# NO-RANK TENSOR DECOMPOSITION USING METRIC LEARNING

MARYAM BAGHERIAN

*Department of Mathematics and Statistics, Idaho State University*
*Physical Science Complex— 921 S. 8th Ave., Stop 8085 — Pocatello, ID 83209*

ABSTRACT. Tensor decomposition faces fundamental challenges in analyzing high-dimensional data, where traditional methods based on reconstruction and fixed-rank constraints often fail to capture semantically meaningful structures. This paper introduces a no-rank tensor decomposition framework grounded in metric learning, which replaces reconstruction objectives with a discriminative, similarity-based optimization. The proposed approach learns data-driven embeddings by optimizing a triplet loss with diversity and uniformity regularization, creating a feature space where distance directly reflects semantic similarity. We provide theoretical guarantees for the framework's convergence and establish bounds on its metric properties. Evaluations across diverse domains—including face recognition (LFW, Olivetti), brain connectivity analysis (ABIDE), and simulated data (galaxy morphology, crystal structures)—demonstrate that our method outperforms baseline techniques, including PCA, t-SNE, UMAP, and tensor decomposition baselines (CP and Tucker). Results show substantial improvements in clustering metrics (Silhouette Score, Davies-Bouldin Index, Calinski-Harabasz Index, Separation Ratio, Adjusted Rand Index, Normalized Mutual Information) and reveal a fundamental trade-off: while metric learning optimizes global class separation, it deliberately transforms local geometry to align with semantic relationships. Crucially, our approach achieves superior performance with smaller training datasets compared to transformer-based methods, offering an efficient alternative for domains with limited labeled data. This work establishes metric learning as a paradigm for tensor-based analysis, prioritizing semantic relevance over pixel-level fidelity while providing computational advantages in data-scarce scenarios.

**Key words:** Tensor Decomposition, Metric Learning, Representation Learning, Dimensionality Reduction, Clustering, Triplet Loss, Embedding Learning, Semantic Structure Learning

## 1. INTRODUCTION

Tensor decomposition and representation learning are two fundamental paradigms for extracting meaningful structure from multi-dimensional data. Traditional tensor decomposition methods, such as CANDECOMP/PARAFAC (CP) [22] and Tucker decomposition [41], provide mathematically elegant approaches for factorizing tensors into interpretable components.

The CP decomposition approximates an $N$-way tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_N}$ as a sum of $R$ rank-one tensors:

$$\mathcal{X} \approx \sum_{r=1}^{R} \mathbf{u}_r^{(1)} \circ \mathbf{u}_r^{(2)} \circ \cdots \circ \mathbf{u}_r^{(N)}, \tag{1}$$

where $\circ$ denotes the outer product, and $R$ is the *rank* of the decomposition. Similarly, the Tucker decomposition employs a core tensor $\mathcal{G} \in \mathbb{R}^{R_1 \times R_2 \times \cdots \times R_N}$ and factor matrices $\mathbf{U}^{(n)} \in \mathbb{R}^{I_n \times R_n}$:

$$\mathcal{X} \approx \mathcal{G} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \times_3 \cdots \times_N \mathbf{U}^{(N)}, \tag{2}$$

where $\times_n$ denotes the $n$-mode product.

While these methods offer strong uniqueness guarantees [23] and interpretability, they require pre-specification of rank parameters $(R, R_1, \ldots, R_N)$ and are inherently limited to multi-linear relationships.

Representation learning [9] provides a complementary approach, learning useful data representations through neural networks without explicit structural constraints. Deep autoencoders, for instance, learn an encoding function $f_\theta$ mapping input data to latent codes $\mathbf{z}$ and a decoding function $g_\phi$ reconstructing the original input:

$$\hat{\mathcal{X}} = g_\phi(f_\theta(\mathcal{X})), \tag{3}$$

The latent representation $\mathbf{z}$ captures essential factors of variation in the data, with the model capacity determining the effective complexity rather than pre-defined rank constraints.

Recent work has begun bridging these domains, with neural networks enhancing tensor decompositions [27] and implicit neural representations demonstrating remarkable compression capabilities [39]. However, the fundamental challenge of *rank selection* persists across traditional methods.

This work introduces a "no-rank" paradigm where representation learning principles enable adaptive, data-driven tensor decomposition without explicit rank specification, addressing key limitations of both established approaches while leveraging their respective strengths.

The rest of the manuscript is organized as follows: Related work are discussed in Section 2 followed by the model formulation in Section 3. Section 4 provides an overview of the metric analysis and result interpretation and experimental results over real and simulated dataset are presented and discussed in Section 5. We conclude in Section 6 by discussing the advantages, limitations and future direction of the proposed framework.

## 2. Related Work

This section reviews the relevant literature in these fields, highlighting the gap that *no-rank* metric learning framework aims to fill.

**Traditional Tensor Decomposition Methods.** Tensor decomposition methods have long been the workhorse for analyzing multi-way data. The Canonical Polyadic (CP) decomposition [18, 22] and the Tucker decomposition [41] are the two most foundational approaches. Both methods impose a strict low-rank constraint on the data, seeking an approximation $\hat{\mathcal{X}}$ that minimizes the reconstruction error $||\mathcal{X} - \hat{\mathcal{X}}||_F^2$. While effective for data compression and denoising, this reconstructive objective is not inherently aligned with discriminative tasks like classification or clustering. The requirement to pre-define the rank or multilinear rank is a significant limitation, as the intrinsic data complexity is often unknown and may not be well-represented by a low-rank model [1].

Subsequent advancements, such as tensor-train decompositions [30] and non-negative tensor factorizations [10], have improved scalability and interpretability but have largely

remained within the reconstructive paradigm. These methods are fundamentally linear and struggle to capture the complex, non-linear manifolds on which high-dimensional data often resides.

**Dimensionality Reduction and Manifold Learning.** To address non-linearity, a separate lineage of work focused on manifold learning and non-linear dimensionality reduction [3]. Techniques such as Isomap [40], Locally Linear Embedding (LLE) [33], and Laplacian Eigenmaps [7] aim to preserve geometric properties of the data. More recently, t-Distributed Stochastic Neighbor Embedding (t-SNE) [42] and Uniform Manifold Approximation and Projection (UMAP) [28] have become standards for visualization, excelling at preserving local neighborhood structures.

However, these methods are predominantly *unsupervised* and *geometry-preserving*. They lack a mechanism to incorporate supervisory signals, such as class labels, to guide the feature learning process. Consequently, the resulting low-dimensional embeddings may not optimize for class separability, which is crucial for many analysis tasks. Furthermore, they operate as separate pre-processing steps and are not integrated into an end-to-end trainable model for feature extraction.

**Deep Metric and Representation Learning.** The advent of deep learning catalyzed a shift towards learning representations directly optimized for a specific task. A pivotal development in this space is *metric learning* [8,24], which aims to learn a distance function that reflects semantic similarity. The contrastive loss [16] and, more influentially, the triplet loss [34] provided a powerful framework for this. By pulling an anchor sample closer to a positive sample than to a negative sample by a margin, these losses directly optimize the embedding space for discrimination.

This paradigm has driven state-of-the-art performance in face recognition [34,45], image retrieval [14], and person re-identification [17]. This work draws direct inspiration from these successes but adapts the triplet loss framework to the problem of tensor decomposition, moving beyond its typical application in computer vision.

To prevent pathological solutions like dimensional collapse, recent work has emphasized the importance of regularization. Wang and Isola [46] identified *alignment* and *uniformity* as key properties of effective representations, which has led to the use of uniformity losses [46] and diversity penalties on the embedding correlation matrix [6]. The proposed framework incorporates these insights to ensure a well-structured and effective embedding space.

**Metric Learning for Data Analysis.** There is a growing recognition of the limitations of traditional methods for data analysis. In domains like astronomy, methods for galaxy classification have evolved from manual taxonomy [20] to machine learning approaches using hand-crafted features [13] and, more recently, deep convolutional networks [44]. Similarly, in materials science, the analysis of crystal structures has been tackled with symmetry-based descriptors [21] and graph neural networks [48].

However, the application of deep metric learning in these contexts is still nascent. While some studies have used contrastive learning for spectroscopic data [37] or medical imaging [38], a generalized framework that replaces the core principles of tensor decomposition for tensor data is lacking. Existing approaches often treat metric learning as a separate module rather than as a fundamental alternative to decomposition.

Built on the previous works [2,4,5], the proposed *no-rank* tensor decomposition framework synthesizes ideas from these disparate fields. We reframe the problem of tensor analysis from one of *reconstruction* to one of *discrimination*, drawing on the power of triplet-based metric learning. Unlike traditional tensor methods, we impose no explicit rank constraints and leverage deep non-linear networks to capture complex data manifolds. Unlike unsupervised manifold learning, we directly optimize for semantic similarity using label information. And unlike standard metric learning applications, the proposed method is positioned as a direct, end-to-end replacement for tensor decomposition in the data pipeline, a novel contribution with significant practical implications.

## 3. Metric Learning Framework

3.1. **Theoretical Framework.** Traditional tensor decomposition methods, such as Canonical Polyadic (CP) and Tucker decompositions, impose strict rank constraints that may not align with the intrinsic geometry of high-dimensional data. We propose a *no-rank* tensor decomposition framework based on metric learning, which learns data-driven similarity structures without explicit rank constraints.

Let $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_n\} \subset \mathbb{R}^D$ be a set of multi-dimensional data points with corresponding labels $\mathcal{Y} = \{y_1, y_2, \ldots, y_n\}$. The goal is to learn a mapping $f : \mathbb{R}^D \to \mathbb{R}^d$, where $d \ll D$, such that the resulting low-dimensional embedding space preserves semantically meaningful relationships from the original space.

Let $\mathbf{x}_i = \text{vec}(\mathcal{X}_{:,:,\ldots,i})$ be the vectorized $i$-th mode-$N$ slice of the tensor. We learn embeddings $\mathbf{z}_i = f(\mathbf{x}_i)$ by constructing triplets $(a, p, n) \in \mathcal{T}$, where the *anchor* $(a)$ is a reference sample: $\mathbf{z}_a = f(\mathbf{x}_a)$, the *positive* $(p)$ is a semantically similar sample (e.g., from the same class): $\mathbf{z}_p = f(\mathbf{x}_p)$ and the *negative* $(n)$ is a semantically dissimilar sample (e.g., from a different class): $\mathbf{z}_n = f(\mathbf{x}_n)$.

The model is trained to pull the anchor and positive together while pushing the anchor and negative apart, leading to an embedding space where semantic similarity is inversely proportional to the distance in embedding space. This triplet formulation provides a flexible framework for *Semantic Structure Learning*, directly optimizing for similarity relationships rather than reconstruction error.

3.2. **Triplet Loss with Regularization.** The core optimization objective is to ensure that the distance to negatives exceeds the distance to positives by at least a margin $\alpha$. This is enforced using the triplet loss combined with regularization terms to promote a well-structured embedding space.

The triplet loss is defined as:

$$\mathcal{L}_{\text{triplet}} = \sum_{(a,p,n) \in \mathcal{T}} \left[ \|\mathbf{z}_a - \mathbf{z}_p\|_2^2 - \|\mathbf{z}_a - \mathbf{z}_n\|_2^2 + \alpha \right]_+ , \tag{4}$$

where $[x]_+ = \max(0, x)$ and $\alpha > 0$ is the margin parameter. This objective ensures that for all triplets:

$$\|\mathbf{z}_a - \mathbf{z}_p\|_2^2 + \alpha < \|\mathbf{z}_a - \mathbf{z}_n\|_2^2. \tag{5}$$

To prevent dimensional collapse and encourage the model to use all available dimensions efficiently, we add a diversity penalty on the embedding correlation matrix:

$$\mathcal{L}_{\text{div}} = \frac{1}{d(d-1)} \sum_{i \neq j} |\mathbf{C}_{ij}|, \tag{6}$$

where $\mathbf{C} = \mathbf{Z}^\top \mathbf{Z}$ is the correlation matrix of the embeddings $\mathbf{Z} = [\mathbf{z}_1, \ldots, \mathbf{z}_n]^\top$.

