VIDROP: VIDEO DENSE REPRESENTATION THROUGH SPATIO-TEMPORAL SPARSITY

Anonymous authors

Paper under double-blind review

ABSTRACT

Self-supervised learning (SSL) has revolutionized image processing, but extending its success to video understanding presents unique challenges due to increased data complexity and computational demands. We introduce ViDROP (Video Dense Representation thrOugh spatio-temporal sParsity), a novel SSL architecture for video understanding that combines token dropping and masking strategies. Our approach eliminates the need for a decoder and enables per-patch loss computation, overcoming limitations of previous video SSL methods. Moreover, we propose a simple yet effective video compression technique using k-means clustering in pixel space, significantly accelerating data loading and facilitating rapid experimentation. ViDROP demonstrates remarkable scalability across model sizes, from ViT-Small to ViT-Huge, when starting from pretrained models (VideoMAE or V-JEPA), achieving significant performance gains. Pushing the boundaries even further, we leverage network expansion techniques to successfully train ViT-Huge from scratch using modest computational resources, achieving comparable accuracy to VideoMAE $25 \times$ faster in training time. This marks a significant breakthrough in large-scale video SSL, enabling the training of state-of-the-art models with limited resources. Extensive experiments show that ViDROP achieves state-of-the-art performance on various video understanding benchmarks, including Kinetics400, SSv2, UCF101, and HMDB51, as well as in temporal action detection (THUMOS14). These results highlight the effectiveness of our finegrained token-level learning strategy in a domain traditionally dominated by finetuned SSL models, while enabling the training of large-scale models with limited computational resources.

033

004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

034 035

1 INTRODUCTION

Self-supervised learning (SSL) has transformed computer vision by enabling models to learn rich
 representations from vast amounts of unlabeled data. While SSL methods like DINOv2 Oquab
 et al. (2023) have achieved remarkable success in image processing, extending these techniques to
 video understanding presents unique challenges due to increased data complexity and computational
 demands.

The temporal dimension in videos captures essential information about object movement and scene changes, crucial for action understanding. However, it also significantly increases the amount of data to be processed, exacerbating computational burdens. Recent works such as VideoMAE Tong et al. (2022) have addressed these challenges by employing sparse encoder and dense decoder architectures, leveraging vision transformers Dosovitskiy et al. (2020) to efficiently process video data split into tokens or patches.

Despite their promise, these approaches face two significant limitations. First, the computational cost of the decoder remains considerable, as it still operates on all or a large portion of video tokens (see *e.g.*, Wang et al. (2023a); Hwang et al. (2022)). Second, the split between encoder and decoder, coupled with a low-level reconstruction objective, can lead to suboptimal representation learning.
The encoder may not capture all relevant high-level features, as some information is relegated to the decoder for reconstruction purposes. This division of representational power can result in less robust or comprehensive features, potentially misaligning with high-level downstream tasks.

054



Figure 1: **Visualization of per-patch representations.** Top row: Original video frames of a dog walking on grass. Middle row: 3-channel PCA visualization Oquab et al. (2023) of patch features from VideoMAE Tong et al. (2022). Bottom row: 3-channel PCA visualization of patch features from ViDROP. The more detailed and coherent colorization in the bottom row demonstrates ViDROP's ability to produce higher-quality per-patch representations, capturing finer details and maintaining better consistency (semantic correspondences) across frames.

072 073 074

096

098

099

102

103

068

069

071

To address these challenges, we propose ViDROP (Video Dense Representation thrOugh spatiotemporal sParsity), a novel SSL architecture for video understanding that combines token dropping and masking strategies. ViDROP employs a sparse encoder with masked tokens and eliminates the need for a decoder, enabling per-patch loss computation for robust representations while maintaining exceptional efficiency. This approach overcomes the limitations of previous methods, providing a more effective and computationally efficient solution for self-supervised learning in video understanding tasks (see Figure 2).

To accelerate training, we introduce a simple yet effective video compression technique using k-means clustering Lloyd (1982) in pixel space. This innovation addresses the data loading bottleneck that emerges as our model accelerates computation, especially for smaller models.

ViDROP achieves state-of-the-art linear probing results using only minimal data augmentations (random resized cropping and flipping), in contrast to top-performing models like SVT Ranasinghe et al. (2021) and ρ BYOL Feichtenhofer et al. (2021) that rely on heavy, potentially dataset-specific augmentations. This approach aims for more robust and generalizable representations. Additionally, ViDROP produces high-quality per-patch representations alongside strong global representations, distinguishing it from traditional contrastive methods that excel in global video clip representations but may struggle with fine-grained, per-patch understanding (see Figure 1). This versatility enables ViDROP to be effective across a wide range of downstream tasks, from those requiring local understanding to those needing global video comprehension.

