# A Cookbook of Self-Supervised Learning

Randall Balestriero[*], Mark Ibrahim[*], *Vlad Sobal*[*]*, Ari Morcos*[*]*, Shashank Shekhar*[*]*, Tom Goldstein*[†]*, Florian Bordes*[*‡]*, Adrien Bardes*[*]*, Gregoire Mialon*[*]*, Yuandong Tian*[*]*, Avi Schwarzschild*[†]*, Andrew Gordon Wilson*[**]*, Jonas Geiping*[†]*, Quentin Garrido*[*§]*, Pierre Fernandez*[**]*, Amir Bar*[*]*, Hamed Pirsiavash*[+]*, Yann LeCun*[*] and Micah Goldblum*[**]*

[*]Meta AI, FAIR
[**]New York University
[†]University of Maryland
[+]University of California, Davis
[‡]Universite de Montreal, Mila
[§]Univ Gustave Eiffel, CNRS, LIGM
[⋆]Univ. Rennes, Inria, CNRS, IRISA
[italic]Equal contributions, randomized ordering

arXiv:2304.12210v1 [cs.LG] 24 Apr 2023

# Contents

# 1 What is Self-Supervised Learning and Why Bother?

*Self-supervised learning*, dubbed "the dark matter of intelligence" [1], is a promising path to advance machine learning. As opposed to *supervised learning*, which is limited by the availability of labeled data, self-supervised approaches can learn from vast unlabeled data [Chen et al., 2020b, Misra and Maaten, 2020]. Self-supervised learning (SSL) underpins deep learning's success in natural language processing leading to advances from automated machine translation to large language models trained on web-scale corpora of unlabeled text [Brown et al., 2020, Popel et al., 2020]. In computer vision, SSL pushed new bounds on data size with models such as SEER trained on 1 billion images [Goyal et al., 2021]. SSL methods for computer vision have been able to match or in some cases surpass models trained on labeled data, even on highly competitive benchmarks like ImageNet [Tomasev et al., 2022, He et al., 2020a, Deng et al., 2009]. SSL has also been successfully applied across other modalities such as video, audio, and time series [Wickstrøm et al., 2022, Liu et al., 2022a, Schiappa et al., 2022a].

Self-supervised learning defines a pretext task based on unlabeled inputs to produce descriptive and intelligible representations [Hastie et al., 2009, Goodfellow et al., 2016]. In natural language, a common SSL objective is to mask a word in the text and predict the surrounding words. This objective of predicting the context surrounding a word encourages the model to capture relationships among words in the text without the need for any labels. The same SSL model representations can be used across a range of downstream tasks such as translating text across languages, summarizing, or even generating text, along with many others. In computer vision, analogous objectives exist with models such as MAE or BYOL learning to predict masked patches of an image or representation [Grill et al., 2020, He et al., 2022]. Other SSL objectives encourage two views of the same image, formed by say adding color or cropping, to be mapped to similar representations.

With the power to train on vast unlabeled data comes many benefits. While traditional supervised learning methods are trained on a specific task often known a priori based on the available labeled data, SSL learns generic representations useful across many tasks. SSL can be especially useful in domains such as medicine where labels are costly or the specific task can not be known a priori [Krishnan et al., 2022, Ciga et al., 2022]. There's also evidence SSL models can learn representations that are more robust to adversarial examples, label corruption, and input perturbations—and are more fair—compared to their supervised counterparts [Hendrycks et al., 2019, Goyal et al., 2022]. Consequently, SSL is a field garnering growing interest. Yet, much like cooking, training SSL methods is a delicate art with a high barrier to entry.

## 1.1 Why a Cookbook for Self-Supervised Learning?

While many components of SSL are familiar to researchers, successfully training a SSL method involves a dizzying set of choices from the pretext tasks to training hyperparameters. SSL research has a high barrier to entry due to (i) its computational cost, (ii) the absence of fully transparent papers detailing the intricate implementations required

---

[1] https://ai.facebook.com/blog/self-supervised-learning-the-dark-matter-of-intelligence/

to fully enable SSL's potential, and (iii) the absence of a unified vocabulary and theoretical view of SSL. As SSL established a distinct paradigm from traditional *reconstruction-based* unsupervised learning methods such as (denoising, variational) Autoencoders [Vincent et al., 2008, 2010, Kingma and Welling, 2013], our vocabulary for understanding SSL in a unified view is limited. In fact, attempts at unifying SSL methods under a single viewpoint have only started to emerge in the last year [HaoChen et al., 2021, Balestriero and LeCun, 2022, Shwartz-Ziv et al., 2022, Garrido et al., 2022b]. Without a common ground to characterize the different components of SSL methods, it's more challenging for researchers to start working on SSL methods. Meanwhile, SSL research is in dire need for new researchers since SSL is now deployed throughout the real-world. Yet, many open research questions remain regarding SSL's generalization guarantees, fairness properties, and robustness to adversarial attacks or even naturally occurring variations. Such questions are crucial to the reliability of SSL methods.

Furthermore, SSL—which is empirically driven—comes with many moving pieces (mostly hyper-parameters) that may impact key properties of the final representations and are not necessarily well-detailed in published work. That is, to start studying SSL methods, one must first exhaustively empirically probe those methods to fully grasp the impact and behaviors of all those moving pieces. Such empirical blind spots are strong limitations as they demand large computational resources and pre-existing hands-on experience. All in all, the co-occurrence of SOTA performances from seemingly different yet overlapping methods, little existing theoretical research, and widespread real-world deployment, make the need for a cookbook unifying the techniques and their recipes essential to lower SSL's research barrier to entry.

Our goal is to lower the barrier to entry into SSL research by laying the foundations and latest SSL recipes in the style of a cookbook. To successfully cook, you must first learn the basic techniques: chopping, sautéing, etc. We begin in Section 2 with the fundamental techniques of self-supervised learning using a common vocabulary. Specifically, we describe the families of methods along with theoretical threads to connect their objectives in a unified perspective. We highlight key concepts such as loss terms or training objectives in concept boxes. Next, a cook must learn to skillfully apply the techniques to form a delicious dish. This requires learning existing recipes, assembling ingredients, and evaluating the dish. In Section 3 we introduce the practical considerations to implementing SSL methods successfully. We discuss common training recipes including hyperparameter choices, how to assemble components such as architectures and optimizers, as well as how to evaluate SSL methods. We also share practical tips from leading researchers on common training configurations and pitfalls. We hope this cookbook serves as a practical foundation for successfully training and exploring self-supervised learning.

## 2    The Families and Origins of SSL

SSL methods have enjoyed a renaissance since 2020, thanks in large part to the availability of extremely large datasets and high-memory GPUs. However, the origins of SSL go back to the very beginning of the deep learning era.

## 2.1 Origins of SSL

Contemporary methods build upon the knowledge we gained from early experiments. In this section, we give a brief overview of the main ideas of SSL prior to 2020. While many of the specific methods have fallen out of mainstream use because they no longer provide state-of-the-art performance on benchmark problems, and they will not be discussed in great detail, the ideas from these papers form the foundation for many of the modern methods. For example, the core objective of restoring missing or distorted parts of an input or contrasting two views of the same image form the foundation for modern SSL methods. Early progress in SSL focused on the development of methods that fell into the following (sometimes overlapping) categories:

**1. Information restoration:** A wide range of methods have been developed that mask or remove something from an image, and then train a neural network to restore the missing information. Colorization-based SSL methods convert an image to grayscale, and then train a network to predict the original RGB values [Zhang et al., 2016, Larsson et al., 2016, Vondrick et al., 2018]. Because colorization requires understanding object semantics and boundaries, colorization was demonstrated as an early SSL method for object segmentation. The most straightforward application of information restoration is to mask, aka remove, a portion of an image and then train a network to inpaint the missing pixel values [Pathak et al., 2016]. This idea evolved into masked auto-encoding methods [He et al., 2022], in which the masked region is a union of image patches that can be predicted using a transformer.

**2. Using temporal relationships in video:** While the focus of this review is on image (and not video) processing, a range of specialized methods have been developed for learning single-image representations by pre-training on videos. Note that information restoration methods are particularly useful for videos, which contain multiple modalities of information that can be masked. Wang and Gupta [2015] pre-train a model using a triplet loss that promotes similarities between representations of an object in two different frames. The resulting model performed well for object detection. Pathak et al. [2017] trains a model to predict the motion of objects in a single frame, and adapts the resulting features to solve single-frame detection problems. Agrawal et al. [2015] predicts the ego-motion of a camera given multiple frames. Owens et al. [2016] propose to remove the audio track from a video, and then predict the missing sound. For specialized applications like depth mapping, self-supervised methods have been proposed that learn monocular depth models from unlabeled image pairs [Eigen et al., 2014] and later the frames from a single-camera video [Zhou et al., 2017]. Such methods remain an active area of research.

**3. Learning spatial context:** This category of methods trains a model to understand the relative positions and orientations of objects within a scene. RotNet [Gidaris et al., 2018] masks the direction of gravity by applying a random rotation and then asks the model to predict the rotation. Doersch et al. [2015] is one of the first SSL methods that simply predicts the relative location of two randomly sampled patches in an image. This strategy was superseded by "jigsaw" methods [Pathak et al., 2016, Noroozi et al., 2018] that break an image into an array of disjoint patches and predict the relative location of each. A different spatial task is learning to count [Noroozi et al., 2017]: the model is trained to output the number of objects in an image in a self-supervised way.

**4. Grouping similar images together:** One can learn rich features by grouping semantically similar images together. K-means clustering is one of the most widely used methods from classical machine learning. A number of studies have adapted k-means to perform SSL with neural models. Deep clustering alternates between assigning labels to images by performing k-means in the feature space, and updating the model to respect these assigned class labels [Caron et al., 2018]. More recent treatments of this approach use mean-shift updates to push features towards their cluster center, and have been shown to complement BYOL, a method based on two networks with the objective to predict pseudo-labels for each sample [Koohpayegani et al., 2021] (discussed in Section 2.3). Other improvements to deep clustering include using optimal transport methods in feature space to create more informative clusters [Asano et al., 2019].

**5. Generative models:** An early influential SSL method is greedy layer-wise pretraining [Bengio et al., 2006], in which layers of a deep network are trained one-at-a-time using an autoencoder loss. An analogous approach from the time used Restricted Boltzman Machines (RBMs), which could be trained layer-wise and stacked to create deep belief nets [Hinton et al., 2006]. While these methods were abandoned in favor of simpler initialization strategies and longer training runs, they were historically impactful uses of SSL, as they enabled the training of the first "deep" networks. Later advancements improved on the representation learning ability of auto-encoders, including denoising autoencoders [Vincent et al., 2008], cross-channel prediction [Zhang et al., 2017], and deep canonically correlated autoencoders [Wang et al., 2015]. Nonetheless, it was ultimately found that representation transferability is better when the auto-encoder is asked to restore a missing part of its input, resulting in the "information restoration" category of SSL methods.

Generative Adversarial Networks (GANs) [Goodfellow et al., 2014] consist of an image generator and a discriminator that differentiates real images from generated images. Both components of this model pair can be trained without supervision, and both potentially contain knowledge useful for transfer learning. Early GANs papers [Salimans et al., 2016] experimented with downstream image classification using GAN components. Specialized feature learning routines have also been developed that modify the discriminator [Springenberg, 2015], add a generator [Dai et al., 2017], or learn additional mappings from image to latent space [Donahue et al., 2017] to improve transfer learning.

**6. Multi-view invariance:** Many modern SSL methods, especially those that we focus on in this article, use contrastive learning to create feature representations that are invariant to simple transforms. The idea of contrastive learning is to encourage a model to represent two augmented versions of an input similarly. A number of methods led the charge in this direction by enforcing invariance in various ways before contrastive learning was widely adopted.

One of the most popular frameworks for learning from unlabeled data is to use a weakly trained network to apply pseudolabels to images, and then train using these labels in a standard supervised fashion [Lee et al., 2013]. This approach was later improved by enforcing invariance to transformations. Virtual adversarial training [Miyato et al., 2018] trains a network on images using their pseudolabels, and additionally performs adversarial training so that learned features are nearly invariant to small perturbations to the input image. Later works focused on maintaining invariance to data augmentation transforms.

Important early methods in this category include MixMatch [Berthelot et al., 2019], which chooses pseudolabels by averaging outputs of a network on several different random augmentations of the training images, resulting in labels that are augmentation invariant. Around the same time, it was discovered that good SSL performance could be achieved by training a network to maximize the mutual information between the representations of an image under different views [Bachman et al., 2019]. These augmentation-based methods formed a bridge between the older methods described above and the contemporary methods that are the focus of this paper.

With these origins, we now turn to categorizing SSL into four broad families: The Deep Metric Learning Family, The Self-Distillation Family, The Canonical Correlation Analysis Family, and the Masked Image Modeling Family.

## 2.2 The Deep Metric Learning Family: SimCLR/NNCLR/MeanSHIFT/SCL

The Deep Metric Learning (DML) family of methods is based on the principle of encouraging similarity between semantically transformed versions of an input. DML originated with the idea of a *contrastive loss*, which transforms this principle into a learning objective. Contrastive loss was first introduced in [Bromley et al., 1993] then more formally defined in [Chopra et al., 2005, Hadsell et al., 2006]. In DML one trains a network to predict whether two inputs are from the same class (or not) by making their embedding close (or far from each other). Since data is without labels, to identify similar inputs, we often form variants of a single input using known semantic preserving transformations. The variants of the inputs are called *positive pairs* or examples; the samples we wish to make dissimilar are called *negatives*. Often there's a margin parameter, $m$, imposing the distance between examples from different classes should be larger than $m$. Similar to the contrastive loss, the Triplet loss [Weinberger and Saul, 2009, Chechik et al., 2010, Schroff et al., 2015] shares a similar spirit, but is composed of triplets: a query, a positive example, and a negative example (see eq. (3)). Compared to contrastive loss, triplet loss only requires the difference of (dis-)similarities between positive and negative examples to the query point to be larger than a margin $m$.

The shift from DML to what is now referred to as SSL might have occurred when Sohn [2016] introduced the (N+1)-tuple loss, a loss similar to the contrastive predictive coding (CPC) loss from [Oord et al., 2018]. The use of other sample positive views as the negative view of other pairs is introduce as an efficient strategy coined *N-pair-mc loss*. Ni et al. [2021b] shows that contrastive learning is a special case of meta-learning, and existing meta-learners can be directly applied to SSL with competitive performance. CPC was extended to images in [Henaff, 2020]. A key ingredient in CPC was the introduction of the InfoNCE loss described in 3 Oh Song et al. [2016], which became central in SSL.

To summarize, the main paradigm shift between DML and Contrastive SSL arises from a few key changes, namely using data-augmentation instead of sampling to obtain the positive/negative pairs, the use of deeper networks, and the use of a predictor network, which we note in Figure 4. One of the most prominent methods coming from the paradigm shift to SSL in the deep learning family is SimCLR.

**SimCLR** learns visual representations by encouraging similarity between two augmented views of an image. In SimCLR the two views are formed by applying a combination

of transformations including random resizing, cropping, color jittering, and random blurring. After encoding each view, SimCLR uses a *projector*, often a MLP (multi-layer perceptron) followed by a ReLU (rectified linear unit) activation, to map the initial embeddings into another space where the constrative loss if applied to encourage similarity between the views. For downstream tasks, extracting the representation before the projector has been shown to improve performance. Further discussions of the role of the projector are in sections 2.6.1 and 3.2.

Another key ingredient along with the InfoNCE loss used in SimCLR is the non-parametric softmax introduced by Wu et al. [2018]. This name is motivated by removing the need to have a "parametrized" linear layer on top of the representation to compute the softmax by instead comparing representations with each others. This loss formulation already contained a *temperature parameter* in the softmax which is responsible for increasing or decreasing the sharpness of events in predictions. Other noteworthy developements include Schroff et al. [2015] use triplet loss with active triplet selection (hard positive, hard negative) either online from the current mini-batch or from a past checkpoint akin to momentum networks (discussed in section 2.3). Weinberger and Saul [2009] introduced push-pull weighting, to push negatives apart while pulling positives together, in a triplet loss to increase the margin of K-NN based models. Tian et al. [2020a] introduced the possibility of many positive views.

Aside from forming positives using semantic preserving transformations, mining positive pairs naturally arising in data is also possible. An iconic triplet loss is based on video frames where the positive pairs come from nearby frames (while negatives are from far away frames) developed in Sermanet et al. [2018] coined Time-Contrastive (TC) Time-Contrastive Learning. Nonlinear ICA [Hyvarinen and Morioka, 2016] introduced a proof that you can learn the log PDF when doing classification tasks. Alexey et al. [2015] trains a classification pretext task by transforming image patches in comparison to different transformations of image patches. One disadvantage is that this setup can involve too many classes leading performance to degrade on downstream tasks. To overcome this, NCE has been successfully employed in Mnih and Teh [2012], Mnih and Kavukcuoglu [2013] to modify the denominator in order not to loop over all classes. This is an alternative to sampling based estimation of the gradient that was found to be less stable [Bengio and Senécal, 2003, 2008]. This introduces the concept of will become momentum encoder by imposing that features maps do not vary quickly referred to as proximal algorithm [Parikh et al., 2014]. One other consideration in SSL motivated by the DML is the idea of "hard negative data mining" where the negative samples are intentionaly selected to be close to but distinct from the positives to form a more challenging learning objective. Next we describe an alternative to deep metric learning based on self-distillation.

## 2.3   The Self-Distillation Family: BYOL/SimSIAM/DINO

Self-distillation methods such as BYOL [Grill et al., 2020], SimSIAM [Chen and He, 2021], DINO [Caron et al., 2021], along with their variants rely on a simple mechanism: feeding two different views to two encoders, and mapping one to the other by means of a predictor. To prevent the encoders from *collapsing* by predicting a constant for any input, various techniques are employed. A common approach to prevent collapse is to update one of the two encoder weights with a running average of the other encoder's weights. We discuss

**Noise Contrastive Estimation: Learning Unnormalized Densities**

- introduced by Gutmann and Hyvärinen [2010] to learn unnormalized probability distributions given i.i.d observations $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_N$ from the distribution $X \sim p_X$. NCE enables approximation of $p_X$ by a parametrized function $f_\theta$ without enforcing $\int f_\theta(\boldsymbol{x}) d\boldsymbol{x} = 1$ during training

- let's first introduce a noise variable $\epsilon \sim p_\epsilon$ and let's consider the following mixture distribution

$$T \sim \mathcal{B}(s), s \in (0,1)$$

$$V \sim X 1_{\{T=1\}} + \varepsilon 1_{\{T=0\}}$$

- using Bayes rule and denoting $\eta = (1-s)/s$ we have

$$p_{T|V}(T = 1 | V = \boldsymbol{v}) = \frac{p_{V|T}(V = \boldsymbol{v} | T = 1)}{p_{V|T}(V = \boldsymbol{v} | T = 1) + \eta p_{V|T}(V = \boldsymbol{v} | T = 0)},$$

- parametrize $p_{V|T}(V = \boldsymbol{v} | T = 1) = f_\theta(\boldsymbol{v}) \exp(c)$ with $f_\theta > 0$ and learnable parameters $\{\theta, c\}$

- minimize the NLL of logistic regression (usual binary classification set-up)

$$\mathcal{L}(\theta, c) = -\mathbb{E}_{(\boldsymbol{v},t) \sim (V,T)} \log[p_{T|V}(T = t | V = \boldsymbol{v})]$$

- the minimum is attained at $f_{\theta^*} \exp(c^*) = p_X$ if $p_X(\boldsymbol{v}) = 0 \implies p_\epsilon(\boldsymbol{v}) > 0$. If $f_\theta$ is powerful enough, one can set $c = 0$ and the model will self-normalize [Mnih and Teh, 2012]

- Ceylan and Gutmann [2018] extends NCE to nonindependent noise realization i.e. $\epsilon$ depends on $X$, Ma and Collins [2018] considers conditional distribution $X|Y$, Dyer [2014] compares NCE and Negative Sampling [Mikolov et al., 2013] (the latter being a special case of the former) both extending Importance Sampling estimation [Bengio and Senécal, 2003] of the partition function (normalization factor)

Figure 1: Noise Contrastive Estimation

the particularities of each method.

**BYOL** (bootstrap your own latent) first introduced self-distillation as a means to avoid collapse. BYOL uses two networks along with a predictor to map the outputs of one network to the other. The network predicting the output is called the *online* or *student* network while the network producing the target is called the *target* or *teacher* network. Each network receives a different view of the same image formed by image transformations including random resizing, cropping, color jittering, and brightness alterations. The student network is updated throughout training using gradient descent. The teacher network is updated with an exponential moving average (EMA) updates of the weights of the online network. The slow updates induced by exponential moving average creates an asymmetry that is crucial to BYOL's success. The loss can be defined as

$$\mathcal{L}_{\text{BYOL}}(\theta_s, \gamma) = \mathbb{E}_{(\boldsymbol{x}, t_1, t_2) \sim (X, T_1, T_2)} \left[ \left\| \text{renorm}(p_\gamma(f_{\theta_s}(t_1(\boldsymbol{x})))) - \text{renorm}(f_{\theta_t}(t_2(\boldsymbol{x}))) \right\|_2^2 \right] \quad (9)$$

where the two vectors in representation space are automatically $\ell_2$-normalized i.e.

$$\text{renorm}(\boldsymbol{v}) = \frac{\boldsymbol{v}}{\max(\|\boldsymbol{v}\|_2 + \epsilon)}, \quad (10)$$

where $\epsilon$ is often set at $1^{-12}$. $f_{\theta_s}$ is the online encoder network often denoted as the *student* parametrized by $\theta_s$, and $p_\gamma$ is the predictor network parameterized by $\gamma$. $\boldsymbol{x} \sim X$ is the

## A Brief History of the infoNCE loss

In the descriptions below $z_i$ denotes the model representation of sample i, $\mathbb{P}$ denotes the set of positive samples, and $\tau$ is the temperature hyperparameter.

- Bromley et al. [1993], Chopra et al. [2005] introduces the **contrastive loss** for Deep Metric Learning

$$\mathcal{L}_{\text{cont}}(\boldsymbol{Z}) = \sum_{(i,j)\in\mathbb{P}} \|\boldsymbol{z}_j - \boldsymbol{z}_i\|_2 + \sum_{(i,j)\notin\mathbb{P}} \text{ReLU}(m - \|\boldsymbol{z}_i - \boldsymbol{z}_j\|_2)^2, m > 0, \tag{1}$$

- Goldberger et al. [2004] introduced **Neighbourhood Component Analysis** to improve maximum margin of NN-classifiers by learning a quadratic distance (Mahalanobis distance is a special case of such a distance) using

$$\mathcal{L}_{\text{NCA}}(\boldsymbol{Z}) = -\sum_{(i,j)\in\mathbb{P}} \frac{e^{-\|\boldsymbol{z}_i - \boldsymbol{z}_j\|_2^2}}{\sum_{(k,l)\in[N]^2} e^{-\|\boldsymbol{z}_k - \boldsymbol{z}_l\|_2^2}}, \tag{2}$$

- Weinberger and Saul [2009], Chechik et al. [2010] extends eq. (1) to a **triplet loss**

$$\mathcal{L}_{\text{triplet}}(\boldsymbol{Z}) = \sum_{(i,j)\in\mathbb{P}} \sum_{(k,l)\notin\mathbb{P}, k=i\}} \text{ReLU}(\|\boldsymbol{z}_i - \boldsymbol{z}_j\|_2 - \|\boldsymbol{z}_i - \boldsymbol{z}_k\| + m), m > 0, \tag{3}$$

- Sohn [2016] extends the triplet and NCA losses to form the **(N+1)-tuple loss**

$$\mathcal{L}_{\text{tuple}}(\boldsymbol{Z}) = -\sum_{(i,j)\in\mathbb{P}} \log\left(\frac{e^{\langle \boldsymbol{z}_i, \boldsymbol{z}_j\rangle}}{\sum_{(k,l)\in\mathbb{P}} e^{\langle \boldsymbol{z}_i, \boldsymbol{z}_l\rangle}}\right) + \beta\|\boldsymbol{Z}\|_F^2, \tag{4}$$

  where the denominator sum only runs through one view of the other samples, and the negative distance is replaced by the inner product <u>and</u> $\ell_2$-penalty of the feature maps $\boldsymbol{Z}$. Explicit normalization was found to be unstable yet introduced (along with a temperature parameter) in Yu and Tao [2019].

