Generative Flow Maps: An overview of the math and methods behind them.
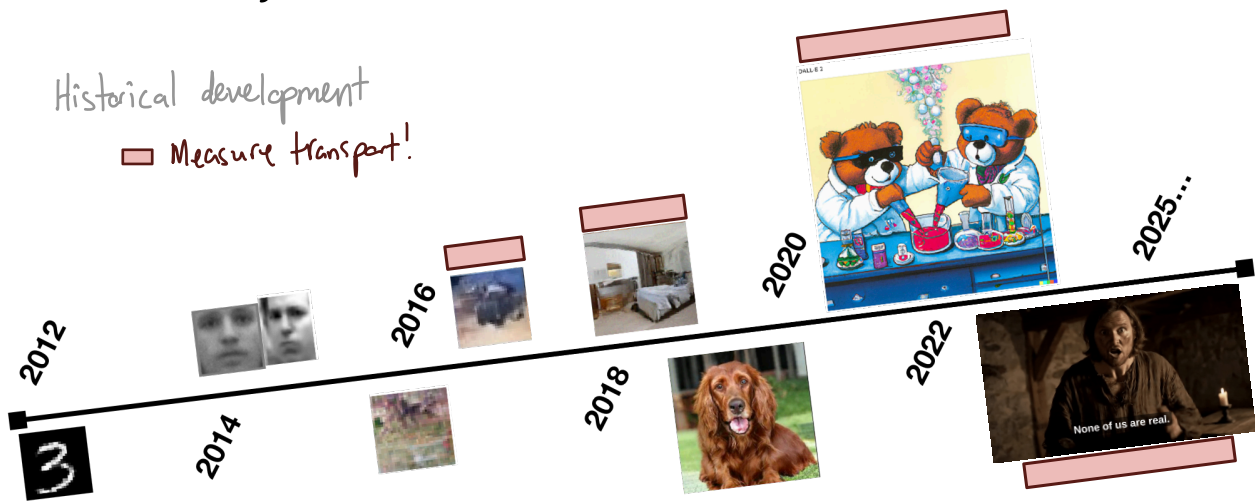
arXiv: 2406.07507 and 2505.18825

Agenda for this talk:

- Introduce dynamical measure transport for generative modeling
- Motivate the flow map as a computationally efficient method
- Illustrate how equations governing this map can be used to learn it, and categorize the recent efforts made in this direction

## Generative Modeling

Goal: Estimate some unknown distribution with density $\rho_1$ through sample data $x_1 \sim \rho_1$.



Historical development
☐ Measure transport!

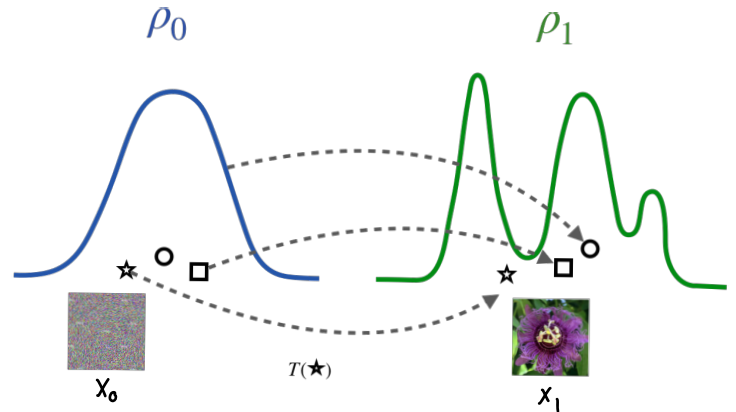2012  2014  2016  2018  2020  2022  2025...