- Our contributions can be summarized as follows:
 - A novel, efficient SSL architecture for video understanding: We introduce ViDROP, combining sparse token processing with masked learning, achieving state-of-the-art performance without a decoder while producing high-quality per-patch and global representations.
 - Scalable and resource-efficient training: We demonstrate exceptional scalability across model sizes (ViT-Small to ViT-Huge) and initialization strategies (VideoMAE Tong et al. (2022), V-JEPA Bardes et al. (2024)). By leveraging network expansion techniques Wang et al. (2023b) and introducing a novel k-means clustering-based video compression method, we train a ViT-Huge model with comparable accuracy to VideoMAE in just 4% of the total training time, enabling large-scale video SSL with limited computational resources.
- Comprehensive evaluation: ViDROP achieves superior performance across a wide range of video understanding tasks, including action recognition, temporal action detection, frame-wise tracking, and copy detection. Our method's versatility extends to image classification, showcasing the robustness and generalizability of the learned representations.

¹⁰⁸ 2 PRIOR WORK

110

Evolution of image SSL: Self-supervised learning (SSL) in computer vision has evolved from image-based to video understanding approaches. Image SSL progressed from contrastive learning van den Oord et al. (2018); Chen et al. (2020;b) to clustering-based methods Caron et al. (2018; 2020; 2021), then to direct feature prediction Grill et al. (2020); Chen & He (2020) and cross-correlation techniques Zbontar et al. (2021). These discriminative approaches have consistently advanced representation learning in static images.

117 Masked modeling and diverse SSL approaches: Complementing these discriminative methods, 118 various other approaches have played a crucial role in SSL. Notably, masked modeling approaches 119 have gained significant traction. BEiT Bao et al. (2021) introduced the concept of infilling in latent 120 space, inspired by BERT Devlin et al. (2019) in natural language processing. Masked Autoencoders 121 (MAE) He et al. (2021) adapted this idea to regress missing pixel patches directly in pixel space, 122 while AIM El-Nouby et al. (2024) advanced these concepts with next token prediction in latent 123 space. Hybrid models such as iBOT Zhou et al. (2021) and DINOv2 Oquab et al. (2023) have further pushed the boundaries, combining generative and discriminative elements to produce high-124 quality representations. These approaches have demonstrated exceptional performance in tasks like 125 KNN classification Fix & Hodges (1989) and unsupervised semantic segmentation Hamilton et al. 126 (2022), showcasing the power of integrating multiple SSL paradigms. 127

128 Video SSL challenges: The transition from image to video SSL introduces unique challenges due to 129 the temporal dimension and increased data scale. Video SSL methods often employ pretext tasks designed to capture both spatial and temporal aspects, such as determining the correct order of shuffled 130 video frames Misra et al. (2016); Jenni et al. (2020); Dorkenwald et al. (2022); Dave et al. (2023). 131 Generative models like VideoMAE Tong et al. (2022) and VideoMAEv2 Wang et al. (2023a) have 132 been adapted to reconstruct video segments or entire frames, while consistency-based approaches 133 like ρ BYOL Feichtenhofer et al. (2021), BRAVE Recasens et al. (2021), and SVT Ranasinghe et al. 134 (2021) ensure feature consistency across different video segments. 135

136 **Computational optimization:** The computational demands of processing videos, especially with 137 Vision Transformers (ViT) Dosovitskiy et al. (2020), have led to various optimization strategies. Factorized attention, as used in SVT Ranasinghe et al. (2021), offers one approach, although com-138 prehensive spatio-temporal attention, as in ViViT Arnab et al. (2021), typically yields superior re-139 sults. Some methods leverage pretrained image or short video clip encoders to derive representations 140 for longer clips Sameni et al. (2022); Lei et al. (2021). Inspired by MAE He et al. (2021), sparse 141 encoder inputs paired with dense decoder outputs have been employed to reduce computational re-142 quirements, with models like VideoMAEv2 Wang et al. (2023a) and EVEREST Hwang et al. (2022) 143 further economizing resources by introducing sparse reconstruction tokens to the decoder. 144

Data processing bottlenecks: Addressing the data bottleneck in video SSL remains a significant 145 challenge. While tools like FFCV Leclerc et al. (2023) and its SSL variant Bordes et al. (2023) are 146 effective for images, they fall short for video processing. Wrappers such as decord and Avion Zhao 147 & Krahenbuhl (2023) offer some improvements but remain relatively slow. Data remasking tech-148 niques Bardes et al. (2024); Feichtenhofer et al. (2022) provide a cost-effective way to alleviate 149 data loading bottlenecks. Another approach involves using precomputed latent representations of 150 videos Wiles et al. (2022); Jaegle et al. (2021), which can be integrated with learned data augmen-151 tations Lee et al. (2023). However, this method introduces higher costs at the inference stage due to 152 the need for an expensive encoder to convert raw pixel videos into the latent representation.

