
Identifiable Shared Component Analysis of Unpaired Multimodal Mixtures

Subash Timilsina
School of EECS
Oregon State University
Corvallis, OR 97331
timilsis@oregonstate.edu

Sagar Shrestha
School of EECS
Oregon State University
Corvallis, OR 97331
shressag@oregonstate.edu

Xiao Fu
School of EECS
Oregon State University
Corvallis, OR 97331
xiao.fu@oregonstate.edu

Abstract

A core task in multi-modal learning is to integrate information from multiple feature spaces (e.g., text and audio), offering modality-invariant essential representations of data. Recent research showed that, classical tools such as *canonical correlation analysis* (CCA) provably identify the shared components up to minor ambiguities, when samples in each modality are generated from a linear mixture of shared and private components. Such identifiability results were obtained under the condition that the cross-modality samples are aligned/paired according to their shared information. This work takes a step further, investigating shared component identifiability from multi-modal linear mixtures where cross-modality samples are unaligned. A distribution divergence minimization-based loss is proposed, under which a suite of sufficient conditions ensuring identifiability of the shared components are derived. Our conditions are based on cross-modality distribution discrepancy characterization and density-preserving transform removal, which are much milder than existing studies relying on independent component analysis. More relaxed conditions are also provided via adding reasonable structural constraints, motivated by available side information in various applications. The identifiability claims are thoroughly validated using synthetic and real-world data.

1 Introduction

The same data entities can often be represented in different feature spaces (e.g., audio, text and image), due to the variety of sensing modalities or domains. Learning common latent components of data from multiple modalities is well-motivated in representation learning. The shared components are considered modality-invariant essential representations of data, which can often enhance performance of downstream tasks by shedding modality-specific noise [1-4] and avoiding over-fitting [5-7].

A prominent theoretical aspect of shared component learning lies in *identifiability* of the components of interest. The literature posed an intriguing theoretical question [1, 2, 8]: If every modality of data is represented by a linear mixture of shared and private components with an unknown mixing system, are the shared components identifiable (up to acceptable ambiguities)? Such component identification problems are often nontrivial due to the ill-posed nature of any linear mixture model (see, e.g., [9-14]). Interestingly, the work [1] showed that using the classical *canonical correlation*

analysis (CCA) provably find the shared components up to rotation and scaling. In fact, shared component identification from multimodal/multiview linear mixtures were considered in various contexts (see, e.g., [15–18]), although some of these works did not model private components. The identifiability results in [1, 2] were generalized to nonlinear mixture models as well [4, 19]. The shared component identification perspective was also related to the success of representation learning in self-supervised learning (SSL) [5–7].

Nonetheless, the treatment in [1, 2] and the related works [15–17] all assumed that the cross-modality data are aligned (i.e., paired) according to their shared components. In many applications, such as cross-language information retrieval [20–22], domain adaptation [23–25], and biological data translation [26, 27], aligned cross-modality data are hard to acquire, if not outright unavailable. A natural question is: *When the multimodal linear mixtures are unaligned, can the shared latent components still be provably identified under reasonably mild conditions?*

Existing Studies. Theoretical characteristics of unaligned multimodal learning were studied under various settings. The work [28] considered a case where one modality is a linear transform of another modality, and showed that the linear transformation is potentially identifiable. The recent work [29] extended this model to a nonlinear transform setting. However, these works did not consider *latent* component models—yet the latter are more versatile in many ways, e.g., facilitating one-to-many cross-domain translations [30, 31]. The work [32] considered unaligned mixtures of shared and private components, but the assumptions (e.g., the availability of a large amount of modalities) to ensure identifiability may not be easy to satisfy. The most related work is perhaps [8]. But their approach also relied on somewhat stringent assumptions, e.g., that all the latent components are element-wise statistically independent with at most one component being Gaussian. This is because their procedure had to invoke the classical *independent component analysis* (ICA) [33].

Contributions. In this work, we provide a suite of sufficient conditions under which the shared components can be provably identified from unaligned multimodal linear mixtures up to reasonable ambiguities. The model and identification problem are referred to as *unaligned shared component analysis* (unaligned SCA) in the sequel.

(i) *An Identifiable Learning Loss for Unaligned SCA.* We propose to tackle the unaligned SCA problem by matching the probability distributions of linearly embedded multi-modal data. We show that under reasonable conditions, the linear transformations identifies the shared components up to the same ambiguities as those in the aligned case [1, 2]. The conditions are considerably milder compared to the existing unaligned SCA work [8].

(ii) *Enhanced Identifiability via Structural Constraints.* We come up with two types of structural constraints, motivated by available side information in applications, to further relax the identifiability conditions. Specifically, we look into cases where the multi-modal data have similar linear mixing systems and cases where a few cross-domain aligned samples available. We show that by adding constraints accordingly, unaligned SCA are identifiable under much milder conditions.

Our contributions primarily lie in identifiability analysis. Nonetheless, we also show the usefulness of our results in real-world applications, namely, *cross-lingual word retrieval*, *genetic information alignment* and *image data domain adaptation*. Particularly, it shows that our succinct multimodal linear mixture model can effectively post-process outputs of pre-trained encoders, e.g., those in [34, 35], to improve data representations and enhance downstream task performance.

Notation. Notation definitions can be found in Appendix A.

2 Background

Generative Model of Interest. Following the classical settings in [1, 2, 15, 16, 18], we consider modeling the multi-modal data as linear mixtures. More specifically, we adopt the model in [1, 2] that splits the latent representation of data into shared components and private components:

$$\mathbf{x}^{(q)} = \mathbf{A}^{(q)} \mathbf{z}^{(q)}, \quad \mathbf{z}^{(q)} = [\mathbf{c}^\top, (\mathbf{p}^{(q)})^\top]^\top, \quad q = 1, 2, \quad (1)$$

where $\mathbf{x}^{(q)} \in \mathbb{R}^{d^{(q)}}$ represents the data from the q th modality, $\mathbf{z}^{(q)} \in \mathbb{R}^{d_c + d_p^{(q)}}$ represents the corresponding latent code, $\mathbf{c} \in \mathbb{R}^{d_c}$ and $\mathbf{p}^{(q)} \in \mathbb{R}^{d_p^{(q)}}$ stand for the shared components and the private components, respectively. The data $\mathbf{x}^{(q)}$'s are assumed to be zero-mean, which can be enforced by

centering. Note that the positions of \mathbf{c} and \mathbf{p}_q are not necessarily arranged as $[\mathbf{c}^\top, (\mathbf{p}^{(q)})^\top]^\top$ (more generally, $\mathbf{z}^{(q)} = \mathbf{\Pi}^{(q)}[\mathbf{c}^\top, (\mathbf{p}^{(q)})^\top]^\top$ with an unknown permutation matrix $\mathbf{\Pi}^{(q)}$). However, the representation in (1) is without loss of generality as one can define $\mathbf{A}^{(q)} := \mathbf{A}^{(q)}(\mathbf{\Pi}^{(q)})^\top$ to reach the representation in (1). For all the domains, we have

$$\mathbf{c} \sim \mathbb{P}_{\mathbf{c}}, \quad \mathbf{p}^{(q)} \sim \mathbb{P}_{\mathbf{p}^{(q)}}, \quad (2)$$

where $\mathbb{P}_{\mathbf{c}}$ and $\mathbb{P}_{\mathbf{p}^{(q)}}$ represent the distributions of the shared components and the domain-private components, respectively. Under (1), the two different range spaces $\text{range}(\mathbf{A}^{(q)})$ for $q = 1, 2$ represent two feature spaces. Then latent $\mathbf{p}^{(q)}$ further distinguishes the modalities and often has interesting physical interpretation. For example, some vision literature use \mathbf{c} to model “content” and $\mathbf{p}^{(q)}$ “style” of the images [31, 36]. In cross-lingual word embedding retrieval [2], \mathbf{c} represents the semantic meaning of the words, while $\mathbf{p}^{(q)}$ represents the language-specific components. The goal of SCA boils down to finding linear operators to recover \mathbf{c} to a reasonable extent.

Aligned SCA: Identifiability of CCA and Extensions. Learning \mathbf{c} without knowing $\mathbf{A}^{(q)}$ is a typical component analysis problem. Learning latent components from *linear mixture models* (LMMs) like $\mathbf{x} = \mathbf{A}\mathbf{z}$ lacks identifiability in general, due to the bilinear nature of the models. This is because one can find an infinite number of invertible matrices \mathbf{B} such that $\mathbf{x} = \mathbf{A}\mathbf{B}\mathbf{B}^{-1}\mathbf{z}$. Then, both (\mathbf{A}, \mathbf{z}) and $(\mathbf{A}\mathbf{B}, \mathbf{B}^{-1}\mathbf{z})$ can fit to the data \mathbf{x} , making the problem ill-posed in terms of solution uniqueness; see, e.g., [9, 37] and more discussions in Sec. 5. Nonetheless, the works [1, 2] studied the identifiability of \mathbf{c} under the model (1), using the assumption that the cross-modality samples share the same \mathbf{c} are aligned. In particular, [1] formulated the \mathbf{c} -identification problem as a CCA problem:

$$\underset{\{\mathbf{Q}^{(q)}\}_{q=1}^2}{\text{minimize}} \quad \mathbb{E} \left[\left\| \mathbf{Q}^{(1)}\mathbf{x}^{(1)} - \mathbf{Q}^{(2)}\mathbf{x}^{(2)} \right\|_2^2 \right] \quad (3a)$$

$$\text{subject to} \quad \mathbf{Q}^{(q)} \mathbb{E} \left[\mathbf{x}^{(q)}(\mathbf{x}^{(q)})^\top \right] (\mathbf{Q}^{(q)})^\top = \mathbf{I} \quad q = 1, 2, \quad (3b)$$

where $\mathbf{Q}^{(q)} \in \mathbb{R}^{d_C \times d^{(q)}}$. The expectation in (3a) is taken from the joint distribution of the *aligned pairs* $\mathbb{P}_{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}}$, where every pair $(\mathbf{x}^{(1)}, \mathbf{x}^{(2)})$ shares the same \mathbf{c} . The formulation aims to find $\mathbf{Q}^{(q)}$ such that the transformed representations of the aligned pairs $\mathbf{Q}^{(1)}\mathbf{x}^{(1)}$ and $\mathbf{Q}^{(2)}\mathbf{x}^{(2)}$ are equal. In [1], it was shown that

$$\hat{\mathbf{Q}}^{(q)}\mathbf{x}^{(q)} = \mathbf{\Theta}\mathbf{c} \quad (4)$$

under mild conditions (see Appendix E.1 for details), where $(\hat{\mathbf{Q}}^{(1)}, \hat{\mathbf{Q}}^{(2)})$ is an optimal solution of the CCA formulation and $\mathbf{\Theta}$ is a certain non-singular matrix. Eq. (4) means that $\hat{\mathbf{Q}}^{(q)}$ finds the range space where \mathbf{c} lives in, i.e., $\text{range}(\mathbf{A}_{1:d_C}^{(q)})$ under our notation.

Unaligned SCA: Existing Result and Theoretical Gap. The work in [8] studied the identifiability of \mathbf{c} under (1) when $\mathbf{x}^{(1)}$ and $\mathbf{x}^{(2)}$ are *unaligned*. Their approach works under the condition that the elements of $\mathbf{z}^{(q)} = [\mathbf{c}^\top, (\mathbf{p}^{(q)})^\top]^\top$ are mutually statistically independent. There, $\hat{\mathbf{z}}^{(q)} = \mathbf{\Pi}^{(q)}\mathbf{\Sigma}^{(q)}\mathbf{z}^{(q)}$ is assumed to have been estimated by ICA, where $\mathbf{\Pi}^{(q)}$ and $\mathbf{\Sigma}^{(q)}$ represent the scaling and permutation ambiguities, respectively, which cannot be removed by ICA. The work [8] assumed $\mathbf{\Sigma}^{(q)} = \mathbf{I}$ by imposing a unit-variance assumption on all the $z_i^{(q)}$ ’s. Then, a cross-domain matching algorithm is used to match the shared elements in $\hat{\mathbf{z}}^{(1)}$ and $\hat{\mathbf{z}}^{(2)}$. The formulation can be summarized as finding d_C pairs of non-repetitive (i, j) such that $\mathbf{e}_i^\top \hat{\mathbf{z}}^{(1)}$ and $\mathbf{e}_j^\top \hat{\mathbf{z}}^{(2)}$ have identical distributions, where \mathbf{e}_i is the i th unit vector. Denote $\hat{\mathbf{c}}_m^{(1)} = \mathbf{e}_{i_m}^\top \hat{\mathbf{z}}^{(1)}$ and $\hat{\mathbf{c}}_m^{(2)} = \mathbf{e}_{j_m}^\top \hat{\mathbf{z}}^{(2)}$ for $m \in [d_C]$. It can be shown that

$$\hat{\mathbf{c}}_m^{(q)} = k\mathbf{c}_{\pi(m)}^{(q)}, \quad m \in [d_C], \quad (5)$$

where $k \in \{+1, -1\}$ and π is a permutation of $\{1, \dots, d_C\}$ (see details in Appendix E.2 summarized from [8]). This method effectively applies ICA to each modality, and thus the ICA identifiability conditions [33] have to met by $\mathbf{x}^{(1)}$ and $\mathbf{x}^{(2)}$ individually. However, if one only aims to extract $\mathbf{\Theta}\mathbf{c}$ as in CCA, these assumptions appear to be overly stringent.

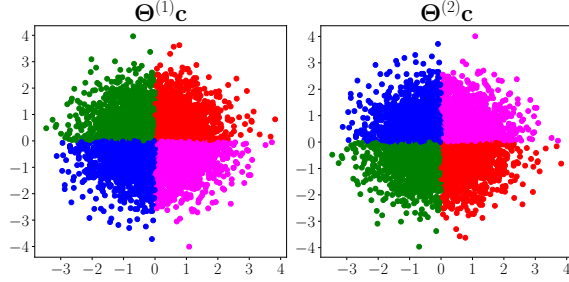


Figure 1: Scatter plots of matched distribution $\Theta^{(1)}c$ (left) and $\Theta^{(2)}c$ (right) when c follows the Gaussian distribution. Colors in the scatter plot represent alignment; same color represent the data are aligned.

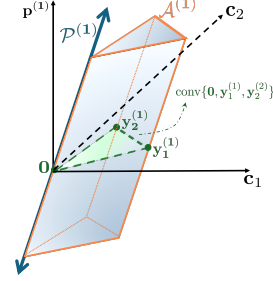


Figure 2: Illustration of $\mathcal{A}^{(1)}$ in Assumption 1 in a case where $d_C = 2$ and $d_P^{(1)} = 1$.

3 Proposed Approach

Unaligned SCA: Problem Formulation We assume that $x^{(q)}$'s are zero-mean. We use the notation from CCA in (3a). However, since no aligned samples are available, we replace the sample-level matching objective with a distribution matching (DM) module, as DM can be carried out without sample level alignment:

$$\text{find } Q^{(q)} \in \mathbb{R}^{d_C \times d^{(q)}}, q = 1, 2, \quad (6a)$$

$$\text{subject to } Q^{(1)}x^{(1)} \stackrel{(d)}{=} Q^{(2)}x^{(2)}, \quad (6b)$$

$$Q^{(q)}\mathbb{E}[x^{(q)}(x^{(q)})^\top](Q^{(q)})^\top = I \quad q = 1, 2. \quad (6c)$$

where “ $u \stackrel{(d)}{=} v$ ” means the distributions of u and v are the same.

The formulation in (6) can be realized using various distribution matching tools, e.g., *maximum mean discrepancy* (MMD) [38] and *Wasserstein distance* [39]. We use the adversarial loss:

$$\min_{Q^{(1)}, Q^{(2)}} \max_f \mathbb{E}_{x^{(1)}} \log(f(Q^{(1)}x^{(1)})) + \mathbb{E}_{x^{(2)}} \log(1 - f(Q^{(2)}x^{(2)})) + \lambda \sum_{q=1}^2 \mathcal{R}(Q^{(q)}), \quad (7)$$

The first and second terms comprise the adversarial loss from GAN [40]. It finds $Q^{(q)}$ to confuse the best-possible discriminator $f: \mathbb{R}^{d_C} \rightarrow \mathbb{R}$, where f is represented by a neural network in practice. It is well known that the minimax optimal point of the first two terms is attained when (6b) is met [40]. We use $\mathcal{R}(Q^{(q)}) = \|Q^{(q)}\mathbb{E}[x^{(q)}(x^{(q)})^\top](Q^{(q)})^\top - I\|_F^2$ to “lift” the constraints. This way, the learning criterion in (7) can be readily handled by any off-the-shelf adversarial learning tools.

Identifiability of Unaligned SCA As we saw in Theorem 4, CCA identifies $\hat{Q}^{(q)}x^{(q)} = \Theta c$ where $\Theta \in \mathbb{R}^{d_C \times d_C}$ under the settings of aligned SCA. Establishing a similar result for unaligned SCA is much more challenging. First, it is unclear if (6b) could disentangle c from $p^{(q)}$. In general, $Q^{(q)}x^{(q)}$ could still be a mixture of c and $p^{(q)}$ yet (6b) still holds (e.g., when both c and $p^{(q)}$ are Gaussian.)

Second, even when the disentanglement is attained via enforcing (6b) and we have $Q^{(q)}x^{(q)} = \Theta^{(q)}c$, in general it does not hold that $\Theta^{(1)} = \Theta^{(2)}$. This is because $\Theta^{(1)}c \stackrel{(d)}{=} \Theta^{(2)}c$ where $\Theta^{(1)} \neq \Theta^{(2)}$ can still be perfectly met (e.g., when $\mathbb{P}_{\Theta^{(q)}c}$ is symmetric Gaussian in Fig. 1). However, $\Theta^{(1)} \neq \Theta^{(2)}$ means that the extracted representations from the two modalities are not matched. This creates challenges for applications like cross-domain information retrieval, language translation, or domain adaptation.