What do we mean by measure transport, and how can we can adapt the equations governing it to create more understandable and performant tools?
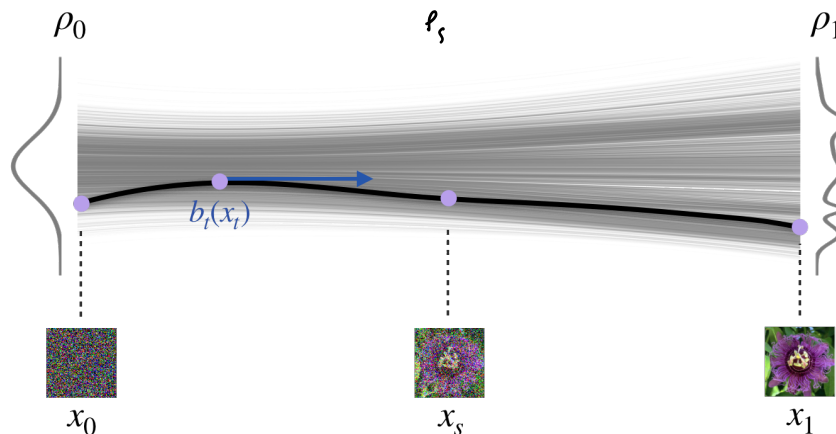
## Measure Transport

Building maps between distributions

- Sample base distribution $x_0 \sim \rho_0$

- Build a <u>map</u> $T : \Omega \to \Omega$

- Produce $x_1 \sim \rho_1$ via $T(x_0) = x_1$



$\rho_0$     $\rho_1$

$T(\star)$

$x_0$     $x_1$

## Dynamical Measure Transport

This map can be constructed as the solution to a dynamical equation. Imagine that $x_0$ continually evolves over time $t \in [0, 1]$ to some $x_1$.



$\rho_0$     $\rho_s$     $\rho_1$

$b_t(x_t)$

$x_0$     $x_s$     $x_1$

**Probability flow ODE**

(1)   $\dot{x}_t = b_t(x_t), \quad x_0 \sim \rho_{t=0}$

**Continuity equation**

(2)   $\partial_t \rho_t + \nabla \cdot (b_t \rho_t) = 0, \quad \rho_{t=0} = \rho_0$

- $b_t$ is a velocity field which defines how $x_t$ should instantaneously evolve

- Equation governing the evolution of $\rho_t$ with $b_t$

## Learning $b_t$ via flow matching/stochastic interpolants

- To construct a $\rho_{t_1}$ stochastically combine $x_0, x_1$ via the interpolant:

$$(3) \quad I_t(x_0, x_1) = \alpha_t x_0 + \beta_t x_1 \quad , \quad (x_0, x_1) \sim p(x_0, x_1) \qquad eg \quad \begin{array}{l} \alpha_t = 1-t \\ \beta_t = t \end{array}$$
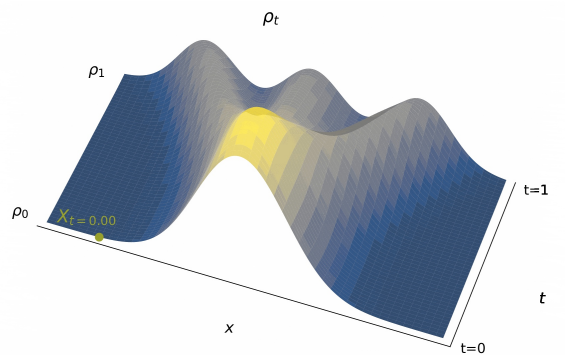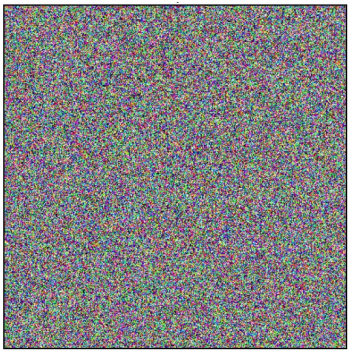
then $\rho_t = \text{Law}(I_t)$ and $b_t$ associated to (1), (2) is given by

$$(4) \qquad b_t = \mathbb{E}\left[\dot{I}_t \,\middle|\, I_t = x\right] \qquad \text{Expectation over } p(x_0, x_1) \text{ conditional on } I_t = x.$$

