# Self-Supervised Learning by Cross-Modal Audio-Video Clustering

**Humam Alwassel**[1][*]
humam.alwassel@kaust.edu.sa

**Dhruv Mahajan**[2]
dhruvm@fb.com

**Bruno Korbar**[2]
bkorbar@fb.com

**Lorenzo Torresani**[2]
torresani@fb.com

**Bernard Ghanem**[1]
bernard.ghanem@kaust.edu.sa

**Du Tran**[2]
trandu@fb.com

[1]King Abdullah University of Science and Technology (KAUST)    [2]Facebook AI
http://humamalwassel.com/publication/xdc

## Abstract

Visual and audio modalities are highly correlated, yet they contain different information. Their strong correlation makes it possible to predict the semantics of one from the other with good accuracy. Their intrinsic differences make cross-modal prediction a potentially more rewarding pretext task for self-supervised learning of video and audio representations compared to within-modality learning. Based on this intuition, we propose *Cross-Modal Deep Clustering (XDC)*, a novel self-supervised method that leverages unsupervised clustering in one modality (*e.g.*, audio) as a supervisory signal for the other modality (*e.g.*, video). This cross-modal supervision helps XDC utilize the semantic correlation and the differences between the two modalities. Our experiments show that XDC outperforms single-modality clustering and other multi-modal variants. XDC achieves state-of-the-art accuracy among self-supervised methods on multiple video and audio benchmarks. Most importantly, our video model pretrained on large-scale unlabeled data significantly outperforms the same model pretrained with full-supervision on ImageNet and Kinetics for action recognition on HMDB51 and UCF101. To the best of our knowledge, XDC is the first self-supervised learning method that outperforms large-scale fully-supervised pretraining for action recognition on the same architecture.

## 1 Introduction

Do we need to explicitly name the actions of "laughing" or "sneezing" in order to recognize them? Or can we learn to visually classify them without labels by associating characteristic sounds with these actions? Indeed, a wide literature in perceptual studies provides evidence that we rely heavily on hearing sounds to make sense of actions and dynamic events in the visual world. For example, objects moving together are perceived as bouncing off each other when the visual stimulus is accompanied by a brief sound [58], and the location and timing of sounds are leveraged as important cues to direct our spatiotemporal visual attention [20, 43]. The influence of hearing sounds in visual perception is also suggested by perceptual studies showing that individuals affected by profound deafness exhibit poorer visual perceptual performance compared to age-matched hearing controls [11, 40].

In this work, we investigate the hypothesis that spatiotemporal models for action recognition can be reliably pretrained from *unlabeled* videos by capturing cross-modal information from audio and video. The motivation for our study stems from two fundamental challenges facing a fully-supervised

---

[*]Work done during an internship at Facebook AI

line of attack to learning video models. The first challenge is the exorbitant cost of scaling up the size of manually-labeled video datasets. The recent creation of large-scale action recognition datasets [5, 16, 26, 27] has undoubtedly enabled a major leap forward in video models accuracies. However, it may be argued that additional significant gains by dataset growth would require scaling up existing labeled datasets by several orders of magnitude. The second challenge is posed by the unclear definition of suitable label spaces for action recognition. Recent video datasets differ substantially in their label spaces, which range from sports actions [26] to verb-noun pairs for kitchen activities [7]. This suggests that the definition of the "right" label space for action recognition, and more generally for video understanding, is still very much up for debate. It also implies that finetuning models pretrained on large-scale labeled datasets is a suboptimal proxy for learning models for small- or medium-size datasets due to the label-space gap often encountered between source and target datasets.

In this paper, we present three approaches for training video models from self-supervised audio-visual information. At a high-level, the idea behind all three frameworks is to leverage one modality (say, audio) as a supervisory signal for the other (say, video). We posit that this is a promising avenue because of the simultaneous synergy and complementarity of audio and video: correlations between these two modalities make it possible to perform prediction from one to the other, while their intrinsic differences make cross-modal prediction an enriching self-supervised task compared to within-modality learning. Specifically, we adapt the single-modality DeepCluster work of Caron *et al.* [6] to our multi-modal setting. DeepCluster was introduced as a self-supervised procedure for learning image representation. It alternates between unsupervised clustering of image features and using these cluster assignments as pseudo-labels to revise the image representation. In our work, the clusters learned from one modality are used as pseudo-labels to refine the representation for the other modality. In two of our approaches—Multi-Head Deep Clustering (MDC) and Concatenation Deep Clustering (CDC)—the pseudo-labels from the second modality are *supplementary*, *i.e.*, they complement the pseudo-labels generated in the first modality. The third approach—Cross-Modal Deep Clustering (XDC)—instead uses the pseudo-labels from the other modality as an *exclusive* supervisory signal. This means that in XDC, the audio clusters drive the learning of the video representation and vice versa. Our experiments support several interesting conclusions:

- All three of our cross-modal methods yield representations that generalize better to the downstream tasks of action recognition and audio classification, compared to their within-modality counterparts.
- XDC (*i.e.*, the cross-modal deep clustering relying on the other modality as an exclusive supervisory signal) outperforms all the other approaches. This underscores the complementarity of audio and video and the benefits of learning label-spaces across modalities.
- Self-supervised cross-modal learning with XDC on a large-scale video dataset yields an action recognition model that achieves higher accuracy when finetuned on HMDB51 or UCF101, compared to that produced by fully-supervised pretraining on Kinetics. To the best of our knowledge, this is the first method to demonstrate that self-supervised video representation learning outperforms large-scale fully-supervised pretraining for action recognition. Moreover, unlike previous self-supervised methods that are only pretrained on curated data (*e.g.*, Kinetics [27] without action labels), we also report results of XDC pretrained on a large-scale uncurated video dataset.

## 2 Related work

**Early unsupervised representation learning.** Pioneering works include deep belief networks [21], autoencoders [22, 67], shift-invariant decoders [53], sparse coding algorithms [33], and stacked ISAs [32]. While these approaches learn by reconstructing the input, our approach learns from a self-supervised pretext task by generating pseudo-labels for supervised learning from unlabeled data.

**Self-supervised representation learning from images and videos.** Several pretext tasks exploit image spatial context, *e.g.*, by predicting the relative position of patches [8] or solving jigsaw puzzles [41]. Others include creating image classification pseudo-labels (*e.g.*, through artificial rotations [13] or clustering features [6]), colorization [81], inpainting [47], motion segmentation [46], and instance counting [42]. Some works have extended image pretext tasks to video [28, 72, 79]. Other video pretext tasks include frame ordering [9, 34, 39, 78], predicting flow or colors [31, 70], exploiting region correspondences across frames [23, 24, 75, 76], future frame prediction [36, 37, 60, 68, 69], and tracking [77]. Unlike this prior work, our model uses two modalities: video and audio.

**Cross-modal learning and distillation.** Several works [2, 17] train a fully-supervised encoder on one modality and distill its discriminative knowledge to an encoder of a different modality. Other works learn from unlabeled data for a specific target task [82, 55]. Unlike these methods, our work is purely self-supervised and aims at learning representations that transfer well to a wide range
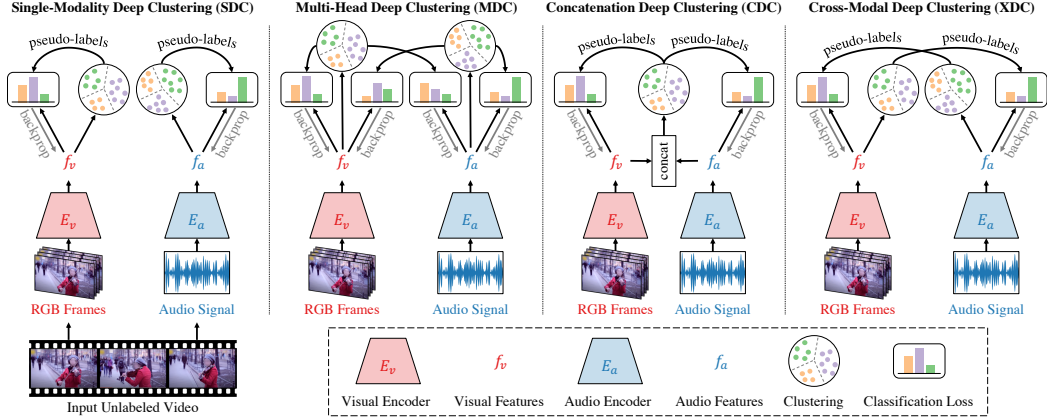
Figure 1: **Overview of our framework.** We present Single-Modality Deep Clustering (SDC) baseline vs. our three multi-modal deep clustering models: Multi-Head Deep Clustering (MDC), Concatenation Deep Clustering (CDC), and Cross-Modal Deep Clustering (XDC). The video and audio encoders ($E_v$ and $E_a$) map unlabeled videos to visual and audio features ($f_v$ and $f_a$). These features, or their concatenations, are clustered using $k$-means. The cluster assignments are then used as pseudo-labels to train the encoders. We start with randomly-initialized encoders, then alternates between clustering to generate pseudo-labels and training to improve the encoders. The four models employ different ways to cluster features and generate self-supervision signals. Illustration video is from [63].

of downstream tasks. Previous cross-modal self-supervised methods most relevant to our work include audio-visual correspondence [1], deep aligned representations [3], audio-visual temporal synchronization [29, 44], contrastive multiview coding [65], and learning image representations using ambient sound [45]. While [1, 3, 45, 65] use only a single frame, we use a video clip. Unlike our method, [45] clusters handcrafted audio features and does not iterate on the pseudo-labels. [29, 44] require constructing positive/negative examples for in-sync and out-of-sync video-audio pairs. This sampling strategy makes these approaches more difficult to scale compared to ours, as many potential out-of-sync pairs can be generated, yielding largely different results depending on the sampling choice [29]. Recent works, such as MIL-NCE [38] and CBT [64], learn from unlabeled instructional videos using text from ASR, while our approach makes use of the audio signal instead.

