

TranSG: Transformer-Based Skeleton Graph Prototype Contrastive Learning with Structure-Trajectory Prompted Reconstruction for Person Re-Identification

– Appendix I Theoretical Hypotheses and Analyses

Haocong Rao Chunyan Miao*

LILY Research Center, Nanyang Technological University, Singapore

School of Computer Science and Engineering, Nanyang Technological University, Singapore

{haocong001, ascmiao}@ntu.edu.sg

The proposed graph prototype contrastive learning (GPC) can be formulated as Expectation-Maximization (EM) solutions. In this appendix, we first provide a theoretical EM modeling for GPC to prove its validity and convergence, and then systematically present the relations and differences between the proposed approach and existing contrastive learning paradigms.

Preliminaries. For clarity and convenience, we adopt a more general notation here, which is different from that used in the paper. Suppose that a training set $S = \{s_i\}_{i=1}^N$ contains N skeleton sequences, where $s_i \in \mathbb{R}^{F \times K}$, $K = J \times 3$, J is the number of body joints with 3D positions, and F is the sequence length (*i.e.*, number of consecutive skeletons). We first represent each sequence with skeleton graphs by: $G(s_i) = \mathbf{x}_i = (\mathcal{G}_1, \dots, \mathcal{G}_F)$, where $G(\cdot)$ is the *pre-defined* graph construction function, $\mathbf{x}_i = (\mathcal{G}_1, \dots, \mathcal{G}_F)$ is the consecutive graphs representing the i^{th} skeleton sequence s_i , and \mathcal{G}_t denotes the t^{th} graph corresponding to the t^{th} skeleton in s_i . In our work, instead of using original skeleton sequences, we utilize their skeleton graph representations as inputs to capture richer skeletal body and motion features. The objective of skeleton graph representation learning is to learn a graph embedding/encoder function f_θ (realized via θ -parameterized neural networks) that maps the skeleton graphs \mathbf{x}_i to $\mathbf{v}_i \in \mathbb{R}^H$, by $\mathbf{v}_i = f_\theta(\mathbf{x}_i)$, such that \mathbf{v}_i can effectively represent latent features of \mathbf{x}_i to perform person re-identification.

Formally, the goal is to find the network parameter θ that maximizes the log-likelihood function of the observed graph representations $\{\mathbf{x}_i\}_{i=1}^N$ of N skeleton sequences as

follows:

$$\begin{aligned} \theta^* &= \arg \max_{\theta} L(\mathbf{x}_1, \dots, \mathbf{x}_N; \theta) \\ &= \arg \max_{\theta} \prod_{i=1}^N p(\mathbf{x}_i; \theta) \\ &\iff \arg \max_{\theta} \sum_{i=1}^N \log p(\mathbf{x}_i; \theta), \end{aligned} \quad (1)$$

where $L(\mathbf{x}_1, \dots, \mathbf{x}_N; \theta)$ denotes the likelihood function of the observed skeleton graph representations *with regard to* θ , and each skeleton graph representation¹ \mathbf{x}_i is hypothetically related to a certain skeleton graph prototype $\mathbf{c}_j \in \mathbb{R}^H$, with $\mathbf{c}_j \in \{\mathbf{c}_j\}_{j=1}^K$ and K is the number of graph prototypes. Under this assumption, we can re-formulate the objective in Eq. (1) as:

$$\begin{aligned} \theta^* &= \arg \max_{\theta} \sum_{i=1}^N \log p(\mathbf{x}_i; \theta) \\ &= \arg \max_{\theta} \sum_{i=1}^N \log \sum_{j=1}^K p(\mathbf{x}_i, \mathbf{c}_j; \theta), \end{aligned} \quad (2)$$

Directly optimizing this function is intractable, thus we consider a lower-bound by using a surrogate function as:

$$\begin{aligned} &\sum_{i=1}^N \log \sum_{j=1}^K p(\mathbf{x}_i, \mathbf{c}_j; \theta) \\ &= \sum_{i=1}^N \log \sum_{j=1}^K Q(\mathbf{c}_j) \frac{p(\mathbf{x}_i, \mathbf{c}_j; \theta)}{Q(\mathbf{c}_j)} \\ &\geq \sum_{i=1}^N \sum_{j=1}^K Q(\mathbf{c}_j) \log \frac{p(\mathbf{x}_i, \mathbf{c}_j; \theta)}{Q(\mathbf{c}_j)}, \end{aligned} \quad (3)$$

