

# **Improving and Generalizing Flow-Based Generative Models with Minibatch Optimal Transport**

Tong et al.

Presented by Beichen Gu, Siddarth Chilukuri

# Problem Setup

- Goal: Generative Modeling
- Learn mapping:  
 $p_0 \rightarrow p_1$
- From simple  $\rightarrow$  complex distributions
  
- Applications:
  - Image generation
  - Scientific data modeling

# Continuous Normalizing Flows (CNFs)

- Key idea: model transformation as ODE

$$dx = u_t(x) dt$$

- $u_t(x)$ : velocity field
- Defines how samples move over time
- Produces continuous, invertible transformations

# From Points to Distributions

- Density evolution follows:

$$\frac{\partial p_t}{\partial t} = -\nabla \cdot (p_t u_t)$$

- Ensures probability is preserved
- Connects: point dynamics  $\rightarrow$  distribution dynamics

# CNFs vs Diffusion models

- **CNF:**
  - Simulating ODE trajectories
  - Computing likelihood via integration
- **Diffusion models:**
  - Simple regression training
  - No simulation needed
  - But Requires many sampling steps → Slow inference

## Flow Matching (Before This Paper)

- Learn velocity field via regression
- Avoid simulation

$$\mathcal{L}_{FM}(\theta) = \mathbb{E}_{t,x \sim p_t} \|v_{\theta}(t, x) - u_t(x)\|^2$$

- Limitation:
- sampling from  $p_t(x)$
- often Gaussian assumptions

## Key Idea: Conditional Decomposition

- Break global problem into conditional subproblems
- Each  $z$ : defines a simple transport task

$$p_t(x) = \int p_t(x|z)q(z)dz$$

- What is condition  $z$ ?

$$z = (x_0, x_1)$$

$x_0$ : source sample,  $x_1$ : target sample

## True Velocity Field (Intractable)

- Problem:
- Requires  $p_t(x)$

$$u_t(x) = \mathbb{E}_{q(z)} \left[ u_t(x|z) \frac{p_t(x|z)}{p_t(x)} \right]$$

## CFM Objective

- Same gradient as original objective

$$\mathcal{L}_{CFM}(\theta) = \mathbb{E}_{t,z,x \sim p_t(x|z)} \|v_\theta(t, x) - u_t(x|z)\|^2$$

$$\mathcal{L}_{FM}(\theta) = \mathbb{E}_{t,x \sim p_t} \|v_\theta(t, x) - u_t(x)\|^2$$

## Design Choice: Coupling

- Independent coupling: simple but inefficient
- Optimal transport coupling: matches points optimally, shorter transport paths

$$z \sim \pi(x_0, x_1)$$

## OT-CFM (Key Contribution)

- OT objective:

$$\min \int p_t(x) \|u_t(x)\|^2$$

- Minimizes transport cost
- Efficient trajectories
- Fewer ODE steps

## Recap: Flow Matching & CFM

- We want to learn a vector field that transports one distribution → another
- Flow Matching (FM):
  - trains CNFs using a regression objective
  - avoids expensive simulation during training
- Conditional Flow Matching (CFM):
  - generalizes FM to more flexible settings
- Key idea:
  - learn velocity field instead of simulating dynamics

# Paper Contributions

- Generalized Conditional Flow Matching (CFM):
  - unified framework for training CNFs without simulation
  - works with arbitrary source distributions
- Introduced OT-CFM:
  - uses optimal transport coupling  $\pi(\mathbf{x}_0, \mathbf{x}_1)$
  - produces straighter, more efficient flows
- Proposed Minibatch OT approximation:
  - scalable alternative to full OT
  - works well in practice despite approximation