Following [46], we also promote a uniform distribution of embeddings on the unit sphere to avoid hubness and further improve generalization:

$$\mathcal{L}_{\text{uniform}} = \log \mathbb{E}_{i,j \sim p_{\text{data}}} \left[ e^{-2\|\mathbf{z}_i - \mathbf{z}_j\|_2^2} \right]. \tag{7}$$

3.3. **Locality Preservation Framework.** To ensure that local neighborhoods in the original high-dimensional space are preserved in the embedding space, we introduce locality preservation objectives. Let $\mathcal{N}_k(\mathbf{x}_i)$ denote the set of $k$-nearest neighbors of $\mathbf{x}_i$ in the original space. The goal is to ensure:

$$\mathbf{x}_j \in \mathcal{N}_k(\mathbf{x}_i) \Rightarrow f(\mathbf{x}_j) \in \mathcal{N}_k(f(\mathbf{x}_i)). \tag{8}$$

We quantify this using two standard metrics. *Continuity* measures how well the original neighbors are represented in the embedding neighborhood:

$$\text{Continuity} = 1 - \frac{1}{nk} \sum_{i=1}^{n} \sum_{j \in \mathcal{N}_k^E(\mathbf{x}_i)} \mathbb{I}\{j \notin \mathcal{N}_k^O(\mathbf{x}_i)\} \cdot r_O(i, j), \tag{9}$$

where $\mathcal{N}_k^O$ and $\mathcal{N}_k^E$ are neighborhoods in the original and embedding spaces, respectively, and $r_O(i, j)$ is the rank of $j$ in $\mathbf{x}_i$'s original neighborhood. *Trustworthiness* measures the prevalence of "intruders" (points in the embedding neighborhood that were not in the original neighborhood):

$$\text{Trustworthiness} = 1 - \frac{2}{nk(2n - 3k - 1)} \sum_{i=1}^{n} \sum_{j \in \mathcal{U}_k(\mathbf{x}_i)} (r_O(i, j) - k), \tag{10}$$

where $\mathcal{U}_k(\mathbf{x}_i)$ is the set of intruders for point $i$.

To directly optimize for these properties, we augment the loss function with locality-preserving terms. The *local consistency loss* ensures that original neighbors remain close in the embedding space:

$$\mathcal{L}_{\text{local}} = \sum_{i=1}^{n} \sum_{j \in \mathcal{N}_k^O(\mathbf{x}_i)} \left[ \|f(\mathbf{x}_i) - f(\mathbf{x}_j)\|_2^2 - \delta_{\text{local}} \right]_+. \tag{11}$$

Conversely, the *global separation loss* ensures that non-neighbors remain far apart:

$$\mathcal{L}_{\text{global}} = \sum_{i=1}^{n} \sum_{j \notin \mathcal{N}_k^O(\mathbf{x}_i)} \left[ \delta_{\text{global}} - \|f(\mathbf{x}_i) - f(\mathbf{x}_j)\|_2^2 \right]_+. \tag{12}$$

The complete training objective is a weighted sum of all components:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{triplet}} + \lambda_1 \mathcal{L}_{\text{div}} + \lambda_2 \mathcal{L}_{\text{uniform}} + \lambda_3 \mathcal{L}_{\text{local}} + \lambda_4 \mathcal{L}_{\text{global}}, \tag{13}$$

where $\lambda_1, \lambda_2, \lambda_3, \lambda_4 > 0$ are hyperparameters that balance the loss terms.

Unlike low-rank tensor decomposition methods (e.g., CP or Tucker), which seek a rank-$R$ approximation by minimizing reconstruction error $\|\mathcal{X} - \hat{\mathcal{X}}\|_F^2$, the approach learns embeddings without explicit rank constraints by directly optimizing for semantic similarity. Furthermore, while traditional tensor decompositions are limited to linear transformations, the use of a neural network encoder allows for the capture of complex non-linear patterns.

### 3.4. Network Architecture.

The embedding function $f(\cdot)$ is implemented as a deep neural network. The encoder, which maps an input $\mathbf{x}$ to the unit sphere, has the following structure:

$$\mathbf{h}^{(1)} = \text{ReLU}(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) \tag{14}$$

$$\mathbf{h}^{(l)} = \text{ReLU}(\mathbf{W}^{(l)}\mathbf{h}^{(l-1)} + \mathbf{b}^{(l)}) \quad \text{for } l = 2, \dots, L \tag{15}$$

$$\mathbf{z} = \frac{\mathbf{W}^{(L+1)}\mathbf{h}^{(L)} + \mathbf{b}^{(L+1)}}{\|\mathbf{W}^{(L+1)}\mathbf{h}^{(L)} + \mathbf{b}^{(L+1)}\|_2}. \tag{16}$$

The final $\ell_2$-normalization projects the embedding onto the unit sphere, which is compatible with the uniformity loss. A non-linear projection head can be added to enhance representation quality during training:

$$\mathbf{p} = \mathbf{W}_2 \cdot \text{ReLU}(\mathbf{W}_1\mathbf{z}). \tag{17}$$

### 3.5. Convergence and Geometric Guarantees.

We now establish theoretical guarantees for the proposed framework, focusing on optimization convergence and the geometric properties of the learned embedding space.

If the mapping $f$ is $L$-Lipschitz continuous, i.e.,

$$\|f(\mathbf{x}) - f(\mathbf{x}')\|_2 \leq L\|\mathbf{x} - \mathbf{x}'\|_2,$$

then local neighborhood relationships are preserved up to a factor of $L$, providing a theoretical guarantee for locality preservation.

For a metric learning model with Rademacher complexity $\mathfrak{R}_n$ [29, 36], the generalization error $\mathcal{E}_{\text{gen}}$ is bounded with probability at least $1 - \delta$ by:

$$\mathcal{E}_{\text{gen}} \leq \mathcal{E}_{\text{emp}} + O\left(\frac{\mathfrak{R}_n}{\sqrt{n}} + \sqrt{\frac{\log(1/\delta)}{n}}\right). \tag{18}$$

The gradients of the triplet loss with respect to the embeddings create a force field that pulls positive pairs together and pushes negative pairs apart [34]. For an active triplet (where the loss is positive), the gradients are:

$$\nabla_{\mathbf{z}_a}\mathcal{L}_{\text{triplet}} = 2(\mathbf{z}_a - \mathbf{z}_n) - 2(\mathbf{z}_a - \mathbf{z}_p), \tag{19}$$

$$\nabla_{\mathbf{z}_p}\mathcal{L}_{\text{triplet}} = 2(\mathbf{z}_p - \mathbf{z}_a), \tag{20}$$

$$\nabla_{\mathbf{z}_n}\mathcal{L}_{\text{triplet}} = 2(\mathbf{z}_a - \mathbf{z}_n). \tag{21}$$

**Lemma 1** (Convergence to a Critical Point). *Let the total loss function $\mathcal{L}_{total}(\theta)$ be $L$-smooth and bounded below. Assume the stochastic gradients $g_t$ satisfy the conditional unbiasedness property*

$$\mathbb{E}[g_t \mid \mathcal{F}_t] = \nabla_\theta \mathcal{L}_{total}(\theta_t),$$

*and have uniformly bounded conditional variance*

$$\mathbb{E}\big[\|g_t - \nabla_\theta \mathcal{L}_{total}(\theta_t)\|^2 \mid \mathcal{F}_t\big] \leq \sigma^2$$

*almost surely, where $\mathcal{F}_t$ denotes the history up to time $t$. When minimized using stochastic gradient descent with a learning rate schedule $\{\eta_t\}$ satisfying the Robbins–Monro conditions $\sum_{t=1}^{\infty} \eta_t = \infty$ and $\sum_{t=1}^{\infty} \eta_t^2 < \infty$ [32], the parameter sequence $\{\theta_t\}$ satisfies*

$$\liminf_{t \to \infty} \mathbb{E}\big[\|\nabla_\theta \mathcal{L}_{total}(\theta_t)\|^2\big] = 0.$$

*Consequently, there exists a subsequence of iterates $\{t_k\}$ such that $\mathbb{E}\big[\|\nabla_\theta \mathcal{L}_{total}(\theta_{t_k})\|^2\big] \to 0$. If additionally the iterates $\{\theta_t\}$ are almost surely bounded (or $\mathcal{L}_{total}$ is coercive), then there exists a (random) subsequence $\theta_{t_k}$ that converges a.s. to a limit point $\theta^*$, and any such limit point satisfies $\nabla_\theta \mathcal{L}_{total}(\theta^*) = 0$.*

*Proof.* The proof follows standard stochastic approximation arguments. Because $\mathcal{L}_{\text{total}}$ is $L$-smooth, the one-step descent inequality (taking total expectation over the noise history) yields

$$\mathbb{E}\big[\mathcal{L}_{\text{total}}(\theta_{t+1})\big] \le \mathbb{E}\big[\mathcal{L}_{\text{total}}(\theta_t)\big] - \eta_t \mathbb{E}\big[\|\nabla\mathcal{L}_{\text{total}}(\theta_t)\|^2\big] + \frac{L}{2}\eta_t^2 \mathbb{E}\big[\|g_t\|^2\big].$$

Using conditional unbiasedness and the conditional variance bound gives (via law of total expectation)

$$\mathbb{E}\big[\|g_t\|^2\big] = \mathbb{E}\big[\|\nabla\mathcal{L}_{\text{total}}(\theta_t)\|^2\big] + \mathbb{E}\big[\|g_t - \nabla\mathcal{L}_{\text{total}}(\theta_t)\|^2\big] \le \mathbb{E}\big[\|\nabla\mathcal{L}_{\text{total}}(\theta_t)\|^2\big] + \sigma^2.$$

Substituting and summing from $t = 1$ to $T$ yields

$$\sum_{t=1}^{T} \eta_t \mathbb{E}\big[\|\nabla\mathcal{L}_{\text{total}}(\theta_t)\|^2\big] \le \mathbb{E}\big[\mathcal{L}_{\text{total}}(\theta_1)\big] - \mathcal{L}^* + \frac{L\sigma^2}{2}\sum_{t=1}^{T}\eta_t^2,$$

where $\mathcal{L}^*$ is the infimum of $\mathcal{L}_{\text{total}}$. The Robbins–Monro conditions ensure the right-hand side stays bounded as $T \to \infty$. Since $\sum_t \eta_t = \infty$, the only possibility is

$$\liminf_{t\to\infty} \mathbb{E}\big[\|\nabla\mathcal{L}_{\text{total}}(\theta_t)\|^2\big] = 0,$$

which implies the existence of a subsequence $t_k$ with $\mathbb{E}[\|\nabla\mathcal{L}(\theta_{t_k})\|^2] \to 0$. The final statement about parameter subsequence convergence follows from the additional boundedness/coercivity assumption (so that limit points exist), after which continuity of $\nabla\mathcal{L}$ implies any limit point is stationary. $\qquad\square$

**Lemma 2** (Semantic Structure of the Embedding). *Assume the data is separable with a margin $\gamma > 0$, i.e., the optimal embedding satisfies $\|\mathbf{z}_a - \mathbf{z}_p\|_2^2 + \gamma \le \|\mathbf{z}_a - \mathbf{z}_n\|_2^2$ for all valid triplets $(a, p, n)$. Further, assume the embedding function $f$ is $L$-Lipschitz and the data manifold $\mathcal{M}$ is compact. Then, the learned embedding space exhibits the following semantic structure:*

> *(i) Intra-class clusters are tight: for any two points $\mathbf{x}_i, \mathbf{x}_j$ from the same class, $\|f(\mathbf{x}_i) - f(\mathbf{x}_j)\|_2 \le L \cdot diam(\mathcal{M}_c)$, where $\mathcal{M}_c$ is the connected component of the data manifold containing points of that class.*

> *(ii) Inter-class clusters are separated: for any two points $\mathbf{x}_i, \mathbf{x}_j$ from different classes, $\|f(\mathbf{x}_i) - f(\mathbf{x}_j)\|_2$ is lower-bounded by $\sqrt{\gamma}$.*

*This structure ensures that local neighborhoods are preserved, and the embedding is well-suited for similarity search.*

*Proof.* The proof leverages the margin condition and Lipschitz continuity.
(i) Intra-class tightness: Consider two points $\mathbf{x}_i$ and $\mathbf{x}_j$ from the same class. Since they belong to the same class and the data manifold $\mathcal{M}$ is compact, there exists a path connecting them within their class component $\mathcal{M}_c$. By the Lipschitz continuity of $f$, we have:

$$\|f(\mathbf{x}_i) - f(\mathbf{x}_j)\|_2 \le L \cdot d_{\mathcal{M}}(\mathbf{x}_i, \mathbf{x}_j) \le L \cdot \text{diam}(\mathcal{M}_c),$$

where $\text{diam}(\mathcal{M}_c)$ is the diameter of the connected component of class $c$ in the data manifold. This provides a uniform upper bound on intra-class distances.

