STREAM-LEVEL FLOW MATCHING FROM A BAYESIAN DECISION THEORETIC PERSPECTIVE

Anonymous authors

Paper under double-blind review

ABSTRACT

Flow matching (FM) is a family of training algorithms for fitting continuous normalizing flows (CNFs). Conditional flow matching (CFM) exploits the fact that the marginal vector field of a CNF can be learned by fitting least-square regression to the so-called conditional vector field specified given one or both ends of the flow path. We show that viewing CFM training from a Bayesian decision theoretic perspective on parameter estimation opens the door to generalizations of CFM algorithms. We propose one such extension by introducing a CFM algorithm based on defining conditional probability paths given what we refer to as "streams," instances of latent stochastic paths that connect pairs of noise and observed data. Further, we advocate the modeling of these latent streams using Gaussian processes (GPs). The unique distributional properties of GPs, and in particular the fact that the velocity of a GP is still a GP, allows drawing samples from the resulting stream-augmented conditional probability path without simulating the actual streams, and hence the "simulation-free" nature of CFM training is preserved. We show that this generalization of the CFM can substantially reduce the variance in the estimated marginal vector field at a moderate computational cost, thereby improving the quality of the generated samples under common metrics. Additionally, we show that adopting the GP on the streams allows for flexibly linking multiple related training data points (e.g., time series) and incorporating additional prior information. We empirically validate our claim through both simulations and applications to image and neural time series data.

031 032

033

034

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

028

029

1 INTRODUCTION

Deep generative models aim to estimate and sample from an unknown probability distribution. Continuous normalizing flows (CNFs, Chen et al. (2018)) construct an invertible and differentiable map-037 ping, using neural ordinary differential equation (ODE), between a source and the target distribution. However, traditionally, it has been difficult to scale CNF training to large datasets (Chen et al., 2018; Grathwohl et al., 2019; Onken et al., 2021). Recently, Lipman et al. (2023); Albergo & Vanden-040 Eijnden (2023); Liu et al. (2023b) showed that CNFs can be trained via a regression objective, and proposed the flow matching (FM) algorithm. The FM exploits the fact that the marginal vector field 041 inducing a desired CNF can be learned through a regression formulation, approximating per-sample 042 conditional vector fields using a smoother such as a deep neural network (Lipman et al., 2023). In 043 the original FM approach, the training objective is conditioned on samples from the target distribu-044 tion, and the source distribution has to be Gaussian. This limitation was later relaxed, allowing the 045 target distribution to be supported on manifolds (Chen & Lipman, 2024) and the source distribution 046 to be non-Gaussian (Pooladian et al., 2023). Tong et al. (2024a) provided a unifying framework 047 with arbitrary transport maps by conditioning on both ends. While their framework is, in principle, 048 general, it does require the induced conditional probability paths be readily sampled from, and as such they considered a few Gaussian probability paths. Moreover, most existing FM methods only consider the inclusion of two endpoints, and hence cannot accommodate data involving multiple 051 related observations, such as time series and other data with a grouping structure. Recently, Albergo et al. (2024) proposed a multimarginal stochastic interpolants, which can jointly learn a multivariate 052 distribution by generalization of the stochastic interpolant framework. However, this approach may be too restrictive for path design, and hence may limit it usage to such as time series data.

054 In this paper, we first view FM from the perspective of Bayesian estimation under squared error loss, 055 which motivates us to go one level deeper in Bayesian hierarchical modeling and specify distribu-056 tional assumptions on *streams*, which are latent stochastic paths connecting the two endpoints. This 057 leads to a class of CFM algorithms that conditions at the "stream" level, which broadens the range 058 of conditional probability paths allowed in CFM training. By endowing the streams with Gaussian process (GP) distributions, these algorithms provide wider sampling coverage over the support of 059 the marginal vector field, leading to reduced variance in the estimated vector field and improved 060 synthetic samples from the target distribution. Furthermore, conditioning on GP streams allows for 061 flexible integration of related observations through placing them along the streams between two end-062 points and for incorporating additional prior information, all while maintaining analytical tractability 063 and computational efficiency of CFM algorithms. 064

065

066

067

068

069

070

071

073

075 076

077

078 079

080

081

082

084

092 093

103 104 105 In summary, the main contributions of this paper are:

- 1. We present a Bayesian decision theoretic perspective on FM algorithms, which provides an additional justification for FM algorithms beyond gradient matching and serves as the foundation for extensions to these algorithms by latent variable modeling on the streams.
- 2. We generalize CFM training by augmenting the specification of conditional probability paths through latent variable modeling on the streams. We show that streams endowed with GP distributions lead to a simple stream-level CFM algorithm that preserves the "simulation-free" training.
 - 3. We demonstrate that adjusting the GP streams can reduce the variance of the estimated marginal vector field with moderate computational cost.
 - 4. We demonstrate how to use GP streams to integrate related observations, thereby taking advantage of the correlation among related samples to enhance the quality of the generated samples from the target distributions.
 - 5. These benefits are illustrated by simulations and applications to image (CIFAR-10, MNIST and HWD+) and neural time series data (LFP), with code for Python implementation in the supplementary materials.

2 A BAYESIAN DECISION THEORETIC PERSPECTIVE ON FLOW MATCHING

We start by viewing FM training from a Bayesian decision theoretic perspective. Consider i.i.d. training observations from an unknown population distribution q_1 over \mathbb{R}^d . A CNF is a timedependent differomorphic map ϕ_t that transforms a random variable $x_0 \in \mathbb{R}^d$ from a source distribution q_0 into a random variable from q_1 . The CNF induces a distribution of $x_t = \phi_t(x_0)$ at each time t, which is denoted by p_t , thereby forming a probability path $\{p_t : 0 \le t \le 1\}$. This probability path should (at least approximately) satisfy the boundary conditions $p_0 = q_0$ and $p_1 = q_1$. It is related to the flow map through the change-of-variable formula or the push-forward equation

$$p_t = [\phi_t]_* p_0$$

FM aims at learning the corresponding vector field $u_t(x)$, which induces the probability path over time by satisfying the continuity equation (Villani, 2008).

The key observation underlying FM algorithms is that the vector field $u_t(x)$ can be written as a conditional expectation involving a conditional vector field $u_t(x|z)$, which induces a conditional probability path $p_t(\cdot|z)$ corresponding to the conditional distribution of $\phi_t(x)$ given z. Here, z is the conditioning latent variable, which can be the target sample x_1 (e.g. Ho et al. (2020); Song et al. (2021); Lipman et al. (2023),) or a pair of (x_0, x_1) on source and target distribution (e.g. Liu et al. (2023b); Tong et al. (2024a)). Specifically, Tong et al. (2024a), generalizing the result from Lipman et al. (2023), showed that

$$u_t(x) = \int u_t(x|z) \frac{p_t(x|z)q(z)}{p_t(x)} dz = \mathbb{E} \left(u_t(x|z) | x_t = x \right),$$

where the expectation is taken over z, which one can recognize is the conditional expectation of $u_t(x|z)$ conditional on the event that $x_t = x$. The integral is with respect to the conditional distribution of z given $x_t = x$.

108 It is well-known in Bayesian decision theory (Berger, 1985) that under squared error loss, the 109 Bayesian estimator, which minimizes both the posterior expected loss (which conditions on the 110 data and integrates out the parameters) and the marginal loss (which integrates out both the param-111 eters and the data), is exactly the posterior expectation of that parameter. This implies immediately 112 that if one considers the conditional vector field $u_t(x|z)$ as the target of "estimation", and the corresponding "data" being the event that $x_t = x$, i.e., that the path goes through x at time t, then the 113 corresponding Bayes estimate for $u_t(x|z)$ will be exactly the marginal vector field $u_t(x)$, as it is now 114 the "posterior mean" of $u_t(x|z)$. We emphasize again that here the "data" differs from the actual 115 training and the generated noise observations, which in fact help form the "prior" distribution. 116

The FM algorithm is motivated from the goal of approximating the marginal vector field $u_t(x)$ through a smoother v_t^{θ} (typically a neural network), via the flow matching (FM) objective

$$\mathcal{L}_{\mathrm{FM}}(\theta) = \mathbb{E}_{t \sim U(0,1), x \sim p_t(x)} \| v_t^{\theta}(x) - u_t(x) \|^2$$

which is not identifiable due to the non-uniqueness of the marginal vector fields that satisfy the boundary conditions without further constraints. In the following, we presume $t \sim U(0, 1)$ and only show random variables to save notations. FM algorithms address this by fitting v_t^{θ} to the conditional vector field $u_t(x|z)$ after further specifying the distribution of q(z) along with the conditional probability path $p_t(x|z)$, through minimizing the finite-sample version of the marginal squared error loss. This approach was referred to as the conditional flow matching (CFM) objective

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

Traditionally, optimizing the CFM objective is justified because it has the same gradients w.r.t. θ to the corresponding FM loss (Lipman et al., 2023; Tong et al., 2024a). The Bayesian decisiontheoretic perspective provides a further validation because approximating the conditional vector field by minimizing the marginal squared error loss can be interpreted as approximating the "posterior expectation" of $u_t(x|z)$, which is exactly $u_t(x)$.