153 **Role of data augmentation:** The role of data augmentation in SSL has been a subject of significant 154 research Zhai et al. (2023); Wagner et al. (2022); Kalibhat et al. (2023). While heavy augmentations 155 have been crucial for the success of many image SSL methods Chen et al. (2020a); Grill et al. (2020), 156 their application to video data is more complex due to the need to preserve temporal coherence Fe-157 ichtenhofer et al. (2021). Some video SSL approaches have adapted image augmentation techniques 158 to the video domain Ranasinghe et al. (2021), while others have developed video-specific augmenta-159 tions Qian et al. (2020). However, the computational cost of these augmentations can be substantial, especially for large-scale video datasets. Moreover, the reliance on carefully designed augmenta-160 tions raises questions about the broader applicability and robustness of SSL methods across diverse 161 video datasets and tasks.



Figure 2: **Comparison of patch reconstruction architectures.** This figure illustrates three different approaches to patch reconstruction. (a) Traditional Encoder/Decoder architecture, where a sparse set of input patches is fed to the encoder, and a small decoder reconstructs all dropped tokens. (b) Our proposed ViDROP architecture, which uniquely combines sparse (dropped) and masked (MSK) tokens in a single encoder. (c) Masked Encoder approach (such as DINOv2 Oquab et al. (2023)) that uses only MSK tokens. In all methods, the input is flattened and patched (image or video), and targets are either pixels or the output of the teacher network (an exponential moving average of the encoder, omitted for clarity).

3 Method

186 187 188

177

178

179

180

181

182

183

184 185

ViDROP introduces a novel approach to self-supervised learning for video understanding, com-189 bining efficient processing with effective representation learning. Our method builds upon the DI-190 NOv2 Oquab et al. (2023) framework, adapting it for video data through a sparse encoder archi-191 tecture that integrates token dropping and masking within a single network. We employ a video 192 sampling and processing strategy using both large and small crops, along with a multi-component 193 loss function ensuring global and local feature consistency. Additionally, we introduce a lossy data 194 loading scheme based on k-means clustering for efficient video compression. In the following sub-195 sections, we detail these components and explain how they work together to overcome the limitations 196 of existing video SSL methods.

197

199

3.1 VIDEO SAMPLING AND PROCESSING

Given a video, we sample multiple clips, creating two large crops and eight small crops. The large crops are fed directly to an exponential moving average (EMA) of the network, serving as the teacher. For the student network, we apply sparsification (token dropping) to these same views and replace some of the remaining tokens with a special MSK token. Small crops maintain dense inputs without sparsification or masking Assran et al. (2022). Both student and teacher networks include a CLS token for global representation.

206 207

3.2 Loss Function

ViDROP's loss function consists of two main components. The first is a DINO-style loss calculated
between the CLS tokens of the student and teacher Caron et al. (2021). The second is a patchlevel loss for the masked tokens in the student, using the corresponding outputs from the teacher
as targets Zhou et al. (2021). We employ a global-global consistency loss to ensure coherence
between different views of the same video and a global-local consistency loss that aligns features
from large crops (global) with those from small crops (local). The patch-level loss enables finegrained representation learning, complemented by the KoLeo regularization adapted from DINOv2.
For detailed formulations of these losses, we refer readers to the DINOv2 paper Oquab et al. (2023).

219 220 221



224

225

226 227

228 229 230

231

232

233

234

235 236



Figure 3: **Comparison of video frame compression methods.** Top row: Original frames sampled from the Kinetics-400 dataset Kay et al. (2017). Middle row: Frames compressed and reconstructed using VQ-GAN Esser et al. (2020). Bottom row: Frames compressed and reconstructed using our proposed k-means clustering method Lloyd (1982). Our approach achieves a balance between compression efficiency and visual quality, maintaining compatibility with pretrained models while significantly accelerating data loading.

Unlike traditional consistency-based SSL approaches, we minimize the use of data augmenta tion Moutakanni et al. (2024), relying primarily on heavy masking and our lossy data compression
 technique as pseudo-augmentations.