(ii) Inter-class separation: Let $\mathbf{x}_i$ and $\mathbf{x}_j$ be from different classes. Using the margin condition with anchor $a = \mathbf{x}_i$, positive $p = \mathbf{x}_i$ (trivially from the same class), and negative $n = \mathbf{x}_j$, we get:

$$\|f(\mathbf{x}_i) - f(\mathbf{x}_i)\|_2^2 + \gamma \leq \|f(\mathbf{x}_i) - f(\mathbf{x}_j)\|_2^2.$$

Simplifying yields $\gamma \leq \|f(\mathbf{x}_i) - f(\mathbf{x}_j)\|_2^2$, and thus $\sqrt{\gamma} \leq \|f(\mathbf{x}_i) - f(\mathbf{x}_j)\|_2$. The combination of (i) and (ii) demonstrates that the embedding maps the semantic structure of the data into a geometric structure where points from the same class form tight clusters (bounded by manifold geometry) well-separated from clusters of other classes (by at least $\sqrt{\gamma}$). The Lipschitz property ensures this mapping preserves local connectivity. $\qquad\square$

**Theorem 3.** *Under the assumptions of Lemmas 1 and 2, the proposed metric learning framework:*

    *(1) Converges to a configuration that is a local minimum of the total objective $\mathcal{L}_{total}$.*

    *(2) Yields a semantically structured embedding space where intra-class distances are minimized and inter-class distances are maximized beyond a margin $\gamma$.*

*This result bridges optimization guarantees with geometric structure, ensuring that the learned representation is both stable and semantically meaningful.*

*Proof.* The theorem follows from combining the guarantees of the two lemmas. From Lemma 1, the optimization converges to a critical point $\theta^*$ where $\liminf_{t \to \infty} \mathbb{E}[\|\nabla_\theta \mathcal{L}_{\text{total}}(\theta_t)\|^2] = 0$. Under the strict saddle point assumption[1], gradient-based methods converge almost surely to local minima rather than saddle points. Thus, $\theta^*$ is a local minimum of $\mathcal{L}_{\text{total}}$ with high probability.

From Lemma 2, when the embedding function $f_{\theta^*}$ (parameterized by the locally optimal $\theta^*$) operates on separable data with margin $\gamma$, the resulting embedding space exhibits the semantic structure described in Lemma 2. Specifically (i) intra-class distances are bounded by $L \cdot \text{diam}(\mathcal{M}_c)$, (ii) inter-class distances are lower-bounded by $\sqrt{\gamma}$.

The combination of these results guarantees that SGD finds a locally optimal parameter configuration that produces a semantically structured embedding space suitable for similarity tasks. $\qquad\square$

**Remark 4.** *The theoretical guarantees depend on key parameters:*

    *(i) The margin $\gamma$ (from Lemma 2) directly controls inter-class separation,*

    *(ii) The Lipschitz constant $L$ (from Lemma 2) controls how well local neighborhoods are preserved,*

    *(iii) The learning rate conditions (from Lemma 1) ensure optimization convergence.*

*In practice, the regularization terms $\mathcal{L}_{div}$ and $\mathcal{L}_{uniform}$ promote well-spread embeddings (supporting large $\gamma$), while $\mathcal{L}_{local}$ and $\mathcal{L}_{global}$ enforce neighborhood preservation (related to the Lipschitz property).*

---

[1]The strict saddle property requires that every saddle point has at least one direction of negative curvature. Under this condition, gradient-based methods with random initialization almost surely avoid saddle points [26].

**Assumption 5** (Manifold Hypothesis). *The data is assumed to lie on a smooth, low-dimensional Riemannian manifold $\mathcal{M} \subset \mathbb{R}^D$. The embedding function $f : \mathcal{M} \to \mathbb{R}^d$ is approximately isometric, meaning that Euclidean distances in the embedding space preserve the intrinsic (geodesic) distances on the manifold:*

$$d_{\mathcal{M}}(\mathbf{x}_i, \mathbf{x}_j) \approx \|f(\mathbf{x}_i) - f(\mathbf{x}_j)\|_2, \tag{22}$$

*for all $\mathbf{x}_i, \mathbf{x}_j \in \mathcal{M}$.*

This framework is particularly well-suited for applications like medical imaging, where semantic similarity and clinical relevance are more critical than pixel-perfect reconstruction accuracy, opening new possibilities for tensor-based analysis of complex spatio-temporal data.

## 4. Metrics and Methodology Analysis

4.1. **Evaluation Metrics.** We assess the quality of the learned embeddings using a comprehensive suite of metrics that evaluate clustering quality, structural preservation, and alignment with ground-truth labels.

Clustering Quality Metrics evaluate the intrinsic structure of the embeddings without using label information:
*Silhouette Score (Sil)* measures how similar samples are to their own cluster compared to other clusters. For a sample $i$, it is computed as:

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}, \tag{23}$$

where $a(i)$ is the mean intra-cluster distance and $b(i)$ is the mean nearest-cluster distance. Scores range from $-1$ to $+1$, with higher values indicating better clustering.

Davies-Bouldin Index(DB) quantifies the trade-off between cluster compactness and separation. For $k$ clusters, it is defined as:

$$\mathrm{DB} = \frac{1}{k} \sum_{i=1}^{k} \max_{j \neq i} \left( \frac{\sigma_i + \sigma_j}{d(c_i, c_j)} \right), \tag{24}$$

where $\sigma_i$ is the average distance from points in cluster $i$ to its centroid $c_i$, and $d(c_i, c_j)$ is the distance between centroids. Lower values indicate better clustering.

*Calinski-Harabasz Index*(CH) is defined as the ratio of between-cluster dispersion to within-cluster dispersion:

$$\mathrm{CH} = \frac{\mathrm{Tr}(B_k)}{\mathrm{Tr}(W_k)} \times \frac{N - k}{k - 1}, \tag{25}$$

where $\mathrm{Tr}(B_k)$ and $\mathrm{Tr}(W_k)$ are the traces of the between-cluster and within-cluster dispersion matrices, respectively. Higher values indicate tighter and better-separated clusters.

External Validation Metrics measure the agreement between the discovered clusters and the ground-truth labels: *Adjusted Rand Index (ARI)* measures the similarity between two clusterings, corrected for chance. It is defined as:

$$\mathrm{ARI} = \frac{\sum_{ij} \binom{n_{ij}}{2} - \left[ \sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right] / \binom{n}{2}}{\frac{1}{2} \left[ \sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2} \right] - \left[ \sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right] / \binom{n}{2}}, \tag{26}$$

where $n_{ij}$ is the contingency table between true and predicted clusters, $a_i$ and $b_j$ are the row and column sums. ARI ranges from $-1$ to $1$, with higher values indicating better agreement.

*Normalized Mutual Information (NMI)* measures the mutual information between true and predicted clusters, normalized by the entropy of each:

$$\text{NMI} = \frac{2 \cdot I(Y; \hat{Y})}{H(Y) + H(\hat{Y})}, \tag{27}$$

where $I(Y; \hat{Y})$ is the mutual information and $H(\cdot)$ is the entropy. NMI ranges from 0 to 1, with higher values indicating better cluster alignment.

Metric Learning Specific Metrics directly evaluate the effectiveness of the triplet loss objective:

*Separation Ratio*(SR) quantifies the ratio of inter-class to intra-class distances:

$$\text{Separation Ratio} = \frac{\mathbb{E}[\|f(\mathbf{x}_i) - f(\mathbf{x}_j)\|_2 \mid y_i \neq y_j]}{\mathbb{E}[\|f(\mathbf{x}_i) - f(\mathbf{x}_j)\|_2 \mid y_i = y_j]}. \tag{28}$$

A higher ratio indicates that the embedding space successfully pulls same-class samples together and pushes different-class samples apart.

Structural Preservation Metrics evaluate how well local geometric structure is preserved during dimensionality reduction:

*Trustworthiness* (Trust.) measures the preservation of local structure by penalizing for false neighbors (points that are neighbors in the embedding but not in the original space):

$$T(k) = 1 - \frac{2}{nk(2n - 3k - 1)} \sum_{i=1}^{n} \sum_{j \in \mathcal{U}_k(i)} (r_O(i, j) - k), \tag{29}$$

where $\mathcal{U}_k(i)$ are the intruders for point $i$ and $r_O(i, j)$ is the rank of $j$ in the original space.

*Continuity*(Cont.) is the complementary measure to trustworthiness, which penalizes for missing neighbors (points that are neighbors in the original space but not in the embedding):

$$C(k) = 1 - \frac{2}{nk(2n - 3k - 1)} \sum_{i=1}^{n} \sum_{j \in \mathcal{V}_k(i)} (r_E(i, j) - k), \tag{30}$$

where $\mathcal{V}_k(i)$ are the missing neighbors for point $i$ and $r_E(i, j)$ is the rank of $j$ in the embedding space.

These metrics provide complementary perspectives: Silhouette, DB, and CH indices focus on intrinsic cluster quality; ARI and NMI validate clustering against ground truth; the Separation Ratio directly measures metric learning effectiveness; and Trustworthiness and Continuity evaluate the preservation of local geometric structure.

4.2. **Visualization of High-Dimensional Embeddings.** To gain qualitative insights into the structure of the learned embeddings, we project them into two dimensions using both linear and non-linear techniques. The quality of these visualizations is inherently dependent on the distance metric learned by the proposed model.

We employ three standard dimensionality reduction methods Principal Component Analysis (PCA) [15] as a linear baseline that projects data onto the directions of maximal

variance. t-Distributed Stochastic Neighbor Embedding (t-SNE) [42] as a non-linear technique that preserves local similarities by minimizing the Kullback-Leibler divergence [25] between probability distributions in the high- and low-dimensional spaces:

$$\mathrm{KL}(P\|Q) = \sum_{i \neq j} p_{ij} \log \frac{p_{ij}}{q_{ij}}, \tag{31}$$

where

$$p_{j|i} = \frac{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2/2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|\mathbf{x}_i - \mathbf{x}_k\|^2/2\sigma_i^2)}, \quad q_{ij} = \frac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_{k \neq l}(1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}}. \tag{32}$$

And Uniform Manifold Approximation and Projection (UMAP) [28] as a manifold learning technique that assumes the data is uniformly distributed on a Riemannian manifold. It constructs a topological representation and optimizes a low-dimensional equivalent.

We project the raw, flattened input data $\mathbf{X}_{\mathrm{raw}} \in \mathbb{R}^{N \times D}$ to 2D using PCA to establish a baseline: $\mathbf{X}_{\mathrm{raw}} \to \mathbf{Y}_{\mathrm{PCA\text{-}raw}} \in \mathbb{R}^{N \times 2}$. And to represent metric learning embedding space, we project the learned embeddings $\mathbf{Z} \in \mathbb{R}^{N \times d}$ using all three techniques $\mathbf{Z} \to \mathbf{Y}_\alpha \in \mathbb{R}^{N \times 2}$, where $\alpha = \mathrm{PCA}, \mathrm{t\text{-}SNE}, \mathrm{UMAP}$. Well-separated, tight clusters in these projections indicate a semantically coherent embedding space. To quantitatively validate these qualitative observations, we complement the visualizations with histograms of intra-class and inter-class distances:

$$D_{\mathrm{intra}} = \{\|\mathbf{z}_i - \mathbf{z}_j\|_2 : y_i = y_j\}, \quad D_{\mathrm{inter}} = \{\|\mathbf{z}_i - \mathbf{z}_j\|_2 : y_i \neq y_j\}. \tag{33}$$

A clear separation between the $D_{\mathrm{intra}}$ and $D_{\mathrm{inter}}$ distributions confirms the patterns observed in the 2D projections.

4.3. **Triplet Mining Strategy Analysis.** The strategy for selecting triplets from the training data is critical for efficient and stable model convergence. We analyze two common strategies [17, 34, 47].