Moreover, this is true for any coherently specified probability model q(z). So long as the conditional probability path $p_t(x|z)$ is tractable, a suitable CFM algorithm can be designed. Therefore one can enrich the specification of q(z) using Bayesian latent variable modeling strategies. This motivates us to generalize CFM training to the stream level, which we describe in the next section.

138 139

140

141

151

3 STREAM-LEVEL FLOW MATCHING

3.1 A PER-STREAM PERSPECTIVE ON FLOW MATCHING

A stream s is a stochastic process $s = \{s_t : 0 \le t \le 1\}$, where each s_t is a random variable in the sample space of the training data. We focus on streams connecting one end x_0 in the source to the other x_1 in the training data. From here on, s will take the space of the latent quantity z.

Instead of defining a conditional probability path and vector field given one endpoint at t = 1(Lipman et al., 2023) or two endpoints at t = 0 and 1 (Tong et al., 2024a), we shall consider defining it given the whole stream connecting the two ends. In order to achieve this, we need to specify a probability model for *s*. This can be separated into two parts—the marginal model on the endpoints $\pi(x_0, x_1)$ and the conditional model for *s* given the two ends. That is

$$(x_0, x_1) \sim \pi$$
 and $s|_{s_0} = x_0, s_1 = x_1 \sim p_s(\cdot|x_0, x_1).$

¹⁵² Our model and algorithm will generally apply to any choice of π that satisfies the boundary condition, including all of the examples considered in Tong et al. (2024a). We defer the description of specific choices of $p_s(\cdot|x_0, x_1)$ to the next section and for now focus on the general framework.

Given a stream s, the "per-stream" vector field $u_t(x|s)$ represents the "velocity" (or derivative) of the stream at time t, conditional on the event that $s_t = x$, i.e, the stream s passes through x at time t. Assuming that the stream is differentiable with in time, the per-stream vector field is

$$u_t(x|\mathbf{s}) := \dot{s}_t = \mathrm{d}s_t/\mathrm{d}t,$$

which is defined only on all pairs of (t, x) that satisfy $s_t = x$. The per-stream view extends previous CFM conditioning on endpoints and provides more flexibility. See Appendix A for more detailed discussion on how the per-stream perspective relates to the per-sample perspective on FM. While the endpoint of the stream $s_1 = x_1$ is an actual observation in the training data, for the task of learning the marginal vector field $u_t(x)$, one can think of our "data" as the event that a stream spasses through a point x at time t, that is $s_t = x$. Under the squared error loss, the Bayes estimate for the per-stream conditional vector field $u_t(x|s)$ will be the "posterior" expectation given the "data", which is exactly the marginal vector field

$$u_t(x) = \mathbb{E}(u_t(x|\mathbf{s})|s_t = x) = \mathbb{E}(\dot{s}_t|s_t = x).$$

$$\tag{1}$$

Following Theorem 3.1 in Tong et al. (2024a), we can show that the marginal vector $u_t(x)$ indeed generates the probability path $p_t(x)$. (See the proof in the Appendix H.1.) The essence of the proof is to check the continuity equation for the (degenerate) conditional probability path $p_t(x|s)$.

A general stream-level CFM loss for learning $u_t(x)$ is then 173

167 168

174

182 183

185 186

191

193

$$\mathcal{L}_{\text{sCFM}}(\theta) = \mathbb{E}_{t,\boldsymbol{s}} \| v_t^{\theta}(s_t) - u_t(\boldsymbol{x}|\boldsymbol{s}) \|^2 = \mathbb{E}_{t,\boldsymbol{s}} \| v_t^{\theta}(s_t) - \dot{s}_t \|^2$$

where the integration over t is again U(0, 1) and that over s is with respect to the marginal distribution of s induced by $\pi(x_0, x_1)$ and $p_s(\cdot|x_0, x_1)$. As in previous research such as Lipman et al. (2023); Tong et al. (2024a), we can show that the gradient of \mathcal{L}_{sCFM} equals that of \mathcal{L}_{FM} with details of proof in Appendix H.2. However, the stream-level CFM can be justified from a Bayesian decision theoretic perspective without gradient matching (Section 2). Because the (population-level) minimizer for the sCFM loss is $u_t(x)$, minimizing the sCFM loss provides a reasonable estimate for the marginal vector field $u_t(x)$. To see this, rewrite the sCFM loss by the law of iterated expectation as

$$\mathcal{L}_{\text{sCFM}}(\theta) = \mathbb{E}_t \mathbb{E}_s \left(\| v_t^{\theta}(s_t) - \dot{s}_t \|^2 | t \right)$$

184 The inner expectation can be further written in terms of another iterated expection:

$$\mathbb{E}_{\boldsymbol{s}}\left(\|v_t^{\theta}(s_t) - \dot{s}_t\|^2 | t\right) = \mathbb{E}_{s_t} \mathbb{E}_{\boldsymbol{s}}\left(\|v_t^{\theta}(s_t) - \dot{s}_t\|^2 | t, s_t\right).$$

For any x, $\mathbb{E}_{s}(||v_{t}^{\theta}(s_{t}) - \dot{s}_{t}||^{2}|t, s_{t} = x) = \mathbb{E}_{s}(||v_{t}^{\theta}(x) - \dot{s}_{t}||^{2}|t, s_{t} = x)$, whose minimizer is the conditional expectation of \dot{s}_{t} given $s_{t} = x$, which is exactly $u_{t}(x)$. Hence, one can estimate $u_{t}(x)$ by minimizing $\mathcal{L}_{sCFM}(\theta)$. This justifies training $u_{t}(x)$ through the sCFM loss without regard to any specific optimization strategy.

192 3.2 CHOICE OF THE STREAM MODEL

Next, we specify the conditional model for the stream given the endpoints $p_s(\cdot|x_0, x_1)$. This model 194 should emit streams differentiable with respect to time, with readily available velocity (either ana-195 lytically or easily computable). Previous methods such as optimal transport (OT) conditional path 196 (Liu et al., 2023b; Lipman et al., 2023; Tong et al., 2024a) can provide rather poor coverage of 197 the (t, x) space, resulting in extensive extrapolation of the estimated vector field $v_t(x)$. It is thus 198 desirable to consider stochastic models for the streams that ensure the smoothness while allowing 199 streams to diverge and provide more spread-out coverage of the (t, x) space. Previous research sug-200 gested that, compared to ODEs, SDEs (Ho et al., 2020; Song et al., 2021) can be more robust in 201 high-dimensional spaces (Tong et al., 2024b; Shi et al., 2023; Liu et al., 2023a), likely due to the robustness arising from the regularization induced by the additional stochasticity. 202

203 To preserve the "simulation-free" nature of CFM, we consider models where the joint distribution 204 of the stream and its velocity is available in closed form. In particular, we further explore the 205 streams following Gaussian processes (GPs). A desirable property of GP is that its velocity is 206 also a GP, with mean and covariance directly derived from original GP (Rasmussen & Williams, 207 2005). This enables efficient joint sampling of (s_t, \dot{s}_t) given observations from a GP in streamlevel CFM training. By adjusting covariance kernels for the joint GP, one can fine-tune the variance 208 level to control the level of regularization, thereby further improving the estimation of the marginal 209 vector field $u_t(x)$ (Section 4.1). The prior path constraints can also be incorporated into the kernel 210 design. Additionally, GP conditioning on the event that the stream passes through a finite number of 211 intermediate locations between two endpoints again leads to a GP with analytic mean and covariance 212 kernel (Section 4.2). This is particularly useful for incorporating multiple related observations. 213

214 Specifically, given M time points $t = (t_1, t_2, ..., t_M)$ with $t_1 = 0$ and $t_M = 1$, we let $s_t = (s_{t_1}, s_{t_2}, ..., s_{t_M})$, and consider a more general conditional model for $p_s(\cdot | s_t = x_{obs})$, where $x_{obs} = (x_{t_1}, x_{t_2}, ..., x_{t_M})$ are a set of "observed values" that we require the statistic process s to

216 pass through at time (t_1, t_2, \ldots, t_M) . Note that this contains the special case of conditioning on 217 two endpoints (i.e., M = 2) described in Section 3.1. We consider a more general construction 218 for $M \ge 2$ because later we will use this to incorporate multiple related observations (such as time 219 series or other measurements from the same subject).

