided search gradient

$$\nabla_{\boldsymbol{\mu}_{t-1}} \mathcal{F}(\boldsymbol{\mu}_{t-1}) = \frac{1}{\sqrt{\beta_t}} \mathbb{E}_{\mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})} \left[\underbrace{f(\boldsymbol{\mu}_{t-1} + \sqrt{\beta_t} \mathbf{z}) \mathbf{z}}_{\text{Reward query}} \right]$$

Obtain rewards for N samples. Update $\boldsymbol{\mu}_{t-1}$ towards high-reward regions

Sample Reward Soups

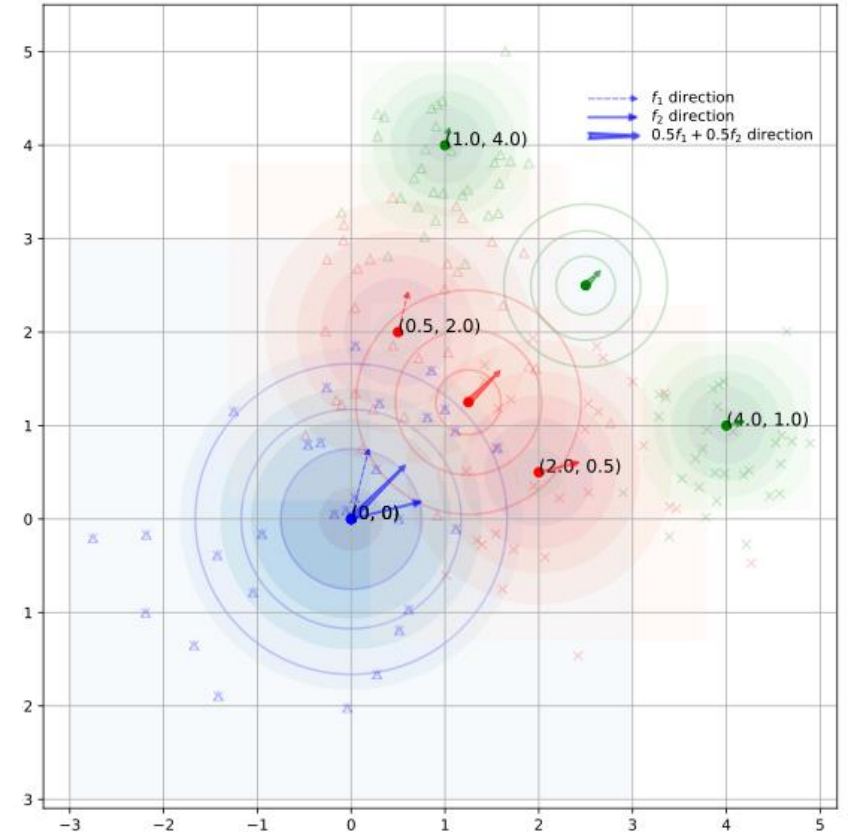
- Obtain search gradient for each reward function