240 241

3.3 Architecture

242 ViDROP employs a sparse encoder with masked tokens, eliminating the need for a separate decoder 243 (Figure 2). This design differs from asymmetric encoder-decoder architectures He et al. (2021); 244 Tong et al. (2022); Liu et al. (2022); Bardes et al. (2024); Assran et al. (2023); Baevski et al. (2022a) 245 by uniquely combining token dropping and masking strategies within a single encoder. Our approach 246 achieves the efficiency of sparse processing while maintaining the rich representational learning of 247 mask-based self-distillation methods Zhou et al. (2021); Oquab et al. (2023); Bao et al. (2021); 248 Xie et al. (2021); Baevski et al. (2022b). Crucially, our design enables per-patch loss computation, 249 allowing fine-grained representation learning without the computational overhead of dense decoders. 250 This bridges the gap between token dropping efficiency and mask-based representational power.

251

253

3.4 LOSSY DATA LOADING

To address the data loading bottleneck in video processing, we introduce a lossy data loading scheme based on k-means clustering Lloyd (1982) of video patches (see Figure 3). By applying k-means clustering directly in pixel space on 10×10 pixel patches, we achieve a compression rate of 150 (using 65536 clusters), significantly accelerating the data loading process.

This approach offers several advantages over VQVAE-based methods van den Oord et al. (2017); Esser et al. (2020); Park et al. (2023); Wiles et al. (2022). It allows for faster processing during both training and inference by avoiding complex encoding and decoding steps. Our method maintains compatibility with pretrained models that operate on raw pixel data, enabling us to leverage existing large-scale pretrained checkpoints. Furthermore, our approach provides a flexible compression scheme that can be easily adjusted to balance between compression rate and representation quality.

264 265

266

268

4 EXPERIMENTS

- 267 4.1 EXPERIMENTAL SETUP AND PROTOCOLS
- We use the training set of Kinetics-400 Kay et al. (2017) for self-supervised training. For evaluation, we use four common action classification datasets (Kinetics-400, Something-Something-

Table 1: ViDROP ablation experiments. We report linear probing (single crop) accuracy (%) with
ViT-B/16 on K400. If not specified, the default is: the loss is iBOT Zhou et al. (2021), the data
augmentation is random resized cropping, the number of small crops is 8, the masking ratio is 85%,
and the pre-training length is 60 epochs. Default settings are marked in gray.

274 275 276	(a) Loss function. Patch loss with MASK tokens improves the performance.	(b) Number of small crops. Even having a few small crops boosts the performance.	(c) Drop pattern. Simple ran- dom token dropping is more ef- fective than complex patterns.		
277 278 279 280	patch losslinearvisible tokens51.7none53.1masked tokens54.9	num. linear 8 54.9 4 54.1 0 50.2	patternlinearrandom tokens54.9random tubes54.5block (vjepa)53.1		
281 282 283 284 285 286	(d) Drop rate . Lower drop rates yield better performance but re- quire longer training times. The model with an 80% drop rate took the longest to train due to reduced batch size.	(e) Number of clusters . Re- ducing clusters from 65k to 16k maintains accuracy, while 2k clusters slightly decrease it. Us- ing a shared head for 65k clus- ters also lowers accuracy.	(f) Masking probability . Increasing the masking probabil- ity range from 10-40% to higher values slightly improves model accuracy (all the models use 16k clusters).		
287	rate time(hh:mm) linear	num. shared linear	min max linear		

287	rate	time(hh:mm)	linear		num.	shared	linear	min	max	linear
288	80%	31:54	55.2	•	65k	X	54.9	0.1	0.4	54.9
289	85%	24:21	54.9		16k	X	54.9	0.1	0.7	54.9
290	90%	23:46	54.4		2k	X	54.3	0.5	0.7	55.5
291	95%	23:19	51.5		65k	\checkmark	53.8	0.7	0.7	55.0

292 293

v2 Goyal et al. (2017), UCF101 Soomro et al. (2012), and HMDB51 Kuehne et al. (2011)) and THUMOS14 Jiang et al. (2014) for temporal action detection. For all evaluations, contrary to common settings in reconstruction-based models Tong et al. (2022); He et al. (2021); Liu et al. (2022); Baevski et al. (2022b); Bardes et al. (2024) but compatible with consistency-based models that rely on heavy data augmentations, we use a **frozen** backbone and only train a linear head on top (except in the case of THUMOS14, where we train a transformer following ActionFormer Zhang et al. (2022)).

For main results, we train ViT-Base, ViT-Large, and ViT-Huge models from pretrained checkpoints.
ViT-Base models are trained for 60 epochs with k-means compressed data (24 hours), while ViT-Large and ViT-Huge are trained for 40 epochs without compression (57 and 100 hours, respectively).
All use a total batch size of 512 (with gradient accumulation).

For the LEMON Wang et al. (2023b) experiment with random initialization, we follow a progressive schedule: ViT-Small (100 epochs, k-means data, 17 hours), ViT-Base (6 epochs, real data, 12 hours), ViT-Large (19 epochs, 52 hours), and ViT-Huge (30 epochs, 62 hours), totaling 155 epochs and 143 hours (almost 6 days or 1144 V100 hours). For comparison, the VideoMAE ViT-Huge model trained for 1600 epochs is estimated to take approximately 28862 V100 hours (25× longer).