- Wu et al. [2018] introduces the **Noise-Contrastive Estimation** (NCE) loss without positive pairs

$$\mathcal{L}_{\text{NCE}} = -\sum_{n=1}^{N} \log\left(\frac{e^{\text{CoSim}(\boldsymbol{z}_i, \boldsymbol{z}_i^{(t-1)})/\tau}}{\sum_{k=1}^{N} e^{\text{CoSim}(\boldsymbol{z}_i, \boldsymbol{z}_k)/\tau}}\right) + \beta\|\boldsymbol{Z} - \boldsymbol{Z}^{(t-1)}\|_F^2, \tag{5}$$

  also coining the term non-parametric softmax. NCE loss introduces explicit normalization, a temperature parameter $\tau$, and the idea of momentum encoder (via proximal optimization method), and employs NCE to approximate the denominator when $N$ is large

- Oord et al. [2018, **CPC]** coins the name **infoNCE** by removing the proximal constraint and using positive pairs

$$\mathcal{L}_{\text{infoNCE}} = -\sum_{(i,j)\in\mathbb{P}} \log\left(\frac{e^{\text{CoSim}(\boldsymbol{z}_i, \boldsymbol{z}_j)/\tau}}{\sum_{k=1}^{N} e^{\text{CoSim}(\boldsymbol{z}_i, \boldsymbol{z}_k)/\tau}}\right), \tag{6}$$

Figure 2: History of infoNCE

input sampled from the data distribution $X$, and $t_1(\boldsymbol{x}), t_2(\boldsymbol{x})$ are two augmented views of $\boldsymbol{x}$ where $t_1 \sim T_1, t_2 \sim T_2$ are two data augmentations. The target network $f_{\theta_t}$ is of the same architecture as the student and is updated by EMA with $\xi$ controlling to what degree the target network preserves its history as in

$$\theta_t \leftarrow \xi\theta_t + (1 - \xi)\theta_s$$

with initialization $\eta = \theta_s$.

- He et al. [2020a, **MoCo**] introduces momentum encoder as an alternative to the memory bank regularization of eq. (5) and introduces a queue to store many negative samples from previous batches; [Chen et al., 2020d, **MoCoV2**] adds a projector, [Chen et al., 2021b, **MoCoV3**] adds ViTs

- Chen et al. [2020b, **SimCLR**] removes the momentum encoder and the $i^{\text{th}}$ term from the denominator coining it **NT-Xent** (Normalized Temperature-scaled cross entropy)

$$\mathcal{L}_{\text{NT}-\text{Xent}}(\boldsymbol{Z}) = -\sum_{(i,j)\in\mathbb{P}} \frac{e^{\text{CoSim}(\boldsymbol{z}_i,\boldsymbol{z}_j)}}{\sum_{k=1}^{N}\mathbf{1}_{\{k\neq i\}}e^{\text{CoSim}(\boldsymbol{z}_i,\boldsymbol{z}_k)}},$$

- Yeh et al. [2021, **DCL**] additionally removes the positive pair in the denominator

$$\mathcal{L}_{\text{DCL}}(\boldsymbol{Z}) = -\sum_{(i,j)\in\mathbb{P}} \frac{e^{\text{CoSim}(\boldsymbol{z}_i,\boldsymbol{z}_j)}}{\sum_{k=1}^{N}\mathbf{1}_{\{k\neq i\wedge(i,k)\neq\mathbb{P}\}}e^{\text{CoSim}(\boldsymbol{z}_i,\boldsymbol{z}_k)}},$$

- Dwibedi et al. [2021, **NNCLR**] uses nearest neighbors from a queue $\mathbb{Q}$

$$\mathcal{L}_{\text{NNCLR}}(\boldsymbol{Z}) = -\sum_{(i,j)\in\mathbb{P}} \frac{e^{\text{CoSim}(\text{NN}(\boldsymbol{z}_i,\mathbb{Q}),\boldsymbol{z}_j)}}{\sum_{(k,l)\in\mathbb{P}}^{N}e^{\text{CoSim}(\text{NN}(\boldsymbol{z}_i,\mathbb{Q}),\boldsymbol{z}_l)}},$$

- Mitrovic et al. [2020, **RELIC**] adds a regularization term to enforce invariance

$$\mathcal{L}_{\text{RELIC}}(\boldsymbol{Z}) = -\sum_{(i,j)\in\mathbb{P}} \frac{e^{\text{CoSim}(\boldsymbol{z}_i,\boldsymbol{z}_j)}}{\sum_{k=1}^{N}\mathbf{1}_{\{k\neq i\}}e^{\text{CoSim}(\boldsymbol{z}_i,\boldsymbol{z}_k)}} + KL(p(\boldsymbol{z}_i),p(\boldsymbol{z}_j)),$$

- Li et al. [2020, **PCL**] uses prototypes

Figure 3: Extensions of the infoNCE loss.

| **Deep Metric Learning** | | **Contrastive SSL** |
|---|---|---|
| positive/negative pairs come from labels or fixed transforms e.g. two halves of an image | $\Longrightarrow$ | positive pairs come from designed DAs that are continuously sampled, negative pairs are all non-positive pairs regardless of class membership |
| Hard-Negative Sampling for each mini-batch | $\Longrightarrow$ | random sampling |
| encoder DN | $\Longrightarrow$ | encoder DN + projector MLP |
| small dataset (N<200k) | $\Longrightarrow$ | large dataset |
| zero-shot k-NN validation | $\Longrightarrow$ | -zero-shot k-NN validation<br>-zero/few-shot/fine-tuning linear probing |

Figure 4: Deep Metric Learning versus Contrastive SSL

**SimSiam** is aimed at understanding which components in BYOL are most important. SimSiam showed that the EMA was not necessary in practice, even if it led to a small boost in performance. This enabled the use of a simplified loss defined by

$$\mathcal{L}_{\text{SimSIAM}}(\theta_{\text{s}},\gamma) = \mathbb{E}_{(\boldsymbol{x},t_1,t_2)}\left[\|\text{renorm}(p_\gamma(f_{\theta_{\text{s}}}(t_1(\boldsymbol{x})))) - \text{sg}(\text{renorm}(f_{\theta_{\text{s}}}(t_2(\boldsymbol{x}))))\|_2^2\right], \qquad (11)$$

<div style="border:1px solid black; border-radius:10px;">

**A Brief History of the Self-Distillation Family**

- Xu et al. [2004], Joulin et al. [2010, **MMC**] searches pseudo-labels so that if a classifier were train on them it would have good margin (on true labels)

- Bojanowski and Joulin [2017, **NaT**] introduces Noise as Targets i.e. $C$ real *frozen* targets $M \triangleq [m_1, \dots, m_N] \in \mathbb{R}^{D \times C}$ with *assignment constraints* of $P \triangleq [p_1, \dots, p_N] \in \{0,1\}^{C \times N}$ with

$$\mathcal{L}_{\mathrm{NaT}} = \min_{P : 1P \leq 1, P^T 1 = 1} - \sum_{n=1}^{N} CosSim(f_\theta(x_n), Mp_n), \tag{7}$$

- Caron et al. [2018, **DeepCluster**] extends NaT by allowing learning of the targets in a K-means fashion with various cluster sampling and reallocation tricks to prevent collapse

$$\mathcal{L}_{\mathrm{DeepCluster}} = CrossEntropy(f_\theta(x), \arg\min_k \|f_\theta(x) - m_k\|_2^2) + K - means(f_\theta(X), M), \tag{8}$$

- YM. et al. [2020, **SLSC**] further prevents collapse in DeepCluster through *constrained clustering membership* using Sinkhorn to infer the cluster membership probabilities

- Grill et al. [2020, **BYOL**] introduces BYOL removing the clustering step, introducing a *predictor* and projector network, defining the continuous targets as the output of a momentum network, renormalize each sample representation by its $\ell_2$-norm and leverage positive pairs. The predictor acts as a whitening operator preventing collapse [Tian et al., 2021], and momentum network can be applied only to the projector [Pham et al., 2022]

- Chen and He [2021, **SimSIAM**] replaces the BYOL moving average encoder by a stop-gradient

- Caron et al. [2021, **DINO**] introduces DINO which extends BYOL and SimSIAM to discrete representations/targets and still relies on momentum encoder

- Zhou et al. [2022a, **iBOT**] and Oquab et al. [2023, **DINOv2**] build upon DINO by combining its objective with a latent space masked-image modeling one, combining the best of both families

</div>

Figure 5: History of Self-Labeling

where for clarity we omit the distribution over which $x, t_1, t_2$ are sampled from. Several works have aimed at understanding how BYOL and SimSiam avoid collapse such as Tian et al. [2021] or Halvagal et al. [2022], where they found that the asymmetry between the two branches is the key, as well the training dynamics which regularize the variance of the embeddings implicitly.

**DINO** performs a centering of the output of the student network using a running mean (to avoid sensitivity to mini-batch size) and discretize (smoothly) the representations by means of a softmax with a temperate $\tau$ usually taken to be around $0.1$ as in

$$\mathcal{L}_{\mathrm{DINO}}(\theta_s, \gamma) = \mathbb{E}_{(x, t_1, t_2)} \left[ \mathrm{CrossEnt} \left( \mathrm{softmax}(f_{\theta_s}(t_1(x))/\tau), \mathrm{sg}(\mathrm{softmax}(\mathrm{center}(f_{\theta_t}(t_2(x)))/\tau)) \right) \right], \tag{12}$$

where akin to BYOL the teacher again has a moving average of the student network's weights, usually with value $\xi$ following a cosine schedule from $0.996$ to $1$ during training. The discretization in DINO caused by the softmax can be interepreted as an online clustering mechanism, where the last layer before the softmax contains the clustering prototypes and its weight. As such, the output of the penultimate layer is clustered using the weights of the last layer.

**iBOT** builds on DINO and combines its objective with a masked image modeling objective applied in latent space directly. Here, the target reconstruction is not the image pixels but the same patches embedded through the teacher network.

**DINOv2** further builds on iBOT and improves its performance significantly in both linear and k-NN evaluations by improving the training recipe, the architecture, and by introducing additional regularizers such as KoLeo [Sablayrolles et al., 2018]. In addition, DINOv2 curates a larger pretraining dataset consisting of 142 million images (further discussion in Section 2.7).

Many other methods belong to this self-distillation family. MoCo is another popular method based on building a dictionary look-up that was shown to in some cases to surpass supervised learning on segmentation and object detection benchmarks He et al. [2020a]. Originally the momentum encoder was introduced as a substitute for a queue in contrastive learning [He et al., 2020a], which extends the result of [Dosovitskiy et al., 2014]. MoCo's moving average uses a relatively large momentum with a default value of $\xi = 0.999$. This higher momentum value works much better than a smaller value of say $\xi = 0.9$. When SimCLR introduced the use of a projector and stronger data-augmentations, MoCoV2 [Chen et al., 2020d] followed suite with stronger data-augmentations and a projector head to boost performance. In a similar spirit, ISD [Tejankar et al., 2021] compares a query distribution to anchors from the student distribution using KL-divergence that relaxes the binary distinction between positive and negative samples. MSF [Koohpayegani et al., 2021] compares a query's nearest neighbor representation to the student target's representation and then minimize the $\ell_2$ distnace between them with renormalization (akin to cosine similarity maximization). Another approach, SSCD builds on the contrastive objective to the task of copy detection outperforming copy detection models and other contrastive methods [Pizzi et al., 2022]. Aside from the widespread use of the contrastive objective, many more methods employ similar running average updates as part of their training mechanism. For example, self-distillation [Hinton et al., 2015, Furlanello et al., 2018], Deep Q Network in reinforcement learning [Mnih et al., 2013], Mean Teacher in semi-supervised learning [Tarvainen and Valpola, 2017], and even model average in supervised and generative modeling [Jean et al., 2014].

## 2.4 The Canonical Correlation Analysis Family: VICReg/BarlowTwins/SWAV/W-MSE

The SSL canonical correlation analysis family originates with the Canonical Correlation Framework (CCA) [Hotelling, 1992]. The high-level goal of CCA is to infer the relationship between two variables by analyzing their cross-covariance matrices. Specifically, let $\boldsymbol{X} \in \mathbb{R}^D$ and $\boldsymbol{Y} \in \mathbb{R}^D$. The CCA framework seeks two transformations $\boldsymbol{U} = f_x(\boldsymbol{X})$ and

$V = f_y(Y)$ such that

$$\mathcal{L} = -\sum_{n=1}^{N} \langle U_n, V_n \rangle,$$

$$\text{such that } \underbrace{\frac{1}{N}\sum_{n=1}^{N} U_n = \frac{1}{N}\sum_{n=1}^{N} V_n = \mathbf{0}}_{\text{zero-mean representations}}, \quad \underbrace{\frac{1}{N}U^T U = \frac{1}{N}V^T V = I}_{\text{identity covariance representations}}, \qquad (13)$$

with $d$ (the dimension of the output mappings) such that $d \leq \min(\dim(X), \dim(Y))$. Linear CCA [Hotelling, 1992] considers the two mappings to be linear in which case the optimal parameters can be found through the SVD of $\Sigma_x^{-\frac{1}{2}}\Sigma_{xy}\Sigma_y^{-\frac{1}{2}}$, involving the covariance matrices of $X, Y$ and their cross-covariance. A major advance in the study nonlinear CCA was achieved by Breiman and Friedman [1985] in the univariate output setting, and by Makur et al. [2015] in the multivariate output setting, by connecting the solution to eq. (13) to the Alternating Conditional Expectation (ACE) method. Painsky et al. [2020] study the link between the optimal representation for nonlinear CCA using the Alternating Conditional Expectation proving new theoretical bounds that lead to further refinements of CCA.

These ideas were extended to deep learning in Deep Canonically Correlated Autoencoders (DCCAE) an autoencoder regularized via CCA. Hsieh [2000] and Andrew et al. [2013] introduce the objective of jointly learning parameters for two networks, $f_1, f_2$, such they their outputs are maximally correlated. The inputs to these networks are two views $X_1$ and $X_2$. Specifically the objective is then to find parameters $\theta_1, \theta_2$ for each network such that

$$(\theta_1^*, \theta_2^*) = \text{argmax}_{(\theta_1, \theta_2)} \text{corr}(f_1(X_1; \theta_1), f_2(X_2; \theta_2)). \qquad (14)$$

This DCCAE objective was extended to multivariate outputs and arbitrary DDNs in Wang et al. [2015].

From these origins, stems SSL methods such as VICReg [Bardes et al., 2021], Barlow Twins [Zbontar et al., 2021], SWAV [Caron et al., 2020], and W-MSE [Ermolov et al., 2021]. **VICReg**, the most recent among these methods, balances three objectives based on co-variance matrices of representations from two views: variance, invariance, co-variance shown in Figure 6. Regularizing the variance along each dimension of the representation prevents collapse, the invariance ensures two views are encoded similarly, and the co-variance encourages different dimensions of the representation to capture different features.

## 2.5  Masked Image Modeling

A number of prominent early self-supervised pre-training algorithms for computer vision applied degradations to training images, such as decolorization [Zhang et al., 2016], noise [Vincent et al., 2008], or shuffling image patches [Noroozi and Favaro, 2016], and taught models to undo these degradations. Context encoders instead mask out large portions of an image and replace their pixel values with white, teaching an autoencoder to inpaint the white patches [Pathak et al., 2016]. This early attempt at masked image modeling does
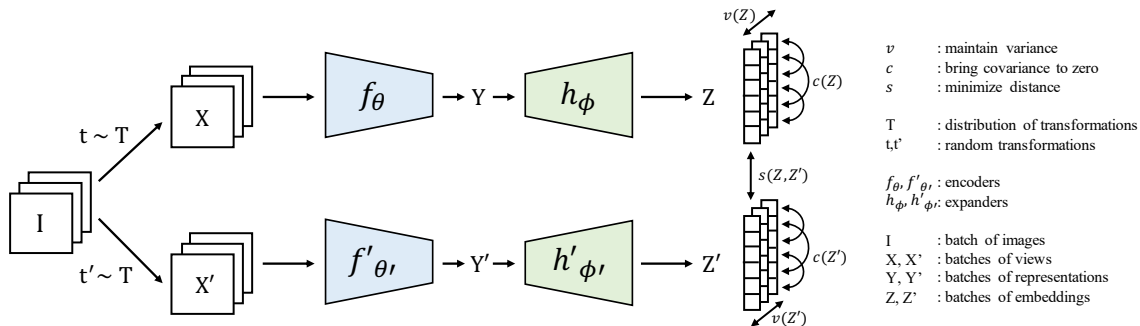
Figure 6: **VICReg**: penalizes variance, invariance, and co-variance terms to learn representations from unlabeled data.

not achieve competitive performance with supervised learning on downstream tasks, and pre-dates vision transformer architectures which modern masked training routines build upon. Subsequently, BERT [Devlin et al., 2019] shook up the natural language processing world by replacing text tokens input to a transformer language model with learnable mask tokens and teaching the model to recover the original text. This paradigm, termed *masked language modeling* (MLM), can also be interpreted as a form of the above strategy, degrading a sample via masking and teaching a model to undo the masking degradation. MLM, along with span-infilling techniques, remains popular as a SSL objective for large language models [Raffel et al., 2020, Wang et al., 2022a, Tay et al., 2022].

We can also similarly mask out portions of an image and teach a model to inpaint them. This pre-training vision strategy is known as masked image modeling (MIM). Inspired by BERT, Dosovitskiy et al. exploit the vision transformer architecture by masking out patch tokens and replacing them with learned mask tokens. They then teach their model to predict pixel values directly, but they find that this pre-training strategy is significantly less effective than supervised pre-training.

Bao et al. [2021a] note that applying the BERT strategy directly to images is difficult because whereas text tokens can only take on a small number of values that can be predicted as a classification problem, image patches can assume considerably more possible values and hence more classes than would be suitable for classification. Instead, the authors cast MIM as a regression problem, first using an autoencoder to encode image patches as discrete tokens, and then pre-training their transformer to predict the discrete token values for masked tokens. BEiT achieves significantly improved performance on downstream image classification and semantic segmentation over previous supervised and self-supervised baselines, but its training pipeline is complex since it requires a powerful autoencoder for converting image patches to discrete tokens.

In order to streamline MIM pre-training, two concurrent works [He et al., 2022, Xie et al., 2022] propose simplified algorithms, masked autoencoders (MAE) and SimMIM respectively, which directly reconstruct masked image patches rather than discrete image tokens extracted from an encoder as in BEiT. Moreover, these simplified pre-training strategies achieve superior performance to BEiT on downstream image classification, semantic segmentation, and object detection tasks. Since then, masked image modeling

has achieved competitive performance on a wide variety of vision tasks [Zhou et al., 2022a, Woo et al., 2023, Oquab et al., 2023] and even vision-language representation learning [Fang et al., 2022a]. The most successful approaches when using a frozen encoder, iBOT [Zhou et al., 2022a] and DINOV2 [Oquab et al., 2023] employ a mix of masked image modeling and more classical approaches such as self-distillation. Howver, their masked image modeling objective reconstructs in latent-space with a teacher network used to provide targets instead of using the original image as the reconstruction target.

Consider that MIM is fundamentally a generative modeling task. Such models are trained to generate missing image parts conditional on the observed ones. Note that BEiT, MAE, and SimMIM are deployed on downstream prediction problems by removing the decoder and replacing it with a prediction head. However, masked image models can also achieve strong generative modeling [Chang et al., 2022], including text-conditional generation [Chang et al., 2023]. Compared to autoregressive models for image generation [Yu et al.] which generate patches sequentially, MIM-based generative models are significantly more efficient, since they can generate patches in parallel.

In Section 3.6, we will discuss various techniques harnessed by state-of-the-art masked image modeling systems to achieve such competitive performance.
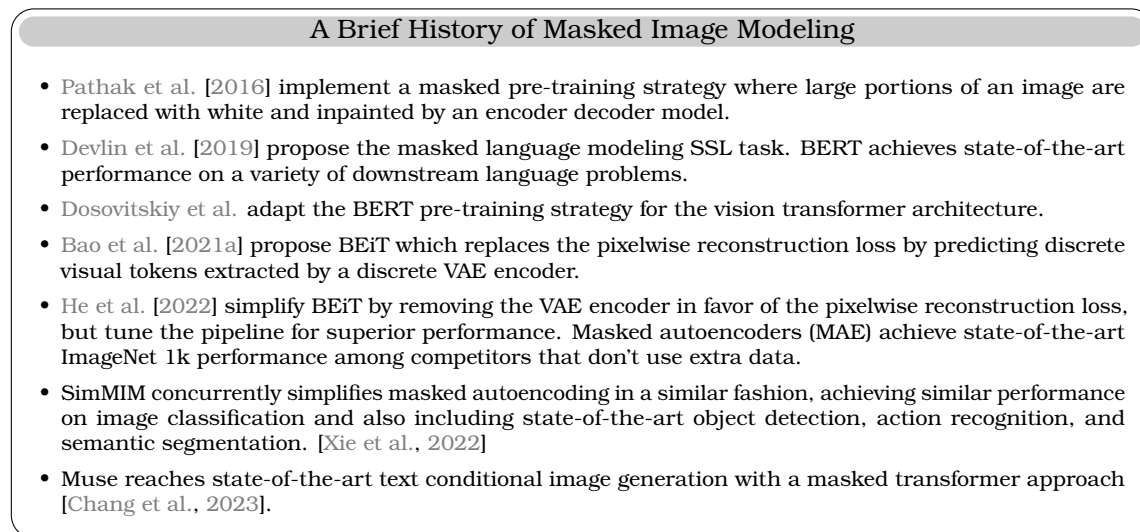
---

**A Brief History of Masked Image Modeling**

- Pathak et al. [2016] implement a masked pre-training strategy where large portions of an image are replaced with white and inpainted by an encoder decoder model.
- Devlin et al. [2019] propose the masked language modeling SSL task. BERT achieves state-of-the-art performance on a variety of downstream language problems.
- Dosovitskiy et al. adapt the BERT pre-training strategy for the vision transformer architecture.
- Bao et al. [2021a] propose BEiT which replaces the pixelwise reconstruction loss by predicting discrete visual tokens extracted by a discrete VAE encoder.
- He et al. [2022] simplify BEiT by removing the VAE encoder in favor of the pixelwise reconstruction loss, but tune the pipeline for superior performance. Masked autoencoders (MAE) achieve state-of-the-art ImageNet 1k performance among competitors that don't use extra data.
- SimMIM concurrently simplifies masked autoencoding in a similar fashion, achieving similar performance on image classification and also including state-of-the-art object detection, action recognition, and semantic segmentation. [Xie et al., 2022]
- Muse reaches state-of-the-art text conditional image generation with a masked transformer approach [Chang et al., 2023].

