
Tree-Based Diffusion Schrödinger Bridge with Applications to Wasserstein Barycenters

Maxence Noble*

CMAP, CNRS, École polytechnique,
Institut Polytechnique de Paris,
91120 Palaiseau, France

Valentin De Bortoli

Computer Science Department,
ENS, CNRS, PSL University

Arnaud Doucet

Department of Statistics,
University of Oxford, UK

Alain Oliviero Durmus

CMAP, CNRS, École polytechnique,
Institut Polytechnique de Paris,
91120 Palaiseau, France

Abstract

Multi-marginal Optimal Transport (mOT), a generalization of OT, aims at minimizing the integral of a cost function with respect to a distribution with some prescribed marginals. In this paper, we consider an entropic version of mOT with a tree-structured quadratic cost, i.e., a function that can be written as a sum of pairwise cost functions between the nodes of a tree. To address this problem, we develop Tree-based Diffusion Schrödinger Bridge (TreeDSB), an extension of the Diffusion Schrödinger Bridge (DSB) algorithm. TreeDSB corresponds to a dynamic and continuous state-space counterpart of the multi-marginal Sinkhorn algorithm. A notable use case of our methodology is to compute Wasserstein barycenters which can be recast as the solution of a mOT problem on a star-shaped tree. We demonstrate that our methodology can be applied in high-dimensional settings such as image interpolation and Bayesian fusion.

1 Introduction

In the last decade, computational Optimal Transport (OT) has shown great success with applications in various fields such as biology (Schiebinger et al., 2019; Bunne et al., 2022), shape correspondence (Su et al., 2015; Feydy et al., 2017; Eisenberger et al., 2020), control theory (Bayraktar et al., 2018; Acciaio et al., 2019) and computer vision (Schmitz et al., 2018; Carion et al., 2020). While OT commonly seeks at computing the transport plan that minimizes the cost of moving between two distributions, it can naturally be extended to the multi-marginal setting (mOT) when considering several distributions. This extension of OT has notably been studied in quantum chemistry (Cotar et al., 2013), clustering (Cuturi & Doucet, 2014) and statistical inference (Srivastava et al., 2018). In particular, a popular application in unsupervised learning of mOT with Euclidean cost consists in computing the Wasserstein barycenter of a set of probability distributions (Agueh & Carlier, 2011; Benamou et al., 2015; Álvarez-Esteban et al., 2016; Peyré et al., 2019).

Interior point methods can be used to solve OT and mOT problems but they come with computational challenges (Pele & Werman, 2009). In order to mitigate these limitations, one often considers an *entropic regularization* of OT, known as Entropic OT (EOT). This regularized formulation can be efficiently solved in discrete state-spaces using the celebrated *Sinkhorn* algorithm (Cuturi, 2013; Knight, 2008; Sinkhorn & Knopp, 1967), which admits a continuous state-space counterpart referred

*Corresponding author. Contact at: maxence.noble-bourillot@polytechnique.edu.

to as the *Iterative Proportional Fitting* (IPF) procedure (Fortet, 1940; Kullback, 1968; Ruschendorf, 1995). In the case of a quadratic cost, EOT is equivalent to the *static* formulation of the Schrödinger Bridge (SB) problem (Schrödinger, 1932). Given a reference diffusion with finite time horizon T and two probability measures, solving SB amounts to finding the closest diffusion to the reference (in terms of Kullback–Leibler divergence on path spaces) with the given marginals at times $t = 0$ and $t = T$. This framework naturally arises in stochastic control (Dai Pra, 1991) where one aims at controlling the marginal distribution of a stochastic process at a fixed time. Recently, De Bortoli et al. (2021) introduced Diffusion Schrödinger Bridge (DSB), an approximation of a *dynamic* version of the IPF scheme on path spaces, see also Vargas et al. (2021); Chen et al. (2022). This methodology leverages advances in the field of denoising diffusion models (Song et al., 2021; Ho et al., 2020) in order to derive a scalable and efficient scheme to solve SB, and thus EOT.

Similarly to OT, mOT admits an entropic regularization (EmOT), which can be solved via a multi-marginal generalization of Sinkhorn/IPF algorithm (Benamou et al., 2015; Marino & Gerolin, 2020). Recently, Haasler et al. (2021) proposed an extension of the *static* SB problem in *discrete* state-space to any multi-marginal tree-based setting. They notably made the correspondence between this formulation and EmOT, when the cost function writes as the sum of interaction energies onto the given tree structure, and introduced an efficient version of Sinkhorn algorithm to solve it.

Motivations and contributions. In this work, we investigate the *continuous* and *dynamic* counterpart of the tree-based framework from Haasler et al. (2021). To be more specific, we present an extension of the static SB formulation in continuous state-space to any multi-marginal tree-based setting, referred to as TreeSB. Then, we establish the equivalence between TreeSB and a formulation of EmOT relying on a (quadratic) tree-structured cost function, analogously to Haasler et al. (2021). Inspired by DSB, we develop TreeDSB, a dynamic counterpart of the multi-marginal IPF (mIPF) to solve it, by operating on path spaces and using score-based diffusion techniques. To bridge gaps in literature, we prove the convergence of mIPF iterations in a *non-compact* setting under mild assumptions, by extending results on IPF convergence (Ruschendorf, 1995). Finally, we illustrate our approach on examples of Wasserstein barycenters from statistical inference and image processing.

Although our approach can be applied to any tree, we focus on *star-shaped trees*. In this setting, we show that TreeSB reduces to a regularized Wasserstein barycenter problem. Our method comes with several benefits compared to existing works. First, it is out-of-sample, *i.e.*, it does not require re-running the full procedure when given a new data point. Second, our formulation of the Wasserstein barycenter problem obtained from TreeSB allows us to avoid numerical issues of having to choose the regularization too small, see Section 5. Finally, to the best of our knowledge, this is the first methodology to extend ideas from diffusion-based models to the computation of Wasserstein barycenters. In particular, we believe that the idea of iterative refinement, *i.e.*, solving the *dynamic* counterpart of a *static* problem, plays a key role in the efficiency and scalability of the method.

Notation. For any measurable space (X, \mathcal{X}) , we denote by $\mathcal{P}(X)$ the space of probability measures defined on (X, \mathcal{X}) . Unless specified, \mathcal{X} is defined as the Borel sets on X . For any $\ell \in \mathbb{N}$, let $\mathcal{P}^{(\ell)} = \mathcal{P}((\mathbb{R}^d)^\ell)$; we denote $\mathcal{P}^{(1)}$ by \mathcal{P} . Assume that $X = (\mathbb{R}^d)^\ell$ for some $\ell \in \mathbb{N}$. For any $x \in X$ and any $m, n \in \{0, \dots, \ell\}$ such that $m \leq n$, let $x_{m:n} = (x_m, x_{m+1}, \dots, x_n)$. Let Leb be the Lebesgue measure. For any non-negative function $f : X \rightarrow \mathbb{R}_+$, such that $\int_X f d\text{Leb} < +\infty$, define $H(f) = -\int_X f \log f d\text{Leb} \in (-\infty, +\infty]$. For any distribution $\mu \in \mathcal{P}(X)$, we define the entropy of μ as $H(\mu) = H(d\mu/d\text{Leb})$ if $\mu \ll \text{Leb}$ and $H(\mu) = +\infty$ otherwise. For any two arbitrary measures μ and ν defined on (X, \mathcal{X}) , define the Kullback–Leibler divergence between μ and ν as $\text{KL}(\mu|\nu) = \int_X \log(d\mu/d\nu) d\mu - \int_X d\mu + \int_X d\nu$ if $\mu \ll \nu$ and $\text{KL}(\mu|\nu) = +\infty$ otherwise. For any $T > 0$, we denote by $C([0, T], \mathbb{R}^d)$ the space of continuous functions from $[0, T]$ to \mathbb{R}^d . For any path measure $\mathbb{P} \in \mathcal{P}(C([0, T], \mathbb{R}^d))$, we denote by $\text{Ext}(\mathbb{P}) \in \mathcal{P}^{(2)}$ the coupling between the *extremal* distributions of \mathbb{P} , *i.e.*, $\text{Ext}(\mathbb{P}) = \mathbb{P}_{0,T}$. Note that, for a given coupling $\pi_{0,T} \in \mathcal{P}^{(2)}$, there may exist several path measures \mathbb{P} verifying $\text{Ext}(\mathbb{P}) = \pi_{0,T}$. For any undirected tree $T = (V, E)$ with vertices V and edges E , we denote by $\{v, v'\}$ (or $\{v', v\}$) the undirected edge between $v \in V$ and $v' \in V$, if it exists. Given $r \in V$, we denote by $T_r = (V, E_r)$ the directed version of T rooted in r , where the directed edges E_r are uniquely defined from the edges E , see Appendix B for further details. In this case, the edge linking $v \in V$ to $v' \in V$ in T_r is denoted by (v, v') . Finally, for any integers $(n, K) \in \mathbb{N} \times \mathbb{N}^*$, we define $n \bmod(K)$ as the remainder of the Euclidean division of n by K .

2 Background and setting

Multi-marginal optimal transport. Let $\ell \in \mathbb{N}^*$. Given a cost function $c : (\mathbb{R}^d)^{\ell+1} \rightarrow \mathbb{R}$, a subset $S \subset \{0, \dots, \ell\}$ and a family of probability measures $\{\mu_i\}_{i \in S} \in \mathcal{P}^{|\mathbb{S}|}$, mOT consists in solving

$$\pi^* = \arg \min \left\{ \int c(x_{0:\ell}) d\pi(x_{0:\ell}) : \pi \in \mathcal{P}^{(\ell+1)}, \pi_i = \mu_i, \forall i \in S \right\}, \quad (\text{mOT})$$

where π_i is the i -th marginal of π , i.e., $\pi_i(A) = \pi(\text{proj}_i^{-1}(A))$ for any $A \in \mathcal{B}(\mathbb{R}^d)$, with $\text{proj}_i : x_{0:\ell} \mapsto x_i$. Given some weights $(w_i)_{i \in \{1, \dots, \ell\}} \in (\mathbb{R}_+)^{\ell}$, the Wasserstein barycenter between the measures $\{\mu_i\}_{i \in S}$ is given by π_0^* in (mOT), in the case where $S = \{1, \dots, \ell\}$ and $c(x_{0:\ell}) = \sum_{i=1}^{\ell} w_i \|x_0 - x_i\|^2$ (Peyré et al., 2019). In particular, when $w_i = 1/\ell$, the distribution π_0^* can be regarded as the Fréchet mean (Karcher, 2014) of the measures $\{\mu_i\}_{i \in S}$ for the Wasserstein distance of order 2. Similarly to OT, (mOT) can be relaxed using the following entropic regularization

$$\pi^* = \arg \min \left\{ \int c(x_{0:\ell}) d\pi(x_{0:\ell}) + \varepsilon \text{KL}(\pi|\nu) : \pi \in \mathcal{P}^{(\ell+1)}, \pi_i = \mu_i, \forall i \in S \right\}, \quad (\text{EmOT})$$

where $\varepsilon > 0$ is a hyperparameter and ν is an arbitrary measure defined on $((\mathbb{R}^d)^{\ell+1}, \mathcal{B}((\mathbb{R}^d)^{\ell+1}))$.

Link with Schrödinger Bridge. We first recall the relationship between Schrödinger Bridge and EOT. Given $T > 0$, \mathbb{Q} a (reference) path measure, i.e., $\mathbb{Q} \in \mathcal{P}(C([0, T], \mathbb{R}^d))$ and two measures $\mu_0, \mu_1 \in \mathcal{P}(\mathbb{R}^d)$, solving the SB problem amounts to finding the path measure \mathbb{P}^* defined by

$$\mathbb{P}^* = \text{argmin} \{ \text{KL}(\mathbb{P}|\mathbb{Q}) : \mathbb{P} \in \mathcal{P}(C([0, T], \mathbb{R}^d)), \mathbb{P}_0 = \mu_0, \mathbb{P}_T = \mu_1 \}. \quad (\text{SB})$$

If \mathbb{Q} is associated with a Stochastic Differential Equation (SDE)², of the form $d\mathbf{X}_t = -a\mathbf{X}_t dt + d\mathbf{B}_t$, with $a \geq 0$, then it can be shown, see (Léonard, 2014, Proposition 1) that $\mathbb{P}_{0,T}^*$ verifies

$$\mathbb{P}_{0,T}^* = \text{argmin} \{ \text{KL}(\pi|\mathbb{Q}_{0,T}) : \pi \in \mathcal{P}^{(2)}, \pi_0 = \mu_0, \pi_1 = \mu_1 \}. \quad (\text{static-SB})$$

This is called the *static* formulation of SB. It can be shown that solving (static-SB) is equivalent to solving EOT with quadratic cost and regularization $\varepsilon = 2 \sinh(aT)/a$ if $a > 0$, $\varepsilon = 2T$ if $a = 0$. Moreover, since $\mathbb{P}^* = \mathbb{P}_{0,T}^* \otimes \mathbb{Q}_{0,T}$, where $\mathbb{Q}_{0,T}$ is the measure \mathbb{Q} conditioned on initial and terminal conditions, solving the *dynamic* problem (SB) is equivalent to solving (static-SB).

Similarly, (EmOT) can be easily rewritten in a *static* multi-marginal SB fashion

$$\pi^* = \text{argmin} \{ \text{KL}(\pi|\pi^0) : \pi \in \mathcal{P}^{(\ell+1)}, \pi_i = \mu_i, \forall i \in S \}, \quad (\text{mSB-like})$$

with $(d\pi^0/d\text{Leb})(x_{0:\ell}) \propto \exp[-c(x_{0:\ell})/\varepsilon](d\nu/d\text{Leb})(x_{0:\ell})$, where π^0 is the *reference* measure.

Diffusion Schrödinger Bridge. Recently, De Bortoli et al. (2021) introduced Diffusion Schrödinger Bridge (DSB), a numerical scheme to solve (SB). It approximates the iterates of a *dynamic* version of the *Iterative Proportional Fitting* (IPF) scheme (Sinkhorn & Knopp, 1967; Knight, 2008; Peyré et al., 2019; Cuturi & Doucet, 2014), which can be described as follows: consider a sequence of path measures $(\mathbb{P}^n)_{n \in \mathbb{N}}$ such that $\mathbb{P}^0 = \mathbb{Q}$ and for any $n \in \mathbb{N}$

$$\mathbb{P}^{2n+1} = \text{argmin} \{ \text{KL}(\mathbb{P}|\mathbb{P}^{2n}) : \mathbb{P}_T = \mu_1 \}, \quad \mathbb{P}^{2n+2} = \text{argmin} \{ \text{KL}(\mathbb{P}|\mathbb{P}^{2n+1}) : \mathbb{P}_0 = \mu_0 \}.$$

This procedure alternatively projects between the measures with fixed initial distribution and the ones with fixed terminal distribution. For the first iteration, we get that $\mathbb{P}^1 = \mu_1 \otimes \mathbb{Q}_{|T}$. Assuming that \mathbb{Q} is given by $d\mathbf{X}_t = f_t(\mathbf{X}_t)dt + d\mathbf{B}_t$, with $f : [0, T] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$, then \mathbb{P}^1 is associated with the *time-reversal* of this SDE initialized at μ_1 . The time-reversal of an SDE has been derived under mild assumptions on the drift and diffusion coefficients (Hausmann & Pardoux, 1986; Cattiaux et al., 2021). In this case, we have $(\mathbf{Y}_{T-t})_{t \in [0, T]} \sim \mathbb{P}^1$, with $\mathbf{Y}_0 \sim \mu_1$ and

$$d\mathbf{Y}_t = \{-f_{T-t}(\mathbf{Y}_t) + \nabla \log p_{T-t}(\mathbf{Y}_t)\}dt + d\mathbf{B}_t,$$

where p_t is the density of \mathbb{P}_t^0 w.r.t. the Lebesgue measure. The score $\nabla \log p_t$ is estimated using score matching techniques (Hyvärinen, 2005; Vincent, 2011). The first iterate of DSB, \mathbb{P}^1 , corresponds to a *denoising diffusion model* (Ho et al., 2020; Song et al., 2021). DSB iterates further and not only parameterizes the backward process but also the forward process. It can therefore be seen as a refinement of diffusion models drawing a bridge between generative modeling and optimal transport.

²We refer to Appendix C for details on solutions of SDEs and associated measures.

Tree-based framework. Consider an undirected tree $T = (V, E)$, with vertices V and edges E , such that V is identified with $\{0, \dots, \ell\}$. Inspired by Haasler et al. (2021), we restrict our study of (EmOT), to the case where the cost function c is the tree-structured *quadratic* cost derived from T

$$c(x_{0:\ell}) = \sum_{\{v,v'\} \in E} w_{v,v'} \|x_v - x_{v'}\|_2^2, \quad (1)$$

where $w_{v,v'}$ is a weight on the edge $\{v, v'\}$, which links v to v' (and v' to v). Furthermore, as in Haasler et al. (2021), we choose S , *i.e.*, the set of vertices of T with constrained marginals, to coincide with the *leaves* of T . This framework recovers important applications, from Wasserstein barycenters to Wasserstein propagation, see Solomon et al. (2014, 2015). We emphasize that it differs from an OT problem defined on the space of graphs (Chen et al., 2016). Here, each node represents a probability measure (observed or to be inferred) and each edge represents a coupling between two distributions.

We consider an arbitrary vertex $r \in V$ and choose ν in (EmOT) such that $(d\nu/d\text{Leb})(x_{0:\ell}) = \varphi_r(x_r)$, where φ_r is a density defined on \mathbb{R}^d . Due to the form of ν and c , the reference measure π^0 in (mSB-like) is therefore a *probability* distribution which factorizes along $T_r = (V, E_r)$, the directed version of T rooted in r . We refer to Appendix B for more details on the notion of directed trees. In this setting, (EmOT) is equivalent to the tree-based problem

$$\begin{aligned} \pi^* = \operatorname{argmin} \{ & \text{KL}(\pi|\pi^0) : \pi \in \mathcal{P}^{(|V|)}, \pi_i = \mu_i, \forall i \in S \}, & (\text{TreeSB}) \\ & \text{with } \pi^0 = \pi_r^0 \otimes_{(v,v') \in E_r} \pi_{v'|v}^0, & (2) \end{aligned}$$

where $\pi_{v'|v}^0(\cdot | x_v) = N(x_v, \varepsilon/(2w_{v,v'})I_d)$ and $\pi_r^0 \ll \text{Leb}$ with density φ_r . In a manner akin to Haasler et al. (2021), we thus establish, in *continuous* state-space, the correspondence between (TreeSB), a *static* tree-based version of SB, and a version of EmOT with tree-structured cost (1). In our work, we make the following assumption on the constrained marginals $\{\mu_i\}_{i \in S}$.

A0. For any $i \in S$, $\mu_i \ll \text{Leb}$ and $H(\mu_i) < \infty$.

In what follows, we define K as the number of leaves of T , denoting $S = \{i_0, \dots, i_{K-1}\}$, and define the horizon times $T_{v,v'} = \varepsilon/(2w_{v,v'})$ for any $\{v, v'\} \in E$. For any $i_k \in S$, we will denote by $T_{k_n} = (V, E_{k_n})$ the directed version of T rooted in the leaf i_k . In the next section, we present our *dynamic* method to solve (TreeSB), called *Tree-based Diffusion Schrödinger Bridge*.

3 Tree-based Diffusion Schrödinger Bridge

In this section, we present a method to solve (TreeSB) in the case where $r \in S$, *i.e.*, r is a leaf of T . We refer to Appendix E for the extension to the case where $r \in V \setminus S$. Without loss of generality, see Appendix E, we assume that $r = i_{K-1}$ and choose $\varphi_r = d\mu_{i_{K-1}}/d\text{Leb}$, such that $\pi_{i_{K-1}}^0 = \mu_{i_{K-1}}$.

Dynamic approach to mIPF. In order to approximate solutions of (TreeSB), we consider the *multi-marginal* extension of the IPF algorithm, denoted by mIPF. Namely, we define a sequence of probability distributions $(\pi^n)_{n \in \mathbb{N}}$ such that for any $n \in \mathbb{N}$

$$\pi^{n+1} = \operatorname{argmin} \{ \text{KL}(\pi|\pi^n) : \pi \in \mathcal{P}^{(|V|)}, \pi_{i_{k_n+1}} = \mu_{i_{k_n+1}} \}, \quad (\text{mIPF})$$

where $k_n = (n-1) \bmod(K)$ and (k_n+1) is identified with $n \bmod(K)$. We define a *mIPF cycle* as a sequence of K consecutive mIPF updates. In particular, each marginal constraint is considered exactly once during one mIPF cycle. In a practical setting, our main aim is to sample from the (mIPF) iterates at the lowest cost. Although these updates can be made explicit, see Marino & Gerolin (2020) for instance, direct sampling is unfeasible in practice when d is large. To overcome this limitation, we suggest to compute these iterates in a *dynamic* fashion with equivalent path measures.

Since π^0 factorizes along T , see (2), one can show that the iterates of (mIPF) also factorize along T , see Section 4. Since these iterates all have a constrained marginal, we obtain the following decomposition for any $n \in \mathbb{N}$: $\pi^n = \mu_{i_{k_n}} \otimes_{(v,v') \in E_{k_n}} \pi_{v'|v}^n$ where E_{k_n} denotes the set of edges of the directed tree T_{k_n} . Then, our approach consists in computing *dynamic* iterates, *i.e.*, path measures, along the edges of T that coincide on their extremal times with the *static* iterates $(\pi^n)_{n \in \mathbb{N}}$. Namely, for any $n \in \mathbb{N}$, for any edge $(v, v') \in E_{k_n}$, we define a path measure $\mathbb{P}_{(v,v')}^n \in \mathcal{P}(C([0, T_{v,v'}], \mathbb{R}^d))$ such that $\text{Ext}(\mathbb{P}_{(v,v')}^n) = \pi_{v,v'}^n$, where $\text{Ext}(\mathbb{P}_{(v,v')}^n)$ stands for the joint distribution of $\mathbb{P}_{(v,v')}^n$ at times 0 and $T_{v,v'}$. In particular, it comes that $\pi_{v'|v}^n = \mathbb{P}_{(v,v'), T_{v,v'}|0}^n$. Using the tree-based form of the (mIPF) iterates, we can thus sample from π^n by (i) following the directed edges of T_{k_n} , (ii) diffusing along them the corresponding path measures $(\mathbb{P}_{(v,v')}^n)_{(v,v') \in E_{k_n}}$ and (iii) picking the samples on the vertices. When T is a *bridge-shaped* tree (2 vertices, 1 edge), it simply reduces to the dynamic reformulation of the IPF scheme. In what follows, we explain how to obtain our *dynamic* sequence.

Definition of the dynamic iterates. We first compute the iterate \mathbb{P}^0 , corresponding to the dynamic version of π^0 defined (2), in Proposition 1. Then, we build the following iterates by recursion on $n \in \mathbb{N}$ and prove their well-posedness in Proposition 2.

Proposition 1. Let $\mathbb{T}_{K-1} = (\mathbb{V}, \mathbb{E}_{K-1})$, the directed tree associated with $\mathbb{T} = (\mathbb{V}, \mathbb{E})$ and root i_{K-1} . Then, for any $(v, v') \in \mathbb{E}_{K-1}$, there exists $\mathbb{P}_{(v, v')}^0 \in \mathcal{P}(\mathcal{C}([0, T_{v, v'}], \mathbb{R}^d))$ with $\text{Ext}(\mathbb{P}_{(v, v')}^0) = \pi_{(v, v')}^0$ and such that $\mathbb{P}_{(v, v')}^0|_0$ is the distribution of $(\mathbf{B}_t)_{t \in [0, T_{v, v'}]}$, recalling that $T_{v, v'} = \varepsilon / (2w_{v, v'})$.

Before deriving the dynamic counterpart of the (mIPF) iterates, we introduce several definitions. For any path measure \mathbb{P} , we denote by \mathbb{P}^R the time-reversal of \mathbb{P} . For any directed tree and any vertex v of this tree, $p(v)$ refers to the (unique) parent of v , and $c(v)$ to the unique child of v when it exists, see Appendix B for more details.

Let $n \in \mathbb{N}$. Assume that we have defined the sequence of our dynamic iterates $(\mathbb{P}_{(v, v')}^m)_{(v, v') \in \mathbb{E}_{k_m}, m \leq n}$ up to stage n .

Consider the path $P_n = \{(v_j, v_{j+1})\}_{j=1}^J$ in the directed tree \mathbb{T}_{k_n} such that $v_1 = i_{k_n}$ and $v_{J+1} = i_{k_{n+1}}$. In particular, for any $(v, v') \in \mathbb{E}_{k_{n+1}}$, either $(v', v) \in P_n$ or $(v, v') \in \mathbb{E}_{k_n} \setminus P_n$. This is illustrated in Figure 1 when $\mathbb{V} = \{0, 1, 2, 3, 4\}$, $S = \{2, 3, 4\}$, $i_k = 3$ and $i_{k+1} = 4$: in this case, $P = \{(3, 1), (1, 0), (0, 4)\}$ and $(1, 2)$ is the only edge common to \mathbb{E}_k and \mathbb{E}_{k+1} .

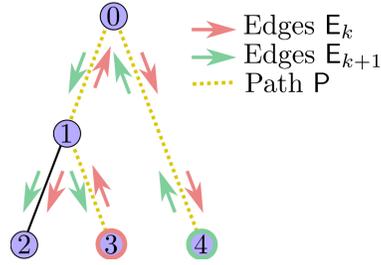


Figure 1: Illustration of the change of root in a toy tree with 5 vertices.

Consider now the directed tree $\mathbb{T}_{k_{n+1}}$. We define the $(n + 1)$ -th iterate of our dynamic sequence by recursion on the edges of this tree, following the breadth-first order. In this order, $(i_{k_{n+1}}, c(i_{k_{n+1}})) = (v_{J+1}, v_J)$ is the first edge considered.

First, we define $\mathbb{P}_{(v_{J+1}, v_J)}^{n+1} = \mu_{i_{k_{n+1}}} \otimes (\mathbb{P}_{(v_J, v_{J+1})}^n)^R$. In the case of a bridge-shaped tree, this is exactly the $(n + 1)$ -th update described in DSB. Then, for any $(v, v') \in \mathbb{E}_{k_{n+1}} \setminus \{(v_{J+1}, v_J)\}$,

- (a) either $(v, v') \in \mathbb{E}_{k_n} \setminus P_n$, and we define $\mathbb{P}_{(v, v')}^{n+1} = \mathbb{P}_{(p(v), v), T_{p(v), v}}^{n+1} \otimes \mathbb{P}_{(v, v')}^n|_0$,
- (b) or $(v', v) \in P_n$, and we define $\mathbb{P}_{(v, v')}^{n+1} = \mathbb{P}_{(p(v), v), T_{p(v), v}}^{n+1} \otimes (\mathbb{P}_{(v', v)}^n)^R|_0$.

Proposition 2. Consider the sequence of dynamic iterates defined by (a) and (b). Then, for any $n \in \mathbb{N}$ and any $(v, v') \in \mathbb{E}_{k_n}$, $\mathbb{P}_{(v, v')}^n \in \mathcal{P}(\mathcal{C}([0, T_{v, v'}], \mathbb{R}^d))$ and we have $\text{Ext}(\mathbb{P}_{(v, v')}^n) = \pi_{(v, v')}^n$.

Proposition 2 highlights the equivalence between the (mIPF) iterates and our dynamic iterates. These path measures are defined iteratively, by following the updates (a) and (b) along the edges of \mathbb{T} . The key observation here is that the computation of each dynamic iterate reduces to a sequence of updates (b) on a path linking two leaves of \mathbb{T} . We emphasize that our iterates could be similarly obtained by directly considering a dynamic formulation of (TreeSB) and introducing the formalism of deterministic time branching processes. We leave the study of this problem for future work. We now get into the details of our practical implementation, which relies on score-based methods.

Approximation of the dynamic iterates. The time-reversal operated in the update (b) can be computed explicitly, see Haussmann & Pardoux (1986) for instance. Indeed, assuming that $\mathbb{P}_{(v', v)}^n$ is associated with $d\mathbf{X}_t = f_{t, v', v}(\mathbf{X}_t)dt + d\mathbf{B}_t$ with $\mathbf{X}_0 \sim \pi_{v'}^n$, then, under mild conditions, its time-reversal $(\mathbb{P}_{(v', v)}^n)^R$ is associated with $d\mathbf{Y}_t = \{-f_{T-t, v', v} + \nabla \log p_{v', v, T-t}\}(\mathbf{Y}_t)dt + d\mathbf{B}_t$ with $\mathbf{Y}_0 \sim \pi_v^{n+1}$, where $p_{v', v, t}$ is the density of $\mathbb{P}_{(v', v), t}^n$ w.r.t. the Lebesgue measure. The score $\nabla \log p_{v', v, T-t}$ can then be approximated using score-matching techniques (Hyvärinen, 2005; Vincent, 2011) which are now ubiquitous in diffusion models (Song et al., 2021) and used in DSB De Bortoli et al. (2021). Therefore, at iteration $(n + 1)$, the update (b) is similar to the one of DSB for each edge on the path joining i_{k_n} and $i_{k_{n+1}}$. In practice, we parameterize the drifts $f_{t, v, v'}$ for any $\{v, v'\} \in \mathbb{E}$ with neural networks $f_{t, \theta_{v, v'}}$ and use the mean-matching loss introduced by De Bortoli et al. (2021). Note that doing so, we obtain $2|\mathbb{E}|$ neural networks. The whole procedure consisting in computing our dynamic iterates using the DSB framework is called Tree-based Diffusion Schrödinger Bridge (TreeDSB) and is summarized in Algorithm 1.