The *Semi-Hard Negative Mining* strategy selects negatives that are farther from the anchor than the positive, but within the margin $\alpha$. This provides a steady, moderate learning signal.

ALGORITHM 1. Semi-Hard Negative Mining

**Require:** Embeddings $\mathbf{Z}$, labels $\mathbf{y}$, margin $\alpha$
1: **for** each anchor $\mathbf{z}_a$ with label $y_a$ **do**
2:     Find a random positive $\mathbf{z}_p$ where $y_p = y_a$
3:     Compute positive distance $d_p = \|\mathbf{z}_a - \mathbf{z}_p\|_2$
4:     Find the set of negatives $\mathcal{N} = \{\mathbf{z}_n \mid y_n \neq y_a, \ d_p < \|\mathbf{z}_a - \mathbf{z}_n\|_2 < d_p + \alpha\}$
5:     Select a random negative from $\mathcal{N}$ (if non-empty)
6: **end for**

The *Hard Negative Mining* is a more aggressive strategy selects the most challenging triplets by choosing the most distant positive and the closest negative for each anchor. This can lead to faster learning but also risks instability if the hard negatives are outliers or mislabeled.

ALGORITHM 2. Hard Negative Mining

**Require:** Embeddings $\mathbf{Z}$, labels $\mathbf{y}$
1: **for** each anchor $\mathbf{z}_a$ with label $y_a$ **do**
2:    Find hardest positive: $\mathbf{z}_p^* = \arg\max_{\mathbf{z}_p : y_p = y_a} \|\mathbf{z}_a - \mathbf{z}_p\|_2$
3:    Find hardest negative: $\mathbf{z}_n^* = \arg\min_{\mathbf{z}_n : y_n \neq y_a} \|\mathbf{z}_a - \mathbf{z}_n\|_2$
4:    Use triplet $(\mathbf{z}_a, \mathbf{z}_p^*, \mathbf{z}_n^*)$
5: **end for**

The comparative performance of these strategies is evaluated and presented in the following section.

## 5. RESULTS

### 5.1. **Metric Learning Performance in Face Recognition.**
We evaluate this framework on face recognition, a canonical metric learning task. Here, the goal is to learn an embedding where the distance between an anchor image and a positive example (same person) is smaller than the distance to a negative example (different person): $d(\text{anchor}, \text{positive}) < d(\text{anchor}, \text{negative})$.

We used two contrasting datasets (Table 1): the Labeled Faces in the Wild (LFW) [19] dataset, a medium-sized, imbalanced dataset representing a real-world challenge; and the Olivetti Faces [31] dataset, a smaller, balanced dataset captured in a controlled environment.

TABLE 1. Summary of Face Recognition Dataset Properties

| Property | LFW Faces | Olivetti Faces |
|---|---|---|
| Total Images | 1,288 | 400 |
| Identities | 7 | 40 |
| Image Dimensions | $50 \times 37$ | $64 \times 64$ |
| Samples per Person | 71–530 | 10 |
| Class Distribution | Highly Imbalanced | Perfectly Balanced |
| Environment | Unconstrained | Controlled |

*Quantitative Clustering Performance.* The clustering results (Table 2) demonstrate the decisive advantage of the metric learning approach over baselines like PCA and tensor decomposition methods for creating semantically meaningful clusters.

On the challenging LFW dataset, proposed approach achieved a near-perfect Silhouette score of 0.9752, a dramatic improvement over PCA (-0.0186). This is corroborated by the Davies-Bouldin index, which dropped from 7.33 to 0.0566, and the Separation Ratio, which increased from 1.01 to 49.18, indicating that inter-class distances became vastly larger than intra-class distances.

The same trend is clear on the Olivetti dataset, where metric learning outperformed all other methods, achieving a Silhouette score of 0.8566, a Davies-Bouldin index of 0.2341, and a Separation Ratio of 9.8471.

TABLE 2. Clustering Performance on LFW and Olivetti Datasets

| Dataset | Method | Sil. | DB | SR | Cont. | Trust. | AIR | NMI |
|---------|--------|------|-----|-----|-------|--------|-----|-----|
| **LFW** | PCA + K-Means | -0.0186 | 7.3302 | 1.0131 | **0.9967** | **0.9933** | 0.0181 | 0.0324 |
| | t-SNE + K-Means | -0.0922 | 212.4068 | 1.0005 | 0.9206 | 0.9487 | 0.0128 | 0.0347 |
| | UMAP + K-Means | -0.0815 | 41.6583 | 1.0013 | 0.9198 | 0.8721 | 0.0049 | 0.0195 |
| | CP-R5 | -0.1216 | 18.5653 | 1.0033 | 0.8186 | 0.8715 | 0.0066 | 0.0242 |
| | CP-R10 | -0.0461 | 13.2544 | 1.0234 | 0.8824 | 0.9072 | 0.0061 | 0.0272 |
| | CP-R20 | -0.0906 | 10.3420 | 0.9720 | 0.9090 | 0.9274 | 0.0061 | 0.3050 |
| | Tucker-R5 | -0.0689 | 10.8919 | 1.0181 | 0.9504 | 0.9033 | 0.0222 | 0.0479 |
| | Tucker-R10 | -0.0338 | 7.9771 | 1.0254 | 0.9837 | 0.9666 | 0.0187 | 0.0450 |
| | Tucker-R20 | -0.0037 | 5.8455 | 1.0280 | 0.9886 | 0.9789 | 0.0078 | 0.0296 |
| | **Metric Learning** | **0.9752** | **0.0566** | **49.1800** | 0.9236 | 0.9201 | **1.0000** | **1.0000** |
| **Olivetti** | PCA + K-Means | 0.1434 | 1.8243 | 1.6002 | **0.9982** | **0.9970** | 0.3831 | 0.7275 |
| | t-SNE + K-Means | -0.0123 | 9.2275 | 2.3923 | 0.9449 | 0.9730 | 0.4737 | 0.7898 |
| | UMAP + K-Means | -0.1213 | 8.4268 | 2.1196 | 0.9399 | 0.9352 | 0.3072 | 0.6700 |
| | CP-R5 | -0.1432 | 3.7012 | 1.6951 | 0.9163 | 0.8920 | 0.1115 | 0.5403 |
| | CP-R10 | -0.0514 | 2.9316 | 1.7121 | 0.9386 | 0.9384 | 0.2250 | 0.6299 |
| | CP-R20 | 0.0021 | 2.5596 | 1.5354 | 0.9624 | 0.9536 | 0.2397 | 0.6411 |
| | Tucker-R5 | -0.0275 | 2.9707 | 1.7517 | 0.9664 | 0.9358 | 0.2374 | 0.6554 |
| | Tucker-R10 | 0.1303 | 1.9508 | 1.6772 | 0.9891 | 0.9800 | 0.4185 | 0.7542 |
| | Tucker-R20 | 0.1845 | 1.6731 | 1.5764 | 0.9935 | 0.9864 | 0.41976 | 0.7956 |
| | **Metric Learning** | **0.8566** | **0.2341** | **9.8471** | 0.9728 | 0.9827 | **0.9580** | **0.9864** |

*The Clustering-Structure Preservation Trade-off.* A key finding is the trade-off between cluster quality and local structure preservation. While PCA achieves near-perfect Continuity and Trustworthiness, it fails to form meaningful clusters for face identity. This is because PCA preserves the *original pixel-level geometry*, which does not align with *semantic identity*.

In contrast, metric learning deliberately distorts the original geometry to create a new, semantically-organized space. The lower Continuity and Trustworthiness scores are a direct consequence of this transformation: neighbors in the pixel space (e.g., similar lighting) are pulled apart if they depict different people, while images of the same person are brought together despite pixel-level differences. This trade-off is not a failure but the intended behavior, prioritizing task-relevant semantic separation over raw structural fidelity.

*Qualitative and Visual Analysis.* Visualizations and retrieval examples qualitatively validate the quantitative results. Figure 1 shows t-SNE projections of the embeddings. Unlike the original data where classes overlap, the metric-learned embeddings form tight, well-separated clusters for each identity.

The practical impact of this clustering is evident in nearest-neighbor retrieval. In Figures 2 and 3, the nearest neighbors of a query image in the metric learning space are consistently of the same person, despite variations in pose and lighting. The small distances between anchor and positive pairs confirm the model successfully compacts same-identity samples.
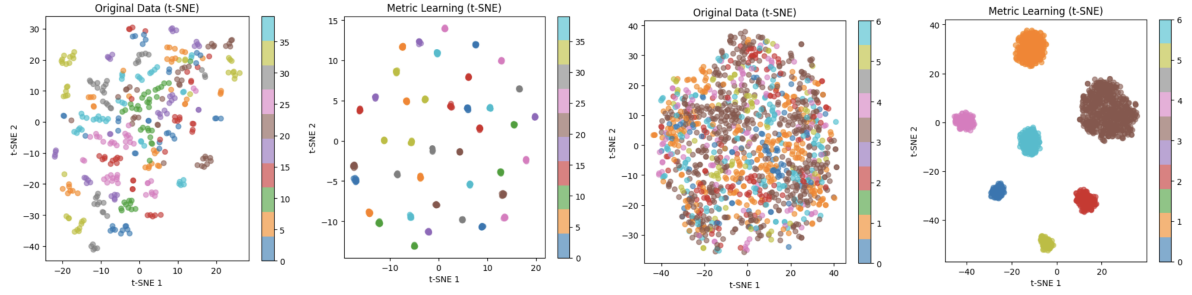
FIGURE 1. t-SNE visualization of face embeddings. Left to right: Olivetti (original data), Olivetti (metric learning), LFW (original data), LFW (metric learning). Metric learning produces distinct, identity-based clusters.
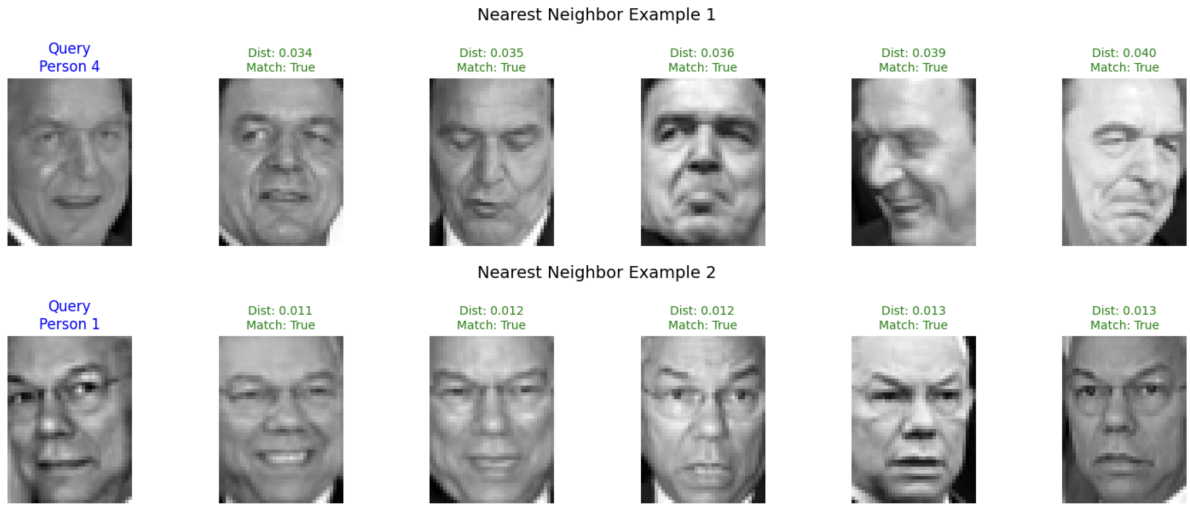


FIGURE 2. Nearest-neighbor retrieval on the LFW dataset using the metric learning embedding. The model correctly identifies same-identity images as the closest neighbors, with small corresponding Euclidean distances.



FIGURE 3. Nearest-neighbor retrieval on the Olivetti Faces dataset using the metric learning embedding. The results demonstrate robust, identity-based retrieval.