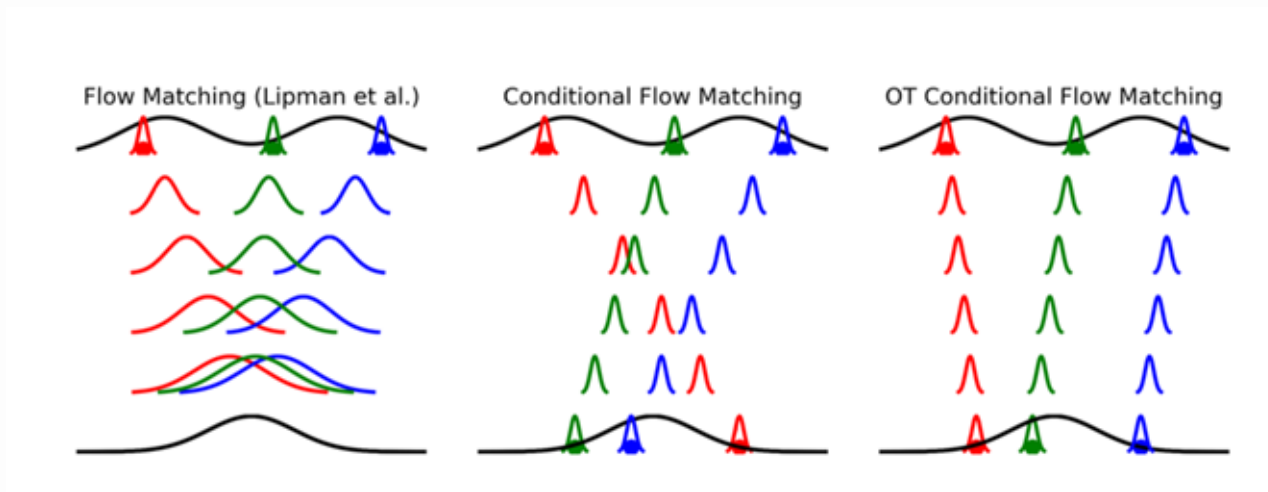
# Paper Contributions

- Extensive experiments:
  - low-dimensional OT tasks
  - CIFAR-10 image generation
  - single-cell dynamics
  - image translation + Schrödinger bridge

## Problem with CFM

- In CFM, we sample pairs  $(x_0, x_1)$
- These pairs define how mass is transported
- Current approach:
  - pairs are sampled independently
- Problem:
  - transport paths are arbitrary
  - flows become inefficient and noisy
  - harder to train, slower to simulate

# Problem with CFM



# Optimal Transport (OT)

- Instead of random pairing, use optimal pairing
- Match points that are closest and most efficient to transport
- Objective:
  - minimize total transport cost
- Produces a transport plan  $\pi(\mathbf{x}_0, \mathbf{x}_1)$
- Gives “straight” and efficient paths

# OT-CFM

- Keep same framework (CFM)
- Change **ONLY** how we sample pairs  $(x_0, x_1)$  :
  - Instead of random pairing,
  - Use optimal transport plan  $\pi(x_0, x_1)$
- Leads to:
  - straighter trajectories
  - smoother vector fields
- This reduces:
  - variance in training
  - complexity of learned dynamics

# Connection to Dynamic Optimal Transport

- OT-CFM uses optimal transport coupling  $\pi(x_0, x_1)$
- Paper shows:
  - learned vector field approximates dynamic OT solution
- As noise  $\rightarrow 0$ :
  - flow minimizes transport energy
- This provides theoretical grounding, not just heuristic improvement

