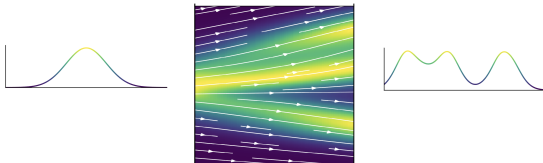


# Improving Flow Matching by Aligning Flow Divergence

Yuhao Huang, Taos Transue, Shih-Hsin Wang, William Feldman  
Hong Zhang & Bao Wang



# Introduction to Flow Matching

Learn a flow map  $\psi_t : [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$  via the following ODE

$$\frac{d}{dt}\psi_t(\mathbf{x}) = \mathbf{u}_t(\psi_t(\mathbf{x})) \approx \mathbf{v}_t(\psi_t(\mathbf{x}); \theta), \quad (1)$$

with IC  $\psi_0(\mathbf{x}) = \mathbf{x} \sim p_{\text{prior}}$ . It **maps prior** distribution  $p_0 = p_{\text{prior}}$  **to data** distribution  $p_1 \approx p_{\text{data}}$  via:

$$p_t(\mathbf{x}) = p_0(\psi_t^{-1}(\mathbf{x})) \det \left[ \frac{\partial \psi_t^{-1}}{\partial \mathbf{x}}(\mathbf{x}) \right], \quad \forall \mathbf{x} \in p_0, \forall t \in [0, 1]. \quad (2)$$

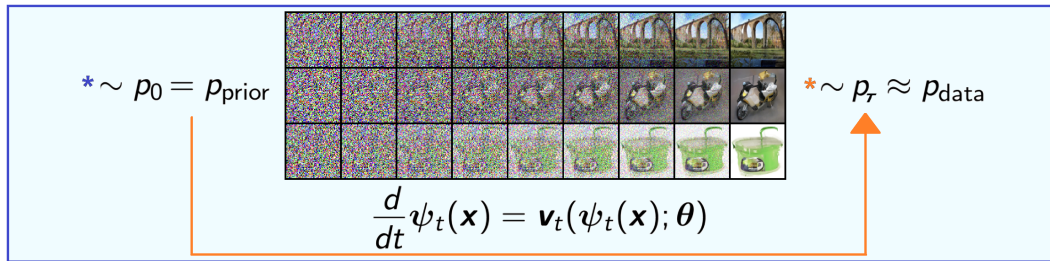


Figure: Generative Flow Map

# Introduction to Flow Matching

## Flow Matching for Generative Model

### Flow Matching Loss:

$$\mathcal{L}_{FM}(\theta) := \mathbb{E}_{t, p_t(\mathbf{x})} [\|\mathbf{v}_t(\mathbf{x}, \theta) - \mathbf{u}_t(\mathbf{x})\|^2] \quad (3)$$

where  $p_t(\mathbf{x})$  and  $\mathbf{u}_t(\mathbf{x})$  are intractable. (Lipman et al.2023) proposed **Conditional Flow Matching** loss:

$$\mathcal{L}_{CFM}(\theta) := \mathbb{E}_{t \sim \mathcal{U}[0,1], p_{\text{data}}(\mathbf{x}_1), p_t(\mathbf{x}|\mathbf{x}_1)} [\|\mathbf{v}_t(\mathbf{x}, \theta) - \mathbf{u}_t(\mathbf{x}|\mathbf{x}_1)\|^2]. \quad (4)$$

where  $p_t(\mathbf{x}|\mathbf{x}_1)$  and  $\mathbf{u}_t(\mathbf{x}|\mathbf{x}_1)$  are **pre-defined**, s.t.,

$$p_t(\mathbf{x}) = \int p_t(\mathbf{x} | \mathbf{x}_1) p_{\text{data}}(\mathbf{x}_1) d\mathbf{x}_1, \quad \mathbf{u}_t(\mathbf{x}) = \int \mathbf{u}_t(\mathbf{x}|\mathbf{x}_1) \frac{p_t(\mathbf{x} | \mathbf{x}_1) p_{\text{data}}(\mathbf{x}_1)}{p_t(\mathbf{x})} d\mathbf{x}_1 \quad (5)$$

