

# Maximum Likelihood Training of Score-Based Diffusion Models

Yang Song\*, Conor Durkan\*, Iain Murray, Stefano Ermon  
Stanford University and University of Edinburgh

## ABSTRACT

- Score-based diffusion models synthesize images by reversing a stochastic process that diffuses data to noise, and are trained by minimizing a weighted combination of score matching losses.
- We show that by choosing a special weighting function, called the likelihood weighting, minimizing the weighted combination of score matching losses amounts to maximum likelihood training.
- Our theoretical results enable ScoreFlow, a continuous normalizing flow model trained with a variational objective, which is much more efficient than neural ODEs. We report the state-of-the-art likelihood on CIFAR-10 and ImageNet 32x32 among all flow models, achieving comparable performance to cutting-edge autoregressive models.

Code:

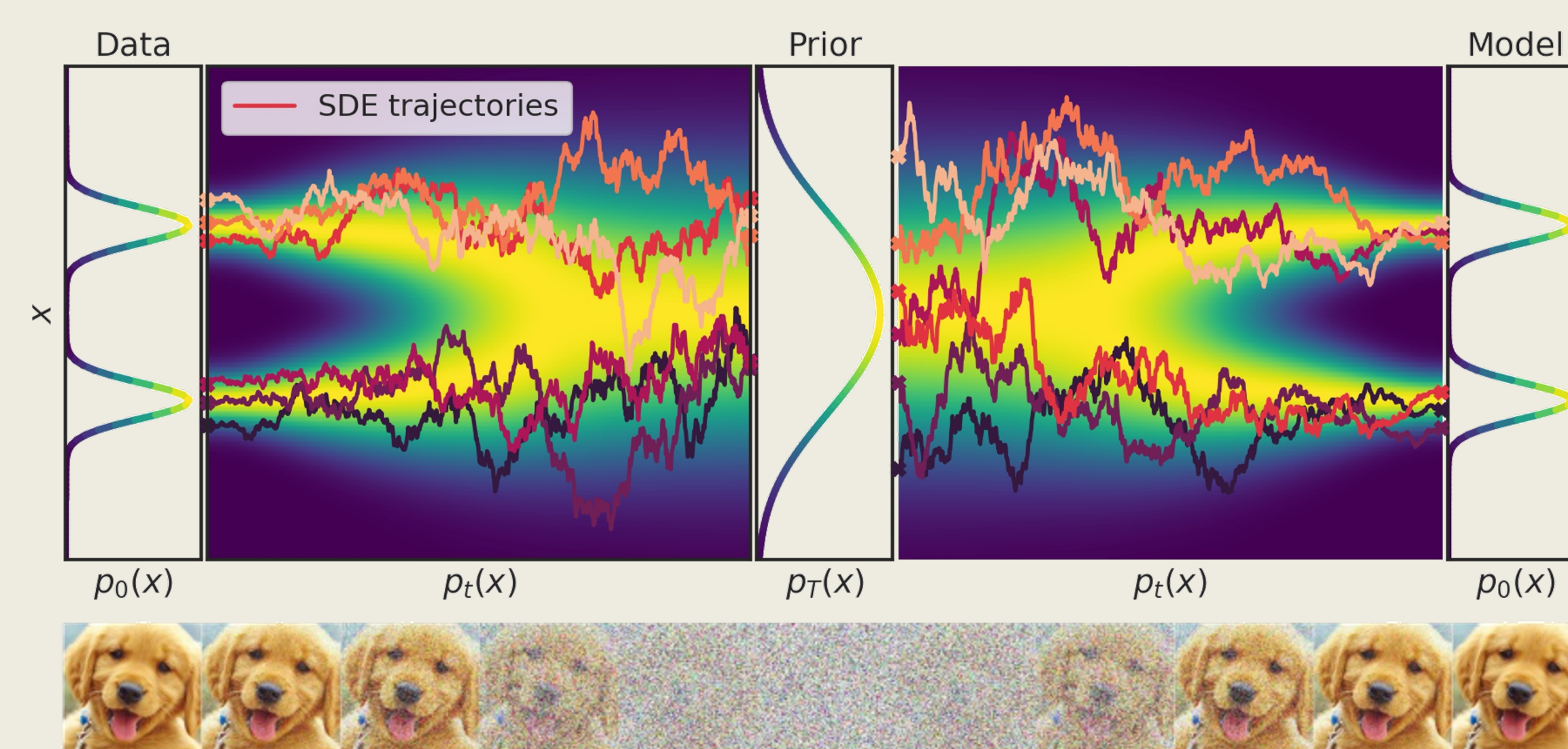