## 3 Technical approach

Here, we briefly discuss previous work on single-modality deep clustering in images [6]. Then, we introduce our three multi-modal deep clustering frameworks for representation learning (Figure 1).

### 3.1 Single-modality deep clustering

Caron *et al.* [6] proposed DeepCluster for self-supervised representation learning from images. DeepCluster iteratively clusters deep features from a single-modality encoder, and then uses the cluster assignments to train the same encoder to improve its representation. Inspired by the simplicity of this work, our paper studies deep clustering in the large-scale multi-modal setting. For completeness, we summarize DeepCluster details. Let $X$ be the set of unlabeled inputs (*e.g.*, images), $E$ be an encoder that maps an input $x \in X$ to a deep feature vector $f \in \mathbb{R}^d$. DeepCluster iterates between clustering the features $F = \{f = E(x) \mid x \in X\}$ and discriminative training to improve $E$ using the clustering assignments as pseudo-labels. The process starts with a randomly-initialized $E$, and only the weights of the classification `fc`-layer are reset between clustering iterations when the supervision-taxonomy is switched. DeepCluster uses a 2D CNN (*e.g.* ResNet-50) for $E$ and clusters the features after each epoch using $k$-means. We refer to DeepCluster as **Single-Modality Deep Clustering (SDC)**.

### 3.2 Multi-modal deep clustering

Contrary to the single-modality case, there exist multiple encoders in a multi-modal setting, each of which encodes a different modality of the input. In our paper, we consider two modalities, the visual and the audio modalities from the unlabeled training video clips. In particular, let $X$ be the set of unlabeled video clips, and $E_v$ and $E_a$ be the visual and audio encoders, respectively. Let $F_v = \{f_v = E_v(x) \in \mathbb{R}^{d_v} \mid x \in X\}$ and $F_a = \{f_a = E_a(x) \in \mathbb{R}^{d_a} \mid x \in X\}$ be the set of visual

and audio deep features produced by the two encoders, respectively. There are different ways we can adapt the deep clustering framework to a multi-modal input. We describe three approaches (MDC, CDC, and XDC) by detailing the steps taken at each deep clustering iteration. Refer to the *supplementary material* for the implementation differences between SDC and our three approaches.

**Multi-Head Deep Clustering (MDC).** This model builds on SDC by adding a second classification head supervised by the other modality. Thus, in this model, each encoder has two classification heads. At each deep clustering iteration, MDC uses the cluster assignments of $F_v$ as pseudo-labels for one head and that of $F_a$ as pseudo-labels for the other head. Thus, each encoder needs to predict the cluster assignments of its own modality (as in SDC), but also those generated by the other modality.

**Concatenation Deep Clustering (CDC).** This model performs clustering of joint visual and audio features. Specifically, at each deep clustering iteration, CDC clusters vectors obtained by concatenating the visual and audio feature vectors, separately $l_2$-normalized. Then, it uses the resulting cluster assignments as pseudo-labels to update the weights of both $E_v$ and $E_a$.

**Cross-Modal Deep Clustering (XDC).** Each encoder in this model relies exclusively on the clusters learned from the other modality as the supervisory signal. At each deep clustering iteration, XDC clusters the audio deep features, $F_a$, and uses their cluster assignments as pseudo-labels to train the visual encoder, $E_v$. Vice versa, XDC supervises $E_a$ with the cluster assignments of $F_v$.

# 4 Experiments

## 4.1 Experimental setup

**Pretraining datasets.** We use four datasets: Kinetics [27], AudioSet [10], IG-Kinetics [12], and IG-Random, which have 240K, 2M, 65M, and 65M training videos, respectively. As our approach is self-supervised, thus the labels from the first three datasets are **not used** during pretraining. While Kinetics and AudioSet are supervised benchmarks for action recognition and audio classification, IG-Kinetics is a weakly-supervised dataset collected from a social media website using tags related to Kinetics actions. IG-Random is an *uncurated* dataset of random videos from the same website. Videos are 10-second long in Kinetics and AudioSet and 10-to-60-second long in IG-Kinetics and IG-Random. We filter out around 7K Kinetics videos that have no audio. Furthermore, we randomly sample 240K videos from AudioSet and denote this subset as AudioSet-240K. We generate this subset to have AudioSet data of the same size as Kinetics, in order to study the effects of pretraining with the same data size but on a different data distribution and domain.

**Downstream datasets.** We evaluate our pretraining performance on three downstream benchmarks: UCF101 [59], HMBD51 [30], and ESC50 [49], which have 13K, 7K, and 2K examples from 101, 51, and 50 classes, respectively. UCF101 and HMDB51 are action recognition datasets, while ESC50 is a sound classification dataset. UCF101 and HMDB51 have 3 official train/test splits, while ESC50 has 5 splits. We conduct our ablation study (Subsection 4.2) using split-1 of each dataset. We also report our average performance over all splits when we compare with state-of-the-art methods in Section 6.

**Baselines.** We consider two baselines: *Scratch* and *Supervised Pretraining (Superv)*. The first is a randomly-initialized model trained from scratch directly on the downstream task, while the second is a model pretrained in a supervised fashion on a large labeled dataset (*e.g.*, Kinetics) and then finetuned on the downstream task. We note that these two baselines are commonly regarded as the lower and upper bounds to gauge the quality of self-supervised representation learning methods [1, 29].

**Backbone architectures.** We employ R(2+1)D [66] and ResNet [19] as $E_v$ and $E_a$, respectively. $E_v$'s input is a $3 \times L \times H \times W$ clip, where 3 refers to the RGB channels, $L$ is the number of frames, and $H$ and $W$ are the frame height and width. $E_a$'s input is a $Q \times P$ spectrogram image extracted from the audio signal, where $Q$ is the number of MEL filters and $P$ is the number of audio frames.

**Pretraining and evaluation details.** We choose the 18-layer variants of R(2+1)D and ResNet encoders. We use clips of $L=8$ frames for pretraining and finetuning our visual encoder $E_v$. We scale frames such that the smallest dimension is 256 pixels and then random crop images of size $224 \times 224$. We extract video clips at 30 fps and employ temporal jittering during training. For the audio input, we sample 2 seconds and use $Q=40$ MEL filters and $P=100$ audio frames. For inference on the downstream tasks, we uniformly sample 10 clips per testing example and average their predictions to make a video-level prediction. We use only one crop per clip: the center $8 \times 224 \times 224$ crop for video and the full $40 \times 100$ crop for audio. We provide more details in the *supplementary material*.

Table 1: **Single-modality vs. multi-modal deep clustering.** We compare the four self-supervised deep clustering models (Section 3) and the three baselines: Scratch, Supervised Pretraining (Superv), and same-modality-XDC (XDC with two encoders defined on the same modality). Models are pretrained via self-supervision on Kinetics and fully finetuned on each downstream dataset. We report the top-1 accuracy on split-1 of each dataset. All multi-modal models significantly outperform the single-modality deep clustering model. We mark in bold the best and underline the second-best models.

| Dataset | Scratch | Superv | SDC | MDC | CDC | XDC | same-modality-XDC | |
| | | | | | | | 2 visual encoders | 2 audio encoders |
|---|---|---|---|---|---|---|---|---|
| UCF101 | 54.5 | 90.9 | 61.8 | 68.4 | <u>72.9</u> | **74.2** | 61.3 | N/A |
| HMDB51 | 24.1 | 58.0 | 31.4 | 37.1 | <u>37.5</u> | **39.0** | 30.5 | N/A |
| ESC50 | 54.3 | 82.3 | 66.5 | 70.3 | <u>74.8</u> | **78.0** | N/A | 66.0 |

