
Low-Rank Optimal Transport through Factor Relaxation with Latent Coupling

Peter Halmos^{1,*}, Xinhao Liu^{1,*}, Julian Gold^{2,*}, and Benjamin J. Raphael¹

¹Department of Computer Science, Princeton University

²Center for Statistics and Machine Learning, Princeton University

Abstract

Optimal transport (OT) is a general framework for finding a minimum-cost transport plan, or coupling, between probability distributions, and has many applications in machine learning. A key challenge in applying OT to massive datasets is the quadratic scaling of the coupling matrix with the size of the dataset. Forrow et al. (2019) introduced a factored coupling for the k -Wasserstein barycenter problem, which Scetbon et al. (2021) adapted to solve the primal low-rank OT problem. We derive an alternative parameterization of the low-rank problem based on the *latent coupling* (LC) factorization previously introduced by Lin et al. (2021) generalizing Forrow et al. (2019). The LC factorization has multiple advantages for low-rank OT including decoupling the problem into three OT problems and greater flexibility and interpretability. We leverage these advantages to derive a new algorithm *Factor Relaxation with Latent Coupling* (FRLC), which uses *coordinate* mirror descent to compute the LC factorization. FRLC handles multiple OT objectives (Wasserstein, Gromov-Wasserstein, Fused Gromov-Wasserstein), and marginal constraints (balanced, unbalanced, and semi-relaxed) with linear space complexity. We provide theoretical results on FRLC, and demonstrate superior performance on diverse applications – including graph clustering and spatial transcriptomics – while demonstrating its interpretability.

1 Introduction

Optimal transport (OT) is a powerful geometric framework for comparing probability distributions. OT problems seek a transport plan P efficiently transforming one distribution (a) into another (b), subject to a ground cost C . The minimum cost yields a distance between a and b , while the optimal transport plan reveals key structural similarities between the distributions. Owing to its versatility – different ground costs result in different ways to compare data – OT has found many applications in machine learning and beyond: from self-attention Tay et al. (2020); Sander et al. (2022); Geshkovski et al. (2023) and domain adaptation Courty et al. (2014); Solomon et al. (2015) to computational biology Schiebinger et al. (2019); Yang et al. (2020); Bunne et al. (2023); Liu et al. (2023).

This versatility is compounded by several variants using different forms of the objective function and/or constraints on the transport plan P . Wasserstein (W) OT Kantorovich (1942) compares distributions over the same space through the expected work of P , while Gromov-Wasserstein (GW) OT Mémoli (2011) compares distributions supported on distinct geometries through the expected metric distortion of P . Fused Gromov-Wasserstein (FGW) Vayer et al. (2020) OT is suited to structured data, taking a convex combination of the former two objectives. Independently, one can relax constraints on the *marginals* of P : in computational applications, P is a matrix whose row-sum $P\mathbf{1}_m$ and column-sum $P^T\mathbf{1}_n$ are called its left and right marginals. *Balanced* OT requires $P\mathbf{1}_m = a$ and $P^T\mathbf{1}_n = b$. *Unbalanced* OT Frogner et al. (2015) replaces these constraints with penalties in the

transport cost, and is more robust to outliers. *Semi-relaxed* OT can be used to understand how one dataset embeds into another by imposing one hard constraint on either the left or right marginal, used for feature transfer Dong et al. (2023), and alignment of spatiotemporal data Halmos et al. (2024).

An important consideration in applying OT is the quadratic space of the transport plan. To address both the quadratic complexity and to provide robustness under sampling noise, Forrow et al. (2019) introduced another variant of OT, optimizing a k -Wasserstein Barycenter proxy for the rank-constrained Wasserstein objective. Their approach factors the transport plan through a small set of anchor points called hubs. Generalizing this approach, Scetbon et al. (2021) introduce the factorization $\mathbf{P} = \mathbf{Q} \text{diag}(1/g) \mathbf{R}^T$ comprised of *sub-coupling* matrices \mathbf{Q} and \mathbf{R} sharing an *inner marginal* \mathbf{g} , meaning $\mathbf{Q}^T \mathbf{1}_n = \mathbf{R}^T \mathbf{1}_m = \mathbf{g}$. Building on this, Scetbon et al. (2021, 2022, 2023) derived algorithms to compute low-rank optimal transport plans for the primal OT problem with general costs, extending low-rank OT to GW and unbalanced problems using factored couplings.

Interestingly, a different factorization of \mathbf{P} was proposed by Lin et al. (2021) in the context of k -Wasserstein barycenters. We call their factorization a *latent coupling* (LC) factorization, given by $\mathbf{P} = \mathbf{Q} \text{diag}(1/g_Q) \mathbf{T} \text{diag}(1/g_R) \mathbf{R}^T$, with *two* inner marginals $\mathbf{g}_Q = \mathbf{Q}^T \mathbf{1}_n$ and $\mathbf{g}_R = \mathbf{R}^T \mathbf{1}_m$ and a general coupling \mathbf{T} . Lin et al. (2021) constrain the transport between \mathbf{a} and \mathbf{b} through two sets of learned anchor points, where the factorization is defined by three transport plans computed from three cost matrices between the points and their anchors. This objective differs from that of Forrow et al. (2019); Scetbon et al. (2021), who seek a minimal rank coupling with respect to a single, fixed cost \mathbf{C} . We observe that factored couplings of Forrow et al. (2019) correspond to LC factorizations with diagonal \mathbf{T} , suggesting the LC factorization of Lin et al. (2021) may provide an alternative parameterization of transport plans for the low-rank OT problem considered in Forrow et al. (2019); Scetbon et al. (2021). To our knowledge, this idea has not yet been explored.

Contributions. We present a new algorithm, Factor Relaxation with Latent Coupling (FRLC, with the informal mnemonic “frolic”), to compute a minimum cost low-rank transport plan using the LC factorization. Parameterizing low-rank transport plans with the LC factorization has a number of advantages. First, optimization of the low-rank OT objective decouples into three OT sub-problems on the LC factors \mathbf{Q} , \mathbf{R} , \mathbf{T} , leading to a simpler optimization algorithm. Second, this decoupling provides straightforward extensions of FRLC to low-rank unbalanced and semi-relaxed OT; similar extensions for factored couplings required additional work Scetbon et al. (2023) beyond the balanced case. Third, the latent coupling \mathbf{T} in the LC factorization provides additional flexibility to model transport between datasets with different numbers of clusters, and to model mass-splitting between these clusters, providing a high-level and interpretable description of \mathbf{P} that differs from the factored couplings of Forrow et al. (2019). FRLC computes the LC factorization using a novel *coordinate* mirror descent scheme, alternating descent steps on variables (\mathbf{Q}, \mathbf{R}) and \mathbf{T} , inspired by the mirror descent approach of Scetbon et al. (2021). We call the descent step on (\mathbf{Q}, \mathbf{R}) *factor relaxation*, as the factors \mathbf{Q} and \mathbf{R} have relaxed inner marginals, allowing FRLC to be solved by OT sub-problems. FRLC handles multiple OT objectives (Wasserstein, Gromov-Wasserstein, Fused Gromov-Wasserstein), and marginal constraints (balanced, unbalanced, and semi-relaxed). We show FRLC performs better than existing state-of-the-art low-rank methods on a range of synthetic and real datasets, retaining the interpretability of Lin et al. (2021), and inheriting the broad applicability of Scetbon et al. (2021); Scetbon & Cuturi (2022); Scetbon et al. (2022, 2023).

2 Background

Wasserstein OT. Let $\{x_1, \dots, x_n\}$ and $\{y_1, \dots, y_m\}$ be datasets in a metric space \mathcal{X} , and let Δ_d be the probability simplex of size d . Through probability vectors $\mathbf{a} \in \Delta_n$ and $\mathbf{b} \in \Delta_m$, each dataset is encoded as a probability measure: $\mu = \sum_{i=1}^n \mathbf{a}_i \delta_{x_i}$ and $\nu = \sum_{j=1}^m \mathbf{b}_j \delta_{y_j}$. Let

$$\Pi_{\mathbf{a}, \cdot} := \{\mathbf{P} \in \mathbb{R}_+^{n \times m} : \mathbf{P} \mathbf{1}_m = \mathbf{a}\}, \quad \Pi_{\cdot, \mathbf{b}} := \{\mathbf{P} \in \mathbb{R}_+^{n \times m} : \mathbf{P}^T \mathbf{1}_n = \mathbf{b}\}, \quad \Pi_{\mathbf{a}, \mathbf{b}} := \Pi_{\mathbf{a}, \cdot} \cap \Pi_{\cdot, \mathbf{b}}.$$

Thus, $\Pi_{\mathbf{a}, \mathbf{b}}$ is the set of transport plans (probabilistic coupling matrices) with marginals \mathbf{a} and \mathbf{b} . Given a cost function $c : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+$, define the cost matrix $\mathbf{C} \in \mathbb{R}_+^{n \times m}$ via $\mathbf{C}_{ij} = c(x_i, y_j)$. The Kantorovich formulation Kantorovich (1942) of discrete OT, also called the Wasserstein problem, seeks a transport plan \mathbf{P} of minimal cost :

$$W(\mu, \nu) := \min_{\mathbf{P} \in \Pi_{\mathbf{a}, \mathbf{b}}} \langle \mathbf{C}, \mathbf{P} \rangle_F. \quad (1)$$

Gromov-Wasserstein OT. In many applications, one wishes to compare datasets $\{x_1, \dots, x_n\} \subset \mathcal{X}$ and $\{y_1, \dots, y_m\} \subset \mathcal{Y}$ across distinct metric spaces \mathcal{X} and \mathcal{Y} . The Gromov-Wasserstein (GW) objective Mémoli (2007, 2011) addresses the absence of a common metric or coordinate system through intra-domain cost functions $c_1 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+$ and $c_2 : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$, leading to intra-domain cost matrices $\mathbf{A}_{ik} = c_1(x_i, x_k)$ and $\mathbf{B}_{jl} = c_2(y_j, y_l)$. The GW objective function $\mathcal{Q}_{\mathbf{A}, \mathbf{B}}(\mathbf{P}) := \sum_{i,j,k,l} (\mathbf{A}_{ik} - \mathbf{B}_{jl})^2 \mathbf{P}_{ij} \mathbf{P}_{kl}$ quantifies the expected metric distortion under \mathbf{P} , leading to the optimization problem:

$$\text{GW}(\mu, \nu) := \min_{\mathbf{P} \in \Pi_{\mathbf{a}, \mathbf{b}}} \mathcal{Q}_{\mathbf{A}, \mathbf{B}}(\mathbf{P}). \quad (2)$$

The Fused Gromov-Wasserstein (FGW) objective function Vayer et al. (2020) is a convex combination of the W and GW objectives, given as $\alpha \langle \mathbf{C}, \mathbf{P} \rangle_F + (1 - \alpha) \mathcal{Q}_{\mathbf{A}, \mathbf{B}}(\mathbf{P})$, for hyperparameter $\alpha \in (0, 1)$.

Relaxed marginal constraints. *Balanced* OT (1) constrains \mathbf{P} to lie in $\Pi_{\mathbf{a}, \mathbf{b}}$. *Unbalanced* OT relaxes constraints $\mathbf{P}\mathbf{1}_m = \mathbf{a}$ and $\mathbf{P}^T\mathbf{1}_n = \mathbf{b}$, replacing them with penalties in the form of KL divergences (or other divergences, see Chizat et al. (2018)):

$$\text{U-W}(\mu, \nu) := \min_{\mathbf{P} \in \mathbb{R}_+^{n \times m}} \langle \mathbf{C}, \mathbf{P} \rangle_F + \tau_L \text{KL}(\mathbf{P}\mathbf{1}_m \| \mathbf{a}) + \tau_R \text{KL}(\mathbf{P}^T\mathbf{1}_n \| \mathbf{b}), \quad (3)$$

where $\tau_L, \tau_R > 0$ control the strength of each penalty. *Semi-relaxed* optimal transport relaxes exactly one of the hard constraints $\mathbf{P}\mathbf{1}_m = \mathbf{a}$ and $\mathbf{P}^T\mathbf{1}_n = \mathbf{b}$ in the same manner. The semi-relaxed version of (1) obtained by relaxing only the ‘‘right’’ marginal constraint on \mathbf{b} is:

$$\text{SR}^R\text{-W}(\mu, \nu) := \min_{\mathbf{P} \in \Pi_{\mathbf{a}, \cdot}} \langle \mathbf{C}, \mathbf{P} \rangle_F + \tau \text{KL}(\mathbf{P}^T\mathbf{1}_n \| \mathbf{b}), \quad (4)$$

while its ‘‘left’’ marginal counterpart $\text{SR}^L\text{-W}(\mu, \nu)$ is defined analogously over $\mathbf{P} \in \Pi_{\cdot, \mathbf{b}}$, using penalty $\tau \text{KL}(\mathbf{P}\mathbf{1}_m \| \mathbf{a})$. Likewise, one can form semi-relaxed or unbalanced GW and FGW problems.

Entropy regularization. The seminal work Cuturi (2013b) introduced the Sinkhorn algorithm to solve an entropy regularized version of (1), $W_\epsilon(\mu, \nu) := \min_{\mathbf{P} \in \Pi_{\mathbf{a}, \mathbf{b}}} \langle \mathbf{C}, \mathbf{P} \rangle_F - \epsilon H(\mathbf{P})$, massively improving the $O(n^3 \log n)$ time complexity of classical techniques Orlin (1997); Tarjan (1997). Above, H is the entropy, $H(\mathbf{P}) = -\sum_{ij} \mathbf{P}_{ij} (\log \mathbf{P}_{ij} - 1)$, and $\epsilon > 0$ is the regularization strength.

Low-rank regularization. The nonnegative rank $\text{rk}_+(M)$ of matrix M is the least number of nonnegative rank-one matrices summing to M . For $r \geq 1$, define

$$\Pi_{\mathbf{a}, \cdot}(r) = \{\mathbf{P} \in \Pi_{\mathbf{a}, \cdot} : \text{rk}_+(\mathbf{P}) \leq r\}, \quad \Pi_{\cdot, \mathbf{b}}(r) = \{\mathbf{P} \in \Pi_{\cdot, \mathbf{b}} : \text{rk}_+(\mathbf{P}) \leq r\}, \quad (5)$$

and let $\Pi_{\mathbf{a}, \mathbf{b}}(r) = \Pi_{\mathbf{a}, \cdot}(r) \cap \Pi_{\cdot, \mathbf{b}}(r)$. To estimate Wasserstein distances with greater stability and accuracy under sampling noise, Forrow et al. (2019) proposed a low-rank regularization on the coupling matrix, factoring the transport through a small set of anchor points. More explicitly, Scetbon et al. (2021) parameterized the set as $\Pi_{\mathbf{a}, \mathbf{b}}(r)$ through the set $\text{FC}_{\mathbf{a}, \mathbf{b}}(r)$ of *factored couplings*,

$$\text{FC}_{\mathbf{a}, \mathbf{b}}(r) := \{(\mathbf{Q}, \mathbf{R}, \mathbf{g}) \in \mathbb{R}_+^{n \times r} \times \mathbb{R}_+^{m \times r} \times (\mathbb{R}_+^*)^r : \mathbf{Q} \in \Pi_{\mathbf{a}, \mathbf{g}}, \mathbf{R} \in \Pi_{\mathbf{b}, \mathbf{g}}\}.$$

The set $\text{FC}_{\mathbf{a}, \mathbf{b}}(r)$ parameterizes $\Pi_{\mathbf{a}, \mathbf{b}}(r)$ through $(\mathbf{Q}, \mathbf{R}, \mathbf{g}) \mapsto \mathbf{Q} \text{diag}(1/\mathbf{g}) \mathbf{R}^T$, as shown by Cohen & Rothblum (1993).

Scetbon et al. (2021) apply this factorization to solve the Wasserstein problem subject to $\mathbf{P} \in \Pi_{\mathbf{a}, \mathbf{b}}(r)$ for general cost matrices:

$$W_r(\mu, \nu) := \min_{\mathbf{P} \in \Pi_{\mathbf{a}, \mathbf{b}}(r)} \langle \mathbf{C}, \mathbf{P} \rangle_F \quad (6)$$

GW, unbalanced and semi-relaxed low-rank OT problems are defined as in (2), (3) and (4), replacing $\mathbb{R}_+^{n \times m}$, $\Pi_{\mathbf{a}, \cdot}$, or $\Pi_{\cdot, \mathbf{b}}$ with rank-constrained counterparts (5). Scetbon & Cuturi (2022); Scetbon et al. (2022, 2023) developed a robust framework for solving all of these problems.

3 Factor Relaxation with Latent Coupling (FRLC) algorithm

3.1 Latent Coupling Factorization

We parameterize low-rank coupling matrices $\mathbf{P} \in \Pi_{\mathbf{a}, \mathbf{b}}(r)$ using a factorization introduced in Lin et al. (2021), which we call the *latent coupling (LC) factorization* (Fig. 1). The key property of

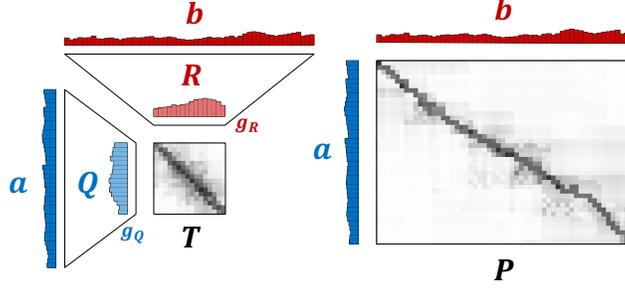


Figure 1: (Left) The LC factorization $P = Q\text{diag}(1/g_Q)T\text{diag}(1/g_R)R^T$ of coupling matrix P with outer marginals a, b , inner marginals g_Q, g_R , factors Q, R , and latent coupling T . (Right) Full-rank coupling matrix P .

this factorization is the presence of a coupling matrix T linking two distinct inner marginals. For simplicity we describe this factorization using an r -dimensional latent space, but we also extend to non-square matrices linking two latent spaces of different dimensions, as demonstrated in the results.

Definition 3.1 (Inner marginals). Given a factorization $P = QR^T$ of a coupling matrix $P \in \Pi_{a,b}(r)$, the *inner marginals* of Q and R are $g_Q := Q^T \mathbf{1}_n$ and $g_R := R^T \mathbf{1}_m$, respectively, where $g_Q, g_R \in \Delta_r$.

To distinguish the different marginals, we refer to a and b as *outer marginals*.

Definition 3.2 (LC factorization). Given a coupling matrix $P \in \Pi_{a,b}(r)$, a *latent coupling (LC) factorization* of P is $P = Q\text{diag}(1/g_Q)T\text{diag}(1/g_R)R^T$, where g_Q and g_R are the inner marginals of Q and R , $Q \in \Pi_{a,\cdot}$, $R \in \Pi_{\cdot,b}$, and $T \in \Pi_{g_Q, g_R}$.

We call the factors Q, R, T in an LC factorization *sub-couplings*. Let $\mathcal{R}_+ := \mathbb{R}_+^{n \times r} \times \mathbb{R}_+^{m \times r} \times \mathbb{R}_+^{r \times r}$. Given probability vectors $a \in \Delta_n, b \in \Delta_m$ and a positive integer rank r , let

$$\text{LC}_{a,b}(r) := \{(Q, R, T) \in \mathcal{R}_+ : Q \in \Pi_{a,\cdot}, R \in \Pi_{\cdot,b}, T \in \Pi_{g_Q, g_R}\},$$

be the set of admissible sub-couplings for the LC factorization. Definition 3.2 gives the following map from $\text{LC}_{a,b}(r)$ to $\Pi_{a,b}(r)$:

$$(Q, R, T) \mapsto Q\text{diag}(1/g_Q)T\text{diag}(1/g_R)R^T =: P_{(Q,R,T)}. \quad (7)$$

Since this map is surjective, the set $\text{LC}_{a,b}(r)$ parameterizes $\Pi_{a,b}(r)$. Surjectivity follows from the fact that $\text{FC}_{a,b}(r)$ maps injectively into $\text{LC}_{a,b}(r)$, through $(Q, R, g) \mapsto (Q, R, \text{diag}(g))$, and $\text{FC}_{a,b}(r)$ maps surjectively onto $\Pi_{a,b}(r)$ via $(Q, R, g) \mapsto Q\text{diag}(1/g)R^T$. Definition 3.1 and Definition 3.2 are readily extended in two directions: the case when the outer marginal constraints are relaxed such that $Q \in \mathbb{R}_+^{n \times m}$ or $R \in \mathbb{R}_+^{n \times m}$, while maintaining the constraint that $T \in \Pi_{g_Q, g_R}$; as well as the case of non-square T .

3.2 The Balanced FRLC Algorithm

We introduce an algorithm Factor Relaxation with Latent Coupling (FRLC), to compute a LC factorization of minimum cost. We first describe the FRLC algorithm for the balanced Wasserstein problem. Extensions to other and marginal constraints are discussed later. The FRLC objective function, for low-rank, balanced Wasserstein OT, is

$$\mathcal{L}_{\text{LC}}(Q, R, T) := \langle C, P_{(Q,R,T)} \rangle_F, \quad (8)$$

where $P_{(Q,R,T)}$ is defined by (7). Since $\text{LC}_{a,b}(r)$ parameterizes $\Pi_{a,b}(r)$, problem (8) is equivalent to low rank problem (6). The FRLC algorithm is built from projections onto convex sets, described by constraints on the outer marginals alone for (Q, R) and by the inner marginals alone for T . Given $(Q, R, T) \in \text{LC}_{a,b}(r)$, sub-couplings Q and R are constrained by:

$$\mathcal{C}_1(a) := \{(Q, R, T) \in \mathcal{R}_+ : Q\mathbf{1}_r = a\}, \quad \mathcal{C}_1(b) := \{(Q, R, T) \in \mathcal{R}_+ : R\mathbf{1}_r = b\}.$$

The convex sets constraining the latent coupling matrix T are

$$\mathcal{C}_2(\mathbf{g}_Q) := \{(\mathbf{Q}, \mathbf{R}, \mathbf{T}) \in \mathcal{R}_+ : \mathbf{T}\mathbf{1}_r = \mathbf{g}_Q\}, \quad \mathcal{C}_2(\mathbf{g}_R) := \{(\mathbf{Q}, \mathbf{R}, \mathbf{T}) \in \mathcal{R}_+ : \mathbf{T}^\top \mathbf{1}_r = \mathbf{g}_R\},$$

where $\mathbf{g}_Q = \mathbf{Q}^\top \mathbf{1}_n$ and $\mathbf{g}_R = \mathbf{R}^\top \mathbf{1}_m$ as per Definition 3.1. Writing $\mathcal{C}_1 = \mathcal{C}_1(\mathbf{a}) \cap \mathcal{C}_1(\mathbf{b})$ and $\mathcal{C}_2 = \mathcal{C}_2(\mathbf{g}_Q) \cap \mathcal{C}_2(\mathbf{g}_R)$, one has $\text{LC}_{\mathbf{a}, \mathbf{b}}(r) = \mathcal{C}_1 \cap \mathcal{C}_2$.

We use *coordinate* mirror descent to optimize (8), building on the mirror descent (MD) approach of Scetbon et al. (2021); Scetbon & Cuturi (2022); Scetbon et al. (2022, 2023) for the low-rank problem. First we take a descent step in the variables (\mathbf{Q}, \mathbf{R}) for a fixed \mathbf{T} , using KL penalties on their inner marginals. These ‘‘soft’’ constraints allow the joint optimization in (\mathbf{Q}, \mathbf{R}) to decouple into two semi-relaxed OT problems, one for each variable. We call this step *factor relaxation* as this allows (\mathbf{Q}, \mathbf{R}) to have relaxed inner marginals \mathbf{g}_Q and \mathbf{g}_R . Next we take a descent step in the latent coupling variable \mathbf{T} , fixing the \mathbf{Q} and \mathbf{R} , equivalent to solving a balanced OT problem. Thus, solving both coordinate descent steps corresponds to solving three OT problems.

We now provide further details on these coordinate descent steps, with the full algorithm given in Algorithm 1. Let $(\gamma_k)_{k=1}^N$ be a sequence of step sizes. As in Scetbon & Cuturi (2022), we choose ℓ^∞ -normalization for the step-sizes. Our coordinate mirror descent in the factor relaxation step is:

$$\begin{aligned} (\mathbf{Q}_{k+1}, \mathbf{R}_{k+1}) \leftarrow \arg \min_{(\mathbf{Q}, \mathbf{R}) : (\mathbf{Q}, \mathbf{R}, \mathbf{T}_k) \in \mathcal{C}_1} & \langle (\mathbf{Q}, \mathbf{R}), \nabla_{(\mathbf{Q}, \mathbf{R})} \mathcal{L}_{\text{LC}} \rangle + \frac{1}{\gamma_k} \text{KL}((\mathbf{Q}, \mathbf{R}) \| (\mathbf{Q}_k, \mathbf{R}_k)) \\ & + \tau \text{KL}((\mathbf{Q}^\top \mathbf{1}_n, \mathbf{R}^\top \mathbf{1}_m) \| (\mathbf{Q}_k^\top \mathbf{1}_n, \mathbf{R}_k^\top \mathbf{1}_m)) \end{aligned}$$

The Sinkhorn kernels for the semi-relaxed OT problems arising from the factor relaxation step are:

$$\begin{aligned} \mathbf{K}_Q^{(k)} &:= \mathbf{Q}_k \odot \exp(-\gamma_k (\mathbf{C} \mathbf{R}_k \mathbf{X}_k^\top - \mathbf{1}_n \text{diag}^{-1}((\mathbf{C} \mathbf{R}_k \mathbf{X}_k^\top)^\top \mathbf{Q}_k \text{diag}(1/\mathbf{g}_{Q_k})))^\top) \\ \mathbf{K}_R^{(k)} &:= \mathbf{R}_k \odot \exp(-\gamma_k (\mathbf{C}^\top \mathbf{Q}_k \mathbf{X}_k - \mathbf{1}_m \text{diag}^{-1}(\text{diag}(1/\mathbf{g}_{R_k}) \mathbf{R}_k^\top \mathbf{C}^\top \mathbf{Q}_k \mathbf{X}_k)^\top)), \end{aligned}$$

introducing the shorthand $\mathbf{X} = \text{diag}(1/\mathbf{g}_Q) \mathbf{T} \text{diag}(1/\mathbf{g}_R)$ and where $\text{diag}^{-1}(\cdot) : \mathbb{R}^{r \times r} \rightarrow \mathbb{R}^r$ denotes the matrix-to-vector extraction of the diagonal. This τ -dependent regularization also allows us to show smoothness of the objective in Proposition E.5, from which the convergence guarantee Proposition 3.3 follows. We derive the semi-relaxed projection Algorithm 2 of the sub-couplings \mathbf{Q} and \mathbf{R} in Appendix G for completeness. We also show in Lemma A.1 that \mathbf{g}_Q and \mathbf{g}_R induced by the semi-relaxed projection are both feasible and locally optimal, not requiring separate optimization.

As $(\mathbf{Q}_{k+1}, \mathbf{R}_{k+1}, \mathbf{T}) \in \mathcal{C}_2$ if and only if $\mathbf{T} \in \Pi_{\mathbf{g}_{Q_{k+1}}, \mathbf{g}_{R_{k+1}}}$, after the factor relaxation step, we next take a coordinate MD step on the latent coupling \mathbf{T} :

$$\mathbf{T}_{k+1} \leftarrow \arg \min_{\mathbf{T} : (\mathbf{Q}_{k+1}, \mathbf{R}_{k+1}, \mathbf{T}) \in \mathcal{C}_2} \langle \mathbf{T}, \nabla_{\mathbf{T}} \mathcal{L}_{\text{LC}} \rangle + \frac{1}{\gamma_k} \text{KL}(\mathbf{T} \| \mathbf{T}_k). \quad (9)$$

This is equivalent to applying Sinkhorn (Algorithm 5) to \mathbf{T} given \mathbf{g}_Q and \mathbf{g}_R with the kernel:

$$\mathbf{K}_T^{(k)} := \mathbf{T}_k \odot \exp(-\gamma_k \text{diag}(1/\mathbf{g}_{Q_{k+1}}) \mathbf{Q}_{k+1}^\top \mathbf{C} \mathbf{R}_{k+1} \text{diag}(1/\mathbf{g}_{R_{k+1}})).$$

After the final iteration of the coordinate-MD scheme, $\mathbf{X} = \text{diag}(1/\mathbf{g}_Q) \mathbf{T} \text{diag}(1/\mathbf{g}_R)$ satisfies $\mathbf{X} \mathbf{g}_R = \mathbf{1}_r$ and $\mathbf{X}^\top \mathbf{g}_Q = \mathbf{1}_r$ as \mathbf{T} is a coupling between \mathbf{g}_Q and \mathbf{g}_R . Thus $\mathbf{P}_r = \mathbf{Q} \mathbf{X} \mathbf{R}^\top \in \Pi_{\mathbf{a}, \mathbf{b}}$ and the iterates $(\mathbf{Q}_k, \mathbf{T}_k, \mathbf{R}_k)$ remain in the intersection of the constraint sets. Thus, in contrast to other approaches Scetbon et al. (2021); Lin et al. (2021); Forrow et al. (2019), we do not require Dykstra projections back into the intersection to maintain feasibility. We note that our implementation of FRLC allows for a non-square latent coupling \mathbf{T} , providing greater interpretability in problem-specific applications. Above, we presented FRLC in the simplest case that \mathbf{T} is square.

3.3 Initialization, convergence, and FRLC extensions

Full-rank random initializations of the sub-coupling matrices. We propose a new initialization of the sub-couplings $(\mathbf{Q}, \mathbf{R}, \mathbf{T})$ for the LC-factorization in Algorithm 6. This generates a full-rank initialization (Proposition F.1) in the set of rank- r couplings $\Pi_{\mathbf{a}, \mathbf{b}}(r)$ and is accomplished by applying Sinkhorn to random matrices. Our approach differs from Scetbon et al. (2021); Scetbon & Cuturi (2022) who use initializations for the diagonal factorization of Forrow et al. (2019), and are not applicable to a latent coupling that is non-diagonal, non-square, or with two distinct inner marginals.

Algorithm 1 Balanced FRLC

Input $\mathbf{C}, r, \mathbf{a}, \mathbf{b}, \tau, \gamma, \delta, \varepsilon$
 Initialize $\mathbf{g}_Q, \mathbf{g}_R = \frac{1}{r} \mathbf{1}_r$
 $\mathbf{Q}_0, \mathbf{R}_0, \mathbf{T}_0 \leftarrow \text{Initialize-Couplings}(\mathbf{a}, \mathbf{b}, \mathbf{g}_Q, \mathbf{g}_R)$ # Alg. 6
 $\mathbf{X}_0 \leftarrow \text{diag}(1/\mathbf{Q}_0^\top \mathbf{1}_n) \mathbf{T}_0 \text{diag}(1/\mathbf{R}_0^\top \mathbf{1}_m)$
while $\Delta((\mathbf{Q}_k, \mathbf{R}_k, \mathbf{T}_k), (\mathbf{Q}_{k-1}, \mathbf{R}_{k-1}, \mathbf{T}_{k-1})) > \varepsilon$ **do** # Δ as in (10)
 $\nabla_Q \leftarrow \mathbf{C} \mathbf{R}_k \mathbf{X}_k^\top - \mathbf{1}_n \text{diag}^{-1}((\mathbf{C} \mathbf{R}_k \mathbf{X}_k^\top)^\top \mathbf{Q}_k \text{diag}(1/\mathbf{g}_Q))^\top$
 $\nabla_R \leftarrow \mathbf{C}^\top \mathbf{Q}_k \mathbf{X}_k - \mathbf{1}_m \text{diag}^{-1}(\text{diag}(1/\mathbf{g}_R) \mathbf{R}_k^\top \mathbf{C}^\top \mathbf{Q}_k \mathbf{X}_k)^\top$
 $\gamma_k \leftarrow \gamma / \max\{\|\nabla_Q\|_\infty, \|\nabla_R\|_\infty\}$ # ℓ^∞ -normalization of Scetbon & Cuturi (2022)
 $\mathbf{K}_Q^{(k)}, \mathbf{K}_R^{(k)} \leftarrow \mathbf{Q}_k \odot \exp(-\gamma_k \nabla_Q), \mathbf{R}_k \odot \exp(-\gamma_k \nabla_R)$
 $\mathbf{Q}_k \leftarrow \text{SR}^R\text{-projection}(\mathbf{K}_Q^{(k)}, \gamma_k, \tau, \mathbf{a}, \mathbf{Q}_{k-1}^\top \mathbf{1}_n, \delta)$ # Semi-relaxed OT, Alg. 2
 $\mathbf{R}_k \leftarrow \text{SR}^R\text{-projection}(\mathbf{K}_R^{(k)}, \gamma_k, \tau, \mathbf{b}, \mathbf{R}_{k-1}^\top \mathbf{1}_m, \delta)$ # Semi-relaxed OT
 $\mathbf{g}_Q, \mathbf{g}_R \leftarrow \mathbf{Q}_k^\top \mathbf{1}_n, \mathbf{R}_k^\top \mathbf{1}_m$
 $\nabla_T = \text{diag}(1/\mathbf{g}_Q) \mathbf{Q}_k^\top \mathbf{C} \mathbf{R}_k \text{diag}(1/\mathbf{g}_R)$
 $\gamma_T = \gamma / \|\nabla_T\|_\infty$ # ℓ^∞ -normalization
 $\mathbf{K}_T^{(k)} \leftarrow \mathbf{T}_k \odot \exp(-\gamma_T \nabla_T)$
 $\mathbf{T}_k \leftarrow \text{Sinkhorn}(\mathbf{K}_T^{(k)}, \mathbf{g}_R, \mathbf{g}_Q, \delta)$ # Balanced OT, Alg. 5
 $\mathbf{X}_k \leftarrow \text{diag}(1/\mathbf{g}_Q) \mathbf{T}_k \text{diag}(1/\mathbf{g}_R)$
end while
 Return $\mathbf{P}_r = \mathbf{Q} \mathbf{X} \mathbf{R}^\top$

Convergence analysis of FRLC. As objective (8) is non-convex, it is important to have convergence guarantees. Our convergence criterion $\Delta(\cdot, \cdot)$ is defined in (10). To prove convergence we require a lower bound on the entries of \mathbf{g}_Q and \mathbf{g}_R . Previous works introduce a lower-bound vector $\alpha \leq \mathbf{g}$ enforced element-wise for stability and smoothness Scetbon et al. (2021). In FRLC the use of semi-relaxed projections naturally enforces a lower-bound. In Appendix E.5, we show that for any $\delta \in (0, \frac{1}{r})$, the FRLC algorithm's τ -weighted regularization on the inner marginals can guarantee a uniform lower-bound of δ on the entries: for sufficiently large τ and $\tilde{O}(m^2/\epsilon)$ iterations for the sub-coupling Pham et al. (2020), one guarantees a lower bound of δ on \mathbf{g}_R and \mathbf{g}_Q . This allows us to show objective smoothness in Proposition E.5. Previous work on low-rank optimal transport Scetbon et al. (2021) use the non-asymptotic convergence criterion of Ghadimi et al. (2014). Following existing works Dang & Lan (2015) establishing convergence rates of coordinate mirror-descent for smooth objectives, we show in Proposition 3.3 this criterion may be extended to coordinate-MD by adapting the block-descent lemma of Beck & Tetruashvili (2013).