220 We construct a conditional (multi-output) GP for s that (approximately) satisfies the boundary con-221 ditions, with differentiable mean function m and covariance kernel k_{11} . Since the derivative of a GP 222 is also a GP, the joint distribution of s and corresponding velocity process $\dot{s} := {\dot{s}_t : t \in [0,1]}$ given s_t is also a GP, with the mean function for \dot{s} be $\dot{m}(t) = dm(t)/dt$ and kernels defined by 224 derivatives of k_{11} . To facilitate the construction of this GP, we consider an auxiliary GP on s with 225 differentiable mean function ξ and covariance kernel c_{11} . Using the property that the conditional 226 distribution of Gaussian remains Gaussian, we can obtain a joint GP model on $(s, \dot{s}) \mid s_t$, which satisfies the boundary conditions. For computational efficiency and ease of implementation, we 227 assume independence of the GP across dimensions of sNotably, while we are modeling streams 228 conditionally given s_t as a GP, the marginal (i.e., unconditional) distribution of s at all time points 229 are allowed to be non-Gaussian, which is necessary for satisfying the boundary condition and for the 230 needed flexibility to model complex distributions. The detailed derivation can be found in Appendix 231 B, and the training algorithm for GP-CFM is summarized in Algorithm 1. 232

Algorithm 1: Gaussian Process Conditional Flow Matching (GP-CFM)

Input : observation distribution $\pi(x_{obs})$, initial network v^{θ} , and a GP defining the conditional distribution $(s_t, \dot{s}_t) \mid s_t = x_{obs} \sim \mathcal{N}(\tilde{\mu}_t, \tilde{\Sigma}_t)$, for $t \in [0, 1]$. **Output:** fitted vector field $v_t^{\theta}(x)$ while Training do $\boldsymbol{x}_{\text{obs}} \sim \pi(\boldsymbol{x}_{\text{obs}}); t \sim U(0, 1)$ $\begin{aligned} (s_t, \dot{s}_t) \mid \boldsymbol{s_t} = \boldsymbol{x_{\text{obs}}} \sim \mathcal{N}(\tilde{\boldsymbol{\mu}}_t, \tilde{\Sigma}_t) \\ \mathcal{L}_{\text{sCFM}}(\boldsymbol{\theta}) \leftarrow \| \boldsymbol{v}_t^{\boldsymbol{\theta}}(s_t) - \dot{s}_t \|^2 \end{aligned}$ $\theta \leftarrow \text{update} (\theta, \nabla_{\theta} \mathcal{L}_{\text{sCFM}}(\theta))$ end

244 245

247

251

233 234

235

236

237

238

239

240

241 242

243

246 Several conditional probability paths considered in previous works are special cases of the general GP representation. For example, if we set $m(t) = tx_1 + (1-t)x_0$ (therefore, $\dot{m}(t) = x_1 - x_0$) and 248 $k_{11}(t,t') = \sigma^2 I_d$, the path reduces to the OT conditional path used in I-CFM with constant variance 249 (Tong et al., 2024a). The I-CFM path can also be induced by conditional GP construction (Appendix 250 B) using a linear kernel for c_{11} , with more details in Appendix C. In the following, we set $\xi(t) = 0$ and use squared exponential (SE) kernel for c_{11} for each dimension (may with additional terms such as in Figure2). The details of SE kernel can be found in Appendix D. 252

253 Probability paths with time-varying variance, such as Song & Ermon (2019); Ho et al. (2020); 254 Lipman et al. (2023), also motivate the adoption of non-stationary GPs whose covariance kernel 255 could vary over t. For example, to encourage samples that display larger deviation from those in 256 the training set (and hence more regularization), one could consider using a kernel producing larger variance as t approaches to ends with finite training samples (Figure 2 and 7). Moreover, because 257 the GP model for s is specified given the two endpoints, both its mean and covariance kernel can 258 be specified as functions of (x_0, x_1) . For example, if x_1 is an outlier of the training data, e.g., from 259 a tail region of q_1 , then one may incorporate a more variable covariance kernel for $p_s(\cdot|x_0, x_1)$ to 260 account for the uncertainty in the "optimal" transport path from x_0 to x_1 . 261

262

263

4 NUMERICAL EXPERIMENTS

264 265

266 In this section, we demonstrate the benefits of GP stream models by several simulation examples. Specifically, we show that using GP stream models can improve the generated sample quality at a 267 moderate cost of training time, through tweaking the variance function to reduce sampling variance 268 of the estimated vector field. Moreover, the GP stream model makes it easy to integrate multiple 269 related observations along the time scale.

287

288

290 291

292

293

295

296

297

298 299

300 301

4.1 Adjusting GP variance for high quality samples

We first show that one can reduce the variance in estimating $u_t(x)$ by incorporating additional stochasticity in the sampling of streams with appropriate GP kernels. As illustrated in Figure 1A, for estimating 2-Gaussian mixtures from standard Gaussian noise, the straight conditional stream used in I-CFM covers a relatively narrow region (gray). For points outside the searching region, there are no "data" and the neural network $v_t^{\theta}(x)$ must be extrapolated. In the sampling stage, this can lead to potential "leaky" or outlying samples that are far from the training observations.

278 For constructing GP conditional streams, we condition on the endpoints but expand the coverage 279 region (red) by tweaking the kernel function (e.g. decrease the SE bandwidth in this case). This provides a layer of protection against extrapolation errors. We then train the I-CFM and GP-I-CFM 280 100 times using a 2-hidden layer multi-layer perceptron (MLP) with 100 training samples at t = 1, 281 and calculate 2-Wasserstein (W2) distance between generated and test samples. For fair comparison, 282 we set $\sigma = 0$ for I-CFM and use noise-free GP-I-CFM. The results are summarized in Table 1B. 283 Empirically, the GP-I-CFM has smaller W2 distance than I-CFM. We further generate 1000 samples 284 and streams for I-CFM and GP-I-CFM with largest W2 distance in Figure 1C, starting with the same 285 points from standard Gaussian. In this example, several outliers are generated from I-CFM. 286



Figure 1: GP streams reduce extrapolation by expanding coverage area. We use a 2-Gaussian 302 mixture distribution as an example. Training observations are shown in red, generated samples in 303 orange, and noise source samples in black. A. FM with straight conditional stream (e.g. I-CFM) may 304 generate "leaky" or outlier samples due to extrapolation errors. The FM method with GP conditional 305 stream has a broader coverage area. B. We train models with I-CFM and GP-I-CFM 100 times and 306 calculate 2-Wasserstein (W2) distance between generated and test samples. Results of 100 seeds 307 are summarized by mean and standard error. C. Among these 100 trained models, generate 1000 308 samples (orange) and streams (blue) for I-CFM and GP-I-CFM with largest W2 distance. 309

310 We can further modify the GP variance function over time to efficiently improve sample quality. 311 Here, we consider the task of estimating and sampling from a 2-Gaussian mixture with 100 training 312 samples at t = 1. For constant noise, diagonal white noise is added to perturb stream locations while 313 retaining the SE kernel. For varying noise, we add a non-stationary dot product kernel to the SE kernel. Specifically, denote the kernel for auxiliary GP on s in dimension i as c_{11}^i , for i = 1, ..., d. Let $c_{11}^i(t, t') = c_{11}^{SE}(t, t') + \alpha tt'$ for increasing variance and $c_{11}^i(t, t') = c_{11}^{SE}(t, t') + \alpha(t-1)(t'-1)$ for decreasing variance, where $\{t, t'\} \in [0, 1]$ and $c_{11}^{SE}(t, t') = \sigma^2 \exp\left(-\frac{(t-t')^2}{2t^2}\right)$. (See Appendix 314 315 316 317 B for additional details.) Some examples of the streams connecting two endpoints under different 318 variance schemes are shown in Figure 2A. We train models 100 times and calculate 2-Wasserstein 319 (W2) distance between generated and test samples, and the results of these 100 seeds are summarized 320 in Table 2B. In this example, with infinite samples at t = 0 (from the standard Gaussian) but only 321 100 samples at t = 1, injecting noise at t = 0 worsens estimation. However, when approaching the target distribution (t = 1), adding noise can improve estimation with small samples (100). This 322 noise perturbs the limited data, encouraging broader exploration and adding regularization to reduce 323 estimation error.



Figure 2: Change variance over time by tweaking the covariance kernel. We revisit the 2-Gaussian mixture distribution as an example. A. Examples of conditional stream between two points, under different variance change scheme. B. We then train models under each variance scheme for 100 times and calculate 2-Wasserstein (W2) distance between generated and test samples for each. The results of 100 seeds are summarized by mean and standard error.

We further consider the transformation between two 2-Gaussian distributions with finite samples (100) at both ends. Results are shown in Appendix E. In this scenario, injecting noise near either endpoint improves estimation.

4.2 INCORPORATING MULTIPLE RELATED TRAINING OBSERVATIONS

Next, we show that GP streams enable the flexible inclusion of related observations along the same
 stream over time. This is particularly useful for generating related samples, such as time series (e.g.,
 videos), where correlations between observations can enhance information sharing and improve
 estimation at each time point.

To illustrate the main idea, we consider 100 paired observations and place the two observations in each pair at t = 0.5 and t = 1 respectively (Figure 3A) while t = 0 still corresponds to a source distribution. Here, we show the generated samples (at t = 0.5 and t = 1) and the corresponding streams for GP-I-CFM and I-CFM. Again, 2-hidden layer MLP is used in this case. The I-CFM strategy employs two separate models with I-CFM algorithms (Figure 3B), whereas GP-I-CFM offers a single unifying model for all observations, resulting in a smooth stream across all time points (Figure 3C).





362

363

364

365 366

367

368

369

370

371

338

339

340

341

342

343 344

348 349

Figure 3: **GP streams can include related points flexibly**. **A**. Paired data with observations on t = 0.5 (red) and t = 1 (orange). **B**. The generated samples (red for t = 0.5 and orange for t = 1) and streams (blue) for I-CFMs. The I-CFMs contain two separate models trained by I-CFM, t = 0 (standard Gaussian noise) to t = 0.5 and t = 0.5 to t = 1. **C**. The generated samples for GP-I-CFM.