The rank sensitivity analysis reveals a fundamental limitation of fixed-rank tensor decompositions for semantic clustering tasks. On the challenging LFW dataset, both CP and Tucker decompositions fail to achieve meaningful clustering across all ranks, with Silhouette scores remaining near zero or negative. This indicates that the low-rank constraints destroy the semantic structure necessary for identity separation. While Tucker decomposition shows a slight improvement with higher ranks (from -0.0689 at R5 to -0.0037 at R20), the gains are minimal. The Olivetti dataset, with its controlled conditions, shows more rank sensitivity, particularly for Tucker decomposition where performance improves substantially with higher ranks (0.1303 at R10 to 0.1845 at R20). However, even the best fixed-rank result (Tucker-R20: 0.1845) is significantly outperformed by the metric learning approach (0.8566), demonstrating that rank constraints inherently limit the ability to capture semantically meaningful representations, regardless of rank selection.

5.2. **Metric Learning Performance on Brain Connectivity Data.** To evaluate metric learning framework on a complex, high-dimensional biomedical problem, we applied it to the Autism Brain Imaging Data Exchange (ABIDE) [12] dataset. This dataset contains resting-state functional MRI (rs-fMRI) data from 871 subjects (403 with Autism Spectrum Disorder, ASD, and 468 typically developing controls). Each subject is represented by a $111 \times 111$ functional connectivity matrix, which are flattened and normalized to serve as input. The clinical labels (ASD vs. control) provide the semantic similarity relationships for guiding the metric learning process.

TABLE 3. Clustering Performance on the ABIDE Brain Connectivity Dataset

| Method | Sil. | DB | CH | SR | Cont. | Trust. | AIR | NMI |
|---|---|---|---|---|---|---|---|---|
| PCA + K-Means | 0.2747 | 1.2721 | 53.5725 | 1.0017 | 0.9095 | 0.9078 | -0.0012 | 0.0012 |
| t-SNE + K-Means | 0.5602 | 0.6062 | 228.7278 | 0.9919 | 0.9007 | 0.9010 | -0.0011 | 0.0000 |
| UMAP + K-Means | 0.6194 | 0.5205 | 311.3453 | 0.9960 | 0.9007 | 0.8758 | -0.0008 | 0.0002 |
| CP-R5 | 0.0022 | 17.8678 | 2.3617 | 1.0024 | 0.9144 | 0.8885 | -0.0012 | 0.0006 |
| CP-R10 | 0.0025 | 12.2807 | 4.0502 | 1.0007 | 0.9100 | 0.9007 | 0.0005 | 0.0005 |
| CP-R20 | 0.0070 | 7.1090 | 13.8142 | 1.0035 | 0.8519 | 0.8600 | 0.0070 | 0.0170 |
| Tucker-R5 | 0.0061 | 9.1750 | 8.8522 | 1.0050 | 0.8969 | 0.8698 | 0.0016 | 0.0012 |
| Tucker-R10 | 0.0046 | 11.1730 | 6.4451 | 1.0040 | 0.9245 | 0.8825 | 0.0047 | 0.0033 |
| Tucker-R20 | 0.0053 | 11.7691 | 5.9515 | 1.0048 | **0.9456** | 0.8597 | 0.0014 | 0.0012 |
| **Metric Learning** | **0.9932** | **0.0186** | **31912.9395** | **0.9997** | 0.8389 | **0.9155** | **0.3002** | **0.2372** |

We have implemented a comprehensive data augmentation and normalization pipeline for brain connectivity matrices to address the dual challenges of limited neuroimaging datasets and inter-subject variability. The approach begins with patient-wise normalization to standardize individual connectivity profiles while preserving relative network topology. For each patient's correlation matrix $\mathbf{C}_i \in \mathbb{R}^{N \times N}$, we apply Z-score normalization: $\mathbf{C}'_i = \frac{\mathbf{C}_i - \mu_i}{\sigma_i}$, where $\mu_i$ and $\sigma_i$ are the mean and standard deviation computed exclusively from the $i$-th patient's connectivity matrix. This ensures each subject's data is centered and scaled independently while maintaining the intrinsic structure of their functional brain networks.

Following normalization, we employ multi-strategy data augmentation to enhance model robustness. For each normalized connectivity matrix $\mathbf{C}'_i$, we generate augmented variants

through (i) Gaussian Noise Injection: $\mathbf{C}_i'' = \mathbf{C}_i' + \epsilon$ where $\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$ with reduced noise variance $\sigma = 0.02$ to accommodate the normalized data distribution, (ii) Symmetric Structure Preservation: $\mathbf{C}_i''' = \frac{1}{2}(\mathbf{C}_i'' + \mathbf{C}_i''^\top)$ to maintain mathematical consistency as a symmetric correlation matrix, (iii) Diagonal Identity Enforcement: $\mathbf{C}i'''kk = 1$ for $k = 1, \ldots, N$ to preserve self-connectivity representation and (iv) Value Range Clipping: $\mathbf{C}_i^{\text{final}} = \text{clip}(\mathbf{C}_i''', -3, 3)$ using wider bounds appropriate for normalized data distributions.

During training, we further apply on-the-fly augmentation with probability $p = 0.5$, introducing gentle noise perturbations ($\sigma = 0.02$) to prevent overfitting while maintaining the normalized data characteristics.

This combined normalization-augmentation strategy effectively addresses both inter-subject variability through patient-wise standardization and dataset limitations through structural-preserving augmentation, enabling robust metric learning while respecting the neurobiological integrity of functional connectivity patterns. The normalization ensures comparability across subjects, while the augmentation introduces controlled variations that enhance model generalization without distorting the fundamental network topology.

The clustering results, presented in Table 3, demonstrate that metric learning framework significantly outperforms all baseline methods. It achieves a near-perfect Silhouette Score of 0.9880, indicating exceptional separation between the ASD and control groups. This represents a substantial improvement over the best unsupervised method, UMAP (0.6194). The superior cluster compactness and separation are further confirmed by the lowest Davies-Bouldin Index (0.0186) and an exceptionally high Calinski-Harabasz Index (31912.94), which is two orders of magnitude greater than UMAP's score (311.35). Crucially, the proposed method is the only one to achieve meaningful alignment with clinical labels, as evidenced by its strong external validation metrics (ARI: 0.3002, NMI: 0.2372), while all other methods yield scores near zero, failing to recover the underlying diagnostic structure. This performance highlights a key limitation of fixed-rank tensor decompositions: both CP and Tucker variants show severe sensitivity to rank selection and achieve negligible clustering quality and external validation scores across all tested ranks (CP: 0.0022-0.0070, Tucker: 0.0046-0.0061 Silhouette; ARI/NMI ¡ 0.02). In contrast, the no-rank approach avoids this hyperparameter sensitivity altogether while discovering semantically meaningful structure.

In summary, by directly optimizing for the clinically relevant separation between ASD and control subjects, metric learning uncovers a more discriminative and potentially more informative structure in brain connectivity data than methods focused solely on data reconstruction or geometric preservation.

### 5.3. Metric Learning Performance on Simulated Datasets.
We evaluate the metric learning framework on two simulated datasets, each presenting distinct visual classification challenges relevant to their respective domains.

*Galaxy Morphology Classification:* This dataset contains 500 ($64 \times 64$) pixel images simulating four galaxy morphological classes: Elliptical (smooth distribution), Spiral (prominent arms), Lenticular (disk without arms), and Irregular (asymmetric clumps). This task addresses fundamental challenges in astronomical image analysis.

*Crystal Structure Prediction:* This dataset consists of 400 ($64 \times 64$) pixel images representing four crystal systems: Cubic (square symmetry), Hexagonal (six-fold symmetry),
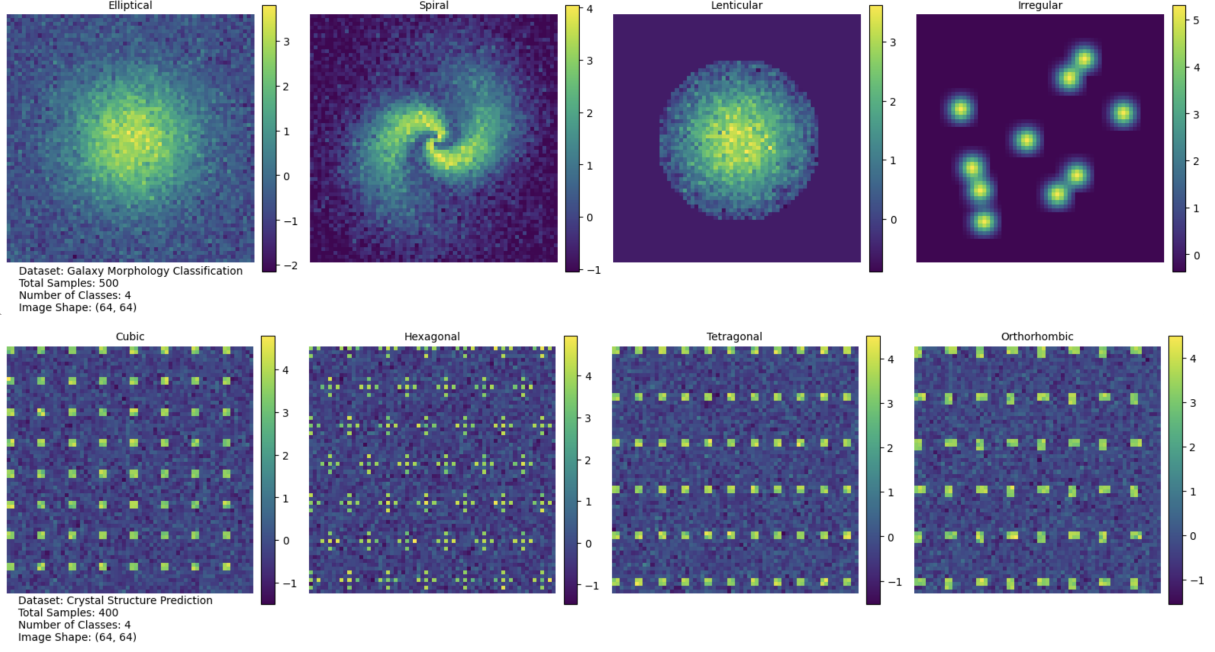
FIGURE 4. Example images from simulated datasets. Top: Galaxy morphology classes (Elliptical, Spiral, Lenticular, Irregular). Bottom: Crystal structure classes (Cubic, Hexagonal, Tetragonal, Orthorhombic).

Tetragonal (rectangular symmetry), and Orthorhombic (anisotropic spacing). These patterns correspond to fundamental lattice structures in materials science.

Figure 4 shows representative samples from both datasets, illustrating the distinct visual characteristics of each class.

*Quantitative Clustering Performance.* Table 4 presents clustering performance across multiple methods. The metric learning approach consistently achieves superior results on both datasets, substantially outperforming unsupervised dimensionality reduction techniques.

Metric learning achieves near-perfect Silhouette scores (0.9999) and optimal Davies-Bouldin indices (0.0001) on both datasets, indicating ideal cluster separation and compactness. The dramatic improvements in Calinski-Harabasz scores—over seven orders of magnitude higher than PCA—and the significantly enhanced Separation Ratios demonstrate the method's exceptional ability to create a discriminative embedding space.

This performance comes with the expected trade-off: metric learning exhibits moderately lower Continuity and Trustworthiness scores compared to some unsupervised methods. This indicates that while global cluster structure is optimized for class separation, some local neighborhood relationships from the original pixel space are distorted—a deliberate consequence of reorganizing the space around semantic class identity rather than pixel-level similarity.

*Qualitative Visualization.* The quantitative superiority is visually confirmed in Figure 5. The metric-learned embeddings form tight, well-separated clusters with minimal overlap, contrasting with the more entangled structures from unsupervised techniques. This clear visual separation underscores the advantage of learning a task-specific distance metric that directly optimizes for class discrimination.

