

Discretize-Optimize Methods for Neural ODEs in Continuous Normalizing Flows

SIAM MDS 2020

Derek Onken

Department of Computer Science

Emory University

donken@emory.edu

derekonken.com



EMORY
UNIVERSITY

Collaborators and Acknowledgments



Samy Wu Fung
UCLA



Xingjian Li
Emory



Lars Ruthotto
Emory

Emory Funding:  DMS 1751636  US-Israel BSF 2018209

UCLA Funding: AFOSR MURI FA9550-18-1-0502 and FA9550-18-1-0167,
ONR N00014-18-1-2527

Special thanks: Organizers and staff of IPAM Long Program MLP 2019.

Overview

- **Continuous ResNet and Neural ODEs**
 - ▶ Discrete neural networks viewed in continuous framework
- **Continuous Normalizing Flows (CNFs)**
 - ▶ Discrete normalizing flows similarly moved to the continuous framework
- **Discretize-Optimize vs Optimize-Discretize**
 - ▶ Comparing approaches in solving CNFs
- **OT-Flow**
 - ▶ Incorporating optimal transport with Discretize-Optimize for fast and accurate CNFs

Motivation: Existing CNFs approaches are prohibitively slow and expensive.



DO and L Ruthotto
*Discretize-Optimize vs. Optimize-Discretize
for Time-Series Regression and CNFs*
arXiv:2005.13420, 2020.



DO, S Wu Fung, X Li, L Ruthotto
*OT-Flow: Fast and Accurate CNFs
via OT*
arXiv:2006.00104, 2020.

Neural ODE

A *neural ODE* is an ordinary differential equation (ODE) with neural network components.

For input $\mathbf{x} \in \mathbb{R}^d$, neural network $f: \mathbb{R}^d \rightarrow \mathbb{R}^d$ models the solution to

$$\partial_t \mathbf{z}(\mathbf{x}, t) = \mathbf{v}(\mathbf{z}(\mathbf{x}, t), t; \boldsymbol{\theta}), \quad \mathbf{z}(\mathbf{x}, 0) = \mathbf{x} \quad (1)$$

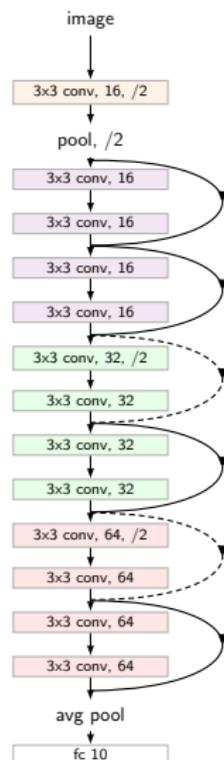
where

- time $t \in [0, T]$
- $\mathbf{v}: \mathbb{R}^d \times [0, T] \rightarrow \mathbb{R}^d$ is a neural network layer with parameters $\boldsymbol{\theta}$
- $\mathbf{z}(\mathbf{x}, t)$ are the features for initial \mathbf{x} at time t
- $f(\mathbf{x}) = \mathbf{z}(\mathbf{x}, T)$

Background

Historically

- Residual Neural Networks (ResNets)¹ perform well on image classification and more.
- ResNets are merely Forward Euler of some Continuous ResNet (1).^{2,3}

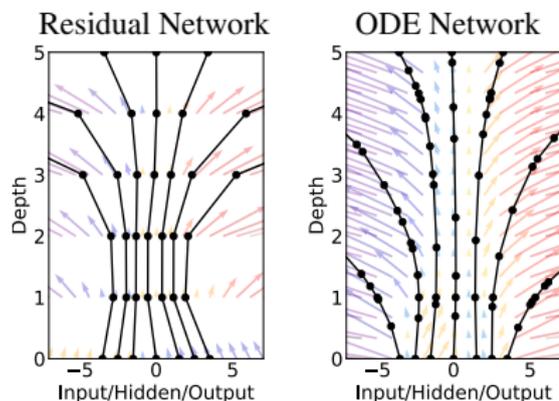


¹He et al. “Deep residual learning for image recognition”. 2016.

²E. “A Proposal on Machine Learning via Dynamical Systems”. 2017.

³Haber and Ruthotto. “Stable Architectures for Deep Neural Networks”. 2017.

Background (cont.)



Neural ODE⁴

Popularized

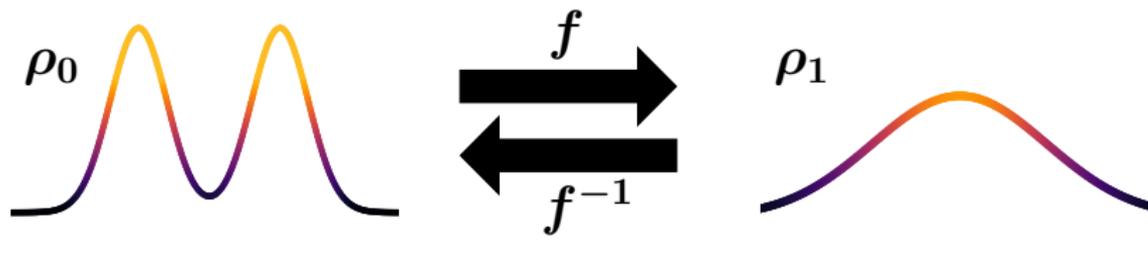
- Incorporate a black-box solver and coin the term *neural ODE*.⁴
- Applied to normalizing flows.⁵

⁴Chen et al. “Neural Ordinary Differential Equations”. 2018.

⁵Grathwohl et al. “FFJORD: Free-form continuous dynamics for scalable reversible generative models”. 2019.

Discrete Normalizing Flows

A normalizing flow^{6,7} is an invertible mapping $f: \mathbb{R}^d \rightarrow \mathbb{R}^d$ between an arbitrary probability distribution and a standard normal distribution whose densities we denote by ρ_0 and ρ_1 , respectively.



