# A Survey on Self-supervised Learning: Algorithms, Applications, and Future Trends

Jie Gui, *Senior Member, IEEE,* Tuo Chen, Jing Zhang, *Senior Member, IEEE,* Qiong Cao, Zhenan Sun, *Senior Member, IEEE,* Hao Luo, Dacheng Tao, *Fellow, IEEE*

**Abstract**—Deep supervised learning algorithms typically require a large volume of labeled data to achieve satisfactory performance. However, the process of collecting and labeling such data can be expensive and time-consuming. Self-supervised learning (SSL), a subset of unsupervised learning, aims to learn discriminative features from unlabeled data without relying on human-annotated labels. SSL has garnered significant attention recently, leading to the development of numerous related algorithms. However, there is a dearth of comprehensive studies that elucidate the connections and evolution of different SSL variants. This paper presents a review of diverse SSL methods, encompassing algorithmic aspects, application domains, three key trends, and open research questions. Firstly, we provide a detailed introduction to the motivations behind most SSL algorithms and compare their commonalities and differences. Secondly, we explore representative applications of SSL in domains such as image processing, computer vision, and natural language processing. Lastly, we discuss the three primary trends observed in SSL research and highlight the open questions that remain. A curated collection of valuable resources can be accessed at https://github.com/guijiejie/SSL.

**Index Terms**—Self-supervised learning, Contrastive learning, Generative model, Representation learning, Transfer learning

✦

## 1 INTRODUCTION

DEEP supervised learning algorithms have demonstrated impressive performance in various domains, including computer vision (CV) and natural language processing (NLP). To address this, models pre-trained on large-scale datasets like ImageNet [1] are commonly employed as a starting point and subsequently fine-tuned for specific downstream tasks (Table 1). This practice is motivated by two primary reasons. Firstly, the parameters acquired from large-scale datasets offer a favorable initialization, enabling faster convergence of models trained on other tasks [2]. Secondly, a network trained on a large-scale dataset has already learned discriminative features, which can be easily transferred to downstream tasks and mitigate the overfitting issue arising from limited training data in such tasks [3], [4].

Unfortunately, numerous real-world data mining and

- J. Gui is with the School of Cyber Science and Engineering, Southeast University and with Purple Mountain Laboratories, Nanjing 210000, China (e-mail: guijie@seu.edu.cn).

- T. Chen is with the School of Cyber Science and Engineering, Southeast University (e-mail: 230219309@seu.edu.cn).

- Jing Zhang is with the School of Computer Science, The University of Sydney, Camperdown, NSW 2050, Australia (e-mail: jing.zhang1@sydney.edu.au).

- D. Tao is with the College of Computing & Data Science at Nanyang Technological University, #32 Block N4 #02a-014, 50 Nanyang Avenue, Singapore 639798 (email: dacheng.tao@gmail.com).

- Q. Cao is with JD Explore Academy (e-mail: mathqiong2012@gmail.com).

- Z. Sun is with the Center for Research on Intelligent Perception and Computing, Chinese Academy of Sciences, Beijing 100190, China (e-mail: znsun@nlpr.ia.ac.cn).

- H. Luo is with Alibaba Group, Hangzhou 310052, China (e-mail: michuan.lh@alibaba-inc.com).

machine learning applications face a common challenge where an abundance of unlabeled training instances coexists with a limited number of labeled ones. The acquisition of labeled examples is frequently costly, arduous, or time-consuming due to the requirement of skilled human annotators with sufficient domain expertise [12], [13]. To illustrate, consider the analysis of web user profiles, where a substantial amount of data can be readily collected. However, the labeling of non-profitable or profitable users necessitates thorough scrutiny, judgment, and sometimes even time-intensive tracing tasks performed by experienced human assessors, resulting in significant expenses. Another instance pertains to the medical field, where unlabeled examples can be easily obtained through routine medical examinations. Nevertheless, assigning diagnoses individually to such a large number of cases places a substantial burden on medical experts. For example, in the case of breast cancer diagnosis, radiologists must label each focus in a vast collection of easily attainable, high-resolution mammograms. This process often proves to be highly inefficient and time-consuming. Additionally, supervised learning methods are susceptible to spurious correlations and generalization errors, and vulnerable to adversarial attacks.

To address the aforementioned limitations of supervised learning, various machine learning paradigms have been introduced, including active learning, semi-supervised learning, and self-supervised learning (SSL). This paper specifically emphasizes SSL. SSL algorithms aim to learn discriminative features from vast quantities of unlabeled instances without relying on human annotations. The general pipeline of SSL is depicted in Fig. 1. In the self-supervised pre-training phase, a pre-defined pretext task is formulated for the deep learning algorithm to solve. Pseudo-labels for the pretext task are automatically generated based on specific attributes of the input data. Once the self-supervised pre-

TABLE 1: Comparison between supervised and self-supervised pre-training and fine-tuning.

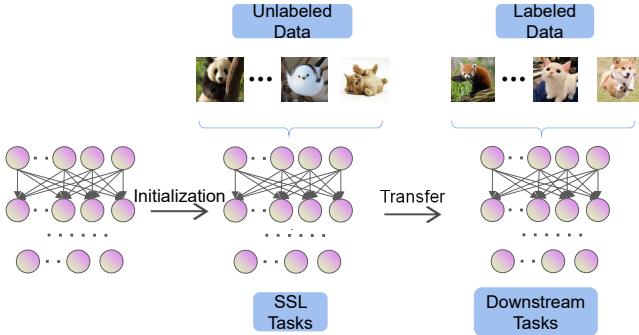| Pre-training | Data | Pre-training Tasks | Downstream Tasks |
|---|---|---|---|
| Supervised | extensive labeled data | image categorization [5] | detection / segmentation / pose estimation / depth estimation, etc. |
| | | video action categorization [6] | action recognition / object tracking, etc. |
| SSL | extensive unlabeled data | Image: rotation [7], jigsaw [8], etc. | detection / segmentation / pose estimation / depth estimation, etc. |
| | | Video: the order of frames [9], playing direction [10], etc. | action recognition / object tracking, etc. |
| | | NLP: masked language modeling [11] | question answering / textual entailment recognition / natural language inference, etc. |



Fig. 1: The general pipeline of applying SSL methods to downstream tasks. The SSL models are first pre-trained on the unlabeled data and then fine-tuned, or directly evaluated, on the labeled data of the downstream tasks.
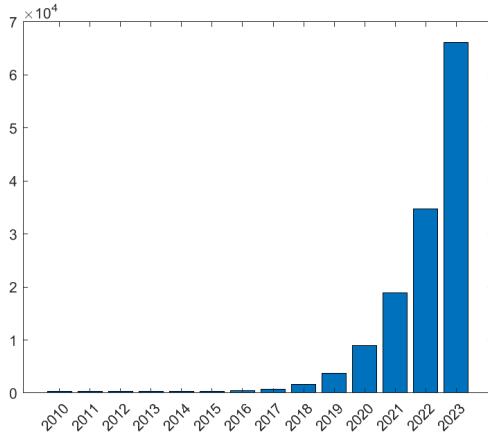


Fig. 2: Google Scholar search results for "self-supervised learning". The vertical and horizontal axes denote the number of SSL publications and the year, respectively.

training process is completed, the acquired model can be transferred to downstream tasks.

One notable advantage of SSL algorithms is their ability to leverage extensive unlabeled data since the generation of pseudo-labels does not necessitate human annotations. By utilizing these pseudo-labels during training, self-supervised algorithms have demonstrated promising outcomes, resulting in a reduced performance disparity compared to supervised algorithms in downstream tasks. Asano et al. [14] demonstrated that SSL can produce generalizable features that exhibit robust generalization even when applied to a single image.

The advancement of SSL [3], [4], [15]–[24] has exhib-

ited rapid progress, capturing significant attention within the research community (Fig. 2), and is recognized as a crucial element for achieving human-level intelligence [25]. Google Scholar reports a substantial volume of SSL-related publications, with approximately 18,900 papers published in 2021 alone. This accounts for an average of 52 papers per day or more than two papers per hour (Fig. 2). To assist researchers in navigating this vast number of SSL papers and to consolidate the latest research findings, we aim to provide a timely and comprehensive survey on this subject.

**Differences from previous work**: Previous works have provided reviews on SSL that cater to specific applications such as recommender systems [26], graphs [27], sequential transfer learning [28], videos [29], and adversarial pre-training of self-supervised deep networks [30]. Besides, Liu et al. [4] primarily focused on papers published before 2020, lacking the latest advancements. Jaiswal et al. [31] centered their survey on contrastive learning (CL). Notably, recent breakthroughs in SSL research within the CV domain are of significant importance. Thus, this review predominantly encompasses recent SSL research derived from the CV community, particularly those influential and classic findings. The primary objectives of this review are to elucidate the concept of SSL, its categories and subcategories, its differentiation and relationship with other machine learning paradigms, as well as its theoretical foundations. We present an extensive and up-to-date review of the frontiers of visual SSL, dividing it into four key areas: context-based, CL, generative, and contrastive generative algorithms, aiming to outline prominent research trends for scholars.

## 2 ALGORITHMS

This section begins by providing an introduction to SSL, followed by an explanation of the pretext tasks associated with SSL and their integration with other learning paradigms.

### 2.1 What is SSL?

The introduction of SSL is attributed to [32] (Fig. 3), who employed this architecture to learn in natural environments featuring diverse modalities. Although the cow image may not warrant a cow label, it is frequently associated with a "moo" sound. The crux lies in the co-occurrence relationship between them.

Subsequently, the machine learning community has advanced the concept of SSL, which falls within the realm of unsupervised learning. SSL involves generating output labels "intrinsically" from input data examples by revealing the relationships between data components or various views of the data. These output labels are derived directly from the data examples. According to this definition, an autoencoder
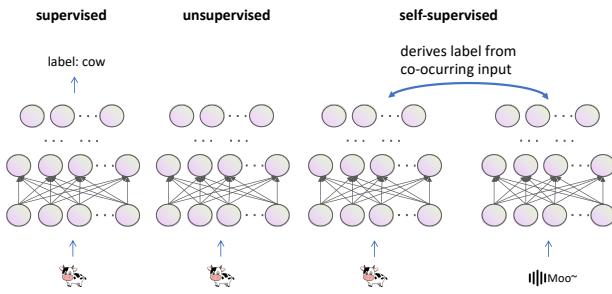
Fig. 3: The differences among supervised learning, unsupervised learning, and SSL. The image is reproduced from [32]. SSL utilizes freely derived labels as supervision instead of manually annotated labels.

(AE) can be perceived as a type of SSL algorithms, where the output labels correspond to the data itself. AEs have gained extensive usage across multiple domains, including dimensionality reduction and anomaly detection.

In the keynote talk at ICLR 2020 [33], Yann LeCun elucidated the concept of SSL as an analogous process to completing missing information (reconstruction). He presented multiple variations as follows: 1) Predict any part of the input from any other part; 2) Predict the future from the past; 3) Predict the invisible from the visible; and 4) Predict any occluded, masked, or corrupted part from all available parts. In summary, a portion of the input is unknown in SSL, and the objective is to predict that particular segment.

Jing et al. [34] expanded the definition of SSL to encompass methods that operate without human-annotated labels. Consequently, any approach devoid of such labels can be categorized under SSL, effectively equating SSL with unsupervised learning. This categorization includes generative adversarial networks (GANs) [35], thereby positioning them within the realm of SSL.

Pretext tasks, also referred to as surrogate or proxy tasks, are a fundamental concept in the field of SSL. The term "pretext" denotes that the task being solved is not the primary objective but serves as a means to generate a robust pre-trained model. Prominent examples of pretext tasks include rotation prediction and instance discrimination, among others. Each pretext task necessitates the use of distinct loss functions to achieve its intended goal. Given the significance of pretext tasks in SSL, we proceed to introduce them in further detail.

## 2.2 Pretext tasks

This section provides a comprehensive overview of the pretext tasks employed in SSL. A prevalent approach in SSL involves devising pretext tasks for networks to solve, where the networks are trained by optimizing the objective functions associated with these tasks. Pretext tasks typically exhibit two key characteristics. Firstly, deep learning methods are employed to learn features that facilitate the resolution of pretext tasks. Secondly, supervised signals are derived from the data itself, a process known as self-supervision. Commonly employed techniques encompass



Fig. 4: Illustration of three common context-based methods: rotation, jigsaw, and colorization.

four categories of pretext tasks: context-based methods, CL, generative algorithms, and contrastive generative methods. In our paper, generative algorithms primarily refer to masked image modeling (MIM) methods.

### 2.2.1 Context-based methods

Context-based methods rely on the inherent contextual relationships among the provided examples, encompassing aspects such as spatial structures and the preservation of both local and global consistency. We illustrate the concept of context-based pretext tasks using rotation as a simple example [36]. Subsequently, we progressively introduce additional tasks (Fig. 4).

**Rotation**: Gidaris et al. [7] trained deep neural networks (DNNs) to learn image representations by recognizing the random geometric transformations. They streamlined image augmentation by introducing rotations of $0°$, $90°$, $180°$, and $270°$ to generate three additional images from each original. This method employs rotation angles as self-supervised labels, using a set of $K = 4$ geometric transformations $G = \{g(\cdot|y)\}_{y=1}^K$. Here, $g(\cdot|y)$ applies a geometric transformation labeled $y$ to an image $X$, resulting in a transformed image $X^y = g(X|y)$.

Gidaris et al. utilized a deep convolutional neural network (CNN), $\mathcal{F}(\cdot)$, to perform rotation prediction through a four-class categorization task. This CNN processes an input image $X^{y^*}$, with $y^*$ being unknown to $\mathcal{F}(\cdot)$, and outputs a probability distribution over possible geometric transformations, expressed as

$$\mathcal{F}\left(X^{y^*}|\theta\right) = \left\{\mathcal{F}^y\left(X^{y^*}|\theta\right)\right\}_{y=1}^K. \tag{1}$$

Here, $\mathcal{F}^y\left(X^{y^*}|\theta\right)$ represents the predicted probability for the geometric transformation labeled as $y$, while $\theta$ denotes the learnable parameters of $\mathcal{F}(\cdot)$.

Given training instances $D = \{X_i\}_{i=1}^N$, the training objective can be formulated as

$$\min_\theta \frac{1}{N} \sum_{i=1}^N \mathcal{L}(X_i, \theta). \tag{2}$$

Here, the loss function is defined as

$$\mathcal{L}(X_i, \theta) = -\frac{1}{K} \sum_{y=1}^K \log(\mathcal{F}^y\left(g\left(X_i|y\right)|\theta\right)). \tag{3}$$

In [37], the relative rotation angle was confined to the interval of $[-30^o, 30^o]$. These rotations were discretized into bins of $3^o$ each, leading to a total of 20 classes (or bins).

**Colorization**: The concept of colorization was initially introduced in [38], and subsequent studies [39]–[41] demonstrated its effectiveness as a pretext task for SSL. Color prediction offers the advantageous feature of requiring freely available training data. In this context, a model can utilize the lightness channel of any color image as input and utilize the corresponding $ab$ color channels in the CIE $Lab$ color space as self-supervised signals. The objective is to predict the $ab$ color channels $Y \in R^{H \times W \times 2}$ given an input lightness channel $X \in R^{H \times W \times 1}$. A commonly employed learning objective is

$$\mathcal{L} = \left\| \hat{Y} - Y \right\|_F^2, \tag{4}$$

where $Y$ and $\hat{Y}$ denote the ground truth and predicted values, respectively.

Besides, [38] utilized the multinomial cross-entropy loss instead of (4) to enhance robustness. Upon completing the training process, the $ab$ color channels would be predicted for any grayscale image. Consequently, the lightness channel and the $ab$ color channels can be concatenated to restore the original grayscale image to a colorful representation.

**Jigsaw**: The jigsaw approach leverages jigsaw puzzles as surrogate tasks, operating under the assumption that a model accomplishes these tasks by comprehending the contextual information embedded within the examples. Specifically, images are fragmented into discrete patches, and their positions are randomly rearranged, with the objective of reconstructing the original order. In [42], the impact of scaling two self-supervised methods, namely jigsaw [8], [43] and colorization, was investigated along three dimensions: data size, model capacity, and problem complexity. The results indicated that transfer performance exhibits a log-linear growth pattern in relation to data size. Furthermore, representation quality was found to improve with higher-capacity models and increased problem complexity.

**Others**: The pretext task employed in [44], [45] involved a conditional motion propagation problem. To enforce a specific constraint on the feature representation process, Noroozi et al. [46] introduced an additional requirement where the sum of feature representations of all image patches should approximate the feature representation of the entire image. While many pretext tasks yield representations that exhibit covariance with image transformations, [47] argued for the importance of semantic representations being invariant to such transformations. In response, they proposed a pretext-invariant representation learning approach that enables the learning of invariant representations through pretext tasks.

### 2.2.2 Contrastive Learning

Numerous SSL methods based on CL have emerged, building upon the foundation of simple instance discrimination tasks [48], [49]. Notable examples include MoCo v1 [50], MoCo v2 [51], SimCLR v1 [52] and SimCLR v2 [53]. Pioneering algorithms, such as MoCo, have significantly enhanced the performance of self-supervised pre-training, reaching a level comparable to that of supervised learning, thus rendering SSL highly pertinent for large-scale applications. Early CL approaches were built upon the concept of utilizing negative examples. However, as CL has progressed, a range of methods have emerged that eliminate the need for negative examples. These methods embrace distinct ideas such as self-distillation and feature decorrelation, yet all adhere to the principle of maintaining positive example consistency. The following section outlines the various CL methods currently available (Fig. 5).