Figure 7: A Brief History of Masked Image Modeling

## 2.6   A Theoretical Unification Of Self-Supervised Learning

### 2.6.1   Theoretical Study of SSL

Numerous works have attempted to unify various SSL methods. In Huang et al. [2021], Barlow Twins' criterion is shown to be linked to an upper bound of a contrastive loss. This suggests a link exists between contrastive and covariance-based methods. This direction was further pursued in Garrido et al. [2022b], where a covariance-based and contrastive
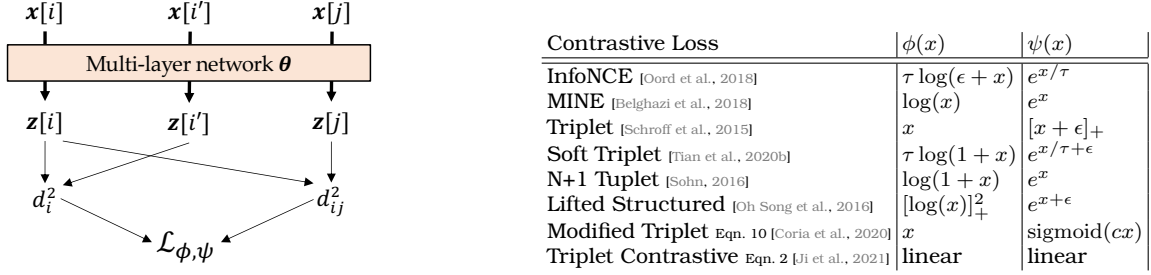
| Contrastive Loss | $\phi(x)$ | $\psi(x)$ |
|---|---|---|
| InfoNCE [Oord et al., 2018] | $\tau \log(\epsilon + x)$ | $e^{x/\tau}$ |
| MINE [Belghazi et al., 2018] | $\log(x)$ | $e^x$ |
| Triplet [Schroff et al., 2015] | $x$ | $[x + \epsilon]_+$ |
| Soft Triplet [Tian et al., 2020b] | $\tau \log(1 + x)$ | $e^{x/\tau+\epsilon}$ |
| N+1 Tuplet [Sohn, 2016] | $\log(1 + x)$ | $e^x$ |
| Lifted Structured [Oh Song et al., 2016] | $[\log(x)]_+^2$ | $e^{x+\epsilon}$ |
| Modified Triplet Eqn. 10 [Coria et al., 2020] | $x$ | $\mathrm{sigmoid}(cx)$ |
| Triplet Contrastive Eqn. 2 [Ji et al., 2021] | linear | linear |

Figure 8: Problem Setting. **Left**: Data points ($i$-th sample $\boldsymbol{x}[i]$ and its augmented version $\boldsymbol{x}[i']$, $j$-th sample $\boldsymbol{x}[j]$) are sent to networks with weights $\boldsymbol{\theta}$, to yield outputs $\boldsymbol{z}[i]$, $\boldsymbol{z}[i']$ and $\boldsymbol{z}[j]$. From the outputs $\boldsymbol{z}$, we compute pairwise squared distance $d_{ij}^2$ between $\boldsymbol{z}[i]$ and $\boldsymbol{z}[j]$ and intra-class squared distance $d_i^2$ between $\boldsymbol{z}[i]$ and $\boldsymbol{z}[i']$ for contrastive learning with a general family of contrastive loss $\mathcal{L}_{\phi,\psi}$ (Eqn. 15). **Right**: Different existing loss functions corresponds to different monotonous functions $\phi$ and $\psi$. Here $[x]_+ := \max(x, 0)$.

criterion are shown to be equivalent up to normalization by deriving the precise gap between the two approaches. These results were further validated empirically as methods were shown to exhibit similar performance and representation properties at ImageNet's scale (1.2 million samples). The similarities among methods was also studied in Tao et al. [2021] where this unification was tackled from a study of the losses' gradients.

**Contrastive Learning** Initially, InfoNCE was suggested as a variational approximation to the mutual information between two views [Aitchison and Ganev, 2023, Wang and Isola, 2020, Oord et al., 2018]. Li et al. [2021a] explains the role of InfoNCE in contrastive learning through the lens of the Hilbert-Schmidt Independence Criterion (HSIC), which was used to present a variational lower bound on the mutual information (MI) between different transformations. Tschannen et al. [2020] shows the performance of InfoNCE cannot be explained only in terms of mutual information. Instead other factors such as the feature extractor and formualtion of the mutual information estimator are important and can lead to drastically different performance [Guo et al., 2022a]. Alternative theories suggest that InfoNCE balances alignment of "positive" examples and uniformity of the overall feature representation [23], or that (under strong assumptions) it can identify the latent structure in a hypothesized data-generating process, akin to nonlinear ICA [Khemakhem et al., 2020]. In Wang and Isola [2020], Theorem 1 shows that contrastive learning with an RBF kernel (an expressive map of features into a higher dimensional space) converges to a uniform distribution on the sphere with matched pairs.

Unified contrastive losses. Tian [2022] unified contrastive losses as minimizing a general family of loss functions $\mathcal{L}_{\phi,\psi}$, where $\phi$ and $\psi$ are monotonously increasing and differentiable scalar functions

$$\min_{\boldsymbol{\theta}} \mathcal{L}_{\phi,\psi}(\boldsymbol{\theta}) = \sum_{i=1}^{N} \phi \left( \sum_{j \neq i} \psi(\|\boldsymbol{z}_i - \boldsymbol{z}_{i'}\|_2^2 - \|\boldsymbol{z}_i - \boldsymbol{z}_j\|_2^2) \right). \tag{15}$$

where $z$ are representations with indices $i$ and $j$ running from 1 to $N$. With different $\phi$ and $\psi$, Eqn. 15 covers many loss functions (Figure 8). In particular, setting $\phi(x) = \tau \log(\epsilon + x)$

and $\psi(x) = \exp(x/\tau)$ gives a generalized version of InfoNCE loss [Oord et al., 2018]:

$$\mathcal{L}_{nce} := -\tau \sum_{i=1}^{N} \log \frac{e^{-\|\boldsymbol{z}_i - \boldsymbol{z}_{i'}\|_2^2 / \tau}}{\epsilon e^{-\|\boldsymbol{z}_i - \boldsymbol{z}_{i'}\|_2^2 / \tau} + \sum_{j \neq i} e^{-\|\boldsymbol{z}_i - \boldsymbol{z}_j\|_2^2 / \tau}} \tag{16}$$

where $\epsilon > 0$ is some constant e.g. $\epsilon = 1$ has been used in He et al. [2020b], Tian et al. [2020a], $\epsilon = 0$ yields a slight variation of SimCLR [Chen et al., 2020b], the DCL loss [Yeh et al., 2021].

Hard negative sampling within a batch uses $\psi = e^{x/\tau}$. As opposed to recent SSL, negative mining has been thoroughly studied in (deep) metric learning. Recently, some works have focused on putting more weight on hard samples [Robinson et al., 2020]. Yet, Kalantidis et al. [2020], Tian [2022] showed that contrastive SSL losses i.e. $\psi = e^{x/\tau}$ already have such mechanisms at the batch level, focusing on *hard-negative pairs* without explicit "hard-negative sampling". This means that *contrastive losses need large batch sizes to ensure that hard negative samples are observed* which occurs at an additional memory cost.

**Study of the projector.** The projector network, first introduction by Chen et al. [2020b], maps the representations into another space where the loss is computed. Despite strong empricial evidence the projector improves performance, few theoretical works attempted to explain its role. Jing et al. [2022] study the role of linear projectors in contrastive learning. More precisely, it is argued that the projector prevents dimensional collapse in the representation space and that it only needs to be diagonal and low-rank to do so. Although the proposed method without a projector outperforms SimCLR with a one layer linear projector, for 2- and 3-MLP projectors, performance remains out of reach. Cosentino et al. [2022] study the interplay of the projector and data augmentations when the augmentations are Lie group transformations, and, as Mialon et al. [2022], provide an explanation on the effect of width and depth of the projector. Further empirical investigations of the role of the projector are presented in section 3.2.

### 2.6.2 Dimensional Collapse of Representations

While the goal of joint self-supervised methods is to learn meaningful representations, a significant part of the approaches suffer from what is called *dimensional collapse*. Dimensional collapse occurs when information encoded across different dimensions of the representation is redundant. In other words in the output of the projector, the embeddings are rank-deficient, which can be approximated via the singular value spectrum of the embeddings, as illustrated in Figure 12.

This phenomenon was first illustrated by Hua et al. [2021] where the use of a whitening batch normalization helped alleviate collapse. Dimensional collapse was also studied from a theoretical point of view by Jing et al. [2022] with a focus on contrastive methods. Several following works linked dimensional collapse to an impact on performance [He and Ozay, 2022, Ghosh et al., 2022, Li et al., 2022a, Garrido et al., 2022a]. Some works focused on unsupervised evaluation [Ghosh et al., 2022, Garrido et al., 2022a] where dimensional collapse was found to be a good proxy for downstream performance. Different measures of dimensional collapse have been introduced such as the entropy of
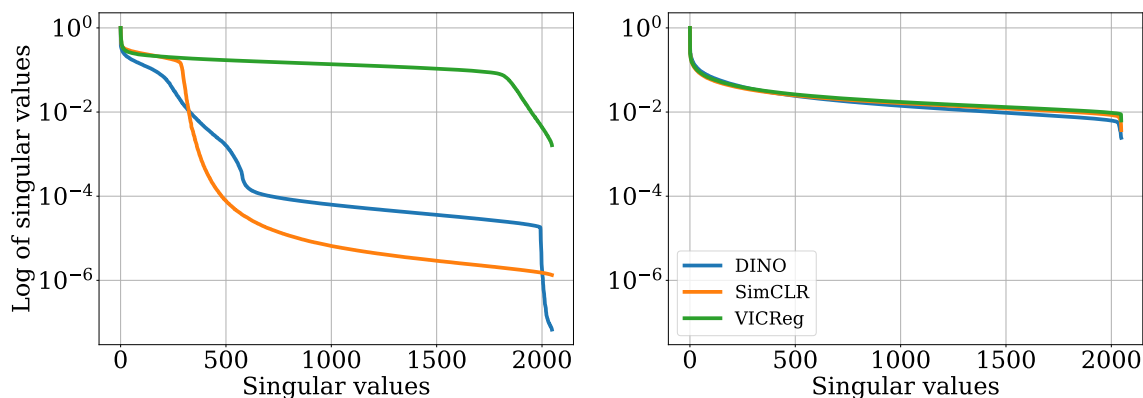
Figure 9: Illustration of dimensional collapse before the projector (Left), and after the projector (Right). Methods suffer from different levels of collapse after the projector; while no such collapse occurs for representations before the projector.

the singular value distribution [Garrido et al., 2022a], the classical rank estimator [Jing et al., 2022], fitting a power law to the singular value distribution [Ghosh et al., 2022] or the AUC of the singular value distribution [Li et al., 2022a]. Nonetheless, all of these measures focus on evaluating the rank of the representations to measure dimensional collapse in the learned representations.

## 2.7 Pretraining Data

**Curated (standard)** : The most common practice is to pretrain SSL models on curated datasets such as ImageNet and alternatives such as PASS Asano et al. [2021]. These datasets tend to generally be class-balanced and contain object-centric images, where the object is prominentely feature often in the center of the photo.

**Training with data from the wild** : Even though ImageNet has been the dataset of choice for pretraining, it is definitely not the only option. Its simplicity (object centric, single object, balanced classes) makes it a very good playground but most datasets in the wild are not as clean. If we want to leverage large uncurated datasets for SSL methods need to translate well outside of ImageNet. To this effect some works have explored pretraining on large uncurated datasets [Goyal et al., 2021], or on datasets that different from ImageNet such as COCO [El-Nouby et al., 2021], or iNaturalist [Assran et al., 2022a]. While these works have shown promising results, ImageNet (or similarly curated dataset) pretraining has remained the norm.

To provide other insights, we pretrained methods on Places205 [Zhou et al., 2014] and iNaturalist18 [Horn et al., 2018] without changing the augmentations strategy but tuning heavily loss related coefficients. The goal is to see if the setups used on ImageNet transfer well to other datasets. Places205 has the advantage of not being object centric, and iNaturalist of having a power law distribution of classes as well as requiring a lot of fine-grained information. We report our results in table 1. As we can see most methods are

| Target Dataset | iNaturalist18 | | | | Places205 | | |
|---|---|---|---|---|---|---|---|
| Method | VICReg | SimCLR | DINO | MSN | VICReg | SimCLR | DINO |
| ImageNet pretraining | 38.8 | 39.2 | 46.3 | 40.5 | 52.6 | 51.8 | 54.4 |
| Target dataset pretraining | 37.0 | 28.6 | 41.9 | 29.1 | 53.4 | 51.6 | 57.2 |

Table 1: Comparison of top-1 accuracy on a target dataset by pretraining on ImageNet or on the target dataset directly. We use the same data augmentation strategy as originally developed on ImageNet to study its transferability and heavily tune loss related hyperparameters. Methods were pretained for the same number of iterations on all datasets.

able to achieve similar performance when pretraining either on ImageNet or on the target dataset. This would suggest that the protocol developed on ImageNet can transfer decently, since we noticed that hyperparameters that were optimal on ImageNet also tended to be on different datasets. There is one visbile exception though, SimCLR and MSN perform poorly on iNaturalist18 when pretraining on it directly. While conclusions are impossible to draw precisely here, it would suggest that certain method exhibit more sensitivity on the pretraining dataset than other.

**Weakly-curated training data** : A successful approach to leverage large uncurated datasets is to perform retrieval in them based on curated data. This means that the dataset will contain images similar to a curated or smaller source dataset such as ImageNet, while being much larger and more diverse. This strategy was used in DINOv2 [Oquab et al., 2023] where LVD-142M was built using a wide variety of small and domain specific datasets. While this does not lead to big performance boosts in classification on ImageNet, it can lead to significant boosts in performance on other tasks such as image retrieval.

# 3   A Cook's Guide to Successful SSL Training and Deployment

## 3.1   Role of Data-Augmentation

Many SSL methods, especially joint embedding methods derived from Chen et al. [2020b], require a way to define positive views from a given image to learn invariances. The proxy used in these SSL methods is to leverage data augmentation to define these invariances. For example, by using different crops of a given images and positive view, the SSL model will be trained to produce a representation that is invariant to these different crops. When using a grayscale operation, or a colorjitter one as positive views, the representation will have to be invariant to the color information. Thus, the deep nature of what is learned by the SSL models is defined by the data augmentation pipeline. It is worth noting that perfect invariance is not achieved thanks to the projector [Bordes et al., 2022a], which helps improve performance on tasks which are not entirely invariant. Chen et al. [2020b] study how much influence have specific data augmentations on SimCLR with respect to the performances over ImageNet. They show that simpler data augmentation such as noise aren't beneficial on ImageNet classification downstream. Instead, cropping and

multiple color jittering operations lead to competitive results with a supervised baseline. This key element of data-augmentation had also been largely used in the following SSL works [Chen et al., 2020d, Bardes et al., 2021, Zbontar et al., 2021] without significant changes. The only variant that is sometimes used is adding smaller crops in addition of bigger crops when learning in-variances. We discuss this use of big and smaller crops, called multi-crop, in the coming subsections.

However this specific combination of data augmentation was specifically designed to reach good performances on ImageNet. Bordes et al. [2023a] study the impact of different choice of data augmentation on different downstream tasks and found that even if the addition of ColorJitter seem beneficial for many classification task it might not always be the case on other downstream tasks. Similarly, Ericsson et al. [2021a] show that different augmentations lead to learning different type of invariances for which some of them are better on some downstream tasks than other. The authors suggest to merge representations learned with different augmentations to improve transferability across a wider range of downstream task. There is also an hidden cost when using a complex pipeline of data augmentation: the data preprocessing time which might slow down significantly the training. Thus, when the training budget matter, it might be preferable to just use random crop along a grayscale operation when training a SSL model. We discuss common approaches for speeding up the training pipeline in Section 3.8.1. Ni et al. [2021b] further show that contrastive learners can benefit from very aggressive data augmentations such as large rotations when explicitly trained not to be invariant to them, as in meta-learning [Ni et al., 2021a].

Another line of work attempts to remove the need for these handcrafted data augmentations. One approach is to use a reconstruction-based objectives such as MAE [He et al., 2022] which uses a reconstruction loss in pixel space to avoid the need for defining precise invariances. Another approach is based on a joint-embedding where based on random parts of an images the goal is to predict the representations of the missing parts of the image in the representation space. An example of such method is I-JEPA [Assran et al., 2023] or Data2Vec2.0 [Baevski et al., 2022] which use a context part of an image to predict missing small parts of the image. Another line of work tries to retain style information about the augmentations to improve downstream performance on tasks requiring style information such as color by predicting style information [Xiao et al., 2020, Dangovski et al., 2021, Gidaris et al., 2018, Scherr et al., 2022]. Encoding true equivariance to augmentations (which requires a mapping between embedding) is an active line of work with approaches such as EquiMod[Dangovski et al., 2021], SEN [Park et al., 2022], or [Marchetti et al., 2022] which also aims at splitting the representations as class and pose. This idea of splitting representations as invariant and equivariant was also explored in SIE [Garrido et al., 2023] and using Lie group formalism in Ibrahim et al. [2022].

### 3.1.1 Role of multi-crop

While works such as MoCo [Meng et al., 2021] are focused on increasing the number or quality of negative pairs, another direction to improve performance is to increase the number of positives for a given image. Multi-crop, which was introduced with SwAV [Caron et al., 2020], tackles this problem by introducing smaller crops ($96 \times 96$) on top of the usual two large ones ($224 \times 224$). Instead of only comparing the two large crops together, or all

pairs of crops, the two large crops are each compared to all other crops (big or small). As such, if we have 2 large crops and $N$ small crops, the invariance loss is computed $2(N-1)$ times, increasing the positive-pair related signal. The use of smaller crops as well as not comparing all pairs of crops helps reduce the computational cost of these additional crops. While the number of additional crops can vary (10 in Mugs [Zhou et al., 2022b] compared to 6 in SwAV), it always lead to an icrease in training time and memory usage if used as is. To mitigate this cost, using $160 \times 160$ large crops and 4 $96 \times 96$ in SwAV helped mitigate the memory cost and only lead to a training time increase of $25\%$ compared to the classical setting using two crops of size $224 \times 224$, while leading to a 4 point performance boost. As such, multi-crop is a very useful strategy to help boost performance for a marginal additional compute cost. It has thus become almost ubiquitous in recent works [Caron et al., 2021, Zhou et al., 2022a,b, Bardes et al., 2022, Oquab et al., 2023]. It is worth pointing out that some works have only noticed minor increases in performance [Wang et al., 2021a] where it only lead to a 0.3 point performance increase.

Other approaches have emerged to negate the computational burden of feeding additional crops to the encoder by using nearest-neighbours in embedding space. While with NNCLR [Dwibedi et al., 2021] the matched positive crop is replaced by its nearest-neighbour in latent space, in MSF [Koohpayegani et al., 2021], a $k$-NN graph is built in embedding space to provide a similar effect as multi-crop and increase positive-pair related signal. This strategy was further employed in UniVCL [Tang et al., 2022] which used augmentation strategies such as edge of node masking in combination with a $k$-NN graph in latent space. All of these approaches show significant performance boosts for a smaller computational cost compared to multi-crop. In MSF, the use of this $k$-NN graph only increases training time by $6\%$.

## 3.2 Role of the Projector

Most SSL with joint embedding methods include a projector (usually 2- or 3-layers MLP with ReLU) after the encoder. The SSL loss is applied to the projector's output, and the projector is usually discarded after training. This crucial component was introduced in SimCLR [Chen et al., 2020b] and, although not responsible for avoiding collapse, allows significant top-1 accuracy gains on ImageNet. For example, in a 100-epochs training, the projector adds around $20\%$ of top-1 accuracy in SimCLR and VICReg (from around $50\%$ to $68\%$ and $48\%$ to $68\%$ respectively).

Bordes et al. [2022a] show that adding a projector is not only useful for SSL but is also highly beneficial in a supervised training setting when there is a misalignment between the training and downstream tasks (which was also demonstrated by Sariyildiz et al. [2022]). In fact, it's well known from Yosinski et al. [2014] that cutting layers in a trained deep neural network is beneficial when doing transfer learning mostly to avoid the training task's overfitting bias. When looking through the lens of transfer learning, it becomes easy to understand why a projector is needed in SSL since the training task is always different from the downstream task. To bridge the gap between the terms used in the SSL and in the transfer learning literature, Bordes et al. [2022a] suggested coining the method of probing intermediate representations or cutting layers as: *Guillotine Regularization* (GR). They also highlight how crucial it is to dissociate GR from the addition of a projector in

SSL because the optimal layer on which one should probe the representation might not always be the backbone (but could be an intermediate projector layer as demonstrated in Chen et al. [2020c]). Lastly, Bordes et al. [2022a] demonstrated that reducing the misalignement between the training and pretext task (by using class label to find the positives pair in contrastive learning) leads to learning a network for which the best linear probe performance on ImageNet are obtained at the last projector layer (instead of the backbone) as shown in Figure 10.



Figure 10: Figure from Bordes et al. [2022a] that show the accuracy difference between the backbone and projector representation across several downstream tasks. When using the tradition SSL positive pairs (in blue) the backbone accuracy is always much higher than the projector accuracy. However when using the class label information to define the positive pairs (in green), thus by reducing the misalignement between the training and downstream task, the projector representation lead to higher accuracy than the backbone representation on ImageNet.

| Projector | Oracle | Top-1 | Top-5 |
|-----------|--------|-------|-------|
| X | X | 50.1 | 75.8 |
| X | V | $56.4^{+6.3}$ | 80.2 |
| V | X | 68.9 | 88.2 |
| V | V | $69.5^{+0.6}$ | 88.8 |

Table 2: The projector may handle noise that originates from random data augmentations. Training VICReg without a projector can benefit from filtering semantically inconsistent augmented views using an oracle. With a projector, using an oracle provide only minor gains. Top-1 and Top-5 correspond to linear probing performance on IN-1k.

**Using a projector to handle noisy image augmentations.** The projector may also be necessary to mitigate the noise of data augmentation. As described in Section 3.1, SSL methods typically randomly augment input images to generate two different views of the same image. In some cases, enforcing invariance over two very different views might be a very strong constraint that could harm the performance, like when the content of the two views is different. To demonstrate how using the projector can mitigate that, we pretrain VICReg [Bardes et al., 2021] with and without projector using image augmentations that are semantically similar according to an "oracle", e.g a ResNet50 pretrained on ImageNet with full supervision. We pretrain for 100 epochs and include the linear probing results of these experiments in Table 2. Without projector and with an oracle, the Top1 performance is 6.3% higher compared to not using an oracle. However, equipped with a projector, using an oracle to remove noisy views only boosts Top1 performance by 0.6%. This might imply that the projector has a role in handling inconsistent or noisy augmented views during the SSL training process.