## 4.2 Ablation study

**Study 1: Single-modality vs. multi-modal deep clustering.** This experiment compares the four models presented in Section 3. We pretrain SDC, MDC, CDC, and XDC on Kinetics and report their performance on the downstream tasks in Table 1. To better understand XDC, we also conduct a new set of baselines, called same-modality-XDC, where XDC is trained with two encoders defined on the *same* modality (either visual or audio). Note that all models in this ablation study use the same visual and audio encoders and only differ in the way they use self-supervision. It takes on average 5 to 6 deep clustering iterations for these models to converge. ***Observations:*** **(I)** The four self-supervised deep clustering models outperform the Scratch baseline on all downstream benchmarks. This shows that our self-supervised pretraining is effective and generalizes well to multiple tasks. **(II)** All multi-modal models (MDC, CDC, and XDC) significantly outperform SDC by up to 12.4%, 7.6%, and 11.5% on UCF101, HMDB51, and ESC50, respectively. This validates the importance of multi-modal modeling compared to single-modality. **(III)** XDC achieves the best performance across all tasks. What distinguishes XDC from the other models is that each modality encoder in XDC is self-supervised purely by the signal from the other modality. The encoders in CDC, MDC, and SDC all employ a self-supervision signal coming from the same modality. Thus, this suggests that encoders learn better when purely supervised by a different modality. We provide the following intuition on why XDC is better than CDC and MDC. XDC groups samples together when they are similar in one of the two modalities (video to supervise the audio encoder, audio to supervise the visual encoder). Instead, CDC groups samples together only if they are similar according to both the audio *and* the video modality (to supervise both encoders). Thus, XDC visual and audio clusters allow for more diversity than those of CDC. We hypothesize that this diversity allows XDC to learn richer representations, which translates into better performance on the downstream tasks. Also, recent work [74] has shown that models trained on different modalities learn and generalize at different speeds, and that training them jointly (as done in MDC which uses two-modality heads) is sub-optimal. We believe that this could contribute to MDC performing worse than XDC, which optimizes for each modality independently. **(IV)** The same-modality-XDC baselines perform similarly to SDC and are 8-12% worse than multi-modal-XDC. This suggests that cross-modality provides a superior supervisory signal for self-supervised learning and that multi-modal-XDC is the best model not because of its optimization strategy but rather because of the use of the other modality for pseudo-labeling. Given the results of this study, we opt to use only XDC in the rest of the experiments. Finally, to show that XDC works for different backbones, we re-do Study 1 with ResNet3D in the *supplementary material*.

**Study 2: The number of $k$-means clusters.** This study explores the effects of changing the hyperparameter $k$ in $k$-means clustering. We pretrain XDC on three datasets, Kinetics, AudioSet-240K, and AudioSet, using $k$=64, 128, 256, 512, and 1024 clusters (Table 2). ***Observations:*** **(I)** The best $k$ value is not sensitive to the number of semantic labels in the downstream datasets. For example, HMDB51 and ESC50 have about the same number of labels but different best $k$ value. **(II)** Similarly, the best $k$ value seems uncorrelated with the number of original semantic labels of the pretraining dataset, *e.g.* 400 in Kinetics. We reiterate here that our approach is self-supervised and **does not use** the labels of the pretraining dataset. **(III)** The best $k$ value tends to get larger as the pretraining data size increases. For example, the best $k$ for HMDB51 shifts from 128 to 256 when moving from pretraining on AudioSet-240K to the full AudioSet. We hypothesize that there is a more diverse sample set to cluster when the pretraining data size increases. Thus, we can have more fine-grained clusters (higher $k$) and make our self-supervised classification problem harder. This aligns with previous self-supervised works [15, 29] that showed benefits from making the pretext task harder.

Table 2: **The number of clusters** ($k$). We show the effect of the number of $k$-means clusters on XDC performance. XDC is pretrained on three large datasets, and then fully finetuned on three downstream tasks. We report the top-1 accuracy on split-1. The best $k$ value increases as the size of the pretraining dataset increases.

| Pretraining Dataset | Downstream Dataset | $k$ | | | | |
|---|---|---|---|---|---|---|
| | | 64 | 128 | 256 | 512 | 1024 |
| Kinetics (240K videos) | UCF101 | 73.8 | 73.1 | **74.2** | <u>74.0</u> | 72.6 |
| | HMDB51 | 36.5 | **39.0** | <u>38.3</u> | 37.7 | 37.7 |
| | ESC50 | **78.0** | <u>76.3</u> | 75.0 | 74.5 | 71.5 |
| AudioSet-240K (240K videos) | UCF101 | **77.4** | <u>77.2</u> | 76.7 | 77.1 | 75.3 |
| | HMDB51 | 41.3 | **42.6** | <u>41.6</u> | 40.6 | 40.7 |
| | ESC50 | **78.5** | <u>77.8</u> | 77.3 | 76.8 | 73.5 |
| AudioSet (2M videos) | UCF101 | 84.1 | 84.3 | **84.9** | <u>84.4</u> | 84.2 |
| | HMDB51 | 47.4 | 47.6 | **48.8** | <u>48.5</u> | 48.4 |
| | ESC50 | 84.8 | **85.8** | <u>85.0</u> | 84.5 | 83.0 |

Table 3: **Pretraining data type and size.** We compare XDC pretrained on five datasets vs. fully-supervised pretrained baselines (Superv). XDC significantly outperforms fully-supervised pretraining on HMDB51.

| | Pretraining | | Downstream Dataset | | |
|---|---|---|---|---|---|
| Method | Dataset | Size | UCF101 | HMDB51 | ESC50 |
| Scratch | None | 0 | 54.5 | 24.1 | 54.3 |
| Superv | ImageNet | 1.2M | 79.9 | 44.5 | NA |
| Superv | Kinetics | 240K | <u>90.9</u> | 58.0 | 82.3 |
| Superv | AudioSet-240K | 240K | 76.6 | 40.8 | 78.3 |
| Superv | AudioSet | 2M | 84.0 | 53.5 | **90.3** |
| XDC | Kinetics | 240K | 74.2 | 39.0 | 78.0 |
| XDC | AudioSet-240K | 240K | 77.4 | 42.6 | 78.5 |
| XDC | AudioSet | 2M | 84.9 | 48.8 | 85.8 |
| XDC | IG-Random | 65M | 88.8 | <u>61.2</u> | <u>86.3</u> |
| XDC | IG-Kinetics | 65M | **91.5** | **63.1** | 84.8 |

**Study 3: Pretraining data type and size.** Here, we investigate the effects of two pretraining characteristics: data size and type. To this end, we pretrain XDC on Kinetics (240K examples), AudioSet-240K (240K examples), AudioSet (2M examples), IG-Kinetics (65M examples), and IG-Random (65M examples). Kinetics and IG-Kinetics videos are collected originally for activity recognition, while AudioSet videos are aimed for audio event classification. IG-Random is an uncurated/unsupervised dataset. In addition to video datasets, we also experiment with ImageNet to understand how much action recognition benefits from supervised pretraining on object classification. For ImageNet, we inflate the images into static video clips (repeating the same frame) and pretrain our video model on this dataset. Table 3 presents the results of this study. ***Observations:*** **(I)** XDC improves across all three downstream tasks as the pretraining data size increases. For example, XDC on HMDB51 improves by $9.8\%$, $22.2\%$, and $24.1\%$ when pretrained on AudioSet, IG-Random, and IG-Kinetics, respectively, compared to the results when pretrained on Kinetics. **(II)** XDC outperforms Kinetics fully-supervised pretraining by $5.1\%$ on HMDB51 and by $0.6\%$ on UCF101. To the best of our knowledge, XDC is the first method to demonstrate that self-supervision can outperform large-scale full-supervision in representation learning for action recognition. **(III)** The performance of the fully-supervised pretrained model is influenced by the taxonomy of the pretraining data more than the size. For example, supervised-pretraining on Kinetics gives better performance on both UCF101 and HMDB51 compared to supervised-pretraining on AudioSet (which is $8$ times larger than Kinetics) and ImageNet. One the other hand, XDC performance is less sensitive to the data type, as it implicitly learns the label space rather than depend on a space manually defined by annotators.

**Study 4: Curated vs. uncurated pretraining data.** The overarching goal of self-supervised representation learning is to learn from the massive amounts of unlabeled data. Previous self-supervised methods have pretrained on videos from supervised (curated) datasets (*e.g.*, Kinetics) without using the labels. However, even without using labels, those videos are still biased due to the sampling distribution (*e.g.*, taxonomy of the curated dataset). To this end, we study the effects of self-supervised representation learning from uncurated data. Table 4 compares XDC pretrained on IG-Kinetics (curated, as videos were tag-retrieved) vs. IG-Random (uncurated) using 1M, 16M, and 65M videos. ***Observations:*** **(I)** Curated pretraining gives better results on UCF101 and HMDB51, while uncurated pretraining is better on ESC50 at large scale. We hypothesize that the bias of IG-Kinetics towards semantics of human actions is the reason behind the positive effect of curation on

Table 4: **Curated vs. uncurated pretraining data.** XDC pretrained on IG-Kinetics (curated) vs. IG-Random (uncurated) using different training set sizes. Uncurated pretraining has better results on ESC at large scale. On UCF and HMDB, the accuracy gap between curated and uncurated pretraining decreases as data size increases.

| Downstream | Pretraining Size | | | | | | | | |
| Dataset | **1M** | | | **16M** | | | **65M** | | |
| | IG-Random | IG-Kinetics | $\Delta$ | IG-Random | IG-Kinetics | $\Delta$ | IG-Random | IG-Kinetics | $\Delta$ |
| UCF101 | 79.6 | **84.2** | -4.6 | 84.1 | **87.6** | -3.5 | 88.8 | **91.5** | -2.7 |
| HMDB51 | 45.1 | **50.3** | -5.2 | 55.2 | **57.3** | -2.1 | 61.2 | **63.1** | -1.9 |
| ESC50 | 77.8 | **79.5** | -1.7 | **84.3** | 82.5 | +1.8 | **86.3** | 84.8 | +1.5 |

Table 5: **Full finetuning vs. learning `fc`-only.** We compare XDC against the supervised pretrained models (Superv) under two transfer-learning schemes: when models are used as features extractor ('`fc`' column) or as a finetuning initialization ('all' column). XDC fixed features outperform several fully-finetuned Superv models.

