# GenURL: A General Framework for Unsupervised Representation Learning

Siyuan Li*, [ID], *Student Member, IEEE*, Zicheng Liu*, [ID], *Student Member, IEEE*,

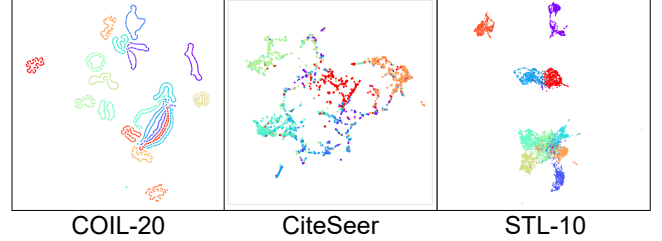Zelin Zang, Di Wu, Zhiyuan Chen, and Stan Z. Li†, [ID], *Fellow, IEEE*

Fig. 1. Illustration of various empirical structures of high-dimensional data. We encode COIL-20 [1], CiteSeer [2], and STL-10 [3] to 2-dim, 128-dim, and 512-dim by GenURL (128-dim and 512-dim latent spaces are then visualized by UMAP [4] in 2-dim). *Left*: we preserve local geometric structures of the circle manifolds in COIL-20 in the DR task. *Middle*: the topological and geometric structures of citation networks in CiteSeer are encoded in the GE task. *Right*: with the instance discriminative proxy task, we learn a discriminative representation in the validation dataset of STL-10.

*Abstract*—**Unsupervised representation learning (URL), which learns compact embeddings of high-dimensional data without supervision, has made remarkable progress recently. However, the development of URLs for different requirements is independent, which limits the generalization of the algorithms, especially prohibitive as the number of tasks grows. For example, dimension reduction methods, t-SNE, and UMAP optimize pair-wise data relationships by preserving the global geometric structure, while self-supervised learning, SimCLR, and BYOL focus on mining the local statistics of instances under specific augmentations. To address this dilemma, we summarize and propose a unified similarity-based URL framework, GenURL, which can smoothly adapt to various URL tasks. In this paper, we regard URL tasks as different implicit constraints on the data geometric structure that help to seek optimal low-dimensional representations that boil down to data structural modeling (DSM) and low-dimensional transformation (LDT). Specifically, DMS provides a structure-based submodule to describe the global structures, and LDT learns compact low-dimensional embeddings with given pretext tasks. Moreover, an objective function, General Kullback-Leibler divergence (GKL), is proposed to connect DMS and LDT naturally. Comprehensive experiments demonstrate that GenURL achieves consistent state-of-the-art performance in self-supervised visual learning, unsupervised knowledge distillation (KD), graph embeddings (GE), and dimension reduction. The code is available at https://github.com/Westlake-AI/openmixup.**

*Index Terms*—**Contrastive learning (CL), dimension reduction (DR), graph embedding (GE), knowledge distillation (KD), self-supervised learning (SSL)**

## I. INTRODUCTION

**L**earning low-dimensional representations from complex data without human supervision, *i.e.*, unsupervised representation learning (URL), is a long-standing problem. As the high-dimensional data is usually highly redundant and non-Euclidean, a widespread assumption is that data lies in a low-dimensional ambient space. However, URL algorithms under different tasks and data structures are designed independently

of each other, yet they have the same ultimate goal of finding the desired embedding space.

URL studies are now broadly categorized into three types of applications. (i) Dimension reduction (DR) and graph embedding (GE) algorithms aim to encode non-Euclidean input data into a latent space $Z$ plainly *without any prior knowledge* of the related domains, as shown in Figure 1 left and middle. (ii) Complementary to DR and GE, another popular path of URL focuses on data-specific augmentations, such as image crop, which leads to a clustering structure and learns discriminative representations, as illustrated in Figure 1 right. (iii) In addition, knowledge distillation (KD) is another approach that can be regarded as an implicit URL method transferring the knowledge from the pretrained model to enhance the student model unsupervised instead of considering geometric structure or special prior information of the target data. Specifically, from the perspective of algorithmic bias, these representative URL algorithms are grouped into two classes: global structure-oriented *e.g.,* t-SNE and UMAP, *etc.*, and individual augmentation-oriented *e.g.,* SimCLR, MoCo, *etc.*, respectively. It is a fact that these two independent algorithms are designed to excel in their respective areas of applicability. **There is thus a natural question if the intrinsic representation of data is determined by both the global data structure and data-specific prior assumptions in a unified framework.**

**This work: A general framework of unsupervised representation learning.** Based on the above URL methods, developing an effective and unified URL framework adaptive to various scenarios is a new trend in the community [5], [6]. In this work, we consider both the global structure and local

discriminative statistics and thus reformulate the URL problem as a non-Euclidean data embedding problem that encodes the structure and content parallelly in a compact low-dimensional space. For instance, we introduce an effective and general framework of unsupervised representation learning called GenURL, containing two steps: data structural modeling (DSM) and low-dimensional Transformation (LDT). To model the global structures, we combine the graph distance calculated on both the *raw feature space* and *predefined graphs*, which define the task-specific knowledge, *e.g.*, the provided graph $\mathcal{G}$ for GE tasks and the proxy task of SSL tasks. To learn the embedding according to data structures, we define a similarity between each sample pair and corresponding latent representations based on the pair-wise distances and design robust loss functions to optimize the encoder $f_\theta$. Unlike previous DR and GE methods, the proposed GenURL can import extra prior knowledge and is robust to highly redundant data; additionally, different from GE and SSL methods, GenURL is agnostic to network structures and predefined proxy tasks. Extensive experiments conducted on benchmarks of four URL tasks (self-supervised visual presentation learning, unsupervised knowledge distillation, graph embedding, and dimension reduction) show that GenURL achieves state-of-the-art performance. We further analyze the relationship between empirical data structures in various tasks and the loss functions and hyper-parameters in GenURL. In short, this paper makes three contributions:

- We propose a unified and general URL framework (GenURL) that encodes the global structure and local discriminative statistics of input data into a compact latent space parallelly.
- We discuss two types of embedding problems based on whether the input distance is well-defined and propose a novel objective function named General Kullback-Leibler divergence (GKL) to connect global and local efficiently.
- We adapt GenURL to various scenarios to verify the effectiveness and explain the relationship between various tasks through extensive experiments.

## II. RELATED WORK

*a) Dimension reduction:* Adopting the manifold assumption in DR, which assumes data lie on a low-dimensional manifold immersed in the high-dimensional space, most DR methods try to preserve intrinsic geometric properties of data [7]–[12]. Another practical branch of DR introduced by t-SNE [13] and UMAP [4] optimizes the pair-wise similarities between latent and input spaces. More recently, deep manifold learning methods can learn more complex manifolds and can be transferred to unseen data by using neural networks. Parametric t-SNE (P-SNE) [14], Parametric UMAP (P-UMAP) [15], DMT [16], and its variant [17] are proposed directly based on t-SNE and UMAP. However, current DR methods model desired data structures only relying on the geometry of input space and might fail with highly redundant data, such as natural images.

*b) Graph Embedding:* GE transfers graph data into a low-dimensional and continuous feature space while preserving most graph structures and topological and geometric structures, such as vertex content. Most early methods

are model-free, which contains four categories: Laplacian eigenmaps-based [18], local similarity-based [19]–[21], matrix factorization-based [22], and nonparametric Bayesian modeling-based methods [23]. Recently, model-based methods [24], [25] utilize graph convolutional networks (GCN) [26] or graph auto-encoders [27], [28] to learn both graph structures and feature information. More recently, some methods [29] take both the geometric and topological structures into consideration. This type of learning, which relies exclusively on the graph structure, eventually leads to the problem of homogenization of representations, so in addition to conveying information through the structure, the nodes themselves need to be bounded by independent a priori information.

*c) Self-supervised Visual Representation Learning:* Early SSL methods design hand-crafted pretext tasks [30]–[33], which rely on somewhat ad-hoc heuristics and have limited abilities to capture practically useful representations. Another popular form is clustering-based methods [34]–[37] learning discriminative representation by offline or online clustering. Recently, contrastive learning (CL) [38]–[41], which discriminates positive pairs against negative pairs, achieved state-of-the-art performance in various vision tasks. Different mechanisms [41]–[44] are proposed to prevent trivial solutions in CL to learn useful representations. To fully utilize negative samples, [45]–[48] explore hard samples in the momentum memory bank. Meanwhile, some efforts have been made on top of contrastive methods to improve pre-training quality for specific downstream tasks [49]–[52]. More recently, with the introduction of the Vision Transformers [53], [54], masked image modeling (MIM) [55]–[58] achieved state-of-the-art performance based on Transformers, which randomly mask out patches in the input image and predict the masked patches with decoder. While SSL can extract highly redundant data features, the loss of geometric structure constraints is an inability to model a priori semantic relationships between clusters in a global context, such as temporal evolutionary order.

*d) Knowledge Distillation:* KD was first proposed by [59], which aims to transfer knowledge from trained neural networks to a smaller one without losing too much generalization power. There are three types of existing KD methods: response-based [59]–[61] and feature-based [62] methods require labels to utilize intermediate-level supervision from the teacher model. Relation-based KD [63]–[67] methods explore the relationships between different layers or data sample pairs which can work without supervision and extend to self-supervised settings [67]–[69]. However, The performance of replicating the teacher model alone is unsatisfactory, and considering both the data structure and the a priori assumptions is a critical step in improving the efficiency of transfer learning. Our method handles the KD task as a special type of DR task without label supervision.

We summarize and analyze the above URL approaches and propose GenURL, a framework that successfully links the two important elements of URLs, data geometry and task hypotheses, based on generalized similarity. GenURL not only takes into account the data geometries that are the focus of DE and GE but also introduces task-relevant data and model hypotheses from CL and KD to enhance the overall quality
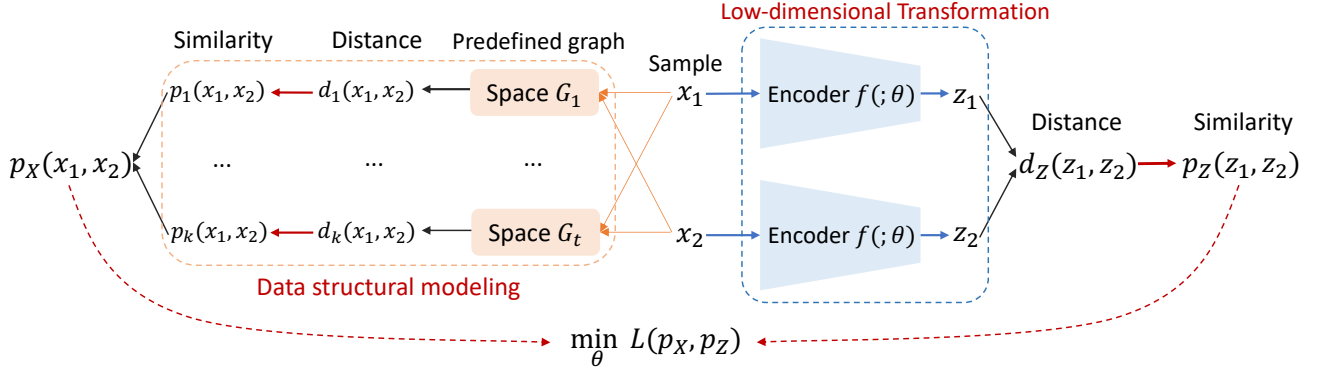
Fig. 2. Illustration of GenURL. The data structures are first modeled as similarity $P_X$ by calculating the graph distance on each predefined graph. Then, the low-dimensional transformation mapping $f_\theta$ is optimized by minimizing $\mathcal{L}$ based on the fixed $P_X$.

of the representations and thus improve the downstream tasks. performance of the downstream tasks.