During training, 85% of the tokens are dropped for the student when fed with two large clips (16 frames of 224×224). For half of the mini-batch, we randomly mask a portion of the remaining tokens. Additionally, 8 local crops of size 96×96 (8 frames) are fed as dense input without token dropping or masking.

314 315

316 4.2 ABLATIONS

Here we perform in-depth ablation studies on ViDROP design choices with a ViT-B. We use 16384
clusters for the DINO loss, both for the CLS token and the per-patch loss, and mask uniformly between 10% to 40% of the tokens. To accelerate training, we start from a pretrained VideoMAE Tong
et al. (2022) checkpoint trained on K400 for 800 epochs.

Loss function. Table 1a shows that patch token loss is crucial for ViDROP. An extra loss on masked
 tokens (iBOT Zhou et al. (2021)) outperforms DINO Caron et al. (2021)'s approach. Using MASK
 tokens for loss calculation proves more effective than calculating loss on all visible tokens.



335

336

337

338

339

Figure 4: **Probing accuracy as a function of training time for dense and sparse models.** The dense model (without token dropping) takes significantly longer to train and achieves lower accuracy compared to the sparse model (with 85% token drop rate), even after training for four times longer.

Table 2: Effect of initialization on ViT-Base model performance. The random model was trained for 240 epochs (4× the epochs of the pretrained models). Initial pretraining was conducted on 64 Tesla V100 GPUs, and our training on 4 RTX 4090 GPUs.† is estimated by multiplying the training time of the VMAE₁₆₀₀ by the relative training throughput of a ViT-Base and ViT-Large (4.24)

Satting	Pretraining Time	Linear
Setting	(V100 Days)	Accuracy
VMAE ₈₀₀	74	41.0%
+ViDROP	+8	55.5%
VMAE ₁₆₀₀	148	43.5%
+ViDROP	+8	57.0%
ViDROP rand	32	53.3%
VMAE ^{Large} ₁₆₀₀	627†	52.5%

Number of small crops. A key success element of DINO Caron et al. (2021) and MSN Assran et al.
(2022) was using small crops. We see a similar pattern in Table 1b, where small crops significantly
boost performance. We use the same setting as DINOv2 Oquab et al. (2023), but observe that fewer
crops are viable. If computational load is an issue (as small crops are not sparsified), the number
can be reduced while maintaining considerable performance gains.

Drop pattern. While previous video reconstruction methods emphasized specific token dropping
patterns Tong et al. (2022); Bardes et al. (2024); Feichtenhofer et al. (2022), Table 1c shows that
random patch dropping outperforms tube dropping Tong et al. (2022); Wang et al. (2023a) and block
masking Bardes et al. (2024). We observe higher self-supervised loss for these patterns compared to
random dropping, suggesting that their difficulty may hinder the model's representational power.

Drop rate. Token dropping is an essential component for reducing the computational load of our model. In Table 1d, we can see that there is a trade-off in terms of training time and quality. Reducing the drop rate improves the quality but at the cost of extra training time. Note that in this experiment, we had to reduce the batch size of the model trained with an 80% token drop rate to be able to train it on our hardware. We further studied this trade-off and trained a dense model (i.e., a model with a drop rate of 0%) for 100 hours (almost 4× the base model) and achieved an accuracy of only 47.1%, which is significantly lower than the 54.9% accuracy of the base model (see Figure 4).

362 Number of clusters. Since we are using Sinkhorn-Knopp centering for our DINO Caron et al. 363 (2021) losses, we are significantly more compute-heavy in the loss (compared to V-JEPA Bardes et al. (2024) and VideoMAE Tong et al. (2022)). Additionally, we can't apply gradient accumulation 364 with many steps, since the centering operation depends on the whole batch. Changing the loss is beyond the scope of this paper, so instead, we reduced the number of clusters both for the global loss 366 and patch loss. Results in Table 1e show that, similar to the findings of MSN Assran et al. (2022), 367 we can significantly reduce the number of clusters and maintain the same performance. We also 368 notice that, similar to the findings of DINOv2 Oquab et al. (2023), having different heads for the 369 two loss terms is beneficial. For the rest of the ablation studies, we used 16k clusters. 370

Masking probability. Following the setting of iBOT Zhou et al. (2021) and DINOv2 Oquab et al. (2023), we apply masking to only half of the large crops. For the other half, we initially randomly masked between 10% to 40% of the remaining tokens (after token dropping). Table 1f shows that we can use larger probabilities and achieve slight improvement. This likely stems from video data's higher redundancy compared to images, allowing for greater sparsity (e.g., 90% for videos Tong et al. (2022) vs 75% for images He et al. (2021)).