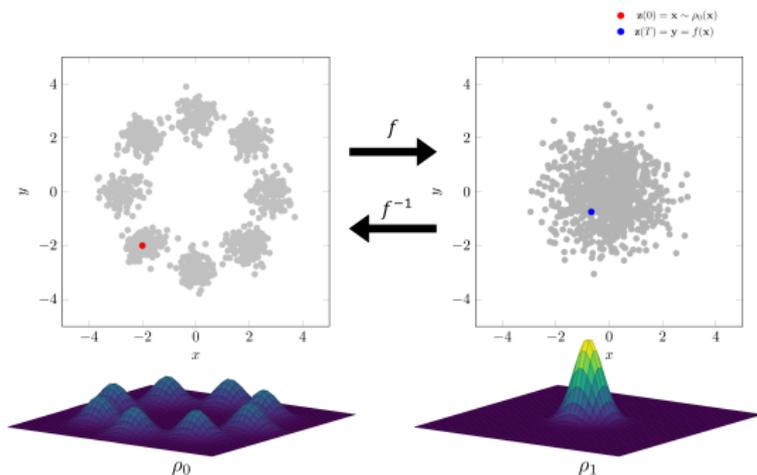
By the change of variables formula, the flow satisfies

$$\log \rho_0(\mathbf{x}) = \log \rho_1(f(\mathbf{x})) + \log |\det \nabla f(\mathbf{x})| \quad \text{for all } \mathbf{x} \in \mathbb{R}^d. \quad (2)$$

⁶Rezende and Mohamed. "Variational Inference with Normalizing Flows". 2015.

⁷Papamakarios et al. "Normalizing Flows for Probabilistic Modeling and Inference". 2019.

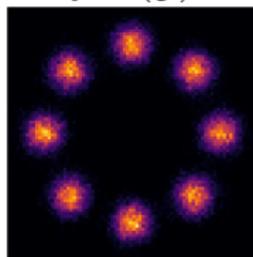
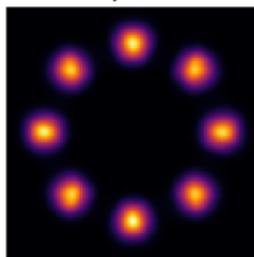
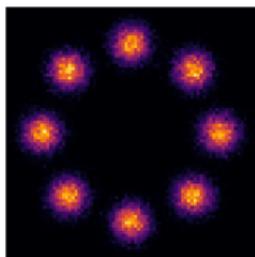
Gaussian Mixture Toy Example



Data
 x

Estimate
 ρ_0

Generation
 $f^{-1}(y)$



Continuous Normalizing Flows (CNFs)

Replace the log-det with a trace

Issue:

- log-determinants cost $\mathcal{O}(d^3)$ FLOPS in general.

Solutions:

- Use specific neural network architectures for \mathbf{v} so the log-det computation is manageable.
- Replace the log-det with a trace computation.

Using the neural ODE f (1) and Jacobi's formula⁸, we can rewrite (2) as

$$\ell(\mathbf{x}, T) := \log \rho_0(\mathbf{x}) - \log \rho_1(f(\mathbf{x})) = \int_0^T \text{tr}(\nabla_{\mathbf{v}}(\mathbf{z}(\mathbf{x}, t), t; \boldsymbol{\theta})) dt. \quad (3)$$

⁸Chen et al. 2018.

CNF Optimization Problem

For expected negative log-likelihood ^{9,10}

$$C(\mathbf{x}, T) := \frac{1}{2} \|\mathbf{z}(\mathbf{x}, T)\|^2 - \ell(\mathbf{x}, T) + \frac{d}{2} \log(2\pi),$$

we optimize

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{\rho_0(\mathbf{x})} C(\mathbf{x}, T)$$

where for a given $\boldsymbol{\theta}$, the trajectory \mathbf{z} satisfies the CNF ¹¹

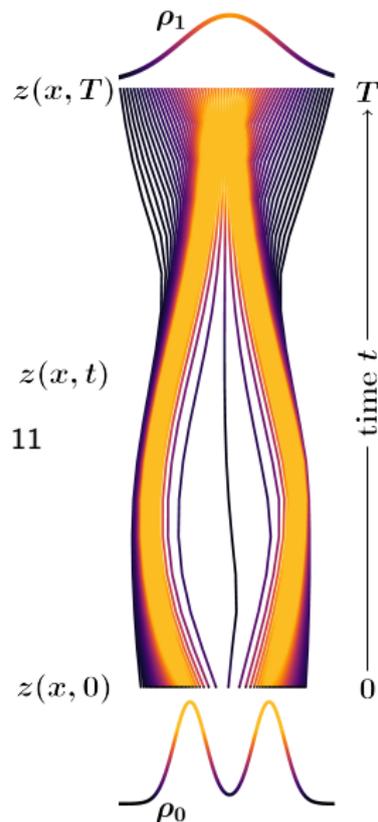
$$\partial_t \begin{bmatrix} \mathbf{z}(\mathbf{x}, t) \\ \ell(\mathbf{x}, t) \end{bmatrix} = \begin{bmatrix} \mathbf{v}(\mathbf{z}(\mathbf{x}, t), t; \boldsymbol{\theta}) \\ \text{tr}(\nabla \mathbf{v}(\mathbf{z}(\mathbf{x}, t), t; \boldsymbol{\theta})) \end{bmatrix},$$
$$\begin{bmatrix} \mathbf{z}(\mathbf{x}, 0) \\ \ell(\mathbf{x}, 0) \end{bmatrix} = \begin{bmatrix} \mathbf{x} \\ 0 \end{bmatrix}.$$

In optimal control, \mathbf{z} is the state and $\boldsymbol{\theta}$ is the control.

⁹Rezende and Mohamed. 2015.

