

---

# Flow-Based Generative Models with Minibatch Optimal Transport

---

**Navin Vincent**  
navinv@kth.se

**Jacques Fürst**  
jfurst@kth.se

**Michał Sitarz**  
sitarz@kth.se

## Abstract

Optimal Transport Conditional Flow Matching (OT-CFM) has recently emerged as a promising framework for learning generative models by combining continuous normalizing flows with optimal transport principles. While OT-CFM demonstrates notable improvements in sample quality, it often demands substantial computational resources during training and inference. In this report, we replicate the results of chosen experiments from the original paper (See Tong et al. 2023) and we evaluate the Light Schrödinger Bridge (Light SB) (inspired from Korotin, Gushchin, and Burnaev 2023) solver as an efficient, simulation-free optimization technique to approximate dynamic optimal transport maps. By parameterizing log-Schrödinger potentials as energy functions, Light SB dramatically reduces training time—often from hours to mere minutes—without compromising quality. Through experiments on two-dimensional toy datasets, image generation on CIFAR-10, and latent-space translation tasks with CelebA, we show that Light SB not only preserves or improves upon the generative performance of OT-CFM but also significantly accelerates inference by forgoing multiple passes of ODE based flow computation. These results highlight the potential of Light SB to enable more scalable and readily deployable generative modeling solutions, paving the way for broader practical applications. Code repository: [https://gits-15.sys.kth.se/sitarz/Group33\\_OT-CFM.git](https://gits-15.sys.kth.se/sitarz/Group33_OT-CFM.git).

## 1 Introduction

This report discusses the evolution and challenges of generative modeling approaches, particularly focusing on continuous normalizing flows (CNFs) and their relationship to diffusion models. Initially, normalizing flows were developed as static transformations mapping between base and target distributions, before evolving into CNFs that use neural ODEs for more flexible mappings. While diffusion models, which use stochastic differential equations, have become state-of-the-art, they share with CNFs the limitation of requiring multiple network passes for integration. Flow matching (FM) emerged as a training method that improved CNF training by regressing the ODE’s drift term, similar to diffusion models’ training approach. The concept was later expanded into conditional flow matching (CFM), allowing for arbitrary transport maps and generalizing both FM and diffusion approaches. The paper we chose introduces optimal transport conditional flow matching (OT-CFM), which addresses the computational efficiency challenges of these models while improving optimal transport flow accuracy. This method, along with its entropic variant, enables efficient training of CNFs to match Schrödinger bridge probability flows and can approximate dynamic optimal transport maps when the true transport plan can be sampled.

In our report, we implemented and evaluated the lightweight Schrödinger Bridge (Light SB) solver proposed by Korotin et al. on the datasets of the original paper. This solver combines the parameterization of Schrödinger potentials with sum-exp quadratic functions and the interpretation of log-Schrödinger potentials as energy functions, resulting in a simple, efficient, and simulation-free optimization framework. We replicated experiments demonstrating its ability to solve SB problems in moderate dimensions within minutes on a CPU, without extensive hyperparameter tuning. We were able to reproduce to a reasonable extent the results presented in the paper using the compute capacity that was available to us. We were able to replicate the experiments performed in Sections 5.1, 5.3 and 5.4 in the original paper. We were also able to extend our experiments to use the light SB framework and were able to record a massive decrease in the compute time required to match the target distribution.

## 2 Related Work

Simulation-free training is widely used in stochastic flow models to address challenges in backpropagation through simulations, which are often numerically unstable and high in variance (Li et al 2020). While diffusion models have demonstrated exceptional generative performance((Song and Ermon 2019, Ho, Jain, and Abbeel 2020), their reliance on costly SDE simulations has prompted efforts to enhance inference efficiency (Lu et al. 2022, Salimans and Ho 2022). Most methods assume Gaussian diffusion processes, with limited exploration of general source distributions, which complicates optimization and inference. Other papers that consider more general source distributions but this would make optimization and inference more difficult as it would require multiple iterations or other artifices to achieve good performance (Bortoli et al. 2023, Wang et al. 2021)

Recent work has advanced simulation-free training of continuous normalizing flows which are equivalent to CFMs using Gaussian sources or independent samples from the source and target distributions.(. (Ben-Hamu et al. 2022, Rozen et al. 2021). There was also a paper concurrent to this work (Pooladian et al. 2023) with very similar results.

Dynamic OT methods typically rely on constrained architectures ((Leygonie et al. 2019, Makkavu et al. 2019) or regularized CNFs ( Huang, et al. 2020, Finlay et al. 2020), which are difficult to optimize. This paper introduces an approach that achieves efficient and accurate OT flows without these constraints, extending the scope of simulation-free modeling and bridging gaps in generalizing source distributions and transport maps.