| Method | Pretraining Dataset | UCF101 | | HMDB51 | | ESC50 | |
| | | `fc` | all | `fc` | all | `fc` | all |
| Random | None | 6.0±1.0 | 54.5 | 7.5±0.6 | 24.1 | 61.3±2.5 | 54.3 |
| Superv | ImageNet | 74.5 | 79.9 | 42.8 | 44.5 | NA | NA |
| Superv | Kinetics | **89.7** | <u>90.9</u> | **61.5** | 58.0 | 79.5 | 82.3 |
| Superv | AudioSet | 80.2 | 84.0 | 51.6 | 53.5 | **88.5** | **90.3** |
| XDC | IG-Random | 80.7 | 88.8 | 49.9 | <u>61.2</u> | <u>84.5</u> | <u>86.3</u> |
| XDC | IG-Kinetics | <u>85.3</u> | **91.5** | <u>56.0</u> | **63.1** | 84.3 | 84.8 |

UCF101 and HMDB51. However, such bias negatively impacts the performance on ESC50. **(II)** The performance gap between the curated and uncurated pretraining shrinks significantly as we increase the data size. For example, the performance gap on HMDB51 drops from $5.2\%$ to $2.1\%$ and $1.9\%$ when the pretraining size increases from 1M to 16M and 65M videos, respectively. This implies that XDC can learn meaningful representations from truly uncurated data. To the best of our knowledge, XDC is the first self-supervised method to study pretraining on large-scale uncurated video data.

**Study 5: Full finetuning vs. learning `fc`-only.** Here, we study two approaches for transferring XDC to downstream tasks. *Full finetuning*: we finetune all parameters of the pretrained encoder on the downstream task. *Learning `fc`-only*: we fix the pretrained encoder and learn a linear classifier for the downstream task, *i.e.*, a fully-connected (`fc`) layer on top of the frozen features. Table 5 compares XDC with the supervised pretrained approaches under these two transfer-learning schemes. ***Observations:*** **(I)** The accuracy of most pretrained models (fully-supervised or self-supervised) degrades, when used as a fixed feature extractor compared to when they are fully-finetuned on the downstream tasks. Nonetheless, the relative performance of XDC compared to supervised pretrained models stays generally the same when fully vs. `fc`-only finetuned on the downstream task. This suggests that XDC pretraining is useful both as a fixed feature extractor and as a pretraining initialization. **(II)** XDC as a fixed feature extractor outperforms many fully-finetuned supervised models. For example, `fc`-only XDC outperforms, by significant margins, the fully-finetuned supervised AudioSet- and ImageNet-pretrained models on both UCF101 and HMDB51. **(III)** We observe that fully-supervised pretraining, followed by `fc`-only finetuning performs well when the pretraining taxonomy is well aligned with that of the downstream task. For example, pretraining on Kinetics by learning `fc`-only on HMDB51 and UCF101 gives the best performance. This is expected as the label spaces of HMBD51 and UCF101 overlap largely with that of Kinetics. This suggests that fully-supervised pretraining is more taxonomy/downstream-task dependent, while our self-supervised XDC is taxonomy-independent.

# 5 Understanding XDC

What does XDC actually learn? What semantic signals does the algorithm use to train its encoders? Here, we try to answer these questions by inspecting the $k$-means clustering results produced by the last iteration of XDC. Figure 2 visualizes some audio and video clusters learned by XDC when trained on Kinetics. These clusters are the top 2 audio clusters (left) and the top 2 video clusters (right) ranked by purity *w.r.t.* Kinetics action labels. More clusters are presented in Table 6. We observe that the top-purity clusters learned from both modalities exhibit strong semantic coherence. For example, the audio 1st and 8th ranked clusters include concepts related to playing musical instruments that have similar sounds, while the 1st ranked video cluster also groups playing-instrument concepts, but mainly because of their appearance, as the cluster is all about guitars. Other interesting clusters include: grouping by motor-engine sounds (audio #10), by different swimming strokes (video #4), by
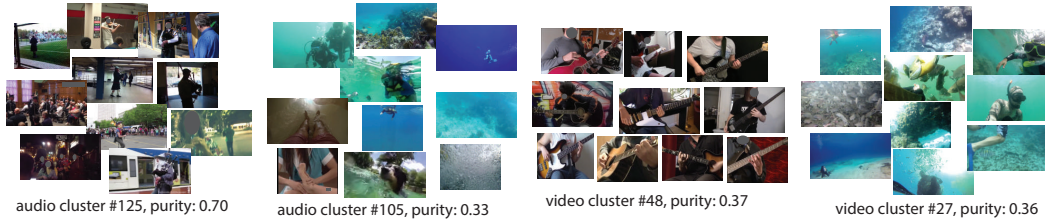
Figure 2: **Visualization of XDC clusters on Kinetics videos**. The top-2 audio clusters (left) and video clusters (right) in terms of purity *w.r.t.* the Kinetics labels. Clusters are represented by the 10 closest videos (shown as frames) to their centroid. Interestingly, XDC learned to group "scuba diving" with "snorkeling" (second left, cluster #105) based on audio features and "scuba diving" with "feeding fish" (rightmost, cluster #27) based on visual features. Please refer to our project website for an interactive visualization of all XDC clusters.

Table 6: **XDC clusters.** Top and bottom audio (left) and video (right) XDC clusters ranked by clustering purity *w.r.t.* Kinetics labels. For each cluster, we list the three concepts with the highest purity (given in parentheses).

| # | Kinetics concepts | # | Kinetics concepts |
|---|---|---|---|
| 1 | play bagpipes(0.70), play harmonica(0.04), play violin(0.03) | 1 | play bass guitar(0.37), play guitar(0.16), tap guitar(0.15) |
| 2 | scuba diving(0.33), snorkeling(0.27), feeding fish(0.11) | 4 | swim backstroke(0.21), breast stroke(0.16), butterfly stroke(0.1) |
| 8 | play cello(0.15), play trombone(0.11), play accordion(0.09) | 5 | golf putting(0.18), golf chipping(0.11), golf driving(0.05) |
| 10 | mowing lawn(0.14), driving tractor(0.09), motorcycling(0.06) | 10 | cook chicken(0.11), barbeque(0.07), fry vegetables(0.06) |
| 127 | abseiling(0.01), grooming horse(0.01), milking cow(0.01) | 63 | pull ups(0.01), gymnastics tumbling(0.01), punching bag(0.01) |
| 128 | washing feet(0.01), motorcycling(0.01), headbanging(0.01) | 64 | capoeira(0.01), riding elephant(0.01), feeding goats(0.01) |

different golf shots (video #5), and different cooking activities (video #10). In the bottom-ranked clusters, although the purity *w.r.t.* Kinetics concepts is low, we still find some coherence, mostly at the scene level: a farm setting in audio #127 ("grooming horse", "milking cow") and gym activities in video #63 ("pull ups", "punching bag"). Many other bottom-ranked clusters appear to lack semantic coherence when viewed through the lens of Kinetics labels. However, one of the motivations behind the design of self-supervised methods is precisely to bypass the hand-design of label spaces, which may not be the optimal ones for general representation learning. Our experiments suggest that the label space learned by XDC yields strong and general audio and video features even though it does not align perfectly with the taxonomies of existing datasets.

# 6 State-of-the-art self-supervised learning comparison

**Experimental setup.** Here, training is similar to our ablations except that we re-train our video encoder on the last clustering assignment using 32-frame clips. Then following [29, 66], we finetune on UCF101 and HMDB51 using 32-frame clips for both XDC and the fully-supervised baselines. Inference is similar to our ablations except for using 32-frame clips. For the audio event classification dataset DCASE [62], we follow [29] and extract `conv_5` features for 60 uniformly-sampled clips per audio sample and learn a linear SVM. We report the average top-1 accuracy over *all splits*.

**Video action recognition.** Table 7(a) compares XDC pretrained on four large-scale datasets against state-of-the-art self-supervised methods, after finetuning on the UCF101 and HMDB51 benchmarks[2]. We also compare against two fully-supervised methods pretrained on ImageNet and Kinetics. ***Results:*** **(I)** XDC pretrained on IG-Kinetics sets new state-of-the-art performance for self-supervised methods on both benchmarks, outperforming Elo [50] by 1.7% on UCF101 and 1.5% on HMDB51. Moreover, XDC significantly outperforms fully-supervised pretraining on Kinetics: by 1.3% on UCF101 and by 3.8% on HMDB51. **(II)** When directly compared on the same R(2+1)D-18 architecture, XDC pretrained on Kinetics slightly outperforms AVTS [29] by 0.6% on UCF101 and 0.3% on HMDB51. However, when both methods are pretrained on AudioSet, XDC outperforms AVTS with larger margins: by 3.9% on UCF101 and 5.6% on HMDB51. This shows that XDC scales better than AVTS. To further verify that XDC scales better, we pretrained AVTS on AudioSet-240K using R(2+1)D-18 and got 76.9% and 40.7% for UCF101 and HMDB51 on split-1, showing a smaller margin between XDC and AVTS than when both are pretrained on the full AudioSet (cf. Table 3).

**Audio event classification.** Table 7(b) compares XDC pretrained on AudioSet and IG-Random against the state-of-the-art self-supervised methods for audio classification. XDC achieves state-of-the-art performance on DCASE and competitive results on ESC50 with only a 1.1% gap with [56].