Our intuition is as follows: If the two distributions $\mathbb{P}_{c, p^{(1)}}$ and $\mathbb{P}_{c, p^{(2)}}$ are very different, then $Q^{(1)}x^{(1)} \stackrel{(d)}{=} Q^{(2)}x^{(2)}$ cannot hold unless $Q^{(q)}A^{(q)} = [\Theta^{(q)}, 0]$. We use the following to characterize such difference between the joint distributions:

Assumption 1 (Modality Variability). *For any two linear subspaces $\mathcal{P}^{(q)} \subset \mathbb{R}^{d_C + d_P^{(q)}}$, $q = 1, 2$, with $\dim(\mathcal{P}^{(q)}) = d_P^{(q)}$, $\mathcal{P}^{(q)} \neq \mathbf{0} \times \mathbb{R}^{d_P^{(q)}}$ and linearly independent vectors $\{\mathbf{y}_i^{(q)} \in \mathbb{R}^{d_C + d_P^{(q)}}\}_{i=1}^{d_C}$, $q = 1, 2$, the sets $\mathcal{A}^{(q)} = \text{conv}\{\mathbf{0}, \mathbf{y}_1^{(q)}, \dots, \mathbf{y}_{d_C}^{(q)}\} + \mathcal{P}^{(q)}$, $q = 1, 2$, are such that if $\mathbb{P}_{\mathbf{c}, \mathbf{p}^{(q)}}[\mathcal{A}^{(q)}] > 0$ for $q = 1$ or $q = 2$, then there exists a $k \in \mathbb{R}$ such that the joint distributions $\mathbb{P}_{\mathbf{c}, \mathbf{p}^{(1)}}[k\mathcal{A}^{(1)}] \neq \mathbb{P}_{\mathbf{c}, \mathbf{p}^{(2)}}[k\mathcal{A}^{(2)}]$, where $k\mathcal{A}^{(q)} = \{k\mathbf{a} \mid \mathbf{a} \in \mathcal{A}^{(q)}\}$.*

The condition in Assumption 1 is a geometric way to characterize the difference between $\mathbb{P}_{\mathbf{c}, \mathbf{p}^{(1)}}$ and $\mathbb{P}_{\mathbf{c}, \mathbf{p}^{(2)}}$ —if the joint distributions have different measures for all possible “stripes”, each being a direct sum of a subspace and a convex hull (see Fig. 2), then $\mathbb{P}_{\mathbf{c}, \mathbf{p}^{(1)}}$ and $\mathbb{P}_{\mathbf{c}, \mathbf{p}^{(2)}}$ must be very different. Note that the difference is contributed by the modality-specific term $\mathbf{p}^{(q)}$, and thus we call this condition “modality variability”. Modality variability is similar to the “domain variability” used in [32, 41]—both characterize the discrepancy of the joint probabilities $\mathbb{P}_{\mathbf{c}, \mathbf{p}^{(1)}}$ and $\mathbb{P}_{\mathbf{c}, \mathbf{p}^{(2)}}$. However, there are key differences: The domain variability was defined in a unified latent domain over arbitrary sets \mathcal{A} , which could be stringent. Instead, we use the fact that (6) relies on linear operations to construct $\mathcal{A}^{(q)}$, which makes the condition defined over a much smaller class of sets—thereby largely relaxing the requirements. Restricting $\mathcal{A}^{(q)}$ to be stripes also makes the modality variability condition much more relaxed compared to the domain variability condition.

We show the following:

Theorem 1. *Under Assumption 1 and the generative model in (1), denote any solution of (6) as $\hat{\mathbf{Q}}^{(q)}$, $q = 1, 2$. Then, if the mixing matrices $\mathbf{A}^{(q)}$ are full column ranks and $\mathbb{E}[\mathbf{c}\mathbf{c}^\top]$ is full rank, we have $\hat{\mathbf{Q}}^{(q)} \mathbf{x}^{(q)} = \Theta^{(q)} \mathbf{c}$. In addition, assume that either of the following is satisfied:*

- (a) *The individual elements of the content components are statistically independent and non-Gaussian. In addition, $c_i \neq kc_j, \forall i \neq j, \forall k \in \mathbb{R}$ and $c_i \neq -c_i, \forall i$, i.e., the marginal distributions of the content elements cannot be matched with each other by mere scaling.*
- (b) *The support of $\mathbb{P}_{\mathbf{c}}$, denoted by \mathcal{C} , is a hyper-rectangle, i.e., $\mathcal{C} = [-a_1, a_1] \times \dots \times [-a_{d_C}, a_{d_C}]$. Further, suppose that $c_i \neq kc_j, \forall i \neq j, \forall k \in \mathbb{R}$ and $c_i \neq -c_i, \forall i$.*

Then, we have $\hat{\mathbf{Q}}^{(q)} \mathbf{x}^{(q)} = \Theta \mathbf{c}$, i.e., $\Theta^{(q)} = \Theta$ for all $q = 1, 2$, where $\Theta^{(q)}$.

In Theorem 1, Assumption 1 is used to guarantee $\hat{\mathbf{Q}}^{(q)} \mathbf{x}^{(q)} = \Theta^{(q)} \mathbf{c}$ and either of conditions (a) or (b) is used to make sure $\Theta^{(1)} = \Theta^{(2)}$. Note that both (a) and (b) are milder than those in [8] (cf. Theorem 5), where the element-wise statistical independence of $\mathbf{z}^{(q)}$ was relied on to find shared representation of $\mathbf{x}^{(1)}$ and $\mathbf{x}^{(2)}$. The proof is in Appendix B.

Numerical Validation. In Fig. 3, the top and bottom rows validate Theorem 1 under the assumptions in (a) and (b), respectively. In the top row, we set $\mathbf{c} \in \mathbb{R}^2$, where c_1 is sampled from Gaussian mixtures with three components and c_2 is sampled from a Gamma distribution (and $c_1 \perp c_2$). We set $p^{(1)}$ and $p^{(2)}$ as one-dimensional Laplacian and uniform distributions. In the bottom row, the dimensions of \mathbf{c} and $\mathbf{p}^{(q)}$ for $q = 1, 2$ are unchanged, but their distributions are replaced in order to satisfy conditions in (b) (see details in Appendix F). One can see that clearly $\hat{\mathbf{c}}^{(q)} = \Theta \mathbf{c}$; i.e., the learned $\hat{\mathbf{c}}^{(q)}$ for $q = 1, 2$ are identically rotated and scaled versions of \mathbf{c} .

A remark is that our framework still allows to identify individual c_i ’s as in [8].

Corollary 1. *Under the conditions in Theorem 1 (a), Assume that at most one c_i for $i \in [d_C]$ is Gaussian. Then, the components of \mathbf{c} are identifiable up to permutation and scaling ambiguities by applying ICA to $\hat{\mathbf{c}}^{(q)} = \hat{\mathbf{Q}}^{(q)} \mathbf{x}^{(q)}$ for either $q = 1$ or $q = 2$.*

The corollary means that to identify individual c_i , using our formulation still enjoys much milder conditions relative to [8]. Specifically, our condition only specifies the independence among elements of \mathbf{c} , but the condition in [8] needs that all the elements in $\mathbf{z}^{(q)} = [\mathbf{c}^\top, (\mathbf{p}^{(q)})^\top]^\top$ are independent.

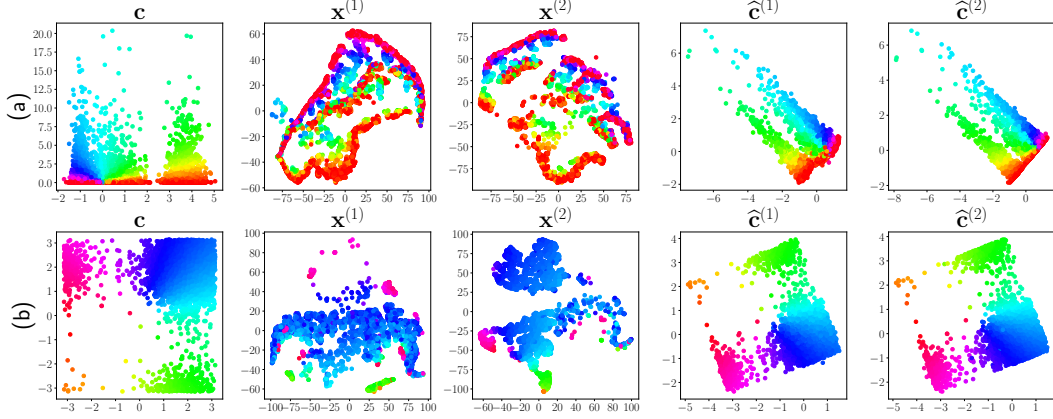


Figure 3: Validation of Theorem 1. Top row: results under assumption (a). Bottom row: results under assumption (b).

4 Enhanced Identifiability via Structural Constraints

Theorem 1 was well-supported by the synthetic data experiments. However, our experiments found that the learning criterion (6) often struggles to produce sensible results in some applications. Our conjecture is that the Assumptions in Theorem 1 (a) and (b) might not have been satisfied by the real data under our tests. Although they are not necessary conditions for identifiability, these conditions do indicate that the requirements to guarantee identifiability of unaligned SCA using (6) are nontrivial to meet. In this section, we explore a couple of structural constraints arising from side information in applications to remove the need for the relatively stringent assumptions on \mathbf{c} .

Homogeneous Domains. The first structural constraint that we consider is $\mathbf{A}^{(q)} = \mathbf{A}$ for $q = 1, 2$. This model is motivated by the fact that advanced representation learning tools, e.g., self-supervised learning tools (e.g., SimCLR [42]) and foundation models (e.g., CLIP [35]), are already capable of mapping the data clusters to a shared linearly separable space—which indicates that the representations share a subspace, i.e., $\mathbf{x}^{(q)} \approx \mathbf{A}\mathbf{z}^{(q)}$. Under such circumstances, the proposed model and method can be used to further process the data by discarding the private components in the latent representation.

Here, we consider the special case of generative process in (1) where,

$$\mathbf{x}^{(q)} = \mathbf{A}[\mathbf{c}^\top, (\mathbf{p}^{(q)})^\top]^\top. \quad (8)$$

Under this model, we look for the shared components by solving (6) with a single $\mathbf{Q} = \mathbf{Q}^{(1)} = \mathbf{Q}^{(2)}$. We use the following version of the modality variability condition:

Assumption 2. For any linear subspace $\mathcal{P} \subset \mathbb{R}^{d_C + d_P}$, $d_P = d_P^{(1)} = d_P^{(2)}$, with $\dim(\mathcal{P}) = d_P$, $\mathcal{P} \neq \mathbf{0} \times \mathbb{R}^{d_P}$ and linearly independent vectors $\{\mathbf{y}_i \in \mathbb{R}^{d_C + d_P}\}_{i=1}^{d_C}$, $q = 1, 2$, the sets $\mathcal{A} = \text{conv}\{\mathbf{0}, \mathbf{y}_1, \dots, \mathbf{y}_{d_C}\} + \mathcal{P}$, $q = 1, 2$, are such that if $\mathbb{P}_{\mathbf{c}, \mathbf{p}^{(q)}}[\mathcal{A}] > 0$ for $q = 1$ or $q = 2$, then the joint distributions $\mathbb{P}_{\mathbf{c}, \mathbf{p}^{(1)}}[k\mathcal{A}] \neq \mathbb{P}_{\mathbf{c}, \mathbf{p}^{(2)}}[k\mathcal{A}]$ for some $k \in \mathbb{R}$.

Theorem 2. Consider the mixture model in (8). Assume that $\text{rank}(\mathbf{A}) = d_C + d_P$ and $\text{rank}(\mathbb{E}[\mathbf{c}\mathbf{c}^\top]) = d_C$, and that Assumption 2 holds. Denote $\hat{\mathbf{Q}}$ as any solution of (6) by constraining $\mathbf{Q} = \mathbf{Q}^{(1)} = \mathbf{Q}^{(2)}$. Then, we have $\hat{\mathbf{Q}}\mathbf{x}^{(q)} = \Theta\mathbf{c}$.

One can see that the conditions (a) and (b) in Theorem 1 are completely removed, if the structure $\mathbf{A}^{(1)} = \mathbf{A}^{(2)}$ is imposed. In fact, the result in Theorem 2 is expected and readily seen from the proof of Theorem 1, as the cause for $\Theta^{(1)} \neq \Theta^{(2)}$ is the use of two different $\mathbf{Q}^{(q)}$'s. Nonetheless, this simple variation will prove useful in a series of real-data experiments.

The Weakly Supervised Case. Another way to add structural constraints is to use available auxiliary information. For example, some datasets have weak annotations and selected pairs; see, e.g., [43, 44].

Assumption 3 (Weak Supervision). There exist a set of available aligned samples $(\mathbf{x}_\ell^{(1)}, \mathbf{x}_\ell^{(2)})$ for $\ell \in \mathcal{L}$ such that $\mathbf{x}_\ell^{(q)} = \mathbf{A}^{(q)}\mathbf{z}_\ell^{(q)}$, $\mathbf{z}_\ell^{(q)} = [\mathbf{c}_\ell^\top, (\mathbf{p}_\ell^{(q)})^\top]^\top$; i.e., $(\mathbf{x}_\ell^{(1)}, \mathbf{x}_\ell^{(2)})$ share the same \mathbf{c}_ℓ .

The condition can be added into our formulation in (6) as a constraint, i.e.,

$$\mathbf{Q}^{(1)}\mathbf{x}_\ell^{(1)} = \mathbf{Q}^{(2)}\mathbf{x}_\ell^{(2)}, \forall \ell \in \mathcal{L}. \quad (9)$$

In the next theorem, we show that the incorporation of aligned samples helps relax conditions (a) and (b) in Theorem 1.

Theorem 3. Assume that Assumption 1 is satisfied, that $|\mathcal{L}| \geq d_C$ paired samples $(\mathbf{x}_\ell^{(1)}, \mathbf{x}_\ell^{(2)})$ are available, that $\mathbf{A}^{(q)}$ for $q = 1, 2$ have full column rank, and that $\mathbb{P}_\mathcal{C}$ is absolutely continuous. Denote $(\hat{\mathbf{Q}}^{(1)}, \hat{\mathbf{Q}}^{(2)})$ as any optimal solution of (6) under the constraint (9). Then, we have $\hat{\mathbf{Q}}^{(q)}\mathbf{x}^{(q)} = \Theta\mathbf{c}$.

The proof and synthetic data validation can be found in Appendices D and F respectively. Note that to realize (9), one only needs to add a regularization term $\beta \sum_{\ell \in \mathcal{L}} \|\mathbf{Q}^{(1)}\mathbf{x}_\ell^{(1)} - \mathbf{Q}^{(2)}\mathbf{x}_\ell^{(2)}\|_2^2$ to the loss in (7), where $\beta \geq 0$ is a tunable parameter. The overall loss is still differentiable and thus can be easily handled by gradient based approaches.

A remark is that our weakly supervised formulation can use as few as d_C pairs of $(\mathbf{x}_\ell^{(1)}, \mathbf{x}_\ell^{(2)})$ to establish identifiability of shared component. In contrast, CCA requires at least $d_C + d_p^{(1)} + d_p^{(2)}$ pairs to attain the same identifiability (cf. Appendix E.1).

Private Component Identifiability. Although our focus is shared component identification, we show that private components are also identifiable with additional assumptions; see Appendix H.

5 Related Works

Identifiability of Component Analysis under Linear Mixture Models. Various component analysis models were studied in the past several decades, e.g., principal component analysis [45], independent component analysis [33], sparse component analysis [10, 12], bounded component analysis [13], simplex component analysis [46, 47], and polytopic component analysis [14]—motivated by their applications in dimensionality reduction, representation learning, and latent variable identification (see, e.g., topic mining [48, 49], hyperspectral unmixing [46, 47], audio/speech separation [33] and community detection [50]). The classical component analysis tools mostly study a single modality. The identifiability results under these models are well developed and documented.

Identifiability of Shared Components from Aligned Modalities. Modeling multimodal data as two or more linear/nonlinear mixtures of latent components was considered in CCA-related works [1, 2, 15, 19], independent vector analysis (IVA) works [17, 18], multiview ICA works [16, 51], and SSL works [5–7, 52]. Partitioning the latent components into shared and private blocks was considered in [1, 2, 4, 5, 7, 52]. Shared component identifiability was established at the block level (see, e.g., [1, 2, 5]) and the individual component level (e.g., [51]) in these works. Nonetheless, they all rely on completely paired/aligned cross-modality samples, which we do not use in this work.

Distribution Matching and Unaligned Multimodal Analysis. Using distribution matching in unaligned multimodal data analytics for different purpose also has a long history; see applications in image-to-image translation [53], domain adaptation [54], cross-platform image super-resolution [55], and cross-domain information retrieval [21]. The recent works [56] and [57] pointed out the identifiability challenge and the existence of density-preserving transforms. The works in [28, 29] started studying the uniqueness issues in distribution matching. However, the latent mixture models were not studied in this line of work.