, e.g., **Optimal transport (OT)** path:

$$p_t(\mathbf{x}|\mathbf{x}_T) = \mathcal{N}(\mathbf{x}|\mu_t(\mathbf{x}_T), \sigma_t(\mathbf{x}_T)^2 \mathbf{I}), \text{ with } \mathbf{u}_t(\mathbf{x}|\mathbf{x}_1) = \frac{\mathbf{x}_1 - (1 - \sigma_{\min})\mathbf{x}}{1 - (1 - \sigma_{\min})t}. \quad (6)$$

, where  $\mu_t(\mathbf{x}) = t\mathbf{x}_T$  and  $\sigma_t(\mathbf{x}) = 1 - (1 - \sigma_{\min})t$ .

# Relation to Continuity Equation

## Continuity Equation

From (Villani, 2009), we can know the vector field  $\mathbf{u}_t$  generates a probability path  $p_t$  satisfies the **continuity equation**:

$$\begin{aligned}\frac{\partial p_t(\mathbf{x})}{\partial t} + \nabla \cdot \left( p_t(\mathbf{x}) \mathbf{u}_t(\mathbf{x}) \right) &= 0 \\ \Leftrightarrow \frac{\partial p_t(\mathbf{x})}{\partial t} &= - \left( \nabla \cdot \mathbf{u}_t(\mathbf{x}) \right) p_t(\mathbf{x}) - \mathbf{u}_t(\mathbf{x}) \cdot \nabla p_t(\mathbf{x})\end{aligned}\tag{7}$$

with  $p_0(\mathbf{x}) = p_{\text{prior}}(\mathbf{x})$  and  $p_1(\mathbf{x}) = p_{\text{data}}(\mathbf{x})$ .

Similarly, the learned  $\mathbf{v}_t(\mathbf{x}; \theta)$  with estimated path  $\hat{p}_t(\mathbf{x})$  and  $p_0(\mathbf{x}) = p_{\text{prior}}(\mathbf{x})$  satisfy

$$\frac{\partial \hat{p}_t(\mathbf{x})}{\partial t} = - \left( \nabla \cdot \mathbf{v}_t(\mathbf{x}; \theta) \right) \hat{p}_t(\mathbf{x}) - \mathbf{v}_t(\mathbf{x}; \theta) \cdot \nabla \hat{p}_t(\mathbf{x})\tag{8}$$

**Notice:**  $|\mathbf{u} - \mathbf{v}|$  is controlled by flow matching loss but  $|\nabla \cdot \mathbf{u} - \nabla \cdot \mathbf{v}|$  not!

# Error Defined by Continuity Equation

## Error of Approximated Probability Path

Let  $\epsilon_t := p_t - \hat{p}_t$  be the error, satisfying the following transport equation

$$\begin{cases} \partial_t \epsilon_t + \nabla \cdot (\epsilon_t \mathbf{v}_t) = L_t, \\ \epsilon_0(x) = 0, \end{cases} \quad (9)$$

where  $L_t = -p_t [\nabla \cdot (\mathbf{u}_t - \mathbf{v}_t) + (\mathbf{u}_t - \mathbf{v}_t) \cdot \nabla \log p_t]$ .

**Duhamel's formula:**

$$\epsilon_t(\phi_t(\mathbf{x})) \cdot \det \nabla \phi_t(\mathbf{x}) = - \int_0^t p_s \left[ (\nabla \cdot (\mathbf{u}_s - \mathbf{v}_s)) + (\mathbf{u}_s - \mathbf{v}_s) \cdot \nabla \log p_s \right] \cdot \det \nabla \phi_s(\mathbf{x}) ds \quad (10)$$

where  $\phi(\mathbf{x})$  is the flow induced by  $\mathbf{v}_t$  in  $\frac{d}{dt} \phi_t(\mathbf{x}) = \mathbf{v}_t(\phi_t(\mathbf{x}); \boldsymbol{\theta})$ , and  $\det \nabla \phi(\mathbf{x})$  denotes the determinant of the Jacobian matrix.