Proposition 3.3. *Suppose one has $f \in C^1(\mathcal{X}, \mathbb{R})$ with block-coordinate Lipschitz gradient and block smoothness constants $(L_i)_{i=1}^p$, and a function $h \in C(\mathcal{X}, \mathbb{R})$ which is α -strongly convex. For $\Phi = f + h$, suppose one performs coordinate mirror descent on Φ minimized over a product of closed convex sets $\mathcal{X} = \prod_{i=1}^p \mathcal{X}_i$. Let the sub-iterates with respect to the i -th block update be $\{\mathbf{x}_k^i\}_{i=0}^p$ where $\mathbf{x}_k := \mathbf{x}_k^0$ for $k \in [N]$ outer iterations. Then one has:*

$$\min_k \Delta(\mathbf{x}_k, \mathbf{x}_{k-1}) \leq \frac{D^2 L}{N(\alpha^2/2L)} = \frac{2D^2 L^2}{N\alpha^2},$$

where D is (36), L is the global smoothness constant, stepsizes $\gamma_{k,i} := \alpha/L$, and convergence criterion $\Delta(\mathbf{x}_k, \mathbf{x}_{k-1})$ is given in (35).

Specialized to the LC-parametrization, the criterion $\Delta_k(\mathbf{x}_k, \mathbf{x}_{k+1})$ is:

$$\Delta_k(\mathbf{x}_k, \mathbf{x}_{k+1}) := \frac{1}{\gamma_k^2} [\|\mathbf{Q}_{k+1} - \mathbf{Q}_k\|_F^2 + \|\mathbf{R}_{k+1} - \mathbf{R}_k\|_F^2 + \|\mathbf{T}_{k+1} - \mathbf{T}_k\|_F^2] \quad (10)$$

for $\mathbf{x}_k = (\mathbf{Q}_k, \mathbf{R}_k, \mathbf{T}_k)$. We show through Propositions 3.3, E.5 the following result:

Proposition 3.4. *The FRLC algorithm with step-sizes $\gamma_k = \alpha/L$ and iterates $\mathbf{x}_k = (\mathbf{Q}_k, \mathbf{R}_k, \mathbf{T}_k)$ has non-asymptotic stationary convergence in the criterion $\Delta(\cdot, \cdot)$ with:*

$$\min_{k \in 1, \dots, N-1} \Delta_k(\mathbf{x}_k, \mathbf{x}_{k+1}) \leq 2D^2 L^2 / N\alpha^2$$

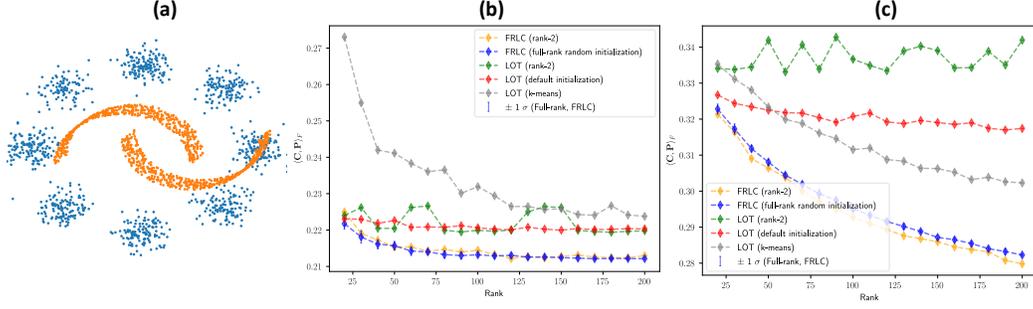


Figure 2: (a) Simulated dataset containing points from two moons (orange) and eight Gaussians (blue). (b) Transport cost $\langle C, P \rangle_F$ achieved by FRLC and LOT Scetbon et al. (2021) for the balanced Wasserstein problem on the dataset in (a) for different ranks and initializations. FLRC full rank (blue curve) is average over 10 random initializations. (c) Results on the 10D mixture of Gaussians dataset.

Where N is the number of iterations, D the optimality-gap as in (36), and $L = \max_{i \in \{1, \dots, 3\}} (L_i)$ the global smoothness for $L_i = \text{poly}(\|C\|_F, n, m, r, \delta)$ the block-wise smoothness constants.

The proof of Proposition 3.4 follows directly from our extension of the non-asymptotic criterion with the block-descent result Proposition 3.3 and the proof that this lemma holds in FRLC Proposition E.5. We also mention two improvements to other low rank approximation results in literature. In Proposition F.2 we show that one can *analytically* solve for the block-optimal g for the factorization of Scetbon et al. (2021), and we improve the bound on the low-rank approximation error in Proposition E.7.

FRLC for other marginal constraints and objectives. The balanced FRLC algorithm can be extended simply to other marginal constraints owing to the decoupling of the coordinate MD scheme. In particular, by using either the semi-relaxed projections (Algorithm 2) or fully-relaxed (unbalanced) projections (Algorithm 3) on sub-couplings Q and R , one can solve the balanced problem, the problem with the left or right marginal relaxed, or the unbalanced problem. As such, *all* variants of marginal constraints can be handled by a single algorithm, given in Algorithm 4.

We also extend the FRLC algorithm to the Gromov-Wasserstein problem. This consists of computing a GW-specific gradient with the appropriate marginal constraints applied to simplify their form, and re-computing Sinkhorn kernels as exponentiations of these gradients. The matrix form of the quadratic GW objective is $\mathbf{1}_m^T \mathbf{P}^T \mathbf{A}^{\odot 2} \mathbf{P} \mathbf{1}_m + \mathbf{1}_n^T \mathbf{P} \mathbf{B}^{\odot 2} \mathbf{P}^T \mathbf{1}_n - 2 \langle \mathbf{A} \mathbf{P} \mathbf{B}, \mathbf{P} \rangle$, where \odot denotes the Hadamard (entrywise) product. Then the GW-specific Sinkhorn kernels are

$$\begin{aligned} K_Q^{(k)} &\leftarrow \exp(2\gamma_k(2\mathbf{A} \mathbf{Q} \mathbf{X} \mathbf{R}^T \mathbf{B} \mathbf{R} \mathbf{X}^T - \mathbf{A}^{\odot 2} \mathbf{Q} \mathbf{1}_r \mathbf{1}_r^T)), \\ K_R^{(k)} &\leftarrow \exp(2\gamma_k(2\mathbf{B} \mathbf{R} \mathbf{X}^T \mathbf{Q}^T \mathbf{A} \mathbf{Q} \mathbf{X} - \mathbf{B}^{\odot 2} \mathbf{R} \mathbf{1}_r \mathbf{1}_r^T)), \\ K_T^{(k)} &\leftarrow \exp(4\gamma_k \text{diag}(g_Q^{-1}) \mathbf{Q}^T \mathbf{A} \mathbf{Q} \mathbf{X} \mathbf{R}^T \mathbf{B} \mathbf{R} \text{diag}(g_R^{-1})). \end{aligned}$$

In Algorithm 4, one can solve the GW-problem by using the kernels above. Here, we present the kernels omitting a rank-1 perturbation, which is given in Appendix D. From the Wasserstein and GW gradients, the FGW gradient is easily taken as a convex combination of the two. In this work, we primarily focus on the LC-factorization for the rank r Wasserstein problem (6).

4 Experimental Results

We compare FRLC to existing low-rank and full-rank optimal transport algorithms on several datasets: simulated datasets previously used in Tong et al. (2023) and Scetbon et al. (2021); a massive spatial-transcriptomics dataset Chen et al. (2022); and a graph partitioning task Chowdhury & Needham (2021). Further details of each experiment (e.g. pre-processing, validation) are in Appendices K, L, and M. In the section below, LOT refers to the works of Scetbon et al. (2021, 2023, 2022) and Latent OT refers to Lin et al. (2021).

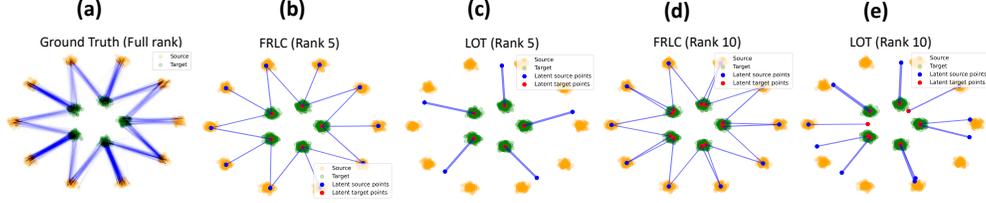


Figure 3: LC-projections of couplings of Gaussians centered on the 5th-roots of unity (green) and 10th roots of unity (yellow). (a) Ground-truth full-rank coupling. (b) Non-square rank-5 latent-coupling of FRLC (c) LC-projection barycenters aligned with rank-5 diagonal coupling of LOT Scetbon et al. (2021). (d) Square rank-10 latent coupling of FRLC. (e) Rank-10 diagonal coupling of LOT .

4.1 Evaluation of Low-rank Approximations for Balanced OT on Synthetic Data

We first compare the balanced OT version of FRLC with the the low-rank balanced OT algorithm LOT of Scetbon et al. (2021) on a synthetic dataset following Tong et al. (2023). The dataset consists of $m = 1000$ points from two moons and $n = 1000$ points sampled from eight 2D Gaussian densities (Fig. 2a). We solve the Wasserstein problem (1) with cost matrix C computed using the Euclidean distance. The full-rank coupling matrix P has rank 1000, and we compute both FRLC and LOT solutions with rank between 20 and 200. For each rank, we initialize FRLC adapting the deterministic rank-2 initialization proposed in Scetbon et al. (2021) and the random initialization of Alg. 6. We initialize LOT using the rank-2 initialization and two other options in `ott-jax` Cuturi et al. (2022).

We find that FRLC obtains lower transport cost $\langle C, P \rangle_F$ with increasing rank (Fig. 2b) and consistently achieves lower transport cost than LOT across all ranks and all initializations. Specifically, starting both methods at the same rank-2 initialization, FRLC consistently achieves a lower cost than LOT for all ranks. Additionally, we observe smooth convergence of FRLC for both rank-2 initialization and the full-rank random initialization of Alg. 6 (Fig. 5).

We also evaluate FRLC and LOT on two datasets of Gaussian mixtures, one in 2-dimensions and one in 10-dimensions, each with $n = m = 5,000$ points from two mixtures of Gaussians, following Scetbon et al. (2021), with further details in Appendix K. We observe the same trend as the previous simulation for both datasets (Fig. 2c, Fig. 7), with FRLC achieving lower transport costs than LOT across all ranks and all initializations. In addition FRLC has half the runtime of LOT (CPU) – including the setup time of FRLC but excluding the setup time of LOT in `ott-jax` – on datasets of $n = m = 1000$ points from all three datasets with rank $r = 100$ (Table 2). At the same time FRLC achieves lower primal cost $\langle C, P \rangle_F$ with tighter marginals $\|P\mathbf{1}_n - \mathbf{a}\|_2$ and $\|P^T\mathbf{1}_m - \mathbf{b}\|_2$. Lin et al. (2021) only solves a proxy for the rank-constrained Wasserstein problem, and thus is not the focus of our comparisons. Nevertheless, we verify that on all synthetic experiments that FRLC achieves significantly lower primal OT cost than Latent OT (Table 5).

4.2 Interpretation of the Latent Coupling and LC-Projection

We demonstrate the interpretability of the latent coupling T in the LC factorization. In both the LC factorization and factored couplings, the sub-couplings Q and R each have associated barycentric projection operators which coarse-grain input datasets $Z^{(1)}, Z^{(2)}$. In particular, the LC projection is defined from the LC factorization as follows.

Definition 4.1 (LC-Projection). Let $Q \text{diag}(1/g_Q)T \text{diag}(1/g_R)R^T$ be an LC factorization of of a coupling matrix $P \in \Pi_{a,b}(r)$ computed from datasets $Z^{(1)} \in \mathbb{R}^{n \times d}, Z^{(2)} \in \mathbb{R}^{m \times d}$, with $T \in \mathbb{R}_+^{r_1 \times r_2}$. The LC-projections $Y^{(1)}$ and $Y^{(2)}$ of $Z^{(1)}$ and $Z^{(2)}$ are $Y^{(1)} := \text{diag}(1/g_Q)Q^T Z^{(1)}$, and $Y^{(2)} := \text{diag}(1/g_R)R^T Z^{(2)}$.

By interpreting any factored coupling (Q, R, g) as an LC factorization $(Q, R, \text{diag}(g))$, Definition 4.1 describes the barycentric projections for both factorizations. We compare the projections of the coupling computed by FRLC to those of LOT Scetbon et al. (2021) on a dataset containing 1000 samples from 2D-Gaussians centered at the 5th-roots of unity and 1000 samples from 2D Gaussians

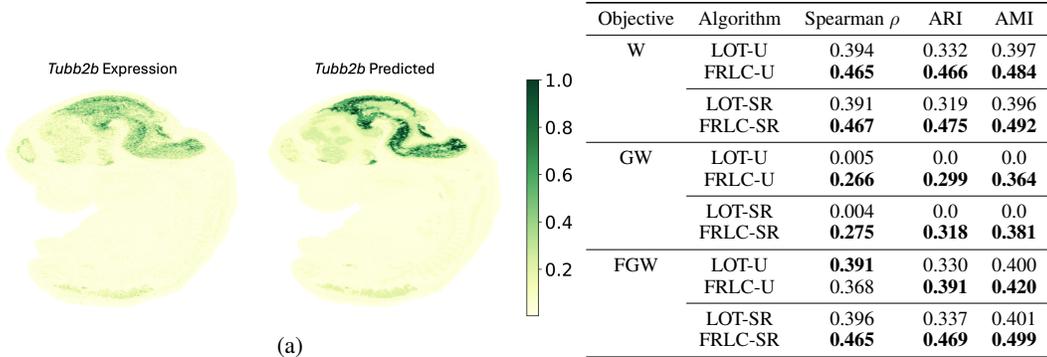


Figure 4: (a) Brain marker gene *Tubb2b* expression and FRLC prediction. (b) Comparison of the low-rank unbalanced (LOT-U) algorithm of Scetbon et al. (2023) and FRLC on aligning spatial transcriptomics data. Bold indicates top performing method for each metric on each objective.

centered at the 10th-roots of unity (the latter scaled by a factor of two, Fig. 3a). In both cases, the latent coupling \mathbf{T} or $\text{diag}(\mathbf{g})$ is visualized as a transport between barycenters. We run FRLC and LOT with ranks $r = 5$ and $r = 10$ to match the number of target and source clusters. In the rank-5 case, FRLC uses a *non-square* latent coupling $\mathbf{T} \in \mathbb{R}_+^{10 \times 5}$ which correctly captures the coupling between clusters (Fig. 3(b)), while the LOT rank-5 projection computes barycenters that are outside of the clusters (Fig. 3c). A similar result is observed for square rank-10 latent couplings computed by FRLC (Fig. 3d) and LOT (Fig. 3e) demonstrating that the LOT barycenters in Fig. 3b) are not an artifact of using the lowest rank. We observe similar results on other simulated datasets (Fig. 11).

4.3 Evaluation on Spatial Transcriptomics Alignment

We compare FRLC and the algorithm (LOT-U) of Scetbon et al. (2023) (which solves unbalanced low-rank Wasserstein, GW, and FGW problems) on the task of computing an alignment between cells from different time points during mouse embryonic development. Specifically, we compute an alignment between a spatial transcriptomics (ST) dataset of an E11.5 stage mouse embryo and an E12.5 stage mouse embryo Chen et al. (2022). Optimal transport is a popular approach to align single-cell Schiebinger et al. (2019) and spatial transcriptomics datasets Zeira et al. (2022); Liu et al. (2023); Klein et al. (2023). In single-cell transcriptomics, one measures a gene expression vector for each cell, and in spatial transcriptomics one additionally measures the 2D location of each cell. The cost matrix \mathbf{C} describes the difference between gene expression vectors and intra-domain cost matrices \mathbf{A} and \mathbf{B} are derived from the 2D coordinates within each slice. Therefore, OT problems of W, GW, and FGW objectives can be solved and the coupling matrix represents the cell-cell alignment (Appendix M). However, computation of a full-rank OT solution is not feasible in our large-scale dataset: the E11.5 slice has about 30,000 cells while the E12.5 slice has about 50,000 cells.

We evaluate the alignments by assessing performance on two prediction tasks from Scetbon et al. (2023): (1) a *gene expression prediction* task where we predict the expression of a gene in E12.5 from expression of the gene in E11.5 using the alignment; (2) a *cell type prediction* task where we predict the cell types of E12.5 from the cell type clustering of E11.5 (Appendix M). We evaluate the accuracy of the gene expression prediction task through the Spearman correlation ρ between the predicted expression and the ground truth expression of 10 test marker genes. We evaluate the accuracy of the cell type prediction task by computing the Adjusted Rand Index (ARI) and Adjusted Mutual Information (AMI) between the predicted cell types and the cell types derived in the original publication Chen et al. (2022). Being a comparison between different objectives, this relies on downstream metrics. For completeness, we validate the efficacy of FRLC on directly minimizing the balanced Wasserstein cost $\langle \mathbf{C}, \mathbf{P} \rangle_F$ against Scetbon et al. (2021) in Figure 8.

For a direct comparison, we use FRLC to solve the same unbalanced problems (denoted FRLC-U). We perform an extensive grid search (Appendix M.3) to pick the best hyperparameters (including rank $\ll 30,000$) for all algorithms. Scetbon et al. (2023) previously showed that unbalanced FGW algorithm has the best performance on ST alignment. We find that unbalanced FRLC achieves comparable or better results than the previous state-of-the-art unbalanced low-rank method on all three objectives (Table 4). We also solve a semi-relaxed version of each problem motivated by the

Method	Factorization	Cost	Variables	Algorithm	Sub-routine for coupling
Factored Coupling Forrow et al. (2019)	Factored coupling	k -Wasserstein barycenter	Anchors & sub-couplings	Lloyd-type	Dijkstra’s
Latent OT Lin et al. (2021)	Latent coupling	Extension of k -Wasserstein barycenter	Anchors & sub-couplings	Lloyd-type	Dijkstra’s
LOT Scetbon et al. (2021)	Factored coupling	Primal OT cost	Sub-couplings & inner marginal	Mirror-descent	Dijkstra’s
FRLC (this work)	Latent coupling	Primal OT cost	Sub-couplings	Coordinate mirror-descent	OT

Table 1: Comparing aspects of low-rank OT methods. Factorization indicates the structure of the inner matrix.

observation that all cells from E12.5 have an ancestor, but not all cells from E11.5 have the same number of descendants due to cell growth and death. Thus the former marginal is tight, and the latter relaxed Halmos et al. (2024). We run both semi-relaxed FRLC (FRLC-SR) and a setting of LOT-U that recovers the semi-relaxed problem (LOT-SR). Semi-relaxed FRLC achieves the best results on all three metrics by a large margin (Table 4). As one example, the expression of *Tubb2b*, a mouse brain marker gene, agreeing with the expression predicted from the semi-relaxed alignment of FRLC (Fig. 4a).

4.4 Additional Experiments

We evaluate FRLC on an unsupervised graph partitioning problem Chowdhury & Needham (2021) on four real-world graph datasets Yang & Leskovec (2012); Yin et al. (2017); Banerjee et al. (2013). We benchmark the performance of the semi-relaxed and GW settings of FRLC against (1) GWL Xu et al. (2019), solving a balanced GW problem; (2) SpecGWL Chowdhury & Needham (2021) using the heat kernel on the graph Laplacian as the cost matrix. We find FRLC achieves the better clustering performance than GWL and SpecGWL on 9/12 and 11/12 of the datasets (Table 3 and Appendix L).

5 Discussion

We provide comparison of existing low-rank solvers in Table 1. The FRLC algorithm has a number of advantages, including (1) coarsening a full-rank plan P to non-diagonal latent coupling T ; (2) minimizing the primal OT problem for general cost C rather than a barycentric problem; (3) optimizing only sub-couplings; and (4) using Sinkhorn alone as the sub-routine for low-rank OT. While we argue these are substantial advantages, FRLC has limitations which warrant follow-up work. In particular, three key limitations of our work, common to the existing low-rank OT algorithms, are: (1) selecting values of the latent coupling ranks; (2) strengthening the convergence criterion; (3) addressing sensitivity to the initialization from non-convexity of the objective. A limitation specific to our work is the selection of the τ hyperparameter controlling the smoothness of the trajectory. These and other limitations are discussed in Section N of the Appendix. Another direction for further investigation is to better understand what structure LC factorizations capture when the optimal plan is known to have full rank, e.g. when the Monge map exists, as has been explored by Liu et al. (2021).

6 Conclusion

We introduce FRLC, an algorithm to compute low-rank optimal transport plan from the latent coupling (LC) factorization. FRLC handles different OT objective costs and relaxations of the marginal constraints. Moreover, the LC factorization provides an interpretable coarse-graining of the full transport plan and its marginals through the mapping $(P, a, b) \rightarrow (T, g_Q, g_R)$. We demonstrate the superior performance of FRLC compared to state-of-the-art low-rank methods on real and synthetic datasets.

Acknowledgments and Disclosure of Funding

This work is supported by NCI grant U24CA248453 to B.J.R. J.G. gratefully acknowledges support from the Schmidt DataX Fund at Princeton University made possible through a major gift from the Schmidt Futures Foundation.

References

- Bakshi, A. and Woodruff, D. Sublinear Time Low-Rank Approximation of Distance Matrices. *Advances in Neural Information Processing Systems*, 31, 2018.
- Banerjee, A., Chandrasekhar, A. G., Duflo, E., and Jackson, M. O. The Diffusion of Microfinance. *Science*, 341(6144):1236498, 2013.
- Bauschke, H. H. and Lewis, A. S. Dykstras algorithm with Bregman projections: A convergence proof. *Optimization*, 48(4):409–427, January 2000. ISSN 1029-4945. doi: 10.1080/02331930008844513. URL <http://dx.doi.org/10.1080/02331930008844513>.
- Beck, A. and Tetruashvili, L. On the Convergence of Block Coordinate Descent Type Methods. *SIAM J. Optim.*, 23:2037–2060, 2013. URL <https://api.semanticscholar.org/CorpusID:6866704>.
- 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, January 2015. ISSN 1095-7197. doi: 10.1137/141000439. URL <http://dx.doi.org/10.1137/141000439>.
- Bregman, L. M. The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR computational mathematics and mathematical physics*, 7(3):200–217, 1967.
- Bunne, C., Stark, S. G., Gut, G., del Castillo, J. S., Levesque, M., Lehmann, K.-V., Pelkmans, L., Krause, A., and Rätsch, G. Learning single-cell perturbation responses using neural optimal transport. *Nature Methods*, 20(11):1759–1768, September 2023. ISSN 1548-7105. doi: 10.1038/s41592-023-01969-x. URL <http://dx.doi.org/10.1038/s41592-023-01969-x>.
- Charikar, M., Chen, B., Ré, C., and Waingarten, E. Fast Algorithms for a New Relaxation of Optimal Transport. In Neu, G. and Rosasco, L. (eds.), *Proceedings of Thirty Sixth Conference on Learning Theory*, volume 195 of *Proceedings of Machine Learning Research*, pp. 4831–4862. PMLR, 12–15 Jul 2023. URL <https://proceedings.mlr.press/v195/charikar23a.html>.
- Chen, A., Liao, S., Cheng, M., Ma, K., Wu, L., Lai, Y., Qiu, X., Yang, J., Xu, J., Hao, S., et al. Spatiotemporal transcriptomic atlas of mouse organogenesis using DNA nanoball-patterned arrays. *Cell*, 185(10):1777–1792, 2022.
- Chen, X. and Price, E. Condition number-free query and active learning of linear families. 2017.
- Chizat, L., Peyré, G., Schmitzer, B., and Vialard, F.-X. Unbalanced Optimal Transport: Dynamic and Kantorovich Formulations. *Journal of Functional Analysis*, 274(11):3090–3123, June 2018. ISSN 0022-1236. doi: 10.1016/j.jfa.2018.03.008. URL <http://dx.doi.org/10.1016/j.jfa.2018.03.008>.
- Chowdhury, S. and Needham, T. Generalized Spectral Clustering via Gromov-Wasserstein Learning. In *International Conference on Artificial Intelligence and Statistics*, pp. 712–720. PMLR, 2021.
- Chung, F. Laplacians and the Cheeger Inequality for Directed Graphs. *Annals of Combinatorics*, 9: 1–19, 2005.
- Cohen, J. E. and Rothblum, U. G. Nonnegative Ranks, Decompositions, and Factorizations of Nonnegative Matrices. *Linear Algebra and its Applications*, 190:149–168, 1993.
- Courty, N., Flamary, R., and Tuia, D. Domain adaptation with regularized optimal transport. In *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2014, Nancy, France, September 15-19, 2014. Proceedings, Part I 14*, pp. 274–289. Springer, 2014.

- Cuturi, M. Sinkhorn Distances: Lightspeed Computation of Optimal Transport. *Advances in neural information processing systems*, 26, 2013a.
- Cuturi, M. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in Neural Information Processing Systems*, pp. 2292–2300, 2013b. URL <https://proceedings.neurips.cc/paper/2013/hash/af21d0c97db2e27e13572cbf59eb343d-Abstract.html>.
- Cuturi, M., Meng-Papaxanthos, L., Tian, Y., Bunne, C., Davis, G., and Teboul, O. Optimal Transport Tools (OTT): A JAX Toolbox for all things Wasserstein. *arXiv preprint arXiv:2201.12324*, 2022.
- Dang, C. D. and Lan, G. Stochastic Block Mirror Descent Methods for Nonsmooth and Stochastic Optimization. *SIAM J. Optim.*, 25(2):856–881, January 2015.
- Dong, S., Pan, Z., Fu, Y., Xu, D., Shi, K., Yang, Q., Shi, Y., and Zhuo, C. Partial Unbalanced Feature Transport for Cross-Modality Cardiac Image Segmentation. *IEEE Transactions on Medical Imaging*, 2023.
- Dykstra, R. L. An Algorithm for Restricted Least Squares Regression. *Journal of the American Statistical Association*, 78(384):837–842, December 1983. ISSN 1537-274X. doi: 10.1080/01621459.1983.10477029. URL <http://dx.doi.org/10.1080/01621459.1983.10477029>.
- Forrow, A., Hütter, J.-C., Nitzan, M., Rigollet, P., Schiebinger, G., and Weed, J. Statistical Optimal Transport via Factored Couplings. In Chaudhuri, K. and Sugiyama, M. (eds.), *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pp. 2454–2465. PMLR, 16–18 Apr 2019. URL <https://proceedings.mlr.press/v89/forrow19a.html>.
- Frieze, A., Kannan, R., and Vempala, S. Fast Monte-Carlo Algorithms for Finding Low-rank Approximations. *J. ACM*, 51(6):1025–1041, nov 2004. ISSN 0004-5411. doi: 10.1145/1039488.1039494. URL <https://doi.org/10.1145/1039488.1039494>.
- Frogner, C., Zhang, C., Mobahi, H., Araya, M., and Poggio, T. A. Learning with a Wasserstein Loss. *Advances in neural information processing systems*, 28, 2015.
- Geshkovski, B., Letrouit, C., Polyanskiy, Y., and Rigollet, P. A mathematical perspective on Transformers. *arXiv preprint arXiv:2312.10794*, 2023.
- Ghadimi, S., Lan, G., and Zhang, H. Mini-batch stochastic approximation methods for nonconvex stochastic composite optimization. *Mathematical Programming*, 155(1–2):267–305, December 2014. ISSN 1436-4646. doi: 10.1007/s10107-014-0846-1. URL <http://dx.doi.org/10.1007/s10107-014-0846-1>.
- Halmos, P., Liu, X., Gold, J., Chen, F., Ding, L., and Raphael, B. J. DeST-OT: Alignment of Spatiotemporal Transcriptomics Data. In *International Conference on Research in Computational Molecular Biology*, pp. 434–437. Springer, 2024.
- Indyk, P., Vakilian, A., Wagner, T., and Woodruff, D. P. Sample-optimal low-rank approximation of distance matrices. In Beygelzimer, A. and Hsu, D. (eds.), *Proceedings of the Thirty-Second Conference on Learning Theory*, volume 99 of *Proceedings of Machine Learning Research*, pp. 1723–1751. PMLR, 25–28 Jun 2019. URL <https://proceedings.mlr.press/v99/indyk19a.html>.
- Kantorovich, L. On the Translocation of Masses: Doklady akademii nauk ussr. 1942.
- Klein, D., Palla, G., Lange, M., Klein, M., Piran, Z., Gander, M., Meng-Papaxanthos, L., Sterr, M., Bastidas-Ponce, A., Tarquis-Medina, M., et al. Mapping cells through time and space with moscot. *bioRxiv*, pp. 2023–05, 2023.
- Lin, C.-H., Azabou, M., and Dyer, E. L. Making transport more robust and interpretable by moving data through a small number of anchor points. *Proceedings of machine learning research*, 139: 6631, 2021.
- Liu, W., Zhang, C., Zheng, N., and Qian, H. Approximating optimal transport via low-rank and sparse factorization. *CoRR*, abs/2111.06546, 2021. URL <https://arxiv.org/abs/2111.06546>.

- Liu, X., Zeira, R., and Raphael, B. J. Partial alignment of multislice spatially resolved transcriptomics data. *Genome Research*, 33(7):1124–1132, 2023.
- Mémoli, F. On the use of Gromov-Hausdorff Distances for Shape Comparison. 2007.
- Mémoli, F. Gromov–Wasserstein Distances and the Metric Approach to Object Matching. *Foundations of computational mathematics*, 11:417–487, 2011.
- Nesterov, Y. Efficiency of Coordinate Descent Methods on Huge-Scale Optimization Problems. *SIAM Journal on Optimization*, 22(2):341–362, 2012.
- Orlin, J. B. A polynomial time primal network simplex algorithm for minimum cost flows. *Mathematical Programming*, 78(2):109–129, Aug 1997. ISSN 1436-4646. doi: 10.1007/BF02614365. URL <https://link.springer.com/content/pdf/10.1007/BF02614365.pdf>.
- Pham, K., Le, K., Ho, N., Pham, T., and Bui, H. H. On Unbalanced Optimal Transport: An Analysis of Sinkhorn Algorithm. In *International Conference on Machine Learning*, 2020. URL <https://api.semanticscholar.org/CorpusID:211068892>.
- Sander, M. E., Ablin, P., Blondel, M., and Peyré, G. Sinkformers: Transformers with Doubly Stochastic Attention. In Camps-Valls, G., Ruiz, F. J. R., and Valera, I. (eds.), *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pp. 3515–3530. PMLR, 28–30 Mar 2022. URL <https://proceedings.mlr.press/v151/sander22a.html>.
- Scetbon, M. and Cuturi, M. Low-rank Optimal Transport: Approximation, Statistics and Debiasing. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K. (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=4btNeXKFAQ>.
- Scetbon, M., Cuturi, M., and Peyré, G. Low-Rank Sinkhorn Factorization. In *International Conference on Machine Learning*, 2021. URL <https://api.semanticscholar.org/CorpusID:232147563>.
- Scetbon, M., Peyré, G., and Cuturi, M. Linear-time Gromov Wasserstein Distances using Low Rank Couplings and Costs. In *International Conference on Machine Learning*, pp. 19347–19365. PMLR, 2022.
- Scetbon, M., Klein, M., Palla, G., and Cuturi, M. Unbalanced Low-rank Optimal Transport Solvers, 2023.
- 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.
- 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.
- Stähl, P. L., Salmén, F., Vickovic, S., Lundmark, A., Navarro, J. F., Magnusson, J., Giacomello, S., Asp, M., Westholm, J. O., Huss, M., et al. Visualization and analysis of gene expression in tissue sections by spatial transcriptomics. *Science*, 353(6294):78–82, 2016.
- Tarjan, R. E. Dynamic trees as search trees via Euler tours, applied to the network simplex algorithm. *Mathematical Programming*, 78(2):169–177, Aug 1997. ISSN 1436-4646. doi: 10.1007/BF02614369. URL <https://link.springer.com/content/pdf/10.1007/BF02614369.pdf>.
- Tay, Y., Bahri, D., Yang, L., Metzler, D., and Juan, D.-C. Sparse Sinkhorn Attention. In III, H. D. and Singh, A. (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 9438–9447. PMLR, 13–18 Jul 2020. URL <https://proceedings.mlr.press/v119/tay20a.html>.
- Tong, A., Malkin, N., Hugué, G., Zhang, Y., Rector-Brooks, J., Fatras, K., Wolf, G., and Bengio, Y. Improving and generalizing flow-based generative models with minibatch optimal transport. In *ICML Workshop on New Frontiers in Learning, Control, and Dynamical Systems*, 2023.

- Vayer, T., Chapel, L., Flamary, R., Tavenard, R., and Courty, N. Fused Gromov-Wasserstein distance for structured objects. *Algorithms*, 13(9):212, August 2020. ISSN 1999-4893. doi: 10.3390/a13090212. URL <http://dx.doi.org/10.3390/a13090212>.
- Vincent-Cuaz, C., Flamary, R., Corneli, M., Vayer, T., and Courty, N. Semi-relaxed Gromov-Wasserstein divergence and applications on graphs. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=RShaMexjc-x>.
- Wolf, F. A., Angerer, P., and Theis, F. J. SCANPY: large-scale single-cell gene expression data analysis. *Genome biology*, 19:1–5, 2018.
- Xu, H., Luo, D., Zha, H., and Duke, L. C. Gromov-Wasserstein Learning for Graph Matching and Node Embedding. In *International conference on machine learning*, pp. 6932–6941. PMLR, 2019.
- Yang, J. and Leskovec, J. Defining and Evaluating Network Communities based on Ground-truth. In *Proceedings of the ACM SIGKDD Workshop on Mining Data Semantics*, pp. 1–8, 2012.
- Yang, K. D., Damodaran, K., Venkatachalapathy, S., Soylemezoglu, A. C., Shivashankar, G., and Uhler, C. Predicting cell lineages using autoencoders and optimal transport. *PLoS computational biology*, 16(4):e1007828, 2020.
- Yin, H., Benson, A. R., Leskovec, J., and Gleich, D. F. Local Higher-Order Graph Clustering. In *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 555–564, 2017.
- Zeira, R., Land, M., Strzalkowski, A., and Raphael, B. J. Alignment and integration of spatial transcriptomics data. *Nature Methods*, 19(5):567–575, 2022.

A Low-rank optimal transport

A.1 Low-rank factorizations

The set of low-rank couplings. Given $M \in \mathbb{R}_+^{n \times m}$, the *nonnegative rank* of M is the least number of nonnegative, rank-1 matrices that sum to M :

$$\text{rk}_+(M) = \min_{r \geq 1} \left\{ M = \sum_{i=1}^r M_i, \text{ such that } \text{rk}(M_i) = 1 \text{ and } M_i \geq 0 \text{ for all } i \right\}.$$

Let $\mathbf{a} \in \Delta_n, \mathbf{b} \in \Delta_m$ be probability vectors, and let $\Pi_{\mathbf{a},\mathbf{b}}(r)$ denote the set of rank- r coupling matrices with marginals \mathbf{a} and \mathbf{b} :

$$\Pi_{\mathbf{a},\mathbf{b}}(r) = \{P \in \mathbb{R}_+^{n \times m} : P^T \mathbf{1}_m = \mathbf{a}, P \mathbf{1}_n = \mathbf{b}, \text{rk}_+(P) \leq r\}.$$

To optimize any cost over $\Pi_{\mathbf{a},\mathbf{b}}(r)$, one requires a parameterization of this set.