## III. METHOD

### A. Preliminaries

The general goal of URL is identical across various URL tasks: Given a finite set of samples $X = [x_1, ..., x_n] \in \mathbb{R}^{n \times D}$, we seek a continuous mapping, $f_\theta(x) : \mathbb{R}^D \to \mathbb{R}^d$, where $d \ll D$, that transfers data into a compact latent space $Z$ while preserving essential structures of $X$ to facilitate most downstream tasks [25], [70], [71]. Generally, the intrinsic structures of data are determined by *both the input data and task-specific assumptions*. However, most research in URL has focused on individual data modalities or specific tasks, resulting in specific designs and different learning objectives. We take the following four widely used URL tasks as examples.

Although there are various task-specific assumptions in URL tasks, we may summarize them as a data embedding problem based on the manifold assumption. Assuming the data $X$ is constrained on a compact low-dimensional manifold $\mathcal{M}$, a neighborhood system for each sample $x_i$ is defined as $\mathcal{N}_i \in \mathbb{R}^d$, and $\mathcal{M}$ is supported on the mixture of these neighborhood systems, $\mathcal{M} = \cup_{i=1}^n \mathcal{N}_i$. The discriminative property of $\mathcal{N}$ facilitates various downstream tasks. Since different neighborhood systems are disjoint, we use an adjacency matrix $A$ to represent neighborhood systems: $A_{ij} = 1$ indicates $x_i$ and $x_j$ are in the same neighborhood system, or $A_{ij} = 0$. $A$ can be built by an undirected graph $\mathcal{G} = (X, E)$ provided by a specific URL task (discussed in Sec. IV). Local (neighborhood systems) and global (manifold) structures are two essential properties

for learning good representations: local geometries describe the discriminative features of instances, while relationships between neighborhood systems reflect the global view of topological structures. Based on $\mathcal{M}$, the similarity between two non-adjacent samples within each $\mathcal{N}_i$ can be approximated by the shortest-path distance. Since most URL algorithms are designed for specific tasks or data, the following two typical issues arise:

**Over-uniformity.** Without the global structure, the optimal solution for constructing $\mathcal{N}$s is to place each $x$ evenly in the embedding space, as shown in the middle of Fig 3. In other words, the decision boundaries are maximized to the discrimination between $x$. However, the manifold structure is then damaged and unorganized in this case, which means the learned representations can not be generalized to other downstream tasks [72], such as data visualization.

**Ill-clustering.** The converse is also true; if we focus excessively on the global structure and ignore instance differences, a local collapse will occur. The result is shown on the right of Fig 3, also called local homogenization. A classical case in the node classification task of graph data is the over-smoothing issue [73]: global-based message passing makes the connected nodes ultimately non-distinguishable.

Therefore, there is a question worth thinking about *whether we can solve both of these issues at the same time in a mutually constraining way*. Motivated by this, we propose a general and unified framework for URL that can be adapted to various URL tasks effectively.

### B. Similarity-preserving Framework

Given a set of $m$ empirical graphs $\mathcal{G} = \{G_t\}_{t=1}^m$ defined on the data $X$, where $\mathcal{G}_t = (X, d_t)$, we calculate the pair-wise distance $d_{X,t}$ based on $X$ and $\mathcal{G}_t$ to model the empirical data structures. To capture the local geometry defined in $\mathcal{G}_t$ and eliminate the scale factor between different distances, we adopt the *similarly* defined in $[0, 1]$ by converting the pair-wise distance $d$ to the similarity with a long-tailed t-distribution kernel function $\kappa(.)$,



Fig. 3. Illustration of two typical issues in unsupervised learning tasks: over-uniformity and ill-clustering.

$$\kappa(d, \nu) = \sqrt{2\pi} \cdot \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\nu\pi}\Gamma\left(\frac{\nu}{2}\right)} \left(1 + \frac{d^2}{\nu}\right)^{-\frac{(\nu+1)}{2}}, \quad (1)$$
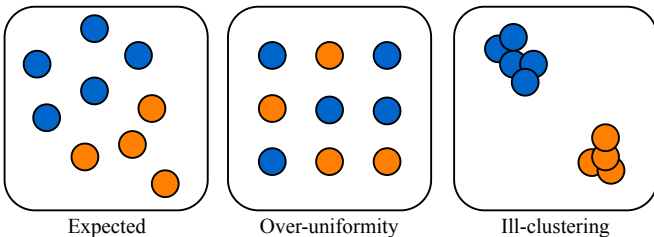
where $\nu$ denotes the degree of t-distribution freedom. Notice that the t-distribution approximates the Gaussian distribution when $\nu \to +\infty$, and approximates the uniform distribution when $\nu \to 0$. The latent space is usually a (normalized) Euclidean space $(Z, d_Z)$. The similarities of input and latent spaces are written as

$$p_X(x_i, x_j) = \alpha_t \sum_{t=1}^{m} A_{ij,t} \kappa(d_{X,t}(x_i, x_j), v_X), \quad (2)$$

$$p_Z(z_i, z_j) = \kappa(d_Z(z_i, z_j), v_Z), \quad (3)$$

where $\alpha_t$ is a balancing hyper-parameter for $d_{X,t}$. Notice that $\sum_{i,j} p_X(x_i, x_j)$ and $\sum_{i,j} p_Z(z_i, z_j)$ are not equal to 1 in a mini-batch.

**Data Structural Modeling.** Since $p_X$ can reflect the reliability of relations between $x$ and other samples, we can control learned representations by the *push* and *pull* forces with various $\nu_X$ and $\nu_Z$. Taking the DR task as an example, we assume $d_X$ is reliable in the original structure of input data while $d_Z$ is likely to be distorted in the extremely low-dimensional space (*e.g.*, 2-dim), as shown in Figure 4: We set $\nu_X \to +\infty$ (*i.e.*, the Gaussian distribution) giving the local sample pair $(x_i, x_j)$ and the disjoint sample pair $(x_i, x_k)$ high and low similarities, respectively, while set $\nu_X = 1$ to make $d_Z(x_i, x_j) \ll d_Z(x_i, x_k)$. As explained in Figure 4, the *push* and *pull* forces enable the learned embedding preserving geometric and topological properties of the input data after optimizing Eq. 6. For practical purposes, we can fix $\nu_X$ and adjust $\nu_Z$ in $[100, 0)$ to control the structure of the latent space based on the characters of URL tasks (detailed in Sec. V-F).

We provide *static* and *dynamic* methods to adaptively model the similarity $p_X$ based on Eq. 2. As for the *static*, we first normalize $d_{X,t}$ as $\tilde{d}_t(x_i, x_j) = \frac{d_t(x_i,x_j)-\mu_{i,t}}{\sigma_t}$, where $\mu_{i,t}$ measures the distance scale of each $x_i$ and $\sigma_t$ controls the extension of local neighborhood systems. We select proper $\mu_{i,t}$ and $\sigma_t$ in terms of population statistics of data (detailed in Sec. IV), such as mean and standard deviation of all samples. We rewrite Eq. 2 and Eq. 3 for the *static* $\tilde{p}_X$ as

$$\tilde{p}_X(x_i, x_j) = \sum_{k=1}^{K} \alpha_k \kappa\left(\frac{d_k(x_i,x_j) - \mu_{i,k}}{\sigma_k}, v_X\right), \quad (4)$$

$$\tilde{p}_Z(z_i, z_j) = \kappa\left(\frac{d_Z(z_i,z_j) - \mu_Z}{\sigma_Z}, v_Z\right). \quad (5)$$

Then, we embed $X$ into the latent space $Z$ by minimizing the dis-similarity between $\tilde{p}_X$ and $\tilde{p}_Z$ by a loss function $\mathcal{L}(.)$,

$$\min_\theta \mathcal{L}(\tilde{p}_X, \tilde{p}_Z). \quad (6)$$

**Low-dimensional Transformation.** Notice that the *static* $\tilde{p}_X$ mainly relies on the balancing weight $\alpha_t$ for each $d_{X,t}$, resulting in sub-optimal performances when some distances are unreliable. Thus, we design the *dynamic* method utilized in the early stage of the encoder $f_\theta$ to achieve reliable low-dimensional transformation, say, the $l$-th layer where $l \in [1, L-1]$. Since the encoder $f_\theta$ will gradually capture data structures by optimizing Eq.6, we regard the latent space of the $l$-th stage $\tilde{p}_{Z,l}$ as the *dynamic* $\tilde{p}_X$, which adaptively combines various $d_{X,t}$. Based on the *static* $\tilde{p}_X$, we define the
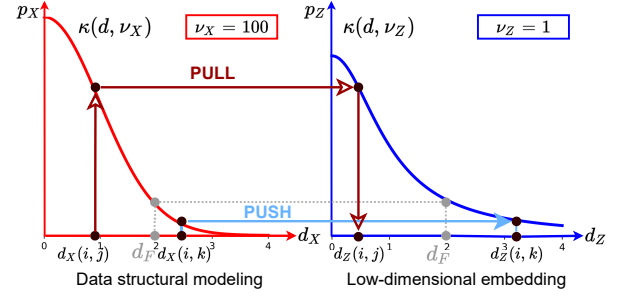


Fig. 4. Visualization of $p_X$ and $p_Z$ using the t-distribution. Let $d_Z^X = \kappa^{-1}(\kappa(d_X, \nu_X), \nu_Z)$ to be the projected distance of $d_X$ to the latent space. When $\nu_X > \nu_Z$, there exists a fix point $d_F$ between the t-distribution with $\nu_X$ and $\nu_Z$, we have $d_X(i, j) - d_F < d_Z(i, j) - d_F$ (*pull* between neighbors) and $d_X(i, k) - d_F < d_Z(i, k) - d_F$ (*push* between disjoint samples).

*dynamic* $\tilde{p}'_X = \beta\tilde{p}_X + \tilde{p}_{Z,l}$, where $\beta$ is a weight which linearly decays from 1 to 0. Notice that $\tilde{p}_{Z,l}$ does not require gradient. Finally, our proposed GenURL is demonstrated in Figure 2. As discussed in Sec. IV and Sec. V, the *static* $\tilde{p}_X$ usually suits URL tasks with well-defined input spaces like DR and KD, while the *dynamic* fits other tasks like CL and GE.