378 In some cases the GP streams may not be well separated and thus may confuse the training of 379 the vector field at crossing points. In Figure 4, we show a time series dataset over 3 time points, 380 where training data at t = 0 and t = 1 on one horizontal side while points at t = 0.5 are on 381 the opposite side (Figure 4A). Therefore, these streams have two crossing regions (marked with 382 blue boxes in Figure 4A), where the training of vector field is deteriorated when simply using the GP-I-CFM (Figure 4B). One easy solution is to further condition the neural net $v_t^{\theta}(x)$ on covariate 383 (subject label) c, such that the optimizing objective is $\mathcal{L}_{cCFM} = E_{t \sim U(0,1), s \sim q(s|c)} \|v_t^{\theta}(s_t, c) - \dot{s}_t\|^2$, 384 where $q(s \mid c)$ represents the distribution of s given c. The validity for approximating the covariate-385 dependent vector field using the above optimizing objective is shown in the Appendix H.3. In this 386 example, similar subjects have close starting points at t = 0, and we let $c = x_0$. By conditioning 387 on c (covariate model), the neural net are separated for different subjects, and hence the training of 388 vector field will not be confused (Figure 4C). 389



Figure 4: Further conditioning on the starting points helps with stream generation. A. Paired data with observations on three time points: t = 0 (black), t = 0.5 (red) and t = 1 (orange). The two stream cross regions are marked with light blue square. **B**. The generated samples and streams for GP-I-CFM (without covariate), where the initial points at t = 0 are generated from noise using a separate I-CFM. **C**. The generated samples and streams for GP-I-CFM with covariate using the same starting points, where the neural network is further conditioning on data at t = 0.

408 409 410

411

402

403

404

405

406

407

5 APPLICATIONS

We apply our GP-based CFM methods to two hand-written image datasets (MNIST and HWD+), CIFAR-10 dataset and local field potential (LFP) dataset from mouse brain to illustrate how GPbased algorithms 1) reduce sampling variance (MNIST and CIFAR-10) and 2) flexibly incorporate multiple related observations (e.g. time series data) and generate smooth transformation across different time points (HWD+ and LFP dataset). The reported running times for the experiments are based on results obtained on a server configured with 2 CPUs, 24 GB RAM, and 2 RTXA5000 GPUs.

418 419

5.1 VARIANCE REDUCTION

420

We explore the empirical benefits of variance reduction using FM with GP conditional streams on CIFAR-10 (Krizhevsky, 2009) the MNIST (Deng, 2012) database. For CIFAR-10, we compare performance for I-CFM and GP-I-CFM. Since OT strategy can be complementary to our GP-stream method to enhance the performance, we consider four algorithms in MNIST application: two linear stream models (I-CFM, OT-CFM) and two GP stream models (GP-I-CFM, GP-OT-CFM). Here, we show the details and results for CIFAR-10 application, the results for MNIST application (with more extensive experiments) can be found in the Appendix F.

We perform an experiment on unconditional CIFAR-10 generation (Krizhevsky, 2009) from a standard Gaussian source, using I-CFM and proposed GP-I-CFM, to evaluate the performance in the high-dimensional image setting. We use the similar setup to that of Tong et al. (2024a), such as time-dependent U-Net (Ronneberger et al., 2015; Nichol & Dhariwal, 2021) with 128 channels, a learning rate of 2×10^{-4} , clipping gradient norm to 1 and exponential moving average with a decay

432 of 0.9999. Besides, we add diagonal white noise 10^{-6} in GP-I-CFM, and set $\sigma = 10^{-3}$ in I-CFM 433 for a fair comparison. The samples of s_t is shown Figure 5B. The models are trained for 400,000 434 epochs, with batch size be 128. The I-CFM runs around 3.6 iterations per second, while GP-I-CFM 435 runs around 3.0 iterations per second. The 64 generated images from I-CFM and GP-I-CFM are 436 shown in Figure 5A, using a DOPRI5 adaptive solver. Visually, images generated by GP-I-CFM are generally sharper and exhibit more details compared to those generated by I-CFM (e.g. first row 437 of Figure 5A). The Fréchet inception distance (FID) (Heusel et al., 2017), calculated by clean-fid 438 library (Parmar et al., 2022) with 2000 samples, along training steps are plotted in Figure 5C, show-439 ing that the GP-I-CFM performs better than I-CFM. To show even a more significant performance 440 improvement of GP-I-CFM, we may use a smaller GP bandwidth. However, this will lead to slower 441 convergence, and it should be chosen by the practitioner to balance computational time (number of 442 iterations) and sample quality. 443



Figure 5: Application to CIFAR-10 dataset. Here, we compare the performance of unconditional image generation from standard Gaussian noise for I-CFM and GP-I-CFM. A. 64 generated samples for I-CFM and GP-I-CFM, starting from the same points. B. The samples of $x_t(s)$ from GP-I-CFM between a pair of two endpoints. C. FID of two algorithms over training steps.

5.2 MULTIPLE TRAINING OBSERVATIONS

444

445

446

447

448 449

450

451

452

453 454

455 456

457

458

459

460 461

462 463

Finally, we demonstrate how our GP stream-level CFM can flexibly incorporate related observations (between two endpoints at t = 0 and t = 1) into a single model and provide smooth transformation across different time points, using the HWD+ dataset (Beaulac & Rosenthal, 2022) and LFP dataset (Steinmetz et al., 2019), where LFP dataset is a time series data for mouse brain. Here, we show results for HWD+ dataset; refer to Appendix G for LFP application.

The HWD+ dataset contains images of handwritten digits along with writer IDs and characteristics, 469 which are not available in MNIST dataset used in SectionF. Here, we consider the task of transform-470 ing from "0" (at t = 0) to "8" (at t = 0.5), and then to "6" (at t = 1). The intermediate image, "8", is 471 placed at t = 0.5 (artificial time) for "symmetric" transformations. All three images have the same 472 number of samples, totaling 1,358 samples (1,086 for training and 272 for testing) from 97 subjects. 473 The U-Nets with 32 channels and 1 residual block are used. Both models with and without covariate 474 (using starting images, as in Figure 4C) are considered. Each model is trained both by I-CFM and 475 GP-I-CFM. The I-CFM transformation contains two separate models trained by I-CFM ("0" to "8" 476 and "8" to "6"). Noise-free GP-I-CFM and I-CFM with $\sigma = 0$ are used for fair comparisons. In each 477 training iteration, we randomly select samples within each writer, to preserve the grouping structure of data. The runtime for all algorithms (I-CFM, GP-I-CFM and corresponding labeled versions) 478 are similar, which take 0.74s for passing all training data once. However, since I-CFMs contain 2 479 separated models, the running time is doubled. 480

The traces for 10 generated samples from each algorithm are shown in Figure 6A, where the starting
images ("0" in the first rows) are generated by an I-CFM from standard Gaussian noise. Visually,
the GP-based algorithms generate higher quality images and smoother transformation compared to
algorithms using linear conditional stream (I-CFM), highlighting the benefit of including correlations across different time points. Additionally, the transformation generally looks smoother when
the CFM training is further conditioned on the starting images.

486 We then quantify the performance of different algorithms by calculating the FID for "0", "8" and 487 "6", and plot them over time for each (Figure 6B). For all FIDs, the GP-based algorithms (green 488 & red) outperform their straight connection (I-) counterparts (blue & orange), especially for the 489 FID for "8" at t = 0.5 and the FID to "6" at t = 1. This also holds for the FID for "0", as 490 the GP-based algorithms are unified and the information is shared across all time points. This aligns with the observation by Albergo et al. (2023) that jointly learning multiple distributions better 491 preserves the original image's characteristics during translation. However, for the I-algorithms, 492 the conditional version (orange) performs worse than unconditional one (blue), as conditioning on 493 the starting images makes the stream more separated, requiring more data to achieve comparable 494 performance. In contrast, the data in GP-based algorithms is more efficiently utilized, as correlations 495 across time points for the same subject are integrate into one model. Therefore, explicitly accounting 496 for the grouping effect by conditioning on starting images (red) further improves performance. 497



Figure 6: Application to HWD+ dataset. We fit models for transforming "0" to "8" and then to "6". Both covariate and non-covariate (on starting images) models are considered, and each model is fitted by both I-CFM and GP-I-CFM. The I-CFM transformation consists of two separate models trained by I-CFM ("0" to "8" and "8" to "6"). A. 10 sample traces for the four trained models. The starting images ("0"s in the first row) are generated by an I-CFM from standard Gaussian noise, and all four trained models use the same starting images. **B**. The corresponding FID to "0", "8" and "6" for these four trained models over time.