---

[2]All XDC pretrained models are publicly released on our project website.

Table 7: **State-of-the-art comparison.** We report the average top-1 accuracy over the official splits for all benchmarks. **(a) Video action recognition:** Comparison between XDC with self-supervised and fully-supervised methods on UCF101 and HMDB51. Not only does XDC set new state-of-the-art performance for self-supervised methods, it also outperforms fully-supervised Kinetics and ImageNet pretraining. * For fair comparison with XDC, we report AVTS performance without dense prediction, *i.e.*, we average the predictions of 10 uniformly-sampled clips at inference. † For direct comparison with XDC, we evaluate AVTS using R(2+1)D-18 and 10 uniformly-sampled clips at inference. **(b) Audio event classification:** We compare XDC with self-supervised methods on ESC50 and DCASE. XDC achieves state-of-the-art performance on DCASE.

(a) Video action recognition.

| Method | Pretraining Architecture | Dataset | Evaluation UCF101 | HMDB51 |
|---|---|---|---|---|
| ClipOrder [79] | R(2+1)D-18 | UCF101 | 72.4 | 30.9 |
| MotionPred [72] | C3D | Kinetics | 61.2 | 33.4 |
| ST-Puzzle [28] | 3D-ResNet18 | Kinetics | 65.8 | 33.7 |
| DPC [18] | 3D-ResNet34 | Kinetics | 75.7 | 35.7 |
| CBT [64] | S3D | Kinetics | 79.5 | 44.6 |
| SpeedNet [4] | S3D | Kinetics | 81.1 | 48.8 |
| AVTS [29]* | MC3-18 | Kinetics | 84.1 | 52.5 |
| AVTS [29]† | R(2+1)D-18 | Kinetics | 86.2 | 52.3 |
| **XDC** (ours) | R(2+1)D-18 | Kinetics | 86.8 | 52.6 |
| AVTS [29]* | MC3-18 | AudioSet | 87.7 | 57.3 |
| AVTS [29]† | R(2+1)D-18 | AudioSet | 89.1 | 58.1 |
| **XDC** (ours) | R(2+1)D-18 | AudioSet | 93.0 | 63.7 |
| MIL-NCE [38] | S3D | HowTo100M | 91.3 | 61.0 |
| ELo [50] | R(2+1)D-50 | YouTube-8M | 93.8 | 67.4 |
| **XDC** (ours) | R(2+1)D-18 | IG-Random | 94.6 | 66.5 |
| **XDC** (ours) | R(2+1)D-18 | IG-Kinetics | **95.5** | **68.9** |
| Fully supervised | R(2+1)D-18 | ImageNet | 84.0 | 48.1 |
| Fully supervised | R(2+1)D-18 | Kinetics | 94.2 | 65.1 |

(b) Audio event classification.

| Method | ESC50 |
|---|---|
| Random Forest [49] | 44.3 |
| Piczak ConvNet [48] | 64.5 |
| SoundNet [2] | 74.2 |
| $L^3$-Net [1] | 79.3 |
| AVTS [29] | 82.3 |
| ConvRBM [56] | **86.5** |
| **XDC** (AudioSet) | 84.8 |
| **XDC** (IG-Random) | 85.4 |

| Method | DCASE |
|---|---|
| RG [52] | 69 |
| LTT [35] | 72 |
| RNH [54] | 77 |
| Ensemble [61] | 78 |
| SoundNet [2] | 88 |
| $L^3$-Net [1] | 93 |
| AVTS [29] | 94 |
| **XDC** (AudioSet) | **95** |
| **XDC** (IG-Random) | **95** |

# 7 XDC for temporal action localization

In this section, we further demonstrate that XDC can be useful beyond video and audio classification. In particular, we employ the recent G-TAD [80] action localization algorithm, where we replace the clip features (originally extracted from a TSN [73] model pretrained on Kinetics) with our XDC features from the R(2+1)D-18 model pretrained on IG-Kinetics or IG-Random. We compare against the features from the R(2+1)D-18 model fully-supervised pretrained on Kinetics. We emphasize that we do not finetune any of the feature extractors used in this experiment. We follow the default hyperparameters setting of G-TAD. Table 8 shows temporal action localization results of G-TAD with different features on THUMOS14 [25] dataset. It reports the mean Average Precision (mAP) results at different temporal Intersection over Union (tIoU) thresholds. Both XDC variants outperform the fully-supervised features across all tIoU thresholds. This confirms the same trend observed in tasks presented in Section 6 and suggests that XDC can be used for other tasks.

Table 8: **Temporal action localization on THUMOS14.** We compare G-TAD algorithm using our XDC features vs. using the fully-supervised Kinetics-pretrained (Superv) features. We report the mean Average Precision (mAP) results at different temporal Intersection over Union (tIoU) thresholds. Both XDC variants outperform the fully-supervised features across all tIoU thresholds.

| Features Type | mAP @ tIoU 0.3 | 0.4 | 0.5 | 0.6 | 0.7 |
|---|---|---|---|---|---|
| Superv (Kinetics) | 50.9 | 44.4 | 36.6 | 28.4 | 19.8 |
| XDC (IG-Random) | **51.5** | 44.8 | 36.9 | 28.6 | **20.0** |
| XDC (IG-Kinetics) | **51.5** | **44.9** | **37.2** | **28.7** | **20.0** |

# 8 Conclusion

We presented Cross-Modal Deep Clustering (XDC), a novel self-supervised model for video and audio. XDC outperforms not only existing self-supervised methods but also fully-supervised ImageNet- and Kinetics-pretraining for action recognition. To the best of our knowledge, XDC is the first to show self-supervision outperforming large-scale full-supervision pretraining for action recognition when pretrained on the same architecture and a larger number of uncurated videos.

## Broader Impact Statement

Video has become a commonplace in society. Its uses range from entertainment, to communication and teaching. Thus, the learning of semantic representations of video has broad and far-reaching potential applications. The authors do not foresee major ethical issues associated to this work. However, as the proposed approach is self-supervised, it will learn the inherent properties and structure of the training data. Thus, the learned model may exhibit biases intrinsically present in the data.

## Acknowledgments

## References

[1] Relja Arandjelovic and Andrew Zisserman. Look, listen and learn. In *ICCV*, 2017. 3, 4, 9

[2] Yusuf Aytar, Carl Vondrick, and Antonio Torralba. Soundnet: Learning sound representations from unlabeled video. In *NeurIPS*, 2016. 2, 9

[3] Yusuf Aytar, Carl Vondrick, and Antonio Torralba. See, hear, and read: Deep aligned representations. *arXiv preprint arXiv:1706.00932*, 2017. 3

[4] Sagie Benaim, Ariel Ephrat, Oran Lang, Inbar Mosseri, William T Freeman, Michael Rubinstein, Michal Irani, and Tali Dekel. Speednet: Learning the speediness in videos. In *CVPR*, 2020. 9

[5] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *CVPR*, 2015. 2

[6] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In *ECCV*, 2018. 2, 3, 14

[7] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. Scaling egocentric vision: The epic-kitchens dataset. In *ECCV*, 2018. 2

[8] Carl Doersch, Abhinav Gupta, and Alexei A Efros. Unsupervised visual representation learning by context prediction. In *ICCV*, 2015. 2

[9] Basura Fernando, Hakan Bilen, Efstratios Gavves, and Stephen Gould. Self-supervised video representation learning with odd-one-out networks. In *CVPR*, 2017. 2

[10] Jort F. Gemmeke, Daniel P. W. Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R. Channing Moore, Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for audio events. In *ICASSP*, 2017. 4, 15

[11] A Gentile and S DiFrancesca. Academic achievement test performance of hearing-impaired students. united states, spring, 1969.(series d, no. 1). washington, dc: Gallaudet university. *Center for Assessment and Demographic Studies*, 1969. 1

[12] Deepti Ghadiyaram, Matt Feiszli, Du Tran, Xueting Yan, Heng Wang, and Dhruv Mahajan. Large-scale weakly-supervised pre-training for video action recognition. In *CVPR*, 2019. 4, 15

[13] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. *ICLR*, 2018. 2

[14] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch SGD: training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017. 14, 15

[15] Priya Goyal, Dhruv Mahajan, Abhinav Gupta, and Ishan Misra. Scaling and benchmarking self-supervised visual representation learning. In *ICCV*, 2019. 5

[16] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The "something something" video database for learning and evaluating visual common sense. In *ICCV*, 2017. 2

[17] Saurabh Gupta, Judy Hoffman, and Jitendra Malik. Cross modal distillation for supervision transfer. In *CVPR*, 2016. 2

[18] Tengda Han, Weidi Xie, and Andrew Zisserman. Video representation learning by dense predictive coding. In *ICCVW*, 2019. 9

[19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 4

[20] Rickye S. Heffner and Henry E. Heffner. *Evolution of Sound Localization in Mammals*, pages 691–715. Springer New York, New York, NY, 1992. 1

[21] Geoffrey E Hinton, Simon Osindero, and Yee-Whye Teh. A fast learning algorithm for deep belief nets. *Neural computation*, 2006. 2

[22] Geoffrey E Hinton and Ruslan R Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 2006. 2