2.2.2.1 Negative example-based CL: Negative examples-based CL adheres to a pretext task known as instance discrimination, which involves generating distinct views of an instance. In negative examples-based CL, views originating from the same instance are treated as positive examples for an anchor sample, while views from different instances serve as negative examples. The underlying principle is to promote proximity between positive examples and maximize the separation between negative examples within the latent space. The definition of positive and negative examples varies depending on factors such as the modality being considered and specific requirements, including spatial and temporal consistency in video understanding or the co-occurrence of modalities in multi-modal learning scenarios. In the context of conventional 2D image CL, image augmentation techniques are utilized to generate diverse views from a single image.

**MoCo**: He et al. [50] framed CL as a dictionary look-up task. In this framework, a query $q$ exists and a set of encoded examples $\{k_0, k_1, k_2, \cdots\}$ serve as the keys in a dictionary. Assuming a single key, denoted as $k_+$ in the dictionary, matches the query $q$, a contrastive loss [57] function is employed. The value of this function is low when $q$ is similar to its positive key $k_+$ and dissimilar to all other negative keys. In the MoCo v1 [50] framework, the InfoNCE loss function [58], a form of contrastive loss, is utilized, *i.e.*,

$$\mathcal{L}_q = -\log \frac{\exp(q \cdot k_+/\tau)}{\sum_{i=0}^{K} \exp(q \cdot k_i/\tau)}, \tag{5}$$

where $\tau$ represents the temperature hyper-parameter and $(\cdot)$ denotes vector product. The summation is computed over one positive example and $K$ negative examples. InfoNCE is derived from noise contrastive estimation (NCE) [59].

MoCo v2 [51] builds upon MoCo v1 [50] and SimCLR v1 [52], incorporating a multilayer perceptron (MLP) projection head and more data augmentations.

**SimCLR**: SimCLR v1 [52] employs a mini-batch sampling strategy with $N$ instances, wherein a contrastive prediction task is formulated on pairs of augmented instances from the mini-batch, generating a total of $2N$ instances. Notably, SimCLR v1 does not explicitly select negative instances. Instead, for a given positive pair, the remaining $2(N-1)$ augmented instances in the mini-batch are treated as negatives. Let $sim(u, v) = u^T v/(\|u\| \|v\|)$ represent the cosine similarity between two instances $u$ and $v$. The loss function of SimCLR v1 for a positive instance pair $(i, j)$ is defined as

$$\mathcal{L}_{i,j} = -\log \frac{\exp(sim(z_i, z_j)/\tau)}{\sum_{k=1}^{2N} 1_{[k \neq i]} \exp(sim(z_i, z_k)/\tau)}, \tag{6}$$

where $1_{[k \neq i]} \in \{0, 1\}$ is an indicator function equal to 1 if $k \neq i$, and $\tau$ denotes the temperature hyper-parameter. The overall loss is computed across all positive pairs, including both $(i, j)$ and $(j, i)$, within the mini-batch.
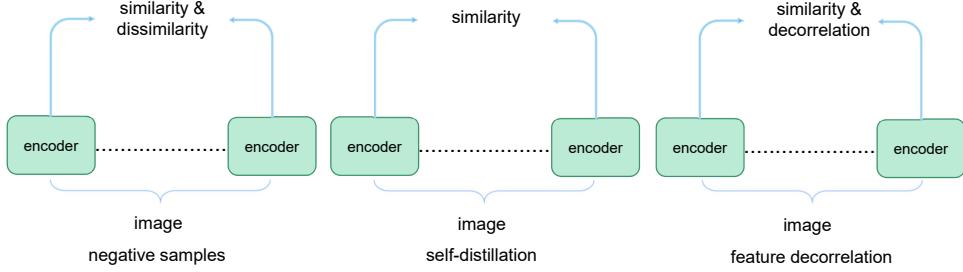
Fig. 5: Illustration of different CL methods: CL based on negative examples (left), CL based on self-distillation (middle), and CL based on feature decorrelation (right). For a demonstration of the concepts of similarity and dissimilarity, one can refer to [52], [54], while for insights into decorrelation, [55], [56] provide a comprehensive overview.

In MoCo, the features generated by the momentum encoder are stored in a feature queue as negative examples. These negative examples do not undergo gradient updates during backpropagation. Conversely, SimCLR utilizes negative examples from the current mini-batch, and all of them are subjected to gradient updates during backpropagation. Both MoCo and SimCLR rely on data augmentation techniques, including cropping, resizing, and color distortion. Notably, SimCLR made a significant contribution by highlighting the crucial role of strong data augmentation in CL, a finding subsequently confirmed by MoCo v2. Additional augmentation methods have also been explored [60]. For instance, in [61], foreground saliency levels were estimated in images, and augmentations were created by selectively copying and pasting image foregrounds onto diverse backgrounds, such as grayscale images with random grayscale levels, texture images, and ImageNet images. Furthermore, views can be derived from various sources, including different modalities such as photos and sounds [62], as well as coherence among different image channels [63].

Minimizing the contrastive loss is known to effectively maximize a lower bound of the mutual information $I(\mathbf{x}_1; \mathbf{x}_2)$ between the variables $\mathbf{x}_1$ and $\mathbf{x}_2$ [58]. Building upon this understanding, [64] proposes principles for designing diverse views based on information theory. These principles suggest that the views should aim to maximize $I(\mathbf{v}_1; \mathbf{y})$ and $I(\mathbf{v}_2; \mathbf{y})$ ($\mathbf{v}_1$, $\mathbf{v}_2$, and $\mathbf{y}$ denoting the first view, the second view, and the label, respectively), representing the amount of information contained about the task label, while simultaneously minimizing $I(\mathbf{v}_1; \mathbf{v}_2)$, indicating the shared information between inputs encompassing both task-relevant and irrelevant details. Consequently, the optimal data augmentation method is contingent on the specific downstream task. In the context of dense prediction tasks, [65] introduces a novel approach for generating different views. This study reveals that commonly employed data augmentation methods, as utilized in SimCLR, are more suitable for categorization tasks rather than dense prediction tasks such as object detection and semantic segmentation. Consequently, the design of data augmentation methods tailored to specific downstream tasks has emerged as a significant area of exploration.

Given the observed benefits of strong data augmentation in enhancing CL performance [52], there has been a growing interest in leveraging more robust augmentation techniques. However, it is worth noting that solely relying on strong data augmentation can actually lead to a decline in performance [64]. The distortions introduced by strong data augmentation can alter the image structure, resulting in a distribution that differs from that of weakly augmented images. This discrepancy poses optimization challenges. To address the overfitting issue arising from strong data augmentation, [66] proposes an alternative approach. Instead of employing a one-hot distribution, they suggest using the distribution generated by weak data augmentation as a mimic. This mitigates the negative impact of strong data augmentation by aligning the distribution of augmented examples with that of weakly augmented examples.

2.2.2.2 Self-distillation-based CL: Bootstrap Your Own Latent (BYOL) [67] is a prominent self-distillation algorithm designed specifically for self-supervised image representation learning, eliminating the need for negative pairs. BYOL employs two identical DNNs, known as Siamese networks, with the same architecture but different weights. One serves as the online network, while the other is the target network. Similar to MoCo [50], BYOL enhances the target network through a gradual averaging of the online network. Siamese networks have emerged as prevalent architectures in contemporary self-supervised visual representation learning models, including SimCLR, BYOL, and SwAV [68]. These models aim to maximize the similarity between two augmented versions of a single image while incorporating specific conditions to mitigate the risk of collapsing solutions.

Simple Siamese (SimSiam) networks, introduced by [69], offers a straightforward approach to learning effective representations in SSL without the need for negative example pairs, large batches, or momentum encoders. Given a data point $x$ and two randomly augmented views $x_1$ and $x_2$, an encoder $f$ and an MLP prediction head $h$ process these views. The resulting outputs are denoted as $p_1 = h(f(x_1))$ and $z_2 = f(x_2)$. The objective of [69] is to minimize their negative cosine similarity:

$$D(p_1, z_2) = -\frac{p_1}{\|p_1\|_2} \frac{z_2}{\|z_2\|_2}. \tag{7}$$

Here, $\|\|_2$ represents the $l_2$-norm. Similar to [67], a symmetric loss [69] is defined as

$$\mathcal{L} = \frac{1}{2}(D(p_1, z_2) + D(p_2, z_1)). \tag{8}$$

This loss is defined based on the example $x$, and the overall loss is the average of all examples. Notably, [69] employs
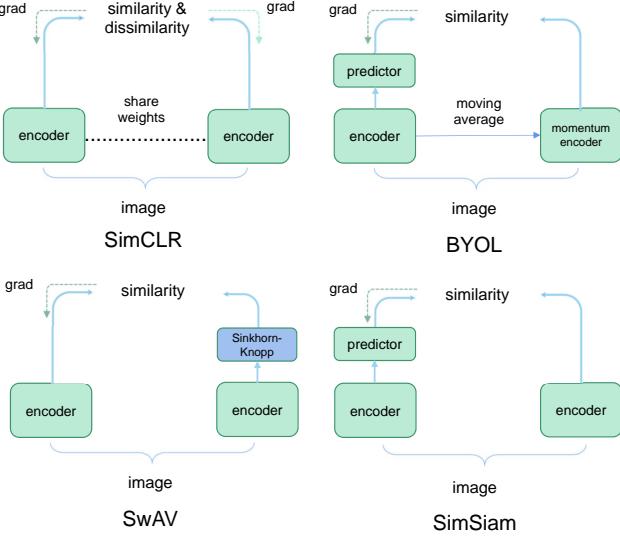
Fig. 6: Comparison among different Siamese architectures. The image is reproduced from [69].

a stop-gradient (*stopgrad*) operation by modifying Eq. (7) as $D\left(p_1, stopgrad\left(z_2\right)\right)$. This implies that $z_2$ is treated as a constant. Similarly, Eq. (8) is revised as

$$\mathcal{L} = \frac{1}{2}\left(D\left(p_1, stopgrad\left(z_2\right)\right) + D\left(p_2, stopgrad\left(z_1\right)\right)\right). \quad (9)$$

Figure 6 illustrates the distinctions among SimCLR, BYOL, SwAV, and SimSiam. The categorization of BYOL and SimSiam as CL methods is a subject of debate due to their exclusion of negative examples. However, to be consistent with [70], this paper considers BYOL and SimSiam to belong to CL methods.

  2.2.2.3  Feature decorrelation-based CL: The objective of feature decorrelation is to learn decorrelated features.

**Barlow Twins**:

Barlow Twins [55] introduced a novel loss function that encourages the similarity of embedding vectors from distorted versions of an example while minimizing redundancy between their components. Similar to other SSL methods such as MoCo [50] and SimCLR [52], Barlow Twins generates two distorted views $Y^A$ and $Y^B$ via a distribution of data augmentations $\mathcal{T}$ for each image in a data batch sampled from a dataset, resulting in batches of embeddings $Z^A$ and $Z^B$. The loss function of Barlow Twins is defined as

$$\mathcal{L}_{BT} = \sum_i \left(1 - C_{ii}\right)^2 + \lambda \sum_i \sum_{j \neq i} C_{ij}^2. \quad (10)$$

Here, $\lambda$ is a hyper-parameter, and $C$ represents the cross-correlation matrix computed between the two batches of embeddings $Z^A$ and $Z^B$, defined as

$$C_{ij} = \frac{\sum_b z_{b,i}^A z_{b,j}^B}{\sqrt{\sum_b \left(z_{b,i}^A\right)^2}\sqrt{\sum_b \left(z_{b,j}^B\right)^2}}, \quad (11)$$

where $b$ indexes batch samples and $i, j$ index the vector dimension of the networks' outputs. $C$ is a square matrix that measures the correlation between the two batches of

embeddings $Z^A$ and $Z^B$. The first term in Eq. (10) encourages the diagonal elements of $C$ to be close to 1, while the second term encourages the off-diagonal elements to be close to 0.

**Variance-Invariance-Covariance Regularization**: Borrowing the covariance regularization from the Barlow Twins method, Variance-invariance-covariance regularization (VICReg) [56] proposes a new self-supervised method for training joint embedding architectures that simultaneously considers variance, invariance, and covariance. Similar to Barlow Twins, VICReg generates two distorted views $Y^A$ and $Y^B$ via a distribution of the data augmentation $\mathcal{T}$ and gets their embeddings $Z^A \in \mathbb{R}^{n \times d}$ and $Z^B \in \mathbb{R}^{n \times d}$. Let the subscript $j$ index the embedding in the batch and $d$, $n$ represent the dimensionality of the vectors in $Z^A$ and the batch size, respectively. The main contribution of VICReg is the variance preservation term, which explicitly prevents a collapse due to a shrinkage of the embedding vectors toward zero. The variance regularization term $v$ in VICReg is defined as a hinge loss function applied to the standard deviation of the embeddings along the batch dimension:

$$v\left(Z^A\right) = \frac{1}{d}\sum_{j=1}^d \max(0, \gamma - S\left(z_j^A, \varepsilon\right)). \quad (12)$$

Here, $z_j^A$ represents the vector composed of each value at dimension $j$ in $Z^A$ and $S$ represents the regularized standard deviation, defined as

$$S(y, \varepsilon) = \sqrt{\text{Var}(y) + \varepsilon}. \quad (13)$$

The constant $\gamma$ determines the standard deviation and is set to 1 in the experiments, while $\varepsilon$ is a small scalar used to prevent numerical instabilities. This criterion encourages the variance within the current batch to be equal to or greater than $\gamma$ for every dimension, thereby preventing collapse scenarios where all data are mapped to the same vector.

The invariance criterion $s$ in VICReg, which captures the similarity between $Z^A$ and $Z^B$, is defined as the mean-squared Euclidean distance between each pair of data without any normalization:

$$s\left(Z^A, Z^B\right) = \frac{1}{n}\sum_{b=1}^n \left\|z_b^A - z_b^B\right\|_2^2. \quad (14)$$

In addition, the covariance criterion $c(Z)$ in VICReg is defined as

$$c\left(Z\right) = \frac{1}{d}\sum_{i \neq j}\left[C(Z)\right]_{i,j}^2, \quad (15)$$

where $C(Z)$ represents the covariance matrix of $Z$. The overall loss of VICReg is a weighted sum of the variance, invariance, and covariance:

$$\begin{aligned}\mathcal{L} = &\, s\left(Z^A, Z^B\right) + \alpha\left(v\left(Z^A\right) + v\left(Z^B\right)\right) \\ &+ \beta\left(C\left(Z^A\right) + C\left(Z^B\right)\right),\end{aligned} \quad (16)$$

where $\alpha$ and $\beta$ are two hyper-parameters. Note that both regularization terms — the variance regularization term and the covariance regularization term — are applied independently to each branch of the architecture. This differs from the Barlow Twins, which uses a cross-correlation matrix between the two branches of the Siamese architecture.
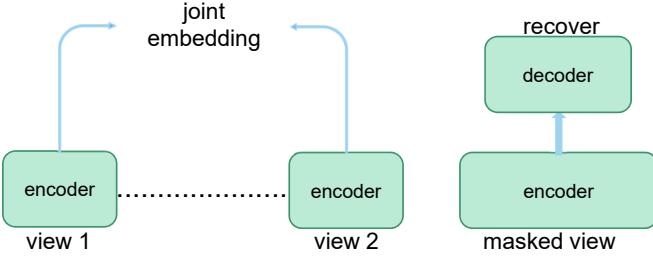
Fig. 7: The broad differences between CL and MIM. Note that the actual differences between their pipelines are not limited to what is shown.

*2.2.2.4  Analysis of CL:* Despite the impressive results achieved by contrastive SSL, the underlying mechanisms remain obscure and not fully understood. Several studies have delved into this area [54], [71]–[81]. Theoretical investigations by [72], [76], [79] have provided support for the value of feature representations generated through CL. In the Appendix, we also provide explanations of the connections between contrastive learning and other concepts, such as Principal Component Analysis, Spectral Clustering, and Supervised Learning.

*2.2.2.5  Others:* Besides the aforementioned works, several other approaches have employed CL. Among them, [82], [83] investigated the utilization of vision transformers (ViTs) as the backbone for contrastive SSL, employing multi-crop and cross-entropy loss [83]. Notably, [83] discovered that the resultant features exhibited exceptional performance as $K$-nearest neighbors ($K$-NN) classifiers and effectively encoded explicit information regarding the semantic segmentation of images. These desirable properties have also motivated specific downstream tasks [84].

In a different study, [85] adopted patches extracted from the same image as a positive pair, while patches from different images served as negative pairs. A mixing operation is further explored in RegionCL [86] to diversify the contrastive pairs. Yang et al. [87] integrated CL and MIM in the context of text recognition, utilizing a weighted objective function.

Numerous CL-based methods are available in the literature [88]–[96]. It should be noted that CL is not restricted solely to SSL, as it can also be used in supervised learning [97].

### 2.2.3  Generative algorithms

For the category of generative algorithms, this study primarily focuses on MIM methods. MIM methods [98] (Fig. 7)—namely, bidirectional encoder representation from image transformers (BEiT) [99], masked AE (MAE) [70], context AE (CAE) [100], and a simple framework for MIM (SimMIM) [101]—have gained significant popularity and pose a considerable challenge to the prevailing dominance of CL. MIM leverages co-occurrence relationships among image patches as supervision signals.

MIM is a variant of the denoising AE (DAE) [16]. Notably, the Bidirectional Encoder Representations from Transformers (BERT) [11] and Generative Pre-trained Transformer (GPT) [102] have emerged as a renowned variant of the DAE

and achieved remarkable success in NLP. Researchers aspire to extend this success to CV by employing BERT-like pre-training strategies. However, it is crucial to acknowledge that BERT's success in NLP can be attributed not only to its large-scale self-supervised pre-training but also to its scalable network architecture. A notable distinction between the NLP and CV communities is their use of different primary models, with transformers being prevalent in NLP and CNNs being widely adopted in CV.