### C. Loss Function

As we formulate the URL problem as Eq. 6 where $\tilde{p}_X$ is regarded as the target, we discuss several similarity functions to achieve optimal embedding. Here, we consider $\tilde{p}_X$ in two cases: (i) generated by incomplete metric spaces where the relationship between distant neighbors is unknown; (ii) generated by well-defined metric spaces. In case (i), we focus on preserving the structures of each neighborhood system (similar pairs). In case (ii), we pay equal attention to dis-similar pairs to capture global relationships. We analyze these losses in various tasks by ablation studies in Sec. V.

*a) Mean squared error (MSE):* Mean squared error is the most commonly used loss function to measure the similarity between $\tilde{p}_X$ and $\tilde{p}_Z$ with $L_2$-norm,

$$\mathcal{L}_{MSE}(\tilde{p}_X, \tilde{p}_Z) = \sum_{i,j} ||\tilde{p}_X(x_i, x_j) - \tilde{p}_Z(z_i, z_j)||_2^2, \quad (7)$$

However, the MSE treats all sample pairs equally, resulting in sub-optimal solutions in both cases.

*b) Kullback-Leibler divergence (KL):* The KL divergence is commonly used to measure the similarity between two probability distributions,

$$\mathcal{L}_{KL}(\tilde{p}_X, \tilde{p}_Z) = -\sum_{i,j} \tilde{p}_X(x_i, x_j) \log \frac{\tilde{p}_X(x_i, x_j)}{\tilde{p}_Z(z_i, z_j)}. \quad (8)$$

When $\log \tilde{p}_X$ is constant, the KL divergence is equal to the cross-entropy between $\tilde{p}_X$ and $\tilde{p}_Z$. We can regard $\tilde{p}_X$ as a re-weight factor from similar sample pairs to dis-similar sample pairs. However, the KL divergence requires $\tilde{p}_Z(z_i, z_j) = 1$ [13]. When $\tilde{p}_Z$ is not a probability distribution [4], [44], it might result in a trivial solution, $\tilde{p}_Z(z_i, z_j) \to 0$ for each $\tilde{p}_X(x_i, x_j)$.

*c) Binary cross-entropy (BCE):* To make up the defect in the KL divergence, the binary cross-entropy loss [4] adds a symmetric term for $\tilde{p}_X(x_i, x_j) \to 0$ (placing higher weights than the KL divergence) to prevent the trivial solution,

$$\mathcal{L}_{BCE}(\tilde{p}_X, \tilde{p}_Z) = -\sum_{i,j} \tilde{p}_X(x_i, x_j) \log \tilde{p}_Z(z_i, z_j) - \sum_{i,j} (1 - \tilde{p}_X(x_i, x_j)) \log(1 - \tilde{p}_Z(z_i, z_j)). \quad (9)$$

The binary cross-entropy loss optimizes the most similar and dis-similar pairs symmetrically, which is suitable to perverse a well-defined metric space in the case (ii).

*d) General Kullback-Leibler divergence (GKL):* Although the BCE loss can solve the case (i) well under ideal conditions, it might suffer from outliers, *e.g.,* false negative samples in SSL and GE tasks [29], [47], [74], resulting in performance degradation in the case (ii). Since we can regard the symmetric term in Eq. 9 as the normalization constrain, $\sum_{i,j} \tilde{p}_X(x_i, x_j) = \sum_{i,j} \tilde{p}_Z(z_i, z_j)$, it is a direct way to prevent the trivial solution in KL divergence. However, the symmetric term emphasizes the importance of the negative samples (dissimilar pairs) with the reweight factor $1 - \tilde{p}_X$. Therefore, we propose a relaxed version of the BCE loss with a relaxed symmetric term,

$$\mathcal{L}_{GKL}(\tilde{p}_X, \tilde{p}_Z) = -\sum_{i,j} \tilde{p}_X(x_i, x_j) \log \tilde{p}_Z(z_i, z_j) + \gamma \sum_{i,j} ||\tilde{p}_X(x_i, x_j) - \tilde{p}_Z(z_i, z_j)||_p, \quad (10)$$

where $\gamma$ is a balancing weight and $||.||_p$ is $L_p$-norm. We adopt $L_1$-norm and set $\gamma = 0.1$ in the GKL loss. Compared to the BCE, the GKL loss is less affected by unreliable samples (usually dis-similar samples) when some $d_{X,t}$ are not reliable.

## IV. INSTANTIATION OF GENURL

GenURL generalizes different tasks by fully utilizing partially available information within corresponding datasets and patching the missing properties of URL modeling.

*a) Dimension Reduction and Graph Embedding:* The goal of GE tasks is to encode the geometric and topological structures of the input data. Given a graph $\mathcal{G} = (V, E, X)$, an adjacency matrix $A_1$ can be defined based on $\mathcal{G}$ and another $A_2$ based on a kNN graph built with the feature space $X$. We caculate the shortest-path distance for the entire graph, $d_1(v_i, v_j) = ||v_i - v_j||_2$ when $A_{1,ij} = 1$, and set $d_1(v_i, v_j)$ to a large constant when $A_{1,ij} = 0$. The distance $d_2$ of the kNN graph is obtained in the same way. Intuitively, we define the input similarity $p_X \triangleq \alpha_1 p_1 + \alpha_2 p_2$, where $\alpha_1 = 1$. To remove the scale effects of distances in different spaces, we calculate $\mu_{i,t} = \min(d_t(x_i, x_0), ..., d_t(x_i, x_n))$ and $\sigma_t = \frac{1}{n} \sum_{i=1}^{n} \sigma_i$ in data $X$. Instead of GCN, used in most GE methods, we use 5-layer MLP using leaky ReLU activation, with the middle latent dimension of 512 and the embedding dimension of 128. As for the DR task, it is regarded as a special case of GE tasks, which only requires a kNN graph. We adopt the *static* $\tilde{p}_X$ for both the tasks. Similarly, 3-layer MLP is adopted for DR tasks with the 2-dim embedding space.

*b) Self-supervised Learning for visual representation:* Unlike the DR and GE tasks, the distance on raw images in the open scenes is unreliable to reflect the desired low-dimensional structures for most discriminative downstream tasks since most images are highly redundant and unstructured. Hence, we import the proxy knowledge of the instance discriminative task in CL [39], [40] as follows. Given a mini-batch of data $X^B = \{x_i\}_{i=1}^{B}$, we apply augmentation $\tau \sim \mathcal{T}$ to each sample as $\tau(x_i)$ to obtain two correlated views $X_a^B = \{x_i^a\}_{i=1}^{B}$ and $X_b^B = \{x_j^b\}_{j=1}^{B}$, and fed to the encoder $f_\theta$ (*e.g.*, ResNet [70]) and a projection MLP neck [40], denoted as $h_\phi$, producing batches of latent representations $Z_a^B, Z_b^B$ and $H_a^B, H_b^B$, where $z_i = f_\theta(x_i)$ and $h_i = h_\phi(z_i)$. The projection neck will be discarded after pre-training. We can convert the proxy knowledge of content invariance into an adjacency matrix $A$: $A_{ii} = 1$ for two different views $x_i^1$ and $x_i^2$ (positive pairs) of the image $x_i$, while $A_{ij} = 0$ for any negative image pair $x_i$ and $x_j$. We adopt the $L_2$-normalized cosine distance of the projection as the latent representation, $d_Z \triangleq \frac{h_i}{||h_i||_2} \cdot \frac{h_j}{||h_j||_2}$. As we discussed in Sec. III-B, we adopt two versions of the input similarity. As for the *static* $\tilde{p}_X$, we use the discrete distance $d_1$ defined by the proxy knowledge and the Euclidean distance $d_2$ defined by kNN graph in $X$, $p_X \triangleq \alpha_1 p_1 + \alpha_2 p_2$, where $\alpha_1 = 1$ and $\alpha_2 = 0.01$ which is linearly decreased to 0. As for the *dynamic*, we calculate the cosine distance $\frac{z_i}{||z_i||_2} \cdot \frac{z_j}{||z_j||_2}$ and define $\tilde{p}_{Z,L}$ on the first latent space. The *dynamic* $\tilde{p}_X' = \beta \tilde{p}_X + \tilde{p}_{Z,L}$.

*c) Unsupervised Knowledge Distillation:* As for the KD task, we regard it as a special type of DR task, *i.e.*, encode the compact latent space of pertained teachers $z^T$ into the lower latent space of a student. Since the input space is already a well-defined Euclidean metric space, where the distance of every sample pair can be measured by $d_{Z^T}(x_i, x_j)$, *i.e.*, the input distance $d_X \triangleq d_{Z^T}$, we adopt the *static* method and use $L_2$-normalized cosine distance for the latent space of both the teacher and student. To fully explore the knowledge in teacher models, we should pay more attention to the global structural relation between distant samples while preserving the local geometry. As discussed in Sec. III-C, the BCE loss is more suitable for KD tasks.

## V. EXPERIMENTS

In this section, we evaluate the effectiveness of GenURL on various unsupervised learning tasks, including self-supervised visual representation (SSL), unsupervised knowledge distillation (KD), graph embedding (GE), and dimension reduction (DR). Meanwhile, we conduct ablation studies of loss functions and hyper-parameters to explore characters of various scenarios.

### A. Experimental Setup.

As for evaluation protocols, we adopt the linear protocol as the standard practice [39], [75], which trains a linear classifier on top of fixed representations. As for self-supervised visual representation, we further follow the semi-supervised classification [40] and evaluate the generation ability of representations by transfer learning [39]. As for dimension reduction, we further adopt trustworthiness (Trust) and continuity (Cont) [11] to measure the distortion between the input data

and representations. We use the following training settings for different tasks unless specified. We use Adam optimizer [76] with a learning rate of $lr \times BatchSize/256$ (linear scaling [77]) and a base $lr = 0.0005$. The batch size is 256 by default. All experiments report the mean of 3 times as default. The best and second results are denoted by **bold** and underlined.

*a) Datasets:* Various types of datasets are used in diverse URL tasks. Image datasets include: (1) MNIST [78] contains gray-scale images of 10 classes in $28 \times 28$ resolutions, 50K for training, and 10K for testing; (2) FashionMNIST [79] contains images of 10 classes of fashion clothing (same setting as MNIST); (3) COIL-20 [1] contains 72 different views (over an interval of 3 degrees) for 20 objects, for a total of 1440 images in $128 \times 128$ resolutions; (4) CIFAR-10/100 [80] contains 50K training images and 10K test images in $32 \times 32$ resolutions, with 10 classes and 100 classes settings; (5) STL-10 [3] consists of 5K labeled training images for 10 classes and 100K unlabelled training images and a test set of 8K images in $96 \times 96$ resolutions; (6) Tiny-ImageNet (Tiny) [81] has 10K training images and 10K validation images of 200 classes in $64 \times 64$ resolutions; (7) ImageNet-1K (IN-1K) [82] contains 1.28M training image and 50K validation images from 1000 classes in $224 \times 224$ resolutions. Datasets (1-3) are used for DR tasks, and (4-7) are used for SSL and KD tasks. Graph datasets for GE tasks include (8) CORA [83] contains binary word vectors of 7 classes with 2708 nodes, 1433 features, and 5429 edges; (9) CiteSeer [2] has binary word vectors of 6 classes with 3327 nodes, 3703 features, and 4732 edges; (10) PubMed [26] are associated with tf-idf weighted word vectors with 19717 nodes, 500 features, and 44338 edges for 3 classes.

