



## Introduction

Large-scale diffusion models have achieved remarkable performance in generative tasks. Beyond their initial training applications, these models have proven their ability to function as versatile plug-and-play priors. For instance, 2D diffusion models can serve as loss functions to optimize 3D implicit models. Rectified Flow, a novel class of generative models, has demonstrated superior performance across various domains. Compared to diffusion-based methods, rectified flow approaches surpass them in terms of generation quality and efficiency. In this work, we present theoretical and experimental evidence demonstrating that rectified flow based methods offer similar functionalities to diffusion models — they can also serve as effective priors. Besides the generative capabilities of diffusion priors, motivated by the unique time-symmetry properties of rectified flow models, a variant of our method can additionally perform image inversion. Experimentally, our rectified flow based priors outperform their diffusion counterparts — the SDS and VSD losses — in text-to-3D generation. Our method also displays competitive performance in image inversion and editing.

**TL;DR: Distill knowledge from 2D flow-matching models for 3D generation and 2D editing.**

## Our Methods

1), Baseline: RFDS – equivalent to Score Distillation Sampling

We are interested in finding the gradient with respect to the input given a trained flow matching model. Starting from the training loss of flow matching:

$$\int_0^1 \mathbb{E} \left[ \left\| v_\phi(\mathbf{x}_t, t) - \dot{\alpha}_t \mathbf{x} - \dot{\sigma}_t \epsilon \right\|^2 \right] dt, \text{ with } \mathbf{x} = g(\theta) \text{ and } \mathbf{x}_t = \alpha_t \mathbf{x} + \sigma_t \epsilon. \quad (1)$$

To find  $\theta$ , the gradients of the above equation with respect to  $\theta$  can be written as:

$$\nabla_\theta \mathcal{L}_{\text{rds}}(\phi, \mathbf{x}, \epsilon, t) = 2 \times \mathbb{E} \left[ \underbrace{(v_\phi(\mathbf{x}_t, t) - \dot{\alpha}_t \mathbf{x} - \dot{\sigma}_t \epsilon)}_{\text{Flow Residual}} \underbrace{\left( \frac{\partial (v_\phi(\mathbf{x}_t, t) - \dot{\alpha}_t \mathbf{x} - \dot{\sigma}_t \epsilon)}{\partial \mathbf{x}} \right)}_{\text{Generator Jacobian}} \underbrace{\left( \frac{\partial \mathbf{x}}{\partial \theta} \right)}_{\text{Generator Jacobian}} \right]. \quad (2)$$

It can be further simplified as:

$$\nabla_\theta \mathcal{L}_{\text{rds}}(\phi, \mathbf{x}, \epsilon, t) = 2 \times \mathbb{E} \left[ \underbrace{(v_\phi(\mathbf{x}_t, t) - \dot{\alpha}_t \mathbf{x} - \dot{\sigma}_t \epsilon)}_{\text{Flow Residual}} \underbrace{\left( \frac{\partial (-\dot{\alpha}_t \mathbf{x} - \dot{\sigma}_t \epsilon)}{\partial \mathbf{x}} \right)}_{-\dot{\alpha}_t} + \underbrace{\left( \frac{\partial v_\phi(\mathbf{x}_t, t)}{\partial \mathbf{x}_t} \right)}_{\text{Network Jacobian}} \underbrace{\left( \frac{\partial \mathbf{x}_t}{\partial \alpha_t} \right)}_{\dot{\alpha}_t} \underbrace{\left( \frac{\partial \mathbf{x}}{\partial \theta} \right)}_{\text{Generator Jacobian}} \right]. \quad (3)$$

Finally:

$$\nabla_\theta \mathcal{L}_{\text{rds}}(\phi, \mathbf{x}, \epsilon, t) \simeq \mathbb{E} \left[ w(t) \underbrace{(v_\phi(\mathbf{x}_t, t) - \dot{\alpha}_t \mathbf{x} - \dot{\sigma}_t \epsilon)}_{\text{Flow Residual}} \underbrace{\left( \frac{\partial \mathbf{x}}{\partial \theta} \right)}_{\text{Generator Jacobian}} \right]. \quad (4)$$

2), inverse RFDS (iRFDS) – find the noise instead of the image – for image editing:

$$\nabla_\epsilon \mathcal{L}_{\text{irfds}}(\phi, \mathbf{x}, \epsilon, t) = 2 \times \mathbb{E} \left[ \underbrace{(v_\phi(\mathbf{x}_t, t) - \dot{\alpha}_t \mathbf{x} - \dot{\sigma}_t \epsilon)}_{\text{Flow Residual}} \underbrace{\left( \frac{\partial (-\dot{\alpha}_t \mathbf{x} - \dot{\sigma}_t \epsilon)}{\partial \epsilon} \right)}_{-\dot{\sigma}_t} + \underbrace{\left( \frac{\partial v_\phi(\mathbf{x}_t, t)}{\partial \mathbf{x}_t} \right)}_{\text{Network Jacobian}} \underbrace{\left( \frac{\partial \mathbf{x}_t}{\partial \sigma_t} \right)}_{\dot{\sigma}_t} \right]. \quad (5)$$

Then

$$\nabla_\epsilon \mathcal{L}_{\text{irfds}}(\phi, \mathbf{x}, \epsilon, t) \simeq \mathbb{E} \left[ w'(t) \underbrace{(v_\phi(\mathbf{x}_t, t) - \dot{\alpha}_t \mathbf{x} - \dot{\sigma}_t \epsilon)}_{\text{Flow Residual}} \right]. \quad (6)$$

3), RFDS Reversal (RFDS-Rev) – improved RFDS

We propose a novel method that improves upon the baseline RFDS by iteratively applying RFDS for image optimization and iRFDS for image inversion. Please refer to our paper for detailed derivations.