**Influence of the projector's output dimension.** Similarly to how large batch sizes were seen as a requirement for contrastive methods, a large output dimension of the projector was seen as a requirement for covariance based methods. This is illustrated by figure 4 in Zbontar et al. [2021], and table 12 in Bardes et al. [2021], where drops of up to 15% in top-1 on ImageNet can be observed. As pointed out in Garrido et al. [2022b] this was due to the projector's intermediate layers scaling with the output dimension as well as loss weights that needed to be scaled as well. By tuning these parameters, VICReg's top-1 accuracy increases from 55.9% to 65.1% with 256 dimensional embeddings. The peak performance is also achieved at 1024 dimensions and plateaus afterwards. While VICReg stays more sensitive to the output dimension of the projector than SimCLR, it is significantly more robust than originally thought and very large output dimensions are not a requirement. Comparable results should be achievable for Barlow Twins due to the similarities between the two methods.
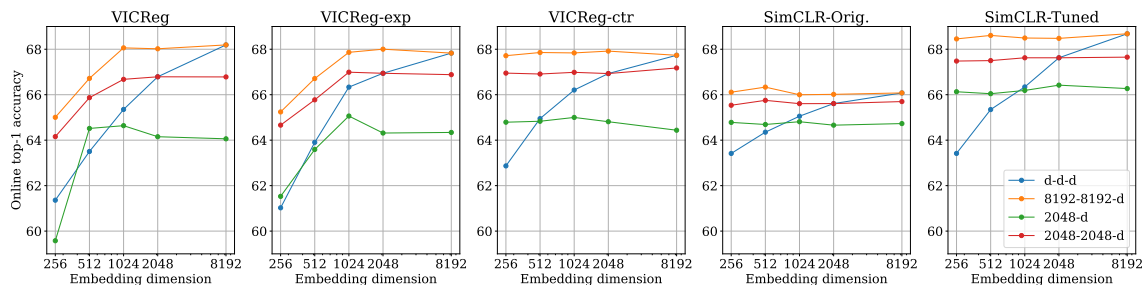


Figure 11: Impact of different projector architectures and output dimension on popular methods.$x - y - z$ denotes a MLP with layers of output dimension $x$,$y$ ad $z$ respectively. From Garrido et al. [2022b].

**Influence of the backbone's output dimension.** Recent works also investigated the effect of the backbone dimension. Dubois et al. [2022] observed that larger backbone

representations lead to better linear probe performance when using CISSL. Bordes et al. [2023b] investigated more deeply the impact of the backbone dimension across common SSL methods like VICReg, SimCLR or BYOL. They show that traditional supervised methods decline in performance when the dimension of the backbone is increased. On the other hand, SSL methods highly benefit from wider backbone representations as shown in Figure 12a. In fact, it is much more beneficial in SSL to increase the backbone size when training a ResNet than increasing the width or depth of the ResNet as illustrated in Figure 12b. This observation highlights that the current architectures used in SSL, which are often the same as the those used in supervised training, might not be optimal.
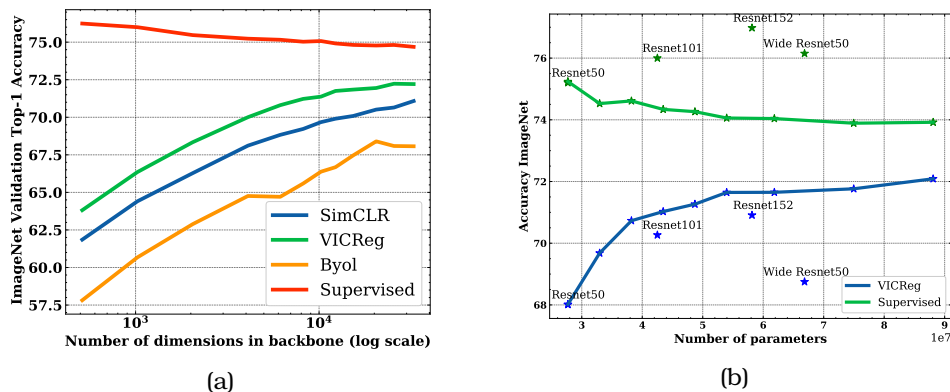


Figure 12: Figures from Bordes et al. [2023b] . a) ImageNet accuracy of various SSL methods with respect to the backbone output dimension. b) ImageNet accuracy with respect to the number of parameters. The dots on the blue and green line are models trained for different backbone output dimension.

**Properties of the representation induced by the projector.** Mialon et al. [2022] argue that the projector enforces pairwise independence of the features in the representation and provide a demonstration for random projectors in the context of VICReg, BarlowTwins and W-MSE [Bardes et al., 2021, Zbontar et al., 2021, Ermolov et al., 2021]. In particular, higher degrees of independence are reached with wider projectors. Pairwise independence, or a soft notion thereof, can be more appropriate to learn unsupervised representations from "real world" datasets such as ImageNet than mutual independence [Li et al., 2019]. Alternatively, seeking alternative SSL regularization to VCReg is needed if mutual independence is sought for. The optimization dynamics resulting from applying VCReg (the anti-collapse term in VICReg) at the projector's output is also worth noting: minimizing VCReg with respect to the projector parameters is not necessary, and VCReg is rather optimized with respect to the encoder parameters. Whether this analysis fully extends to other SSL methods is an open question.

### 3.3 The Uniform Prior in SSL or the Failure of SSL on Unbalanced Data

Despite their recent successes, there is an important limitation of SSL methods: poor performance on unbalanced datasets. Since real world data is imbalanced, such a limitation is an important factor that made the use of SSL methods on vast amount of uncurated data challenging. Assran et al. [2022a] explains such a limitation by the use of an hidden uniform prior that is common to many SSL methods. By distributing the data uniformly in the representation space, SSL methods learn to find the most discriminative features in a given mini batch. When data is uniformly distributed across classes labels, the most discriminative features that the model will learn will be class specific. However, when using imbalanced data, the most discriminative features inside the mini batch might not be the class anymore but more low level information which decrease the performances on downstream classification tasks. To alleviate this issue, Assran et al. [2022a] introduce the use of an additional regularization term on the SSL method MSN [Assran et al., 2022c] to change the distribution of the SSL clustering.

### 3.4 Teacher-Student Architecture Specific Tricks

#### 3.4.1 Role of the Moving Average Teacher

While the original BYOL method is based on exponential moving average (EMA) updates of the weights for the target (teacher) network, it was later confirmed that EMA is not necessary (i.e., the online and target networks can be identical). This is also confirmed with SimSiam [Chen and He, 2021], as long as the predictor is updated more often or has larger learning rate compared to the backbone. In the case of DQN, the target network with EMA is shown to remove bias Fan et al. [2020] and Piché et al. [2021] showed that the EMA could be removed from the target network by using the correct regularizer. For BYOL, a stop gradient of the online network, meaning the decay rate is 0 for the target network, collapses as shown in Table 5 of Grill et al. [2020]. Pham et al. [2022] shows the idea of exponential moving averages provide training stability that can even be used in non student-teacher frameworks such as SimCLR. Specifically, they show applying EMA updates to the projector of SimCLR can boost performance.

#### 3.4.2 Role of the Predictor in Self-Labeling SSL

The predictor network plays a central role in BYOL's success by predicting the representation of the teacher network from the student networks' representation. Shi et al. [2020] shows removing the predictor leads to a performance drop from 68% to 21% top-1 accuracy on ImageNet (compared to the original two-layer MLP predictor in BYOL). In Figure 1 of Shi et al. [2020], they demonstrate even a linear predictor leads to good performance and can recover from poor initialization in 10-20 epochs of training. For SimSiam, Table 1 of Chen and He [2021] shows removing the predictor in SimSiam also leads to collapse with a top-1 accuracy of < 1% on ImageNet. Is it possible to train a SSL without a projector? Yes. Tian et al. [2021], whose implementation can be found[2], proposed a contrastive

---

[2]https://github.com/facebookresearch/luckmatters/tree/main/ssl

method, DirectCLR, that does not require a trainable projector. They show regularizing the representation in DirectCLR by applying the InfoNCE SimCLR objective on sub-vectors of the representation without a trainable projector is sufficient to outperform SimCLR with a linear projector in terms of ImagNet top-1 accuracy.

## 3.5 Role of Standard Hyper-Parameters

A common issue in SSL research is that each method has different configurations of hyper-parameters. Hence comparisons directly between different SSL methods or models is often challenging. In this section, we present and describe the impact of each hyper-parameters to help SSL practitioners identify which are most important depending of their setup.

### 3.5.1 Role of Mini-Batch Size

It was originally thought that contrastive methods such as SimCLR or MoCo require large batch sizes or memory banks to work. This turns out to be misleading as both methods can be made to work at small batch sizes. A square root scaling of the learning rate was discussed in the appendix of Chen et al. [2020b] which already gave a significant increase in performance of up to 5 points in top-1 accuracy on ImageNet for a 100 epochs training. Similarly, Bordes et al. [2023a] investigated the impact of the learning rate with small batch sizes and found how one can train SimCLR on ImageNet using a single gpu without an important drop in performances. Furthermore, some works such as DCL [Yeh et al., 2021] show that you can reach top performance with a batch size of 256 or more for SimCLR, and a queue size of only 256 or more for MoCo, by simply removing the positive pair from the denominator of the softmax and with more careful hyperparameter tuning. Similarly, it was shown by Zhang et al. [2022a] that by decomposing the dictionary in MoCo and by using different temperatures for the positive and negative pairs it is possible to increase the robustness to the dictionary dimension.

### 3.5.2 Role of Learning Rate (Schedulers) and Optimizers

Here we overview typical standard settings for learning rate schedulers and optimizers across methods. To determine the learning rate, methods often scale a base learning rate based on the batch size according to the heuristic by Goyal et al. [2017]: learning rate = $\frac{\text{batch size}}{256} *$ base learning rate. For ImageNet pretraining, VICREg, Barlow Twins, BYOL, and SimCLR use a base learning rate of $0.2 - 0.3$ with the LARS optimizer [You et al., 2017]. Additionally for some methods such as Barlow twins, a much smaller learning rate (0.0048) is used to update the bias terms and batch norm parameters. Other methods such as MAE, DINO, and iBot use the AdamW optimizer [Loshchilov and Hutter, 2017] with a smaller base learning rate of $1e - 5 - 5e - 4$. For a discussion of weight decay see Section 3.5.3. The most common training schedule involves a warmup period, usually 10 epochs, where the learning rate is linearly increased to its base value. After the warmup period, most methods use cosine decay.

### 3.5.3   Role of Weight-Decay

Weight-decay is an important component of backprogagation for many SSL methods. Table 15 in BYOL [Grill et al., 2020] indicates that no weight decay may lead to unstable results. A recent blogpost[3] also mentions using weight decay leads to stable learning in BYOL. In Figure 4 of Tian et al. [2020b] the effect of weight decay is explained in terms of its effect on memory of the initial conditions. The hypothesis is that weight decay allows the online network and predictor to better model invariance to augmentations regardless of the initial condition. For further reading, Zhang et al. [2022b] provides a good review of SimSIAM collapse understanding and Shi et al. [2020] does the same for BYOL.

### 3.5.4   Vision Transformers Considerations

Training Vision Transformers (ViT) [Dosovitskiy et al.] requires special care. They are more prone to collapse and instability, and are more sensitive to the setting of hyperparameters [Touvron et al., 2021a].

**Batch size.** [Chen et al., 2021b] found that large batch (e.g., 4096) training for joint-embedding ViT SSL methods can be unstable. This instability does not reflect as a large drop in the final accuracy, but appears as drops in kNN probe accuracy during training when the $L_\infty - norm$ of the gradient spikes. Using a random (versus a learned) patch projection layer to embed pixel patches into input tokens for ViT stabilizes training for MoCo-V3, SimCLR, and BYOL and also improves the final accuracy. A learning rate warm-up period of 10k iterations [Goyal et al., 2017, Dosovitskiy et al.] also improves training stability. On the other hand, Caron et al. [2021] noted a drop in final k-NN accuracy when training with very small batch sizes (128). So, a batch size of 1024 or 2048 seems to be the sweet spot for SSL pre-training of ViTs.

While the ViT architecture does not have any BatchNorm layers, training a MoCo-V3 model with BN layers in the projector heads improved the linear probing accuracy of the ViT [Chen et al., 2021b]. Note that for joint embedding methods, batching can be done either together for all samples and crops in one batch, or separately for each batch of crops. SimCLR adopts the former, while BYOL and MoCo-V3 adopt the latter.

**Patch size.** [Caron et al., 2021] found that training with smaller patch sizes ($5 \times 5$, or $8 \times 8$ instead of $16 \times 16$) leads to improved linear probing accuracy on DINO ViT pre-training. Note that while increasing patch sizes leads to a reduction in running time, it also increases memory usage (which makes it hard to train on patches smaller than $8 \times 8$).

**Stochastic depth** [Huang et al., 2016] originated from NLP and was subsequently used in vision models [Touvron et al., 2021b] to train deeper models. It randomly drops blocks of the ViT as a regularization. The per-layer drop-rate may depend linearly on the layer depth or uniformly as suggested in recent works [Touvron et al., 2021b]. It has huge importance when training larger models (ViT-L, ViT-H, etc.). For instance Touvron et al. [2022] use $0.5$ drop path rate for ViT-H models. Conversely, when training smaller models like ViT-B, such regularization usually hurts the performance [Steiner et al., 2021].

**LayerDecay** [Clark et al., 2020] decreases the learning rate geometrically along the layers. Put differently, the last layer is not affected, while the first has very small learning rate. In SSL vision models, LayerDecay increases performance when fine-tuning on

---

[3]https://generallyintelligent.ai/blog/2020-08-24-understanding-self-supervised-contrastive-learning/

downstream tasks [Bao et al., 2021b, Zhou et al., 2022a, He et al., 2022]. Depending on the model size, the parameter is set between $0.65$ to $0.85$ – larger models usually need higher values because there are more layers. The underlying principle is that SSL builds strong model backbones, therefore we only need to fine-tune the shallowest layers.

**LayerScale** [Touvron et al., 2021b] is a per-channel multiplication of the vector produced by each residual block of the transformer. It increases the stability of the optimization and permits deeper ViT (larger than ViT-B).

`[cls] token.` When it is not explicitly needed by the method, using the average of the patch tokens instead of the class token saves memory without much change on the accuracies of the network [Zhai et al., 2022a].

## 3.6   Techniques for High Performance Masked Image Modeling

While there are several approaches to masked pretraining, the state-of-the-art systems that employ them tend to pair MIM with other techniques. For example, the ConvNextV2 architecture, which was state of the art on ImageNet (for models trained with only public data) when released, employs MAE pretraining [Woo et al., 2023]. Interestingly, the authors point out that simply pretraining a ConvNextV2 with the MAE framework is subpar. They propose adding a novel normalization layer, called "global response normalization," that proves vital to reaching state-of-the-art results [Woo et al., 2023].

In other works that claim state-of-the-art performance on image classification and semantic segmentation, MIM pretraining is paired with distillation. While some MIM routines involve reconstructing the masked portion of the input in pixel space, another option is to use a teacher network to generate target representations of the unmasked image. Zhou et al. [2022a] propose iBOT, which uses ViTs for both the teacher and the student in distillation-based MIM and outperforms prior methods on ImageNet classification. Subsequently, Liu et al. [2022b] propose dBOT, an updated distillation-based MIM approach which also achieves state-of-the-art results on image classification and semantic segmentation. A major finding in their work is that the choice of the teacher model does not have to be chosen carefully if the distillation is done in stages. This is where the teacher is updated periodically to match the student's weights and the student is reinitialized. Oquab et al. [2023] employ similar distillations to train smaller models from a ViT-g teacher with much better performance than training from scratch. This line of work highlights that pairing distillation with MIM is extremely effective.

For object detectors that utilize MIM to outperform prior work, techniques that allow MIM to work with recent and high performing pyramid ViTs like Swin are critical. Since pyramid ViTs collapse patches, random masking can leave some local windows with no information. Li et al. [2022d] propose an approach to masking that accounts for the hierarchical structure of these models called "uniform masking." They constrain the masking to hide equal amounts of information in each local window ensuring that each has some information intact. This technique helps self-supervised models (on ImageNet1K) outperform supervised models (even on ImageNet22K) on object detection benchmarks Li et al. [2022d].

## 3.7  Evaluating Your SSL Models

### 3.7.1  Evaluation with labels

Self-supervised pre-training is mainly evaluated on image classification, since it has been at the core of computer vision for decades. The three main common protocols are referred to as $k$-nearest neighbors (KNN), linear and full fine-tuning evaluations (ranked by order of complexity). They are offline evaluations, meaning that they are done independently of the self-supervised training procedure, conversely to online evaluation, which are performed during training. While online evaluation can provide a useful signal of downstream performance, because it's optimized alongside the varying self-supervised learning objective, it can be misleading. In addition, to these procedures which require labels for the downstream task, more recently, RankMe [Garrido et al., 2022a] has appeared as a viable alternative to costly evaluations, and is used as an oracle to final accuracy without having to do any training.

**KNN**  is one of the best known algorithms of machine learning and has been extensively used throughout the fields. With regards to image classification, a KNN classifier determines the label of a data point from the labels of its neighbors.

Formally speaking, the model is first used to extract frozen features $\mathcal{X} = x_1, ..., x_n$ (often $l_2$-normalized), from all the images in the training dataset. To classify a new image, we extract its feature representation $x'$, and retrieve its $k$ nearest-neighbors. They are the $k$ vectors of the training set $\mathcal{X}$ that have highest cosine-similarity with $x'$. Then, the vanilla approach applies a majority voting scheme: every neighbor counts $+1$ in its corresponding label, and we choose the label with most votes at the end. More sophisticated approaches use a weigthed voting scheme. Instead of counting $+1$ in its corresponding label, every neighbor counts a weight $w = f(x^T x')$, for instance DINO implemenation employs $w = e^{x^T x'/T}$ [Caron et al., 2021]. This allows to account for imbalanced training set, not i.i.d. features, and usually gives more accurate results, at the cost of introducing an additional hyperparameter $T$.

K-NN classifiers have the great advantage of not relying on many hyperparameters, being fast and light to deploy, without requiring any domain adaptation.

**Linear**  In the context of SSL evaluation, training a linear classifier on top of pre-trained feature representations, a.k.a. linear-probing evaluation, was introduced by Zhang et al. [2016, 2017]. It is the most popular protocol for several reasons: it achieves high-accuracy, its performance heavily rely on the quality of the representation since their discriminative power is low, it imitates how the features can be used in practice, and last but not least, it is not very computationally expensive.

Most of the time, it is done simply by appending a linear layer at the end of the frozen backbone, and optimizing its parameters for a few epochs (around $100$). Sometimes, as introduced by Bao et al. [2021b], we can benefit from the fact that the linear evaluation is lightweight and evaluate multiple linear heads at the same time, to test many hyperparameters at the same time (learning rate, averaging features or using a class-token for ViT-like architectures, number of features, etc.). A linear probe can also be trained online, by simply cutting gradient from the representations. Though only an approximation, an
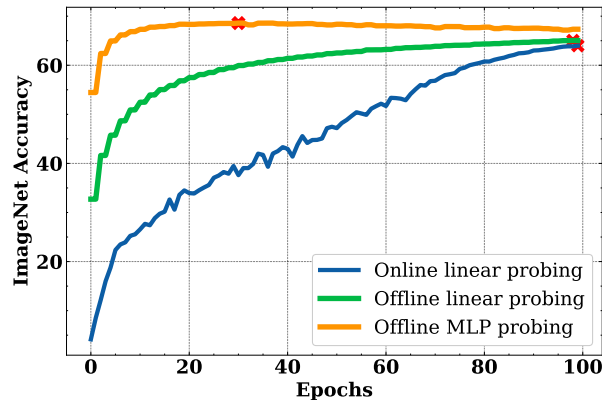
Figure 13: Figure from Bordes et al. [2023a]. Depiction of the classifier probe trained to predict the Imagenet-1k labels from the output of a Resnet50 backbone during SimCLR training (**online**) and post-training (**offline**) using a linear or MLP classifiers. The cross in red corresponds to the best accuracy. In the offline setting no data-augmentation is employed. We observe clearly that (i) when employing an MLP only a few epochs are needed and regularization or early-stopping should be employed, however, in the popular linear case, we clearly see that there is limited differences between the online and offline performances, and that over-fitting never occurs during either of the training cases.

online linear probe is extremely cheap as it reuses the computations for the SSL pretraining, and gives a good indication of downstream performance, as shown in Figure 13.

**MLP**   Instead of a simple linear probing, a multi-layer perceptron (with two or three layers) could also be used to extract which information is learned in a SSL model. Non linear evaluation is rarely present in work around SSL, however it is needed when the learned features are not linearly separable, and when it is too difficult to extract information present in features with a linear model. In fact, comparing results with a linear and a non linear probe, can give us some ideas about how well structured a representation is. Bordes et al. [2023a] present some results that compare different evaluation regime using a linear or a non linear probe. In Figure 13, one can observe that it's possible to get some gain in accuracy when using a multilater layer perceptron instead of a linear probe. However, the main issue with adding capacity into the probe is one related to overfitting: the best MLP head might not be the ones you get after 100 epochs, as showed in Figure 13.

**Full Fine-tuning**   The Masked Auto-encoders (MAE) paper [He et al., 2022] re-introduced fine-tuning as the main evaluation metrics. The main arguments are that linear-probing is uncorrelated with fine-tuning and transfer learning performances, and that small MLP heads do not evaluate the strength of the method to create strong but non-linear features. The majority of works that followed [Bao et al., 2021b, Zhou et al., 2022a, Dong et al., 2021] focused on this type of evaluation (and sometimes do not report linear/MLP

| Dataset | Method | VICReg | | SimCLR | DINO | |
|---|---|---|---|---|---|---|
| | | cov. | inv. | temp. | t-temp. | s-temp. |
| ImageNet | ImageNet Oracle | 68.2 | 68.2 | 68.5 | 72.3 | 72.4 |
| | $\alpha$-ReQ | **67.9** | 67.5 | 63.5 | 71.7 | 66.2 |
| | RankMe | 67.8 | **67.9** | **67.1** | **72.2** | **72.4** |
| OOD | ImageNet Oracle | 68.7 | 68.7 | 68.7 | 71.9 | 72.5 |
| | $\alpha$-ReQ | **68.1** | 67.8 | 65.1 | **71.8** | 68.5 |
| | RankMe | 67.7 | **68.3** | **67.6** | **71.8** | **72.5** |

Table 3: Hyperparameter selection using the common supervised linear probe strategy (ImageNet oracle), RankMe and $\alpha$-ReQ. OOD indicates the average performance over iNaturalist18, Places 205, Sun397, EuroSat, StanfordCars, CIFAR-10, CIFAR-100, Pascal VOC2007, CLEVR-cnt and FOOD101. Without any label, optimization or parameters, RankMe recovers most of the performance obtained by using ImageNet validation set, highlighting its strength as a hyperparameter selection tool. From Garrido et al. [2022a]

results). It has been shown that contrasting methods show inferior performance than masked image modeling with regards to fine-tuning, because they are less "optimization friendly" [Wei et al., 2022] - which explains the overall interest over MIM. It is by far the most computationally expensive of the evaluation methods, since it needs to re-train the whole network. The most common benchmark on ImageNet runs the optimization over $100$ epochs for ViT smaller than base, and for $50$ epochs for larger models [He et al., 2022]. Other works [Bao et al., 2021b, Peng et al., 2022, Wang et al., 2022b] first fine-tune on ImageNet-21k for $60$ epochs, and further fine-tune on ImageNet-1k, which represents between $1/5$ to $2$ times the cost of the pre-training phase.

### 3.7.2 Evaluation without labels

As we just discussed, most evaluations rely on the use of labels and training an auxiliary model. This can make evaluations expensive and sensitive to hyperparameters or their optimizations. To help alleviate these issues multiple methods have been proposed to evaluate or help tune hyperparameters of methods without relying on labels. Using a pretext-task such as rotation prediction can facilitate performance evaluation without labels, as demonstrated in Reed et al. [2021] for data augmentation policy selection. However, a drawback of this approach is the requirement for training the classifier for the pretext-task and the assumption that rotations were not part of the pretraining augmentations, or the model would be invariant to it. The eigenspectrum of representations is used in conjunction with the loss value to evaluate performance in Li et al. [2022a]. While a correlation with performance is shown, it requires training a performance classifier with the rank and loss value, making it hard to use for unsupervised evaluation. In Agrawal et al. [2022] $\alpha$-ReQ is introduced to evaluate methods by looking at the eigenspectrum decay of representations before the projector.