## 3 Methods

### 3.1 Main Paper

CNFs define a time-dependent vector field  $u_t(x)$  that transforms samples between distributions via the ODE  $\frac{dx}{dt} = u_t(x)$ , with corresponding density evolution governed by the continuity equation  $\frac{\partial p}{\partial t} + \nabla \cdot (pu_t) = 0$ . Flow matching optimizes a neural network  $v_\theta(t, x)$  to approximate this vector field by minimizing  $\mathcal{L}_{\text{FM}}(\theta) = \mathbb{E}_{t \sim U(0,1), x \sim p_t} |v_\theta(t, x) - u_t(x)|^2$ .

The framework connects to optimal transport (OT), which finds minimal-cost mappings between distributions through both static ( $W_2^2(q_0, q_1) = \inf_{\pi \in \Pi(q_0, q_1)} \int_{\mathbb{R}^d \times \mathbb{R}^d} |x - y|^2 d\pi(x, y)$ ) and dynamic ( $W_2^2(q_0, q_1) = \inf_{(p_t, u_t)} \int_0^1 \int_{\mathbb{R}^d} p_t(x) |u_t(x)|^2 dx dt$ ) formulations. The Schrödinger Bridge (SB) problem extends this by adding entropy regularization, solving  $\pi^* = \arg \min_{\pi \in \Pi(q_0, q_1)} \text{KL}(\pi \| \pi_{\text{ref}})$  with respect to a reference Brownian motion. These concepts are unified through Conditional Flow Matching (CFM), which minimizes  $\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{t, z \sim q(z), x \sim p_t(x|z)} |v_\theta(t, x) - u_t(x|z)|^2$ , where conditioning variable  $z$  couples source and target distributions (ie  $q(z) := \pi(x_0, x_1)$ ). It is important to note that  $\nabla_\theta \mathcal{L}_{\text{FM}} = \nabla_\theta \mathcal{L}_{\text{CFM}}$ , ensures that optimizing either objective yields the same result. This is critical because CFM allows flexibility in conditioning on arbitrary couplings  $z$ , while retaining the computational stability of FM. OT-CFM extends CFM by incorporating optimal transport couplings as the conditioning variable  $z = (x_0, x_1)$ , sampled from the OT plan  $\pi$ . This ensures the learned flows approximate dynamic OT. Similarly, SB-CFM incorporates entropy-regularized OT couplings, aligning the learned flows with the Schrödinger Bridge solutions. Please see Tong et al. 2023 as the theory is too extensive to describe fully given the page limit.

### 3.2 Light Schrödinger Bridge

Going beyond the paper’s contents, we implement the ‘Light Schrödinger Bridge’ Korotin, Gushchin, and Burnaev 2023 method and evaluate its performance in all three experimental setups reproduced from the original paper. This method takes advantage of the inherent stochasticity of the OT problem. Just as in SB-CFM, they aim to minimize the KL-divergence between the Brownian motion from  $p_0$  to  $p_1$  and the transport map  $\pi$ . The key difference lies in the parameterization of the variables. Their approach is parameterizing the Gaussian mixture density of  $p_1$  and its regularizer  $c$ . Here, the loss is constructed as the distance between the source and target distribution, rather than between the learned and optimal transport map. Again, please see Korotin, Gushchin, and Burnaev 2023 for full details.

## 4 Data

Firstly, pairs of various 2D toy datasets were used in Section 5.1: (1) standard gaussian, (2) 8gaussian, (3) scurve, (4) moons, and (5) swiss roll. They datasets are generated using NumPy and SkLearn, to fit ranges presented in (Tong et al. 2023; Korotin, Gushchin, and Burnaev 2023). Secondly, in Section 5.2 we use CIFAR-10 images, from TorchVision datasets. Pre-processing included normalization and random flipping. Finally, in Section 5.3, CelebA dataset was used (obtained from Liu et al. 2015), and it contains 200,000 images has 40 binary notation (like “smiling” or not).

Pre-processing included center cropping, resizing to size 64, and normalizing. Furthermore, a VAE model was trained to encode them into 128-dimensional latent vectors (see Section 5.3).

## 5 Experiments and Findings

This section evaluates the empirical difference on various tasks, between methods utilizing Optimal Transport plans, such as OT-CFM and SB-CFM, and the ones that do not, including FM and I-CFM,

### 5.1 Toy Datasets

The following set of experiments values the different algorithms as generative models in low-dimensional space, working with toy datasets (as described in Section 4). A simple Multilayer Perceptron (MLP) is used for learning the vector field.