[23] Phillip Isola, Daniel Zoran, Dilip Krishnan, and Edward H Adelson. Learning visual groups from co-occurrences in space and time. *ICLR*, 2015. 2

[24] Dinesh Jayaraman and Kristen Grauman. Slow and steady feature analysis: higher order temporal coherence in video. In *CVPR*, 2016. 2

[25] Y. Jiang, J. Liu, A. Zamir, G. Toderici, I. Laptev, M. Shah, and R. Sukthankar. THUMOS challenge: Action recognition with a large number of classes, 2014. 9

[26] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. Large-scale video classification with convolutional neural networks. In *CVPR*, 2014. 2

[27] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijaya-narasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, Mustafa Suleyman, and Andrew Zisserman. The kinetics human action video dataset. *CoRR*, abs/1705.06950, 2017. 2, 4, 15

[28] Dahun Kim, Donghyeon Cho, and In So Kweon. Self-supervised video representation learning with space-time cubic puzzles. In *AAAI*, 2019. 2, 9

[29] Bruno Korbar, Du Tran, and Lorenzo Torresani. Cooperative learning of audio and video models from self-supervised synchronization. In *NeurIPS*, 2018. 3, 4, 5, 8, 9

[30] H. Kuehne, H. Jhuang, E. Garrote, T. Poggio, and T. Serre. HMDB: a large video database for human motion recognition. In *ICCV*, 2011. 4

[31] Zihang Lai and Weidi Xie. Self-supervised learning for video correspondence flow. *BMVC*, 2019. 2

[32] Quoc V Le, Will Y Zou, Serena Y Yeung, and Andrew Y Ng. Learning hierarchical invariant spatio-temporal features for action recognition with independent subspace analysis. In *CVPR*, 2011. 2

[33] Honglak Lee, Alexis Battle, Rajat Raina, and Andrew Y Ng. Efficient sparse coding algorithms. In *NeurIPS*, 2007. 2

[34] Hsin-Ying Lee, Jia-Bin Huang, Maneesh Singh, and Ming-Hsuan Yang. Unsupervised representation learning by sorting sequences. In *ICCV*, 2017. 2

[35] David Li, Jason Tam, and Derek Toub. Auditory scene classification using machine learning techniques. *AASP Challenge*, 2013. 9

[36] William Lotter, Gabriel Kreiman, and David Cox. Deep predictive coding networks for video prediction and unsupervised learning. *ICLR*, 2017. 2

[37] Michael Mathieu, Camille Couprie, and Yann LeCun. Deep multi-scale video prediction beyond mean square error. *ICLR*, 2016. 2

[38] Antoine Miech, Jean-Baptiste Alayrac, Lucas Smaira, Ivan Laptev, Josef Sivic, and Andrew Zisserman. End-to-end learning of visual representations from uncurated instructional videos. In *CVPR*, 2020. 3, 9

[39] Ishan Misra, C Lawrence Zitnick, and Martial Hebert. Shuffle and learn: unsupervised learning using temporal order verification. In *ECCV*, 2016. 2

[40] Helmer R Myklebust. *The psychology of deafness: Sensory deprivation, learning, and adjustment.* Grune & Stratton, 1960. 1

[41] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *ECCV*, 2016. 2

[42] Mehdi Noroozi, Hamed Pirsiavash, and Paolo Favaro. Representation learning by learning to count. In *ICCV*, 2017. 2

[43] Risto Näätänen. *Attention and Brain Function*. Lawrence Erlbaum Associates, Inc, 1992. 1

[44] Andrew Owens and Alexei A Efros. Audio-visual scene analysis with self-supervised multisensory features. In *ECCV*, 2018. 3

[45] Andrew Owens, Jiajun Wu, Josh H McDermott, William T Freeman, and Antonio Torralba. Ambient sound provides supervision for visual learning. In *ECCV*, 2016. 3

[46] Deepak Pathak, Ross Girshick, Piotr Dollár, Trevor Darrell, and Bharath Hariharan. Learning features by watching objects move. In *CVPR*, 2017. 2

[47] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *CVPR*, 2016. 2

[48] Karol J. Piczak. Environmental sound classification with convolutional neural networks. *MLSP*, 2015. 9

[49] Karol J. Piczak. Esc: Dataset for environmental sound classification. In *ACM Multimedia*, 2015. 4, 9

[50] AJ Piergiovanni, Anelia Angelova, and Michael S Ryoo. Evolving losses for unsupervised video representation learning. In *CVPR*, 2020. 8, 9

[51] Jordi Pons and Xavier Serra. Randomly weighted cnns for (music) audio classification. In *ICASSP*, 2019. 14

[52] Alain Rakotomamonjy and Gilles Gasso. Histogram of gradients of time-frequency representations for audio scene classification. *IEEE/ACM Trans. Audio, Speech and Lang. Proc.*, 2015. 9

[53] Marc'aurelio Ranzato, Fu Jie Huang, Y-Lan Boureau, and Yann LeCun. Unsupervised learning of invariant feature hierarchies with applications to object recognition. In *CVPR*, 2007. 2

[54] Guido Roma, Waldo Nogueira, and Perfecto Herrera. Recurrence quantification analysis features for environmental sound recognition. *WASPAA*, 2013. 9

[55] Andrew Rouditchenko, Hang Zhao, Chuang Gan, Josh McDermott, and Antonio Torralba. Self-supervised audio-visual co-segmentation. In *ICASSP*, 2019. 2

[56] Hardik B. Sailor, Dharmesh M Agrawal, and Hemant A Patil. Unsupervised filterbank learning using convolutional restricted boltzmann machine for environmental sound classification. In *INTERSPEECH*, 2017. 8, 9

[57] Andrew M Saxe, Pang Wei Koh, Zhenghao Chen, Maneesh Bhand, Bipin Suresh, and Andrew Y Ng. On random weights and unsupervised feature learning. In *ICML*, 2011. 14

[58] Ladan Shams and Robyn Kim. Crossmodal influences on visual perception. *Physics of Life Reviews*, 7(3):269–284, 2010. 1

[59] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. UCF101: A dataset of 101 human action classes from videos in the wild. In *CRCV-TR-12-01*, 2012. 4

[60] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhudinov. Unsupervised learning of video representations using lstms. In *ICML*, 2015. 2

[61] Dan Stowell, Dimitrios Giannoulis, Emmanouil Benetos, Mathieu Lagrange, and Mark D. Plumbley. Detection and classification of acoustic scenes and events. *TM*, 2015. 9

[62] D. Stowell, D. Giannoulis, E. Benetos, M. Lagrange, and M. D. Plumbley. Detection and classification of acoustic scenes and events. *IEEE Transactions on Multimedia*, 17(10):1733–1746, Oct 2015. 8

[63] Street Performers. This violinist plays beautiful music St Petersburg, Russia. *https://www.youtube.com/watch?v=yLw2Sq8Bz78*, https://creativecommons.org/licenses/by/3.0. 3

[64] Chen Sun, Fabien Baradel, Kevin Murphy, and Cordelia Schmid. Learning video representations using contrastive bidirectional transformer. *arXiv preprint arXiv:1906.05743*, 2019. 3, 9

[65] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. In *ECCV*, 2020. 3

[66] Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A closer look at spatiotemporal convolutions for action recognition. In *CVPR*, 2018. 4, 8

[67] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *ICML*, 2008. 2

[68] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Anticipating visual representations from unlabeled video. In *CVPR*, 2016. 2

[69] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating videos with scene dynamics. In *NeurIPS*, 2016. 2

[70] Carl Vondrick, Abhinav Shrivastava, Alireza Fathi, Sergio Guadarrama, and Kevin Murphy. Tracking emerges by colorizing videos. In *ECCV*, 2018. 2

[71] Heng Wang and Cordelia Schmid. Action recognition with improved trajectories. In *ICCV*, 2013. 14

[72] Jiangliu Wang, Jianbo Jiao, Linchao Bao, Shengfeng He, Yunhui Liu, and Wei Liu. Self-supervised spatio-temporal representation learning for videos by predicting motion and appearance statistics. In *CVPR*, 2019. 2, 9

[73] Limin Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, Dahua Lin, Xiaoou Tang, and Luc Van Gool. Temporal segment networks: Towards good practices for deep action recognition. In *ECCV*, 2016. 9

[74] Weiyao Wang, Du Tran, and Matt Feiszli. What makes training multi-modal classification networks hard? In *CVPR*, 2020. 5

[75] Xiaolong Wang and Abhinav Gupta. Unsupervised learning of visual representations using videos. In *ICCV*, 2015. 2

[76] Xiaolong Wang, Allan Jabri, and Alexei A Efros. Learning correspondence from the cycle-consistency of time. In *CVPR*, 2019. 2

[77] Xiaolong Wang, Allan Jabri, and Alexei A. Efros. Learning correspondence from the cycle-consistency of time. In *CVPR*, 2019. 2

[78] Donglai Wei, Joseph J Lim, Andrew Zisserman, and William T Freeman. Learning and using the arrow of time. In *CVPR*, 2018. 2

[79] Dejing Xu, Jun Xiao, Zhou Zhao, Jian Shao, Di Xie, and Yueting Zhuang. Self-supervised spatiotemporal learning via video clip order prediction. In *CVPR*, 2019. 2, 9

[80] Mengmeng Xu, Chen Zhao, David S. Rojas, Ali Thabet, and Bernard Ghanem. G-tad: Sub-graph localization for temporal action detection. In *CVPR*, 2020. 9

[81] Richard Zhang, Phillip Isola, and Alexei A Efros. Colorful image colorization. In *ECCV*, 2016. 2

[82] Hang Zhao, Chuang Gan, Andrew Rouditchenko, Carl Vondrick, Josh McDermott, and Antonio Torralba. The sound of pixels. In *ECCV*, 2018. 2

# Supplementary Material

## A    Optimization challenges

In this section, we give the details of the full optimization cycle and discuss differences between the single-modality baseline and our multi-modal models.