¹⁰Papamakarios et al. 2019.

¹¹Grathwohl et al. 2019.



Solving ODE-constrained Optimization Problems

Two Predominant Approaches:

- Discretize-Optimize (DO)
 - ▶ Discretize the ODE, then optimize on that discretization.
 - ▶ Typical machine learning approach: set up architecture with N layers, the optimize on that discretization (propagate forward, calculate loss, backpropagate)
 - ▶ ANODE¹²
- Optimize-Discretize (OD)
 - ▶ Optimize in the continuous space, then discretize.
 - ▶ Use the Karush-Kuhn-Tucker (KKT) conditions or the adjoint equations to optimize, then choose a discretization.
 - ▶ Neural ODEs paper¹³ and FFJORD¹⁴

¹²Gholaminejad, Keutzer, and Biros. “ANODE: Unconditionally Accurate Memory-Efficient Gradients for Neural ODEs”. 2019.

¹³Chen et al. 2018.

¹⁴Grathwohl et al. 2019.

Solving ODE-constrained Optimization Problems

Popular methods use Optimize-Discretize.

- Adaptive solver `dopri5` for the forward propagation
- Adjoint-based backpropagation recomputes the intermediate gradients
- **Drawback:** inaccurate gradients when the adjoint equation is not solved well enough.^{15,16}

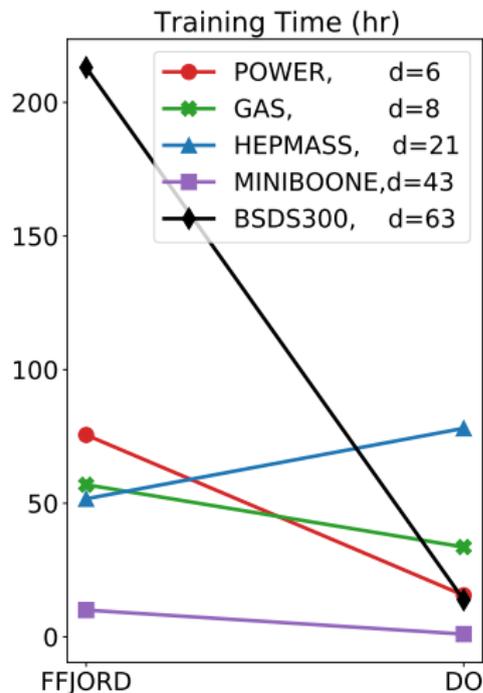
We choose Discretize-Optimize.

- Same discretization for the forward and backpropagation.
 - ▶ Use automatic differentiation (AD) for the backpropagation.
 - ▶ The gradients are accurate.
- We use Runge-Kutta 4 with a fixed step size.
- **Drawback:** have to tune a sufficiently small step size for the solver

¹⁵Li et al. "Maximum principle based algorithms for deep learning". 2017.

¹⁶Gholaminejad, Keutzer, and Biros. 2019.

Results



For all five data sets, DO and FFJORD¹⁷ (OD) achieve similar results with different training time.

DO has an average speed-up of 6.4x even with slower training on HEPMASS.

Reasons

- Fewer function evaluations (RK4 instead of dopri5).
- Intermediate gradients are stored (with AD) rather than recomputed.
- DO has more accurate gradients.



DO and L Ruthotto

Discretize-Optimize vs. Optimize-Discretize for Time-Series Regression and CNFs
arXiv:2005.13420, 2020.

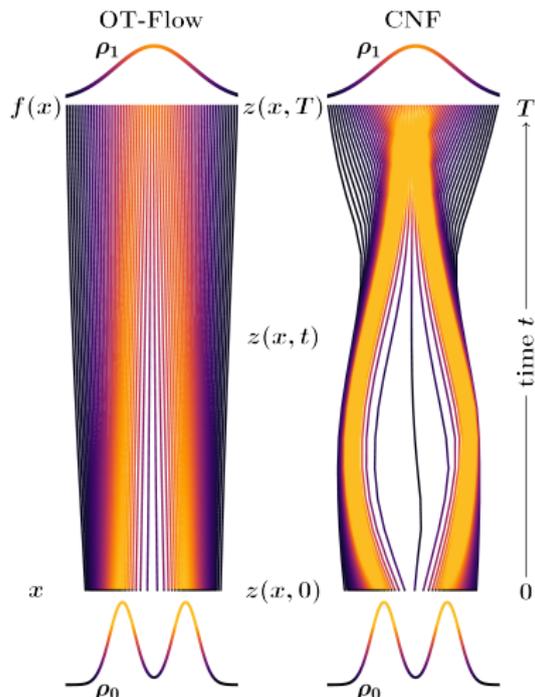
¹⁷Grathwohl et al. 2019.

Can we solve even faster?

$$\partial_t \begin{bmatrix} \mathbf{z}(\mathbf{x}, t) \\ \ell(\mathbf{x}, t) \end{bmatrix} = \begin{bmatrix} \mathbf{v}(\mathbf{z}(\mathbf{x}, t), t; \boldsymbol{\theta}) \\ \text{tr}(\nabla \mathbf{v}(\mathbf{z}(\mathbf{x}, t), t; \boldsymbol{\theta})) \end{bmatrix}, \quad \begin{bmatrix} \mathbf{z}(\mathbf{x}, 0) \\ \ell(\mathbf{x}, 0) \end{bmatrix} = \begin{bmatrix} \mathbf{x} \\ 0 \end{bmatrix}$$

What makes CNFs slow?

- Trajectories can be complicated leading to high number of function evaluations.
- Trace computation costs $\mathcal{O}(d^2)$ FLOPS in general.
 - ▶ FFJORD uses Hutchinson's estimator for $\mathcal{O}(d)$ FLOPS in training.



Straight Trajectories

Include some optimal transport (OT)

In OT, a unique mapping exists.