Algorithm 1 TreeDSB (Training)

```

1: Input:  $\mathbb{T} = (V, E)$ ,  $\{\mu_i\}_{i \in S}$ ,  $\{\theta_{v,v'}\}_{\{v,v'\} \in E}$ ,  $N \in \mathbb{N}$ 
2: for  $n = 0, \dots, N$  do
3:   Let  $k_n = (n - 1) \bmod(K)$ 
4:   Get path between  $i_{k_n}$  and  $i_{k_n+1}$ ,  $P_n = \{v_j, v_{j+1}\}_{j=1}^J$ 
5:   while not converged do
6:     for  $j = 1, \dots, J$  do
7:       Sample from  $\mathbb{P}_{v_j, v_{j+1}}^n$  (Euler-Maruyama)
8:       Compute mean matching loss  $\ell(\theta_{v_{j+1}, v_j})$ 
9:        $\theta_{v_{j+1}, v_j} \leftarrow$  Gradient Step( $\ell(\theta_{v_{j+1}, v_j})$ )
10:      Update  $f_{t, \theta_{v_{j+1}, v_j}}$ 
11:     end for
12:   end while
13: end for
14: Output:  $\{\theta_{v,v'}\}_{\{v,v'\} \in E}$ 

```

The algorithm is initialized with $f_{t, \theta_{v,v'}} = 0$ for all $\{v, v'\} \in E$. This corresponds to Brownian motion dynamics when sampling at the first iteration of TreeDSB, see Proposition 1. Note that in Algorithm 1, when we sample from $\mathbb{P}_{(v_j, v_{j+1})}^n$, we update $f_{t, \theta_{v_{j+1}, v_j}}$ which will be used to sample from $\mathbb{P}_{(v_{j+1}, v_j)}^{n+1}$ in the next iterations. In order to sample from the dynamics $\mathbb{P}_{(v_j, v_{j+1})}^n$, we consider its Euler–Maruyama discretization, see Appendix F for more details. We describe the different steps of the algorithm in the case of a toy example below, see Figure 2 for an illustration.

TreeDSB on a toy tree. We consider a star-shaped tree with three leaves denoted $\{1, 2, 3\}$ and its central node $\{0\}$. Following (2), we define π^0 with $r = 3$ and $\varphi_r = (d\mu_3/d\text{Leb})$. During the first iteration of TreeDSB, \mathbb{T} is rooted at vertex 3 and we compute samples from the *forward* path $P_0 = \{(3, 0), (0, 1)\}$ with Brownian motions, see Proposition 1, in order to learn the *backward* path $\{(1, 0), (0, 3)\}$. In the next iteration, we re-root the tree \mathbb{T} at vertex 1 and consider the *forward* path $P_1 = \{(1, 0), (0, 2)\}$, where the edges $(1, 0)$ and $(0, 2)$ are respectively given by the first iteration and the initialisation. This highlights that *TreeDSB does not require to update the whole tree*. The following iterations are done similarly. At each iteration $n \in \mathbb{N}$, we sample from π^n by first sampling from μ_{k_n} at leaf i_{k_n} and then following the parameterized SDEs on the directed edges of \mathbb{T}_{k_n} .

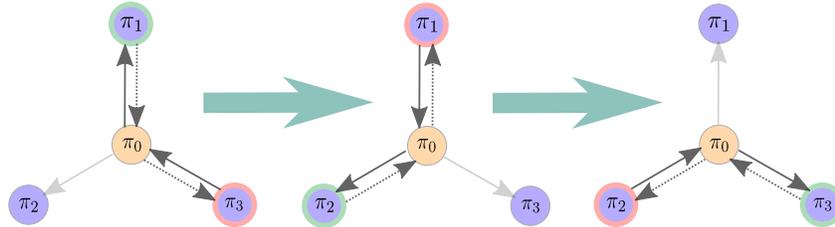


Figure 2: Illustration of one mIPF cycle solved by TreeDSB for a toy star-shaped tree. At each iteration, our method learns the *backward* stochastic process (dotted arrows) that goes from the target leaf (green-circled), corresponding to the constrained marginal, to the current root of the tree (red-circled) by using samples from the *forward* stochastic process (solid arrows).

4 Theoretical properties of mIPF

In this section, we study some of the theoretical properties of the *static* iterates $(\pi^n)_{n \in \mathbb{N}}$, that are equivalent to our *dynamic* iterates according to Proposition 2. In the case where the cost function c is bounded in (EmOT), results of convergence of (mIPF) exist (Marino & Gerolin, 2020; Carlier, 2022). However, our setting does not satisfy their assumptions, since our transport cost is quadratic and the measures are defined on \mathbb{R}^d . In what follows, we provide the first non-quantitative convergence results for (mIPF) in a *non-compact* setting.

For the rest of the section, we consider a static formulation of the multi-marginal Schrödinger bridge problem which is more general than (TreeSB), defined as

$$\pi^* = \operatorname{argmin}\{\text{KL}(\pi | \pi^0) : \pi \in \mathcal{P}^{(\ell+1)}, \pi_i = \mu_i, \forall i \in S\}, \quad (\text{static-mSB})$$

where $S \subset \{0, \dots, \ell\}$, $\pi^0 \in \mathcal{P}$, $\{\mu_i\}_{i \in S} \in \mathcal{P}^{|S|}$. We consider the following set of assumptions.

A1. *There exists a family of measures $\{\nu_i\}_{i \in \{0, \dots, \ell\}}$ defined on $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$ such that $\pi^0 \ll \bigotimes_{i=0}^{\ell} \nu_i$ with density $h = d\pi^0 / (d \bigotimes_{i=0}^{\ell} \nu_i)$ and $\mu_i \ll \nu_i$ with density $r_i = d\mu_i / d\nu_i$ for any $i \in S$.*

A2. $\{\pi \in \mathcal{P}^{(\ell+1)} : \text{KL}(\pi | \pi^0) < \infty, \pi_i = \mu_i, \forall i \in S\} \neq \emptyset$.

A3. *There exists a family of probability measures $\{\tilde{\mu}_j\}_{j \in \{0, \dots, \ell\} \setminus S}$ such that $\pi^0 \sim \tilde{\pi}^0$, where $\tilde{\pi}^0 = \bigotimes_{i \in S} \mu_i \bigotimes_{j \in \{0, \dots, \ell\} \setminus S} \tilde{\mu}_j$.*

In particular, (static-mSB) recovers (TreeSB) by considering $\nu_i = \text{Leb}$ for any $i \in \{0, \dots, \ell\}$ and $h(x_{0:\ell}) = \varphi_r(x_r) \exp[-c(x_{0:\ell})/\varepsilon]$ in **A1**. We detail in Appendix **D** how **A2** and **A3** can be met in (TreeSB). Under these assumptions, the multi-marginal Schrödinger Bridge exists.

Proposition 3. *Assume **A1** and **A2**. Then, there exists a unique solution π^* to (static-mSB). In addition, assume **A3**. Then, there exists a family $\{\psi_i^*\}_{i \in S}$ of measurable functions $\psi_i^* : \mathbb{R}^d \rightarrow \mathbb{R}$ such that*

$$(d\pi^*/d\pi^0) = \exp\left[\bigoplus_{i \in S} \psi_i^*\right] \quad \pi^0\text{-a.s.}$$

In order to establish the existence and uniqueness result of Proposition 3, we extend results from Nutz (2021) to the multi-marginal setting. A consequence of Proposition 3 is that the iterates of (mIPF) can be described using potentials.

Corollary 4. *Assume **A1**, **A2** and **A3**. Let $(\pi^n)_{n \in \mathbb{N}}$ be the sequence given by (mIPF). Then, for any $n \in \mathbb{N}^*$ with $k_n = (n - 1) \bmod(K)$ and $q_n \in \mathbb{N}$ such that $n = q_n K + k_n + 1$, there exists a family of measurable functions $\{\psi_{i_0}^{q_n+1}, \dots, \psi_{i_{k_n}}^{q_n+1}, \psi_{i_{k_n+1}}^{q_n}, \dots, \psi_{i_{K-1}}^{q_n}\}$ such that*

$$(d\pi^n/d\pi^0)(x_{0:\ell}) = \exp\left[\bigoplus_{j=0}^{k_n} \psi_{i_j}^{q_n+1}(x_{i_j}) \bigoplus_{j=k_n+1}^{K-1} \psi_{i_j}^{q_n}(x_{i_j})\right] \quad \pi^0\text{-a.s.}$$

In the tree-based setting, Corollary 4 explains why the (mIPF) iterations preserve the tree-based Markovian nature of π^0 . We now prove that the marginal π_i^n converges to μ_i for any $i \in S$, as n goes to infinity, i.e., we have marginal convergence on the leaves of T .

Proposition 5. *Assume **A1** and **A2**. Let $(\pi^n)_{n \in \mathbb{N}}$ be the sequence given by (mIPF). Then, we have $\lim_{n \rightarrow \infty} \|\pi_i^n - \mu_i\|_{\text{TV}} = 0$ for any $i \in S$.*

The previous result does not ensure the convergence of $(\pi^n)_{n \in \mathbb{N}}$ to the solution to (static-mSB). In particular, Proposition 5 does not provide the convergence of the marginals on the nodes $v \in V \setminus S$, which is key to compute regularized Wasserstein barycenters with TreeDSB. Relying on additional assumptions, we now derive the convergence of (mIPF).

A4. $\bigoplus_{i \in S} L^1(\mu_i) \subset L^1(\pi^*)$ is closed.

A5. There exist $\bar{c} \in (0, \infty)$ such that $\exp(\psi_{i_k}^n - \psi_{i_k}^{n+1}) \leq \bar{c}$, for any $n \in \mathbb{N}$, any $k \in \{0, \dots, K - 2\}$.

These assumptions can be seen as multi-marginal extensions of the ones of Ruschendorf (1995), see Appendix **D** for a discussion and examples.

Proposition 6. *Assume **A1**, **A2**, **A3**, **A4** and **A5**. Let $(\pi^n)_{n \in \mathbb{N}}$ be the sequence given by (mIPF). Then, we have $\lim_{n \rightarrow \infty} \|\pi^n - \pi^*\|_{\text{TV}} = 0$, where π^* is given in Proposition 3.*

To the best of our knowledge, Proposition 6 is the first convergence result of (mIPF) without assuming that the space is compact or that the cost is bounded. We highlight that traditional techniques to prove the convergence of IPF cannot be easily extended to the multi-marginal setting as pointed by Carlier (2022). In the case of bounded cost, quantitative results exist (Marino & Gerolin, 2020; Carlier, 2022). We leave the study of such results in the unbounded cost setting for future work.

5 Application to Wasserstein barycenters

Although Algorithm 1 can be applied to trees T with fixed marginals on the leaves, one case of particular interest is star-shaped trees, i.e., trees with a central node, denoted by index 0, and such that $S = \{1, \dots, \ell\}$ (see Figure 2 for an illustration with $\ell = 3$). In this section, we draw a link between (TreeSB) and regularized Wasserstein barycenters. We recall the definition of the Wasserstein distance of order 2 with ε -entropic regularization between μ and ν (Peyré et al., 2019, Chapter 4)

$$W_{2,\varepsilon}^2(\mu, \nu) = \inf\left\{\int \|x_1 - x_0\|^2 d\pi(x_0, x_1) - \varepsilon H(\pi) : \pi \in \mathcal{P}^{(2)}, \pi_0 = \mu, \pi_1 = \nu\right\}. \quad (3)$$

In this work, we consider the $(\ell\varepsilon, (\ell - 1)\varepsilon)$ -doubly-regularized Wasserstein-2 barycenter problem (Chizat, 2023) defined as follows

$$\mu_\varepsilon^* = \arg \min\left\{\sum_{i=1}^{\ell} w_i W_{2,\varepsilon/w_i}^2(\mu, \mu_i) + (\ell - 1)\varepsilon H(\mu) : \mu \in \mathcal{P}\right\}, \quad (\text{regWB})$$

where $(w_i)_{i \in \{1, \dots, \ell\}} \in (0, +\infty)^\ell$. The following proposition shows the equivalence between the barycenter problem (regWB) and the multi-marginal Schrödinger bridge problem (TreeSB) over T . In particular, it allows us to use TreeDSB to estimate the solution μ_ε^* of (regWB).

Proposition 7. Let $\varepsilon > 0$. Assume **A0**. Also assume that T is a star-shaped tree with central node indexed by 0, and that the reference measure of **(TreeSB)** defined in (2) verifies $r = i_{K-1}$ and $\varphi_r = d\mu_{i_{K-1}}/d\text{Leb} > 0$. Under **A2**, **(regWB)** has a unique solution π_0^* , where π^* solves **(TreeSB)**.

The proof of this result is postponed to Appendix D. More generally, we show in Appendix D that, for any tree T , **(TreeSB)** is equivalent to a regularized version of the Wasserstein propagation problem (Solomon et al., 2014, 2015). Moreover, we present in Appendix E an extension of Proposition 7 in the case where the chosen root r is not a leaf of T . We finally emphasize that the formulation of **(regWB)** leads to a *minimization* of the entropy of the barycenter. In particular, this allows us to choose ε reasonably large in TreeDSB, which is a stability advantage compared to other regularized methods which do not consider this further regularization.

6 Related work

Diffusion Schrödinger Bridge. Schrödinger Bridges (Schrödinger, 1932) have been extensively studied using tools from stochastic control and probability theory (Léonard, 2014; Dai Pra, 1991; Chen et al., 2021). More recently, algorithms were proposed to efficiently approximate such bridges in the context of machine learning. In particular, De Bortoli et al. (2021) proposed DSB while Vargas et al. (2021); Chen et al. (2022) developed related algorithms. In Chen et al. (2023), the authors study a multi-marginal version of DSB in a linear tree-based setting, where the set of observed nodes is the whole set of vertices. However, contrary to our setting, Chen et al. (2023) introduced a momentum variable. This allows for smoother trajectories which are desirable for single-cell trajectories applications and correspond to some spline interpolation in the space of probability measures (Chen et al., 2018). A general framework for tree-based static Schrödinger Bridges on discrete state-spaces was given in Haasler et al. (2021). In this work, we extend their formulation to a dynamic and continuous setting, see Appendix D for more a thorough comparison.

Wasserstein barycenters. The notion of Wasserstein barycenter was first introduced in Rabin et al. (2012) and then later studied in Agueh & Carlier (2011). The algorithms to solve this problem can be split into two families: the in-sample based approaches and the parametric ones. In-sample approaches require access to all the measures μ_i which are assumed to be empirical measures (Cuturi & Doucet, 2014; Benamou et al., 2015; Solomon et al., 2015). Related to this class of algorithms is the semi-discrete approach, which aims at computing a Wasserstein barycenter between continuous distribution but rely on a discretization of the barycenter (Claici et al., 2018; Staib et al., 2017; Mi et al., 2020). Most recent approaches do not rely on a discrete representation of the barycenter, but instead parameterize it using neural networks. These approaches can be further split into two categories. First, *measure-based optimization* approaches parameterize the measures using a neural network. This is the case of Cohen et al. (2020), where the barycenter is given by a generative model, which is then optimized. Fan et al. (2020) introduce an optimization procedure which relies on a *min-max-min* problem using the framework of Makuva et al. (2020). More recently, Korotin et al. (2022) considered a fixed point-based algorithm introduced in Álvarez-Esteban et al. (2016) to update a generative model parametrizing the barycenter. On the one hand, *potential-based methods* rely on a dual formulation of the barycenter. Korotin et al. (2021) parameterized the dual potentials using Input Convex Neural Network and considered regularizing losses imposing conjugacy and congruency. On the other hand, Li et al. (2020) consider a dual version of the *regularized* Wasserstein barycenter problem contrary to other works. Our approach applied to star-shaped trees also approximates a *regularized* Wasserstein barycenter. However, contrary to Li et al. (2020), we do not consider a parameterization of the potentials in the *static* setting but instead, parameterize the drift of an associated *dynamic* formulation using Schrödinger bridges. To the best of our knowledge TreeDSB is the first approach leveraging DSB-like algorithms to compute Wasserstein barycenters.

7 Experiments

In our experiments³, we illustrate the performance of TreeDSB to compute entropic regularized Wasserstein barycenters for various tasks. We choose to compare our method with state-of-the-art regularized algorithms: fast free-support Wasserstein barycenter (fsWB) (Cuturi & Doucet, 2014), and continuous regularized Wasserstein barycenter (crWB) (Li et al., 2020). In all of our settings, we consider a star-shaped tree with K leaves and edge weights that are equal to $1/K$, resulting in a

³Code available at <https://github.com/maxencenoble/tree-diffusion-schrodinger-bridge>.

sequential training procedure over $2K$ neural networks. The initial diffusion is always a Brownian motion parameterized as explained in Proposition 1. Hence, the time horizon on each edge is defined by $T = K\varepsilon/2$. The order of the leaves is randomly shuffled between the mIPF cycles. We consider 50 steps for the time discretization on $[0, T]$. We refer to Appendix G for details on the choice of the schedule, the architecture of the neural networks and the settings of our experiments.

Synthetic two dimensional datasets. We first illustrate TreeDSB in a synthetic two dimensional setting. We consider three different datasets *Swiss-roll* (vertex 0, starting node r), *Circle* (vertex 2) and *Moons* (vertex 3) and compute their Wasserstein barycenter (vertex 1) by running TreeDSB for 50 mIPF cycles with $\varepsilon = 0.1$. In Figure 3, we show the estimated densities of the datasets on the leaves of the tree (we emphasize that the distributions plotted on each leaf are generated from the central barycenter measure). In Figure 4, we observe the consistency between the barycenters generated from the different leaves. In Appendix G, we present additional results for this setting.

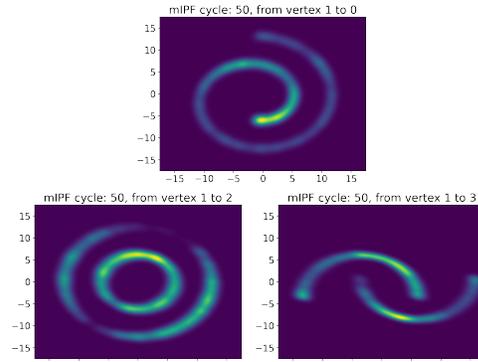


Figure 3: Estimated densities on the leaves.

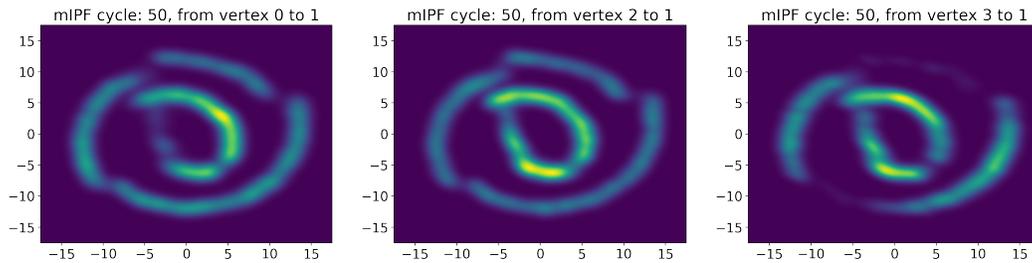


Figure 4: From left to right: barycenter estimated from the leaves *Swiss-roll*, *Circle* and *Moons*.

Synthetic Gaussian datasets. Next, we consider three independent Gaussian distributions with zero mean and random non-diagonal covariance matrices whose conditional number is less than 10, following Fan et al. (2020). In this case, the non-regularized barycenter can be exactly computed. To evaluate the performance of the algorithms, we use the Bures-Wasserstein Unexplained Variance Percentage (UVP), following (Korotin et al., 2021, Section 5). Given a target distribution $\mu^* \in \mathcal{P}$ and some approximation $\mu \in \mathcal{P}$, we define

$$\text{BW}_2^2\text{-UVP}(\mu, \mu^*) = 100 \cdot 2 \text{BW}_2^2(\mu, \mu^*) / \text{Var}(\mu^*)\%,$$

where $\text{BW}_2^2(\mu, \mu^*) = W_2^2(\mathcal{N}(\mathbb{E}[\mu], \text{Cov}(\mu)), \mathcal{N}(\mathbb{E}[\mu^*], \text{Cov}(\mu^*)))$.

Method	$d = 2$	$d = 16$	$d = 64$	$d = 128$	$d = 256$
fsWB (Cuturi & Doucet, 2014)	0.06 ± 0.01	2.86 ± 0.06	11.12 ± 0.06	14.47 ± 0.07	17.41 ± 0.05
crWB (Li et al., 2020)	0.02 ± 0.01	1.52 ± 0.11	11.41 ± 0.73	5.75 ± 0.02	18.27 ± 0.54
Tree DSB	0.63 ± 0.26	1.07 ± 0.58	1.39 ± 0.07	1.92 ± 0.02	2.62 ± 0.07

Table 1: Gaussian setting: comparison with the regularized methods crWB and fsWB.

In this setting, we choose μ^* to be the non-regularized barycenter and assess the dependency w.r.t. the dimension of the algorithms using the $\text{BW}_2^2\text{-UVP}$ metric. In Table 1, we compare ourselves with the two regularized methods Li et al. (2020) (L_2 -reg. equal to 10^{-4}) and Cuturi & Doucet (2014). We run TreeDSB for 10 mIPF cycles with $\varepsilon = 0.1$. Bold numbers represent the best values up to statistical significance. While Li et al. (2020) and Cuturi & Doucet (2014) enjoy better performance in low dimensions ($d = 2$), TreeDSB outperforms these methods as the dimension increases.

MNIST Wasserstein barycenter. We then turn to an image experiment using MNIST dataset (LeCun, 1998). Here, an image is not considered as a 2D-dimensional distribution as in Cuturi & Doucet (2014) and Li et al. (2020), but as a sample from a high-dimensional probability measure ($d = 784$). We aim at computing a Wasserstein barycenter between the digits 2,4 and 6. To do so, we

run TreeDSB for 10 mIPF cycles with r that corresponds to the digit 6 and $\varepsilon = 0.5$. In Figure 5, we display samples from the estimated marginals on the leaves, to assess the reconstruction of the digits 2, 4 and 6, and samples from the barycenter, obtained by diffusing from the leaf corresponding to the digit 6. Our results prove the scalability of TreeDSB to the high-dimensional setting, compared to state-of-the-art regularized methods. Additional results on MNIST dataset are given in Appendix G.

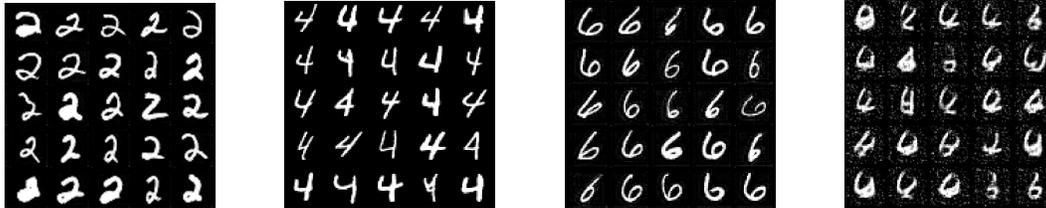


Figure 5: Samples from the estimated MNIST 2-4-6 marginals and from their Wasserstein barycenter.

Subset posterior aggregation. Finally, we evaluate TreeDSB in the context of *Bayesian fusion* (Srivastava et al., 2018), also called posterior aggregation. Given a Bayesian model and a dataset partitioned into several shards, this task aims at recovering the full data posterior distribution from the posterior distributions computed on each shard.

Method	Without het.	With het.
fsWB (Cuturi & Doucet, 2014)	12.95 \pm 0.35	14.43 \pm 0.51
crWB (Li et al., 2020)	20.66 \pm 0.71	23.06 \pm 0.12
Tree DSB	8.69\pm0.12	8.90\pm0.68

Table 2: Bayesian fusion setting: comparison with the regularized methods crWB and fsWB.

In particular, it has been proved that the barycenter of the subdataset posteriors is close to the full data posterior under mild assumptions (Srivastava et al., 2018). Here, we consider a logistic regression model applied to the wine dataset⁴ ($d = 42$) and proceed as follows. We first split this dataset into 3 subsets, with or without heterogeneity, and estimate the posterior parameters on each shard. Then, we draw samples from the obtained logistic distributions to define μ_1, μ_2, μ_3 . Then, we compute the Wasserstein barycenter of these measures, and compare it to the posterior computed on the full dataset. As in the synthetic Gaussian experiment, we run TreeDSB for 10 mIPF cycles $\varepsilon = 0.1$ and we compare ourselves with Li et al. (2020) (L_2 -reg. equal to 10^{-4}) and Cuturi & Doucet (2014). We evaluate the methods using the BW_2^2 -UVP metric, where μ^* is the estimated full data posterior, and report the results in Table 2. In both settings, we observe that our method outperforms existing regularized methods to compute Wasserstein barycenters.

Limitations. One of the main limitation of entropic regularized OT approach is that their behavior is usually badly conditioned as $\varepsilon \rightarrow 0$. In our setting, we observe that if ε , or equivalently T , is too low then the algorithm becomes less stable as the training of the models slows down. In the future, we plan to mitigate this issue by incorporating fixed point techniques like the one used in Korotin et al. (2022). Finally, since our algorithm is based on DSB (De Bortoli et al., 2021), it suffers from the same limitations. In particular, training different neural networks iteratively incurs some bias in the SDE which is harmful for large number of mIPF iterations.

8 Discussion

In this paper, we introduced Tree-based Diffusion Schrödinger Bridge (TreeDSB) a scalable scheme to approximate solutions of entropic-regularized multi-marginal Optimal Transport (mOT) problems. Our methodology leverages tools from the diffusion model literature and extends Diffusion Schrödinger Bridge (De Bortoli et al., 2021). In particular, it approximates the iterates of the multi-marginal Iterative Proportional Fitting (mIPF) algorithm, for which we prove its convergence under mild assumptions. We illustrate the efficiency of TreeDSB for image processing and Bayesian fusion, using the link between mOT and Wasserstein barycenters. In future work, we would like to study quantitative convergence bounds for mIPF in the *unbounded* cost setting. Another line of work would be to scale TreeDSB to higher dimensional problems building on recent developments in the diffusion model and flow matching community (Lipman et al., 2023; Peluchetti, 2023; Shi et al., 2023).

⁴<https://archive.ics.uci.edu/ml/datasets/wine>

Acknowledgments

We thank James Thornton for the DSB codebase⁵ and useful discussions. AD acknowledges support from the Lagrange Mathematics and Computing Research Center. AD and MN would like to thank the Isaac Newton Institute for Mathematical Sciences for support and hospitality during the programme *The mathematical and statistical foundation of future data-driven engineering* when work on this paper was undertaken. MN acknowledges funding from the grant SCAI (ANR-19-CHIA-0002).

References

- Acciaio, B., Backhoff-Veraguas, J., and Carmona, R. Extended mean field control problems: stochastic maximum principle and transport perspective. *SIAM journal on Control and Optimization*, 57(6):3666–3693, 2019.
- Agueh, M. and Carlier, G. Barycenters in the Wasserstein space. *SIAM Journal on Mathematical Analysis*, 43(2):904–924, 2011.
- Álvarez-Esteban, P. C., Del Barrio, E., Cuesta-Albertos, J., and Matrán, C. A fixed-point approach to barycenters in Wasserstein space. *Journal of Mathematical Analysis and Applications*, 441(2): 744–762, 2016.
- Bayraktar, E., Cox, A. M., and Stoev, Y. Martingale optimal transport with stopping. *SIAM Journal on Control and Optimization*, 56(1):417–433, 2018.
- Benamou, J.-D., Carlier, G., Cuturi, M., Nenna, L., and Peyré, G. Iterative Bregman projections for regularized transportation problems. *SIAM Journal on Scientific Computing*, 37(2):A1111–A1138, 2015.
- Bunne, C., Papaxanthos, L., Krause, A., and Cuturi, M. Proximal optimal transport modeling of population dynamics. In *International Conference on Artificial Intelligence and Statistics*, pp. 6511–6528. PMLR, 2022.
- Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., and Zagoruyko, S. End-to-end object detection with transformers. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part 1 16*, pp. 213–229. Springer, 2020.
- Carlier, G. On the linear convergence of the multimarginal Sinkhorn algorithm. *SIAM Journal on Optimization*, 32(2):786–794, 2022.
- Cattiaux, P., Conforti, G., Gentil, I., and Léonard, C. Time reversal of diffusion processes under a finite entropy condition. *arXiv preprint arXiv:2104.07708*, 2021.
- Chen, T., Liu, G.-H., and Theodorou, E. A. Likelihood training of Schrödinger bridge using forward-backward SDEs theory. *International Conference on Learning Representations*, 2022.
- Chen, T., Liu, G.-H., Tao, M., and Theodorou, E. A. Deep momentum multi-marginal Schrödinger bridge. 2023. doi: 10.48550/ARXIV.2303.01751.
- Chen, Y., Georgiou, T., Pavon, M., and Tannenbaum, A. Robust transport over networks. *IEEE Transactions on Automatic Control*, 62(9):4675–4682, 2016.
- Chen, Y., Conforti, G., and Georgiou, T. T. Measure-valued spline curves: An optimal transport viewpoint. *SIAM Journal on Mathematical Analysis*, 50(6):5947–5968, 2018.
- Chen, Y., Georgiou, T. T., and Pavon, M. Optimal transport in systems and control. *Annual Review of Control, Robotics, and Autonomous Systems*, 4, 2021.
- Chizat, L. Doubly regularized entropic Wasserstein barycenters. *arXiv preprint arXiv:2303.11844*, 2023.
- Claici, S., Chien, E., and Solomon, J. Stochastic Wasserstein barycenters. In *International Conference on Machine Learning*, pp. 999–1008. PMLR, 2018.

