

---

# Self-supervised visual representation learning: SimCLR & Deep Clustering

---

**Dzodzoenyenye Adjowa Senanou, Prince Mensah, Najlaa Mohamed, Leonard Sanya**

African Masters for Machine Intelligence (AMMI)

dsenanou@aimsammi.org, pmensah@aimsammi.org

nmohamed@aimsammi.org, lsanya@aimsammi.org

## Abstract

1 Manual annotation of large collections of unlabeled data is extremely expensive.  
2 These large-scale datasets come with a variety of complexity and higher dimensionality.  
3 Several self-supervised methods have been developed to face challenges  
4 in analyzing unlabelled data. In this work, we provided a comprehensive review of  
5 two methods; Simple Contrastive Learning Representations (SimCLR) and Deep  
6 Clustering.

7 SimCLR learns robust data representations by simultaneously optimizing for contrastive  
8 similarity among augmented views of the same instance and dissimilarity  
9 between views of different instances.

10 Deep Clustering is a combination of Deep Learning and clustering approaches. So it  
11 generally retrieve information about features and perform clustering simultaneously  
12 basing on them. With the efficiency of deep learning, the accuracy of the clusters  
13 outperformed the standard clustering approach.

## 1 Introduction

15 Building machine learning models typically requires large-scale annotated datasets, which demand  
16 significant manual annotation effort, time, and expense. To alleviate the burden of manual data  
17 annotation and reduce memory costs, self-supervised learning techniques have emerged as a viable  
18 solution [1]. These techniques leverage unsupervised learning capabilities to perform tasks tradi-  
19 tionally requiring supervised learning. Unlike supervised learning, where ground truth labels are  
20 provided by humans, self-supervised learning generates these labels during the learning process from  
21 unstructured or unlabeled data.

22 Self-supervised learning is a notable form of unsupervised learning where "pseudo-labels" are derived  
23 directly from raw input data and used in place of human-annotated labels through pretext tasks [2].  
24 This approach reduces the dependency on labeled data, making it scalable to large datasets. The  
25 technique is easily generalized across various tasks and domains, and the learned representations can  
26 be fine-tuned for specific tasks, showcasing its transfer learning capabilities.

27 Self-supervised learning typically involves two main stages: the pretext task (representation learning)  
28 and the downstream task [3]. The pretext task focuses on learning useful representations from data  
29 without needing labeled examples, while the downstream task applies these learned representations  
30 to specific applications, often incorporating labeled data. Multiple techniques exist for training  
31 self-supervised models, primarily divided into discriminative and generative methods. Contrastive  
32 learning, a discriminative approach, trains models to differentiate between pairs of data points,  
33 enabling the discovery of valuable features by comparing data pairs. In contrast, self-predictive  
34 techniques, a generative approach, train models to generate labels from the data itself, predicting  
35 certain aspects of the input data based on other elements.

36

37 In this work , we focus on **SimCLR** and **Deep Clustering**.

38 **2 SimCLR: A simple framework for contrastive learning of visual  
39 representations**

40 SimCLR is a contrastive learning approach that maximises agreement across various augmented  
41 views of a given data example through a contrastive loss in a latent space, thereby learning represen-  
42 tations.

43 **2.1 SimCLR Architecture**

44 There are four main parts to the SimCLR architecture.

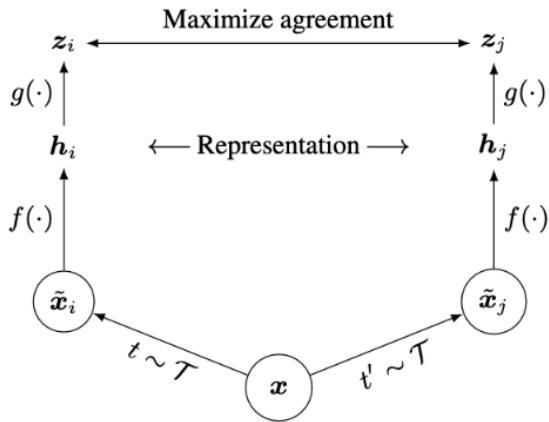


Figure 1: SimCLR Architecture [4]