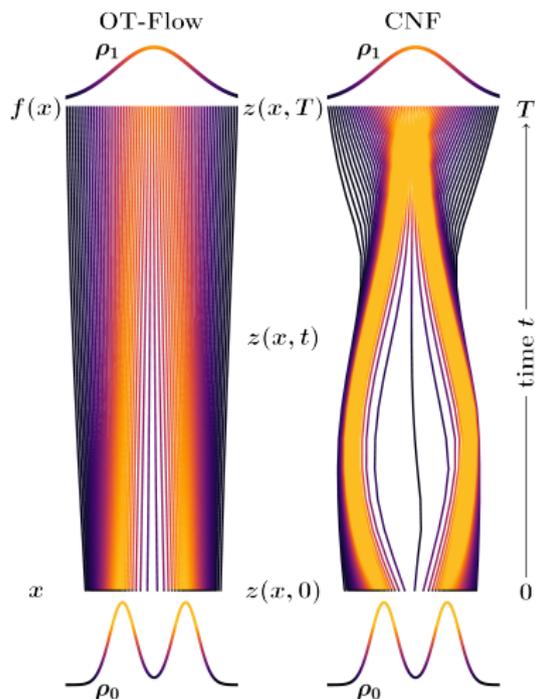
We regularize the optimization problem

$$\min_{\theta} \mathbb{E}_{\rho_0(\mathbf{x})} \left\{ C(\mathbf{x}, T) + L(\mathbf{x}, T) \right\} \quad (4)$$

subject to (1).

The L_2 transport costs are given by

$$L(\mathbf{x}, T) = \int_0^T \frac{1}{2} \|\mathbf{v}(\mathbf{z}(\mathbf{x}, t), t; \theta)\|^2 dt.$$



More OT

Potential Function Φ

Apply the Pontryagin maximum principle¹⁸
to (4)

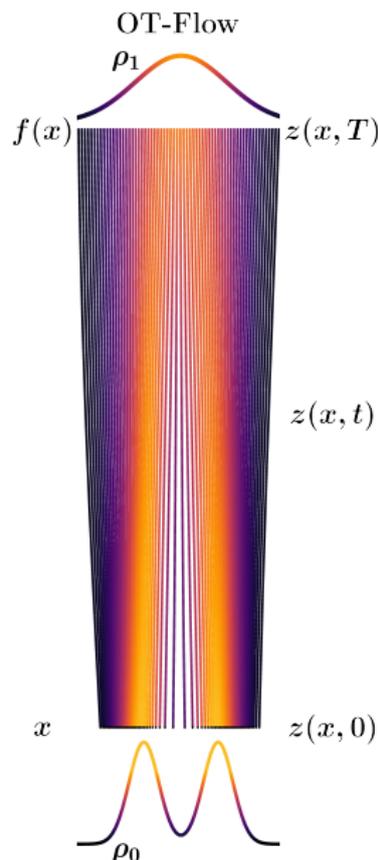
There exists a scalar potential function

$\Phi: \mathbb{R}^d \times [0, T] \rightarrow \mathbb{R}$ such that

$$\mathbf{v}(\mathbf{x}, t; \boldsymbol{\theta}) = -\nabla \Phi(\mathbf{x}, t; \boldsymbol{\theta}). \quad (5)$$

Analogous to classical physics, samples move in a manner to minimize their potential.

We parametrize potential Φ instead of \mathbf{v} .



¹⁸Evans. *An Introduction to Mathematical Optimal Control Theory Version 0.2*. 2013.

More OT

HJB Equation

The optimality conditions of (4) lead to another regularizer.

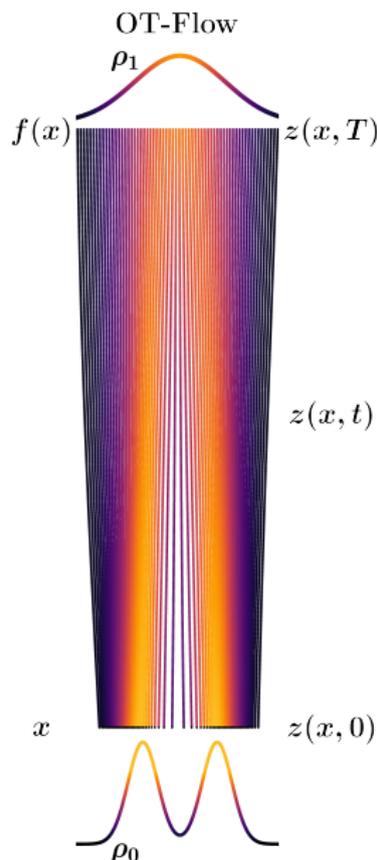
Potential Φ satisfies the Hamilton-Jacobi-Bellman (HJB) equation¹⁹

$$-\partial_t \Phi(\mathbf{x}, t) + \frac{1}{2} \|\nabla \Phi(\mathbf{z}(\mathbf{x}, t), t)\|^2 = 0, \quad (6)$$
$$\Phi(\mathbf{x}, T) = G(\mathbf{x})$$

where

$$G(\mathbf{z}(\mathbf{x}, T)) = 1 + \log(\rho_0(\mathbf{x})) - \log(\rho_1(\mathbf{z}(\mathbf{x}, T))) - \ell(\mathbf{x}, T) \quad (7)$$

Terminal condition G derives from the variational derivative or the Kullback-Leibler (KL) divergence.



¹⁹Evans. "Partial differential equations and Monge-Kantorovich mass transfer". 1997.

