

1 We sincerely thank the reviewers for their valuable comments. We proofread and fixed the mentioned errors.

2 **Broader Impact:** Since our method improves subspace clustering, it advances learning from unannotated data.  
3 Improving the learning process and providing more accurate similarity matrices for unannotated data can positively  
4 impact accountability, transparency and explainability of AI methods. However, if not controlled, providing the  
5 opportunity to learn from big unannotated datasets could increase the concerns about violating the privacy of individuals.

6 **Related Work:** Thank you for the additional references. We will include and discuss them in the revised version.

7 **Publishing codes:** Upon the acceptance of our paper, we will publicly release the source codes.

8 **Equation 9:** While in MLRDSC (Eq. (9)),  $Q^T$  is the same as  $\hat{Q}^T$ , in our method,  $Q^T$  relates to  $\hat{Q}^T$  by Eq. (4).

9 **Replacing  $Q_{Ens}^t$  with  $\hat{Q}$**  in Eq. (7) could also make an alternative consistency loss. However, using  $Q_{Ens}^t$  has the  
10 advantages of: 1) Ensembling multiple subspace membership predictions which results in a better estimate. 2) Providing  
11 smoother transitions for subspace membership predictions that could result in better learning.

12 **Motivations behind augmentations for subspace clustering:** Additional unlabeled data can only improve the repre-  
13 sentation learning in encoder and decoder, and not necessarily the self-expressiveness layer’s parameters (which scales  
14 quadratically with additional data). Augmentations in our model regularize parameters in both encoder and decoder  
15 as well as the self-expressiveness layer. Intuitively, if two samples are similar, our method makes sure they remain  
16 similar even if they undergo transformations that do not change their label. Adding additional data does not provide  
17 such property.

18 **Rationale Behind the Augmentation Search:** Note that our augmentations search algorithm does not minimize  
19 Eq. (8). It looks for augmentation policies that yield the highest mean Silhouette coefficient. The rationale is similar  
20 to supervised models. Augmentations may initially increase the value of the loss function, but they lead to improved  
21 learning by regularizing the training.

22 **Relation to metric learning and stochastic subspace clustering:** Our model shares more similarities to stochastic  
23 subspace clustering (reference [17] in the submission) models (SSSC) than to metric learning methods. In SSSC,  
24 stochastic transformations are applied to the inputs of the subspace clustering algorithm. In addition, many subspace  
25 algorithms apply nonlinear functions to deal with nonlinearity (references [8,9,10,11] in the submission). We add the  
26 discussion in the final paper.

27 **Silhouette coefficient:** Silhouette coefficient measures how similar an object is to its own cluster compared to other  
28 clusters. As we assume in a clustering task we do not have any labels, an *external evaluation* such as Silhouette can be  
29 used as an estimation of clustering performance to avoid bad augmentations.

30 **Experiments:** We include the following experiments in the revision:

- 31 • **Runtime:** In Table 1, we report the runtime for the policy search algorithm on the COIL-20 dataset.
- 32 • **Silhouette Coefficient v.s. Ground-truth:** We tested our augmentation algorithm with accuracy as the Score on the  
33 COIL-20 dataset, and observed that for a set of 9 best augmentation policies it achieves the same clustering error rate  
34 of 1.79. Our found augmentations in Table 1 of the main paper are also among these best augmentation policies.
- 35 • **Search Baselines:** In Table 2, we compare the clustering error rate of training our method with the augmentations  
36 found by 1) our algorithm, 2) uniform sampling, and 3) coordinate descent. We test on the Extended Yale-B dataset.
- 37 • **Downstream tasks:** We use the latent space features from our network and from MLRDSC to perform a classification  
38 task. We randomly set 50% of the learned latent space features for ORL samples as test set and use the remaining  
39 samples in training an SVM model. Table 3 shows that features learned via our method are better for classification.
- 40 • **Training deeper networks:** We add two more layers to DSC network (denoted by DSC-5layer). Table 4 compares  
41 the clustering error rate of DSC-5layer with and without augmentations applied to the ORL dataset.

Table 1: Runtime.	Table 2: Error rates on Yale-B. Search baselines.			Table 3: Downstream task: Accuracy of SVM on ORL.		Table 4: Error rates on ORL. Deeper networks: DSC-5layer.	
COIL-20	Ours	Uniform	Coordinate Descent	Ours	MLRDSC	Without Aug.s	With Aug.s
261 mins	<b>0.82</b>	3.72	0.95	<b>95.5</b>	93	22.50	<b>13.25</b>

42 **Details:** We include the following details in the revision:

- 43 • **EMA decay:** We used the silhouette coefficient as an evaluation metric in cross-validation.
- 44 • **Statistical Significance:** In all the conducted experiments, we reported 5-fold averages.
- 45 • **Minibatches:** Similar to other DSC methods, we input the whole dataset as a batch.
- 46 • **Policies with  $< \ell_{max}$  sub-policies** were also considered as augmentation candidates.