⁵https://github.com/JTT94/diffusion_schrodinger_bridge

- Cohen, S., Arbel, M., and Deisenroth, M. P. Estimating barycenters of measures in high dimensions. *arXiv preprint arXiv:2007.07105*, 2020.
- Cotar, C., Friesecke, G., and Klüppelberg, C. Density functional theory and optimal transportation with Coulomb cost. *Communications on Pure and Applied Mathematics*, 66(4):548–599, 2013.
- Csiszár, I. I-divergence geometry of probability distributions and minimization problems. *The Annals of Probability*, pp. 146–158, 1975.
- Cuturi, M. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in Neural Information Processing Systems*, 26, 2013.
- Cuturi, M. and Doucet, A. Fast computation of Wasserstein barycenters. In *International Conference on Machine Learning*, pp. 685–693. PMLR, 2014.
- Dai Pra, P. A stochastic control approach to reciprocal diffusion processes. *Applied Mathematics and Optimization*, 23(1):313–329, 1991.
- De Bortoli, V., Thornton, J., Heng, J., and Doucet, A. Diffusion Schrödinger bridge with applications to score-based generative modeling. *Advances in Neural Information Processing Systems*, 34: 17695–17709, 2021.
- Dupuis, P. and Ellis, R. S. *A weak convergence approach to the theory of large deviations*. John Wiley & Sons, 2011.
- Eisenberger, M., Toker, A., Leal-Taixé, L., and Cremers, D. Deep shells: Unsupervised shape correspondence with optimal transport. *Advances in Neural Information Processing Systems*, 33: 10491–10502, 2020.
- Fan, J., Taghvaei, A., and Chen, Y. Scalable computations of Wasserstein barycenter via input convex neural networks. *arXiv preprint arXiv:2007.04462*, 2020.
- Feydy, J., Charlier, B., Vialard, F.-X., and Peyré, G. Optimal transport for diffeomorphic registration. In *Medical Image Computing and Computer Assisted Intervention- MICCAI 2017: 20th International Conference, Quebec City, QC, Canada, September 11-13, 2017, Proceedings, Part I 20*, pp. 291–299. Springer, 2017.
- Flamary, R., Courty, N., Gramfort, A., Alaya, M. Z., Boisbunon, A., Chambon, S., Chapel, L., Corenflos, A., Fatras, K., Fournier, N., et al. POT: Python optimal transport. *The Journal of Machine Learning Research*, 22(1):3571–3578, 2021.
- Fortet, R. Résolution d’un système d’équations de M. Schrödinger. *Journal de Mathématiques Pures et Appliquées*, 19(1-4):83–105, 1940.
- Haasler, I., Ringh, A., Chen, Y., and Karlsson, J. Multimarginal Optimal Transport with a tree-structured cost and the Schrödinger bridge problem. *SIAM Journal on Control and Optimization*, 59(4):2428–2453, 2021.
- Hausmann, U. G. and Pardoux, E. Time reversal of diffusions. *The Annals of Probability*, pp. 1188–1205, 1986.
- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- Hyvärinen, A. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(4), 2005.
- Karcher, H. Riemannian center of mass and so called Karcher mean. *arXiv preprint arXiv:1407.2087*, 2014.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Knight, P. A. The Sinkhorn–Knopp algorithm: convergence and applications. *SIAM Journal on Matrix Analysis and Applications*, 30(1):261–275, 2008.

- Kober, H. A theorem on Banach spaces. *Compositio Mathematica*, 7:135–140, 1940.
- Koller, D. and Friedman, N. *Probabilistic Graphical Models: Principles and Techniques*. MIT press, 2009.
- Korotin, A., Li, L., Solomon, J., and Burnaev, E. Continuous Wasserstein-2 barycenter estimation without minimax optimization. *arXiv preprint arXiv:2102.01752*, 2021.
- Korotin, A., Egiazarian, V., Li, L., and Burnaev, E. Wasserstein iterative networks for barycenter estimation. *arXiv preprint arXiv:2201.12245*, 2022.
- Kullback, S. Probability densities with given marginals. *The Annals of Mathematical Statistics*, 39(4):1236–1243, 1968.
- LeCun, Y. The MNIST database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- Léonard, C. A survey of the Schrödinger problem and some of its connections with optimal transport. *Discrete & Continuous Dynamical Systems-A*, 34(4):1533–1574, 2014.
- Li, L., Genevay, A., Yurochkin, M., and Solomon, J. M. Continuous regularized Wasserstein barycenters. *Advances in Neural Information Processing Systems*, 33:17755–17765, 2020.
- Lipman, Y., Chen, R. T., Ben-Hamu, H., Nickel, M., and Le, M. Flow matching for generative modeling. *International Conference on Learning Representations*, 2023.
- Luo, Z.-Q. and Tseng, P. On the convergence of the coordinate descent method for convex differentiable minimization. *Journal of Optimization Theory and Applications*, 72(1):7–35, 1992.
- Makkuva, A., Taghvaei, A., Oh, S., and Lee, J. Optimal transport mapping via input convex neural networks. In *International Conference on Machine Learning*, pp. 6672–6681. PMLR, 2020.
- Marino, S. D. and Gerolin, A. An optimal transport approach for the Schrödinger bridge problem and convergence of Sinkhorn algorithm. *Journal of Scientific Computing*, 85(2):27, 2020.
- Mi, L., Yu, T., Bento, J., Zhang, W., Li, B., and Wang, Y. Variational Wasserstein barycenters for geometric clustering. *arXiv preprint arXiv:2002.10543*, 2020.
- Minsker, S., Srivastava, S., Lin, L., and Dunson, D. Scalable and robust Bayesian inference via the median posterior. In *International Conference on Machine Learning*, pp. 1656–1664. PMLR, 2014.
- Nichol, A. Q. and Dhariwal, P. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pp. 8162–8171, 2021.
- Nutz, M. Introduction to entropic optimal transport, 2021.
- Pele, O. and Werman, M. Fast and robust Earth mover’s distances. In *2009 IEEE 12th International Conference on Computer Vision*, pp. 460–467. IEEE, 2009.
- Peluchetti, S. Diffusion bridge mixture transports, Schrödinger bridge problems and generative modeling. *arXiv preprint arXiv:2304.00917*, 2023.
- Peyré, G., Cuturi, M., et al. Computational optimal transport: With applications to data science. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019.
- Rabin, J., Peyré, G., Delon, J., and Bernot, M. Wasserstein barycenter and its application to texture mixing. In *Scale Space and Variational Methods in Computer Vision: Third International Conference, SSVM 2011, Ein-Gedi, Israel, May 29–June 2, 2011, Revised Selected Papers 3*, pp. 435–446. Springer, 2012.
- Ruschendorf, L. Convergence of the Iterative Proportional Fitting procedure. *The Annals of Statistics*, pp. 1160–1174, 1995.
- Schiebinger, G., Shu, J., Tabaka, M., Cleary, B., Subramanian, V., Solomon, A., Gould, J., Liu, S., Lin, S., Berube, P., et al. Optimal-transport analysis of single-cell gene expression identifies developmental trajectories in reprogramming. *Cell*, 176(4):928–943, 2019.

- Schmitz, M. A., Heitz, M., Bonneel, N., Ngole, F., Coeurjolly, D., Cuturi, M., Peyré, G., and Starck, J.-L. Wasserstein dictionary learning: Optimal transport-based unsupervised nonlinear dictionary learning. *SIAM Journal on Imaging Sciences*, 11(1):643–678, 2018.
- Schrödinger, E. Sur la théorie relativiste de l'électron et l'interprétation de la mécanique quantique. *Annales de l'Institut Henri Poincaré*, 2(4):269–310, 1932.
- Shi, Y., De Bortoli, V., Campbell, A., and Doucet, A. Diffusion Schrödinger bridge matching. *arXiv preprint arXiv:2303.16852*, 2023.
- Sinkhorn, R. and Knopp, P. Concerning nonnegative matrices and doubly stochastic matrices. *Pacific Journal of Mathematics*, 21(2):343–348, 1967.
- Solomon, J., Rustamov, R., Guibas, L., and Butscher, A. Wasserstein propagation for semi-supervised learning. In *International Conference on Machine Learning*, pp. 306–314. PMLR, 2014.
- Solomon, J., De Goes, F., Peyré, G., Cuturi, M., Butscher, A., Nguyen, A., Du, T., and Guibas, L. Convolutional Wasserstein distances: Efficient optimal transportation on geometric domains. *ACM Transactions on Graphics (ToG)*, 34(4):1–11, 2015.
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., and Poole, B. Score-based generative modeling through stochastic differential equations. *International Conference on Learning Representations*, 2021.
- Srivastava, S., Li, C., and Dunson, D. B. Scalable Bayes via barycenter in Wasserstein space. *The Journal of Machine Learning Research*, 19(1):312–346, 2018.
- Staib, M., Claici, S., Solomon, J. M., and Jegelka, S. Parallel streaming Wasserstein barycenters. *Advances in Neural Information Processing Systems*, 30, 2017.
- Stroock, D. W. and Varadhan, S. S. *Multidimensional diffusion processes*, volume 233. Springer Science & Business Media, 1997.
- Su, Z., Wang, Y., Shi, R., Zeng, W., Sun, J., Luo, F., and Gu, X. Optimal mass transport for shape matching and comparison. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(11):2246–2259, 2015.
- Valiente, G. *Algorithms on Trees and Graphs*, volume 112. Springer, 2002.
- Vargas, F., Thodoroff, P., Lamacraft, A., and Lawrence, N. Solving Schrödinger bridges via maximum likelihood. *Entropy*, 23(9):1134, 2021.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 2017.
- Vincent, P. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674, 2011.

Appendix organization

First, additional notation is introduced in Appendix A. Then, we briefly recall some notions on undirected and directed trees in Appendix B. Similarly, martingale problems are introduced in Appendix C. The proofs of the main manuscript and additional theoretical results on Tree Schrödinger Bridges are given in Appendix D. Additional details on our consideration of the tree-based static SB problem are described in Appendix E. Details on the implementation of TreeDSB are given in Appendix F and the experiments are investigated in Appendix G.

A Additional notation

For any finite set E , we equivalently refer to the cardinal of E as $\text{card}(E)$ or $|E|$. Let (X, \mathcal{X}) be a measurable space. For any $x \in (\mathbb{R}^d)^{\ell+1}$ and any $m \in \{0, \dots, \ell\}$, let $x_{-m} = (x_0, \dots, x_{m-1}, x_{m+1}, \dots, x_\ell)$. For any family of measures $\{\nu_j\}_{j \in \{0, \dots, \ell\}}$ defined on (X, \mathcal{X}) and any $i \in \{0, \dots, \ell\}$, let $\nu_{-i} = \bigotimes_{j \in \{0, \dots, \ell\} \setminus \{i\}} \nu_j$. Let $I = \{i_1, \dots, i_q\} \subset \{1, \dots, \ell\}$ and $\mu \in \mathcal{P}^{(\ell)}$ such that $\mu \ll \text{Leb}$. We define $I^c = \{1, \dots, \ell\} \setminus I$ and denote it by $\{i_1^c, \dots, i_{\bar{q}}^c\}$ where $\bar{q} = \ell - q$. We denote the marginal of μ along I by μ_I , i.e., $\mu_I \in \mathcal{P}^{(q)}$ and we have for any $A \in \mathcal{B}((\mathbb{R}^d)^q)$, $\mu_I(A) = \int_X \mu(x) \prod_{j=1}^q \delta_{x_{i_j}}(A_j) dx$. In addition, note that $\mu_I \ll \text{Leb}$. We denote the conditional distribution of μ given I by $\mu_{|I}(\cdot|\cdot)$, i.e., $\mu_{|I}(\cdot|\cdot) \in \mathcal{P}^{(\bar{q})} \times (\mathbb{R}^d)^q$ and we have for any $y \in (\mathbb{R}^d)^q$ and any $A \in \mathcal{B}((\mathbb{R}^d)^{\bar{q}})$, $\mu_{|I}(A|y) = \int_X \mu(x) / \mu_I(y) \prod_{j=1}^q \delta(x_{i_j} - y_j) \prod_{j'=1}^{\bar{q}} \delta_{x_{i_{j'}^c}}(A_{j'}) dx$. Remark that for any $y \in (\mathbb{R}^d)^q$, $\mu_{|I}(\cdot|y) \ll \text{Leb}$. For any subset $J \subset I^c$ with $\text{card}(J) = q_J$, we also define $\mu_{|J}(\cdot|\cdot) \in \mathcal{P}^{(q_J)} \times (\mathbb{R}^d)^q$ such that for any $y \in (\mathbb{R}^d)^q$, $\mu_{|J}(\cdot|y) = \{\mu_{|I}(\cdot|y)\}_J$. For a collection of functions $\{f_i\}_{i \in I}$, with $I \subset \{1, \dots, n\}$ and $n \in \mathbb{N}$ such that $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$, we define $\bigoplus_{i \in I} f_i : (\mathbb{R}^d)^n \rightarrow \mathbb{R}$ such that for any $x = (x_1, \dots, x_n) \in (\mathbb{R}^d)^n$, $\bigoplus_{i \in I} f_i(x) = \sum_{i \in I} f_i(x_i)$.

B Introduction to trees

Undirected tree. An undirected graph $T = (V, E)$, with vertices V and edges E , is said to be an *undirected tree* if it is *acyclic* and *connected* (Valiente, 2002, Definition 1.19). In particular, we have $\text{card}(E) = \text{card}(V) - 1$. The undirected edge between two nodes v_1 and v_2 is similarly denoted by $\{v_1, v_2\}$ or $\{v_2, v_1\}$. We say that $T' = (V', E')$ is a *sub-tree* of T if T' is an undirected tree with vertices $V' \subset V$ and edges $E' \subset E$. For any vertex $v \in V$, we define the set of its *neighbours* N_v as the set of vertices $v' \in V$ such that $\{v, v'\} \in E$. The integer $\text{card}(N_v)$ is referred to as the *degree* of v . The vertices with degree 1 are called *leaves*, and we denote the set of leaves by $V_L \subset V$. The (unique) *path* in T between two vertices v and v' is the sequence of two-by-two distinct edges $\{\{v_i, v_{i+1}\}\}_{i=1}^n$ (with $n \geq 1$) such that $v_k = v_{k+1}$ for any $k \in \{1, \dots, n\}$ such that $k \equiv 0 \pmod{2}$, $v_1 = v$ and $v_{n+1} = v'$. This path can be seen as a linear sub-tree of T , and we define n as the *length* of this path. We say that T is *weighted* if there exists a map $w : E \mapsto \mathbb{R}_+$; in this case, $w(\{v_1, v_2\})$, or equivalently $w(\{v_2, v_1\})$ (also denoted by w_{v_1, v_2} or w_{v_2, v_1}) is called the *weight* of the edge $\{v_1, v_2\}$. The tree T is said to be *rooted* in $r \in V$ if r defines a partial ordering $\leq_{T,r} \subset V \times V$ such that for any $v_1, v_2 \in V$, $v_1 \leq_{T,r} v_2$ if the node v_1 lies on the unique path between r and v_2 .

Directed tree. Consider a directed graph $T_r = (V, E_r)$ rooted in $r \in V$. Any directed edge $e \in E_r$ from $v_1 \in V$ to $v_2 \in V$ is denoted by (v_1, v_2) . T_r is said to be a *directed tree* rooted in r if (i) the underlying undirected graph $T = (V, E)$ is an undirected tree rooted in r and (ii) any $(v_1, v_2) \in E_r$ is directed according to the partial ordering $\leq_{T,r}$, i.e., $\{v_1, v_2\} \in E$ and $v_1 \leq_{T,r} v_2$. For any vertices $(v, v') \in V \times V$ such that $v \leq_{T,r} v'$, the (unique) *path* in T_r from v to v' , denoted by $\text{path}_{T_r}(v, v')$, is defined as the directed version of the path in T between v and v' (viewed as a sub-tree of T), which is rooted in v . We say that T_r is *weighted*, if T is weighted and the edges of T_r have the same weights as the corresponding undirected edges of T . For any $(v_1, v_2) \in E_r$, we denote this weight by w_{v_1, v_2} . We say that T_r is the (unique) *directed version* of T rooted in r . It is endowed with a canonical vertex numbering $\zeta : V \rightarrow \{0, \dots, \text{card}(V) - 1\}$, corresponding to a depth-first traversal of its nodes, starting from the root r (Valiente, 2002, Definition 3.1.). This numbering is consistent with the partial ordering on T , i.e., if $v_1 \leq_{T,r} v_2$, $\zeta(v_1) \leq \zeta(v_2)$, and satisfies $\zeta(r) = 0$. In the rest of the paper, we will write in an equivalent manner v or $\zeta(v)$.

For any vertices $(v_1, v_2) \in E \times E$ such that $v_1 \leq_{T,r} v_2$, $\text{path}_{T_r}(v_1, v_2)$ corresponds to the ordered set of edges in E_r which define the ordered path between two vertices v_1 and v_2 . For any vertex $v \in V$, we define:

- (a) the set of its *children* C_v as the set of vertices $v' \in V$ such that $(v, v') \in E_r$. In particular, for any $v \in V_L$, the set of leaves, one has $C_v = \emptyset$.
- (b) its *parent* as the unique vertex $p(v)$ such that $(p(v), v) \in E_r$, if $v \neq r$ (the parent of the root is not defined).

Note that $N_r = C_r$ and, for any vertex $v \in V \setminus \{r\}$, $N_v = \{p(v)\} \cup C_v$.

Definition 8 (Tree-structured directed Probabilistic Graphical Model (PGM)). Consider a directed tree $T_r = (V, E_r)$. The directed PGM induced by T_r (Koller & Friedman, 2009, Definition 3.4.), denoted by \mathcal{P}_{T_r} , is the family of distributions $\pi \in \mathcal{P}^{(|V|)}$ which have a Markovian factorization along T_r , i.e.,

$$\mathcal{P}_{T_r} = \{ \pi \in \mathcal{P}^{(|V|)} : \pi = \pi_r \otimes_{(v,v') \in E_r} \pi_{v'|v} \}.$$

Lemma 9. Consider an undirected tree $T = (V, E)$. Let $(r, r') \in V \times V$. Let T' be a sub-tree of T with vertices V' such that $r' \in V'$. Denote by $T'_{r'}$ the directed version of T' rooted in r' . Then, for any $\pi \in \mathcal{P}_{T_r}$, we have $\pi_{V'} \in \mathcal{P}_{T'_{r'}}$.

Proof. Let $(r, r') \in V \times V$. We denote by $T_r = (V, E_r)$, respectively $T_{r'} = (V, E_{r'})$, the directed version of T rooted in r , respectively r' . We define the paths $P_{r,r'} = \text{path}_{T_r}(r, r') \subset E_r$ and $P_{r',r} = \text{path}_{T_{r'}}(r', r) \subset E_{r'}$. It is easy to see that

- (a) $E_r \setminus P_{r,r'} = E_{r'} \setminus P_{r',r}$,
- (b) $P_{r,r'} = \{(v_2, v_1) : (v_1, v_2) \in P_{r',r}\}$,
- (c) $P_{r',r} = \{(v_2, v_1) : (v_1, v_2) \in P_{r,r'}\}$.

Let $\pi \in \mathcal{P}_{T_r}$. First note that for any $(v_1, v_2) \in E_r$, we have by Bayes decomposition $\pi_{v_1} \pi_{v_2|v_1} = \pi_{v_2} \pi_{v_1|v_2} = \pi_{v_1, v_2}$. Then it comes

$$\begin{aligned} \pi &= \pi_r \otimes_{(v_1, v_2) \in E_r} \pi_{v_2|v_1} \\ &= \pi_r \otimes_{(v_1, v_2) \in P_{r,r'}} \pi_{v_2|v_1} \otimes_{(v_1, v_2) \in E_r \setminus P_{r,r'}} \pi_{v_2|v_1} \\ &= \pi_r \otimes_{(v_2, v_1) \in P_{r',r}} \pi_{v_2|v_1} \otimes_{(v_1, v_2) \in E_r \setminus P_{r',r}} \pi_{v_2|v_1} \\ &= \pi_{r'} \otimes_{(v_1, v_2) \in P_{r',r}} \pi_{v_2|v_1} \otimes_{(v_1, v_2) \in E_{r'} \setminus P_{r',r}} \pi_{v_2|v_1} \\ &= \pi_{r'} \otimes_{(v_1, v_2) \in E_{r'}} \pi_{v_2|v_1}, \end{aligned}$$

and therefore, we have $\pi \in \mathcal{P}_{T_{r'}}$.

Let T' be a sub-tree of T with vertices V' such that $r' \in V'$. First note that $E'_{r'} \subset E_{r'}$. Using the previous computation, we have for any $A \in \mathcal{B}((\mathbb{R}^d)^{|V'|})$,

$$\begin{aligned} \pi_{V'}(A) &= \int_{(\mathbb{R}^d)^{|V|}} \pi_{r'}(x_{r'}) \otimes_{(v_1, v_2) \in E_{r'}} \pi_{v_2|v_1}(x_{v_2}|x_{v_1}) \prod_{v' \in V'} \delta_{x_{v'}}(A_{v'}) dx \\ &= \int_{(\mathbb{R}^d)^{|V| - |V'|} } \{ \pi_{r'}(A_{r'}) \otimes_{(v_1, v_2) \in E'_{r'}} \pi_{v_2|v_1}(A_{v_2}|x_{v_1}) \} \otimes_{(v_1, v_2) \in E_{r'} \setminus E'_{r'}} \pi_{v_2|v_1}(x_{v_2}|x_{v_1}) dx_{V \setminus V'} \\ &= \{ \pi_{r'} \otimes_{(v_1, v_2) \in E'_{r'}} \pi_{v_2|v_1} \}(A), \end{aligned}$$

which proves that $\pi_{V'} \in \mathcal{P}_{T'_{r'}}$. □

Discretized undirected tree. Let $N \geq 1$. Consider an undirected tree $T = (V, E)$ with weights w . We say that $T^{(N)} = (V^{(N)}, E^{(N)})$ is a N -discretized version of T if it is an undirected tree with weights $w^{(N)}$ such that

- (a) $V^{(N)} = V \bigsqcup \bigcup_{\substack{e \in E, \\ k \in \{1, \dots, N-1\}}} \{v_e^k\}$,

(b) $E^{(N)} = \cup_{e \in E} \cup_{k=0, \dots, N-1} \{v_e^k, v_e^{k+1}\}$ with the convention that the vertices v_e^N and v_e^0 are defined such that $\{v_e^0, v_e^N\} = e$,

(c) $\sum_{e \in \text{path}_T(v, v')} 1/w_e^{(N)} = 1/w_{v, v'}$, if $\{v, v'\} \in E$.

Remark that the leaves of $T^{(N)}$ are exactly the original leaves of T and that $T^{(1)} = T$. The non-uniqueness of $T^{(N)}$ comes from the freedom of choice on the weights of its edges.

Discretized directed tree. Let $N \geq 1$. Consider a directed tree $T_r = (V, E_r)$ rooted in $r \in V$ with weights w . We say that $T_r^{(N)} = (V^{(N)}, E_r^{(N)})$ is a N -discretized version of T_r if it is the directed version of $T^{(N)}$ rooted in r , where $T^{(N)}$ is a N -discretized version of the underlying undirected tree of T_r .

C Background on martingale problems

In this section, we introduce the background on Stochastic Differential Equations (SDEs) and weak solutions of SDEs following the framework of (Stroock & Varadhan, 1997, Section 10.1, page 249). We recall that $C_0^\infty(\mathbb{R}^d)$ is the space of infinitely differentiable real-valued functions which vanish at infinity. In addition, we have that S_+^d is the space of $d \times d$, symmetric, non-negative matrices.

Definition 10. Let $T > 0$ or $T = +\infty$, $\sigma : [0, T) \times \mathbb{R}^d \rightarrow S_+^d$ and $b : [0, T) \times \mathbb{R}^d \rightarrow \mathbb{R}^d$, locally bounded measurable functions. We define the infinitesimal generator, \mathcal{A} , given for any $f \in C_0^\infty(\mathbb{R}^d)$, $t \in [0, T)$ and $x \in \mathbb{R}^d$ by

$$\mathcal{A}_t(f)(x) = \langle b_t(x), \nabla f(x) \rangle + \frac{1}{2} \langle \sigma_t(x) \sigma_t(x)^\top, \nabla^2 f(x) \rangle. \quad (4)$$

We say that a probability measure \mathbb{P} satisfies the martingale problem for \mathcal{A} if for any $t \in [0, T)$ and $f \in C_0^\infty(\mathbb{R}^d)$, we have that $(f(\mathbf{X}_t) - \int_0^t \mathcal{A}_s(f)(\mathbf{X}_s) ds)_{s \in [0, t]}$ is a \mathbb{P} -martingale.

In the main document, see Section 2, we say that “a path measure \mathbb{P} is associated with $d\mathbf{X}_t = b(t, \mathbf{X}_t)dt + \sigma(t, \mathbf{X}_t)d\mathbf{B}_t$ with $(\mathbf{B}_t)_{t \geq 0}$ a d -dimensional Brownian motion” if \mathbb{P} solves the martingale problem associated with \mathcal{A} given by (4). Unless specified, we always assume that such a path measure exists and is unique. Below, we recall the following theorem, see (Stroock & Varadhan, 1997, Theorem 10.2.2), which gives sufficient conditions for the existence and uniqueness of solutions to the martingale problem.

Theorem 11. Assume that for any $x \in \mathbb{R}^d$ we have

$$\inf\{\langle \theta, \sigma \sigma^\top(s, x) \theta \rangle : \theta \in \mathbb{R}^d, \|\theta\| = 1, s \in [0, T]\} > 0, \\ \limsup_{y \rightarrow x} \{\|\sigma(s, x) - \sigma(s, y)\|\} = 0.$$

In addition, assume that there exists $C > 0$ such that for any $x \in \mathbb{R}^d$

$$\sup\{\|\sigma \sigma^\top(t, x)\| : s \in [0, T]\} + \sup\{\langle x, b(t, x) \rangle : s \in [0, T]\} \leq C(1 + \|x\|^2).$$

Then, there exists a unique solution to the martingale problem with initialization $x_0 \in \mathbb{R}^d$.

D Theoretical results on Tree Schrödinger Bridges

We respectively provide in Appendix D.1, Appendix D.2 and Appendix D.3 the proofs of the results of the main manuscript presented in Section 3, Section 4 and Section 5. Finally, we make a detailed comparison between our setting and the framework of Haasler et al. (2021) in Appendix D.4. In the rest of this section, we consider an undirected tree $T = (V, E)$, where $|V| = \ell + 1$, and some subset $S \subset V$ which we denote by $S = \{i_0, \dots, i_{K-1}\}$. We define $S^c = V \setminus S$.

D.1 Proofs of Section 3

Proposition 1 is straightforward to obtain by combining the definition of the Brownian motion with the definition of π^0 given in (2). The following lemma details the recursion relation between the (mIPF) iterates, which is key to prove Proposition 2.

Lemma 12. Let $(\pi^n)_{n \in \mathbb{N}}$ be the sequence given by (mIPF). Let $n \in \mathbb{N}$, $k_n = (n - 1) \bmod(K)$, $k_n + 1 = n \bmod(K)$. Denote by \mathbb{T}_{k_n} , respectively \mathbb{T}_{k_n+1} with edges E_{k_n+1} , the directed version of \mathbb{T} rooted in i_{k_n} , respectively in i_{k_n+1} . We have:

(i) $\pi^n \in \mathcal{P}_{\mathbb{T}_{k_n}}$,

(ii) $\pi^{n+1} = \mu_{i_{k_n+1}} \otimes_{(v,v') \in E_{k_n+1}} \pi_{v'|v}^n$. In particular, for any $(v, v') \in E_{k_n+1}$, $\pi_{v'|v}^{n+1} = \pi_{v'|v}^n$.

Proof. We show the result (i) by recursion on $n \in \mathbb{N}$, and will deduce (ii) from the proof. Using (2), we first have $\pi^0 \in \mathcal{P}_{\mathbb{T}_r}$, where r is chosen as i_{K-1} , see Section 3. Thus, we obtain the result (i) at step $n = 0$. Assume now that $\pi^n \in \mathcal{P}_{\mathbb{T}_{k_n}}$ for some $n \in \mathbb{N}$.

Consider the paths $P_n = \text{path}_{\mathbb{T}_{k_n}}(i_{k_n}, i_{k_n+1})$ and $P_{n+1} = \text{path}_{\mathbb{T}_{k_n+1}}(i_{k_n+1}, i_{k_n})$. Note that these two paths have the same length, denoted by J , and contain the same vertices, denoted by V_n . Let $\pi \in \mathcal{P}^{(\ell+1)}$ such that $\text{KL}(\pi|\pi^n) < +\infty$. We have the following decomposition

$$\text{KL}(\pi|\pi^n) = \text{KL}(\pi_{V_n}|\pi_{V_n}^n) + \int_{(\mathbb{R}^d)^{J+1}} \text{KL}(\pi|_{V_n}|\pi_{V_n}^n) d\pi_{V_n}(x_{V_n}).$$

Hence, the $(n + 1)$ -th iterate of (mIPF) is given by $\pi^{n+1} = \pi_{V_n}^{n+1} \otimes \pi_{V_n}^n$, with

$$\pi_{V_n}^{n+1} = \text{argmin}\{\text{KL}(\pi|\pi_{V_n}^n) : \pi \in \mathcal{P}^{(J+1)}, \pi_{i_{k_n+1}} = \mu_{i_{k_n+1}}\}.$$

Since $\pi^n \in \mathcal{P}_{\mathbb{T}_{k_n}}$, we have (i) $\pi_{V_n}^n = \otimes_{(v,v') \in E_{k_n} \setminus P_n} \pi_{v'|v}^n$ and (ii) $\pi_{V_n}^n \in \mathcal{P}_{P_{n+1}}$ by Lemma 9, where P_{n+1} is viewed as a directed tree rooted in i_{k_n+1} . Defining $V_{n+1} = V_n \setminus \{i_{k_n+1}\}$, we thus have $\pi_{V_n}^n = \pi_{i_{k_n+1}}^n \otimes \pi_{V_{n+1}|i_{k_n+1}}^n$ where $\pi_{V_{n+1}|i_{k_n+1}}^n = \otimes_{(v,v') \in P_{n+1}} \pi_{v'|v}^n$.