More OT

HJB regularizer R

Penalize deviations from the HJB equation

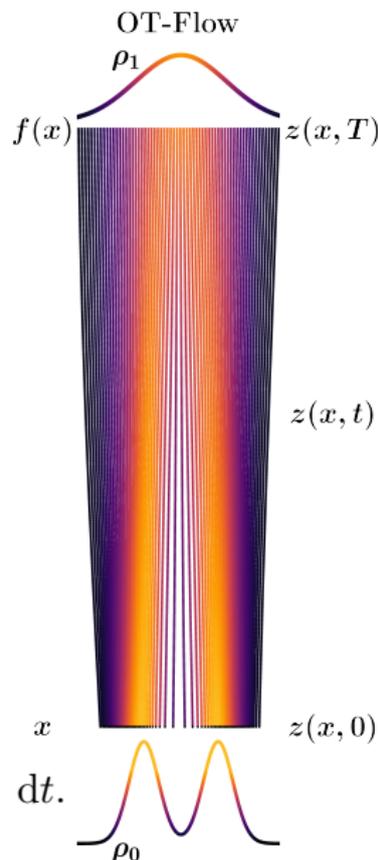
We add another regularizer, so the optimization problem is

$$\min_{\theta} \mathbb{E}_{\rho_0(\mathbf{x})} \left\{ C(\mathbf{x}, T) + L(\mathbf{x}, T) + R(\mathbf{x}, T) \right\}$$

subject to (1).

The HJB regularizer is computed as

$$R(\mathbf{x}, T) = \int_0^T \left| \partial_t \Phi(\mathbf{z}(\mathbf{x}, t), t) - \frac{1}{2} \|\nabla \Phi(\mathbf{z}(\mathbf{x}, t), t)\|^2 \right| dt.$$



OT-Flow Formulation

We incorporate the time integration in the ODE solver.

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{\rho_0(\mathbf{x})} \left\{ C(\mathbf{x}, T) + L(\mathbf{x}, T) + R(\mathbf{x}, T) \right\}$$

subject to

$$\partial_t \begin{pmatrix} \mathbf{z}(\mathbf{x}, t) \\ \ell(\mathbf{x}, t) \\ L(\mathbf{x}, t) \\ R(\mathbf{x}, t) \end{pmatrix} = \begin{pmatrix} -\nabla\Phi(\mathbf{z}(\mathbf{x}, t), t; \boldsymbol{\theta}) \\ -\text{tr}(\nabla^2\Phi(\mathbf{z}(\mathbf{x}, t), t; \boldsymbol{\theta})) \\ \frac{1}{2}\|\nabla\Phi(\mathbf{z}(\mathbf{x}, t), t; \boldsymbol{\theta})\|^2 \\ \left| \partial_t\Phi(\mathbf{z}(\mathbf{x}, t), t; \boldsymbol{\theta}) - \frac{1}{2}\|\nabla\Phi(\mathbf{z}(\mathbf{x}, t), t; \boldsymbol{\theta})\|^2 \right| \end{pmatrix}$$

with initial conditions

$$\mathbf{z}(\mathbf{x}, 0) = \mathbf{x} \quad \text{and} \quad \ell(\mathbf{x}, 0) = L(\mathbf{x}, 0) = R(\mathbf{x}, 0) = 0$$

Other OT approaches in CNFs

Table: Comparison of flow formulations.

Model	ODE (1)	Φ	L_2 cost	HJB reg.	$\ \nabla \mathbf{v}\ _F^2$
FFJORD ²⁰	✓	✗	✗	✗	✗
RNODE ²¹	✓	✗	✓	✗	✓
Monge-Ampère Flows ²²	✓	✓	✗	✗	✗
Potential Flow Gen. ²³	✓	✓	✗	✓	✗
OT-Flow	✓	✓	✓	✓	✗

²⁰Grathwohl et al. 2019.

²¹Finlay et al. "How to train your neural ODE". 2020.

²²Zhang, E, and Wang. "Monge-Ampère Flow for Generative Modeling". 2018.

²³Yang and Karniadakis. "Potential Flow Generator with L_2 Optimal Transport Regularity for Generative Models". 2019.

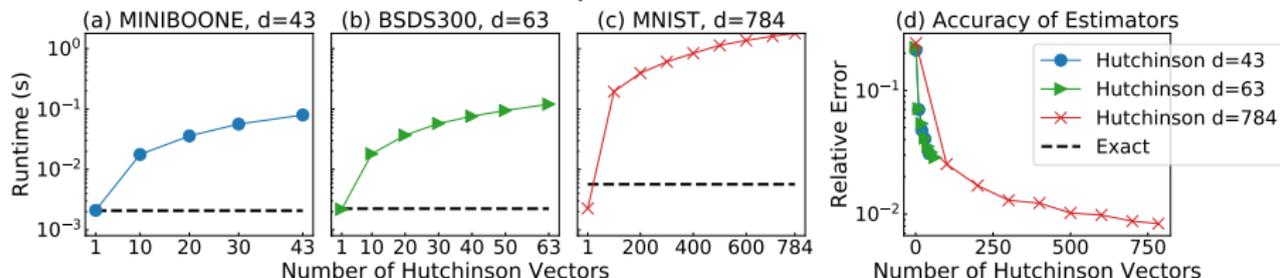
Improving the Trace Computation

General Trace Computation: $\mathcal{O}(d^2)$ FLOPS

Trace Estimators used in state-of-the-art: $\mathcal{O}(d)$ FLOPS

Our Exact Trace in OT-Flow $\mathcal{O}(d)$ FLOPS

In runtime, our exact trace is competitive with the estimators



Exact Trace Computation

Our model

Neural Network

$$\Phi(\mathbf{s}; \boldsymbol{\theta}) = \mathbf{w}^\top N(\mathbf{s}; \boldsymbol{\theta}_N) + \frac{1}{2} \mathbf{s}^\top (\mathbf{A}^\top \mathbf{A}) \mathbf{s} + \mathbf{b}^\top \mathbf{s} + c, \quad (8)$$

where $\boldsymbol{\theta} = (\mathbf{w}, \boldsymbol{\theta}_N, \mathbf{A}, \mathbf{b}, c)$