Factored couplings. The *factored coupling* parameterization of $\Pi_{\mathbf{a},\mathbf{b}}(r)$ introduced in Forrow et al. (2019), and used by Scetbon et al. (2021); Scetbon & Cuturi (2022); Scetbon et al. (2022, 2023) is

$$\text{FC}_{\mathbf{a},\mathbf{b}}(r) := \{(\mathbf{Q}, \mathbf{R}, \mathbf{g}) \in \mathbb{R}_+^{n \times r} \times \mathbb{R}_+^{m \times r} \times (\mathbb{R}_+^*)^r : \mathbf{Q} \in \Pi_{\mathbf{a},\mathbf{g}}, \mathbf{R} \in \Pi_{\mathbf{b},\mathbf{g}}\}.$$

Cohen & Rothblum (1993) show that any $P \in \Pi_{\mathbf{a},\mathbf{b}}(r)$ may be decomposed as $P = \mathbf{Q} \text{diag}(1/\mathbf{g}) \mathbf{R}^T$ for some triple $(\mathbf{Q}, \mathbf{R}, \mathbf{g}) \in \text{FC}$. Thus, for cost matrix $C \in \mathbb{R}^{n \times m}$, the general low-rank optimal transport problem is equivalent to an optimization over factored couplings:

$$\min_{P \in \Pi_{\mathbf{a},\mathbf{b}}(r)} \langle C, P \rangle_F = \min_{(\mathbf{Q}, \mathbf{R}, \mathbf{g}) \in \text{FC}_{\mathbf{a},\mathbf{b}}(r)} \langle C, \mathbf{Q} \text{diag}(1/\mathbf{g}) \mathbf{R}^T \rangle_F. \quad (11)$$

Latent coupling factorization. The *latent coupling* parameterization of $\Pi_{\mathbf{a},\mathbf{b}}(r)$ introduced in Lin et al. (2021), and used in the present work is

$$\text{LC}_{\mathbf{a},\mathbf{b}}(r) := \{(\mathbf{Q}, \mathbf{R}, \mathbf{T}) \in \mathbb{R}_+^{n \times r} \times \mathbb{R}_+^{m \times r} \times \mathbb{R}_+^{r \times r} : \mathbf{Q} \in \Pi_{\mathbf{a},\cdot}, \mathbf{R} \in \Pi_{\cdot,\mathbf{b}}, \mathbf{T} \in \Pi_{\mathbf{g}_Q, \mathbf{g}_R}\},$$

where $\mathbf{g}_Q, \mathbf{g}_R$ are the inner marginals of \mathbf{Q} and \mathbf{R} .

Latent coupling diagonalization. The LC-factorization recovers the factorization of Forrow et al. (2019) as a sub-case. While the diagonal factorization of previous works cannot be directly converted to the LC-factorization, the LC-factorization can easily recover the diagonal factorization. In particular, taking $\mathbf{Q}' \leftarrow \mathbf{Q} \text{diag}(1/\mathbf{g}_Q) \mathbf{T}$ one can refactor

$$P_r = \mathbf{Q} \text{diag}(1/\mathbf{g}_Q) \mathbf{T} \text{diag}(1/\mathbf{g}_R) \mathbf{R}^T = \mathbf{Q}' \text{diag}(1/\mathbf{g}_R) \mathbf{R}^T$$

or alternatively taking $\mathbf{R}' = \mathbf{R} \text{diag}(1/\mathbf{g}_R) \mathbf{T}^T$ may refactor as

$$P_r = \mathbf{Q} \text{diag}(1/\mathbf{g}_Q) \mathbf{T} \text{diag}(1/\mathbf{g}_R) \mathbf{R}^T = \mathbf{Q} \text{diag}(1/\mathbf{g}_Q) (\mathbf{R}')^T$$

So that instead of returning $(\mathbf{Q}, \mathbf{R}, \mathbf{T})$ one may alternatively return $(\mathbf{Q}, \mathbf{R}, \mathbf{T}) \rightarrow (\mathbf{Q}', \mathbf{R}, \text{diag}(\mathbf{g}_R))$ or $(\mathbf{Q}, \mathbf{R}, \mathbf{T}) \rightarrow (\mathbf{Q}, \mathbf{R}', \text{diag}(\mathbf{g}_Q))$ to recover the Forrow et al. (2019) factorization. An example of this diagonal-conversion is offered in Figure 12.

A.2 Balanced low-rank optimal transport

The FRLC optimization problem Our optimization problem is over the variables $(\mathbf{Q}, \mathbf{R}, \mathbf{T})$ and defined as follows:

$$\min_{(\mathbf{Q}, \mathbf{R}, \mathbf{T}) \in \text{LC}_{\mathbf{a},\mathbf{b}}(r)} \mathcal{L}_{\text{LC}}(\mathbf{Q}, \mathbf{R}, \mathbf{T}), \quad (12)$$

where our objective function \mathcal{L}_{LC} is

$$\mathcal{L}_{\text{LC}}(\mathbf{Q}, \mathbf{R}, \mathbf{T}) = \langle C, \mathbf{Q} (\text{diag}(1/\mathbf{Q}^T \mathbf{1}_n)) \mathbf{T} (\text{diag}(1/\mathbf{R}^T \mathbf{1}_m)) \mathbf{R}^T \rangle \quad (13)$$

Given $(\mathbf{Q}, \mathbf{R}, \mathbf{T}) \in \text{LC}_{\mathbf{a},\mathbf{b}}(r)$, sub-couplings \mathbf{Q} and \mathbf{R} are constrained by:

$$\mathcal{C}_1(\mathbf{a}) := \{(\mathbf{Q}, \mathbf{R}, \mathbf{T}) \in \mathcal{R}_+ : \mathbf{Q} \mathbf{1}_r = \mathbf{a}\}, \quad \mathcal{C}_1(\mathbf{b}) := \{(\mathbf{Q}, \mathbf{R}, \mathbf{T}) \in \mathcal{R}_+ : \mathbf{R} \mathbf{1}_r = \mathbf{b}\},$$

while the convex sets constraining the latent coupling matrix \mathbf{T} are

$$\mathcal{C}_2(\mathbf{g}_Q) := \{(\mathbf{Q}, \mathbf{R}, \mathbf{T}) \in \mathcal{R}_+ : \mathbf{T}\mathbf{1}_r = \mathbf{g}_Q\}, \quad \mathcal{C}_2(\mathbf{g}_R) := \{(\mathbf{Q}, \mathbf{R}, \mathbf{T}) \in \mathcal{R}_+ : \mathbf{T}^\top \mathbf{1}_n = \mathbf{g}_R\},$$

where $\mathbf{g}_Q = \mathbf{Q}^\top \mathbf{1}_n$ and $\mathbf{g}_R = \mathbf{R}^\top \mathbf{1}_m$ as per Definition 3.1. Under these definitions, $\mathcal{L}_{\mathcal{C}_a, \mathcal{C}_b}(r) = \mathcal{C}_1 \cap \mathcal{C}_2$, where

$$\mathcal{C}_1 = \mathcal{C}_1(\mathbf{a}) \cap \mathcal{C}_1(\mathbf{b}) \quad \text{and} \quad \mathcal{C}_2 = \mathcal{C}_2(\mathbf{g}_Q) \cap \mathcal{C}_2(\mathbf{g}_R) \quad (14)$$

To solve (12), we separate the variables into two ‘‘blocks’’ of variables, (\mathbf{Q}, \mathbf{R}) and \mathbf{T} , and perform two block updates per iteration, as follows. Let $(\gamma_k)_{k \geq 0}$ be a positive sequence of stepsizes. Suppose we have $(\mathbf{Q}_k, \mathbf{R}_k, \mathbf{T}_k) \in \mathcal{C}$. We update the first variable block (\mathbf{Q}, \mathbf{R}) by taking a locally optimal (mirror descent) update step, while the second variable block \mathbf{T} is held fixed:

$$(\mathbf{Q}_{k+1}, \mathbf{R}_{k+1}) \leftarrow \arg \min_{(\mathbf{Q}, \mathbf{R}) : (\mathbf{Q}, \mathbf{R}, \mathbf{T}_k) \in \mathcal{C}_1} \langle (\mathbf{Q}, \mathbf{R}), \nabla_{(\mathbf{Q}, \mathbf{R})} \mathcal{L}_{\text{LC}} \rangle + \frac{1}{\gamma_k} \text{KL}((\mathbf{Q}, \mathbf{R}) \| (\mathbf{Q}_k, \mathbf{R}_k)) \quad (15)$$

Here, we slightly abused the notation by putting (\mathbf{Q}, \mathbf{R}) inside an inner product. The triple $(\mathbf{Q}_{k+1}, \mathbf{R}_{k+1}, \mathbf{T}_k)$ produced by this update lies in \mathcal{C}_1 . Next, we update \mathbf{T} , the second variable block, by taking another locally optimal (mirror descent) step, while the first variable block is held fixed.

$$\mathbf{T}_{k+1} \leftarrow \arg \min_{\mathbf{T} : (\mathbf{Q}_{k+1}, \mathbf{R}_{k+1}, \mathbf{T}) \in \mathcal{C}_2} \langle \mathbf{T}, \nabla_{\mathbf{T}} \mathcal{L}_{\text{LC}} \rangle + \frac{1}{\gamma_k} \text{KL}(\mathbf{T} \| \mathbf{T}_k). \quad (16)$$

By construction, the triple $(\mathbf{Q}_{k+1}, \mathbf{R}_{k+1}, \mathbf{T}_{k+1}) \in \mathcal{C}_2$. However, because the set \mathcal{C}_1 only constrains the first variable block $\mathbf{Q}_{k+1}, \mathbf{R}_{k+1}$, we have that $(\mathbf{Q}_{k+1}, \mathbf{R}_{k+1}, \mathbf{T}_{k+1}) \in \mathcal{C}_2$, and hence $(\mathbf{Q}_{k+1}, \mathbf{R}_{k+1}, \mathbf{T}_{k+1}) \in \mathcal{C}$. Thus, each iteration produces a feasible triple $(\mathbf{Q}_{k+1}, \mathbf{R}_{k+1}, \mathbf{T}_{k+1})$ through a pair of locally optimal block updates.

The LOT optimization problem Scetbon et al. (2021) For comparison, recall the optimization problem that is solved in the LOT framework:

$$\min_{(\mathbf{Q}, \mathbf{R}, \mathbf{g}) \in \tilde{\mathcal{C}}} \mathcal{L}_{\text{LOT}}, \quad (17)$$

where the objective function \mathcal{L}_{LOT} is

$$\mathcal{L}_{\text{LOT}} := \langle \mathbf{C}, \mathbf{Q} \text{diag}(1/\mathbf{g}) \mathbf{R}^\top \rangle \quad (18)$$

and where $\tilde{\mathcal{C}} = \tilde{\mathcal{C}}_1 \cap \tilde{\mathcal{C}}_2$ with:

$$\begin{aligned} \tilde{\mathcal{C}}_1 &:= \{(\mathbf{Q}, \mathbf{R}, \mathbf{g}) \in \mathbb{R}_+^{n \times r} \times \mathbb{R}_+^{m \times r} \times (\mathbb{R}_+^*)^r : \mathbf{Q}\mathbf{1}_r = \mathbf{a}, \mathbf{R}\mathbf{1}_r = \mathbf{b}\} \\ \tilde{\mathcal{C}}_2 &:= \{(\mathbf{Q}, \mathbf{R}, \mathbf{g}) \in \mathbb{R}_+^{n \times r} \times \mathbb{R}_+^{m \times r} \times \mathbb{R}_+^r : \mathbf{Q}^\top \mathbf{1}_n = \mathbf{g} = \mathbf{R}^\top \mathbf{1}_m\}. \end{aligned}$$

Here, there are also three optimization variables $(\mathbf{Q}, \mathbf{R}, \mathbf{g})$, but they are updated *together* in each iteration of LOT. That is, given feasible $(\mathbf{Q}_k, \mathbf{R}_k, \mathbf{g}_k) \in \tilde{\mathcal{C}}$, an iteration of LOT updates this triple via

$$(\mathbf{Q}_{k+1}, \mathbf{R}_{k+1}, \mathbf{g}_{k+1}) := \arg \min_{(\mathbf{Q}, \mathbf{R}, \mathbf{g}) \in \tilde{\mathcal{C}}} \langle (\mathbf{Q}, \mathbf{R}, \mathbf{g}), \nabla_{(\mathbf{Q}, \mathbf{R}, \mathbf{g})} \mathcal{L}_{\text{LOT}} \rangle + \frac{1}{\gamma_k} \text{KL}((\mathbf{Q}, \mathbf{R}, \mathbf{g}) \| (\mathbf{Q}_k, \mathbf{R}_k, \mathbf{g}_k)),$$

where $(\gamma_k)_{k \geq 0}$ is again a positive sequence of stepsizes. Scetbon et al. (2021) then compute the unconstrained argmin across all variables to yield a set of unconstrained kernels $(\mathbf{K}_Q, \mathbf{K}_R, \mathbf{k}_g)$, using Dykstra to jointly project the unconstrained update onto the intersection $\tilde{\mathcal{C}}$ of the constraint sets.

OT subroutine in FRLC To see why we do not need Dykstra in the FRLC scheme, observe that the update (15) of variables (\mathbf{Q}, \mathbf{R}) can be equivalently expressed as

$$(\mathbf{Q}_{k+1}, \mathbf{R}_{k+1}) \leftarrow \arg \min_{\mathbf{Q}, \mathbf{R} : (\mathbf{Q}, \mathbf{R}, \mathbf{T}_k) \in \mathcal{C}_1} \langle (\mathbf{Q}, \mathbf{R}), \nabla_{(\mathbf{Q}, \mathbf{R})} \mathcal{L}_{\text{LC}} \rangle + \frac{1}{\gamma_k} \text{KL}((\mathbf{Q}, \mathbf{R}) \| (\mathbf{Q}_k, \mathbf{R}_k)),$$

Thus, even though the pair (\mathbf{Q}, \mathbf{R}) is being updated at this step, solving for $(\mathbf{Q}_{k+1}, \mathbf{R}_{k+1})$ above is equivalent to updating each individually because \mathbf{Q} and \mathbf{R} do not share an inner marginal:

$$\mathbf{Q}_{k+1} \leftarrow \arg \min_{\mathbf{Q} : \mathbf{Q}\mathbf{1}_r = \mathbf{a}} \langle \mathbf{Q}, \nabla_{\mathbf{Q}} \mathcal{L}_{\text{LC}} \rangle + \frac{1}{\gamma_k} \text{KL}(\mathbf{Q} \| \mathbf{Q}_k) \quad (19)$$

$$\mathbf{R}_{k+1} \leftarrow \arg \min_{\mathbf{R}: \mathbf{R}\mathbf{1}_r = \mathbf{b}} \langle \mathbf{R}, \nabla_{\mathbf{R}} \mathcal{L}_{\text{LC}} \rangle + \frac{1}{\gamma_k} \text{KL}(\mathbf{R} \| \mathbf{R}_k) \quad (20)$$

We choose to add the regularization $\tau \text{KL}((\mathbf{Q}^{\text{T}} \mathbf{1}_n, \mathbf{R}^{\text{T}} \mathbf{1}_m) \| (\mathbf{Q}_k^{\text{T}} \mathbf{1}_n, \mathbf{R}_k^{\text{T}} \mathbf{1}_m))$ to turn each update here into an entropy-regularized semi-relaxed optimal transport problem, and to ensure β -smoothness:

$$\mathbf{Q}_{k+1} \leftarrow \arg \min_{\mathbf{Q}: \mathbf{Q}\mathbf{1}_r = \mathbf{a}} \langle \mathbf{Q}, \nabla_{\mathbf{Q}} \mathcal{L}_{\text{LC}} \rangle + \frac{1}{\gamma_k} \text{KL}(\mathbf{Q} \| \mathbf{Q}_k) + \tau \text{KL}(\mathbf{Q}^{\text{T}} \mathbf{1}_n \| \mathbf{Q}_k^{\text{T}} \mathbf{1}_n) \quad (21)$$

$$\mathbf{R}_{k+1} \leftarrow \arg \min_{\mathbf{R}: \mathbf{R}\mathbf{1}_r = \mathbf{b}} \langle \mathbf{R}, \nabla_{\mathbf{R}} \mathcal{L}_{\text{LC}} \rangle + \frac{1}{\gamma_k} \text{KL}(\mathbf{R} \| \mathbf{R}_k) + \tau \text{KL}(\mathbf{R}^{\text{T}} \mathbf{1}_m \| \mathbf{R}_k^{\text{T}} \mathbf{1}_m) \quad (22)$$

After updating \mathbf{Q} and \mathbf{R} , the update on \mathbf{T} then follows a similar form:

$$\mathbf{T}_{k+1} \leftarrow \arg \min_{\mathbf{T}: (\mathbf{Q}_{k+1}, \mathbf{R}_{k+1}, \mathbf{T}) \in \mathcal{C}_2} \langle \mathbf{T}, \nabla_{\mathbf{T}} \mathcal{L}_{\text{LC}} \rangle + \frac{1}{\gamma_k} \text{KL}(\mathbf{T} \| \mathbf{T}_k), \quad (23)$$

leading to balanced constraints on \mathbf{T} of the form $\mathbf{T}^{\text{T}} \mathbf{1}_r = \mathbf{Q}_{k+1}^{\text{T}} \mathbf{1}_n$ and $\mathbf{T} \mathbf{1}_r = \mathbf{R}_{k+1}^{\text{T}} \mathbf{1}_m$, allowing the problem to be solved by Sinkhorn.

Importantly, there are no constraints in \mathcal{C}_1 involving both \mathbf{Q} and \mathbf{R} , which is what allows the optimization to split in this way. If there were such constraints, we would have needed to use Dykstra to update the pair (\mathbf{Q}, \mathbf{R}) . Because our update scheme is equivalent to the three updates of individual variables given in (21), (22), (23), we can solve for each update using optimal transport. As \mathbf{Q} and \mathbf{R} are not required to match the inner marginals exactly, the OT problems associated to \mathbf{Q} and \mathbf{R} are semi-relaxed by construction.

The separation of our block updates into a step where $(\mathbf{Q}, \mathbf{R}, \mathbf{T}_k) \in \mathcal{C}_1$ and $(\mathbf{Q}_{k+1}, \mathbf{R}_{k+1}, \mathbf{T}) \in \mathcal{C}_2$ allows us to entirely remove the optimization over inner marginals \mathbf{g}_Q and \mathbf{g}_R , as done in all previous works on low-rank optimal transport which optimize \mathbf{g} explicitly as both a variable and a constraint of the optimization Scetbon et al. (2021, 2023); Scetbon & Cuturi (2022). If one were to introduce an extended loss in the style of previous works which adds \mathbf{g}_Q and \mathbf{g}_R as variables in the form $\mathcal{H}(\mathbf{Q}, \mathbf{R}, \mathbf{T}, \mathbf{g}_Q, \mathbf{g}_R) = \langle \mathbf{Q} \text{diag}(1/\mathbf{g}_Q) \mathbf{T} \text{diag}(1/\mathbf{g}_R) \mathbf{R}^{\text{T}}, \mathbf{C} \rangle_F$, one observes an equivalence to simply taking a semi-relaxed projection.

Lemma A.1. *Define the function $\mathcal{H}(\mathbf{Q}, \mathbf{R}, \mathbf{T}, \mathbf{g}_Q, \mathbf{g}_R) = \langle \mathbf{C}, \mathbf{Q} \text{diag}(1/\mathbf{g}_Q) \mathbf{T} \text{diag}(1/\mathbf{g}_R) \mathbf{R}^{\text{T}} \rangle_F$ and let $\mathcal{L}_{\text{LC}}(\mathbf{Q}, \mathbf{R}, \mathbf{T})$ be as in (13). One has the following equivalence:*

$$\min_{\mathbf{g}_R \in \Delta_r, \mathbf{g}_Q \in \Delta_r, \mathbf{Q} \in \Pi_{\mathbf{a}, \mathbf{g}_Q}, \mathbf{R} \in \Pi_{\mathbf{b}, \mathbf{g}_R}} \mathcal{H}(\mathbf{Q}, \mathbf{R}, \mathbf{T}_k, \mathbf{g}_Q, \mathbf{g}_R) = \min_{(\mathbf{Q}, \mathbf{R}, \mathbf{T}_k) \in \mathcal{C}_1} \mathcal{L}_{\text{LC}}(\mathbf{Q}, \mathbf{R}, \mathbf{T}_k) \quad (24)$$

Thus the semi-relaxed projections yield locally optimal inner marginals.

Proof. To see why this is true, notice that, so long as the outer marginals are tightly satisfied $\mathbf{Q}\mathbf{1}_r = \mathbf{a}$ and $\mathbf{R}\mathbf{1}_r = \mathbf{b}$ for $\mathbf{Q} \geq \mathbf{0}_{n \times r}$, $\mathbf{R} \geq \mathbf{0}_{m \times r}$, one has

$$\sum_i \mathbf{g}_{R,i} = \sum_i \langle \mathbf{R}_{i,\cdot}^{\text{T}}, \mathbf{1}_m \rangle = \sum_i \langle \mathbf{R}_{\cdot,i}, \mathbf{1}_m \rangle = \sum_{\ell} \sum_i \mathbf{R}_{\ell,i} = \sum_{\ell} \mathbf{b}_{\ell} = 1$$

and

$$\sum_i \mathbf{g}_{Q,i} = \sum_i \langle \mathbf{Q}_{i,\cdot}^{\text{T}}, \mathbf{1}_n \rangle = \sum_i \langle \mathbf{Q}_{\cdot,i}, \mathbf{1}_n \rangle = \sum_{\ell} \sum_i \mathbf{Q}_{\ell,i} = \sum_{\ell} \mathbf{a}_{\ell} = 1$$

Therefore, the inner marginals $\mathbf{g}_Q^* = (\mathbf{Q}^*)^{\text{T}} \mathbf{1}_n$ and $\mathbf{g}_R^* = (\mathbf{R}^*)^{\text{T}} \mathbf{1}_m$ induced by the optimal \mathbf{Q}^* , \mathbf{R}^* of the optimization problem on the right hand side satisfy the constraints $\mathbf{g}_R^* \in \Delta_r$, $\mathbf{g}_Q^* \in \Delta_r$ on the left hand side, so the two minimums coincide. \square

This implies that an extra optimization for \mathbf{g}_Q and \mathbf{g}_R is unnecessary in a coordinate update which alternates (\mathbf{Q}, \mathbf{R}) and \mathbf{T} . If we did not perform a block-update, in the form of standard MD in the objective described above, we would encounter some difficulty. In particular \mathbf{g}_Q and \mathbf{g}_R would be optimized with the constraint $\mathbf{g}_Q \in \Delta_r$ and $\mathbf{g}_R \in \Delta_r$, and would concurrently constrain all of the other variables as $\mathbf{Q}^{\text{T}} \mathbf{1}_n = \mathbf{g}_Q$, $\mathbf{T} \mathbf{1}_r = \mathbf{g}_Q$, and $\mathbf{R}^{\text{T}} \mathbf{1}_m = \mathbf{g}_R$, $\mathbf{T}^{\text{T}} \mathbf{1}_r = \mathbf{g}_R$.

B Block-Coordinate steps for the OT sub-problems

We use a latent non-diagonal coupling instead of an inner diagonal coupling $\text{diag}(\mathbf{g})$ of the form of Forrow et al. (2019). This allows us to loosen the constraint that the inner marginals have to be joined by a common coupling $\mathbf{Q}^\top \mathbf{1}_n = \mathbf{R}^\top \mathbf{1}_m = \mathbf{g}$. The fundamental advantage of this choice is that we can decouple the convex-optimization problem for $(\mathbf{Q}, \mathbf{R}, \mathbf{T})$ entirely. One can simply solve for the optimal \mathbf{Q} and \mathbf{R} *independently*, yield the associated inner marginals for each $\mathbf{Q}^\top \mathbf{1}_n = \mathbf{g}_Q$ and $\mathbf{R}^\top \mathbf{1}_m = \mathbf{g}_R$, and then find the optimal \mathbf{T} which links the two. This link is provided by the aforementioned form of the problem, where:

$$\mathbf{P} = \mathbf{Q}\mathbf{X}\mathbf{R}^\top$$

For \mathbf{Q}, \mathbf{R} in either the appropriate set of couplings or a relaxation thereof (which we will describe shortly). \mathbf{X} is related to \mathbf{T} by:

$$\mathbf{X} = \text{diag}(1/\mathbf{g}_Q)\mathbf{T}\text{diag}(1/\mathbf{g}_R)$$

And, $\mathbf{T} \in \Pi_{\mathbf{g}_Q, \mathbf{g}_R}$ consistently for all cases. As the semi-relaxed case is intermediate between fully-relaxed and balanced, it has ideas which generalize to both directly. As such, we use it as the leading example again. As in Scetbon et al. (2021), we take proximal-steps of the form:

$$\min_{\zeta} \langle \nabla \mathcal{L}(\zeta) |_{\zeta^k}, \zeta \rangle_F + \frac{1}{\gamma_k} \text{KL}(\zeta \| \mathbf{K}^{(k)})$$

Where these steps are now in block-wise fashion on (\mathbf{Q}, \mathbf{R}) and \mathbf{T} , rather than joint. One may identify for each sub-factor in (\mathbf{Q}, \mathbf{R}) and \mathbf{T} a linearized gradient as before, which yields a set of objectives which each solve an independent optimal-transport for the sub-factors. In particular, we have that:

$$\begin{aligned} \langle \mathbf{Q}\mathbf{X}\mathbf{R}^\top, \mathbf{C} \rangle_F &= \text{Tr} [\mathbf{Q}\mathbf{X}\mathbf{R}^\top \mathbf{C}^\top] = \langle \mathbf{C}\mathbf{R}\mathbf{X}^\top, \mathbf{Q} \rangle_F \\ \langle \mathbf{Q}\mathbf{X}\mathbf{R}^\top, \mathbf{C} \rangle_F &= \text{Tr} [\mathbf{Q}\mathbf{X}\mathbf{R}^\top \mathbf{C}^\top] = \langle \mathbf{C}^\top \mathbf{Q}\mathbf{X}, \mathbf{R} \rangle_F \\ \langle \mathbf{Q}\mathbf{X}\mathbf{R}^\top, \mathbf{C} \rangle_F &= \text{Tr} [\mathbf{R}^\top \mathbf{C}^\top \mathbf{Q}\mathbf{X}] = \langle \mathbf{Q}^\top \mathbf{C}\mathbf{R}, \mathbf{X} \rangle_F \end{aligned}$$

A linearization in the left-slot of the inner product as $\langle \mathbf{Q}, \mathbf{C}\mathbf{R}\mathbf{X}(\mathbf{Q}_k)^\top \rangle := \langle \mathbf{Q}, \mathbf{C}\mathbf{R}\mathbf{X}^\top \rangle$ or $\langle \mathbf{C}^\top \mathbf{Q}\mathbf{X}(\mathbf{R}_k), \mathbf{R} \rangle_F := \langle \mathbf{C}^\top \mathbf{Q}\mathbf{X}, \mathbf{R} \rangle_F$ is common practice for quadratic problems. In this case, the directional derivative of \mathbf{Q} and \mathbf{R} in the matrix-direction \mathbf{V} are respectively:

$$\begin{aligned} D\langle \mathbf{C}\mathbf{R}\mathbf{X}^\top, \mathbf{Q} \rangle_F \circ (\mathbf{V}) &= \langle \mathbf{C}\mathbf{R}\mathbf{X}^\top, \mathbf{V} \rangle_F \implies \nabla_{\mathbf{Q}} \mathcal{L} = \mathbf{C}\mathbf{R}\mathbf{X}^\top \\ D\langle \mathbf{C}^\top \mathbf{Q}\mathbf{X}, \mathbf{R} \rangle_F \circ (\mathbf{V}) &= \langle \mathbf{C}^\top \mathbf{Q}\mathbf{X}, \mathbf{V} \rangle_F \implies \nabla_{\mathbf{R}} \mathcal{L} = \mathbf{C}^\top \mathbf{Q}\mathbf{X} \end{aligned}$$

Without this linearization assumption on \mathbf{X} , the full gradient may be evaluated as well. In particular, we note that for $\text{diag}^{-1}(\cdot)$ the matrix-to-vector extraction of the diagonal, the directional derivative on \mathbf{Q} is:

$$\begin{aligned} D\langle \mathbf{C}, \mathbf{Q}\mathbf{X}\mathbf{R}^\top \rangle_F \circ \mathbf{V} &= \langle \mathbf{C}\mathbf{R}\mathbf{X}^\top, \mathbf{V} \rangle_F + \langle \mathbf{C}, \mathbf{Q}D\mathbf{X}^\top \circ (\mathbf{V})\mathbf{R}^\top \rangle_F \\ &= \langle \mathbf{C}\mathbf{R}\mathbf{X}^\top - \mathbf{1}_n \text{diag}^{-1}((\mathbf{C}\mathbf{R}\mathbf{X}^\top)^\top \mathbf{Q} \text{diag}(1/\mathbf{g}_Q))^\top, \mathbf{V} \rangle_F \end{aligned}$$

Thus, without the linearization assumption one may use product rule on \mathbf{X} as an implicit function of \mathbf{Q} (resp. \mathbf{R}) to take the total derivative:

$$\nabla_{\mathbf{Q}} \mathcal{L}_{\text{FRLC}} = \mathbf{C}\mathbf{R}\mathbf{X}^\top - \mathbf{1}_n \text{diag}^{-1}((\mathbf{C}\mathbf{R}\mathbf{X}^\top)^\top \mathbf{Q} \text{diag}(1/\mathbf{g}_Q))^\top$$

Likewise for \mathbf{R} ,

$$\begin{aligned} D\langle \mathbf{C}, \mathbf{Q}\mathbf{X}\mathbf{R}^\top \rangle_F \circ \mathbf{V} &= \langle \mathbf{C}^\top \mathbf{Q}\mathbf{X}, \mathbf{V} \rangle_F + \langle \mathbf{C}, \mathbf{Q}D\mathbf{X} \circ (\mathbf{V})\mathbf{R}^\top \rangle_F \\ &= \langle \mathbf{C}^\top \mathbf{Q}\mathbf{X} - \mathbf{1}_m \text{diag}^{-1}(\text{diag}(1/\mathbf{g}_R)\mathbf{R}^\top \mathbf{C}^\top \mathbf{Q}\mathbf{X})^\top, \mathbf{V} \rangle_F, \end{aligned}$$

and so

$$\nabla_{\mathbf{R}} \mathcal{L}_{\text{FRLC}} = \mathbf{C}^\top \mathbf{Q}\mathbf{X} - \mathbf{1}_m \text{diag}^{-1}(\text{diag}(1/\mathbf{g}_R)\mathbf{R}^\top \mathbf{C}^\top \mathbf{Q}\mathbf{X})^\top$$

Both the W and GW problem have these rank-one perturbations of the gradient from the derivative with respect to \mathbf{X} .