Let $\pi \in \mathcal{P}^{(J+1)}$ such that $\pi_{i_{k_n+1}} = \mu_{i_{k_n+1}}$ and $\text{KL}(\pi|\pi_{V_n}^n) < +\infty$. Similarly to the previous computation, we have the following decomposition

$$\begin{aligned} \text{KL}(\pi|\pi_{V_n}^n) &= \text{KL}(\pi_{i_{k_n+1}}|\pi_{i_{k_n+1}}^n) + \int_{\mathbb{R}^d} \text{KL}(\pi|_{i_{k_n+1}}|\pi_{V_{n+1}|i_{k_n+1}}^n) d\pi_{i_{k_n+1}}(x_{i_{k_n+1}}) \\ &= \text{KL}(\mu_{i_{k_n+1}}|\pi_{i_{k_n+1}}^n) + \int_{\mathbb{R}^d} \text{KL}(\pi|_{i_{k_n+1}}|\pi_{V_{n+1}|i_{k_n+1}}^n) d\mu_{i_{k_n+1}}(x_{i_{k_n+1}}). \end{aligned}$$

Therefore, we obtain

$$\pi_{V_n}^{n+1} = \mu_{i_{k_n+1}} \otimes \pi_{V_{n+1}|i_{k_n+1}}^n = \mu_{i_{k_n+1}} \otimes_{(v,v') \in P_{n+1}} \pi_{v'|v}^n.$$

Noting that $E_{k_n} \setminus P_n = E_{k_n+1} \setminus P_{n+1}$ and recalling that $\pi^{n+1} = \pi_{V_n}^{n+1} \otimes \pi_{V_n}^n$, it finally comes

$$\pi^{n+1} = \mu_{i_{k_n+1}} \otimes_{(v,v') \in P_{n+1}} \pi_{v'|v}^n \otimes_{(v,v') \in E_{k_n+1} \setminus P_{n+1}} \pi_{v'|v}^n = \mu_{i_{k_n+1}} \otimes_{(v,v') \in E_{k_n+1}} \pi_{v'|v}^n. \quad (5)$$

Therefore, $\pi^{n+1} \in \mathcal{P}_{\mathbb{T}_{k_n+1}}$, which achieves the recursion for (i), and we obtain (ii) by (5). \square

Hence, Lemma 12 shows that the (mIPF) iterates admit a Markovian factorization on \mathbb{T} , and can be defined recursively using the edges of \mathbb{T} . We now provide the proof of Proposition 2.

Proof of Proposition 2. We will prove this result by recursion on $n \in \mathbb{N}$. Observe that the initialisation is directly given by Proposition 1. Assume now that the result of Proposition 2 stands for some $n \in \mathbb{N}$. Let $k_n = (n - 1) \bmod(K)$, $k_n + 1 = n \bmod(K)$. Denote by \mathbb{T}_{k_n} with edges E_{k_n} , respectively \mathbb{T}_{k_n+1} with edges E_{k_n+1} , the directed version of \mathbb{T} rooted in i_{k_n} , respectively in i_{k_n+1} . For any vertex v of \mathbb{T}_{k_n+1} , we define $p(v)$ as the (unique) parent of v and $c(v)$ as the unique child of v when it exists. Consider the $(n + 1)$ -th dynamic iterate defined by (a) and (b), i.e., $(\mathbb{P}_{(v,v') \in E_{k_n+1}}^{n+1})_{(v,v') \in E_{k_n+1}}$. To prove that this iterate has the properties stated in Proposition 2, we proceed by recursion on the edges of \mathbb{T}_{k_n+1} , following the bread-first order in \mathbb{T}_{k_n+1} . In this order, the edge $(i_{k_n+1}, c(i_{k_n+1}))$ is the first to be considered. Remark that $c(i_{k_n+1})$ is well defined since i_{k_n+1} is a leaf of \mathbb{T} .

Here, we denote $T_{c(i_{k_n+1}), i_{k_n+1}}$ by T . By construction, we have $\mathbb{P}_{(i_{k_n+1}, c(i_{k_n+1}))}^{n+1} = \mu_{i_{k_n+1}} \otimes (\mathbb{P}_{(c(i_{k_n+1}), i_{k_n+1})}^n)^R|_0$. By recursion assumption, $\mathbb{P}_{(c(i_{k_n+1}), i_{k_n+1})}^n \in \mathcal{P}(C([0, T], \mathbb{R}^d))$ since

$(c(i_{k_n+1}), i_{k_n+1}) \in E_{k_n}$. Then, $\mathbb{P}_{(i_{k_n+1}, c(i_{k_n+1}))}^{n+1}$ is a well defined path measure on $[0, T]$. By definition of the (mIPF) sequence, we have $\mu_{i_{k_n+1}} = \pi_{i_{k_n+1}}^{n+1}$. By recursion assumption, we also have that $\text{Ext}(\mathbb{P}_{(c(i_{k_n+1}), i_{k_n+1})}^n) = \pi_{c(i_{k_n+1}), i_{k_n+1}}^n$. Hence, it comes that $(\mathbb{P}_{(c(i_{k_n+1}), i_{k_n+1})}^n)^R_{T|0} = \pi_{c(i_{k_n+1}), i_{k_n+1}}^n = \pi_{c(i_{k_n+1})|i_{k_n+1}}^{n+1}$, where the last equality comes from Lemma 12. Finally, we obtain that $\text{Ext}(\mathbb{P}_{(i_{k_n+1}, c(i_{k_n+1}))}^{n+1}) = \pi_{i_{k_n+1}, c(i_{k_n+1})}^{n+1}$, which proves the initialisation.

Assume now that \mathbb{P}^{n+1} is well defined and has the right properties, up to some edge in \mathbb{T}_{k_n+1} . Consider the following edge, denoted by $(v, v') \in E_{k_n+1}$, in the breadth-first order. By edge recursion, we have that $\text{Ext}(\mathbb{P}_{(p(v), v)}^{n+1}) = \pi_{p(v), v}^{n+1}$, and thus $\mathbb{P}_{(p(v), v), T_{p(v), v}}^{n+1} = \pi_v^{n+1}$. Define the path $P_n = \text{path}_{\mathbb{T}_{k_n}}(i_{k_n}, i_{k_n+1})$. Then, we face two cases.

(i) Either $(v, v') \in E_{k_n} \setminus P_n$. Then, we have by (a) that

$$\mathbb{P}_{(v, v')}^{n+1} = \mathbb{P}_{(p(v), v), T_{p(v), v}}^{n+1} \otimes \mathbb{P}_{(v, v')|0}^n = \pi_v^{n+1} \otimes \mathbb{P}_{(v, v')|0}^n$$

In particular, $\mathbb{P}_{(v, v')}^{n+1}$ is a well defined path measure on $[0, T_{v, v'}]$. Since $(v, v') \in E_{k_n}$, $\text{Ext}(\mathbb{P}_{(v, v')}^n) = \pi_{v, v'}^n$ by recursion assumption. In particular, $\mathbb{P}_{(v, v'), T_{v, v'}|0}^n = \pi_{v'|v}^n = \pi_{v'|v}^{n+1}$ where the last equality comes from Lemma 12. We thus have $\text{Ext}(\mathbb{P}_{(v, v')}^{n+1}) = \pi_{v, v'}^{n+1}$.

(ii) Or $(v', v) \in P_n$. Then, we have by (b) that

$$\mathbb{P}_{(v, v')}^{n+1} = \mathbb{P}_{(p(v), v), T_{v, v'}}^{n+1} \otimes (\mathbb{P}_{(v', v)}^n)^R_{|0} = \pi_v^{n+1} \otimes (\mathbb{P}_{(v', v)}^n)^R_{|0}$$

In particular, $\mathbb{P}_{(v, v')}^{n+1}$ is a well defined path measure on $[0, T_{v, v'}]$. Here, $(v', v) \in E_{k_n}$ and thus, $\text{Ext}(\mathbb{P}_{(v', v)}^n) = \pi_{v', v}^n$ by recursion assumption. In particular, $(\mathbb{P}_{(v', v)}^n)^R_{T_{v', v}|0} = \pi_{v'|v}^n = \pi_{v'|v}^{n+1}$ where the last equality comes from Lemma 12. We thus have $\text{Ext}(\mathbb{P}_{(v, v')}^{n+1}) = \pi_{v, v'}^{n+1}$.

This achieves the recursion. \square

D.2 Proofs of Section 4

Remark on assumption A1. Although A1 is not needed to establish the result of Proposition 3, Corollary 4 and Proposition 5, it is however crucial in the proof of convergence of (mIPF) stated in Proposition 6. Nevertheless, we choose to keep A1 as an assumption in the statement of every theoretical result presented in Section 4 for sake of clarity.

Additional definitions. We define the set $\mathcal{P}_S = \cap_{i \in S} \mathcal{P}_i$, where $\mathcal{P}_i = \{\pi \in \mathcal{P}^{(\ell+1)} : \pi_i = \mu_i\}$, i.e., \mathcal{P}_S is the set of all probability measures $\pi \in \mathcal{P}^{(\ell+1)}$ which verify

$$\int_{(\mathbb{R}^d)^{\ell+1}} f_i(x_i) d\pi(x_{0:\ell}) = \int_{\mathbb{R}^d} f_i(x_i) d\mu_i(x_i),$$

for any family of bounded measurable functions $\{f_i\}_{i \in S} \in C(\mathbb{R}^d, \mathbb{R})^K$. Since \mathbb{R}^d is separable, there exists a dense family of functions $\{f_i^k\}_{k \in \mathbb{N}^*, i \in S}$, with $f_i^k \in L^\infty(\mu_i)$ for any $k \in \mathbb{N}^*$ and any $i \in S$, such that $\pi \in \mathcal{P}_S$ if and only if

$$\int_{(\mathbb{R}^d)^{\ell+1}} f_i^k(x_i) d\pi(x_{0:\ell}) = \int_{\mathbb{R}^d} f_i^k(x_i) d\mu_i(x_i)$$

or equivalently, upon centering f_i^k ,

$$\int_{(\mathbb{R}^d)^{\ell+1}} f_i^k(x_i) d\pi(x_{0:\ell}) = 0.$$

In the rest of the section, we consider such family $\{f_i^k\}_{k \in \mathbb{N}^*, i \in S}$.

For any $n \in \mathbb{N}^*$, we also define $\mathcal{P}_S^n = \cap_{i \in S} \mathcal{P}_i^n$, where $\mathcal{P}_i^n = \{\pi \in \mathcal{P}^{(\ell+1)} : \int_{(\mathbb{R}^d)^{\ell+1}} f_i^k(x_i) d\pi(x_{0:\ell}) = 0, \forall k \in \{1, \dots, n\}\}$. In particular, we have

$$\mathcal{P}_S = \cap_{n \in \mathbb{N}^*} \mathcal{P}_S^n. \quad (6)$$

Finally, (static-mSB) can be rewritten as

$$\pi^* = \text{argmin}\{\text{KL}(\pi | \pi^0) : \pi \in \mathcal{P}_S\}. \quad (7)$$

Proof of Proposition 3 and Corollary 4. In this part of the section, we present an extension of the theoretical results from Nutz (2021) to the multi-marginal setting. We first present two technical results, Lemma 13 and Lemma 14, which are respectively adapted from (Nutz, 2021, Lemma 2.10.) and (Nutz, 2021, Lemma 2.11.).

Lemma 13. Let $\{\tilde{\mu}_j\}_{j \in S^c}$ be a family of probability measures defined on $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$. We define $\tilde{\pi}^0 = \bigotimes_{i \in S} \mu_i \bigotimes_{j \in S^c} \tilde{\mu}_j$. Let $A \in \bigotimes_{m=0}^{\ell} \mathcal{B}(\mathbb{R}^d)$ such that $\tilde{\pi}^0(A) = 1$. Then, for $\tilde{\pi}^0$ -almost any $x^* \in A$, there exists a family of sets $\{X_m^0\}_{m=0}^{\ell} \subset (\mathbb{R}^d)^{\ell+1}$ such that

(a) $\mu_i(X_i^0) = 1$ for any $i \in S$, and $\tilde{\mu}_j(X_j^0) = 1$ for any $j \in S^c$,

(b) $A^0 = A \cap (\prod_{m=0}^{\ell} X_m^0)$ satisfies $x^* \in A^0$ and

$$(x_0^*, \dots, x_{m-1}^*, x_m, x_{m+1}^*, \dots, x_{\ell}^*) \in A^0, \forall x \in A^0, \forall m \in \{0, \dots, \ell\}.$$

Proof. Consider such set A . We define for any $m \in \{0, \dots, \ell\}$ the set

$$X_m = \{u \in \mathbb{R}^d : \tilde{\pi}_{-m}^0(A_m^u) = 1\},$$

where $A_m^u = \{y \in (\mathbb{R}^d)^{\ell} : (y_0, \dots, y_{m-1}, u, y_m, \dots, y_{\ell-1}) \in A\}$.

Take $i \in S$. Assume that $\mu_i(X_i) < 1$. We recall that $\tilde{\pi}^0 = \tilde{\pi}_{-i}^0 \otimes \mu_i$. Using Fubini's theorem and that $\int_{A_i^{x_i}} d\tilde{\pi}_{-i}^0(x_{-i}) < 1$ for any $x_i \notin X_i$, we have

$$\begin{aligned} 1 = \tilde{\pi}^0(A) &= \int_A d\tilde{\pi}_{-i}^0(x_{-i}) \otimes d\mu_i(x_i) \\ &= \int_{\mathbb{R}^d} \left\{ \int_{A_i^{x_i}} d\tilde{\pi}_{-i}^0(x_{-i}) \right\} d\mu_i(x_i) \\ &= \int_{X_i} \left\{ \int_{A_i^{x_i}} d\tilde{\pi}_{-i}^0(x_{-i}) \right\} d\mu_i(x_i) + \int_{X_i^c} \left\{ \int_{A_i^{x_i}} d\tilde{\pi}_{-i}^0(x_{-i}) \right\} d\mu_i(x_i) \\ &< \mu_i(X_i) + \mu_i(X_i^c) = 1, \end{aligned}$$

which is absurd. Therefore, we obtain $\mu_i(X_i) = 1$, and similarly, we have $\tilde{\mu}_j(X_j) = 1$ for any $j \in S^c$. For any $y \in (\mathbb{R}^d)^{\ell}$, any $m \in \{0, \dots, \ell\}$, we define the set

$$\bar{A}_m^y = \{u \in \mathbb{R}^d : (y_0, \dots, y_{m-1}, u, y_m, \dots, y_{\ell-1}) \in A\}.$$

Let $i \in S$. We have by Fubini's theorem

$$\begin{aligned} 1 = \tilde{\pi}^0(A) &= \int_A d\mu_i(x_i) \otimes d\tilde{\pi}_{-i}^0(x_{-i}) \\ &= \int_{(\mathbb{R}^d)^{\ell}} \left\{ \int_{\bar{A}_i^{x_{-i}}} d\mu_i(x_i) \right\} d\tilde{\pi}_{-i}^0(x_{-i}) \\ &= \int_{\prod_{\substack{m=0 \\ m \neq i}}^{\ell} X_i} \left\{ \int_{\bar{A}_i^{x_{-i}}} d\mu_i(x_i) \right\} d\tilde{\pi}_{-i}^0(x_{-i}), \end{aligned}$$

where the last equality comes from the fact that $\mu_i(X_i) = 1$ for any $i \in S$, $\tilde{\mu}_j(X_j) = 1$ for any $j \in S^c$ and that $\tilde{\pi}^0 = \bigotimes_{i \in S} \mu_i \bigotimes_{j \in S^c} \tilde{\mu}_j$. Consequently, there exists a measurable set $A_{-i} \subset \prod_{\substack{m=0 \\ m \neq i}}^{\ell} X_i$ such that the following properties hold: (a) $\mu_i(\bar{A}_i^y) = 1$ for any $y \in A_{-i}$, (b) $\tilde{\pi}_{-i}^0(A_{-i}) = 1$. Similarly, this result holds for any $j \in S^c$, i.e., there exists a measurable set $A_{-j} \subset \prod_{\substack{m=0 \\ m \neq j}}^{\ell} X_i$ such that the following properties hold: (a) $\tilde{\mu}_j(\bar{A}_j^y) = 1$ for any $y \in A_{-j}$, (b) $\tilde{\pi}_{-j}^0(A_{-j}) = 1$. We consider such sets $\{A_{-m}\}_{m=0}^{\ell}$ for the rest of the proof and finally define the set

$$\tilde{A} = \bigcap_{m=0}^{\ell} \tilde{A}_m,$$

where $\tilde{A}_m = A_{-m} \times \{u \in \bar{A}_m^y : y \in A_{-m}\}$. By definition, we have $\tilde{A} \subset A \cap \prod_{m=0}^{\ell} X_m$, using the fact that $\tilde{A}_m \subset A$ for any $m \in \{0, \dots, \ell\}$. In addition, for any $i \in S$, we get by Fubini's theorem

$$\tilde{\pi}^0(\tilde{A}_i) = \int_{\tilde{A}_i} d\mu_i(x_i) \otimes d\tilde{\pi}_{-i}^0(x_{-i}) = \int_{A_{-i}} \left\{ \int_{\bar{A}_i^{x_{-i}}} d\mu_i(x_i) \right\} d\tilde{\pi}_{-i}^0(x_{-i}) = \tilde{\pi}_{-i}^0(A_{-i}) = 1,$$

and similarly, we get $\tilde{\pi}^0(\tilde{A}_j) = 1$ for any $j \in S^c$. We can deduce that $\tilde{\pi}^0(\tilde{A}) = 1$ since $\tilde{\pi}^0(\tilde{A}^c) \leq \sum_{m=0}^{\ell} \tilde{\pi}^0(\tilde{A}_m^c) = 0$.

Let $x^* \in \tilde{A}$. In particular, $x^* \in A$. We define the set $A^0 = A \cap (\prod_{m=0}^{\ell} X_m^0)$, where $X_m^0 = X_m \cap \bar{A}_m^{x^*-m}$ for any $m \in \{0, \dots, \ell\}$. We now establish the result of Lemma 13.

We first prove (a). Let $i \in S$. Since $x^* \in \tilde{A}$, we have $x^* \in \tilde{A}_i$ and therefore $x_{-i}^* \in A_{-i}$. By definition of A_{-i} , we obtain that $\mu_i(\bar{A}_i^{x_{-i}^*}) = 1$ and thus,

$$\mu_i(\{X_i^0\}^c) \leq \mu_i(X_i^c) + \mu_i(\{\bar{A}_i^{x_{-i}^*}\}^c) = 0,$$

which gives $\mu_i(X_i^0) = 1$, and similarly, we have $\tilde{\mu}_j(X_j^0) = 1$ for any $j \in S^c$.

We now prove (b). Let $m \in \{0, \dots, \ell\}$. Since $x^* \in \tilde{A} \subset A$, we get $x_m^* \in \bar{A}_m^{x^*-m}$. Using that $\tilde{A} \subset A \cap \prod_{m=0}^{\ell} X_m$, we get $x^* \in A^0$. Let $x \in A^0$. We denote $x^m = (x_0^*, \dots, x_{m-1}^*, x_m, x_{m+1}^*, \dots, x_{\ell}^*)$. We need to show that $x^m \in A$ and $x^m \in \prod_{j=1}^{\ell} X_j^0 = \prod_{j=1}^{\ell} (X_j \cap \bar{A}_j^{x^*-m})$. First, since $x_j^m = x_j$ or x_j^* for any $j \in \{0, \dots, \ell\}$, and $x \in A^0$ and $x^* \in A^0$, we get that for any $j \in \{0, \dots, \ell\}$, $x_j^m \in X_j$. Similarly, for any $j \in \{0, \dots, \ell - 1\}$, $x_j^m \in \bar{A}_j^{x^*-m}$. Therefore, we get that $x^m \in \prod_{j=1}^{\ell} (X_j \cap \bar{A}_j^{x^*-m})$. Since $x_m \in A_m^{x^*-m}$ (because $x \in \prod_{j=1}^{\ell} (X_j \cap \bar{A}_j^{x^*-m})$), we get that $x \in A$, which concludes the proof. \square

Lemma 14. Let $A^0 \subset (\mathbb{R}^d)^{\ell+1}$. For any $m \in \{0, \dots, \ell\}$, we denote $X_m^0 = \text{proj}_m(A^0)$. We make the following assumptions.

(a) Assume there exists $x^* \in A^0$ such that for any $x \in A^0$, for any $m \in \{0, \dots, \ell\}$, we have $(x_0^*, \dots, x_{m-1}^*, x_m, x_{m+1}^*, \dots, x_{\ell}^*) \in A^0$.

(b) Assume there exists a family of functions $\{\varphi_{i_k}^n\}_{n \in \mathbb{N}^*, k \in \{0, \dots, K-1\}}$ with $\varphi_{i_k}^n : X_{i_k}^0 \rightarrow [-\infty, +\infty)$ such that for any $n \in \mathbb{N}^*$ and any $k \in \{0, \dots, K-2\}$, we have $\varphi_{i_k}^n(x_{i_k}^*) = 0$.

(c) Denote $F^n(x) = \sum_{k=0}^{K-1} \varphi_{i_k}^n(x_{i_k})$ for any $x \in A^0$. Assume that for any $x \in A^0$, $F(x) = \lim_{n \rightarrow \infty} F^n(x)$ exists and is such that $F(x) \in [-\infty, +\infty)$ with $F(x^*) \in \mathbb{R}$.

Then, for any $i \in S$, for any $x_i \in X_i^0$, $\varphi_i(x_i) = \lim_{n \rightarrow \infty} \varphi_i^n(x_i)$ exists and is such that $\varphi_i(x_i) \in [-\infty, +\infty)$.

Proof. Consider $A^0 \subset (\mathbb{R}^d)^{\ell+1}$ such that assumptions (a), (b) and (c) hold. Remark that we have $F^n(x^*) = \varphi_{i_{K-1}}^n(x_{i_{K-1}}^*)$.

Let $x \in A^0$. We denote $x^m = (x_1^*, \dots, x_{m-1}^*, x_m, x_{m+1}^*, \dots, x_{\ell}^*)$ for any $m \in \{0, \dots, \ell\}$. In particular, we have $x^m \in A^0$ by assumption (a). Let us define

$$\begin{aligned} \varphi_{i_k}(x_{i_k}) &= F(x^{i_k}) - F(x^*), \quad \forall k \in \{0, \dots, K-2\}, \\ \varphi_{i_{K-1}}(x_{i_{K-1}}) &= F(x^{i_{K-1}}). \end{aligned}$$

Using assumption (c), we have $\varphi_i(x_i) \in [-\infty, +\infty)$ for any $i \in S$. Let $k \in \{0, \dots, K-2\}$. We have by definition of F^n ,

$$\varphi_{i_k}^n(x_{i_k}) = F^n(x^{i_k}) - \sum_{\substack{m=0 \\ m \neq k}}^{K-1} \varphi_{i_m}^n(x_{i_m}^*) = F^n(x^{i_k}) - F^n(x^*),$$

where we used assumption (b) in the last equality. Since $x^{i_k} \in A^0$ and $x^* \in A^0$, we have by assumption (c),

$$\lim_{n \rightarrow \infty} \varphi_{i_k}^n(x_{i_k}) = F(x^{i_k}) - F(x^*) = \varphi_{i_k}(x_{i_k}).$$

Furthermore, by combining the definition of F^n with assumption (b), we have

$$\lim_{n \rightarrow \infty} \varphi_{i_{K-1}}^n(x_{i_{K-1}}) = F(x^{i_{K-1}}) = \varphi_{i_{K-1}}(x_{i_{K-1}}),$$

which concludes the proof. \square

In what follows, before proving Proposition 3, we respectively show in Proposition 15 and Proposition 16 how A2 and A3 can be satisfied in the case where $\pi^0 \in \mathcal{P}_{T_r}$, as in (2), that is

$$\pi^0 = \pi_r^0 \otimes_{(v,v') \in E_r} \pi_{v'|v}^0.$$

Proposition 15. Let $\pi^0 \in \mathcal{P}_{T_r}$. Assume that $\pi_r^0 = \mu_r$ if $r \in S$ or $\pi_r^0 = N(m_r, \sigma_r \text{Id})$, with $m_r \in \mathbb{R}^d$ and $\sigma_r > 0$ if $r \in S^c$. In addition, assume that for any $(v, v') \in E_r$, $\pi_{v'|v}^0(\cdot | x_v) = N(x_v, \sigma_{v,v'} \text{Id})$ with $\sigma_{v,v'} > 0$. Finally, assume that for any $i \in S$, $\int_{\mathbb{R}^d} \|x\|^2 d\mu_i(x) < +\infty$ and $H(\mu_i) < +\infty$. Then A2 is satisfied.

Proof. Let $\pi = \otimes_{i \in S} \mu_i \otimes_{i \in S^c} \nu_i$ with ν_i any Gaussian measure with positive definite covariance matrix. First, we have that

$$\text{KL}(\pi | \pi^0) = \text{KL}(\pi_r | \pi_r^0) + \sum_{(v,v') \in E_r} \int_{\mathbb{R}^d} \text{KL}(\pi_{v'|v} | \pi_{v'|v}^0) d\pi_v.$$

For any $(v, v') \in E_r$, there exists $C_{v,v'} \geq 0$ such that

$$\begin{aligned} \int_{\mathbb{R}^d} \text{KL}(\pi_{v'|v} | \pi_{v'|v}^0) d\pi_v &\leq C_{v,v'} - H(\pi_{v'}) + \int_{\mathbb{R}^d \times \mathbb{R}^d} \|x_v - x_{v'}\|^2 / (2\sigma_{v,v'}^2) d\pi_v \otimes \pi_{v'}(x_v, x_{v'}) \\ &\leq C_{v,v'} - H(\pi_{v'}) + (1/\sigma_{v,v'}^2) \int_{\mathbb{R}^d} \|x_v\|^2 d\pi_v(x_v) + (1/\sigma_{v,v'}^2) \int_{\mathbb{R}^d} \|x_{v'}\|^2 d\pi_{v'}(x_{v'}) < +\infty. \end{aligned}$$

We conclude the proof upon remarking that $\text{KL}(\pi_r | \pi_r^0) < +\infty$. \square

Proposition 16. Let $\pi^0 \in \mathcal{P}_{T_r}$. Assume that $\pi_r^0 = \mu_r$ if $r \in S$ or $\pi_r^0 = N(m_r, \sigma_r \text{Id})$, with $m_r \in \mathbb{R}^d$ and $\sigma_r > 0$ if $r \in S^c$. In addition, assume that for any $(v, v') \in E_r$, $\pi_{v'|v}^0(\cdot | x_v) = N(x_v, \sigma_{v,v'} \text{Id})$ with $\sigma_{v,v'} > 0$. Finally, assume that for any $i \in S$, μ_i admits a positive density w.r.t. the Lebesgue measure. Then A3 is satisfied.

Proof. We have that π^0 admits a positive density w.r.t. the Lebesgue measure. Letting $\tilde{\pi}^0 = \otimes_{i \in S} \mu_i \otimes_{j \in S^c} \tilde{\mu}_j$ where $\tilde{\mu}_j$ which admits a positive density w.r.t. the Lebesgue measure for any $j \in S^c$, we get that $\tilde{\pi}^0$ admits a positive density w.r.t. the Lebesgue measure and therefore $\pi^0 \sim \tilde{\pi}^0$, which concludes the proof. \square

Using the preliminary results presented above, we are now ready to prove Proposition 3.

Proof of Proposition 3. Assume A1 and A2. Since \mathcal{P}_S is convex and closed in total-variation norm, there exists a probability distribution π^* solution to (7), or equivalently to (static-mSB), by using A2 with (Csiszár, 1975, Theorem 2.1.). Moreover, this solution is unique by strict convexity of $\text{KL}(\cdot | \pi^0)$.

We now turn to the proof of existence of potentials defining $(d\pi^*/d\pi^0)$, by adapting the arguments of (Nutz, 2021, Section 2.3.). Define $\nu^n = \text{argmin}\{\text{KL}(\pi | \pi^0) : \pi \in \mathcal{P}_S^n\}$ for any $n \in \mathbb{N}^*$. Since $\{\mathcal{P}_S^n\}_{n \in \mathbb{N}^*} \subset \mathcal{P}^{(\ell+1)}$ is a decreasing sequence of sets that are convex and closed in total-variation norm such that (6) holds, we get from (Nutz, 2021, Proposition 1.17.) with A2 that

$$\lim_{n \rightarrow \infty} \|\nu^n - \pi^*\|_{\text{TV}} = 0,$$

or equivalently

$$\lim_{n \rightarrow \infty} \|(d\nu^n/d\pi^0) - (d\pi^*/d\pi^0)\|_{L^1(\pi^0)} = 0. \quad (8)$$

Following (Nutz, 2021, Example 1.18), there exists a family of bounded measurable functions $\{\varphi_i^n\}_{n \in \mathbb{N}^*, i \in S}$ with $\varphi_i^n : \mathbb{R}^d \rightarrow \mathbb{R}$ such that for any $n \in \mathbb{N}^*$

$$(d\nu^n/d\pi^0) = \exp[\bigoplus_{i \in S} \varphi_i^n]. \quad (9)$$

We consider such family $\{\varphi_i^n\}_{n \in \mathbb{N}^*, i \in S}$ for the rest of the proof. By combining (8) and (9), we obtain, up to extraction,

$$(d\pi^*/d\pi^0) = \lim_{n \rightarrow \infty} \exp[\bigoplus_{i \in S} \varphi_i^n] \quad \pi^0\text{-a.s.} \quad (10)$$

We now define the following sets

$$\begin{aligned} A^* &= \{x \in (\mathbb{R}^d)^{\ell+1} : \lim_{n \rightarrow \infty} \bigoplus_{i \in S} \varphi_i^n(x_i) \in [-\infty, +\infty)\}, \\ B^* &= \{x \in (\mathbb{R}^d)^{\ell+1} : \lim_{n \rightarrow \infty} \bigoplus_{i \in S} \varphi_i^n(x_i) > -\infty\} \subset A^* \end{aligned}$$

Using (10), we have $\pi^0(A^*) = 1$. Using **A3**, it comes $\tilde{\pi}^0(A^*) = 1$. Moreover, we also get that $\pi^*(B^*) = 1$ by (10). Thus, it comes $\pi^0(B^*) > 0$, and $\tilde{\pi}^0(B^*) > 0$ using **A3**.

We then apply Lemma 13 to $\tilde{\pi}^0$ and $A = A^*$. Since $\tilde{\pi}^0(B^*) > 0$, it implies that there exists $x^* \in B^*$ and a measurable set $A^0 \subset B^*$ verifying the properties (a) and (b). Following (Nutz, 2021, Corollary 2.12), we may assume without loss of generality in the statement of Lemma 13 that the sets X_m^0 are measurable with $\prod_{m=0}^{\ell} X_m^0 \subset A$. In this case, we obtain that $\mu_i(\text{proj}_i(A^0)) = 1$ for any $i \in S$.

