
Batch Curation for Unsupervised Contrastive Representation Learning

Michael C. Welle^{*1} Petra Poklukar^{*1} Danica Kragic¹

Abstract

The state-of-the-art unsupervised contrastive visual representation learning methods that have emerged recently (SimCLR, MoCo, SwAV) all make use of data augmentations in order to construct a pretext task of instant discrimination consisting of similar and dissimilar pairs of images. Similar pairs are constructed by randomly extracting patches from the same image and applying several other transformations such as color jittering or blurring, while transformed patches from different image instances in a given batch are regarded as dissimilar pairs. We argue that this approach can result similar pairs that are *semantically* dissimilar. In this work, we address this problem by introducing a *batch curation* scheme that selects batches during the training process that are more inline with the underlying contrastive objective. We provide insights into what constitutes beneficial similar and dissimilar pairs as well as validate *batch curation* on CIFAR10 by integrating it in the SimCLR model.

1. Introduction

Extracting data representation that are compact, expressive, and meaningful is a long term goal in the field of machine learning (Bengio et al., 2013). In recent years, most of the efforts have been put on learning general data representations that are suitable for a variety of different downstream task without using human annotation. In computer vision, prior work addressed this by solving pretext tasks defined in a self-supervised manner such as denoising autoencoders (Vincent et al., 2008), solving jigsaw puzzles (Noroozi & Favaro, 2016), or motion segmentation (Pathak et al., 2017). A specific pretext task, namely instance discrimination (Wu et al., 2018), has recently emerged as the state-of-the-art for unsupervised visual representation learning, spawning

a polyhedral of new models based on contrastive learning such as SimCLR (Chen et al., 2020), MoCo (He et al., 2020) or SwAV (Caron et al., 2020).

In these methods, the pretext tasks are constructed by applying various data transformations, for example, cropping, color jittering or gaussian blur, in order to obtain similar and dissimilar pairs in an unsupervised way. Thanks to the thorough ablation study performed by the respective methods, it is established that *random resized cropping*, i.e., extracting a patch of random size from the original image and resizing it back to the size of the original image, is one of the crucial data augmentation for learning useful representations (Chen et al., 2020). The intuition behind it is that patches drawn from the same image contain similar semantic information and should therefore be encoded close to each other. However, we hypothesise that this approach can produce patches containing semantically dissimilar information which decelerate learning as they violate this core similarity assumption.

In this work, we first take a closer look at the random resize cropping and investigate how the size of the patches as well as their relative position to each other in the original image influences the learning process. Based on our findings, we then propose a *batch curation* scheme, which during training constructs similar and dissimilar pairs that are aligned with the contrastive objective. Our scheme is simple as it uses the given model in training to self-curate batches such that similar patches are semantically close as judged by the model itself. It can be easily adapted to any contrastive learning setting with arbitrary number of similar and dissimilar pairs, and can therefore be integrated into any of the current state-of-the-art methods. We validate our batch curation scheme in a small scale experiment by integrating it into the SimCLR (Chen et al., 2020) model trained using CIFAR-10 (Krizhevsky et al., 2009) dataset. We show an improvement of $\approx 1.5\%$ on the K-NN classification evaluation as well as competitive performance on the linear classification protocol. To the best of our knowledge, we are the first to introduce a procedure for providing training batches of higher quality for unsupervised visual contrastive learning methods.

^{*}Equal contribution ¹RPL, KTH Royal Institute of Technology, Sweden. Correspondence to: Michael C. Welle <mwelle@kth.se>.