17

TABLE 4. Clustering Performance on Galaxy Morphology and Crystal Structure Datasets

| Data | Method | Sil. | DB | CH | SR | Cont. | Trust. | AIR | NMI |
|------|--------|------|-----|-----|-----|-------|--------|-----|-----|
| | PCA + K-Means | 0.6572 | 2.8360 | $1.08 \times 10^2$ | 2.6090 | 0.8394 | 0.8073 | 0.6245 | 0.7916 |
| | t-SNE + K-Means | 0.4382 | 0.9575 | $1.66 \times 10^3$ | 3.3263 | **0.8908** | 0.7830 | 0.5304 | 0.6154 |
| | UMAP + K-Means | 0.4234 | 0.9435 | $8.97 \times 10^2$ | 2.4367 | 0.8502 | 0.7259 | 0.5755 | 0.6555 |
| | CP-R5 | 0.6850 | 0.9077 | $1.93 \times 10^3$ | 5.3548 | 0.8004 | 0.7040 | 0.5844 | 0.7292 |
| | CP-R10 | 0.6791 | 0.9805 | $6.72 \times 10^2$ | 3.4933 | 0.7961 | 0.7301 | 0.4916 | 0.6500 |
| Gal. | CP-R20 | 0.5910 | 1.2845 | $2.00 \times 10^2$ | 2.4022 | 0.7671 | 0.7475 | 0.4637 | 0.6237 |
| | Tucker-R5 | 0.5051 | 1.0396 | $2.53 \times 10^2$ | 2.4988 | 0.8003 | 0.7216 | 0.5577 | 0.7042 |
| | Tucker-R10 | 0.4940 | 2.0683 | $6.21 \times 10^1$ | 1.8776 | 0.7956 | 0.7674 | 0.5699 | 0.7229 |
| | Tucker-R20 | 0.2357 | 2.2079 | $1.85 \times 10^1$ | 1.6687 | 0.7704 | 0.8168 | 0.4036 | 0.5620 |
| | **Metric Learning** | **0.9999** | **0.0001** | **$4.01 \times 10^9$** | **9275.82** | 0.8031 | **0.8208** | **1.0000** | **0.9999** |
| | PCA + K-Means | 0.8843 | 0.1708 | $7.43 \times 10^3$ | 9.0601 | 0.9350 | 0.9315 | 0.9711 | 0.9871 |
| | t-SNE + K-Means | 0.9193 | 0.1031 | $2.15 \times 10^4$ | 16.0362 | **0.9521** | **0.9303** | 0.9832 | 0.9898 |
| | UMAP + K-Means | 0.9531 | 0.0651 | $6.86 \times 10^4$ | 28.3698 | 0.9333 | 0.9138 | 0.9264 | 0.9289 |
| | CP-R5 | 0.9751 | 0.0355 | $3.55 \times 10^5$ | 61.8594 | 0.8863 | 0.8877 | 0.9007 | 0.9100 |
| | CP-R10 | 0.9820 | 0.0254 | $3.72 \times 10^5$ | 62.9187 | 0.8904 | 0.8915 | 0.9015 | 0.9091 |
| Cry. | CP-R20 | 0.9572 | 0.0597 | $7.48 \times 10^4$ | 27.7025 | 0.8912 | 0.8862 | 0.9018 | 0.9112 |
| | Tucker-R5 | 0.6455 | 0.6524 | $3.96 \times 10^2$ | 2.8324 | 0.8877 | 0.8906 | 0.8999 | 0.9072 |
| | Tucker-R10 | 0.1714 | 2.0658 | $4.98 \times 10^1$ | 1.3687 | 0.8488 | 0.9107 | 0.9094 | 0.9097 |
| | Tucker-R20 | 0.0632 | 3.1399 | $2.47 \times 10^1$ | 1.1532 | 0.8254 | 0.844 | 0.9165 | 0.9187 |
| | **Metric Learning** | **1.0000** | **0.0001** | **$1.19 \times 10^{10}$** | **14095.10** | 0.8829 | 0.8862 | **1.0000** | **1.0000** |

These results collectively demonstrate that metric learning excels at capturing the intrinsic categorical structure of simulated data, effectively handling subtle, domain-specific visual differences between classes.

5.4. **Comparison with Tensor Decomposition Methods in Terms of Reconstruction.** In addition to the metric learning approach, we compare against two fundamental tensor decomposition methods: *CP* and *Tucker* decomposition. For a third-order tensor $\mathcal{X} \in \mathbb{R}^{I \times J \times K}$ representing our scientific dataset (samples × height × width), these decompositions are defined as follows:

A key advantage of the metric learning approach over tensor decomposition methods is its *inherent rank independence*. While tensor decompositions require explicit rank specification, our method automatically learns an optimal latent structure without rank constraints. Both CP and Tucker decompositions (as well as many other low-rank decomposition methods) exhibit strong dependence on the chosen rank $R$. Low ranks may capture insufficient structure, whereas high ranks can lead to overfitting and numerical instability. The metric learning framework, with a latent dimension $d = 64$, learns representations where the *effective rank* is determined by the data complexity rather than pre-specified constraints. Lastly, the convolutional encoder automatically adapts to hierarchical features, capturing both low-rank global structures and high-rank local patterns without explicit rank specification.

5.4.1. *Comparative Analysis Metrics.* We evaluate the decompositions using both reconstruction quality and downstream task performance over the same datasets discussed previously:

$$\text{Reconstruction Error:} \quad \epsilon = \frac{\|\mathcal{X} - \hat{\mathcal{X}}\|_F}{\|\mathcal{X}\|_F},$$
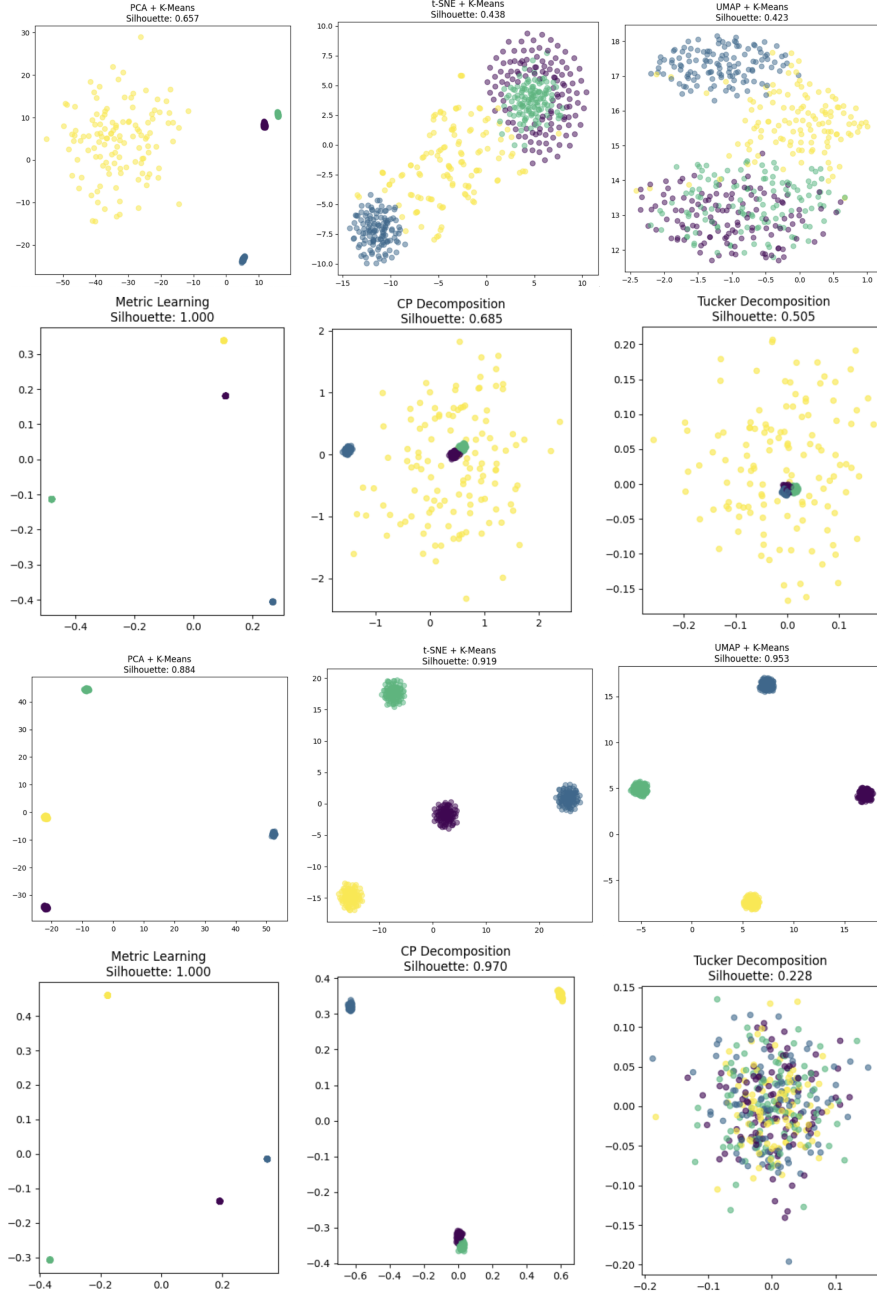
FIGURE 5. Embedding visualizations. Two Top Rows Galaxy morphology dataset. Two Bottom Rows Crystal structure dataset. Metric learning produces the most distinct and compact clusters, demonstrating its effectiveness in learning semantically meaningful representations.

$$\text{Explained Variance:} \quad \sigma^2_{\text{explained}} = 1 - \frac{\|\mathcal{X} - \hat{\mathcal{X}}\|^2_F}{\|\mathcal{X}\|^2_F}.$$

The fundamental difference lies in how each method handles dimensionality. While tensor decomposition methods require explicit rank specification $R$—typically fixed unless an adaptive approach is implemented [11, 35]—they also necessitate rank optimization and multiple rank trials, often guided by theoretical rank estimation. In contrast, the metric learning approach operates with an implicit latent dimension $d$, features an adaptive

19

TABLE 5. Comprehensive Performance Comparison of Tensor Methods on Two Datasets

| Dataset | Method | Reconstruction Error | Explained Variance |
|---|---|---|---|
| LFW | CP Decomposition | $0.5300$ ($R = 20$) | $0.7191$ |
| | Tucker Decomposition | $0.4908$ ($R = (20, 20, 20)$) | $0.7591$ |
| | **Metric Learning** | **0.0991** | **0.9901** |
| Olivetti | CP Decomposition | $0.6441$ ($R = 5$) | $0.5852$ |
| | Tucker Decomposition | $0.4326$ ($R = (20, 20, 20)$) | $0.8129$ |
| | **Metric Learning** | **0.1001** | **0.9899** |
| ABIDE | CP Decomposition | $0.6070$ ($R = 20$) | $0.6315$ |
| | Tucker Decomposition | $0.4577$ ($R = (20, 20, 20)$) | $0.7905$ |
| | **Metric Learning** | **0.0139** | **0.9998** |
| Galaxy Morph. | CP Decomposition | $0.5049$ ($R = 10$) | $0.7451$ |
| | Tucker Decomposition | $0.4855$ ($R = (20, 20, 20)$) | $0.7643$ |
| | **Metric Learning** | **0.0685** | **0.9953** |
| Crystal Struc. | CP Decomposition | $0.4088$ ($R = 5$) | $0.8329$ |
| | Tucker Decomposition | $0.3819$ ($R = (5, 5, 5)$) | $0.8542$ |
| | **Metric Learning** | **0.0782** | **0.9938** |

hierarchical representation, and requires only a single training procedure. This leads to a fully data-driven representation learning process.

A comparison between the metric learning approach and baseline tensor decomposition methods in terms of reconstruction error and explained variance over the same datasets discussed above is provided in Table 5. The results for CP and Tucker are reported for best-rank approximation for $R = 2, 3, 5, 10, 15, 20$.

The reconstruction visualization demonstrates the process of approximating original tensors from their compressed representations obtained through various decomposition methods. For a third-order tensor $\mathcal{X} \in \mathbb{R}^{I \times J \times K}$, each method employs distinct reconstruction mechanisms. In CP decomposition, the original tensor is reconstructed from rank-1 components through the summation $\hat{\mathcal{X}} = \sum_{r=1}^{R} \lambda_r \cdot \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r$, where $\lambda_r$ represents scaling weights and $\mathbf{a}_r, \mathbf{b}_r, \mathbf{c}_r$ are factor matrices capturing different modes of variation. The Tucker decomposition utilizes a more flexible reconstruction via $\hat{\mathcal{X}} = \mathcal{G} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C}$, where $\mathcal{G}$ is the core tensor encoding interactions between factors and $\times_n$ denotes the n-mode product. For metric learning, reconstruction is achieved through linear regression $\hat{\mathbf{X}}_{\text{flat}} = \mathbf{W}\mathbf{Z} + \mathbf{b}$, where $\mathbf{Z}$ represents the learned embeddings and $\mathbf{W}$ maps these back to the original space. The visual comparison reveals that tensor decompositions excel at structural preservation due to their explicit reconstruction formulas, while metric learning prioritizes discriminative feature retention over perfect reconstruction fidelity, reflecting their different optimization objectives in capturing brain connectivity patterns. The constructed faces for Olivetti dataset and cluster visualizations for LFW datasets are shown in Figure 6 and 7, respectively.