We now aim at applying Lemma 14 to the set A^0 . Remark that A^0 directly satisfies assumption (a). For any $n \in \mathbb{N}^*$, consider the following transformation of the functions $\{\varphi_i^n\}_{i \in S}$

$$\begin{aligned} \varphi_{i_k}^n &\leftarrow \varphi_{i_k}^n - \varphi_{i_k}^n(x_{i_k}^*), \quad \forall k \in \{0, \dots, K-2\}, \\ \varphi_{i_{K-1}}^n &\leftarrow \varphi_{i_{K-1}}^n + \sum_{k=0}^{K-2} \varphi_{i_k}^n(x_{i_k}^*). \end{aligned}$$

For any $i \in S$, we restrict $\varphi_{i_k}^n$ to $X_{i_k}^0$, so that the family $\{\varphi_i^n\}_{n \in \mathbb{N}^*, i \in S}$ now verifies assumption (b). Finally, since $A^0 \subset A^*$ and $x^* \in B^*$, we directly obtain assumption (c).

Therefore, Lemma 14 may be applied. It provides us with the family of functions $\{\varphi_i\}_{i \in S}$ defined by $\varphi_i : X_i^0 \rightarrow [-\infty, +\infty)$ with $\varphi_i = \lim_{n \rightarrow \infty} \varphi_i^n$ μ_i -a.s. for any $i \in S$. Since $\mu_i(\text{proj}_i(A^0)) = 1$ for any $i \in S$, we may extend the functions φ_i to \mathbb{R}^d . In particular, we can find a family of functions $\{\psi_i^*\}_{i \in S}$ with $\psi_i^* : \mathbb{R}^d \rightarrow [-\infty, +\infty)$ such that $\psi_i^* = \varphi_i$ μ_i -a.s. Note that these functions are measurable as limits of measurable functions.

Since $\pi^0 \sim \tilde{\pi}^0$ by **A3**, (10) turns into

$$(d\pi^*/d\pi^0) = \exp\left[\bigoplus_{i \in S} \psi_i^*\right] \quad \pi^0\text{-a.s.} \quad (11)$$

Finally, we show that the functions ψ_i^* are μ_i -a.s. finite. Let $i \in S$. Let us define $A_i = \{x_i \in \mathbb{R}^d : \psi_i^*(x_i) = -\infty\}$. Using (11), we obtain $(d\pi^*/d\pi^0)(A_i \times (\mathbb{R}^d)^\ell) = 0$. Since $\pi_i^* = \mu_i$, we have

$$\mu_i(A_i) = \pi^*(A_i \times (\mathbb{R}^d)^\ell) = \int_{A_i \times (\mathbb{R}^d)^\ell} (d\pi^*/d\pi^0) d\pi^0 = 0,$$

which gives the result. \square

We now turn to the proof of Corollary 4, which states that the iterates of (mIPF) can be expressed via potentials, in the same manner as the solution π^* to (static-mSB).

Proof of Corollary 4. Assume **A1**, **A2** and **A3**. We prove the result of this corollary by recursion on $n \in \mathbb{N}^*$. First take $n = 1$. In this case, the first iteration of (mIPF) is a multi-marginal SB problem of the form (static-mSB) where $S = \{i_0\}$ with reference measure π^0 . Therefore, using **A2** and **A3**, we can apply Proposition 3 and obtain existence of $\psi_{i_0}^1 : \mathbb{R}^d \rightarrow \mathbb{R}$ such that

$$(d\pi^1/d\pi^0) = \exp[\psi_{i_0}^1] \quad \pi^0\text{-a.s.}$$

By taking $\psi_{i_k}^0 = 0$ for $k \in \{1, \dots, K-1\}$, we thus obtain the result at step $n = 1$.

Now assume that the result is verified for some $n \in \mathbb{N}^*$, with $k_n = (n-1) \bmod(K)$. We define $k_n + 1 = n \bmod(K)$ and $q_n \in \mathbb{N}$ as the quotient of the Euclidean division of n by K . In this case, the $(n+1)$ -th iteration of (mIPF) is a multi-marginal SB problem of the form (static-mSB) where $S = \{i_{k_n+1}\}$ with reference measure π^n . Using (13), we have that **A2** is satisfied for this new (static-mSB) problem. **A1** and **A3** are satisfied for this problem, given the form of π^n . Therefore, we can apply Proposition 3 and obtain existence of $\psi_{i_{k_n+1}}^{q_n+1} : \mathbb{R}^d \rightarrow \mathbb{R}$ such that

$$(d\pi^{n+1}/d\pi^n) = \exp[\psi_{i_{k_n+1}}^{q_n+1}] \quad \pi^n\text{-a.s.} \quad (12)$$

By assumption, we have that $\pi^n \ll \pi^0$. Hence, we obtain $\pi^{n+1} \ll \pi^0$ and thus,

$$(d\pi^{n+1}/d\pi^0) = (d\pi^{n+1}/d\pi^n)(d\pi^n/d\pi^0) \quad \pi^0\text{-a.s.}$$

By combining (12) with the result of the recursion at step n , we directly obtain the result at step $n+1$, which achieves the proof. \square

Proofs of Proposition 5 and Proposition 6. In this part of the section, we establish the proofs of results related to the convergence of (mIPF), respectively Proposition 5 and Proposition 6, which can be seen as a natural extension of (Ruschendorf, 1995, Proposition 2.1.) and (Ruschendorf, 1995, Theorem 3.1.).

Proof of Proposition 5. Under A1 and A2, we obtain by Proposition 3 existence and uniqueness of a solution to (static-mSB), which we denote by π^* . Since $\pi^* \in \mathcal{P}_S$, using recursively (Csiszár, 1975, Theorem 3.12.), the fact that $\{\pi_{i_k} = \mu_{i_k} : \pi \in \mathcal{P}^{(|V|)}\}$ is convex for any $k \in \{0, \dots, K-1\}$ and (mIPF), we obtain

$$\text{KL}(\pi^* | \pi^0) = \text{KL}(\pi^* | \pi^n) + \sum_{i=1}^n \text{KL}(\pi^i | \pi^{i-1}). \quad (13)$$

Therefore, we have $\sum_{i=1}^\infty \text{KL}(\pi^i | \pi^{i-1}) \leq \text{KL}(\pi^* | \pi^0) < \infty$ and thus,

$$\lim_{i \rightarrow +\infty} \text{KL}(\pi^i | \pi^{i-1}) = 0. \quad (14)$$

Let $n \in \mathbb{N}^*$ with $n > 2K$, $k \in \{0, \dots, K-1\}$ and let $q_n \in \mathbb{N}$ be the quotient of the Euclidean division of $n-1$ by K . We define $n_k = q_n K + k + 1$ with $(n_k - 1) = k \bmod(K)$ if $n_k \leq n$. Otherwise, we set $n_k = (q_n - 1)K + k + 1$ with $(n_k - 1) = k \bmod(K)$. Note that we always have $|n - n_k| \leq 2K$. In particular, we have $\pi_{i_k}^{n_k} = \mu_{i_k}$ by definition of (mIPF). Therefore, we obtain

$$\begin{aligned} \|\pi_{i_k}^n - \mu_{i_k}\|_{\text{TV}} &\leq \|\pi^n - \pi^{n_k}\|_{\text{TV}} \\ &\leq \|\pi^n - \pi^{n-1}\|_{\text{TV}} + \dots + \|\pi^{n_k+1} - \pi^{n_k}\|_{\text{TV}} \quad (\text{triangle inequality}) \\ &\leq (2\text{KL}(\pi^n | \pi^{n-1}))^{1/2} + \dots + (2\text{KL}(\pi^{n_k+1} | \pi^{n_k}))^{1/2}, \quad (\text{ Pinsker's inequality}) \end{aligned}$$

where each term goes to 0 as $n \rightarrow +\infty$ in the last inequality by (14), which achieves the proof. \square

We now turn to the proof Proposition 6, which requires several preliminary technical results. For the rest of this section, we define, for any $n \in \mathbb{N}$, q_n as the quotient of the Euclidean division of $n-1$ by K (in particular, $q_0 = -1$).

Schrödinger equations. Under A1, A2 and A3, we know from Proposition 3 that the unique solution π^* to (static-mSB) can be π^0 -a.s. written as $(d\pi^*/d\pi^0) = \exp[\bigoplus_{i \in S} \psi_i^*]$, where $\{\psi_i^*\}_{i \in S}$ are measurable potentials, referred to as *Schrödinger potentials*. These functions are determined by the fixed-point *Schrödinger equations*

$$\psi_i^*(x_i) = \log[r_i(x_i) / \int_{(\mathbb{R}^d)^\ell} \exp[\sum_{j \in S \setminus \{i\}} \psi_j^*(x_j)] h(x_{0:\ell}) d\nu_{-i}(x_{-i})] \quad \mu_i\text{-a.s.}, \quad \forall i \in S,$$

which are obtained by marginalising π^* along its constrained marginals. This family of potentials is not unique. Indeed, for any family of real numbers $\{\lambda_{i_k}\}_{k \in \{0, \dots, K-2\}}$, we have

$$(d\pi^*/d\pi^0) = \exp[\bigoplus_{i \in S} \tilde{\psi}_i],$$

where $\tilde{\psi}_{i_k} = \psi_{i_k}^* + \tilde{\lambda}_{i_k}$ for any $k \in \{0, \dots, K-1\}$ with $\tilde{\lambda}_{i_k} = \lambda_{i_k}$ if $k \in \{0, \dots, K-2\}$ and $\tilde{\lambda}_{i_{K-1}} = -\sum_{i=0}^{K-2} \lambda_{i_k}$.

Remark on the initialisation of (mIPF). Consider a probability measure $\bar{\pi}^0 \in \mathcal{P}^{(\ell+1)}$ of the form

$$(d\bar{\pi}^0/d\pi^0) = \exp[\bigoplus_{i \in S} \psi_i^0], \quad (15)$$

where $\{\psi_i^0\}_{i \in S}$ is a family of measurable potentials with $\psi_i^0 : \mathbb{R}^d \rightarrow \mathbb{R}$ such that $|\int_{\mathbb{R}^d} \psi_i^0 d\mu_i| < \infty$ for any $i \in S$. Then, for any $\pi \in \mathcal{P}_S$, we have

$$\text{KL}(\pi | \pi^0) = \text{KL}(\pi | \bar{\pi}^0) + \int_{(\mathbb{R}^d)^K} \bigoplus_{i \in S} \psi_i^0 d\pi = \text{KL}(\pi | \bar{\pi}^0) + \sum_{i \in S} \int_{\mathbb{R}^d} \psi_i^0 d\mu_i.$$

Hence, (static-mSB) is equivalent to the multi-marginal SB problem

$$\text{argmin}\{\text{KL}(\pi | \bar{\pi}^0) : \pi \in \mathcal{P}^{(\ell+1)}, \pi_i = \mu_i, \forall i \in S\}.$$

We refer to (Peyré et al., 2019, Proposition 4.2) for the EOT counterpart of this result. This means that the solutions of the multi-marginal Schrödinger Bridge problem are invariant by multiplication of the reference measure by potentials on the *fixed* marginals. Consequently, the initialisation of the (mIPF) sequence may be chosen as $\bar{\pi}^0$ instead of π^0 .

For sake of clarity, we now refer to the reference probability measure of (static-mSB) as $\bar{\pi}$ or π^{-1} and to the initialisation of the (mIPF) iterates as π^0 .

Solving (mIPF) with potentials. To prove the convergence of the (mIPF) iterates to the solution π^* given by Proposition 3, we first rewrite these iterates with potentials, following the form of π^* .

To do so, we recursively define the sequence of potentials $\{\psi_i^n\}_{n \in \mathbb{N}, i \in S}$ by

$$\begin{aligned} \psi_{i_0}^0 &= \dots = \psi_{i_{K-2}}^0 = 0, \\ \psi_{i_{K-1}}^0(x_{i_{K-1}}) &= \log(r_{i_{K-1}}(x_{i_{K-1}}) / \int_{(\mathbb{R}^d)^\ell} h(x_{0:\ell}) d\nu_{-i_{K-1}}(x_{-i_{K-1}})), \end{aligned} \tag{16}$$

and for any $n \in \mathbb{N}^*$ and $k \in \{0, \dots, K-1\}$

$$\begin{aligned} \psi_{i_k}^{q_n+1}(x_{i_k}) &= \log[r_{i_k}(x_{i_k}) / \int_{(\mathbb{R}^d)^\ell} \exp[\bigoplus_{\ell=0}^k \psi_{i_\ell}^{q_n+1}(x_{i_\ell}) \bigoplus_{m=k+1}^{K-1} \psi_{i_m}^{q_n}(x_{i_m})] \\ &\quad \times h(x_{0:\ell}) d\nu_{-i_k}(x_{-i_k})], \end{aligned} \tag{17}$$

recalling that q_n is the quotient of the Euclidean division of $n-1$ by K .

We now define the sequence of probability measures $\{\pi^n\}_{n \in \mathbb{N}}$ by

$$d\pi^n / d\bar{\pi} = \exp[\bigoplus_{\ell=0}^{k_n} \psi_{i_\ell}^{q_n+1} \bigoplus_{m=k_n+1}^{K-1} \psi_{i_m}^{q_n}], \quad k_n = (n-1) \bmod(K), \quad n = q_n K + k_n + 1. \tag{18}$$

In particular, we have $(d\pi^0 / d\bar{\pi}) = \exp[\bigoplus_{\ell=0}^{K-1} \psi_{i_\ell}^0] = \exp[\psi_{i_{K-1}}^0]$, and thus $\int_{\mathbb{R}^d} \psi_{i_{K-1}}^0 d\mu_{i_{K-1}} = \text{KL}(\mu_{i_{K-1}} | \bar{\pi}_{i_{K-1}})$. Consequently, π^0 can be chosen as the initialisation of (mIPF), following the previous remark, if we assume that $\text{KL}(\mu_{i_{K-1}} | \bar{\pi}_{i_{K-1}}) < \infty$. In (TreeSB) with $r = i_{K-1}$, the latter assumption is directly verified since we choose $\bar{\pi}_{i_{K-1}} = \mu_{i_{K-1}}$.

Let $n \in \mathbb{N}$, with $k_n = (n-1) \bmod(K)$, $k_n + 1 = n \bmod(K)$. Using (16) and (17), we get that $\pi_{i_{k_n}}^n = \mu_{i_{k_n}}$. Moreover, we have

$$d\pi^n / d\pi^{n-1} = \exp[\psi_{i_{k_n}}^{q_n+1} - \psi_{i_{k_n}}^{q_n}], \tag{19}$$

with the convention that $\psi_{i_{K-1}}^{-1} = 0$. In particular, we obtain that $\pi_{i_{k_n+1}}^{n+1} = \pi_{i_{k_n+1}}^n$.

In conclusion, the sequence $\{\pi^n\}_{n \in \mathbb{N}}$ defined in (18) verifies $\pi^{n+1} = \mu_{i_{k_n+1}} \pi_{i_{k_n+1}}^n$ for any $n \in \mathbb{N}$. By decomposition property of the Kullback-Leibler divergence, this sequence solves (mIPF) with initialisation π^0 . We consider such iterates in the following.

Since $\pi_{i_{k_n}}^n = \mu_{i_{k_n}}$, we have that

$$\text{KL}(\pi^n | \pi^{n-1}) = \int_{\mathbb{R}^d} (\psi_{i_{k_n}}^{q_n+1} - \psi_{i_{k_n}}^{q_n}) d\mu_{i_{k_n}}. \tag{20}$$

Before proving a multi-marginal counterpart to (Ruschendorf, 1995, Lemma 4.1), we state and prove the following result.

Proposition 17. Let π_0, π_1 two probability measures on \mathbb{R}^d such that $\pi_0 \ll \pi_1$. Then, denoting $f = d\pi_0 / d\pi_1$, the following assertions are equivalent:

- (a) $\text{KL}(\pi_0 | \pi_1) < +\infty$
- (b) $\int_{\mathbb{R}^d} |\log(f)(x)| d\pi_0(x) < +\infty$
- (c) $\int_{\mathbb{R}^d} \log(f)(x) \mathbb{1}_{f(x) > 1} d\pi_0(x) < +\infty$

If one of these conditions is satisfied then $\int_{\mathbb{R}^d} |\log(f)(x)| d\pi_0 \leq \text{KL}(\pi_0 | \pi_1) + 2/e$.

Proof. First, note that

$$\int_{\mathbb{R}^d} |\log(f)(x)| \mathbb{1}_{f < 1} d\pi_0(x) \leq \int_{\mathbb{R}^d} |\log(f)(x) f(x)| \mathbb{1}_{f < 1} d\pi_1(x) \leq 1/e, \tag{21}$$

where we have used that for any $u \in [0, 1]$, $|u \log(u)| \leq 1/e$. We have that (b) implies (c). Using the previous result we have that (c) implies (b). Hence (c) and (b) are equivalent. In addition, it is clear that (b) implies (a). Finally (this is more of a convention), we have that $\text{KL}(\pi_0 | \pi_1) = \int_{\mathbb{R}^d} \log(f)(x) \mathbb{1}_{f(x) > 1} d\pi_0(x) + \int_{\mathbb{R}^d} \log(f)(x) \mathbb{1}_{f(x) < 1} d\pi_0(x) < +\infty$. Using (21) this implies (c). Finally, we have

$$\begin{aligned} \int_{\mathbb{R}^d} |\log(f)(x)| d\pi_0(x) &= \int_{\mathbb{R}^d} \log(f)(x) d\pi_0(x) - 2 \int_{\mathbb{R}^d} \log(f)(x) \mathbb{1}_{f(x) < 1} d\pi_0(x) \\ &\leq \text{KL}(\pi_0 | \pi_1) + 2/e, \end{aligned}$$

which concludes the proof. \square

We begin with the following lemma which controls the integral of the potentials uniformly w.r.t. $n \in \mathbb{N}$. It can be seen as the *multi-marginal* counterpart of (Ruschendorf, 1995, Lemma 4.1).

Lemma 18. *Assume A4. There exist $\{c_i\}_{i \in S} \in (0, +\infty)^K$ such that for any function $f : (\mathbb{R}^d)^{\ell+1} \rightarrow \mathbb{R}$ of the form $f = \bigoplus_{i \in S} f_i$, we have*

$$c_i \|f\|_{L^1(\pi^*)} \geq \|f_i\|_{L^1(\mu_i)}, \quad \forall i \in S. \quad (22)$$

For any $n \in \mathbb{N}^*$, we have

- (a) $\sum_{i \in S} \int_{\mathbb{R}^d} \psi_i^n d\mu_i \leq \text{KL}(\pi^* | \bar{\pi}) < \infty$,
- (b) $\int_{(\mathbb{R}^d)^{\ell+1}} (\bigoplus_{i \in S} \psi_i^* - \bigoplus_{i \in S} \psi_i^n) d\pi^* \leq \text{KL}(\pi^* | \bar{\pi}) < \infty$,
- (c) $\sup_{n \in \mathbb{N}} \int_{\mathbb{R}^d} |\psi_i^n| d\mu_i < \infty, \forall i \in S$.

Proof. First, we have that (22) is a direct consequence of (Kober, 1940, Theorem 1) and A4. Let us now prove (a). Using (20), we have

$$\begin{aligned} \sum_{m=0}^{K^n} \text{KL}(\pi^m | \pi^{m-1}) &= \sum_{\ell=0}^{n-1} \sum_{k=0}^{K-1} \text{KL}(\pi^{\ell K+k+1} | \pi^{\ell K+k}) + \text{KL}(\pi^0 | \pi^{-1}) \\ &= \sum_{\ell=0}^{n-1} \sum_{i \in S} \int_{\mathbb{R}^d} (\psi_i^{\ell+1} - \psi_i^\ell) d\mu_i + \int_{\mathbb{R}^d} (\psi_{i_{K-1}}^0 - \psi_{i_{K-1}}^{-1}) d\mu_{i_{K-1}} \\ &= \sum_{i \in S} \sum_{\ell=0}^{n-1} \int_{\mathbb{R}^d} (\psi_i^{\ell+1} - \psi_i^\ell) d\mu_i + \int_{\mathbb{R}^d} (\psi_{i_{K-1}}^0 - \psi_{i_{K-1}}^{-1}) d\mu_{i_{K-1}} \\ &= \sum_{i \in S} \int_{\mathbb{R}^d} (\psi_i^n - \psi_i^0) d\mu_i + \int_{\mathbb{R}^d} (\psi_{i_{K-1}}^0 - \psi_{i_{K-1}}^{-1}) d\mu_{i_{K-1}} \\ &= \sum_{i \in S} \int_{\mathbb{R}^d} \psi_i^n d\mu_i \leq \text{KL}(\pi^* | \bar{\pi}). \end{aligned}$$

where the last inequality follows the proof of Proposition 5.

Since the first term in the inequality of (b) is equal to $\text{KL}(\pi^* | \pi^{nK})$, we obtain (b) using that $\text{KL}(\pi^* | \pi^{nK}) \leq \text{KL}(\pi^* | \bar{\pi})$ following the proof of Proposition 5.

Let us now prove (c). Since $\text{KL}(\pi^* | \bar{\pi}) < \infty$, using Proposition 17, we have that $\bigoplus_{i \in S} \psi_i^* \in L^1(\pi^*)$. From (b) and Proposition 17, we also get that $\bigoplus_{i \in S} (\psi_i^* - \psi_i^n) \in L^1(\pi^*)$, and thus $\int_{(\mathbb{R}^d)^{\ell+1}} |\bigoplus_{i \in S} (\psi_i^* - \psi_i^n)| d\pi^* \leq C_0$ with $C_0 > 0$. Therefore, we have

$$\int_{(\mathbb{R}^d)^{\ell+1}} |\bigoplus_{i \in S} \psi_i^n| d\pi^* \leq \int_{(\mathbb{R}^d)^{\ell+1}} |\bigoplus_{i \in S} \psi_i^*| d\pi^* + \int_{(\mathbb{R}^d)^{\ell+1}} |\bigoplus_{i \in S} (\psi_i^* - \psi_i^n)| d\pi^* \leq 2C_0.$$

Using (22), we conclude with A4 that for any $i \in S$, we have

$$\int_{\mathbb{R}^d} |\psi_i^n| d\mu_i \leq 2c_i C_0,$$

which concludes the proof of (c). □

The next lemma gives an explicit expression for $\text{KL}(\pi^n | \bar{\pi})$. It can be seen as the *multi-marginal* counterpart of (Ruschendorf, 1995, Lemma 4.2).

Lemma 19. *For any $n \in \mathbb{N}$, with $k_n = (n - 1) \bmod(K)$, we have*

$$\begin{aligned} \text{KL}(\pi^n | \bar{\pi}) &= \int_{\mathbb{R}^d} \psi_{i_{k_n}}^{q_n+1} d\mu_{i_{k_n}} + \sum_{\ell=0}^{k_n-1} \int_{\mathbb{R}^d} \psi_{i_\ell}^{q_n+1} \exp[\psi_{i_\ell}^{q_n+1} - \psi_{i_\ell}^{q_n+2}] d\mu_{i_\ell} \\ &\quad + \sum_{m=k_n+1}^{K-1} \int_{\mathbb{R}^d} \psi_{i_m}^{q_n} \exp[\psi_{i_m}^{q_n} - \psi_{i_m}^{q_n+1}] d\mu_{i_m}. \end{aligned}$$

Proof. Let $n \in \mathbb{N}$, with $k_n = (n - 1) \bmod(K)$. Using (18), we have

$$\text{KL}(\pi^n | \bar{\pi}) = \int_{\mathbb{R}^d} \psi_{i_{k_n}}^{q_n+1} d\mu_{i_{k_n}} + \sum_{\ell=0}^{k_n-1} \int_{\mathbb{R}^d} \psi_{i_\ell}^{q_n+1} d\pi_{i_\ell}^n + \sum_{m=k_n+1}^{K-1} \int_{\mathbb{R}^d} \psi_{i_m}^{q_n} d\pi_{i_m}^n. \quad (23)$$

Consider $m \in \{k_n + 1, \dots, K - 1\}$. Let m_n be the closest integer to n such that $m_n > n$ and $m = (m_n - 1) \bmod(K)$. By (19), we have

$$d\pi^n = \exp[\bigoplus_{j=k_n+1}^m \psi_{i_j}^{q_n} - \psi_{i_j}^{q_n+1}] d\pi^{m_n}.$$

Using (19) recursively, we obtain

$$d\pi_{i_m}^n = \exp[\psi_{i_m}^{q_n} - \psi_{i_m}^{q_n+1}] d\pi_{i_m}^{m_n}, \quad (24)$$

where we recall that $\pi_{i_m}^{m_n} = \mu_{i_m}$.

Consider now $\ell \in \{0, \dots, k_n - 1\}$. Let ℓ_n be the closest integer to n such that $\ell_n > n$ and $\ell = (\ell_n - 1) \bmod(K)$. By (19), we have

$$d\pi^n = \exp\left[\bigoplus_{j=k_n+1}^{K-1} \{\psi_{i_j}^{q_n} - \psi_{i_j}^{q_n+1}\} \bigoplus_{j'=0}^{\ell} \{\psi_{i_{j'}}^{q_n+1} - \psi_{i_{j'}}^{q_n+2}\}\right] d\pi^{\ell_n},$$

and using (19) recursively, we obtain

$$d\pi_{i_\ell}^n = \exp[\psi_{i_\ell}^{q_n+1} - \psi_{i_\ell}^{q_n+2}] d\pi_{i_\ell}^{\ell_n}, \quad (25)$$

where we recall that $\pi_{i_\ell}^{\ell_n} = \mu_{i_\ell}$. We conclude the proof upon combining (23), (24) and (25). \square

We are now ready to prove a *uniform integrability* result which is the multi-marginal counterpart of (Ruschendorf, 1995, Lemma 4.4). Before stating Lemma 21, we prove the following well-known lemma. We recall that a sequence $(\Psi_n)_{n \in \mathbb{N}}$ such that for any $n \in \mathbb{N}$, $\Psi_n \in L^1(\mu)$, is *uniformly integrable* w.r.t. μ if (i) $\sup_{n \in \mathbb{N}} \int_{\mathbb{R}^d} |\Psi_n| d\mu < +\infty$ and (ii) for any $\varepsilon > 0$, there exists $K > 0$ such that for any $n \in \mathbb{N}$, $\int_{\mathbb{B}(0,K)^c} |\Psi_n| d\mu \leq \varepsilon$.

Lemma 20. *Let $f : \mathbb{R} \rightarrow \mathbb{R}$, convex and non-decreasing on $[A, +\infty)$ with $A > 0$ and $\lim_{x \rightarrow +\infty} f(x)/x = +\infty$. Assume that $\sup_{n \in \mathbb{N}} \int_{\mathbb{R}^d} f(|\Psi_n|) d\mu < +\infty$. Then, $(\Psi_n)_{n \in \mathbb{N}}$ is uniformly integrable w.r.t. μ .*

Proof. Since f is convex, using Jensen's inequality, we get that $\sup_{n \in \mathbb{N}} f(\int_{\mathbb{R}^d} |\Psi_n| d\mu) < +\infty$ and since $\lim_{x \rightarrow +\infty} f(x)/x = +\infty$ we have $\sup_{n \in \mathbb{N}} \int_{\mathbb{R}^d} |\Psi_n| d\mu < +\infty$. Let $\varepsilon > 0$, there exists $K > 0$ such that for any $x > K$, $x \leq \varepsilon f(x)/B$ with $B = \sup_{n \in \mathbb{N}} \int_{\mathbb{R}^d} f(|\Psi_n|) d\mu < +\infty$. Therefore, we have for any $n \in \mathbb{N}$

$$\int_{\mathbb{B}(0,K)^c} |\Psi_n| d\mu \leq (\varepsilon/B) \int_{\mathbb{B}(0,K)^c} f(|\Psi_n|) d\mu \leq \varepsilon,$$

which concludes the proof. \square

Lemma 21. *Assume A4 and A5. Then, $\{\exp[\bigoplus_{i \in S} \psi_i^n]\}_{n \in \mathbb{N}}$ is uniformly integrable w.r.t. $\bar{\pi}$.*

Proof. It is enough to show that the sequence $\{f(\exp[\bigoplus_{i \in S} \psi_i^n])\}_{n \in \mathbb{N}}$ is bounded in $L^1(\bar{\pi})$, where $f : u \mapsto u \log(u)$ is continuous, convex and such that $\lim_{u \rightarrow \infty} f(u)/u = +\infty$, see Lemma 20. Let $n \in \mathbb{N}$. We have

$$\begin{aligned} \int_{(\mathbb{R}^d)^{\ell+1}} f(\exp[\bigoplus_{i \in S} \psi_i^n]) d\bar{\pi} &= \text{KL}(\pi^{nK} \mid \bar{\pi}) \\ &= \int_{\mathbb{R}^d} \psi_{i_{K-1}}^n d\mu_{i_{K-1}} + \sum_{k=0}^{K-2} \int_{\mathbb{R}^d} \psi_{i_k}^n \exp[\psi_{i_k}^n - \psi_{i_k}^{n+1}] d\mu_{i_k} && \text{(Lemma 19)} \\ &= \sum_{k=0}^{K-1} \int_{\mathbb{R}^d} \psi_{i_k}^n d\mu_{i_k} + \sum_{k=0}^{K-2} \int_{\mathbb{R}^d} \psi_{i_k}^n \{\exp[\psi_{i_k}^n - \psi_{i_k}^{n+1}] - 1\} d\mu_{i_k} \\ &\leq \text{KL}(\pi^* \mid \bar{\pi}) + (\bar{c} + 1) \sum_{k=0}^{K-2} \int_{\mathbb{R}^d} \psi_{i_k}^n d\mu_{i_k} && \text{(Lemma 18-(a), A5)} \\ &\leq \text{KL}(\pi^* \mid \bar{\pi}) + (\bar{c} + 1) \sum_{k=0}^{K-2} \sup_{n \in \mathbb{N}} \int_{\mathbb{R}^d} |\psi_{i_k}^n| d\mu_{i_k} < \infty. && \text{(Lemma 18-(c))} \end{aligned}$$

\square

With the preliminary results stated above, we are now ready to prove Proposition 6.