The landscape changed significantly with the introduction of the original ViT [5], which marked a pivotal moment. Alexey Dosovitskiy et al. conducted pioneering research on applying MIM to CV, drawing inspiration from BERT's masked image prediction paradigm. Their smaller ViT-B/16 model achieved 79.9% accuracy on ImageNet [1] through self-supervised pre-training, an impressive 2% improvement over training from scratch. However, it still fell short of the accuracy attained by supervised pre-training. iGPT [103] further employs the GPT-style next token prediction, but it received limited attention due to its subpar accuracy and computational efficiency. Beyond ViTs, a separate early investigation adopted context encoders [104], employing a concept akin to MAE, *i.e.*, image inpainting.

However, the differences between natural language and visual signals limit the effectiveness of naive paradigm transfer. BEiT introduces a tailored MIM task for visual pre-training, *i.e.*, an extra tokenization procedure which breaks down the input image into visual tokens, and then predicts randomly masked subset of the image tokens. To address the challenge of tokenization, the authors leveraged a discrete variational autoencoder (dVAE) [105] to create a predefined visual vocabulary. In contrast to BEiT, MAE does not utilize image tokens; instead, it approaches the problem from the perspective of image signal sparsity. MAE identifies a significant amount of redundancy in image signals, necessitating a higher masking rate, such as 75%.

Here, we define

$$\text{MIM} := \mathcal{L}\left(\mathcal{D}\left(\mathcal{E}\left(\mathcal{T}_1\left(I\right)\right)\right), \mathcal{T}_2\left(I\right)\right), \qquad (17)$$

where $\mathcal{E}$ denotes the encoder, $\mathcal{D}$ denotes the decoder, $\mathcal{T}_1$ represents the transformation applied to the input before it is fed into the network, and $\mathcal{T}_2$ represents the transformation used to derive the target label. It is noteworthy that this representation is provided for the sake of clarity and ease of understanding rather than serving as a strict definition.

The primary distinction between BEiT and MAE lies in their choice of $\mathcal{T}$. While BEiT employs the token output from the pre-trained tokenizer as its target, MAE directly uses the original pixels as its target. BEiT adopts a two-stage approach, initially training a tokenizer to convert images into visual tokens, followed by BERT-style training. On the other hand, MAE is a one-stage end-to-end approach, incorporating a decoder to decode the encoder-derived representation into the original pixels. The two representative MIM approaches BEiT and MAE, showcase different architectural designs, with subsequent MIM methods often following one of these techniques. A central challenge in MIM lies in the selection of the target representation $\mathcal{T}_2$, which leads to the categorization of MIM methods, as presented in Table 2.

Following the introduction of BEiT and MAE, several variants have been proposed. iBOT [98] is an "online tok-

TABLE 2: Categorization of MIM methods based on the reconstruction target. The second and third rows denote MIM methods and reconstructing targets, respectively.

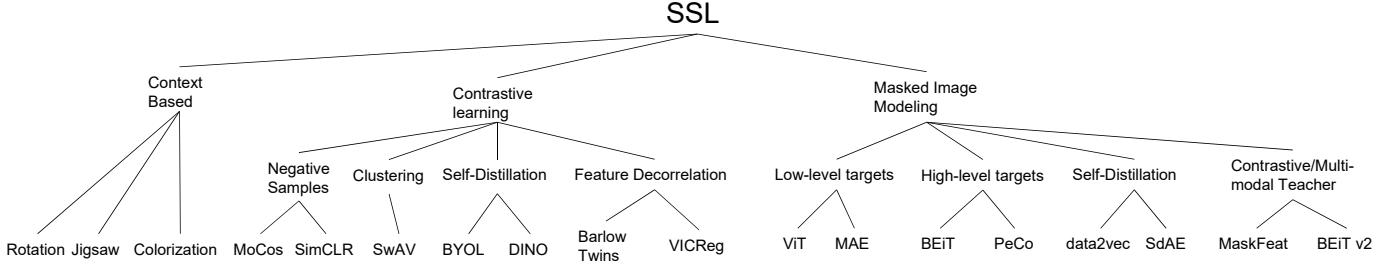| | Low-Level Targets | | | | High-Level Targets | | | Self-Distillation | | Contrastive / Multi-modal Teacher | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Algorithm | ViT [5] | MAE [70] | SimMIM [101] | Maskfeat [106] | BEiT [99] | CAE [100] | PeCo [107] | data2vec [108] | SdAE [109] | MimCo [110] | BEiT v2 [111] |
| Target | Raw Pixel | | | HOG | VQ-VAE | | VQ-GAN | self | | MoCo v3 | CLIP |



Fig. 8: Several representative pretext tasks of SSL.

enizer" adaptation of BEiT, aiming to address the limitation of dVAE in capturing only low-level semantics within local details. The CAE introduces an alignment constraint to encourage masked patch representations (predicted by a "latent contextual regressor") to lie in the encoded representation space. This decoupling of the representation learning task and pretext task enhances the model's capacity for representation learning. Furthermore, MAE has been extended to other modalities beyond images [112]–[114].

Generative pre-training has also evolved in the video domain. BEVT [115] decouples video representation learning into spatial representation learning and temporal dynamics learning. It first undertakes masked image modeling on image data, followed by a joint approach of masked image modeling and masked video modeling on video data. This accelerates training and achieves results comparable to those of strongly-supervised baselines. Similarly, VideoMAE [116] extends the MAE to videos and discovers that an extremely high proportion of masking ratio (90% to 95%) is permissible in video mask modeling. Moreover, it remains effective even on very small datasets, consisting of only 3,000 to 4,000 videos. OmniMAE [117] demonstrates that a unified model can be concurrently trained across multiple visual modalities, breaking the paradigm of previously studying different modes in isolation. This significantly streamlines the training process, enabling more efficient development of large-scale model architectures. SiamMAE [118] indicates that, contrary to images that are (approximately) isotropic, the temporal dimension is unique, necessitating an asymmetric approach to processing temporal and spatial information, as not all spatiotemporal orientations are equally probable.

MIM has demonstrated significant potential in pre-training vision transformers [119]–[121]. However, in prior works, the random masking of image patches led to an underutilization of valuable semantic information essential for effective visual representation learning. Liu et al. [122] introduced an attention-driven masking strategy to explore improvements over random masking for insufficient semantic utilization.

### 2.2.4 Contrastive Generative Methods

As stated in [123], contrastive models tend to be data-hungry and vulnerable to overfitting issues, whereas generative models encounter data-filling challenges and exhibit inferior data scaling capabilities when compared to contrastive models. While contrastive models often focus on global views [83], overlooking internal structures within images, MIM primarily models local relationships. The divergent characteristics and challenges encountered in contrastive self-supervised learning and generative self-supervised learning have motivated researchers to explore the combination of these two kinds of approaches.

To elaborate further, let us compare the challenges faced by contrastive self-supervised methods and generative self-supervised methods. Generative self-supervised methods are characterized as data-filling approaches [124]. For a model of a certain size, when the dataset reaches a certain magnitude, further scaling of the data does not lead to significant performance gains in generative self-supervised methods. In contrast, recent studies have revealed the potential of data scaling to enhance the performance of CL [125]. As data increases, CL shows substantial performance improvements, demonstrating remarkable generalization without additional fine-tuning on downstream tasks. However, the scenario differs in low-data regimes. Contrastive models may find shortcuts with trivial representations that overfit the limited data [50], thus leading to inconsistent improvements in generalization performance for downstream tasks using pre-trained models with contrastive self-supervised methods [123]. On the other hand, generative methods are more adept at handling low-data scenarios and can even achieve notable performance improvements when data is extremely scarce, such as with only 10 images [126].

Several endeavors have sought to integrate both types of algorithms [123], [127]. In [127], GANs are employed for online data augmentation in CL. The study devises a contrastive module that learns view-invariant features for generation and introduces a view-invariant loss function to facilitate learning between original and generated views. On the other hand, [98] draws inspiration from both BEiT and DINO [83]. It modifies the tokenizer of BEiT to an online distilled teacher while integrating cross-view distillation from

the DINO framework. As a result, iBOT [98] significantly enhances linear probing accuracy compared to the MIM method. RePre [128] integrates local feature learning into self-supervised vision transformers through reconstructive pre-training, an approach that enhances contrastive frameworks. This is achieved by incorporating an additional branch dedicated to reconstructing raw image pixels, which operates concurrently with the established contrastive objective. CMAE [129] concurrently performs CL and MIM tasks. To align CL with MIM effectively, CMAE introduces two novel components: pixel shifting for generating plausible positive views, and a feature decoder for enhancing the features of contrastive pairs. This approach significantly improves the quality of representation and transfer performance compared to its MIM-only counterparts. SiameseIM [130] does not simply merge the objectives of CL and MIM, but rather utilizes the views generated by CL as the target for MIM reconstruction in the latent space.

Despite attempts to combine both types of approaches, naive combinations may not always yield performance gains and can even perform worse than the generative model baseline, thereby exacerbating the issue of representation over-fitting [123]. The performance degradation could be attributed to the disparate properties of CL and generative methods. For instance, CL methods typically exhibit longer attention distances, whereas generative methods tend to favor local attention [131]. In light of this challenge, RECON [123] emerges as a solution by training generative modeling to guide CL, thereby leveraging the benefits of both paradigms.

### 2.2.5 Summary

As described above, numerous pretext tasks for SSL have been devised, with several significant milestone variants depicted in Fig. 8.

Several other pretext tasks are available [132], [133], encompassing diverse approaches such as relative patch location [134], noise prediction [135], feature clustering [136]–[138], cross-channel prediction [139], and combining different cues [140]. Kolesnikov et al. [141] conducted a comprehensive investigation of previously proposed SSL pretext tasks, yielding significant insights. Besides, Krähenbühl et al. [142] proposed an alternative approach to pretext tasks and demonstrated the ease of obtaining data from video games.

It has been observed that context-based approaches exhibit limited applicability due to their inferior performance. In the realm of visual SSL, two dominant types of algorithms are CL and MIM. While visual CL may encounter overfitting issues, CL algorithms that incorporate multi-modality, exemplified by CLIP [2], have gained popularity.

### 2.3 Combinations with other learning paradigms

It is essential to acknowledge that the advancements in SSL did not occur in isolation; instead, they have been the result of continuous development over time. In this section, we provide a comprehensive list of relevant learning paradigms that, when combined with SSL, contribute to a clearer understanding of their collective impact.

#### 2.3.1 GANs

GANs represent classical unsupervised learning methods and were among the most successful approaches in this domain before the surge of SSL techniques. The integration of GANs with SSL offers various avenues, with self-supervised GANs (SS-GAN) serving as one such example. The GANs' objective function [35], [143] is given as

$$\min_G \max_D V(G, D) = E_{x \sim p_{data}(x)} [\log D(x)] \\ + E_{z \sim p_z(z)} [\log(1 - D(G(z)))]. \quad (18)$$

The SS-GAN [144] is defined by combining the objective functions of GANs with the concept of rotation [7]:

$$L_G(G, D) = -V(G, D) \\ - \alpha \mathbb{E}_{x \sim p_G} \mathbb{E}_{r \sim R} [\log Q_D(R = r | x^r)], \quad (19)$$

$$L_D(G, D) = V(G, D) \\ - \beta \mathbb{E}_{x \sim p_{data}} \mathbb{E}_{r \sim R} [\log Q_D(R = r | x^r)], \quad (20)$$

where $V(G, D)$ represents the objective function of GANs as given in Eq. (18), and $r \sim R$ refers to a rotation selected from a set of possible rotations, similar to the concept presented in [7]. Here, $x^r$ denotes an image $x$ rotated by $r$ degrees, and $Q(R | x^r)$ corresponds to the discriminator's predictive distribution over the angles of rotation for a given example $x$. Notably, rotation [7] serves as a classical SSL method. The SS-GAN incorporates rotation invariance into the GANs' generation process by integrating the rotation prediction task during training.

#### 2.3.2 Semi-supervised learning

SSL and semi-supervised learning are contrasting paradigms that can be effectively combined. One notable example of this combination is self-supervised semi-supervised learning (S$^4$L) [145]. In S$^4$L, the objective function is given by

$$\mathcal{L} = \min_\theta \mathcal{L}_l(D_l, \theta) + w\mathcal{L}_u(D_u, \theta). \quad (21)$$

This means optimizing the corresponding loss objectives on a labeled dataset $D_l$ and an unlabeled dataset $D_u$. $\mathcal{L}_l$ is the categorization loss (e.g., cross-entropy) and $\mathcal{L}_u$ stands for the self-supervised loss (e.g., rotation task in Eq. (3)). $\theta$ is the learnable parameters.

Incorporating SSL as an auxiliary task is a well-established approach in semi-supervised learning. Another classical method to leverage SSL within this context involves implementing SSL on unlabeled data, followed by fine-tuning the resultant model on labeled data, as demonstrated in the SimCLR.

To demonstrate the robustness of self-supervision against adversarial perturbations, Hendrycks et al. [146] proposed an overall loss function as a linear combination of supervised and self-supervised losses:

$$\mathcal{L}(x, y, \theta) = \mathcal{L}_{CE}(y, p(y | PGD(x)), \theta) \\ + \lambda \mathcal{L}_{SS}(PGD(x), \theta), \quad (22)$$

where $x$ is the example, $y$ is the one-hot vector of ground-truth and $\theta$ denotes the model parameters. The adversarial example is generated from $x$ by projected gradient descent (PGD) and adversarial training is implemented by cross-entropy loss $\mathcal{L}_{CE}$. $\mathcal{L}_{SS}$ is the self-supervised loss.

### 2.3.3 Multi-instance learning (MIL)

Miech et al. [13] introduced an extension of the InfoNCE loss (5) for MIL and termed it MIL-NCE:

$$\max_{f,g} \sum_{i=1}^{n} \log \left( \frac{\sum\limits_{(x,y)\in P_i} e^{f(x)^T g(y)}}{\sum\limits_{(x,y)\in P_i} e^{f(x)^T g(y)} + \sum\limits_{(x',y')\in N_i} e^{f(x')^T g(y')}} \right), \quad (23)$$

where $x$ and $y$ represent a video clip and a narration, respectively. The functions $f$ and $g$ generate embeddings of $x$ and $y$, respectively. For a specific example indexed by $i$, $P_i$ denotes the set of positive video/narration pairs, while $N_i$ corresponds to the set of negative video/narration pairs.

### 2.3.4 Multi-view/multi-modal(ality) learning

Observation plays a vital role in infants' acquisition of knowledge about the world. Notably, they can grasp the concept of apples through observational and comparative processes, which distinguishes their learning approach from traditional supervised algorithms that rely on extensive labeled apple data. This phenomenon was demonstrated by Orhan et al. [22], who gathered perceptual data from infants and employed an SSL algorithm to model how infants learn the concept of "apple". Moreover, infants' learning about the world extends to multi-view and multi-modal(ality) learning [2], encompassing various sensory inputs such as video and audio. Hence, SSL and multi-view/multi-modal(ality) learning converge naturally in infants' learning mechanisms as they explore and comprehend the workings of the world.

2.3.4.1 Multiview CL: The objective function in standard multiview CL, as proposed by Tian et al. [64], is given by

$$\mathcal{L}_{NCE} = E[L_q], \quad (24)$$

where $L_q$ corresponds to Eq. (5). Multiview CL treats different views of the same sample as positive examples for contrastive learning. Tian et al. [64] introduced both unsupervised and semi-supervised multiview learning based on adversarial learning. Let $\hat{X}$ denote $g(X)$, i.e., $\hat{X} = g(X)$. Two encoders, $f_1$ and $f_2$, were trained to maximize $I_{NCE}(\hat{X}_1, \hat{X}_{2:3})$ as stated in Eq. (24). A flow-based model $g$ was trained to minimize $I_{NCE}(\hat{X}_1, \hat{X}_{2:3})$ and $\{X_1, X_{2:3}\}$ is obtained from image splitting over its channels. Formally, the objective function for unsupervised view learning can be expressed as

$$\min_{g} \max_{f_1, f_2} I_{NCE}^{f_1, f_2}(g(X)_1, g(X)_{2:3}). \quad (25)$$

In the context of semi-supervised view learning, when several labeled examples are available, the objective function is formulated as

$$\min_{g,c_1,c_2} \max_{f_1,f_2} I_{NCE}^{f_1,f_2}(g(X)_1, g(X)_{2:3}) \\ + L_{ce}(c_1(g(X)_1), y) + L_{ce}(c_2(g(X)_{2:3}), y), \quad (26)$$

where $y$ represents the labels, $c_1$ and $c_2$ are classifiers, and $L_{ce}$ denotes the cross-entropy. Further relevant works can be found in [63], [64], [147]. Table 3 summarizes different SSL losses.

2.3.4.2 Images and text: In the study conducted by Gomez et al. [148], the authors employed a topic modeling framework to project the text of an article into the topic probability space. This semantic-level representation was then utilized as the self-supervised signal for training CNN models on images. On a similar note, CLIP [2] leverages a CL-style pre-training task to predict the correspondence between captions and images. Benefiting from the CL paradigm, CLIP is capable of training models from scratch on an extensive dataset comprising 400 million image-text pairs collected from the internet. Consequently, CLIP's advancements have significantly propelled multimodal learning to the forefront of research attention.

2.3.4.3 Point clouds and other modalities: Several SSL methods have been proposed for joint learning of 3D point cloud features and 2D image features by leveraging cross-modality and cross-view correspondences through triplet and cross-entropy losses [149]. Additionally, there are efforts to jointly learn view-invariant and mode-invariant characteristics from diverse modalities, such as images, point clouds, and meshes, using heterogeneous networks for 3D data [150]. SSL has also been employed for point cloud datasets, with approaches including CL and clustering based on graph CNNs [151]. Furthermore, AEs have been used for point clouds in works like [113], [114], [152], [153], while capsule networks have been applied to point cloud data in [154].