Initialization. To demonstrate the general applicability of our model, we trained it from scratch and compared it to models initialized with different pretrained weights, as shown in Table 2. The

378 Table 3: Comparison of different compression methods. 379 Evaluation of various methods for compressing and decom-380 pressing 0.5M frames on 4 GPUs. KMeans-based methods offer a good balance between quality (PSNR), processing 381 time, and compression factor. 382

Table 4: Effect of KMeans on training speed and accuracy with ViT-Small. Using KMeans compression achieves a $5.37 \times$ speedup with a minor performance cost. Training on pixel data for the same duration results in lower performance.

Time

(hh:mm)

9:42

52:09

9:42

*K*Means

./

Х

Х

Linear

Accuracy

44.4

46.1

29.5

Method	PSNR	Time (mm:ss)	Compression Factor	
SD-XL	30.29	40:29	24	
TinyAE-XL _{byte}	26.44	7:25	48	
VQGAN-f16	23.36	25:59	384	
$\widetilde{KMeans}_{16 \times 16}$	24.28	1:23	384	
$KMeans_{10\times 10}$	25.75	1:41	150	

Table 5: Training throughput of various SSL methods with ViT models. Comparison of training speeds for different SSL methods using ViT-Base, ViT-Large, and ViT-Huge architectures on a single NVIDIA RTX 4090 GPU. DINOv2 refers to the dense version of our model without sparsity. While ViDROP shows lower raw throughput compared to reconstruction-based methods, it significantly outperforms consistency-based models (DINOv2 and SVT) and achieves superior sample efficiency (see Table 2).

Mathad	Vi1	-Base	ViT	-Large	ViT-Huge		
Method	Max B.S. Throughput		Max B.S.	Throughput	Max B.S.	Throughput	
SVT	4	8.4	1	1.5	OOM	N/A	
VMAE	30	139.6	8	32.9	4	17.0	
VMAEv2	15	58.8	8	29.7	4	16.1	
VJEPA	37	74.4	22	37.9	8	16.2	
DINOv2	12	23.4	4	7.8	1	3.1	
ViDROP	33	43.1	11	17.4	4	7.8	

407 408

391 392

393

394

395

396

397

409

410 randomly initialized model was trained for $4\times$ the epochs of the pretrained models but still required 411 significantly less total training time when considering the pretraining duration of the checkpoints.

412 Data compression method. To accelerate training for ViT-Small and ViT-Base models, we em-413 ployed KMeans-based data compression. Table 3 shows that KMeans compression strikes a good 414 balance between quality (measured by PSNR), encoding/decoding time, and compression rate. 415 Other methods rely on large models, complicating inference and disallowing the use of pretrained 416 checkpoints. For all ablation experiments, data was compressed once, taking around 40 hours for 60 417 epochs of data. We used a patch size of 10×10 and 65536 clusters.

418 Table 4 demonstrates that using KM eans data results in a $5.37 \times$ speedup with a small performance 419 cost when training a ViT-Small model. Training the same model on pixel data for the same duration 420 yields worse performance. ViT-Base sees a $2.71 \times$ speedup, and ViT-Large a $1.24 \times$ speedup (better 421 GPUs can lead to greater speedup, as data becomes the bottleneck).

422 **Training throughput.** Table 5 compares the training speed of various self-supervised learning 423 (SSL) methods for video ViT models on a single NVIDIA RTX 4090 GPU (24 GB VRAM). We 424 used official code and configurations for each model: VJEPA Bardes et al. (2024) with repeated 425 masking of 2, VideoMAEv2 Wang et al. (2023a) with 4, and VideoMAE Tong et al. (2022) without 426 repeated masking. DINOv2 represents the dense version of our model without sparsity. Measure-427 ments exclude data loading time, focusing on forward and backward passes. While ViDROP shows 428 lower raw throughput compared to reconstruction-based methods due to its combination of small 429 and large crops and more complex loss calculation, it significantly outperforms consistency-based models like SVT and DINOv2. Crucially, as demonstrated in Table 2, ViDROP achieves superior 430 sample efficiency, reaching higher performance with fewer training iterations, thus offsetting its 431 lower per-iteration speed.

Table 6: **Comprehensive performance comparison on action classification tasks.** We report linear probing accuracies (%) for all datasets and KNN accuracies for smaller datasets. Numbers in parentheses indicate evaluation clips. For K400, we also present low-shot learning results with varying amounts of labeled data. For temporal action detection on THUMOS14, we report the average mean Average Precision (mAP) across different temporal Intersection over Union (tIoU) thresholds. ViDROP variants consistently outperform their baselines across all tasks and metrics. Superscript numbers indicate performance gap compared to the respective baseline.