## Score-Based Diffusion Models

**Stochastic differential equation (SDE):**

$$dx = \underbrace{f(x, t) dt}_{\text{Deterministic drift}} + \underbrace{g(t) dw}_{\text{Stochastic diffusion}} \quad \text{Brownian motion}$$

Perturbing data with a fixed SDE, and reverse it for generative modeling



**The reverse-time SDE:**

$$dx = [f(x, t) - g^2(t) \nabla_x \log p_t(x)] dt + g(t) dw$$

↓  
Score function of  $p_t(x)$

- Must be solved in the reverse time direction
- Requires estimating score functions at all time steps.

**Estimating the score function for the reverse SDE:**

- Time-dependent score-based model**

$$s_{\theta}(\cdot, t) : \mathbb{R}^d \rightarrow \mathbb{R}^d$$

- Goal:**  $s_{\theta}(\mathbf{x}, t) \approx \nabla_x \log p_t(\mathbf{x})$

- Training objective:**

$$\min_{\theta} \mathbb{E}_{t \sim \mathcal{U}(0, T)} \mathbb{E}_{\mathbf{x} \sim p_t(\mathbf{x})} [\lambda(t) \|\mathbf{s}_{\theta}(\mathbf{x}, t) - \nabla_x \log p_t(\mathbf{x})\|_2^2]$$

- Score matching (Hyvarinen 2005)**

- Denosing score matching (Vincent 2010)
- Sliced score matching (Song et al., 2019)