Proof of Proposition 6. Using A4 and A5, we have, by Lemma 21, uniform integrability of $\{\exp[\bigoplus_{i \in S} \psi_i^n]\}_{n \in \mathbb{N}}$ in $L^1(\bar{\pi})$. Therefore, the sequence $\{\pi^{nK}\}_{n \in \mathbb{N}}$ is relatively compact with respect to the weak topology of $\sigma(L^1(\bar{\pi}), L^\infty(\bar{\pi}))$, denoted as the τ -topology. We recall that $\lim_{n \rightarrow \infty} \text{KL}(\pi^{nK+1} \mid \pi^{nK}) = 0$. This implies that $\{\pi^{nK+1}\}_{n \in \mathbb{N}}$ is also relatively τ -compact. By trivial recursion, we obtain that the sequences $\{\pi^{nK+k}\}_{n \in \mathbb{N}}$, where $k \in \{2, \dots, K-1\}$ are also relatively τ -compact. Therefore, $\{\pi^n\}_{n \in \mathbb{N}}$ is relatively τ -compact and τ -sequentially compact.

We consider an increasing function $\Phi : \mathbb{N} \rightarrow \mathbb{N}$ such that $\{\pi^m\}_{m \in \Phi(\mathbb{N})}$ is a τ -convergent subsequence, and we denote by $\bar{\pi}$ its limit for this topology. In particular, $\bar{\pi} \in \mathcal{P}_S$ by Proposition 5. We assume without loss of generality that $\Phi(\mathbb{N}) \subset K\mathbb{N}$.

Using the lower semi-continuity of the Kullback-Leibler divergence (Dupuis & Ellis, 2011, Lemma 1.4.3), we get

$$\text{KL}(\tilde{\pi} \mid \bar{\pi}) \leq \liminf \text{KL}(\pi^m \mid \bar{\pi}) \leq \limsup \text{KL}(\pi^m \mid \bar{\pi}).$$

Consider $k \in \{0, \dots, K-2\}$. By (19), we have

$$\frac{d\mu_{i_k}^{nK+k}}{d\pi_{i_k}^{nK+k}} = \frac{d\pi_{i_k}^{nK+k+1}}{d\pi_{i_k}^{nK+k}} = \frac{d\pi_{i_k}^{nK+k+1}}{d\pi_{i_k}^{nK+k}} = \exp[\psi_{i_k}^{n+1} - \psi_{i_k}^n],$$

and thus,

$$\|\mu_{i_k} - \pi_{i_k}^{nK+k}\|_{\text{TV}} = (1/2) \int_{\mathbb{R}^d} |d\pi_{i_k}^{nK+k} / d\mu_{i_k} - 1| d\mu_{i_k} = (1/2) \int_{\mathbb{R}^d} |\exp[\psi_{i_k}^n - \psi_{i_k}^{n+1}] - 1| d\mu_{i_k}.$$

With Proposition 5, we obtain that $\{\exp[\psi_{i_k}^n - \psi_{i_k}^{n+1}]\}_{n \in \mathbb{N}}$ converges to 1 in $L^1(\mu_{i_k})$. In addition using the uniform integrability of $\{\psi_{i_k}^n\}_{n \in \mathbb{N}}$ and A5, we get

$$\limsup_{n \rightarrow +\infty} \int_{\mathbb{R}^d} \psi_{i_k}^n \exp[\psi_{i_k}^n - \psi_{i_k}^{n+1}] d\mu_{i_k} = \limsup_{n \rightarrow +\infty} \int_{\mathbb{R}^d} \psi_{i_k}^n d\mu_{i_k}.$$

We denote $m = K\ell$. Since $\text{KL}(\pi^m \mid \bar{\pi}) = \int_{\mathbb{R}^d} \psi_{i_{K-1}}^\ell d\mu_{i_{K-1}} + \sum_{k=0}^{K-2} \int_{\mathbb{R}^d} \psi_{i_k}^\ell \exp[\psi_{i_k}^\ell - \psi_{i_k}^{\ell+1}] d\mu_{i_k}$ by Lemma 19, we finally have

$$\text{KL}(\tilde{\pi} \mid \bar{\pi}) \leq \limsup \left\{ \sum_{k=0}^{K-1} \int_{\mathbb{R}^d} \psi_{i_k}^\ell d\mu_{i_k} \right\} \leq \text{KL}(\pi^* \mid \bar{\pi})$$

where the last inequality comes from Lemma 18.

Since $\tilde{\pi}_i = \mu_i$ for any $i \in \mathcal{S}$, using Proposition 5, we have $\tilde{\pi} = \pi^*$ by uniqueness of π^* . Hence, π^* is the only limit point of $\{\pi^n\}_{n \in \mathbb{N}}$ in the τ -topology. In particular, $\text{KL}(\pi^n \mid \bar{\pi}) \rightarrow \text{KL}(\pi^* \mid \bar{\pi})$. Since $\mathcal{P}_{\mathcal{S}}$ is convex, this last result implies $\|\pi^* - \pi^n\|_{\text{TV}} \rightarrow 0$, see the proof of Theorem 2.1 in Csiszár (1975). \square

We finish this section by highlighting that A5 is stronger than (Ruschendorf, 1995, B1). A natural extension of the latter assumption would consist of having a guarantee on the $(K-1)$ first potentials given by (17), as presented below.

A6. There exist $0 < \underline{c} \leq \bar{c}$ such that for any $k \in \{0, \dots, K-2\}$, we have $\underline{c} \leq \exp(-\psi_{i_k}^1) \leq \bar{c}$.

Under A6, (Ruschendorf, 1995, Lemma 4.3) can be adapted as written below.

Lemma 22. Assume A6. Then, for any $n \in \mathbb{N}^*$

(a) for any $k \in \{0, \dots, K-2\}$, there exists $\alpha_{n,k} \in \mathbb{N}$ such that

$$\underline{c} \cdot (\underline{c}/\bar{c})^{\alpha_{n,k}(K-2)} \leq \exp[\psi_{i_k}^{n-1} - \psi_{i_k}^n] \leq \bar{c} \cdot (\bar{c}/\underline{c})^{\alpha_{n,k}(K-2)}$$

(b) there exists $\alpha_{n,K-1} \in \mathbb{N}$ such that

$$1/\bar{c}^{K-1} \cdot (\underline{c}/\bar{c})^{\alpha_{n,K-1}(K-2)} \leq \exp[\psi_{i_{K-1}}^{n-1} - \psi_{i_{K-1}}^n] \leq 1/\underline{c}^{K-1} \cdot (\bar{c}/\underline{c})^{\alpha_{n,K-1}(K-2)}$$

where $\{\alpha_{n,k}\}_{n \in \mathbb{N}^*, k \in \{0, \dots, K-1\}}$ is a strictly increasing sequence that can be explicitly defined.

Proof. We prove the result by recursion on $n \in \mathbb{N}^*$.

Take $n = 1$. Let $k \in \{0, \dots, K-2\}$. We define $\alpha_{1,k} = 0$ and directly obtain (a) by A5 since $\psi_{i_k}^0 = 0$. Let us prove (b). We have by (17)

$$\begin{aligned} \exp[\psi_{i_{K-1}}^0 - \psi_{i_{K-1}}^1] &= \frac{\int_{(\mathbb{R}^d)^\ell} \exp[\bigoplus_{k=0}^{K-2} \psi_{i_k}^1] h d\nu_{-i_{K-1}}}{\int_{(\mathbb{R}^d)^\ell} \exp[\bigoplus_{k=0}^{K-2} \psi_{i_k}^0] h d\nu_{-i_{K-1}}} \\ &= \frac{\int_{(\mathbb{R}^d)^\ell} \exp[\bigoplus_{k=0}^{K-2} \{\psi_{i_k}^1 - \psi_{i_k}^0\} + \bigoplus_{k=0}^{K-2} \psi_{i_k}^0] h d\nu_{-i_{K-1}}}{\int_{(\mathbb{R}^d)^\ell} \exp[\bigoplus_{k=0}^{K-2} \psi_{i_k}^0] h d\nu_{-i_{K-1}}}. \end{aligned}$$

Using (a) at rank $n = 1$, we have

$$1/\bar{c}^{K-1} \leq \exp[\bigoplus_{k=0}^{K-2} \{\psi_{i_k}^1 - \psi_{i_k}^0\}] \leq 1/\underline{c}^{K-1},$$

and therefore, we obtain (b) by taking $\alpha_{1,K-1} = 0$. Let us assume that the result is verified for some $n \in \mathbb{N}^*$. We have

$$\begin{aligned} \exp[\psi_{i_0}^n - \psi_{i_0}^{n+1}] &= \frac{\int \exp[\bigoplus_{k=1}^{K-1} \psi_{i_k}^n] h d\nu_{-i_0}}{\int \exp[\bigoplus_{k=1}^{K-1} \psi_{i_k}^{n-1}] h d\nu_{-i_0}} \\ &= \frac{\int \exp[\bigoplus_{k=1}^{K-2} \{\psi_{i_k}^n - \psi_{i_k}^{n-1}\} \oplus \{\psi_{i_{K-1}}^n - \psi_{i_{K-1}}^{n-1}\} + \bigoplus_{k=1}^{K-1} \psi_{i_k}^{n-1}] h d\nu_{-i_0}}{\int \exp[\bigoplus_{k=1}^{K-1} \psi_{i_k}^{n-1}] h d\nu_{-i_0}} \end{aligned}$$

Using (a) and (b) at rank n , we have

$$\begin{aligned} 1/\bar{c}^{K-2} \cdot (\underline{c}/\bar{c})^{(K-2) \sum_{k=1}^{K-2} \alpha_{n,k}} &\leq \exp[\bigoplus_{k=1}^{K-2} \{\psi_{i_k}^n - \psi_{i_k}^{n-1}\}] \\ &\leq 1/\underline{c}^{K-2} \cdot (\bar{c}/\underline{c})^{(K-2) \sum_{k=1}^{K-2} \alpha_{n,k}}, \\ \underline{c}^{K-1} \cdot (\underline{c}/\bar{c})^{\alpha_{n,K-1}(K-2)} &\leq \exp[\psi_{i_{K-1}}^n - \psi_{i_{K-1}}^{n-1}] \leq \bar{c}^{K-1} \cdot (\bar{c}/\underline{c})^{\alpha_{n,K-1}(K-2)}. \end{aligned}$$

Therefore, we obtain

$$\begin{aligned} \underline{c} \cdot (\underline{c}/\bar{c})^{(K-2) \sum_{k=1}^{K-1} \alpha_{n,k}} &\leq \exp[\bigoplus_{k=1}^{K-2} \{\psi_{i_k}^n - \psi_{i_k}^{n-1}\} \oplus \{\psi_{i_{K-1}}^n - \psi_{i_{K-1}}^{n-1}\}] \\ &\leq \bar{c} \cdot (\bar{c}/\underline{c})^{(K-2) \sum_{k=1}^{K-1} \alpha_{n,k}}, \\ \underline{c} \cdot (\underline{c}/\bar{c})^{(K-2) \sum_{k=1}^{K-1} \alpha_{n,k}} &\leq \exp[\psi_{i_0}^n - \psi_{i_0}^{n+1}] \leq \bar{c} \cdot (\bar{c}/\underline{c})^{(K-2) \sum_{k=1}^{K-1} \alpha_{n,k}}. \end{aligned}$$

Now, we define $\alpha_{n+1,0} = \sum_{k=1}^{K-1} \alpha_{n,k}$ to obtain (a) for $k = 0$. Consider now $k \in \{1, \dots, K-2\}$. Following the same steps as above, we recursively define

$$\alpha_{n+1,k} = \sum_{j=0}^{k-1} \alpha_{n+1,j} + \sum_{j'=k+1}^{K-1} \alpha_{n,j'},$$

which gives (a) at rank $n+1$. Let us now prove (b) at rank $n+1$. We have

$$\begin{aligned} \exp[\psi_{i_{K-1}}^n - \psi_{i_{K-1}}^{n+1}] &= \frac{\int \exp[\bigoplus_{k=0}^{K-2} \psi_{i_k}^{n+1}] h d\nu_{-i_{K-1}}}{\int \exp[\bigoplus_{k=0}^{K-2} \psi_{i_k}^n] h d\nu_{-i_{K-1}}} \\ &= \frac{\int \exp[\bigoplus_{k=0}^{K-2} \{\psi_{i_k}^{n+1} - \psi_{i_k}^n\} + \bigoplus_{k=0}^{K-2} \psi_{i_k}^n] h d\nu_{-i_{K-1}}}{\int \exp[\bigoplus_{k=0}^{K-2} \psi_{i_k}^n] h d\nu_{-i_{K-1}}}. \end{aligned}$$

Using (a) at rank $n+1$, we obtain

$$\begin{aligned} 1/\bar{c}^{K-1} \cdot (\underline{c}/\bar{c})^{(K-2) \sum_{k=0}^{K-2} \alpha_{n+1,k}} &\leq \exp[\bigoplus_{k=0}^{K-2} \{\psi_{i_k}^{n+1} - \psi_{i_k}^n\}] \\ &\leq 1/\underline{c}^{K-1} \cdot (\bar{c}/\underline{c})^{(K-2) \sum_{k=0}^{K-2} \alpha_{n+1,k}}. \end{aligned}$$

Therefore, by taking $\alpha_{n+1,K-1} = \sum_{k=0}^{K-2} \alpha_{n+1,k}$, we obtain (b), which concludes the proof. \square

Unfortunately, Lemma 22 only yields non-vacuous bounds in the case $K = 2$. Indeed, when $K > 2$, the sequence $\{\alpha_{n,k}\}_{n \in \mathbb{N}^*, k \in \{0, \dots, K-1\}}$ leads to increase the bounds on the quantities $\exp[\psi_{i_k}^{n-1} - \psi_{i_k}^n]$, which motivates the use of A5.

D.3 Proof of Section 5

For the rest of this section, we consider the multi-marginal Schrödinger bridge problem given by (TreeSB) and establish in Proposition 24 the correspondence with the regularized Wasserstein propagation problem presented in Solomon et al. (2014, 2015). We first state a technical result.

Lemma 23. *Let $\varepsilon > 0$. Assume that π^0 is given by (2), where $r \in \mathbb{V}$ is chosen arbitrarily. Then, for any $\pi \in \mathcal{P}_{\mathbb{T}_r}$, we have*

$$\begin{aligned} \varepsilon \text{KL}(\pi \mid \pi^0) &= \sum_{(v,v') \in E_r} \{w_{v,v'} \mathbb{E}_{\pi_{v,v'}} [\|X_v - X_{v'}\|^2] - \varepsilon \mathbb{H}(\pi_{v,v'})\} \\ &\quad + \varepsilon \sum_{v \in \mathbb{V}} \text{card}(C_v) \mathbb{H}(\pi_v) + \varepsilon \text{KL}(\pi_r \mid \pi_r^0), \end{aligned}$$

where we recall that $C_v = \{v' \in \mathbb{V} : (v, v') \in E_r\}$.

Proof. Since $\pi, \pi^0 \in \mathcal{P}_{T_r}$, we obtain the following decomposition

$$\begin{aligned} & \text{KL}(\pi \mid \pi^0) \\ &= \text{KL}(\pi_r \prod_{(v,v') \in E_r} \pi_{v'|v} \mid \pi_r^0 \prod_{(v,v') \in E_r} \pi_{v'|v}^0) \\ &= \text{KL}(\pi_r \mid \pi_r^0) + \sum_{(v,v') \in E_r} \int_{\mathbb{R}^d} \text{KL}(\pi_{v'|v}(\cdot \mid x_v) \mid \pi_{v'|v}^0(\cdot \mid x_v)) d\pi_v(x_v) \\ &= \text{KL}(\pi_r \mid \pi_r^0) - \sum_{(v,v') \in E_r} \int_{\mathbb{R}^d \times \mathbb{R}^d} \log \pi_{v'|v}^0 d\pi_{v,v'} - \sum_{(v,v') \in E_r} \int_{\mathbb{R}^d} \text{H}(\pi_{v'|v}(\cdot \mid x_v)) d\pi_v(x_v). \end{aligned}$$

We finally obtain the result by using the definition of π^0 and noticing that $\int_{\mathbb{R}^d} \text{H}(\pi_{v'|v}(\cdot \mid x_v)) d\pi_v(x_v) = \text{H}(\pi_{v,v'}) - \text{H}(\pi_v)$ for any $(v, v') \in E_r$. \square

Proposition 24. Let $\varepsilon > 0$ and $\mu_0 \in \mathcal{P}$ such that $\mu_0 \ll \text{Leb}$. Assume that π^0 is given by (2), where $r \in V$ is chosen arbitrarily, and that $\varphi_r = d\mu_0/d\text{Leb}$. Also assume **A2**. Then, the set of marginals of the solution to **(TreeSB)** is exactly the solution to the entropic-regularized Wasserstein Propagation problem (Solomon et al., 2014, 2015) defined by

$$\begin{aligned} & \arg \min \left\{ \sum_{(v,v') \in E_r} w_{v,v'} W_{2,\varepsilon/w_{v,v'}}^2(\nu_v, \nu_{v'}) + \varepsilon \sum_{v \in V} \text{card}(C_v) \text{H}(\nu_v) + \varepsilon \text{KL}(\nu_r \mid \mu_0) : \text{(WP)} \right. \\ & \left. \{ \nu_v \}_{v \in V} \in \mathcal{P}^{\ell+1}, \nu_i = \mu_i, \forall i \in S \right\}, \end{aligned}$$

where we recall that $C_v = \{v' \in V : (v, v') \in E_r\}$.

Proof. Assume that π^0 is given by (2), where $r \in V$ is chosen arbitrarily, and that $\varphi_r = d\mu_0/d\text{Leb}$. In particular, we have $\pi_r^0 = \mu_0$. Moreover, it is clear that π^0 verifies **A1**, and **A3** by Proposition 16.

Let $\{\nu_v\}_{v \in V} \in \mathcal{P}^{\ell+1}$ and $\{\nu^{(v,v')}\}_{(v,v') \in E_r} \in (\mathcal{P}^{(2)})^{|E_r|}$. We define

$$\begin{aligned} F(\{\nu_v\}) &= \sum_{(v,v') \in E_r} w_{v,v'} W_{2,\varepsilon/w_{v,v'}}^2(\nu_v, \nu_{v'}) + \varepsilon \sum_{v \in V} \text{card}(C_v) \text{H}(\nu_v) + \varepsilon \text{KL}(\nu_r \mid \mu_0), \\ G(\nu_r, \{\nu^{(v,v')}\}) &= \sum_{(v,v') \in E_r} \{ w_{v,v'} \mathbb{E}_{\nu^{(v,v')}} [\|X_v - X_{v'}\|^2] - \varepsilon \text{H}(\nu^{(v,v')}) \} \\ & \quad + \varepsilon \sum_{(v,v') \in E_r} \text{H}(\nu_v^{(v,v')}) + \varepsilon \text{KL}(\nu_r \mid \mu_0). \end{aligned}$$

By definition of the regularized Wasserstein distance given in (3), we have for any $\{\nu_v\}_{v \in V} \in \mathcal{P}^{\ell+1}$

$$F(\{\nu_v\}) = \min \{ G(\nu_r, \{\nu^{(v,v')}\}) : \nu^{(v,v')} \in \mathcal{P}^{(2)}, \nu_v^{(v,v')} = \nu_v, \nu_{v'}^{(v,v')} = \nu_{v'}, \forall (v, v') \in E_r \}. \quad (26)$$

In particular, we have $F(\{\pi_v\}) \leq G(\pi_r, \{\pi_{v,v'}\})$ for any $\pi \in \mathcal{P}^{(\ell+1)}$. We now prove the result of Proposition 24 in two steps denoted by **Step 1** and **Step 2**.

Step 1. Let us not assume **A2** for now. In this case, we prove in **Step 1.a** and **Step 1.b** that solving **(WP)** is equivalent to solving a modified version of **(TreeSB)** given by

$$\pi^* = \arg \min \{ \text{KL}(\pi \mid \pi^0) : \pi \in \mathcal{P}_{T_r}, \pi_i = \mu_i, \forall i \in S \}. \quad (T_r\text{-TreeSB})$$

Remark that any solution to **(T_r-TreeSB)** is a solution to **(TreeSB)**, but the converse result may not be true.

Step 1.a: (WP) \implies (T_r-TreeSB). Consider a solution $\{\nu_v^*\}_{v \in V}$ to **(WP)**. For any $(v, v') \in E_r$, $W_{2,\varepsilon/w_{v,v'}}^2(\nu_v^*, \nu_{v'}^*)$ is well defined and thus, there exists $\nu^{(v,v')} \in \Pi(\nu_v^*, \nu_{v'}^*)$ such that

$$\nu^{(v,v')} \in \arg \min \{ \mathbb{E}_\pi [\|X_v - X_{v'}\|^2] - (\varepsilon/w_{v,v'}) \text{H}(\pi) : \pi \in \Pi(\nu_v^*, \nu_{v'}^*) \}. \quad (27)$$

Using the gluing lemma, we build the probability measure $\pi^* = \nu_r^* \prod_{(v,v') \in E_r} \nu_{v'|v}^{(v,v')}$ such that (i) $\pi^* \in \mathcal{P}_{T_r}$, and (ii) $\pi_{v,v'}^*$ and $\nu^{(v,v')}$ have the same distribution for any $(v, v') \in E_r$. In particular, we have $\pi_i^* = \mu_i$ for any $i \in S$.

Let us show now that π^* is a solution to **(T_r-TreeSB)**. Let $\pi \in \mathcal{P}_{T_r}$ such that $\pi_i = \mu_i$ for any $i \in S$. We have

$$\begin{aligned} \epsilon\text{KL}(\pi \mid \pi^0) &= G(\pi_r, \{\pi_{v,v'}\}) && \text{(Lemma 23)} \\ &\geq F(\{\pi_v\}) \\ &\geq F(\{\nu_v^*\}) && \text{(definition of } \nu^*) \\ &= G(\nu_r^*, \{\nu^{(v,v')}\}) && \text{(see (27))} \\ &= G(\pi_r^*, \{\pi_{(v,v')}^*\}) && \text{(definition of } \pi^*) \\ &= \epsilon\text{KL}(\pi^* \mid \pi^0). && \text{(Lemma 23)} \end{aligned}$$

Therefore, π^* is a solution to **(T_r-TreeSB)**.

Step 1.b: **(T_r-TreeSB)** \implies **(WP)**. Consider now a solution π^* to **(T_r-TreeSB)**. Since $\pi^* \in \mathcal{P}_{T_r}$, we have $\pi^* = \pi_r^* \prod_{(v,v') \in E_r} \pi_{v'|v}^*$ and $\pi_i^* = \mu_i$ for any $i \in S$.

Let us show that $\{\pi_v^*\}_{v \in V}$ is a solution to **(WP)**. Let $\{\nu_v\}_{v \in V} \in \mathcal{P}^{\ell+1}$ such that $\nu_i = \mu_i$ for any $i \in S$.

Let $\{\nu^{(v,v')}\}_{(v,v') \in E_r}$ be a family of probability measures such that $\nu^{(v,v')} \in \mathcal{P}^{(2)}$, $\nu_v^{(v,v')} = \nu_v$, $\nu_{v'}^{(v,v')} = \nu_{v'}$ for any $(v, v') \in E_r$.

Using the gluing lemma, we build the probability measure $\pi = \nu_r \prod_{(v,v') \in E_r} \nu_{v'|v}^{(v,v')}$, such that (i) $\pi \in \mathcal{P}_{T_r}$ and (ii) $\pi_{v,v'}$ and $\nu^{(v,v')}$ have the same distribution for any $(v, v') \in E_r$. We have

$$\begin{aligned} \epsilon\text{KL}(\pi \mid \pi^0) &= G(\pi_r, \{\pi_{(v,v')}\}) && \text{(Lemma 23)} \\ &= G(\nu_r, \{\nu^{(v,v')}\}) && \text{(definition of } \pi) \\ &\geq \epsilon\text{KL}(\pi^* \mid \pi^0) && \text{(definition of } \pi^*) \\ &= G(\pi_r^*, \{\pi_{(v,v')}^*\}) && \text{(Lemma 23)} \end{aligned}$$

By taking the infimum in the previous inequality over the families $\{\nu^{(v,v')}\}_{(v,v') \in E_r}$, we obtain by **(26)** that

$$F(\{\nu_v\}) \geq G(\pi_r^*, \{\pi_{(v,v')}^*\}) \geq F(\{\pi_v^*\}),$$

and therefore, $\{\pi_v^*\}_{v \in V}$ is a solution to **(WP)**.

Step 2. We now assume **A2**. By Proposition 3, there exists a unique solution $\pi^* \in \mathcal{P}^{(\ell+1)}$ to **(TreeSB)** such that we π^0 -a.s. have $(d\pi^*/d\pi^0) = \exp[\bigoplus_{i \in S} \psi_i^*]$, where $\{\psi_i^*\}_{i \in S}$ are measurable potentials with $\psi^* : \mathbb{R}^d \rightarrow \mathbb{R}$. Since $\pi^0 \in \mathcal{P}_{T_r}$, we also have $\pi^* \in \mathcal{P}_{T_r}$, i.e., the potentials $\{\psi_i^*\}_{i \in S}$ do not modify the Markovian nature of π^0 . Therefore, π^* is also the unique solution to **(T_r-TreeSB)**. Using the equivalence between **(T_r-TreeSB)** and **(WP)** established in **Step 1**, we finally obtain the result of Proposition 24. \square

In particular, Proposition 7 directly derives from Proposition 24 by taking $r = i_{K-1}$ and $\mu_0 = \mu_{i_{K-1}}$.

D.4 Comparison with Haasler et al. (2021)

In their work, Haasler et al. (2021) study the *static* and *discrete-state* counterpart of our approach. Given a state space X such that $|X| = n + 1$ with $n \in \mathbb{N}$, they establish a correspondence between multi-marginal EOT with a general tree-based cost and discrete-time multi-marginal static Schrödinger bridge, and provide an efficient method to solve these problems. In this section, we provide details on their framework and give a precise comparison between our theory and their results.

To be coherent with the setting of Haasler et al. (2021), we adapt here some of our notation. Let us define $Z^{(q)} = \mathbb{R}_+^{(n+1)^q}$. For any $q \in \mathbb{N}^*$, the set of probability measures on X^q is defined as $\mathcal{P}^{(q)} = \{M \in Z^{(q)} : \langle M, \mathbf{1} \rangle = 1\}$. We denote $\mathcal{P} = \mathcal{P}^{(1)}$. For any tensors $M, P \in Z^{(q)}$, the Kullback-Leibler divergence between M and P is defined as $\text{KL}(M \mid P) = \langle M \log(M/P) - M + P, \mathbf{1} \rangle$ and the

entropy of M is defined as $H(M) = -\text{KL}(M \mid \mathbf{1})$, where the operations are meant componentwise. In the rest of the section, we consider an undirected tree $T = (V, E)$ with $|V| = \ell + 1$ such that V may be identified with $\{0, \dots, \ell\}$.

Details on the results of Haasler et al. (2021). In their paper, the authors consider a cost tensor $C \in \mathbb{Z}^{(\ell+1)}$ that factorizes along T , i.e., for any $\{j_0, \dots, j_\ell\}$ with for any $i \in \{0, \dots, \ell\}$, $j_i \in \{0, \dots, n\}$, we have

$$C_{j_0, \dots, j_\ell} = \sum_{(v, v') \in E} C_{j_v, j_{v'}}^{\{v, v'\}},$$

where $C^{\{v, v'\}} \in \mathbb{Z}^{(2)}$ is a cost matrix for transportation between the marginals at vertices v and v' , see (Haasler et al., 2021, Eq. (3.1)). In particular, this cost can be seen as the discrete counterpart of the tree-based cost introduced in (1) in the quadratic setting.

Given a subset $S \subset V$ with $|S| = K$ and a set of marginals $\{\mu_i\}_{i \in S} \in \mathcal{P}^K$, Haasler et al. (2021) study the EOT problem associated to T , see (Haasler et al., 2021, Eq. (2.4)), which is given by

$$\operatorname{argmin}\{\langle C, M \rangle - \varepsilon H(M) : M \in \mathcal{P}^{(\ell+1)}, \operatorname{proj}_i(M) = \mu_i, \forall i \in S\}. \quad (\text{discrete-EmOT})$$

This problem may be solved with Sinkhorn algorithm (Cuturi, 2013; Knight, 2008; Sinkhorn & Knopp, 1967), for which the authors provide an efficient implementation adapted to the tree-based setting, see (Haasler et al., 2021, Algorithm 3.1). Moreover, they state the convergence of their method in (Haasler et al., 2021, Theorem 3.5), as a direct consequence of the results presented in Luo & Tseng (1992).

In (Haasler et al., 2021, Section 4.2), it is assumed that S corresponds to the set of the leaves of T , as we do, and it is shown an equivalence between (discrete-EmOT) and the discrete-state static SB problem stated in (Haasler et al., 2021, Eq 4.2), which is given by

$$\operatorname{argmin}\{\sum_{(v, v') \in E_r} \text{KL}(M^{(v, v')} \mid \operatorname{diag}(\nu_v)A^{(v, v')}) : \quad (\text{discrete-TreeSB}) \\ M^{(v, v')} \in \mathcal{P}^{(2)}, \{\nu_v\}_{v \in V} \in \mathcal{P}^{\ell+1}, M^{(v, v')} \mathbf{1} = \nu_v, M^{(v, v')}^\top \mathbf{1} = \nu_{v'}, \nu_i = \mu_i, \forall i \in S\},$$

where $T_r = (V, E_r)$ is the directed version of T rooted in an arbitrary vertex $r \in S$, and $A^{(v, v')} = \exp(-C^{(v, v')}/\varepsilon) \in \mathbb{Z}^{(2)}$ for any $(v, v') \in E_r$. Remark that $A^{(v, v')}$ may not necessarily be a transition probability matrix.