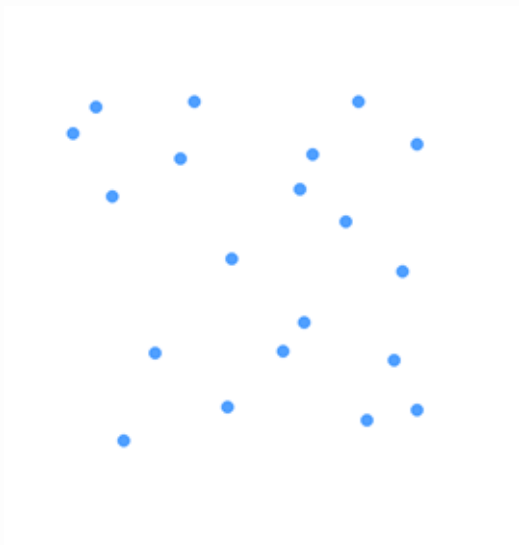
# OT Computation Cost

- Exact OT:
  - cubic time complexity
  - quadratic memory
- Not scalable for large datasets

## **Solution: Minibatch OT**

- Compute OT within batches
- Instead of global OT:
  - approximate locally
- For each batch:
  - Sample batch from source and target
  - compute OT matching within batch
  - train model on matched pairs

# How Minibatch OT Works

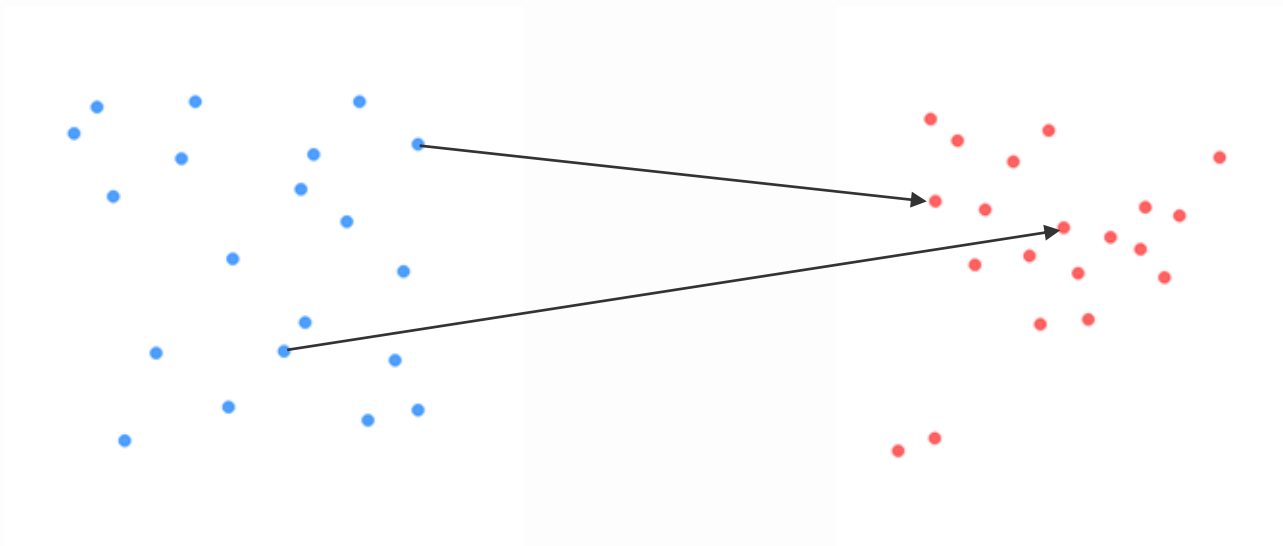


Source distribution  
Sample source batch  $\mathbf{x}_0$



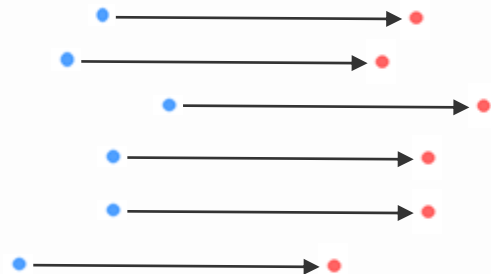
Target distribution  
Sample target batch  $\mathbf{x}_1$