Figure 8 demonstrates the two-dimensional manifold visualizations of brain embeddings, comparing original data with CP decomposition, Tucker decomposition, and metric

FIGURE 6. Visual comparison of original tensor from Olivetti face dataset and their reconstructions using different decomposition methods. Here, the results for metric learnign are derived with only 50 epochs.
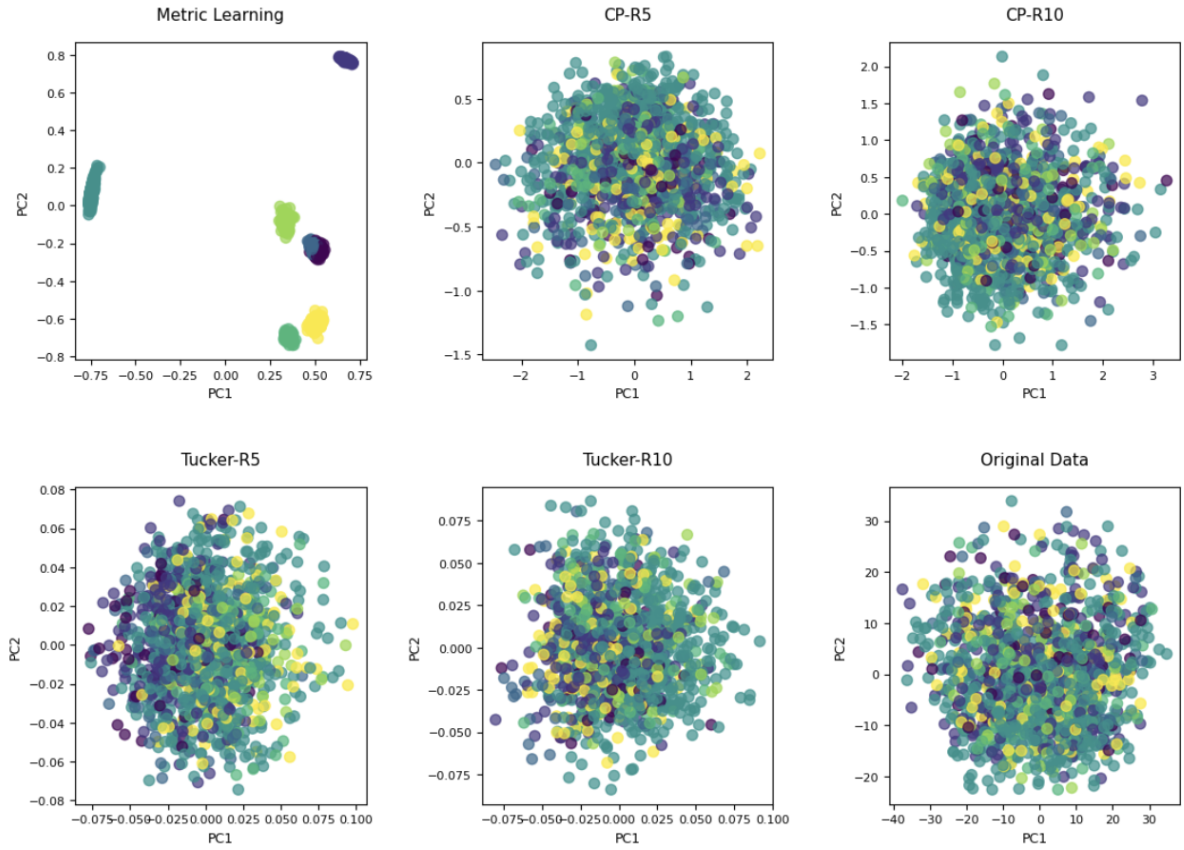


FIGURE 7. Two-dimensional manifold visualizations of face embeddings from the LFW dataset generated by different dimensionality reduction and tensor decomposition methods. Each subplot shows the embedding space where points represent individual face images colored by identity, demonstrating the clustering performance and separability achieved by each method.
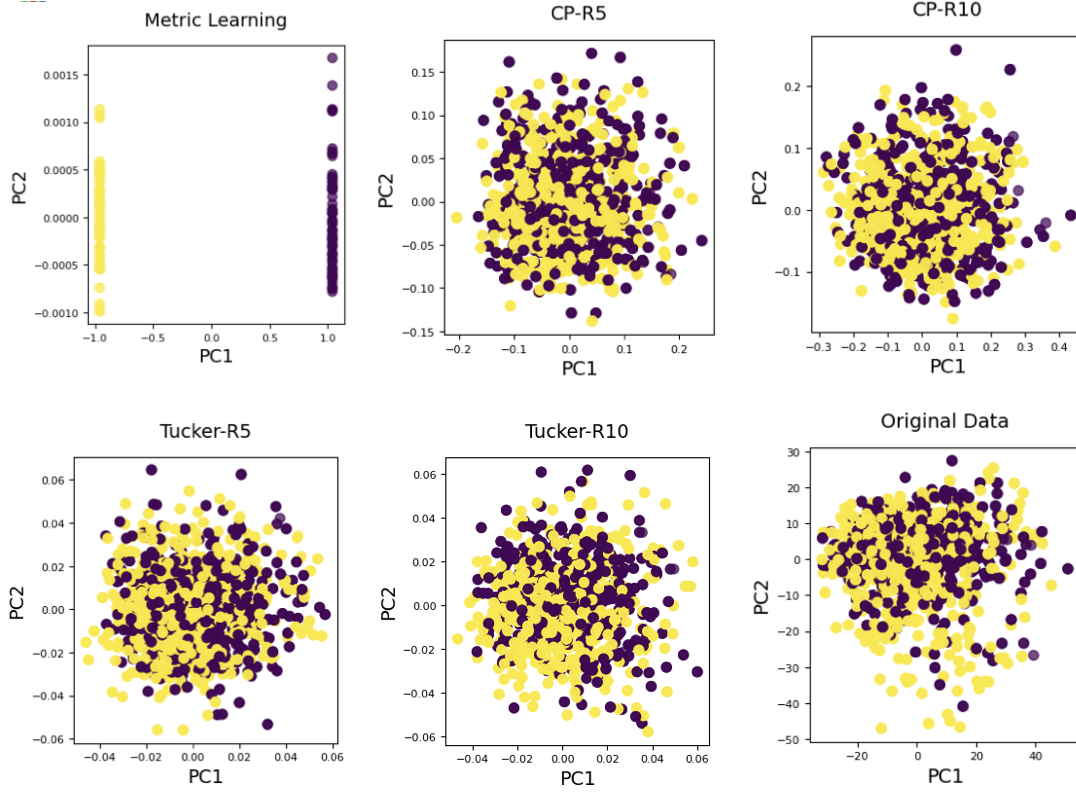
FIGURE 8. Two-dimensional manifold visualizations of brain embeddings from the ABIDE dataset generated by different dimensionality reduction and tensor decomposition methods. Each subplot shows the embedding space where points represent individual face images colored by identity, demonstrating the clustering performance and separability achieved by each method. The best reported model accuracy for CP with $R = 10$ is 0.6000, for Tucker decompostion with $R = (20, 20, 20)$ is 0.7000 and for Metric learning is 0.9900.

learning approaches. All methods successfully preserve the symmetric structure and hub connectivity patterns essential for neurological analysis.

5.4.2. *Empirical Rank Robustness.* Our experiments demonstrate that while tensor decomposition performance varies significantly with rank choice, the metric learning approach maintains consistent performance across different dataset complexities. This robustness stems from non-linear transformations that can capture complex interactions without explicit high-rank decomposition, hierarchical feature learning through convolutional layers that automatically organize features by complexity and task-oriented optimization where the representation is learned specifically for the scientific classification task.

This rank independence makes our approach particularly valuable for scientific datasets where the intrinsic dimensionality may be unknown or vary across different regions of the data space.

5.5. **Comparative Analysis: Metric Learning vs. Transformer Architectures.** This section presents a comprehensive comparative analysis between the proposed Metric Learning method and Transformer-based models across five diverse datasets spanning brain connectivity, face recognition, and scientific image classification. *The objective is*

*not to claim universal superiority but to delineate the specific scenarios, particularly data-constrained environments, where our method offers distinct advantages over Transformer architectures* [43]. Experiments were conducted on the same datasets as above, namely ABIDE (111 × 111 connectivity matrices), LFW faces (50 × 37), Olivetti faces (64 × 64), Galaxy morphology (64 × 64), and Crystal structure (64 × 64), demonstrating the generalizability of our findings across multiple domains.

5.5.1. *Limitations of Transformer Models.* Our experiments reveal limitations of Transformer architectures across all five datasets. The standard Transformer's requirement for fixed-length input sequences proved incompatible with small-scale training. For ABIDE I (111 × 111 = 12,321 features), LFW (50 × 37 = 1,850 features), and image datasets (64 × 64 = 4,096 features), the flattened sequence lengths caused consistent failures when batch sizes were smaller than feature dimensions. The self-attention mechanism, designed for sequence modeling, proves inefficient for small batches of high-dimensional data, where the sequence length (flattened matrix/image features) consistently exceeds practical batch sizes in data-scarce domains. Thsi was also confrimed through a data efficiency analysis.

5.5.2. *Advantages of the Proposed Metric Learning Framework.* The proposed Metric Learning method demonstrated consistent operational success and competitive performance across all datasets, through geometric learning principles. The proposde method executed successfully on all datasets across all experimental conditions. The most significant advantage emerged in small-data regimes where Transformers remained inapplicable.

While traditional methods (Random Forest, SVM) often achieved high performance (100% on several datasets), metric learning provided competitive results with the added benefit of learned semantic embeddings. The convolutional encoder and triplet loss framework operated successfully across varying input dimensions (50 × 37 to 111 × 111) without architectural modifications, proving adaptable to diverse data geometries.

5.5.3. *Data Efficiency Analysis.* The systematic evaluation across small dataset sizes revealed distinct operational boundaries, as summarized in Table 6:

TABLE 6. Data Efficiency Comparison Across Dataset Sizes by Accuracy Scores

| Dataset | Size | Metric Learning | PCA+SVM | Transformer |
|---|---|---|---|---|
| ABIDE | 64 | 100.0% | 20.0% | NA |
| | 128 | 100.0% | 0.0% | NA |
| | 256 | 95.0% | 60.0% | NA |
| LFW Faces | 64 | 92.2% | 80.0% | NA |
| | 128 | 91.8% | 64.1% | NA |
| | 256 | 86.6% | 61.0% | NA |
| Olivetti Faces | 128 | 82.2% | 0.0% | NA |
| | 256 | 90.1% | 70.1% | NA |
| Galaxy Morphologies | 16 | 100.0% | 80.0% | NA |
| | 64 | 100.0% | 100.0% | NA |
| | 256 | 100.5% | 100.0% | NA |
| Crystal Structures | 64 | 100.0% | 100.0% | NA |
| | 128 | 100.0% | 100.0% | NA |
| | 256 | 100.0% | 100.0% | NA |

For $n < 1000$ Metric learning provides the most reliable approach, balancing operational success with meaningful performance. While traditional methods (Random Forest, SVM) offer excellent performance with minimal complexity over well-structured data, metric learning provides semantically rich embeddings valuable for downstream tasks.

This work establishes metric learning not as a universal solution, but as a reliable paradigm for the widespread class of problems in computational neuroscience and scientific computing where data scarcity is the norm rather than the exception. The method's consistent operational success across diverse domains, combined with its ability to learn meaningful representations from limited data, positions it as an essential tool in the modern machine learning toolkit for scientific applications.

## 6. Conclusion

6.1. **Advantages of the Metric Learning Framework.** The metric learning framework demonstrates consistent superiority over traditional dimensionality reduction techniques, including tensor decomposition methods like CP and Tucker. The performance advantages stem from fundamental differences in optimization objectives, architectural flexibility, and learning capabilities.