*Corresponding author

¹For simplicity of presentation, we use the “skeleton graph representation” to denote the graph representation of a skeleton sequence.

where $Q(c_j)$ represents some distribution over $\{c_j\}_{j=1}^K$ and $\sum_{j=1}^K Q(c_j) = 1$. We apply Jensen's inequality to derive the last step of Eq. (3), where the equality can be achieved under the condition that $\frac{p(\mathbf{x}_i, c_j; \theta)}{Q(c_j)}$ is a constant. To realize this equality, we have:

$$Q(c_j) = \frac{p(\mathbf{x}_i, c_j; \theta)}{\sum_{m=1}^K p(\mathbf{x}_i, c_m; \theta)} = \frac{p(\mathbf{x}_i, c_j; \theta)}{p(\mathbf{x}_i; \theta)} = p(c_j; \mathbf{x}_i, \theta), \quad (4)$$

where $Q(c_j)$ is a posterior probability related to c_j, \mathbf{x}_i , and θ . Different from [1, 2] that employ clustering to estimate prototypes and representation distributions under the fixed θ , we exploit the graph feature centroid of each *ground-truth* identity as a different graph prototype. In particular, when given the θ -parameterized encoder to encode all skeleton graphs (\mathbf{x}_i) at the *Expectation step*, their corresponding graph prototypes (c_j) are assumed to follow the distribution of ground-truth classes in the dataset. The graph prototype distribution is hence constant and can be computed by $Q(c_j) = p(c_j; \mathbf{x}_i, \theta)$. We can re-write Eq. (3) as:

$$\sum_{i=1}^N \sum_{j=1}^K (Q(c_j) \log p(\mathbf{x}_i, c_j; \theta) - Q(c_j) \log Q(c_j)), \quad (5)$$

where the constant $-\sum_{i=1}^N \sum_{j=1}^K Q(c_j) \log Q(c_j)$ can be ignored and we need to maximize:

$$\sum_{i=1}^N \sum_{j=1}^K Q(c_j) \log p(\mathbf{x}_i, c_j; \theta). \quad (6)$$

For the **Expectation (E)-step**, $p(c_j; \mathbf{x}_i, \theta)$ (see Eq. (4)) is estimated by the ground-truth class distribution. In our approach, the number of graph prototypes (K) is identical to the number of different classes, and we generate skeleton graph prototypes $\{c_j\}_{j=1}^K$ by computing the feature centroids of encoded skeleton graph representations \mathbf{v}_i in different classes. We use $\{C_j\}_{j=1}^K$ to denote the groups² of graph representations (referred to as “sample groups”) belonging to different prototypes. Then, we compute $p(c_j; \mathbf{x}_i, \theta) = \mathbb{1}(\mathbf{x}_i \in C_j)$, where $\mathbb{1}(\mathbf{x}_i \in C_j) = 1$ if \mathbf{x}_i belongs to the j^{th} sample group C_j (i.e., corresponding to graph prototype c_j); otherwise $\mathbb{1}(\mathbf{x}_i \in C_j) = 0$.

Assumption 1 Prototype-Class Consistency. *The global distribution of graph prototypes is consistent with the distribution of class feature centroids, i.e., all samples belonging*

²The graph prototypes $\{c_j\}_{j=1}^K$ for $\{\mathbf{x}_i\}_{i=1}^N$ are generated based on their encoded features $\{\mathbf{v}_i\}_{i=1}^N$, while $\{C_j\}_{j=1}^K$ are groups of $\{\mathbf{x}_i\}_{i=1}^N$ belonging to different graph prototypes.

to a ground-truth class explicitly correspond to the sample group of a certain prototype. In the *E-step*, we adopt this assumption to generate skeleton graph prototypes and derive $p(c_j; \mathbf{x}_i, \theta) = \mathbb{1}(\mathbf{x}_i \in C_j)$.