Lastly, owing to the block-coordinate updates which fix (\mathbf{Q}, \mathbf{R}) preceding the update on \mathbf{T} , the derivative on \mathbf{T} follows directly by chain rule on \mathbf{X} :

$$\begin{aligned} D\langle \mathbf{Q}^\top \mathbf{C} \mathbf{R}, \mathbf{X} \rangle_{F \circ} (\mathbf{V}) &= \langle \mathbf{Q}^\top \mathbf{C} \mathbf{R}, \mathbf{V} \rangle_F \implies \nabla_{\mathbf{X}} \mathcal{L} = \mathbf{Q}^\top \mathbf{C} \mathbf{R} \\ \nabla_{\mathbf{T}} \mathcal{L} &= \text{diag}(1/g_{\mathbf{Q}}) \mathbf{Q}^\top \mathbf{C} \mathbf{R} \text{diag}(1/g_{\mathbf{R}}) \end{aligned}$$

Let (γ_k) be a sequence of step sizes and consider the first-order conditions required for the proximal step as before. From these we have the updated proximal-step updates:

$$\begin{aligned} \mathbf{K}_{\mathbf{Q}}^{(k)} &\leftarrow \min_{\mathbf{Q}} \langle \mathbf{C} \mathbf{R} \mathbf{X}^\top, \mathbf{Q} \rangle_F + \frac{1}{\gamma_k} \text{KL}(\mathbf{Q}_k \| \mathbf{K}_{\mathbf{Q}}^{(k)}) \\ \mathbf{K}_{\mathbf{R}}^{(k)} &\leftarrow \min_{\mathbf{R}} \langle \mathbf{C}^\top \mathbf{Q} \mathbf{X}, \mathbf{R} \rangle_F + \frac{1}{\gamma_k} \text{KL}(\mathbf{R}_k \| \mathbf{K}_{\mathbf{R}}^{(k)}) \\ \mathbf{K}_{\mathbf{T}}^{(k)} &\leftarrow \min_{\mathbf{T}} \langle \mathbf{Q}^\top \mathbf{C} \mathbf{R}, \mathbf{X} \rangle_F + \frac{1}{\gamma_k} \text{KL}(\mathbf{T}_k \| \mathbf{K}_{\mathbf{T}}^{(k)}), \end{aligned}$$

with kernels $\mathbf{K}_{\zeta_j}^{(k)}$, for $j = 1, 2, 3$ given by

$$\begin{aligned} \mathbf{K}_{\mathbf{Q}}^{(k)} &:= \mathbf{Q}_k \odot \exp(-\gamma_k \mathbf{C} \mathbf{R}_k \mathbf{X}_k^\top) \\ \mathbf{K}_{\mathbf{R}}^{(k)} &:= \mathbf{R}_k \odot \exp(-\gamma_k \mathbf{C}^\top \mathbf{Q}_k \mathbf{X}_k) \\ \mathbf{K}_{\mathbf{T}}^{(k)} &:= \mathbf{T}_k \odot \exp(-\gamma_k \text{diag}(g_{\mathbf{Q}}^{-1}) \mathbf{Q}^\top \mathbf{C} \mathbf{R} \text{diag}(g_{\mathbf{R}}^{-1})) \end{aligned}$$

Or, dropping the linearization assumption on (\mathbf{Q}, \mathbf{R}) the updates are: $\mathbf{K}_{\zeta_j}^{(k)}$, for $j = 1, 2$ given by

$$\begin{aligned} \mathbf{K}_{\mathbf{Q}}^{(k)} &:= \mathbf{Q}_k \odot \exp(-\gamma_k (\mathbf{C} \mathbf{R}_k \mathbf{X}_k^\top - \mathbf{1}_n \text{diag}^{-1}((\mathbf{C} \mathbf{R}_k \mathbf{X}_k^\top)^\top \mathbf{Q}_k \text{diag}(1/g_{\mathbf{Q}_k})))^\top) \\ \mathbf{K}_{\mathbf{R}}^{(k)} &:= \mathbf{R}_k \odot \exp(-\gamma_k (\mathbf{C}^\top \mathbf{Q}_k \mathbf{X}_k - \mathbf{1}_m \text{diag}^{-1}(\text{diag}(1/g_{\mathbf{R}}) \mathbf{R}_k^\top \mathbf{C}^\top \mathbf{Q}_k \mathbf{X}_k)^\top)) \end{aligned}$$

The first projection is onto the set that satisfies the marginal constraint $\mathbf{R} \mathbf{1}_r = \mathbf{b}$. In particular, following the discussion above, one has the coordinate-MD step:

$$\begin{aligned} \min_{(\mathbf{Q}, \mathbf{R}, \mathbf{T})} \quad & \frac{1}{\gamma_k} \text{KL}((\mathbf{Q}, \mathbf{R}, \mathbf{T}) \| (\mathbf{K}_{\mathbf{Q}}, \mathbf{K}_{\mathbf{R}}, \mathbf{K}_{\mathbf{T}})) + \tau \text{KL}(\mathbf{Q} \mathbf{1}_r \| \mathbf{a}) \\ \text{s.t.} \quad & \mathbf{R} \mathbf{1}_r = \mathbf{b} \end{aligned} \quad (25)$$

As before, there is no difference from the previous case, where one takes the unconstrained projection $\mathbf{R} = \text{diag}(\mathbf{b}/\mathbf{K}_{\mathbf{R}} \mathbf{1}_r) \mathbf{K}_{\mathbf{R}}$. To generalize this, we also consider adding a soft-constraint on the inner marginal of \mathbf{R} to be near that of the previous iteration. In particular, we consider the problem:

$$\begin{aligned} \min_{(\mathbf{Q}, \mathbf{R}, \mathbf{T})} \quad & \frac{1}{\gamma_k} \text{KL}((\mathbf{Q}, \mathbf{R}, \mathbf{T}) \| (\mathbf{K}_{\mathbf{Q}}, \mathbf{K}_{\mathbf{R}}, \mathbf{K}_{\mathbf{T}})) \\ & + \tau \text{KL}(\mathbf{Q} \mathbf{1}_r \| \mathbf{a}) + \tau \text{KL}(\mathbf{R}^\top \mathbf{1}_m \| \mathbf{g}_{\mathbf{R}}^{(k-1)} \equiv \mathbf{R}_{k-1}^\top \mathbf{1}_m) \\ \text{s.t.} \quad & \mathbf{R} \mathbf{1}_r = \mathbf{b} \end{aligned} \quad (26)$$

Which yields the relaxed solution of $\mathbf{R} = \text{SR}^{\mathbf{R}}\text{-projection}(\mathbf{K}_{\mathbf{R}}, \gamma_k, \tau, \mathbf{b}, \mathbf{g}_{\mathbf{R}}^{(k-1)})$ which generalizes the original projection and is equivalent to it for $\tau = 0$. This regularization is essential, as it ensures β -smoothness of the objective. For \mathbf{Q} , as the constraint on \mathbf{g} is fully relaxed, the Lagrange multiplier $\lambda_1 = \mathbf{0}$ entirely, such that the problem 40 now becomes fully-unconstrained:

$$\inf_{\mathbf{Q}} \left(\frac{1}{\gamma_k} \text{KL}(\mathbf{Q} \| \mathbf{K}_{\mathbf{Q}}) + \tau \text{KL}(\mathbf{Q} \mathbf{1}_r \| \mathbf{a}) \right) \quad (27)$$

To generalize this solution, we again consider adding a soft-regularization on the inner marginal of \mathbf{Q} , where we consider the alternate problem:

$$\inf_{\mathbf{Q}} \left(\frac{1}{\gamma_k} \text{KL}(\mathbf{Q} \| \mathbf{K}_{\mathbf{Q}}) + \tau \text{KL}(\mathbf{Q} \mathbf{1}_r \| \mathbf{a}) + \tau \text{KL}(\mathbf{Q}^\top \mathbf{1}_n \| \mathbf{g}_{\mathbf{Q}}^{(k-1)} \equiv \mathbf{Q}_{k-1}^\top \mathbf{1}_n) \right) \quad (28)$$

Which trivially recovers the original for $\tau = 0$. This form has a solution given by an unbalanced optimal transport with kernel \mathbf{K}_Q . We see these two, convex problems in sequence give *independent* optimal solutions for \mathbf{Q} and \mathbf{R} as the two matrices are not required to share an inner marginal. The last step is to link them via \mathbf{T} , corresponding to the projection of $(\mathbf{Q}_k, \mathbf{R}_k, \mathbf{T})$ onto the set of valid rank- r couplings $\min_{(\mathbf{Q}_k, \mathbf{R}_k, \mathbf{T}) \in \mathcal{C}} \mathcal{L}_{\text{LC}}(\mathbf{Q}_k, \mathbf{R}_k, \mathbf{T}) := \min_{\mathbf{T} \in \Pi(\mathbf{g}_Q, \mathbf{g}_R)} \mathcal{L}_{\text{LC}}(\mathbf{Q}_k, \mathbf{R}_k, \mathbf{T})$. We verify in C that if $\mathbf{T} \in \Pi_{\mathbf{g}_Q = \mathbf{Q}^T \mathbf{1}_n, \mathbf{g}_R = \mathbf{R}^T \mathbf{1}_m}$ then $\mathbf{P} \in \Pi_{\mathbf{a}, \mathbf{b}}$. As such, one does not require any projection onto the intersection of convex sets, as done in Scetbon et al. (2021, 2023) via the Dykstra projection algorithm Dykstra (1983). Alternating a coordinate-MD step in (\mathbf{Q}, \mathbf{R}) and a step on \mathbf{T} , one not only minimizes the objective in an alternating fashion but remains in the feasible set without the need for projection algorithms beyond Sinkhorn. As such, the final linking step is done via:

$$\begin{aligned} \min_{\mathbf{T}} \quad & \frac{1}{\gamma_k} \text{KL}(\mathbf{T} \| \mathbf{K}_T) \\ \text{s.t.} \quad & \mathbf{T} \mathbf{1}_r = \mathbf{g}_Q, \mathbf{T}^T \mathbf{1}_r = \mathbf{g}_R \end{aligned} \quad (29)$$

This formulation amounts to solving a *balanced* Sinkhorn problem on \mathbf{T} with respect to the proximal step kernel matrix.

Algorithm 2 SR^R-projection (*semi-relaxed OT, right marginal relaxed*)

Input $\mathbf{K}, \gamma, \tau, \mathbf{a}, \mathbf{b}, \delta$
 $\mathbf{u} \leftarrow \mathbf{1}_n$
 $\mathbf{v} \leftarrow \mathbf{1}_r$
repeat
 $\tilde{\mathbf{u}} \leftarrow \mathbf{u}$
 $\tilde{\mathbf{v}} \leftarrow \mathbf{v}$
 $\mathbf{u} \leftarrow (\mathbf{a} / \mathbf{K} \mathbf{v})$
 $\mathbf{v} \leftarrow (\mathbf{b} / \mathbf{K}^T \mathbf{u})^{\tau / (\tau + \gamma^{-1})}$
until $\gamma^{-1} \max\{\|\log \tilde{\mathbf{u}} / \mathbf{u}\|_\infty, \|\log \tilde{\mathbf{v}} / \mathbf{v}\|_\infty\} < \delta$
return $\text{diag}(\mathbf{u}) \mathbf{K} \text{diag}(\mathbf{v})$

Algorithm 3 U-projection (*unbalanced OT*)

Input $\mathbf{K}, \gamma, \tau, \mathbf{a}, \mathbf{b}, \delta$
 $\mathbf{u} \leftarrow \mathbf{1}_n$
 $\mathbf{v} \leftarrow \mathbf{1}_r$
repeat
 $\tilde{\mathbf{u}} \leftarrow \mathbf{u}$
 $\tilde{\mathbf{v}} \leftarrow \mathbf{v}$
 $\mathbf{u} \leftarrow (\mathbf{a} / \mathbf{K} \mathbf{v})^{\tau / (\tau + \gamma^{-1})}$
 $\mathbf{v} \leftarrow (\mathbf{b} / \mathbf{K}^T \mathbf{u})^{\tau / (\tau + \gamma^{-1})}$
until $\gamma^{-1} \max\{\|\log \tilde{\mathbf{u}} / \mathbf{u}\|_\infty, \|\log \tilde{\mathbf{v}} / \mathbf{v}\|_\infty\} < \delta$
return $\text{diag}(\mathbf{u}) \mathbf{K} \text{diag}(\mathbf{v})$

C The Latent Coupling Matrix

To solve the balanced form (and generalize the principle to the relaxed problems), we consider an alternative parametrization of the inner matrix. In particular, previous works Scetbon et al. (2021, 2023) consider $\text{diag}(1/\mathbf{g})$ to be the inner matrix, with marginals $\mathbf{Q}^T \mathbf{1}_n = \mathbf{R}^T \mathbf{1}_m = \mathbf{g}$ to ensure that the outer conditions $\mathbf{Q} \mathbf{1}_r = \mathbf{a}$ and $\mathbf{R} \mathbf{1}_r = \mathbf{b}$ hold. We instead relax the constraint that $\mathbf{Q}^T \mathbf{1}_n$ and $\mathbf{R}^T \mathbf{1}_m$ are equal, by allowing $\mathbf{Q}^T \mathbf{1}_n = \mathbf{g}_Q$ and $\mathbf{R}^T \mathbf{1}_m = \mathbf{g}_R$ to vary arbitrarily, and considering a non-diagonal inner matrix $\mathbf{X} \in \mathbb{R}^{r \times r}$ in the place of $\text{diag}(1/\mathbf{g})$ where we have the conditions:

$$\mathbf{X} \mathbf{g}_R = \mathbf{X} \mathbf{R}^T \mathbf{1}_m = \mathbf{1}_r$$

And

$$\mathbf{X}^T \mathbf{g}_Q = \mathbf{X} \mathbf{Q}^T \mathbf{1}_n = \mathbf{1}_r$$

Thus, if one considers the coupling formed as $P_r = QXR^T$, one maintains that $P_r \in \Pi_{a,b}$ from the condition of \mathcal{C}_1 , defined in the balanced case as in (14). We clearly have that:

$$P_r \mathbf{1}_m = QXR^T \mathbf{1}_m = Q\mathbf{1}_r = \mathbf{a}$$

And:

$$P_r^T \mathbf{1}_n = RX^T Q\mathbf{1}_n = R\mathbf{1}_r = \mathbf{b}$$

We first consider two approaches to optimizing for such a X , given that it is not a coupling. First, we consider the appropriate proximal step for X

$$\min_{\zeta} \langle \nabla \mathcal{L}(\zeta) |_{\xi_k}, \zeta \rangle_F + \frac{1}{\gamma_k} \text{KL}(\zeta \| \mathbf{K}^{(k)})$$

We first note that the gradient of our loss with respect to X is given as $Q^T C R$. If one supposes that X is invertible with $X^{-1} = T$, we have that for any such T :

$$X^{-1} \mathbf{1}_r = T \mathbf{1}_r = \mathbf{g}_R$$

And

$$X^{-T} \mathbf{1}_r = T^T \mathbf{1}_r = \mathbf{g}_Q$$

This implies that this inverse matrix T is a coupling such that $T \in \Pi_{\mathbf{g}_R, \mathbf{g}_Q}$, which also suggests one might be able to update it using Sinkhorn. In fact, being a density in $\mathbb{R}_+^{r \times r}$ it represents a transition matrix between the latent r -dimensional variables. In particular, writing the proximal step in full, we have:

$$\min_T \langle Q^T C R, T^{-1} \rangle_F + \frac{1}{\gamma_k} \text{KL}(T_k \| \mathbf{K}_T^{(k)})$$

Noting the derivative $D(X^{-1}) \circ V = -X^{-1} V X^{-1}$, we have from the first-order condition that:

$$-T^{-T} Q^T C R T^{-T} + \frac{1}{\gamma_k} \log \left[\frac{T_k}{\mathbf{K}_T^{(k)}} \right] = \mathbf{0}$$

This implies the kernel matrix update:

$$\mathbf{K}_T^{(k)} = T_k \odot \exp\{+\gamma_k T^{-T} Q^T C R T^{-T}\}$$

Where one then takes the Sinkhorn projection J onto the set $\Pi_{\mathbf{g}_R, \mathbf{g}_Q}$ as $\mathcal{P}_{\Pi_{\mathbf{g}_R, \mathbf{g}_Q}}(\mathbf{K}_T^{(k)})$ using the Sinkhorn algorithm Cuturi (2013b). However, a more stable and inversion-free update exists which ensures X remains positive by a diagonal re-scaling in the form introduced by Lin et al. (2021). In particular, if one takes $X = \text{diag}(1/\mathbf{g}_Q) T \text{diag}(1/\mathbf{g}_R)$, then

$$X \mathbf{g}_R = \text{diag}(1/\mathbf{g}_Q) T \text{diag}(1/\mathbf{g}_R) \mathbf{g}_R = \text{diag}(1/\mathbf{g}_Q) T \mathbf{1}_r = \mathbf{1}_r$$

and likewise

$$X^T \mathbf{g}_Q = \text{diag}(1/\mathbf{g}_R) T^T \text{diag}(1/\mathbf{g}_Q) \mathbf{g}_Q = \text{diag}(1/\mathbf{g}_R) T^T \mathbf{1}_r = \mathbf{1}_r$$

so that X necessarily satisfies $X \mathbf{g}_R = \mathbf{1}_r$ and $X^T \mathbf{g}_Q = \mathbf{1}_r$. Thus $T \mathbf{1}_r = \mathbf{g}_Q$ and $T^T \mathbf{1}_r = \mathbf{g}_R$ and $T \in \Pi_{\mathbf{g}_Q, \mathbf{g}_R}$. With analogous reasoning to before, one has a step for the coupling T in the form:

$$\min_T \langle Q^T C R, \text{diag}(1/\mathbf{g}_Q) T \text{diag}(1/\mathbf{g}_R) \rangle_F + \frac{1}{\gamma_k} \text{KL}(T_k \| \mathbf{K}_T^{(k)})$$

Which yields the kernel matrix:

$$\mathbf{K}_T^{(k)} = T_k \odot \exp\{-\gamma_k \text{diag}(\mathbf{g}_Q)^{-1} Q^T C R \text{diag}(\mathbf{g}_R)^{-1}\}$$

Which is likewise projected onto $\Pi_{\mathbf{g}_Q, \mathbf{g}_R}$ using the Sinkhorn algorithm. From this $T \in \Pi_{\mathbf{g}_Q, \mathbf{g}_R}$, one takes $X = \text{diag}(1/\mathbf{g}_Q) T \text{diag}(1/\mathbf{g}_R)$ as the inner matrix that corresponds the unequal marginals \mathbf{g}_Q and \mathbf{g}_R which ensuring $P_r \in \Pi_{a,b}$.

Algorithm 4 FRLC (*General marginal constraint low-rank optimal transport*)

Input $\mathbf{C}, r, r_2, \mathbf{a}, \mathbf{b}, \tau, \tau_2, \gamma, \delta, \varepsilon$
Initialize $\mathbf{g}_Q, \mathbf{g}_R = \frac{1}{r} \mathbf{1}_r, \frac{1}{r_2} \mathbf{1}_{r_2}$
 $\mathbf{Q}_0, \mathbf{R}_0, \mathbf{T}_0 \leftarrow \text{Initialize-Couplings}(\mathbf{a}, \mathbf{b}, \mathbf{g}_Q, \mathbf{g}_R)$
if $r = r_2$ **then**
 $\mathbf{X}_0 \leftarrow \mathbf{T}_0^{-1}$ # Invertible case
else
 $\mathbf{X}_0 \leftarrow \text{diag}(1/\mathbf{Q}_0^T \mathbf{1}_n) \mathbf{T}_0 \text{diag}(1/\mathbf{R}_0^T \mathbf{1}_m)$ # General case
end if
while $\Delta((\mathbf{Q}_k, \mathbf{R}_k, \mathbf{T}_k), (\mathbf{Q}_{k-1}, \mathbf{R}_{k-1}, \mathbf{T}_{k-1})) > \varepsilon$ **do**
 $\nabla_Q \leftarrow \mathbf{C} \mathbf{R}_k \mathbf{X}_k^T - \mathbf{1}_n \text{diag}^{-1}((\mathbf{C} \mathbf{R}_k \mathbf{X}_k^T)^T \mathbf{Q}_k \text{diag}(1/\mathbf{g}_Q))^T$
 $\nabla_R \leftarrow \mathbf{C}^T \mathbf{Q}_k \mathbf{X}_k - \mathbf{1}_m \text{diag}^{-1}(\text{diag}(1/\mathbf{g}_R) \mathbf{R}_k^T \mathbf{C}^T \mathbf{Q}_k \mathbf{X}_k)^T$
 $\gamma_k \leftarrow \gamma / \max\{\|\nabla_Q\|_\infty, \|\nabla_R\|_\infty\}$ # ℓ^∞ -normalization of Scetbon & Cuturi (2022)
 $\mathbf{K}_Q^{(k)}, \mathbf{K}_R^{(k)} \leftarrow \mathbf{Q}_k \odot \exp(-\gamma_k \nabla_Q), \mathbf{R}_k \odot \exp(-\gamma_k \nabla_R)$
 if Balanced **then**
 $\mathbf{Q}_k \leftarrow \text{SR}^R\text{-projection}(\mathbf{K}_Q^{(k)}, \gamma_k, \tau, \mathbf{a}, \mathbf{Q}_{k-1}^T \mathbf{1}_n, \delta)$ # Semi-relaxed OT
 $\mathbf{R}_k \leftarrow \text{SR}^R\text{-projection}(\mathbf{K}_R^{(k)}, \gamma_k, \tau, \mathbf{b}, \mathbf{R}_{k-1}^T \mathbf{1}_m, \delta)$ # Semi-relaxed OT
 else if Unbalanced **then**
 $\mathbf{Q}_k \leftarrow \text{U-projection}(\mathbf{K}_Q^{(k)}, \gamma_k, \tau, \mathbf{a}, \mathbf{Q}_{k-1}^T \mathbf{1}_n, \delta)$ # Unbalanced OT
 $\mathbf{R}_k \leftarrow \text{U-projection}(\mathbf{K}_R^{(k)}, \gamma_k, \tau, \mathbf{b}, \mathbf{R}_{k-1}^T \mathbf{1}_m, \delta)$ # Unbalanced OT
 else if Semi-Relaxed Left **then**
 $\mathbf{Q}_k \leftarrow \text{U-projection}(\mathbf{K}_Q^{(k)}, \gamma_k, \tau, \mathbf{a}, \mathbf{Q}_{k-1}^T \mathbf{1}_n, \delta)$ # Unbalanced OT
 $\mathbf{R}_k \leftarrow \text{SR}^R\text{-projection}(\mathbf{K}_R^{(k)}, \gamma_k, \tau, \mathbf{b}, \mathbf{R}_{k-1}^T \mathbf{1}_m, \delta)$ # Semi-relaxed OT
 else if Semi-Relaxed Right **then**
 $\mathbf{Q}_k \leftarrow \text{SR}^R\text{-projection}(\mathbf{K}_Q^{(k)}, \gamma_k, \tau, \mathbf{a}, \mathbf{Q}_{k-1}^T \mathbf{1}_n, \delta)$ # Semi-relaxed OT
 $\mathbf{R}_k \leftarrow \text{U-projection}(\mathbf{K}_R^{(k)}, \gamma_k, \tau, \mathbf{b}, \mathbf{R}_{k-1}^T \mathbf{1}_m, \delta)$ # Unbalanced OT
 end if
 $\mathbf{g}_Q, \mathbf{g}_R \leftarrow \mathbf{Q}_k^T \mathbf{1}_n, \mathbf{R}_k^T \mathbf{1}_m$
 $\nabla_T = \text{diag}(\mathbf{g}_Q)^{-1} \mathbf{Q}_k^T \mathbf{C} \mathbf{R}_k \text{diag}(\mathbf{g}_R)^{-1}$
 $\gamma_T = \gamma / \|\nabla_T\|_\infty$ # ℓ^∞ -normalization
 $\mathbf{K}_T^{(k)} \leftarrow \mathbf{T}_k \odot \exp\{-\gamma_T \nabla_T\}$
 $\mathbf{T}_k \leftarrow \text{Sinkhorn}(\mathbf{K}_T^{(k)}, \mathbf{g}_R, \mathbf{g}_Q, \delta)$ # Balanced OT
 $\mathbf{X}_k \leftarrow \text{diag}(1/\mathbf{g}_Q) \mathbf{T}_k \text{diag}(1/\mathbf{g}_R)$
end while
Return $\mathbf{P}_r = \mathbf{Q} \mathbf{X} \mathbf{R}^T$

Algorithm 5 Sinkhorn Algorithm (*Cuturi (2013b), balanced OT*)

Input $\mathbf{K}, \mathbf{a}, \mathbf{b}, \delta$
 $\mathbf{u} \leftarrow \mathbf{1}_n$
 $\mathbf{v} \leftarrow \mathbf{1}_m$
while $\|\text{diag}(\mathbf{u}) \mathbf{K} \mathbf{v} - \mathbf{a}\|_1 + \|\text{diag}(\mathbf{v}) \mathbf{K}^T \mathbf{u} - \mathbf{b}\|_1 > \delta$ **do**
 $\mathbf{u}^{(l+1)} \leftarrow \mathbf{a} / \mathbf{K} \mathbf{v}^{(l)}$
 $\mathbf{v}^{(l+1)} \leftarrow \mathbf{b} / \mathbf{K}^T \mathbf{u}^{(l+1)}$
end while
Return $\text{diag}(\mathbf{u}) \mathbf{K} \text{diag}(\mathbf{v})$

D Gromov-Wasserstein (GW)

As defined in, the general Gromov-Wasserstein problem concerns a minimization of the energy:

$$\mathcal{Q}_{A,B}(P) = \sum_{i,j,k,l} (\mathbf{A}_{ik} - \mathbf{B}_{jl})^2 P_{ij} P_{kl} \quad (30)$$

Where the minimization is over the set of all couplings $\Pi_{a,b}$:

$$\text{GW}(\mu, \nu) := \min_{P \in \Pi_{a,b}} \mathcal{Q}_{A,B}(P) \quad (31)$$

We consider extending the semi-relaxed framework to the GW problem under the low-rank restriction on P . This extension has thus far been considered for the balanced and unbalanced case in two previous works Scetbon et al. (2023, 2022). In both of these works, the algorithms for the Wasserstein case extend trivially to the Gromov-Wasserstein (GW) problem. In particular, each kernel with variable ζ has an update of the form $K \leftarrow \zeta \odot \exp(-\gamma_k \nabla_{\zeta} \mathcal{L}(\zeta))$ for $\mathcal{L}(\zeta)$ heretofore taken to be the Wasserstein loss $\langle P(\zeta), C \rangle_F$ where the coupling $P = QXR^T$ is interpreted as a function of the low-rank sub-factor variables $\zeta \in \{Q, R, X\}$. Taking $\mathcal{L} := \mathcal{Q}_{A,B}(P(\zeta))$ one can simply extend the gradient through the GW-loss and directly use it in place of the Wasserstein gradient in the update $K \leftarrow \zeta \odot \exp(-\gamma_k \nabla_{\zeta} \mathcal{L}(\zeta))$ for $\mathcal{L}(\zeta)$ of each algorithm. The matrix-form of the GW-cost is expressed as:

$$\mathcal{Q}_{A,B}(P) = \mathbf{1}_m^T P^T A^{\odot 2} P \mathbf{1}_m + \mathbf{1}_n^T P B^{\odot 2} P^T \mathbf{1}_n - 2\langle APB, P \rangle$$

Which, using the constraints of \mathcal{C}_2 reduces the cost as a function of Q, R, X to:

$$\mathcal{Q}_{A,B}(Q, R, X) = \mathbf{1}_r^T Q^T A^{\odot 2} Q \mathbf{1}_r + \mathbf{1}_r^T R^T B^{\odot 2} R \mathbf{1}_r - 2\langle QXR^T, AQXR^T B \rangle_F$$

$$\nabla_Q \mathcal{Q}_{A,B}(Q, R, X) = 2A^{\odot 2} Q \mathbf{1}_r \mathbf{1}_r^T - 4AQXR^T BRX^T$$

Which is proportional to $\nabla_Q \mathcal{Q}_{A,B}(Q, R, X) \propto -4AQXR^T BRX^T$ for the balanced and right-marginal semi-relaxed case. And:

$$\nabla_R \mathcal{Q}_{A,B}(Q, R, X) = 2B^{\odot 2} R \mathbf{1}_r \mathbf{1}_r^T - 4BRX^T Q^T AQX$$

Which likewise can be reduced in proportionality to $\nabla_R \mathcal{Q}_{A,B}(Q, R, g) \propto -4BRX^T Q^T AQX$ in the balanced and left-marginal semi-relaxed case. The gradients, as presented above, assume a linearization in $X \leftarrow X_k$. If one does not make this assumption and takes $X(Q, R) = \text{diag}(1/Q^T \mathbf{1}_n) T \text{diag}(1/R^T \mathbf{1}_m)$, a rank-one perturbation must be added to the Q and R gradient of the form:

$$\begin{aligned} \nabla_Q^{(2)} &= 4\mathbf{1}_n \text{diag}^{-1}(XR^T B(QXR^T)^T AQ \text{diag}(1/g_Q))^T \\ \nabla_R^{(2)} &= 4\mathbf{1}_m \text{diag}^{-1}(X^T Q^T AQXR^T BR \text{diag}(1/g_R))^T \end{aligned}$$

Analogously, for the gradient on T one can simply take the gradient with respect to X , and subsequently T as $X(T) = \text{diag}(g_Q)^{-1} T \text{diag}(g_R)^{-1}$. The gradient with respect to X is given as:

$$\nabla_X \mathcal{Q}_{A,B}(Q, R, X) = -4Q^T AQXR^T BR$$

And thus, the directional derivative with respect to X in the direction $V_X = \text{diag}(g_Q)^{-1} V_T \text{diag}(g_R)^{-1}$ and thus by the chain rule V_T is:

$$\begin{aligned} D\mathcal{Q}_{A,B}(Q, R, X) \circ (V_X) &= -4\langle Q^T AQXR^T BR, V_X \rangle_F \\ &= -4\langle Q^T AQXR^T BR, \text{diag}(1/g_Q) V_T \text{diag}(1/g_R) \rangle_F \end{aligned}$$

So that the gradient with respect to the coupling matrix T is given as:

$$\nabla_T \mathcal{Q}_{A,B}(Q, R, T) = -4 \text{diag}(1/g_Q) Q^T AQXR^T BR \text{diag}(1/g_R)$$

E Convergence Analysis and Other Proofs

E.1 Convergence and Smoothness of the Objective

We show in Proposition 3.3 that directly applying the block-descent lemma of Beck & Tetrushvili (2013) to the template of Ghadimi et al.'s proof Ghadimi et al. (2014) is sufficient to show the non-asymptotic stationary convergence of a coordinate mirror descent procedure in Ghadimi's criterion. The non-asymptotic guarantee of Ghadimi et al. (2014) follows directly in the case of coordinate mirror descent using Lemma E.1 and Lemma E.2 below. For completeness, we define all notation used, describe the coordinate mirror descent algorithm in general, and discuss a few relevant preliminaries.

Suppose that the vector of n variables $\mathbf{x} \in \mathbb{R}^n$ is partitioned into p blocks, $\mathbf{x} = (\mathbf{x}(1), \dots, \mathbf{x}(p))$, where $\mathbf{x}(i) \in \mathbb{R}^{n_i}$. Here, n_1, \dots, n_p are positive integers summing to n . Following the notation of Beck & Tetrushvili (2013); Nesterov (2012), we define matrices $\mathbf{U}_i \in \mathbb{R}^{n \times n_i}$ such that $\mathbf{x}(i) = \mathbf{U}_i^T \mathbf{x}$ for all $i = 1, \dots, p$. This also implies $\mathbf{x} = \sum_{i=1}^p \mathbf{U}_i \mathbf{x}(i)$. This allows us to define the vector of partial derivatives corresponding to each block of variables $\mathbf{x}(i)$:

$$\nabla_i f(\mathbf{x}) := \mathbf{U}_i^T \nabla f(\mathbf{x}).$$

In Beck & Tetrushvili (2013), the gradient of f is assumed to be block-coordinate-wise Lipschitz, with L_i the smoothness constant associated to the i -th block of variables: for all $\mathbf{h}_i \in \mathbb{R}^{n_i}$, one has

$$\|\nabla_i f(\mathbf{x} + \mathbf{U}_i \mathbf{h}_i) - \nabla_i f(\mathbf{x})\| \leq L_i \|\mathbf{h}_i\|, \quad (32)$$

and for such functions we denote by $L := \max_i L_i$ the (global) smoothness constant of ∇f . To be clear, a *smoothness constant* associated to f is a Lipschitz constant of its gradient.

Lemma E.1 (Block descent lemma, Beck & Tetrushvili (2013), Lemma 3.2.). *Suppose $f \in C^1(\mathbb{R}^n, \mathbb{R})$ is a continuously differentiable function over \mathbb{R}^n whose gradient is block-coordinate-wise Lipschitz (32) for L_i the smoothness constant associated to the i -th block of variables $\mathbf{x}(i)$. Let \mathbf{u}, \mathbf{v} be two vectors differing only in the i -th block: there exists $\mathbf{h}_i \in \mathbb{R}^{n_i}$ such that $\mathbf{v} - \mathbf{u} = \mathbf{U}_i \mathbf{h}_i$. Then,*

$$f(\mathbf{v}) \leq f(\mathbf{u}) + \langle \nabla f(\mathbf{u}), \mathbf{v} - \mathbf{u} \rangle + \frac{L_i}{2} \|\mathbf{u} - \mathbf{v}\|^2. \quad (33)$$

Lemma E.1 is central to adapting the proof of Theorem 1 of Ghadimi et al. (2014) to our case. Their Theorem 1 concerns the non-asymptotic convergence of mirror descent for objectives of the form

$$\min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) + h(\mathbf{x}),$$

where \mathcal{X} is a closed, convex subset of \mathbb{R}^n , $f \in C^1(\mathcal{X}, \mathbb{R})$ is a possibly non-convex objective, and where h is an α -strongly convex function. Using notation similar to Ghadimi et al. (2014), we write $\Phi(\mathbf{x}) = f(\mathbf{x}) + h(\mathbf{x})$. Additionally, Ghadimi et al. (2014) assume ∇f is L -Lipschitz for some $L > 0$.

To unify our assumptions on f , we suppose that $\nabla f \in C^1(\mathcal{X}, \mathbb{R})$ is block-coordinate Lipschitz (32) with block Lipschitz constants $(L_i)_{i=1}^p$, and that \mathcal{X} itself decomposes as a product $\mathcal{X} = \prod_{i=1}^p \mathcal{X}_i$, where each \mathcal{X}_i is a closed convex set constraining the block variables $\mathbf{x}(i)$.

The proof of Ghadimi relies on β -smoothness of the objective in all variables. We show that component-wise smoothness in each block is sufficient to achieve an analogous convergence result for a coordinate mirror descent. To provide context for Proposition 3.3, we now describe in general (1) mirror descent, and (2) block-coordinate mirror descent.

We again follow the notation of Ghadimi et al. (2014). A function $\omega : \mathcal{X} \rightarrow \mathbb{R}$ is a *distance generating function* with modulus $\alpha > 0$, with respect to the Euclidean norm $\|\cdot\|$, if ω is continuously differentiable and strongly convex, so that

$$\langle \mathbf{x} - \mathbf{z}, \nabla \omega(\mathbf{x}) - \nabla \omega(\mathbf{z}) \rangle \geq \alpha \|\mathbf{x} - \mathbf{z}\|^2, \quad \text{for all } \mathbf{x}, \mathbf{z} \in \mathcal{X}.$$

The *prox-function* (or Bregman divergence) associated with ω is then

$$V(\mathbf{x}, \mathbf{z}) = \omega(\mathbf{x}) - \omega(\mathbf{z}) - \langle \nabla \omega(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle,$$

and from this prox-function, $\gamma > 0$, and some $\mathbf{g} \in \mathbb{R}^n$, we define the *generalized projection*

$$\mathbf{x}^+ := \arg \min_{\mathbf{u} \in \mathcal{X}} \left(\langle \mathbf{g}, \mathbf{u} \rangle + \frac{1}{\gamma} V(\mathbf{u}, \mathbf{x}) + h(\mathbf{u}) \right). \quad (34)$$

To describe the projected gradient descent algorithm (which coincides with mirror descent in our case of interest), we first define the *generalized projected gradient of Φ at \mathbf{x}* :

$$P_{\mathcal{X}}(\mathbf{x}, \mathbf{g}, \gamma) := \frac{1}{\gamma} (\mathbf{x} - \mathbf{x}^+).$$

The *mirror descent (MD) algorithm* is as follows: given initial point $\mathbf{x}_0 \in \mathcal{X}$, a total number of iterations N , and positive stepsizes $(\gamma_k)_{k=1}^N$, at step k , the $(k+1)$ -st iterate is computed via

$$\mathbf{x}_{k+1} \leftarrow \arg \min_{\mathbf{u} \in \mathcal{X}} \left(\langle \nabla f(\mathbf{x}_k), \mathbf{u} \rangle + \frac{1}{\gamma_k} V(\mathbf{u}, \mathbf{x}_k) + h(\mathbf{u}) \right).$$

Among all iterates \mathbf{x}_k , the MD algorithm outputs the one at which the generalized projected gradient is of least norm. Concretely, \mathbf{x}_R is the output of the MD algorithm, where

$$R := \arg \min_{k=0, \dots, N} \|\mathbf{g}_{\mathcal{X},k}\|^2,$$

and where $\mathbf{g}_{\mathcal{X},k}$ is

$$\mathbf{g}_{\mathcal{X},k} := P_{\mathcal{X}}(\mathbf{x}_k, \nabla f(\mathbf{x}_k), \gamma_k).$$

Having described the MD algorithm, let us consider a block-coordinate variant; we assume \mathbf{x} admits the block-coordinate structure described above. To simplify the presentation, we suppose that in a given iteration k , the block variables are updated sequentially from $i = 1, \dots, p$. This leads to doubly-indexed iterates (\mathbf{x}_k^i) with $k = 1, \dots, N$ indexing each full iteration through all variables, and $i = 0, 1, \dots, p$ indexing the sub-iterations which update one block of variables at a time.

The *coordinate mirror descent (CMD) algorithm* takes as input an initial point $\mathbf{x}_0 \in \mathcal{X}$, a number of iterations N , and a sequence of positive stepsizes $(\gamma_{k,i})_{k=1, i=1}^{N,p}$. We set $\mathbf{x}_0^0 = \mathbf{x}_0$, and for $k = 0, \dots, N-1$ and $i = 1, \dots, p$, we compute \mathbf{x}_k^i from \mathbf{x}_k^{i-1} as follows:

$$\begin{aligned} \mathbf{x}_k^i(i) &\leftarrow \arg \min_{\mathbf{u}_i \in \mathcal{X}_i} \left(\langle \nabla_i f(\mathbf{x}_k^{i-1}), \mathbf{u}_i \rangle + \frac{1}{\gamma_{k,i}} V_i(\mathbf{u}_i, \mathbf{U}_i^T \mathbf{x}_k^{i-1}) + h_i(\mathbf{u}_i) \right) \\ \mathbf{x}_k^i(j) &\leftarrow \mathbf{x}_k^{i-1}(j) \quad \text{for } j \neq i \end{aligned}$$

Here, we have assumed that V can be written as a composite function of the block variables,

$$V(\mathbf{x}, \mathbf{z}) = \sum_{i=1}^p V_i(\mathbf{x}(i), \mathbf{z}(i)),$$

as is the case for the KL divergence. We also have assumed that h has this composite structure (as with entropy):

$$h(\mathbf{x}) = \sum_{i=1}^p h_i(\mathbf{x}(i)).$$

Lastly, for $k = 0, \dots, N-1$, we set $\mathbf{x}_{k+1}^0 := \mathbf{x}_k^p$. We define $\mathbf{g}_{\mathcal{X},k} := (\mathbf{g}_{\mathcal{X},k,1}, \dots, \mathbf{g}_{\mathcal{X},k,p})$ to be the collection of block-wise differences, where by definition

$$\mathbf{g}_{\mathcal{X},k,i} = P_{\mathcal{X}_i}(\mathbf{x}_k^{i-1}, \nabla_i f(\mathbf{x}_k^{i-1}), \gamma_{k,i}) = \frac{1}{\gamma_{k,i}} (\mathbf{x}_k^{i-1} - \mathbf{x}_k^i) = \mathbf{U}_i^T \mathbf{g}_{\mathcal{X},k}.$$

The convergence criterion Δ we use in Proposition 3.3 is the one used by Ghadimi et al. (2014) summed across blocks. In particular, the CMD algorithm returns iterate \mathbf{x}_R , where

$$\begin{aligned} R &:= \arg \min_{k=0, \dots, N} \Delta(\mathbf{x}_k, \mathbf{x}_{k-1}), \\ \Delta(\mathbf{x}_k, \mathbf{x}_{k-1}) &:= \|\mathbf{g}_{\mathcal{X},k}\|^2 = \sum_{i=1}^p \|\mathbf{g}_{\mathcal{X},k,i}\|^2. \end{aligned} \tag{35}$$

For $\Phi = f + h$ as above (f has global smoothness constant L), let $\Phi^* = \arg \min_{\mathbf{x} \in \mathcal{X}} \Phi(\mathbf{x})$, and define

$$D := \left(\frac{\Phi(x_0) - \Phi^*}{L} \right)^{1/2}. \tag{36}$$

The other lemma used in the proof of Proposition 3.3 is as follows.