*b) Implementation of contrastive learning:* We follow MoCo.v2 [84] for contrastive learning (CL) pre-training, which adopts ResNet [70] encoder with a two-layer MLP projector based on OpenMixup codebase [85]. All contrastive learning methods adopt the same network and augmentation settings, while other methods use default settings in their paper. The data augmentation setting in MoCo.v2 is as follows: Geometric augmentations include *RandomResizedCrop* with the scale in $[0.2, 1.0]$ and *RandomHorizontalFlip*. Color augmentation include *ColorJitter* with {brightness, contrast, saturation, hue} strength of $\{0.4, 0.4, 0.4, 0.1\}$ with an applying probability of 0.8, and *RandomGrayscale* with an applying probability of 0.2. Blurring augmentation uses a Gaussian kernel of size $23 \times 23$ with a standard deviation uniformly sampled in $[0.1, 2.0]$. As shown in Table I and Table II, we use the GKL loss with $\nu_X = \nu_Z = 100$ and $\sigma_X = \sigma_Z = 0.1$ for GenURL on CIFAR-10, CIFAR-100, STL-10, Tiny ImageNet, and ImageNet-1k datasets.

*c) Implementation of knowledge distillation:* In KD tasks, GenURL follows the settings of the current-proposed contrastive-based KD method SEED [67], which adopts the non-linear projector network and data augmentations used in MoCo.v2. Note that MoCo.v2 pre-trained ResNet-50 is adopted as the teacher model. Similar to the SSL task, GenURL uses the BCE loss with $\nu_X = \nu_Z = 100$ and $\sigma_X = 1$.

*d) Implementation of graph embedding:* In GE tasks, we adopt $L_2$ distance with $\sigma_Z = 1$ and tune various hyper-parameters as follows. As for $\mu$ and $\sigma$, we perform a grid search

of $\mu_Z$ and $\sigma_Z$ for the latent space in $\{0.01, 0.1, 1, 10, 100\}$ on the validation set. As for $\mu_{i,2}$ and $\sigma_2$ of the raw attribute space, we use a binary search (requires $O(n^2)$) with 5 nearest neighbors for each data point, *i.e.*, the optimal hyper-parameters guarantee the 5 nearest neighbors of $x_i$ have a large similarity score. There are similar practices in UMAP [4] and t-SNE [13]. If the dataset is too large, we set $\mu_{i,2}$ and $\sigma_2$ to the statistic mean and std of the whole dataset.

*e) Implementation of dimension reduction:* GenURL performs DR tasks with the BCE loss and the kNN graph built on the input $X$ following UMAP [4] and DMT [16] based on DMT implementation. Similar to the setting of GE tasks, we conduct a grid search of $\nu_Z$, $\mu_Z$, and $\sigma_Z$. We use $\nu_Z = 0.001$ and $\sigma_X = 5$ for MNIST and FMNIST datasets while using $\nu_Z = 0.01$ and $\sigma_X = 20$ for COIL-20.

TABLE I
**LINEAR EVALUATION AND SEMI-SUPERVISED LEARNING ON STL-10.**
TOP-1 ACCURACY (%) IS REPORTED WITH VARIOUS TRAINING EPOCHS
BASED ON RESNET-50.

| method | Linear | | | Semi-supervised | | |
|---|---|---|---|---|---|---|
| | 400ep | 800ep | 1600ep | 400ep | 800ep | 1600ep |
| Supervised | - | - | - | 71.21 | 72.70 | 72.89 |
| Related loc [30] | 60.19 | 64.20 | 64.37 | 86.49 | 87.93 | 88.41 |
| Rotation [32] | 76.70 | 73.14 | 72.15 | 89.91 | 90.43 | 90.05 |
| NPID [38] | 82.51 | 84.64 | 89.88 | 88.31 | 90.06 | 92.86 |
| ODC [35] | 73.43 | 75.47 | 76.20 | 80.88 | 82.04 | 85.80 |
| SimCLR [40] | 86.92 | 87.25 | 88.75 | 89.88 | 90.25 | 91.30 |
| MoCo.v2 [84] | 84.89 | 89.68 | 91.78 | 89.66 | 92.53 | **93.65** |
| BYOL [41] | 81.17 | 88.74 | 91.41 | 85.38 | 91.71 | 92.69 |
| SwAV* [37] | 84.35 | 88.79 | 91.02 | 86.57 | 92.05 | 92.63 |
| BarlowTwins [43] | 85.74 | 88.90 | 91.23 | 86.35 | 91.82 | 92.78 |
| **GenURL** | **88.35** | **90.82** | **91.85** | **90.88** | **92.58** | 93.55 |

TABLE II
**LINEAR EVALUATION ON CIFAR-10/100, TINY IMAGENET, AND
IMAGENET-1K.** TOP-1 ACCURACY (%) ARE REPORTED. RESNET-18
(R-18) IS USED AS THE ENCODER FOR CIFAR-10/100 AND TINY
IMAGENET TRAINING 400 AND 800 EPOCHS. R-18 AND RESNET-50 (R-50)
ARE USED FOR IMAGENET-1K TRAINING 200 EPOCHS.

| method | CIFAR-10 | | CIFAR-100 | | Tiny ImageNet | | ImageNet-1K | |
|---|---|---|---|---|---|---|---|---|
| | 400ep | 800ep | 400ep | 800ep | 400ep | 800ep | R-18 | R-50 |
| Supervised | 94.55 | 94.89 | 78.07 | 78.09 | 50.68 | 51.01 | 69.87 | 76.56 |
| Rotation [32] | 74.81 | 76.32 | 45.52 | 47.82 | 23.46 | 25.17 | 39.85 | 48.25 |
| NPID [38] | 79.53 | 82.70 | 54.52 | 57.16 | 36.86 | 38.24 | 43.10 | 58.87 |
| ODC [35] | 78.23 | 79.91 | 48.04 | 52.17 | 27.30 | 28.79 | 45.17 | 53.40 |
| SimCLR [40] | 86.22 | 88.24 | 56.45 | 57.42 | 37.64 | 38.46 | 51.03 | 66.67 |
| MoCo.v2 [84] | 82.41 | 88.62 | 56.65 | 61.48 | 33.00 | 37.49 | 52.87 | 67.85 |
| BYOL [41] | 82.61 | 88.15 | 58.32 | **64.40** | 33.93 | **38.81** | 54.62 | 71.88 |
| BarlowTwins [43] | 82.28 | 88.36 | 56.72 | 61.92 | 33.27 | 38.34 | 53.23 | 71.66 |
| **GenURL** | **88.27** | **88.95** | **59.01** | 63.51 | **37.81** | 38.68 | **55.12** | **72.15** |

### B. Self-supervised Visual Representation

In this subsection, we compare GenURL with three types of existing self-supervised methods, including head-craft, clustering-based, and contrastive learning methods. For a fair comparison, we apply the same augmentation settings described in MoCo.v2 [84] to all contrastive learning methods and follow hyper-parameters described in their original papers. We remove the Gaussian blur augmentation in CIFAR experiments [41], [42]. We perform unsupervised pre-training using ResNet-50 on STL-10 and ImageNet-1K, and using ResNet-18 on CIFAR-10/100, Tiny ImageNet, and ImageNet-1K. GenURL adopts the GKL loss and the *dynamic* method in the SSL task.

*a) Evaluation protocols:* As for linear evaluation, we follow the experiment settings in [39] to train a linear classifier for 100 epochs and use different base $lr$ for different datasets. We set the base $lr = 1.0$ for STL-10 and CIFAR-100, $lr = 0.1$ for CIFAR-10 and Tiny, and $lr = 30$ for IN-1K. The learning rate decays by 0.1 at 60 and 80 epochs. As for semi-supervised evaluation, we fine-tune the whole pre-trained model for 20 epochs on STL-10 with a step schedule at 12 and 16 epochs, and batch size is 256. We perform grid search for each test methods on base $lr = \{0.1, 0.01, 0.001\}$ and parameter-wise $lr\_mul = \{1, 10, 100\}$ of the $fc$ layer. Both experiments use the SGD optimizer with the weight decay of 0 for linear evaluation and 0.0001 for semi-supervised. Top-1 and top-5 accuracy are reported on the validation set.

*b) Linear evaluation results:* We first compare with existing methods in terms of different training epochs on STL-10, as shown in Table I. The proposed GenURL achieves the highest accuracy among all settings. It not only converges faster than other algorithms under 400-epoch pre-training but gains better performance when training longer. Then, we compare various methods on CIFAR-10/100 and Tiny, as shown in Table II. GenURL achieves the top performance on three datasets under 400-epoch pre-training and achieves second-best on CIFAR-100 and Tiny-ImageNet when training 800 epochs. Different from the existing contrastive-based methods, GenURL takes more pair-wise relations between samples into consideration, which might help GenURL convergence faster. For example, given a mini-batch of $N$ samples, BYOL utilizes $2N$ sample pairs, and MoCo requires $K + N$ sample pairs ($K$ is the memory bank size), while GenURL optimizes $N^2$ sample pairs (similar to SimCLR [40]).

*c) Semi-supervised evaluation results:* In Table I, we fine-tune a ResNet-50 pre-trained with various methods on the labeled training set of STL-10. GenURL outperforms other methods under 400-epoch and 800-epoch pre-training, which reflects its fast convergence speeds while maintaining the second-best classification accuracy with longer training.

*d) Transferring Features:* The main goal of unsupervised learning is to learn transferrable features. In Table V, we compare the representation quality of unsupervised pre-training on STL-10 by transferring to the classification task. We adopt linear evaluation on the CIFAR-10 in 64×64 resolutions with 1600-epoch pre-trained ResNet-50 on STL-10, and other settings are the same as Sec. V-B. GenURL achieves the highest accuracy among all methods: +3.36%/+3.29%/+2.99% for GenURL pre-training 400/800/1600 epochs over the second-best method.

*e) Ablation studies for SSL tasks:* We first ablate the loss functions used in visual SSL tasks. Since the input distance in SSL tasks is only well-defined for positive pairs where it can be optimized explicitly, and the relationship between negative pairs can be implicitly modeled by the *dynamic* method. As shown in Table III, GenURL prefers the GKL loss when using the *static* $\tilde{p}_X$, while the BCE loss yields better performance when using the *dynamic* to mine the relation between negative pairs. Then, we analyze hyperparameters of GenURL in Figure 5. We find that GenURL prefers $\nu_Z = 100$ and $\sigma = 1$ (approximating a standard Gaussian kernel) with a small batch size of 256.