522 523 524

525 526

527

531

516

517

518

519

520

521

CONCLUSION 6

We have presented a Bayesian decision theoretic perspective to CFM training, which motivates an extension to CFM algorithms based on latent variable modeling. In particular, we adopt GP models 528 on the latent streams. Our GP-CFM algorithm preserves the "simulation-free" feature of CFM 529 training by exploiting distributional properties of GPs. This generalization not only reduces the 530 sampling variance by expanding coverage of the sampling space in CFM training, but also allows easy integration of multiple related observations to achieve borrowing of strength. 532

There are some potential improvements either under GP-CFM frameworks or generally motivated 533 by Bayesian decision theoretic perspective. For example, under GP-CFM framework, current im-534 plementations require the complete observations for all time points, which can be rare in time series 535 applications. To deal with the missingness as well as the potential high-dimensionality of the train-536 ing data, we may fit the GP-CFM in some latent space as in latent diffusion models (Rombach et al., 537 2022) and latent flow matching (Dao et al., 2023). 538

We believe that the Bayesian decision theoretic perspective and GP-CFM generalization so motivated open the door to various further improvements of CFM training of CNFs.

540 REFERENCES 541

548

561

562

563 564

565

566

567

568

569

570

579

581

582

583

- Michael S. Albergo, Nicholas M. Boffi, and Eric Vanden-Eijnden. Stochastic interpolants: A uni-542 fying framework for flows and diffusions, 2023. URL https://arxiv.org/abs/2303. 543 08797. 544
- Michael Samuel Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic 546 interpolants. In The Eleventh International Conference on Learning Representations, 2023. URL 547 https://openreview.net/forum?id=li7qeBbCR1t.
- Michael Samuel Albergo, Nicholas Matthew Boffi, Michael Lindsey, and Eric Vanden-Eijnden. 549 Multimarginal generative modeling with stochastic interpolants. In *The Twelfth International* 550 Conference on Learning Representations, 2024. URL https://openreview.net/forum? 551 id=FHqAzWl2wE. 552
- 553 David G Amaral, Helen E Scharfman, and Pierre Lavenex. The dentate gyrus: fundamental neu-554 roanatomical organization (dentate gyrus for dummies). Progress in brain research, 163:3-22, 2007. ISSN 0079-6123 (Print). doi: 10.1016/S0079-6123(07)63001-5. 555
- 556 Cédric Beaulac and Jeffrey S Rosenthal. Introducing a New High-Resolution Handwritten Digits Data Set with Writer Characteristics. SN Computer Science, 4(1):66, 2022. ISSN 558 2661-8907. doi: 10.1007/s42979-022-01494-2. URL https://doi.org/10.1007/ 559 s42979-022-01494-2.
 - J.O. Berger. Statistical Decision Theory and Bayesian Analysis. Springer Series in Statistics. Springer, 1985. ISBN 9780387960982. URL https://books.google.com/books?id= oY_x7dE15_AC.
 - Mikołaj Bińkowski, Dougal J. Sutherland, Michael Arbel, and Arthur Gretton. Demystifying MMD GANs. In International Conference on Learning Representations, 2018. URL https: //openreview.net/forum?id=r11UOzWCW.
 - Ricky T. Q. Chen and Yaron Lipman. Flow matching on general geometries. In The Twelfth International Conference on Learning Representations, 2024. URL https://openreview.net/ forum?id=g7ohDlTITL.
- 571 Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary dif-572 ferential equations. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, 573 and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 31. Cur-574 ran Associates, Inc., 2018. URL https://proceedings.neurips.cc/paper_files/ paper/2018/file/69386f6bb1dfed68692a24c8686939b9-Paper.pdf. 575
- 576 Quan Dao, Hao Phung, Binh Nguyen, and Anh Tran. Flow matching in latent space, 2023. URL 577 https://arxiv.org/abs/2307.08698. 578
- Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE* Signal Processing Magazine, 29(6):141–142, 2012. 580
 - Will Grathwohl, Ricky T. Q. Chen, Jesse Bettencourt, and David Duvenaud. Scalable reversible generative models with free-form continuous dynamics. In International Conference on Learning Representations, 2019. URL https://openreview.net/forum?id=rJxgknCcK7.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 585 Gans trained by a two time-scale update rule converge to a local nash equilibrium. In Proceedings 586 of the 31st International Conference on Neural Information Processing Systems, NIPS'17, pp. 6629–6640, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964. 588
- 589 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS '20, 591 Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546. 592
- Alex Krizhevsky. Learning multiple layers of features from tiny images. pp. 32–33, 2009. URL 593 https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf.

- 594 Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow 595 matching for generative modeling. In The Eleventh International Conference on Learning Repre-596 sentations, 2023. URL https://openreview.net/forum?id=PqvMRDCJT9t. 597 Guan-Horng Liu, Arash Vahdat, De-An Huang, Evangelos Theodorou, Weili Nie, and Anima 598 Anandkumar. I²sb: Image-to-image schrödinger bridge. In International Conference on Machine Learning (ICML), July 2023a. 600 601 Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and 602 transfer data with rectified flow. In The Eleventh International Conference on Learning Repre-603 sentations, 2023b. URL https://openreview.net/forum?id=XVjTT1nw5z. 604 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic mod-605 els. In Marina Meila and Tong Zhang (eds.), Proceedings of the 38th International Confer-606 ence on Machine Learning, volume 139 of Proceedings of Machine Learning Research, pp. 607 8162-8171. PMLR, 18-24 Jul 2021. URL https://proceedings.mlr.press/v139/ 608 nichol21a.html. 609 610 Derek Onken, Samy Wu Fung, Xingjian Li, and Lars Ruthotto. OT-Flow: Fast and accurate contin-611 uous normalizing flows via optimal transport. In AAAI Conference on Artificial Intelligence, vol-612 ume 35, pp. 9223-9232, May 2021. URL https://ojs.aaai.org/index.php/AAAI/ 613 article/view/17113. 614 Gaurav Parmar, Richard Zhang, and Jun-Yan Zhu. On aliased resizing and surprising subtleties in 615 gan evaluation. In CVPR, 2022. 616 617 Aram-Alexandre Pooladian, Heli Ben-Hamu, Carles Domingo-Enrich, Brandon Amos, Yaron Lip-618 man, and Ricky T. Q. Chen. Multisample flow matching: straightening flows with minibatch 619 couplings. In Proceedings of the 40th International Conference on Machine Learning, ICML'23. 620 JMLR.org, 2023. 621 Carl Edward Rasmussen and Christopher K. I. Williams. Gaussian Processes for Machine Learning. 622 The MIT Press, 11 2005. ISBN 9780262256834. doi: 10.7551/mitpress/3206.001.0001. URL 623 https://doi.org/10.7551/mitpress/3206.001.0001. 624 625 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and BjĶrn Ommer. High-626 resolution image synthesis with latent diffusion models. In Proceedings of the IEEE Conference 627 on Computer Vision and Pattern Recognition (CVPR), 2022. URL https://github.com/ 628 CompVis/latent-diffusionhttps://arxiv.org/abs/2112.10752. 629 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomed-630 ical image segmentation. In Nassir Navab, Joachim Hornegger, William M. Wells, and Alejan-631 dro F. Frangi (eds.), Medical Image Computing and Computer-Assisted Intervention – MICCAI 632 2015, pp. 234–241, Cham, 2015. Springer International Publishing. ISBN 978-3-319-24574-4. 633 634 Yuyang Shi, Valentin De Bortoli, Andrew Campbell, and Arnaud Doucet. Diffusion schrödinger 635 bridge matching. In Thirty-seventh Conference on Neural Information Processing Systems, 2023. 636 URL https://openreview.net/forum?id=qy070HsJT5. 637 Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribu-638 tion. Curran Associates Inc., Red Hook, NY, USA, 2019. 639 640 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben 641 Poole. Score-based generative modeling through stochastic differential equations. In Interna-642 tional Conference on Learning Representations, 2021. URL https://openreview.net/ 643 forum?id=PxTIG12RRHS. 644 Nicholas A Steinmetz, Peter Zatka-Haas, Matteo Carandini, and Kenneth D Harris. Distributed 645 coding of choice, action and engagement across the mouse brain. Nature, 576(7786):266-273, 646 2019. ISSN 1476-4687. doi: 10.1038/s41586-019-1787-x. URL https://doi.org/10. 647
 - 12

1038/s41586-019-1787-x.

653

654

655

656

657

658

659

660

661

662 663 664

665

688

689

696 697

699 700

701

- Alexander Tong, Kilian Fatras, Nikolay Malkin, Guillaume Huguet, Yanlei Zhang, Jarrid Rector-Brooks, Guy Wolf, and Yoshua Bengio. Improving and generalizing flow-based generative models with minibatch optimal transport. *Transactions on Machine Learning Research*, 2024a. ISSN 2835-8856. URL https://openreview.net/forum?id=CD9Snc73AW. Expert Certification.
 - Alexander Y. Tong, Nikolay Malkin, Kilian Fatras, Lazar Atanackovic, Yanlei Zhang, Guillaume Huguet, Guy Wolf, and Yoshua Bengio. Simulation-free Schrödinger bridges via score and flow matching. In Sanjoy Dasgupta, Stephan Mandt, and Yingzhen Li (eds.), Proceedings of The 27th International Conference on Artificial Intelligence and Statistics, volume 238 of Proceedings of Machine Learning Research, pp. 1279–1287. PMLR, 02–04 May 2024b. URL https:// proceedings.mlr.press/v238/tong24a.html.
 - C. Villani. *Optimal Transport: Old and New*. Grundlehren der mathematischen Wissenschaften. Springer Berlin Heidelberg, 2008. ISBN 9783540710509. URL https://books.google. com/books?id=hV8o5R7_5tkC.