Method	K4 Linear (5×3)	00 Attn. (1×1)	SSv2 Linear (2×3)	$ UCH \\ Linear \\ (5 \times 3)$	F-101 KNN (1×1)	HMI Linear (5×3)	DB-51 KNN (1×1)	5% K4	00 Low-s 10% inear (5×	shot 50% 3)	THU. Avg. mAP
hoBYOL SVT	71.5 68.1	N/A N/A	25.3 20.3	89.6 91.3	85.2 87.2	61.2 63.1	49.7 51.8	- -	-	-	-
VMAE ^{Large} +ViDROP ^{kmeans} +ViDROP	$\begin{vmatrix} 60.7 \\ 63.4^{+2.7} \\ 74.3^{+13.6} \end{vmatrix}$	$68.671.9^{+3.3}72.6^{+4.0}$	27.9 32.9 ^{+5.0} 38.4 ^{+10.5}	84.5 86.6 ^{+2.1} 94.2 ^{+9.7}	49.1 73.0 ^{+23.9} 88.0^{+38.9}	60.3 60.9 ^{+0.6} 69.6^{+9.3}	29.2 45.5 ^{+16.3} 55.2 ^{+26.0}	$\begin{vmatrix} 43.3 \\ 50.4^{+7.1} \\ 59.7^{+16.4} \end{vmatrix}$	50.4 55.8 ^{+5.4} 64.2 ^{+13.8}	59.6 64.7 ^{+5.1} 71.8 ^{+12.2}	15.0 50.2 ^{+35.2} 57.1 ^{+42.1}
VJEPA ^{Large} +ViDROP _{noaug} +ViDROP	62.7 72.4 ^{+9.7} 74.8^{+12.1}	73.7 71.1 ^{-2.6} 72.7 ^{-1.0}	43.2 33.8 ^{-9.4} 38.7 ^{-4.5}	92.0 91.1 ^{-0.9} 93.7 ^{+1.7}	81.2 81.3 ^{+0.1} 84.8 ^{+3.6}	66.9 66.0 ^{-0.9} 69.6 ^{+2.7}	54.7 51.0 ^{-3.7} 55.3^{+0.6}	49.6 58.5 ^{+8.9} 61.2 ^{+11.6}	54.8 62.5 ^{+7.7} 65.4^{+10.6}	58.0 69.8 ^{+11.8} 72.5^{+14.5}	$\begin{vmatrix} 20.1 \\ 49.9^{+29.8} \\ 56.2^{+36.1} \end{vmatrix}$
VMAE ^{Huge} +ViDROP ViDROP ^{Huge} noaug	67.2 74.8 ^{+7.6} 66.3 ^{-0.9}		29.3 39.3^{+10.0} 30.1 ^{+0.8}	84.3 94.2 ^{+9.9} 85.2 ^{+0.9}	48.6 85.9 ^{+37.3} 75.7 ^{+27.1}	59.7 69.4 ^{+9.7} 59.6 ^{-0.1}	30.1 53.4 ^{+23.3} 45.4 ^{+15.3}	$\begin{vmatrix} 46.8 \\ 60.3^{+13.5} \\ 50.1^{+3.3} \end{vmatrix}$	52.9 64.4 ^{+11.5} 55.4 ^{+2.5}	64.5 72.2 ^{+7.7} 63.7 ^{-0.8}	41.7 55.7 ^{+14.0} 47.2 ^{+5.5}

4.3 RESULTS

456

For our main results, we train seven different models to comprehensively evaluate the performance of 457 ViDROP under various conditions. We begin with two ViT-Large models based on VideoMAE Tong 458 et al. (2022), one with minimal data augmentation and another incorporating heavier augmentations 459 used in DINO Caron et al. (2021). These are followed by two similarly configured ViT-Large models 460 based on VJEPA Bardes et al. (2024), demonstrating the compatibility and versatility of our method 461 across different pretraining paradigms. To showcase scalability, we include two ViT-Huge mod-462 els: one initialized from a VideoMAE checkpoint and another trained using LEMON Wang et al. 463 (2023b) expansion techniques from random initialization, highlighting the efficiency and potential of 464 our method for training large-scale models with limited computational resources. All VideoMAE-465 based models are initialized with checkpoints pretrained for 1600 epochs on Kinetics400, while VJEPA-based models use the original VJEPA checkpoint. For one VideoMAE-based ViT-Large 466 model, we utilize KMeans compressed data. While KMeans compression doesn't offer significant 467 speed benefits for ViT-Large models, we include this configuration to demonstrate the robustness 468 of our findings. By using compressed data, we establish a lower bound on performance (as shown 469 in Table 4), yet still outperform the original VideoMAE model, underscoring the effectiveness of 470 ViDROP even under potentially suboptimal data conditions. Additionally, to facilitate direct com-471 parisons with academic papers, we separately train a ViT-Base model based on VideoMAE trained 472 for 800 epochs on Kinetics400. 473