TABLE III
**LOSS FUNCTION ANALYSIS ON SELF-SUPERVISED LEARNING.** WE EVALUATE THE LOSS FUNCTIONS PROPOSED IN SEC. 3 ON STL-10, CIFAR-10/100, AND TINY IMAGENET. TOP-1 ACCURACY (%) UNDER LINEAR EVALUATION IS REPORTED.

| Loss | $\tilde{p}_X$ setting | STL-10 | CIFAR-10 | CIFAR-100 | Tiny |
|------|------|--------|----------|-----------|------|
| MSE | *static* | 88.72 | 84.26 | 57.31 | 36.82 |
| BCE | *static* | 91.05 | 88.35 | 60.08 | 38.19 |
| BCE | *dynamic* | 91.60 | 88.87 | 61.16 | **38.85** |
| GKL | *static* | 91.21 | 88.63 | 61.27 | 38.07 |
| GKL | *dynamic* | **91.85** | **88.95** | **61.51** | 38.48 |

Moreover, we compare learned representations of GenURL with other visual SSL methods on STL-10 by UMAP [4] visualization in Figure 6.
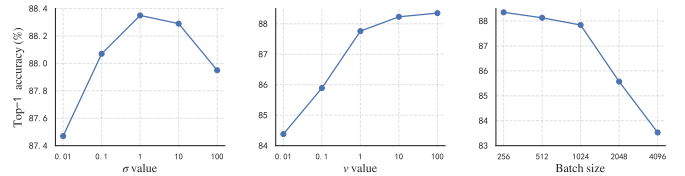


Fig. 5. Ablation of $\nu_Z$, $\sigma$ and batch size of GenURL for visual SSL tasks on STL-10. GenURL is pre-trained 800-epoch with ResNet-50.
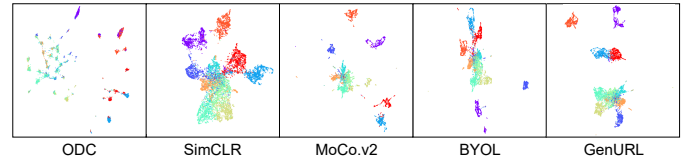


Fig. 6. Visualization of learned representations of CL methods with ResNet-50 on STL-10. We visualize the 2048-dim embeddings by UMAP. Compared to ODC and BYOL, the local structures of clusters are well-preserved, while each cluster is discriminative.

### C. Unsupervised Knowledge Distillation

We evaluate the KD tasks based on self-supervised learning on STL-10 dataset. In this experiment, we adopt MoCo.v2 with ResNet-50 under 1600-epoch pre-training. We choose multiple smaller networks with fewer parameters as the student network: ResNet-18 [70], MobileNet.v2 [86], ShuffleNet.v1 [87]. Similar to the pre-training for the teacher network, we add one additional MLP layer on the basis of the student network. Follow the linear evaluation protocols in Sec. V-B, we compare the existing relation-based KD methods including RKD [65], PKT [64], SP [66], SSKD [68], CRD [69], and SEED [67]. We adopt the BCE loss for GenURL in the KD task.

*a) Linear evaluation results:* From the view of different student models, as shown in Table IV, we notice that smaller networks perform rather worse and also benefit more from distillation than larger networks. From the perspective of various KD loss functions, the results clearly demonstrate that the proposed GenURL with the BCE loss achieves the best results, which is mainly because the BCE loss optimizes both local geometries and global structures.

*b) Ablation studies for KD tasks:* In contrast to SSL tasks, the input distance in KD tasks is assumed to be well-defined for both positive and negative sample pairs. Thus, we adopt the *static* $\tilde{p}_X$ and try to mine relationships among different sub-manifolds by momentum memory bank ($M$) for negative

TABLE IV
**UNSUPERVISED KNOWLEDGE DISTILLATION.** TOP-1 ACCURACY (%) UNDER LINEAR EVALUATION ON
STL-10. THE TEACHER MODEL IS RESNET-50 PRE-TRAINED BY MOCO.V2. † INDICATES USING A
MOMENTUM ENCODER AS MOCO.V2. SSL DENOTES THE INFONCE LOSS. KD DENOTES THE KNOWLEDGE
DISTILLATION LOSS. H+AW DENOTES THE HUBER LOSS AND ANGLE-PRESERVING LOSS IN RKD.

| KD methods | KD loss | ResNet-18 | | | MobileNet.v2 | | | ShuffleNet.v1 | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | SSL | KD | KD+SSL | SSL | KD | KD+SSL | SSL | KD | KD+SSL |
| RKD [65] | H+AW | | 86.48 | 86.76 | | 85.89 | 86.20 | | 84.31 | 85.01 |
| PKT [64] | KL | | 86.89 | 87.12 | | 86.14 | 86.48 | | 84.25 | 84.82 |
| SP [66] | MSE | | 86.53 | 86.74 | | 85.96 | 86.13 | | 84.22 | 84.76 |
| SSKD [68] | KL+InfoNCE | 81.51 | - | 87.78 | 79.96 | - | 86.80 | 77.26 | - | 85.23 |
| CRD [69] | KL+InfoNCE | | - | 87.24 | | - | 86.39 | | - | 84.98 |
| SEED† [67] | KL+InfoNCE | | - | 87.36 | | - | 86.44 | | - | 85.02 |
| **GenURL** | BCE | | 88.05 | 88.13 | | 86.61 | 86.85 | | 84.67 | 85.10 |
| **GenURL†** | BCE | | **88.26** | **88.39** | | **87.28** | **87.47** | | **85.05** | **85.38** |

TABLE V
**TRANSFER LEARNING ON CIFAR-10 CLASSIFICATION.** TOP-1 ACCURACY (%) UNDER LINEAR EVALUATION IS REPORTED.

| method | 400 ep | 800 ep | 1600 ep |
|---|---|---|---|
| Related loc [30] | 66.41 | 69.34 | - |
| Rotation [32] | 71.01 | 64.29 | - |
| NPID [38] | 71.15 | 63.72 | 65.30 |
| ODC [35] | 68.59 | 66.13 | 70.51 |
| SimCLR [40] | 75.97 | 75.08 | 76.86 |
| MoCo.v2 [84] | 74.46 | 76.54 | 75.61 |
| BYOL [41] | 74.04 | 76.83 | 75.55 |
| SwAV* [37] | 74.17 | 76.28 | 76.34 |
| BarlowTwins [43] | 74.63 | 76.71 | 76.12 |
| **GenURL** | **80.22** | **80.12** | **79.85** |

samples, which is similar to dark knowledge [59] in supervised KD tasks. Then, we ablate hyperparameters of GenURL for KD tasks in Figure 7. We find that using the BCE loss and $M$ with a large batch size achieves the best performance, and the GKL loss with $M$ yields the best result when using a small batch size.

### D. Unsupervised Graph Embedding

*a) Setups and results:* Unsupervised graph embedding experiments are conducted on three graph network datasets (Cora, CiteSeer, and PubMed), and we evaluate the learned embeddings by the node classification task. We compare GE methods that utilize both features and graph structures, including AGC [25], AGE [28], GIC [88], and ARGA [27]. The learned node embeddings are passed to logistic regression, and we report the mean and std of linear classification accuracy (Acc) of comparison methods. Table VII shows that GenURL (BCE) using *both* graphs and attributed features achieves new state-of-the-art performances on three GE datasets and improves previous GE methods by at least 0.5%, 0.4%, and 1.0% top-1 accuracy on CORA, CiteSeer, and PubMed datasets.

*b) Ablation and analysis:* We first adopt two ablation studies of the loss functions used in GenURL on GE tasks in Table VII: (i) when the input is *both* (using both graphs

TABLE VII
**NODE CLASSIFICATION.** TOP-1 ACCURACY (%) UNDER LINEAR EVALUATION ON CORA, CITESEER, AND PUBMED. *both* INDICATES USING BOTH THE GRAPH STRUCTURE AND ATTRIBUTED FEATURES.

| Method | Input | CORA | CiteSeer | PubMed |
|---|---|---|---|---|
| DeepWalk [20] | *graph* | $43.68_{\pm0.04}$ | $36.57_{\pm0.05}$ | $60.29_{\pm0.03}$ |
| AGC [25] | *both* | $75.50_{\pm0.02}$ | $67.37_{\pm0.01}$ | $67.70_{\pm0.02}$ |
| AGE [28] | *both* | $77.77_{\pm0.02}$ | $61.30_{\pm0.02}$ | $71.15_{\pm0.0}$ |
| GIC [88] | *both* | $77.27_{\pm0.01}$ | $41.88_{\pm0.03}$ | $76.70_{\pm0.02}$ |
| ARGA [27] | *both* | $72.52_{\pm0.02}$ | $55.21_{\pm0.01}$ | $64.24_{\pm0.02}$ |
| **GenURL(BCE)** | *feature* | $60.33_{\pm0.03}$ | $43.78_{\pm0.08}$ | $70.47_{\pm0.02}$ |
| **GenURL(GKL)** | *both* | $78.17_{\pm0.01}$ | $67.63_{\pm0.02}$ | $77.54_{\pm0.02}$ |
| **GenURL(BCE)** | *both* | $\mathbf{78.32}_{\pm0.02}$ | $\mathbf{67.70}_{\pm0.01}$ | $\mathbf{77.75}_{\pm0.02}$ |

and attributed features with the *dynamic* $\tilde{p}_X$), the BCE loss shows better performance than the GKL loss; (ii) when using the BCE loss, using *both* with the *dynamic* $\tilde{p}_X$ outperforms *feature* (only using the attributed features) with the *static* $\tilde{p}_X$. Then, we visualize the learned embedding on CiteSeer by UMAP in Figure 8. We find that GenURL separates sub-graphs of different classes into different clusters while maintaining the local geometric details of each sub-graph. As shown in Figure 9, we empirically show that $\nu_Z$ can control the balance between local geometric and global topological structures, and $\nu_Z = 0.005$ yields the best visualization results.

TABLE VI
**LOSS FUNCTION ANALYSIS ON UNSUPERVISED KNOWLEDGE DISTILLATION.** WE EVALUATE THE LOSS FUNCTIONS PROPOSED IN SEC. 3 ON STL-10. TOP-1 ACCURACY (%) UNDER LINEAR EVALUATION IS REPORTED. M DENOTES USING THE MOMENTUM ENCODER.

| batch size | MSE | | BCE | | GKL | |
|---|---|---|---|---|---|---|
| | w/o $M$ | w/ $M$ | w/o $M$ | w/ $M$ | w/o $M$ | w/ $M$ |
| 256 | 85.98 | 86.10 | 84.11 | 87.37 | 86.33 | **88.02** |
| 2048 | 86.14 | 86.43 | 88.05 | **88.26** | 85.97 | 85.91 |



Fig. 8. Visualization of the learned representations of GE methods on Wiki. We visualize the last latent space of the encoder by UMAP. The result of GenURL contains both the topology and the local geometric structures.