- $b_t$ can be learned over neural networks by minimizing

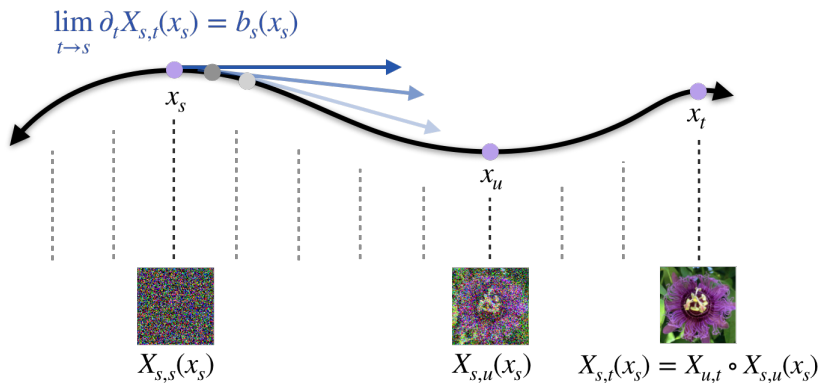$$(5) \quad L_b[\hat{b}] = \int_0^1 \mathbb{E}_{x_0, x_1}\left[|\hat{b}_t(I_t) - \dot{I}_t|^2\right] dt$$

Then use $b_t$ coming from (5) to generate samples by numerically solving (1)



$\rho_t$

$\rho_1$

$\rho_0$    $X_{t=0.00}$

$t=1$

$t$

$x$    $t=0$

**Powerful! But limitation:** Sampling requires many evaluations of $b_t$ to solve (1). How can we avoid this?

# The flow map

Instead of solving (1), we may be interested in learning an arbitrary integrator for the equation in terms of a flow map:

$$\lim_{t \to s} \partial_t X_{s,t}(x_s) = b_s(x_s)$$



$X_{s,s}(x_s)$  $\qquad$ $X_{s,u}(x_s)$  $\qquad$ $X_{s,t}(x_s) = X_{u,t} \circ X_{s,u}(x_s)$

$$X_{s,t}(x_s) = x_t \qquad (6)$$

"Takes steps of arbitrary size $t-s$ along trajectories of the probability flow"

## Properties of the flow map

Semigroup property :
$$X_{u,t}\left(X_{s,u}(x_s)\right) = X_{s,t}(x_s) = x_t \qquad (7)$$

$\hookrightarrow$ Invertibility $\quad X_{s,t}(X_{t,s}(x_t)) = x_t \qquad (8)$

Lagrangian eqn : $\partial_t X_{s,t}(x_s) = \dot{x}_t = \underset{(1)}{b_t(x_t)} = b_t(X_{s,t}(x))$

$\hookrightarrow$ $$\partial_t X_{s,t}(x) = b_t\left(X_{s,t}(x)\right) \qquad (9)$$
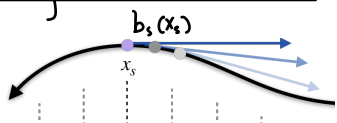
<u>Eulerian eqn</u> : Take a total derivative of (8)

$$\frac{d}{ds} X_{s,t}(X_{t,s}(x)) = \frac{\partial}{\partial s} X_{s,t}(X_{t,s}(x)) + \nabla X_{s,t}(X_{t,s}(x)) \cdot \underbrace{\frac{\partial}{\partial s} X_{t,s}(x)}_{b_s(X_{t,s}(x))} = 0$$

Evaluate it at $X_{t,s} = y$, so that

$$\frac{\partial}{\partial s} X_{s,t}(x) + \nabla X_{s,t}(x) \cdot b_s(x) = 0 \qquad (10)$$