Identifiability of Unaligned SCA. The works in [32, 41] investigated the shared component identifiability when the multimodal data are nonlinear mixtures of content and style (which are shared and private components, respectively) under the same mixing system. Hence, our identical linear mixing case in Theorem 2 can be understood as a special case of theirs. But their analysis relies on the assumption that all the latent components are statistically independent, which is much stronger than our conditions in Theorem 2. Their results also require that there are a large amount of modalities available. But our proof works for just two modalities. The most related work is [8], which uses the model in (1) in the context of multi-view causal graph learning. As discussed before, their assumptions on the latent components are much stronger than ours (see Corollary 1 and Appendix E.2).

Table 1: Classification accuracy on the target domain of *office-31* dataset (ResNet50 embedding).

source \rightarrow target	ResNet	DANN	MDD	MCC	SDAT	ELS	Proposed
A \rightarrow W	85.2 \pm 0.2	86.3 \pm 0.3	86.4 \pm 0.4	88.3 \pm 0.3	88.6 \pm 0.4	87.2 \pm 0.3	90.4 \pm 0.4
D \rightarrow W	97.5 \pm 0.1	97.4 \pm 0.3	97.7 \pm 0.1	96.9 \pm 0.1	97.6 \pm 0.1	97.7 \pm 0.1	97.8 \pm 0.2
W \rightarrow D	99.5 \pm 0.3	98.7 \pm 0.2	99.7 \pm 0.1	97.4 \pm 0.2	99.1 \pm 0.2	99.3 \pm 0.2	99.5 \pm 0.3
A \rightarrow D	89.4 \pm 0.2	84.3 \pm 0.4	89.9 \pm 0.2	87.4 \pm 0.5	86.3 \pm 0.4	87.1 \pm 0.2	90.1 \pm 0.3
D \rightarrow A	71.4 \pm 0.3	71.7 \pm 0.4	70.6 \pm 0.3	74.9 \pm 0.4	72.3 \pm 0.4	71.6 \pm 0.3	71.9 \pm 0.1
W \rightarrow A	73.1 \pm 0.2	73.5 \pm 0.2	72.3 \pm 0.4	73.0 \pm 0.4	73.6 \pm 0.3	73.7 \pm 0.3	74.6 \pm 0.1
Average	86.0 \pm 0.2	85.3 \pm 0.3	86.1 \pm 0.2	86.3 \pm 0.3	86.2 \pm 0.3	86.1 \pm 0.2	87.3 \pm 0.2

6 Numerical Validation

More Synthetic-Data Validation. We first validate our proposed method on synthetic data that follows our model; see Appendix F for details.

Application (i) - Domain Adaptation. We first test the proposed methods over a number of domain adaptation (DA) tasks. In DA, we have the source domain data $\{x^{(1)}\}$ and the target domain $\{x^{(2)}\}$, respectively. Only the source domain data have labels and the two domains are unaligned. We hope to use our method to find shared representations of source and target, and thus the classifier trained using source data can also work well on the target data.

Dataset: We use two standard benchmarks of DA, i.e., *Office-31* [58] and *Office-Home* [59]. The *Office-31* dataset has 4652 images and 31 categories from three domains, namely, Amazon images (A), Webcam images (W) and DSLR images (D). The *Office-Home* dataset contains 15,500 images with 65 object classes from four domains, i.e., Artistic images (Ar), Clip art images (Cl), Product images (Pr), and Real-world images (Rw).

Setup: We first test the homogeneous domain model in Sec. 4. The images are pre-processed using a ResNet50-based image encoder pre-trained over ImageNet1k [42]. As mentioned, it was observed that self-supervised representation encoders find embeddings that are linearly separable [42], which justifies the use of the model $x^{(q)} \approx Az^{(q)}$ in the embedding domain. After pre-processing, each image is represented by $d^{(q)} = 2048$ features for $q = 1, 2$. We set $d_C = 256$ for *Office-31* and $d_C = 512$ for *Office-Home*. More detailed settings are in Appendix G.

Baselines and Training Setup: The baselines are representative DA methods, namely, DANN [25], MDD [60], MCC [61], SDAT [62], and ELS [63]. All the baselines use the same encoder-produced embeddings as inputs; see Appendix G.1 for their configurations. We also use ResNet encoder’s outputs as an extra baseline as it learns informative and transferable features from the ImageNet-1K dataset. We follow the training strategies adopted by the baselines [25, 60, 62] to learn a classifier jointly with the shared latent components. This strategy arguably regularizes towards more classification-friendly geometry of the shared features. Therefore we append a cross-entropy (CE) based classifier training module to our loss in (7) that learns our feature extractor Q . More details are in Appendix G.1.

Metric: The evaluation metric is the classification accuracy in the target domain $\{x^{(2)}\}$. The classifier is trained with the projected source domain $\hat{Q}x^{(1)}$ and the associated labels.

Result: Table 1 and Table 2 show the classification accuracy (mean \pm std) on *Office-31* and *Office-Home*, respectively. The results are averaged over 5 runs. One can observe that the proposed method offers the best and second best performance in most of the cases. In some tasks (e.g., “A \rightarrow W”, “Ar \rightarrow Cl”, “Ar \rightarrow Pr” and “Rw \rightarrow Cl”), the proposed method outperforms the best-performing baselines by at least 2% in accuracy.

More results on the DA task can be found in Appendix G.1.

Application (ii) - Single Cell Sequence Analysis. In biomedical research, it is desired to fuse measurements from multiple sensorial modalities of the same cells, in order to have better characterizations of the cells. However, obtaining multimodal data of the same cells simultaneously is almost impossible, due to the sensing limitations. Therefore, many methods are proposed in the literature for aligning unpaired multi-modal single cell data [27, 64, 65]. We focus on the following two modalities of single-cell data [66]: (1) the RNA sequences $\{x^{(1)}\}$ and (2) the ATAC sequences $\{x^{(2)}\}$.

Table 2: Classification accuracy on the target domain of *office-Home* dataset (ResNet50 embedding).

source \rightarrow target	ResNet	DANN	MDD	MCC	SDAT	ELS	Proposed
Ar \rightarrow Cl	42.0 \pm 0.2	46.7 \pm 0.2	47.4 \pm 0.3	44.4 \pm 0.3	47.3 \pm 0.4	48.5 \pm 0.2	51.0 \pm 0.3
Ar \rightarrow Pr	69.2 \pm 0.1	70.2 \pm 0.4	72.8 \pm 0.4	72.4 \pm 0.2	71.1 \pm 0.3	71.0 \pm 0.3	75.8 \pm 0.1
Ar \rightarrow Rw	80.2 \pm 0.3	81.2 \pm 0.4	81.2 \pm 0.1	80.3 \pm 0.3	80.5 \pm 0.1	80.8 \pm 0.4	82.5 \pm 0.2
Cl \rightarrow Ar	60.7 \pm 0.4	60.8 \pm 0.3	62.4 \pm 0.1	59.2 \pm 0.4	57.6 \pm 0.2	59.8 \pm 0.1	62.7 \pm 0.4
Cl \rightarrow Pr	71.0 \pm 0.1	69.8 \pm 0.3	70.0 \pm 0.4	71.1 \pm 0.4	66.5 \pm 0.1	68.5 \pm 0.2	72.5 \pm 0.3
Cl \rightarrow Rw	74.8 \pm 0.2	73.3 \pm 0.1	74.1 \pm 0.1	76.2 \pm 0.2	70.7 \pm 0.1	71.7 \pm 0.1	75.8 \pm 0.1
Pr \rightarrow Ar	60.6 \pm 0.2	62.2 \pm 0.1	64.3 \pm 0.1	59.2 \pm 0.1	62.5 \pm 0.4	60.9 \pm 0.2	64.4 \pm 0.3
Pr \rightarrow Cl	44.8 \pm 0.1	48.8 \pm 0.1	48.0 \pm 0.3	46.2 \pm 0.2	49.0 \pm 0.3	49.6 \pm 0.3	50.4 \pm 0.1
Pr \rightarrow Rw	79.6 \pm 0.1	80.3 \pm 0.4	79.6 \pm 0.3	80.3 \pm 0.2	80.0 \pm 0.1	79.2 \pm 0.1	81.7 \pm 0.2
Rw \rightarrow Ar	70.1 \pm 0.2	71.5 \pm 0.1	71.4 \pm 0.3	67.8 \pm 0.2	71.6 \pm 0.4	71.3 \pm 0.4	72.6 \pm 0.1
Rw \rightarrow Cl	45.8 \pm 0.2	50.9 \pm 0.2	50.3 \pm 0.1	50.0 \pm 0.2	51.4 \pm 0.1	50.7 \pm 0.1	53.2 \pm 0.1
Rw \rightarrow Pr	80.7 \pm 0.1	80.6 \pm 0.4	81.1 \pm 0.1	81.2 \pm 0.1	80.7 \pm 0.1	79.8 \pm 0.3	82.9 \pm 0.3
Average	64.9 \pm 0.1	66.3 \pm 0.2	66.8 \pm 0.2	65.6 \pm 0.2	65.7 \pm 0.2	65.9 \pm 0.2	68.7 \pm 0.2

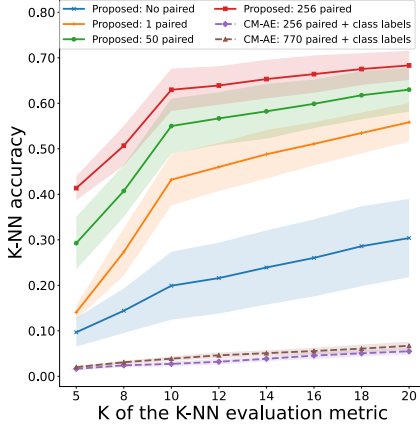
Figure 4: k -NN accuracy for single-cell sequence alignment.

Table 3: Average precision P@1 of cross-language information retrieval.

source \rightarrow target	Adv - NN	proposed - NN	Adv - CSLS	proposed - CSLS
en \rightarrow es	61.3	66.4	70.2	74.9
es \rightarrow en	55.4	65.3	67.6	75.6
en \rightarrow it	48.2	54.4	60.8	67.7
it \rightarrow en	55.2	51.9	63.8	66.0
en \rightarrow fr	63.6	60.2	72.6	73.7
fr \rightarrow en	55.4	58.4	64.1	71.4
en \rightarrow de	51.4	56.7	59.3	67.6
de \rightarrow en	42.5	57.0	51.0	59.3
en \rightarrow ru	32.7	34.9	38.6	41.4
ru \rightarrow en	27.6	41.6	35.0	50.8
en \rightarrow ar	12.6	22.7	16.7	29.1
ar \rightarrow en	15.7	26.9	20.1	35.6
en \rightarrow vi	2.1	10.4	7.7	22.8
vi \rightarrow en	2.7	17.3	4.4	33.0
Average	37.6	44.5	45.1	54.9

Dataset: We use human lung adenocarcinoma A549 cells data from [66]. The dataset contains 1,874 samples of RNA sequences $\{x^{(1)}\}$ and ATAC sequences $\{x^{(2)}\}$. Each data set is split into 1534 training samples and 340 testing samples as in [27]. The data have labeled associations between the two domains—part of which will be used to test our weakly supervised formulation. For this experiment, features of RNA sequence and the ATAC sequence have dimensions of $d^{(1)} = 815$ and $d^{(2)} = 2613$, respectively. We set $d_C = 256$. We use our weakly supervised formulation as shown in (9). We uniformly sampled a set of indices from the training set to serve as \mathcal{L} .

Baseline and Metric: We use weakly supervised algorithm, namely, cross-modal autoencoder (CM-AE) work in [27], as a baseline, which also learns the shared representation between unaligned RNA and ATAC sequences. We use the K -nearest neighbor (k -NN) accuracy to evaluate the performance as suggested in [27].

Result: The plot in Fig. 4 shows the k -NN accuracy of the methods on the test set. Results show the mean and standard deviation over 10 runs, each having a different random initialization. For the proposed method, we vary the number of available paired samples from 0 (cf. Theorem 1) to $d_C = 256$ (cf. Theorem 3). Note that the baseline uses more (i.e., 256 and 770) paired samples. It also needs additional class labels, i.e., $y_i^{(q)}$ for the i th sample $x_i^{(q)}$. Here, $y_i^{(q)}$ represents the number of hours (0, 1 or 3) of cell treatment [27, 66]. The proposed method without any supervision (i.e., 0 paired samples) already exhibits around 3 times greater k -NN accuracy compared to the baseline for all k . Moreover, including just one paired sample boosts the k -NN accuracy of the proposed method to around 5 times higher than the baseline for all k . Finally, one can observe a steadily increasing k -NN accuracy with respect to the number of available paired samples. This corroborates with our Theorem 3.

Application (iii) - Multi-lingual Information Retrieval. We also evaluated our method on a word embedding association problem from the natural language processing literature [20, 21]. This task aims to associate high-dimensional word embeddings across different languages according to their semantic meaning. The word embeddings in two languages are represented using two sets of vectors, i.e., $\{\mathbf{x}_i^{(1)}\}_{i=1}^I$ and $\{\mathbf{x}_j^{(2)}\}_{j=1}^J$. The postulate is that if $\mathbf{x}_i^{(1)}$ and $\mathbf{x}_j^{(2)}$ have the same meaning (e.g., both representing “cat”) in two languages (e.g., English and German), they should share a latent components \mathbf{c} .

Dataset: We use the word embeddings from the MUSE dataset (<https://github.com/facebookresearch/MUSE>) [21]. These monolingual word embeddings are generated using fastText [67] and has dimensions of $d^{(q)} = 300$ for $q = 1, 2$. The training dataset include 200,000 word embeddings in each language. In our experiment we set $d_C = 256$. We follow the generative model under [8] and run the formulation in [7] to learn the linear transformation \mathbf{Q} .

Baseline: We use Adv [21] as the baseline which also uses distribution matching between two language domains. Unlike our method, Adv does not use linear mixture models.

Metric: We follow [21] to use the average precision score calculated based on *nearest neighbor* (NN) and *cross domain similarity local scaling* (CSLS). Precision at k (“ k precision”) is computed by the number of times that one of the correct translations of source word is retrieved at top- k results ($k = \{1, 5, 10\}$). The final score is normalized to be in the range of 0 to 100, with 100 being the highest score indicating the best performance. To evaluate the performance, we use the same test data as in [21]. For each source and target language pair, this dataset includes 1,500 source word embeddings. The source embeddings are used to retrieve corresponding embeddings from a pool of 200,000 target word embeddings.

Result: Table 3 reports the P@1 scores over the test data calculated for each source and target language pair. The languages are denoted as **en** - English, **es** - Spanish, **it** - Italian, **fr** - French, **de** - Germany, **ru** - Russian, **ar** - Arabic and **vi** - Vietnamese. One can observe that the proposed method exhibits a better precision performance than that of Adv in most of the translation tasks. In particular, the proposed method significantly outperforms the baseline on the tasks **en**→**ar**, **ar**→**en**, **en**→**vi** and **vi**→**en**, showing at least 10% precision gains. Similarly, our method shows at least 5% improvements in both NN and CSLS based precision metrics in **en**→**es** and **es**→**en** tasks.

More details and additional experiments can be found in Appendix G.3.

7 Conclusion

In this work, we considered the problem of identifying shared components from unaligned multi-domain mixtures. We proposed a learning loss that matches the distributions of linearly transformed data. Based on this loss, we came up with a suite of sufficient conditions to ensure the identifiability of shared components. Furthermore, we proposed modified models and losses that enjoy more relaxed conditions for shared component identifiability. This was achieved via introducing structural constraints, namely, the homogeneity of the mixing systems and the existence of weak supervision. Our theoretical claims were validated with both synthetic and real-world data, demonstrating soundness of the theorems and usefulness of the models/algorithms.

Limitations. First, our conditions for shared component identification are sufficient. The necessary conditions are not underpinned, but necessary conditions assist understanding the limitations of the models and algorithms. Second, our methods were developed under the linear mixture model, which has limited expressiveness, and thus often requires pre-processing to approximately meet the model specification. We expect that results with similar flavors to be derived for nonlinear models in the future. Third, the results were derived under an unlimited data assumption. It would be interesting have a finite sample analysis. Finally, optimizing GAN-based losses is sensitive to hyperparameter settings. Back-propagation based minimax optimization occasionally fails to converge. More optimization-friendly losses and more stable algorithms are desirable in the context of distribution matching.

Acknowledgment

This work is supported in part by the Army Research Office (ARO) under Project ARO W911NF-21-1-0227, and in part by the National Science Foundation (NSF) CAREER Award ECCS-2144889.