In the **Maximization (M)-step**, we combine Eq. (4) to maximize the lower-bound in Eq. (6) after the E-step:

$$\begin{aligned} & \sum_{i=1}^N \sum_{j=1}^K Q(c_j) \log p(\mathbf{x}_i, c_j; \theta) \\ &= \sum_{i=1}^N \sum_{j=1}^K p(c_j; \mathbf{x}_i, \theta) \log p(\mathbf{x}_i, c_j; \theta) \\ &= \sum_{i=1}^N \sum_{j=1}^K \mathbb{1}(\mathbf{x}_i \in C_j) \log p(\mathbf{x}_i, c_j; \theta). \end{aligned} \quad (7)$$

Assumed that each class is equally important and the sample number of each class is approximately identical in learning, each graph prototype c_j can have a uniform prior probability $p(c_j; \theta) = \frac{1}{K}$. We have:

$$\begin{aligned} p(\mathbf{x}_i, c_j; \theta) &= p(\mathbf{x}_i; c_j, \theta) p(c_j; \theta) \\ &= \frac{1}{K} \cdot p(\mathbf{x}_i; c_j, \theta), \end{aligned} \quad (8)$$

where the distribution of samples around each graph prototype is assumed to be an isotropic Gaussian, leading to:

$$p(\mathbf{x}_i; c_j, \theta) = \frac{\exp\left(\frac{-(\mathbf{v}_i - \mathbf{c}_p)^2}{2\sigma_p^2}\right)}{\sum_{j=1}^K \exp\left(\frac{-(\mathbf{v}_i - \mathbf{c}_j)^2}{2\sigma_j^2}\right)}, \quad (9)$$

where $\mathbf{v}_i = f_\theta(\mathbf{x}_i)$ and \mathbf{c}_p is the graph prototype for the sample group C_p containing \mathbf{x}_i , i.e., $\mathbf{x}_i \in C_p$. We apply ℓ_2 -normalization to both \mathbf{v} and \mathbf{c} to have $(\mathbf{v} - \mathbf{c})^2 = 2 - 2\mathbf{v} \cdot \mathbf{c}$. Then combining this with Eq. (2), (3), (6), (7), (8), and (9), we can get the maximum log-likelihood estimation with:

$$\begin{aligned} \theta^* &= \arg \min_{\theta} \sum_{i=1}^N -\log \frac{\exp(\mathbf{v}_i \cdot \mathbf{c}_p / \tau_p)}{\sum_{j=1}^K \exp(\mathbf{v}_i \cdot \mathbf{c}_j / \tau_j)} \iff \\ \theta^* &= \arg \min_{\theta} \sum_{k=1}^K \sum_{i=1}^{N_k} -\log \frac{\exp(\mathbf{v}_i^k \cdot \mathbf{c}_k / \tau_k)}{\sum_{j=1}^K \exp(\mathbf{v}_i^k \cdot \mathbf{c}_j / \tau_j)}, \end{aligned} \quad (10)$$

where \mathbf{v}_i^k denotes the encoded features of i^{th} sample (i.e., skeleton graph representation) belonging to the k^{th} graph prototype \mathbf{c}_k , N_k is the number of samples in the k^{th} sample group, and τ is related to the distribution of encoded graph representations around different graph prototypes.

Assumption 2 Maximum Homogeneous Similarity. *The homogeneous samples, which are defined as samples within*

the same-prototype sample group, should share higher inherent similarity than heterogeneous samples between different groups. In other words, the graph prototype of each sample group can represent the unique skeleton concepts and attributes of a certain identity, and the same group's samples possess the homogeneity of features corresponding to this prototype [3]. According to Assumption 1, it can be equivalent to the objective that the representation of each sample should be maximally similar to the corresponding prototype and be minimally similar to other prototypes. In the M -step, we maximize the probability that each sample representation belongs to its unique prototype (see Eq. (9)) based on this assumption. The equivalent formulation of this objective in Eq. (10) after applying feature ℓ_2 -normalization can be further interpreted as to maximize the dot-product based similarity between samples and their prototypes while maximizing the dissimilarity to other prototypes.