Lemma E.2 (Ghadimi et al. (2014), Lemma 1). *Let $\mathbf{x}^+ = \arg \min_{\mathbf{u} \in \mathcal{X}} \{ \langle \mathbf{g}, \mathbf{u} \rangle + \frac{1}{\gamma} V(\mathbf{u}, \mathbf{x}) + h(\mathbf{u}) \}$ and $P_{\mathcal{X}}(\mathbf{x}, \mathbf{g}, \gamma) = \frac{1}{\gamma} (\mathbf{x} - \mathbf{x}^+)$. Then for all $\mathbf{x} \in \mathbb{R}^n$, all $\mathbf{g} \in \mathbb{R}^n$, and $\gamma > 0$, one has:*

$$\langle \mathbf{g}, P_{\mathcal{X}}(\mathbf{x}, \mathbf{g}, \gamma) \rangle \geq \alpha \|P_{\mathcal{X}}(\mathbf{x}, \mathbf{g}, \gamma)\|^2 + \frac{1}{\gamma} (h(\mathbf{x}^+) - h(\mathbf{x})).$$

Proposition E.3 (Proposition 3.3). *Suppose one has $f \in C^1(\mathcal{X}, \mathbb{R})$ whose gradient is block-coordinate Lipschitz, with block smoothness constants $(L_i)_{i=1}^p$, and a function $h \in C(\mathcal{X}, \mathbb{R})$ which is α -strongly convex. For $\Phi = f + h$, suppose one performs a coordinate mirror descent on Φ minimized over a product of closed convex sets $\mathcal{X} = \prod_{i=1}^p \mathcal{X}_i$. Let the sub-iterates with respect to the i -th block update be $\{\mathbf{x}_k^i\}_{i=0}^p$ where $\mathbf{x}_k := \mathbf{x}_k^0$ for $k \in [N]$ outer iterations. Then one has:*

$$\min_k \Delta(\mathbf{x}_k, \mathbf{x}_{k-1}) \leq \frac{D^2 L}{N(\alpha^2/2L)} = \frac{2D^2 L^2}{N\alpha^2},$$

where D is (36), L is the global smoothness constant of f , and convergence criterion $\Delta(\mathbf{x}_k, \mathbf{x}_{k-1})$ is given in (35). Above, the stepsizes $\gamma_{k,i}$ in the coordinate mirror descent are $\gamma_{k,i} := \alpha/L$.

Proof. As f satisfies the hypotheses of the block descent lemma, Lemma E.1, we apply (33) to obtain:

$$f(\mathbf{x}_k^i) \leq f(\mathbf{x}_k^{i-1}) + \langle \nabla_i f(\mathbf{x}_k^{i-1}), \mathbf{x}_k^i - \mathbf{x}_k^{i-1} \rangle + \frac{L_i}{2} \|\mathbf{x}_k^i - \mathbf{x}_k^{i-1}\|^2.$$

Noting the definition $\mathbf{g}_{\mathcal{X},k,i} = \frac{1}{\gamma_{k,i}} (\mathbf{x}_k^{i-1} - \mathbf{x}_k^i)$, one has

$$f(\mathbf{x}_k^i) \leq f(\mathbf{x}_k^{i-1}) - \gamma_{k,i} \langle \nabla_i f(\mathbf{x}_k^{i-1}), \mathbf{g}_{\mathcal{X},k,i} \rangle + \frac{L_i}{2} \gamma_{k,i}^2 \|\mathbf{g}_{\mathcal{X},k,i}\|^2.$$

Lemma 1 of Ghadimi et al. (2014) (stated as Lemma E.2 above) applies identically through block-wise optimality on \mathcal{X}_i because $\mathbf{g}_{\mathcal{X},k,i} = P_{\mathcal{X}_i}(\mathbf{x}_k^{i-1}, \nabla_i f(\mathbf{x}_k^{i-1}), \gamma_{k,i})$. Thus for any value $\nabla_i f(\mathbf{x}_k^{i-1})$ takes,

$$f(\mathbf{x}_k^i) \leq f(\mathbf{x}_k^{i-1}) - [\alpha\gamma_{k,i} \|\mathbf{g}_{\mathcal{X},k,i}\| + h(\mathbf{x}_k^i) - h(\mathbf{x}_k^{i-1})] + \frac{L_i}{2} \gamma_{k,i}^2 \|\mathbf{g}_{\mathcal{X},k,i}\|^2,$$

and thus,

$$f(\mathbf{x}_k^i) + h(\mathbf{x}_k^i) \leq f(\mathbf{x}_k^{i-1}) + h(\mathbf{x}_k^{i-1}) - \left[\alpha\gamma_{k,i} - \frac{L_i}{2} \gamma_{k,i}^2 \right] \|\mathbf{g}_{\mathcal{X},k,i}\|^2.$$

The right-hand side above only becomes larger, taking $L_i = L$ to be the global smoothness constant. Introducing a sum over sub-iterates and total iterates, one has

$$\begin{aligned} \sum_{k,i}^{N,p} f(\mathbf{x}_k^i) + h(\mathbf{x}_k^i) &\leq \sum_{k,i}^{N,p} f(\mathbf{x}_k^{i-1}) + h(\mathbf{x}_k^{i-1}) - \sum_{k,i}^{N,p} \left[\alpha\gamma_{k,i} - \frac{L}{2} \gamma_{k,i}^2 \right] \|\mathbf{g}_{\mathcal{X},k,i}\|^2, \\ \sum_{k,i}^{N,p} \Phi(\mathbf{x}_k^i) &\leq \sum_{k,i}^{N,p} \Phi(\mathbf{x}_k^{i-1}) - \sum_{k,i}^{N,p} \left[\alpha\gamma_{k,i} - \frac{L}{2} \gamma_{k,i}^2 \right] \|\mathbf{g}_{\mathcal{X},k,i}\|^2. \end{aligned}$$

Noting the end-point condition $\Phi(\mathbf{x}_k^p) = \Phi(\mathbf{x}_{k+1}^0)$, one may cancel all intermediate terms:

$$\Phi^* \leq \Phi(\mathbf{x}_N) \leq \Phi(\mathbf{x}_0) - \sum_{k,i}^{N,p} \left[\alpha\gamma_{k,i} - \frac{L}{2} \gamma_{k,i}^2 \right] \|\mathbf{g}_{\mathcal{X},k,i}\|^2$$

Thus one finds the upper bound in terms of $\Phi(\mathbf{x}_0)$ and the minimum value $\Phi^* = \min_{\mathbf{x} \in \mathcal{X}} \Phi(\mathbf{x})$ of Φ :

$$\sum_k^N \sum_i^p \left[\alpha\gamma_{k,i} - \frac{L}{2} \gamma_{k,i}^2 \right] \|\mathbf{g}_{\mathcal{X},k,i}\|^2 \leq \Phi(\mathbf{x}_0) - \Phi^*. \quad (37)$$

Taking $\gamma_{k,i} = \alpha/L$ as in Ghadimi et al. (2014), the bracketed term directly above becomes $\alpha^2/2L$, and one has:

$$\begin{aligned} \sum_k^N \left[\frac{\alpha^2}{2L} \right] \left(\min_k \Delta(\mathbf{x}_k, \mathbf{x}_{k-1}) \right) &= \sum_k^N \left[\frac{\alpha^2}{2L} \right] \left(\min_k \sum_i^p \|\mathbf{g}_{\mathcal{X},k,i}\|^2 \right) \\ &\leq \sum_k^N \left[\frac{\alpha^2}{2L} \right] \sum_i^p \|\mathbf{g}_{\mathcal{X},k,i}\|^2 \\ &\leq \Phi(\mathbf{x}_0) - \Phi^* = D^2 L, \end{aligned}$$

where D is as in (36), and where we used (37) to obtain the last line. Thus,

$$\min_k \Delta(\mathbf{x}_k, \mathbf{x}_{k-1}) \leq \frac{D^2 L}{N(\alpha^2/2L)} = \frac{2D^2 L^2}{N\alpha^2},$$

completing the proof. \square

Definition E.4 (Relative smoothness). Let $\beta > 0$ and let $g \in \mathcal{C}^1(\mathbb{R}^n, \mathbb{R})$ be continuously differentiable. Additionally, let ω be a distance generating function with associated prox-function V . The function g is β -smooth relative to ω if the following holds:

$$g(\mathbf{y}) \leq g(\mathbf{x}) + \langle \nabla \omega(\mathbf{x}), \mathbf{x} - \mathbf{y} \rangle + \beta V(\mathbf{y}, \mathbf{x})$$

Proposition E.5. Let $\epsilon > 0$ be a predefined error tolerance, $r > 0$ a small rank parameter, $\delta \in (0, \frac{1}{r})$ a lower-bound parameter, and N the number of inner iterations required for the semi-relaxed projection to converge for each iteration k as $\|\mathbf{Q}^T \mathbf{1}_n - \mathbf{g}_Q^{(k-1)}\|_2 < \epsilon = \frac{1}{N} (\frac{1}{r} - \delta)$ for $\mathbf{g}_Q^{(0)} = \frac{1}{r} \mathbf{1}_r$. Then, the FRLC objective

$$\mathcal{L}_{\text{LC}}(\mathbf{Q}, \mathbf{R}, \mathbf{T}) = \langle \mathbf{Q} \text{diag}(1/\mathbf{Q}^T \mathbf{1}_n) \mathbf{T} \text{diag}(1/\mathbf{R}^T \mathbf{1}_m) \mathbf{R}^T, \mathbf{C} \rangle_F$$

is component-wise smooth with respect to the variables $\mathbf{Q}, \mathbf{R}, \mathbf{T}$ with smoothness constants $\{\beta_i\}_{i=1}^3$ where $\beta_i = \text{poly}(\|\mathbf{C}\|_F, n, m, r, \delta)$.

Proof. The addition of a pair of regularizations on the inner marginals $\tau \text{KL}(\mathbf{Q}^T \mathbf{1}_n \| \mathbf{Q}_k^T \mathbf{1}_n)$ and $\tau \text{KL}(\mathbf{R}^T \mathbf{1}_m \| \mathbf{R}_k^T \mathbf{1}_m)$ in 21 and 22 ensures that we may use standard results on relaxed optimal-transport which bound how far the marginal $\mathbf{Q}_k^T \mathbf{1}_n$ deviates across iterations.

In particular, for all $\epsilon > 0$ there exists τ and N (number of iterations) sufficiently large, so that $\|\mathbf{R}^T \mathbf{1}_m - \mathbf{g}_R\|_2^2 < \epsilon$ and $\|\mathbf{Q}^T \mathbf{1}_n - \mathbf{g}_Q\|_2^2 < \epsilon$. In particular, Pham et al. (2020) shows that one can attain convergence to any ϵ for the unbalanced problem in $N = \tilde{O}(m^2/\epsilon)$ iterations (hiding logarithmic and τ -factors) for each sub-problem solving for \mathbf{R}_k in Algorithm 4 and analogously $N = \tilde{O}(n^2/\epsilon)$ for \mathbf{Q} . Under the uniform initialization of \mathbf{g}_R , and after N iterations, we have:

$$\|\mathbf{g}_R^{(0)} - \mathbf{g}_R^{(N)}\|_2 = \left\| \frac{1}{r} \mathbf{1}_r - \mathbf{g}_R^{(N)} \right\|_2 \leq \sum_{k=1}^N \|\mathbf{g}_R^{(k)} - \mathbf{g}_R^{(k-1)}\|_2.$$

With sufficiently large τ and $N = \tilde{O}(m^2/\epsilon)$ sub-iterations, one can guarantee that:

$$\|\mathbf{g}_R^{(k)} - \mathbf{g}_R^{(k-1)}\|_2 < \epsilon = \frac{1}{N} \left(\frac{1}{r} - \delta \right).$$

This implies, for all iterations of the algorithm and all indices i , that

$$(\mathbf{g}_{R_k})_i > \delta, \quad \text{and analogously,} \quad (\mathbf{g}_{Q_k})_i > \delta. \quad (38)$$

Thus, by adding the regularization on the inner marginal, one may guarantee a lower-bound on the entries of \mathbf{g}_R and \mathbf{g}_Q . This is essential for demonstrating smoothness of the objective.

First, we consider smoothness in \mathbf{Q} . We note that the gradient in \mathbf{Q} splits into two terms:

$$\begin{aligned} \nabla_{\mathbf{Q}} \mathcal{L}_{\text{LC}}(\mathbf{Q}, \mathbf{R}, \mathbf{T}) &= \nabla_{\mathbf{Q}}^{(A)} \mathcal{L}_{\text{LC}} + \nabla_{\mathbf{Q}}^{(B)} \mathcal{L}_{\text{LC}} \\ &= \mathbf{C} \mathbf{R} \mathbf{X}^T - \mathbf{1}_n \text{diag}^{-1}((\mathbf{C} \mathbf{R} \mathbf{X}^T)^T \mathbf{Q} \text{diag}(1/\mathbf{g}_Q))^T \\ &= \nabla_{\mathbf{Q}}^{(A)} \mathcal{L}_{\text{LC}} - \mathbf{1}_n \text{diag}^{-1}((\nabla_{\mathbf{Q}}^{(A)} \mathcal{L}_{\text{LC}})^T \mathbf{Q} \text{diag}(1/\mathbf{g}_Q))^T \end{aligned}$$

Where

$$\begin{aligned} &\|\nabla_{\mathbf{Q}} \mathcal{L}_{\text{LC}}(\mathbf{Q}_{k+1}, \mathbf{R}_k, \mathbf{T}_k) - \nabla_{\mathbf{Q}} \mathcal{L}_{\text{LC}}(\mathbf{Q}_k, \mathbf{R}_k, \mathbf{T}_k)\|_F \\ &\leq \|\nabla_{\mathbf{Q}}^{(A)} \mathcal{L}_{\text{LC}}(\mathbf{Q}_{k+1}, \mathbf{R}_k, \mathbf{T}_k) - \nabla_{\mathbf{Q}}^{(A)} \mathcal{L}_{\text{LC}}(\mathbf{Q}_k, \mathbf{R}_k, \mathbf{T}_k)\|_F \\ &\quad + \|\nabla_{\mathbf{Q}}^{(B)} \mathcal{L}_{\text{LC}}(\mathbf{Q}_{k+1}, \mathbf{R}_k, \mathbf{T}_k) - \nabla_{\mathbf{Q}}^{(B)} \mathcal{L}_{\text{LC}}(\mathbf{Q}_k, \mathbf{R}_k, \mathbf{T}_k)\|_F. \end{aligned} \quad (39)$$

Starting with the first term on the right side of (39), one has:

$$\begin{aligned} & \|\nabla_{\mathbf{Q}}^{(A)} \mathcal{L}_{\text{LC}}(\mathbf{Q}_{k+1}, \mathbf{R}_k, \mathbf{T}_k) - \nabla_{\mathbf{Q}}^{(A)} \mathcal{L}_{\text{LC}}(\mathbf{Q}_k, \mathbf{R}_k, \mathbf{T}_k)\|_F \\ &= \|\mathbf{C}\mathbf{R}_k(\text{diag}(1/\mathbf{g}_{\mathbf{Q}_k})\mathbf{T}_k\text{diag}(1/\mathbf{g}_{\mathbf{R}_k}))^T - \mathbf{C}\mathbf{R}_k(\text{diag}(1/\mathbf{g}_{\mathbf{Q}_{k-1}})\mathbf{T}_k\text{diag}(1/\mathbf{g}_{\mathbf{R}_k}))^T\|_F \\ &\leq \|\text{diag}(1/\mathbf{g}_{\mathbf{R}_k})\|_F \|\mathbf{C}\|_F \|\mathbf{R}_k\|_F \|\mathbf{T}_k\|_F \|\text{diag}(1/\mathbf{g}_{\mathbf{Q}_k}) - \text{diag}(1/\mathbf{g}_{\mathbf{Q}_{k-1}})\|_F. \end{aligned}$$

Note that $\|\mathbf{R}_k\|_F^2 = \sum_{i,j} (\mathbf{R}_k)_{i,j}^2 < \sum_{i,j} (\mathbf{R}_k)_{i,j} = 1$, as \mathbf{R}_k has marginals which sum to one. The same bound holds for $\|\mathbf{T}_k\|_F^2$, which is also a coupling. Invoking the lower-bound (38) of δ on the entries of the inner marginals, and continuing from the above display,

$$\begin{aligned} & \|\nabla_{\mathbf{Q}}^{(A)} \mathcal{L}_{\text{LC}}(\mathbf{Q}_{k+1}, \mathbf{R}_k, \mathbf{T}_k) - \nabla_{\mathbf{Q}}^{(A)} \mathcal{L}_{\text{LC}}(\mathbf{Q}_k, \mathbf{R}_k, \mathbf{T}_k)\|_F \\ &\leq \frac{\|\mathbf{C}\|_F}{\delta} \|\text{diag}(1/\mathbf{g}_{\mathbf{Q}_k}) - \text{diag}(1/\mathbf{g}_{\mathbf{Q}_{k-1}})\|_F \\ &= \frac{\|\mathbf{C}\|_F}{\delta} \|\text{diag}(1/\mathbf{g}_{\mathbf{Q}_{k-1}})\text{diag}(1/\mathbf{g}_{\mathbf{Q}_k})(\text{diag}(\mathbf{g}_{\mathbf{Q}_{k-1}}) - \text{diag}(\mathbf{g}_{\mathbf{Q}_k}))\|_F \\ &\leq \frac{\|\mathbf{C}\|_F}{\delta^3} \|\text{diag}(\mathbf{g}_{\mathbf{Q}_{k-1}}) - \text{diag}(\mathbf{g}_{\mathbf{Q}_k})\|_F. \end{aligned}$$

To further bound the right-hand side above, consider:

$$\|\text{diag}(\mathbf{g}_{\mathbf{Q}_k}) - \text{diag}(\mathbf{g}_{\mathbf{Q}_{k-1}})\|_F^2 = \|\mathbf{Q}_k^T \mathbf{1}_n - \mathbf{Q}_{k-1}^T \mathbf{1}_n\|_2^2 = \sum_{i=1}^r \left(\sum_{j=1}^r (\mathbf{Q}_k)_{i,j} - (\mathbf{Q}_{k-1})_{i,j} \right)^2$$

While can easily be upper-bounded by an application of Jensen's inequality as

$$\begin{aligned} &= \sum_{i=1}^r r^2 \left(\sum_{j=1}^r \frac{1}{r} \left((\mathbf{Q}_k)_{i,j} - (\mathbf{Q}_{k-1})_{i,j} \right) \right)^2 \\ &\leq \sum_{i=1}^r r^2 \left(\sum_{j=1}^r \frac{1}{r} \left((\mathbf{Q}_k)_{i,j} - (\mathbf{Q}_{k-1})_{i,j} \right)^2 \right) \\ &= r \sum_{i=1}^r \sum_{j=1}^r \left((\mathbf{Q}_k)_{i,j} - (\mathbf{Q}_{k-1})_{i,j} \right)^2 \\ &= r \|\mathbf{Q}_k - \mathbf{Q}_{k-1}\|_F^2. \end{aligned}$$

Likewise, we have that $\|\text{diag}(\mathbf{g}_{\mathbf{R}_k}) - \text{diag}(\mathbf{g}_{\mathbf{R}_{k-1}})\|_F^2 \leq r \|\mathbf{R}_k - \mathbf{R}_{k-1}\|_F^2$. Thus, it holds that

$$\frac{\|\mathbf{C}\|_F}{\delta^3} \|\mathbf{g}_{\mathbf{Q}_k} - \mathbf{g}_{\mathbf{Q}_{k-1}}\|_2 \leq \frac{\|\mathbf{C}\|_F \sqrt{r}}{\delta^3} \|\mathbf{Q}_k - \mathbf{Q}_{k-1}\|_F.$$

Next, we focus on the $\nabla_{\mathbf{Q}}^{(B)}$ term. Observe that

$$\begin{aligned} \|\mathbf{1}_n \text{diag}^{-1} \mathbf{X}\|_F^2 &= \text{Tr}(\mathbf{1}_n \text{diag}^{-1} \mathbf{X})^T (\mathbf{1}_n \text{diag}^{-1} \mathbf{X}) \\ &= n \|\text{diag}^{-1} \mathbf{X}\|_2^2 \leq n \|\mathbf{X}\|_F^2. \end{aligned}$$

Thus:

$$\|\nabla_{\mathbf{Q}_{k+1}}^{(B)} - \nabla_{\mathbf{Q}_k}^{(B)}\|_F \leq \sqrt{n} \|(\nabla_{\mathbf{Q}_{k+1}}^{(A)})^T \mathbf{Q}_{k+1} \text{diag}(1/\mathbf{g}_{\mathbf{Q}_{k+1}}) - (\nabla_{\mathbf{Q}_k}^{(A)})^T \mathbf{Q}_k \text{diag}(1/\mathbf{g}_{\mathbf{Q}_k})\|_F$$

Adding and subtracting terms in the norm and applying triangle inequality:

$$\begin{aligned} & \sqrt{n} \|(\nabla_{\mathbf{Q}_{k+1}}^{(A)})^T \mathbf{Q}_{k+1} \text{diag}(1/\mathbf{g}_{\mathbf{Q}_{k+1}}) - (\nabla_{\mathbf{Q}_k}^{(A)})^T \mathbf{Q}_{k+1} \text{diag}(1/\mathbf{g}_{\mathbf{Q}_{k+1}}) \\ & \quad + (\nabla_{\mathbf{Q}_k}^{(A)})^T \mathbf{Q}_{k+1} \text{diag}(1/\mathbf{g}_{\mathbf{Q}_{k+1}}) - (\nabla_{\mathbf{Q}_k}^{(A)})^T \mathbf{Q}_k \text{diag}(1/\mathbf{g}_{\mathbf{Q}_k})\|_F \\ & \leq \sqrt{n} \|\nabla_{\mathbf{Q}_{k+1}}^{(A)} - \nabla_{\mathbf{Q}_k}^{(A)}\|_F \|\mathbf{Q}_{k+1}\|_F \|\text{diag}(1/\mathbf{g}_{\mathbf{Q}_{k+1}})\|_F \\ & \quad + \sqrt{n} \|\nabla_{\mathbf{Q}_k}^{(A)}\|_F \|(\mathbf{Q}_{k+1} \text{diag}(1/\mathbf{g}_{\mathbf{Q}_{k+1}}) - \mathbf{Q}_k \text{diag}(1/\mathbf{g}_{\mathbf{Q}_k}))\|_F \end{aligned}$$

Invoking the lower-bound on the marginal, continuing from the above display,

$$\begin{aligned} &\leq \frac{\sqrt{n}}{\delta} \|\nabla_{\mathbf{Q}_{k+1}}^{(A)} - \nabla_{\mathbf{Q}_k}^{(A)}\|_F \\ &\quad + \sqrt{n} \|\nabla_{\mathbf{Q}_k}^{(A)}\|_F \|(\mathbf{Q}_{k+1} \text{diag}(1/g_{\mathbf{Q}_{k+1}}) - \mathbf{Q}_k \text{diag}(1/g_{\mathbf{Q}_k}))\|_F \end{aligned}$$

Let us consider $\|\nabla_{\mathbf{Q}_k}^{(A)}\|_F \leq \|\mathbf{X}_k\|_F \|\mathbf{C}\|_F \|\mathbf{R}_k\|_F \leq \|\mathbf{X}_k\|_F \|\mathbf{C}\|_F$. We have that:

$$\|\mathbf{X}_k\|_F^2 = \sum_{i,j} \frac{1}{g_{(\mathbf{Q}_k)_i}^2} \mathbf{T}_{ij}^2 \frac{1}{g_{(\mathbf{R}_k)_j}^2}$$

Where we always have that $\mathbf{T}_{ij} \leq g_{\mathbf{Q}_i}$ and $\mathbf{T}_{ij} \leq g_{\mathbf{R}_j}$ by definition of \mathbf{T} as a coupling. As such:

$$\leq \sum_{ij} \frac{1}{g_{\mathbf{Q}_i}^2} (g_{\mathbf{Q}_i} g_{\mathbf{R}_j}) \frac{1}{g_{\mathbf{R}_j}^2} = \sum_{ij} \frac{1}{g_{\mathbf{Q}_i} g_{\mathbf{R}_j}} = \left\langle g_{\mathbf{Q}}^{-1} g_{\mathbf{R}}^{-T}, \mathbf{1}_m \mathbf{1}_r^T \right\rangle_F \leq \frac{mr}{\delta^2}$$

Thus $\|\nabla_{\mathbf{Q}_k}^{(A)}\|_F \leq \frac{\sqrt{mr}}{\delta} \|\mathbf{C}\|_F$, and the bound above reduces to

$$\begin{aligned} &\leq \frac{\sqrt{n}}{\delta} \|\nabla_{\mathbf{Q}_{k+1}}^{(A)} - \nabla_{\mathbf{Q}_k}^{(A)}\|_F \\ &\quad + \frac{\sqrt{nmr}}{\delta} \|\mathbf{C}\|_F \|(\mathbf{Q}_{k+1} \text{diag}(1/g_{\mathbf{Q}_{k+1}}) - \mathbf{Q}_k \text{diag}(1/g_{\mathbf{Q}_k}))\|_F \end{aligned}$$

Further bounding the last term in the norm

$$\begin{aligned} &\frac{\sqrt{nmr}}{\delta} \|\mathbf{C}\|_F \|(\mathbf{Q}_{k+1} \text{diag}(1/g_{\mathbf{Q}_{k+1}}) - \mathbf{Q}_k \text{diag}(1/g_{\mathbf{Q}_{k+1}}) \\ &\quad + \mathbf{Q}_k \text{diag}(1/g_{\mathbf{Q}_{k+1}}) - \mathbf{Q}_k \text{diag}(1/g_{\mathbf{Q}_k}))\|_F \\ &\leq \frac{\sqrt{nmr}}{\delta} \|\mathbf{C}\|_F (\|\text{diag}(1/g_{\mathbf{Q}_{k+1}})\|_F \|\mathbf{Q}_{k+1} - \mathbf{Q}_k\|_F \\ &\quad + \|\mathbf{Q}_k\|_F \|\text{diag}(1/g_{\mathbf{Q}_{k+1}}) - \text{diag}(1/g_{\mathbf{Q}_k})\|_F) \\ &\leq \frac{\sqrt{nmr}}{\delta} \|\mathbf{C}\|_F \left(\frac{1}{\delta} + \frac{\sqrt{r}}{\delta^2} \right) \|\mathbf{Q}_{k+1} - \mathbf{Q}_k\|_F \end{aligned}$$

Thus, the final bound on the $\nabla_{\mathbf{Q}}^{(B)}$ term is:

$$\begin{aligned} &\leq \frac{\sqrt{n}}{\delta} \|\nabla_{\mathbf{Q}_{k+1}}^{(A)} - \nabla_{\mathbf{Q}_k}^{(A)}\|_F + \frac{\sqrt{nmr}}{\delta} \|\mathbf{C}\|_F \|(\mathbf{Q}_{k+1} \text{diag}(1/g_{\mathbf{Q}_{k+1}}) - \mathbf{Q}_k \text{diag}(1/g_{\mathbf{Q}_k}))\|_F \\ &\leq \left(\frac{\|\mathbf{C}\|_F \sqrt{nr}}{\delta^4} + \frac{\sqrt{nmr}}{\delta^2} \|\mathbf{C}\|_F \left(1 + \frac{\sqrt{r}}{\delta} \right) \right) \|\mathbf{Q}_{k+1} - \mathbf{Q}_k\|_F \end{aligned}$$

The total component-wise smoothness bound on \mathbf{Q} is then

$$\begin{aligned} &\|\nabla_{\mathbf{Q}} \mathcal{L}_{\text{LC}}(\mathbf{Q}_{k+1}, \mathbf{R}_k, \mathbf{T}_k) - \nabla_{\mathbf{Q}} \mathcal{L}_{\text{LC}}(\mathbf{Q}_k, \mathbf{R}_k, \mathbf{T}_k)\|_F \\ &\leq \frac{\|\mathbf{C}\|_F}{\delta^2} \left(\left(\frac{\sqrt{nr}}{\delta^2} + \sqrt{nmr} \left(1 + \frac{\sqrt{r}}{\delta} \right) \right) + \frac{\sqrt{r}}{\delta} \right) \|\mathbf{Q}_{k+1} - \mathbf{Q}_k\|_F \end{aligned}$$

Identical reasoning applies for $\nabla_{\mathbf{R}} \mathcal{L}_{\text{LC}}$, where we similarly have the gradient split into two terms:

$$\begin{aligned} \nabla_{\mathbf{R}} \mathcal{L}_{\text{LC}}(\mathbf{Q}, \mathbf{R}, \mathbf{T}) &= \nabla_{\mathbf{R}}^{(A)} \mathcal{L}_{\text{LC}} \nabla_{\mathbf{R}}^{(B)} \mathcal{L}_{\text{LC}} \\ &= \mathbf{C}^T \mathbf{Q} \mathbf{X} - \mathbf{1}_m \text{diag}^{-1}(\text{diag}(1/g_{\mathbf{R}}) \mathbf{R}^T \mathbf{C}^T \mathbf{Q} \mathbf{X})^T \\ &= \nabla_{\mathbf{R}}^{(A)} \mathcal{L}_{\text{LC}} - \mathbf{1}_m \text{diag}^{-1}(\text{diag}(1/g_{\mathbf{R}}) \mathbf{R}^T \nabla_{\mathbf{R}}^{(A)} \mathcal{L}_{\text{LC}})^T \end{aligned}$$

As before, we may first show smoothness in $\nabla_{\mathbf{R}}^{(A)}$ using the same steps as for \mathbf{Q}

$$\begin{aligned}
& \|\nabla_{\mathbf{R}}^{(A)} \mathcal{L}_{\text{LC}}(\mathbf{Q}_{k+1}, \mathbf{R}_{k+1}, \mathbf{T}_k) - \nabla_{\mathbf{R}}^{(A)} \mathcal{L}_{\text{LC}}(\mathbf{Q}_{k+1}, \mathbf{R}_k, \mathbf{T}_k)\|_F \\
&= \|\mathbf{C}^T \mathbf{Q}_{k+1} \text{diag}(1/\mathbf{g}_{\mathbf{Q}_{k+1}}) \mathbf{T}_k \text{diag}(1/\mathbf{g}_{\mathbf{R}_{k+1}}) - \mathbf{C}^T \mathbf{Q}_{k+1} \text{diag}(1/\mathbf{g}_{\mathbf{Q}_{k+1}}) \mathbf{T}_k \text{diag}(1/\mathbf{g}_{\mathbf{R}_k})\|_F \\
&\leq \|\mathbf{C}\|_F \|\text{diag}(1/\mathbf{g}_{\mathbf{Q}_{k+1}})\|_F \|\mathbf{Q}_{k+1}\|_F \|\mathbf{T}_k\|_F \|\text{diag}(1/\mathbf{g}_{\mathbf{R}_{k+1}}) - \text{diag}(1/\mathbf{g}_{\mathbf{R}_k})\|_F \\
&\leq \frac{\|\mathbf{C}\|_F}{\delta} \|\text{diag}(1/\mathbf{g}_{\mathbf{R}_{k+1}}) - \text{diag}(1/\mathbf{g}_{\mathbf{R}_k})\|_F \\
&\leq \frac{\|\mathbf{C}\|_F}{\delta^3} \|\mathbf{g}_{\mathbf{R}_{k+1}} - \mathbf{g}_{\mathbf{R}_k}\|_2 \\
&\leq \frac{\|\mathbf{C}\|_F \sqrt{r}}{\delta^3} \|\mathbf{R}_{k+1} - \mathbf{R}_k\|_F.
\end{aligned}$$

For $\nabla_{\mathbf{R}}^{(B)}$, one may use the same reasoning as before to find:

$$\begin{aligned}
& \|\nabla_{\mathbf{R}}^{(B)} \mathcal{L}_{\text{LC}}(\mathbf{Q}_{k+1}, \mathbf{R}_{k+1}, \mathbf{T}_k) - \nabla_{\mathbf{R}}^{(B)} \mathcal{L}_{\text{LC}}(\mathbf{Q}_{k+1}, \mathbf{R}_k, \mathbf{T}_k)\|_F \\
&\leq \|\mathbf{1}_m \left[\text{diag}^{-1}(\text{diag}(1/\mathbf{g}_{\mathbf{R}_{k+1}}) \mathbf{R}_{k+1}^T \nabla_{\mathbf{R}_{k+1}}^{(A)} \mathcal{L}_{\text{LC}}) - \text{diag}^{-1}(\text{diag}(1/\mathbf{g}_{\mathbf{R}_k}) \mathbf{R}_k^T \nabla_{\mathbf{R}_k}^{(A)} \mathcal{L}_{\text{LC}}) \right]^T\|_F \\
&\leq \sqrt{m} \|\text{diag}(1/\mathbf{g}_{\mathbf{R}_{k+1}}) \mathbf{R}_{k+1}^T \nabla_{\mathbf{R}_{k+1}}^{(A)} \mathcal{L}_{\text{LC}} - \text{diag}(1/\mathbf{g}_{\mathbf{R}_k}) \mathbf{R}_k^T \nabla_{\mathbf{R}_k}^{(A)} \mathcal{L}_{\text{LC}}\|_F
\end{aligned}$$

As before, one may apply three rounds of triangle inequality inside the norm to bound this directly in terms of $\|\nabla_{\mathbf{R}_{k+1}}^{(A)} - \nabla_{\mathbf{R}_k}^{(A)}\|_F$, $\|\text{diag}(1/\mathbf{R}_{k+1}) - \text{diag}(1/\mathbf{R}_k)\|_F$, and $\|\mathbf{R}_{k+1} - \mathbf{R}_k\|_F$. Each of these terms is smooth in \mathbf{R} by the lower-bound argument, so that smoothness in \mathbf{R} holds analogously to \mathbf{Q} . The remainder of the proof for smoothness in \mathbf{R} thus follows identically to that of \mathbf{Q} above.