# How Minibatch OT Works



OT Matching  
Compute OT Matching  $\pi_{batch}$

# How Minibatch OT Works



Matched pairs  
Form  $(x_0, x_1)$



Sample Interpolation  
 $x_t$

# How Minibatch OT Works

Sample Interpolation  
 $x_t$



Train velocity model  
 $v_\theta(t, x_t)$

## Why Minibatch OT is enough

- Exact OT not needed
- Even with approximate matching, improves training significantly
- Works well empirically

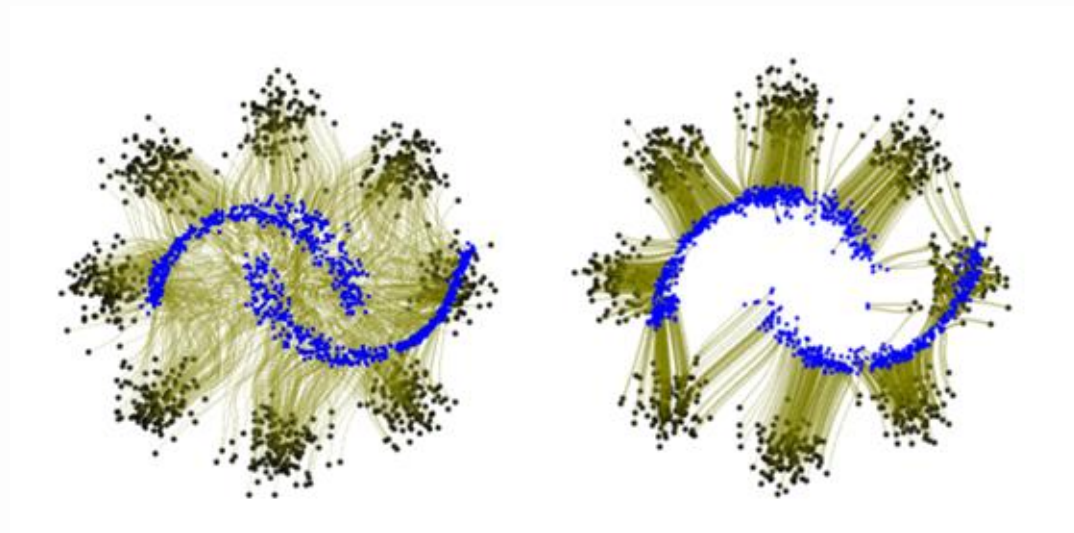
## Summary

- CFM  $\rightarrow$  regression-based training
- OT-CFM  $\rightarrow$  better pairing using OT
- Minibatch OT  $\rightarrow$  scalable version
- Outcome:
  - better flows
  - faster training
  - faster inference

# Experimental Setup

- Evaluated on:
  - low-dimensional datasets (Gaussian, moons, etc.)
  - Schrödinger bridge tasks
  - single-cell data
  - CIFAR-10 image generation
  - CelebA translation
- Metrics:
  - Wasserstein distance
  - path energy
  - FID (images)
  - MMD (translation)

# Low-Dimensional Results



OT-CFM achieves lowest normalized path energy, closest to true optimal transport solution