Another simple way to evaluate SSL methods, called RankMe, was introduced by Garrido et al. [2022a]. The idea is to use the effective rank of representations, defined as the

entropy of the singular value distribution of the embeddings. It can be computed as:

$$\text{RankMe}(\boldsymbol{Z}) = \exp\left(-\sum_{k=1}^{\min(N,K)} p_k \log p_k\right),\; p_k = \frac{\sigma_k(\boldsymbol{Z})}{\|\sigma(\boldsymbol{Z})\|_1} + \epsilon \tag{17}$$

It is shown to be a necessary condition for good performance, though you can achieve full rank representations with degenerate results (e.g. a random matrix with entries sampled i.i.d. from a Gaussian distribution). While this cannot be used to evaluate different methods, it works well for hyperparameter selection, as shown in Table 3.

### 3.7.3 Going beyond classification

While classification is a commonly used performance metric for evaluating self-supervised learning models, it is important to consider other types of vision tasks as well. Tasks such as object detection and semantic segmentation have gained popularity as they require models to learn more complex representations of visual information. Recent works Caron et al. [2021], Zhou et al. [2022a], Bardes et al. [2022] have demonstrated the effectiveness of self-supervised learning for these tasks. However, a limitation is that there is currently no standardized protocol for evaluating self-supervised models on these tasks. Various evaluation methods exist, such as finetuning the encoder on a downstream task or using the encoder as a feature extractor. Further research is needed to establish a standardized evaluation protocol for these tasks in the context of self-supervised learning.

### 3.7.4 Visual Evaluation

Another way to evaluate what information is contained or not in a representation is to use a decoder over the representation that is able to map back this information to pixel space. Some methods like [He et al., 2022] are built with a specific decoder which make such visual analysis easy, however most SSL methods aren't shipped with a decoder. To alleviate this issue and to allow researchers to visualize what can be learned by any type of SSL method, Bordes et al. [2022b] suggest training a conditional generative diffusion model using a SSL representation as conditioning. By analyzing which information remains constant across different generated samples using a specific conditioning and what information does not remain constant (because of the stochasticity in the generative model), one can get some hints about what information is contained in the representation. If a representation encodes every information about each pixel, the conditional generative model would exploit every bit of this information to perform a perfect reconstruction which will lead to no variance across different samples. If the representation encodes only the class information, the conditional generative model will only be able to use that to reconstruct the image belonging to this class, which means that when generating different samples, the object class will remain constant but the background/context/color would change across samples. In Figure 14, we show how RCDM was used by Bordes et al. [2022b] to compare the representations learned at the projector level versus the representations learned at the backbone level. In this Figure, we observe that the representations at the

Figure 14: Figure from Bordes et al. [2022b]. RCDM visualization of **what is encoded inside various representations?** First to fourth rows show our samples conditioned on the usual resnet50 backbone representation (size 2048) while fifth to eigth rows show samples conditionned on the projector/head representation of various ssl models. (Note that a separate our generative model was trained specifically for each representation). *Common/stable aspects* among a set of generated images reveal *what is encoded* in the conditioning representation. *Aspects that vary* show *what is not encoded* in the representation. We clearly see that the projector representation only keeps global information and not its context, contrary to the backbone representation. This indicates that invariances in SSL models are mostly achieved in the projector representation, not the backbone. Furthermore, it also confirms the linear classification results of Table a) which show that backbone representation are better for classifications since they contain more information about an input than the ones at the projector level.

projector level are much more invariant since the color/background information does not remain constant across different samples while this is not the case at the backbone level.

## 3.8 Speeding up Training

### 3.8.1 Distributed Training

Training self-supervised models often requires large batch sizes [Chen et al., 2020b, He et al., 2020b], or can be considerably speed up by increasing the batch size, which is ultimately limited by the memory capacity of the device the model is trained on. Distributed training divides batches across several devices that run in parallel, which increases the overall size of the batch. This is mainly done with DDP: Distributed Data Parallel or FSDP: Fully Sharded Data Parallel, available in libraries like FairScale [FairScale, 2021] or Apex [NVidia, 2021]. However some self-supervised methods rely on the statistics of the current batch for the computation of their loss value [Chen et al., 2020b, Zbontar et al., 2021, Bardes et al., 2021], which has to be taken into account when distributing the training across multiple devices. In this section, we present the elements that need to be taken into account in order to correctly distribute the training of common self-supervised learning methods. We call effective batch size, the size of the full batch distributed on the devices, and per device batch size, the size of each sub-batch on a single device.

**Synchronized batch normalization.** Batch normalization is one a the most common technique for stabilizing neural network training, as well as improving the performance of the network. It is present in most convolutional backbones used in self-supervised learning, in particular in ResNet. Batch norm uses the statistics from the current batch, which need to be aggregated for distributed training. This can be done easily in PyTorch by wrapping your distributed model the following way: `model = torch.nn.SyncBatchNorm.convert_sync_batchnorm(model)` This will replace all the BatchNorm modules in the network by a custom BatchNorm class that aggregates the statistics automatically.

**Aggregate batches for exact loss computation.** Batch norm is not the only operation that operate on batches, multiple self-supervised loss functions do as well, such as SimCLR [Chen et al., 2020b] that uses the examples in the current batch as negative example for its contrastive loss, or VICReg [Bardes et al., 2021] that computes the covariance matrix of its embeddings. In these cases the batches from each device need to be aggregated into the full batch manually. This can be done using the all_gather operation from PyTorch, however this operation does not allow back-propagation through it. We therefore implement a custom gather operation that does, the code is provided below:

We use an `all_reduce` operation on the gradient, which sums them, because DDP will divide them later by the number of devices. One can use the operation by simply calling: `FullGatherLayer.apply(x)` on the input x. Practically, for the methods above, this needs to be done on the embeddings just before the computation of the loss.

**Additional tricks.** We advise to always use the effective batch size as argument to the training script, as well as for comparing runs. The `DataLoader` class takes the per device batch size as, argument, which can be obtained by dividing the effective batch size by the number of devices which is `world_size` in PyTorch. We also advise to use an adaptive learn-

**Algorithm 1:**

```
1  class GatherLayer(torch.autograd.Function):
2      """
3      Gather tensors from all process and support backward propagation
4      for the gradients across processes.
5      """
6
7      @staticmethod
8      def forward(ctx, x):
9          output = [torch.zeros_like(x) for _ in range(dist.get_world_size())]
10         dist.all_gather(output, x)
11         return tuple(output)
12
13     @staticmethod
14     def backward(ctx, *grads):
15         all_gradients = torch.stack(grads)
16         dist.all_reduce(all_gradients)
17         return all_gradients[dist.get_rank()]
```

ing rate scaled with the effective batch_size, for example using `effective_lr = base_lr * effective_batch_size / 256` where is the `base_lr` is the argument of the training script. This reduce the learning rate search range when changing batch_size. When using small batch size, it is recommended by Chen et al. [2020b] to use `effective_lr = base_lr * ` $\sqrt{}$ `(effective_batch_size) / 256`.

### 3.8.2 Even Faster Training with FFCV and Other Speedups

Since most join-embedding SSL methods requires different set of handcrafted data augmentation, data processing can become a real bottleneck when training SSL models. Some approaches[4] have used DALI as an alternative data loader to pytorch vision while some other have relied on FFCV-SSL[5] which is based on the FFCV library [Leclerc et al., 2022]. FFCV-SSL [Bordes et al., 2023a] shows that one can train SimCLR on ImageNet in less than 2 days on a single GPU or in just a few hours using 8 GPUs (Figure 15).

### 3.8.3 Speeding Up Training of Vision Transformers

Training ViT can be made more efficient for two reasons. First, it is made easy for ViTs not to process all patches. This is especially helpful when using masked prediction pre-training objectives such as MAE [He et al., 2022] or Masked Siamese Networks [Assran et al., 2022b]. For instance, with ViT and such objectives, Data2vec 2.0 [Baevski et al., 2022] achieves 84% top-1 accuracy after pre-training for only 3 hours on 32 GPUs.

The second reason is linked to the architecture. Since transformers [Vaswani et al., 2017] are employed in almost all domains of computer science, many works aim to reduce the compute and memory requirements of the attention mechanism. One approach is with low-rank and/or sparse approximation mechanisms [Kitaev et al., 2020, Choromanski et al., 2020, Wang et al., 2020a, Chen et al., 2021a, Zaheer et al., 2020]. For instance, Li et al. [2022b] use sparse self-attention to improve efficiency in the context of SSL vision

---

[4]https://github.com/vturrisi/solo-learn
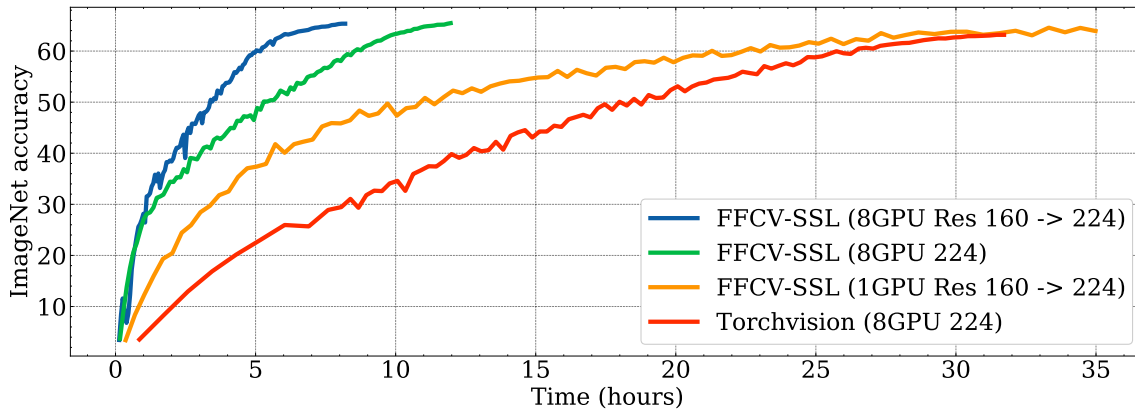[5]https://github.com/facebookresearch/FFCV-SSL

Figure 15: Figure from Bordes et al. [2023a]. ImageNet validation accuracy (y-axis) during training of SimCLR with respect to the training time (x-axis). FFCV-SSL is a library that is specifically optimized for Self-Supervised Learning, and that extends the original FFCV library [Leclerc et al., 2022]. FFCV-SSL allows a 3x time speed up with respect to torchvision and enables the training of SSL model in less than 2 days on a single gpu.

models. Another approach, is to resort to IO-aware optimizations [Ivanov et al., 2021], the most known one perhaps being FlashAttention [Dao et al., 2022].

These speed-ups are available in open-source libraries: Fairseq [Ott et al., 2019], FairScale [FairScale, 2021], XFormers [Lefaudeux et al., 2022], Apex [NVidia, 2021], etc.

Another simple way to speed up the training of vision transformers is to use Pytorch bfloat16 which allow faster training while keeping the same precision range as float32 (this is useful to avoid the usual numerical instability issues one can encounter when training vision transformers in float16).

# 4 Extending Self-Supervised Learning Beyond Images and Classification

## 4.1 Strategies for Other Data Domains

Pre-training large models with self-supervision objectives is popular not only for vision systems, but also for audio, text, and tabular data as well. The performance of existing SSL methods varies across these domains – yielding state-of-the-art language models but limited success on tabular data – which may either reflect better suitability of self-supervision or alternatively the wildly differing amount of attention which has been paid to the various domains in the SSL literature.

Applying SSL techniques to any of these data domains requires care as unique challenges arise in each domain which necessitate special considerations. For example, SSL for vision often revolves around data augmentations that may not naturally apply to speech signals. The 'positive pairs' available for contrastive learning varies from slightly different views of the same image to totally different segments of an audio recording. Nonetheless,

both contrastive and generative objectives can be applied to these other data domains. One generically useful technique across data types is masking. Whether predicting missing words in a sentence, pixels in an image, or entries of a row in a table, masking is an effective component of SSL approaches across domains.

This section is not intended as a thorough survey of self-supervision for other data modalities, as each of those fields is vast. Domain-specific surveys can be found in Liu et al. [2022a] (audio), Schiappa et al. [2022b] (video), Min et al. [2021] (text), and Rubachev et al. [2022] (tabular data). Rather, this section provides a discussion of the interesting similarities and differences in how SSL is applied to audio, text, and tabular data.

**Audio data.** Audio signals, both raw audio and mel spectrograms, have a lot in common with images. As inputs to a neural network, there are strong similarities. For example convolutions can be useful [Oord et al., 2016, Schneider et al., 2019, Baevski et al., 2021]. But as data for SSL, major differences arise. For example, horizontally flipping an image does not usually change the semantic meaning of an image (and is a wildly popular data augmentation), but for speech recordings this would completely distort the data. Similarly, while masking images is often done with random pixels, the two dimensions of a spectrogram represent time and frequency and masking with horizontal and/or vertical bands is more effective [Wang et al., 2020b]. Additionally, the existence of tones other than speech (background noise, room tone) presents a unique challenge when looking for positive pairs for contrastive learning, which is to prevent the learned representations from over fitting to the noise withing a given clip [Oord et al., 2018, Wang et al., 2020b]. In fact, the high frequency noisy artifacts, which are generally unrelated to the semantic meaning, mean reconstruction in input space is more complicated than other domains (e.g. text). Multi-modal models, on the other hand, can consider a soundbite and its text [Sermanet et al., 2018, Chung et al., 2018] or some frames of a video and the corresponding sound clip [Zhao et al., 2018, Alwassel et al., 2020] as different views to be used as positive pairs for contrastive learning.

**Video data.** Most of SSL images methods have a counter-part video SSL method. For instance, Feichtenhofer et al. [2021] have generalized SimCLR, MoCo, SwAV and BYOL to space-time video data. Indeed, in all these methods it is possible to incorporate the notion of similarity between different temporal clips of the same video. More recently, masked auto-encoding objectives for video have been built around the same idea as images, but by masking patches/ tubes of patches in the temporal axis as well [Feichtenhofer et al., 2022, Tong et al., 2022, Girdhar et al., 2022]. Besides, it is common practice to use SSL vision pre-trained models for video downstream tasks like action recognition. With ViT for instance, the patch embedding convolutional layer can be transferred from 2D to 3D by repeating the weights along the temporal axis [Feichtenhofer et al., 2022]. Vision models can then transfer to video models by using them as initialization for fine-tuning on video tasks [Fang et al., 2022a]. The frame features can also be used directly, by appending a linear layer on top of the features [Radford et al., 2019], or by using more complex heads [Ni et al., 2022, Arnab et al., 2021]. In this case, the visual system is frozen and the temporal information is learned after.

**Text data.** In contrast to audio data, text is a relatively clean input signal and representations that are useful for reconstruction do not over fit to a noisy part of the signal. In fact, the most popular large language models are all trained with reconstruction objective as opposed to contrastive objectives popular in other data domains [Radford et al., 2018,

2019, Brown et al., 2020, Devlin et al., 2018]. The Word2Vec objective [Mikolov et al., 2013] predicts a masked out portion of the training text has served as a foundational objective for self-supervised learning in natural language. While uncommon, language modeling can be done with contrastive learning for word or or character representations [Chen et al., 2022]. One other difference between text and images is that the masked token prediction for text is done over an entire dictionary. This approach is not the dominant one for images but it has been tried at the pixel level [Chen et al., 2020a]. While there are few augmentations for language data that do not change the semantic meaning, large scale systems generally use enough data and various types of masking to overcome this. Specifically, next token prediction [Radford et al., 2018, 2019, Brown et al., 2020] is akin to masking the last token in a string, while bidirectional encoders mask tokens anywhere in the string [Devlin et al., 2018] or fill larger spans of missing text [Raffel et al., 2020, Tay et al., 2022]. This choice of unidirectional next-token prediction versus bidirectional approaches leads to meaningful differences in downstream text applications [Artetxe et al., 2022]. For contrastive learning, positive pairs often come from masking and/or cropping input sequences [Meng et al., 2021, Giorgi et al., 2021]. They can also be generated using dropout so that one input has two different latent representations [Gao et al., 2021]. Additionally, some methods for both contrastive and reconstructive pretraining corrupt the input with several other augmentations including document rotation, sentence permutation, and token deletion [Lewis et al., 2020, Raffel et al., 2020, Wu et al., 2018].

**Tabular data.** Unlike text, audio, and images, classical machine learning tools are still popular for processing tabular data. However, while deep learning for tabular data is comparatively a small field, finding sensible data augmentation strategies is already a much studied topic. Several SSL methods for tabular data utilize masking in various ways and some techniques creatively employ other augmentations developed for images, like mixup [Zhang et al., 2018]. As with images and audio, some algorithms aim to generate the missing or corrupted values while others employ contrastive learning.

The masked reconstruction approaches account for a variety of masking tactics. Furthermore, it is common with tabular data to predict mask vectors as a pretext task [Yoon et al., 2020, Iida et al., 2021]. Since the predicting mask itself is part of the pretraining objective, the masked entries in the input must be filled, and typically this is done by sampling from the empirical distribution of that column or feature.

With the same augmentation, i.e. masking and sampling from the empirical marginal distribution, Bahri et al. [2021] propose pretraining with a contrastive loss. Specifically, they propose using the InfoNCE loss [Gutmann and Hyvärinen, 2010, Ceylan and Gutmann, 2018] to compare the representations of the clean and corrupted inputs.

Several other works outline ways to augment the data for a combination of generation and contrastive learning. For example, tabular data can be split into groups of columns so each sample (row) has several views available Ucar et al. [2021]. Borrowing from vision systems, a combination of CutMix [Yun et al., 2019] in input space and mixup [Zhang et al., 2018] in embedding space is also an effective augmentation for tabular data [Somepalli et al., 2021]. These methods generate augmented views that are used along with the clean input for contrastive learning. However, contrastive pretraining for both the SAINT model [Somepalli et al., 2021] and SubTab [Ucar et al., 2021] seems to work best when this is paired with a reconstructive loss term.

In their work focusing on comparing the SSL methods for tabular data, Rubachev

et al. [2022] find that pretraining objectives generally do help boost the performance of tabular models. But more specifically, they find that pretraining objectives that use the labels are best, implying that SSL for tabular data has yet to be the state of the art in its domain [Rubachev et al., 2022]. Similarly, Levin et al. [2023] show that unlike in computer vision, existing SSL pre-training routines yield less transferable features than supervised pre-training.

**Reinforcement learning.** SSL has been used to improve reinforcement learning (RL) on visual inputs. This setting is similar to video, except apart from the sequence of images, we also have access to the sequence of actions. The most common approach to apply SSL here is to use contrastive learning to train a model to match current state representation and the next time step's representation, or to match representations of the same state but with different augmentations applied. One of the earliest examples is CURL [Srinivas et al., 2020]. Recently, SSL has been used to improve sample efficiency on a challenging Atari100k benchmark [Kaiser et al., 2020]. Recent works have modified BYOL [Grill et al., 2020] or Barlow Twins [Zbontar et al., 2021] by feeding images of consecutive timesteps' observations to the two branches of the siamese network: SGI [Schwarzer et al., 2021b] and Barlow Balance [Zhang et al., 2022c] did this for offline pretraining, while SPR [Schwarzer et al., 2021a] uses it as an additional objective in the online setting. The best-performing method doing this is EfficientZero [Ye et al., 2021], which modifies MuZero [Schrittwieser et al., 2020] by, among other modifications, adding the SimSiam [Chen and He, 2020] objective to train the encoder and the forward model, and sets the new state of the art on Atari100k. Parisi et al. [2022] propose PVR, a method based on MoCo [He et al., 2020a] that improves sample efficiency on control tasks. Eysenbach et al. [2022] show that contrastive learning in RL setting is directly linked to goal-conditioned RL, and demonstrate that a method based on InfoNCE [Oord et al., 2018] achieves great performance on robotic arm control tasks.

There are a few additional challenges when applying SSL to RL. First, if the data is recorded on-line, individual observations are highly correlated with each other and are not IID (independent and identically distributed), so sampling from replay buffer should be done carefully. One failure mode of SSL objectives when applied to RL agents' data is the proclivity to latch on 'slow features' [Sobal et al., 2022]. The contrastive objective may learn for example to only look at the cloud patterns in the sky to tell apart frames in a self-driving dataset, so one must be careful to design augmentations in a way to remove useless static features in the image, or to sample data accordingly.

SSL has been used not only to improve sample efficiency, but also to improve exploration. Guo et al. [2022b] propose BYOL-Explore which uses BYOL [Grill et al., 2020] to learn the encoder and the forward model, and use the forward model disagreement as the exploration objective. The follow-up work by **?** address the problem of BYOL-Explore latching on a noisy TV. Yarats et al. [2021] proposed using a clustering method akin to SwAV [Caron et al., 2020] to do unsupervised exploration, i.e. exploration with only intrinsic rewards.

A few works have explored using vast natural videos data available to pre-train representations for RL agents. Xiao et al. [2022] introduce MVP, which uses masked-autoencoder to pre-train the transformer encoder for robotic control, while Ma et al. [2022] propose VIP, a method to learn universal features for RL using ResNet-50 backbone and the objective based the time between frames in the observations as the supervision signal. Another

method for training foundation models for RL, R3M [Nair et al., 2022], combines time-contrastive and video-language alignment objectives. VIP and R3M are trained on the large Ego4D dataset [Grauman et al., 2022], while MVP combines Imagenet, Ego4D, and additional hand manipulation data. Majumdar et al. propose VC-1, a method based on masked auto-encoding. The authors test the proposed method and other foundation models on the new test suite called CortexBench. The benchmark includes control, object manipulation, and navigation tasks, with different methods excelling at different parts of the benchmark.

There are also unsupervised methods for learning representations that are specific to RL and are not commonly used for images: e.g. Laplacian eigenmaps [Machado et al., 2017], forward-backward representations [Touati et al., 2023]. Zhang et al. [2021] propose to learn representations by making representations the same for states that lead to the same rewards, and different otherwise.

## 4.2  Incorporating Multiple Modalities into SSL Training

Self-supervised learning need not be based on a single modality. Especially multimodal vision-language have recently demonstrated this to great effect. Contrastive Language–Image Pre-training (CLIP) [Radford et al., 2021], and ALIGN [Jia et al., 2021] are self-supervised learning approaches that use image-caption pairs to learn a joint embedding space for images and captions. The objective here is contrastive, given an image and its caption are fed through separate encoder models that encode each modality into a fixed-length embedding vector. The embeddings of the training data image-caption pair are aligned, whereas other combinations in a batch are repelled.

This approach is especially interesting in comparison to contrastive SSL based on pure vision, as discussed in Section 2.6.1. The use of a second modality, here text, anchors the entire SSL training. It is no longer necessary to generate multiple augmented views to form a notion of robust representation as the joint approach learns semantically meaningful representations simply by observing similar captions re-occurring with similar images.

As a result, image encoders arising from such a joint pre-training are especially robust to visual changes that leave semantic meaning unchanged, such as sketches of objects as evaluated in ImageNet-Sketch [Wang et al., 2019, Radford et al., 2021], and are strong on out-of-domain generalization tasks. Yet, this is not always a desired representation, as visualizations in Ghiasi et al. [2022] show that these models also group features that are visually dissimilar, but semantically, or literally, alike. This can be mitigated, and overall performance, e.g. in linear probing, can even be improved by combining both image-text and image-image SSL as done in Mu et al. [2022], who combine CLIP and SimCLR [Radford et al., 2021, Chen et al., 2020b].

Recent work has pushed these vision-language systems to larger scales [Ding et al., 2021, Yuan et al., 2021, Singh et al., 2022, Wang et al., 2022c, Fang et al., 2022b], based on freely available image-caption pairs collected from the internet, such as in [Schuhmann et al., 2022]. These modern SSL models are capable of representing both vision and text, and can be used in a number of applications that are multimodal, from visual-question answering to multimodal generation [Alayrac et al., 2022, Li et al., 2022c, Nichol et al., 2022, Rao et al., 2022].