**Trivial solutions.** As discussed in [6], SDC may converge to trivial solutions, corresponding to empty clusters or encoder parameterizations, where the classifier predicts the same label regardless of the input. DeepCluster proposes workarounds to tackle these issues, involving reassigning empty cluster centers and sampling training images uniformly over the cluster assignments. While these strategies mitigate the issues, they do not fix the main cause of the problem: SDC learns a discriminative classifier on the *same* input from which it learns the labels. On the other hand, our multi-modal deep clustering models are less prone to trivial solutions because they learn the discriminative classifier on one modality and obtain the labels from a *different* modality. In our training, we never encountered the issue of empty clusters or few-class predictions for any of our multi-modal clustering approaches.

**Initialization and convergence.** Our initial pseudo-labels come from clustering features of randomly-initialized encoders. Such pseudo-labels are "good enough" to capture some weak similarities between the input samples as features from randomly-weighted networks have shown decent results on image and audio classification [51, 57]. Another potential option involves generating the initial pseudo-labels by clustering hand-crafted features, *e.g.* iDT [71] and audio spectrograms. Hand-crafted features capture low-level semantics that may help the encoders learn better or faster. Indeed, in small-scale experiments, we observed that clustering handcrafted features in the initial iteration reduces the number of clustering iterations needed to learn a well-performing encoder. However, we decided to not pursue this further, since these features are computationally expensive to extract and thus are not suitable for large-scale training on millions of examples. Furthermore, handcrafted features may bias the learning to reflect the design choices behind these manually-engineered descriptors.

**Clustering and optimization schedule.** Following previous work [6], we cluster the deep features using the $k$-means algorithm primarily for its desirable properties of efficiency and scalability. The number of $k$-means clusters is a key hyperparameter in our framework. Intuitively, using more clusters makes the pretext task harder, as it increases the number of pseudo-classes the classifier must recognize. On the other hand, the diversity of samples to cluster effectively dictates the maximum $k$, for which the grouping is still sensible. Taking into account these factors, we explore the effects of $k$ in our ablation study in Subsection 4.2 of the main manuscript. Another important hyperparameter of our framework is the number of training epochs for the encoders, before re-clustering the learned features. DeepCluster re-clusters after each epoch, which is an expensive design choice when scaling to millions of training samples. Thus, we choose to fix the pseudo-labels and train the encoders until the validation loss for predicting the pseudo-labels saturates. Then, we re-cluster the newly learned features, reassign pseudo-labels, reset the classification layer, and repeat the same process. We find this strategy to be more efficient, as it reduces the number of times we need to invoke $k$-means.

## B    Learning using audio rather than text from ASR

We note that while our approach was demonstrated by leveraging audio, the method is general and is easy to adapt to other modalities, including text. While video and text are semantically correlated, audio and video are temporally correlated. Thus, these two form of correlations are likely to provide different forms of self-supervision, potentially leading to further gains when used in combination. A disadvantage of text from ASR is that it is only available for videos with speech. Audio provides information about environmental sounds beyond speech (*e.g.* walking steps, playing guitar, and dog barking) and allows us to train on uncurated datasets of arbitrary Web videos, as we demonstrated with IG-Random.

## C    Hyperparameters and training details

**Training**. We train our models using caffe2 with distributed SGD on a GPU cluster, and employ the warmup scheme proposed in [14]. The main training parameters are presented in Table 9. We note that the epoch size can be different from the actual number of videos. This is because the total

Table 9: **Training parameter definitions.** The abbreviations and descriptions of each training parameters.

| Abv. | Name | Description |
|------|------|-------------|
| es | epoch size | The total number of examples the model trains on in one epoch. |
| bs | batch size | The size of a mini-batch. |
| lr | base lr | The initial learning rate. |
| we | warmup epoch | The number of epochs used for warmup [14]. |
| se | step epoch | Every se epochs, the learning rate |
| $\gamma$ | lr decay | is decayed by multiplying with $\gamma$. |
| te | total epoch | The training lasts for te epochs. |
| wd | weight decay | The weight decay used in SGD. |
| e-stop | early stop | Stop training when validation loss is increased in 3 consecutive epochs. |

Table 10: **Pretraining parameters.** We use early-stopping for Kinetics and AudioSet since we observe some overfiting on the pretext tasks. For the last iteration of XDC on IG-Kinetics and IG-Random, we pretrain XDC 3x longer (iteration denoted as IG-Kinetics* and IG-Random* in this table). $\gamma$ is set to 0.01 for all settings.

| method | dataset | es | bs | lr | we/se/te | wd | e-stop |
|--------|---------|-----|-----|------|----------|---------|--------|
| Superv | Kinetics | 1M | 32 | 0.01 | 10/10/45 | $10^{-4}$ | no |
| Superv | AudioSet | 2M | 32 | 0.04 | 10/20/45 | $10^{-5}$ | no |
| All DCs | Kinetics | 1M | 32 | 0.01 | 10/10/30 | $10^{-4}$ | yes |
| All DCs | AudioSet | 2M | 32 | 0.01 | 10/10/45 | $10^{-4}$ | yes |
| All DCs | IG-Kinetics & IG-Random | 10M | 32 | 0.01 | 1/3/10 | $10^{-4}$ | no |
| All DCs | IG-Kinetics* & IG-Random* | 10M | 32 | 0.01 | 0/9/30 | $10^{-4}$ | no |

number of clips the model sees during training (with temporal jittering) can be larger than the number of videos.

**Pretraining parameters**. We pretrain XDC and other baselines using the parameters described in Table 10. Early stopping is used for pretraining on small datasets such as Kinetics [27] and AudioSet [10] to stop before the model starts overfitting on the pretext task. For IG-Kinetics [12] and IG-Random, we do not observe overfitting. We pretrain XDC on IG-Kinetics and IG-Random longer in the last deep clustering iteration (denoted as IG-Kinetics* and IG-Random* in Table 10). When pretraining our R(2+1)D on longer clips (*e.g.* 32 frames), due to the GPU memory limit, we reduce the mini-batch size to 8 (instead of 32) and the base learning rate to 0.0025 (instead of 0.01).

**Finetuning parameters**. We provide finetuning hyperparameters in Table 11. Different pretraining methods may have different optimal base learning rate when finetuned on downstream tasks. Thus to make a fair comparison, we cross-validate the finetuning using the same set of base learning rates (presented in Table 12) and report the best result for each pretraining method. As we observed that higher learning rates tend to be beneficial when learning FC-only, we use a wider set of learning rates to cross-validate FC-only models. As done during pretraining, when finetuning R(2+1)D on longer clips (*i.e.* 32 frames), we reduce the mini-batch size to 8 and reduce the base learning rate to $1/4$ of its original rate.

## D  XDC using a different backbone architecture

We pretrain XDC on Kinetics with ResNet3D-18 as the visual backbone and keep the same audio encoder (ResNet-18). The results are compared with those of baselines in Table 13. XDC with the ResNet3D-18 backbone outperforms the training from scratch baseline by good margins on three downstream tasks.

## E  Additional qualitative results

**XDC clusters.** Tables 14 and 15 present the top and bottom 10 audio and video clusters learned with XDC on Kinetics, ranked by their purity with respect to Kinetics labels. We list the 5 most frequent concepts of each cluster.

Table 11: **Finetuning parameters.** Different pretraining methods have different ranges of optimal base learning rate when finetuning on downstream tasks. Thus, we cross-validate all methods with the same set of base learning rates and report the best result for each method. $\gamma$ is set to 0.01 for all settings.

| dataset | es | bs | we/se/te | wd | e-stop |
|---------|------|----|----------|-------|--------|
| HMDB51  | 40K  | 32 | 2/2/8    | 0.005 | no     |
| UCF101  | 106K | 32 | 2/2/8    | 0.005 | no     |
| ESC50   | 20K  | 32 | 2/2/8    | 0.005 | no     |

Table 12: **Finetuning base learning rates.** For a fair comparison, we cross-validate all pretraining methods with the same set of base learning rates. We report the best finetuning result for each method. Learning FC-only benefits from cross-validation with a wider range of base learning rates.

| Setup   | Base learning rates |
|---------|---------------------|
| Full    | $0.001, 0.002, 0.004, 0.006, 0.008, 0.01$ |
| FC only | $0.001, 0.002, 0.004, 0.006, 0.008, 0.01, 0.02, 0.04$ |

Table 13: **XDC using a different backbone.** We present the results of XDC on a different backbone, ResNet3D-18, for the visual encoder. We compare XDC pretrained on Kinetics vs. the two baselines: Scratch and fully-supervised Kinetics-pretraining (Superv) for the same backbone. We report the top-1 accuracy on split-1 of each dataset.

| Method | UCF101 | HMDB51 | ESC50 |
|--------|--------|--------|-------|
| Scratch (ResNet3D-18) | 60.1 | 25.7 | 54.3 |
| Superv (ResNet3D-18)  | 87.5 | 54.5 | 82.3 |
| XDC (ResNet3D-18)     | 68.0 | 36.3 | 75.5 |

**XDC filters.** Figure 3 visualizes and compares `conv_1` spatial and temporal filters of R(2+1)D learned by self-supervised XDC pretraining on IG-Kinetics versus fully-supervised pretraining on Kinetics. We observe some differences in both spatial and temporal filters between XDC and fully-supervised pretraining. In particular, XDC learns a more diverse set of motion filters.