The first experiment was to verify whether OT-CFM approximates the dynamic Optimal Plan, and working with a low dimensional dataset allows to verify whether it works as anticipated. Figure 3 shows a clear difference between using I-CFM and OT-CFM when it comes to trajectory generation. The trajectories were created using an rk4 solver with 100 steps. Furthermore, we used the 2-Wasserstein difference on the endpoints of the generated trajectories to verify how well the source distribution maps to the goal distribution, with the trained vector field. The results in Table 1 show that with a large batch size, the two OT plan methods outperform I-CFM which does not use it, which supports the findings of the original paper. However, we specifically decided to show results on the *swiss roll* dataset which is slightly more complex than the toy datasets that were used in the original paper. An interesting finding that we found was that the batch size seems to impact the training much more on the OT methods, rather than on the ones not using OT. In the original paper the authors claimed that "OT-CFM requires surprisingly small batches to approximate the OT map well", but our training methods struggle to perform well with smaller batch sizes (similar findings are shown in Section 5.2). Furthermore, we had to extend the training time as the length of training proposed in the original paper was not sufficient enough for the model to learn the vector field accurately (we trained for 20 thousand training steps instead of a thousand).

The second experiment that we tried to reproduce involved testing if OT-CFM yields faster training and inference, and the results can be observed in Figure 6. It can be concluded from the left graph that OT-CFM performs better for the same number of steps than I-CFM on the validation set. The overall findings show that OT-CFM is easier to train, which most likely comes from variance reductions of the conditional flow. For the other graphs, we can see that the quality of samples, as we decrease the NFEs (number of function evaluations), decreases for FM and I-CFM, but retains most of the quality for the algorithms with OT plan. This can once again be attributed to the generated paths we showed in Figure 3.

Thirdly, we reproduced the experiment for SB-CFM reproduction of Schrödinger bridge flows, focusing exclusively on SB-CFM to assess its performance. Unlike the original experiment, we did not include the Diffusion Schrödinger Bridge (DSB) method, as our goal was to evaluate SB-CFM independently. Our results closely matched those reported in the study, demonstrating similar levels of accuracy based on the average 2-Wasserstein distance to ground truth Schrödinger bridge samples over 18 time steps ( see Table 2). This consistency validates the robustness of SB-CFM as a method for reproducing Schrödinger bridge flows.

Finally, as shown in Table 1 we did some extra modifications, hyperparameter searches, and ablations. The main hyperparameter involved is  $\sigma$ , and on the toy datasets, we found that it does not lead to any noticeable difference in training ( $\sigma \in \{0.0, 0.01, 0.1\}$ ). A major component of training is the model that is trained to estimate the vector field, and the authors used the same one as Lipman et al. 2022, so we thought it would be interesting to experiment slightly with it and observing how it affects OT-CFM. We tried networks that are: wider (128 instead of 64 neurons per layer), deeper (extra 64 neuron layer), shallower (removed one layer), and have a different activation function (ReLU instead of SELU). The results showed that all the modifications decreased the performance, except for a very tiny boost for a wider network. However, these results show that the networked proposed initially by Lipman et al. 2022 seems to work the best for this task, and batch size is much more important than the choice of a network.

### 5.2 CIFAR-10 dataset

The experiment aims to examine the performance of OT-CFM on high-dimensional data (CIFAR-10) generation from Gaussian noise. For learning the vector field a UNet model was used. We compare the results of OT-CFM with I-CFM and FM. To evaluate the performance at every training step, we use Fréchet inception distance (FID) with a *dopri5* evaluation solver. Figure 2 summarizes the results. The original paper found that they achieve significantly better results

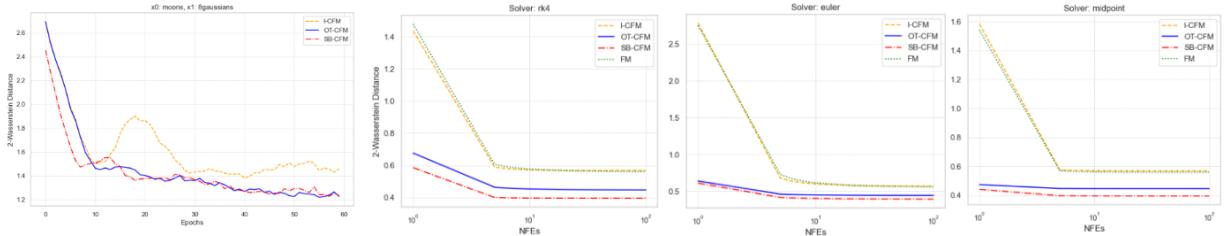


Figure 1: The difference in how fast the Optimal Plan methods train in terms of validation set error in comparison to ones without (left); Error based on the number of function evaluation (NFEs) for different ODE integrators (right).