A DISCUSSION ON PER-STREAM PERSPECTIVE ON FLOW MATCHING

It is helpful to recognize the relationship between the per-stream vector field and the conditional 667 vector field given one or both endpoints introduced previously in the literature. Specifically, the 668 per-sample vector field in Lipman et al. (2023) corresponds to marginalizing out s given the end 669 point x_1 , that is, $u_t(x \mid x_1) = \mathbb{E}(u_t(x \mid s) \mid s_t = x, s_1 = x_1)$. Similarly, the conditional vector 670 field of Tong et al. (2024a), corresponds to marginalizing out s given both x_0 and x_1 , that is $u_t(x)$ 671 $x_0, x_1 = \mathbb{E}\left(u_t(x \mid s) \mid s_t = x, s_0 = x_0, s_1 = x_1\right)$. Furthermore, when $p_s(\cdot \mid x_0, x_1)$ is simply a unit-point mass (Dirac) concentrated on the optimal transport (OT) path, i.e., a straight line that 672 connects two endpoints x_0 and x_1 , then $u_t(x \mid s) = u_t(x \mid x_1) = u_t(x \mid x_0, x_1)$ for all (s, t, x)673 tuples that satisfy $s_0 = x_0, s_1 = x_1, s_t = x$. Intuitively, when the stream connecting two ends is 674 unique, conditioning on the two ends is equivalent to conditioning on the corresponding stream s. 675 In this case, our stream-level FM algorithm (Section 3.2) coincides with those previous algorithms. 676 More generally, however, this equivalence does not hold when $p_s(\cdot \mid x_0, x_1)$ is non-degenerate. 677

The per-stream view affords additional modeling flexibility and alleviates the practitioners from 678 the burden of directly sampling from the conditional probability paths given one (Lipman et al., 679 2023) or both endpoints (Tong et al., 2024a). While the per-stream vector field induces a degenerate 680 unit-point mass conditional probability path, we will attain non-degenerate marginal and conditional 681 probability paths that satisfy the boundary conditions after marginalizing out the streams. Sampling 682 the streams in essence provides a data-augmented Monte Carlo alternative to sampling directly from 683 the conditional probability paths, which can then allow estimation of the marginal vector field $u_t(x)$ 684 when direct sampling from the conditional probability path is challenging. Additionally, as we will 685 demonstrate later, by approaching FM at the stream level, one could more readily incorporate prior 686 knowledge or other external features into the design of the stream distribution $p_s(\cdot \mid x_0, x_1)$. 687

B DERIVATION OF JOINT CONDITIONAL MEAN AND COVARIANCE

For computational efficiency and ease of implementation, we assume independent GPs across dimensions and present the derivation dimension-wise throughout the Appendices. We use s_t^i to denote the location of stream s at time t in dimension i, for i = 1, ..., d. Suppose each dimension of stream s follows a Gaussian process with a differentiable mean function ξ^i and covariance kernel c_{11}^i . Then, the joint distribution of $s_{t_1,...,t_g}^i = (s_{t_1}^i, ..., s_{t_g}^i)'$ and $\dot{s}_{t_1,...,t_g}^i = (\dot{s}_{t_1}^i, ..., \dot{s}_{t_g}^i)'$ at g time points is

$$\begin{pmatrix} s_{t_1,\dots,t_g}^i \\ \dot{s}_{t_1,\dots,t_g}^i \end{pmatrix} \sim \mathcal{N}\left(\begin{pmatrix} \xi_{t_1,\dots,t_g}^i \\ \dot{\xi}_{t_1,\dots,t_g}^i \end{pmatrix}, \begin{pmatrix} \Sigma_{11}^i & \Sigma_{12}^i \\ \Sigma_{12}^i & \Sigma_{22}^i \end{pmatrix} \right), \tag{2}$$

where $\xi_t^i = \xi^i(t), \dot{\xi}_t^i = d\xi_t^i/dt, \xi_{t_1,\dots,t_g}^i = (\xi_{t_1}^i,\dots,\xi_{t_g}^i)', \dot{\xi}_{t_1,\dots,t_g}^i = (\dot{\xi}_{t_1}^i,\dots,\dot{\xi}_{t_g}^i)'$ and covariance Σ_{jl}^i is determined by kernel c_{jl}^i . The kernel function for the covariance between s and \dot{s} in dimension

 $i ext{ is } c_{12}^i(t,t') = \frac{\partial c_{11}^i(t,t')}{\partial t'}, ext{ and the kernel defining covariance of } \dot{s} ext{ is } c_{22}^i = \frac{\partial^2 c_{11}^i(t,t')}{\partial t \partial t'} ext{ (Rasmussen & Williams (2005) Chapter 9.4). The conditional distribution of <math>(s, \dot{s})$ in dimension i given M observations $s_t^i = x_{obs}^i$ is also a (bivariate) Gaussian process. In particular, for $t \in [0, 1]$, let $\mu_t^i = (\xi_t^i, \dot{\xi}_t^i)'$ and $\mu_{obs}^i = (\xi_{t_1}^i, \dots, \xi_{t_s}^i)$, the joint distribution is

$$\begin{pmatrix} s_t^i, \dot{s}_t^i, {\boldsymbol{x}_{\text{obs}}^i}' \end{pmatrix}' \sim \mathcal{N} \left(\begin{pmatrix} \boldsymbol{\mu}_t^i \\ \boldsymbol{\mu}_{\text{obs}}^i \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_t^i & \boldsymbol{\Sigma}_{t,\text{obs}}^i \\ \boldsymbol{\Sigma}_t^i \boldsymbol{\tau} & \boldsymbol{\Sigma}_{\text{obs}}^i \end{pmatrix} \right)$$

where $\Sigma_t^i = \operatorname{Cov}(s_t^i, \dot{s}_t^i)$ and $\Sigma_{obs}^i = \operatorname{Cov}(\boldsymbol{x}_{obs}^i)$. Accordingly, the conditional distribution $(s_t^i, \dot{s}_t^i) | \boldsymbol{x}_{obs}^i \sim \mathcal{N}(\tilde{\boldsymbol{\mu}}_t^i, \tilde{\Sigma}_t^i)$, where $\tilde{\boldsymbol{\mu}}_t^i = \boldsymbol{\mu}_t^i + \Sigma_{t,obs}^i \Sigma_{obs}^{i^{-1}}(\boldsymbol{x}_{obs}^i - \boldsymbol{\mu}_{obs}^i)$ and $\tilde{\Sigma}_t^i = \Sigma_t^i - \Sigma_{t,obs}^i \Sigma_{obs}^{i^{-1}} \Sigma_{t,obs}^{i^{-1}}$.

C OPTIMAL TRANSPORT PATH FROM CONDITIONAL GP CONSTRUCTION

In this section, we show how to derive the path in I-CFM (Tong et al., 2024a) from the conditional GP construction (Appendix B) using a linear kernel. Without loss of generality, we present the derivation of "noise-free" path with $\sigma^2 = 0$ (i.e., the rectified flow, Liu et al. (2023b)).

Let
$$\mathbf{x}_{obs}^{i} = (x_{0}^{i}, x_{1}^{i})', \xi_{t}^{i} = \xi_{t}^{i} = 0 \text{ and } c_{11}^{i}(t, t') = \sigma_{a}^{2} + \sigma_{b}^{2}(t-1)(t'-1), \text{ such that}$$

$$\Sigma_{t}^{i} = \begin{pmatrix} \sigma_{a}^{2} + \sigma_{b}^{2}(t-1)^{2} & \sigma_{b}^{2}(t-1) \\ \sigma_{b}^{2}(t-1) & \sigma_{b}^{2} \end{pmatrix}, \qquad \Sigma_{t,obs}^{i} = \begin{pmatrix} \sigma_{a}^{2} - \sigma_{b}^{2}(t-1) & \sigma_{a}^{2} \\ -\sigma_{b}^{2} & 0 \end{pmatrix},$$

$$\Sigma_{obs}^{i} = \begin{pmatrix} \sigma_{a}^{2} + \sigma_{b}^{2} & \sigma_{a}^{2} \\ \sigma_{a}^{2} & \sigma_{a}^{2} \end{pmatrix}, \qquad \Sigma_{obs}^{i}^{-1} = \frac{1}{\sigma_{b}^{2}} \begin{pmatrix} 1 & -1 \\ -1 & 1 + \frac{\sigma_{b}^{2}}{\sigma_{a}^{2}} \end{pmatrix}$$