<u>Tangent Condition</u> :



$$\lim_{s \to t} \partial_t X_{s,t}(x) = b_t(x) \qquad (11)$$

## Parameterizing the Flow map

Choose $\hat{X}_{s,t}(x) = x + (t-s)\hat{v}_{s,t}(x)$ (12), then, using (11)  $\boxed{v_{t,t}(x) = b_t(x)}$ (13)

*will be learned as an NN*

---

Proposition: If $X_{s,t}$ is given by (12) and $v_{s,t}$ satisfies (13),

**A: Lagrangian**
$$\partial_t X_{s,t}(x) = v_{t,t}(X_{s,t}(x))$$

**B: Eulerian**
$$\frac{\partial}{\partial s} X_{s,t}(x) + \nabla X_{s,t}(x) \cdot b_s(x) = 0$$

**C: consistency**
$$X_{u,t}(X_{s,u}(x_s)) = X_{s,t}(x_s)$$

each characterize the flow map !

---

Let's use them, along with (13), to learn $X_{s,t}$ directly!

Objective Function:

Learn $v_{t,t} = b_t$ on the diagonal using (5)

Self-distillation

Learn the flow map on the off-diagonal w/ A, B, C or any combination

$$L_{SD}(\hat{v}) = L_b(\hat{v}) + L_D(\hat{v}) \qquad (14)$$

Example: Lagrangian Self-Distillation

- $L_b(\hat{v}) = \int_0^1 \mathbb{E}_{X_0, X_1}\left[| \hat{v}_{t,t}(I_t) - \dot{I}_t |^2\right] dt$  ·  Rewriting of (5) w/ $v_{t,t}$

- $L_D^{LSD}(\hat{v}) = \int_0^1 \int_0^t \mathbb{E}_{X_0, X_1}\left[| \partial_t \hat{X}_{s,t}(I_s) - \hat{v}_{t,t}(\hat{X}_{s,t}(I_s)) |^2\right] ds\, dt$  ·  PINN enforcing A

Others

- $L^{ESD}(\hat{v}) = \int_0^1 \int_0^t \mathbb{E}_{X_0, X_1}\left[| \partial_s \hat{X}_{s,t}(I_s) + \nabla \hat{X}_{s,t}(I_s) \hat{v}_{s,s}(I_s) |^2\right] ds\, dt$  ·  PINN enforcing B

- $L^{PSD}(\hat{v}) = \int_0^1 \int_0^t \int_s^t \mathbb{E}_{X_0, X_1}\left[| \hat{X}_{s,t}(I_s) - \hat{X}_{u,t}(\hat{X}_{s,u}(I_s)) |^2\right] du\, ds\, dt$  ·  Enforcing C
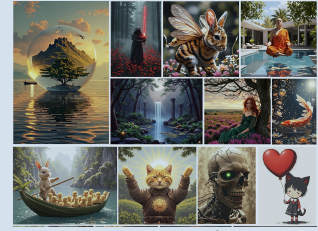
## Connection w/ the literature

Distill from a known velocity field $b_t(x)$:

$$\mathbb{E} \; \frac{1}{2} \left| \partial_S X_{S,t}(I_S) + \underbrace{b_S(I_S) \cdot \nabla X_{S,t}(I_S)}_{} \right|^2$$

stop gradient on these terms



"Align your flow"
Saban, Fidler, Kreis

Thm 3.2/3.3

Take $\nabla_\theta$ gradient of objective:

$$\mathbb{E} \; \nabla_\theta \frac{1}{2} \left[ \partial_S X_{S,t}(I_S) + \text{stopgrad}\left( b_S(I_S) \cdot \nabla X_{S,t}(I_S) \right) \right]^2$$

$$= \mathbb{E} \; \left[ \nabla_\theta \partial_S X_{S,t}(I_S) \right] \left( \partial_S X_S(I_S) + b_S(I_S) \cdot \nabla_x X_{S,t}(I_S) \right)$$
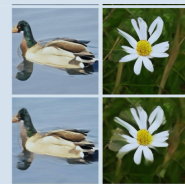
Eulerian Map Distillation with stopgrad