## Two Ways to Define Likelihoods

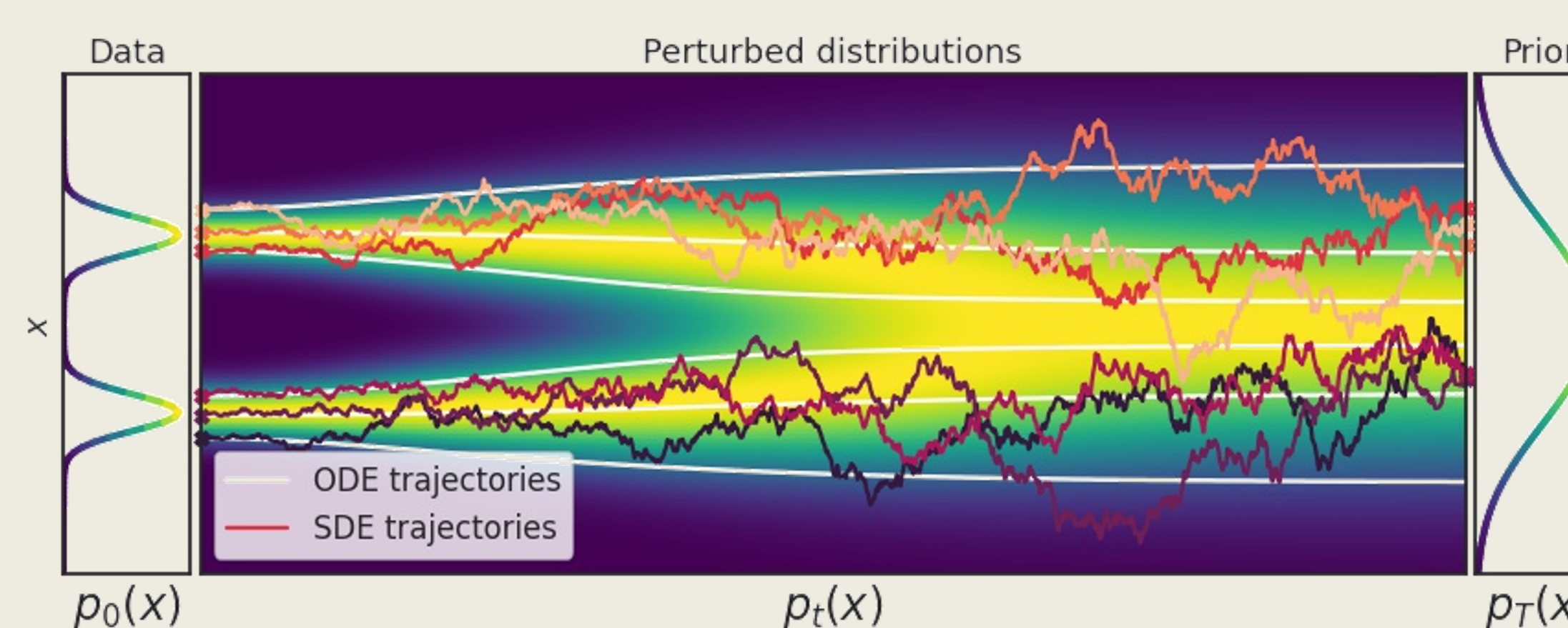
**Probability flow ODE:**

$$dx = \left[ f(\mathbf{x}, t) - \frac{1}{2} g^2(t) \nabla_x \log p_t(\mathbf{x}) \right] dt$$

$$\{p_t(\mathbf{x})\}_{t \in [0, T]}$$

$$dx = [f(\mathbf{x}, t) - g^2(t) \nabla_x \log p_t(\mathbf{x})] dt + g(t) dw$$

Can sample from the same distribution by solving the ODE instead of the SDE.



**Prior distribution:**  $\pi(\mathbf{x}) \approx p_T(\mathbf{x})$

**Model 1:** Reversing the probability flow ODE

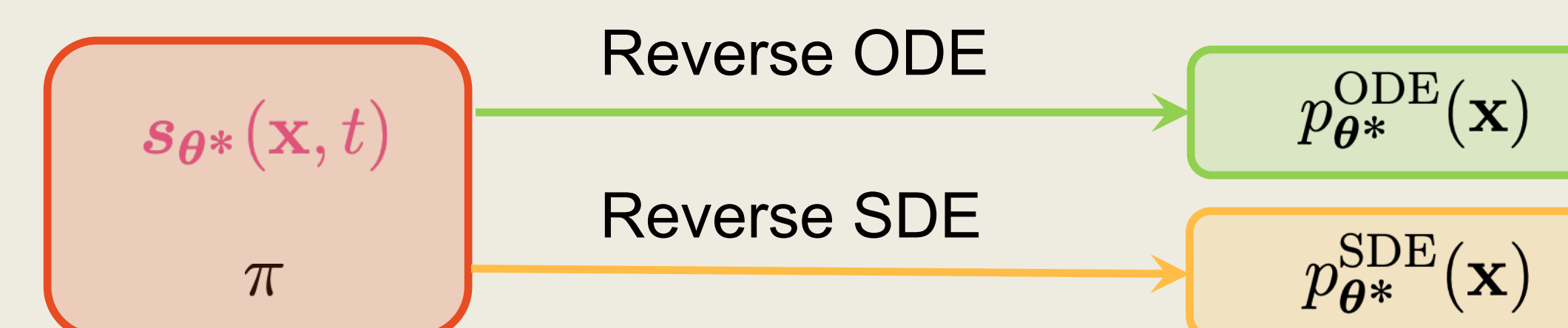
$$d\tilde{\mathbf{x}} = \left[ f(\tilde{\mathbf{x}}, t) - \frac{1}{2} g(t)^2 s_{\theta^*}(\tilde{\mathbf{x}}, t) \right] dt, \quad \tilde{\mathbf{x}}(T) \sim \pi$$