$$\nabla_{\mathbf{c}_{t-1}^m} \mathcal{F}(\mathbf{c}_{t-1}^m, e_m)$$

- Interpolate search gradients

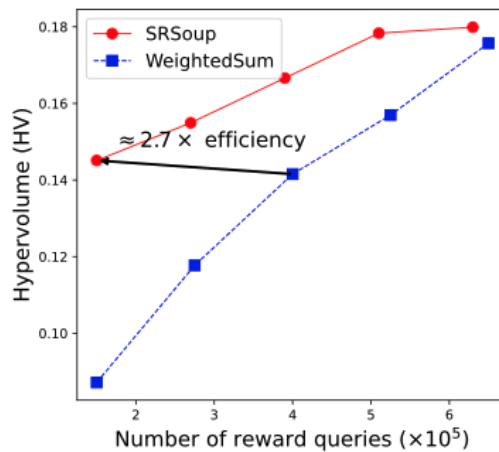
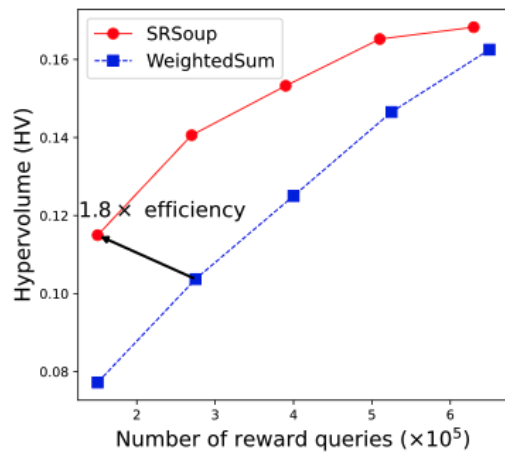
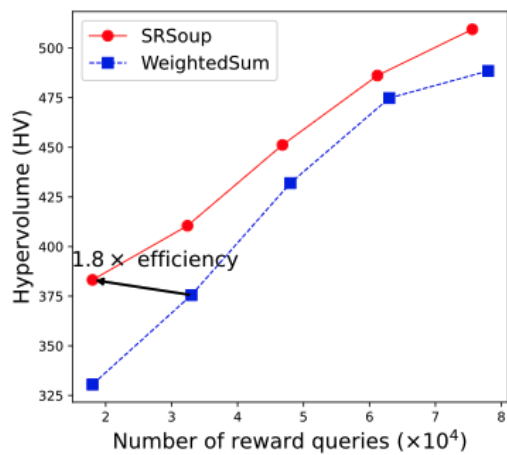
$$\begin{aligned} \nabla_{\boldsymbol{\mu}_{t-1}^l} \mathcal{F}(\boldsymbol{\mu}_{t-1}^l, w_{1:M}^l) &= \sum_{m=1}^M w_m^l \left(\nabla_{\mathbf{c}_{t-1}^m} \mathcal{F}(\mathbf{c}_{t-1}^m, e_m) \right. \\ &\quad \left. + \frac{1}{\sqrt{\beta_t} N} \sum_{n=1}^N f_m(\mathbf{x}_{t-1}^{m,n}) (\mathbf{c}_{t-1}^m - \boldsymbol{\mu}_{t-1}^l) \right) \end{aligned}$$

- Hybrid schedule: interpolate gradients at the first K steps, weighted-sum update for later steps

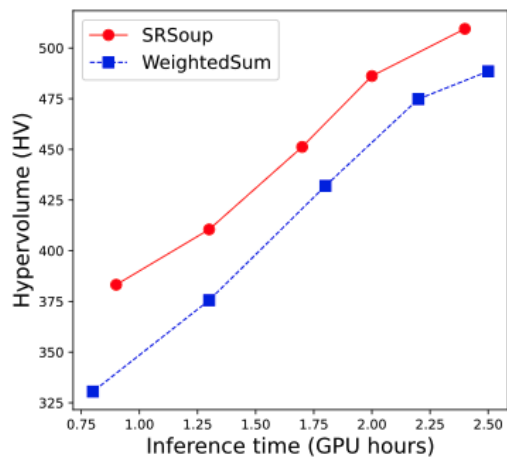


Insight: Early denoising distributions for different rewards are still close, so sample rewards can be shared