Therefore,

$$\tilde{\boldsymbol{\mu}}_{t}^{i} = \Sigma_{t,\text{obs}}^{i} \Sigma_{\text{obs}}^{i}^{-1} \begin{pmatrix} x_{0}^{i} \\ x_{1}^{i} \end{pmatrix} = \begin{pmatrix} 1-t & t \\ -1 & 1 \end{pmatrix} \begin{pmatrix} x_{0}^{i} \\ x_{1}^{i} \end{pmatrix} = \begin{pmatrix} (1-t)x_{0}^{i} + tx_{1}^{i} \\ x_{1}^{i} - x_{0}^{i} \end{pmatrix}$$
$$\tilde{\Sigma}_{t}^{i} = \Sigma_{t}^{i} - \Sigma_{t,\text{obs}}^{i} \Sigma_{\text{obs}}^{i}^{-1} \Sigma_{t,\text{obs}}^{i\mathsf{T}} = \boldsymbol{O}$$

D COVARIANCE UNDER SQUARED EXPONENTIAL KERNEL

Throughout this paper, we adopted the squared exponential (SE) kernel, with the same hyperparameters for each dimension. The kernel defining block covariance for s, (s, \dot{s}) and \dot{s} in dimension i from Equation 2 are as follows:

$$c_{11}^{i}(t,t') = \alpha \exp\left(-\frac{(t-t')^{2}}{2l^{2}}\right) \qquad c_{12}^{i}(t,t') = \frac{\alpha}{l^{2}}(t-t')\exp\left(-\frac{(t-t')^{2}}{2l^{2}}\right)$$
$$c_{21}^{i}(t,t') = -c_{12}^{i}(t,t') \qquad c_{22}^{i}(t,t') = \frac{\alpha}{l^{4}}\left[l^{2} - (t-t')^{2}\right]\exp\left(-\frac{(t-t')^{2}}{2l^{2}}\right).$$

E A SUPPLEMENTARY EXAMPLE FOR VARIANCE CHANGING OVER TIME

Here, instead of generating data from standard Gaussian noise, we consider 100 training (unpaired) samples from a 2-Gaussian to another 2-Gaussian (Figure 7A). The example streams connecting two points under different variance schemes are shown in Figure 7B, again using additional nugget noise for constant noise, and a dot product kernel for decreasing and increasing noise, as described in Section 4.1. We then fit 100 independent models and calculate the W2 distance between generated and test samples at t = 1. The results are summarized in Figure 7C. Now, since both ends have finite samples, injecting noise (a.k.a. adding regularization) at both ends helps.

F APPLICATION TO MNIST DATABASE

755 We explore the empirical benefits of variance reduction by using FM with GP conditional streams on the MNIST database (Deng, 2012). Four algorithms are considered: two linear stream models



Figure 7: Supplementary Example for Variance Change over Time.A. The 100 observations in training data at t = 0 and t = 1. B. Examples of streams between two points, under different variance change scheme. C. Train models 100 times and calculate 2-Wasserstein (W2) distance between generated and test samples for each. The results of these 100 seeds are summarized by mean and standard error.

778 (I-CFM, OT-CFM) and two GP stream models (GP-I-CFM, GP-OT-CFM). For a fair comparison, 779 we set $\sigma = 0$ for linear stream models and use noise-free GP stream models. For all models, U-Nets (Ronneberger et al., 2015; Nichol & Dhariwal, 2021) with 32 channels and 1 residual block are used. It takes around 50s, 51s, 52s and 53s for I-CFM, OT-I-CFM, GP-I-CFM and GP-OT-781 CFM to pass through all training dataset once. Figure 8A shows the 10 generated images for each 782 trained model, starting from the same standard Gaussian noise. Compared to I-CFM, the OT version 783 jointly samples two endpoints by 2-Wasserstein optimal transport (OT) map π (Tong et al., 2024a). 784 Here, we demonstrate how much the GP stream-level CFM can further improve the estimation. We 785 train each algorithm 100 times, and calculate the kernel inception distance (KID) (Bińkowski et al., 786 2018) and Fréchet inception distance (FID) (Heusel et al., 2017). The histograms in Figure 8B show 787 distribution of these 100 KIDs and FIDs, with results summarized in Figure 8C. According to KID 788 and FID, the independent sampling algorithms (I-algorithms) are comparable to optimal transport 789 sampling algorithms (OT-algorithms). However, algorithms using GP conditional stream exhibit 790 lower standard error and fewer extreme values for KID and FID, thereby reducing the occurrence of 791 outlier samples, as illustrated in Figure 1).

792 793 794

777

G APPLICATION TO LFP DATASET

795 In this section, to illustrate the usage of proposed GP-CFM for time series data, we apply the labeled-796 GP-I-CFM to a session of local field potential (LFP) data from a mouse brain. In the LFP dataset, the 797 neural activity across multiple brain regions is recorded when the mice perform a task on choosing 798 the side with highest contrast for visual gratings. The data contains 39 sessions from 10 mice, and 799 each session contains multiple trials. Time bins for all measurements are 10 ms, starting 500 ms 800 before stimulus onset. Here, we study LFP from stimulus onset to 500ms after stimulus, and hence each trial contains data from 50 time points. See Steinmetz et al. (2019) for more details of the LFP 801 dataset. 802

Here, we choose recordings from a mouse in one session, where the trial is repeated 214 times. For each single trial, the data contains a time series from 7 brain regions. To illustrate the temporal smoothness over time in a visually significant way, we subset the data so that there are 5 evenlyspaced time points. In summary, the training data has 214 observations, and dimension for each observation is 5×7 . Here, we fit the data by covariate GP-I-CFM, using the starting point as covariates, and generate 1000 LFP time series for each region (the starting LFP is generated from an I-CFM). For each second, the algorithm can run around 100 iterations per second (if we use all 50 time points, it runs around 2.5 iterations per second, and it take longer time to converge). The



Figure 8: **Application to MNIST dataset**. We compare the performance of four algorithms (I-CFM, OT-CFM, GP-I-CFM and GP-OT-CFM) on fitting MNIST dataset. **A**. The 10 images generated from each trained model. Fit the models 100 times for each, and evaluate the quality of the samples by KID and FID. **B**. The histograms of KID and FID. **C**. The mean and standard error for KID and FID.

results are shown in Figure 9. The generated time series can further be used to study neural activity in different brain regions. For example, the mean trajectories in Figure 9A suggest that the LFPs in Cornu Ammonis region 3 (CA3) and dentate gyrus (DG) are highly correlated, which is consistent with the experiment fact that the rat DG does not project to any brain region other than the CA3 field of the hippocampus (Amaral et al., 2007). Besides this, we can use the generated samples to make more scientific insightful and concrete conclusions. But this is beyond the scope of this paper.



Figure 9: **Application to LFP data**. We apply the GP-I-CFM with covariate (on starting point) to a session of local field potential (LFP) data from 7 regions of mouse brain. In the training dataset, there are 214 observations (repeated trials). For each observation, it is a time series data of 5 time points from 7 brain regions. Here, we generate 1000 LFP time series for each region, where the starting LFP is generated from an I-CFM. **A.** The mean trajectories over 1000 samples. **B.** The generated 1000 time series for CA3 and DG.

H PROOF OF PROPOSITIONS

In this section, we provide proofs for several propositions in the main text. All these proofs are adapted from Lipman et al. (2023); Tong et al. (2024a).

- 861 H.1 PROOF FOR CONDITIONAL FM ON STREAM
- **Proposition 1.** The marginal vector field over stream $u_t(x)$ generates the marginal probability path $p_t(x)$ from initial condition $p_0(x)$.

Reference Proof. Denote probability over stream as $q(s) = \int p_s(s \mid x_0, x_1) \pi(x_0, x_1) d(x_0, x_1)$ and $p_t(x \mid s) = \delta(x - s_t)$, then

$$\frac{d}{dt}p_t(x) = \frac{d}{dt}\int p_t(x \mid \boldsymbol{s})q(\boldsymbol{s})d\boldsymbol{s}$$

Assume the regularity condition holds, such that we can exchange limit and integral (and differentiation and integral) by dominated convergence theorem (DCT). Therefore,

$$= \int \frac{d}{dt} p_t(x \mid \boldsymbol{s}) q(\boldsymbol{s}) d\boldsymbol{s}$$

To handle the derivative on zero measure, define s_t -centered Gaussian conditional path and corresponding flow map as

$$p_{\sigma,t}(x \mid \boldsymbol{s}) := \mathcal{N}(x \mid s_t, \sigma^2 I)$$

$$\psi_{\sigma,t}(z \mid \boldsymbol{s}) := \sigma z + s_t,$$

for $z \sim N(0, I)$, such that $\lim_{\sigma \to 0} p_{\sigma,t}(x \mid s) = p_t(x \mid s)$. Then by Theorem 3 of Lipman et al. (2023), the unique vector field defining $\psi_{\sigma,t}(z \mid s)$ (and hence generating $p_{\sigma,t}(x \mid s)$) is $u_t^*(x \mid s) = ds_t/dt = u_t(s_t \mid s)$, for all (t, x). Note that $u_t^*(x \mid s)$ extends $u_t(x \mid s)$ by defining on all x, and they are equivalent when $s_t = x$. Since $u_t^*(\cdot \mid s)$ generates $p_{\sigma,t}(\cdot \mid s)$, by continuity equation,

$$\begin{aligned} \frac{d}{dt}p_t(x) &= \int \frac{d}{dt} \lim_{\sigma \to 0} p_{\sigma,t}(x \mid \boldsymbol{s})q(\boldsymbol{s})d\boldsymbol{s} \\ &= \int -\lim_{\sigma \to 0} \operatorname{div}(u_t^*(x \mid \boldsymbol{s})p_{\sigma,t}(x \mid \boldsymbol{s}))q(\boldsymbol{s})d\boldsymbol{s} \end{aligned}$$

Then by DCT,

$$= -\lim_{\sigma \to 0} \operatorname{div} \left(\int u_t^*(x \mid \boldsymbol{s}) p_{\sigma,t}(x \mid \boldsymbol{s}) q(\boldsymbol{s}) d\boldsymbol{s} \right)$$
$$= -\operatorname{div} \left(\int u_t^*(x \mid \boldsymbol{s}) \lim_{\sigma \to 0} p_{\sigma,t}(x \mid \boldsymbol{s}) q(\boldsymbol{s}) d\boldsymbol{s} \right)$$
$$= -\operatorname{div} \left(\mathbb{E} \left(u_t(x \mid \boldsymbol{s}) \mid s_t = x \right) p_t(x) \right)$$

By definition in equation 1,

$$= -\mathrm{div}\left(u_t(x)p_t(x)\right),$$

which shows that $p_t(\cdot)$ and $u_t(\cdot)$ satisfy the continuity equation, and hence $u_t(x)$ generates $p_t(x)$.