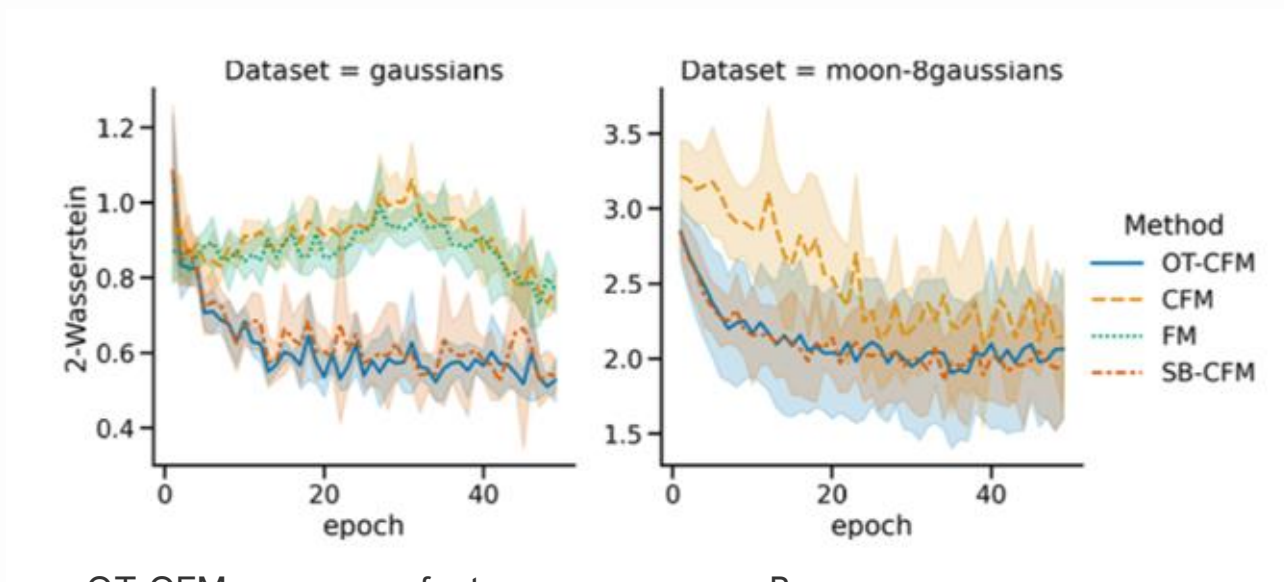
# Low-Dimensional Results

Table 2: Comparison of neural optimal transport methods over four distribution pairs ( $\mu \pm \sigma$  over five seeds) in terms of fit (2-Wasserstein), optimal transport performance (normalized path energy), and runtime. ‘—’ indicates a method that requires a Gaussian source. Best in **bold**. CFM and RF models are trained on a single CPU core, other baselines are trained with a GPU and two CPUs.