References

- [1] M. S. Ibrahim, A. S. Zamzam, A. Konar, and N. D. Sidiropoulos, “Cell-edge detection via selective cooperation and generalized canonical correlation,” *IEEE Transactions on Wireless Communications*, vol. 20, no. 11, pp. 7431–7444, 2021.
- [2] M. Sørensen, C. I. Kanatsoulis, and N. D. Sidiropoulos, “Generalized canonical correlation analysis: A subspace intersection approach,” *IEEE Transactions on Signal Processing*, vol. 69, pp. 2452–2467, 2021.
- [3] P. Rastogi, B. Van Durme, and R. Arora, “Multiview LSA: Representation learning via generalized CCA,” in *Proc. Conference of the North American chapter of the Association for Computational Linguistics: human language technologies*, 2015, pp. 556–566.
- [4] Q. Lyu and X. Fu, “Nonlinear multiview analysis: Identifiability and neural network-assisted implementation,” *IEEE Transactions on Signal Processing*, vol. 68, pp. 2697–2712, 2020.
- [5] J. Von Kügelgen, Y. Sharma, L. Gresele, W. Brendel, B. Schölkopf, M. Besserve, and F. Locatello, “Self-supervised learning with data augmentations provably isolates content from style,” in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 34, 2021, pp. 16 451–16 467.
- [6] Q. Lyu, X. Fu, W. Wang, and S. Lu, “Understanding latent correlation-based multiview learning and self-supervision: An identifiability perspective,” in *Proc. International Conference on Learning Representations (ICLR)*, 2022.
- [7] I. Daunhawer, A. Bizeul, E. Palumbo, A. Marx, and J. E. Vogt, “Identifiability results for multimodal contrastive learning,” in *Proc. International Conference on Learning Representations (ICLR)*, 2022.
- [8] N. Sturma, C. Squires, M. Drton, and C. Uhler, “Unpaired multi-domain causal representation learning,” in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 34, 2021.
- [9] P. Comon and C. Jutten, *Handbook of Blind Source Separation: Independent component analysis and applications*. Academic press, 2010.
- [10] J. Hu and K. Huang, “Global identifiability of ℓ_1 -based dictionary learning via matrix volume optimization,” in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 36, 2024.
- [11] N. Gillis, “The why and how of nonnegative matrix factorization,” *Regularization, optimization, kernels, and support vector machines*, vol. 12, no. 257, pp. 257–291, 2014.
- [12] J. Sun, Q. Qu, and J. Wright, “Complete dictionary recovery over the sphere i: Overview and the geometric picture,” *IEEE Transactions on Information Theory*, vol. 63, no. 2, pp. 853–884, 2016.
- [13] A. T. Erdogan, “A class of bounded component analysis algorithms for the separation of both independent and dependent sources,” *IEEE Transactions on Signal Processing*, vol. 61, no. 22, pp. 5730–5743, 2013.
- [14] G. Tatli and A. T. Erdogan, “Polytopic matrix factorization: Determinant maximization based criterion and identifiability,” *IEEE Transactions on Signal Processing*, vol. 69, pp. 5431–5447, 2021.
- [15] F. R. Bach and M. I. Jordan, “A probabilistic interpretation of canonical correlation analysis,” 2005.
- [16] H. Richard, P. Ablin, B. Thirion, A. Gramfort, and A. Hyvarinen, “Shared independent component analysis for multi-subject neuroimaging,” in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 34, 2021, pp. 29 962–29 971.
- [17] Y.-O. Li, T. Adali, W. Wang, and V. D. Calhoun, “Joint blind source separation by multiset canonical correlation analysis,” *IEEE Transactions on Signal Processing*, vol. 57, no. 10, pp. 3918–3929, 2009.

- [18] M. Anderson, G.-S. Fu, R. Phlypo, and T. Adali, “Independent vector analysis: Identification conditions and performance bounds,” *IEEE Transactions on Signal Processing*, vol. 62, no. 17, pp. 4399–4410, 2014.
- [19] P. A. Karakasis and N. D. Sidiropoulos, “Revisiting deep generalized canonical correlation analysis,” *IEEE Transactions on Signal Processing*, vol. 71, pp. 4392–4406, 2023.
- [20] G. Lample, A. Conneau, L. Denoyer, and M. Ranzato, “Unsupervised machine translation using monolingual corpora only,” in *Proc. International Conference on Learning Representations (ICLR)*, 2018.
- [21] G. Lample, A. Conneau, M. Ranzato, L. Denoyer, and H. Jégou, “Word translation without parallel data,” in *Proc. International Conference on Learning Representations (ICLR)*, 2018.
- [22] E. Grave, A. Joulin, and Q. Berthet, “Unsupervised alignment of embeddings with wasserstein procrustes,” in *Proc. International Conference on Artificial Intelligence and Statistics (AISTATS)*. PMLR, 2019, pp. 1880–1890.
- [23] S. Ben-David, J. Blitzer, K. Crammer, and F. Pereira, “Analysis of representations for domain adaptation,” in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 19, 2006.
- [24] M. Long, Y. Cao, J. Wang, and M. Jordan, “Learning transferable features with deep adaptation networks,” in *Proc. International Conference on Machine Learning (ICML)*. PMLR, 2015, pp. 97–105.
- [25] Y. Ganin and V. Lempitsky, “Unsupervised domain adaptation by backpropagation,” in *Proc. International Conference on Machine Learning (ICML)*. PMLR, 2015, pp. 1180–1189.
- [26] S. Waaijenborg, P. C. Verselwele de Witt Hamer, and A. H. Zwinderman, “Quantifying the association between gene expressions and dna-markers by penalized canonical correlation analysis,” *Statistical applications in genetics and molecular biology*, vol. 7, no. 1, 2008.
- [27] K. D. Yang, A. Belyaeva, S. Venkatachalapathy, K. Damodaran, A. Katcoff, A. Radhakrishnan, G. Shivashankar, and C. Uhler, “Multi-domain translation between single-cell imaging and sequencing data using autoencoders,” *Nature communications*, vol. 12, no. 1, p. 31, 2021.
- [28] I. Gulrajani and T. Hashimoto, “Identifiability conditions for domain adaptation,” in *Proc. International Conference on Machine Learning (ICML)*, 2022, pp. 7982–7997.
- [29] S. Shrestha and X. Fu, “Towards identifiable unsupervised domain translation: A diversified distribution matching approach,” in *Proc. International Conference on Learning Representations (ICLR)*, 2024.
- [30] Y. Choi, Y. Uh, J. Yoo, and J.-W. Ha, “StarGAN v2: Diverse image synthesis for multiple domains,” in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 8188–8197.
- [31] X. Huang, M.-Y. Liu, S. Belongie, and J. Kautz, “Multimodal unsupervised image-to-image translation,” in *Proc. European Conference on Computer Vision (ECCV)*, 2018, pp. 172–189.
- [32] S. Xie, L. Kong, M. Gong, and K. Zhang, “Multi-domain image generation and translation with identifiability guarantees,” in *Proc. International Conference on Learning Representations (ICLR)*, 2023.
- [33] P. Comon, “Independent component analysis, a new concept?” *Signal processing*, vol. 36, no. 3, pp. 287–314, 1994.
- [34] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [35] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, “Learning transferable visual models from natural language supervision,” in *Proc. International Conference on Machine Learning (ICML)*. PMLR, 2021, pp. 8748–8763.
- [36] Y. Choi, M. Choi, M. Kim, J.-W. Ha, S. Kim, and J. Choo, “StarGAN: Unified generative adversarial networks for multi-domain image-to-image translation,” in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 8789–8797.
- [37] X. Fu, K. Huang, N. D. Sidiropoulos, and W.-K. Ma, “Nonnegative matrix factorization for signal and data analytics: Identifiability, algorithms, and applications,” *IEEE Signal Processing Magazine*, vol. 36, no. 2, pp. 59–80, 2019.

- [38] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola, “A kernel two-sample test,” *The Journal of Machine Learning Research*, vol. 13, no. 1, pp. 723–773, 2012.
- [39] C. Villani *et al.*, *Optimal transport: old and new*. Springer, 2009, vol. 338.
- [40] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial networks,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2014.
- [41] L. Kong, S. Xie, W. Yao, Y. Zheng, G. Chen, P. Stojanov, V. Akinwande, and K. Zhang, “Partial disentanglement for domain adaptation,” in *Proc. International Conference on Machine Learning (ICML)*, 2022, pp. 11 455–11 472.
- [42] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, “A simple framework for contrastive learning of visual representations,” in *Proc. International Conference on Machine Learning (ICML)*. PMLR, 2020, pp. 1597–1607.
- [43] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra *et al.*, “Matching networks for one shot learning,” in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 29, 2016.
- [44] E. Triantafillou, T. Zhu, V. Dumoulin, P. Lamblin, U. Evci, K. Xu, R. Goroshin, C. Gelada, K. Swersky, P.-A. Manzagol *et al.*, “Meta-dataset: A dataset of datasets for learning to learn from few examples,” in *Proc. International Conference on Learning Representations (ICLR)*, 2019.
- [45] S. Wold, K. Esbensen, and P. Geladi, “Principal component analysis,” *Chemometrics and intelligent laboratory systems*, vol. 2, no. 1-3, pp. 37–52, 1987.
- [46] J. Li and J. M. Bioucas-Dias, “Minimum volume simplex analysis: A fast algorithm to unmix hyperspectral data,” in *IEEE International Geoscience and Remote Sensing Symposium*, vol. 3, 2008, pp. III–250.
- [47] X. Fu, K. Huang, B. Yang, W.-K. Ma, and N. D. Sidiropoulos, “Robust volume minimization-based matrix factorization for remote sensing and document clustering,” *IEEE Transactions on Signal Processing*, vol. 64, no. 23, pp. 6254–6268, 2016.
- [48] S. Arora, R. Ge, Y. Halpern, D. Mimno, A. Moitra, D. Sontag, Y. Wu, and M. Zhu, “A practical algorithm for topic modeling with provable guarantees,” in *Proc. International Conference on Machine Learning (ICML)*. PMLR, 2013, pp. 280–288.
- [49] K. Huang, X. Fu, and N. D. Sidiropoulos, “Anchor-free correlated topic modeling: Identifiability and algorithm,” *Advances in Neural Information Processing Systems*, vol. 29, 2016.
- [50] X. Mao, P. Sarkar, and D. Chakrabarti, “Estimating mixed memberships with sharp eigenvector deviations,” *Journal of the American Statistical Association*, vol. 116, no. 536, pp. 1928–1940, 2021.
- [51] L. Gresele, P. K. Rubenstein, A. Mehrjou, F. Locatello, and B. Schölkopf, “The incomplete Rosetta stone problem: Identifiability results for multi-view nonlinear ICA,” in *Proc. Uncertainty in Artificial Intelligence*, 2020, pp. 217–227.
- [52] C. Eastwood, J. von Kügelgen, L. Ericsson, D. Bouchacourt, P. Vincent, M. Ibrahim, and B. Schölkopf, “Self-supervised disentanglement by leveraging structure in data augmentations,” in *Causal Representation Learning Workshop at NeurIPS 2023*, 2023.
- [53] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, “Unpaired image-to-image translation using cycle-consistent adversarial networks,” in *Proc. the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 2223–2232.
- [54] M. Long, H. Zhu, J. Wang, and M. I. Jordan, “Deep transfer learning with joint adaptation networks,” in *Proc. International Conference on Machine Learning (ICML)*. PMLR, 2017, pp. 2208–2217.
- [55] W. Wang, H. Zhang, Z. Yuan, and C. Wang, “Unsupervised real-world super-resolution: A domain adaptation perspective,” in *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2021, pp. 4318–4327.
- [56] T. Galanti, L. Wolf, and S. Benaïm, “The role of minimal complexity functions in unsupervised learning of semantic mappings,” in *Proc. International Conference on Learning Representations (ICLR)*, 2018.

- [57] N. Moriaikov, J. Adler, and J. Teuwen, “Kernel of CycleGAN as a principle homogeneous space,” in *Proc. International Conference on Learning Representations (ICLR)*, 2020.
- [58] K. Saenko, B. Kulis, M. Fritz, and T. Darrell, “Adapting visual category models to new domains,” in *Proc. European Conference on Computer Vision (ECCV)*. Springer, 2010, pp. 213–226.
- [59] H. Venkateswara, J. Eusebio, S. Chakraborty, and S. Panchanathan, “Deep hashing network for unsupervised domain adaptation,” in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 5018–5027.
- [60] Y. Zhang, T. Liu, M. Long, and M. Jordan, “Bridging theory and algorithm for domain adaptation,” in *Proc. International Conference on Machine Learning (ICML)*, 2019.
- [61] Y. Jin, X. Wang, M. Long, and J. Wang, “Minimum class confusion for versatile domain adaptation,” in *Proc. IEEE International Conference on Computer Vision (ICCV)*. Springer, 2020, pp. 464–480.
- [62] H. Rangwani, S. K. Aithal, M. Mishra, A. Jain, and V. B. Radhakrishnan, “A closer look at smoothness in domain adversarial training,” in *Proc. International Conference on Machine Learning (ICML)*. PMLR, 2022, pp. 18 378–18 399.
- [63] Y. Zhang, J. Liang, Z. Zhang, L. Wang, R. Jin, T. Tan *et al.*, “Free lunch for domain adversarial training: Environment label smoothing,” in *Proc. International Conference on Learning Representations (ICLR)*, 2023.
- [64] L. Eyring, D. Klein, T. Uscidda, G. Palla, N. Kilbertus, Z. Akata, and F. Theis, “Unbalancedness in neural monge maps improves unpaired domain translation,” *arXiv preprint arXiv:2311.15100*, 2023.
- [65] M. Amodio and S. Krishnaswamy, “MAGAN: Aligning biological manifolds,” in *Proc. International Conference on Machine Learning (ICML)*. PMLR, 2018, pp. 215–223.
- [66] J. Cao, D. A. Cusanovich, V. Ramani, D. Aghamirzaie, H. A. Pliner, A. J. Hill, R. M. Daza, J. L. McFaline-Figueroa, J. S. Packer, L. Christiansen *et al.*, “Joint profiling of chromatin accessibility and gene expression in thousands of single cells,” *Science*, vol. 361, no. 6409, pp. 1380–1385, 2018.
- [67] A. Joulin, E. Grave, P. Bojanowski, M. Douze, H. Jégou, and T. Mikolov, “Fasttext.zip: Compressing text classification models,” *arXiv preprint arXiv:1612.03651*, 2016.
- [68] K. B. Athreya and S. N. Lahiri, *Measure theory and probability theory*. Springer, 2006, vol. 19.
- [69] J.-F. Cardoso, “Super-symmetric decomposition of the fourth-order cumulant tensor. blind identification of more sources than sensors,” in *Proc. International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, vol. 91, 1991, pp. 3109–3112.
- [70] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *Proc. ICLR*, 2015.
- [71] A. L. Maas, A. Y. Hannun, and A. Y. Ng, “Rectifier nonlinearities improve neural network acoustic models,” in *Proc. International Conference on Machine Learning (ICML)*, vol. 30, no. 1, 2013, p. 3.
- [72] K. Ghasedi Dizaji, A. Herandi, C. Deng, W. Cai, and H. Huang, “Deep clustering via joint convolutional autoencoder embedding and relative entropy minimization,” in *Proc. the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 5736–5745.
- [73] J. Xie, R. Girshick, and A. Farhadi, “Unsupervised deep embedding for clustering analysis,” in *Proc. International Conference on Machine Learning (ICML)*. PMLR, 2016, pp. 478–487.
- [74] W. Wang, X. Yan, H. Lee, and K. Livescu, “Deep variational canonical correlation analysis,” *arXiv preprint arXiv:1610.03454*, 2016.
- [75] A. Gretton, K. Fukumizu, C. Teo, L. Song, B. Schölkopf, and A. Smola, “A kernel statistical test of independence,” in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 20, 2007.
- [76] T. M. Cover, *Elements of information theory*. John Wiley & Sons, 1999.

Supplementary Material of “Identifiable Shared Component Analysis of Unpaired Multimodal Mixtures”

A Notation

The notations used throughout the paper are summarized in the Table 4:

Table 4: Definition of notations.

Notation	Definition
$x, \mathbf{x}, \mathbf{X}$	scalar, vector and matrix
$\mathbf{x}^{(q)}$	variable from q-th domain
x_i, \mathbf{x}_i	both represents i-th element of vector \mathbf{x}
\mathbf{X}_{ij}	represents the element of i-th row and j-th column of matrix \mathbf{X}
$\mathbf{x}^\top, \mathbf{X}^\top$	transpose of \mathbf{x} and \mathbf{X}
$ \mathcal{X} $	represents the cardinality of set \mathcal{X}
$\text{Null}(\mathbf{X})$	represents the null space of matrix \mathbf{X}
$\text{conv}(\cdot)$	returns the convex hull of the given set
$\dim(\mathcal{X})$	denotes the dimension of subspace \mathcal{X}
$k\mathcal{A}$	$\{k\mathbf{a} \mid \mathbf{a} \in \mathcal{A}, k \in \mathbb{R}\}$
$\mathbf{x} + \mathcal{X}$	$\{\mathbf{x} + \mathbf{z} \mid \mathbf{z} \in \mathcal{X}\}$
$\mathcal{X} + \mathcal{Y}$	$\{\mathbf{x} + \mathbf{y} \mid \mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y}\}$
$\mathbf{A}_{\text{PreImg}}(\mathcal{X})$	preimage of \mathcal{X} ; $\{\mathbf{x} \mid \mathbf{A}\mathbf{x} \in \mathcal{X}\}$
$[N]$	set of whole numbers up to N ; $\{1 \dots N\}$
\mathbf{I}	identity matrix
$\mathbb{P}_{\mathbf{x}}$	probability distribution of random variable \mathbf{x}
$\mathbb{P}_{\mathbf{x}, \mathbf{y}}$	joint probability distribution of random variable \mathbf{x} and \mathbf{y}
$\mathbb{E}[\cdot]$	expectation
$\mathbf{x} \stackrel{(d)}{=} \mathbf{y}$	\mathbf{x} and \mathbf{y} random vectors have the same distribution
$\mathbf{x} \stackrel{(d)}{\neq} \mathbf{y}$	\mathbf{x} and \mathbf{y} random vectors have different distributions
$\mathbf{x} \perp\!\!\!\perp \mathbf{y}$	\mathbf{x} and \mathbf{y} random vectors are statistically independent
$[a, b]$	represents continuous interval between a and b
$\mathcal{N}(\mu, \sigma^2)$	normal distribution with mean μ and variance σ^2
$\text{Uniform}[a, b]$	uniform distribution with interval a and b
$\text{Gamma}(\alpha, \theta)$	gamma distribution with the shape parameter α and scale parameter θ
$\text{Laplace}(\mu, b)$	Laplace distribution with location μ and diversity or scale parameter b
$\text{VonMises}(\mu, \kappa)$	von Mises distribution with location μ and κ concentration parameter.
$\text{Beta}(\alpha, \beta)$	beta distribution with the shape parameters α and β

B Proof of Theorem 1

We restate the theorem here:

Theorem 1 Under Assumption 1 and the generative model in (1), denote any solution of (6) as $\hat{\mathbf{Q}}^{(q)}$ $q = 1, 2$. Then, if the mixing matrices $\mathbf{A}^{(q)}$ are full column ranks and $\mathbb{E}[\mathbf{c}\mathbf{c}^\top]$ is full rank, we have $\hat{\mathbf{Q}}^{(q)} \mathbf{x}^{(q)} = \Theta^{(q)} \mathbf{c}$. In addition, assume that either of the following is satisfied:

- (a) The individual elements of the content components are statistically independent and non-Gaussian. In addition, $c_i \stackrel{(d)}{\neq} kc_j, \forall i \neq j, \forall k \in \mathbb{R}$ and $c_i \stackrel{(d)}{\neq} -c_i, \forall i$, i.e., the marginal distributions of the content elements cannot be matched with each other by mere scaling.
- (b) The support \mathcal{C} is a hyper-rectangle, i.e., $\mathcal{C} = [-a_1, a_1] \times \cdots \times [-a_{d_C}, a_{d_C}]$. Further, suppose that $c_i \stackrel{(d)}{\neq} kc_j, \forall i \neq j, \forall k \in \mathbb{R}$ and $c_i \stackrel{(d)}{\neq} -c_i, \forall i$.