### 2.3.5 Test time training

Sun et al. [155] introduced "test time training (TTT) with self-supervision" to enhance the performance of predictive models when the training and test data come from distinct distributions. TTT converts an individual unlabeled test example into an SSL problem, enabling model parameter updates before making predictions. Recently, Gandelsman et al. [156] combined TTT with MAE for improved performance. They argued that by treating TTT as a one-sample learning problem, optimizing a model for each test input could be addressed using the MAE as

$$h_0 = \arg\min_{h} \frac{1}{n} \sum_{i=1}^{n} \mathcal{L}_m(h \circ f_0(x_i), y_i), \quad (27)$$

$$f_x, g_x = \arg\min_{f,g} \mathcal{L}_s(g \circ f(\text{mask}(x)), x). \quad (28)$$

Here, $f$ and $g$ refer to the encoder and decoder of MAE, and $h$ denotes the main task head, respectively.

TTT achieves an improved bias-variance tradeoff under distribution shifts. A static model heavily depends on training data that may not accurately represent the new test distribution, leading to bias. On the other hand, training a new model from scratch for each test input, ignoring all training data, is undesirable. This approach results in an unbiased representation for each test input but exhibits high variance due to its singularity.

### 2.3.6 Summary

The evolution of SSL is characterized by its dynamic and interconnected nature. Analyzing the amalgamation of various methods allows for a clearer grasp of SSL's developmental trajectory. An exemplar of this success is evident

TABLE 3: Different losses of SSL.

| Category | | Method | Loss | Equation |
|---|---|---|---|---|
| | Context-Based | Rotation [7] | Rotation Prediction | (3) |
| | | MoCo v1 [50] | InfoNCE | (5) |
| | | SimCLR v1 [52] | InfoNCE | (6) |
| Pretext | CL | SimSiam [69] | Cosine Similarity | (9) |
| | | Barlow Twins [55] | Invariance, and Covariance | (10) |
| | | VICReg [56] | Variance, Invariance, and Covariance | (16) |
| Combinations | | SS-GAN [144] | GAN loss + Rotation Prediction | (19 & 20) |
| with Other | | S⁴L [145] | Supervised and Unsupervised Loss | (21) |
| Learning Paradigms | | SSL improving robustness [146] | Supervised and Self-supervised Adversarial Training Loss | (22) |
| | | unsupervised multi-view learning [64] | Self-supervised Loss on Multiple Views | (25) |

in CLIP, which effectively combines CL with multi-modal learning, leading to remarkable achievements. SSL has been extensively integrated with various machine learning tasks, showcasing its versatility and potential. It has been combined with clustering [68], semi-supervised learning [145], multi-task learning [157], [158], transfer learning [159]–[161], graph NNs [147], [162], [163], reinforcement learning [164]–[166], few-shot learning [167], [168], neural architecture search [169], robust learning [146], [170]–[172], and meta-learning [173], [174]. This diverse integration underscores the widespread applicability and impact of SSL in the machine learning domain.

## 3 APPLICATIONS

SSL initially emerged in the context of vowel class recognition [175], and subsequently, it was extended to encompass object extraction tasks [176]. SSL has found widespread applications in diverse domains, including CV, NLP, medical image analysis, and remote sensing (RS).

### 3.1 CV

Sharma et al. [177] introduced a fully convolutional volumetric AE for unsupervised deep embeddings learning of object shapes. In addition, SSL has been extensively applied to various aspects of image processing and CV: image inpainting [104], human parsing [178], [179], scene deocclusion [180], semantic image segmentation [181], [182], monocular vision [183], person reidentification (re-ID) [184], [185], visual odometry [186], scene flow estimation [187], knowledge distillation [188], optical flow prediction [189], vision-language navigation [190], physiological signal estimation [191], [192], image denoising [193], [194], object detection [195]–[197], super-resolution [198], [199], voxel prediction from 2D images [200], and ego-motion [201], [202]. These applications highlight the broad impact and relevance of SSL in the realm of image processing and CV.

#### 3.1.1 SSL models for videos

SSL has garnered widespread usage across various applications, including video representation learning [203]–[205] and video retrieval [206].

3.1.1.1 Temporal information in videos: Various forms of temporal information in videos can be employed, encompassing frame order, video playback direction, video playback speed, and future prediction information [207], [208]. 1) The order of the frames. Several studies have explored the significance of frame order in videos. Misra et

al. [9] introduced a method for learning visual representations from raw spatiotemporal signals and determining the correct temporal sequence of frames extracted from videos. Fernando et al. [209] proposed a novel self-supervised CNN pre-training approach called "odd-one-out learning," where the objective is to identify the unrelated or odd element within a set of related elements. This odd element corresponds to a video subsequence with an incorrect temporal frame order, while the related elements maintain the correct temporal order. Lee et al. [210] employed temporally shuffled frames, presented in a non-chronological order, as inputs to train a CNN for predicting the correct order of the shuffled sequences, effectively using temporal coherence as a self-supervised signal. Building upon this work, Xu et al. [211] utilized temporally shuffled clips as inputs instead of individual frames, training 3D CNNs to sort these shuffled clips. 2) Video playback direction. Temporal direction analysis in videos, as studied by Wei et al. [10], involves discerning the arrow of time to determine if a video sequence progresses in the forward or backward direction. 3) Video playback speed. Video playback speed has been a subject of investigation in several studies. Benaim et al. [212] focused on predicting the speeds of moving objects in videos, determining whether they moved faster or slower than the normal speed. Yao et al. [213] leveraged playback rates and their corresponding video content as self-supervision signals for video representation learning. Additionally, Wang et al. [214] addressed the challenge of self-supervised video representation learning through the lens of video pace prediction.

3.1.1.2 Motions of objects in videos: Diba et al. [215] focused on SSL of motions in videos by employing dynamic motion filters to enhance motion representations, particularly for improving human action recognition. The concept of SSL with videos (CoCLR) [216] bears similarities to SimCLR [52].

3.1.1.3 Multi-modal(ality) data in videos: The auditory and visual components in a video are intrinsically interconnected. Leveraging this correlation, Korbar et al. [217] employed a self-supervised temporal synchronization approach to learn comprehensive and effective models for both video and audio analysis. Similarly, other methodologies [62], [218] are also founded on joint video and audio modalities while certain studies [219]–[221] incorporated both video and text modalities. Moreover, Alayrac et al. [222] explored a tri-modal approach involving vision, audio, and language in videos. On a different note, Sermanet et al. [223] proposed a self-supervised technique for learning rep-

resentations and robotic behaviors from unlabeled videos captured from various viewpoints.

    3.1.1.4 Spatial-temporal coherence of objects in videos: Wang et al. [224] introduced a self-supervised algorithm for learning visual correspondence in unlabeled videos by utilizing cycle consistency in time as a self-supervised signal. Extensions of this work have been explored by Li et al. [225] and Jabri et al. [226]. Lai et al. [227] presented a memory-augmented self-supervised method that enables generalizable and accurate pixel-level tracking. Zhang et al. [228] employed spatial-temporal consistency of depth maps to mitigate forgetting during the learning process. Zhao et al. [229] proposed a novel self-supervised algorithm named the "video cloze procedure (VCP)," which facilitates learning rich spatial-temporal representations for videos. Feichtenhofer et al. [112] extended the MAE to video representation learning and demonstrated that leveraging the naive ViT along with the spatiotemporal co-occurrence of videos can outperform the vanilla supervised training. Gupta et al. [118] demonstrated the importance of asymmetrically modeling the spatiotemporal information of videos.

### 3.1.2 Universal sequential SSL models for image processing and CV

Contrastive predictive coding (CPC) [58] operates on the fundamental concept of acquiring informative representations through latent space predictions of future data using robust autoregressive models. While initially applied to sequential data like speech and text, CPC has also found applicability to images [230].

    Drawing inspiration from the accomplishments of GPT [102], [231] in NLP, iGPT [103] investigates whether similar models can effectively learn representations for images. iGPT explores two training objectives, namely autoregressive prediction and a denoising objective, thereby sharing similarities with BERT [11]. In high-resolution scenarios, this approach [103] competes favorably with other self-supervised methods on ImageNet [1]. Similar to iGPT, ViT [5] also adopts a transformer architecture for vision tasks. By applying a pure transformer to sequences of image patches, ViT has demonstrated outstanding performance in image recognition tasks. The transformer architecture has been further extended to various vision-related applications, as evidenced by [70], [82], [83], [99], [232].

### 3.2 NLP

In the realm of NLP, pioneering works for performing SSL on word embeddings include the continuous bag-of-words model and the continuous skip-gram model [233]. They can be considered as belonging to generative self-supervised learning algorithms, which have long dominated the field of NLP. Despite their diverse forms, these algorithms are fundamentally based on language models that employ maximum likelihood estimation. Discriminative algorithms (e.g., contrastive learning) were initially deemed ineffective due to the distinct semantics inherent in language. Some discriminative algorithms aim to challenge the conventions, among which ELECTRA [234] stands out as a pioneer. ELECTRA employs the Replaced Token Detection (RTD) task and draws upon the structure and ideas of GANs

(notably, without adopting GAN's training paradigm) to pre-train a language model. [235] demonstrated that supervised contrastive pretraining enables zero-shot prediction of unseen text classes and enhances few-shot performance. A series of subsequent works have demonstrated that task-agnostic self-supervised contrastive pre-training has been shown to improve language modeling [236]–[238]. However, the automatic creation of textual input augmentation remains a significant challenge, as a single token can reverse the meaning of a sentence. Generative SSL algorithms continue to dominate NLP, from early works such as BERT and GPT to recent trillion-scale large language models.

### 3.3 Other fields

Within the medical field [239], the availability of labeled data is typically limited, while a vast amount of unlabeled data exists. This natural scenario makes SSL a compelling approach, which has been effectively employed for various tasks like medical image segmentation [240] and 3D medical image analysis [241]. Recently, SSL has also found applications in the remote sensing domain, benefiting from the abundance of large-scale unlabeled data that remains largely unexplored. For example, SeCo [242] leverages seasonal changes in RS images to construct positive pairs and perform CL. On the other hand, RVSA [243] introduces a novel rotated varied-size window attention mechanism that advances the plain vision transformer to serve as a fundamental model for various remote sensing tasks. Notably, it is pre-trained using the generative SSL method MAE [70] on the large-scale MillionAID dataset.

## 4 PERFORMANCE COMPARISON

Once a pre-trained model is obtained through SSL, the assessment of its performance becomes necessary. The conventional approach involves gauging the achieved performance on downstream tasks to ascertain the quality of the extracted features. However, this evaluation metric does not provide insights into what the network has specifically learned during self-supervised pre-training. To delve into the interpretability of self-supervised features, alternative evaluation metrics, such as network dissection [245] and other unsupervised methods [246], can be employed. In this section, we aim to present a clear demonstration of the performance comparison. We summarize the pre-trained dataset performance and transfer learning efficacy of typical SSL methods on well-established datasets. Note that SSL can technically be applied to diverse modalities. However, for the sake of simplicity, we narrow our focus to SSL in the vision domain.

### 4.1 Comprehensive comparison

We present the results in Table 4 and 5. In cases where a method reproduced from another subsequent work achieves superior accuracy compared to the original paper, we report the results with the higher one. Please note that although we have endeavored to align the experimental settings, minor variations in hyper-parameters can still affect the performance. Refer to the original paper if necessary. The experimental results are obtained according to the default

TABLE 4: Experimental results of the tested algorithms for linear classification and transfer learning tasks. DB denotes the default batch size. The symbol "-" indicates the absence or unavailability of the data point in the respective paper. The subscripts A, R, and V represent AlexNet, ResNet-50, and ViT-B, respectively. The superscript "e" indicates the utilization of extra data, specifically VOC2012.

| Methods | Linear Probe | Fine-Tuning | VOC_det | VOC_seg | COCO_det | COCO_seg | ADE20K_seg | DB |
|---|---|---|---|---|---|---|---|---|
| **Random:** | | | | | | | | |
| Random | $17.1_A$ [8] | - | $60.2_R^e$ [69] | $19.8_A$ [8] | $36.7_R$ [50] | $33.7_R$ [50] | - | - |
| R50 Sup | 76.5 [68] | 76.5 [68] | $81.3^e$ [69] | 74.4 [67] | 40.6 [50] | 36.8 [50] | - | - |
| ViT-B Sup | 82.3 [70] | 82.3 [70] | - | - | 47.9 [70] | 42.9 [70] | 47.4 [70] | - |
| **Context-Based:** | | | | | | | | |
| Jigsaw [8] | $45.7_R$ [68] | 54.7 | $61.4_R$ [42] | 37.6 | - | - | - | 256 |
| Colorization [38] | $39.6_R$ [68] | 40.7 [7] | 46.9 | 35.6 | - | - | - | - |
| Rotation [7] | 38.7 | 50.0 | 54.4 | 39.1 | - | - | - | 128 |
| **CL Based on Negative Examples:** | | | | | | | | |
| Examplar [132] | 31.5 [48] | - | - | - | - | - | - | - |
| Instdisc [48] | 54.0 | - | 65.4 | - | - | - | - | 256 |
| MoCo v1 [50] | 60.6 | - | 74.9 | - | 40.8 | 36.9 | - | 256 |
| SimCLR [52] | $73.9_V$ [82] | - | $81.8^e$ [69] | - | 37.9 [69] | 33.3 [69] | - | 4096 |
| MoCo v2 [51] | 72.2 [69] | - | $82.5^e$ | - | 39.8 [56] | 36.1 [56] | - | 256 |
| MoCo v3 [82] | 76.7 | 83.2 | - | - | 47.9 [70] | 42.7 [70] | 47.3 [70] | 4096 |
| **CL Based on Clustering:** | | | | | | | | |
| SwAV [68] | 75.3 | - | $82.6^e$ [56] | - | 41.6 | 37.8 [56] | - | 4096 |
| **CL Based on Self-distillation:** | | | | | | | | |
| BYOL [67] | 74.3 | - | $81.4^e$ [69] | 76.3 | 40.4 [56] | 37.0 [56] | - | 4096 |
| SimSiam [69] | 71.3 | - | $82.4^e$ [69] | - | 39.2 | 34.4 | - | 512 |
| DINO [83] | 78.2 | 83.6 [98] | - | - | 46.8 [100] | 41.5 [100] | 44.1 [99] | 1024 |
| **CL Based on Feature Decorrelation:** | | | | | | | | |
| Barlow Twins [55] | 73.2 | - | $82.6^e$ [56] | - | 39.2 | 34.3 | - | 2048 |
| VICReg [56] | 73.2 | - | $82.4^e$ | - | 39.4 | 36.4 | - | 2048 |
| **Masked Image Modeling (ViT-B by default):** | | | | | | | | |
| Context Encoder [104] | $21.0_A$ [7] | - | $44.5_A$ [7] | $30.0_A$ | - | - | - | - |
| BEiT v1 [99] | 56.7 [111] | 83.4 [98] | - | - | 49.8 [70] | 44.4 [70] | 47.1 [70] | 2000 |
| MAE [70] | 67.8 | 83.6 | - | - | 50.3 | 44.9 | 48.1 | 4096 |
| SimMIM [101] | 56.7 | 83.8 | - | - | $52.3_{Swin-B}$ [244] | - | $52.8_{Swin-B}$ [244] | 2048 |
| PeCo [107] | - | 84.5 | - | - | 43.9 | 39.8 | 46.7 | 2048 |
| iBOT [98] | 79.5 | 84.0 | - | - | 51.2 | 44.2 | 50.0 | 1024 |
| MimCo [110] | - | 83.9 | - | - | 44.9 | 40.7 | 48.91 | 2048 |
| CAE [100] | 70.4 | 83.9 | - | - | 50 | 44 | 50.2 | 2048 |
| data2vec [108] | - | 84.2 | - | - | - | - | - | 2048 |
| SdAE [109] | 64.9 | 84.1 | - | - | 48.9 | 43.0 | 48.6 | 768 |
| BEiT v2 [111] | 80.1 | 85.5 | - | - | - | - | 53.1 | 2048 |

backbone specified in the original papers, such as ResNet-50 or ViT-B/16. Additionally, results from alternative backbones were provided in instances where data using the default backbone was not available, and marked accordingly.

**Setup**. Upon the 2D image, the model was pre-trained on ImageNet-1k [1] and evaluated on semantic segmentation tasks on PASCAL VOC [247], COCO [248], and ADE20k [249], [250], as well as on object detection tasks on VOC and COCO, and on classification tasks on ImageNet-1k. Upon the video, the model is pre-trained on the Kinetics [251] or Something Something-v2 (SSv2) [252] datasets, and its performance is evaluated on action detection tasks on the Kinetics, SSv2, and AVA [253] datasets.

The evaluation of object detection on the PASCAL VOC dataset employs mean average precision (mAP), specifically $AP_{50}$. By default, the object detection task on PASCAL VOC employs VOC2007 for training. However, certain methods employ the combined 07+12 dataset and are annotated with a superscript "e". As for the object detection and instance segmentation tasks on COCO, we adopt the bounding-box AP ($AP_{bb}$) and mask AP ($AP_{mk}$) metrics, in accordance with [50]. The results on video understanding are evaluated using fine-tuned Top-1 accuracy as the metric.

## 4.2 Summary

First, the linear probe performance of contrastive learning models typically surpasses that of other algorithms, and contrastive learning approaches tend to regard the linear probe as a significant performance metric. This superiority is attributed to contrastive learning generating well-structured latent spaces, wherein distinct categories are effectively separated, and similar categories are appropriately clustered.