Finally, Haasler et al. (2021) provide two main numerical experiments. In (Haasler et al., 2021, Section 5.2), they consider a tree with 15 vertices, 14 edges and 8 leaves, combined to the state-space $X = \{0, 1\}^{50 \times 50}$, and solve the corresponding (discrete-EmOT) problem for the quadratic cost. In (Haasler et al., 2021, Section 6), they apply their methodology to estimate ensemble flows on a hidden Markov chain. Given $\tau \in \mathbb{N}^*$, they consider a tree T with τ internal vertices (modeling the distribution of N agents at time $t \in \{1, \dots, \tau\}$), that are linearly linked, and such that each of these vertices is independently linked to S leaves of T (modeling observations at time $t \in \{1, \dots, \tau\}$). In this setting, the state space is given by $X = \{1, \dots, 100\}^N$. They solve the formulation (discrete-TreeSB) where the reference measure is chosen as a random walk.

Comparison with our results. We now establish remarks on the main differences between our methodology and the work of Haasler et al. (2021).

First of all, the continuous state-space counterpart of (discrete-TreeSB) is given by

$$\operatorname{argmin}\{\text{KL}(\pi \mid \pi^0) : \pi \in \mathcal{P}_{T_r}, \pi_i = \mu_i, \forall i \in S\}, \quad (28)$$

where π^0 is a reference measure which factorizes along T_r . In this case, $\pi_{v, v'}$, π_v and $\pi_{v'|v}^0$ in (28) respectively correspond to the continuous version of $M^{(v, v')}$, ν_v and $A^{(v, v')}$ in (discrete-TreeSB). In contrast, our formulation of the multi-marginal Tree Schrödinger Bridge problem given in (TreeSB) is a minimization problem over all probability measures $\pi \in \mathcal{P}^{(\ell+1)}$, and is not restricted to the distributions that admit a Markovian factorization along T as in (28). Hence, our framework may be considered more general. Remark that under A1, A2 and A3, Proposition 3 states that (TreeSB) admits a unique solution $\pi^* \ll \pi^0$ such that $(d\pi^*/d\pi^0)$ can be written with potentials. Then, $\pi^* \in \mathcal{P}_{T_r}$ since $\pi^0 \in \mathcal{P}_{T_r}$, and (TreeSB) is then equivalent to (28).

Furthermore, (EmOT) is more general than the continuous version of (discrete-EmOT), which we can recover by taking any measure ν of the form $(d\nu/d\text{Leb}) = \exp[\bigoplus_{i \in S} \varphi_i]$ in (EmOT), where $\{\varphi_i\}_{i \in S}$ is a family of potentials such that $|\int_{\mathbb{R}^d} \varphi_i d\mu_i| < \infty$ for any $i \in S$. As a consequence, our setting allows us to choose the root $r \in V \setminus S$ for the SB problem, whereas Haasler et al. (2021) only consider the case where $r \in S$. In the latter case, we establish in Appendix E that r can be chosen arbitrarily, as stated by (Haasler et al., 2021, Corollary 4.3).

Finally, TreeDSB deeply differs from the framework of Haasler et al. (2021) due its *dynamic* nature. Although we solve the same tree-based static SB problem (up to continuous/discrete state-space consideration), our approach consists in computing dynamic iterates (*i.e.*, path measures) using diffusion-based methods instead of static iterates (*i.e.*, distributions) using Sinkhorn algorithm. This paradigm is at the core of the DSB (De Bortoli et al., 2021) methodology, and offers an efficient approach to tackle high-dimensional settings, where Sinkhorn algorithm would fail.

Here, we present some advantages of the method proposed by Haasler et al. (2021) compared to ours. First, Haasler et al. (2021) may choose any kind of tree-based cost in practice, while our methodology only holds for the quadratic cost. This limitation is shared with all approaches based on the DSB (De Bortoli et al., 2021) methodology. Indeed, since the cost is determined by the reference path measure, we often choose quadratic costs associated with Brownian motions or Ornstein-Uhlenbeck processes. Moreover, Haasler et al. (2021) may consider various inhomogeneous (discrete) state spaces for the vertices of T , as presented in their numerical experiments. In our case, this approach is not compatible with our diffusion-based method. Finally, unlike Haasler et al. (2021), our method is not scalable with the number of vertices or edges in T due to computational limits. This limitation is common to all multi-marginal approaches which rely on neural networks to parameterize the potential and/or the distributions of the multi-marginal OT method, see Li et al. (2020); Fan et al. (2020); Korotin et al. (2022, 2021) for instance.

E Further results on TreeSB

Choice of the root r in (TreeSB). We recall that the reference measure π^0 considered in (TreeSB), which is defined in (2), verifies $\pi^0 \in \mathcal{P}_{T_r}$ for some fixed root $r \in V$ and $\pi_r^0 \ll \text{Leb}$ with density φ_r . Moreover, we have $\pi_{v'|v}^0(\cdot | x_v) = N(x_v, \varepsilon/(2w_{v,v'})I_d)$ for any $(v, v') \in E_r$, and thus, π^0 is entirely determined by the choice of the root r and the density on the corresponding vertex φ_r .

As presented in Appendix D.1, we recall that (TreeSB) is equivalent to any multi-marginal Tree-SB problem with a reference measure $\bar{\pi}^0$ given by (15), *i.e.*, $\bar{\pi}^0$ writes as $(d\bar{\pi}^0/d\pi^0) = \exp[\bigoplus_{i \in S} \psi_i^0]$, where $\{\psi_i^0\}_{i \in S}$ is a family of measurable potentials with $\psi_i^0 : \mathbb{R}^d \rightarrow \mathbb{R}$ such that $|\int_{\mathbb{R}^d} \psi_i^0 d\mu_i| < \infty$ for any $i \in S$. In the case where r is chosen as a leaf of T , this result implies that (TreeSB) is unchanged if

- (a) $\varphi_r = d\nu/d\text{Leb}$ where $\nu \in \mathcal{P}$ is such that $\text{KL}(\mu_r|\nu) < \infty$,
- (b) r is replaced by $r' \in S$, as long as $H(\mu_r) < \infty$ and $H(\mu_{r'}) < \infty$.

Therefore, under A0, the setting chosen in Section 3 is equivalent to any other setting where r is arbitrarily chosen in S and $\varphi_r = d\nu/d\text{Leb}$ where $\text{KL}(\mu_r|\nu) < \infty$.

Consider now the case where $r \in S^c$, *i.e.*, r is not a leaf of T . Then, the choice of φ_r can not be made arbitrarily anymore, since it determines a further regularization on the r -th marginal of the solution to (TreeSB). In this setting, the sequence defined by (mIPF) is unchanged. Hence, TreeDSB proceeds in the same manner as presented in Section 3, except for the first iteration, which we detail now.

Let us define $P = \text{path}_{T_{i_0}}(i_0, r)$, where $T_{i_0} = (V, E_{i_0})$ is the directed version of T rooted in i_0 . We recall that first iterate of (mIPF) is defined by

$$\pi^1 = \text{argmin}\{\text{KL}(\pi | \pi^0) : \pi \in \mathcal{P}^{(\ell+1)}, \pi_{i_0} = \mu_{i_0}\}.$$

Following the proof of Lemma 12, it is clear that

$$\pi^1 = \mu_{i_0} \otimes_{(v,v') \in P} \pi_{v'|v}^0 \otimes_{(v,v') \in E_{i_0} \setminus P} \pi_{v'|v}^0 = \mu_{i_0} \otimes_{(v,v') \in E_{i_0}} \pi_{v'|v}^0,$$

where we emphasize that $P = \{(v, v') \in E_{i_0} : (v', v) \in E_r\}$. Therefore, Proposition 2 still applies between r and i_0 , by considering r instead of i_{K-1} . In practice, this means that the first iteration of TreeDSB consists in computing the time reversal of the path measures $\mathbb{P}_{(v',v)}^0$ for any $(v, v') \in P$.

Extension of the regularized Wasserstein barycenter problem (regWB). Consider the regularized Wasserstein-2 barycenter problem defined as follows

$$\mu_\varepsilon^* = \arg \min \left\{ \sum_{i=1}^\ell w_i W_{2,\varepsilon/w_i}^2(\mu, \mu_i) + \ell \varepsilon H(\mu) + \varepsilon \text{KL}(\mu \mid \mu_0) : \mu \in \mathcal{P} \right\}, \quad (\mu_0\text{-regWB})$$

where $(w_i)_{i \in \{1, \dots, \ell\}} \in (0, +\infty)^\ell$ and $\mu_0 \in \mathcal{P}$ is a reference measure. This formulation admits a further regularization compared to (regWB), which tends to make μ_ε^* closer to μ_0 . In particular, given a Wasserstein barycenter problem onto a star-shaped tree, the formulation (μ_0 -regWB) may be more adapted than (regWB) if we have an *a priori* on the form of the regularized barycenter. In the case where $\mu_0 = N(0, \sigma_0^2 \text{Id})$, letting $\sigma_0 \rightarrow \infty$, we recover the $(\ell\varepsilon, (\ell-1)\varepsilon)$ doubly-regularized Wasserstein barycenter problem (regWB). In the same spirit as Proposition 7, we can derive the following result from Proposition 24, which proves that (μ_0 -regWB) can be solved with TreeDSB.

Proposition 25. *Let $\varepsilon > 0$ and $\mu_0 \in \mathcal{P}$ such that $\mu_0 \ll \text{Leb}$. Assume A0. Also assume that \mathbb{T} is a star-shaped tree with central node indexed by 0, and that the reference measure of (TreeSB) defined in (2) verifies $r = 0$ and $\varphi_r = d\mu_0/d\text{Leb} > 0$. Under A2, (μ_0 -regWB) has a unique solution π_0^* , where π^* is the solution to (TreeSB).*

Below, we provide practical guidelines to parameterize μ_0 when it is chosen as a Gaussian distribution.

Gaussian design of μ_0 in (μ_0 -regWB). Consider an undirected star-shaped tree \mathbb{T} with $K+1$ vertices and leaves $\{1, \dots, K\}$. In order to incorporate the marginal constraints in the penalization brought by μ_0 when it is a Gaussian distribution, we set its mean to $\sum_{i=1}^K \mathbb{E}[\mu_i]/K$ and its diagonal covariance matrix as $\alpha \times (\sum_{i=1}^K \text{diag}(\text{Cov}[\mu_i])^{-1}/K)^{-1}$, where the inverse operation is component-wise and α is a positive hyperparameter. This choice of variance helps to correctly explore the state-space at the very first iteration of TreeDSB, which is key to ensure numerical stability. In this setting, (TreeSB) verifies A2 and A3, by Proposition 15 and Proposition 16. In particular, we use this approach for two of our experiments: synthetic Gaussian datasets and Bayesian fusion, see Appendix G.

F Algorithmic techniques

Time discretization in TreeDSB. Denote $k_n = (n-1) \bmod(K)$ for any $n \in \mathbb{N}$. Let $\mathbb{T} = (V, E)$ be a weighted undirected tree and consider the multi-marginal Schrödinger bridge problem (TreeSB) associated to this tree. We recall that for any $\{v, v'\} \in E$, we define $T_{v,v'} = \varepsilon/(2w_{v,v'})$.

Consider the path measures $\{\mathbb{P}_{(v,v')}^n\}_{n \in \mathbb{N}, (v,v') \in E_{k_n}}$ recursively defined by (a) and (b). By combining Proposition 1, Proposition 2 and results on time reversal theory (Haussmann & Pardoux, 1986), we obtain by recursion that for any $n \in \mathbb{N}$, any $(v, v') \in E_{k_n}$, $\mathbb{P}_{(v,v')}^n$ is associated with a Stochastic Differential Equation on $[0, T_{v,v'}]$ given by

$$d\mathbf{X}_t = f_{t,v,v'}^n(\mathbf{X}_t)dt + d\mathbf{B}_t, \quad \mathbf{X}_0 \sim \pi_v^n. \quad (29)$$

Let $N \in \mathbb{N}^*$. In order to sample from the dynamics (29) at iteration $n \in \mathbb{N}$, we consider its Euler-Maruyama discretization on $(N+1)$ time steps,

$$X_{m+1} = X_m + \gamma_{m+1} f_{t_m,v,v'}^n(X_m) + \sqrt{\gamma_{m+1}} Z_{m+1}, \quad X_0 \sim \pi_v^n, \quad (30)$$

where $Z_m \sim N(0, \text{Id})$ for any $m \in \{1, \dots, N\}$, $t_m = \sum_{i=1}^m \gamma_i$, and $\{\gamma_m\}_{m=1}^N \in (0, \infty)^N$ is a time schedule such that $\sum_{m=1}^N \gamma_m = T_{v,v'}$. This results in approximating the path measure $\mathbb{P}_{(v,v')}^n$ by the joint distribution $\pi_{(v,v')}^{n,N} \in \mathcal{P}^{(N+1)}$ defined by

$$\pi_{(v,v')}^{n,N} = \pi_v^n \otimes_{m=0}^{N-1} \pi_{(v,v'),m+1|m}^{n,N},$$

where $\pi_{(v,v'),m+1|m}^{n,N}(\cdot | x_m) = N(x_m + \gamma_{m+1} f_{t_m,v,v'}^n(x_m), \gamma_{m+1} \text{Id})$ for any $m \in \{0, \dots, N-1\}$. If N is chosen large enough, then $\pi_{(v,v'),m}^{n,N}$ and $\mathbb{P}_{(v,v'),t_m}^n$ have approximately the same distribution

for any $m \in \{0, \dots, N\}$. Consequently, $(\mathbb{P}_{(v,v')}^n)^R$ is naturally approximated by the joint distribution $\tilde{\pi}_{(v,v')}^{n,N} \in \mathcal{P}^{(N+1)}$ defined by

$$\tilde{\pi}_{(v,v')}^{n,N} = \pi_{v'}^n \otimes_{m=0}^{N-1} \pi_{(v,v'),N-m-1|N-m}^{n,N}.$$

If N is chosen large enough, we obtain that

$$\pi_{(v,v'),N-m-1|N-m}^{n,N}(\cdot|x_{N-m}) = \mathbb{N}(x_{N-m} - \gamma_{N-m} f_{t_{N-m},v,v'}^n(x_{N-m}) + \gamma_{N-m} \nabla \log p_{v,v',t_{N-m}}(x_{N-m}), \gamma_{N-m} \mathbf{I}_d),$$

where $p_{v,v',t}$ is the density of $\mathbb{P}_{(v,v'),t}^n$ w.r.t. the Lebesgue measure.

Following the construction of our dynamic iterates, we now explain how the sequence $\{\pi_{(v,v')}^n\}_{n \in \mathbb{N}^*, (v,v') \in E_{k_n}}$ is recursively defined. Let $n \in \mathbb{N}$, $k_n = (n-1) \bmod(K)$. Define the path $P_n = \text{path}_{T_{i_{k_n}}}(i_{k_n}, i_{k_n+1})$. Then, for any $(v, v') \in E_{k_n+1}$,

- (a) if $(v, v') \in E_{k_n} \setminus P_n$, then $\pi_{(v,v')}^{n+1,N} = \pi_v^{n+1} \otimes_{m=0}^{N-1} \pi_{(v,v'),m+1|m}^{n,N}$,
- (b) if $(v', v) \in P_n$, then $\pi_{(v,v')}^{n+1,N} = \pi_v^{n+1} \otimes_{m=0}^{N-1} \pi_{(v',v),N-m-1|N-m}^{n,N}$.

These computations may be obtained by considering the sequence given by (mIPF) to solve the multi-marginal Tree-SB problem associated to $T^{(N)} = (V^{(N)}, E^{(N)})$, the N -discretized version of T (see Appendix B) with weights $w_{e_m}^{(N)} = 2\gamma_m/\varepsilon$, which is given by

$$\pi^* = \text{argmin}\{\text{KL}(\pi|\pi^{0,N}) : \pi \in \mathcal{P}^{(V^{(N)})}, \pi_i = \mu_i, \forall i \in S\},$$

with $\pi^{0,N} = \pi_r^0 \otimes_{(v,v') \in E_r} \pi_{(v,v'),1:N|0}^{0,N}$.

To approximate the IPF recursion given by (a) and (b), we use on each edge of T the score-matching approach of De Bortoli et al. (2021), which avoids heavy computations of score approximations. The next proposition is direct adaptation of (De Bortoli et al., 2021, Proposition 3).

Proposition 26. Assume that for any $n \in \mathbb{N}$, any $(v, v') \in E_{k_n}$ with $k_n = (n-1) \bmod(K)$, we have

$$\pi_{(v,v'),m+1|m}^{n,N}(\cdot|x_m) = \mathbb{N}(F_{m,v,v'}^n(x_m), \gamma_m \mathbf{I}_d).$$

Let $n \in \mathbb{N}$. Consider the path $P_n = \text{path}_{T_{i_{k_n}}}(i_{k_n}, i_{k_n+1})$. Let $(v, v') \in E_{k_n+1}$. Define $p^n = \pi_{(v,v')}^{n,N}$ and $m_N = N - m - 1$. Then, if $(v', v) \in P_n$, we have

$$F_{m,v,v'}^{n+1} = \text{argmin}_{F \in L^2(\mathbb{R}^d, \mathbb{R}^d)} \mathbb{E}_{P_{m_N, m_N+1}^n} [\|F(X_{m_N+1}) - (X_{m_N+1} + F_{m_N, v', v}^n(X_{m_N}) - F_{m_N, v', v}^n(X_{m_N+1}))\|^2], \quad (31)$$

otherwise, we have $F_{m,v,v'}^{n+1} = F_{m,v,v'}^n$.

In practice, we use two neural networks per edge $\{v, v'\} \in E$, one for each possible direction of the edge, such that $F_{v,v'}(\theta_{v,v'}^n, m, x) \approx F_{m,v,v'}^n(x)$ and $F_{v',v}(\theta_{v',v}^n, m, x) \approx F_{m,v',v}^n(x)$. For any $\{v, v'\} \in E$, the parameter $\theta_{v,v'}^n$ is updated at iteration n via the score matching loss defined by (31) in Proposition 26 if $(v, v') \in \text{path}_{T_{i_{k_n}}}(i_{k_n}, i_{k_n+1})$, see Algorithm 1.

G Additional experimental results and details

The numerical experiments presented in Section 7 are obtained by our own Pytorch implementation, which is inspired from the code⁶ provided by De Bortoli et al. (2021). We first provide information on the general setting of our experiments in Appendix G.1, and then give details on each of them in Appendix G.2 along with additional results. We recall that a mIPF cycle is defined as a subset of K consecutive iterations of (mIPF) and that the order of the leaves given by $\{i_0, \dots, i_{K-1}\}$ is randomly shuffled at each new mIPF cycle.

⁶https://github.com/JTT94/diffusion_schrodinger_bridge

G.1 General experimental setup

Implementation of Algorithm 1 in practice. Let $n \in \mathbb{N}$, with $k_n = (n - 1) \bmod(K)$, $k_n + 1 = n \bmod(K)$. Consider the path $P_n = \text{path}_{\mathbb{T}_{i_{k_n}}}(i_{k_n}, i_{k_n+1})$. Assume that we are provided with a dataset $D_{i_{k_n}}$, which contains M samples from $\pi_{i_{k_n}}^n$. Following Lines 7-9 in Algorithm 1, we apply processes (a) and (b) recursively on the edges $(v, v') \in P_n$.

- (a) **Sampling step (Line 7).** For any $x_0 \in D_v$, we sample from the diffusion trajectory (30) given by the Euler Maruyama discretization of $\mathbb{P}_{v,v'}^n$ starting from x_0 . This gives us $M \times N$ trajectory samples. We then store the last iterate of each trajectory in a new dataset $D_{v'}$, which thus approximates $\pi_{v'}^n$.
- (b) **Training step (Lines 8-9).** In order to avoid heavy computation, we approximate the *mean-matching* loss (31) by an unbiased estimator obtained by subsampling b elements from the *full* trajectories computed in the sampling process, see (De Bortoli et al., 2021, Eq. (97)-(98)). Here, b refers to the *batch-size* parameter of the neural networks. Then, we perform gradient descent to optimize the parameter $\theta_{v',v}$, which parameterizes the *backward* drift on the edge (v, v') .

To avoid any bias issue, the whole trajectories obtained at process (a) are refreshed at a certain frequency over the training iterations of the neural networks by once again simulating the diffusion (30). In our experiments, this refresh occurs each 500 iterations.

Setting of the time discretization. The number of time-steps N in the time discretization of the diffusions is chosen to be even and identical for each of the edges of the tree. Let $\{v, v'\} \in \mathbf{E}$. We now give details on the design of the time schedule $\{\gamma_k\}_{k=1}^N$ related to the edge $\{v, v'\}$, see Appendix F. Following De Bortoli et al. (2021), we choose this sequence to be invariant by time reversal and consider $\gamma_k = \gamma_0 + (2k/N)(\bar{\gamma} - \gamma_0)$ for any $k \in \{0, \dots, N/2\}$ (the rest of the sequence being obtained by symmetry) where γ_0 is a free parameter and $\bar{\gamma}$ is determined by $\sum_{k=1}^N \gamma_k = T_{v,v'}$. In our experiments, we set $N = 50$ and $\gamma_0 = 10^{-5}$.

Sampling improvement. In our code, we implemented the corrector scheme of Song et al. (2021) and the *probability flow*-based sampling approach detailed in (De Bortoli et al., 2021, Section H.3), but did not observe any significant improvement in our experiments using one of these techniques.

Choice of the architectures of the neural networks. In the case of the experiments related to synthetic datasets (two-dimensional toy datasets, Gaussian distributions) and to the subset posterior aggregation task, we implement the same architecture as presented in (De Bortoli et al., 2021, Figure 3). We refer to this model as “Basic Model” and detail it in Figure 6. In the “Basic Model”, the PositionalEncoding block applies the sine transform described in Vaswani et al. (2017), with output dimension equal to 32, and each MLP Block represents a Multilayer Perceptron Network. In particular, MLPBlock (1a) has shape $(d, 128, \max(256, 2d))$, MLPBlock (1b) has shape $(32, 128, \max(256, 2d))$, and MLPBlock (2) has shape $(2 \times \max(256, 2d), \max(256, 2d), \max(128, d), d)$, where d denotes the dimension of input data. We optimize the networks with ADAM (Kingma & Ba, 2014) with learning rate 10^{-4} and momentum 0.9. For each of the networks, we set the batch size to 4,096 and the number of iterations to 10,000 for the synthetic datasets and 15,000 for the subset posterior aggregation task. Our experiments ran on 1 Intel Xeon CPU Gold 6230 20 cores @ 2.1 Ghz CPU.

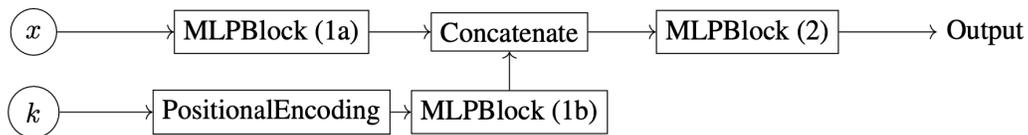


Figure 6: Architecture of the “Basic Model”.

In the case of the experiments related to MNIST dataset, we use a reduced UNET architecture based on Nichol & Dhariwal (2021), where we set the number of channels to 64 rather than 128. We implement an exponential moving average of network parameters across training iterations, with rate 0.999. We optimize the networks with ADAM (Kingma & Ba, 2014) with learning rate 10^{-4} and momentum 0.9. Finally, we set the batch size to 256 and the number of training iterations to 30,000. Our experiments ran using 1 Nvidia A100.

Details on regularized state-of the art methods. We run the fsWB algorithm (Cuturi & Doucet, 2014) with the implementation provided by Flamary et al. (2021). For each experiment, we run 100 Sinkhorn iterations with 1500 samples for each dataset (*i.e.*, the maximum number of samples that it can generate) and set the regularization parameter ε to its lowest value such that the algorithm is stable. Finally, for sake of fairness with our method, we initialise the barycenter measure with π_r^0 when solving the problem (μ_0 -regWB) for synthetic Gaussian datasets and Bayesian fusion. To run the crWB algorithm (Li et al., 2020), we use the code provided by the authors. We consider the quadratic regularization, which is shown to be empirically more stable than entropic regularization. Following Fan et al. (2020), we choose the potential networks to be fully connected neural networks with 3 hidden layers of shape $(\max(128, 2d), \max(128, 2d), \max(128, 2d))$. The activation functions are ReLu. We optimize the networks with ADAM (Kingma & Ba, 2014) with learning rate 10^{-4} for the subset posterior aggregation task and 10^{-3} for the Gaussian experiment. Finally, we set the batch size to 4,096 and the number of training iterations to 50,000. We highlight that fsWB and crWB solve a regularized Wasserstein barycenter problem, which does not contain an additional *penalization* term on the entropy of the barycenter, contrary to TreeDSB.

G.2 Details on the experiments

Synthetic Gaussian datasets. For each dimension that we consider, we generate three different triplets of random non-diagonal covariance matrices whose condition number is less than 10. We then run the algorithms on each triplet and aggregate the obtained results. The Gaussian datasets contain 1,500 samples for fsWB, and 10,000 samples for crWB and TreeDSB. We run fsWB with the following settings $(d, \varepsilon) \in \{(2, 0.1), (16, 0.2), (64, 0.5), (128, 1.0), (256, 2.0)\}$. We run TreeDSB for 10 mIPF cycles with regularization parameter $\varepsilon = 0.1$, starting from the central node initialized to a Gaussian distribution μ_0 chosen as detailed in Appendix F with $\alpha = 1$. Thus, we solve the regularized Wasserstein barycenter problem (μ_0 -regWB), which contains an additional regularization with respect to μ_0 . This choice is justified, since the non-regularized barycenter is known to be a Gaussian distribution, and μ_0 can be seen as an *a priori* for the regularized barycenter. For each of the three settings, we keep the best result among the 30 mIPF iterations. In this setting, TreeDSB and crWB have roughly the same training time.

Subset posterior aggregation. When considering a dataset splitted into several subdatasets, a common paradigm in bayesian inference consists in running Monte Carlo Markov Chain methods separately on these subdatasets, and then merge the obtained posteriors to recover the full posterior. The barycenter of these subdataset posteriors is proved to be close to the full data posterior under mild assumptions (Srivastava et al., 2018). In our setting, we consider the posterior aggregation problem for the logistic regression model associated to the wine dataset⁷ ($d = 42$) with 3 subdatasets. We consider here two splitting methods: (i) either, data is uniformly splitted between 3 subdatasets with respect to the label distribution, denoted by wine-homogeneous, or (ii) data is splitted with some heterogeneity according to a Dirichlet distribution whose parameter is randomly chosen, denoted by wine-heterogeneous. Following Korotin et al. (2021), we use the stochastic approximation trick so that the subset posterior samples do not vary consistently from the full posterior in covariance (Minsker et al., 2014). We implement the Unadjusted Langevin Algorithm (ULA) to sample from each subdataset posterior and from the full posterior. In each case, we run ULA for $5.5 \cdot 10^6$ iterations with a well chosen step-size, and obtain 9,900 samples after applying a *burn-in* of order 10% and then a *thinning* of size 500. We provide in Figure 7 some metrics which assess the quality of this sampling process. We recall that the the full posterior samples serve as ground truth in this experiment.

The results presented in Table 2 were computed as follows. For fsWB, we first subsample 1,500 samples out of the 9,900 samples from each posterior, and then run the algorithm with $\varepsilon = 0.5$. We repeat three times this procedure and then aggregate the results. In the case of crWB and TreeDSB, we run the algorithms three times with various seeds. Similarly to the Gaussian setting, we run TreeDSB for 10 mIPF cycles with regularization parameter $\varepsilon = 0.1$. We start from the central node with a Gaussian distribution μ_0 chosen as detailed in Appendix F with $\alpha = 1$, and thus solve the barycenter formulation (μ_0 -regWB). For each of the three settings, we keep the best result among the 30 IPF iterations. In this setting, TreeDSB and crWB have roughly the same training time.

⁷<https://archive.ics.uci.edu/ml/datasets/wine>

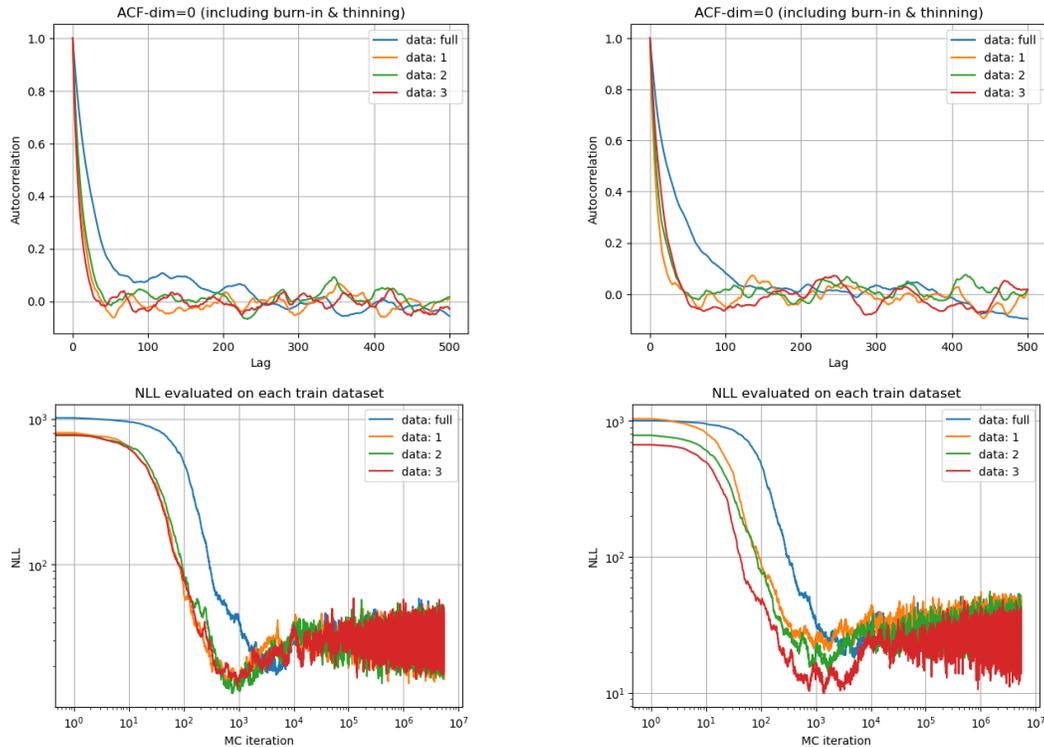


Figure 7: Evaluation of the sampling process for wine-homogeneous (left) and wine-heterogeneous (right). We display the Autocorrelation function on 500 lags (above) and the evolution over the iterations of ULA of the negative log-likelihood (NLL) evaluated on each training dataset (below). In particular, the samples are decorrelated and the NLL has a satisfying profile.