$$\tilde{\mathbf{x}}(0) \sim p_{\theta^*}^{\text{ODE}}$$

**Model 2:** Reversing the SDE

$$d\hat{\mathbf{x}} = [f(\hat{\mathbf{x}}, t) - g(t)^2 s_{\theta^*}(\hat{\mathbf{x}}, t)] dt + g(t) dw, \quad \hat{\mathbf{x}}(T) \sim \pi$$

$$\hat{\mathbf{x}}(0) \sim p_{\theta^*}^{\text{SDE}}$$



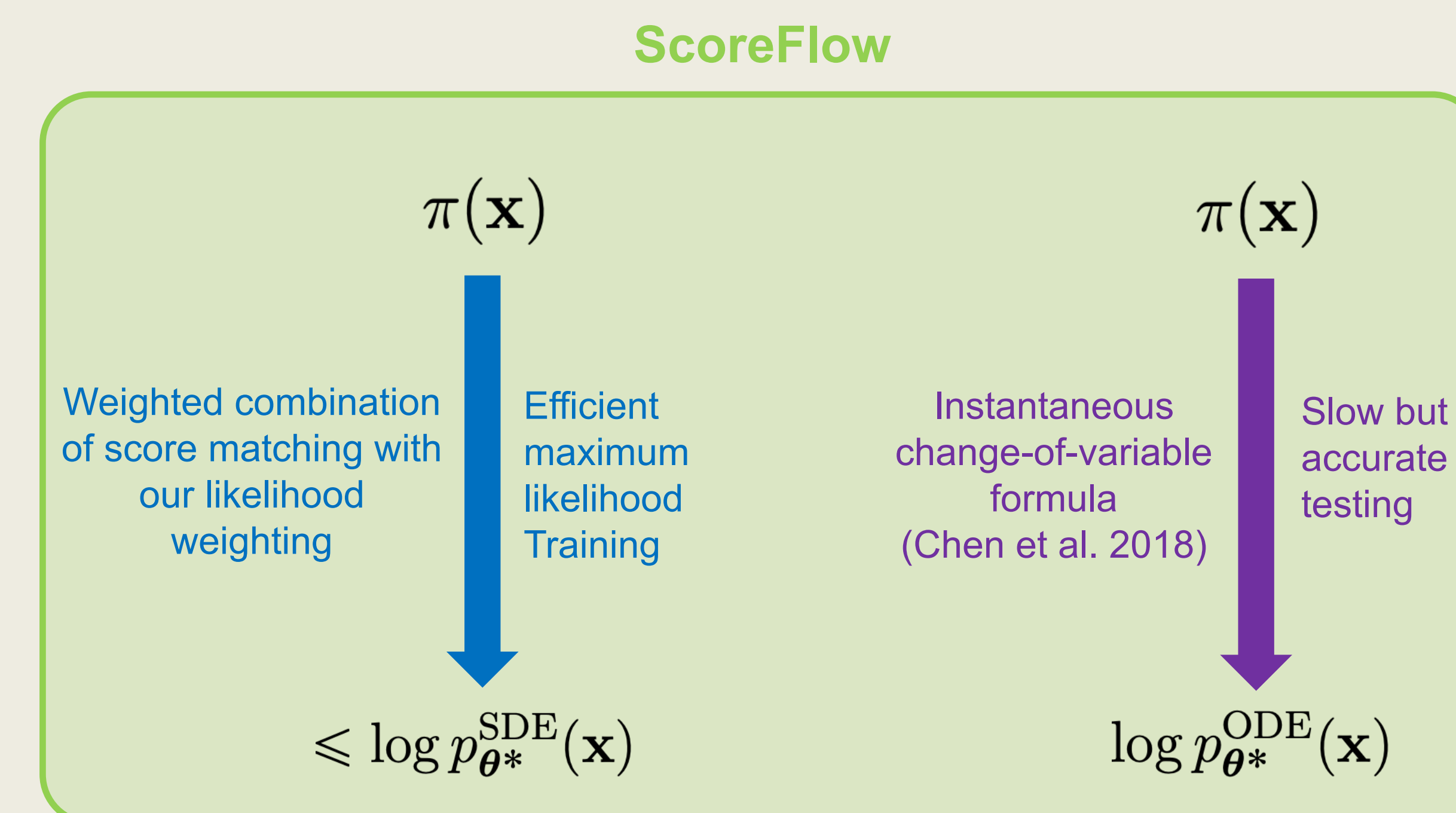
## Maximum Likelihood Training

**Theorem:** KL divergence is upper-bounded by the weighted combination of score matching.

$$\text{KL}(p_{\text{data}} \parallel p_{\theta^*}^{\text{SDE}}) \leq \frac{1}{2} \mathbb{E}_{t \sim \text{Uniform}[0, T]} [\sigma(t)^2 \mathbb{E}_{p_t(\mathbf{x})} [\|\nabla_x \log p_t(\mathbf{x}) - s_{\theta^*}(\mathbf{x}, t)\|_2^2]] + \text{KL}(p_T \parallel \pi) \approx 0$$

↑ "Likelihood weighting"

**Theorem:** There exists an efficiently computable variational lower bound to  $\log p_{\theta^*}^{\text{SDE}}(\mathbf{x})$ , analogous to the evidence lower bound of variational inference or Variational Auto-Encoders (VAEs).



**Remarks:**

- Under ideal conditions:  $\log p_{\theta^*}^{\text{SDE}}(\mathbf{x}) \approx \log p_{\theta^*}^{\text{ODE}}(\mathbf{x})$
- Importance sampling w.r.t. variable  $t$  to reduce the variance when estimating expectations.
- Variational dequantization for comparing with models trained on discrete data.

## Empirical Results

Negative log-likelihood (bits/dim) and sample quality (FID scores) on CIFAR-10 and ImageNet 32x32.