Table 14: **XDC audio clusters.** Top and bottom 10 XDC audio clusters ranked by clustering purity *w.r.t.* Kinetics labels. For each, we list the 5 concepts with the highest purity (given in parentheses).

| # | Kinetics concepts |
|---|---|
| 1 | playing bagpipes(0.70), playing 2harmonica(0.04), playing violin(0.03), playing accordion(0.02), marching(0.01) |
| 2 | scuba diving(0.33), snorkeling(0.27), feeding fish(0.11), canoeing or kayaking(0.02), jumping into pool(0.02) |
| 3 | playing cymbals(0.21), playing drums(0.17), marching(0.03), air drumming(0.02), drumming fingers(0.02) |
| 4 | passing American football(0.17), play kickball(0.06), catching or throwing softball(0.05), kick field goal(0.02), sled dog racing(0.02) |
| 5 | presenting weather forecast(0.17), playing poker(0.05), testifying(0.03), tying knot (not on a tie)(0.02), golf putting(0.02) |
| 6 | hurling (sport)(0.17), swimming backstroke(0.05), skiing slalom(0.04), vault(0.03), ski jumping(0.02) |
| 7 | presenting weather forecast(0.15), news anchoring(0.05), filling eyebrows(0.02), braiding hair(0.02), tossing salad(0.02) |
| 8 | playing cello(0.15), playing trombone(0.11), playing accordion(0.09), playing harp(0.07), playing clarinet(0.06) |
| 9 | playing recorder(0.14), playing violin(0.12), playing trumpet(0.08), playing harmonica(0.07), tapping guitar(0.06) |
| 10 | mowing lawn(0.14), driving tractor(0.09), motorcycling(0.06), blowing leaves(0.04), water skiing(0.04) |
| 119 | side kick(0.02), front raises(0.01), dunking basketball(0.01), smoking(0.01), high kick(0.01) |
| 120 | clay pottery making(0.02), crawling baby(0.02), brushing teeth(0.01), playing harmonica(0.01), eating spaghetti(0.01) |
| 121 | pushing cart(0.01), hula hooping(0.01), high kick(0.01), blowing out candles(0.01), bench pressing(0.01) |
| 122 | shot put(0.01), feeding birds(0.01), squat(0.01), push up(0.01), high jump(0.01) |
| 123 | opening present(0.01), petting cat(0.01), pushing cart(0.01), washing dishes(0.01), punching bag(0.01) |
| 124 | trimming or shaving beard(0.01), petting cat(0.01), front raises(0.01), massaging back(0.01), tai chi(0.01) |
| 125 | feeding birds(0.01), tobogganing(0.01), riding elephant(0.01), feeding goats(0.01), jumping into pool(0.01) |
| 126 | climbing tree(0.01), writing(0.01), archery(0.01), brushing hair(0.01), shining shoes(0.01) |
| 127 | abseiling(0.01), grooming horse(0.01), milking cow(0.01), feeding goats(0.01), juggling balls(0.01) |
| 128 | washing feet(0.01), motorcycling(0.01), headbanging(0.01), cheerleading(0.01), krumping(0.01) |

Table 15: **XDC video clusters.** Top and bottom 10 XDC video clusters ranked by clustering purity *w.r.t.* Kinetics labels. For each, we list the 5 concepts with the highest purity (given in parentheses).

| # | Kinetics concepts |
|---|---|
| 1 | playing bass guitar(0.37), playing guitar(0.16), tapping guitar(0.15), strumming guitar(0.09), playing ukulele(0.09) |
| 2 | scuba diving(0.36), snorkeling(0.32), feeding fish(0.10), diving cliff(0.02), jumping into pool(0.02) |
| 3 | presenting weather forecast(0.26), playing poker(0.10), news anchoring(0.05), testifying(0.03), giving or receiving award(0.02) |
| 4 | swimming backstroke(0.21), swimming breast stroke(0.16), swimming butterfly stroke(0.10), play ice hockey(0.04), jump into pool(0.04) |
| 5 | golf putting(0.18), golf chipping(0.11), golf driving(0.05), hitting baseball(0.03), archery(0.03) |
| 6 | hurling (sport)(0.17), passing American football (in game)(0.06), skiing slalom(0.04), playing ice hockey(0.03), vault(0.03) |
| 7 | filling eyebrows(0.13), braiding hair(0.05), massaging back(0.05), curling hair(0.05), dying hair(0.03) |
| 8 | playing cello(0.12), playing harp(0.12), playing trombone(0.06), playing piano(0.06), playing accordion(0.05) |
| 9 | windsurfing(0.12), jetskiing(0.10), water skiing(0.09), surfing water(0.08), kitesurfing(0.06) |
| 10 | cooking chicken(0.11), barbequing(0.07), frying vegetables(0.06), cooking sausages(0.04), making pizza(0.04) |
| 55 | yoga(0.02), folding napkins(0.02), doing nails(0.02), cutting watermelon(0.01), writing(0.01) |
| 56 | eating spaghetti(0.02), making pizza(0.02), brushing teeth(0.02), blowing out candles(0.02), reading book(0.02) |
| 57 | answering questions(0.02), tai chi(0.02), dancing ballet(0.02), dunking basketball(0.02), sign language interpreting(0.01) |
| 58 | trimming or shaving beard(0.02), barbequing(0.02), kissing(0.02), dining(0.01), playing poker(0.01) |
| 59 | punching bag(0.02), blowing out candles(0.02), pumping fist(0.02), dancing gangnam style(0.02), opening present(0.01) |
| 60 | feeding goats(0.02), blowing out candles(0.02), milking cow(0.02), arm wrestling(0.02), finger snapping(0.02) |
| 61 | air drumming(0.02), pumping fist(0.02), pushing cart(0.02), brushing teeth(0.02), eating ice cream(0.01) |
| 62 | clean and jerk(0.01), robot dancing(0.01), bench pressing(0.01), side kick(0.01), punching bag(0.01) |
| 63 | pull ups(0.01), gymnastics tumbling(0.01), punching bag(0.01), cracking neck(0.01), eating ice cream(0.01) |
| 64 | capoeira(0.01), riding elephant(0.01), feeding goats(0.01), feeding birds(0.01), crawling baby(0.01) |

a) conv1 spatial and temproal filtes learned by Kinetics fully supervision.



b) conv1 spatial and temporal filters learned by IG65M self-supervised XDC.
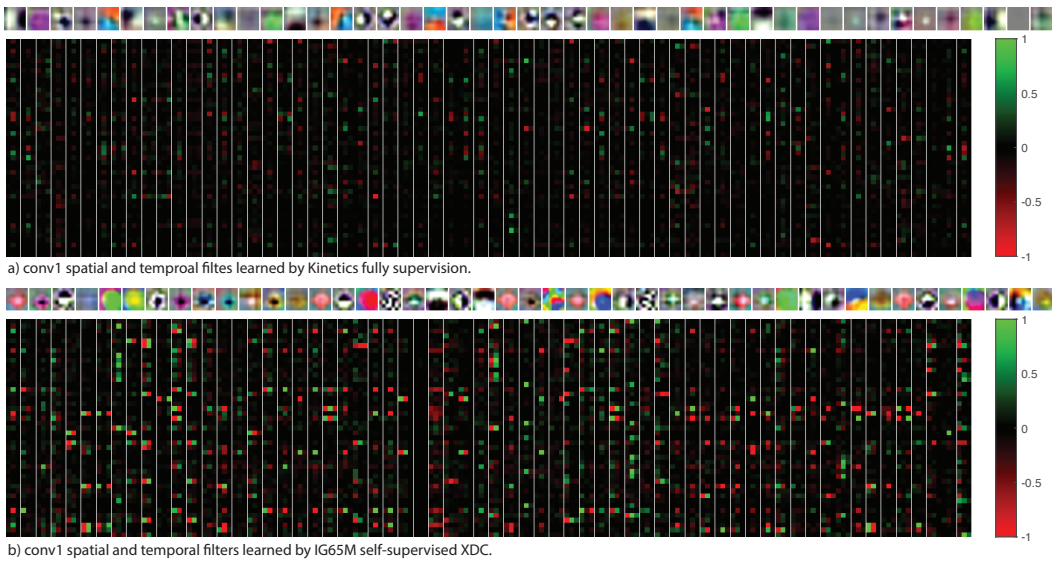
Figure 3: **R(2+1)D filters learned with self-supervised XDC vs. fully-supervised training.** (a) R(2+1)D `conv_1` filters learned by fully-supervised training on Kinetics. (b) The same filters learned by self-supervised XDC pretraining on IG-Kinetics. XDC learns a more diverse set of temporal filters compared to fully-supervised pretraining.