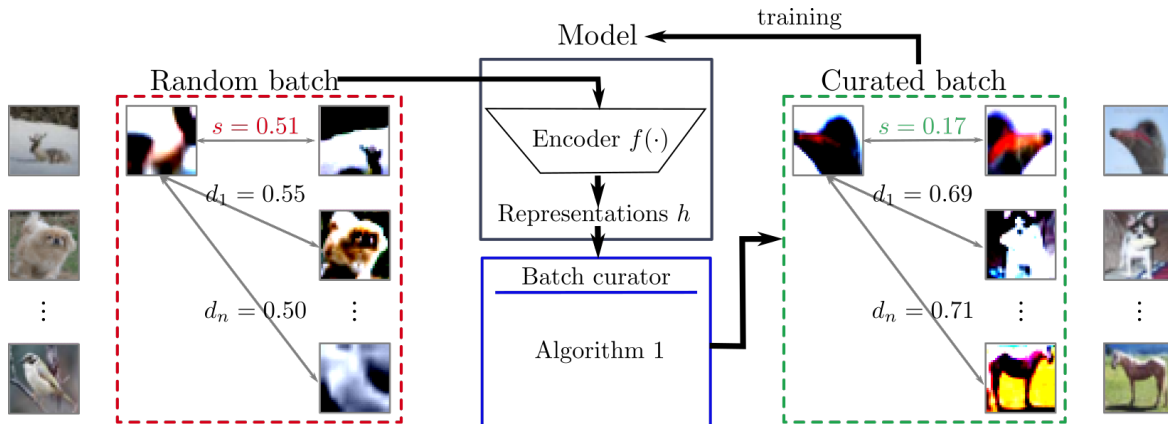


Figure 1. Overview of the proposed batch curation scheme. The curator (middle, blue) uses the model in training (middle, black) to curate a randomly sampled batch (left, red) of images such that the distance between the similar pairs s is smaller than the distance d_i among any dissimilar pair. The curated batch (right, green) is then used for training of the model. The original images from which patches in the batches are extracted are shown on the left and right. To obtain meaningful approximation of the distances among similar and dissimilar pairs, the curation scheme is activated after a *warm-up* period. See Section 3 for more information.

2. Investigation of Patch Configurations

In this section, we perform a deeper investigation on how the patch size and their pairwise relative position in the original image influences the learning process. Our analysis is based on extracting two patches from the same image similar to the approach taken in (Chen et al., 2020). All the investigations in this paper are performed on the CIFAR-10 (Krizhevsky et al., 2009) dataset containing images of dimension 32×32 .

2.1. Patch Configurations

We first take a closer look into possible relative positions that two patches extracted from an image can attain. In Fig. 2 we show the three possible configurations: a) *global-local view* where one patch is completely contained in the other, b) *adjacent view* where patches have no intersection, and c) *intersection view* where patches have a non-empty intersection but are not contained in one another.

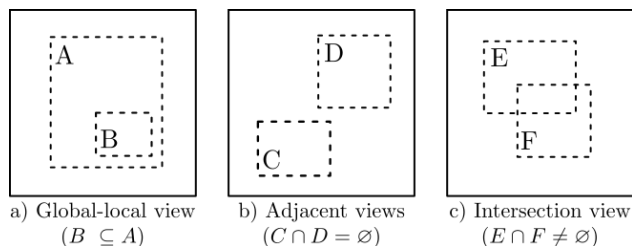


Figure 2. Three different configurations of two patches drawn from an image. From left to right: a) global-local view, b) adjacent view, and c) intersection view.

We perform random resized cropping using default scale $[0.08, 1]$ and range $[0.75, 1.33]$ parameters in Pytorch implementation¹ one million times on CIFAR-10 images, and show the occurrence of the three different configurations in the right of Fig. 3. We can see that the majority of patch configurations are intersection views (81.33%), while the global-local views occur in 17.27% of cases. The adjacent views occur in only 1.4% of cases which is expected since the average patch size is 49% of the original image. On the left of Fig. 3 we display a normalised coverage heatmap of the sampled patches onto an image where pixels in blue denote the ones that are covered the most. As expected, we can observe a clear bias towards the middle of the image, which might be more beneficial for object centered datasets such as CIFAR-10 and ImageNet (Russakovsky et al., 2015).

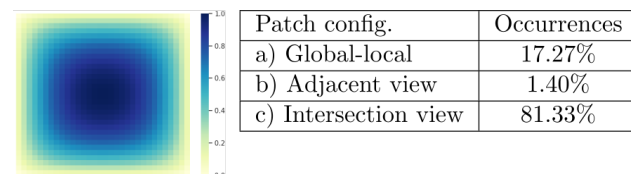


Figure 3. Normalised coverage heatmap of one million randomly drawn patches (left) and their patch configurations (right).

¹<https://pytorch.org/vision/stable/transforms.html#torchvision.transforms.RandomResizedCrop>

2.2. Performance of Different Patch Configurations on SimCLR

Given the insights of patch configuration occurrences, we investigate the performance of a SimCLR when limiting the patches we sample to either *i)* only global-local views, *ii)* only adjacent views, *iii)* only intersection views, and *iv)* a version where all configurations occur with equal probability. We trained all models for 500 epochs using batch size 128, feature dimension 128, and temperature 0.5. We evaluated them using K-NN evaluation where $K = 200$, and linear evaluation protocol where one linear layer is trained on the obtained representations for 100 epochs.

We show in Table 1 that the best performing models are the one using only intersection views and the default SimCLR (both obtaining $\approx 90.5\%$ on linear evaluation and 86.5% on K-NN evaluation). As we showed in Fig. 3, the latter naturally contains a high percentage of intersection views. Enforcing just global-local patches as well as using an equal amount of each configuration performs slightly worse, while the model employing only adjacent views performs the worst (only $\approx 65\%$ on both linear and K-NN evaluations). We observed that both the model containing equal amount of each configuration and using adjacent views have smaller average patch size.