Model	SDE	CIFAR-10				ImageNet 32x32					
		Uni. deq. NLL↓	Var. deq. Bound↓	NLL↓	Bound↓	FID↓	Uni. deq. NLL↓	Var. deq. Bound↓	FID↓		
Baseline	VP	3.16	3.28	3.04	3.14	3.98	3.90	3.96	3.84	3.91	8.34
Baseline + LW	VP	3.06	3.18	2.94	3.03	5.18	3.91	3.96	3.86	3.92	17.75
Baseline + LW + IS	VP	2.95	3.08	2.83	2.94	6.03	3.86	3.92	3.80	3.88	11.15
Deep	VP	3.13	3.25	3.01	3.10	3.09	3.89	3.95	3.84	3.90	8.40
Deep + LW	VP	3.06	3.17	2.93	3.02	7.88	3.91	3.96	3.86	3.92	17.73
Deep + LW + IS	VP	2.93	3.06	2.80	2.92	5.34	3.85	3.92	3.79	3.88	11.20
Baseline	subVP	2.99	3.09	2.88	2.98	3.20	3.87	3.92	3.82	3.88	8.71
Baseline + LW	subVP	2.97	3.07	2.86	2.96	7.33	3.87	3.92	3.82	3.88	12.99
Baseline + LW + IS	subVP	2.94	3.05	2.84	2.94	5.58	3.84	3.91	3.79	3.87	10.57
Deep	subVP	2.96	3.06	2.85	2.95	2.86	3.86	3.91	3.81	3.87	8.87
Deep + LW	subVP	2.95	3.05	2.85	2.94	6.57	3.88	3.93	3.83	3.88	16.55
Deep + LW + IS	subVP	2.90	3.02	2.81	2.90	5.40	3.82	3.90	3.76	3.86	10.18

NLLs on CIFAR-10 and ImageNet 32x32 without any data augmentation

Model	CIFAR-10	ImageNet
FFJORD [15]	3.40	-
Flow++ [18]	3.08	3.86
Gated PixelCNN [35]	3.03	3.83
VFlow [4]	2.98	3.83
PixelCNN++ [40]	2.92	-
NVAE [54]	2.91	3.92
Image Transformer [36]	2.90	3.77
Very Deep VAE [8]	2.87	3.80
PixelSNAIL [7]	2.85	3.80
δ-VAE [38]	2.83	3.77
Sparse Transformer [9]	2.80	-
ScoreFlow (Ours)	2.83	3.76