Relations to Existing Contrastive Losses [1,2,4-7]:

1. The InfoNCE loss [4] re-formulated in MoCo [5] and SimCLR [6] can be interpreted as special cases of the maximum log-likelihood estimation in Eq. (10), where the prototype c_p for a feature v_i is replaced by the augmented feature v'_i generated from different views of augmentation of the same instance (*i.e.*, $c_p = v'_i$) and τ is fixed as a temperature for contrastive learning.
2. The masked prototype contrastive (MPC) loss in [1] and skeleton prototype contrastive (SPC) loss in [2] can be viewed as *unsupervised* generalized versions of the objective in Eq. (10). Both of them leverage unsupervised clustering algorithms (*e.g.*, DBSCAN [8]) to generate *class-agnostic* prototypes c_p . The SPC loss exploits the multi-scale graph features of a skeleton sequence as v_i , while the MPC loss replaces it with the features of random skeleton subsequences of a sequence. However, the instability (*e.g.*, varying cluster numbers caused by over-clustering) or/and unreliability of the used *identity-agnostic* prototypes (*e.g.*, lower confidence to characterize a ground-truth identity) largely limit their practical performance.
3. The ProtoNCE loss used in PCL [7] is a combination of momentum-based contrastive (MoCo) learning [5] and unsupervised prototype estimation with k -means clustering. It has a similar form as Eq. (10), where τ is estimated with the assumption that the distribution of feature representations around each prototype varies in different clusters. However, PCL estimates the feature distribution under the Euclidean distance metric used in the k -means clustering. Such estimation could be inapplicable (*e.g.*, can not be generalized) to models that employ different clustering algorithms (*e.g.*,

density-based DBSCAN [8]) or/and different distance metrics (*e.g.*, Jaccard metric), thus failing to getting satisfactory performance in practice [1].

Temperatures. In our work, we adopt a generic form following the common practice [5,6,9], *i.e.*, setting a global temperature τ for the proposed approach. By assuming a uniform feature distribution around each instance (*i.e.*, $\tau = \tau_k = \tau_j$), we encourage the model to learn representations with higher global uniformity, which could improve the quality of contrastive representation learning as theoretically and empirically proved in [1,10,11].

In the proposed approach, each ground-truth identity is represented with a unique graph prototype, which is generated by computing the class centroid of encoded graph representations. We combine both *sequence-level* and *skeleton-level* graph representations to perform the graph prototype contrastive (GPC) learning, so as to learn more identity-associated graph semantics from different levels. The proposed sequence-level ($\mathcal{L}_{\text{GPC}}^{\text{seq}}$) and skeleton-level GPC loss ($\mathcal{L}_{\text{GPC}}^{\text{ske}}$) can be formulated based on Eq. (10) as:

$$\mathcal{L}_{\text{GPC}}^{\text{seq}} = \sum_{k=1}^K \sum_{i=1}^{N_k} -\log \frac{\exp(v_i^k \cdot c_k / \tau_1)}{\sum_{j=1}^K \exp(v_i^k \cdot c_j / \tau_1)}, \quad (11)$$

$$\mathcal{L}_{\text{GPC}}^{\text{ske}} = \sum_{k=1}^K \sum_{i=1}^{N_k} \sum_{t=1}^F -\log \frac{\exp(\hat{v}_{i,t}^k \cdot \hat{c}_k / \tau_2)}{\sum_{j=1}^K \exp(\hat{v}_{i,t}^k \cdot \hat{c}_j / \tau_2)}, \quad (12)$$

where c_j , \hat{c}_j denotes the j^{th} graph prototype and its linear projection, *i.e.*, $\hat{c}_j = \mathcal{F}(c_j)$, v_i^k (equivalent to $\mathcal{S}_{k,j}$ in the paper) denotes the graph representation of the i^{th} skeleton sequence belonging to the k^{th} class, $\hat{v}_{i,t}^k$ (equivalent to $\mathcal{F}_1(s_{k,j}^t)$ in the paper) represents the linear projection of graph representation of the t^{th} skeleton in the i^{th} sequence belonging to the k^{th} identity, τ_1 and τ_2 represent the global temperatures for sequence-level and skeleton-level contrastive learning, and $\mathcal{F}_1(\cdot)$, $\mathcal{F}(\cdot)$ are linear projection heads to transform skeleton-level graph representations and graph prototypes into the same contrastive space \mathbb{R}^d . It is worth noting that the graph prototypes are generated from higher level (*i.e.*, sequence-level) representations and the learnable linear projection in Eq. (11) can be viewed as integrating related graph features from both levels for contrastive learning. Overall, the proposed GPC loss can be viewed as a generalization of existing skeleton prototype contrastive losses [1,2] with: (1) Full skeletal relation encoding, which exploits Skeleton Graph Transformer (SGT) to *simultaneously* captures structural and actional relations from both adjacent and non-adjacent body joints (see Sec. 3.2 of the paper); (2) Fine-grained skeleton semantics learning, which combines graph prototype contrastive learning

(sequence and skeleton level semantics, see Sec. 3.3 of the paper) and graph structure-trajectory prompted reconstruction (graph and node level semantics, see Sec. 3.4 of the paper).