Secondly, it is observed that pre-trained models using MIM can be fine-tuned to achieve superior performance in most cases. Conversely, pre-trained models based on CL lack this property. One primary reason for this discrepancy lies in the increased susceptibility of CL-based models to overfitting [66], [262], [263]. This observation also extends to the fine-tuning of pre-trained models for downstream tasks. MIM-based approaches consistently exhibit substantial performance enhancements in downstream tasks, while CL-based methods offer comparatively limited assistance.

Thirdly, CL-based methods tend to employ resource-intensive techniques like momentum encoders, memory queues, and multi-crop, significantly increasing the demands on computing, storage, and communication resources. In contrast, MIM-based methods have a more efficient resource utilization, possibly attributed to the absence of example interactions. This advantageous property allows MIM-based algorithms to easily scale up models and

TABLE 5: Performance comparison of SSL methods for video.

| | | | | Downstream Dataset | | | |
| | | | | UCF101 [254] | | HMDB51 [255] | |
| Contrastive Methods | | | | | | | |
| Method | Pre-training Dataset | Backbone | Linear Probe | Linear | Fine-tune | Linear | Fine-tune |
|---|---|---|---|---|---|---|---|
| DSM [256] | K400 | R3D34 | - | - | 78.2 | - | 52.8 |
| TCE [257] | K400 | R50 | - | - | 71.2 | - | 36.6 |
| CoCRL [216] | K400 | S3D-G | - | 74.5 [258] | 87.9 | 46.1 [258] | 54.6 |
| CoCRL | K400 | 2×S3D-G | - | - | 90.6 | - | 62.9 |
| VTHCL [259] | K400 | R3D50 | 37.8 [260] | - | 82.1 | - | 49.2 |
| CVRL [261] | K400 | R3D50 | 66.1 | 89.2 | 92.2 | 57.3 | 66.7 |
| CVRL | K600 | R3D50 | 70.4 | 90.6 | 93.4 | 59.7 | 68.0 |
| $\rho$BYOL [260] | K400 | R3D50 | 71.5 | - | 95.5 | - | 73.6 |
| $\rho$BYOL | K400 | S3D-G | - | - | 96.3 | - | 75.0 |
| BraVe [258] | K400 | R3D50 | - | 90.6 | 93.7 | 65.1 | 72.0 |
| BraVe | K600 | R3D50 | 69.1 | 91.9 | 94.4 | 67.6 | 73.9 |

| | | | Downstream Dataset | | |
| Masked Image Modeling Methods | | | | | |
| Method | Pre-training Dataset | Backbone | K400 [251] | SSv2 [252] | AVA [253] |
|---|---|---|---|---|---|
| MaskFeat [106] | K400 | MViTv2-L/312 | 86.4 | 74.4 | 37.5 |
| BEVT [115] | K400 | Swin-B | 76.2 | 67.1 | - |
| BEVT | IN1K + K400 | Swin-B | 80.6 | 70.6 | - |
| VidelMAE [116] | K400 | ViT-B | 80.0 | 68.5 | 26.7 |
| VidelMAE | SSv2 | ViT-B | 69.6 | 79.6 | - |
| VidelMAE | SSv2 | ViT-L | - | 75.4 | 34.3 |
| MAE-ST [112] | K400 | ViT-L | 84.8 | 72.1 | 32.3 |
| OmniMAE [117] | IN1K + K400 | ViT-B | 80.8 | 69.0 | - |
| OmniMAE | IN1K + SSv2 | ViT-B | 80.6 | 69.5 | - |
| OmniMAE | IN1K + SSv2 | ViT-L | 84.0 | 74.2 | - |

data, efficiently leveraging modern GPUs for high parallel computing. We compared the computational complexity of different SSL methods in Table 1 of the Appendix. Note that the primary sources of time complexity and memory consumption are the neural network other than SSL components, e.g., the calculation of the cross-correlation matrix in Barlow Twins.

# 5 CONCLUSIONS, FUTURE TRENDS, AND OPEN QUESTIONS

In summary, this comprehensive review offers essential insights into contemporary SSL research, providing newcomers with an overall picture of the field. The paper presents a thorough survey of SSL from three main perspectives: algorithms, applications, and future trends. We focus on mainstream visual SSL algorithms, classifying them into four major types: context-based methods, generative methods, contrastive methods, and contrastive generative methods. Furthermore, we investigate the correlation between SSL and other learning paradigms. Lastly, we will delve into future trends and open problems as outlined below.

**Main trends**: Firstly, the theoretical cloud still looms over SSL. How can we understand different SSL algorithms and unify them in the same way physics seeks to unify the four fundamental forces? [54] analyzed the key properties of contrastive learning based on negative samples, enhancing the understanding of representation distributions. [78] rethought contrastive learning from the perspective of spectral decomposition, providing a high-level understanding of why contrastive learning is effective. [264] showed

practical properties, with InfoMin [64] indicating that the design of views should consider downstream tasks. [265] investigated why distillation-based methods do not collapse. [266] demonstrated the duality between negative example-based contrastive learning and covariance regularization-based methods such as Barlow Twins, indicating the latter can be seen as contrastive between the dimensions of the embeddings instead of between the samples. [267] demonstrated that introducing discrete sparse overcomplete representations for SSL can improve generalization. [268] presented the connections and distinctions among various SSL methods from the perspective of gradients. We anticipate that new theoretical studies will aid in comprehending and unifying various SSL approaches, particularly in harmonizing CL-based methods with MIM-based methods.

Secondly, a crucial question arises concerning the automatic design of an optimal pretext task to enhance the performance of a fixed downstream task. Various methods have been proposed to address this challenge, including the pixel-to-propagation consistency method [65] and dense contrastive learning [269]. However, this problem remains insufficiently resolved, and further theoretical investigations are warranted in this direction.

Thirdly, there is a pressing need for a unified SSL paradigm that encompasses multiple modalities. MIM has demonstrated remarkable progress in vision tasks, akin to the success of masked language model in NLP, suggesting the possibility of unifying learning paradigms. Additionally, the ViT architecture bridges the gap between visual and verbal modalities, enabling the construction of a unified transformer model for both CV and NLP tasks. Recent en-

deavors [108], [270] have sought to unify SSL models, yielding impressive results in downstream tasks and showing broad applicability. Nevertheless, NLP has advanced further in leveraging SSL models, prompting the CV community to draw inspiration from NLP approaches to effectively harness the potential of pre-trained models.

**Open problems**: Can SSL effectively leverage vast amounts of unlabeled data? How does it consistently benefit from additional unlabeled data, and how can we determine the theoretical inflection point?

Secondly, it is pertinent to explore the interconnection between SSL and multi-modality learning, as both methodologies share resemblances with the cognitive processes observed in infants. Consequently, a critical inquiry arises: how can these two approaches be synergistically integrated to forge a robust and comprehensive learning model?

Thirdly, determining the most optimal or recommended SSL algorithm poses a challenge as there is no universally applicable solution. The ideal selection of an algorithm should align with the specific problem structure, yet practical situations often complicate this process. Consequently, the development of a checklist to aid users in identifying the most suitable method under particular circumstances warrants investigation and should be pursued as a promising avenue for future research.

Fourthly, the assumption that unlabeled data invariably leads to improved outcomes warrants scrutiny. Our hypothesis challenges this notion, especially concerning semi-supervised learning methods, as the *no free lunch* theorem comes into play. Performance degradation can arise when model assumptions fail to align effectively with the underlying problem structure. For instance, if a model assumes a substantial separation between decision boundaries and regions of high data density, it may perform poorly when faced with data originating from heavily overlapping Cauchy distributions, as the decision boundary would traverse through dense areas. However, preemptively identifying such mismatches remains intricate and an unresolved matter. Consequently, this topic merits further research to shed light on the matter.

## ACKNOWLEDGMENTS

## REFERENCES

[1] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 248–255, 2009.

[2] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, *et al.*, "Learning transferable visual models from natural language supervision," in *Int. Conf. Mach. Learn.*, pp. 8748–8763, 2021.

[3] L. Ericsson, H. Gouk, and T. M. Hospedales, "How well do self-supervised models transfer?," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 5414–5423, 2021.

[4] X. Liu, F. Zhang, Z. Hou, L. Mian, Z. Wang, J. Zhang, and J. Tang, "Self-supervised learning: Generative or contrastive," *IEEE T. Knowl. Data Eng.*, 2022.

[5] 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 *Int. Conf. Learn. Represent.*, 2021.

[6] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, "Learning spatiotemporal features with 3d convolutional networks," in *IEEE Int. Conf. Comput. Vis.*, pp. 4489–4497, 2015.

[7] S. Gidaris, P. Singh, and N. Komodakis, "Unsupervised representation learning by predicting image rotations," in *Int. Conf. Learn. Represent.*, pp. 1–14, 2018.

[8] M. Noroozi and P. Favaro, "Unsupervised learning of visual representations by solving jigsaw puzzles," in *Eur. Conf. Comput. Vis.*, pp. 69–84, 2016.

[9] I. Misra, C. L. Zitnick, and M. Hebert, "Shuffle and learn: unsupervised learning using temporal order verification," in *Eur. Conf. Comput. Vis.*, pp. 527–544, 2016.

[10] D. Wei, J. J. Lim, A. Zisserman, and W. T. Freeman, "Learning and using the arrow of time," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 8052–8060, 2018.

[11] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.

[12] X. Zeng, Y. Pan, M. Wang, J. Zhang, and Y. Liu, "Realistic face reenactment via self-supervised disentangling of identity and pose," in *AAAI Conf.Artif. Intell.*, pp. 12154–12163, 2020.

[13] A. Miech, J.-B. Alayrac, L. Smaira, I. Laptev, J. Sivic, and A. Zisserman, "End-to-end learning of visual representations from uncurated instructional videos," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 9879–9889, 2020.

[14] Y. M. Asano, C. Rupprecht, and A. Vedaldi, "A critical analysis of self-supervision, or what we can learn from a single image," in *Int. Conf. Learn. Represent.*, 2020.

[15] G. E. Hinton and R. R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *Science*, vol. 313, no. 5786, pp. 504–507, 2006.

[16] P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol, "Extracting and composing robust features with denoising autoencoders," in *Int. Conf. Mach. Learn.*, pp. 1096–1103, 2008.

[17] L. Pinto and A. Gupta, "Supersizing self-supervision: Learning to grasp from 50k tries and 700 robot hours," in *IEEE Int. Conf. Robot. Autom.*, pp. 3406–3413, 2016.

[18] Y. Li, M. Paluri, J. M. Rehg, and P. Dollár, "Unsupervised learning of edges," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 1619–1627, 2016.

[19] D. Li, W.-C. Hung, J.-B. Huang, S. Wang, N. Ahuja, and M.-H. Yang, "Unsupervised visual representation learning by graph-based consistent constraints," in *Eur. Conf. Comput. Vis.*, pp. 678–694, 2016.

[20] H. Lee, S. J. Hwang, and J. Shin, "Rethinking data augmentation: Self-supervision and self-distillation," *arXiv preprint arXiv:1910.05872*, 2019.

[21] B. Zoph, G. Ghiasi, T.-Y. Lin, Y. Cui, H. Liu, E. D. Cubuk, and Q. Le, "Rethinking pre-training and self-training," in *Neural Inf. Process. Syst.*, pp. 1–13, 2020.

[22] A. E. Orhan, V. V. Gupta, and B. M. Lake, "Self-supervised learning through the eyes of a child," in *Neural Inf. Process. Syst.*, pp. 9960–9971, 2020.

[23] J. Mitrovic, B. McWilliams, J. Walker, L. Buesing, and C. Blundell, "Representation learning via invariant causal mechanisms," in *Int. Conf. Learn. Represent.*, pp. 1–19, 2021.

[24] T. Hua, W. Wang, Z. Xue, S. Ren, Y. Wang, and H. Zhao, "On feature decorrelation in self-supervised learning," in *IEEE Int. Conf. Comput. Vis.*, pp. 9598–9608, 2021.

[25] VentureBeat, "Yann LeCun, Yoshua Bengio: Self-supervised learning is key to human-level intelligence." https://cacm.acm.org/news/244720-yann-lecun-yoshua-bengio-self-supervised-learning-is-key-to-human-level-intelligence/fulltext.

[26] J. Yu, H. Yin, X. Xia, T. Chen, J. Li, and Z. Huang, "Self-supervised learning for recommender systems: A survey," *arXiv preprint arXiv:2203.15876*, 2022.

[27] Y. Liu, M. Jin, S. Pan, C. Zhou, Y. Zheng, F. Xia, and P. Yu, "Graph self-supervised learning: A survey," *IEEE T. Knowl. Data Eng.*, 2022.

[28] H. H. Mao, "A survey on self-supervised pre-training for sequential transfer learning in neural networks," *arXiv preprint arXiv:2007.00800*, 2020.

[29] M. C. Schiappa, Y. S. Rawat, and M. Shah, "Self-supervised learning for videos: A survey," *arXiv preprint arXiv:2207.00419*, 2022.

[30] G.-J. Qi and M. Shah, "Adversarial pretraining of self-supervised deep networks: Past, present and future," *arXiv preprint arXiv:2210.13463*, 2022.

[31] A. Jaiswal, A. R. Babu, M. Z. Zadeh, D. Banerjee, and F. Makedon, "A survey on contrastive self-supervised learning," *Technologies*, vol. 9, no. 1, pp. 1–22, 2020.

[32] V. R. de Sa, "Learning classification with unlabeled data," in *Neural Inf. Process. Syst.*, pp. 112–119, 1994.

[33] Y. LeCun and Y. Bengio, "Reflections from the turing award winners." https://iclr.cc/virtual_2020/speaker_7.html.

[34] L. Jing and Y. Tian, "Self-supervised visual feature learning with deep neural networks: A survey," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 43, no. 11, pp. 4037–4058, 2021.

[35] J. Gui, Z. Sun, Y. Wen, D. Tao, and J. Ye, "A review on generative adversarial networks: Algorithms, theory, and applications," *IEEE T. Knowl. Data Eng.*, 2022.

[36] T. Nathan Mundhenk, D. Ho, and B. Y. Chen, "Improvements to context based self-supervised learning," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 9339–9348, 2018.

[37] P. Agrawal, J. Carreira, and J. Malik, "Learning to see by moving," in *IEEE Int. Conf. Comput. Vis.*, pp. 37–45, 2015.

[38] R. Zhang, P. Isola, and A. A. Efros, "Colorful image colorization," in *Eur. Conf. Comput. Vis.*, pp. 649–666, 2016.

[39] G. Larsson, M. Maire, and G. Shakhnarovich, "Learning representations for automatic colorization," in *Eur. Conf. Comput. Vis.*, pp. 577–593, 2016.

[40] R. Zhang, J.-Y. Zhu, P. Isola, X. Geng, A. S. Lin, T. Yu, and A. A. Efros, "Real-time user-guided image colorization with learned deep priors," *arXiv preprint arXiv:1705.02999*, 2017.

[41] G. Larsson, M. Maire, and G. Shakhnarovich, "Colorization as a proxy task for visual understanding," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 6874–6883, 2017.

[42] P. Goyal, D. Mahajan, A. Gupta, and I. Misra, "Scaling and benchmarking self-supervised visual representation learning," in *IEEE Int. Conf. Comput. Vis.*, pp. 6391–6400, 2019.

[43] U. Ahsan, R. Madhok, and I. Essa, "Video jigsaw: Unsupervised learning of spatiotemporal context for video action recognition," in *Proc. Winter Conf. Appl. Comput. Vis.*, pp. 179–189, 2019.

[44] X. Zhan, X. Pan, Z. Liu, D. Lin, and C. C. Loy, "Self-supervised learning via conditional motion propagation," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 1881–1889, 2019.

[45] K. Wang, L. Lin, C. Jiang, C. Qian, and P. Wei, "3d human pose machines with self-supervised learning," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 42, no. 5, pp. 1069–1082, 2019.

[46] M. Noroozi, H. Pirsiavash, and P. Favaro, "Representation learning by learning to count," in *IEEE Int. Conf. Comput. Vis.*, pp. 5898–5906, 2017.

[47] I. Misra and L. v. d. Maaten, "Self-supervised learning of pretext-invariant representations," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 6707–6717, 2020.

[48] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin, "Unsupervised feature learning via non-parametric instance discrimination," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 3733–3742, 2018.

[49] N. Zhao, Z. Wu, R. W. Lau, and S. Lin, "What makes instance discrimination good for transfer learning?," in *Int. Conf. Learn. Represent.*, pp. 1–11, 2021.

[50] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 9729–9738, 2020.

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

[52] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *Int. Conf. Mach. Learn.*, pp. 1597–1607, 2020.

[53] T. Chen, S. Kornblith, K. Swersky, M. Norouzi, and G. Hinton, "Big self-supervised models are strong semi-supervised learners," in *Neural Inf. Process. Syst.*, pp. 1–13, 2020.

[54] T. Wang and P. Isola, "Understanding contrastive representation learning through alignment and uniformity on the hypersphere," in *Int. Conf. Mach. Learn.*, pp. 9929–9939, 2020.

[55] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny, "Barlow twins: Self-supervised learning via redundancy reduction," in *Int. Conf. Mach. Learn.*, 2021.

[56] A. Bardes, J. Ponce, and Y. LeCun, "Vicreg: Variance-invariance-covariance regularization for self-supervised learning," in *Int. Conf. Learn. Represent.*, pp. 1–12, 2022.

[57] R. Hadsell, S. Chopra, and Y. LeCun, "Dimensionality reduction by learning an invariant mapping," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 1735–1742, 2006.

[58] A. v. d. Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2019.

[59] M. Gutmann and A. Hyvärinen, "Noise-contrastive estimation: A new estimation principle for unnormalized statistical models," in *Int. Conf. Artif. Intell. Statist.*, pp. 297–304, 2010.