Synthetic two-dimensional datasets. In this setting, we consider three different datasets (*Swiss-roll*, *Circle* and *Moons*) that each contain 10,000 samples. Since we do not have an *a priori* on the shape of the barycenter between these datasets, we consider the regularized Wasserstein barycenter problem (**regWB**), *i.e.*, r is chosen as a leaf and corresponds to one of the input datasets. We emphasize that this experiment is not intended to demonstrate the superiority of TreeDSB to compute 2D Wasserstein barycenters, but is rather meant to illustrate that (a) the marginals of the leaves are well recovered by the algorithm, see Figure 3, and that (b) the obtained barycenter is consistent when diffusing from the different leaves, see Figure 4. In all our experiments on 2D datasets, we observed that (a) was persistently verified without difficulty. In this section, we rather aim at illustrating (b) by providing additional results which assess the quality of the barycenter obtained by TreeDSB with respect to the choice of the starting leaf r and to the choice of the regularization parameter ε .

To do so, we consider three different choices of regularization in TreeDSB: (i) $\varepsilon = 0.2$ (50 mIPF cycles), see Figure 8, (ii) $\varepsilon = 0.1$ (50 mIPF cycles), see Figure 9 and (iii) $\varepsilon = 0.05$ (60 mIPF cycles), see Figure 10. For each of these settings, we run TreeDSB with the starting leaf r chosen as *Swiss-roll* (first row), *Circle* (second row) or *Moons* (third row), and display the final barycenter obtained by diffusing from *Swiss-roll* (first column), *Circle* (second column) and *Moons* (third column). Note that the vertex 0 always corresponds to the starting leaf, the vertex 1 to the barycenter node and that Figure 4 corresponds to the first row of Figure 9.

We can make the following observations. First, the estimated barycenter is always coherent within each row, which assesses the convergence of our method. Then, for each value of ε , the TreeDSB barycenter is rather consistent between the rows, *i.e.*, the choice of the starting leaf does not have a meaningful impact on our method. Finally, as expected, we observe that the support of the barycenter is less and less diffuse as long as ε decreases.

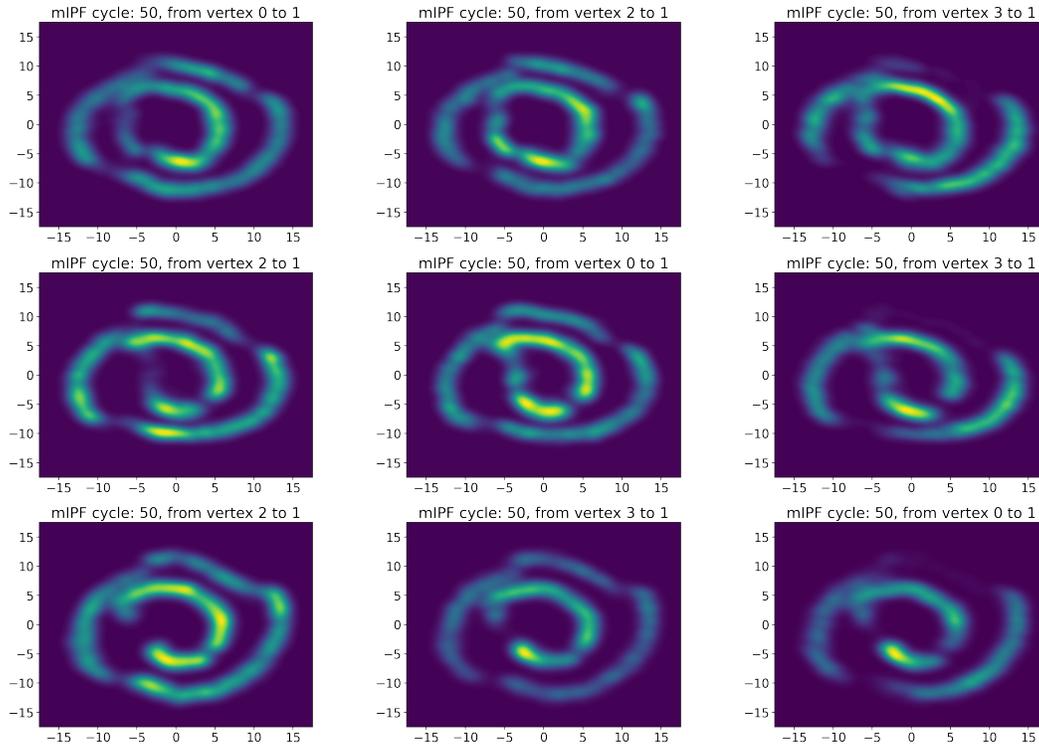


Figure 8: Estimated 2D barycenter obtained by TreeDSB with $\varepsilon = 0.2$ (50 mIPF cycles). First row: starting from *Swiss-roll*. Second row: starting from *Circle*. Third row: starting from *Moons*.

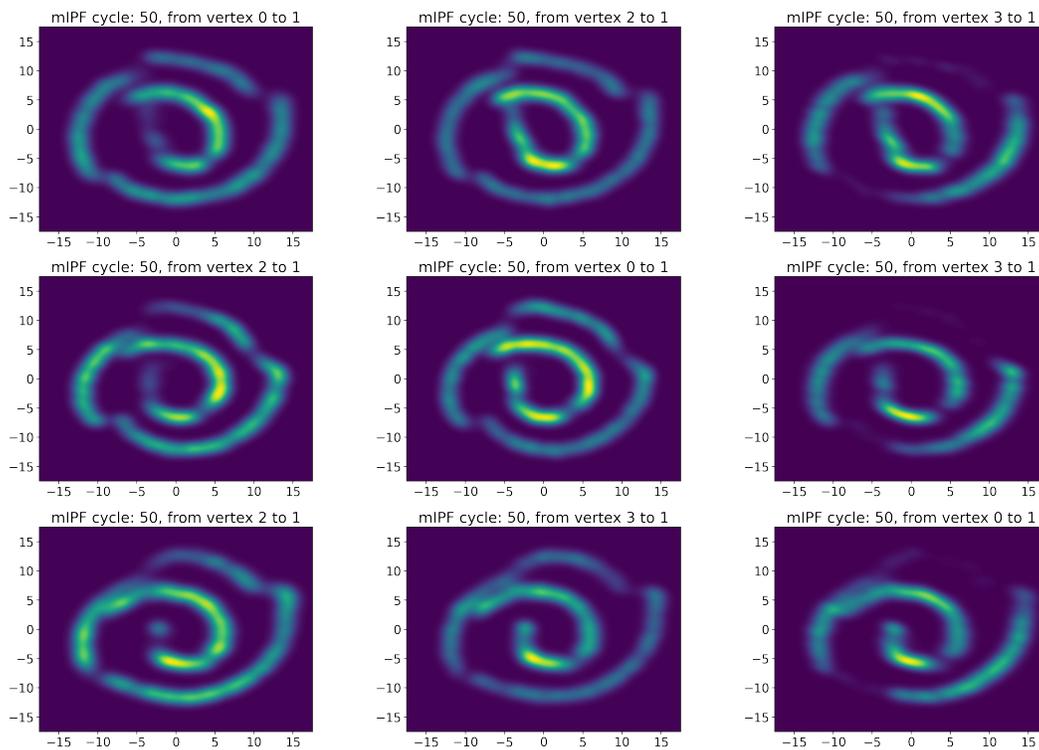


Figure 9: Estimated 2D barycenter obtained by TreeDSB with $\varepsilon = 0.1$ (50 mIPF cycles). First row: starting from *Swiss-roll*. Second row: starting from *Circle*. Third row: starting from *Moons*.

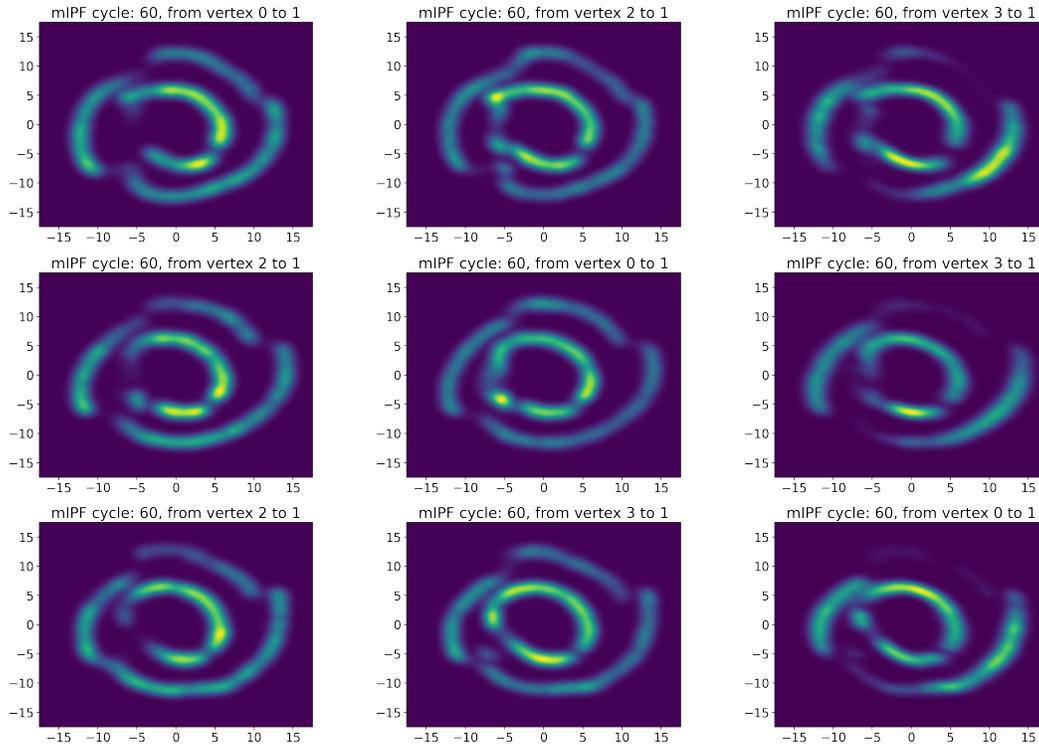


Figure 10: Estimated 2D barycenter obtained by TreeDSB with $\varepsilon = 0.05$ (60 mIPF cycles). First row: starting from *Swiss-roll*. Second row: starting from *Circle*. Third row: starting from *Moons*.

For purpose of illustration, we provide in Figure 11 the barycenter obtained by state-of-the-art two-dimensional *in-sample* methods that are available in POT library (Flamary et al., 2021): (i) non-regularized free-support Wasserstein barycenter (Cuturi & Doucet, 2014), (ii) entropic-regularized free-support Wasserstein barycenter (fsWB) with $\varepsilon = 0.5$ (Cuturi & Doucet, 2014) and (iii) entropic-regularized convolutional Wasserstein barycenter with $\varepsilon = 5.10^{-4}$ (Solomon et al., 2015), which is specifically designed for images. We notably observe that TreeDSB cannot capture the full complexity of the 2D barycenter compared to these methods. We infer that this gap comes from the *dynamic* nature of TreeDSB, since increasing the number of training iterations per IPF iteration or improving the complexity of the neural networks did not bring any significant change in our results. Finally, we recall that the methods (i), (ii) and (iii) do not scale well with dimension, and have to be completely run again when new data samples are available, contrary to TreeDSB.

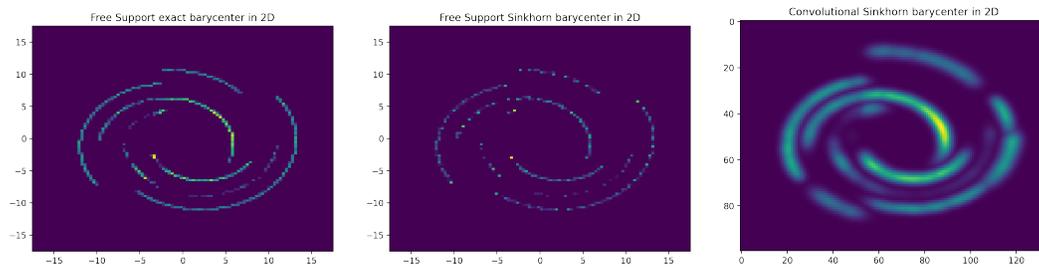


Figure 11: Estimated 2D barycenter obtained by *in-sample* algorithms. From left to right: Cuturi & Doucet (2014) (non-regularized), Cuturi & Doucet (2014) (regularized), Solomon et al. (2015).

MNIST Wasserstein barycenter. This setting can be qualified as *high-dimensional*, since the data dimension is $d = 784$. Here, each digit dataset contains 1,000 samples. As in the two-dimensional setting, we do not have an *a priori* on the shape of the barycenter between MNIST digits, and thus consider the formulation (**regWB**), where the root r is chosen as a leaf. We propose below several experiments to assess the scalability of TreeDSB to this setting.

Digits 0 and 1. In Figure 12, we report the results obtained by running TreeDSB on MNIST digits 0 and 1, for 15 mIPF cycles with $\varepsilon = 0.5$, starting from the leaf MNIST-0. We display 25 samples from the reconstructed MNIST-0 marginal (first column), from the reconstructed MNIST-1 marginal (fourth column), from the estimated barycenter by diffusing from MNIST-0 (second column) and diffusing from MNIST-1 (third column). We notably observe that the digits are well recovered and that the barycenter samples are consistent. We draw the reader’s attention to the fact that TreeDSB showed numerical instability with a regularization value ε lower than 0.5. For purpose of illustration, we display in Figure 13 the Wasserstein barycenter obtained by *non-regularized* methods from Fan et al. (2020) and Korotin et al. (2021), and by the *regularized* approach from Li et al. (2020).

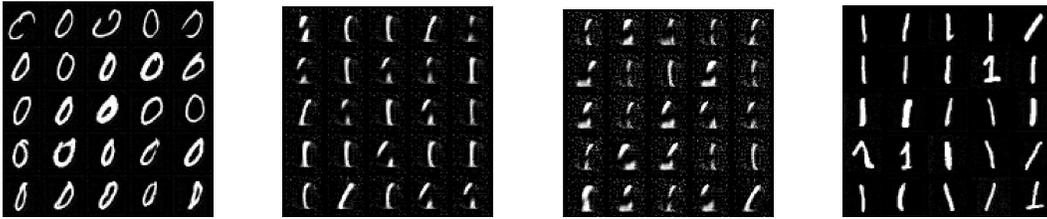


Figure 12: Tree DSB results for MNIST digits 0 and 1, after 15 mIPF cycles with $\varepsilon = 0.5$.

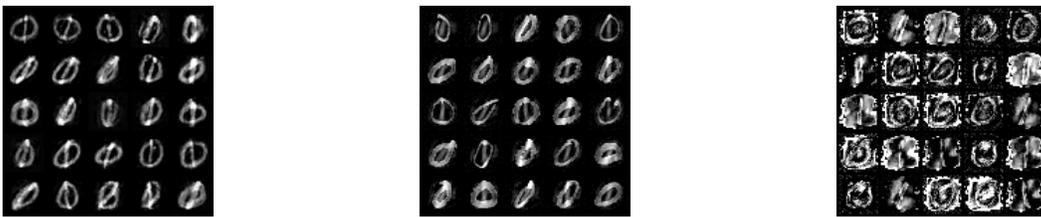


Figure 13: From left to right: Fan et al. (2020), Korotin et al. (2021) and Li et al. (2020).

Digits 2,4 and 6. In Figure 14, we report the results obtained by running TreeDSB on MNIST digits 2,4 and 6, for 10 mIPF cycles with $\varepsilon = 0.5$. Here, we consider three settings which differ by the starting leaf r in the algorithm: MNIST-2 (first row), MNIST-4 (second row), or MNIST-6 (third row). For each of these settings, we display 30 samples from the estimated barycenter by diffusing from MNIST-2 (first column), diffusing from MNIST-4 (third column) and diffusing from MNIST-6 (third column). We notably observe a global consistency of the barycenter samples across the various settings. In Figure 15, we report the results obtained by running TreeDSB on MNIST digits 2,4 and 6, for 10 mIPF cycles with $\varepsilon = 0.2$, starting from MNIST-6. We display 30 samples from the reconstructed marginals (first row), from the estimated barycenter (second row) by diffusing from MNIST-2 (first column), diffusing from MNIST-4 (second column) and diffusing from MNIST-6 (third column). As expected, we observe less noisy barycenter samples compared to Figure 14, while still well recovering MNIST digits.

Digits 0,1 and 4. In Figure 16, we report the results obtained by running TreeDSB on MNIST digits 0,1 and 4, for 10 mIPF cycles with $\varepsilon = 0.5$. We consider two settings which differ by the starting leaf r in the algorithm: MNIST-0 (second row) and MNIST-1 (first/third rows), for which we display samples from the reconstructed measures (first row). In Figure 17, we report the results obtained by running TreeDSB on MNIST digits 0,1 and 4, for 10 mIPF cycles with $\varepsilon = 0.2$. We consider two settings which differ by the starting leaf r in the algorithm: MNIST-0 (first/second row), for which display samples from the reconstructed measures (first row), and MNIST-1 (third row). For all of these settings, we display 30 samples from the estimated barycenter by diffusing from MNIST-0 (first column), diffusing from MNIST-1 (third column) and diffusing from MNIST-4 (third column). Similarly to the digits 2-4-6, we observe consistency within the barycenter samples, unconditionally to the starting leaf, and less noise as ε decreases. Note that the reconstructed MNIST digits are less truthful to the original datasets when ε is low.

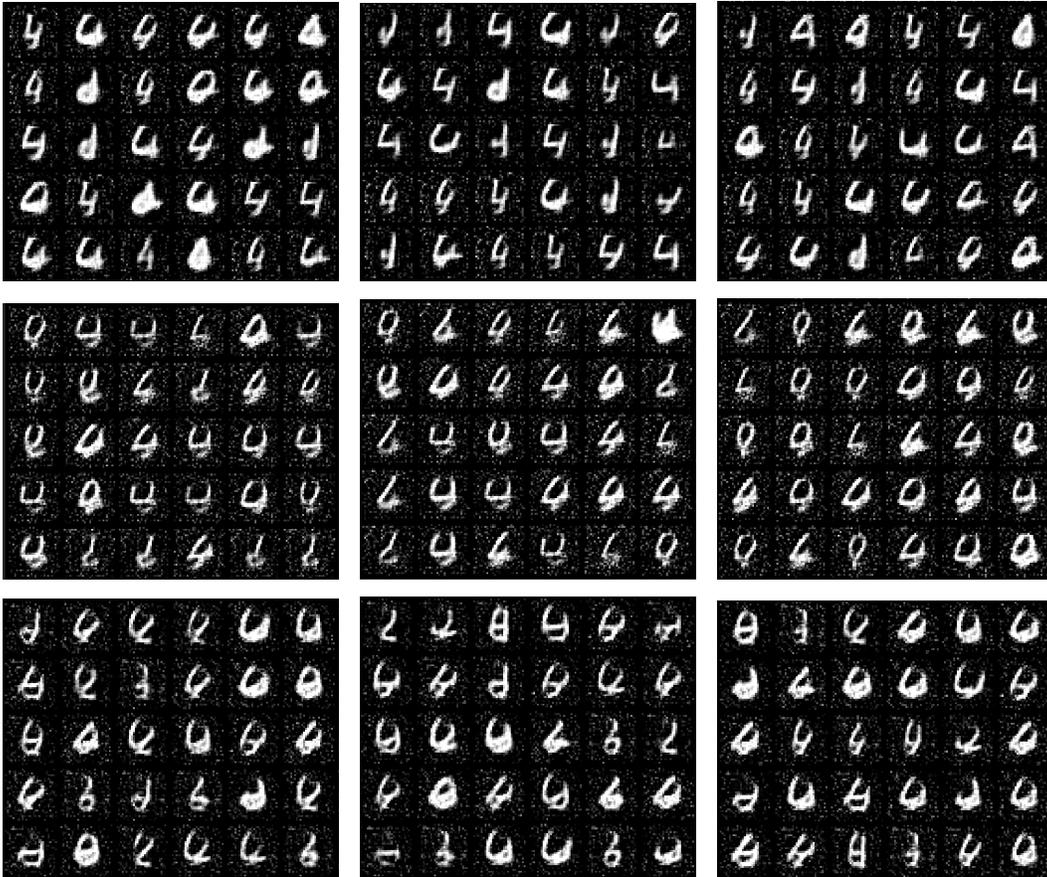


Figure 14: Tree DSB results for MNIST digits 2,4 and 6, after 10 mIPF cycles with $\varepsilon = 0.5$. First row: starting from MNIST-2. Second row: starting from MNIST-4. Third row: starting from MNIST-6.

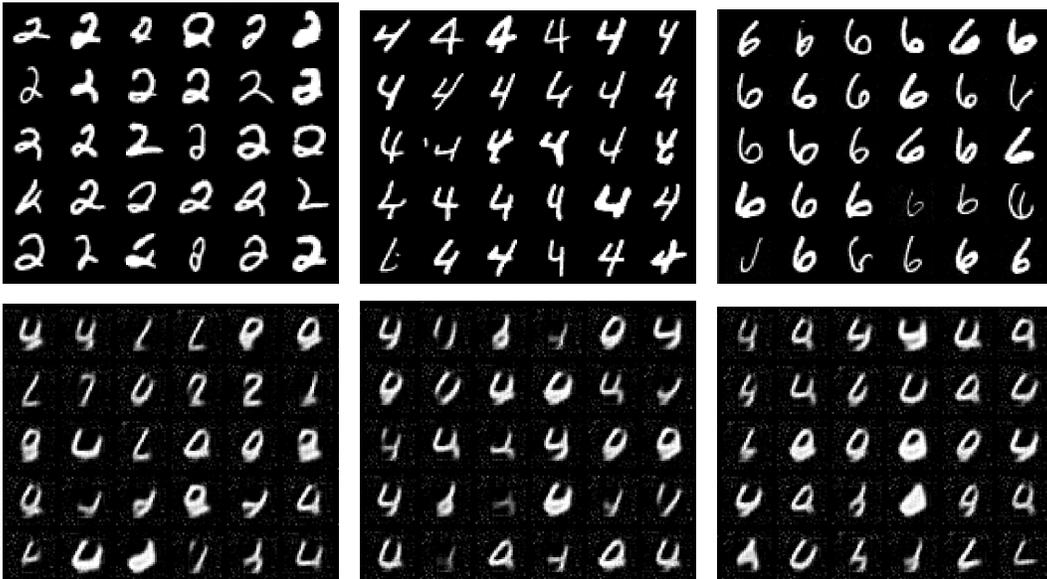


Figure 15: Tree DSB results for MNIST digits 2,4 and 6, after 10 mIPF cycles with $\varepsilon = 0.2$, starting from MNIST-6. First row: samples from the reconstructed marginals. Second row: samples from the estimated barycenter.

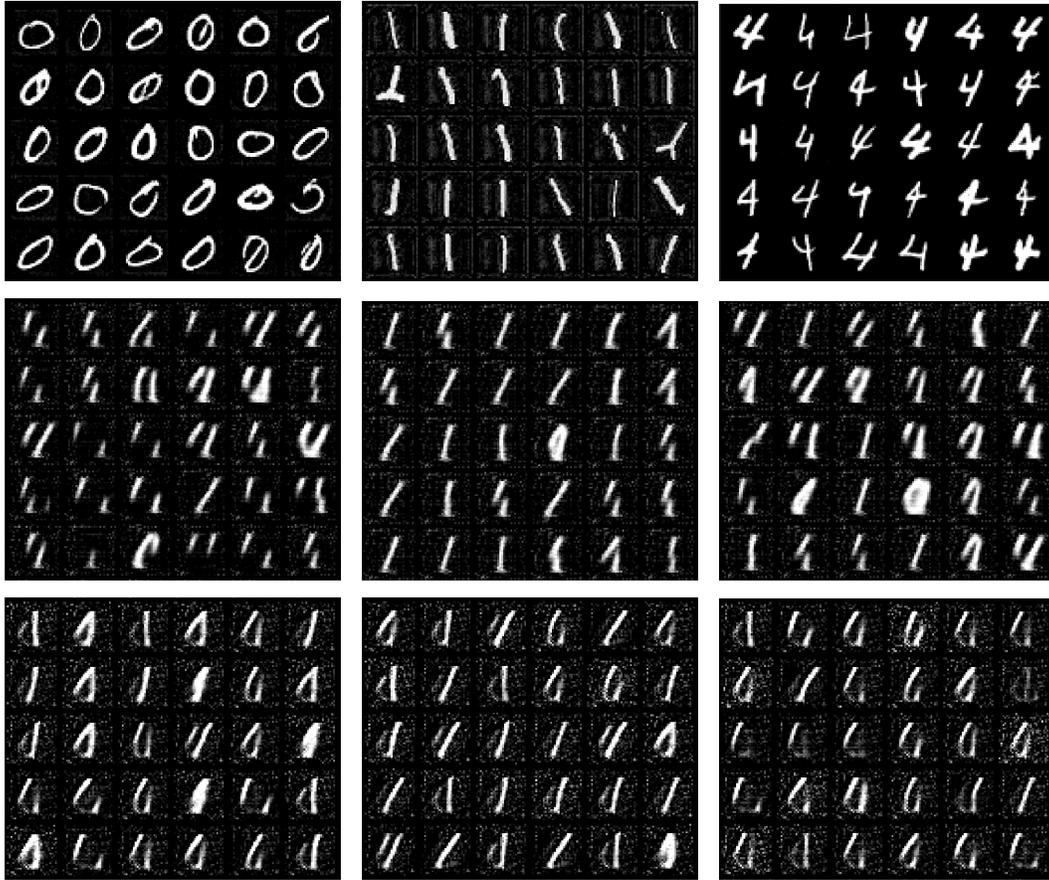


Figure 16: Tree DSB results for MNIST digits 0,1 and 4, after 10 mIPF cycles with $\varepsilon = 0.5$. First row: samples from the reconstructed marginals, starting from MNIST-1. Second row: samples from the estimated barycenter, starting from MNIST-0. Third row: samples from the estimated barycenter, starting from MNIST-1.

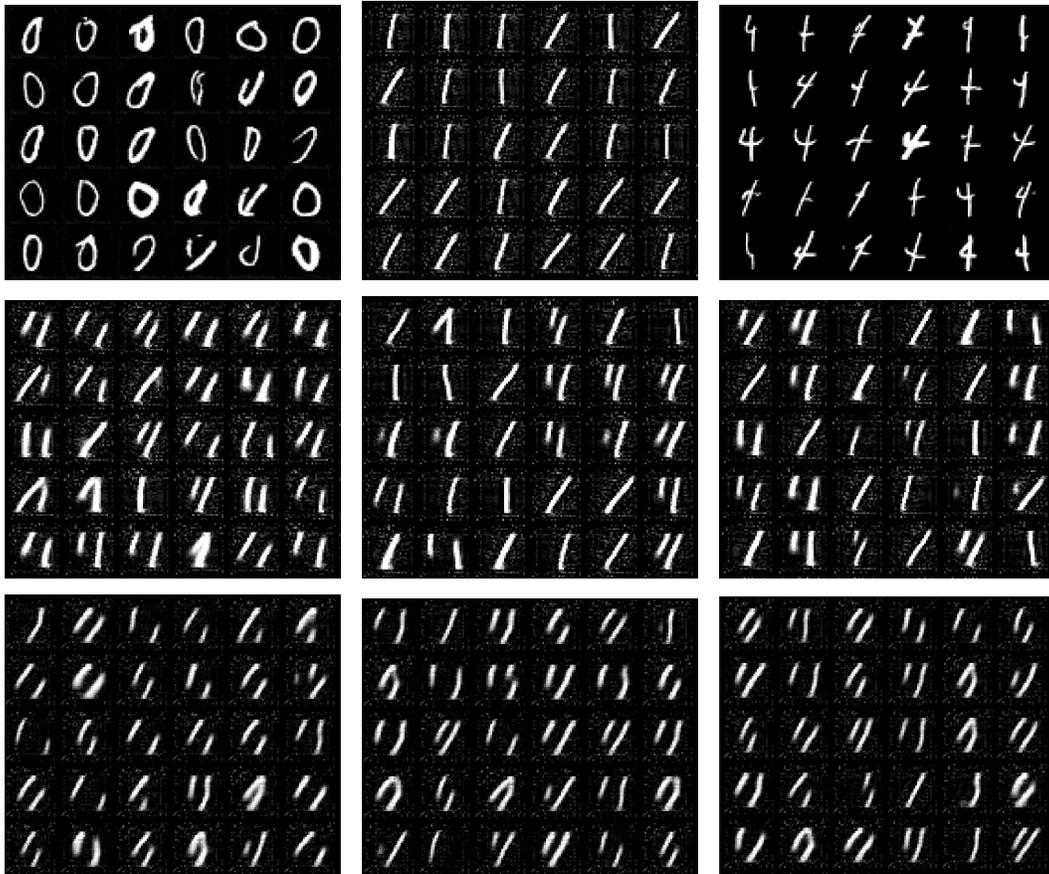


Figure 17: Tree DSB results for MNIST digits 0,1 and 4, after 10 mIPF cycles with $\varepsilon = 0.2$. First row: samples from the reconstructed marginals, starting from MNIST-0. Second row: samples from the estimated barycenter, starting from MNIST-0. Third row: samples from the estimated barycenter, starting from MNIST-1.