Then, we have $\hat{\mathbf{Q}}^{(q)} \mathbf{x}^{(q)} = \Theta \mathbf{c}$, i.e., $\Theta^{(q)} = \Theta$ for all $q = 1, 2$, where $\Theta^{(q)}$.

We will prove the theorem in following two steps. For the first step we will prove $\hat{\mathbf{Q}}^{(q)} \mathbf{x}^{(q)} = \Theta^{(q)} \mathbf{c}$ and for second step we will employ either assumption (a) or (b) to prove that $\Theta^{(q)} = \Theta$, $\forall q = 1, 2$.

B.1 Linearly transformed content identification

Let us define

$$\mathbf{H}^{(q)} = \mathbf{Q}^{(q)} \mathbf{A}^{(q)} \in \mathbb{R}^{d_C \times (d_C + d_P^{(q)})}.$$

We want to show that

$$\text{Null}(\mathbf{H}^{(q)}) = \mathbf{0} \times \mathbb{R}^{d_P^{(q)}}, \quad (10)$$

since this will imply that $\mathbf{H}^{(q)}$ does not depend upon the style component. Combined with the fact that $\text{rank}(\mathbf{H}^{(q)}) = d_C$, this will imply that $\mathbf{H}^{(q)}$ is an invertible function of the content component. To that end, consider the following line of arguments.

Since the objective in (6) matches the distribution for latent random variables $\hat{\mathbf{c}}^{(1)} = \mathbf{Q}^{(1)} \mathbf{x}^{(1)}$ and $\hat{\mathbf{c}}^{(2)} = \mathbf{Q}^{(2)} \mathbf{x}^{(2)}$, the following holds for any $\mathcal{R}_c \subseteq \mathbb{R}^{d_C}, \forall k \in \mathbb{R}$,

$$\begin{aligned} \mathbb{P}_{\hat{\mathbf{c}}^{(1)}}[k\mathcal{R}_c] &= \mathbb{P}_{\hat{\mathbf{c}}^{(2)}}[k\mathcal{R}_c], \\ \stackrel{(a)}{\iff} \mathbb{P}_{\mathbf{z}^{(1)}}[\mathbf{H}_{\text{PreImg}}^{(1)}(k\mathcal{R}_c)] &= \mathbb{P}_{\mathbf{z}^{(2)}}[\mathbf{H}_{\text{PreImg}}^{(2)}(k\mathcal{R}_c)] \\ \stackrel{(b)}{\iff} \mathbb{P}_{\mathbf{z}^{(1)}}[k\mathbf{H}_{\text{PreImg}}^{(1)}(\mathcal{R}_c)] &= \mathbb{P}_{\mathbf{z}^{(2)}}[k\mathbf{H}_{\text{PreImg}}^{(2)}(\mathcal{R}_c)], \end{aligned} \quad (11)$$

where, $\mathbf{H}_{\text{PreImg}}^{(q)}(\mathcal{R}_c) := \{\mathbf{z}^{(q)} \mid \mathbf{H}^{(q)} \mathbf{z}^{(q)} \in \mathcal{R}_c\}$ is the pre-image of $\mathbf{H}^{(q)}$. (a) follows because $\mathbb{P}_{\hat{\mathbf{c}}^{(q)}}[k\mathcal{R}_c] = \mathbb{P}_{\mathbf{H}^{(q)} \mathbf{z}^{(q)}}[k\mathcal{R}_c] = \mathbb{P}_{\mathbf{z}^{(q)}}[\mathbf{H}_{\text{PreImg}}^{(q)}(k\mathcal{R}_c)]$ [68, Section 2.2]. (b) follows because $\mathbf{H}^{(q)}$ is a linear operator.

Although (11) holds for any \mathcal{R}_c , we will see that it is sufficient to consider a special \mathcal{R}_c to prove (10). To that end, take $\mathcal{R}_c = \text{conv}\{\mathbf{0}, \mathbf{a}_1, \dots, \mathbf{a}_{d_C}\}$, where $\mathbf{a}_i \in \mathbb{R}^{d_C}$ such that $\mathbb{P}_{\hat{\mathbf{c}}^{(q)}}[\mathcal{R}_c] > 0$. Let us take $\mathbf{y}_i^{(q)} \in \mathbb{R}^{d_C + d_P^{(q)}}$, such that $\mathbf{H}^{(q)} \mathbf{y}_i^{(q)} = \mathbf{a}_i$. For reasons that will be clear later, we hope to show that

$$\mathbf{H}_{\text{PreImg}}^{(q)}(\mathcal{R}_c) = \text{conv}\{\mathbf{0}, \mathbf{y}_1^{(q)}, \dots, \mathbf{y}_{d_C}^{(q)}\} + \text{Null}(\mathbf{H}^{(q)}).$$

To that end, observe that for any $\mathbf{r} \in \mathcal{R}_c$, we can represent \mathbf{r} as,

$$\mathbf{r} = \frac{1}{d_C + 1} \sum_{i=1}^{d_C} w_i \mathbf{a}_i, \text{ for some } \{w_i\}_{i=1}^{d_C} \text{ s.t. } \sum_{i=1}^{d_C} w_i \leq 1, \forall i.$$

For both view $q = 1, 2$, we get,

$$\begin{aligned} \mathbf{r} &= \frac{1}{d_C + 1} \sum_{i=1}^{d_C} w_i \mathbf{H}^{(q)} \mathbf{y}_i^{(q)} \\ \implies \mathbf{r} &= \mathbf{H}^{(q)} \left(\frac{1}{d_C + 1} \sum_{i=1}^{d_C} w_i \mathbf{y}_i^{(q)} \right) \\ \mathbf{H}_{\text{PreImg}}^{(q)} \left(\frac{1}{d_C + 1} \sum_{i=1}^{d_C} w_i \mathbf{a}_i \right) &= \frac{1}{d_C + 1} \sum_{i=1}^{d_C} w_i \mathbf{y}_i^{(q)} + \text{Null}(\mathbf{H}^{(q)}) \end{aligned} \quad (12)$$

We can write,

$$\mathbf{H}_{\text{PreImg}}^{(q)}(\mathcal{R}_c) = \text{conv}\{\mathbf{0}, \mathbf{y}_1^{(q)}, \dots, \mathbf{y}_{d_C}^{(q)}\} + \text{Null}(\mathbf{H}^{(q)}) \quad (13)$$

We have that $\text{Null}(\mathbf{H}^{(q)}) \subset \mathbb{R}^{d_C + d_P^{(q)}}$ is a linear subspace with $\dim(\text{Null}(\mathbf{H}^{(q)})) = d_P^{(q)}$. Let $\mathcal{A}^{(q)} = \mathbf{H}_{\text{PreImg}}^{(q)}(\mathcal{R}_c)$. Note that $\mathbb{P}_{\mathbf{z}^{(1)}}[k\mathcal{A}^{(1)}] = \mathbb{P}_{\mathbf{z}^{(2)}}[k\mathcal{A}^{(2)}]$, $\forall k \in \mathbb{R}$ (from [\(11\)](#)), and $\mathbb{P}_{\mathbf{z}^{(q)}}[\mathcal{A}^{(q)}] > 0$ (by the construction of \mathcal{R}_c). Further, the set $\mathcal{A}^{(q)}$ is of the form

$$\text{conv}\{\mathbf{0}, \mathbf{y}_1^{(q)}, \dots, \mathbf{y}_{d_C}^{(q)}\} + \mathcal{P}^{(q)},$$

because $\text{Null}(\mathbf{H}^{(q)})$ is a subspace of dimension $d_P^{(q)}$, hence it satisfies the definition of $\mathcal{P}^{(q)}$. Hence, Assumption [1](#) implies that

$$\text{Null}(\mathbf{H}^{(q)}) = \mathbf{0} \times \mathbb{R}^{d_P^{(q)}}.$$

Denoting the N th to M th columns of $\mathbf{H}^{(q)}$ by $\mathbf{H}^{(q)}(N : M)$, the above is equivalent to saying

$$\mathbf{H}^{(q)}(d_C + 1 : d_C + d_P^{(q)}) = \mathbf{0}. \quad (14)$$

Denote,

$$\mathbf{\Theta}^{(q)} = \mathbf{H}^{(q)}(1 : d_C) \quad \forall q = 1, 2.$$

Then, we can write,

$$\mathbf{Q}^{(q)} \mathbf{x}^{(q)} = \mathbf{\Theta}^{(q)} \mathbf{c}, \quad \forall q = 1, 2. \quad (15)$$

Next, we use Assumption (a) or (b) to show that $\mathbf{\Theta}^{(1)} = \mathbf{\Theta}^{(2)} = \mathbf{\Theta}$. To that end, note that the distribution matching constraint implies that

$$\begin{aligned} \mathbf{\Theta}^{(1)} \mathbf{c} &\stackrel{(d)}{=} \mathbf{\Theta}^{(2)} \mathbf{c} \\ \implies \mathbf{c} &\stackrel{(d)}{=} (\mathbf{\Theta}^{(1)})^{-1} \mathbf{\Theta}^{(2)} \mathbf{c}. \end{aligned}$$

Hence $\mathbf{M} = (\mathbf{\Theta}^{(1)})^{-1} \mathbf{\Theta}^{(2)}$ is an invertible matrix such that $\mathbf{c} \stackrel{(d)}{=} \mathbf{M} \mathbf{c}$. However, in the following, we will show that if either Assumption (a) or (b) is satisfied, then $\mathbf{M} = \mathbf{I}$, and thus $\mathbf{\Theta}^{(1)} = \mathbf{\Theta}^{(2)}$.

B.2 Considering Assumption (a)

We want to show that when Assumption (a) is satisfied, if $\mathbf{M} \mathbf{c} \stackrel{(d)}{=} \mathbf{c}$ for any invertible \mathbf{M} , then $\mathbf{M} = \mathbf{I}$.

Note that $\mathbf{M} \mathbf{c} = [\mathbf{m}_1 \dots \mathbf{m}_{d_C}] \begin{bmatrix} c_1 \\ \vdots \\ c_{d_C} \end{bmatrix}$. By Assumption (a), we have that the components of content are statistically independent $c_i \perp\!\!\!\perp c_j$, $i \neq j$, non-Gaussian, and has non-zero kurtosis. Then, according to cumulant multilinearity and additivity properties, the fourth order cumulant tensor $\text{Cum}(\mathbf{M} \mathbf{c})$ of $\mathbf{M} \mathbf{c}$ has the following unique decomposition [\[69\]](#),

$$\text{Cum}(\mathbf{M} \mathbf{c}) = \sum_{i=1}^{d_C} \kappa_i \mathbf{m}_i \circ \mathbf{m}_i \circ \mathbf{m}_i \circ \mathbf{m}_i \quad (16)$$

where \circ is the outer product, κ_i is the kurtosis of component c_i , and \mathbf{m}_i , $i \in [d_C]$ are the columns of \mathbf{M} .

Since $\mathbf{M}\mathbf{c} \stackrel{(d)}{=} \mathbf{c}$, the following should hold

$$\text{Cum}(\mathbf{M}\mathbf{c}) = \text{Cum}(\mathbf{c}) = \text{Cum}(\mathbf{I}\mathbf{c}) \quad (17)$$

$$\implies \sum_{d=1}^{d_C} \kappa_d \mathbf{m}_d \circ \mathbf{m}_d \circ \mathbf{m}_d \circ \mathbf{m}_d = \sum_{d=1}^{d_C} \kappa_d \mathbf{e}_d \circ \mathbf{e}_d \circ \mathbf{e}_d \circ \mathbf{e}_d, \quad (18)$$

\mathbf{e}_i is the i th column of identity matrix \mathbf{I} .

Because of statistical independence of components of \mathbf{c} , the CP-decomposition of $\text{Cum}(\mathbf{M}\mathbf{c}) = \text{Cum}(\mathbf{I}\mathbf{c})$ is unique [69] upto permutation and scaling ambiguities, i.e., \mathbf{M} should be a permutation scaling matrix.

Let $\mathbf{M} = \mathbf{\Pi}\mathbf{\Sigma}$ where, $\mathbf{\Pi} \in \mathbb{R}^{d_C \times d_C}$ is a permutation matrix and $\mathbf{\Sigma} = \text{Diag}(r_1, \dots, r_{d_C}) \in \mathbb{R}^{d_C \times d_C}$ is a diagonal scaling matrix.

Finally, since $c_i \stackrel{(d)}{\neq} kc_j, \forall i \neq j, \forall k \in \mathbb{R}$, \mathbf{M} has to be identity matrix. To see the reason, for the sake of contradiction, suppose that either (i) there exist $i, j \in [d_C] \times [d_C]$ and $k \in \mathbb{R}$, with $i \neq j$ such that $[\mathbf{M}\mathbf{c}]_i = kc_j$, or (ii) $\exists i \in [d_C]$ such that $[\mathbf{M}\mathbf{c}]_i = kc_i$ for some $k \in \mathbb{R}, k \neq 1$.

For case (i), since $\mathbf{M}\mathbf{c} \stackrel{(d)}{=} \mathbf{c}$, $[\mathbf{M}\mathbf{c}]_i \stackrel{(d)}{=} c_i, \forall i$. Hence,

$$\begin{aligned} [\mathbf{M}\mathbf{c}]_i &= kc_j \\ \implies [\mathbf{M}\mathbf{c}]_i &\stackrel{(d)}{=} kc_j \\ \implies c_i &\stackrel{(d)}{=} kc_j, \end{aligned}$$

which is a contradiction to the assumption $c_i \stackrel{(d)}{\neq} kc_j$.

For case (ii), $[\mathbf{M}\mathbf{c}]_i = kc_j$ implies that $c_i \stackrel{(d)}{=} kc_j$. First, $k \neq \pm 1$, cannot hold because it will mean that $\text{var}(c_i) = k^2 \text{var}(c_j)$ which cannot hold for $k \neq \pm 1$ since $\text{var}(c_i) > 0$. Hence, the only possible option is $k = -1$, which is already ruled out by the assumption that $c_i \stackrel{(d)}{\neq} -c_i$. Hence \mathbf{M} is an identity matrix. This concludes the proof.

B.3 Considering Assumption (b)

Let

$$\mathbf{e}_i = [0, 0, \dots, \underset{i\text{th location}}{1}, 0, 0] \in \mathbb{R}^{d_C}$$

denote the standard basis vector in \mathbb{R}^{d_C} . Let vertex of hyper-rectangle $\mathbf{v}_i = a_i \mathbf{e}_i = \mathbf{\Lambda} \mathbf{e}_i$, where

$$\mathbf{\Lambda} = \text{Diag}([a_1, \dots, a_{d_C}]^T),$$

where $\text{Diag}(\cdot)$ represents the diagonal matrix formed by the given vector.

If $\mathbf{M}\mathbf{c} \stackrel{(d)}{=} \mathbf{c}$, then the supports of $\mathbf{M}\mathbf{c}$ and \mathbf{c} should match, i.e.,

$$\mathbf{M}(\mathcal{C}) = \mathcal{C},$$

where $\mathbf{M}(\mathcal{C}) = \{\mathbf{M}\mathbf{c} \mid \mathbf{c} \in \mathcal{C}\}$.

Note that $\forall \mathbf{c} \in \mathcal{C}$, the set of points \mathbf{v}_i satisfy the following property

$$\begin{aligned} \mathbf{c} &= \sum_{i=1}^{d_C} \alpha_i \mathbf{v}_i, \quad \text{for some } -1 \leq \alpha_i \leq 1 \\ \implies \mathbf{M}\mathbf{c} &= \sum_{i=1}^{d_C} \alpha_i (\mathbf{M}\mathbf{v}_i). \end{aligned} \quad (19)$$

Since the support of $M\mathbf{c}$ is \mathcal{C} , this implies that $\forall \mathbf{c} \in \mathcal{C}$,

$$\mathbf{c} = \sum_{i=1}^{d_C} \alpha_i (M\mathbf{v}_i), \quad \text{for some } -1 \leq \alpha_i \leq 1$$

The last equation implies that the set of points $M\mathbf{v}_i, \forall i \in [d_C]$ also satisfy property (19). Hence, for each $i \in [d_C]$, $M\mathbf{v}_i = \pm \mathbf{v}_j$ for some unique $j \in [d_C]$. Note that j should be unique for each i because M is invertible, hence M cannot map two orthogonal vectors \mathbf{v}_i and $\mathbf{v}_k, i \neq k$, to the same vector $\pm \mathbf{v}_j$ with same or different signs.