Method ↓	Batch Size (b)			Network (b = 128)			
	16	128	512	Wide	Deep	Shallow	Activation
I-CFM	$0.27 \pm 0.03$	$0.26 \pm 0.02$	$0.26 \pm 0.04$	$0.20 \pm 0.01$	$0.20 \pm 0.02$	$0.23 \pm 0.03$	$0.26 \pm 0.02$
OT-CFM	$0.26 \pm 0.00$	$0.20 \pm 0.02$	$0.20 \pm 0.01$	$0.19 \pm 0.02$	$0.23 \pm 0.02$	$0.21 \pm 0.02$	$0.24 \pm 0.02$
SB-CFM	$0.26 \pm 0.05$	$0.19 \pm 0.02$	$0.19 \pm 0.10$	$0.17 \pm 0.02$	$0.19 \pm 0.01$	$0.20 \pm 0.02$	$0.22 \pm 0.03$
LightSB	$0.18 \pm 0.06$	$0.19 \pm 0.03$	$0.19 \pm 0.00$	-	-	-	-

Table 1: 2-Wasserstein distance ( $\downarrow$ ) with different training modifications when training on the *Swiss Roll* dataset, evaluated with 512 validation samples after 20k training steps (for the different network types see Section 5.1). For each of the metrics,  $\mu$  and  $\sigma$  were computed after running the experiments 5 times.

for OT-CFM over the other two methods. However, as discussed in Section 6, we were not able to run the batch size of 128, but other than that we have the same setup as Tong et al. 2023. Because of that, our results differ greatly. Firstly, the batch size has a massive impact on the FID score, where it can be concluded that by increasing the batch size from 16 to 64, it improves the model training enormously. Furthermore, just like it was discussed in Section 5.1, OT plan does not work as well as I-CFM with a smaller batch size.

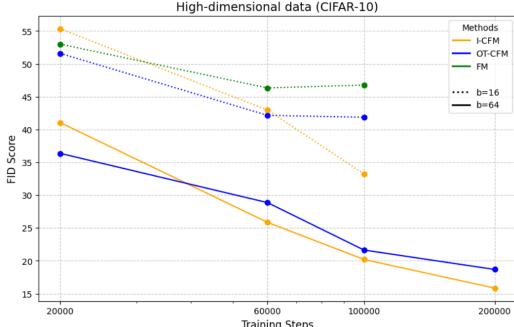


Figure 2: FID scores on CIFAR-10 across training steps for different

Dataset ↓	Mean Distance	Std Dev
gaussian → 8gaussians	0.655	$\pm 0.049$
moons → 8gaussians	1.025	$\pm 0.025$
gaussian → moons	0.472	$\pm 0.048$
gaussian → scurve	0.579	$\pm 0.051$

Table 2: FM methods trained with different batch sizes Schrödinger bridge flow comparison, showing average error over flow time to ground truth averaged over 5 models for SB-CFM

### 5.3 Celeb-A Image Translation

We replicated the experiment 5.4 described in the original study, utilizing a custom-trained VAE to encode images from the Celeb-A into a 128-dimensional latent space. This approach allowed us to perform unsupervised translation by learning a flow in the VAE-latent space, that maps between the embeddings of images with and without specific attributes. After the CNF is learned, we push forward a held-out set of negative vectors by the CNF and compare them to the held-out positive vectors and vice versa. Using maximum mean discrepancy (MMD) with a broad Gaussian kernel as the metric for divergence, our results in terms of MMD closely aligned with those reported in the original experiment (see Table 3). We can see that OT-CFM discovers a better mapping than other methods.

The training for the flow-matching was done on a CPU with 2 cores (from Colab). Each flow matching run took on average 3 hours.

Algorithm ↓	$\sigma = 0.1$	$\sigma = 0.3$	$\sigma = 1$
Identity	0.0117726	0.0117726	0.0117726
I-CFM	0.0085874	0.0057337	0.0018896
OT-CFM	0.0066700	0.0052683	0.0013703

Table 3: MMD values between target and transformed source samples of CelebA latent vectors for each algorithm and  $\sigma$  value. ‘Identity’ refers to performing no translation and treating source samples as approximate samples from the target.

## 5.4 Light Schrödinger Bridge Experiments

This section dives into our extensions of the experiments with LightSB.

### 5.4.1 FID scores on CIFAR-10 (Experiment 5.3)

Looking at the results in Table 4 (CIFAR-10), we can deduce that the images reconstructed from the VAE and the images sampled from the learned light SB latent space still differ significantly in their general distribution. That might be due to the small amount of images generated (640, so 10 batches of 64). Still, the scores give an indication that the light SB model learned the latent representation to a satisfactory degree.

Furthermore, we can observe that higher epsilon values do have a notable impact on the FID score, where smaller epsilon values lead to stronger regularization and thus a more stringent adherence of the trajectories to the optimal transport map minimizing the 2-Wasserstein distance. Hence, a higher epsilon leads to more stochasticity in the process which might have led to a different minimum in this case that is further away from the target distribution.