[60] M. Zheng, S. You, F. Wang, C. Qian, C. Zhang, X. Wang, and C. Xu, "Ressl: Relational self-supervised learning with weak augmentation," *arXiv preprint arXiv:2107.09282*, 2021.

[61] N. Zhao, Z. Wu, R. W. Lau, and S. Lin, "Distilling localization for self-supervised representation learning," in *AAAI Conf. Artif. Intell.*, pp. 10990–10998, 2021.

[62] R. Arandjelovic and A. Zisserman, "Objects that sound," in *Eur. Conf. Comput. Vis.*, pp. 435–451, 2018.

[63] Y. Tian, D. Krishnan, and P. Isola, "Contrastive multiview coding," in *Eur. Conf. Comput. Vis.*, pp. 776–794, 2020.

[64] Y. Tian, C. Sun, B. Poole, D. Krishnan, C. Schmid, and P. Isola, "What makes for good views for contrastive learning," in *Neural Inf. Process. Syst.*, pp. 1–13, 2020.

[65] Z. Xie, Y. Lin, Z. Zhang, Y. Cao, S. Lin, and H. Hu, "Propagate yourself: Exploring pixel-level consistency for unsupervised visual representation learning," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 16684–16693, 2021.

[66] X. Wang and G.-J. Qi, "Contrastive learning with stronger augmentations," *IEEE Trans. Pattern Anal. Mach. Intell.*, pp. 1–12, 2022.

[67] 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 *Neural Inf. Process. Syst.*, pp. 1–14, 2020.

[68] M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin, "Unsupervised learning of visual features by contrasting cluster assignments," in *Neural Inf. Process. Syst.*, 2020.

[69] X. Chen and K. He, "Exploring simple siamese representation learning," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 15750–15758, 2021.

[70] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, "Masked autoencoders are scalable vision learners," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 16000–16009, 2022.

[71] M. Tschannen, J. Djolonga, P. K. Rubenstein, S. Gelly, and M. Lucic, "On mutual information maximization for representation learning," in *Int. Conf. Learn. Represent.*, pp. 1–12, 2020.

[72] N. Saunshi, O. Plevrakis, S. Arora, M. Khodak, and H. Khandeparkar, "A theoretical analysis of contrastive unsupervised representation learning," in *Int. Conf. Mach. Learn.*, pp. 5628–5637, 2019.

[73] Y. Yang and Z. Xu, "Rethinking the value of labels for improving class-imbalanced learning," in *Neural Inf. Process. Syst.*, 2020.

[74] Y.-H. H. Tsai, Y. Wu, R. Salakhutdinov, and L.-P. Morency, "Self-supervised learning from a multi-view perspective," *arXiv preprint arXiv:2006.05576*, 2020.

[75] C.-Y. Chuang, J. Robinson, L. Yen-Chen, A. Torralba, and S. Jegelka, "Debiased contrastive learning," in *Int. Conf. Learn. Represent.*, 2020.

[76] J. D. Lee, Q. Lei, N. Saunshi, and J. Zhuo, "Predicting what you already know helps: Provable self-supervised learning," *arXiv preprint arXiv:2008.01064*, 2020.

[77] S. Chen, G. Niu, C. Gong, J. Li, J. Yang, and M. Sugiyama, "Large-margin contrastive learning with distance polarization regularizer," in *Int. Conf. Mach. Learn.*, pp. 1673–1683, 2021.

[78] J. Z. HaoChen, C. Wei, A. Gaidon, and T. Ma, "Provable guarantees for self-supervised deep learning with spectral contrastive loss," in *Neural Inf. Process. Syst.*, Nov. 2021.

[79] C. Tosh, A. Krishnamurthy, and D. Hsu, "Contrastive learning, multi-view redundancy, and linear models," in *Algorithmic Learning Theory*, pp. 1179–1206, 2021.

[80] C. Wei, K. Shen, Y. Chen, and T. Ma, "Theoretical analysis of self-training with deep networks on unlabeled data," in *Int. Conf. Learn. Represent.*, pp. 1–15, 2021.

[81] Y. Tian, "Deep contrastive learning is provably (almost) principal component analysis," *arXiv preprint arXiv:2201.12680*, 2022.

[82] X. Chen, S. Xie, and K. He, "An empirical study of training self-supervised visual transformers," in *IEEE Int. Conf. Comput. Vis.*, pp. 9640–9649, 2021.

[83] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin, "Emerging properties in self-supervised vision transformers," in *IEEE Int. Conf. Comput. Vis.*, pp. 9650–9660, 2021.

[84] Y. Wang, X. Shen, S. X. Hu, Y. Yuan, J. L. Crowley, and D. Vaufreydaz, "Self-supervised transformers for unsupervised object discovery using normalized cut," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 14543–14553, 2022.

[85] E. Hoffer, I. Hubara, and N. Ailon, "Deep unsupervised learning through spatial contrasting," *arXiv preprint arXiv:1610.00243*, 2016.

[86] Y. Xu, Q. Zhang, J. Zhang, and D. Tao, "Regioncl: exploring contrastive region pairs for self-supervised representation learning," in *Eur. Conf. Comput. Vis.*, pp. 477–494, Springer, 2022.

[87] M. Yang, M. Liao, P. Lu, J. Wang, S. Zhu, H. Luo, Q. Tian, and X. Bai, "Reading and writing: Discriminative and generative modeling for self-supervised text recognition," *arXiv preprint arXiv:2207.00193*, 2022.

[88] R. Zhu, B. Zhao, J. Liu, Z. Sun, and C. W. Chen, "Improving contrastive learning by visualizing feature transformation," in *IEEE Int. Conf. Comput. Vis.*, pp. 10306–10315, 2021.

[89] M. Yang, Y. Li, Z. Huang, Z. Liu, P. Hu, and X. Peng, "Partially view-aligned representation learning with noise-robust contrastive loss," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 1134–1143, 2021.

[90] A. Islam, C.-F. Chen, R. Panda, L. Karlinsky, R. Radke, and R. Feris, "A broad study on the transferability of visual representations with contrastive learning," in *IEEE Int. Conf. Comput. Vis.*, pp. 8845–8855, 2021.

[91] J. Li, C. Xiong, and S. C. Hoi, "Learning from noisy data with robust representation learning," in *IEEE Int. Conf. Comput. Vis.*, pp. 9485–9494, 2021.

[92] L. Jing, P. Vincent, Y. LeCun, and Y. Tian, "Understanding dimensional collapse in contrastive self-supervised learning," in *Int. Conf. Learn. Represent.*, pp. 1–11, 2022.

[93] J. Zhang, X. Xu, F. Shen, Y. Yao, J. Shao, and X. Zhu, "Video representation learning with graph contrastive augmentation," in *ACM Int. Conf. Multimedia*, pp. 3043–3051, 2021.

[94] Q. Hu, X. Wang, W. Hu, and G.-J. Qi, "Adco: Adversarial contrast for efficient learning of unsupervised representations from self-trained negative adversaries," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, 2021.

[95] Y. Kalantidis, M. B. Sariyildiz, N. Pion, P. Weinzaepfel, and D. Larlus, "Hard negative mixing for contrastive learning," in *Neural Inf. Process. Syst.*, pp. 1–12, 2020.

[96] S. Purushwalkam and A. Gupta, "Demystifying contrastive self-supervised learning: Invariances, augmentations and dataset biases," in *Neural Inf. Process. Syst.*, pp. 1–12, 2020.

[97] P. Khosla, P. Teterwak, C. Wang, A. Sarna, Y. Tian, P. Isola, A. Maschinot, C. Liu, and D. Krishnan, "Supervised contrastive learning," in *Neural Inf. Process. Syst.*, pp. 18661–18673, 2020.

[98] J. Zhou, C. Wei, H. Wang, W. Shen, C. Xie, A. Yuille, and T. Kong, "ibot: Image bert pre-training with online tokenizer," in *Int. Conf. Learn. Represent.*, pp. 1–12, 2022.

[99] H. Bao, L. Dong, S. Piao, and F. Wei, "Beit: Bert pre-training of image transformers," in *Int. Conf. Learn. Represent.*, pp. 1–13, 2022.

[100] X. Chen, M. Ding, X. Wang, Y. Xin, S. Mo, Y. Wang, S. Han, P. Luo, G. Zeng, and J. Wang, "Context autoencoder for self-supervised representation learning," *arXiv preprint arXiv:2202.03026*, 2022.

[101] 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," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 9653–9663, 2022.

[102] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, *et al.*, "Language models are few-shot learners," *arXiv preprint arXiv:2005.14165*, 2020.

[103] M. Chen, A. Radford, R. Child, J. Wu, H. Jun, P. Dhariwal, D. Luan, and I. Sutskever, "Generative pretraining from pixels," in *Int. Conf. Mach. Learn.*, pp. 1691–1703, 2020.

[104] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros, "Context encoders: Feature learning by inpainting," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 2536–2544, 2016.

[105] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever, "Zero-shot text-to-image generation," in *Int. Conf. Mach. Learn.*, pp. 8821–8831, 2021.

[106] C. Wei, H. Fan, S. Xie, C.-Y. Wu, A. Yuille, and C. Feichtenhofer, "Masked feature prediction for self-supervised visual pre-training," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 14668–14678, 2022.

[107] X. Dong, J. Bao, T. Zhang, D. Chen, W. Zhang, L. Yuan, D. Chen, F. Wen, and N. Yu, "Peco: Perceptual codebook for bert pre-training of vision transformers," *arXiv preprint arXiv:2111.12710*, 2021.

[108] 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.

[109] Y. Chen, Y. Liu, D. Jiang, X. Zhang, W. Dai, H. Xiong, and Q. Tian, "Sdae: Self-distillated masked autoencoder," in *Eur. Conf. Comput. Vis.*, pp. 108–124, 2022.

[110] Q. Zhou, C. Yu, H. Luo, Z. Wang, and H. Li, "Mimco: Masked image modeling pre-training with contrastive teacher," in *ACM Int. Conf. Multimedia*, pp. 4487–4495, 2022.

[111] Z. Peng, L. Dong, H. Bao, Q. Ye, and F. Wei, "Beit v2: Masked image modeling with vector-quantized visual tokenizers," *arXiv preprint arXiv:2208.06366*, 2022.

[112] C. Feichtenhofer, H. Fan, Y. Li, and K. He, "Masked autoencoders as spatiotemporal learners," *arXiv preprint arXiv:2205.09113*, 2022.

[113] Y. Liang, S. Zhao, B. Yu, J. Zhang, and F. He, "Meshmae: Masked autoencoders for 3d mesh data analysis," in *Eur. Conf. Comput. Vis.*, pp. 37–54, 2022.

[114] Y. Pang, W. Wang, F. E. Tay, W. Liu, Y. Tian, and L. Yuan, "Masked autoencoders for point cloud self-supervised learning," in *Eur. Conf. Comput. Vis.*, pp. 604–621, 2022.

[115] R. Wang, D. Chen, Z. Wu, Y. Chen, X. Dai, M. Liu, Y.-G. Jiang, L. Zhou, and L. Yuan, "Bevt: Bert pretraining of video transformers," in *Proceedings of the IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 14733–14743, 2022.

[116] Z. Tong, Y. Song, J. Wang, and L. Wang, "Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training," *Neural Inf. Process. Syst.*, vol. 35, pp. 10078–10093, 2022.

[117] R. Girdhar, A. El-Nouby, M. Singh, K. V. Alwala, A. Joulin, and I. Misra, "Omnimae: Single model masked pretraining on images and videos," in *Proceedings of the IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 10406–10417, June 2023.

[118] A. Gupta, J. Wu, J. Deng, and L. Fei-Fei, "Siamese masked autoencoders," in *Neural Inf. Process. Syst.*, Nov. 2023.

[119] Z. Liu, H. Hu, Y. Lin, Z. Yao, Z. Xie, Y. Wei, J. Ning, Y. Cao, Z. Zhang, L. Dong, *et al.*, "Swin transformer v2: Scaling up capacity and resolution," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 12009–12019, 2022.

[120] Y. Li, H. Mao, R. Girshick, and K. He, "Exploring plain vision transformer backbones for object detection," in *Eur. Conf. Comput. Vis.*, pp. 280–296, 2022.

[121] Y. Xu, J. Zhang, Q. Zhang, and D. Tao, "Vitpose: Simple vision transformer baselines for human pose estimation," in *Neural Inf. Process. Syst.*, pp. 38571–38584, 2022.

[122] Z. Liu, J. Gui, and H. Luo, "Good helper is around you: Attention-driven masked image modeling," in *AAAI Conf. Artif. Intell.*, pp. 1799–1807, 2023.

[123] Z. Qi, R. Dong, G. Fan, Z. Ge, X. Zhang, K. Ma, and L. Yi, "Contrast with reconstruct: Contrastive 3d representation learning guided by generative pretraining," *arXiv preprint arXiv:2302.02318*, 2023.

[124] Z. Xie, Z. Zhang, Y. Cao, Y. Lin, Y. Wei, Q. Dai, and H. Hu, "On data scaling in masked image modeling," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 10365–10374, 2023.

[125] M. Oquab, T. Darcet, T. Moutakanni, H. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haziza, F. Massa, A. El-Nouby, *et al.*, "Dinov2: Learning robust visual features without supervision," *arXiv preprint arXiv:2304.07193*, 2023.

[126] X. Kong and X. Zhang, "Understanding masked image modeling via learning occlusion invariant feature," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 6241–6251, 2023.

[127] H. Chen, Y. Wang, B. Lagadec, A. Dantcheva, and F. Bremond, "Joint generative and contrastive learning for unsupervised person re-identification," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 2004–2013, 2021.

[128] L. Wang, F. Liang, Y. Li, H. Zhang, W. Ouyang, and J. Shao, "Repre: Improving self-supervised vision transformer with reconstructive pre-training," Jan. 2022.

[129] Z. Huang, X. Jin, C. Lu, Q. Hou, M.-M. Cheng, D. Fu, X. Shen, and J. Feng, "Contrastive masked autoencoders are stronger vision learners," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–13, 2023.

[130] C. Tao, X. Zhu, W. Su, G. Huang, B. Li, J. Zhou, Y. Qiao, X. Wang, and J. Dai, "Siamese image modeling for self-supervised vision representation learning," in *Proceedings of the IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 2132–2141, 2023.

[131] Z. Xie, Z. Geng, J. Hu, Z. Zhang, H. Hu, and Y. Cao, "Revealing the dark secrets of masked image modeling," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 14475–14485, 2023.

[132] A. Dosovitskiy, J. T. Springenberg, M. Riedmiller, and T. Brox, "Discriminative unsupervised feature learning with convolutional neural networks," in *Neural Inf. Process. Syst.*, pp. 766–774, 2014.

[133] A. Dosovitskiy, P. Fischer, J. T. Springenberg, M. Riedmiller, and T. Brox, "Discriminative unsupervised feature learning with exemplar convolutional neural networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 9, pp. 1734–1747, 2015.

[134] C. Doersch, A. Gupta, and A. A. Efros, "Unsupervised visual representation learning by context prediction," in *IEEE Int. Conf. Vis.*, pp. 1422–1430, 2015.

[135] P. Bojanowski and A. Joulin, "Unsupervised learning by predicting noise," in *Int. Conf. Mach. Learn.*, 2017.

[136] J. Xie, R. Girshick, and A. Farhadi, "Unsupervised deep embedding for clustering analysis," in *Int. Conf. Mach. Learn.*, pp. 478–487, 2016.

[137] J. Yang, D. Parikh, and D. Batra, "Joint unsupervised learning of deep representations and image clusters," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 5147–5156, 2016.

[138] M. Caron, P. Bojanowski, A. Joulin, and M. Douze, "Deep clustering for unsupervised learning of visual features," in *Eur. Conf. Comput. Vis.*, pp. 132–149, 2018.

[139] R. Zhang, P. Isola, and A. A. Efros, "Split-brain autoencoders: Unsupervised learning by cross-channel prediction," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 1058–1067, 2017.

[140] X. Wang, K. He, and A. Gupta, "Transitive invariance for self-supervised visual representation learning," in *IEEE Int. Conf. Vis.*, pp. 1329–1338, 2017.

[141] A. Kolesnikov, X. Zhai, and L. Beyer, "Revisiting self-supervised visual representation learning," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 1920–1929, 2019.

[142] P. Krähenbühl, "Free supervision from video games," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 2955–2964, 2018.

[143] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Neural Inf. Process. Syst.*, pp. 2672–2680, 2014.

[144] T. Chen, X. Zhai, M. Ritter, M. Lucic, and N. Houlsby, "Self-supervised gans via auxiliary rotation loss," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 12154–12163, 2019.

[145] X. Zhai, A. Oliver, A. Kolesnikov, and L. Beyer, "S4l: Self-supervised semi-supervised learning," in *IEEE Int. Conf. Comput. Vis.*, pp. 1476–1485, 2019.

[146] D. Hendrycks, M. Mazeika, S. Kadavath, and D. Song, "Using self-supervised learning can improve model robustness and uncertainty," in *Neural Inf. Process. Syst.*, pp. 15663–15674, 2019.

[147] K. Hassani and A. H. Khasahmadi, "Contrastive multi-view representation learning on graphs," in *Int. Conf. Mach. Learn.*, 2020.

[148] L. Gomez, Y. Patel, M. Rusiñol, D. Karatzas, and C. Jawahar, "Self-supervised learning of visual features through embedding images into text topic spaces," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 4230–4239, 2017.

[149] L. Jing, Y. Chen, L. Zhang, M. He, and Y. Tian, "Self-supervised feature learning by cross-modality and cross-view correspondences," *arXiv preprint arXiv:2004.05749*, 2020.

