

Table 1: CIFAR-10 posterior sampling results for CNF prior. We report expected classifier log probability and FID scores for the class posteriors, averaged over all 10 classes.

| Model | Sampler                 | $\mathbb{E}[\log p(\mathbf{y} \mid \mathbf{x})](\uparrow)$ | FID ( $\downarrow$ ) | $\log Z(\mathbf{y})(\uparrow)$ |
|-------|-------------------------|--|----------------------|--------------------------------|
| I-CFM | Prior                   | -5.88  | 84.79                | -24.04                         |
|       | DPS                     | -2.22  | 84.96                | -                              |
|       | Latent HMC              | -2.80  | 46.69                | -                              |
|       | RTB                     |  |                      |                                |
|       | Adj. Matching           | -3.09  | 19.45                | -17.23                         |
|       | <b>Outsourced Diff.</b> | -3.35  | 34.28                | -20.36                         |

Table 2: SD 1.5 fine-tuning results, averaged across three prompts used in [Venkatraman et al.]. DPOK, DDPO and RTB results taken from the same paper.<sup>2</sup>

| Sampler                 | $\mathbb{E}[\log r(\mathbf{x}, \mathbf{y})](\uparrow)$ | CLIP diversity ( $\uparrow$ ) |
|-------------------------|--|-------------------------------|
| Prior                   | -0.17  | 0.18                          |
| DDPO                    | 1.37   | 0.09                          |
| DPOK                    | 1.23   | 0.13                          |
| RTB                     | 1.4  | 0.11                          |
| <b>Outsourced Diff.</b> | 1.26   | 0.14                          |

<sup>2</sup>"A green rabbit.", "A cat and a dog.", and "Four roses."

---

**Algorithm 1** Training loop for Outsourced Diffusion Sampler

---

- 1: **Initialize:** deterministic prior function  $f$  (e.g. CNF integrators, GAN generator), randomly initialized noise posterior model  $p_F^\phi$ , randomly initialized  $Z^\phi(\mathbf{y})$  (scalar for fixed  $\mathbf{y}$ ), VP-SDE backward policy  $p_B$ , log reward function  $\log r(\mathbf{x}, \mathbf{y})$ , on-policy update fraction  $p$ .
  - 2: **for** each step  $n = 1, 2, \dots, N$  **do**
  - 3:   Sample a batch of trajectories:  $\{\tau^{(i)}\}_{i=1}^B \sim p_F^\phi(\tau \mid \mathbf{y})$
  - 4:   **for**  $i = 1, \dots, B$  **do**
  - 5:     Compute log density:  $\log R^{(i)} \leftarrow \log \mathcal{N}(\mathbf{z}^{(i)}; \mathbf{0}, \mathbf{I}) + r(f(\mathbf{z}^{(i)}), \mathbf{y})$
  - 6:     Store experience  $(\tau^{(i)}, \log R^{(i)})$  in replay buffer  $\mathcal{D}$
  - 7:   **end for**
  - 8:   Draw  $u \sim \text{Uniform}(0, 1)$
  - 9:   **if**  $u \leq p$  **then**
  - 10:     Keep on-policy batch  $\{(\tau^{(i)}, \log R^{(i)})\}_{i=1}^B$
  - 11:   **else**
  - 12:     Sample off-policy batch  $\{(\tau^{(i)}, \log R^{(i)})\}_{i=1}^B \sim \mathcal{D}$
  - 13:   **end if**
  - 14:   Compute  $\mathcal{L}_{\text{TB}}(\tau; \mathbf{y}, \phi)$  for batch using TB loss eq(4).
  - 15:   Update  $p_F^\phi$ ,  $Z^\phi(\mathbf{y})$  using  $\nabla_\phi \mathcal{L}_{\text{TB}}(\tau; \mathbf{y}, \phi)$ .
  - 16: **end for**
-