### 5.4.2 Celeb-A Image Translation in Latent Space (Exp. 5.4)

The results in Table 4 (CelebA) show how the choice of volatility  $\epsilon$  and the number of potentials  $N_{potentials}$  influences the performance and training time of the Light Schrödinger Bridge (Light SB) solver on the CelebA image translation task. Using fewer potentials (e.g.,  $N_{potentials} = 2$ ) generally reduces training time but somewhat increases the MMD compared to the 10-potentials setup. This suggests that more potentials can improve the quality of the learned transport map, but at the expense of longer training times. Notably, both  $\epsilon = 0.1$  and a very small  $\epsilon = 2e^{-3}$  yield similar MMD values (around  $1.109 \times 10^{-3}$ ), indicating that some flexibility in choosing  $\epsilon$  exists without severely compromising performance.

On the other hand, increasing  $N_{potentials}$  to 100 does not yield a substantial improvement over  $N_{potentials} = 10$ , with results remaining around  $1.12 \times 10^{-3}$  MMD. However, training time escalates significantly for  $N_{potentials} = 100$ , ballooning to nearly 37 minutes, compared to about 3-4 minutes for  $N_{potentials} = 10$ . Thus, the results suggest a trade-off: while additional potentials and tuning  $\epsilon$  can marginally improve the learned transport map, doing so may come with diminishing returns in performance relative to the increased computational cost.

Compared to OT-CFM, which typically takes around three hours of training, the Light Schrödinger Bridge approach achieves a remarkable reduction in computation time—down to a few minutes. This approximately 60-fold speedup can be attributed to the fact that Light SB does not require iterating through large mini-batches and optimizing a complex neural network to directly model the transport velocity field. Instead, it learns compact parameterizations of potentials governing the transport cost, allowing for efficient simulation-free optimization that circumvents the heavy overhead of expensive gradient calculations inherent in the OT-CFM framework.

$\epsilon$	$N_{potentials}$	CIFAR-10		CelebA	
		FID	Training time (min)	MMD (1e-3)	Training time (min)
0.1	10	37.649	3:00	$1.109 \pm 0.04$	3:44
10	10	41.369	2:53	$1.707 \pm 0.05$	4:35
0.001	10	37.127	2:56	$1.109 \pm 0.01$	3:30
0.1	2	37.662	1:28	$1.353 \pm 0.04$	1:53
0.1	100	38.306	1:00:41	$1.12 \pm 0.03$	36:59

Table 4: Ablation study: **(CIFAR-10)** Comparison of different epsilon values and numbers of potentials on the FID score. **(CelebA)** Comparison of different epsilon values and number of sources on Mean Maximum Discrepancy between target distribution and  $\pi(x_1|x_0)$  obtained from LSB for the ‘smiling’ attribute, 3 trials. Training on CPU (2 cores).

## 6 Challenges

The following section summarizes the challenges that we faced when trying to reimplement the methods and when we tried to reproduce the results for different experiments.

### 6.1 Toy Datasets

When trying to reproduce the results from "Low-dimensional data: Optimal transport and faster convergence" in Tong et al. 2023, it was hard to directly reproduce the results from their Figure 2, and in the end, we only managed to replicate them to an extent (see Section 5.1). The confusion arose because of the Optimal Transport plan, because on one hand the trajectories with the OT plan were as expected, but the metrics did not improve over other methods (especially for the datasets with the Gaussian as the source distribution). After running the authors code without seeing any difference, we concluded that the methods work properly and differences might be due to hyperparameter differences.

Another problem we faced was concerning SB-CFM and the use of the Sinkhorn method when computing the plan. The library we were using, was *POT: Python Optimal Transport*, and the problem was that we were getting numerical errors with the calculations, which resulted in some values collapsing to 0. We tried amending the cost matrix and changing the regularization hyperparameter, but we were not able to remove the error. However, the authors specified in their code that they ended up using the Earth's Movers Distance method (the same one as for OT-CFM) for SB-CFM, as they "found this to perform better (more accurate and faster) in practice for reasonable batch sizes" (Tong et al. 2023).

### 6.2 CIFAR-10 dataset

Unlike the aforementioned problem, in this case, the experiment was clearly defined by the authors, but the main issue was that we did not have the computation resources to run it with the same batch size and for the same length as they did. We had to scale down the experiment to a batch size of 16 (rather than 128) and run it only for 100,000 training steps, instead of 500,000. The resulting metrics were less than satisfactory, and hence, we tried to repeat the experiment on a dedicated GPU with batch size 64 and 200,000 training steps, but we still were not able to match their performance. Given the time frame and the available resources, rising to their proposed batch size would have taken over 50 hours per 100,000 training steps per model, which is unfeasible in our case.

Running the CIFAR experiments with Light SB posed the challenge of converting the images into a digestible form. In the end, we used VAE embeddings based on weights from a paper that scored fairly well on CIFAR image generation (the reference can be found in the relevant notebook). We ended up training light SB to learn the latent representation of the VAE and compare 600 images reconstructed from the original VAE latent space and the same amount sampled from the learnt light SB representation. Still, using these weights naturally lead to a lower FID score due to the lowered quality within the VAE's latent space.