### E. Dimension Reduction

*a) Setups and results:* We perform DR experiments on MNIST, FMNIST, and COIL-20 datasets. We compare the current leading methods, including non-parametric methods (t-SNE [13] and UMAP [4]) and parametric methods (P-UMAP [15], GRAE [89], TopoAE [11], and DMT [16]).
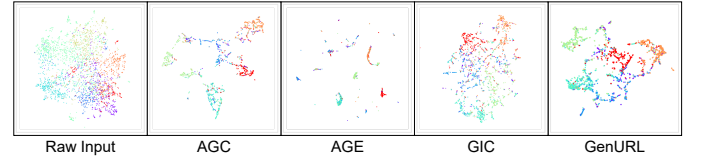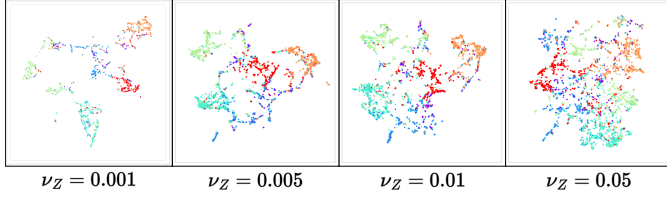


Fig. 7. Ablation of $\nu_Z$, $\sigma$ and batch size of GenURL for KD tasks on STL-10.

Fig. 9. Ablation of various $\nu_Z$ in GenURL for the GE task on CiteSeer. The learned embedding is visualized by UMAP to 2-dim space.

Besides the linear classification top-1 accuracy (Acc) with logistic regression, we evaluate the qualities of the low-dimensional representation in terms of the input space with Trust and Cont. Since the DR task only relies on the input space, GenURL adopts the BCE loss and the *static* $\tilde{p}_X$. As shown in Table VIII, we compare GenURL (BCE) with existing DR methods and find that GenURL yields comparable performance in terms of Trust and Acc, which indicates that GenURL (BCE) keeps the balance between local geometric structures (achieving better Trust and Cont) and the distinction of different sub-manifolds (achieving better linear classification accuracy).

*b) Ablation and analysis:* We first conduct the ablation of the BCE or GKL losses in Table VIII: GenURL (BCE) always outperforms GenURL (GKL) on three DR datasets because we adopt a large batch size as DMT and P-UMAP. Then, we provide DR results on COIL-20 in Figure 10 and find that GenURL captures more geometric structures than previous methods, especially UMAP and TopoAE (focusing on topological structures). Moreover, we ablate hyperparameters of GenURL (BCE) on MNIST in Figure 11 and find that GenURL prefers $\sigma = 1$, $\nu_Z = 0.001$, and the batch size of 2048. Similar to GE tasks, we verify whether $\nu_Z$ can control the balance of local structures and global topology by providing visualization of GenURL with various $\nu_Z$ on COIL-20. In Figure 12, we find that $\nu_Z = 0.01$ yields the best balance between local geometric and global topological structures.

TABLE VIII
**DIMENSION REDUCTION.** TRUST, CONT, AND TOP-1 ACCURACY (%) ARE
REPORTED ON MNIST, FMNIST, AND COIL-20.

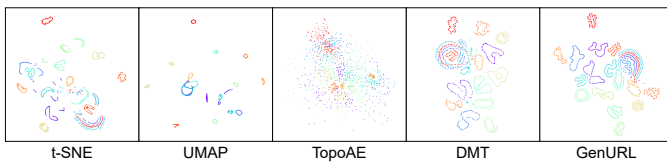| method | MNIST | | | FMNIST | | | COIL-20 | | |
|---|---|---|---|---|---|---|---|---|---|
| | Trust | Cont | Acc | Trust | Cont | Acc | Trust | Cont | Acc |
| t-SNE [13] | 0.872 | 0.855 | 83.4 | 0.949 | 0.934 | 69.2 | **0.925** | 0.867 | 89.0 |
| UMAP [4] | 0.887 | 0.837 | 96.6 | 0.947 | 0.938 | 68.7 | 0.922 | 0.849 | 87.1 |
| P-UMAP [15] | 0.890 | 0.834 | <u>96.7</u> | 0.951 | 0.937 | 68.9 | <u>0.924</u> | 0.843 | 89.3 |
| GRAE [89] | 0.876 | <u>0.861</u> | 75.0 | 0.949 | <u>0.943</u> | 58.2 | 0.920 | 0.887 | 86.5 |
| TopoAE [11] | 0.881 | **0.876** | 75.5 | 0.952 | **0.970** | 60.6 | 0.877 | 0.898 | 85.8 |
| DMT [16] | 0.896 | 0.840 | <u>96.7</u> | <u>0.958</u> | 0.939 | <u>70.0</u> | 0.916 | <u>0.928</u> | 90.2 |
| **GenURL(GKL)** | **0.897** | 0.849 | 96.4 | <u>0.958</u> | 0.937 | <u>70.0</u> | 0.923 | 0.926 | 90.3 |
| **GenURL(BCE)** | 0.898 | 0.842 | **96.8** | **0.959** | 0.940 | **70.2** | <u>0.924</u> | **0.930** | 90.4 |



Fig. 10. Visualization of COIL-20 by DR methods to 2-dim space. Compared to UMAP and DMT, GenURL preserves more local geometric details while capturing global topological structures.
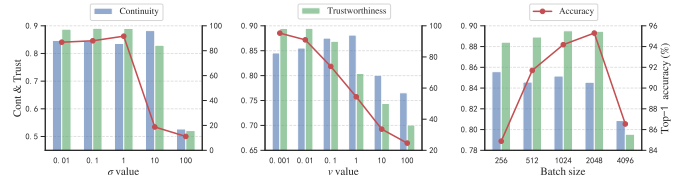


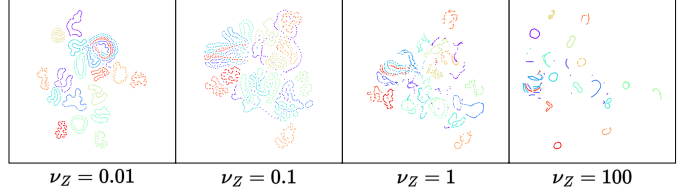Fig. 11. Ablation of $\nu_Z$, $\sigma$ and batch size of GenURL (BCE) for DR tasks on MNIST.



Fig. 12. Ablation of various $\nu_Z$ in GenURL for the DR task on CiteSeer. The learned embedding is visualized by UMAP to 2-dim space.

*F. Analysis and Discussion*

We provide an empirical analysis of the hyper-parameters and loss functions in GenURL on different URL tasks to demonstrate the characteristics of various tasks. We compare the results using different batch sizes, $\nu_Z$, $\sigma$, and loss functions used in GenURL to demonstrate the relationship of SSL, KD, and DR tasks.

*a) Relationship between DR and GE:* As shown in Figure 9 and Figure 12, GenURL prefers smaller $\nu_Z$ for both the GE and DR tasks because the large $\nu_Z$ yields crowd embedding while the small $\nu_Z$ conducts separable results. Therefore, we consider DR as a special type of GE task. The main difference between DR and GE tasks is GE takes both the node and edge features into consideration.

*b) Relationship between SSL and KD:* Then, we compare how GenURL deals with the negative samples in SSL and KD tasks. In Figure 7, we find that GenURL prefers the similar $\nu_Z$ and $\sigma$ for both SSL and KD tasks, which indicates using $\nu_Z = 100$ and $\sigma = 1$ (the standard Gaussian kernel) is suitable for $L_2$-normalized cosine distance. Notice that GenURL prefers small batch sizes like 256 for the SSL task (suffering performance drops when the batch size increases) while prefers larger batch sizes for the KD task. It might be because pair-wise similarities of negative samples in SSL tasks are unreliable and can be regarded as dark knowledge in the KD task [59], [67]. The gradient from negative pairs might overwhelm positive samples at the early training stage of the SSL task, while negative samples are well-defined by the teacher model in the KD task. Meanwhile, Table VI shows that using the *dynamic* version and the GKL loss in SSL tasks yield the best performance while using the large batch size, and the BCE loss in KD tasks performs better. Therefore, We conclude that the GKL with the *dynamic* structural modeling can alleviate the harmful effects of unreliable metric spaces of SSL tasks.

*c) Relationship between SSL and DR:* We further discuss the relation between SSL and DR tasks (the *dynamic* and *static* $\tilde{p}_X$) with GenURL to explain the desired structures of data in URL tasks. Meanwhile, we find that the instance discrimination prior knowledge in SSL tasks is more useful for downstream tasks like clustering and classification of complex scenarios, as
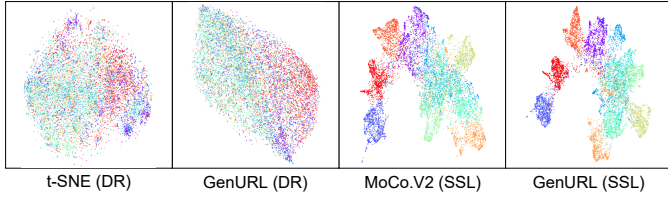
Fig. 13. Visualization of learned representation in DR and SSL tasks on CIFAR-10. The SSL representation is visualized by UMAP to 2-dim space.

shown in Figure 13. Therefore, we can choose a proper URL task for various scenarios: the DR task is suitable for MNIST and FMNIST datasets where the input spaces are discriminative and reliable, while the SSL task suits more complex datasets like CIFAR-10/100 where the input spaces are unreliable.

*d) Hyperparameters:* Firstly, we compare the effects of hyper-parameters in GenURL for DR and SSL tasks. As shown in Figure 12, GenURL prefers smaller $\nu_Z$, *i.e.,* using $\nu = 0.01$ to balance the local and global structures. Figure 5 shows that GenURL prefers $\nu_Z = 100$ for representations with strong discriminative abilities. Then, we find that the prior knowledge of instance discrimination in SSL tasks is more useful for downstream tasks like clustering and classification of complex scenarios. As shown in Figure 13, we summarize the learned representations with DR and SSL methods on CIFAR-10 and find that the results of SSL (adopting the instance discrimination prior knowledge) are more reliable and useful to downstream tasks like clustering and classification than DR. Empirically, DR methods are employed on "simple" datasets (*e.g.,* MNIST and COIL-20) with reliable geometric structures rather than "complex" datasets (*e.g.,* CIFAR-10 and ImageNet) with high-dimensional and redundant features. We hypothesize that this might depend on the property of the dataset. To verify our assumption, we compute the pair-wise distance between raw input samples and the latent space of the SSL model, as shown in Figure 14, to show the difference between DR and SSL tasks. We find that the input space is discriminative and reliable enough on MNIST and FMNIST for the DR task, while it is unreliable on CIFAR-10. Therefore, we can conclude as follows: when the dataset is reliable, GenURL can employ the BCE loss and the *static* $\tilde{p}_X$ to perform DR tasks (or GE tasks on the graph data); when the dataset contains high-dimensional redundant features, GenURL can adopt the GKL loss and the *dynamic* $\tilde{p}_X$ with the prior knowledge to conduct SSL tasks.

*e) Complexities:* GenURL has a constant algorithmic complexity in the four URL tasks, thanks to its property of unification. It improves performance without adding extra complexity. Unlike GNN and GCN, the proposed GenURL doesn't require neighborhood aggregation operations, making the complexity agnostic to the network architectures and the

kNN graph in the input space. In Sec. III, we build an undirected graph before training and calculate the adjacency matrix $A$ with $O(n^2)$. The pre-computed results will be saved to get $p_X(x_i, x_j)$ in Eq. 3. In each iteration, we calculate the latent space pair-wise similarity $p_X(z_i, z_j)$ with $O(n^2)$, assuming the batch size is n (*i.e.,* the whole dataset in GE and DR tasks). In CL and KD tasks, we will practically use a mini-batch of $b$ to reduce the computational complexity to $O(b^2)$.