SimCLR model version	Linear acc.	K-NN acc.	avg. p. size
Different patch configurations			
default	90.67%	86.29%	49%
global-local	87.32%	81.95%	51%
adjacent view	64.59%	65.96%	17%
intersection view	90.72%	86.63%	49%
equal configuration	88.60%	83.95%	39%
Different patch sizes			
big patches	89.17%	85.99%	70%
small patches	23.24%	76.31%	29%
global-local big patches	86.50%	82.87%	73%

Table 1. Best accuracy of SimCLR models trained with different patch configurations on a test set using K-NN evaluation (left) as well as linear evaluation (middle). Right column shows the average patch size in relation to the original image size.

In order to further investigate the effect of different patch sizes, we evaluate three additional models: (i) the default SimCLR, (ii) a global-local version both with scale $[0.5, 1]$ producing larger patches, and (iii) the default SimCLR with scale $[0.08, 0.5]$ producing smaller patches. Fig. 4 shows the K-NN Top-1 accuracy obtained during training for all models. We observe that bigger patches lead to a faster increase in performance but do not yield a better performance in the long run. We also see that the big patch version of global-local views performs slightly better than the default global-local in the K-NN evaluation but slightly worse in the linear evaluation protocol (see Table 1). The default SimCLR with small patches performs significantly worse

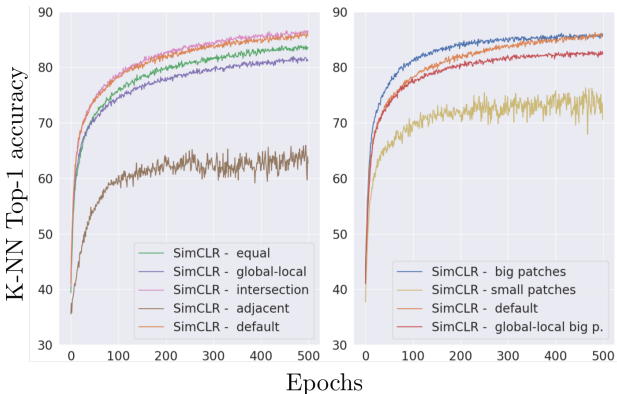


Figure 4. Top-1 accuracy of SimCLR models obtained during training with different patch configurations on a test set using K-NN evaluation.

obtaining only $\approx 23\%$ accuracy on the linear evaluation. In summary, we find that:

- Adjacent views seem to negatively affect the resulting representations possibly due to their small size.
- Larger patches on average give a boost at the beginning of the training but do not increase the overall performance.
- Intersecting views are more beneficial than global-local or adjacent views for the SimCLR model.

3. Batch Curation

Motivated by the differences in the performance of the SimCLR models when varying the type of patch configuration presented in the previous section, we hypothesise more generally that the quality of batches used during training of any unsupervised contrastive visual representation learning model is directly correlated to the quality of the learned representations.

More precisely, from the poor performance of the adjacent configuration of patches in Table 1, we hypothesise that sampling small patches that do not contain any relevant semantic information hinders the learning process of the model in training. We conjecture that this results in an ill-defined pretext task where patches that are labeled similar are in fact not semantically similar. To this end, we propose a *batch curation* scheme presented in the next section.

3.1. The Batch Curator

The proposed *batch curation scheme* curates a randomly sampled batch such that similar pairs of patches are always closer together than any of the dissimilar pairs. We visualize its principles in Fig. 1 using a contrastive learning setting

Algorithm 1 General batch curation scheme.

Require: Set of similar images \mathcal{S}

Require: Set of dissimilar images \mathcal{D}

```

1:  $d_S = d_D \leftarrow 0$ 
2: while not  $d_S < d_D$  do
3:    $h_S \leftarrow f(\mathcal{S}), h_D \leftarrow f(\mathcal{D})$  [extract representations]
4:    $M_S \leftarrow \text{compute distance}(h_S, h_S)$  [distance among similar pairs]
5:    $M_D \leftarrow \text{compute distance}(h_S, h_D)$  [distance among dissimilar pairs]
6:    $d_S \leftarrow \max(M_S)$ 
7:    $d_D \leftarrow \min(M_D)$ 
8:   resample images in  $\mathcal{S}$  and  $\mathcal{D}$  violating  $d_S < d_D$ 
9: end while
10: return  $\mathcal{S}, \mathcal{D}$ 

```

with one similar pair and several dissimilar pairs².

The details of our procedure are outlined in Algorithm 1. Firstly, we encode the randomly sampled patches using the model trained in the previous epoch to obtain their representations. Secondly, we use these representations to calculate the distances M_S, M_D between similar and dissimilar pairs, respectively. We accept the batch if the largest distance between similar pairs d_S is smaller than the smallest distance among the dissimilar pairs d_D . If this criteria is not fulfilled, we re-sample pairs that violate it until we obtain an adequate batch. The curated batch is then used to perform one step of the gradient-descent in the training of the unsupervised contrastive learning model. In order to obtain somewhat meaningful distances in the learned representation space, we start the batch curation scheme after an initial warm-up period of w epochs. Since the curation is simply performed by the model itself, it can be easily integrated to any unsupervised contrastive representation learning method.