### 6.3 Celeb-A Image Translation

In contrast to the original paper's large convolutional VAE (7M parameters) trained on 128×128 CelebA faces for 5000 batches of size 256, we implemented a smaller architecture using 64×64 images, 70% of the dataset, 10 epochs, and batch size 128. Training was conducted on a default Colab CPU (2 cores) environment, with each flow matching model (I-CFM and OT-CFM) requiring approximately three hours to train. Due to these computational constraints, experiments focused solely on the "smiling" attribute rather than multiple attributes. While this setup reduced computational demands and training time at the cost of VAE representational capacity, the reliability of results remains intact through MMD baseline comparisons in Table 3.

## 7 Conclusion

In this report, we reproduced several key experiments from the original OT-CFM paper, verifying that conditional flow matching with optimal transport constraints can yield accurate and high-quality generative models. By carefully tuning hyperparameters and adjusting the experimental setup, we were able to achieve results reasonably close to those presented, despite operating with more limited computational resources. Our ablation studies on toy datasets, CIFAR-10 and Celeb A datasets showed that training with fewer steps or smaller batch sizes still captures the essential behavior and merits of OT-CFM models, albeit with somewhat reduced performance.

Most notably, incorporating the Light Schrödinger Bridge (Light SB) approach proved to be a major asset, substantially reducing training times from hours to minutes while still delivering results comparable to those obtained via OT-CFM. The compact parameterization and simulation-free optimization inherent in Light SB offer a practical solution to the

previously time-consuming ODE integrations during inference and gradient computations during training. As a result, Light SB can make OT-based flow matching methods more accessible and scalable in real-world applications, enabling efficient experimentation and deployment without sacrificing generative performance.

## 8 Ethical consideration, Societal impact, Alignment with UN SDG targets

As for the ethical concerns of this research, there is very little concern for any form of data breaches or similar issues, since only toy datasets were used in the experiments. A research ethical concern that may arise in any work of this kind is the question of plagiarism, which should be taken care of by the authors clearly stating what sources their work is based on and to which degree.

Looking at the paper's societal impacts, OT-CFM is an improvement of generative modeling which has several applications in generating art, text and so forth. Looking specifically at the single-cell interpolation experiment, one could use this to study the progression of cancer cells in the human body, or how single cells respond to specific drugs. Hence, there is potential for exploring medical applications of Flow Matching which can impact the healthcare system significantly which is in line with 'SDG 3:Good Health and Well-being'.

## 9 Limitations

Generally speaking, the paper is well-written, theoretically sound and fairly easy to reproduce due to an elaborate code base that is provided by the authors. Furthermore, most hyperparameters are mentioned in the appendix of the paper. Notably, the authors mention a key limitation of CFM which is the fact that it is impossible to regularize  $u_t(x)$  based on prior information which is needed in some use cases. Furthermore, mini-batch approximation of OT can incur errors in high dimensions.

One limitation in reproducing the results lies in the computational cost that is involved in solving the ODEs is quite high for high dimensional data. Seeing the runtime that is needed for CIFAR-10 data that is comparatively low-dimensional due to a smaller image size, one might question the applicability to data of much higher dimensionality. Still, this can be argued for generative models in general.

## 10 Self Assessment

We not only reproduced the main experiments of the original paper but also achieved results that were closely aligned with the paper's. Through careful analysis of hyperparameters, architectural adjustments, batch size constraints, and computational resource trade-offs, we validated the reproducibility of both OT-CFM and SB-CFM in low-dimensional settings, as well as more complex tasks like CIFAR-10 image generation and CelebA latent translation. Moreover, we brought a novel perspective by integrating the Light Schrödinger Bridge solver, demonstrating that it significantly accelerates training/inference time without sacrificing performance, thus illustrating a meaningful and useful extension beyond the original study.

By combining thorough reproduction of experiments, careful consideration of best practices in machine learning, and a substantial extension leveraging new methodologies, we have surpassed the criteria for merely reproducing the results. Our work not only confirms the original findings, but also provides a deeper understanding of the factors influencing model performance and scalability. Consequently, we strongly believe our efforts deserve the highest grade category (A), reflecting our success in going beyond the paper's scope and advancing the research field with improved efficiency, broader applicability, and refined experimental insights.

## References

- Korotin, Alexander, Nikita Gushchin, and Evgeny Burnaev (2023). "Light Schrodinger Bridge". In: *arXiv preprint arXiv:2310.01174*.
- Lipman, Yaron et al. (2022). "Flow matching for generative modeling". In: *arXiv preprint arXiv:2210.02747*.
- Liu, Ziwei et al. (Dec. 2015). "Deep Learning Face Attributes in the Wild". In: *Proceedings of International Conference on Computer Vision (ICCV)*.
- Tong, Alexander et al. (2023). "Improving and generalizing flow-based generative models with minibatch optimal transport". In: *arXiv preprint arXiv:2302.00482*.