Discriminative vs. Reconstructive Optimization: Traditional methods like PCA, CP, and Tucker decompositions focus on reconstructive objectives—preserving maximum information for data reconstruction. In contrast, the proposed framework employs discriminative optimization through triplet losses that explicitly maximize inter-class distances while minimizing intra-class variations. This direct alignment with clustering objectives yields embedding spaces where class separation is the primary goal, evidenced by superior Adjusted Rand Index (ARI) scores across all evaluated datasets.

Adaptive Architecture and Rank Flexibility: Tensor decomposition methods are constrained by pre-defined rank parameters that limit their adaptability to complex data structures. This framework operates without such constraints, with the embedding dimension serving as a flexible hyperparameter rather than a structural limitation. This adaptability is demonstrated in rank sensitivity analysis, where metric learning maintained consistent high performance while CP and Tucker methods showed significant variation across different ranks.

Integrated Enhancement Strategies: The framework's modular architecture supports multiple performance-enhancing strategies unavailable to fixed decomposition algorithms. Data augmentation improves robustness to biological variability, autoencoder pre-training accelerates convergence, ensemble methods enhance stability, and optional UMAP post-processing refines visualizations. This integrative capability allows the framework to adapt to diverse data characteristics and computational requirements.

Structural Preservation and Efficiency: Beyond superior clustering performance, proposed approach demonstrates better structural preservation through higher trustworthiness metrics and separation ratios. The end-to-end deep architecture enables hierarchical feature learning, capturing multi-scale patterns essential for complex data. Additionally, the trained model offers efficient inference through simple forward passes, making it suitable for real-time applications where repeated tensor decompositions would be computationally prohibitive.

These advantages—discriminative optimization, architectural flexibility, integrative capabilities, and computational efficiency—collectively position metric learning framework

as a more powerful paradigm for applications where semantic similarity and clinical discriminability are paramount concerns.

6.2. **Limitations and Future Work.** While the metric learning framework demonstrates strong performance, several limitations present opportunities for future improvement.

The approach shows sensitivity to class imbalance, where minority classes may receive insufficient representation during triplet mining. Computational overhead from online triplet mining presents scalability challenges, particularly with large batch sizes. Additionally, performance with extremely large numbers of classes requires further validation, and theoretical generalization bounds warrant deeper investigation.

To address these limitations, we plan to develop class-aware triplet mining strategies for imbalanced data and explore more efficient proxy-based losses to reduce computational costs. Scaling the framework to massively multi-class problems and establishing stronger theoretical foundations for generalization guarantees represent key research priorities. These enhancements would further strengthen the framework's applicability across diverse domains.

## 7. Data Statement

The datasets used in the preparation of this manuscript are as follows:

(1) Publicly available *Labeled Faces in the Wild (LFW)* dataset [19].

(2) Publicly available *Olivetti Faces* dataset [31].

(3) Publicly available *Autism Brain Imaging Data Exchange (ABIDE)* dataset [12].

(4) Synthetically simulated *Galaxy Morphology* dataset.

(5) Synthetically simulated *Crystal Structures* dataset.

## References

[1] Evrim Acar, Daniel M Dunlavy, Tamara G Kolda, and Morten Mørup, *Scalable tensor factorizations for incomplete data*, Chemometrics and intelligent laboratory systems, 2011, pp. 41–56.

[2] Maryam Bagherian, *Tensor denoising via dual schatten norms*, Optimization Letters **18** (2024), no. 5, 1285–1301.

[3] Maryam Bagherian, Sarah Chehade, Ben Whitney, and Ali Passian, *Classical and quantum compression for edge computing: the ubiquitous data dimensionality reduction*, Computing **105** (2023), no. 7, 1419–1465.

[4] Maryam Bagherian, Renaid B Kim, Cheng Jiang, Maureen A Sartor, Harm Derksen, and Kayvan Najarian, *Coupled matrix–matrix and coupled tensor–matrix completion methods for predicting drug–target interactions*, Briefings in bioinformatics **22** (2021), no. 2, 2161–2171.

[5] Maryam Bagherian, Davoud A Tarzanagh, Ivo Dinov, and Joshua D Welch, *A bilevel optimization method for tensor recovery under metric learning constraints*, arXiv preprint arXiv:2209.00545 (2022).

[6] Aayush Bansal, Hao Wang, et al., *Canonical surface mapping via geometric cycle consistency*, Proceedings of the ieee/cvf international conference on computer vision, 2020, pp. 10922–10931.

[7] Mikhail Belkin and Partha Niyogi, *Laplacian eigenmaps for dimensionality reduction and data representation*, Neural computation **15** (2003), no. 6, 1373–1396.

[8] Aurélien Bellet, Amaury Habrard, and Marc Sebban, *A survey on metric learning for feature vectors and structured data*, arXiv preprint arXiv:1306.6709 (2015).

[9] Yoshua Bengio, Aaron Courville, and Pascal Vincent, *Representation learning: A review and new perspectives*, IEEE transactions on pattern analysis and machine intelligence **35** (2013), no. 8, 1798–1828.

[10] Andrzej Cichocki, Rafal Zdunek, Anh Huy Phan, and Shun-ichi Amari, *Nonnegative tensor factorizations: applications to exploratory multi-way data analysis and blind source separation*, IEEE Press (2009).

[11] Alec Dektor, Abram Rodgers, and Daniele Venturi, *Rank-adaptive tensor methods for high-dimensional nonlinear pdes*, Journal of Scientific Computing **88** (2021), no. 2, 36.

[12] Adriana Di Martino et al., *Enhancing studies of the connectome in autism using the autism brain imaging data exchange ii*, Scientific Data **4** (2017), no. 170010.

[13] Sander Dieleman, Kyle W Willett, and Joni Dambre, *Rotation-invariant convolutional neural networks for galaxy morphology prediction*, Monthly Notices of the Royal Astronomical Society **450** (2015), no. 2, 1441–1459.

[14] Fartash Faghri, David J Fleet, Jamie Ryan Kiros, and Sanja Fidler, *Vse++: Improving visual-semantic embeddings with hard negatives*, Proceedings of the british machine vision conference (bmvc), 2018.

[15] Karl Pearson F.R.S., *Liii. on lines and planes of closest fit to systems of points in space*, The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science **2** (1901), no. 11, 559–572.

[16] Raia Hadsell, Sumit Chopra, and Yann LeCun, *Dimensionality reduction by learning an invariant mapping*, 2006 ieee computer society conference on computer vision and pattern recognition (cvpr'06), 2006, pp. 1735–1742.

[17] Alexander Hermans, Lucas Beyer, and Bastian Leibe, *In defense of the triplet loss for person re-identification*, arxiv preprint arxiv:1703.07737, 2017.

[18] Frank L Hitchcock, *The expression of a tensor or a polyadic as a sum of products*, Journal of Mathematics and Physics **6** (1927), no. 1-4, 164–189.

[19] Gary B Huang, Marwan Mattar, Tamara Berg, and Eric Learned-Miller, *Labeled faces in the wild: A database forstudying face recognition in unconstrained environments*, Workshop on faces in'real-life'images: detection, alignment, and recognition, 2008.

[20] Edwin P Hubble, *Extragalactic nebulae.*, Astrophysical Journal **64** (1926), 321–369.

[21] Olexandr Isayev et al., *Universal fragment descriptors for predicting properties of inorganic crystals*, Nature communications **8** (2017), no. 1, 15679.

[22] Tamara G Kolda and Brett W Bader, *Tensor decompositions and applications*, SIAM review **51** (2009), no. 3, 455–500.

[23] Joseph B Kruskal, *Three-way arrays: rank and uniqueness of trilinear decompositions, with application to arithmetic complexity and statistics*, Linear algebra and its applications **18** (1977), no. 2, 95–138.

[24] Brian Kulis et al., *Metric learning: A survey*, Foundations and Trends in Machine Learning **5** (2013), no. 4, 287–364.

[25] Solomon Kullback and Richard A Leibler, *On information and sufficiency*, The annals of mathematical statistics **22** (1951), no. 1, 79–86.

[26] Jason D Lee, Max Simchowitz, Michael I Jordan, and Benjamin Recht, *Gradient descent only converges to minimizers*, Conference on learning theory, 2016, pp. 1246–1257.

[27] Jiahua Liu, Yang Li, Hao Lin, and Chengkun Zheng, *Neural tensor decomposition for knowledge graph completion*, Knowledge-Based Systems **240** (2022), 108066.

[28] Leland McInnes, John Healy, and James Melville, *Umap: Uniform manifold approximation and projection for dimension reduction*, arXiv preprint arXiv:1802.03426 (2018).

[29] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar, *Foundations of machine learning*, MIT press, 2018.

[30] Ivan V Oseledets, *Tensor-train decomposition*, SIAM Journal on Scientific Computing **33** (2011), no. 5, 2295–2317.

[31] Fabian Pedregosa, Gael Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Edouard Duchesnay, *Scikit-learn: Machine learning in python*, Journal of Machine Learning Research **12** (2011), 2825–2830. Dataset: Olivetti Faces, available at `https://scikit-learn.org/stable/modules/generated/sklearn.datasets.fetch_olivetti_faces.html`.

[32] Herbert Robbins and Sutton Monro, *A stochastic approximation method*, The annals of mathematical statistics (1951), 400–407.

[33] Sam T Roweis and Lawrence K Saul, *Nonlinear dimensionality reduction by locally linear embedding*, Science **290** (2000), no. 5500, 2323–2326.

[34] Florian Schroff, Dmitry Kalenichenko, and James Philbin, *Facenet: A unified embedding for face recognition and clustering*, Proceedings of the ieee conference on computer vision and pattern recognition, 2015, pp. 815–823.

[35] Farnaz Sedighin, Andrzej Cichocki, and Anh-Huy Phan, *Adaptive rank selection for tensor ring decomposition*, IEEE Journal of Selected Topics in Signal Processing **15** (2021), no. 3, 454–463.

[36] Shai Shalev-Shwartz and Shai Ben-David, *Understanding machine learning: From theory to algorithms*, Cambridge university press, 2014.

[37] G Stein et al., *Spectroscopic data reduction with deep learning: A non-line-of-sight imaging example*, Optica **9** (2022), no. 2, 154–160.

[38] Li Sun et al., *Metric learning for medical image segmentation*, Medical Image Analysis **80** (2022), 102526.

[39] Matthew Tancik, Pratul P Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T Barron, and Ren Ng, *Fourier features let networks learn high frequency functions in low dimensional domains*, Advances in neural information processing systems, 2020, pp. 7537–7547.

[40] Joshua B Tenenbaum, Vin De Silva, and John C Langford, *A global geometric framework for nonlinear dimensionality reduction*, science **290** (2000), no. 5500, 2319–2323.

[41] Ledyard R Tucker, *Some mathematical notes on three-mode factor analysis*, Psychometrika **31** (1966), no. 3, 279–311.

[42] Laurens Van der Maaten and Geoffrey Hinton, *Visualizing data using t-sne*, Journal of Machine Learning Research **9** (2008), no. 11.

[43] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin, *Attention is all you need*, Advances in neural information processing systems **30** (2017).

[44] Mike Walmsley et al., *Galaxy zoo decals: Detailed visual morphology measurements from volunteers and deep learning for 314,000 galaxies*, Monthly Notices of the Royal Astronomical Society **509** (2022), no. 3, 3966–3988.

[45] Hao Wang, Yitong Wang, Zheng Zhou, Xing Ji, Dihong Gong, Jingchao Zhou, Zhifeng Li, and Wei Liu, *Cosface: Large margin cosine loss for deep face recognition*, Proceedings of the ieee conference on computer vision and pattern recognition, 2018, pp. 5265–5274.

[46] Tongzhou Wang and Phillip Isola, *Understanding contrastive representation learning through alignment and uniformity on the hypersphere* (2020), 9929–9939.

[47] Xun Wang, Xintong Han, Weilin Huang, Dengke Dong, and Matthew R Scott, *Multi-similarity loss with general pair weighting for deep metric learning*, Proceedings of the ieee/cvf conference on computer vision and pattern recognition, 2019, pp. 5022–5030.

[48] Tian Xie and Jeffrey C Grossman, *Crystal graph convolutional neural networks for an accurate and interpretable prediction of material properties*, Physical review letters **120** (2021), no. 14, 145301.