Gradient

$$\nabla_{\mathbf{s}} \Phi(\mathbf{s}; \boldsymbol{\theta}) = \nabla_{\mathbf{s}} N(\mathbf{s}; \boldsymbol{\theta}_N) \mathbf{w} + (\mathbf{A}^\top \mathbf{A}) \mathbf{s} + \mathbf{b} \quad (9)$$

where

- space-time inputs $\mathbf{s} = (\mathbf{x}, t) \in \mathbb{R}^{d+1}$
- neural network $N(\mathbf{s}; \boldsymbol{\theta}_N): \mathbb{R}^{d+1} \rightarrow \mathbb{R}^m$ (we choose ResNet)
- $\boldsymbol{\theta}$ consists of all the trainable weights:
 $\mathbf{w} \in \mathbb{R}^m$, $\boldsymbol{\theta}_N \in \mathbb{R}^p$, $\mathbf{A} \in \mathbb{R}^{r \times (d+1)}$, $\mathbf{b} \in \mathbb{R}^{d+1}$, $c \in \mathbb{R}$
where $r = \min(10, d)$

Exact Trace Computation

Analytic Gradient Computation

N is an $(M + 1)$ -layer ResNet

Forward propagation

Compute $N(\mathbf{s}; \boldsymbol{\theta}_N) = \mathbf{u}_M$.

$$\mathbf{u}_0 = \sigma(\mathbf{K}_0 \mathbf{s} + \mathbf{b}_0)$$

$$\mathbf{u}_1 = \mathbf{u}_0 + h\sigma(\mathbf{K}_1 \mathbf{u}_0 + \mathbf{b}_1)$$

$$\vdots \quad \quad \quad \vdots$$

$$\mathbf{u}_M = \mathbf{u}_{M-1} + h\sigma(\mathbf{K}_M \mathbf{u}_{M-1} + \mathbf{b}_M)$$

where

- fixed step size $h > 0$
- ResNet weights $\boldsymbol{\theta}_N$ are
 - ▶ $\mathbf{K}_0 \in \mathbb{R}^{m \times (d+1)}$
 - ▶ $\mathbf{K}_1, \dots, \mathbf{K}_M \in \mathbb{R}^{m \times m}$
 - ▶ $\mathbf{b}_0, \dots, \mathbf{b}_M \in \mathbb{R}^m$
- $\sigma(\mathbf{x}) = \log(\exp(\mathbf{x}) + \exp(-\mathbf{x}))$
 - ▶ the antiderivative of hyperbolic tangent
 - ▶ so, $\sigma'(\mathbf{x}) = \tanh(\mathbf{x})$

Exact Trace Computation

Analytic Gradient Computation

N is an $(M + 1)$ -layer ResNet

Forward propagation

Compute $N(\mathbf{s}; \boldsymbol{\theta}_N) = \mathbf{u}_M$.

$$\mathbf{u}_0 = \sigma(\mathbf{K}_0 \mathbf{s} + \mathbf{b}_0)$$

$$\mathbf{u}_1 = \mathbf{u}_0 + h\sigma(\mathbf{K}_1 \mathbf{u}_0 + \mathbf{b}_1)$$

$$\vdots \quad \quad \quad \vdots$$

$$\mathbf{u}_M = \mathbf{u}_{M-1} + h\sigma(\mathbf{K}_M \mathbf{u}_{M-1} + \mathbf{b}_M) \quad \quad \quad \vdots \quad \quad \quad \vdots$$

Backpropagation

Compute $\nabla_{\mathbf{s}} N(\mathbf{s}; \boldsymbol{\theta}_N) \mathbf{w} = \mathbf{z}_0$

$$\mathbf{z}_{M+1} = \mathbf{w}$$

$$\mathbf{z}_M = \mathbf{z}_{M+1}$$

$$+ h \mathbf{K}_M^\top \text{diag}(\sigma'(\mathbf{K}_M \mathbf{u}_{M-1} + \mathbf{b}_M)) \mathbf{z}_{M+1}$$

$$\mathbf{z}_1 = \mathbf{z}_2 + h \mathbf{K}_1^\top \text{diag}(\sigma'(\mathbf{K}_1 \mathbf{u}_0 + \mathbf{b}_1)) \mathbf{z}_2$$

$$\mathbf{z}_0 = \mathbf{K}_0^\top \text{diag}(\sigma'(\mathbf{K}_0 \mathbf{s} + \mathbf{b}_0)) \mathbf{z}_1$$

Exact Trace Computation

Laplacian of the Potential

$$\text{tr}(\nabla^2 \Phi(\mathbf{s}; \boldsymbol{\theta})) = \text{tr}(\mathbf{E}^\top (\nabla_{\mathbf{s}}^2 (N(\mathbf{s}; \boldsymbol{\theta}_N) \mathbf{w}) + \mathbf{A}^\top \mathbf{A}) \mathbf{E}) \quad \text{for } \mathbf{E} = \text{eye}(d+1, d)$$

Exact Trace Computation

Laplacian of the Potential

$$\text{tr}(\nabla^2 \Phi(\mathbf{s}; \boldsymbol{\theta})) = \text{tr}(\mathbf{E}^\top (\nabla_{\mathbf{s}}^2 (N(\mathbf{s}; \boldsymbol{\theta}_N) \mathbf{w}) + \mathbf{A}^\top \mathbf{A}) \mathbf{E}) \quad \text{for } \mathbf{E} = \text{eye}(d+1, d)$$

Focus on the nontrivial part (the ResNet)

$$\text{tr}(\mathbf{E}^\top \nabla_{\mathbf{s}}^2 (N(\mathbf{s}; \boldsymbol{\theta}_N) \mathbf{w}) \mathbf{E}) = t_0 + h \sum_{i=1}^M t_i,$$