SimCLR model version	Linear acc.	K-NN acc.	avg. p. size
default w/o batch curator	90.67%	86.29%	49%
with batch curator	90.81%	87.63%	49%

Table 2. Best accuracy of SimCLR models trained with batch curation scheme (bottom) and without it (top) on a test set using K-NN evaluation (left) as well as linear evaluation (right). Right column shows the average patch size in relation to the image size.

3.2. Performance of the Batch Curator on SimCLR

We validate our approach by comparing the performance of a SimCLR model trained with the proposed batch curation scheme (SimCLR-batch curation) and one without it (SimCLR-default). We trained and evaluated both models as described in Section 2.2. The evaluation results are shown

²Note that this can be easily adapted to arbitrary number of similar or dissimilar pairs.

in Table 2, where we can observe a 1.5% improvement in the K-NN evaluation and a slight improvement in the linear evaluation. In Fig 5, we show the changes in K-NN Top-1 accuracy of the test set during training. We can clearly see the improvement in the performance when batch curation is started, showing that it leads to improved representations.

We emphasise that the presented experimental results do not match the state-of-the-art 92% accuracy on the linear evaluation on CIFAR-10 as reported in (Chen et al., 2020) due to the limited hardware setup. Note that even when using batch size 128, it takes more than 20h to train each SimCLR model using a single GPU. In future, we plan to validate our batch curation scheme on more complex datasets such as ImageNet as well as integrate it into other state-of-the-art unsupervised contrastive representation learning methods such as MoCo and SWaV. However, we leave the extensive scaling, for example as performed in (Goyal et al., 2021) using up to 1.3 billion model parameters and 512 GPUs, to other research groups with larger computational resources.

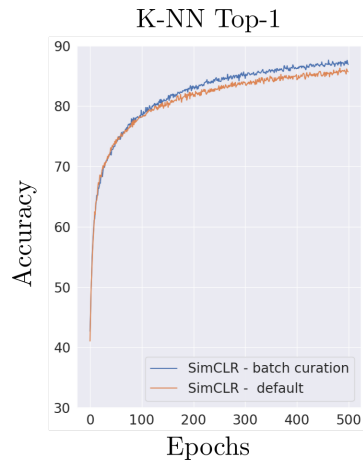


Figure 5. Top-1 accuracy of SimCLR model trained with batch curation scheme (blue) and without it (orange) using K-NN evaluation.

4. Conclusion and Future Work

In this work, we investigated the common transformation applied to images in the unsupervised contrastive visual representation learning domain, namely random resize cropping. We showed that the most beneficial patch configuration consists of intersecting patches, while smaller patches seem to be detrimental. We hypothesise that the latter occurs because patches do not contain enough semantic information. To this end, we proposed a simple batch curator scheme, which curates the training batches to be more aligned with the contrastive learning objective. We show that our scheme improves the end performance when integrated to the SimCLR model on CIFAR-10 dataset. As future work, we plan to perform extensive fine-tuning of our scheme, scaling it to larger datasets such as ImageNet and integrating it to other existing methods.

References

- Bengio, Y., Courville, A., and Vincent, P. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8): 1798–1828, 2013.
- Caron, M., Misra, I., Mairal, J., Goyal, P., Bojanowski, P., and Joulin, A. Unsupervised learning of visual features by contrasting cluster assignments. *arXiv preprint arXiv:2006.09882*, 2020.
- Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PMLR, 2020.
- Goyal, P., Caron, M., Lefaudeaux, B., Xu, M., Wang, P., Pai, V., Singh, M., Liptchinsky, V., Misra, I., Joulin, A., et al. Self-supervised pretraining of visual features in the wild. *arXiv preprint arXiv:2103.01988*, 2021.
- He, K., Fan, H., Wu, Y., Xie, S., and Girshick, R. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9729–9738, 2020.
- Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. 2009.
- Noroozi, M. and Favaro, P. Unsupervised learning of visual representations by solving jigsaw puzzles. In *European conference on computer vision*, pp. 69–84. Springer, 2016.
- Pathak, D., Girshick, R., Dollár, P., Darrell, T., and Hariharan, B. Learning features by watching objects move. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2701–2710, 2017.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3): 211–252, 2015.
- Vincent, P., Larochelle, H., Bengio, Y., and Manzagol, P.-A. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning*, pp. 1096–1103, 2008.
- Wu, Z., Xiong, Y., Yu, S. X., and Lin, D. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3733–3742, 2018.