The future of vision-language pre-training, as an alternative to robust visual represen-

tations learned on vision alone, remains to be further explored. While its advantages in vision-language downstream applications are evident [Shen et al., 2022, Dou et al., 2022], shared embedding spaces can also be constructed by training solely the vision encoder first, fixing it, and then training a matching language encoder, as described in [Zhai et al., 2022b]. Ultimately, vision-language models are only the first step to self-supervised learning from multiple modalities at scale. Prototypes, such as Reed et al. [2022], train self-supervised on arbitrary input streams, ranging from vision and text to tables and agent actions, and so learn re-usable representations that are helpful for general tasks.

## 4.3 Building Feature Extractors with Localization for Dense Prediction Tasks

Aside from semantic understanding, popular computer vision tasks from object detection to segmentation to depth estimation require models which extract localized features, in other words ones which contain information indicating the locations of objects within the input image. Self-supervised learning may be particularly valuable for these dense prediction tasks since collecting segmentation masks or bounding box annotations for training images is significantly more expensive than classification labels. However, learning frameworks which are carefully tuned on image classification benchmarks may lack traits which are valuable for such dense prediction tasks. Several works, which we note perform their experiments in different settings and on different architectures and learning algorithms, express seemingly contradictory findings, namely that existing self-supervised learning strategies are or are not effective for downstream dense prediction tasks [Goyal et al., 2019, Purushwalkam and Gupta, 2020, Zhao et al., 2021, Ericsson et al., 2021b, Shwartz-Ziv et al.]. We now delve further into this discussion.

**Limitations of self-supervised learners for localization.** SSL approaches which rely on augmented views or jigsaw transformations, such as MoCo [He et al., 2020b] and PIRL [Misra and Maaten, 2020], learn occlusion invariance since they are trained with random crops on ImageNet where foreground objects are often large so that different crops contain different parts of the same object [Purushwalkam and Gupta, 2020]. On the other hand, they lack viewpoint invariance and category-instance invariance. Further, Zhao et al. [2021] argue that self-supervised learners also lack localization information because the models are able to use all parts of the image, both foreground and background, to make their predictions. The above works conduct experiments principally on convolutional architectures. It is worth noting that Ericsson et al. [2021b] suggest that the best among the popular SSL algorithms they test on are CNNs, which can still achieve competitive performance with their supervised learning counterparts in some detection and segmentation settings. Interestingly, older pretext tasks such as `jigsaw` or `colorization`, which predate the recent SSL craze sparked by MoCo and SimCLR, can also achieve competitive performance compared to supervised learning backbones when the pretext task is made "hard" enough [Goyal et al., 2019].

**CNNs or ViTs?** Recent works suggest that vision transformers (ViTs) contain superior localization information in their learned representations compared to convolutional architectures [Caron et al., 2021]. Whereas CNNs require specially designed segmentation pipelines to extract localization information from their features, this information arises naturally in the patchwise features of ViTs. Existing SSL methods designed specifically

42

for transformers confirm that the trained models are effective for downstream detection and segmentation tasks, especially when fine-tuned [Li et al., 2021b, He et al., 2022]. However, it should be noted that these SSL algorithms explicitly demand localization in their objective functions, for example via masked autoencoding where patch features should contain information regarding the contents of the corresponding section of the image [He et al., 2022]. More recently, masked autoencoding pre-training strategies have been adapted for convolutional architectures to great effect, where they achieve competitive performance on downstream object detection and instance segmentation [Woo et al., 2023]. Moreover, we will see below that a variety of pre-training strategies designed specifically for localization can be effective on transformers and convolutional networks alike.

**So how do we learn localized features without annotations?**  In order to tailor representations for downstream dense prediction tasks, numerous works propose modifying SSL routines specifically to enhance the localization in their features. Since these SSL pre-training algorithms do not use segmentation or detection annotations, they instead rely on carefully chosen unsupervised object priors.

One style of object prior enforces relationships between features extracted from locations within a single image, just as self-supervised learning procedures often enforce relationships between distinct images. One such prior uses the fact that adjacent ViT patches often contain the same objects. Unlike popular contrastive objectives which encourage augmented views of an image to produce similar features, SelfPatch encourages adjacent patches within a single image to produce similar features [Yun et al., 2022]. A related method, DenseCL [Wang et al., 2021b], matches the most similar pixel-wise features extracted from augmented samples to automatically handle the case in which augmentations move objects around in an image, and we only want to match features corresponding to the same object. More recently, VICRegL [Bardes et al., 2022] applies a similar principle by combining geometric and learned matching, with a non-contrastive criterion. Just as clustering-based methods cluster related images, Leopart [Ziegler and Asano, 2022] fine-tunes a pre-trained model to cluster patch-level features.

In addition to modifying the training loss to improve localization, we can also augment the data with this objective in mind by placing an object in multiple settings so that resulting models extract the same features from an object irrespective of its location. Instance Localization [Yang et al., 2021] leverages RoIAlign [He et al., 2017], an algorithm designed for object detectors which extracts features corresponding to a specific image patch. To this end, Instance Localization pastes a randomly chosen patch cut from the foreground of one image onto two other images and extracts features corresponding to only the pasted foreground patch, using a contrastive loss to ensure that the foreground patch generates similar features regardless of the background present and regardless of its location within an image. A competing approach estimates the location of an object within the training image using saliency maps and then cuts and pastes these objects onto background and optimizes a similar objective [Zhao et al., 2021]. Instead of using augmentations to move objects around, Purushwalkam and Gupta [2020] notes that nearby video frames contain the same object but in different positions or from different viewpoints so that contrastive learning on video data can serve much the same purpose.

Recently, UP-DETR [Dai et al., 2021] and DETReg [Bar et al., 2022] proposed an end-to-end SSL pretraining of the DETR family detectors. UP-DETR proposes to detect the bounding boxes of randomly selected patch regions in images conditioned on their pixel

values while predicting their corresponding SwAV [Caron et al., 2020] embedding. In DETReg, detection targets are obtained using the Selective Search algorithm, which does not require human annotations. Similarly, the detector predicts an associated SwAV [Caron et al., 2018] embedding for each target bounding box.

**Vision-language models for dense prediction tasks.** In Section 4.2, we saw that vision-language models extract semantically meaningful features. These features are also leveraged by recent works for open-vocabulary object detection [Kamath et al., 2021, Gu et al., 2021, Zareian et al., 2021, Minderer et al., 2022]. These works leverage vision and language backbones pre-trained as previously discussed on captioned image databases and fine-tune on object detection data. Crucially, pre-trained language models, paired with image feature extractors, allow open-vocabulary object detectors to detect new objects never seen during their fine-tuning stage simply by querying the language model with an appropriate prompt.

# 5   Conclusion

*Self-supervised learning* (SSL) established a new paradigm for advancing machine intelligence. Despite many successes, SSL remains a daunting field with a dizzying array of methods each with intricate implementations. Due to the fast moving research and the breadth of SSL methods, it remains a challenge to navigate the field. This becomes an issue for researchers and practitioners who joined the field only recently, in turn creating a high barrier to entry for SSL research and deployment. We hope our cookbook will help lower these barriers by enabling the curious researcher of any background to navigate the terrain of methods, understand the role of the various knobs, and gain the know-how required to be successful with SSL.

# References

K. K. Agrawal, A. K. Mondal, A. Ghosh, and B. A. Richards. $\alpha$-req : Assessing representation quality in self-supervised learning by measuring eigenspectrum decay. In A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=ii9X4vtZGTZ. 32

P. Agrawal, J. Carreira, and J. Malik. Learning to see by moving. In *Proceedings of the IEEE international conference on computer vision*, pages 37–45, 2015. 5

L. Aitchison and S. Ganev. Infonce is a variational autoencoder, 2023. 17

J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican, M. Reynolds, R. Ring, E. Rutherford, S. Cabi, T. Han, Z. Gong, S. Samangooei, M. Monteiro, J. Menick, S. Borgeaud, A. Brock, A. Nematzadeh, S. Sharifzadeh, M. Binkowski, R. Barreira, O. Vinyals, A. Zisserman, and K. Simonyan. Flamingo: A Visual Language Model for Few-Shot Learning. *arxiv:2204.14198[cs]*, Nov. 2022. doi: 10.48550/arXiv.2204.14198. URL http://arxiv.org/abs/2204.14198. 41

D. Alexey, P. Fischer, J. Tobias, M. R. Springenberg, and T. Brox. Discriminative unsupervised feature learning with exemplar convolutional neural networks. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 99, 2015. 8

H. Alwassel, D. Mahajan, B. Korbar, L. Torresani, B. Ghanem, and D. Tran. Self-supervised learning by cross-modal audio-video clustering. *Advances in Neural Information Processing Systems*, 33:9758–9770, 2020. 38

G. Andrew, R. Arora, J. Bilmes, and K. Livescu. Deep canonical correlation analysis. In *International conference on machine learning*, pages 1247–1255. PMLR, 2013. 14

A. Arnab, M. Dehghani, G. Heigold, C. Sun, M. Lučić, and C. Schmid. Vivit: A video vision transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6836–6846, 2021. 38

M. Artetxe, J. Du, N. Goyal, L. Zettlemoyer, and V. Stoyanov. On the Role of Bidirectionality in Language Model Pre-Training. *arxiv:2205.11726[cs]*, May 2022. doi: 10.48550/arXiv.2205.11726. URL http://arxiv.org/abs/2205.11726. 39

Y. M. Asano, C. Rupprecht, and A. Vedaldi. Self-labelling via simultaneous clustering and representation learning. *arXiv preprint arXiv:1911.05371*, 2019. 6

Y. M. Asano, C. Rupprecht, A. Zisserman, and A. Vedaldi. Pass: An imagenet replacement for self-supervised pretraining without humans. *arXiv preprint arXiv:2109.13228*, 2021. 19

M. Assran, R. Balestriero, Q. Duval, F. Bordes, I. Misra, P. Bojanowski, P. Vincent, M. Rabbat, and N. Ballas. The hidden uniform cluster prior in self-supervised learning, 2022a. URL https://arxiv.org/abs/2210.07277. 19, 26

M. Assran, M. Caron, I. Misra, P. Bojanowski, F. Bordes, P. Vincent, A. Joulin, M. Rabbat, and N. Ballas. Masked siamese networks for label-efficient learning. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXI*, pages 456–473. Springer, 2022b. 36

M. Assran, M. Caron, I. Misra, P. Bojanowski, F. Bordes, P. Vincent, A. Joulin, M. Rabbat, and N. Ballas. Masked siamese networks for label-efficient learning. In *Computer Vision – ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXI*, page 456–473, Berlin, Heidelberg, 2022c. Springer-Verlag. ISBN 978-3-031-19820-5. doi: 10.1007/978-3-031-19821-2_26. URL https://doi.org/10.1007/978-3-031-19821-2_26. 26

M. Assran, Q. Duval, I. Misra, P. Bojanowski, P. Vincent, M. Rabbat, Y. LeCun, and N. Ballas. Self-supervised learning from images with a joint-embedding predictive architecture, 2023. 21

P. Bachman, R. D. Hjelm, and W. Buchwalter. Learning representations by maximizing mutual information across views. *Advances in neural information processing systems*, 32, 2019. 7

A. Baevski, W.-N. Hsu, A. Conneau, and M. Auli. Unsupervised speech recognition. *Advances in Neural Information Processing Systems*, 34:27826–27839, 2021. 38

A. Baevski, A. Babu, W.-N. Hsu, and M. Auli. Efficient self-supervised learning with contextualized target representations for vision, speech and language. *arXiv preprint arXiv:2212.07525*, 2022. 21, 36

D. Bahri, H. Jiang, Y. Tay, and D. Metzler. Scarf: Self-supervised contrastive learning using random feature corruption. *arXiv preprint arXiv:2106.15147*, 2021. 39

R. Balestriero and Y. LeCun. Contrastive and non-contrastive self-supervised learning recover global and local spectral embedding methods. *arXiv preprint arXiv:2205.11508*, 2022. 4

H. Bao, L. Dong, S. Piao, and F. Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021a. 15, 16

H. Bao, L. Dong, and F. Wei. BEiT: BERT pre-training of image transformers. 2021b. 29, 30, 31, 32

A. Bar, X. Wang, V. Kantorov, C. J. Reed, R. Herzig, G. Chechik, A. Rohrbach, T. Darrell, and A. Globerson. Detreg: Unsupervised pretraining with region priors for object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14605–14615, 2022. 43

A. Bardes, J. Ponce, and Y. LeCun. Vicreg: Variance-invariance-covariance regularization for self-supervised learning. *International Conference on Learning Representations*, 2021. 14, 21, 24, 25, 35

A. Bardes, J. Ponce, and Y. LeCun. Vicregl: Variance-invariance-covariance regularization for self-supervised learning. *Advances in neural information processing systems*, 2022. 22, 33, 43

M. I. Belghazi, A. Baratin, S. Rajeshwar, S. Ozair, Y. Bengio, A. Courville, and D. Hjelm. Mutual information neural estimation. In *International Conference on Machine Learning*, pages 531–540. PMLR, 2018. 17

Y. Bengio and J.-S. Senécal. Quick training of probabilistic neural nets by importance sampling. In *International Workshop on Artificial Intelligence and Statistics*, pages 17–24. PMLR, 2003. 8, 9

Y. Bengio and J.-S. Senécal. Adaptive importance sampling to accelerate training of a neural probabilistic language model. *IEEE Transactions on Neural Networks*, 19(4): 713–722, 2008. 8

Y. Bengio, P. Lamblin, D. Popovici, and H. Larochelle. Greedy layer-wise training of deep networks. *Advances in neural information processing systems*, 19, 2006. 6

D. Berthelot, N. Carlini, I. Goodfellow, N. Papernot, A. Oliver, and C. A. Raffel. Mixmatch: A holistic approach to semi-supervised learning. *Advances in neural information processing systems*, 32, 2019. 7

P. Bojanowski and A. Joulin. Unsupervised learning by predicting noise. In *International Conference on Machine Learning*, pages 517–526. PMLR, 2017. 12

F. Bordes, R. Balestriero, Q. Garrido, A. Bardes, and P. Vincent. Guillotine regularization: Improving deep networks generalization by removing their head. 2022a. 20, 22, 23

F. Bordes, R. Balestriero, and P. Vincent. High fidelity visualization of what your self-supervised representation knows about. *Transactions on Machine Learning Research*, 2022b. URL https://openreview.net/forum?id=urfWb7VjmL. 33, 34

F. Bordes, R. Balestriero, and P. Vincent. Towards democratizing joint-embedding self-supervised learning, 2023a. URL https://arxiv.org/abs/2303.01986. 21, 27, 31, 36, 37

F. Bordes, S. Lavoie, R. Balestriero, N. Ballas, and P. Vincent. A surprisingly simple technique to control the pretraining bias for better transfer: Expand or narrow your representation, 2023b. 25

L. Breiman and J. H. Friedman. Estimating optimal transformations for multiple regression and correlation. *Journal of the American statistical Association*, 80(391):580–598, 1985. 14

J. Bromley, I. Guyon, Y. LeCun, E. Säckinger, and R. Shah. Signature verification using a" siamese" time delay neural network. *Advances in neural information processing systems*, 6, 1993. 7, 10

T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei. Language models are few-shot learners, 2020. URL https://arxiv.org/abs/2005.14165. 3, 39

M. Caron, P. Bojanowski, A. Joulin, and M. Douze. Deep clustering for unsupervised learning of visual features. In *Proceedings of the European conference on computer vision (ECCV)*, pages 132–149, 2018. 6, 12, 44

M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in Neural Information Processing Systems*, 33:9912–9924, 2020. 14, 21, 40, 44

M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9650–9660, 2021. 8, 12, 22, 28, 30, 33, 42

C. Ceylan and M. U. Gutmann. Conditional noise-contrastive estimation of unnormalised models. In *International Conference on Machine Learning*, pages 726–734. PMLR, 2018. 9, 39

H. Chang, H. Zhang, L. Jiang, C. Liu, and W. T. Freeman. Maskgit: Masked generative image transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11315–11325, 2022. 16

H. Chang, H. Zhang, J. Barber, A. Maschinot, J. Lezama, L. Jiang, M.-H. Yang, K. Murphy, W. T. Freeman, M. Rubinstein, et al. Muse: Text-to-image generation via masked generative transformers. *arXiv preprint arXiv:2301.00704*, 2023. 16

G. Chechik, V. Sharma, U. Shalit, and S. Bengio. Large scale online learning of image similarity through ranking. *Journal of Machine Learning Research*, 11(3), 2010. 7, 10

B. Chen, T. Dao, E. Winsor, Z. Song, A. Rudra, and C. Ré. Scatterbrain: Unifying sparse and low-rank attention. *Advances in Neural Information Processing Systems*, 34:17413–17426, 2021a. 36

B. Chen, H. Tang, J. Bu, K. Zhang, J. Wang, Q. Wang, H.-T. Zheng, W. Wu, and L. Yu. Clower: A pre-trained language model with contrastive learning over word and character representations. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 3098–3108, 2022. 39

M. Chen, A. Radford, R. Child, J. Wu, H. Jun, D. Luan, and I. Sutskever. Generative pretraining from pixels. In *International conference on machine learning*, pages 1691–1703. PMLR, 2020a. 39

T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020b. 3, 11, 18, 20, 22, 27, 35, 36, 41

T. Chen, S. Kornblith, K. Swersky, M. Norouzi, and G. Hinton. Big self-supervised models are strong semi-supervised learners. *arXiv preprint arXiv:2006.10029*, 2020c. 23

X. Chen and K. He. Exploring simple siamese representation learning. (arXiv:2011.10566), Nov 2020. URL http://arxiv.org/abs/2011.10566. arXiv:2011.10566 [cs]. 40

X. Chen and K. He. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15750–15758, 2021. 8, 12, 26

X. Chen, H. Fan, R. Girshick, and K. He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020d. 11, 13, 21

X. Chen, S. Xie, and K. He. An empirical study of training self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9640–9649, 2021b. 11, 28

S. Chopra, R. Hadsell, and Y. LeCun. Learning a similarity metric discriminatively, with application to face verification. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, volume 1, pages 539–546. IEEE, 2005. 7, 10

K. Choromanski, V. Likhosherstov, D. Dohan, X. Song, A. Gane, T. Sarlos, P. Hawkins, J. Davis, A. Mohiuddin, L. Kaiser, et al. Rethinking attention with performers. *arXiv preprint arXiv:2009.14794*, 2020. 36

Y.-A. Chung, W.-H. Weng, S. Tong, and J. Glass. Unsupervised cross-modal alignment of speech and text embedding spaces. *Advances in neural information processing systems*, 31, 2018. 38

O. Ciga, T. Xu, and A. L. Martel. Self supervised contrastive learning for digital histopathology. *Machine Learning with Applications*, 7:100198, 2022. ISSN 2666-8270. doi: https://doi.org/10.1016/j.mlwa.2021.100198. URL https://www.sciencedirect.com/science/article/pii/S2666827021000992. 3

K. Clark, M.-T. Luong, Q. V. Le, and C. D. Manning. Electra: Pre-training text encoders as discriminators rather than generators. *arXiv preprint arXiv:2003.10555*, 2020. 28

J. M. Coria, H. Bredin, S. Ghannay, and S. Rosset. A comparison of metric learning loss functions for end-to-end speaker verification. In *International Conference on Statistical Language and Speech Processing*, pages 137–148. Springer, 2020. 17

R. Cosentino, A. Sengupta, S. Avestimehr, M. Soltanolkotabi, A. Ortega, T. Willke, and M. Tepper. Toward a geometrical understanding of self-supervised contrastive learning. *arXiv preprint arXiv:2205.06926*, 2022. 18

Z. Dai, Z. Yang, F. Yang, W. W. Cohen, and R. R. Salakhutdinov. Good semi-supervised learning that requires a bad gan. *Advances in neural information processing systems*, 30, 2017. 6

Z. Dai, B. Cai, Y. Lin, and J. Chen. Up-detr: Unsupervised pre-training for object detection with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1601–1610, 2021. 43

R. Dangovski, L. Jing, C. Loh, S. Han, A. Srivastava, B. Cheung, P. Agrawal, and M. Soljačić. Equivariant contrastive learning. *arXiv preprint arXiv:2111.00899*, 2021. 21

T. Dao, D. Y. Fu, S. Ermon, A. Rudra, and C. Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. *arXiv preprint arXiv:2205.14135*, 2022. 37

J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255. IEEE, 2009. 3

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. 39

J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, 2019. 15, 16

M. Ding, Z. Yang, W. Hong, W. Zheng, C. Zhou, D. Yin, J. Lin, X. Zou, Z. Shao, H. Yang, and J. Tang. CogView: Mastering Text-to-Image Generation via Transformers. In *Advances in Neural Information Processing Systems*, volume 34, pages 19822–19835. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper/2021/hash/a4d92e2cd541fca87e4620aba658316d-Abstract.html. 41

C. Doersch, A. Gupta, and A. A. Efros. Unsupervised visual representation learning by context prediction. In *Proceedings of the IEEE international conference on computer vision*, pages 1422–1430, 2015. 5

J. Donahue, P. Krähenbühl, and T. Darrell. Adversarial feature learning. In *International Conference on Learning Representations*, 2017. 6

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. 31

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 *International Conference on Learning Representations*. 15, 16, 28

A. Dosovitskiy, J. T. Springenberg, M. Riedmiller, and T. Brox. Discriminative unsupervised feature learning with convolutional neural networks. *Advances in neural information processing systems*, 27, 2014. 13

Z.-Y. Dou, Y. Xu, Z. Gan, J. Wang, S. Wang, L. Wang, C. Zhu, P. Zhang, L. Yuan, N. Peng, Z. Liu, and M. Zeng. An Empirical Study of Training End-to-End Vision-and-Language Transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18166–18176, 2022. URL https://openaccess.thecvf.com/content/CVPR2022/html/Dou_An_Empirical_Study_of_Training_End-to-End_Vision-and-Language_Transformers_CVPR_2022_paper.html. 42

Y. Dubois, T. Hashimoto, S. Ermon, and P. Liang. Improving Self-Supervised Learning by Characterizing Idealized Representations, Dec. 2022. URL http://arxiv.org/abs/2209.06235. arXiv:2209.06235 [cs, stat]. 24

D. Dwibedi, Y. Aytar, J. Tompson, P. Sermanet, and A. Zisserman. With a little help from my friends: Nearest-neighbor contrastive learning of visual representations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9588–9597, 2021. 11, 22

C. Dyer. Notes on noise contrastive estimation and negative sampling. *arXiv preprint arXiv:1410.8251*, 2014. 9

D. Eigen, C. Puhrsch, and R. Fergus. Depth map prediction from a single image using a multi-scale deep network. *Advances in neural information processing systems*, 27, 2014. 5

A. El-Nouby, G. Izacard, H. Touvron, I. Laptev, H. Jegou, and E. Grave. Are large-scale datasets necessary for self-supervised pre-training? *arXiv preprint arXiv:2112.10740*, 2021. 19

L. Ericsson, H. Gouk, and T. M. Hospedales. Why do self-supervised models transfer? investigating the impact of invariance on downstream tasks, 2021a. 21

L. Ericsson, H. Gouk, and T. M. Hospedales. How well do self-supervised models transfer? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5414–5423, 2021b. 42

A. Ermolov, A. Siarohin, E. Sangineto, and N. Sebe. Whitening for self-supervised representation learning. In *International Conference on Machine Learning*, pages 3015–3024. PMLR, 2021. 14, 25