- 45 **1. Data augmentation:** This module converts an example of data into two views of the same  
46 example that are related to one another. Operations for enhancing data include rotation,  
47 Gaussian blur, random colour distortion, random cropping, and more. After augmentation  
48 with operations  $t$  and  $t'$ ,  $x$  is shown in Figure 1 in two views:  $\hat{x}_i$  and  $\hat{x}_j$ . Once augmentation  
49 is applied, we will have  $2N$  training examples.
- 50 **2. Neural networks base encoder  $f(\cdot)$ :** Features that are representative of the input images  
51 are extracted. There are no restrictions when selecting a neural network architecture based  
52 on an encoder module thanks to the SimCLR framework. ResNet is the neural network  
53 architecture that is most frequently utilised. The results of applying the encoder are  $h_i$  and  
54  $h_j$ , where  $h = f(\hat{x}_i)$  and  $h = f(\hat{x}_j)$  respectively, and are d-dimensional outputs following  
55 average pooling.
- 56 **3. Small neural networks projection head  $g(\cdot)$ :** Maps the representations to a new space  
57 where the contrastive loss function is applied. They used a one hidden layer MLP to obtain  
58  $z_i = g(h_i) = W^{(2)}\sigma(W^{(1)}h_i)$  where  $\sigma$  is a ReLU nonlinearity.
- 59 **4. Contrastive loss function:** The contrastive prediction task tries to find  $\hat{x}_j$  in  $\{\hat{x}_k\}_{k \neq i}$  for  
60 a given  $\hat{x}_i$ , the collection  $\{\hat{x}_k\}$  includes a positive pair of samples  $\hat{x}_i$  and  $\hat{x}_j$ . The other  
61  $2(N - 1)$  augmented examples in a minibatch are treated as negative examples given a  
62 positive pair.  
63 Denote the dot product between  $l_2$  normalised  $u$  and  $v$  (i.e., cosine similarity) as  $\text{sim}(u, v) =$   
64  $u^T \cdot v / \|u\| \|v\|$ . Next, we define the loss function for a positive pair of samples  $(i, j)$  as  
65 follows:

$$\mathcal{L}_{i,j} = -\log \frac{\exp(\text{sim}(z_i, z_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}[k \neq i] \exp(\text{sim}(z_i, z_k)/\tau)},$$

66 where  $1[k \neq i] \in \{0, 1\}$  is an indicator function evaluating to 1 if  $k \neq i$ , and  $\tau$  denotes a  
 67 temperature parameter. The final loss is computed across all positive pairs, both  $(i, j)$  and  
 68  $(j, i)$ , in a mini-batch.

---

**Algorithm 1:** SimCLR's main learning algorithm [4]
 

---

```

input: batch size  $N$ , constant  $\tau$ , structure of  $f, g, T$ 
for sampled minibatch  $\{x_k\}_{k=1}^N$  do
  for all  $k \in \{1, \dots, N\}$  do
    draw two augmentation functions  $t \sim T, t_0 \sim T$ ;      // the first augmentation
     $\tilde{x}_{2k-1} = t(x_k)$ ;                                // augmented example
     $h_{2k-1} = f(\tilde{x}_{2k-1})$ ;                            // representation
     $z_{2k-1} = g(h_{2k-1})$ ;                              // projection
     $\tilde{x}_{2k} = t_0(x_k)$ ;                               // the second augmentation
     $h_{2k} = f(\tilde{x}_{2k})$ ;                            // representation
     $z_{2k} = g(h_{2k})$ ;                                // projection
  end
  for all  $i \in \{1, \dots, 2N\}$  and  $j \in \{1, \dots, 2N\}$  do
     $s_{i,j} = \frac{z_i^\top z_j}{k \|z_i\| \|z_j\|}$ ;           // pairwise similarity
  end
  Define  $\mathcal{L}(i, j) = -\log \frac{\exp(s_{i,j}/\tau)}{\sum_{k=1}^{2N} 1[k \neq i] \exp(s_{i,k}/\tau)}$ ;          // positive pair loss
   $L = \frac{1}{2N} \sum_{k=1}^N [\mathcal{L}(2k-1, 2k) + \mathcal{L}(2k, 2k-1)]$ ;          // total loss
  update networks  $f$  and  $g$  to minimize  $L$ 
end
return encoder network  $f(\cdot)$ , and discard  $g(\cdot)$ 
  