| Dataset →<br>Algorithm ↓ Metric → | $\mathcal{N} \rightarrow 8\text{gaussians}$ |                          | moons $\rightarrow 8\text{gaussians}$ |                          | $\mathcal{N} \rightarrow \text{moons}$ |                          | $\mathcal{N} \rightarrow \text{scurve}$ |                          | Avg. train time          |
|-----------------------------------|---|--------------------------|---------------------------------------|--------------------------|--|--------------------------|---|--------------------------|--------------------------|
|                                   | $W_2^2$                                     | NPE                      | $W_2^2$                               | NPE                      | $W_2^2$                                | NPE                      | $W_2^2$                                 | NPE                      | ( $\times 10^3$ s)       |
| OT-CFM                            | 1.262 $\pm$ 0.348                           | <b>0.018</b> $\pm$ 0.014 | <b>1.923</b> $\pm$ 0.391              | <b>0.053</b> $\pm$ 0.035 | <b>0.239</b> $\pm$ 0.048               | 0.087 $\pm$ 0.061        | <b>0.264</b> $\pm$ 0.093                | <b>0.027</b> $\pm$ 0.026 | 1.129 $\pm$ 0.335        |
| I-CFM                             | 1.284 $\pm$ 0.384                           | 0.222 $\pm$ 0.032        | 1.977 $\pm$ 0.266                     | 2.738 $\pm$ 0.181        | 0.338 $\pm$ 0.109                      | 0.841 $\pm$ 0.148        | 0.333 $\pm$ 0.060                       | 0.867 $\pm$ 0.117        | <b>0.630</b> $\pm$ 0.365 |
| 2-RF (Liu, 2022)                  | 1.436 $\pm$ 0.344                           | 0.069 $\pm$ 0.027        | 2.211 $\pm$ 0.423                     | 0.149 $\pm$ 0.101        | 0.278 $\pm$ 0.026                      | <b>0.076</b> $\pm$ 0.067 | 0.395 $\pm$ 0.111                       | 0.112 $\pm$ 0.085        | 0.862 $\pm$ 0.166        |
| 3-RF (Liu, 2022)                  | 1.337 $\pm$ 0.367                           | 0.055 $\pm$ 0.043        | 2.700 $\pm$ 0.587                     | 0.123 $\pm$ 0.112        | 0.305 $\pm$ 0.026                      | 0.084 $\pm$ 0.051        | 0.395 $\pm$ 0.082                       | 0.129 $\pm$ 0.075        | 0.954 $\pm$ 0.116        |
| FM (Lipman et al., 2023)          | 1.062 $\pm$ 0.196                           | 0.174 $\pm$ 0.030        | —                                     | —                        | 0.246 $\pm$ 0.077                      | 0.778 $\pm$ 0.144        | 0.377 $\pm$ 0.099                       | 0.772 $\pm$ 0.081        | 0.708 $\pm$ 0.370        |
| Reg. CNF (Finlay et al., 2020)    | 1.144 $\pm$ 0.075                           | 0.274 $\pm$ 0.060        | —                                     | —                        | 0.376 $\pm$ 0.040                      | 0.620 $\pm$ 0.088        | 0.581 $\pm$ 0.195                       | 0.586 $\pm$ 0.503        | 8.021 $\pm$ 3.288        |
| CNF (Chen et al., 2018)           | <b>1.055</b> $\pm$ 0.059                    | 0.151 $\pm$ 0.064        | —                                     | —                        | 0.387 $\pm$ 0.065                      | 2.937 $\pm$ 1.973        | 0.645 $\pm$ 0.343                       | 10.548 $\pm$ 8.100       | 18.810 $\pm$ 12.677      |
| ICNN (Makkuva et al., 2020)       | 1.771 $\pm$ 0.398                           | 0.747 $\pm$ 0.029        | 2.193 $\pm$ 0.136                     | 0.832 $\pm$ 0.004        | 0.532 $\pm$ 0.046                      | 0.267 $\pm$ 0.010        | 0.753 $\pm$ 0.068                       | 0.344 $\pm$ 0.045        | 2.912 $\pm$ 0.626        |

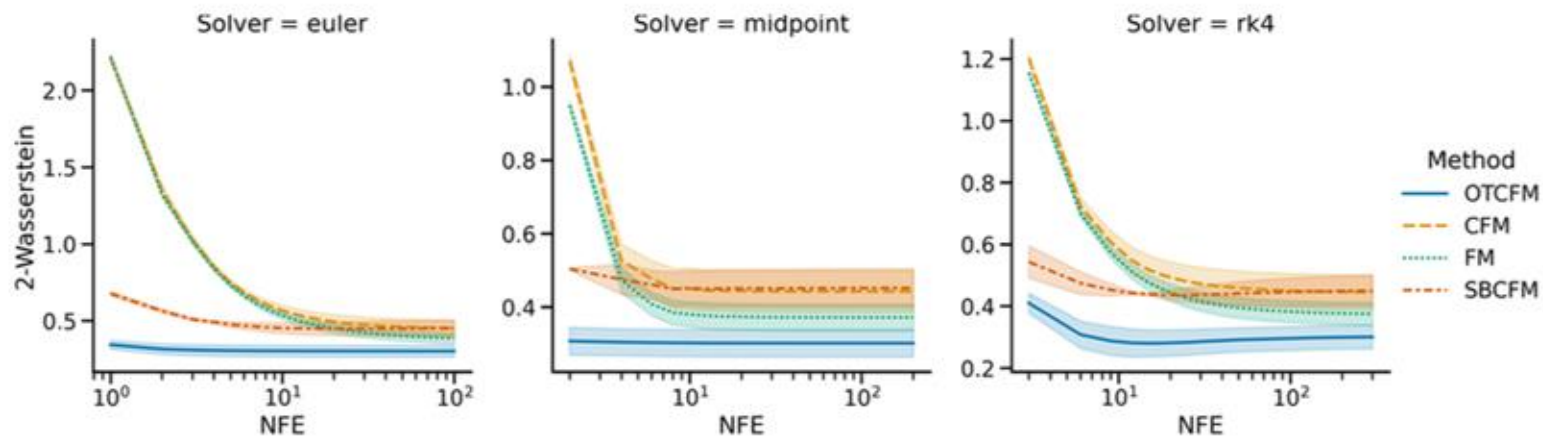
- OT-CFM outperforms I-CFM and other baselines
- Especially strong on moons  $\leftrightarrow$  8 Gaussians (hard case)