[150] L. Jing, Y. Chen, L. Zhang, M. He, and Y. Tian, "Self-supervised modal and view invariant feature learning," *arXiv preprint arXiv:2005.14169*, 2020.

[151] L. Zhang and Z. Zhu, "Unsupervised feature learning for point cloud understanding by contrasting and clustering using graph convolutional neural networks," in *International Conference on 3D Vision*, pp. 395–404, 2019.

[152] Y. Yang, C. Feng, Y. Shen, and D. Tian, "Foldingnet: Point cloud auto-encoder via deep grid deformation," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 206–215, 2018.

[153] M. Gadelha, R. Wang, and S. Maji, "Multiresolution tree networks for 3d point cloud processing," in *Eur. Conf. Comput. Vis.*, pp. 103–118, 2018.

[154] Y. Zhao, T. Birdal, H. Deng, and F. Tombari, "3d point capsule networks," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 1009–1018, 2019.

[155] Y. Sun, X. Wang, Z. Liu, J. Miller, A. A. Efros, and M. Hardt, "Test-time training with self-supervision for generalization under distribution shifts," in *Int. Conf. Mach. Learn.*, 2020.

[156] Y. Gandelsman, Y. Sun, X. Chen, and A. A. Efros, "Test-time training with masked autoencoders," *arXiv preprint arXiv:2209.07522*, 2022.

[157] J. J. Sun, A. Kennedy, E. Zhan, D. J. Anderson, Y. Yue, and P. Perona, "Task programming: Learning data efficient behavior representations," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 2876–2885, 2021.

[158] Z. Ren and Y. Jae Lee, "Cross-domain self-supervised multi-task feature learning using synthetic imagery," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 762–771, 2018.

[159] K. Saito, D. Kim, S. Sclaroff, and K. Saenko, "Universal domain adaptation through self supervision," in *Neural Inf. Process. Syst.*, pp. 1–11, 2020.

[160] Y. Sun, E. Tzeng, T. Darrell, and A. A. Efros, "Unsupervised domain adaptation through self-supervision," *arXiv preprint arXiv:1909.11825*, 2019.

[161] M. Noroozi, A. Vinjimoor, P. Favaro, and H. Pirsiavash, "Boosting self-supervised learning via knowledge transfer," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 9359–9367, 2018.

[162] Z. Hu, Y. Dong, K. Wang, K.-W. Chang, and Y. Sun, "Gpt-gnn: Generative pre-training of graph neural networks," in *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1857–1867, 2020.

[163] Y. Rong, Y. Bian, T. Xu, W. Xie, Y. Wei, W. Huang, and J. Huang, "Self-supervised graph transformer on large-scale molecular data," in *Neural Inf. Process. Syst.*, 2020.

[164] U. Buchler, B. Brattoli, and B. Ommer, "Improving spatiotemporal self-supervision by deep reinforcement learning," in *Eur. Conf. Comput. Vis.*, pp. 770–786, 2018.

[165] D. Guo, B. A. Pires, B. Piot, J.-b. Grill, F. Altché, R. Munos, and M. G. Azar, "Bootstrap latent-predictive representations for multitask reinforcement learning," *arXiv preprint arXiv:2004.14646*, 2020.

[166] N. Hansen, Y. Sun, P. Abbeel, A. A. Efros, L. Pinto, and X. Wang, "Self-supervised policy adaptation during deployment," *arXiv preprint arXiv:2007.04309*, 2020.

[167] S. Gidaris, A. Bursuc, N. Komodakis, P. Pérez, and M. Cord, "Boosting few-shot visual learning with self-supervision," in *IEEE Int. Conf. Comput. Vis.*, pp. 8059–8068, 2019.

[168] J.-C. Su, S. Maji, and B. Hariharan, "Boosting supervision with self-supervision for few-shot learning," *arXiv preprint arXiv:1906.07079*, 2019.

[169] C. Li, T. Tang, G. Wang, J. Peng, B. Wang, X. Liang, and X. Chang, "Bossnas: Exploring hybrid cnn-transformers with block-wisely self-supervised neural architecture search," in *IEEE Int. Conf. Comput. Vis.*, 2021.

[170] L. Fan, S. Liu, P.-Y. Chen, G. Zhang, and C. Gan, "When does contrastive learning preserve adversarial robustness from pre-training to finetuning?," in *Neural Inf. Process. Syst.*, 2021.

[171] M. Kim, J. Tack, and S. J. Hwang, "Adversarial self-supervised contrastive learning," in *Neural Inf. Process. Syst.*, pp. 1–12, 2020.

[172] T. Chen, S. Liu, S. Chang, Y. Cheng, L. Amini, and Z. Wang, "Adversarial robustness: From self-supervised pre-training to fine-tuning," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 699–708, 2020.

[173] Y. Lin, X. Guo, and Y. Lu, "Self-supervised video representation learning with meta-contrastive network," in *IEEE Int. Conf. Comput. Vis.*, pp. 8239–8249, 2021.

[174] Y. An, H. Xue, X. Zhao, and L. Zhang, "Conditional self-supervised learning for few-shot classification," in *Int. Joint Conf. Artif. Intell.*, pp. 2140–2146, 2021.

[175] S. Pal, A. Datta, and D. D. Majumder, "Computer recognition of vowel sounds using a self-supervised learning algorithm," *Journal of the Anatomical Society of India*, pp. 117–123, 1978.

[176] A. Ghosh, N. R. Pal, and S. K. Pal, "Self-organization for object extraction using a multilayer neural network and fuzziness mearsures," *IEEE Transactions on Fuzzy Systems*, pp. 54–68, 1993.

[177] A. Sharma, O. Grau, and M. Fritz, "Vconv-dae: Deep volumetric shape learning without object labels," in *Eur. Conf. Comput. Vis.*, pp. 236–250, 2016.

[178] K. Gong, X. Liang, D. Zhang, X. Shen, and L. Lin, "Look into person: Self-supervised structure-sensitive learning and a new benchmark for human parsing," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 932–940, 2017.

[179] X. Liang, K. Gong, X. Shen, and L. Lin, "Look into person: Joint body parsing & pose estimation network and a new benchmark," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 4, pp. 871–885, 2018.

[180] X. Zhan, X. Pan, B. Dai, Z. Liu, D. Lin, and C. C. Loy, "Self-supervised scene de-occlusion," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 3784–3792, 2020.

[181] D. Pathak, R. Girshick, P. Dollár, T. Darrell, and B. Hariharan, "Learning features by watching objects move," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 2701–2710, 2017.

[182] Y. Wang, J. Zhang, M. Kan, S. Shan, and X. Chen, "Self-supervised equivariant attention mechanism for weakly supervised semantic segmentation," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 12275–12284, 2020.

[183] Z. Chen, X. Ye, L. Du, W. Yang, L. Huang, X. Tan, Z. Shi, F. Shen, and E. Ding, "Aggnet for self-supervised monocular depth estimation: Go an aggressive step furthe," in *ACM Int. Conf. Multimedia*, pp. 1526–1534, 2021.

[184] H. Chen, B. Lagadec, and F. Bremond, "Ice: Inter-instance contrastive encoding for unsupervised person re-identification," in *IEEE Int. Conf. Comput. Vis.*, pp. 14960–14969, 2021.

[185] T. Isobe, D. Li, L. Tian, W. Chen, Y. Shan, and S. Wang, "Towards discriminative representation learning for unsupervised person re-identification," in *IEEE Int. Conf. Comput. Vis.*, pp. 8526–8536, 2021.

[186] S. Li, X. Wang, Y. Cao, F. Xue, Z. Yan, and H. Zha, "Self-supervised deep visual odometry with online adaptation," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 6339–6348, 2020.

[187] W. Wu, Z. Y. Wang, Z. Li, W. Liu, and L. Fuxin, "Pointpwc-net: Cost volume on point clouds for (self-) supervised scene flow estimation," in *Eur. Conf. Comput. Vis.*, 2020.

[188] G. Xu, Z. Liu, X. Li, and C. C. Loy, "Knowledge distillation meets self-supervision," *arXiv preprint arXiv:2006.07114*, 2020.

[189] J. Walker, A. Gupta, and M. Hebert, "Dense optical flow prediction from a static image," in *IEEE Int. Conf. Comput. Vis.*, pp. 2443–2451, 2015.

[190] F. Zhu, Y. Zhu, X. Chang, and X. Liang, "Vision-language navigation with self-supervised auxiliary reasoning tasks," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 10012–10022, 2020.

[191] X. Niu, S. Shan, H. Han, and X. Chen, "Rhythmnet: End-to-end heart rate estimation from face via spatial-temporal representation," *IEEE Trans. Image Process.*, vol. 29, pp. 2409–2423, 2020.

[192] X. Niu, Z. Yu, H. Han, X. Li, S. Shan, and G. Zhao, "Video-based remote physiological measurement via cross-verified feature disentangling," in *Eur. Conf. Comput. Vis.*, 2020.

[193] Y. Xie, Z. Wang, and S. Ji, "Noise2same: Optimizing a self-supervised bound for image denoising," in *Neural Inf. Process. Syst.*, 2020.

[194] T. Huang, S. Li, X. Jia, H. Lu, and J. Liu, "Neighbor2neighbor: Self-supervised denoising from single noisy images," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, 2021.

[195] C. Yang, Z. Wu, B. Zhou, and S. Lin, "Instance localization for self-supervised detection pretraining," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 3987–3996, 2021.

[196] I. Croitoru, S.-V. Bogolin, and M. Leordeanu, "Unsupervised learning from video to detect foreground objects in single images," in *IEEE Int. Conf. Comput. Vis.*, pp. 4335–4343, 2017.

[197] E. Xie, J. Ding, W. Wang, X. Zhan, H. Xu, Z. Li, and P. Luo, "Detco: Unsupervised contrastive learning for object detection," *arXiv preprint arXiv:2102.04803*, 2021.

[198] G. Wu, J. Jiang, X. Liu, and J. Ma, "A practical contrastive learning framework for single image super-resolution," *arXiv preprint arXiv:2111.13924*, 2021.

[199] S. Menon, A. Damian, S. Hu, N. Ravi, and C. Rudin, "Pulse: Self-supervised photo upsampling via latent space exploration of generative models," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 2437–2445, 2020.

[200] R. Girdhar, D. F. Fouhey, M. Rodriguez, and A. Gupta, "Learning a predictable and generative vector representation for objects," in *Eur. Conf. Comput. Vis.*, pp. 484–499, 2016.

[201] D. Jayaraman and K. Grauman, "Learning image representations tied to ego-motion," in *IEEE Int. Conf. Comput. Vis.*, pp. 1413–1421, 2015.

[202] Z. Yin and J. Shi, "Geonet: Unsupervised learning of dense depth, optical flow and camera pose," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 1983–1992, 2018.

[203] L. Huang, Y. Liu, B. Wang, P. Pan, Y. Xu, and R. Jin, "Self-supervised video representation learning by context and motion decoupling," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 13886–13895, 2021.

[204] K. Hu, J. Shao, Y. Liu, B. Raj, M. Savvides, and Z. Shen, "Contrast and order representations for video self-supervised learning," in *IEEE Int. Conf. Comput. Vis.*, pp. 7939–7949, 2021.

[205] M. Tschannen, J. Djolonga, M. Ritter, A. Mahendran, N. Houlsby, S. Gelly, and M. Lucic, "Self-supervised learning of video-induced visual invariances," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 13806–13815, 2020.

[206] X. He, Y. Pan, M. Tang, Y. Lv, and Y. Peng, "Learn from unlabeled videos for near-duplicate video retrieval," in *International Conference on Research on Development in Information Retrieval*, pp. 1–10, 2022.

[207] T. Han, W. Xie, and A. Zisserman, "Video representation learning by dense predictive coding," in *ICCV Workshops*, 2019.

[208] T. Han, W. Xie, and A. Zisserman, "Memory-augmented dense predictive coding for video representation learning," in *Eur. Conf. Comput. Vis.*, 2020.

[209] B. Fernando, H. Bilen, E. Gavves, and S. Gould, "Self-supervised video representation learning with odd-one-out networks," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 3636–3645, 2017.

[210] H.-Y. Lee, J.-B. Huang, M. Singh, and M.-H. Yang, "Unsupervised representation learning by sorting sequences," in *IEEE Int. Conf. Comput. Vis.*, pp. 667–676, 2017.

[211] D. Xu, J. Xiao, Z. Zhao, J. Shao, D. Xie, and Y. Zhuang, "Self-supervised spatiotemporal learning via video clip order prediction," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 10334–10343, 2019.

[212] S. Benaim, A. Ephrat, O. Lang, I. Mosseri, W. T. Freeman, M. Rubinstein, M. Irani, and T. Dekel, "Speednet: Learning the speediness in videos," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 9922–9931, 2020.

[213] Y. Yao, C. Liu, D. Luo, Y. Zhou, and Q. Ye, "Video playback rate perception for self-supervised spatio-temporal representation learning," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 6548–6557, 2020.

[214] J. Wang, J. Jiao, and Y.-H. Liu, "Self-supervised video representation learning by pace prediction," in *Eur. Conf. Comput. Vis.*, 2020.

[215] A. Diba, V. Sharma, L. V. Gool, and R. Stiefelhagen, "Dynamonet: Dynamic action and motion network," in *IEEE Int. Conf. Comput. Vis.*, pp. 6192–6201, 2019.

[216] T. Han, W. Xie, and A. Zisserman, "Self-supervised co-training for video representation learning," in *Neural Inf. Process. Syst.*, pp. 1–12, 2020.

[217] B. Korbar, D. Tran, and L. Torresani, "Cooperative learning of audio and video models from self-supervised synchronization," in *Neural Inf. Process. Syst.*, pp. 7763–7774, 2018.

[218] R. Arandjelovic and A. Zisserman, "Look, listen and learn," in *IEEE Int. Conf. Comput. Vis.*, pp. 609–617, 2017.

[219] C. Sun, A. Myers, C. Vondrick, K. Murphy, and C. Schmid, "Videobert: A joint model for video and language representation learning," in *IEEE Int. Conf. Comput. Vis.*, pp. 7464–7473, 2019.

[220] A. Nagrani, C. Sun, D. Ross, R. Sukthankar, C. Schmid, and A. Zisserman, "Speech2action: Cross-modal supervision for action recognition," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 10317–10326, 2020.

[221] J. C. Stroud, D. A. Ross, C. Sun, J. Deng, R. Sukthankar, and C. Schmid, "Learning video representations from textual web supervision," *arXiv preprint arXiv:2007.14937*, 2020.

[222] J.-B. Alayrac, A. Recasens, R. Schneider, R. Arandjelović, J. Ramapuram, J. De Fauw, L. Smaira, S. Dieleman, and A. Zisserman,

"Self-supervised multimodal versatile networks," *arXiv preprint arXiv:2006.16228*, 2020.

[223] P. Sermanet, C. Lynch, Y. Chebotar, J. Hsu, E. Jang, S. Schaal, and S. Levine, "Time-contrastive networks: Self-supervised learning from video," in *IEEE Int. Conf. Robot. Autom.*, pp. 1134–1141, 2018.

[224] X. Wang, A. Jabri, and A. A. Efros, "Learning correspondence from the cycle-consistency of time," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 2566–2576, 2019.

[225] X. Li, S. Liu, S. De Mello, X. Wang, J. Kautz, and M.-H. Yang, "Joint-task self-supervised learning for temporal correspondence," in *Neural Inf. Process. Syst.*, pp. 318–328, 2019.

[226] A. Jabri, A. Owens, and A. A. Efros, "Space-time correspondence as a contrastive random walk," in *Neural Inf. Process. Syst.*, pp. 19545–19560, 2020.

[227] Z. Lai, E. Lu, and W. Xie, "Mast: A memory-augmented self-supervised tracker," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 6479–6488, 2020.

[228] Z. Zhang, S. Lathuiliere, E. Ricci, N. Sebe, Y. Yan, and J. Yang, "Online depth learning against forgetting in monocular videos," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 4494–4503, 2020.

[229] D. Luo, C. Liu, Y. Zhou, D. Yang, C. Ma, Q. Ye, and W. Wang, "Video cloze procedure for self-supervised spatio-temporal learning," in *AAAI Conf.Artif. Intell.*, pp. 11701–11708, 2020.

[230] O. J. Hénaff, A. Srinivas, J. De Fauw, A. Razavi, C. Doersch, S. Eslami, and A. v. d. Oord, "Data-efficient image recognition with contrastive predictive coding," in *Int. Conf. Mach. Learn.*, 2020.

[231] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever, "Improving language understanding by generative pre-training," 2018.

[232] C. Li, J. Yang, P. Zhang, M. Gao, B. Xiao, X. Dai, L. Yuan, and J. Gao, "Efficient self-supervised vision transformers for representation learning," *arXiv preprint arXiv:2106.09785*, 2021.

[233] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *Neural Inf. Process. Syst.*, pp. 3111–3119, 2013.

[234] K. Clark, M.-T. Luong, Q. V. Le, and C. D. Manning, "Electra: Pre-training text encoders as discriminators rather than generators," in *Int. Conf. Learn. Represent.*, 2020.

[235] N. Pappas and J. Henderson, "Gile: A generalized input-label embedding for text classification," *Transactions of the Association for Computational Linguistics*, vol. 7, pp. 139–155, 2019.

[236] K. Clark, M.-T. Luong, Q. V. Le, and C. D. Manning, "Pre-training transformers as energy-based cloze models," *arXiv preprint arXiv:2012.08561*, 2020.

[237] Z. Wu, S. Wang, J. Gu, M. Khabsa, F. Sun, and H. Ma, "Clear: Contrastive learning for sentence representation," *arXiv preprint arXiv:2012.15466*, 2020.