```

---

### 70 3 Deep Clustering

71 This approach is a clustering based approach of unsupervised learning that is being applied and  
 72 studied in computer vision. Deep clustering technique simultaneously learns a neural network's  
 73 parameters and clusters according to the generated features. Using the traditional clustering approach  
 74 k-means, DeepCluster groups the features iteratively and updates the network's weights using the  
 75 succeeding assignments as supervision.

#### 76 3.1 DeepCluster Process Illustration

77 We used CIFAR-10 dataset which consists of 60,000 32x32 color images in 10 different classes, in  
 78 the input as shown in Figure 2 passed through Convnet model (AlexNet Backbone) feature extraction.  
 79 Principal Component Analysis (PCA) is then used to reduce the dimension of these characteristics in  
 80 order to improve the efficiency of the clustering process and lower noise levels in the data. K-Means  
 81 is used to cluster the reduced features, allocating a picture to a cluster.

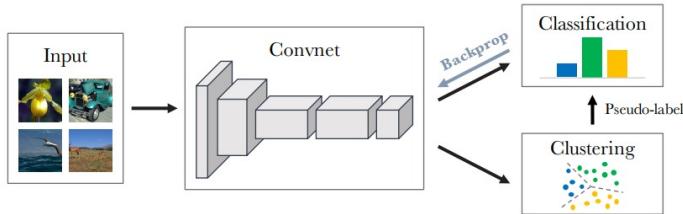


Figure 2: DeepCluster Process Illustration [2]

82 The convnet's characteristics  $f_\theta(x_n)$  are sent into the k-means algorithm, which sorts them into  $k$   
 83 different groups according to a geometric criterion. In other words, it solves the following problem to  
 84 jointly learn the cluster assignments  $y_n$  of each picture  $n$  and a  $d \times k$  centroid matrix  $C$ .

$$\min_{C \in \mathbb{R}^{d \times k}} \frac{1}{N} \sum_{n=1}^N \min_{y_n \in \{0,1\}^k} \|f_\theta(x_n) - Cy_n\|_2^2 \quad \text{such that} \quad y_n^\top 1_k = 1$$

85 For each  $n \leq N$ , solving this problem yields a centroid matrix  $C^*$  and a set of optimal assignments  
 86  $y_n^*$  [2]. As pseudo-labels, these cluster assignments are simply labels that are inferred from the  
 87 data itself without any human annotation. These pseudo-labels are used to refine the Convnet model  
 88 during training. A loss function, cross-entropy is used to test the model's predictions against the  
 89 pseudo-labels, and the model weights are adjusted to minimise this loss.

$$\min_{\theta, W} \frac{1}{N} \sum_{n=1}^N \ell(g_W(f_\theta(x_n)), y_n),$$

90 where  $\ell$  is the multinomial logistic loss, also known as the negative log-softmax function.

## 91 4 Experiments

### 92 4.1 SimCLR Experiment 1:

93 In our experiments, we first trained the SimCLR model for 50 epochs and then for 100 epochs using  
 94 this configuration: ResNet18 as the feature extractor, a projection head with two linear layers and a  
 95 ReLU activation, a contrastive loss function, and the LARS optimizer. Despite extending training  
 96 to 100 epochs, the results remained similar, suggesting that our batch size of 128 was insufficient  
 97 specially that LARS optimizer performs the best with large batch size. This limitation due to memory  
 98 limitations likely hindered the model from generating distinctive embeddings due to the lack of  
 99 diverse negative samples.