Let $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_{d_C}]^T$. Then one can write

$$M\mathbf{V} = \mathbf{V}\Sigma\Pi,$$

where Σ is some diagonal matrix with diagonal entries from $\{+1, -1\}$ and Π is a permutation matrix. Then the above implies

$$\begin{aligned} M\Lambda\mathbf{I} &= \Lambda\mathbf{I}\Sigma\Pi \\ \implies M &= \Lambda\Sigma\Pi\Lambda^{-1}. \end{aligned}$$

Hence M is a permutation and scaling matrix.

Finally, by the same argument presented in last paragraph of Sec. B.2 (i.e., proof with Assumption (a)), we conclude that M is an identity matrix.

C Proof of Theorem 2

We restate the theorem here:

Theorem 2 Consider the mixture model in (8). Assume that $\text{rank}(\mathbf{A}) = d_C + d_P$ and $\text{rank}(\mathbb{E}[\mathbf{c}\mathbf{c}^T]) = d_C$, and that Assumption 2 holds. Denote $\hat{\mathbf{Q}}$ as any solution of (6) by constraining $\mathbf{Q} = \mathbf{Q}^{(1)} = \mathbf{Q}^{(2)}$. Then, we have $\hat{\mathbf{Q}}\mathbf{x}^{(q)} = \Theta\mathbf{c}$.

One can follow the same argument as in the step 1 of proof in B.

Let us define

$$\mathbf{H} = \mathbf{Q}\mathbf{A} \in \mathbb{R}^{d_C \times (d_C + d_P)}.$$

We want to show that

$$\text{Null}(\mathbf{H}) = \mathbf{0} \times \mathbb{R}^{d_P}, \quad (20)$$

since this will imply that \mathbf{H} does not depend upon the style component. Combined with the fact that $\text{rank}(\mathbf{H}) = d_C$, this will imply that \mathbf{H} is an invertible function of the content component. To that end, consider the following line of arguments.

Since the objective in (6) matches the distribution for latent random variables $\hat{\mathbf{c}}^{(1)} = \mathbf{Q}\mathbf{x}^{(1)}$ and $\hat{\mathbf{c}}^{(2)} = \mathbf{Q}\mathbf{x}^{(2)}$, the following holds for any $\mathcal{R}_c \subseteq \mathbb{R}^{d_C}, \exists k \in \mathbb{R}$

$$\mathbb{P}_{\hat{\mathbf{c}}^{(1)}}[k\mathcal{R}_c] = \mathbb{P}_{\hat{\mathbf{c}}^{(2)}}[k\mathcal{R}_c],$$

$$\stackrel{(a)}{\iff} \mathbb{P}_{\mathbf{z}^{(1)}}[\mathbf{H}_{\text{PreImg}}(k\mathcal{R}_c)] = \mathbb{P}_{\mathbf{z}^{(2)}}[\mathbf{H}_{\text{PreImg}}(k\mathcal{R}_c)] \quad (21)$$

$$\stackrel{(b)}{\iff} \mathbb{P}_{\mathbf{z}^{(1)}}[k\mathbf{H}_{\text{PreImg}}(\mathcal{R}_c)] = \mathbb{P}_{\mathbf{z}^{(2)}}[k\mathbf{H}_{\text{PreImg}}(\mathcal{R}_c)], \quad (22)$$

where, $\mathbf{H}_{\text{PreImg}}(\mathcal{R}_c) := \{\mathbf{z} \mid \mathbf{H}\mathbf{z} \in \mathcal{R}_c\}$ is the pre-image of \mathbf{H} . (a) follows because $\mathbb{P}_{\hat{\mathbf{c}}^{(q)}}[k\mathcal{R}_c] = \mathbb{P}_{\mathbf{H}\mathbf{z}^{(q)}}[k\mathcal{R}_c] = \mathbb{P}_{\mathbf{z}^{(q)}}[\mathbf{H}_{\text{PreImg}}(k\mathcal{R}_c)]$ [68, Section 2.2]. (b) follows because \mathbf{H} is a linear operation.

Although (21) holds for any \mathcal{R}_c , we will see that it is sufficient to consider a special \mathcal{R}_c to prove (20). To that end, take $\mathcal{R}_c = \text{conv}\{\mathbf{0}, \mathbf{a}_1, \dots, \mathbf{a}_{d_C}\}$, where $\mathbf{a}_i \in \mathbb{R}^{d_C}$ such that $\mathbb{P}_{\hat{\mathbf{c}}^{(q)}}[\mathcal{R}_c] > 0$. Let us take $\mathbf{y}_i \in \mathbb{R}^{d_C + d_P}$, such that $\mathbf{H}\mathbf{y}_i = \mathbf{a}_i$. For reasons that will be clear later, we hope to show that

$$\mathbf{H}_{\text{PreImg}}(\mathcal{R}_c) = \text{conv}\{\mathbf{0}, \mathbf{y}_1, \dots, \mathbf{y}_{d_C}\} + \text{Null}(\mathbf{H}).$$

To that end, observe that for any $\mathbf{r} \in \mathcal{R}_c$, we can represent \mathbf{r} as,

$$\mathbf{r} = \frac{1}{d_C + 1} \sum_{i=1}^{d_C} w_i \mathbf{a}_i, \text{ for some } \{w_i\}_{i=1}^{d_C} \text{ s.t. } \sum_{i=1}^{d_C} w_i \leq 1, \forall i.$$

For both view $q = 1, 2$, we get,

$$\begin{aligned} \mathbf{r} &= \frac{1}{d_C + 1} \sum_{i=1}^{d_C} w_i \mathbf{H} \mathbf{y}_i \\ \Rightarrow \mathbf{r} &= \mathbf{H} \left(\frac{1}{d_C + 1} \sum_{i=1}^{d_C} w_i \mathbf{y}_i \right) \\ \mathbf{H}_{\text{PreImg}} \left(\frac{1}{d_C + 1} \sum_{i=1}^{d_C} w_i \mathbf{a}_i \right) &= \frac{1}{d_C + 1} \sum_{i=1}^{d_C} w_i \mathbf{y}_i + \text{Null}(\mathbf{H}) \end{aligned} \quad (23)$$

We can write,

$$\mathbf{H}_{\text{PreImg}}(\mathcal{R}_c) = \text{conv}\{\mathbf{0}, \mathbf{y}_1, \dots, \mathbf{y}_{d_C}\} + \text{Null}(\mathbf{H}) \quad (24)$$

We have that $\text{Null}(\mathbf{H}) \subset \mathbb{R}^{d_C + d_P}$ is a linear subspace with $\dim(\text{Null}(\mathbf{H})) = d_P$. Let $\mathcal{A} = \mathbf{H}_{\text{PreImg}}(\mathcal{R}_c)$. Note that $\mathbb{P}_{\mathbf{z}^{(1)}}[k\mathcal{A}] = \mathbb{P}_{\mathbf{z}^{(2)}}[k\mathcal{A}], \forall k \in \mathbb{R}$ (from (21)), and $\mathbb{P}_{\mathbf{z}^{(q)}}[\mathcal{A}] > 0$ (by the construction of \mathcal{R}_c). Further, the set \mathcal{A} is of the form

$$\text{conv}\{\mathbf{0}, \mathbf{y}_1, \dots, \mathbf{y}_{d_C}\} + \mathcal{P},$$

because $\text{Null}(\mathbf{H})$ is a subspace of dimension d_P , hence it satisfies the definition of \mathcal{P} . Hence, Assumption 2 implies that

$$\text{Null}(\mathbf{H}) = \mathbf{0} \times \mathbb{R}^{d_P}.$$

Denoting the N th to M th columns of \mathbf{H} by $\mathbf{H}(N : M)$, the above is equivalent to saying

$$\mathbf{H}(d_C + 1 : d_C + d_P) = \mathbf{0}. \quad (25)$$

Denote,

$$\boldsymbol{\Theta} = \mathbf{H}(1 : d_C) \forall v = 1, 2.$$

Then, we can write,

$$\mathbf{Q} \mathbf{x}^{(q)} = \boldsymbol{\Theta} \mathbf{c}, \forall v = 1, 2. \quad (26)$$

This concludes the proof.

D Proof of Theorem 3

We restate the theorem here:

Theorem 3 Assume that Assumption 1 is satisfied, that $|\mathcal{L}| \geq d_C$ paired samples $(\mathbf{x}_\ell^{(1)}, \mathbf{x}_\ell^{(2)})$ are available, that $\mathbf{A}^{(q)}, q = 1, 2$ have full column rank, and that \mathbb{P}_c is absolutely continuous. Denote $(\hat{\mathbf{Q}}^{(1)}, \hat{\mathbf{Q}}^{(2)})$ as any optimal solution of (6) under the constraint (9). Then, we have $\hat{\mathbf{Q}}^{(q)} \mathbf{x}^{(q)} = \boldsymbol{\Theta} \mathbf{c}$.

From our objective in (6), we obtain

$$\mathbf{Q}^{(1)} \mathbf{x}^{(1)} \stackrel{(d)}{=} \mathbf{Q}^{(2)} \mathbf{x}^{(2)}. \quad (27)$$

Using Assumption 1 and following the proof of step 1 in Theorem B, we can obtain:

$$\mathbf{Q}^{(q)} \mathbf{x}^{(q)} = \boldsymbol{\Theta}^{(q)} \mathbf{c}, \forall q = 1, 2,$$

for some invertible matrices $\Theta^{(q)}, \forall q$. Hence,

$$\Theta^{(1)} \mathbf{c} \stackrel{(d)}{=} \Theta^{(2)} \mathbf{c} \quad (28)$$

$$\implies \mathbf{c} \stackrel{(d)}{=} (\Theta^{(1)})^{-1} \Theta^{(2)} \mathbf{c}. \quad (29)$$

Hence we can have linear transformation $\mathbf{M} := (\Theta^{(1)})^{-1} \Theta^{(2)}$ which has same probability density as $\mathbb{P}_{\mathbf{c}}$. However, the sample matching constraint [9], for ℓ -th sample implies that

$$\begin{aligned} \mathbf{Q}^{(1)} \mathbf{x}_{\ell}^{(1)} &= \mathbf{Q}^{(2)} \mathbf{x}_{\ell}^{(2)} \\ \implies \Theta^{(1)} \mathbf{c}_{\ell} &= \Theta^{(2)} \mathbf{c}_{\ell} \\ \implies \mathbf{c}_{\ell} &= (\Theta^{(1)})^{-1} \Theta^{(2)} \mathbf{c}_{\ell} \\ \implies \mathbf{c}_{\ell} &= \mathbf{M} \mathbf{c}_{\ell}. \end{aligned}$$

Let $\mathbf{C} = [\mathbf{c}_1 \dots \mathbf{c}_{N_p}]$. Then the above implies:

$$\begin{aligned} \mathbf{C} &= \mathbf{M} \mathbf{C} \\ \implies (\mathbf{M} - \mathbf{I}) \mathbf{C} &= \mathbf{0}. \end{aligned}$$

Now we show that \mathbf{C} is a full row rank matrix, which implies that $\mathbf{M} - \mathbf{I} = \mathbf{0} \implies \mathbf{M} = \mathbf{I}$. To that end, note that random variables $\mathbf{x}_i^{(1)}$ and $\mathbf{x}_i^{(2)}$ being i.i.d implies that $\mathbf{c}^{(i)}$ are i.i.d from $\mathbb{P}_{\mathbf{c}}$. This implies that for any $1 \leq i \leq |\mathcal{L}|$,

$$\Pr[\mathbf{c}_i \in \text{span}(\{\mathbf{c}_{n_1}, \dots, \mathbf{c}_{n_{d_C-1}}\})] = 0. \quad (30)$$

This is because $\text{span}(\{\mathbf{c}_{n_1}, \dots, \mathbf{c}_{n_{d_C-1}}\})$ for $n_j \in [|\mathcal{L}|]$, is a lower dimensional subspace in \mathbb{R}^{d_C} , which has zero probability under absolutely continuous distribution $\mathbb{P}_{\mathbf{c}}$. Hence any d_C out of $|\mathcal{L}|$ column vectors in \mathbf{C} are linearly independent with probability 1.

This concludes the proof.

E Detailed Identifiability Conditions of Existing Results

E.1 Identifiability of CCA

Theorem 4 (Identifiability of Aligned SCA via CCA [1]). *Under (1), assume that every aligned pair $(\mathbf{x}^{(1)}, \mathbf{x}^{(2)})$ share the same \mathbf{c} , and that $\mathbf{A}^{(q)}$ has full column rank. Also assume that there exists an N -sample set $\{\ell_1, \dots, \ell_N\}$ such that $[\mathbf{C}^{\top}, (\mathbf{P}^{(1)})^{\top}, (\mathbf{P}^{(2)})^{\top}] \in \mathbb{R}^{N \times (d_C + d_P^{(1)} + d_P^{(2)})}$ has full column rank, where $\mathbf{C} = [\mathbf{c}_{\ell_1}, \dots, \mathbf{c}_{\ell_N}] \in \mathbb{R}^{d_C \times N}$ and $\mathbf{P}^{(q)} = [\mathbf{p}_{\ell_1}^{(q)} \dots \mathbf{p}_{\ell_N}^{(q)}] \in \mathbb{R}^{d_P^{(q)} \times N}$ for $q = 1, 2$. Denote $(\hat{\mathbf{Q}}^{(1)}, \hat{\mathbf{Q}}^{(2)})$ as an optimal solution of the CCA formulation. Then, we have*

$$\hat{\mathbf{Q}}^{(q)} \mathbf{x}^{(q)} = \Theta \mathbf{c},$$

where Θ is nonsingular.

In the above theorem, one can see that $N \geq (d_C + d_P^{(1)} + d_P^{(2)})$ is a necessary condition for the identifiability of $\Theta \mathbf{c}$. Hence, CCA needs at least $d_C + d_P^{(1)} + d_P^{(2)}$ paired samples for identifiability.

E.2 Identifiability of Unaligned SCA in [8]

We summarize the result in [8] in the following

Theorem 5 (Identifiability of Unaligned SCA via ICA [8]). *Under (1), assume that the following are met: (i) The conditions for ICA identifiability [33] is met by each modality, including that the components of $\mathbf{z}^{(q)} = [\mathbf{c}^{\top}, (\mathbf{p}^{(q)})^{\top}]^{\top}$ are mutually statistically independent and contain at most one Gaussian variable. In addition, each $\mathbf{z}_i^{(q)}$ has unit variance; (ii) $\mathbb{P}_{\mathbf{z}_i^{(q)}} \neq \mathbb{P}_{\mathbf{z}_j^{(q)}}, \mathbb{P}_{\mathbf{z}_i^{(q)}} \neq \mathbb{P}_{-\mathbf{z}_j^{(q)}} \forall i, j \in [d_C + d_P^{(q)}], i \neq j$. Then, assume that (i_m, j_m) are obtained by ICA followed by cross domain matching (see the part on Unaligned SCA in Section 2) for $m = 1, \dots, d_C$.*

Denote $\hat{c}_m^{(1)} = \mathbf{e}_{i_m}^\top \hat{\mathbf{z}}^{(1)}$ and $\hat{c}_m^{(2)} = \mathbf{e}_{j_m}^\top \hat{\mathbf{z}}^{(2)}$. We have the following:

$$\hat{c}_m^{(q)} = k c_{\pi(m)}^{(q)}, \quad m \in [d_C], \quad (31)$$

where $k \in \{+1, -1\}$ and π is a permutation of $\{1, \dots, d_C\}$.

F Additional Synthetic Data Experiments

Hyperparameter Settings: We use Adam optimizer [70] to solve (7) and learn matrices $\mathbf{Q}^{(q)}$, $q = 1, 2$ and the discriminator f . We set the initial learning rate of matrix and discriminator to be 0.009 and 0.00008 respectively. We set the $\lambda = 0.1$ in (7) to enforce (6c). For weak supervision experiment in [E], we set $\beta = 0.01$ in (9). We generate total of 100,000 samples in each domain. For our experiment we set the batch size to be 1,000 and run (7) for 50 epochs. Our discriminator is a 6-layer multilayer perceptron (MLP) with hidden units $\{1024, 521, 512, 256, 128, 64\}$ in each layer. All the layers use leaky ReLU activation functions [71] with a slope of 0.2 except for the last layer which has sigmoid activations. We include a label smoothing coefficient of 0.2 in the discriminator predictions as suggested in [40].

Additional Details for Validation of Theorem 1 in Sec. 3: Here we explain the data generation details of the result shown in Fig. 3. For the result in top row, we sample c_1 from a Gaussian mixture with three Gaussian components. Each component follows a normal distribution $\mathcal{N}(\mu, 2)$ where $\mu \sim \mathcal{N}(0, 10)$. The second component, i.e., c_2 , is independently sampled from the gamma distribution $\text{Gamma}(1, 3)$. The private components are sampled from $p^{(1)} \sim \text{Laplace}(1.0, 6.5)$ and $p^{(2)} \sim \text{Uniform}[-10, 10]$, both only having one dimension. In the bottom row, we sample $c \in \mathbb{R}^2 \sim \text{VonMises}(2.5, 2.0)$ distribution. The private components satisfy $p^{(1)} \sim \text{Laplace}(1.0, 6.5)$ and $p^{(2)} \sim \text{Gamma}(0.5, 3.0)$. Each element of mixing matrices are sampled from $\mathbf{A}_{ij}^{(q)} \sim \mathcal{N}(0, 1)$, $q = 1, 2$. The readers are referred to Table 4 for the definition of notations used for distributions.