For \mathbf{T} , the component-wise bound of

$$\|\nabla_{\mathbf{T}_{k+1}} \mathcal{L}_{\text{LC}} - \nabla_{\mathbf{T}_k} \mathcal{L}_{\text{LC}}\|_F = \|\mathbf{Q}_{k+1} \mathbf{C} \mathbf{R}_{k+1}^T - \mathbf{Q}_{k+1} \mathbf{C} \mathbf{R}_k^T\|_F \leq L_T \|\mathbf{T}_{k+1} - \mathbf{T}_k\|_F$$

holds trivially for any $L_T > 0$ as the gradient is uniquely determined by \mathbf{Q} and \mathbf{R} alone. Thus there exist $L_Q, L_R, L_T > 0$ as component-wise smoothness constants for $\mathcal{L}_{\text{LC}}(\mathbf{Q}, \mathbf{R}, \mathbf{T})$. \square

We next prove Proposition 3.4, restated just below for convenience.

Proposition E.6 (Proposition 3.4). *Consider the FRLC objective (8). The FRLC algorithm, Algorithm 4, yields β -smooth iterates for $\beta = \text{poly}(\|\mathbf{C}\|_F, m, r, \delta)$, where δ denotes the lower-bound on the entries of $\mathbf{g}_{\mathbf{Q}}, \mathbf{g}_{\mathbf{R}}$. Consider the convergence metric of 3.3 adapted from Ghadimi et al. (2014), given as:*

$$\Delta_k(\mathbf{x}_k, \mathbf{x}_{k+1}) = \sum_{i=1}^p \|\mathbf{g}_{\mathcal{X}, k, i}\|^2 = \frac{1}{\gamma_k^2} [\|\mathbf{Q}_k - \mathbf{Q}_{k-1}\|_F^2 + \|\mathbf{R}_k - \mathbf{R}_{k-1}\|_F^2 + \|\mathbf{T}_k - \mathbf{T}_{k-1}\|_F^2]$$

for $\mathbf{x}_k = (\mathbf{Q}_k, \mathbf{T}_k, \mathbf{R}_k)$. Define the gap to the optimal solution D as in (36), and let $L = \sup_i(L_i)$ to be the global smoothness constant across all components. Then for $\gamma_k = \alpha/L$ as defined in 3.3 the FRLC algorithm has the non-asymptotic stationary convergence guarantee that:

$$\min_{k \in \{1, \dots, N-1\}} \Delta_k \leq \frac{2D^2 L^2}{N\alpha^2}$$

Proof. The proof of the non-asymptotic stationary convergence of mirror descent of Ghadimi et al. (2014), adapted for coordinate mirror descent using the block-descent lemma in 3.3, only requires component-wise smoothness in $(\mathbf{Q}, \mathbf{R}, \mathbf{T})$. The proof of this for FRLC is given in ??, and the guarantee follows directly for this value of $L = \max(L_Q, L_R, L_T) = \text{poly}(\|\mathbf{C}\|_F, m, r, \delta)$. \square

Proposition E.7. Low-rank Approximation Error. *Let SR-W_r^* denote the optimal rank- r approximation for the semi-relaxed low-rank optimal transport problem, and let SR-W^* denote the optimal solution for the full-rank semi-relaxed optimal transport problem.*

Additionally, suppose $\mathbf{c}_b = \sum_{j=1}^m \mathbf{b}_j$ denotes the sum of the entries of the second marginal ($\mathbf{c}_b = 1$ if a probability measure). Then we have the following upper-bound on the objective error:

$$|\text{SR-W}_r^*(\mu_b) - \text{SR-W}^*(\mu_b)| \leq \mathbf{c}_b \left(\max_{p,q} \{C_{pq}\} - \min_{p,q} \{C_{pq}\} \right) \ln(\min\{n, m\}/(r-1))$$

We note that this bound also applies for the standard balanced optimal transport case, giving:

$$|W_r^*(\mu_a, \mu_b) - W^*(\mu_a, \mu_b)| \leq \left(\max_{p,q} \{C_{pq}\} - \min_{p,q} \{C_{pq}\} \right) \ln(\min\{n, m\}/(r-1))$$

and improves the previous bound of $|W_r^*(\mu_a, \mu_b) - W^*(\mu_a, \mu_b)| \leq \|C\|_\infty \ln(\min\{n, m\}/(r-1))$ as the distance matrix C contains only non-negative entries.

Proof. We adapt the proof from Scetbon & Cuturi (2022) for the balanced case, which previously gave the bound:

$$|W_r^*(\mu_a, \mu_b) - W^*(\mu_a, \mu_b)| \leq \|C\|_\infty \ln(\min\{n, m\}/(r-1))$$

In particular, for $z = \min\{m, n\}$, there exists an optimal $\text{rank}_+(\mathbf{P}^*) \leq z$ where one may express the optimal solution for the non-negative coupling matrix \mathbf{P} as a sum of z rank-one, non-negative outer products $\tilde{\mathbf{q}}_k \tilde{\mathbf{r}}_k^\top \succcurlyeq 0$:

$$\mathbf{P}^* = \sum_{k=1}^z \tilde{\mathbf{q}}_k \tilde{\mathbf{r}}_k^\top = \sum_{k=1}^z \lambda_k \mathbf{q}_k \mathbf{r}_k^\top$$

Where we write this sum in terms of normalized vectors $\mathbf{q}_k = \tilde{\mathbf{q}}_k / \|\tilde{\mathbf{q}}_k\|_1$, $\mathbf{r}_k = \tilde{\mathbf{r}}_k / \|\tilde{\mathbf{r}}_k\|_1$, and $\lambda_k = \|\tilde{\mathbf{r}}_k\|_1 \|\tilde{\mathbf{q}}_k\|_1$. Without any loss of generality, $(\lambda_k)_{k=1}^z$ is ordered in terms of decreasing value such that $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_z$. Note that we have a fixed constraint for the sum of the entries of \mathbf{P} for the semi-relaxed case, assuming \mathbf{b} is a general positive measure, as $\mathbf{1}_n^\top \mathbf{P} \mathbf{1}_m = (\mathbf{P}^\top \mathbf{1}_n)^\top \mathbf{1}_m = \mathbf{b}^\top \mathbf{1}_m = \sum_{j=1}^m \mathbf{b}_j := c_b$ (where $c_b = 1$ if \mathbf{b} is chosen to be a probability measure, i.e. in the balanced case). Moreover, it is simple to observe for these ordered values that $\lambda_k \leq (c_b/k)$. As in Scetbon & Cuturi (2022), define the weighted average of the bottom $z - r + 1$ vectors of the decomposition to be:

$$\begin{aligned} \boldsymbol{\alpha}_r &= \frac{\sum_{i=r}^z \lambda_i \mathbf{q}_i}{\sum_{i=r}^z \lambda_i} \\ \boldsymbol{\beta}_r &= \frac{\sum_{i=r}^z \lambda_i \mathbf{r}_i}{\sum_{i=r}^z \lambda_i} \end{aligned}$$

And take the rank- r approximation using the optimal $r-1$ vectors of OPT and this weighted average of the bottom to be:

$$\tilde{\mathbf{P}}_r = \sum_{i=1}^{r-1} \lambda_i \mathbf{q}_i \mathbf{r}_i^\top + \left(\sum_{i=r}^z \lambda_i \right) \boldsymbol{\alpha}_r \boldsymbol{\beta}_r^\top$$

Where, by the assumption that $\mathbf{P}^* \in \Pi_b$ is feasible:

$$\tilde{\mathbf{P}}_r^\top \mathbf{1}_n = \sum_{i=1}^{r-1} \lambda_i \mathbf{r}_i \mathbf{q}_i^\top \mathbf{1}_n + \left(\sum_{i=r}^z \lambda_i \right) \boldsymbol{\beta}_r \boldsymbol{\alpha}_r^\top \mathbf{1}_n = \sum_{i=1}^{r-1} \lambda_i \mathbf{r}_i + \left(\sum_{i=r}^z \lambda_i \right) \boldsymbol{\beta}_r = \sum_{i=1}^z \lambda_i \mathbf{r}_i = \mathbf{b}$$

Thus $\tilde{\mathbf{P}}_r \in \Pi_{b,r}$ is a feasible rank- r solution by feasibility of \mathbf{P}^* . One can verify that if $\mathbf{P}^* \in \Pi_{a,b}$, then $\tilde{\mathbf{P}}_r \mathbf{1}_m = \mathbf{a}$ and the solution is again feasible. From this, we observe that the difference between this solution and \mathbf{P}^* is an upper-bound to the difference between \mathbf{P}^* and the optimal rank- r solution:

$$\begin{aligned} |\text{SR-}W_r^*(\mu_a, \mu_b) - \text{SR-}W^*(\mu_a, \mu_b)| &= |\langle \tilde{\mathbf{P}}_r, \mathbf{C} \rangle_F - \langle \mathbf{P}^*, \mathbf{C} \rangle_F| \leq \langle \tilde{\mathbf{P}}_r, \mathbf{C} \rangle_F - \langle \mathbf{P}^*, \mathbf{C} \rangle_F \\ &= \left\langle \sum_{i=1}^{r-1} \lambda_i \mathbf{q}_i \mathbf{r}_i^\top + \left(\sum_{i=r}^z \lambda_i \right) \boldsymbol{\alpha}_r \boldsymbol{\beta}_r^\top - \sum_{i=1}^z \lambda_i \mathbf{q}_i \mathbf{r}_i^\top, \mathbf{C} \right\rangle_F \\ &= \left\langle \left(\sum_{i=r}^z \lambda_i \right) \boldsymbol{\alpha}_r \boldsymbol{\beta}_r^\top - \sum_{i=r}^z \lambda_i \mathbf{q}_i \mathbf{r}_i^\top, \mathbf{C} \right\rangle_F \end{aligned}$$

Noting that α_r, β_r and q_i, r_i are unit normalized positive vectors, the sum of the entries of the outer product $\mathbf{1}_n^T q_i r_i^T \mathbf{1}_m = 1$, and likewise for $\alpha_r \beta_r^T$. Thus, continuing from the above display:

$$\begin{aligned}
&= \left(\sum_{i=r}^z \lambda_i \langle \alpha_r \beta_r^T, \mathbf{C} \rangle_F - \sum_{i=r}^z \lambda_i \langle q_i r_i^T, \mathbf{C} \rangle_F \right) \\
&\leq \left(\sum_{i=r}^z \lambda_i \langle \alpha_r \beta_r^T, \mathbf{C} \rangle_F - \left(\sum_{i=r}^z \lambda_i \right) \min_{p,q} \{C_{pq}\} \right) \\
&\leq \left(\max_{p,q} \{C_{pq}\} - \min_{p,q} \{C_{pq}\} \right) \sum_{i=r}^z \lambda_i \\
&\leq c_b \left(\max_{p,q} \{C_{pq}\} - \min_{p,q} \{C_{pq}\} \right) \sum_{i=r}^z \frac{1}{i} \\
&\leq c_b \left(\max_{p,q} \{C_{pq}\} - \min_{p,q} \{C_{pq}\} \right) \ln(z/(r-1))
\end{aligned}$$

Concluding the proof. As discussed, this directly applies to the balanced case (for $c_b = 1$). \square

F Initialization

We propose a new initialization of the sub-couplings Q, R, T for the LC-factorization. Algorithm 6 generates a random full-rank initial condition in the set of couplings $\Pi_{a,b}$ which still satisfies the marginal constraints. It accomplishes this by sampling random matrices which are full-rank and applying the Sinkhorn algorithm to each of them. Scetbon et al. (2021) proposed an initialization which represents an improvement over the rank-1 product measure which is rank-2. Follow-up work proposed initialization using k-means Scetbon & Cuturi (2022). However, this assumes the previous diagonal factorization and is thus not application for generating a latent coupling which may be non-diagonal, non-square, and with two distinct inner marginals. Our initialization is tailored to the LC-factorization, is effective, and has a full-rank guarantee. In particular, higher-rank initializations may exhibit better convergence properties by allowing the gradient to explore a larger set of directions immediately in the optimization. This initialization is given in Algorithm 6.

Algorithm 6 Initialize-Couplings

Input $\mathbf{a} \in \Delta_n, \mathbf{b} \in \Delta_m, \mathbf{g}_Q \in \Delta_r, \mathbf{g}_R \in \Delta_r$
 $\mathbf{C}_Q \sim [0, 1]^{n \times r}, \mathbf{C}_R \sim [0, 1]^{m \times r}, \mathbf{C}_T \sim [0, 1]^{r \times r}$
 $\mathbf{K}_Q \leftarrow e^{\mathbf{C}_Q}, \mathbf{K}_R \leftarrow e^{\mathbf{C}_R}, \mathbf{K}_T \leftarrow e^{\mathbf{C}_T}$
 $\mathbf{Q} \leftarrow \text{Sinkhorn}(\mathbf{K}_Q, \mathbf{a}, \mathbf{g}_Q)$
 $\mathbf{R} \leftarrow \text{Sinkhorn}(\mathbf{K}_R, \mathbf{b}, \mathbf{g}_R)$
 $\mathbf{T} \leftarrow \text{Sinkhorn}(\mathbf{K}_T, \mathbf{g}_Q = \mathbf{Q}^T \mathbf{1}_n, \mathbf{g}_R = \mathbf{R}^T \mathbf{1}_m)$
Return $(\mathbf{Q}, \mathbf{R}, \mathbf{T})$

Proposition F.1. *Suppose one samples an initial condition on the optimal transport coupling using Algorithm 6, where we assume $C_{ij} \sim \text{Unif}(0, 1)$ such that $\mathbb{P}(\text{rank}(\mathbf{C}) < \min\{n, m\}) = 0$. Additionally, suppose that $\mathbf{a}, \mathbf{b} > \mathbf{0}$ holds elementwise for both marginals $\mathbf{a} \in \Delta_n, \mathbf{b} \in \Delta_r$. Then the elementwise exponential $\exp\{\mathbf{C}\}$ (or $\exp\{-\mathbf{C}\}$) has full-rank and the return $\text{Sinkhorn}(e^{-\mathbf{C}}, \mathbf{a}, \mathbf{b})$ has full-rank.*

Proof. It is established that a random matrix $\mathbf{C} \sim [0, 1]^{n \times m}$ has full-rank with probability one. For $\mathbf{K} = \exp\{\mathbf{C}\}$, it holds that the matrix must be entry-wise positive with $K_{ij} \geq 0$. If columns $\mathbf{C}_{\cdot,i} \neq \mathbf{C}_{\cdot,j}$ then clearly $\mathbf{C}_{\cdot,i}^{\odot k} \neq \mathbf{C}_{\cdot,j}^{\odot k}$, and if $\mathbf{C}_{\cdot,i}, \mathbf{C}_{\cdot,j} \succcurlyeq \mathbf{0}$ and are independent remain so under element-wise powers. One may easily show this by contrapositive. Suppose there exist constants c_1, c_2 such that:

$$c_1 \mathbf{C}_{\cdot,i}^{\odot k} + c_2 \mathbf{C}_{\cdot,j}^{\odot k} = \mathbf{0}$$

As $C_{.,i}, C_{.,j} \succcurlyeq \mathbf{0}$, without loss of generality one may assume $c_1 > 0$ and $c_2 < 0$. Then:

$$c_1 C_{.,i}^{\odot k} = c_1 \mathbf{1} \odot C_{.,i}^{\odot k} = -c_2 \mathbf{1} \odot C_{.,j}^{\odot k} = -c_2 C_{.,j}^{\odot k} \implies \left(-\frac{c_1}{c_2}\right)^{1/k} \mathbf{1} = c_1 = \frac{C_{.,j}}{C_{.,i}}$$

So clearly one has that $C_{.,j} - cC_{.,i} = \mathbf{0}$ for $c > 0$. This implies the columns $C_{.,j}$ and $C_{.,i}$ are dependent. Thus it is clear that elementwise powers of entrywise positive independent vectors preserve independence. The same principle extends trivially to exponentiation of the columns, where if one assumes by contradiction that $c_1 e^{C_{.,i}} + c_2 e^{C_{.,j}} = \mathbf{0}$ for $c_1 > 0, c_2 < 0$, one finds $\log\left(-\frac{c_1}{c_2}\right) \mathbf{1} = C_{.,j} - C_{.,i}$. Without loss of generality, assume $0 < -c_2 \leq c_1$ and $C_{.,j} > C_{.,i} > \mathbf{0}$, so that $\delta = \log\left(-\frac{c_1}{c_2}\right) \geq 0$. Then, considering constants q_1, q_2 :

$$q_2 C_{.,j} + q_1 C_{.,i} = (q_1 + q_2) C_{.,i} + \delta q_2 \mathbf{1} = \mathbf{0}$$

Assuming $C_{.,i}$ has greater than one unique entry, which we assume as the entries are sampled densely in \mathbb{R} , the two vectors are dependent if and only if $q_1 = -q_2$ and $\delta = 0$, implying $C_{.,i} = C_{.,j}$. Thus, for the set of independent column vectors of C , given as $\{C_{.,i}\}_{i=1}^m$, the set $\{e^{C_{.,i}}\}_{i=1}^m$ is also linearly independent. This holds analogously for the row vectors. As C is full-rank and $\text{span}(\{C_{.,i}\}) = \mathbb{R}^{\min\{m,n\}}$, we have that $\text{span}(\mathbf{K}) = \mathbb{R}^{\min\{m,n\}}$ as $\mathbf{K} = e^C = \sum_{k=0}^{\infty} \frac{C^{\odot k}}{k!}$ (analogously e^{-C}) and remains full-rank.

Sinkhorn expresses each variable as

$$\mathbf{X} = \text{diag}(\mathbf{u}) \mathbf{K} \text{diag}(\mathbf{v})$$

where $\text{rank}(\mathbf{X}) = \text{rank}(\text{diag}(\mathbf{u}) \mathbf{K} \text{diag}(\mathbf{v}))$. As Cuturi (2013b) updates the vectors $\mathbf{u} \leftarrow \mathbf{a}/\mathbf{K}\mathbf{v}$ and $\mathbf{v} \leftarrow \mathbf{b}/\mathbf{K}^T \mathbf{u}$ from $\mathbf{u}_0 = \mathbf{1}_n$ and $\mathbf{v}_0 = \mathbf{1}_m$, if $\mathbf{a}, \mathbf{b} > \mathbf{0}$ holds element-wise, one has that $\mathbf{u}, \mathbf{v} > \mathbf{0}$ elementwise as well. Then, one has that $\text{null}(\text{diag}(\mathbf{v})) = \{\mathbf{0}\}$ and $\text{null}(\text{diag}(\mathbf{u})) = \{\mathbf{0}\}$, implying that $\text{rank}(\text{diag}(\mathbf{u}) \mathbf{K} \text{diag}(\mathbf{v})) = \text{rank}(\mathbf{K}) = \min\{n, m\}$.

Thus, our initialization returns a random coupling matrix $\mathbf{X} \in \Pi_{\mathbf{a}, \mathbf{b}}$ of full-rank. \square

In the next proposition, we show that one can *analytically* solve for the block-optimal weights \mathbf{g} for the factorization of the coupling matrix \mathbf{P} as $\mathbf{P} = \mathbf{Q} \text{diag}(1/\mathbf{g}) \mathbf{R}^T$ Forrow et al. (2019); Scetbon et al. (2021).

Proposition F.2. *For the minimization problem expressed as*

$$\min_{\mathbf{g} \in \Delta_r} \langle \mathbf{Q} \text{diag}(1/\mathbf{g}) \mathbf{R}^T, \mathbf{C} \rangle_F$$

One has the closed-form minimizer of \mathbf{g}^ defined entrywise as:*

$$g_i^* = \frac{\sqrt{\omega_i}}{\sum_{j=1}^r \sqrt{\omega_j}}$$

For $\omega = \text{diag}^{-1}(\mathbf{Q}^T \mathbf{C} \mathbf{R})$, when $\omega \geq \mathbf{0}$ holds entrywise.

Proof. As $\omega \geq \mathbf{0}$, we consider the simplex condition $\sum_{j=1}^r g_j = 1$ alone. Writing out the Lagrangian associated to our objective, with $\lambda \in \mathbb{R}$ our equality-condition dual variable, we have:

$$\mathcal{L}(\mathbf{g}, \lambda) = \langle \mathbf{Q} \text{diag}(1/\mathbf{g}) \mathbf{R}^T, \mathbf{C} \rangle_F + \lambda(1 - \sum_j g_j)$$

Let us consider a rewriting of the inner product term:

$$\begin{aligned} \langle \mathbf{Q} \text{diag}(1/\mathbf{g}) \mathbf{R}^T, \mathbf{C} \rangle_F &= \sum_{i=1}^n \sum_{j=1}^m C_{ij} \sum_{k=1}^r Q_{ik} \left(\frac{1}{g_k}\right) R_{kj}^T \\ &= \sum_{k=1}^r \left(\frac{1}{g_k}\right) \sum_{i=1}^n \sum_{j=1}^m Q_{ki}^T C_{ij} R_{jk} \\ &= \sum_{k=1}^r \left(\frac{1}{g_k}\right) (\mathbf{Q}^T \mathbf{C} \mathbf{R})_{k,k} \\ &= \sum_{k=1}^r \frac{\omega_k}{g_k} = \omega^T (1/\mathbf{g}) \end{aligned}$$

Where division is interpreted element-wise in the last line. Thus, one can interpret the problem as minimizing the weighted sum of reciprocals of a density. As a result, we can simplify our Lagrangian's Froebenius inner product to a vector dot-product as:

$$\mathcal{L}(\mathbf{g}, \lambda) = \boldsymbol{\omega}^T (1/\mathbf{g}) - \lambda(1 - \sum_{j=1}^r \mathbf{g}_j)$$

Thus, the first order condition tells us that the value of the coupling weight \mathbf{g}_j is related to λ as:

$$\partial_{\mathbf{g}_j} \mathcal{L}(\mathbf{g}, \lambda) = -\frac{\boldsymbol{\omega}_j}{\mathbf{g}_j^2} + \lambda = 0 \implies \mathbf{g}_j = \sqrt{\frac{\boldsymbol{\omega}_j}{\lambda}}$$

And by relying on the summation condition on the probability density \mathbf{g} , yields the Langrange multiplier as

$$\sum_{j=1}^r \mathbf{g}_j = 1 = \sum_{j=1}^r \sqrt{\frac{\boldsymbol{\omega}_j}{\lambda}}$$

so that one finds

$$1 = \frac{1}{\sqrt{\lambda}} \sum_{j=1}^r \sqrt{\boldsymbol{\omega}_j} \implies \lambda = \left(\sum_{j=1}^r \sqrt{\boldsymbol{\omega}_j} \right)^2$$

Plugging our Lagrange-multiplier into the above expression yields:

$$\mathbf{g}_j = \sqrt{\frac{\boldsymbol{\omega}_j}{\lambda}} = \sqrt{\frac{\boldsymbol{\omega}_j}{\left(\sum_{i=1}^r \sqrt{\boldsymbol{\omega}_i}\right)^2}} = \frac{\sqrt{\boldsymbol{\omega}_j}}{\sum_{i=1}^r \sqrt{\boldsymbol{\omega}_i}}$$

As the Hessian $\nabla_{\mathbf{g}}^2 \mathcal{L} = \text{diag}\left(\frac{\boldsymbol{\omega}}{\mathbf{g}^3}\right) \succcurlyeq \mathbf{0}$, we conclude that this value of \mathbf{g} indeed minimizes the loss over Δ_r . \square

G Alternating updates on the dual variables

For the problem:

$$\inf_{(\mathbf{Q}, \mathbf{R}_k, \mathbf{g}_k) \in \mathcal{C}_1 \cap \mathcal{C}_2} \left(\frac{1}{\gamma_k} \text{KL}(\mathbf{Q} \| \mathbf{K}_Q) + \tau \text{KL}(\mathbf{Q} \mathbf{1}_r \| \mathbf{a}) - \boldsymbol{\lambda}_1^T \mathbf{Q}^T \mathbf{1}_n \right) \quad (40)$$

one can find a simple set of semi-relaxed updates for the coupling matrix. We note the primal-dual relationship of Sinkhorn, $\mathbf{Q} = \text{diag}(e^{\gamma_k \mathbf{f}_1}) \mathbf{K}_Q \text{diag}(e^{\gamma_k \mathbf{h}_1})$, and consider the entry-wise first-order condition required for the sub-coupling \mathbf{Q} :

$$\begin{aligned} 0 &= \gamma_k^{-1} \log \left(\frac{\mathbf{Q}_{ij}}{(\mathbf{K}_Q)_{ij}} \right) + \tau \log \left(\frac{\langle \mathbf{Q}_{i,\cdot}, \mathbf{1}_r \rangle}{\mathbf{a}_i} \right) - \boldsymbol{\lambda}_{1,i} \\ \implies \log \mathbf{Q}_{ij} &= \tau \gamma_k \log \left(\frac{\mathbf{a}_i}{\langle \mathbf{Q}_{i,\cdot}, \mathbf{1}_r \rangle} \right) + \log (\mathbf{K}_Q)_{ij} - \boldsymbol{\lambda}_{1j} \gamma_k \end{aligned} \quad (41)$$

Thus:

$$\mathbf{Q}_{ij} = \left(\frac{\mathbf{a}_i}{\mathbf{Q}_{i,\cdot}^T \mathbf{1}_r} \right)^{\tau \gamma_k} (\mathbf{K}_Q)_{ij} e^{-\boldsymbol{\lambda}_{1j} \gamma_k}$$

And in matrix-form, this yields:

$$\begin{aligned} \mathbf{Q} &= \text{diag} \left(\frac{\mathbf{a}}{\mathbf{Q} \mathbf{1}_r} \right)^{\tau \gamma_k} \mathbf{K}_Q \text{diag}(e^{-\gamma_k \boldsymbol{\lambda}_1}) \\ &= \text{diag}(e^{\gamma_k \mathbf{f}_1}) \mathbf{K}_Q \text{diag}(e^{\gamma_k \mathbf{h}_1}) \end{aligned}$$

And expanding the $\mathbf{Q}\mathbf{1}_r$ term explicitly, noting that $\mathbf{X} \text{diag}(\mathbf{v})\mathbf{1} = \mathbf{X}\mathbf{v}$, we have:

$$\begin{aligned} & \text{diag}\left(\frac{\mathbf{a}}{\mathbf{Q}\mathbf{1}_r}\right)^{\tau\gamma_k} \mathbf{K}_Q \text{diag}(e^{-\gamma_k\lambda_1}) \\ &= \text{diag}\left(\frac{\mathbf{a}}{\text{diag}(e^{\gamma_k\mathbf{f}_1})\mathbf{K}_Q \text{diag}(e^{\gamma_k\mathbf{h}_1})\mathbf{1}_r}\right)^{\tau\gamma_k} \mathbf{K}_Q \text{diag}(e^{-\gamma_k\lambda_1}) \\ &= \text{diag}\left(\frac{\mathbf{a}}{e^{\gamma_k\mathbf{f}_1} \odot \mathbf{K}_Q e^{\gamma_k\mathbf{h}_1}}\right)^{\tau\gamma_k} \mathbf{K}_Q \text{diag}(e^{-\gamma_k\lambda_1}) \end{aligned}$$

Thus, we identify $e^{\gamma_k\mathbf{h}_1} = e^{-\gamma_k\lambda_1}$ as the right dual vector, and identify the following relationship in terms of the left dual vector:

$$\left(\frac{\mathbf{a}}{e^{\gamma_k\mathbf{f}_1} \odot \mathbf{K}_Q e^{\gamma_k\mathbf{h}_1}}\right)^{\tau\gamma_k} = e^{\gamma_k\mathbf{f}_1} \implies e^{\gamma_k\mathbf{f}_1} = \left(\frac{\mathbf{a}}{\mathbf{K}_Q e^{\gamma_k\mathbf{h}_1}}\right)^{\frac{\tau}{\tau+1/\gamma_k}}$$

From 40, the condition that $(\mathbf{Q}, \mathbf{R}_k, \mathbf{g}_k) \in \mathcal{C}_1 \cap \mathcal{C}_2$ implies that $\mathbf{Q}^T \mathbf{1}_m = \mathbf{g}_k := \mathbf{g}$. As such, we find that:

$$\mathbf{Q}^T \mathbf{1}_m = \text{diag}(e^{\gamma_k\mathbf{h}_1}) \mathbf{K}_Q^T \text{diag}(e^{\gamma_k\mathbf{f}_1}) \mathbf{1}_m = \text{diag}(e^{\gamma_k\mathbf{h}_1}) \mathbf{K}_Q^T e^{\gamma_k\mathbf{f}_1} = \mathbf{g}$$

Implying an update for $e^{\gamma_k\mathbf{h}_1}$ in the form:

$$e^{\gamma_k\mathbf{h}_1} = \left(\frac{\mathbf{g}}{\mathbf{K}_Q^T e^{\gamma_k\mathbf{f}_1}}\right)$$

Analogous reasoning applies for a relaxation of the other marginal, yielding the SR^R -projection and SR^L -projection (i.e. semi-relaxed OT).

Algorithm 7 SR^L -projection (semi-relaxed OT, left marginal relaxed)

Input $\mathbf{K}, \gamma, \tau, \mathbf{a}, \mathbf{b}, \delta$
 $\mathbf{u} \leftarrow \mathbf{1}_n$
 $\mathbf{v} \leftarrow \mathbf{1}_r$
repeat
 $\tilde{\mathbf{u}} \leftarrow \mathbf{u}$
 $\tilde{\mathbf{v}} \leftarrow \mathbf{v}$
 $\mathbf{u} \leftarrow (\mathbf{a}/\mathbf{K}\mathbf{v})^{\tau/(\tau+\gamma^{-1})}$
 $\mathbf{v} \leftarrow (\mathbf{b}/\mathbf{K}^T\mathbf{u})$
until $\gamma^{-1} \max\{\|\log \tilde{\mathbf{u}}/\mathbf{u}\|_\infty, \|\log \tilde{\mathbf{v}}/\mathbf{v}\|_\infty\} < \delta$
return $\text{diag}(\mathbf{u})\mathbf{K} \text{diag}(\mathbf{v})$

H Discussion of Complexity

For $(\mathbf{Q}, \mathbf{R}, \mathbf{T}) \in \mathbb{R}_+^{n \times r_1} \times \mathbb{R}_+^{m \times r_2} \times \mathbb{R}_+^{r_1 \times r_2}$, the space complexity $O(nr_1 + r_1r_2 + mr_2)$ is linear if the ranks $r_1, r_2 = o(1)$ are taken to be small constants. The time-complexity of Algorithm 4 is $O(BLr^2(n+m))$ for B the number of inner Sinkhorn iterations, L the number of mirror-descent steps, n, m the number of samples in the first and second dataset, and $r = \max\{r_1, r_2, d\}$ for r_1, r_2 the ranks of the latent coupling and d the rank of the factorized distance matrix \mathbf{C} (generally chosen to be a constant near r_1, r_2). Each matrix-multiplication is of max order $(n+m)r^2$, which happens a constant number of times in the computation of each gradient ∇_i , and for the respective Sinkhorn matrix-vector multiplications $\mathbf{K}\mathbf{v}$ and $\mathbf{K}^T\mathbf{u}$. The L outer steps follow from the mirror-descent convergence rate and the number of iterations B required for each projection follow from the convergence of Sinkhorn. In particular, for ε a fixed error tolerance and η the entropy constant, one finds a $\pm\varepsilon D$ approximation for D the diameter of the data in $B = \text{poly}(1/\eta\varepsilon)$ iterations using the Sinkhorn algorithm Charikar et al. (2023); Cuturi (2013a).

I Review of Background Material

I.1 Low-Rank Approximation of Pairwise Distance Matrices

As mentioned previously, works such as Charikar et al. (2023) have developed algorithms with linear $O((n+m)^{1+o(1)}\text{poly}(1/\epsilon))$ time-complexity and $O((n+m)d)$ space-complexity for sketching the optimal transport cost *value*. Recent works on low-rank factorization of the optimal transport coupling matrix \mathbf{P} (the matrix associated to the coupling $\gamma \in \Pi(\mu, \nu)$) Scetbon & Cuturi (2022); Scetbon et al. (2023, 2021) have achieved per-iteration time-complexities of $O(T(n+m)dr)$ for some constant non-negative rank $r \geq 1$, d the dimension of the metric space, and T the number of iterations. By the JL-lemma one can simply embed the points in dimension $d = O(\log(nm)/\epsilon^2)$ while preserving pairwise distances, however, currently no proofs exist which offer the number of iterations T until convergence to some tolerance ϵ . This is partially due to how recent these works are, and also to the non-convexity of the objective which is sensitive to initial conditions Scetbon et al. (2021). However, the space complexity of the algorithm is $O((n+m)dr)$, which is noteworthy for being *linear* in the number of points and avoids storing the potentially intractable $O(nm)$ coupling matrix \mathbf{P} . To accomplish this, however, these works rely on a low-rank approximation of the pairwise distance matrix. A number of works by Indyk and Woodruff have concerned algorithms for finding low-rank approximations for such distance matrices. A seminal work Bakshi & Woodruff (2018) developed an algorithm which, given two point sets $\{\mathbf{x}_i\}_{i=1}^n$ and $\{\mathbf{y}_j\}_{j=1}^m$ in some metric space \mathcal{X} , finds a rank r approximation in $O((n+m)^{1+\gamma}\text{poly}(r, 1/\epsilon))$ for $\gamma > 0$ an arbitrarily small constant and $\epsilon > 0$ an error parameter. A more recent work Indyk et al. (2019) improves on this one, by reading a sample-optimal $O((n+m)r/\epsilon)$ entries of the input matrix with a run-time which removes dependence on γ that is merely $O((n+m)\text{poly}(r, 1/\epsilon))$. This algorithm is used by all of the low-rank optimal transport works. These works, by finding a low-rank approximation to the coupling matrix $\mathbf{P} \approx \mathbf{A}\mathbf{B}^T \in \mathbb{R}_+^{n \times m}$ due to space limitations on the coupling, necessarily cannot store the full distance matrix $\mathbf{C} \in \mathbb{R}_+^{n \times m}$ of the same size in memory either. As such, it must also be approximated as $\mathbf{C} \approx \mathbf{V}\mathbf{U}$ before input to the algorithm, where we necessarily require very effective approximations of \mathbf{U} and \mathbf{V} to tolerate the additional source of error from coarse-graining the distance matrix to be low-rank. As such, we present some of the details and algorithm of Indyk et al. (2019) as an essential component of the existing low-rank optimal transport solvers. We begin by summarizing the main theorems in Indyk et al. (2019), which provide an algorithm (upper-bound) on the low-rank distance-matrix approximation problem and a lower-bound on the number of entries which must be read.

Theorem I.1. *Indyk et al. (2019) There is a randomized algorithm that, given a distance matrix $\mathbf{C} \in \mathbb{R}^{n \times m}$, reads $O((n+m)r/\epsilon)$ entries of \mathbf{C} , runs in time $\tilde{O}(n+m) \cdot \text{poly}(r, 1/\epsilon)^1$ and computes low-rank factors $\mathbf{V} \in \mathbb{R}^{n \times r}$, $\mathbf{U} \in \mathbb{R}^{r \times m}$ that with probability 0.99 satisfy:*

$$\|\mathbf{C} - \mathbf{V}\mathbf{U}\|_F^2 \leq \|\mathbf{C} - \mathbf{C}_r\|_F^2 + \epsilon\|\mathbf{C}\|_F^2 \quad (42)$$

For \mathbf{C}_r the optimal rank- r approximation of \mathbf{C} .

This is a remarkable result, especially in light of the next theorem.