B. Eysenbach, T. Zhang, R. Salakhutdinov, and S. Levine. Contrastive learning as goal-conditioned reinforcement learning. (arXiv:2206.07568), Jun 2022. URL http://arxiv.org/abs/2206.07568. number: arXiv:2206.07568 arXiv:2206.07568 [cs]. 40

FairScale. Fairscale: A general purpose modular pytorch library for high performance and large scale training. https://github.com/facebookresearch/fairscale, 2021. 35, 37

J. Fan, Z. Wang, Y. Xie, and Z. Yang. A theoretical analysis of deep q-learning. In *Learning for Dynamics and Control*, pages 486–489. PMLR, 2020. 26

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. *arXiv preprint arXiv:2211.07636*, 2022a. 16, 38

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. *arxiv:2211.07636[cs]*, Dec. 2022b. doi: 10.48550/arXiv.2211.07636. URL http://arxiv.org/abs/2211.07636. 41

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/CVF Conference on Computer Vision and Pattern Recognition*, pages 3299–3309, 2021. 38

C. Feichtenhofer, H. Fan, Y. Li, and K. He. Masked autoencoders as spatiotemporal learners. *arXiv preprint arXiv:2205.09113*, 2022. 38

T. Furlanello, Z. Lipton, M. Tschannen, L. Itti, and A. Anandkumar. Born again neural networks. In *International Conference on Machine Learning*, pages 1607–1616. PMLR, 2018. 13

T. Gao, X. Yao, and D. Chen. Simcse: Simple contrastive learning of sentence embeddings. In *EMNLP (1)*, 2021. 39

Q. Garrido, R. Balestriero, L. Najman, and Y. Lecun. Rankme: Assessing the downstream performance of pretrained self-supervised representations by their rank. *arXiv preprint arXiv:2210.02885*, 2022a. 18, 19, 30, 32

Q. Garrido, Y. Chen, A. Bardes, L. Najman, and Y. Lecun. On the duality between contrastive and non-contrastive self-supervised learning. *arXiv preprint arXiv:2206.02574*, 2022b. 4, 16, 24

Q. Garrido, L. Najman, and Y. Lecun. Self-supervised learning of split invariant equivariant representations. *arXiv preprint arXiv:2302.10283*, 2023. 21

A. Ghiasi, H. Kazemi, E. Borgnia, S. Reich, M. Shu, M. Goldblum, A. G. Wilson, and T. Goldstein. What do Vision Transformers Learn? A Visual Exploration. *arxiv:2212.06727[cs]*, Dec. 2022. doi: 10.48550/arXiv.2212.06727. URL http://arxiv.org/abs/2212.06727. 41

A. Ghosh, A. K. Mondal, K. K. Agrawal, and B. Richards. Investigating power laws in deep representation learning. *arXiv preprint arXiv:2202.05808*, 2022. 18, 19

S. Gidaris, P. Singh, and N. Komodakis. Unsupervised representation learning by predicting image rotations. *arXiv preprint arXiv:1803.07728*, 2018. 5, 21

J. Giorgi, O. Nitski, B. Wang, and G. Bader. Declutr: Deep contrastive learning for unsupervised textual representations. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 879–895, 2021. 39

R. Girdhar, A. El-Nouby, M. Singh, K. V. Alwala, A. Joulin, and I. Misra. Omnimae: Single model masked pretraining on images and videos. *arXiv preprint arXiv:2206.08356*, 2022. 38

J. Goldberger, G. E. Hinton, S. Roweis, and R. R. Salakhutdinov. Neighbourhood components analysis. *Advances in neural information processing systems*, 17, 2004. 10

I. Goodfellow, Y. Bengio, A. Courville, and Y. Bengio. *Deep learning*, volume 1. MIT Press, 2016. 3

I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial networks. In *stat*, 2014. 6

P. Goyal, P. Dollár, R. Girshick, P. Noordhuis, L. Wesolowski, A. Kyrola, A. Tulloch, Y. Jia, and K. He. Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017. 27, 28

P. Goyal, D. Mahajan, A. Gupta, and I. Misra. Scaling and benchmarking self-supervised visual representation learning. In *Proceedings of the ieee/cvf International Conference on computer vision*, pages 6391–6400, 2019. 42

P. Goyal, M. Caron, B. Lefaudeux, M. Xu, P. Wang, V. Pai, M. Singh, V. Liptchinsky, I. Misra, A. Joulin, et al. Self-supervised pretraining of visual features in the wild. *arXiv preprint arXiv:2103.01988*, 2021. 3, 19

P. Goyal, Q. Duval, I. Seessel, M. Caron, M. Singh, I. Misra, L. Sagun, A. Joulin, and P. Bojanowski. Vision models are more robust and fair when pretrained on uncurated images without supervision. *arXiv preprint arXiv:2202.08360*, 2022. 3

K. Grauman, A. Westbury, E. Byrne, Z. Chavis, A. Furnari, R. Girdhar, J. Hamburger, H. Jiang, M. Liu, X. Liu, M. Martin, T. Nagarajan, I. Radosavovic, S. K. Ramakrishnan, F. Ryan, J. Sharma, M. Wray, M. Xu, E. Z. Xu, C. Zhao, S. Bansal, D. Batra, V. Cartillier, S. Crane, T. Do, M. Doulaty, A. Erapalli, C. Feichtenhofer, A. Fragomeni, Q. Fu, A. Gebreselasie, C. Gonzalez, J. Hillis, X. Huang, Y. Huang, W. Jia, W. Khoo, J. Kolar, S. Kottur, A. Kumar, F. Landini, C. Li, Y. Li, Z. Li, K. Mangalam, R. Modhugu, J. Munro, T. Murrell, T. Nishiyasu, W. Price, P. R. Puentes, M. Ramazanova, L. Sari, K. Somasundaram, A. Southerland, Y. Sugano, R. Tao, M. Vo, Y. Wang, X. Wu, T. Yagi, Z. Zhao, Y. Zhu, P. Arbelaez, D. Crandall, D. Damen, G. M. Farinella, C. Fuegen, B. Ghanem, V. K. Ithapu, C. V. Jawahar, H. Joo, K. Kitani, H. Li, R. Newcombe, A. Oliva, H. S. Park, J. M. Rehg, Y. Sato, J. Shi, M. Z. Shou, A. Torralba, L. Torresani, M. Yan, and J. Malik. Ego4d: Around the world in 3,000 hours of egocentric video. (arXiv:2110.07058), Mar 2022. doi: 10.48550/arXiv.2110.07058. URL http://arxiv.org/abs/2110.07058. arXiv:2110.07058 [cs]. 41

J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. Richemond, E. Buchatskaya, C. Doersch, B. Avila Pires, Z. Guo, M. Gheshlaghi Azar, et al. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020. 3, 8, 12, 26, 28, 40

X. Gu, T.-Y. Lin, W. Kuo, and Y. Cui. Open-vocabulary object detection via vision and language knowledge distillation. In *International Conference on Learning Representations*, 2021. 44

Q. Guo, J. Chen, D. Wang, Y. Yang, X. Deng, L. Carin, F. Li, J. Huang, and C. Tao. Tight mutual information estimation with contrastive fenchel-legendre optimization, 2022a. 17

Z. D. Guo, S. Thakoor, M. Pîslar, B. A. Pires, F. Altché, C. Tallec, A. Saade, D. Calandriello, J.-B. Grill, Y. Tang, M. Valko, R. Munos, M. G. Azar, and B. Piot. Byol-explore: Exploration

by bootstrapped prediction. (arXiv:2206.08332), Jun 2022b. URL http://arxiv.org/abs/2206.08332. arXiv:2206.08332 [cs, stat]. 40

M. Gutmann and A. Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 297–304. JMLR Workshop and Conference Proceedings, 2010. 9, 39

R. Hadsell, S. Chopra, and Y. LeCun. Dimensionality reduction by learning an invariant mapping. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 2, pages 1735–1742. IEEE, 2006. 7

M. S. Halvagal, A. Laborieux, and F. Zenke. Predictor networks and stop-grads provide implicit variance regularization in byol/simsiam. *arXiv preprint arXiv:2212.04858*, 2022. 12

J. Z. HaoChen, C. Wei, A. Gaidon, and T. Ma. Provable guarantees for self-supervised deep learning with spectral contrastive loss. In *Advances in Neural Information Processing Systems*, 2021. 4

T. Hastie, R. Tibshirani, and J. Friedman. Overview of supervised learning. In *The elements of statistical learning*, pages 9–41. Springer, 2009. 3

B. He and M. Ozay. Exploring the gap between collapsed & whitened features in self-supervised learning. In *International Conference on Machine Learning*, pages 8613–8634. PMLR, 2022. 18

K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 43

K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020a. 3, 11, 13, 40

K. He, H. Fan, Y. Wu, S. Xie, and R. B. Girshick. Momentum contrast for unsupervised visual representation learning. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9726–9735, 2020b. 18, 35, 42

K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022. 3, 5, 15, 16, 21, 29, 31, 32, 33, 36, 43

O. Henaff. Data-efficient image recognition with contrastive predictive coding. In *International conference on machine learning*, pages 4182–4192. PMLR, 2020. 7

D. Hendrycks, M. Mazeika, S. Kadavath, and D. Song. Using self-supervised learning can improve model robustness and uncertainty. *Advances in neural information processing systems*, 32, 2019. 3

G. Hinton, O. Vinyals, J. Dean, et al. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2(7), 2015. 13

G. E. Hinton, S. Osindero, and Y.-W. Teh. A fast learning algorithm for deep belief nets. *Neural computation*, 18(7):1527–1554, 2006. 6

G. V. Horn, O. M. Aodha, Y. Song, Y. Cui, C. Sun, A. Shepard, H. Adam, P. Perona, and S. Belongie. The inaturalist species classification and detection dataset. In *CVPR*, 2018. 19

H. Hotelling. Relations between two sets of variates. In *Breakthroughs in statistics*, pages 162–190. Springer, 1992. 13, 14

W. W. Hsieh. Nonlinear canonical correlation analysis by neural networks. *Neural Networks*, 13(10):1095–1105, 2000. 14

T. Hua, W. Wang, Z. Xue, S. Ren, Y. Wang, and H. Zhao. On feature decorrelation in self-supervised learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9598–9608, 2021. 18

G. Huang, Y. Sun, Z. Liu, D. Sedra, and K. Q. Weinberger. Deep networks with stochastic depth. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pages 646–661. Springer, 2016. 28

W. Huang, M. Yi, and X. Zhao. Towards the generalization of contrastive self-supervised learning. *arXiv preprint arXiv:2111.00743*, 2021. 16

A. Hyvarinen and H. Morioka. Unsupervised feature extraction by time-contrastive learning and nonlinear ica. *Advances in neural information processing systems*, 29, 2016. 8

M. Ibrahim, D. Bouchacourt, and A. Morcos. Robust self-supervised learning with lie groups. *arXiv preprint arXiv:2210.13356*, 2022. 21

H. Iida, D. Thai, V. Manjunatha, and M. Iyyer. Tabbie: Pretrained representations of tabular data. *arXiv preprint arXiv:2105.02584*, 2021. 39

A. Ivanov, N. Dryden, T. Ben-Nun, S. Li, and T. Hoefler. Data movement is all you need: A case study on optimizing transformers. *Proceedings of Machine Learning and Systems*, 3:711–732, 2021. 37

D. Jarrett, C. Tallec, F. Altché, T. Mesnard, R. Munos, and M. Valko. Curiosity in hindsight. (arXiv:2211.10515), Nov 2022. URL http://arxiv.org/abs/2211.10515. arXiv:2211.10515 [cs, stat].

S. Jean, K. Cho, R. Memisevic, and Y. Bengio. On using very large target vocabulary for neural machine translation. *arXiv preprint arXiv:1412.2007*, 2014. 13

W. Ji, Z. Deng, R. Nakada, J. Zou, and L. Zhang. The power of contrast for feature learning: A theoretical analysis. *arXiv preprint arXiv:2110.02473*, 2021. 17

C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. Le, Y.-H. Sung, Z. Li, and T. Duerig. Scaling Up Visual and Vision-Language Representation Learning With Noisy Text Supervision. In *Proceedings of the 38th International Conference on Machine Learning*, pages 4904–4916. PMLR, July 2021. URL https://proceedings.mlr.press/v139/jia21b.html. 41

L. Jing, P. Vincent, Y. LeCun, and Y. Tian. Understanding dimensional collapse in contrastive self-supervised learning. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=YevsQ05DEN7. 18, 19

A. Joulin, F. Bach, and J. Ponce. Discriminative clustering for image co-segmentation. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 1943–1950. IEEE, 2010. 12

L. Kaiser, M. Babaeizadeh, P. Milos, B. Osinski, R. H. Campbell, K. Czechowski, D. Erhan, C. Finn, P. Kozakowski, S. Levine, A. Mohiuddin, R. Sepassi, G. Tucker, and H. Michalewski. Model-based reinforcement learning for atari. *arXiv:1903.00374 [cs, stat]*, Feb 2020. URL http://arxiv.org/abs/1903.00374. arXiv: 1903.00374. 40

Y. Kalantidis, M. B. Sariyildiz, N. Pion, P. Weinzaepfel, and D. Larlus. Hard negative mixing for contrastive learning. *Advances in Neural Information Processing Systems*, 33: 21798–21809, 2020. 18

A. Kamath, M. Singh, Y. LeCun, G. Synnaeve, I. Misra, and N. Carion. Mdetr-modulated detection for end-to-end multi-modal understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1780–1790, 2021. 44

I. Khemakhem, D. P. Kingma, R. P. Monti, and A. Hyvärinen. Variational autoencoders and nonlinear ica: A unifying framework, 2020. 17

D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 4

N. Kitaev, Ł. Kaiser, and A. Levskaya. Reformer: The efficient transformer. *arXiv preprint arXiv:2001.04451*, 2020. 36

S. A. Koohpayegani, A. Tejankar, and H. Pirsiavash. Mean shift for self-supervised learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10326–10335, 2021. 6, 13, 22

R. Krishnan, P. Rajpurkar, and E. J. Topol. Self-supervised learning in medicine and healthcare. *Nature Biomedical Engineering*, pages 1–7, 2022. 3

G. Larsson, M. Maire, and G. Shakhnarovich. Learning representations for automatic colorization. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pages 577–593. Springer, 2016. 5

G. Leclerc, A. Ilyas, L. Engstrom, S. M. Park, H. Salman, and A. Madry. ffcv. https://github.com/libffcv/ffcv/, 2022. commit xxxxxxx. 36, 37

D.-H. Lee et al. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *Workshop on challenges in representation learning, ICML*, volume 3, page 896, 2013. 6

B. Lefaudeux, F. Massa, D. Liskovich, W. Xiong, V. Caggiano, S. Naren, M. Xu, J. Hu, M. Tintore, S. Zhang, P. Labatut, and D. Haziza. xformers: A modular and hackable transformer modelling library. https://github.com/facebookresearch/xformers, 2022. 37

R. Levin, V. Cherepanova, A. Schwarzschild, A. Bansal, C. B. Bruss, T. Goldstein, A. G. Wilson, and M. Goldblum. Transfer learning with deep tabular models. *International Conference on Learning Representations (ICLR)*, 2023. 40

M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, and L. Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, 2020. 39

A. C. Li, A. A. Efros, and D. Pathak. Understanding collapse in non-contrastive siamese representation learning. In *European Conference on Computer Vision*, pages 490–505. Springer, 2022a. 18, 19, 32

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. In *International Conference on Learning Representations*, 2022b. URL https://openreview.net/forum?id=fVu3o-YUGQK. 36

J. Li, P. Zhou, C. Xiong, and S. C. Hoi. Prototypical contrastive learning of unsupervised representations. *arXiv preprint arXiv:2005.04966*, 2020. 11

J. Li, D. Li, C. Xiong, and S. Hoi. BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation. In *Proceedings of the 39th International Conference on Machine Learning*, pages 12888–12900. PMLR, June 2022c. URL https://proceedings.mlr.press/v162/li22n.html. 41

X. Li, W. Wang, L. Yang, and J. Yang. Uniform masking: Enabling mae pre-training for pyramid-based vision transformers with locality. *arXiv preprint arXiv:2205.10063*, 2022d. 29

Y. Li, R. Pogodin, D. J. Sutherland, and A. Gretton. Self-supervised learning with kernel dependence maximization. *Advances in Neural Information Processing Systems*, 34: 15543–15556, 2021a. 17

Z. Li, Y. Tang, W. Li, and Y. He. Learning disentangled representation with pairwise independence. In *Proceedings of the AAAI conference on artificial intelligence*, 2019. 25

Z. Li, Z. Chen, F. Yang, W. Li, Y. Zhu, C. Zhao, R. Deng, L. Wu, R. Zhao, M. Tang, et al. Mst: Masked self-supervised transformer for visual representation. *Advances in Neural Information Processing Systems*, 34:13165–13176, 2021b. 43

S. Liu, A. Mallol-Ragolta, E. Parada-Cabeleiro, K. Qian, X. Jing, A. Kathan, B. Hu, and B. W. Schuller. Audio self-supervised learning: A survey. *arXiv preprint arXiv:2203.01205*, 2022a. 3, 38

X. Liu, J. Zhou, T. Kong, X. Lin, and R. Ji. Exploring target representations for masked autoencoders. *arXiv preprint arXiv:2209.03917*, 2022b. 29

I. Loshchilov and F. Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 27

Y. J. Ma, S. Sodhani, D. Jayaraman, O. Bastani, V. Kumar, and A. Zhang. Vip: Towards universal visual reward and representation via value-implicit pre-training. (arXiv:2210.00030), Sep 2022. URL http://arxiv.org/abs/2210.00030. arXiv:2210.00030 [cs]. 40

Z. Ma and M. Collins. Noise contrastive estimation and negative sampling for conditional models: Consistency and statistical efficiency. *arXiv preprint arXiv:1809.01812*, 2018. 9

M. C. Machado, M. G. Bellemare, and M. Bowling. A laplacian framework for option discovery in reinforcement learning. *arXiv:1703.00956 [cs]*, Jun 2017. URL http://arxiv.org/abs/1703.00956. arXiv: 1703.00956. 41

A. Majumdar, K. Yadav, S. Arnaud, Y. J. Ma, C. Chen, S. Silwal, A. Jain, V.-P. Berges, P. Abbeel, J. Malik, D. Batra, Y. Lin, O. Maksymets, A. Rajeswaran, and F. Meier. Where are we in the search for an artificial visual cortex for embodied intelligence? 41

A. Makur, F. Kozynski, S.-L. Huang, and L. Zheng. An efficient algorithm for information decomposition and extraction. In *2015 53rd Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, pages 972–979. IEEE, 2015. 14

G. L. Marchetti, G. Tegnér, A. Varava, and D. Kragic. Equivariant representation learning via class-pose decomposition. *arXiv preprint arXiv:2207.03116*, 2022. 21

Y. Meng, C. Xiong, P. Bajaj, P. Bennett, J. Han, X. Song, et al. Coco-lm: Correcting and contrasting text sequences for language model pretraining. *Advances in Neural Information Processing Systems*, 34:23102–23114, 2021. 21, 39

G. Mialon, R. Balestriero, and Y. Lecun. Variance-covariance regularization enforces pair-wise independence in self-supervised representations. *arXiv preprint arXiv:2209.14905*, 2022. 18, 25

T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*, 26, 2013. 9, 39

B. Min, H. Ross, E. Sulem, A. P. B. Veyseh, T. H. Nguyen, O. Sainz, E. Agirre, I. Heinz, and D. Roth. Recent advances in natural language processing via large pre-trained language models: A survey. *arXiv preprint arXiv:2111.01243*, 2021. 38

M. Minderer, A. Gritsenko, A. Stone, M. Neumann, D. Weissenborn, A. Dosovitskiy, A. Mahendran, A. Arnab, M. Dehghani, Z. Shen, et al. Simple open-vocabulary object detection with vision transformers. *arXiv preprint arXiv:2205.06230*, 2022. 44

I. Misra and L. v. d. Maaten. Self-supervised learning of pretext-invariant representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6707–6717, 2020. 3, 42

J. Mitrovic, B. McWilliams, J. Walker, L. Buesing, and C. Blundell. Representation learning via invariant causal mechanisms. *arXiv preprint arXiv:2010.07922*, 2020. 11

T. Miyato, S.-i. Maeda, M. Koyama, and S. Ishii. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):1979–1993, 2018. 6

A. Mnih and K. Kavukcuoglu. Learning word embeddings efficiently with noise-contrastive estimation. *Advances in neural information processing systems*, 26, 2013. 8

A. Mnih and Y. W. Teh. A fast and simple algorithm for training neural probabilistic language models. *arXiv preprint arXiv:1206.6426*, 2012. 8, 9

V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013. 13

N. Mu, A. Kirillov, D. Wagner, and S. Xie. SLIP: Self-supervision Meets Language-Image Pre-training. In S. Avidan, G. Brostow, M. Cissé, G. M. Farinella, and T. Hassner, editors, *Computer Vision – ECCV 2022*, Lecture Notes in Computer Science, pages 529–544, Cham, 2022. Springer Nature Switzerland. ISBN 978-3-031-19809-0. doi: 10.1007/978-3-031-19809-0_30. 41

S. Nair, A. Rajeswaran, V. Kumar, C. Finn, and A. Gupta. R3m: A universal visual representation for robot manipulation. (arXiv:2203.12601), Nov 2022. URL http://arxiv.org/abs/2203.12601. arXiv:2203.12601 [cs]. 41

B. Ni, H. Peng, M. Chen, S. Zhang, G. Meng, J. Fu, S. Xiang, and H. Ling. Expanding language-image pretrained models for general video recognition. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part IV*, pages 1–18. Springer, 2022. 38

R. Ni, M. Goldblum, A. Sharaf, K. Kong, and T. Goldstein. Data augmentation for meta-learning. In *International Conference on Machine Learning*, pages 8152–8161. PMLR, 2021a. 21

R. Ni, M. Shu, H. Souri, M. Goldblum, and T. Goldstein. The close relationship between contrastive learning and meta-learning. In *International Conference on Learning Representations*, 2021b. 7, 21

A. Nichol, P. Dhariwal, A. Ramesh, P. Shyam, P. Mishkin, B. McGrew, I. Sutskever, and M. Chen. GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models. *arxiv:2112.10741[cs]*, Mar. 2022. doi: 10.48550/arXiv.2112.10741. URL http://arxiv.org/abs/2112.10741. 41

M. Noroozi and P. Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *European conference on computer vision*, pages 69–84. Springer, 2016. 14

M. Noroozi, H. Pirsiavash, and P. Favaro. Representation learning by learning to count. In *Proceedings of the IEEE international conference on computer vision*, pages 5898–5906, 2017. 5

M. Noroozi, A. Vinjimoor, P. Favaro, and H. Pirsiavash. Boosting self-supervised learning via knowledge transfer. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9359–9367, 2018. 5

NVidia. Apex. https://github.com/nvidia/apex, 2021. 35, 37

H. Oh Song, Y. Xiang, S. Jegelka, and S. Savarese. Deep metric learning via lifted structured feature embedding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4004–4012, 2016. 7, 17

A. v. d. Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016. 38

A. v. d. Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 7, 10, 17, 18, 38, 40

M. Oquab, T. Darcet, T. Moutakanni, H. V. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haziza, F. Massa, A. El-Nouby, R. Howes, P.-Y. Huang, H. Xu, V. Sharma, S.-W. Li, W. Galuba, M. Rabbat, M. Assran, N. Ballas, G. Synnaeve, I. Misra, H. Jegou, J. Mairal, P. Labatut, A. Joulin, and P. Bojanowski. Dinov2: Learning robust visual features without supervision, 2023. 12, 16, 20, 22, 29

M. Ott, S. Edunov, A. Baevski, A. Fan, S. Gross, N. Ng, D. Grangier, and M. Auli. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of NAACL-HLT 2019: Demonstrations*, 2019. 37