Efficiency Gains

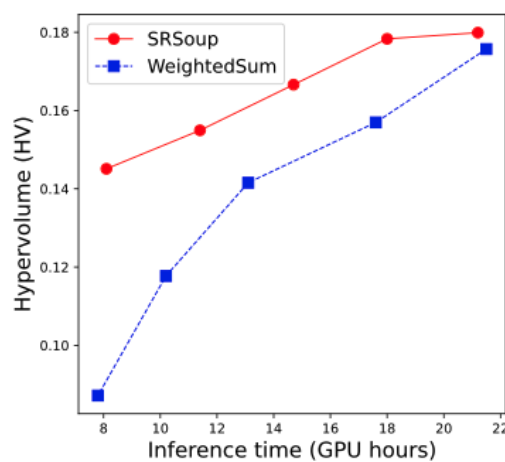


Aesthetic & Compressibility



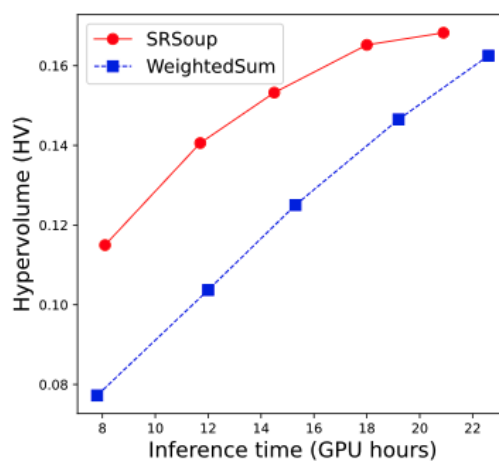
Aesthetic & Compressibility

Aesthetic & PickScore



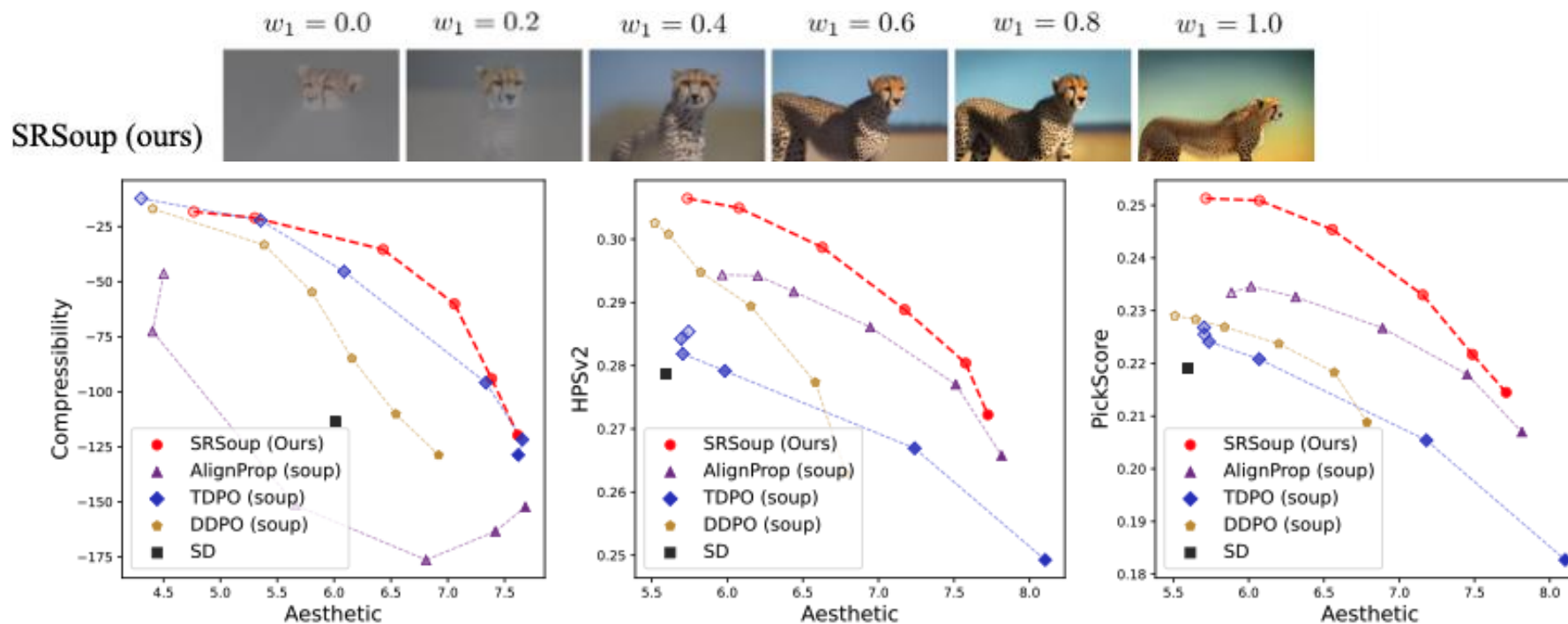
Aesthetic & PickScore

Aesthetic & HPSv2



Aesthetic & HPSv2

Compared with Model Soup Baselines



Conclusion

- First **inference-time soup strategy (gradient soups)** for multi-reward diffusion alignment.
- Enables **training-free Pareto sampling** over the full preference space.
- Achieves **1.8X-2.7X** better query efficiency than weighted-sum guidance.

Code Link:



Thank you!