## VI. LIMITATION AND FUTURE WORK

As for the societal impacts of GenURL, it can be regarded as a unified framework for the unsupervised representation learning (URL) problem that bridges the gap between various methods. The ablation studies of basic hyper-parameters can reflect the relationship between different URL tasks. The core idea of GenURL is to explore intrinsic structures of the data (the raw input space or empirical metric space) and preserve these structures in the latent space, which might inspire some improvements in various URL tasks. For example, the *dynamic* $\tilde{p}_X$ is similar to the hard negative mining problem in the SSL task [47], [74].

As for the limitations of GenURL, we can conclude three aspects: (i) the proposed framework relies on offline hyper-parameter tuning to adapt to new URL tasks, which makes it tough to handle more than two input similarities, (ii) GenURL cannot deal with the case of discrete empirical spaces well, *e.g.,* the SSL tasks, and the *dynamic* $\tilde{p}_X$ should be improved in the future work, (iii) the performance of GenURL is still limited by negative samples (sensitive to the size of datasets and the batch size). In future work, we plan to tackle the aforementioned limitations and design a dynamic framework that can learn to optimize the data structural modeling and low-dimensional embedding in an end-to-end manner. Moreover, GenURL will be used in more URL application scenarios.

## VII. CONCLUSIONS

We propose a simple but efficient similarity-based framework for unsupervised representation learning (URL), called GenURL, that encodes essential structures from the input data and optional prior knowledge. Specifically, we discuss the loss functions for embedding learning in GenURL and proposed binary cross-entropy loss and general Kullback-Leibler divergence loss. Combined with a specific pretext task, we can adapt GenURL to various URL scenarios in a unified manner and achieve state-of-the-art performances, including self-supervised visual representation learning, unsupervised knowledge distillation, graph embeddings, and dimension reduction. Moreover, ablation studies reflect the relationship between data characters and hyper-parameter settings of GenURL.
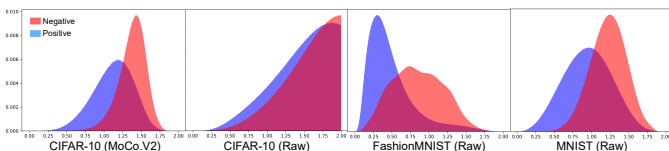
Fig. 14. Histograms of the pair-wise cosine distance of the positive (blue) and negative pairs (red) on STL-10, CIFAR-10, and MNIST datasets. Note that MoCo.v2 and raw denote the distance between the latent space (MoCo.v2) and raw feature spaces.

## REFERENCES

[1] S. A. Nene, S. K. Nayar, and H. Murase, "Columbia object image library (coil-20)," Columbia University, Tech. Rep., 1996. [Online]. Available: https://www.cs.columbia.edu/CAVE/software/softlib/coil-20.php 1, 6

[2] C. L. Giles, K. D. Bollacker, and S. Lawrence, "Citeseer: An automatic citation indexing system," in *Proceedings of the Third ACM Conference on Digital Libraries*, 1998, p. 89–98. [Online]. Available: https://doi.org/10.1145/276675.276685 1, 6

[3] A. Coates, A. Ng, and H. Lee, "An analysis of single-layer networks in unsupervised feature learning," in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2011, pp. 215–223. 1, 6

[4] L. McInnes, J. Healy, and J. Melville, "Umap: Uniform manifold approximation and projection for dimension reduction," *arXiv preprint arXiv:1802.03426*, 2018. 1, 2, 4, 5, 6, 7, 8, 9

[5] N. Mu, A. Kirillov, D. Wagner, and S. Xie, "Slip: Self-supervision meets language-image pre-training," *arXiv preprint arXiv:2112.12750*, 2021. 1

[6] A. Baevski, W.-N. Hsu, Q. Xu, A. Babu, J. Gu, and M. Auli, "data2vec: A general framework for self-supervised learning in speech, vision and language," *arXiv preprint arXiv:2202.03555*, 2022. 1

[7] J. B. Tenenbaum, V. De Silva, and J. C. Langford, "A global geometric framework for nonlinear dimensionality reduction," *science*, pp. 2319–2323, 2000. 2

[8] S. T. Roweis and L. K. Saul, "Nonlinear dimensionality reduction by locally linear embedding," *science*, pp. 2323–2326, 2000. 2

[9] Z. Zhang and H. Zha, "Principal manifolds and nonlinear dimensionality reduction via tangent space alignment," *SIAM journal on scientific computing*, pp. 313–338, 2004. 2

[10] D. L. Donoho and C. Grimes, "Hessian eigenmaps: Locally linear embedding techniques for high-dimensional data," *Proceedings of the National Academy of Sciences*, pp. 5591–5596, 2003. 2

[11] M. Moor, M. Horn, B. Rieck, and K. Borgwardt, "Topological autoencoders," in *Proceedings of the International Conference on Machine Learning (ICML)*, 2020, pp. 7045–7054. 2, 5, 8, 9

[12] S. Melacci and M. Gori, "Unsupervised learning by minimal entropy encoding," *IEEE transactions on neural networks and learning systems*, vol. 23, no. 12, pp. 1849–1861, 2012. 2

[13] L. v. d. Maaten and G. Hinton, "Visualizing data using t-sne," *Journal of machine learning research (JMLR)*, pp. 2579–2605, 2008. 2, 4, 6, 8, 9

[14] L. van der Maaten, "Learning a parametric embedding by preserving local structure." *Journal of Machine Learning Research (JMLR)*, vol. 5, pp. 384–391, 01 2009. 2

[15] T. Sainburg, L. McInnes, and T. Q. Gentner, "Parametric umap embeddings for representation and. semi-supervised learning," 2021. 2, 8, 9

[16] S. Z. Li, Z. Zang, and L. Wu, "Deep manifold transformation for nonlinear dimensionality reduction," *arXiv preprint arXiv:2010.14831*, 2021. 2, 6, 8, 9

[17] S. Li, H. Lin, Z. Zang, L. Wu, J. Xia, and S. Z. Li, "Invertible manifold learning for dimension reduction," in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases (ECML-PKDD)*, 2021, pp. 713–728. 2

[18] M. E. J. Newman, "Finding community structure in networks using the eigenvectors of matrices," *Physical Review E*, no. 3, Sep. 2006, arXiv: physics/0605087. 2

[19] Q. Le and T. Mikolov, "Distributed representations of sentences and documents," in *Proceedings of the International Conference on Machine Learning (ICML)*, 2014, pp. 1188—-1196. 2

[20] B. Perozzi, R. Al-Rfou, and S. Skiena, "Deepwalk: Online learning of social representations," in *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, 2014, pp. 701—-710. 2, 8

[21] R. Zhang, Y. Zhang, and X. Li, "Unsupervised feature selection via adaptive graph learning and constraint," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 33, no. 3, pp. 1355–1362, 2020. 2

[22] Y. Li and C. Sha, "Community Detection in Attributed Graphs: An Embedding Approach," in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, New Orleans, Louisiana, USA, 2018, pp. 338–345. 2

[23] A. Bojchevski and S. Günnemann, "Bayesian robust attributed graph clustering: Joint learning of partial anomalies and group structure," in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2018, pp. 265–271. 2

[24] T. N. Kipf and M. Welling, "Variational graph auto-encoders," *arXiv preprint arXiv:1611.07308*, 2016. 2

[25] X. Zhang and H. Liu, "Attributed graph clustering via adaptive graph convolution," *arXiv preprint arXiv:1906.01210*, 2019. 2, 3, 8

[26] N. K. Thomas and W. Max, "Semi-supervised classification with graph convolutional networks," *arXiv preprint arXiv:1609.02907*, 2017. 2, 6

[27] S. Pan, R. Hu, S.-F. Fung, G. Long, J. Jiang, and C. Zhang, "Learning graph embedding with adversarial training methods," *IEEE Transactions on Cybernetics*, p. 2475–2487, Jun 2020. 2, 8

[28] G. Cui, J. Zhou, C. Yang, and Z. Liu, "Adaptive graph encoder for attributed graph embedding," in *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, 2020, p. 976–985. 2, 8

[29] Z. Zang, S. Li, D. Wu, J. Guo, Y. Xu, and S. Z. Li, "Unsupervised deep manifold attributed graph embedding," *ArXiv*, vol. abs/2104.13048, 2021. 2, 5

[30] C. Doersch, A. Gupta, and A. A. Efros, "Unsupervised visual representation learning by context prediction," in *Proceedings of the International Conference on Computer Vision (ICCV)*, 2015. 2, 6, 8

[31] R. Zhang, P. Isola, and A. A. Efros, "Colorful image colorization," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016, pp. 41–57. 2

[32] S. Gidaris, P. Singh, and N. Komodakis, "Unsupervised representation learning by predicting image rotations," in *International Conference on Learning Representations (ICLR)*, 2018. 2, 6, 8

[33] D. He, C. Liang, C. Huo, Z. Feng, D. Jin, L. Yang, and W. Zhang, "Analyzing heterogeneous networks with missing attributes by unsupervised contrastive learning," *IEEE Transactions on Neural Networks and Learning Systems*, 2022. 2

[34] M. Caron, P. Bojanowski, A. Joulin, and M. Douze, "Deep clustering for unsupervised learning of visual features," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 132–149. 2

[35] X. Zhan, J. Xie, Z. Liu, Y. S. Ong, and C. C. Loy, "Online deep clustering for unsupervised representation learning," in *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 6, 8

[36] Y. M. Asano, C. Rupprecht, and A. Vedaldi, "Self-labelling via simultaneous clustering and representation learning," in *International Conference on Learning Representations (ICLR)*, 2020. 2

[37] M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin, "Unsupervised learning of visual features by contrasting cluster assignments," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2020, p. 9912–9924. 2, 6, 8

[38] Z. Wu, Y. Xiong, S. Yu, and D. Lin, "Unsupervised Feature Learning via Non-Parametric Instance-level Discrimination," *ArXiv:1805.01978 [cs]*, 2018. [Online]. Available: http://arxiv.org/abs/1805.01978 2, 6, 8

[39] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 5, 7

[40] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," *arXiv preprint arXiv:2002.05709*, 2020. 2, 5, 6, 7, 8

[41] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. H. Richemond, E. Buchatskaya, C. Doersch, B. A. Pires, Z. D. Guo, M. G. Azar *et al.*, "Bootstrap your own latent: A new approach to self-supervised learning," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2020, p. 21271–21284. 2, 6, 8

[42] X. Chen and K. He, "Exploring simple siamese representation learning," in *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2, 6

[43] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny, "Barlow twins: Self-supervised learning via redundancy reduction," in *International Conference on Machine Learning (ICML)*. PMLR, 2021, pp. 12 310–12 320. 2, 6, 8

[44] Z. Zang, S. Li, D. Wu, G. Wang, L. Shang, B. Sun, H. Li, and S. Z. Li, "Dlme: Deep local-flatness manifold embedding," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2022, p. 576–592. 2, 4