## A Visualized trajectories on toy datasets

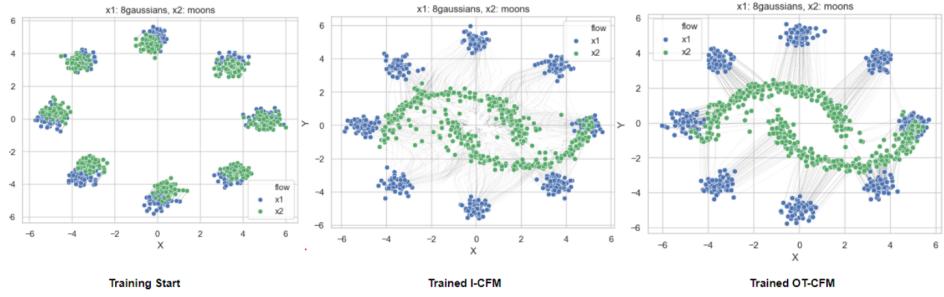


Figure 3: Difference in generated trajectories without (middle) and with (right) an Optimal Transport plan

## B Extended Results

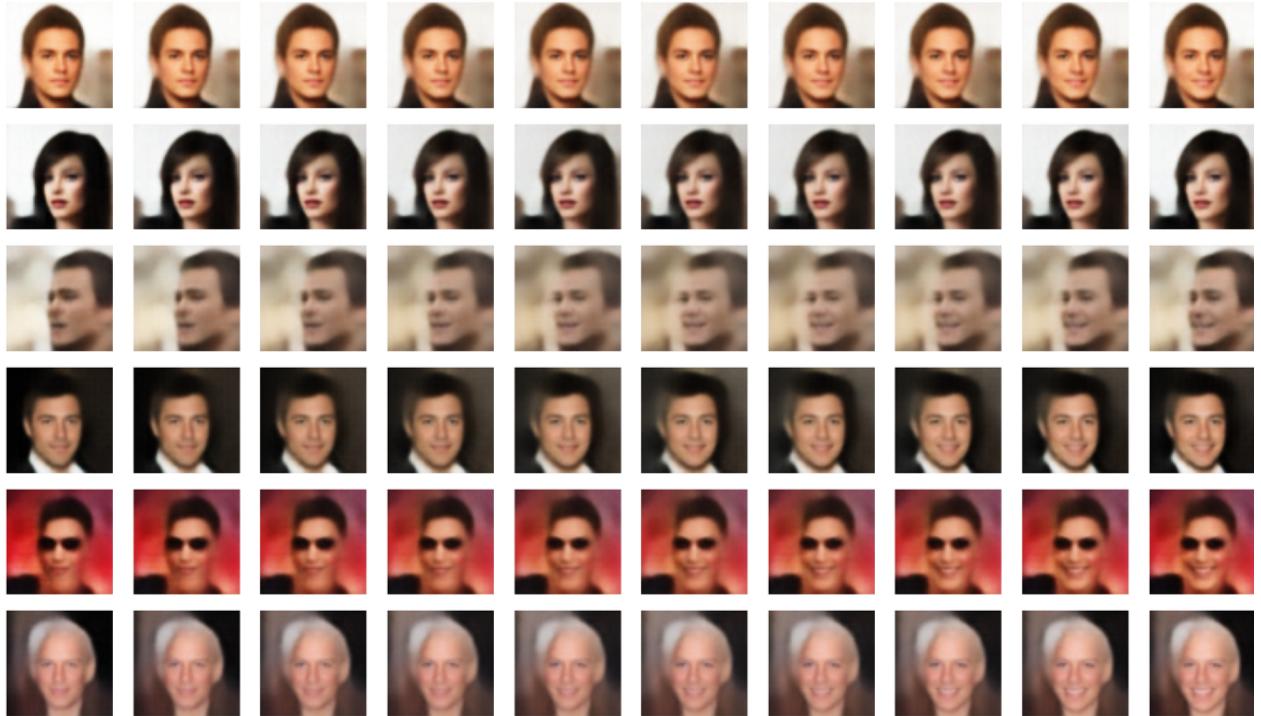