Exact Trace Computation

Laplacian of the Potential

$$\text{tr}(\nabla^2 \Phi(\mathbf{s}; \boldsymbol{\theta})) = \text{tr}(\mathbf{E}^\top (\nabla_{\mathbf{s}}^2(N(\mathbf{s}; \boldsymbol{\theta}_N)\mathbf{w}) + \mathbf{A}^\top \mathbf{A}) \mathbf{E}) \quad \text{for } \mathbf{E} = \text{eye}(d+1, d)$$

Focus on the nontrivial part (the ResNet)

$$\text{tr}(\mathbf{E}^\top \nabla_{\mathbf{s}}^2(N(\mathbf{s}; \boldsymbol{\theta}_N)\mathbf{w}) \mathbf{E}) = t_0 + h \sum_{i=1}^M t_i,$$

With one pass, we calculate the trace of each layer

$$\begin{aligned} t_0 &= \text{tr} \left(J_{i-1}^\top \nabla_{\mathbf{s}} (\mathbf{K}_i^\top \text{diag}(\sigma''(\mathbf{K}_i \mathbf{u}_{i-1}(\mathbf{s}) + \mathbf{b}_i)) \mathbf{z}_{i+1}) J_{i-1} \right) \\ &= (\sigma''(\mathbf{K}_0 \mathbf{s} + \mathbf{b}_0) \odot \mathbf{z}_1)^\top ((\mathbf{K}_0 \mathbf{E}) \odot (\mathbf{K}_0 \mathbf{E})) \mathbf{1} \quad \text{and} \end{aligned}$$

$$\begin{aligned} t_i &= \text{tr} \left(J_{i-1}^\top \nabla_{\mathbf{s}} (\mathbf{K}_i^\top \text{diag}(\sigma''(\mathbf{K}_i \mathbf{u}_{i-1}(\mathbf{s}) + \mathbf{b}_i)) \mathbf{z}_{i+1}) J_{i-1} \right) \\ &= (\sigma''(\mathbf{K}_i \mathbf{u}_{i-1} + \mathbf{b}_i) \odot \mathbf{z}_{i+1})^\top ((\mathbf{K}_i J_{i-1}) \odot (\mathbf{K}_i J_{i-1})) \mathbf{1}. \end{aligned}$$

where Jacobian $J_{i-1} = \nabla_{\mathbf{s}} \mathbf{u}_{i-1}^\top \in \mathbb{R}^{m \times d}$, \odot denotes Hadamard product, and $\mathbf{1} = \text{ones}(d, 1)$.

Exact Trace Computation

Efficiency

Update and overwrite $J_{i-1} = \nabla_{\mathbf{s}} \mathbf{u}_{i-1}^\top \in \mathbb{R}^{m \times d}$ in the forward pass as

$$\begin{aligned}\nabla_{\mathbf{s}} \mathbf{u}_i^\top &= \nabla_{\mathbf{s}} \mathbf{u}_{i-1} + h \sigma'(\mathbf{K}_i \mathbf{u}_{i-1} + \mathbf{b}_i) \mathbf{K}_i^\top \nabla_{\mathbf{s}} \mathbf{u}_{i-1} \\ J &\leftarrow J + h \sigma'(\mathbf{K}_i \mathbf{u}_{i-1} + \mathbf{b}_i) \mathbf{K}_i^\top J\end{aligned}$$

Overall Cost is $\mathcal{O}(m^2 \cdot d)$ FLOPS.

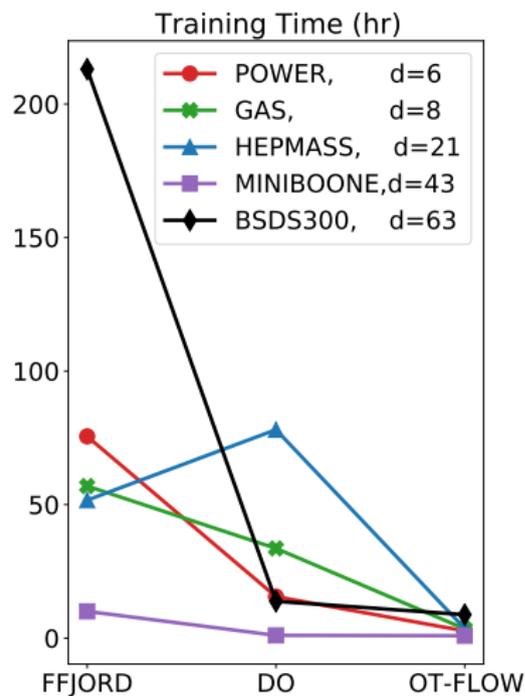
Recall: $\mathbf{K}_0 \in \mathbb{R}^{m \times (d+1)}$ and $\mathbf{K}_1, \dots, \mathbf{K}_M \in \mathbb{R}^{m \times m}$



L Ruthotto, S Osher, W Li, L Nurbekyan, S Wu Fung
A ML Framework for Solving High-Dimensional MFG and MFC
PNAS 117 (17), 9183-9193, 2020.

Results

Fast Training



OT-Flow has 19x training speed-up on average

Reasons

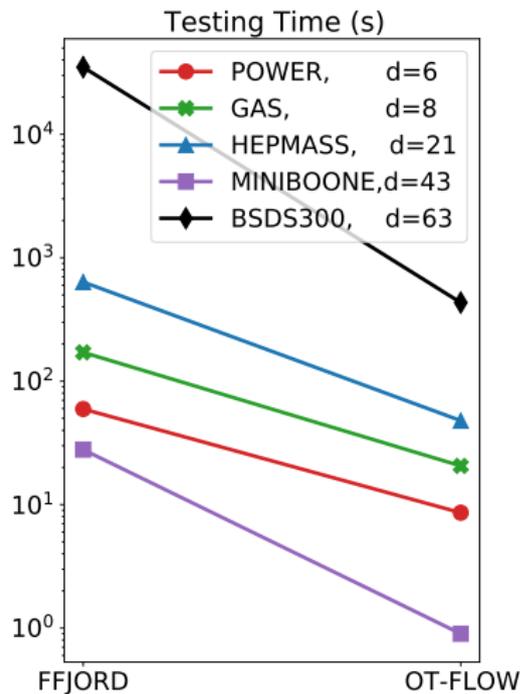
- OT-inspired regularization leads to straight trajectories that are inexpensive to integrate.
- The trace computation is efficient and exact.
- The potential flows approach results in fewer weights and a smaller model.