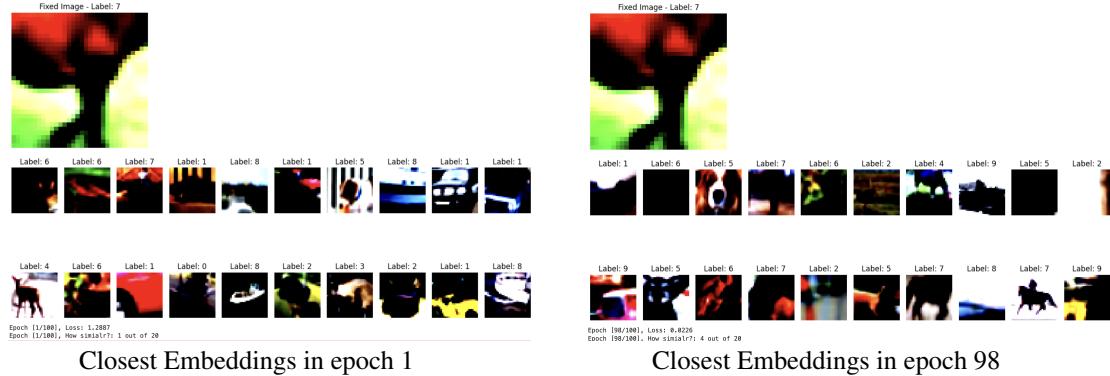


Figure 3: Closest 20 embeddings during training

### 100 4.2 SimCLR Experiment 2:

101 In the second experiment, T-SNE was used to visualize the embedding space throughout the 100  
 102 epochs. The embeddings showed minimal change, indicating that the small batch size constrained the  
 103 model's ability to learn effectively.