H.2 PROOF FOR GRADIENT EQUIVALENCE ON STREAM

Recall

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

$$\mathcal{L}_{\text{sCFM}}(\theta) = \mathbb{E}_{t,s} \| v_t^{\theta}(s_t) - u_t(x \mid s) \|^2,$$

where $x \sim p_t(x)$, $s \sim q(s)$ and $q(s) = \int p_s(s \mid x_0, x_1) \pi(x_0, x_1) d(x_0, x_1)$. Proposition 2. $\nabla_{\theta} \mathcal{L}_{FM}(\theta) = \nabla_{\theta} \mathcal{L}_{sCFM}(\theta)$.

Proof. To ensure existence of all integrals and to allow the changes of integral (Fubini's Theorem), 913 we assume that q(s) are decreasing to zero at a sufficient speed as $||s|| \to \infty$ and that $u_t, v_t, \nabla_{\theta} v_t$ 914 are bounded. To facilitate proof writing, let $p_t(x \mid s) = \delta(x - s_t)$.

915 The L-2 error in the expectation ca be re-written as

916
917
$$\|v_t^{\theta}(x) - u_t(x)\|^2 = \|v_t^{\theta}(x)\|^2 + \|u_t(x)\|^2 - 2\langle v_t^{\theta}(x), u_t(x) \rangle$$

$$\|v_t^{\theta}(s_t) - u_t(x \mid \mathbf{s})\|^2 = \|v_t^{\theta}(s_t)\|^2 + \|u_t(x \mid \mathbf{s})\|^2 - 2\langle v_t^{\theta}(s_t), u_t(x \mid \mathbf{s}) \rangle$$

Thus, it's sufficient to prove the result by showing the expectations of terms including θ are equiva-lent.

First,

$$\begin{split} \mathbb{E}_{x} \|v_{t}^{\theta}(x)\|^{2} &= \int \|v_{t}^{\theta}(x)\|^{2} p_{t}(x) dx \\ &= \int \int \|v_{t}^{\theta}(x)\|^{2} p_{t}(x \mid \boldsymbol{s}) q(\boldsymbol{s}) dx d\boldsymbol{s} \\ &= \mathbb{E}_{\boldsymbol{s}} \int \|v_{t}^{\theta}(x)\|^{2} \delta(x - s_{t}) dx \\ &= \mathbb{E}_{\boldsymbol{s}} \|v_{t}^{\theta}(s_{t})\|^{2} \end{split}$$

Second,

$$\begin{split} \mathbb{E}_x \langle v_t^{\theta}(x), u_t(x) \rangle &= \int \langle v_t^{\theta}(x), u_t(x) \rangle p_t(x) dx \\ &= \int \langle v_t^{\theta}(x), \frac{\int u_t(x \mid \boldsymbol{s}) p_t(x \mid \boldsymbol{s}) q(\boldsymbol{s}) d\boldsymbol{s}}{p_t(x)} \rangle p_t(x) dx \\ &= \int \langle v_t^{\theta}(x), \int u_t(x \mid \boldsymbol{s}) p_t(x \mid \boldsymbol{s}) q(\boldsymbol{s}) d\boldsymbol{s} \rangle dx \\ &= \int \int \langle v_t^{\theta}(x), u_t(x \mid \boldsymbol{s}) \rangle \delta(x - s_t) q(\boldsymbol{s}) d\boldsymbol{s} dx \\ &= \mathbb{E}_{\boldsymbol{s}} \langle v_t^{\theta}(s_t), u_t(x \mid \boldsymbol{s}) \rangle \end{split}$$

These two holds for all t, and hence $\nabla_{\theta} \mathcal{L}_{FM}(\theta) = \nabla_{\theta} \mathcal{L}_{sCFM}(\theta)$

H.3 PROOF FOR GRADIENT EQUIVALENCE CONDITIONING ON COVARIATES

Let x be response, c be covariates, and s be the stream connecting two endpoints (x_0, x_1) . Given covariate c, denote the conditional distribution of s as $q(s \mid c) = \int p_s(s \mid x_0, x_1, c) \pi(x_0, x_1) d(x_0, x_1)$ and marginal conditional probability path as $p_t(x \mid c)$. Further, let

> $\mathcal{L}_{cFM}(\theta) = \mathbb{E}_{t,x} \| v_t^{\theta}(x,c) - u_t(x \mid c) \|^2,$ $\mathcal{L}_{\text{cCFM}}(\theta) = \mathbb{E}_{t,s} \| v_t^{\theta}(s_t, c) - u_t(x \mid s) \|^2,$

where $x \sim p_t(x|c)$ and $s \sim q(s | c)$

Proposition 3. $\nabla_{\theta} \mathcal{L}_{cFM}(\theta) = \nabla_{\theta} \mathcal{L}_{cCFM}(\theta).$

Proof. To ensure existence of all integrals and to allow the changes of integral (Fubini's Theorem), we assume that $q(\cdot \mid c)$ decreases to zero at a sufficient speed as $||s|| \to \infty$ and that $v_t^{\theta}, u_t, \nabla_{\theta} v_t^{\theta}$ are bounded. To facilitate proof writing, let $p_t(x \mid s) = \delta(x - s_t)$.

The L-2 error in the expectation ca be re-written as

$$\|v_t^{\theta}(x,c)$$

$$\|v_t^{\theta}(s_t, c) - u_t(x \mid \boldsymbol{s})\|^2 = \|v_t^{\theta}(s_t, c)\|^2 + \|u_t(x \mid \boldsymbol{s})\|^2 - 2\langle v_t^{\theta}(s_t, c), u_t(x \mid \boldsymbol{s})\rangle$$

 $-u_t(x \mid c)\|^2 = \|v_t^{\theta}(x, c)\|^2 + \|u_t(x \mid c)\|^2 - 2\langle v_t^{\theta}(x, c), u_t(x \mid c)\rangle$

Thus, it's sufficient to prove the result by showing the expectations of terms including θ are equivalent.

First,

$$\mathbb{E}_{x} \|v_{t}^{\theta}(x,c)\|^{2} = \int \|v_{t}^{\theta}(x,c)\|^{2} p_{t}(x \mid c) dx$$

$$= \int \int ||v_t^{\theta}(x,c)||^2 p_t(x \mid s) q(s \mid c) dx ds$$

$$= \mathbb{E}_s \int ||v_t^{\theta}(x,c)||^2 \delta(x-s_t) dx$$
971

$$= \mathbb{E}_{\boldsymbol{s}} \| v_t^{\boldsymbol{\theta}}(s_t, c) \|^2$$

Second, $\mathbb{E}_x \langle v_t^{\theta}(x,c), u_t(x \mid c) \rangle = \int \langle v_t^{\theta}(x,c), u_t(x \mid c) \rangle p_t(x \mid c) dx$ $= \int \langle v_t^{\theta}(x,c), \frac{\int u_t(x \mid \boldsymbol{s}) p_t(x \mid \boldsymbol{s}) q(\boldsymbol{s} \mid c) d\boldsymbol{s}}{p_t(x \mid c)} \rangle p_t(x \mid c) dx$ $= \int \langle v_t^{\theta}(x,c), \int u_t(x \mid \boldsymbol{s}) p_t(x \mid \boldsymbol{s}) q(\boldsymbol{s} \mid c) d\boldsymbol{s} \rangle dx$ $= \int \int \langle v_t^{\theta}(x,c), u_t(x \mid \boldsymbol{s}) \rangle \delta(x-s_t) q(\boldsymbol{s} \mid c) d\boldsymbol{s} dx$ $= \mathbb{E}_{\boldsymbol{s}} \langle v_t^{\theta}(s_t, c), u_t(x \mid \boldsymbol{s}) \rangle$ These two holds for all t, and hence $\nabla_{\theta} \mathcal{L}_{cFM}(\theta) = \nabla_{\theta} \mathcal{L}_{cCFM}(\theta)$.