DO, S Wu Fung, X Li, L Ruthotto
OT-Flow: Fast and Accurate CNFs via OT
arXiv:2006.00104, 2020.

Results

Fast Inference



OT-Flow has 28x testing speed-up on average

Inference-Specific Reasons

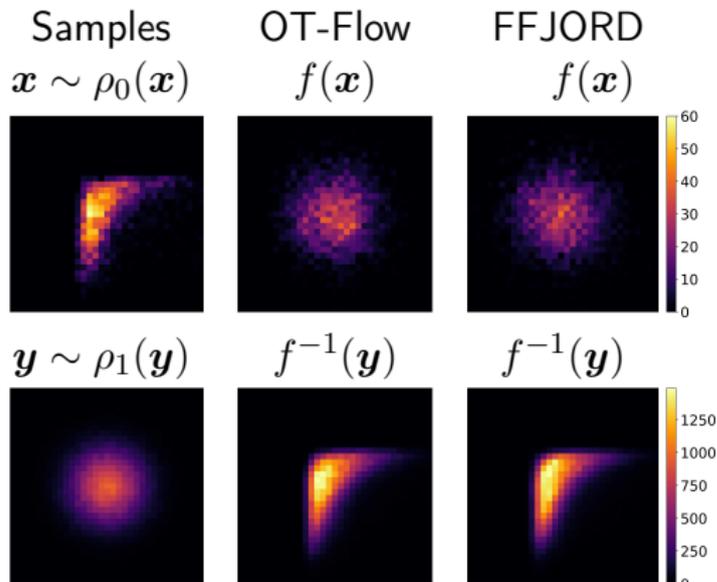
- Inference uses exact trace (no estimates)
 - ▶ State-of-the-art approaches use AD to obtain exact trace with $\mathcal{O}(d^2)$
 - ▶ Meanwhile, our exact trace is $\mathcal{O}(d)$



DO, S Wu Fung, X Li, L Ruthotto
OT-Flow: Fast and Accurate CNFs via OT
arXiv:2006.00104, 2020.

More Results

Details in paper



Two of the 43 dimensions in the MINIBOONE CNF.



MNIST synthetic generation.
Original images boxed in red.



DO, S Wu Fung, X Li, L Ruthotto
*OT-Flow: Fast and Accurate CNFs
via OT*
arXiv:2006.00104, 2020.

Conclusions

Discretize-Optimize

- DO often converges faster than OD when used in neural ODEs
- For CNFs, DO provides 6x training speedup

OT-Flow

- OT regularization \Rightarrow well-posed and efficient time integration
- Potential flow \Rightarrow smaller model
- OT-Flow achieves 19x training speedup and 28x inference speedup over same baseline



DO and L Ruthotto
*DO vs. OD for Time-Series
Regression and CNFs*
arXiv:2005.13420, 2020.



DO, S Wu Fung, X Li, L Ruthotto
OT-Flow: Fast and Accurate CNFs via OT
arXiv:2006.00104, 2020.

Code: github.com/EmoryMLIP/OT-Flow

References I

- Chen, Tian Qi et al. (2018). “Neural Ordinary Differential Equations”. In: *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 6571–6583.
- E, Weinan (2017). “A Proposal on Machine Learning via Dynamical Systems”. In: *Comm. Math. Statist.* 5.1, pp. 1–11.
- Evans, Lawrence C (1997). “Partial differential equations and Monge-Kantorovich mass transfer”. In: *Current developments in mathematics* 1997.1, pp. 65–126.
- (2013). *An Introduction to Mathematical Optimal Control Theory Version 0.2*.
- Finlay, Chris et al. (2020). “How to train your neural ODE”. In: *arXiv:2002.02798*.
- Gholaminejad, Amir, Kurt Keutzer, and George Biros (2019). “ANODE: Unconditionally Accurate Memory-Efficient Gradients for Neural ODEs”. In: *International Joint Conference on Artificial Intelligence (IJCAI)*, pp. 730–736.

References II

- Grathwohl, Will et al. (2019). “FFJORD: Free-form continuous dynamics for scalable reversible generative models”. In: *International Conference on Learning Representations (ICLR)*.
- Haber, Eldad and Lars Ruthotto (2017). “Stable Architectures for Deep Neural Networks”. In: *Inverse Probl.* 34 (1).
- He, Kaiming et al. (2016). “Deep residual learning for image recognition”. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778.
- Li, Qianxiao et al. (2017). “Maximum principle based algorithms for deep learning”. In: *The Journal of Machine Learning Research (JMLR)* 18.1, pp. 5998–6026.
- Papamakarios, George et al. (2019). “Normalizing Flows for Probabilistic Modeling and Inference”. In: *arXiv:1912.02762*.
- Rezende, Danilo Jimenez and Shakir Mohamed (2015). “Variational Inference with Normalizing Flows”. In: *International Conference on Machine Learning (ICML)*. Lille, France, pp. 1530–1538.

References III

- Yang, Liu and George Em Karniadakis (2019). “Potential Flow Generator with L_2 Optimal Transport Regularity for Generative Models”. In: *arXiv:1908.11462*.
- Zhang, Linfeng, Weinan E, and Lei Wang (2018). “Monge-Ampère Flow for Generative Modeling”. In: *arXiv:1809.10188*.