A. Owens, J. Wu, J. H. McDermott, W. T. Freeman, and A. Torralba. Ambient sound provides supervision for visual learning. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part I 14*, pages 801–816. Springer, 2016. 5

A. Painsky, M. Feder, and N. Tishby. Nonlinear canonical correlation analysis: A compressed representation approach. *Entropy*, 22(2):208, 2020. 14

N. Parikh, S. Boyd, et al. Proximal algorithms. *Foundations and trends® in Optimization*, 1(3):127–239, 2014. 8

S. Parisi, A. Rajeswaran, S. Purushwalkam, and A. Gupta. The unsurprising effectiveness of pre-trained vision models for control. (arXiv:2203.03580), Aug 2022. URL http://arxiv.org/abs/2203.03580. arXiv:2203.03580 [cs]. 40

J. Y. Park, O. Biza, L. Zhao, J. W. van de Meent, and R. Walters. Learning symmetric embeddings for equivariant world models. *arXiv preprint arXiv:2204.11371*, 2022. 21

D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2536–2544, 2016. 5, 14, 16

D. Pathak, R. Girshick, P. Dollár, T. Darrell, and B. Hariharan. Learning features by watching objects move. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2701–2710, 2017. 5

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. 32

T. Pham, C. Zhang, A. Niu, K. Zhang, and C. D. Yoo. On the pros and cons of momentum encoder in self-supervised visual representation learning. *arXiv preprint arXiv:2208.05744*, 2022. 12, 26

A. Piché, J. Marino, G. M. Marconi, C. Pal, and M. E. Khan. Beyond target networks: Improving deep $q$-learning with functional regularization. *arXiv preprint arXiv:2106.02613*, 2021. 26

E. Pizzi, S. D. Roy, S. N. Ravindra, P. Goyal, and M. Douze. A self-supervised descriptor for image copy detection, 2022. 13

M. Popel, M. Tomkova, J. Tomek, Ł. Kaiser, J. Uszkoreit, O. Bojar, and Z. Žabokrtský. Transforming machine translation: a deep learning system reaches news translation quality comparable to human professionals. *Nature communications*, 11(1):4381, 2020. 3

S. Purushwalkam and A. Gupta. Demystifying contrastive self-supervised learning: Invariances, augmentations and dataset biases. *Advances in Neural Information Processing Systems*, 33:3407–3418, 2020. 42, 43

A. Radford, K. Narasimhan, T. Salimans, I. Sutskever, et al. Improving language understanding by generative pre-training. 2018. 38, 39

A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019. 38, 39

A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever. Learning Transferable Visual Models From Natural Language Supervision. *arxiv:2103.00020[cs]*, Feb. 2021. doi: 10.48550/arXiv.2103.00020. URL http://arxiv.org/abs/2103.00020. 41

C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020. 15, 39

Y. Rao, W. Zhao, G. Chen, Y. Tang, Z. Zhu, G. Huang, J. Zhou, and J. Lu. DenseCLIP: Language-Guided Dense Prediction With Context-Aware Prompting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18082–18091, 2022. URL https://openaccess.thecvf.com/content/CVPR2022/html/Rao_DenseCLIP_Language-Guided_Dense_Prediction_With_Context-Aware_Prompting_CVPR_2022_paper.html. 41

C. J. Reed, S. Metzger, A. Srinivas, T. Darrell, and K. Keutzer. Selfaugment: Automatic augmentation policies for self-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2674–2683, 2021. 32

S. Reed, K. Zolna, E. Parisotto, S. G. Colmenarejo, A. Novikov, G. Barth-maron, M. Giménez, Y. Sulsky, J. Kay, J. T. Springenberg, T. Eccles, J. Bruce, A. Razavi, A. Edwards, N. Heess, Y. Chen, R. Hadsell, O. Vinyals, M. Bordbar, and N. de Freitas. A Generalist Agent. *Transactions on Machine Learning Research*, Nov. 2022. ISSN 2835-8856. URL https://openreview.net/forum?id=1ikK0kHjvj. 42

J. Robinson, C.-Y. Chuang, S. Sra, and S. Jegelka. Contrastive learning with hard negative samples. *arXiv preprint arXiv:2010.04592*, 2020. 18

I. Rubachev, A. Alekberov, Y. Gorishniy, and A. Babenko. Revisiting pretraining objectives for tabular deep learning. *arXiv preprint arXiv:2207.03208*, 2022. 38, 39, 40

A. Sablayrolles, M. Douze, C. Schmid, and H. Jégou. Spreading vectors for similarity search. *arXiv preprint arXiv:1806.03198*, 2018. 13

T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. Improved techniques for training gans. *Advances in neural information processing systems*, 29, 2016. 6

M. B. Sariyildiz, Y. Kalantidis, K. Alahari, and D. Larlus. Improving the generalization of supervised models, 2022. 22

F. Scherr, Q. Guo, and T. Moraitis. Self-supervised learning through efference copies. In A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=DotEQCtY67g. 21

M. C. Schiappa, Y. S. Rawat, and M. Shah. Self-supervised learning for videos: A survey. *ACM Computing Surveys*, 2022a. 3

M. C. Schiappa, Y. S. Rawat, and M. Shah. Self-supervised learning for videos: A survey. *ACM Computing Surveys*, 2022b. 38

S. Schneider, A. Baevski, R. Collobert, and M. Auli. wav2vec: Unsupervised pre-training for speech recognition. *arXiv preprint arXiv:1904.05862*, 2019. 38

J. Schrittwieser, I. Antonoglou, T. Hubert, K. Simonyan, L. Sifre, S. Schmitt, A. Guez, E. Lockhart, D. Hassabis, T. Graepel, T. Lillicrap, and D. Silver. Mastering atari, go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, Dec 2020. ISSN 0028-0836, 1476-4687. doi: 10.1038/s41586-020-03051-4. arXiv: 1911.08265. 40

F. Schroff, D. Kalenichenko, and J. Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823, 2015. 7, 8, 17

C. Schuhmann, R. Beaumont, R. Vencu, C. Gordon, R. Wightman, M. Cherti, T. Coombes, A. Katta, C. Mullis, M. Wortsman, P. Schramowski, S. Kundurthy, K. Crowson, L. Schmidt, R. Kaczmarczyk, and J. Jitsev. LAION-5B: An open large-scale dataset for training next generation image-text models. *arxiv:2210.08402[cs]*, Oct. 2022. doi: 10.48550/arXiv.2210.08402. URL http://arxiv.org/abs/2210.08402. 41

M. Schwarzer, A. Anand, R. Goel, R. D. Hjelm, A. Courville, and P. Bachman. Data-efficient reinforcement learning with self-predictive representations. *arXiv:2007.05929 [cs, stat]*, May 2021a. URL http://arxiv.org/abs/2007.05929. arXiv: 2007.05929. 40

M. Schwarzer, N. Rajkumar, M. Noukhovitch, A. Anand, L. Charlin, D. Hjelm, P. Bachman, and A. Courville. Pretraining representations for data-efficient reinforcement learning. *arXiv:2106.04799 [cs]*, Jun 2021b. URL http://arxiv.org/abs/2106.04799. arXiv: 2106.04799. 40

P. Sermanet, C. Lynch, Y. Chebotar, J. Hsu, E. Jang, S. Schaal, S. Levine, and G. Brain. Time-contrastive networks: Self-supervised learning from video. In *2018 IEEE international conference on robotics and automation (ICRA)*, pages 1134–1141. IEEE, 2018. 8, 38

S. Shen, L. H. Li, H. Tan, M. Bansal, A. Rohrbach, K.-W. Chang, Z. Yao, and K. Keutzer. How Much Can CLIP Benefit Vision-and-Language Tasks? In *International Conference on Learning Representations*, Jan. 2022. URL https://openreview.net/forum?id=zf_Ll3HZWgy. 42

H. Shi, D. Luo, S. Tang, J. Wang, and Y. Zhuang. Run away from your teacher: Understanding byol by a novel self-supervised approach. *arXiv preprint arXiv:2011.10944*, 2020. 26, 28

R. Shwartz-Ziv, M. Goldblum, H. Souri, S. Kapoor, C. Zhu, Y. LeCun, and A. G. Wilson. Pre-train your loss: Easy bayesian transfer learning with informative priors. In *Advances in Neural Information Processing Systems*. 42

R. Shwartz-Ziv, R. Balestriero, and Y. LeCun. What do we maximize in self-supervised learning? *arXiv preprint arXiv:2207.10081*, 2022. 4

A. Singh, R. Hu, V. Goswami, G. Couairon, W. Galuba, M. Rohrbach, and D. Kiela. FLAVA: A Foundational Language and Vision Alignment Model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15638–15650, 2022. URL https://openaccess.thecvf.com/content/CVPR2022/html/Singh_FLAVA_A_Foundational_Language_and_Vision_Alignment_Model_CVPR_2022_paper.html. 41

V. Sobal, J. S V, S. Jalagam, N. Carion, K. Cho, and Y. LeCun. Joint embedding predictive architectures focus on slow features. (arXiv:2211.10831), Nov 2022. URL http://arxiv.org/abs/2211.10831. arXiv:2211.10831 [cs]. 40

K. Sohn. Improved deep metric learning with multi-class n-pair loss objective. *Advances in neural information processing systems*, 29, 2016. 7, 10, 17

G. Somepalli, M. Goldblum, A. Schwarzschild, C. B. Bruss, and T. Goldstein. Saint: Improved neural networks for tabular data via row attention and contrastive pre-training. *arXiv preprint arXiv:2106.01342*, 2021. 39

J. T. Springenberg. Unsupervised and semi-supervised learning with categorical generative adversarial networks. *arXiv preprint arXiv:1511.06390*, 2015. 6

A. Srinivas, M. Laskin, and P. Abbeel. Curl: Contrastive unsupervised representations for reinforcement learning. *arXiv:2004.04136 [cs, stat]*, Sep 2020. URL http://arxiv.org/abs/2004.04136. arXiv: 2004.04136. 40

A. Steiner, A. Kolesnikov, X. Zhai, R. Wightman, J. Uszkoreit, and L. Beyer. How to train your vit? data, augmentation, and regularization in vision transformers. *arXiv preprint arXiv:2106.10270*, 2021. 28

S. Tang, F. Zhu, L. Bai, R. Zhao, C. Wang, and W. Ouyang. Unifying visual contrastive learning for object recognition from a graph perspective. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXVI*, pages 649–667. Springer, 2022. 22

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. *arXiv preprint arXiv:2112.05141*, 2021. 17

A. Tarvainen and H. Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. *Advances in neural information processing systems*, 30, 2017. 13

Y. Tay, M. Dehghani, V. Q. Tran, X. Garcia, D. Bahri, T. Schuster, H. S. Zheng, N. Houlsby, and D. Metzler. Unifying Language Learning Paradigms. *arxiv:2205.05131[cs]*, May 2022. URL http://arxiv.org/abs/2205.05131. 15, 39

A. Tejankar, S. A. Koohpayegani, V. Pillai, P. Favaro, and H. Pirsiavash. Isd: Self-supervised learning by iterative similarity distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9609–9618, 2021. 13

Y. Tian. Understanding deep contrastive learning via coordinate-wise optimization. *CoRR*, abs/2201.12680, 2022. URL https://arxiv.org/abs/2201.12680. 17, 18

Y. Tian, D. Krishnan, and P. Isola. Contrastive multiview coding. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XI 16*, pages 776–794. Springer, 2020a. 8, 18

Y. Tian, L. Yu, X. Chen, and S. Ganguli. Understanding self-supervised learning with dual deep networks. *arXiv preprint arXiv:2010.00578*, 2020b. 17, 28

Y. Tian, X. Chen, and S. Ganguli. Understanding self-supervised learning dynamics without contrastive pairs. In *International Conference on Machine Learning*, pages 10268–10278. PMLR, 2021. 12, 26

N. Tomasev, I. Bica, B. McWilliams, L. Buesing, R. Pascanu, C. Blundell, and J. Mitrovic. Pushing the limits of self-supervised resnets: Can we outperform supervised learning without labels on imagenet? *arXiv preprint arXiv:2201.05119*, 2022. 3

Z. Tong, Y. Song, J. Wang, and L. Wang. Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training. *arXiv preprint arXiv:2203.12602*, 2022. 38

A. Touati, J. Rapin, and Y. Ollivier. Does zero-shot reinforcement learning exist? 2023. 41

H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jégou. Training data-efficient image transformers & distillation through attention. In *International conference on machine learning*, pages 10347–10357. PMLR, 2021a. 28

H. Touvron, M. Cord, A. Sablayrolles, G. Synnaeve, and H. Jégou. Going deeper with image transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 32–42, 2021b. 28, 29

H. Touvron, M. Cord, and H. Jégou. Deit iii: Revenge of the vit. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXIV*, pages 516–533. Springer, 2022. 28

M. Tschannen, J. Djolonga, P. K. Rubenstein, S. Gelly, and M. Lucic. On mutual information maximization for representation learning, 2020. 17

T. Ucar, E. Hajiramezanali, and L. Edwards. Subtab: Subsetting features of tabular data for self-supervised representation learning. *Advances in Neural Information Processing Systems*, 34:18853–18865, 2021. 39

A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 36

P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning*, pages 1096–1103, 2008. 4, 6, 14

P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, P.-A. Manzagol, and L. Bottou. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *Journal of machine learning research*, 11(12), 2010. 4

C. Vondrick, A. Shrivastava, A. Fathi, S. Guadarrama, and K. Murphy. Tracking emerges by colorizing videos. In *Proceedings of the European conference on computer vision (ECCV)*, pages 391–408, 2018. 5

G. Wang, K. Wang, G. Wang, P. H. Torr, and L. Lin. Solving inefficiency of self-supervised representation learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9505–9515, 2021a. 22

H. Wang, S. Ge, Z. Lipton, and E. P. Xing. Learning Robust Global Representations by Penalizing Local Predictive Power. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper/2019/hash/3eefceb8087e964f89c2d59e8a249915-Abstract.html. 41

S. Wang, B. Z. Li, M. Khabsa, H. Fang, and H. Ma. Linformer: Self-attention with linear complexity. *arXiv preprint arXiv:2006.04768*, 2020a. 36

T. Wang and P. Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *International Conference on Machine Learning*, pages 9929–9939. PMLR, 2020. 17

T. Wang, A. Roberts, D. Hesslow, T. L. Scao, H. W. Chung, I. Beltagy, J. Launay, and C. Raffel. What Language Model Architecture and Pretraining Objective Work Best for Zero-Shot Generalization? *arXiv:2204.05832 [cs, stat]*, Apr. 2022a. URL http://arxiv.org/abs/2204.05832. 15

W. Wang, R. Arora, K. Livescu, and J. Bilmes. On deep multi-view representation learning. In *International conference on machine learning*, pages 1083–1092. PMLR, 2015. 6, 14

W. Wang, Q. Tang, and K. Livescu. Unsupervised pre-training of bidirectional speech encoders via masked reconstruction. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6889–6893. IEEE, 2020b. 38

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*, 2022b. 32

X. Wang and A. Gupta. Unsupervised learning of visual representations using videos. In *Proceedings of the IEEE international conference on computer vision*, pages 2794–2802, 2015. 5

X. Wang, R. Zhang, C. Shen, T. Kong, and L. Li. Dense contrastive learning for self-supervised visual pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3024–3033, 2021b. 43

Z. Wang, J. Yu, A. W. Yu, Z. Dai, Y. Tsvetkov, and Y. Cao. SimVLM: Simple Visual Language Model Pretraining with Weak Supervision. In *International Conference on Learning Representations*, Jan. 2022c. URL https://openreview.net/forum?id=GUrhfTuf_3. 41

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. 32

K. Q. Weinberger and L. K. Saul. Distance metric learning for large margin nearest neighbor classification. *Journal of machine learning research*, 10(2), 2009. 7, 8, 10

K. Wickstrøm, M. Kampffmeyer, K. Ø. Mikalsen, and R. Jenssen. Mixing up contrastive learning: Self-supervised representation learning for time series. *Pattern Recognition Letters*, 155:54–61, 2022. 3

S. Woo, S. Debnath, R. Hu, X. Chen, Z. Liu, I. S. Kweon, and S. Xie. Convnext v2: Co-designing and scaling convnets with masked autoencoders. *arXiv preprint arXiv:2301.00808*, 2023. 16, 29, 43

Z. Wu, Y. Xiong, S. X. Yu, and D. Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3733–3742, 2018. 8, 10, 39

T. Xiao, X. Wang, A. A. Efros, and T. Darrell. What should not be contrastive in contrastive learning. *arXiv preprint arXiv:2008.05659*, 2020. 21

T. Xiao, I. Radosavovic, T. Darrell, and J. Malik. Masked visual pre-training for motor control. (arXiv:2203.06173), Mar 2022. URL http://arxiv.org/abs/2203.06173. arXiv:2203.06173 [cs]. 40

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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9653–9663, 2022. 15, 16

L. Xu, J. Neufeld, B. Larson, and D. Schuurmans. Maximum margin clustering. *Advances in neural information processing systems*, 17, 2004. 12

C. Yang, Z. Wu, B. Zhou, and S. Lin. Instance localization for self-supervised detection pretraining. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3987–3996, 2021. 43

D. Yarats, R. Fergus, A. Lazaric, and L. Pinto. Reinforcement learning with prototypical representations. *arXiv:2102.11271 [cs]*, Jul 2021. URL http://arxiv.org/abs/2102.11271. arXiv: 2102.11271. 40

W. Ye, S. Liu, T. Kurutach, P. Abbeel, and Y. Gao. Mastering atari games with limited data. *arXiv:2111.00210 [cs]*, Dec 2021. URL http://arxiv.org/abs/2111.00210. arXiv: 2111.00210. 40

C.-H. Yeh, C.-Y. Hong, Y.-C. Hsu, T.-L. Liu, Y. Chen, and Y. LeCun. Decoupled contrastive learning. *arXiv preprint arXiv:2110.06848*, 2021. 11, 18, 27

A. YM., R. C., and V. A. Self-labelling via simultaneous clustering and representation learning. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=Hyx-jyBFPr. 12

J. Yoon, Y. Zhang, J. Jordon, and M. van der Schaar. Vime: Extending the success of self-and semi-supervised learning to tabular domain. *Advances in Neural Information Processing Systems*, 33:11033–11043, 2020. 39

J. Yosinski, J. Clune, Y. Bengio, and H. Lipson. How transferable are features in deep neural networks? In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014. URL https://proceedings.neurips.cc/paper/2014/file/375c71349b295fbe2dcdca9206f20a06-Paper.pdf. 22

Y. You, I. Gitman, and B. Ginsburg. Large batch training of convolutional networks. *arXiv preprint arXiv:1708.03888*, 2017. 27

B. Yu and D. Tao. Deep metric learning with tuplet margin loss. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6490–6499, 2019. 10

J. Yu, Y. Xu, J. Y. Koh, T. Luong, G. Baid, Z. Wang, V. Vasudevan, A. Ku, Y. Yang, B. K. Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *Transactions on Machine Learning Research*. 16

L. Yuan, D. Chen, Y.-L. Chen, N. Codella, X. Dai, J. Gao, H. Hu, X. Huang, B. Li, C. Li, C. Liu, M. Liu, Z. Liu, Y. Lu, Y. Shi, L. Wang, J. Wang, B. Xiao, Z. Xiao, J. Yang, M. Zeng, L. Zhou, and P. Zhang. Florence: A New Foundation Model for Computer Vision. *arxiv:2111.11432[cs]*, Nov. 2021. doi: 10.48550/arXiv.2111.11432. URL http://arxiv.org/abs/2111.11432. 41

S. Yun, D. Han, S. J. Oh, S. Chun, J. Choe, and Y. Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6023–6032, 2019. 39

S. Yun, H. Lee, J. Kim, and J. Shin. Patch-level representation learning for self-supervised vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8354–8363, 2022. 43

M. Zaheer, G. Guruganesh, K. A. Dubey, J. Ainslie, C. Alberti, S. Ontanon, P. Pham, A. Ravula, Q. Wang, L. Yang, et al. Big bird: Transformers for longer sequences. *Advances in neural information processing systems*, 33:17283–17297, 2020. 36

A. Zareian, K. D. Rosa, D. H. Hu, and S.-F. Chang. Open-vocabulary object detection using captions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14393–14402, 2021. 44

J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny. Barlow twins: Self-supervised learning via redundancy reduction. In *International Conference on Machine Learning*, pages 12310–12320. PMLR, 2021. 14, 21, 24, 25, 35, 40

X. Zhai, A. Kolesnikov, N. Houlsby, and L. Beyer. Scaling vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12104–12113, 2022a. 29

X. Zhai, X. Wang, B. Mustafa, A. Steiner, D. Keysers, A. Kolesnikov, and L. Beyer. LiT: Zero-Shot Transfer With Locked-Image Text Tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18123–18133, 2022b. URL https://openaccess.thecvf.com/content/CVPR2022/html/Zhai_LiT_

Zero-Shot_Transfer_With_Locked-Image_Text_Tuning_CVPR_2022_paper.html.
42

A. Zhang, R. McAllister, R. Calandra, Y. Gal, and S. Levine. Learning invariant representations for reinforcement learning without reconstruction. (arXiv:2006.10742), Apr 2021. URL http://arxiv.org/abs/2006.10742. arXiv:2006.10742 [cs, stat]. 41

C. Zhang, K. Zhang, T. X. Pham, A. Niu, Z. Qiao, C. D. Yoo, and I. S. Kweon. Dual temperature helps contrastive learning without many negative samples: Towards understanding and simplifying moco. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14441–14450, 2022a. 27

C. Zhang, K. Zhang, C. Zhang, T. X. Pham, C. D. Yoo, and I. S. Kweon. How does simsiam avoid collapse without negative samples? a unified understanding with self-supervised contrastive learning. *arXiv preprint arXiv:2203.16262*, 2022b. 28

H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations*, 2018. 39

R. Zhang, P. Isola, and A. A. Efros. Colorful image colorization. In *European conference on computer vision*, pages 649–666. Springer, 2016. 5, 14, 30

R. Zhang, P. Isola, and A. A. Efros. Split-brain autoencoders: Unsupervised learning by cross-channel prediction. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1058–1067, 2017. 6, 30

W. Zhang, A. GX-Chen, V. Sobal, Y. LeCun, and N. Carion. Light-weight probing of unsupervised representations for reinforcement learning. (arXiv:2208.12345), Aug 2022c. URL http://arxiv.org/abs/2208.12345. arXiv:2208.12345 [cs]. 40

H. Zhao, C. Gan, A. Rouditchenko, C. Vondrick, J. McDermott, and A. Torralba. The sound of pixels. In *Proceedings of the European conference on computer vision (ECCV)*, pages 570–586, 2018. 38

N. Zhao, Z. Wu, R. W. Lau, and S. Lin. Distilling localization for self-supervised representation learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 10990–10998, 2021. 42, 43

B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva. Learning deep features for scene recognition using places database. In *NeurIPS*, 2014. 19

J. Zhou, C. Wei, H. Wang, W. Shen, C. Xie, A. Yuille, and T. Kong. ibot: Image bert pre-training with online tokenizer. *International Conference on Learning Representations (ICLR)*, 2022a. 12, 16, 22, 29, 31, 33

P. Zhou, Y. Zhou, C. Si, W. Yu, T. K. Ng, and S. Yan. Mugs: A multi-granular self-supervised learning framework. *arXiv preprint arXiv:2203.14415*, 2022b. 22

T. Zhou, M. Brown, N. Snavely, and D. G. Lowe. Unsupervised learning of depth and ego-motion from video. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1851–1858, 2017. 5

A. Ziegler and Y. M. Asano. Self-supervised learning of object parts for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14502–14511, 2022. 43