[45] C.-Y. Chuang, J. Robinson, Y.-C. Lin, A. Torralba, and S. Jegelka, "Debiased contrastive learning," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2020, pp. 8765–8775. 2

[46] Y. Kalantidis, M. B. Sariyildiz, N. Pion, P. Weinzaepfel, and D. Larlus, "Hard negative mixing for contrastive learning," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2020, pp. 21798–21809. 2

[47] J. D. Robinson, C.-Y. Chuang, S. Sra, and S. Jegelka, "Contrastive learning with hard negative samples," in *International Conference on Learning Representations (ICLR)*, 2021. 2, 5, 10

[48] S. Li, Z. Liu, Z. Wang, D. Wu, Z. Liu, and S. Z. Li, "Boosting discriminative visual representation learning with scenario-agnostic mixup," *arXiv preprint arXiv:2111.15454*, 2021. 2

[49] E. Xie, J. Ding, W. Wang, X. Zhan, H. Xu, Z. Li, and P. Luo, "Detco: Unsupervised contrastive learning for object detection," in *Proceedings of the International Conference on Computer Vision (ICCV)*, 2021. 2

[50] T. Xiao, C. J. Reed, X. Wang, K. Keutzer, and T. Darrell, "Region similarity representation learning," *arXiv preprint arXiv:2103.12902*, 2021. 2

[51] R. R. Selvaraju, K. Desai, J. Johnson, and N. Naik, "Casting your model: Learning to localize improves self-supervised representations," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 11 058–11 067. 2

[52] D. Wu, S. Li, Z. Zang, and S. Z. Li, "Exploring localization for self-supervised fine-grained contrastive learning," in *Proceedings of the British Machine Vision Conference (BMVC)*, 2022. 2

[53] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," in *International Conference on Learning Representations (ICLR)*, 2021. 2

[54] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 2

[55] H. Bao, L. Dong, and F. Wei, "Beit: Bert pre-training of image transformers," in *International Conference on Learning Representations (ICLR)*, 2022. 2

[56] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, "Masked autoencoders are scalable vision learners," *arXiv preprint arXiv:2111.06377*, 2021. 2

[57] Z. Xie, Z. Zhang, Y. Cao, Y. Lin, J. Bao, Z. Yao, Q. Dai, and H. Hu, "Simmim: A simple framework for masked image modeling," *arXiv preprint arXiv:2111.09886*, 2021. 2

[58] S. Li, D. Wu, F. Wu, Z. Zang, K. Wang, L. Shang, B. Sun, H. Li, and Stan.Z.Li, "Architecture-agnostic masked image modeling - from vit back to cnn," *ArXiv*, vol. abs/2205.13943, 2022. 2

[59] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, 2015. 2, 8, 9

[60] Y. Zhang, T. Xiang, T. M. Hospedales, and H. Lu, "Deep mutual learning," *arXiv preprint arXiv:1706.00384*, 2017. 2

[61] Q. Xie, M.-T. Luong, E. Hovy, and Q. V. Le, "Self-training with noisy student improves imagenet classification," in *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2

[62] S. Zagoruyko and N. Komodakis, "Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer," in *International Conference on Learning Representations (ICLR)*, 2017. 2

[63] Y. Junho, J. Donggyu, B. Jihoon, and K. Junmo, "A gift from knowledge distillation: Fast optimization, network minimization and transfer learning," in *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2

[64] N. Passalis and A. Tefas, "Learning deep representations with probabilistic knowledge transfer," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, p. 283–299. 2, 7, 8

[65] W. Park, D. Kim, Y. Lu, and M. Cho, "Relational knowledge distillation," in *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 3962–3971. 2, 7, 8

[66] F. Tung and G. Mori, "Similarity-preserving knowledge distillation," in *Proceedings of the International Conference on Computer Vision (ICCV)*, 2019. 2, 7, 8

[67] Z. Fang, J. Wang, L. Wang, L. Zhang, Y. Yang, and Z. Liu, "Seed: Self-supervised distillation for visual representation," in *International Conference on Learning Representations (ICLR)*, 2021. 2, 6, 7, 8, 9

[68] G. Xu, Z. Liu, X. Li, and C. C. Loy, "Knowledge distillation meets self-supervision," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020, p. 588–604. 2, 7, 8

[69] Y. Tian, D. Krishnan, and P. Isola, "Contrastive representation distillation," in *International Conference on Learning Representations (ICLR)*, 2020. 2, 7, 8

[70] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2016, pp. 770–778. 3, 5, 6, 7

[71] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *Proceedings of the International Conference on Computer Vision (ICCV)*, 2017. 3

[72] T. Wang and P. Isola, "Understanding contrastive representation learning through alignment and uniformity on the hypersphere," in *International Conference on Machine Learning*. PMLR, 2020, pp. 9929–9939. 3

[73] C. Cai and Y. Wang, "A note on over-smoothing for graph neural networks," *arXiv preprint arXiv:2006.13318*, 2020. 3

[74] Z. Shen, Z. Liu, Z. Liu, M. Savvides, T. Darrell, and E. Xing, "Unmix: Rethinking image mixtures for unsupervised visual representation learning," in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2021, pp. 2216–2224. 5, 10

[75] P. Goyal, D. Mahajan, A. Gupta, and I. Misra, "Scaling and benchmarking self-supervised visual representation learning," in *Proceedings of the International Conference on Computer Vision (ICCV)*, 2019. 5

[76] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *International Conference on Learning Representations (ICLR)*, 2015. 6

[77] P. Goyal, P. Dollár, R. Girshick, P. Noordhuis, L. Wesolowski, A. Kyrola, A. Tulloch, Y. Jia, and K. He, "Accurate, large minibatch sgd: Training imagenet in 1 hour," *arXiv preprint arXiv:1706.02677*, 2020. 6

[78] Y. LeCun, L. Bottou, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, pp. 2278–2324, 1998. 6

[79] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms," *arXiv preprint arXiv:1708.07747*, 2017. 6

[80] A. Krizhevsky, "Learning multiple layers of features from tiny images," *Master's thesis, University of Tront*, 2009. 6

[81] P. Chrabaszcz, I. Loshchilov, and F. Hutter, "A downsampled variant of imagenet as an alternative to the cifar datasets," *arXiv preprint arXiv:1707.08819*, 2017. 6

[82] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2012, pp. 1097–1105. 6

[83] A. McCallum, K. Nigam, J. D. M. Rennie, and K. Seymore, "Automating the construction of internet portals with machine learning," *Information Retrieval*, pp. 127–163, 2004. 6

[84] X. Chen, H. Fan, R. Girshick, and K. He, "Improved baselines with momentum contrastive learning," *arXiv preprint arXiv:2003.04297*, 2020. 6, 8

[85] S. Li, Z. Wang, Z. Liu, D. Wu, C. Tan, and D. W. S. Z. Li, "Openmixup: A comprehensive mixup benchmark for visual classification," *ArXiv*, vol. abs/2209.04851, 2022. 6

[86] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 4510–4520. 7

[87] X. Zhang, X. Zhou, M. Lin, and J. Sun, "Shufflenet: An extremely efficient convolutional neural network for mobile devices," *arXiv preprint arXiv:1707.01083*, 2017. 7

[88] C. Mavromatis and G. Karypis, "Graph infoclust: Leveraging cluster-level node information for unsupervised graph representation learning," *arXiv preprint arXiv:2009.06946*, 2020. 8

[89] A. F. Duque, S. Morin, G. Wolf, and K. R. Moon, "Extendable and invertible manifold learning with geometry regularized autoencoders," *arXiv preprint arXiv:2007.07142*, 2020. 8, 9

**Siyuan Li** (Student Member, IEEE) received the B.S. degree from the Department of Computer Science and Technology, Nanjing University, Nanjing, China, in 2021. He is currently pursuing the Ph.D. degree in the School of Engineering at Westlake University and the Department of Computer Science and Technology at Zhejiang University, supervised by Prof. Stan Z. Li (Fellow, IEEE). His main research interests include self-supervised learning, data augmentation, network architecture design in computer vision, and biological application.
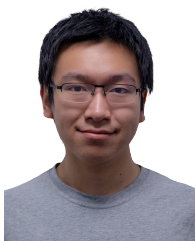
**Zicheng Liu** (Student Member, IEEE) received the B.S. degree from the Department of Information and Computing Science, University of Liverpool, Liverpool, U.K., in 2020. He is currently pursuing the Ph.D. degree with the School of Engineering, Westlake University and Zhejiang University, supervised by Prof. Stan Z. Li (Fellow, IEEE). His research interests include data augmentation, network architecture design, and biological and computer vision applications.

**Zelin Zang** (Student Member, IEEE) received the M.Eng. degree from the Zhejiang University of Technology, Hangzhou, China, in 2020. He is currently pursuing the Ph.D. degree at Westlake University and Zhejiang University, supervised by Prof. Stan Z. Li (Fellow, IEEE). His research interests include manifold learning, dimension reduction, self-supervised learning in computer vision, and biological applications.

**Di Wu** received the B.S. degree from the Department of Computer Science, Harbin Institute of Technology, Harbin, China, in 2016. He received the M.S. degree in electrical and computer engineering from Boston University, Boston, MA, USA, in 2018. He is currently working toward the Ph.D. degree with the Center of Excellence in Biomedical Research on Advanced Integrated-on-chips Neurotechnologies, School of Engineering, Westlake University, Hangzhou, China. His research interests include self-supervised learning and efficient deep learning for neurophysiological applications.

**Zhiyuan Chen** received the B.S. degree from the Australian National University, Australia, in 2022. We worked as a research assistant in the StarBridge program & a research intern at the Star of Tomorrow program under the guidance of Dr. Pan Deng at Microsoft Research Asia & Microsoft Research AI for Science from 2021-2022. He is currently a researcher at Deep Potential, Beijing, China, working on macro-molecules. His research interests include AI4Science, computer vision, and deep learning.

**Stan Z. Li** (Fellow, IEEE) received the B.Eng. degree from Hunan University, Changsha, China, in 1982, the M.Eng. degree from the National University of Defense Technology, Changsha, in 1985, and the Ph.D. degree from the University of Surrey, Guildford, U.K, in 1991. He was the Director of the Center for Biometrics and Security Research (CBSR), Chinese Academy of Sciences, Beijing, China, from 2004 to 2019. He worked at Microsoft Research Asia, Beijing, as a Research Lead from 2000 to 2004. Prior to that, he was an Associate Professor (Tenure) at Nanyang Technological University, Singapore. He joined Westlake University, Hangzhou, China, as a Chair Professor of artificial intelligence in February 2019. He has published over 400 articles in international journals and conferences and authored and edited 10 books, with over 50,000 Google Scholar citations. Among these are Markov Random Field Models in Image Analysis (Springer), Handbook of Face Recognition(Springer), and Encyclopedia of Biometrics (Springer). His research interests include fundamental research in machine learning, data science, and applied research in multiple AI-related interdisciplinary fields (computer vision, smart sensors, life science, material science, and environmental science). Dr. Li served as an Associate Editor for IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE and organized more than 100 international conferences or workshops.