# Error Defined by Continuity Equation

## TV Error from Continuity Equation

**Theorem** Under mild assumptions, there exists a constant  $C > 0$  such that

$$\text{TV}(p_t, \hat{p}_t) \leq \frac{1}{2} \mathbb{E}_{t, p_t(\mathbf{x})} \left[ \left| \nabla \cdot \mathbf{u}_t(\cdot) - \nabla \cdot \mathbf{v}_t(\cdot; \boldsymbol{\theta}) \right| \right] + \frac{C}{2} \mathbb{E}_{t, p_t(\mathbf{x})} \left[ \left| \mathbf{u}_t(\cdot) - \mathbf{v}_t(\cdot; \boldsymbol{\theta}) \right| \right] \quad (11)$$

**Note:** Second term is already flow matching loss.

# Divergence Loss

## Flow Divergence Matching Loss

The TV error bound gives the following divergence loss:

$$\mathcal{L}_{\text{DM}}(\boldsymbol{\theta}) := \mathbb{E}_{t, p_t(\mathbf{x})} \left[ \left| \nabla \cdot (\mathbf{u}_t - \mathbf{v}_t) + (\mathbf{u}_t - \mathbf{v}_t) \cdot \nabla \log p_t \right| \right] \quad (12)$$

which is also intractable.

## Conditional Flow Divergence Loss

$$\begin{aligned} \mathcal{L}_{\text{CDM}}(\boldsymbol{\theta}) := \mathbb{E}_{t, p_t(\mathbf{x}|\mathbf{x}_1), p(\mathbf{x}_1)} & \left[ \left| \left( \nabla \cdot \mathbf{u}_t(\mathbf{x}|\mathbf{x}_1) - \nabla \cdot \mathbf{v}_t(\mathbf{x}, \theta) \right) \right. \right. \\ & \left. \left. + \left( \mathbf{u}_t(\mathbf{x}|\mathbf{x}_1) - \mathbf{v}_t(\mathbf{x}, \theta) \right) \cdot \nabla \log p_t(\mathbf{x}|\mathbf{x}_1) \right| \right]. \end{aligned} \quad (13)$$

which is an upper bound of  $\mathcal{L}_{\text{DM}}(\boldsymbol{\theta})$ .

# Flow Matching with Divergence Loss

## Improve Flow Matching with Aligning Flow Divergence

We propose the flow and divergence matching (FDM) loss:

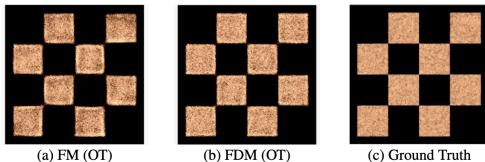
$$\mathcal{L}_{\text{FDM}} = \lambda_1 \mathcal{L}_{\text{CFM}} + \lambda_2 \mathcal{L}_{\text{CDM}}, \quad (14)$$

where  $\lambda_1, \lambda_2 > 0$  are hyperparameters.



# Density Modeling

## Synthetic Data – Checkerboard



Model	FM (OT)	FDM (OT)	FM (VP)	FDM (VP)
Likelihood ( $\uparrow$ )	$2.38 \pm .02$	<b><math>2.53 \pm .02</math></b>	$2.34 \pm .02$	$2.46 \pm .02$

*Table 1.* Likelihood estimation of models on the checkerboard test set. Here, “OT” denotes the optimal transport path and “VP” denotes the variance-preserving path. Unit:  $\times 10^{-2}$

## Image Data – CIFAR10

Model	NLL( $\downarrow$ )	FID( $\downarrow$ )
FM(OT)	2.99	6.35
FDM(OT)	2.85	5.62

*Table 2.* Negative log-likelihood and sample quality (FID scores) estimation on CIFAR-10.