Validation of Theorem 1 under different sample sizes and imbalanced data: Here we observe the shared component identification performance of the proposed method numerically. We conducted two experiments in different settings. First, we vary the sample sizes in both modalities, but the two modalities have the same sample size. Second, we only vary the sample size of modality 2 while keeping the sample size of modality 1 fixed. This way, we create the data imbalance between modalities. Note that the shared components are identified if the following two conditions are met:

1. $\mathbf{Q}^{(q)} \mathbf{A}^{(q)} = [\Theta, \mathbf{0}]$, i.e., $\hat{\Theta}^{(1)} = \hat{\Theta}^{(2)} = \Theta$ and
2. $\|\mathbf{Q}^{(q)} \mathbf{A}^{(q)}(d_C : d_C + d_P^{(q)})\|_F = \mathbf{0}$.

Therefore, we use the above as our performance evaluation metrics. For the following experiments (Table 5 and 6), we generate the data for the two modalities by sampling a two-dimensional content $c \sim \text{VonMises}(2.5, 2.0)$ and private components from $p^{(1)} \sim \text{Laplace}(1.0, 6.5)$ and $p^{(2)} \sim \text{Gamma}(0.5, 3.0)$. The elements of the mixing matrices are sampled as $\mathbf{A}_{ij}^{(q)} \sim \mathcal{N}(0, 1)$, $q = 1, 2$. We report the mean and standard deviation of $\|\hat{\Theta}^{(1)}(1 : d_C) - \hat{\Theta}^{(2)}(1 : d_C)\|_F$ and $1/2 \sum_{q=1}^2 \|\hat{\Theta}^{(q)}(d_C : d_C + d_P^{(q)})\|_F$ obtained from 5 different runs.

Table 5 shows the performance of SCA and CCA under different sample sizes (i.e., N). One can see that the proposed method (SCA) clearly identifies the shared components even when only 100 samples are available. The performance starts to deteriorate when $N \leq 50$, probably because the min-max optimization problem is difficult to solve with very few samples. CCA does not really work under this setting as it needs aligned cross-domain samples.

Table 6 shows the performance of SCA in the cases where two modalities have unbalanced data sizes. The number of samples in the first modality is fixed to 100,000 while the second modality's data size varies from 10,000 to 10 samples. The data generation process remains the same as in the previous experiment. One can see that even under obvious cross-domain data size imbalance (e.g., 100,000 to 1,000), the proposed method performs reasonably well in terms of shared component identification.

Table 5: Shared component identification performance over different N .

N	$\ \widehat{\Theta}^{(1)}(1 : d_C) - \widehat{\Theta}^{(2)}(1 : d_C)\ _F$		$1/2 \sum_{q=1}^2 \ \widehat{\Theta}^{(q)}(d_C : d_C + d_P^{(q)})\ _F$	
	SCA	CCA	SCA	CCA
100,000	0.015 ± 0.020	1.623 ± 0.273	0.031 ± 0.009	0.232 ± 0.033
10,000	0.021 ± 0.006	1.667 ± 0.240	0.030 ± 0.002	0.267 ± 0.031
1,000	0.018 ± 0.011	1.572 ± 0.474	0.042 ± 0.059	0.280 ± 0.070
100	0.053 ± 0.014	2.224 ± 0.525	0.083 ± 0.096	0.364 ± 0.153
50	0.132 ± 0.118	1.469 ± 0.299	0.142 ± 0.132	1.470 ± 1.520
20	1.373 ± 0.626	2.084 ± 0.661	0.490 ± 0.321	0.546 ± 0.269

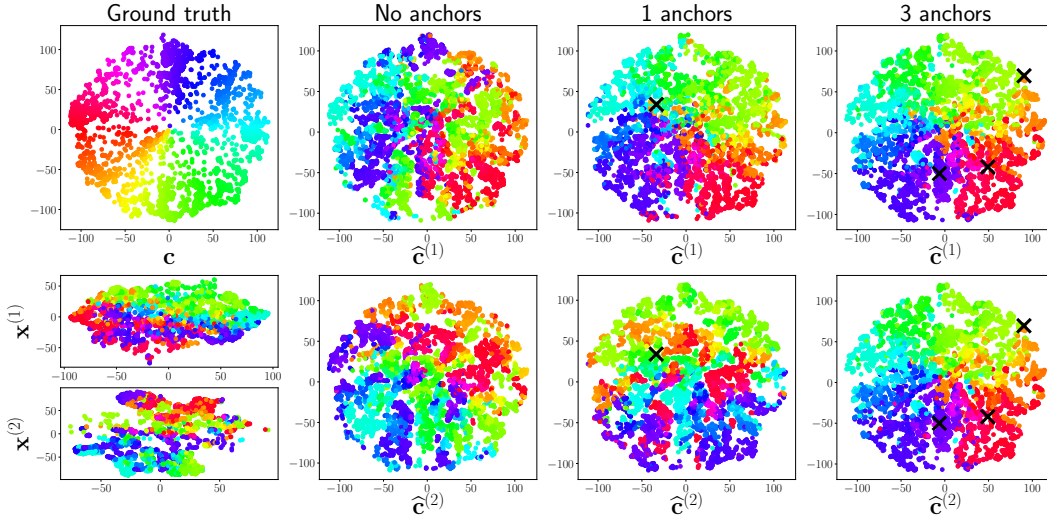
Table 6: Shared component identification performance under imbalanced multi-modal data sizes.

# samples in modality 2	$\ \widehat{\Theta}^{(1)}(1 : d_C) - \widehat{\Theta}^{(2)}(1 : d_C)\ _F$	$\frac{1}{2} \sum_{q=1}^2 \ \widehat{\Theta}^{(q)}(d_C : d_C + d_P^{(q)})\ _F$
10,000	0.020 ± 0.018	0.018 ± 0.005
1,000	0.065 ± 0.029	0.026 ± 0.015
100	0.145 ± 0.051	0.081 ± 0.049
10	1.290 ± 0.239	0.293 ± 0.064

Validation of Theorem 3. Fig. 5 presents numerical validation for Theorem 3.

Data Generation: We set $d_C = 3$ and $d_P^{(q)} = 1$ for $q = 1, 2$. We sample each component of shared component $c_i \sim \text{Laplace}(0.0, 6.5)$ $i = 1, 2, 3$, $p^{(1)} \in \mathbb{R}^1 \sim \text{Uniform}[-10, 10]$ and $p^{(2)} \sim \text{Gamma}(0.5, 3.0)$. Although c satisfies component-wise independence assumption, it does not satisfy the condition that $c_i \stackrel{(d)}{\neq} kc_j, \forall i \neq j$ because $c_i \stackrel{(d)}{=} c_j, \forall i, j \in [3]$. Therefore, Theorem 1 does not cover this case. Nonetheless, this case falls under the jurisdiction of Theorem 3.

Result: Fig. 5 corroborates with our Theorem 3. That is, one needs at least $|\mathcal{L}| \geq d_C = 3$ pairs of “anchors” (i.e., aligned cross domain pairs) to ensure identifiability of $\widehat{c}^{(q)} = \Theta c$ for $q = 1, 2$.

Figure 5: Validation of Theorem 3 $d_C = 3$ and $d_P^{(1)} = 1$.

G Real Data Experiment Settings and Additional Results

G.1 Domain Adaptation

Hyperparameter Settings: The domain adaptation task follows the hyperparameter settings described in Table. 7.

Table 7: Hyperparameter settings for domain adaptation.

Parameter	Value
Optimizer	Adam
Learning rate of Q	0.0002
Learning rate of f	0.00002
Learning rate of classifier	0.02
Learning rate decay of classifier	0.75
λ (see Eq. (7))	1.0
γ (see Eq. (32))	0.1
Batch size	64
Number of epochs	20
Discriminator; f architecture	6 layers, hidden units {1024, 521, 512, 256, 128, 64}
Activation functions of f	Leaky ReLU (slope 0.2), Sigmoid (final layer)
Label smoothing coefficient in f	0.2

Table 8: Classification accuracy on the target domain of *office-31* dataset using CLIP embeddings.

source \rightarrow target	CLIP	DANN	MDD	MCC	SDAT	SDAT+MCC	ELS	ELS+MCC	Proposed	Proposed+MCC
A \rightarrow W	93.4	93.7	94.1	95.9	95.0	98.1	96.8	98.7	95.3	98.3
D \rightarrow W	99.1	100.0	99.3	100.0	100.0	100.0	100.0	100.0	99.7	100.0
W \rightarrow D	100.0	99.5	99.5	98.4	99.5	99.5	99.5	100.0	100.0	99.9
A \rightarrow D	91.9	92.1	94.2	97.7	95.7	97.7	95.0	97.7	93.7	99.1
D \rightarrow A	81.4	81.9	79.2	85.7	81.2	84.6	81.3	83.0	83.9	85.9
W \rightarrow A	81.7	83.0	82.2	84.7	84.7	86.7	82.6	86.3	85.8	87.1
Average	91.2	91.6	91.4	93.7	92.6	94.4	92.5	94.2	93.0	95.0

Baselines and Training Setup: The baselines are representative DA methods, namely, DANN [25], MDD [60], MCC [61], SDAT [62], and ELS [63]. We use the implementations of DANN, MDD, and MCC from the <https://github.com/thuml/Transfer-Learning-Library>, while SDAT and ELS are taken from <https://github.com/yfzhang114/Environment-Label-Smoothing>. In all the baselines, the classifier is jointly optimized with the feature extractor Q which arguably regularizes towards more classification-friendly geometry of the shared features; see [72, 73]. Following their training strategy, we also append a cross-entropy (CE) based classifier training module to our loss in (7) (which learns our feature extractor Q). The CE part uses $Qx^{(1)}$ and the labels of the sources as inputs to learn the classifier, i.e.,

$$\mathcal{L}_{\text{CE}} = -\gamma \sum_{\ell=1}^N \sum_{k=1}^K \mathbb{I}[y_{\ell} = k] \log r_{\theta}([Qx_{\ell}^{(1)}]_k), \quad (32)$$

where $r_{\theta}(\cdot) : \mathbb{R}^{d_c} \rightarrow \mathbb{R}^K$ is the classifier that aims to map the learned feature vector $Qx_{\ell}^{(1)}$ to a K -dimensional probability mass function (i.e., the distribution of the ground-truth label over K classes), $y_{\ell} \in [K]$ represents the label of the ℓ th sample in source domain, and the indicator function $\mathbb{I}[y_{\ell} = k] = 1$ only when the event $y_{\ell} = k$ happens (otherwise $\mathbb{I}[y_{\ell} = k] = 0$). The $\gamma \geq 0$ is the tunable parameter. The joint loss is still differentiable, and thus we still use the Adam optimizer to jointly optimize Q and θ .

Additional domain adaptation experiment using CLIP features: In this experiment, we use CLIP as an image encoder as it learns informative and transferable features from very large datasets [35]. Table 8 and Table 9 show the results on *Office-31* and *Office-Home* datasets, respectively, using CLIP embeddings. Compared to the results on ResNet50 embeddings in Table 1 and Table 2, one can observe that all the methods, including proposed method, gains an advantage. This is likely because CLIP was trained on a large and diverse dataset [35], which may have include similar content to the *Office-31* and *Office-Home* datasets.

The results show that, as a foundation model, CLIP can already unify the embeddings of the source and target domains to a reasonable extent. In addition, our model and algorithm when combined with regularization techniques like MCC, can still further enhance performance, even with simple post-processing of CLIP embeddings.

Table 9: Classification accuracy on the target domain of *office-Home* dataset using CLIP embeddings.

source \rightarrow target	CLIP	DANN	MDD	MCC	SDAT	SDAT+MCC	ELS	ELS+MCC	Proposed	Proposed+MCC
Ar \rightarrow Cl	78.0	80.4	80.2	80.9	79.6	80.7	80.0	81.3	82.0	83.2
Ar \rightarrow Pr	88.7	91.7	88.9	93.3	89.4	94.3	91.2	93.9	91.4	95.2
Ar \rightarrow Rw	90.6	90.2	91.0	92.8	90.1	92.1	89.4	92.1	91.9	93.8
Cl \rightarrow Ar	85.2	83.2	85.1	87.4	83.1	86.1	84.4	87.2	85.4	87.7
Cl \rightarrow Pr	89.0	89.7	90.1	93.4	90.2	93.5	89.7	93.5	91.1	94.9
Cl \rightarrow Rw	89.8	88.1	89.4	89.3	87.9	90.5	88.3	90.6	90.4	92.0
Pr \rightarrow Ar	78.2	80.4	81.8	83.7	81.0	85.0	81.8	86.1	83.0	86.6
Pr \rightarrow Cl	72.7	75.8	75.8	78.4	75.4	78.5	75.7	78.3	77.7	81.3
Pr \rightarrow Rw	89.0	90.4	90.3	92.6	90.8	92.1	90.0	92.3	91.0	93.6
Rw \rightarrow Ar	86.6	84.9	85.9	85.3	85.3	86.3	85.4	87.0	87.7	88.2
Rw \rightarrow Cl	78.1	79.4	79.8	79.0	78.6	79.8	79.1	79.8	81.3	81.8
Rw \rightarrow Pr	94.3	94.6	93.9	95.9	94.8	95.4	94.0	95.2	94.7	96.0
Average	85.0	85.7	86.0	87.6	85.5	87.8	85.7	88.1	87.3	89.5

Visualization Result: Fig. 6 shows the 2-dimensional visualization of the CLIP-learned features ($d = 256$) from two domains, namely, DSLR and Amazon images (*Office-31*), using t-SNE. One can see that CLIP could roughly group the same classes from the two domains together. But the proposed method can further pull the circles and the triangle markers together—meaning that the Q really learns shared representations of the same data in the DSLR and Amazon domains.

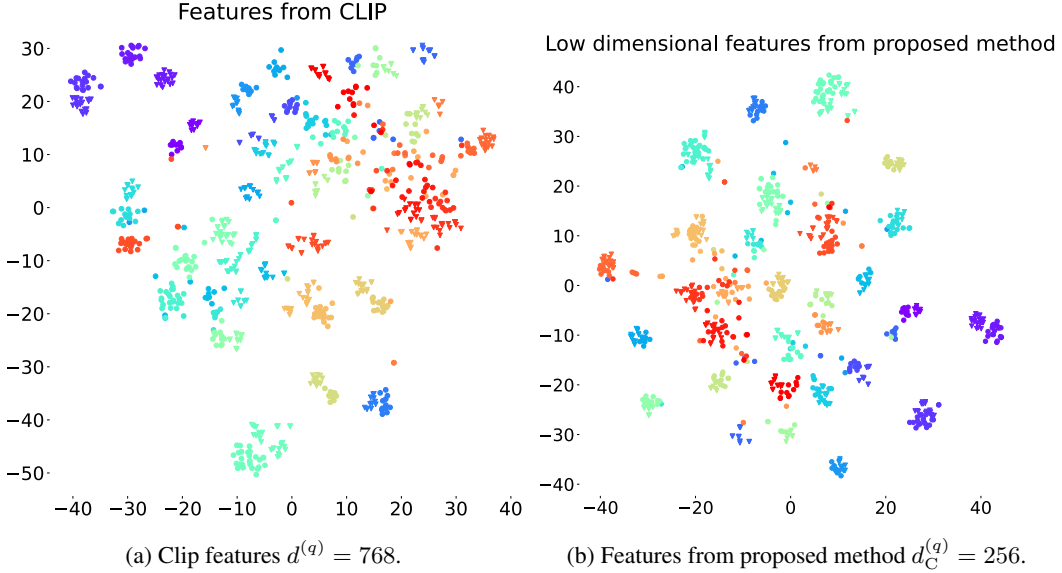


Figure 6: *Office-31* dataset: DSLR images features represented as circle markers, Amazon images features represented as triangle markers. Different color represent different classes.

G.2 Single-cell Sequence Analysis

Hyperparameter Settings: The hyperparameter settings for single-cell sequence analysis is presented in Table. 10.

Baseline: For more details on baseline refer to the implementation in <https://github.com/uhlerlab/cross-modal-autoencoders>

G.3 Multi-lingual Information Retrieval

Hyperparameter Settings: The hyperparameter settings for multi-lingual information retrieval is described in the Table. 11.

Table 10: Hyperparameter settings for single-cell sequence analysis.

Parameter	Value
Optimizer	Adam
Learning rate of $Q^{(q)}$	0.001
Learning rate of f	0.0001
λ (see Eq. (7))	1.0
β (see Eq. (9))	10.0
Batch size	32
Number of epochs	75
Discriminator; f architecture	6 layers, hidden units {1024, 512, 512, 256, 128, 64}
Activation functions of f	Leaky ReLU (slope 0.2), Sigmoid (final layer)
Label smoothing coefficient in f	0.2

Table 11: Hyperparameter settings for multi-lingual information retrieval.

Parameter	Value
Optimizer	Adam
Learning rate of Q	0.0001
Learning rate of f	0.00001
λ (see Eq. (7))	1.0
Batch size	32
Number of epochs	5
Discriminator; f (similar as in [21])	2 layers, 2048 hidden units each
Activation functions of f	Leaky ReLU (slope 0.2), Sigmoid (final layer)
Dropout rate (Input) in f	0.1
Label smoothing coefficient in f	0.2

Additional Results: Table [12] reports the P@5 and P@10 scores over the test data, calculated for different source and target language pairs. It can be observed that the proposed method achieves higher precision than Adv in most of the translation tasks (e.g., by at least 1% in the **en**→**es** and **es**→**en** tasks) when considering both P@5 and P@10 scores.