[238] J. Giorgi, O. Nitski, B. Wang, and G. Bader, "Declutr: Deep contrastive learning for unsupervised textual representations," *arXiv preprint arXiv:2006.03659*, 2020.

[239] H.-Y. Zhou, C. Lu, S. Yang, X. Han, and Y. Yu, "Preservational learning improves self-supervised medical image models by reconstructing diverse contexts," in *IEEE Int. Conf. Comput. Vis.*, pp. 3499–3509, 2021.

[240] K. Chaitanya, E. Erdil, N. Karani, and E. Konukoglu, "Contrastive learning of global and local features for medical image segmentation with limited annotations," in *Neural Inf. Process. Syst.*, 2020.

[241] J. Zhu, Y. Li, Y. Hu, K. Ma, S. K. Zhou, and Y. Zheng, "Rubik's cube+: A self-supervised feature learning framework for 3d medical image analysis," *Medical Image Analysis*, p. 101746, 2020.

[242] O. Manas, A. Lacoste, X. Giró-i Nieto, D. Vazquez, and P. Rodriguez, "Seasonal contrast: Unsupervised pre-training from uncurated remote sensing data," in *IEEE Int. Conf. Comput. Vis.*, pp. 9414–9423, 2021.

[243] D. Wang, Q. Zhang, Y. Xu, J. Zhang, B. Du, D. Tao, and L. Zhang, "Advancing plain vision transformer toward remote sensing foundation model," *IEEE Trans. Geoscience and Remote Sensing*, vol. 61, pp. 1–15, 2022.

[244] J. Liu, X. Huang, Y. Liu, and H. Li, "Mixmim: Mixed and masked image modeling for efficient visual representation learning," *arXiv preprint arXiv:2205.13137*, 2022.

[245] D. Bau, B. Zhou, A. Khosla, A. Oliva, and A. Torralba, "Network dissection: Quantifying interpretability of deep visual representations," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 6541–6549, 2017.

[246] Q. Garrido, R. Balestriero, L. Najman, and Y. Lecun, "Rankme: Assessing the downstream performance of pretrained self-supervised representations by their rank," in *Int. Conf. Mach. Learn.*, pp. 10929–10974, PMLR, July 2023.

[247] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *Int. J. Comput. Vis.*, vol. 88, pp. 303–338, 2010.

[248] T.-Y. Lin, M. Maire, S. Belongie, L. Bourdev, R. Girshick, J. Hays, P. Perona, D. Ramanan, C. L. Zitnick, and P. Dollár, "Microsoft coco: Common objects in context," 2015.

[249] B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, and A. Torralba, "Scene parsing through ade20k dataset," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017.

[250] B. Zhou, H. Zhao, X. Puig, T. Xiao, S. Fidler, A. Barriuso, and A. Torralba, "Semantic understanding of scenes through the ade20k dataset," *Int. J. Comput. Vis.*, vol. 127, no. 3, pp. 302–321, 2019.

[251] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev, *et al.*, "The kinetics human action video dataset," *arXiv preprint arXiv:1705.06950*, 2017.

[252] R. Goyal, S. Ebrahimi Kahou, V. Michalski, J. Materzynska, S. Westphal, H. Kim, V. Haenel, I. Fruend, P. Yianilos, M. Mueller-Freitag, *et al.*, "The" something something" video database for learning and evaluating visual common sense," in *IEEE Int. Conf. Comput. Vis.*, pp. 5842–5850, 2017.

[253] C. Gu, C. Sun, D. A. Ross, C. Vondrick, C. Pantofaru, Y. Li, S. Vijayanarasimhan, G. Toderici, S. Ricco, R. Sukthankar, *et al.*, "Ava: A video dataset of spatio-temporally localized atomic visual actions," in *IEEE Conf, ComputVis.Pattern Recognit.*, pp. 6047–6056, 2018.

[254] K. Soomro, A. R. Zamir, and M. Shah, "Ucf101: A dataset of 101 human actions classes from videos in the wild," *arXiv preprint arXiv:1212.0402*, 2012.

[255] H. Kuehne, H. Jhuang, E. Garrote, T. Poggio, and T. Serre, "Hmdb: a large video database for human motion recognition," in *IEEE Int. Conf. Comput. Vis.*, pp. 2556–2563, IEEE, 2011.

[256] J. Wang, Y. Gao, K. Li, J. Hu, X. Jiang, X. Guo, R. Ji, and X. Sun, "Enhancing unsupervised video representation learning by decoupling the scene and the motion," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, pp. 10129–10137, 2021.

[257] J. Knights, B. Harwood, D. Ward, A. Vanderkop, O. Mackenzie-Ross, and P. Moghadam, "Temporally coherent embeddings for self-supervised video representation learning," in *2020 25th International Conference on Pattern Recognition (ICPR)*, pp. 8914–8921, IEEE, 2021.

[258] A. Recasens, P. Luc, J.-B. Alayrac, L. Wang, F. Strub, C. Tallec, M. Malinowski, V. Pătrăucean, F. Altché, M. Valko, *et al.*, "Broaden your views for self-supervised video learning," in *IEEE Int. Conf. Comput. Vis.*, pp. 1255–1265, 2021.

[259] C. Yang, Y. Xu, B. Dai, and B. Zhou, "Video representation learning with visual tempo consistency," *arXiv preprint arXiv:2006.15489*, 2020.

[260] C. Feichtenhofer, H. Fan, B. Xiong, R. Girshick, and K. He, "A large-scale study on unsupervised spatiotemporal representation learning," in *Proceedings of the IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 3299–3309, 2021.

[261] R. Qian, T. Meng, B. Gong, M.-H. Yang, H. Wang, S. Belongie, and Y. Cui, "Spatiotemporal contrastive video representation learning," in *Proceedings of the IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 6964–6974, 2021.

[262] J. Robinson, L. Sun, K. Yu, K. Batmanghelich, S. Jegelka, and S. Sra, "Can contrastive learning avoid shortcut solutions?," in *Neural Inf. Process. Syst.*, pp. 4974–4986, 2021.

[263] Y. Wei, H. Hu, Z. Xie, Z. Zhang, Y. Cao, J. Bao, D. Chen, and B. Guo, "Contrastive learning rivals masked image modeling in fine-tuning via feature distillation," *arXiv preprint arXiv:2205.14141*, 2022.

[264] T. Chen, C. Luo, and L. Li, "Intriguing properties of contrastive losses," in *Neural Inf. Process. Syst.*, vol. 34, pp. 11834–11845, Curran Associates, Inc., 2021.

[265] Y. Tian, X. Chen, and S. Ganguli, "Understanding self-supervised learning dynamics without contrastive pairs," in *Int. Conf. Mach. Learn.*, pp. 10268–10278, 2021.

[266] Q. Garrido, Y. Chen, A. Bardes, L. Najman, and Y. LeCun, "On the

duality between contrastive and non-contrastive self-supervised learning," in *Int. Conf. Learn. Represent.*, 2023.

[267] S. Lavoie, C. Tsirigotis, M. Schwarzer, A. Vani, M. Noukhovitch, K. Kawaguchi, and A. Courville, "Simplicial embeddings in self-supervised learning and downstream classification," in *Int. Conf. Learn. Represent.*, 2023.

[268] C. Tao, H. Wang, X. Zhu, J. Dong, S. Song, G. Huang, and J. Dai, "Exploring the equivalence of siamese self-supervised learning via a unified gradient framework," in *Proceedings of the IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 14431–14440, 2022.

[269] X. Wang, R. Zhang, C. Shen, T. Kong, and L. Li, "Dense contrastive learning for self-supervised visual pre-training," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 3024–3033, 2021.

[270] W. Wang, H. Bao, L. Dong, J. Bjorck, Z. Peng, Q. Liu, K. Aggarwal, O. K. Mohammed, S. Singhal, S. Som, *et al.*, "Image as a foreign language: Beit pretraining for all vision and vision-language tasks," *arXiv preprint arXiv:2208.10442*, 2022.

[271] A. Newell and J. Deng, "How useful is self-supervised pretraining for visual tasks?," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 7345–7354, 2020.

[272] D. Chicco, "Siamese neural networks: An overview," *Artificial neural networks*, pp. 73–94, 2021.

[273] Y. Cao, Z. Xie, B. Liu, Y. Lin, Z. Zhang, and H. Hu, "Parametric instance classification for unsupervised visual feature learning," in *Neural Inf. Process. Syst.*, pp. 1–11, 2020.

[274] Z. Hou, F. Sun, Y.-K. Chen, Y. Xie, and S.-Y. Kung, "Milan: Masked image pretraining on language assisted representation," *arXiv preprint arXiv:2208.06049*, 2022.

[275] Y. Fang, W. Wang, B. Xie, Q. Sun, L. Wu, X. Wang, T. Huang, X. Wang, and Y. Cao, "Eva: Exploring the limits of masked visual representation learning at scale," in *Proceedings of the IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 19358–19369, 2023.

# APPENDIX

**Connection to principal component analysis**: Tian [81] demonstrated that CL with loss functions like InfoNCE can be formulated as a max-min problem. The max function aims to maximize the contrast between feature representations, while the min function assigns weights to pairs of examples with similar representations. In the context of deep linear networks, Tian showed that the max function in representation learning is equivalent to principal component analysis (PCA), and most local minima correspond to global minima, thus recovering optimal PCA solutions. Experimental results revealed that this formulation, when extended to include new contrastive losses beyond InfoNCE, achieves comparable or even superior performance on datasets like STL-10 and CIFAR10. Furthermore, Tian extended his theoretical analysis to 2-layer rectified linear unit (ReLU) networks, emphasizing the substantial differences between linear and nonlinear scenarios and highlighting the essential role of data augmentation during the training process. It is noteworthy that PCA aims to maximize the inter-example distances within a low-dimensional subspace, making it a specific instance of instance discrimination.

**Connection to spectral clustering**: Chen et al. [78] established a connection between CL and spectral clustering, showing that the representations obtained from CL correspond to embeddings of a positive pair graph in spectral clustering. Specifically, the authors introduced a population augmentation graph, where nodes represent augmented data from the population distribution, and the presence of an edge between nodes is determined by whether they originate from the same original example. Their key assumption is that different classes exhibit only a limited number of connections, resulting in a sparser partition for such a graph. Empirical evidence has confirmed this characteristic, illustrating the data continuity within the same class [80].

Specifically, spectral decomposition is employed on the adjacency matrix to construct a matrix, where each row denotes the representation of an example. Through a linear transformation, they demonstrated that the corresponding feature extractor could be retrieved by minimizing an unconventional contrastive loss given as

$$\mathcal{L}(f) = -2 \cdot \mathbb{E}_{x,x^+}\left[f(x)^\top f\left(x^+\right)\right] \\ + \mathbb{E}_{x,x'}\left[\left(f(x)^\top f\left(x'\right)\right)^2\right], \tag{29}$$

where $(x, x^+)$ is a pair of augmentations of the same data, $(x, x')$ is a pair of independently random augmented data, and $f$ is a parameterized function from augmented data to $\mathbb{R}^k$. It is worth noting that in cases where the dimensionality of the representation surpasses the maximum count of disjoint subgraphs, the utilization of learned representations in linear classification is guaranteed to yield minimal error.

**Connection to supervised learning**: Recent research has highlighted the remarkable efficacy of self-supervised pre-training using CL for downstream tasks involving categorization. However, its effectiveness may vary when applied to other task domains. Thus, there is a compelling need to investigate the potential of contrastive pre-training in augmenting supervised learning, particularly in terms of surpassing the accuracy achieved through traditional supervised learning.

Newell et al. [271] conducted a comprehensive investigation into the potential effects of pre-training on model performance. Their study explored three key hypotheses as follows. Firstly, whether pre-training consistently leads to performance improvements. Secondly, whether pre-training achieves higher accuracy when faced with limited labeled data, but eventually levels off at a performance comparable to the baseline when sufficient labeled data is available. Thirdly, whether pre-training converges to baseline performance before reaching its plateau in accuracy. To address these hypotheses, the authors conducted experiments on the synthetic COCO dataset with rendering, allowing for the availability of a large number of labels. The results revealed that self-supervised pre-training adheres to the assumption outlined in the third hypothesis. This suggests that SSL does not surpass supervised learning in terms of learning capability, but does perform effectively when dealing with limited labeled data.

Table 6 shows the computational complexity of different SSL methods. Note that the primary sources of time complexity and memory consumption are the neural network other than SSL components, e.g., the calculation of the cross-correlation matrix in Barlow Twins. Hence, we categorize the SSL methods into different groups based on the architecture.

TABLE 6: Summary of the computational complexity of different SSL methods. We categorize the SSL methods into different groups based on the architecture. We denote methods that require large batch training or employ multi-crop with an asterisk (*), which demands substantial memory resources. Here, **Encoder Only** refers to models with only an encoder, **Encoder & Decoder** denotes models employing a decoder with significant computational load, **Encoder & Tokenizer** indicates the use of an additional encoder for obtaining semantic tokens, and **Momentum Encoder** refers to one branch of the siamese model being an exponentially moving average version of the other. **Siamese Model** essentially refers to two identical neural networks [272]. In terms of computational load, the hierarchy is as follows: Encoder Only < Encoder & Decoder < Momentum Encoder $\approx$ Encoder & Tokenizer < Siamese Model.

| Complexity | Tower Type | Model Type | Methods |
|---|---|---|---|
| ↓ | one-tower model | Encoder Only | Jigsaw [8] Colorization [38] Rotation [7] Examplar [132] Instdisc [48] PIC [273] SimMIM [101] |
| | | Encoder & Decoder | MAE [70] |
| | dual-tower model | Encoder & Tokenizer | BEiT [99] BEiT v2 [111] MILAN [274] EVA [275] |
| | | Momentum Encoder | MoCo v1 [50] BYOL [67] DINO* [83] MoCo v2 [51] |
| | | Siamese Model [272] | SimCLR* [52] MoCo v3* [82] SwAV* [68] SimSiam [69] Barlow Twins [55] VICReg [56] data2vec [108] iBOT [98] |

**Jie Gui** (SM'16) is currently a professor at the School of Cyber Science and Engineering, Southeast University. He received a BS degree in Computer Science from Hohai University, Nanjing, China, in 2004, an MS degree in Computer Applied Technology from the Hefei Institutes of Physical Science, Chinese Academy of Sciences, Hefei, China, in 2007, and a PhD degree in Pattern Recognition and Intelligent Systems from the University of Science and Technology of China, Hefei, China, in 2010. He has published more than 60 papers in international journals and conferences such as IEEE TPAMI, IEEE TNNLS, IEEE TCYB, IEEE TIP, IEEE TCSVT, IEEE TSMCS, KDD, and ACM MM. He is the Area Chair, Senior PC Member, or PC Member of many conferences such as NeurIPS and ICML. He is an Associate Editor of IEEE Transactions on Circuits and Systems for Video Technology (T-CSVT), Artificial Intelligence Review, Neural Networks, and Neurocomputing. His research interests include machine learning, pattern recognition, and image processing.

**Zhenan Sun** (SM'18) received a B.E. degree in industrial automation from Dalian University of Technology, Dalian, China, in 1999, an M.S. degree in system engineering from Huazhong University of Science and Technology, Wuhan, China, in 2002, and a PhD degree in pattern recognition and intelligent systems from the Institute of Automation, Chinese Academy of Sciences (CASIA), Beijing, China, in 2006.

Since 2006, he has been a Faculty Member with the National Laboratory of Pattern Recognition, CASIA, and he is currently a professor with the Center for Research on Intelligent Perception and Computing. He has authored/coauthored over 200 technical papers. His current research interests include biometrics, pattern recognition, and CV.

Prof. Sun is an Associate Editor of IEEE Transactions on Biometrics, Behavior, and Identity Science. He is a member of the IEEE Computer Society and IEEE Signal Processing Society and a fellow of IAPR.

**Tuo Chen** is a PhD student with the Department of Electronic Information, Southeast University. He received his bachelor's degree from the Department of Information Security, Lanzhou University. His main research interests include Self-supervised learning, representation learning, and adversarial robustness.

**Hao Luo** received B.S. and PhD degrees from Zhejiang University, China, in 2015 and 2020, respectively. He is currently working at the Alibaba DAMO Academy. His research interests include person re-identification, vision transformer, self-supervised, computer vision, and deep learning.

**Jing Zhang** (Senior Member, IEEE) is currently a Research Fellow at the School of Computer Science, The University of Sydney. He has authored over 80 papers in prestigious conferences and journals, including CVPR, ICCV, ECCV, NeurIPS, ICLR, IEEE TPAMI, and IJCV. His research focuses on computer vision and deep learning. Additionally, he is an Area Chair for ICPR, a Senior Program Committee member for AAAI and IJCAI, and a guest editor for IEEE TBD. He regularly reviews for numerous prestigious journals and conferences.

**Dacheng Tao** (Fellow, IEEE) is currently a Distinguished University Professor in the College of Computing & Data Science at Nanyang Technological University. He mainly applies statistics and mathematics to artificial intelligence and data science, and his research is detailed in one monograph and over 200 publications in prestigious journals and proceedings at leading conferences, with best paper awards, best student paper awards, and test-of-time awards. His publications have been cited over 112K times and he has an h-index 160+ in Google Scholar. He received the 2015 and 2020 Australian Eureka Prize, the 2018 IEEE ICDM Research Contributions Award, and the 2021 IEEE Computer Society McCluskey Technical Achievement Award. He is a Fellow of the Australian Academy of Science, AAAS, ACM and IEEE.

**Qiong Cao** is a Research Scientist at JD Explore Academy. Before that, she was a Senior Researcher at Tencent. Prior to joining Tencent, she was a Postdoctoral Researcher at the Department of Engineering Science, University of Oxford. She obtained her PhD in Computer Science from the University of Exeter.