Thm 3.2 in AYF Paper

----

## Shortcut Models is Progressive Self Distillation

$$L^{PSD}(\hat{v}) = \int_0^1 \int_0^t \int_S^t \mathbb{E}_{x_0, x_1} \left[ \left| \hat{X}_{S,t}(I_S) - \hat{X}_{u,t}\left( \hat{X}_{S,u}(I_S) \right) \right|^2 \right] du\,ds\,dt$$



"Shortcut Models"
Frans et al

- Take the Eulerian Self-Distillation term again

$$\min_{v} \int_{[0,1]^2} \iint \frac{1}{2} \left| \underbrace{\partial_s X_{s,t}(I_s) + V_{ss}(I_s) \cdot \nabla X_{s,t}(I_s)}_{\LARGE \star} \right|^2 ds\, dt$$

- Plug in $X_{s,t}(x) = x + (t-s) V_{s,t}(x)$, which means that some terms simplify:

$$\partial_s X_{s,t}(x) = - V_{s,t}(x) + (t-s) \partial_s V_{s,t}(x)$$

- Plug this into $\star$ :

$$\iint \frac{1}{2} \left| -V_{s,t}(I_s) + \underbrace{(t-s)\partial_s V_{s,t}(I_s) + V_{ss}(I_s) \cdot \nabla X_{s,t}(I_s)}_{\substack{\text{stopgrad this} \\ \text{whole term}}} \right|^2$$

- Then the gradient reads

$$\iint - \nabla_\theta V_{s,t}(I_s) \cdot \left[ -V_{s,t}(I_s) + (t-s)\partial_s V_{s,t}(I_s) + V_{ss}(I_s) \cdot \nabla X_{s,t}(I_s) \right]$$

- Expand the last gradient $\nabla X_{s,t}$ using the definition of $X_{s,t}$ :

$$b_t = \iint [I_t | I_t]$$

$$= \iint - \nabla_\theta V_{s,t}(I_s) \cdot \left[ -V_{s,t}(I_s) + (t-s)\partial_s V_{s,t}(I_s) + V_{ss}(I_s) + (t-s)V_{ss}(I_s) \cdot \nabla V_{s,t}(I_s) \right]$$

- Now, because this is linear in $V_{ss}(I_s)$, it can be replaced with $I_s$ here

$$= \iint - \nabla_\theta V_{s,t}(I_s) \cdot \left[ -V_{s,t}(I_s) + (t-s)\partial_s V_{s,t}(I_s) + V_{ss}(I_s) + (t-s) I_s \cdot \nabla V_{s,t}(I_s) \right]$$

- And this means finally that terms in blue can be collected as $\frac{d}{ds} V_{s,t}(I_s)$:

$$= \frac{d}{dt} - \nabla_\theta V_{s,t}(I_s) \cdot \left[ -V_{s,t}(I_s) + (t-s) \frac{d}{ds} V_{s,t}(I_s) \right]$$

Directly learning the flow map solely in terms of $V_{s,t}$ !

## Why? (4)

Take derivative of $f_t(x) = \int_{\mathbb{R}^d \times \mathbb{R}^d} \delta(x - I_t)\, p(x_0, x_1)\, dx_0 dx_1$

$$\partial_t f_t(x) = -\nabla \cdot \int \dot{I}_t\, \delta(x - I_t)\, p(x_0, x_1)\, dx_0 dx_1$$

$$= -\nabla \cdot j_t = -\nabla \cdot (b_t\, f_t)$$

$\hookleftarrow$ current $j_t$

So $b_t = \dfrac{\int \dot{I}_t\, \delta(x - I_t)\, p(x_0, x_1)\, dx_0 dx_1}{\int \delta(x - I_t)\, p(x_0, x_1)\, dx_0 dx_1} = \mathbb{E}\left[\dot{I}_t \mid I_t = x\right]$