G.4 Computation resources

All the experiments were run on Nvidia H100 GPU. The approximate runtime for a single run of the algorithm is 20 minutes for multi-lingual information retrieval, 15 minutes for domain adaptation, and 3 minutes for single-cell sequence analysis.

Complexity Analysis:

Since the proposed objective is tackled using stochastic gradient (SG)-based first-order iterative method, the computational complexity of the proposed algorithm depends upon the per-iteration complexity.

For each sample, the per-iteration complexity is composed of a forward pass and a backward pass.

Note that the problem size depends upon $d^{(q)}$ (the data dimension), d_C , and the batch size B . We assume that the network architecture of f (the number of layers and hidden units in each layer) is fixed, represented by $f = \sigma \circ F_L \circ \dots \circ \sigma \circ F_1$, where F_ℓ and σ are the linear layer (matrix) and activation function corresponding to the ℓ th layer. Only the input dimension d_C of first matrix F_1 , varies with the problem size.

The forward pass involves computing $\hat{c}^{(q)} = Q^{(q)}x^{(q)}$ and $f^{(q)}(\hat{c}^{(q)})$, both of which scale linearly with d_C , $d^{(q)}$ and the batch size B . Hence, the forward pass time complexity is $O(Bd_C(d^{(1)} + d^{(2)}))$.

Similarly, the backward pass requires computing of $\frac{\partial L}{\partial \hat{c}_i^{(q)}}$, $\forall i \in [d_C]$ and $\frac{\partial \hat{c}_i^{(q)}}{\partial Q^{(q)}_{jk}}$, $\forall i \in [d_C], j, k \in [d_C \times d^{(q)}]$, where L is the loss function. The first gradient computation is linear in Bd_C , while the second gradient computation has a complexity of $O(Bd_C(d^{(1)} + d^{(2)}))$. Hence the computational complexity of our method is $O(Bd_C(d^{(1)} + d^{(2)}))$.

Table 12: Average precision P@k of cross-language information retrieval

P@k	source \rightarrow target	Adv - NN	proposed - NN	Adv - CSLS	proposed - CSLS
P@5	en \rightarrow es	77.9	78.8	83.6	85.2
	es \rightarrow en	73.2	79.0	83.2	86.0
	en \rightarrow it	68.0	70.8	76.8	82.4
	it \rightarrow en	79.0	77.6	71.2	66.9
	en \rightarrow fr	79.2	75.2	85.7	85.2
	fr \rightarrow en	69.8	73.5	77.2	83.5
P@10	en \rightarrow es	82.4	82.6	87.0	87.8
	es \rightarrow en	79.0	82.3	87.2	88.8
	en \rightarrow it	74.1	75.6	81.7	85.6
	it \rightarrow en	71.5	75.8	82.0	82.9
	en \rightarrow fr	83.4	79.1	88.0	88.4
	fr \rightarrow en	73.8	77.4	80.6	86.6

The memory complexity involves storing the network parameters and the aforementioned gradients. Hence, only the size of $\mathbf{Q}^{(q)}$, \mathbf{F}_1 , and $\mathbf{c}^{(q)}$ changes with the problem dimension. The size of $\mathbf{Q}^{(q)}$, \mathbf{F}_1 , and $\mathbf{c}^{(q)}$ are $d_C d^{(q)}$, $O(d_C)$ and d_C , respectively. Therefore, the space complexity is $O(Bd_C(d^{(1)} + d^{(2)}))$.

In summary, both the memory and computational complexities of the proposed method scales linearly with d_C .

H Extension: Private Component Identification

Theorems 1-3 are concerned with learning the shared component \mathbf{c} . The goal, there, was to ensure that $\mathbf{Q}_C^{(q)} \mathbf{x}^{(q)} \Theta \mathbf{c}, \forall q$. In some cases, the private components $\mathbf{p}^{(q)}$ is also of interest [6, 31, 74]. To learn $\mathbf{p}^{(q)}$, we propose to solve the following learning criterion:

$$\text{find } \mathbf{Q}_C^{(q)} \in \mathbb{R}^{d_C \times d^{(q)}}, \mathbf{Q}_P^{(q)} \in \mathbb{R}^{d_P^{(q)} \times d^{(q)}} \quad q = 1, 2, \quad (33a)$$

$$\text{subject to } \mathbf{Q}_C^{(1)} \mathbf{x}^{(1)} \stackrel{(d)}{=} \mathbf{Q}_C^{(2)} \mathbf{x}^{(2)}, \quad (33b)$$

$$\mathbf{Q}_C^{(q)} \mathbf{x}^{(q)} \perp\!\!\!\perp \mathbf{Q}_P^{(q)} \mathbf{x}^{(q)} \quad q = 1, 2, \quad (33c)$$

$$\mathbf{Q}_C^{(q)} \mathbb{E} [\mathbf{x}^{(q)} (\mathbf{x}^{(q)})^\top] (\mathbf{Q}_C^{(q)})^\top = \mathbf{I} \quad q = 1, 2, \quad (33d)$$

$$\mathbf{Q}_P^{(q)} \mathbb{E} [\mathbf{x}^{(q)} (\mathbf{x}^{(q)})^\top] (\mathbf{Q}_P^{(q)})^\top = \mathbf{I} \quad q = 1, 2, \quad (33e)$$

where $\mathbf{u} \perp\!\!\!\perp \mathbf{v}$ means that the random vectors \mathbf{u} and \mathbf{v} are independent with each other.

For implementation we use following criterion,

$$\begin{aligned} & \min_{\mathbf{Q}_C^{(1)}, \mathbf{Q}_C^{(2)}, \mathbf{Q}_P^{(1)}, \mathbf{Q}_P^{(2)}} \max_f \mathbb{E}_{\mathbf{x}^{(1)}} \log \left(f(\mathbf{Q}_C^{(1)} \mathbf{x}^{(1)}) \right) + \mathbb{E}_{\mathbf{x}^{(2)}} \log \left(1 - f(\mathbf{Q}_C^{(2)} \mathbf{x}^{(2)}) \right) \\ & + \lambda \sum_{q=1}^2 \mathcal{R}(\mathbf{Q}_C^{(q)}) + \omega \sum_{q=1}^2 \mathcal{R}(\mathbf{Q}_P^{(q)}) + \rho \sum_{q=1}^2 \text{HSIC}(\mathbf{Q}_C^{(q)} \mathbf{x}^{(q)}, \mathbf{Q}_P^{(q)} \mathbf{x}^{(q)}), \end{aligned} \quad (34)$$

where, first two term are adversarial loss for distribution matching. The constraint on (33d) and (33e) are enforced as $\mathcal{R}(\mathbf{Q}_C^{(q)})$ and $\mathcal{R}(\mathbf{Q}_P^{(q)})$ respectively, where $\mathcal{R}(\mathbf{Q}^{(q)}) = \|\mathbf{Q}^{(q)} \mathbb{E} [\mathbf{x}^{(q)} (\mathbf{x}^{(q)})^\top] (\mathbf{Q}^{(q)})^\top - \mathbf{I}\|_F^2$. The constraint on (33c) is realized with Hilbert-Schmidt Independence Criterion (HSIC) [75]. HSIC measures the independence between two distribution. So, we minimize HSIC between estimated shared component and estimated private component to promote independence between shared and private components.

We show that under some reasonable conditions the block $\mathbf{p}^{(q)}$ can also be learned up to a matrix multiplication:

Theorem 6. Assume that the blocks \mathbf{c} , $\mathbf{p}^{(1)}$ and $\mathbf{p}^{(2)}$ are statistically independent, i.e., $p(\mathbf{c}, \mathbf{p}^{(1)}, \mathbf{p}^{(2)}) = p(\mathbf{c})p(\mathbf{p}^{(1)})p(\mathbf{p}^{(2)})$. Then, if one of the following holds:

- (i) Assumption 1 and assumptions in Theorem 1 are satisfied, and (33) is solved yielding solutions $\hat{\mathbf{Q}}_C^{(q)}$ and $\hat{\mathbf{Q}}_P^{(q)}$
- (ii) Assumption 2 is satisfied and has same mixing matrix $\mathbf{A}^{(q)} = \mathbf{A}$ and (33) with $\mathbf{Q}_P^{(q)} = \mathbf{Q}_P$ and $\mathbf{Q}_C^{(q)} = \mathbf{Q}_C$ is solved yielding $\hat{\mathbf{Q}}_C^{(q)}$ and $\hat{\mathbf{Q}}_P^{(q)}$ as the solutions.
- (iii) Assumption 1 is satisfied and d_C paired samples $(\mathbf{x}_\ell^{(1)}, \mathbf{x}_\ell^{(2)})$ are available (weak supervision), and denote $\hat{\mathbf{Q}}_C^{(q)}$ and $\hat{\mathbf{Q}}_P^{(q)}$ as the solutions after solving (33).

Then, we have $\hat{\mathbf{Q}}_C^{(q)} \mathbf{x}^{(q)} = \Theta \mathbf{c}$ and $\hat{\mathbf{Q}}_P^{(q)} \mathbf{x}^{(q)} = \Xi^{(q)} \mathbf{p}^{(q)}$, for some invertible $\Xi^{(q)}$ for all $q = 1, 2$.

Proof. For each case in Theorem 6 (i) - (iii), we can prove

$$\hat{\mathbf{c}}^{(q)} = \hat{\mathbf{Q}}_C^{(q)} \mathbf{x}^{(q)} = \Theta \mathbf{c}, \quad q = 1, 2 \quad (35)$$

using Theorems 1-3. The proofs are referred to Appendix B-D.

Let us denote

$$\hat{\mathbf{p}}^{(q)} = \hat{\mathbf{Q}}_P^{(q)} \mathbf{x}^{(q)} = \hat{\mathbf{Q}}_P^{(q)} \mathbf{A}^{(q)} \begin{bmatrix} \mathbf{c} \\ \mathbf{p}^{(q)} \end{bmatrix} = \mathbf{H}^{(q)} \begin{bmatrix} \mathbf{c} \\ \mathbf{p}^{(q)} \end{bmatrix}, \quad (36)$$

where $\mathbf{H}^{(q)} = \hat{\mathbf{Q}}_P^{(q)} \mathbf{A}^{(q)} \in \mathbb{R}^{d_P^{(q)} \times (d_C + d_P^{(q)})}$. Note that the constraint (33c) implies that the mutual information between $\hat{\mathbf{p}}^{(q)}$ and $\hat{\mathbf{c}}^{(q)}$ is zero, i.e.,

$$I(\hat{\mathbf{p}}^{(q)}; \hat{\mathbf{c}}^{(q)}) = 0.$$

Note that $\hat{\mathbf{p}}^{(q)} \rightarrow \hat{\mathbf{c}}^{(q)} \rightarrow \Theta^{-1} \hat{\mathbf{c}}^{(q)} = \mathbf{c}$ is a Markov chain. This is because when conditioned on $\hat{\mathbf{c}}^{(q)}$, $\Theta^{-1} \hat{\mathbf{c}}^{(q)}$ becomes constant, making it independent of $\hat{\mathbf{p}}^{(q)}$. This allows us to use the data processing inequality [76, Theorem 2.8.1], which results in the following:

$$I(\hat{\mathbf{p}}^{(q)}; \hat{\mathbf{c}}^{(q)}) \geq I(\hat{\mathbf{p}}^{(q)}; \Theta^{-1} \hat{\mathbf{c}}^{(q)}) = I(\hat{\mathbf{p}}^{(q)}; \mathbf{c}).$$

Since mutual information is always non-negative, the above implies that $I(\hat{\mathbf{p}}^{(q)}; \mathbf{c}) = 0$. This implies that $\hat{\mathbf{p}}^{(q)} = \mathbf{H}^{(q)} \begin{bmatrix} \mathbf{c} \\ \mathbf{p}^{(q)} \end{bmatrix}$ is independent of \mathbf{c} . Hence, $\mathbf{H}^{(q)}[1 : d_C] = 0, \forall q$.

Therefore $\hat{\mathbf{p}}^{(q)} = \mathbf{H}^{(q)}[d_C + 1 : d_C + d_P^{(q)}] \mathbf{p}^{(q)} = \Xi^{(q)} \mathbf{p}^{(q)}, \forall q$, where $\Xi^{(q)} = \mathbf{H}^{(q)}[d_C + 1 : d_C + d_P^{(q)}]$. Note that $\mathbf{H}^{(q)}$ is full row-rank because of constraint (33e). This implies that $\Xi^{(q)}, q = 1, 2$ are invertible matrices.

This concludes the proof. □

H.1 Validation of Theorem 6

Fig. 7 presents numerical validation for Theorem 6.

Hyperparameter Setting The hyperparameter setting is the same as mentioned in Appendix F. We solve (34) to obtain $\hat{\mathbf{Q}}_C^{(q)}$ and $\hat{\mathbf{Q}}_P^{(q)}$ to recover the shared and private components, respectively. For learning $\mathbf{Q}_P^{(q)}$, we use Adam optimizer and set initial learning rates to be 0.001. Also we set the regularization parameter $\omega = 10.0$ and $\rho = 50.0$.

Data Generation: We set $d_C = 2$ and $d_P^{(q)} = 1$ as in the previous synthetic experiments. We sample $\mathbf{c} \sim \text{VonMises}(2.5, 2.0)$. The private components are sampled from $p^{(1)} \sim \text{Beta}(1.0, 3.0)$ and $p^{(2)} \sim \text{Gamma}(0.5, 3.0)$ distributions. Each element of mixing matrices are sampled from $A_{ij}^{(q)} \sim \mathcal{N}(0, 1)$, $q = 1, 2$.

Result: Fig. 7 shows the result for proposed method for private component identification. The first column shows the data domain, the second column shows the true and extracted shared component, and the third and fourth columns shows the true and extracted private components. Especially, the last row of the third and fourth columns shows the plot of ground truth $\mathbf{p}^{(q)}$ on x -axis and $\hat{\mathbf{p}}^{(q)}$ on the y -axis. The plot is approximately a straight line which indicates that the estimated private components $\hat{\mathbf{p}}^{(q)}$ are scaled version (i.e., invertible linear transformations) of ground truth private components. This verifies our Theorem 6.

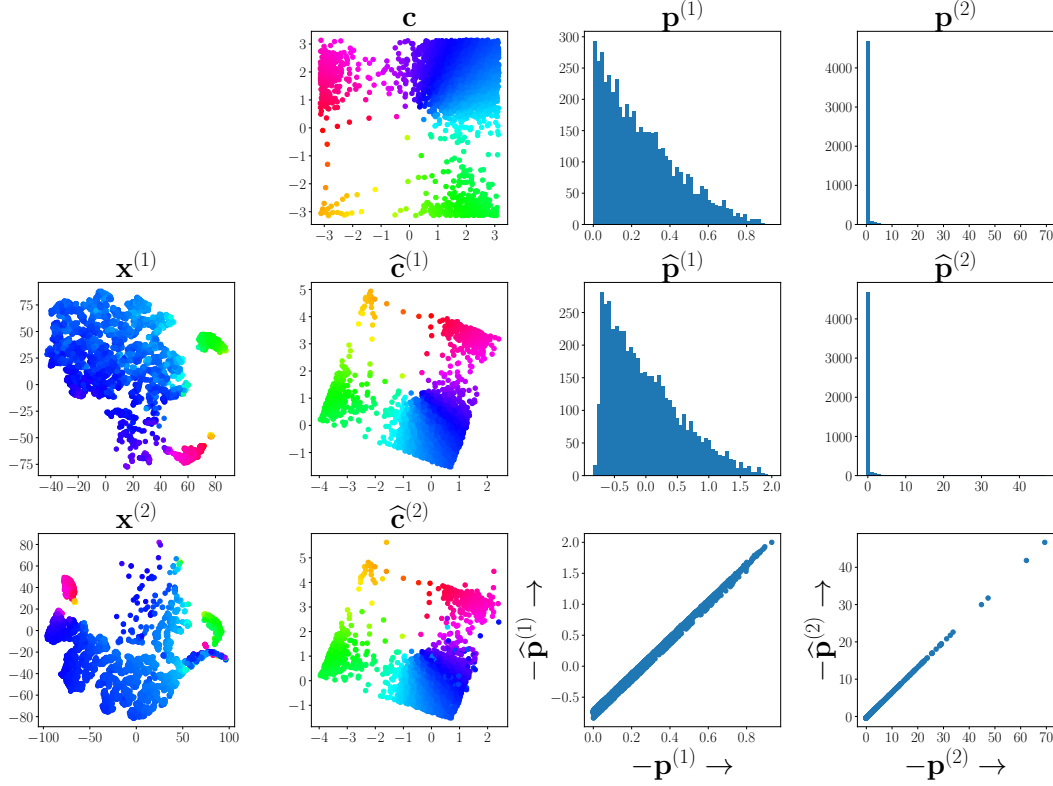


Figure 7: Validation of Theorem 6 $d_C = 2$ and $d_P^{(1)} = 1$.

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: See Section [3](#), [4](#) and [6](#)

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: See Section [7](#)

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: See Appendix [B](#), [C](#), [D](#) and [H](#).

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: See Appendix [F](#) and [G](#).

- 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: Yes the code is provided in the supplemental material.

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: See Appendix [F](#) and [G](#).

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: We have provided the error bars for the experiment that is computationally less demanding.

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: See Appendix [G.4](#).

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: We have followed the NeurIPS Code of Ethics.

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: [NA]

Justification: The contribution of this paper is on theoretical aspects of machine learning. We don't foresee any immediate societal impacts.

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: Theoretical paper. So not applicable in our case.

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: See Appendix [G](#).

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 donot release any 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: Not applicable.

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: Not Applicable.

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.