Convergence Proof

We prove the convergence of GPC under modeling the maximum log-likelihood estimation (see Eq. (10)). Recall Eq. (2) and (3) and let

$$\begin{aligned}\ell(\theta) &= \sum_{i=1}^N \log p(\mathbf{x}_i; \theta) \\ &= \sum_{i=1}^N \log \sum_{j=1}^K p(\mathbf{x}_i, \mathbf{c}_j; \theta) \\ &= \sum_{i=1}^N \log \sum_{j=1}^K Q(\mathbf{c}_j) \frac{p(\mathbf{x}_i, \mathbf{c}_j; \theta)}{Q(\mathbf{c}_j)} \\ &\geq \sum_{i=1}^N \sum_{j=1}^K Q(\mathbf{c}_j) \log \frac{p(\mathbf{x}_i, \mathbf{c}_j; \theta)}{Q(\mathbf{c}_j)}.\end{aligned}\quad (13)$$

The above inequality holds with equality when $Q(\mathbf{c}_j) = p(\mathbf{c}_j; \mathbf{x}_i, \theta)$ is a constant (see Eq. (4)).

In the t^{th} E-step, we have estimated the constant value $Q^{(t)}(\mathbf{c}_j) = p(\mathbf{c}_j; \mathbf{x}_i, \theta^{(t)})$ based on the ground-truth class distribution. Then we have:

$$\ell(\theta^{(t)}) = \sum_{i=1}^N \sum_{j=1}^K Q^{(t)}(\mathbf{c}_j) \log \frac{p(\mathbf{x}_i, \mathbf{c}_j; \theta^{(t)})}{Q^{(t)}(\mathbf{c}_j)}.\quad (14)$$

For the t^{th} M-step, we fix $Q^{(t)}(\mathbf{c}_j) = p(\mathbf{c}_j; \mathbf{x}_i, \theta^{(t)})$ and train model parameters θ to maximize Eq. (14). In this way, we can always have:

$$\begin{aligned}\ell(\theta^{(t+1)}) &\geq \sum_{i=1}^N \sum_{j=1}^K Q^{(t)}(\mathbf{c}_j) \log \frac{p(\mathbf{x}_i, \mathbf{c}_j; \theta^{(t+1)})}{Q^{(t)}(\mathbf{c}_j)} \\ &\geq \sum_{i=1}^N \sum_{j=1}^K Q^{(t)}(\mathbf{c}_j) \log \frac{p(\mathbf{x}_i, \mathbf{c}_j; \theta^{(t)})}{Q^{(t)}(\mathbf{c}_j)} \\ &= \ell(\theta^{(t)}).\end{aligned}\quad (15)$$

The above result that $\ell(\theta^{(t)})$ monotonically increases with more iterations suggests the convergence of the algorithm.

The detailed convergence properties of the EM algorithm are discussed in [12, 13]. Here we only discuss the general case, and follow [13] to make the assumptions for the EM algorithm:

- **(a)** Ω is a subset in the r -dimensional Euclidean space \mathbb{R}^r .

- **(b)** $\Omega_{\theta^{(0)}} = \{\theta \in \Omega : \ell(\theta) \geq \ell(\theta^{(0)})\}$ is compact for any $\ell(\theta^{(0)}) > -\infty$.
- **(c)** $\ell(\cdot)$ is continuous in Ω and differentiable in the interior of Ω .

Under the assumptions of **(a)**, **(b)**, and **(c)**³, we have:

- **(d)** $\{\ell(\theta^{(t)})\}_{t \geq 0}$ is bounded above for any $\theta^{(0)} \in \Omega$. As a consequence of **(d)** and the inequality (15) (i.e., $\ell(\theta^{(t+1)}) \geq \ell(\theta^{(t)})$), $\ell(\theta^{(t)})$ converges monotonically to some ℓ^* .