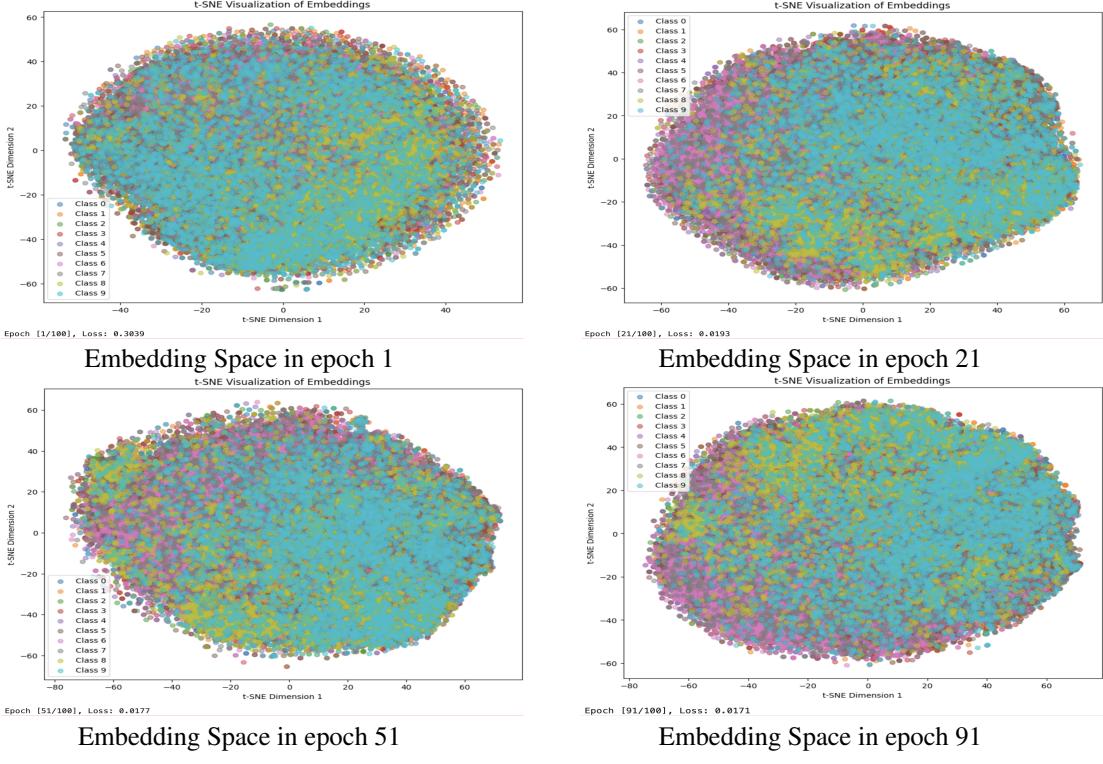


Figure 4: Closest 20 embeddings during 100 epochs of training.

### 104 4.3 Deep Clustering

105 Our deep clustering experiment findings are compiled in Table 1. We found that the performance  
 106 of our model was not up to par after training it on the CIFAR-10 dataset. Despite not being overly  
 107 complex, CIFAR-10 nevertheless needs a more advanced model architecture in order to achieve  
 108 greater generalisation. As such, we choose to train our model using the MNIST dataset, which  
 109 is a simpler dataset. Our model performed noticeably better on MNIST, according to the results,  
 110 demonstrating its usefulness on simpler datasets.

Method	Accuracy	
	MNIST	CIFAR10
Supervised	0.982	0.7409
Self-supervised	0.9193	0.3536

Table 1: Performance of the

## 111 5 Applications of Self-Supervised Learning

- 112 Computer Vision: Enhances image and video understanding tasks such as object detection,  
 113 image classification, and segmentation [5].
- 114 Medical Imaging: Improves analysis of medical images where labeled data is scarce [6].
- 115 Robotics: Enhances visual perception in autonomous systems [7].

## 116 6 Conclusion

117 In this project we provided a solid review of SimCLR and a scalable clustering approach for the  
 118 unsupervised.

- 119 The SimCLR paper introduces a solution to the challenge of limited labeled data, enabling improved  
120 performance on tasks like image classification, object detection, and image retrieval, through the use  
121 of contrastive learning. In our SimCLR implementation, the use of a batch size of 128 compromised  
122 performance, leading to poor representation learning. A T-SNE analysis over 100 epochs showed  
123 minimal changes in the embedding space with this batch size, indicating the need for longer training  
124 periods.
- 125 The DeepCluster alternates between updating the weights of the convolutional network by predicting  
126 the cluster assignments as pseudo-labels and clustering the features generated by the network using  
127 k-means. This method achieves better performance than the previous state-of-the-art on every  
128 standard transfer learning task.
- 129
- 130 Find here the link of the repository for SimCLR and DeepCluster implementation.

131 **References**

- 132 [1] Xiao Liu, Fanjin Zhang, Zhenyu Hou, Li Mian, Zhaoyu Wang, Jing Zhang, and Jie Tang. Self-  
133 supervised learning: Generative or contrastive. *IEEE Transactions on Knowledge and Data  
134 Engineering*, 35(1):857–876, 2023.
- 135 [2] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for  
136 unsupervised learning of visual features. In *Proceedings of the European conference on computer  
137 vision (ECCV)*, pages 132–149, 2018.
- 138 [3] Andrew Ng. Cs229 lecture notes. cs229 lecture notes, 1(1):1–3. 2000.
- 139 [4] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for  
140 contrastive learning of visual representations. In *International conference on machine learning*,  
141 pages 1597–1607. PMLR, 2020.
- 142 [5] Zibei Wang. Self-supervised learning in computer vision: a review. In *International Conference  
143 on Computer Engineering and Networks*, pages 1112–1121. Springer, 2022.
- 144 [6] Saeed Shurrah and Rehab Duwairi. Self-supervised learning methods and applications in medical  
145 imaging analysis: A survey. *PeerJ Computer Science*, 8:e1045, 2022.
- 146 [7] Yongqiang Huang, Juan Wilches, and Yu Sun. Robot gaining accurate pouring skills through  
147 self-supervised learning and generalization. *Robotics and Autonomous Systems*, 136:103692,  
148 2021.