Figure 4: Samples of results for ICFM with  $\sigma = 1$

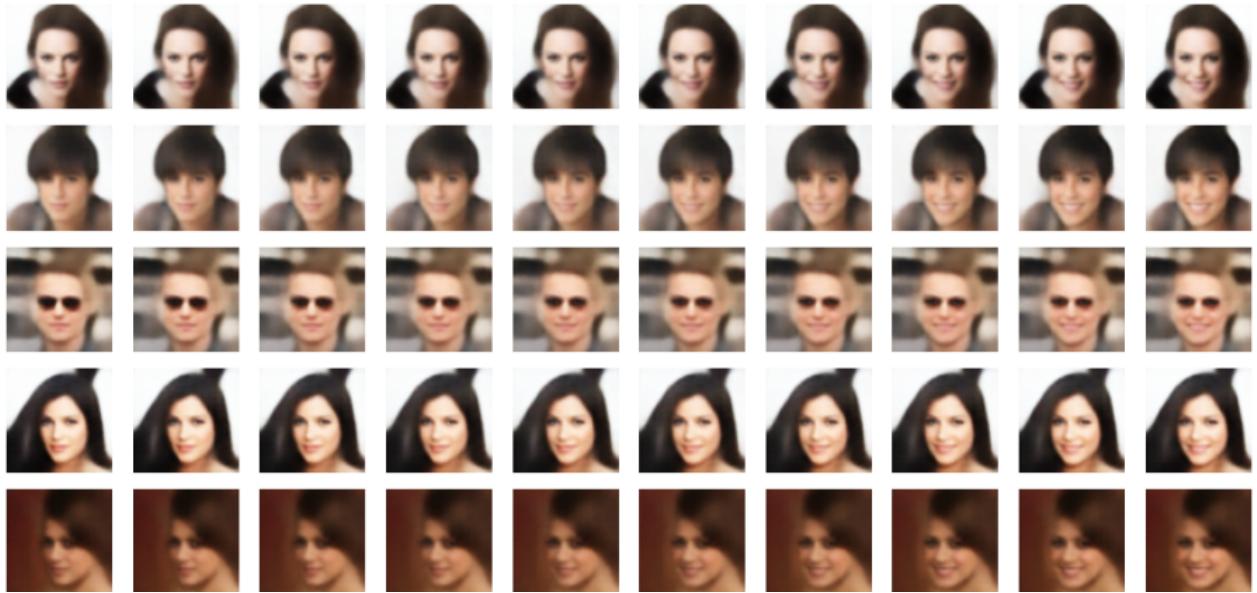


Figure 5: Samples of results for OT CFM with  $\sigma = 1$

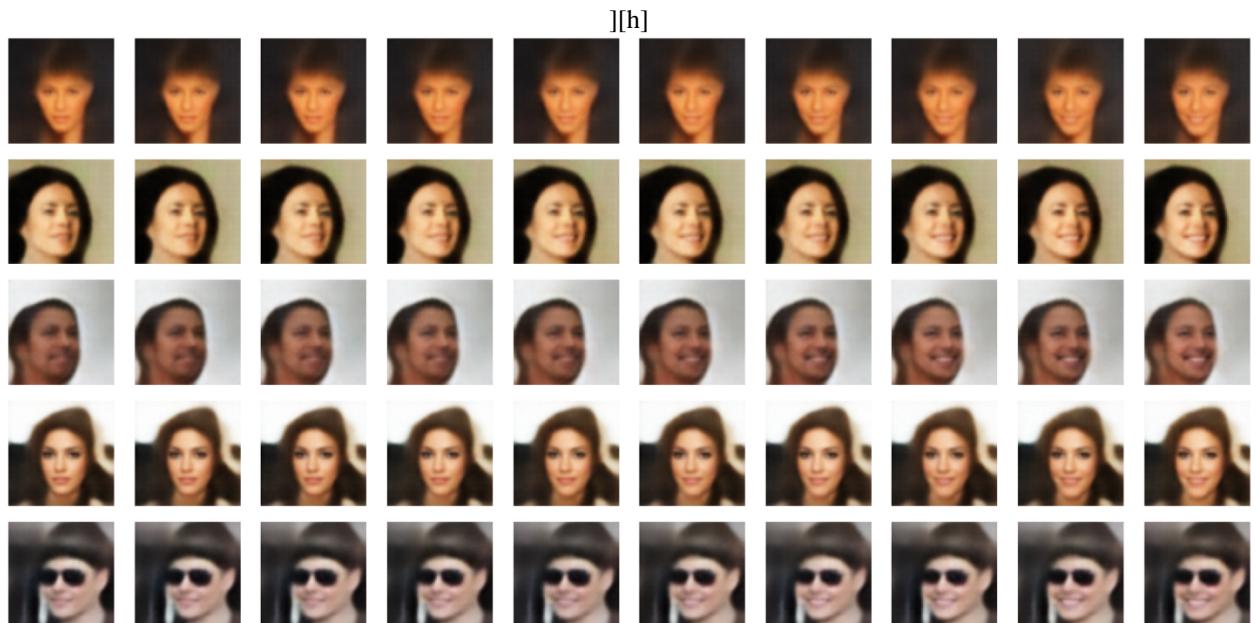


Figure 6: Samples of results for LSB with  $\epsilon = 0.01$ ,  $N_{pots} = 10$