It is worth noting that there is no guarantee that ℓ^* is the global maximum of $\ell(\cdot)$ over Ω . As reported in previous works [13–17], if the log-likelihood function $\ell(\cdot)$ has several (local or global) maxima and stationary points, the convergence of the EM sequence $\{\ell(\theta^{(t)})\}$ to either type of point depends on the choice of starting point. Readers can refer to [12, 13] for more details about different convergence cases of the EM algorithm.

The aforementioned fact may account for the performance changes (i.e., small variations) of our model on the same dataset, as different random initializations of model parameters could change $\theta^{(0)}$ (i.e., the starting point) hence the final convergence result. In practice, we follow [1, 18] to train the model with different random initializations on each dataset and report its average performance, which helps estimate a more stable EM convergence result with different initialized starting points.

References

- [1] H. Rao and C. Miao, “SimMC: Simple masked contrastive learning of skeleton representations for unsupervised person re-identification,” in *International Joint Conference on Artificial Intelligence (IJCAI)*, pp. 1290–1297, 2022. 2, 3, 4
- [2] H. Rao and C. Miao, “Skeleton prototype contrastive learning with multi-level graph relation modeling for unsupervised person re-identification,” *arXiv preprint arXiv:2208.11814*, 2022. 2, 3
- [3] M. Nauta, J. Trienes, S. Pathak, E. Nguyen, M. Peters, Y. Schmitt, J. Schlöterer, M. van Keulen, and C. Seifert, “From anecdotal evidence to quantitative evaluation methods: A systematic review on evaluating explainable AI,” *arXiv preprint arXiv:2201.08164*, 2022. 3
- [4] A. van den Oord, Y. Li, and O. Vinyals, “Representation learning with contrastive predictive coding,” *arXiv preprint arXiv:1807.03748*, 2018. 3
- [5] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, “Momentum contrast for unsupervised visual representation learning,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9729–9738, 2020. 3

³These assumptions can be satisfied in most practical situations. As the related proofs/discussions are out of the scope of this work, readers can refer to [13] for more details.

- [6] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, “A simple framework for contrastive learning of visual representations,” in *International Conference on Machine Learning (ICML)*, pp. 1597–1607, 2020. 3
- [7] J. Li, P. Zhou, C. Xiong, and S. Hoi, “Prototypical contrastive learning of unsupervised representations,” in *International Conference on Learning Representation (ICLR)*, 2021. 3
- [8] M. Ester, H.-P. Kriegel, J. Sander, X. Xu, *et al.*, “A density-based algorithm for discovering clusters in large spatial databases with noise,” in *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, vol. 96, pp. 226–231, 1996. 3
- [9] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin, “Unsupervised feature learning via non-parametric instance discrimination,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3733–3742, 2018. 3
- [10] T. Wang and P. Isola, “Understanding contrastive representation learning through alignment and uniformity on the hypersphere,” in *International Conference on Machine Learning (ICML)*, pp. 9929–9939, 2020. 3
- [11] T. Gao, X. Yao, and D. Chen, “SimCSE: Simple contrastive learning of sentence embeddings,” in *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 6894–6910, 2021. 3
- [12] G. J. McLachlan and T. Krishnan, *The EM algorithm and extensions*. John Wiley & Sons, 2007. 4
- [13] C. J. Wu, “On the convergence properties of the EM algorithm,” *The Annals of statistics*, pp. 95–103, 1983. 4
- [14] V. Hasselblad, “Estimation of finite mixtures of distributions from the exponential family,” *Journal of the American Statistical Association*, vol. 64, no. 328, pp. 1459–1471, 1969. 4
- [15] J. H. Wolfe, “Pattern clustering by multivariate mixture analysis,” *Multivariate behavioral research*, vol. 5, no. 3, pp. 329–350, 1970. 4
- [16] N. Laird, “Nonparametric maximum likelihood estimation of a mixing distribution,” *Journal of the American Statistical Association*, vol. 73, no. 364, pp. 805–811, 1978. 4
- [17] D. B. Rubin and D. T. Thayer, “EM algorithms for ml factor analysis,” *Psychometrika*, vol. 47, no. 1, pp. 69–76, 1982. 4
- [18] H. Rao, S. Wang, X. Hu, M. Tan, Y. Guo, J. Cheng, X. Liu, and B. Hu, “A self-supervised gait encoding approach with locality-awareness for 3D skeleton based person re-identification,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, no. 01, pp. 1–1, 2021. 4