# Faster Training



- OT-CFM converges faster (Tong et al., 2023)
  - Lower validation error earlier
- Because:
- OT pairing reduces variance in training targets
  - easier optimization

# Faster Inference



- OT-CFM needs fewer function evaluations (Tong et al., 2023)
- Same quality with fewer steps

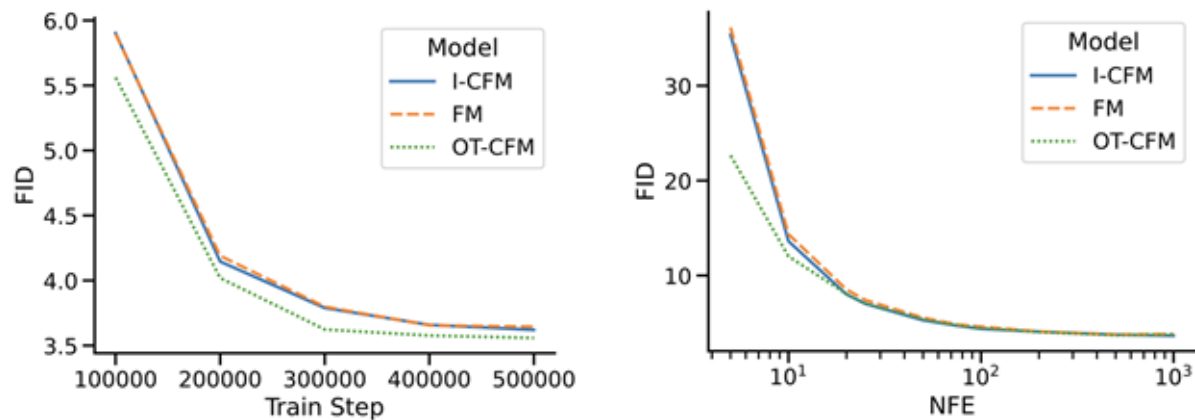
Reason:

- straighter flows = easier ODE integration

# High-Dimensional Results (CIFAR-10)

- Evaluated on CIFAR-10 image generation
- Model setup:
  - U-Net architecture (similar to diffusion models)
- Metrics:
  - FID (Fréchet Inception Distance)
  - NFE (Number of Function Evaluations)
- Key results:
  - OT-CFM achieves competitive or better FID scores
  - Requires fewer function evaluations (NFE) for similar quality
  - Performs especially well under limited compute budgets

# High-Dimensional Results (CIFAR-10)



- OT-based flows are straighter → easier to integrate
- Leads to faster and more efficient generation

# Applications of OT-CFM

- Single-cell dynamics
  - Task: interpolate cell distributions over time
  - Metric: Earth mover's distance (EMD)
- OT-CFM achieves lowest error across datasets (Tong et al., 2023)

Table 4: Single-cell comparison over three datasets averaged over leaving out intermediate time-points measuring EMD to left out distribution following Tong et al. (2020). \*Indicates values taken from aforementioned work.