# Sequential Data Sampling with Guidance – DNA Sequence

Method	MSE ( $\downarrow$ )
Bit Diffusion (One-hot Encoding)(Albergo et al., 2023)	3.95E-2
DDSM (Albergo et al., 2023)	3.34E-2
Large Language Model (Stark et al., 2024)	3.33E-2
Linear FM (Stark et al., 2024)	$2.82 \pm 0.02$ E-2
Linear FDM (ours)	$2.78 \pm 0.01$ E-2
Dirichlet FM (Stark et al., 2024)	$2.68 \pm 0.01$ E-2
Dirichlet FDM (ours)	<b><math>2.59 \pm 0.02</math>E-2</b>

Table 4. Evaluation of transcription profile guided promoter DNA sequence design of different models.

- Train the models **guided** by a profile by providing it as additional input to the vector field;
- Evaluate generated sequences using mean- squared error (MSE) between their predicted and original regulatory activity.

# Spatiotemporal Data – Dynamical System

Model	Lorenz		FitzHugh-Nagumo	
	$p(\mathbf{x}_1) (\downarrow)$	$p(\mathbf{x}_1 E) (\downarrow)$	$p(\mathbf{x}_1) (\downarrow)$	$p(\mathbf{x}_1 E) (\downarrow)$
Diffusion	0.0314	0.1001	0.0277	0.1192
FM	0.0348	0.0972	0.0314	0.2164
FDM	<b>0.0306</b>	<b>0.0914</b>	<b>0.0266</b>	<b>0.1168</b>

*Table 5.* TV distances of the models from the trajectory distribution  $p(\mathbf{x}_1)$  and from the distribution conditioned on an event  $p(\mathbf{x}_1|E)$ . Here, Diffusion results follow from (Finzi et al., 2023), while FM and FDM are based on our implementation, which builds on the code provided by Finzi et al. (2023).

Model	Lorenz		FitzHugh-Nagumo	
	$p(\mathbf{x}_1)$	$p(\mathbf{x}_1 E)$	$p(\mathbf{x}_1)$	$p(\mathbf{x}_1 E)$
Diffusion	0.0056	0.2774	0.0260	0.3011
FM	0.0081	<b>0.2560</b>	0.0280	0.3468
FDM	<b>0.0049</b>	0.3045	<b>0.0280</b>	<b>0.2084</b>

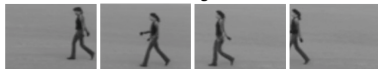
*Table 7.* KL divergence between the histograms of the event constraint value  $C(\mathbf{x}_1)$  for event trajectories  $\mathbf{x}_1$  in the dataset of trajectories computed by an ODE solver and event trajectories sampled with event guidance from the models.

Sample trajectories in dynamical systems using the initial states from the first few steps as additional inputs to the model, with or without event guidance.

# Spatiotemporal Data – Video Prediction



Walking; a-FM



Walking; a-FDM



Boxing; b-FM



Boxing; b-FDM



Hand Waving; c-FM



Hand Waving; c-FDM

Method	FVD(↓)	PSNR(↑)	Time(s/iter)
SRVP (Franceschi et al., 2020)	222	29.7	–
SLAMP (Akan et al., 2021)	228	29.4	–
Latent FM (Davtyan et al., 2023)	180	30.4	0.18
Latent FDM (ours)	<b>155.5<math>\pm</math>5</b>	<b>31.2</b>	0.27

Table 8. KTH dataset evaluation. The evaluation protocol is to predict the next 30 frames given the first 10 frames.

Method	FVD(↓)	MEM(GB)	Time(hours)
TriVD-GAN-FP (Luc et al., 2020)	103	1024	280
Video Transformer (Weissenborn et al., 2019)	94	512	336
LVT (Rakhimov et al., 2020)	126	128	48
RaMViD (Diffusion) (Höppe et al., 2022)	84	320	72
Latent FM (Davtyan et al., 2023)	146	24.2	25
Latent FDM (ours)	<b>123<math>\pm</math>4.5</b>	35	36

Table 9. BAIR dataset evaluation. We adopt the standard evaluation setup, where the model predicts 15 future frames conditioned on a single initial frame. MEM stands for peak memory footprint.

Autoregressive next-frame generation (prediction) guided by several preceding frames, provided as additional inputs to the flow matching vector regressor.

**Thank you!**