**Algorithm 1:** The RFDS-Rev Algorithm.

```

1 Initialize the learnable parameter  $\theta$ 
2 while Not Converge do
3   Sample random timestep  $t$ 
4   Sample random noise  $\epsilon$ 
5   for  $n$  steps do
6     Freeze  $\theta$ , optimize  $\epsilon$  with iRFDS (Eq. 6)
7   Optimize  $\theta$  with optimized  $\epsilon$  based on RFDS (Eq. 4)
8 RETURN  $\theta$ 

```

## Our Results



Figure 1. Demonstration of potential use cases of the proposed rectified flow priors. Our methods can be used on 2D inversion and editing and text-to-3D generation.

## Additional Results



Figure 2. 3D Comparison. Text Prompts: “An antique wooden rocking horse”, “A bear holding a sign that says Hello world” and “Hot popcorn jump out from the red striped popcorn maker”.

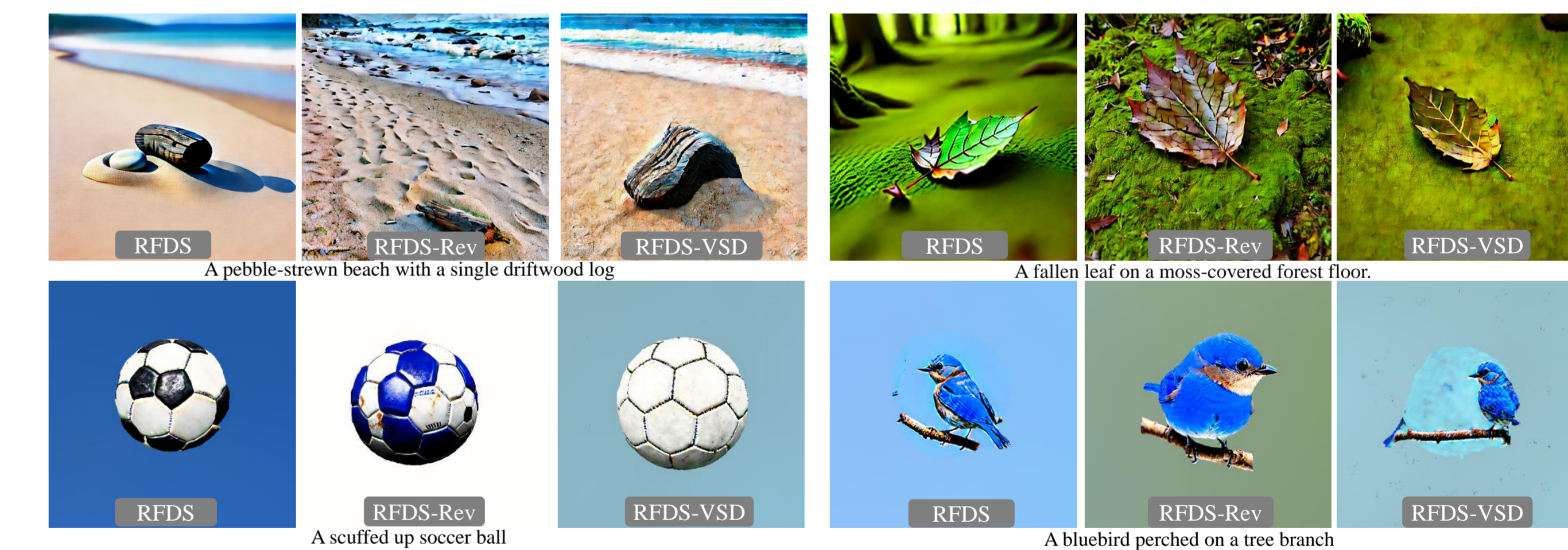


Figure 3. Ablation experiments of RFDS-Rev vs. RFDS-VSD. Top: 2D case. Bottom: 3D case.

Table 1. Results of Text-to-3D on T<sup>3</sup>Bench dataset. Our proposed methods achieve state-of-the-art performance among 2D lifting methods, beating the diffusion based SDS and VSD priors.

Method	Dreamfusion (SDS)	ProlificDreamer (VSD)	RFDS	RFDS-Rev	RFDS	RFDS-Rev
Base Model	Diffusion	Diffusion	InstaFlow	InstaFlow	SD3	SD3
Single Object	24.9	51.1	46.9	57.6	49.4	<b>59.8</b>
Surroundings	19.3	42.5	34.5	44.6	39.3	<b>54.5</b>
Multiple Objects	17.3	45.7	28.3	44.6	42.9	<b>55.2</b>
Average	20.5	46.4	36.5	48.9	43.9	<b>56.5</b>