| Algorithm ↓ Dataset →              | Cite                 | EB                   | Multi                |
|------------------------------------|----------------------|----------------------|----------------------|
| TrajectoryNet (Tong et al., 2020)* | —                    | 0.848 ± —            | —                    |
| Reg. CNF (Finlay et al., 2020)*    | —                    | 0.825 ± —            | —                    |
| DSB (De Bortoli et al., 2021)      | 0.953 ± 0.140        | 0.862 ± 0.023        | 1.079 ± 0.117        |
| I-CFM                              | 0.965 ± 0.111        | 0.872 ± 0.087        | 1.085 ± 0.099        |
| SB-CFM                             | 1.067 ± 0.107        | 1.221 ± 0.380        | 1.129 ± 0.363        |
| OT-CFM                             | <b>0.882 ± 0.058</b> | <b>0.790 ± 0.068</b> | <b>0.937 ± 0.054</b> |

# Applications of OT-CFM

- Image generation
  - Dataset: CIFAR-10
  - Metrics: FID + number of function evaluations (NFE)
- OT-CFM achieves competitive FID.
- Requires fewer evaluations -> faster generation (Tong et al., 2023)

# Applications of OT-CFM

- Unsupervised image translation
  - Task: map between attributes (CelebA latent space)
  - Metric: Maximum Mean Discrepancy (MMD)
- OT-CFM achieves better alignment between distributions (Tong et al., 2023)

Table 6: MMD (in units of  $10^{-3}$ ) between target and transformed source samples of CelebA latent vectors. Mean and standard deviation over 40 attributes and both translation directions ( $- \leftrightarrow +$ ) for each attribute. ‘Identity’ refers to performing no translation and treating source samples as approximate samples from the target.

| Algorithm ↓ | $\sigma = 0.1$  | $\sigma = 0.3$  | $\sigma = 1$    |
|-------------|-----------------|-----------------|-----------------|
| Identity    | $9.17 \pm 5.68$ | $9.17 \pm 5.68$ | $9.17 \pm 5.68$ |
| I-CFM       | $4.85 \pm 5.09$ | $3.44 \pm 2.03$ | $1.59 \pm 0.83$ |
| OT-CFM      | $2.81 \pm 2.62$ | $1.91 \pm 1.30$ | $1.04 \pm 0.60$ |

# Applications of OT-CFM

- Schrödinger Bridge
  - Learns probability flow between distributions
- accurately approximates SB solutions, faster than diffusion-based approaches

Table 3: Schrödinger bridge flow comparison, showing average error over flow time to ground truth averaged over 5 models for SB-CFM and 5 dynamics from DSB (De Bortoli et al., 2021).

| Dataset ↓ Alg. →                            | SB-CFM               | DSB           |
|---|----------------------|---------------|
| $\mathcal{N} \rightarrow 8\text{gaussians}$ | <b>0.454 ± 0.164</b> | 1.440 ± 0.720 |
| moons $\rightarrow 8\text{gaussians}$       | <b>1.377 ± 0.229</b> | 2.407 ± 1.025 |
| $\mathcal{N} \rightarrow \text{moons}$      | <b>0.283 ± 0.048</b> | 0.333 ± 0.129 |
| $\mathcal{N} \rightarrow \text{scurve}$     | <b>0.297 ± 0.064</b> | 0.383 ± 0.134 |

## Paper Strengths

- Simulation-free training
- Supports general source distributions
- Theoretically grounded (OT + SB)
- Faster training + inference
- Strong empirical results

## Paper Weaknesses

- Requires closed-form conditional paths
- Minibatch OT is approximate
- OT still adds computational overhead
- Performance depends on batch quality

## Open Questions

- Can we learn OT maps directly?
- How does minibatch OT scale to very high dimensions?
- Can we combine this with diffusion models?
- Better approximations to OT?

## Final Takeaways

- CFM = simulation-free training for CNFs
- OT-CFM = smarter pairing using optimal transport
- Minibatch OT = makes it scalable
- Result:
  - better flows
  - faster training
  - faster inference