Theorem I.2. *Indyk et al. (2019) Let $r \leq m \leq n$ and $\epsilon > 0$ such that $r/\epsilon = O(\min\{m, n^{1/3}\})$. Any randomized and possibly adaptive algorithm that given a distance matrix $\mathbf{C} \in \mathbb{R}^{n \times m}$ computes $\mathbf{V} \in \mathbb{R}^{n \times r}$, $\mathbf{U} \in \mathbb{R}^{r \times m}$ satisfying $\|\mathbf{C} - \mathbf{V}\mathbf{U}\|_F^2 \leq \|\mathbf{C} - \mathbf{C}_r\|_F^2 + \epsilon\|\mathbf{C}\|_F^2$ must read $\Omega((n+m)r/\epsilon)$ entries of \mathbf{C} in expectation. This lower bound also holds for symmetric distance matrices $\mathbf{C} \in \mathcal{S}_n$.*

The lower-bound follows from a difficult argument which involves constructing a hard distribution over distance matrices, involving the use of random matrix theory. The upper-bound, however, follows relatively straightforwardly from a number of previous algorithms and their associated guarantees, along with the algorithm presented below. We introduce the algorithm and also offer a proof of I.1 for completeness.

¹Where $\tilde{O}(\cdot)$ hides poly-log factors.

Algorithm 8 Low-Rank approximation for distance matrix C

Input point sets $\{\mathbf{x}_i\}_{i=1}^n, \{\mathbf{y}_j\}_{j=1}^M$ in metric space \mathcal{X} and metric d
 Pick indices $i^* \in [n], j^* \in [m]$ uniformly at random
for $i = 1$ to n **do**
 Update sample probability $p_i = d(\mathbf{x}_i, \mathbf{y}_{j^*})^2 + d(\mathbf{x}_{i^*}, \mathbf{y}_{j^*})^2 + \frac{1}{m} \sum_{j=1}^m d(\mathbf{x}_{i^*}, \mathbf{y}_j)^2$
end for
 Sample $O(r/\varepsilon)$ rows $C_{i,\cdot} \sim \text{Categorical}\left(\frac{p_i}{\sum_i p_i}\right)$
 Compute U using Frieze et al. (2004)
 Compute V using Chen & Price (2017)
 return V, U

This algorithm relies on two previous works, whose main results we summarize here.

Theorem I.3. *Frieze et al. (2004)* Let $C \in \mathbb{R}^{n \times m}$. For a sample of $O(r/\varepsilon)$ rows according to a distribution $\mathbf{p} \in \Delta_n$ which satisfies $p_i \geq \Omega(1) \|C_{i,\cdot}\|_2^2 / \|C\|_F^2$ for $i \in [n]$. Then in $O(mr/\varepsilon + \text{poly}(r, 1/\varepsilon))$ time one may compute a matrix $U \in \mathbb{R}^{r \times m}$ from this sample which satisfies:

$$\|C - CU^T U\|_F^2 \leq \|C - C_k\|_F^2 + \varepsilon \|C\|_F^2$$

With probability 0.99.

Thus, to compute the first low-rank factor U , we need to ensure the p_i generated from the algorithm satisfies this requirement and offer the (short) proof below.

Proof. First, it is helpful to note $d(x, y)^2 \leq (d(x, z) + d(z, y))^2 = d(x, z)^2 + 2d(x, z)d(z, y) + d(y, z)^2 \leq 2(d(x, z)^2 + d(y, z)^2)$. Where in the last step one uses AM-GM where $\prod_i a_i^{1/n} \leq \frac{\sum_i a_i}{n}$. Rewriting the norm of row i we have:

$$\begin{aligned}
 \|C_{i,\cdot}\|_2^2 &= \sum_{j=1}^m d(\mathbf{x}_i, \mathbf{y}_j)^2 \leq 2 \sum_{j=1}^m d(\mathbf{x}_i, \mathbf{y}_{j^*})^2 + d(\mathbf{y}_{j^*}, \mathbf{y}_j)^2 \\
 &= 2md(\mathbf{x}_i, \mathbf{y}_{j^*})^2 + 2 \sum_{j=1}^m d(\mathbf{y}_{j^*}, \mathbf{y}_j)^2 \\
 &\leq 2md(\mathbf{x}_i, \mathbf{y}_{j^*})^2 + 4 \sum_{j=1}^m d(\mathbf{y}_{j^*}, \mathbf{x}_{i^*})^2 + d(\mathbf{y}_j, \mathbf{x}_{i^*})^2 \\
 &= 2md(\mathbf{x}_i, \mathbf{y}_{j^*})^2 + 4md(\mathbf{y}_{j^*}, \mathbf{x}_{i^*})^2 + 4 \sum_{j=1}^m d(\mathbf{y}_j, \mathbf{x}_{i^*})^2 \\
 &= 4m \left(\frac{1}{2} d(\mathbf{x}_i, \mathbf{y}_{j^*})^2 + d(\mathbf{y}_{j^*}, \mathbf{x}_{i^*})^2 + \frac{1}{m} \sum_{j=1}^m d(\mathbf{y}_j, \mathbf{x}_{i^*})^2 \right) \leq 4mp_i
 \end{aligned}$$

As we have the re-normalization $p_i \leftarrow \frac{p_i}{\sum_i p_i}$ before sampling, we need to consider the value of the expectation $\mathbb{E}[\sum_i p_i]$ to conclude.

$$\begin{aligned}
\mathbb{E} \left[\sum_{i=1}^n p_i \right] &= \sum_{i=1}^n \mathbb{E} \left[d(\mathbf{x}_i, \mathbf{y}_{j^*})^2 + d(\mathbf{x}_{i^*}, \mathbf{y}_{j^*})^2 + \frac{1}{m} \sum_{j=1}^m d(\mathbf{x}_{i^*}, \mathbf{y}_j)^2 \right] \\
&= \sum_{i=1}^n \mathbb{E}_{j^* \sim [m]} [d(\mathbf{x}_i, \mathbf{y}_{j^*})^2] + \mathbb{E}_{i^* \sim [n], j^* \sim [m]} [d(\mathbf{x}_{i^*}, \mathbf{y}_{j^*})^2] \\
&\quad + \frac{1}{m} \sum_{j=1}^m \mathbb{E}_{i^* \sim [n]} [d(\mathbf{x}_{i^*}, \mathbf{y}_j)^2] \\
&= \sum_{i=1}^n \frac{1}{m} \sum_{j=1}^m d(\mathbf{x}_i, \mathbf{y}_j) + n \sum_{i=1}^n \sum_{j=1}^m \frac{1}{nm} d(\mathbf{x}_i, \mathbf{y}_j) + \frac{n}{m} \sum_{j=1}^m \frac{1}{n} \sum_{i=1}^n d(\mathbf{x}_i, \mathbf{y}_j) \\
&= \frac{3}{m} \sum_{i=1}^n \sum_{j=1}^m d(\mathbf{x}_i, \mathbf{y}_j) = \frac{3}{m} \|\mathbf{C}\|_F^2
\end{aligned}$$

By an application of Markov's inequality we have that:

$$P \left[\left(\sum_{i=1}^n p_i \right)^{-1} \geq \left(\frac{3\|\mathbf{C}\|_F^2}{m\delta} \right)^{-1} \right] = P \left[\sum_{i=1}^n p_i \leq \frac{3\|\mathbf{C}\|_F^2}{m\delta} \right] \geq 1 - \frac{\mathbb{E}[\sum_{i=1}^n p_i]}{\frac{3\|\mathbf{C}\|_F^2}{m\delta}} = 1 - \delta$$

Thus with probability $1 - \delta$ we have:

$$\frac{p_i}{\sum_{i'=1}^n p_{i'}} \geq \frac{mp_i\delta}{3\|\mathbf{C}\|_F^2} = \frac{4mp_i\delta}{12\|\mathbf{C}\|_F^2} \geq \frac{\|\mathbf{C}_{i..}\|_2^2\delta}{12\|\mathbf{C}\|_F^2} = \Omega(\delta) \frac{\|\mathbf{C}_{i..}\|_2^2}{\|\mathbf{C}\|_F^2}$$

This indicates the algorithm presented has probabilities with an appropriate bound for using the algorithm of Frieze et al. (2004) to sample the $O(r/\varepsilon)$ rows of \mathbf{C} and generate a rank- r factor \mathbf{U} . \square

To conclude the result of Indyk et al. (2019) requires reference to an additional work which solves a regression problem for \mathbf{V} given \mathbf{C} and \mathbf{U} .

Theorem I.4. *Chen & Price (2017) There is a randomized algorithm \mathcal{A} , given matrices $\mathbf{C} \in \mathbb{R}^{n \times m}$, $\mathbf{U} \in \mathbb{R}^{r \times m}$ reads only $O(r/\varepsilon)$ columns of \mathbf{C} with time-complexity $O(mr) + \text{poly}(r/\varepsilon)$ and returns $\mathbf{V} \in \mathbb{R}^{n \times r}$ which satisfies*

$$\|\mathbf{C} - \mathbf{V}\mathbf{U}\|_F^2 \leq (1 + \varepsilon) \min_{\mathbf{Z} \in \mathbb{R}^{n \times r}} \|\mathbf{C} - \mathbf{Z}\mathbf{U}\|_F^2$$

with probability 0.99.

Thus, using the result of Chen and Price Chen & Price (2017), one may find a satisfying \mathbf{V} easily for a fixed \mathbf{U} , \mathbf{C} . In particular, Indyk et al. (2019) concludes by tying together the low-rank distance matrix algorithm and the guarantees of the algorithms from Frieze et al. (2004) Chen & Price (2017) as follows

$$\begin{aligned}
\|\mathbf{C} - \mathbf{V}\mathbf{U}\|_F^2 &\leq (1 + \varepsilon) \min_{\mathbf{Z}} \|\mathbf{C} - \mathbf{Z}\mathbf{U}\|_F^2 \\
&\leq (1 + \varepsilon) \|\mathbf{C} - \mathbf{C}\mathbf{U}^T\mathbf{U}\|_F^2 \\
&\leq (1 + \varepsilon) (\|\mathbf{C} - \mathbf{C}_r\|_F^2 + \varepsilon \|\mathbf{C}\|_F^2) \\
&= \|\mathbf{C} - \mathbf{C}_r\|_F^2 + \varepsilon \|\mathbf{C} - \mathbf{C}_r\|_F^2 + (1 + \varepsilon)\varepsilon \|\mathbf{C}\|_F^2 \\
&\leq \|\mathbf{C} - \mathbf{C}_r\|_F^2 + (1 + 2\varepsilon)\varepsilon \|\mathbf{C}\|_F^2 \leq \|\mathbf{C} - \mathbf{C}_r\|_F^2 + \tilde{\varepsilon} \|\mathbf{C}\|_F^2
\end{aligned}$$

Which achieves the bound up to a constant scaling of ε and shows the result of Indyk et al. (2019). We next investigate low-rank optimal transport solvers, which assume the result and algorithm of Indyk et al. (2019) to tractably scale to massive datasets.

J Connection of Optimal Transport to Projection Problems

Before discussing works which have address the problem of finding low-rank decompositions of the coupling matrix \mathbf{P} , we discuss a few relevant preliminaries. The Bregman-divergence of some function $F(\mathbf{x})$, defined by the first-order Taylor expansion of F :

$$D_F(\mathbf{x}, \mathbf{y}) = F(\mathbf{x}) - F(\mathbf{y}) + \nabla F(\mathbf{y})^\top (\mathbf{x} - \mathbf{y})$$

For the negative entropy function $-\mathbb{H}(\boldsymbol{\xi})$, this corresponds to the KL-divergence $\text{KL}(\zeta \mid \boldsymbol{\xi})$. We introduce the iterative-Bregman projection algorithm in the context of the KL-divergence owing to the direct connection with entropically-regularized optimal transport.

Definition J.1. Iterative Bregman Projections Bregman (1967)

Suppose $\mathcal{C} = \bigcap_{l=1}^L \mathcal{C}_l$ is an intersection of closed convex sets $\{\mathcal{C}_l\}_{l=1}^L$. For $n > L$, let the indexing be L -periodic as $\mathcal{C}_n := \mathcal{C}_{n \bmod L}$. Suppose we want to find a minimizer ζ of the KL-divergence with some positive vector $\boldsymbol{\xi} \in \mathbb{R}_+^{n \times m}$ such that $\zeta \in \mathcal{C}$. This is to say, we hope to solve the problem of minimizing a distance subject to the condition that one remains in this intersection of convex sets:

$$\min_{\zeta \in \mathcal{C}} \text{KL}(\zeta \parallel \boldsymbol{\xi})$$

Where the projection of $\boldsymbol{\xi}$ onto the set \mathcal{C} is denoted by

$$\zeta^* = \arg \min_{\zeta \in \mathcal{C}} \text{KL}(\zeta \parallel \boldsymbol{\xi}) := \mathcal{P}_{\mathcal{C}}^{KL}(\boldsymbol{\xi})$$

Supposing that each \mathcal{C}_l forms a quotient space $\mathcal{C}_l = V/U$ for a subspace $U \subset V$, defined as $V/U = \{\mathbf{v} + U \mid \mathbf{v} \in V\}^2$, the iterative Bregman projection algorithm alternates projections onto each set \mathcal{C}_n as

$$\zeta^{(n)} \leftarrow \mathcal{P}_{\mathcal{C}_n}^{KL}(\zeta^{(n-1)}) \quad (43)$$

starting from $\zeta^{(0)} = \boldsymbol{\xi}$.

One may show Bregman (1967) the convergence of Bregman projections to the unique minimizer in \mathcal{C} , ζ^* , where we have the guarantee that $\zeta^{(n)} \rightarrow \zeta^*$ as $n \rightarrow \infty$. These iterative projections only have convergence guarantees when the constraint sets are quotient spaces—this is clearly not the case for the constraints of the optimal transport LP, and a few more notions are required.

Definition J.2. Dykstra’s Algorithm Dykstra (1983)

Given a point $\mathbf{x}_0 \in E$ for E a Euclidean space³, to find the unique point in $\mathcal{C} = \bigcap_{l=1}^L \mathcal{C}_l$ for closed, convex sets \mathcal{C}_l that minimize the distance to \mathbf{x}_0 as

$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathcal{C}} \|\mathbf{x} - \mathbf{x}_0\|_2$$

One may initialize the residuals $\mathbf{q}_{-(L-1)} = \dots = \mathbf{q}_0 = 0$ and apply the algorithm:

$$\begin{aligned} \mathbf{x}_n &= \mathcal{P}_{\mathcal{C}_n}(\mathbf{x}_{n-1} + \mathbf{q}_{n-1}) \\ \mathbf{q}_n &= (\mathbf{x}_{n-1} - \mathbf{x}_n) + \mathbf{q}_{n-L} \end{aligned}$$

Where $\mathcal{C}_n := \mathcal{C}_{n \bmod L}$ as before and $\mathcal{P}_{\mathcal{C}_n}$ denotes the projection operator onto the convex set \mathcal{C}_n .

One can note that Dykstra’s algorithm for projections onto intersections of convex sets no longer relies on the assumption that the set is a quotient space, and applies in the case that it is merely closed under convex-combinations. To generalize Dykstra’s for more general functions than the ℓ_2 -norm, one may define the projection with respect to the Bregman divergence of a cost function F and define the projection by the minimization: $\mathcal{P}_{\mathcal{C}_n} := \arg \min_{\mathbf{x}} D_F(\mathbf{x}, \mathbf{y})$. It was proven in Bauschke & Lewis (2000) that the generalized form of Dykstra’s iterations are given as:

$$\begin{aligned} \mathbf{x}_n &= \mathcal{P}_{\mathcal{C}_n}(\nabla F^*(\nabla F(\mathbf{x}_{n-1}) + \mathbf{q}_{n-1})) \\ \mathbf{q}_n &= (\nabla F(\mathbf{x}_{n-1}) - \nabla F(\mathbf{x}_n)) + \mathbf{q}_{n-L} \end{aligned}$$

For F^* denoting the Fenchel-conjugate of F , which we define later in connection to the low-rank dual problem. Notably, Bauschke & Lewis (2000) also provided guarantees of convergence to the optimal solution which extend to the case that F is the negative entropy and D_F the KL-divergence. These constitute Dykstra’s algorithm with cyclic Bregman projections, and project a point to the closest point in the intersection of convex sets $\mathcal{C} = \bigcap_{l=1}^L \mathcal{C}_l$ for an arbitrary cost function F and its associated Bregman-divergence D_F .

²This is to say, $\mathbf{x}, \mathbf{y} \in \mathcal{C}_l \implies c_1 \mathbf{x} + c_2 \mathbf{y} \in V/U$ for any $c_1, c_2 \in \mathbb{R}$.

³e.g. $\mathbb{R}, \mathbb{R}^d, \mathbb{R}^{n \times m}$, etc.

J.1 Dykstra's algorithm with cyclic Bregman projections

As such, we can see the updates for F being the negative entropy and D_F the KL-divergence. Without proof, the conjugate of the negative entropy is simply given as $F^*(\zeta) = \exp\{\zeta - \mathbf{1}\}$. Thus, for the minimization problem:

$$\zeta^* = \arg \min_{\zeta \in \mathcal{C}} \text{KL}(\zeta \|\xi) := \mathcal{P}_{\mathcal{C}}^{\text{KL}}(\xi)$$

Letting $\log \mathbf{q}_{-(L-1)} = \dots = \log \mathbf{q}_0 = 0$, one may combine the two algorithms above to solve this minimization using generalized Dykstra's iterations. In particular, we have:

$$\begin{aligned} \zeta^{(n)} &= \mathcal{P}_{\mathcal{C}_n} \left(\nabla F^* (\nabla F(\zeta^{(n-1)}) + \log \mathbf{q}_{n-1}) \right) \\ &= \mathcal{P}_{\mathcal{C}_n} \left(\exp \left(\nabla F(\zeta^{(n-1)}) + \log \mathbf{q}_{n-1} - \mathbf{1} \right) \right) \\ &= \mathcal{P}_{\mathcal{C}_n} \left(\exp \left(\log \zeta^{(n-1)} + \mathbf{1} + \log \mathbf{q}_{n-1} - \mathbf{1} \right) \right) \\ &= \mathcal{P}_{\mathcal{C}_n} \left(\zeta^{(n-1)} \odot \mathbf{q}_{n-1} \right) \end{aligned}$$

And:

$$\begin{aligned} \log \mathbf{q}_n &= (\nabla F(\zeta^{(n-1)}) - \nabla F(\zeta^{(n)})) + \log \mathbf{q}_{n-L} \\ &= \left(\log \zeta^{(n-1)} + \mathbf{1} - (\log \zeta^{(n)} + \mathbf{1}) + \log \mathbf{q}_{n-L} \right) \\ &= \log \left(\mathbf{q}_{n-L} \odot \frac{\zeta^{(n-1)}}{\zeta^{(n)}} \right) \end{aligned}$$

So that:

$$\zeta^{(n)} \leftarrow \mathcal{P}_{\mathcal{C}_n} \left(\zeta^{(n-1)} \odot \mathbf{q}_{n-1} \right) \quad (44)$$

$$\mathbf{q}_n \leftarrow \mathbf{q}_{n-L} \odot \frac{\zeta^{(n-1)}}{\zeta^{(n)}} \quad (45)$$

Where division is interpreted to be element-wise, \odot refers to the Hadamard product, and the logarithm is applied elementwise.

J.2 Connection to Sinkhorn distances

Interestingly, the Sinkhorn algorithm described in Algorithm 5 can be alternatively derived in the context of Bregman iterations Benamou et al. (2015) as a minimization of the form:

$$W_\epsilon(\mu, \nu) = \epsilon \min_{\mathbf{P} \in \Pi(\mu, \nu)} \text{KL}(\mathbf{P} \|\xi)$$

Where ξ is the kernel $\xi = e^{-C/\epsilon}$ and $\Pi(\mu, \nu) = \mathcal{C}_1 \cap \mathcal{C}_2$ for the convex constraint sets $\mathcal{C}_1 = \{\mathbf{P} : \mathbf{P}\mathbf{1}_m = \mathbf{a}\}$ and $\mathcal{C}_2 = \{\mathbf{P} : \mathbf{P}^T\mathbf{1}_n = \mathbf{b}\}$. To cast this into the Bregman-projection framework, one alternates between the two updates:

$$\begin{aligned} \mathbf{P}^{(l)} &= \mathcal{P}_{\mathcal{C}_1}(\mathbf{P}^{(l-1)}) \\ \mathbf{P}^{(l+1)} &= \mathcal{P}_{\mathcal{C}_2}(\mathbf{P}^{(l)}) \end{aligned}$$

Where for the first projection one has the following first-order KKT condition:

$$\begin{aligned} \nabla \left(\epsilon \text{KL}(\mathbf{P} \|\mathbf{P}^{(l-1)}) + \lambda^T (\mathbf{P}\mathbf{1}_m - \mathbf{a}) \right) &= \epsilon \log \left(\frac{\mathbf{P}}{\mathbf{P}^{(l-1)}} \right) + \lambda \mathbf{1}_m^T = \mathbf{0} \\ \implies \mathbf{P} &= \text{diag}(e^{-\lambda/\epsilon}) \mathbf{P}^{(l-1)} \end{aligned}$$

With the constraint of \mathcal{C}_1 that $\mathbf{P}\mathbf{1}_m = \mathbf{a}$, this implies $\text{diag}(\mathbf{a}/\mathbf{P}^{(l-1)}\mathbf{1}_m) = \text{diag}(e^{-\lambda/\epsilon})$ and recovers the first update of Sinkhorn $\mathbf{P}^{(l)} \leftarrow \text{diag}(\mathbf{a}/\mathbf{P}^{(l-1)}\mathbf{1}_m) \mathbf{P}^{(l-1)}$. An analogous argument gives the second, where all iterates satisfy $\mathbf{P}^{(l)} = \text{diag}(\mathbf{u}^{(l)}) e^{-C/\epsilon} \text{diag}(\mathbf{v}^{(l)})$ for \mathbf{u}, \mathbf{v} as defined in Algorithm 5.

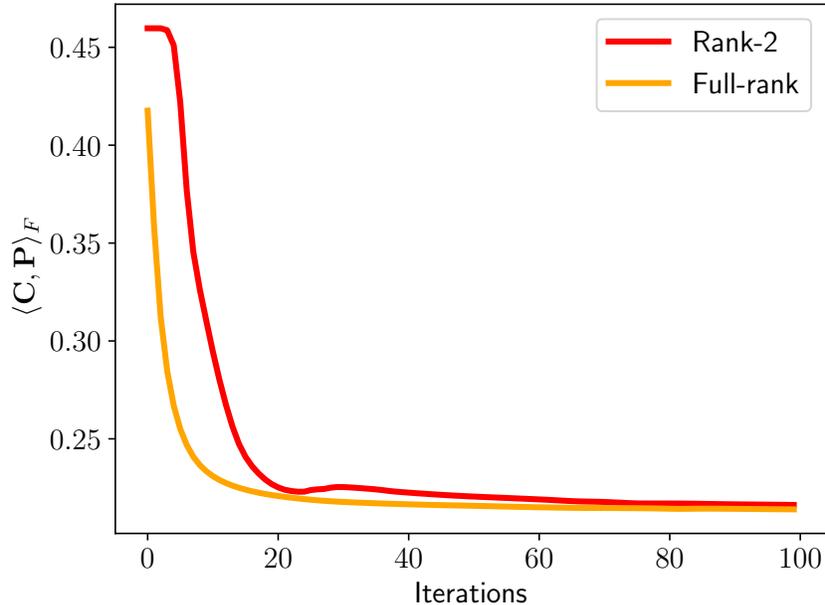


Figure 5: Transport cost $\langle C, P \rangle_F$ against number of iterations for FRLC with rank 200 on the synthetic dataset of two moons and eight Gaussians. Smooth convergence is observed for both rank-2 and full-rank random initialization.

K Additional Simulations

We tested on two additional synthetic datasets, both used as benchmarking datasets in Scetbon et al. (2021). We follow exactly the parameters provided in Scetbon et al. (2021) to simulate these datasets. For the first one, we simulated $n = m = 10,000$ points from two Gaussian mixtures in 2D (Fig. 6). The first Gaussian mixture is a mixture of three Gaussian distributions with means $(0, 0)$, $(0, 1)$, $(1, 1)$ respectively. The mixture proportion is $\frac{1}{3}$ and the covariance is $0.05 \times \text{identity}$ for each Gaussian. The second Gaussian mixture is a mixture of two Gaussian distributions with means $(0.5, 0.5)$, $(-0.5, 0.5)$ respectively. The mixture proportions is $\frac{1}{2}$ and the covariance is $0.05 \times \text{identity}$ for each Gaussian.

For the second dataset, we simulated $n = m = 5,000$ points from two Gaussian mixtures in 10D. The first Gaussian mixture is a mixture of three Gaussian distributions with means $(0, 0, 0, \dots, 0)$, $(0, 1, 0, \dots, 0)$, $(1, 1, 0, \dots, 0)$ respectively. The mixture proportions is $\frac{1}{3}$ and the covariance is $0.05 \times \text{identity}$ for each Gaussian. The second Gaussian mixture is a mixture of two Gaussian distributions with means $(0.5, 0.5, 0, \dots, 0)$, $(-0.5, 0.5, 0, \dots, 0)$ respectively. The mixture proportions is $\frac{1}{2}$ and the covariance is $0.05 \times \text{identity}$ for each Gaussian.

For each dataset, we repeat the same procedure as in § 4.1, running FRLC and LOT with Euclidean distance as cost to find a low-rank coupling matrix between the two Gaussian mixtures with rank between 50 and 200. We observe the same pattern as Fig. 2b. FRLC obtains lower transport cost with increasing rank, and achieves lower cost for each rank than LOT under all initializations (Fig. 23c, Fig. 7).

L Graph Partitioning

L.1 Evaluation on a Graph Partitioning Task

We next evaluate FRLC on an unsupervised graph partitioning (node clustering) problem described by Chowdhury & Needham (2021). Specifically, given a graph $G = (V, E)$ of n nodes, we represent the graph as $G = (\mathbf{A}, \mathbf{h})$, where $\mathbf{A} \in \mathbb{R}^{n \times n}$ encodes the intra-graph node relationship (e.g. adjacency matrix) and $\mathbf{h} \in \Delta_n$ is a uniform measure. We cluster the nodes of G by estimating a GW coupling

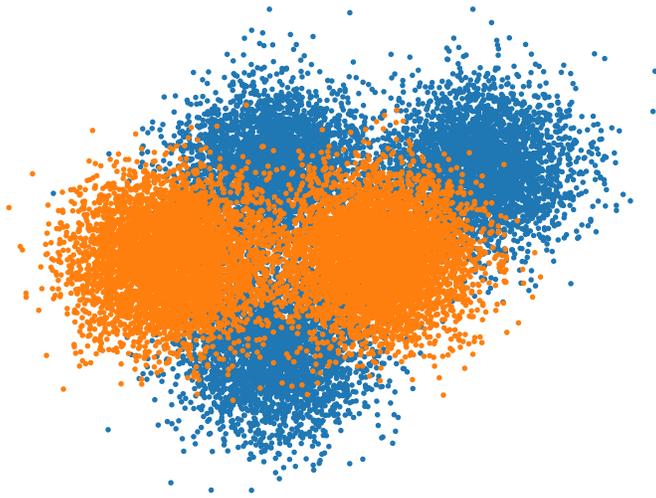


Figure 6: Plot of the two simulated mixtures of Gaussians in 2D, following the same parameters as Scetbon et al. (2021).

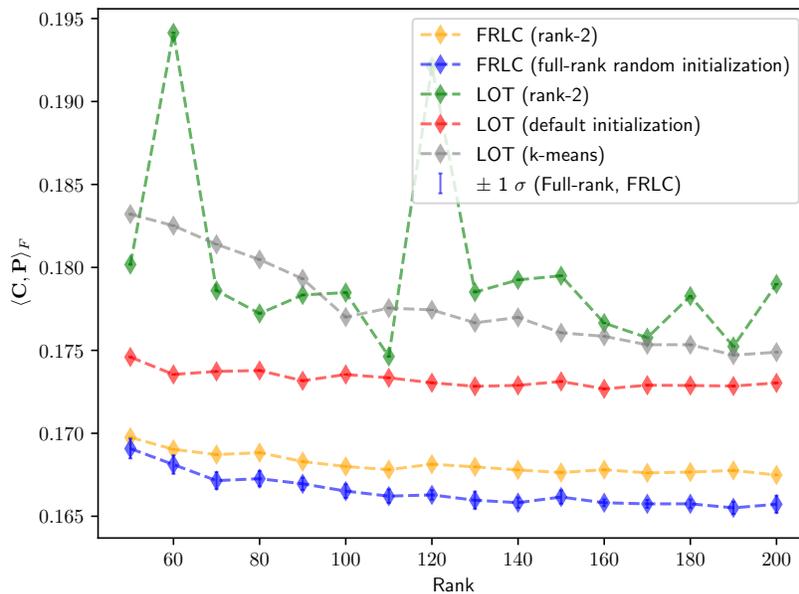


Figure 7: Transport cost $\langle C, P \rangle_F$ achieved by LOT Scetbon et al. (2021) and FRLC across different ranks and different initializations on the Wasserstein problem on the synthetic dataset of two mixtures of Gaussians in 2D.

Method	Dataset	Time in seconds (CPU)	Cost Value $\langle C, P \rangle_F$	$\ P\mathbf{1}_n - \mathbf{a}\ _2$	$\ P^T\mathbf{1}_m - \mathbf{b}\ _2$
FRLC	2-moons,	0.379	0.207	1.56E-5	6.1E-18
LOT	8-Gaussians	0.751	0.210	1.90E-5	1.89E-5
FRLC	Gaussian mixture	0.354	0.178	1.05E-5	7.0E-18
LOT	(2D)	0.735	0.181	1.65E-5	1.70E-5
FRLC	Gaussian mixture	0.323	0.294	2.48E-6	8.9E-18
LOT	(10D)	0.677	0.307	1.39E-5	1.46E-5

Table 2: Runtime of FRLC and LOT on the synthetic datasets of 1000 samples, as well as cost value $\langle C, P \rangle_F$ and marginal tightness for context. This was done with the FRLC setting `max_inneriters_balanced=1000`, `max_inneriters_relaxed=50`, `min_iter=7` and rank $r = 100$. This time excludes the extra time incurred for `ott-jax` problem setup and includes the setup time for FRLC.

between G and a smaller graph $\bar{G} = (\bar{B}, \bar{h})$, where $\bar{B} \in \mathbb{R}^{m \times m}$ represents the relationship between each of the m clusters and $\bar{h} \in \Delta_m$ is the proportion of nodes of G in each cluster. Without a priori knowledge, \bar{h} is set to be uniform and \bar{B} is set as the identity matrix. Vincent-Cuaz et al. (2022) notes that instead of solving a balanced GW problem, semi-relaxed GW with the right marginal \bar{h} relaxed learns \bar{h} from data and leads to more accurate node clustering.

Dataset		GWL	FRLC SR-GW	SpecGWL	SpecFRLC SR-GW
Wikipedia (1998 nodes, 2700 edges)	sym, raw	0.314	0.387	0.372	0.444
	sym, noisy	0.250	0.361	0.293	0.400
	asym, raw	0.263	0.276	0.194	0.304
	asym, noisy	0.208	0.201	0.141	0.177
EU-email (1005 nodes, 25571 edges)	sym, raw	0.434	0.464	0.009	0.040
	sym, noisy	0.392	0.422	0.009	0.014
	asym, raw	0.388	0.398	0.012	0.028
	asym, noisy	0.385	0.348	0.008	0.012
Amazon (1501 nodes, 4626 edges)	raw	0.322	0.338	0.505	0.479
	noisy	0.274	0.257	0.438	0.453
Village (1991 nodes, 8423 edges)	raw	0.531	0.710	0.553	0.579
	noisy	0.413	0.536	0.397	0.829

Table 3: Performance (measured using Adjusted Mutual Information (AMI)) in graph partitioning for *full-rank* OT algorithms GWL, SpecGWL and full-rank semi-relaxed FRLC. The top performing method for each dataset is highlighted in bold.

We benchmark FRLC for semi-relaxed GW on four real-world graphs: a Wikipedia hyperlink network with 15 webpage categories Yang & Leskovec (2012); an email interaction network within a European institute with 42 departments Yin et al. (2017); an Amazon product network with 12 product categories Yang & Leskovec (2012); and a network of interactions between 12 Indian villages Banerjee et al. (2013). We also test on the symmetric and noisy versions of each graph provided by Chowdhury & Needham (2021). We compare with two OT-based methods: (1) GWL Xu et al. (2019), which solves a balanced GW problem between G and \bar{G} with the adjacency matrix of G as the intra-domain cost matrix A ; (2) SpecGWL Chowdhury & Needham (2021) which uses the heat kernel on the graph Laplacian as A . We similarly run our FRLC algorithm using with both the adjacency matrix (denoted FRLC-SR-GW) and heat kernel (denoted SpecFRLC-SR-GW). Since the number of clusters in each dataset is not large, we compute the full-rank coupling matrix in each case.

Overall, FRLC achieves the best clustering performance on 9 out of 12 datasets (Table 3). When using the adjacency matrix, our semi-relaxed algorithm achieves better clustering result than GWL on 9 out of 12 datasets. When using the heat kernel, our semi-relaxed algorithm achieves better

clustering result than SpecGWL on 11 out of 12 datasets. These results show the importance of semi-relaxed OT on real-world problems, as well as the accuracy of FRLC.

L.2 Problem Statement

As discussed in Vincent-Cuaz et al. (2022); Chowdhury & Needham (2021), it is possible to achieve unsupervised graph partitioning (node clustering) through Gromov-Wasserstein (GW) OT. Given a graph (V, E) of n nodes, we encode it as $G = (\mathbf{A}, \mathbf{h})$, where $\mathbf{A} \in \mathbb{R}^{n \times n}$ encodes the intra-graph node relationship (e.g. adjacency matrix, Laplacian) and $\mathbf{h} \in \Delta n$ is a uniform measure over the nodes. If we want to cluster the nodes of G into m clusters, we can define a new graph $\bar{G} = (\bar{\mathbf{B}}, \bar{\mathbf{h}})$, where $\bar{\mathbf{B}} \in \mathbb{R}^{m \times m}$ is a diagonal matrix representing the cluster’s connections and $\bar{\mathbf{h}}$ is a distribution over clusters estimating the proportion of nodes in G in each cluster. Since we don’t know the density of clusters a priori, we can set $\bar{\mathbf{h}}$ to be uniform. Usually $\bar{\mathbf{B}}$ is set as the identity matrix.

To cluster the nodes in G , we can solve a GW problem matching nodes in G with nodes in \bar{G} , with intra-domain cost matrices \mathbf{A} and $\bar{\mathbf{B}}$:

$$\begin{aligned} \min \quad & \sum_{ij'kl'} (\mathbf{A}_{ik} - \bar{\mathbf{B}}_{j'l'})^2 \mathbf{P}_{ij'} \mathbf{P}_{kl'} \\ \text{s.t.} \quad & \mathbf{P} \in \Pi_{\mathbf{h}, \bar{\mathbf{h}}} \end{aligned} \tag{46}$$

The cluster assignment of each node in G can then be recovered from \mathbf{P} by finding the node in \bar{G} mapped to it with the maximum weight.

However, since the proportion of nodes in G in each cluster is not known a priori, solving a balanced GW problem fixing the marginal of \mathbf{P} on \bar{G} to be a uniform $\bar{\mathbf{h}}$ significantly constrains the expressivity of the algorithm. Therefore, as proposed in Vincent-Cuaz et al. (2022), we can instead solve a semi-relaxed GW problem with the right marginal relaxed, fixing the marginal on G to be \mathbf{h} but allowing the marginal on \bar{G} to deviate from $\bar{\mathbf{h}}$:

$$\begin{aligned} \min \quad & \sum_{ij'kl'} (\mathbf{A}_{ik} - \bar{\mathbf{B}}_{j'l'})^2 \mathbf{P}_{ij'} \mathbf{P}_{kl'} + \tau \text{KL}(\mathbf{P}^T \mathbf{1}_n \mid \bar{\mathbf{h}}) \\ \text{s.t.} \quad & \mathbf{P} \in \Pi_{\mathbf{h}, \cdot} \end{aligned} \tag{47}$$

The learned $\bar{\mathbf{h}}$ from semi-relaxed GW estimates the posterior proportion of nodes in G in each of the m clusters.

L.3 Dataset and Preprocessing

We run our semi-relaxed FRLC algorithm on four real-world graph datasets: a Wikipedia hyperlink network with 1998 nodes and 15 clusters Yang & Leskovec (2012), a directed graph of email interactions in a European research institute with 1005 nodes and 42 clusters Yin et al. (2017), an Amazon product network with 1501 nodes and 12 clusters Yang & Leskovec (2012), and a network of interactions between Indian villages with 1991 nodes and 12 clusters Banerjee et al. (2013). The Wikipedia and EU-email graphs are directed, so we also use undirected versions of them. We also use noisy version of each graph by adding up to 10% additional edges following Chowdhury & Needham (2021), leading to a total of 12 different graphs to cluster.

L.4 Experiment Settings

We compare our algorithm with two baseline methods, GWL and SpecGWL. GWL Xu et al. (2019) solves the entropy-regularized version of the balanced GW problem of (46) with the adjacency matrix of G as \mathbf{A} . We set \mathbf{h} such that the density of each node is proportional to its degree. We set $\bar{\mathbf{h}}$ to be a distribution estimated by sorting the weights of \mathbf{h} and sampling m values via linear interpolation, following Chowdhury & Needham (2021). We set $\bar{\mathbf{B}} = \text{diag}(\bar{\mathbf{h}})$. We set the entropy regularization

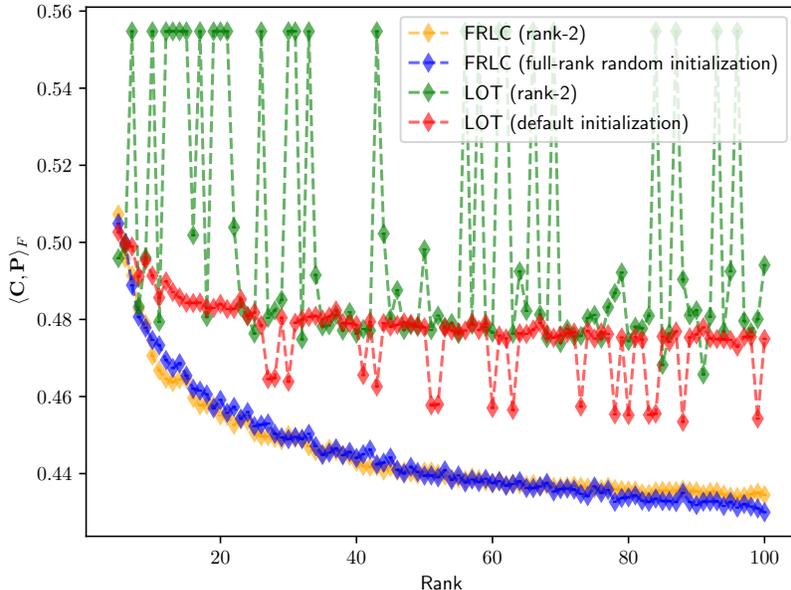


Figure 8: FRLC achieves lower primal cost $\langle C, P \rangle_F$ for $P \in \Pi_{a,b}$ than Scetbon et al. (2021) on a spatial-transcriptomics dataset of mouse embryonic development Chen et al. (2022). FRLC demonstrates a more robust trend of improved performance with higher rank.

$\eta = 10^{-6}$. SpecGWL Chowdhury & Needham (2021) solves the same problem as GWL, but instead of using the adjacency matrix as A , it uses the heat kernel on the normalized graph Laplacians Chung (2005). We set the heat parameter $t = 10$.

For our method, we solve the semi-relaxed GW problem of (47) with full rank solution. We use both adjacency matrix and heat kernel as C and label the result of the two representations as FRLC-SR-GW and SpecFRLC-SR-GW. We set $\tau = 0.01$ as to minimize the conformation to the right marginal. Since our method depends on random initialization, we run our method 10 times on each dataset and report the mean performance. We evaluate the resulting clusterings of all methods by computing the Adjusted Mutual Information (AMI) between the computed clustering and the ground truth clustering. This experiment, and the experiments on mouse embryo spatial transcriptomics, were conducted on cluster GPUs.

M Spatial Transcriptomics Alignment

M.1 Problem Statement

In this problem, we use FRLC to find a low-rank alignment matrix between cells from two spatial transcriptomics (ST) Ståhl et al. (2016) slices, collected at two timepoints, then use the computed alignment matrix for two downstream prediction tasks. An ST experiment on a 2D tissue slice yields a pair (X, Z) . $X \in \mathbb{R}^{n \times p}$ is the gene expression matrix, where n is the number of cells on the slice and p is the number of genes measured. $X_{ij} \in \mathbb{R}$ is the gene expression level of gene j in cell i , where a higher number indicates stronger expression. $Z \in \mathbb{R}^{n \times 2}$ is the spatial coordinate matrix, where each row i stores the x-y coordinate of cell i on the slice. Therefore, each cell on the slice has a gene expression vector of length p , which encodes the feature of the cell, as well as a coordinate vector of length two, which encodes the geometrical information of the cell on the slice.

Our input data is a pair of ST slices $\mathcal{S}_0 = (X^0, Z^0), \mathcal{S}_1 = (X^1, Z^1)$, with n and m cells, of the same tissue region. We assume \mathcal{S}_0 is collected at timepoint $t = 0$ and \mathcal{S}_1 is collected at timepoint $t = 1$, hence the transition from \mathcal{S}_0 to \mathcal{S}_1 reflects the biological development of the tissue during the time period. We would like to find the ancestor-descendant relationship between cells from \mathcal{S}_0 and \mathcal{S}_1 by computing an optimal transport coupling matrix between cells from \mathcal{S}_0 and cells from \mathcal{S}_1 . The state-of-the-art spatial transcriptomics alignment method moscot Klein et al. (2023) claims that

unbalanced OT is the most appropriate setup for this problem. Specifically, given discrete uniform measure \mathbf{a} and \mathbf{b} over cells from \mathcal{S}_0 and \mathcal{S}_1 , we solve the unbalanced Wasserstein problem

$$\min_{\mathbf{P} \in \mathbb{R}_+^{n \times m}} \langle \mathbf{C}, \mathbf{P} \rangle + \tau \text{KL}(\mathbf{P} \mathbf{1}_m \| \mathbf{a}) + \tau \text{KL}(\mathbf{P}^T \mathbf{1}_n \| \mathbf{b}) \quad (48)$$

$\mathbf{C} \in \mathbb{R}_+^{n \times m}$ has entries $C_{ij} = c(X_i^0, X_j^1)$, where c is the Euclidean distance between the features of cell i from \mathcal{S}_0 and cell j from \mathcal{S}_1 .

The above formulation only considers the feature information of the two slices, but not the geometrical information. Therefore, we can also solve the unbalanced GW problem

$$\min_{\mathbf{P} \in \mathbb{R}_+^{n \times m}} \sum_{i,j',k,l'} (\mathbf{A}_{ik} - \mathbf{B}_{j'l'})^2 P_{ij'} P_{kl'} + \tau \text{KL}(\mathbf{P} \mathbf{1}_m \| \mathbf{a}) + \tau \text{KL}(\mathbf{P}^T \mathbf{1}_n \| \mathbf{b}) \quad (49)$$

where $\mathbf{A} \in \mathbb{R}^{n \times n}$ is the Euclidean distance matrix between the spatial location of cells within \mathcal{S}_0 , and $\mathbf{B} \in \mathbb{R}^{m \times m}$ is the Euclidean distance matrix between the spatial location of cells within \mathcal{S}_1 . Similarly, we can combine the above two formulations to solve the unbalanced FGW problem.

Notice the inherent asymmetry stemming from the temporal nature of the problem: all cells on \mathcal{S}_1 should have an ancestor from \mathcal{S}_1 , but not all cells from \mathcal{S}_1 need to have a descendent in \mathcal{S}_0 because of cell death. Therefore, the most natural OT task for this problem is semi-relaxed OT with the left marginal (the marginal on the first/ancestor slice) relaxed. Specifically, we can also solve the semi-relaxed Wasserstein problem

$$\min_{\mathbf{P} \in \Pi_{\cdot, \mathbf{b}}} \langle \mathbf{C}, \mathbf{P} \rangle + \tau \text{KL}(\mathbf{P} \mathbf{1}_m \| \mathbf{a}) \quad (50)$$

as well as semi-relaxed GW and FGW problem.

Gene expression prediction task Given the alignment matrix \mathbf{P} linking cells from \mathcal{S}_0 to \mathcal{S}_1 , we can predict properties of cells in \mathcal{S}_1 from properties of cells in \mathcal{S}_0 . Let the expression of a gene j in \mathcal{S}_0 be a vector $\mathbf{f}_j \in \mathbb{R}^n$, such that f_{ji} is the expression level of gene j in cell i , we can predict the expression of gene j in \mathcal{S}_1 as $\tilde{\mathbf{f}}_j = m \times \mathbf{P}^T \times \mathbf{f}_j \in \mathbb{R}^m$. The accuracy of the prediction can be measured by the Spearman correlation between the predicted expression and the ground truth expression $\bar{\mathbf{f}}_j$ of gene j in \mathcal{S}_1 : $\rho(\tilde{\mathbf{f}}_j, \bar{\mathbf{f}}_j)$. In this work, we test the prediction accuracy on 10 test genes: *Tubb2b*, *Pantr1*, *Actc1*, *Tnni1*, *Afp*, *Hbb-bh1*, *Fez1*, *Crabp1*, *Crabp2*, *Col3a1*, which are markers genes for various cell types in mouse embryo.

Cell type prediction task We can also use the cell type labels of cells in \mathcal{S}_0 to predict the cell type labels of cells in \mathcal{S}_1 . Specifically, for each cell j in \mathcal{S}_1 , we can assign it the type of the cell $\text{argmax}_i P_{ij}$ in \mathcal{S}_0 . We can measure the accuracy of the cell type prediction by computing the Adjusted Rand Index (ARI) and Adjusted Mutual Information (AMI) between the predicted clustering of cells in \mathcal{S}_1 and the ground truth clustering.

M.2 Dataset and Preprocessing

We use the large-scale real-world dataset of mouse embryo development Chen et al. (2022), consisting of eight timepoints of ST slices during the whole process of mouse embryo development. In this work, we align the pair of adjacent timepoints of E11.5 and E12.5 embryos, consisting of 30,124 cells and 51,365 cells, respectively. We preprocess the dataset using the standard SCANPY Wolf et al. (2018) pipeline. We first filter the two slices to have the same set of genes, resulting in 26,436 genes for all cells from both slices. We then log-normalize the gene expression of all cells from the two slices, and apply Principle Component Analysis (PCA) to reduce the dimensionality of gene expressions to 30. We take the Euclidean distance between the gene expression in the PCA space as the cost matrix \mathbf{C} . We take the Euclidean distance between the 2D coordinate of each cell within each slice as the intra-domain cost matrices \mathbf{A} and \mathbf{B} .

Fig. 9 visualizes the two slices in this dataset, with each cell annotated with a cell type from the original publication.

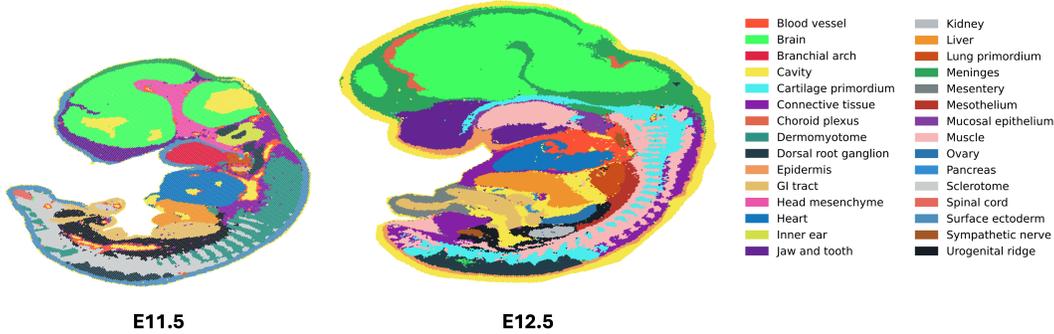


Figure 9: Visualization of the E11.5 and E12.5 mouse embryos, with each cell colored by the cell type annotated by Chen et al. (2022).

M.3 Experiment Settings

We compare with the unbalanced low-rank optimal transport algorithm of Scetbon et al. (2023), the backbone of the spatial transcriptomics alignment method moscot Klein et al. (2023), which was shown by Scetbon et al. (2023); Klein et al. (2023) to achieve state-of-the-art performance on spatial transcriptomics alignment. We use Scetbon et al. (2023) to solve three unbalanced problems for spatial transcriptomics alignment: the unbalanced Wasserstein problem of (48) (result denoted as LOT-U-W), the unbalanced GW problem of (49) (LOT-U-GW), and the unbalanced FGW problem with a convex combination of the previous two costs (LOT-U-FGW). We use FRLC to solve the same three unbalanced problems (results denoted as FRLC-U-W, FRLC-U-GW, FRLC-U-FGW). Since the two slices contain $> 30,000$ cells and $> 50,000$ cells respectively, a full-rank solution is not feasible, hence we solve for low-rank solutions with the rank validated. We also solve semi-relaxed versions of Wasserstein (result denoted as FRLC-SR-W), GW (FRLC-SR-GW), FGW (FRLC-SR-FGW) problems using our FRLC algorithm, as well as using a particular setting of LOT-U that is equivalent to semi-relaxed solver (results denoted as LOT-SR-W, LOT-SR-GW, LOT-SR-FGW).

We perform extensive grid search to find the best hyperparameter combinations for each method and each problem. The grid of hyperparameters searched for each method is reported in Table. 6. The best performing hyperparameter combination for each method is reported in Table. 7 along with the performance on the validation genes. We pick the best hyperparameters using the Spearman correlation on the gene expression prediction task for 10 validation genes: *Ckb*, *Fabp7*, *My14*, *Tnnt2*, *Apoa2*, *Hba-x*, *Tubb3*, *Epha7*, *Ldha*, *Col11a2*, which are marker genes of various cell types as well. We report the performance of the alignment computed by each method using the Spearman correlation on the gene expression prediction task for 10 test genes, as well as the ARI and AMI on the cell type prediction task. Fig. 10 visualizes the ground truth cell type classification versus the classification predicted by our method FRLC-SR-W.

M.4 Runtime

We report the runtime and OT cost of FRLC and LOT Scetbon et al. (2021) on this dataset (mouse embryo E11.5-12.5) as well as two other datasets (E9.5-10.5, E10.5-11.5) from Chen et al. (2022) in Table 4. For all three datasets, FRLC achieves a better OT cost in a shorter time.

Dataset	LOT (seconds)	FRLC (seconds)	LOT (OT Cost)	FRLC (OT Cost)
Mouse embryo (E9.5–10.5)	2.545	1.112	0.440	0.385
Mouse embryo (E10.5–11.5)	4.209	1.190	0.371	0.344
Mouse embryo (E11.5–12.5)	8.667	1.889	0.478	0.439

Table 4: Comparison of methods on Stereo-Seq mouse embryo spatial transcriptomics datasets using GPU and default settings: $\text{min_iter}=10$, $\text{max_iter}=100$, rank $r = 50$.

N n^{th} -roots of unity

N.1 Problem Statement

To test if FRLC can effectively coarse-grain transport between two datasets with obvious cluster structure, we generate a pair of two-dimensional datasets $\mathbf{Z}^{(1)}$ and $\mathbf{Z}^{(2)}$, with ten and five clusters respectively. We run FRLC with $\mathbf{T} \in \mathbb{R}_+^{10 \times 5}$ to see if the barycentric projections for the LC factorization (Definition 4.1) can recover the cluster structure. These projections induce 10 and 5 barycenters for the first and second dataset, and are defined by

$$\mathbf{Y}^{(1)} := \text{diag}(1/g_Q)\mathbf{Q}^T\mathbf{Z}^{(1)}, \quad \mathbf{Y}^{(2)} := \text{diag}(1/g_R)\mathbf{R}^T\mathbf{Z}^{(2)}.$$

We examine, visually, whether these barycenters are good representatives of the clusters in each dataset, and whether the latent coupling depicts a reasonable transfer of mass between barycenters. We also run FRLC with $\mathbf{T} \in \mathbb{R}_+^{10 \times 10}$ and plot the barycenters from the resulting factorization for comparison. As discussed in § 4.2, the barycentric projections defined above, and in Definition 4.1 can be applied to factored couplings Forrow et al. (2019); Scetbon et al. (2021), yielding projections of the form:

$$\mathbf{Y}^{(1)} := \text{diag}(1/g)\mathbf{Q}^T\mathbf{Z}^{(1)}, \quad \mathbf{Y}^{(2)} := \text{diag}(1/g)\mathbf{R}^T\mathbf{Z}^{(2)}.$$

Thus, we also ran the method of Scetbon et al. (2021) on this data, called LOT throughout the experiment, with $\mathbf{g} \in (\mathbb{R}_+^*)^5$ and $\mathbf{g} \in (\mathbb{R}_+^*)^{10}$, plotting its barycenters in each case along with the diagonal latent coupling $\text{diag}(\mathbf{g})$.

N.2 Dataset and Preprocessing

We instantiate a pair of two-dimensional datasets $\mathbf{Z}^{(1)}$ and $\mathbf{Z}^{(2)}$ as follows. Let \mathcal{U}_n denote the n^{th} -roots of unity:

$$\mathcal{U}_n := \{e^{2\pi ik/n} : k = 0, \dots, n\}.$$

These complex numbers can be equivalently expressed as ordered pairs on the unit circle:

$$\mathbf{c}_k = \begin{pmatrix} \text{Re}(e^{2\pi ik/n}) \\ \text{Im}(e^{2\pi ik/n}) \end{pmatrix} = \begin{pmatrix} \cos(2\pi k/n) \\ \sin(2\pi k/n) \end{pmatrix}.$$

We consider n uniformly weighted mixtures of Gaussians, where for each sample X , we first sample a root of unity uniformly,

$$k \sim \text{Uniform}(n),$$

and conditionally on k , we sample X from an isotropic Gaussian centered at this root of unity

$$X \sim \mathcal{N}(\mathbf{c}_k, \sigma^2 \text{Id}_2),$$

using a standard deviation $\sigma = 0.1$. Samples are generated with the `make_blobs` function from `sklearn.datasets`. We generate two datasets in this way, $\mathbf{Z}^{(1)}$ for $n = 10$, and $\mathbf{Z}^{(2)}$ for $n = 5$. To generate dataset $\mathbf{Z}^{(1)}$, we first homogeneously scale the roots of unity to lie on a circle of radius 3 before sampling, so that the two datasets do not overlap but still use the same standard deviation. We do not scale the centers used for $\mathbf{Z}^{(2)}$, so they all lie on the unit circle. We generate $\mathbf{Z}^{(1)} = \{\mathbf{z}_1^{(1)}, \dots, \mathbf{z}_{1000}^{(1)}\}$ and $\mathbf{Z}^{(2)} = \{\mathbf{z}_1^{(2)}, \dots, \mathbf{z}_{1000}^{(2)}\}$ using 1000 samples each and form the empirical measures

$$\mu := \sum_{i=1}^{1000} \frac{1}{1000} \delta_{\mathbf{z}_i^{(1)}}, \quad \nu := \sum_{j=1}^{1000} \frac{1}{1000} \delta_{\mathbf{z}_j^{(2)}}.$$

supported on $\mathbf{Z}^{(1)}$ and $\mathbf{Z}^{(2)}$, corresponding to uniform probability vectors $\mathbf{a}, \mathbf{b} \in \Delta_{1000}$.

N.3 Experiment Settings

We used default hyperparameter settings for both methods. In particular, FRLC sets $\tau = 75$ and $\gamma = 90$, with a maximum of 200 iterations and a minimum of 25 iterations, subject to the stopping criterion Δ given in (10). The LOT default settings are $\gamma = 10$ (there is no hyperparameter analogous to τ in LOT). Both methods use a random initialization and were run on CPU.

O Discussion of differences between FRLC and existing low-rank optimal transport algorithms

We provide an extensive comparison of FRLC against the parametrization and objective of Scetbon et al. (2021) in Appendix A. As Latent-OT of Lin et al. (2021) is an extension of the k -Wasserstein barycenter problem, it has a distinct objective and thus performs worse on primal OT cost (Table 5), making a direct experimental comparison on primal cost minimization only appropriate relative to the works of Scetbon et al. (2021, 2023, 2022). This stated, we still list a number of distinctions between FRLC and Latent-OT – noting that most differences between FRLC and Forrow et al. (2019) transfer from this discussion since Lin et al. (2021) extends the k -Barycenter problem of Forrow et al. (2019). (1) Lin et al. optimize two sets of variables: sub-couplings $(\mathbf{Q}, \mathbf{R}, \mathbf{T})$ and anchor points (z^x, z^y) on which the sub-couplings depend. FRLC only has $(\mathbf{Q}, \mathbf{R}, \mathbf{T})$ as variables of the optimization. (2) Cost matrices used in Lin et al. (2021) are built from distances between each dataset and its representative anchor points for \mathbf{Q}, \mathbf{R} , or the distances between the two sets of anchor points for \mathbf{T} . In contrast, ground costs used in FRLC to update $(\mathbf{Q}, \mathbf{R}, \mathbf{T})$ are always derived from the distance matrix \mathbf{C} in the Wasserstein objective $\langle \mathbf{C}, \mathbf{P} \rangle_F$. Specifically, the cost matrices used by Lin et al. are:

$$[\mathbf{C}_{\mathbf{Q}}]_{ik} = \|x_i - z_k^x\|_2^2, \quad [\mathbf{C}_{\mathbf{R}}]_{j\ell} = \|y_j - z_\ell^y\|_2^2, \quad [\mathbf{C}_{\mathbf{T}}]_{k\ell} = \|z_k^x - z_\ell^y\|_2^2,$$

optionally using a Wasserstein distance for the entries of $\mathbf{C}_{\mathbf{T}}$. (3) FRLC costs are given in the exponents of the Gibbs kernels written above and below Equation 9. These are derived directly from the rank- r Wasserstein problem $\min_{\mathbf{P} \in \Pi_r(a,b)} \langle \mathbf{C}, \mathbf{P} \rangle_F$ and differ substantially from those of the proxy objective in Lin et al. (2021); Forrow et al. (2019). (4) The different objectives and variables lead to very different algorithms: Lin et al. alternate updates to the sub-couplings $(\mathbf{Q}, \mathbf{R}, \mathbf{T})$ using Dykstra, with updates to the latent anchor points (z^x, z^y) using first-order conditions. In contrast, FRLC alternates semi-relaxed OT to update (\mathbf{Q}, \mathbf{R}) and balanced OT to update \mathbf{T} . (5) Because FRLC does not require anchor points to define costs, FRLC can handle cost matrices which are not simple functions of distance. For example, if \mathbf{C}_{ij} is the price of transporting good i to warehouse j one may not be able to re-evaluate a price $c(x_i, z_k^x)$ between x_i and latent anchor z_k^x . In such situations, while finding a low-rank plan may make sense (e.g. to approximate an assignment for a massive dataset), an “anchor” may not have clear definition in the setting of general cost matrices. (6) The Lin et al. objective is only a proxy for a Wasserstein-type loss, and Lin et al. do not explore extensions to Gromov-Wasserstein (GW), or Fused GW, which FRLC readily generalizes to. A summary of the existing low-rank OT algorithms and key distinctions between them is given in Table 1.

For completeness, we offer a compare against the work Latent OT Lin et al. (2021), which solves a variation of the k -Wasserstein barycenter problem. As discussed, while their factorization is similar, their problem is distinct from FRLC as it does not solve the primal OT problem for general cost. We report the cost obtained by FRLC and by Latent OT on various simulated datasets in Table 5.

Dataset	OT-cost (FRLC)	OT-cost (Lin et al.)
5 th and 10 th roots of unity (rank $r_1, r_2 = 5, 10$)	1.174	2.124
Two-moons and 8-Gaussians (rank $r = 20$)	2.716	4.291
2D Gaussian mixture (rank $r = 20$)	0.552	0.922
10D Gaussian mixture (rank $r = 20$)	1.038	1.298

Table 5: Comparison against Lin et al. (2021) in primal OT-cost $\langle \mathbf{C}, \mathbf{P} \rangle_F$.

P Limitations

Our method introduces an additional hyperparameter τ relative to previous approaches Scetbon et al. (2021), controlling the strength of the KL penalty on the inner marginals when updating \mathbf{Q} and \mathbf{R} .

Empirically, we found FRLC to be robust to different choices of τ , but applying the method optimally requires this additional hyperparameter in any grid-search.

We also note that the non-asymptotic criterion $\Delta(\cdot, \cdot)$ is weak relative to stronger notions of convergence, and that often users might prefer to simply run the method up to some number of maximal iterations by setting the parameter for whether $\Delta(\cdot, \cdot)$ is used to `False`. The `W` optimization empirically converges to minima smoothly, so for the most part there is not much of a need for Δ except for early stopping. We recommend that users plot the loss over iterations and use it to set the tolerance parameter `tol`, and the minimum and maximum iteration parameters for the time-being. The needs for these parameters might vary widely—the minimum number of iterations should be very low (around 5) for simple datasets and substantially higher for high-dimensional, structured ones.

Although we demonstrate strong performance already, there is massive room for improvement as our implementation is preliminary and not at the level of a high-performance library like `ott-jax`. We use lightweight vanilla implementations of Sinkhorn as a sub-routine, not taking advantage of the momentum-based techniques which could accelerate it massively. Thus, one can imagine that the potential scalability and speed of this method could be much higher than reported in this document.

Q Broader impacts

FRLC is general enough to be used modularly within any ML algorithm using OT as a subroutine to help with scalability. We also note that the LC factorization is similar to a PCA in the context of OT, yielding an optimal low-rank coupling with an interpretable latent coupling factor.

Hyperparameter	Values
rank (Both)	50, 100, 200
τ (Ours)	30, 50, 100
τ (LOT-U)	0.99, 0.9, 0.7
ϵ (LOT-U)	0.001, 0.01, 0.1

Table 6: Hyperparameter grid considered in hyperparameter search for validation. Scetbon et al. (2023); Klein et al. (2023) scales $\tau' = \frac{\tau}{1-\tau}$ and their τ' are in the same range.

Solver	Rank	τ (Ours)	τ (UL)	ϵ (UL)	Spearman ρ (Validation)
FRLC-SR-W (Ours)	200	30	-	-	0.465
FRLC-SR-GW (Ours)	100	100	-	-	0.288
FRLC-SR-FGW (Ours)	200	50	-	-	0.465
FRLC-U-W (Ours)	200	100	-	-	0.471
FRLC-U-GW (Ours)	50	30	-	-	0.282
FRLC-U-FGW (Ours)	200	30	-	-	0.353
LOT-U-W	200	-	0.9	0.001	0.394
LOT-U-GW	100	-	0.99	0.01	0.001
LOT-U-FGW	200	-	0.7	0.001	0.393
LOT-SR-W	200	-	0.7	0.001	0.394
LOT-SR-GW	200	-	0.7	0.1	0.003
LOT-SR-FGW	200	-	0.99	0.001	0.399

Table 7: The best performing hyperparameters for each solver and the performance on the validation genes.

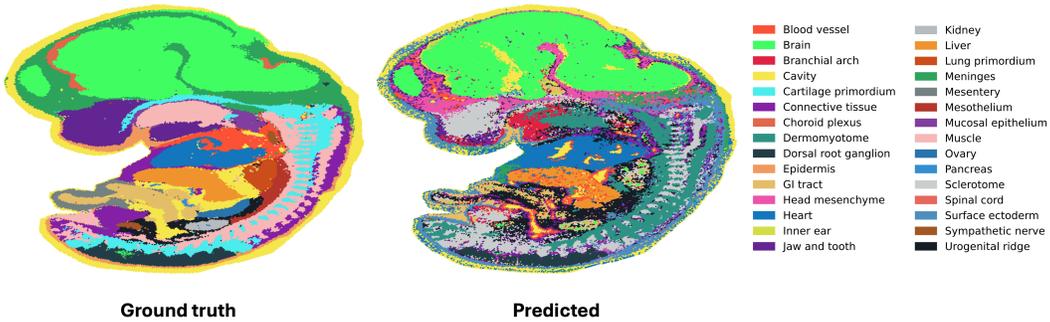


Figure 10: Ground truth and the predicted cell type classification of the E12.5 embryo.

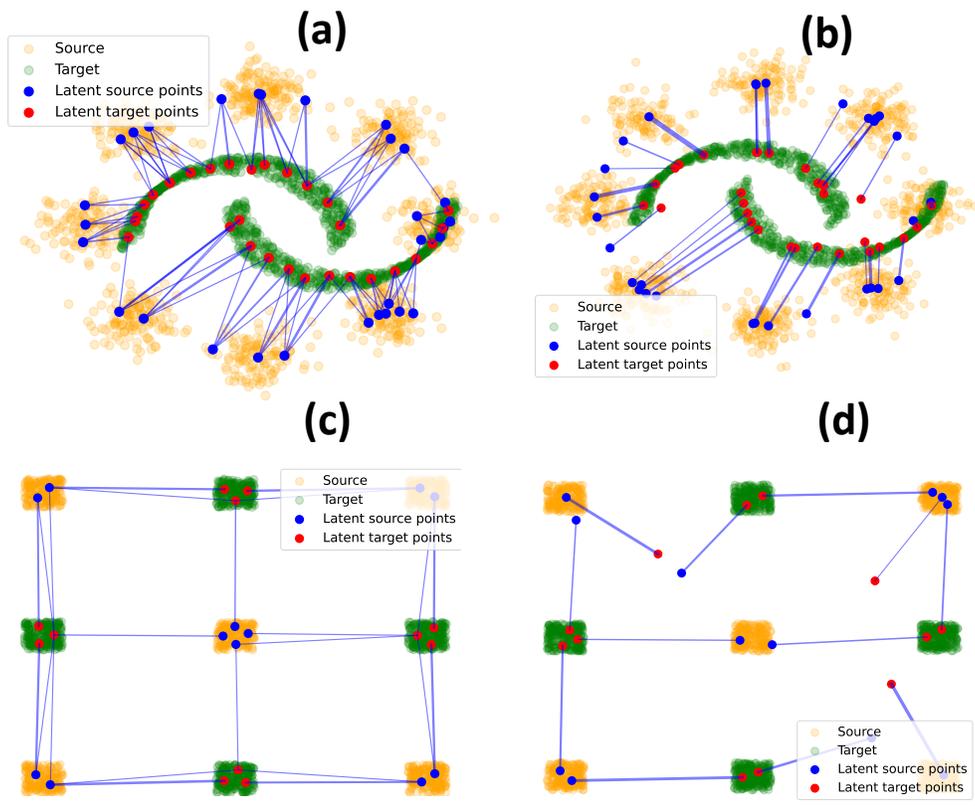


Figure 11: (a) LC-projection barycenters aligned with FRLC latent-coupling T on (a) two moons and eight Gaussians ($r = 30$), (b) LC-projection barycenters aligned with $\text{diag}(g)$ from Scetbon et al. (2021) ($r = 30$), (c) the checkerboard dataset with FRLC latent coupling aligned barycenters ($r = 12$), and (d) with diagonal alignment Scetbon et al. (2021) ($r = 12$). We show in A.1 that the output of FRLC can be diagonalized to the factorization of Forrow et al. (2019) (Figure 12).

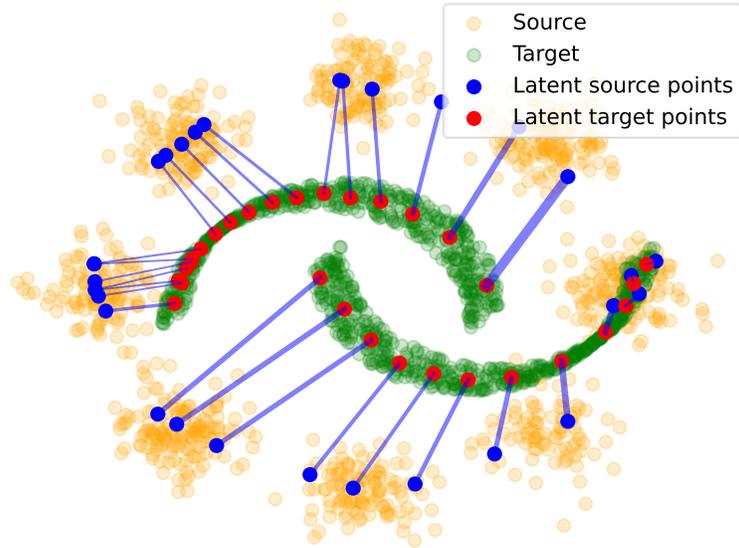


Figure 12: As discussed in A.1, one may recover the factorization of Forrow et al. (2019) as a sub-case of the LC-factorization. Shown is the factorization found by diagonalizing the output of FRLC from $(\mathbf{Q}, \mathbf{R}, \mathbf{T}) \mapsto (\mathbf{Q}', \mathbf{R}, \text{diag}(\mathbf{g}_R))$.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The abstract and introduction exactly describe the background preceding this paper and placing it into context, and exactly describe the contribution and scope of the paper to the field.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We describe the limitations of our method in Section P. These include an additional hyperparameter that FRLC introduces relative to previous work, and the substantial room for code optimization.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: All theoretical results state their assumptions clearly and all propositions are followed by thorough line-by-line proofs with careful justifications made for each step. They have checked over by all of the authors. All proofs are provided in the supplement with detail and are referenced in the main body where appropriate.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We provide all code used for generating the results of our paper. This not only includes the source code for the optimization, but the experimental code used for benchmarking. The algorithm is implemented exactly as described in the paper, and reviewers can freely consult the code we provide to them.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).

- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We will upload all code, and all code required to generate the synthetic data experiments used. Any real data experiments using especially large-scale data have publicly available and easily accessible datasets which we provide references for. We provide comprehensive descriptions of how the data was pre-processed, and provide experimental code which can be followed easily.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We introduce all hyper-parameters with rigorous justification of their utility. The values of the default hyper-parameters are accessible in the code we upload and any experiment which does not use the default hyperparameter (e.g. for a validation search) has the full table provided and the code used for the experiment.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: The optimization is deterministic, so generally there's no need for reporting statistical significance. One small exception is that one of our proposed initializations is randomized, and for experiments which use it we do include $\pm 1\sigma$ error bars.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We indicate when experiments are run on CPU versus GPU in the experimental details.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: This research was conducted ethically and conforms to the guidelines of the code of ethics. We do not anticipate any major negative societal impact to result from this work on optimal transport.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.

- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We address broader impacts in Section Q: FRLC can be used for scaling an interpretability in any ML method using OT as a subroutine.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: We do not generate any new datasets other than a toy dataset of Gaussians centered at n^{th} roots of unity, which has no societal impact. We do not believe our optimal transport work has significant risk of misuse.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: All data that we use and all work that we build on are cited and credited heavily.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: We provide no new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: This paper involves no research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.