Action classification results. Table 6 presents our comprehensive results on various action class-474 sification tasks, including standard action recognition, low-shot learning, and temporal action de-475 tection. For action recognition, ViDROP variants consistently improve linear probing accuracies 476 across all datasets, with the exception of VJEPA-based models on some datasets due to potential 477 forgetting effects from extra training data of VJEPA Bardes et al. (2024). Notably, attention prob-478 ing for VideoMAE-based models shows improvement, while VJEPA-based models exhibit lower 479 attention probing accuracies compared to linear probing, reminiscent of DINOv2's behavior on Im-480 ageNet. KNN accuracies see substantial improvements, particularly for VideoMAE-based models. 481 The incorporation of data augmentation further enhances performance across metrics. In low-shot 482 settings on K400, our ViDROP models often match or surpass the full data baseline using only 10% of the labeled data, demonstrating strong data efficiency. For temporal action detection on 483 THUMOS14 Jiang et al. (2014), ViDROP shows remarkable improvement over baselines, with 484 the average mAP more than doubling, indicating more temporally precise feature representations. 485 The LEMON-trained model, starting from random initialization, achieves comparable results to Table 7: **Performance comparison of ViT-Base models.** Linear probing accuracies (%) on various datasets, DAVIS tracking $((\mathcal{J}\&\mathcal{F})_m \text{ metric})$, and copy detection (μ AP). SIGMA_{DINO} Salehi et al. (2024) uses a teacher model pretrained on extra data.

489		-						
490	Method	SSv2	K400	UCF	HMDB	IN-1K	DAVIS	Copy Det.
491	SIGMA _{MLP}	19.9	30.7	73.8	45.0	24.1	-	-
492	SIGMA _{DINO}	20.8	47.5	80.7	52.3	45.0	-	-
493	VMAE ^{Base}	17.5	20.7	58.6	37.7	20.2	40.8	2.0
494	+ViDROP noaug	25.4	59.5	81.7	54.6	44.6	54.8	22.4
495								

496 497

498 VMAE $_{1600}^{Huge}$, with slightly lower performance on full K400 but improved results in low-shot scenar-499 ios, highlighting the potential of progressive training for large-scale video SSL.

These comprehensive results demonstrate the effectiveness and scalability of ViDROP across various
 model sizes, datasets, and action understanding tasks. Our method consistently achieves state-of the-art performance in linear probing, shows strong transferability to smaller datasets and low-shot
 scenarios, and exhibits remarkable generalization to temporal action detection, underscoring its po tential for a wide range of video understanding applications.

505 **Comparison with small scale methods.** To facilitate fair comparisons with small scale approaches, 506 we trained a ViT-Base model following the setting of SIGMA Salehi et al. (2024) (current SotA 507 method in small scale models). Table 7 presents the results of this comparison across various 508 datasets and tasks. For ImageNet (IN-1K) evaluation, we follow the common practice of repeat-509 ing each image twice to create a pseudo-video input. Our method significantly outperforms both SIGMA variants and the VideoMAE baseline across all datasets, demonstrating its effectiveness 510 even without additional data augmentation. Notably, ViDROP achieves comparable or superior per-511 formance to SIGMA_{DINO}, which benefits from a teacher model pretrained on extra imagenet data. 512 The results also showcase ViDROP's versatility on DAVIS tracking Pont-Tuset et al. (2017) and 513 copy detection Pizzi et al. (2023) tasks. For DAVIS tracking, we adopt the DINO approach of 514 processing frames individually. The findings highlight ViDROP's substantial improvements over 515 VideoMAE baselines, particularly in copy detection where ViDROP achieves a remarkable 22.4 516 μ AP, compared to 2.0 for the VideoMAE baseline. These results underscore ViDROP's ability to 517 learn rich, transferable representations across diverse video understanding tasks.

- 518
- 519 520

5 CONCLUSION

521 522 523

We presented ViDROP, a novel self-supervised learning approach for video understanding that 524 combines token dropping and masking strategies within a single encoder. Our method achieves 525 state-of-the-art performance across various video understanding tasks, including action recognition, 526 temporal action detection, and low-shot learning, while maintaining computational efficiency. Key 527 innovations include a sparse encoder architecture that eliminates the need for a separate decoder, en-528 abling fine-grained representation learning without computational overhead, and a lossy data loading 529 scheme based on k-means clustering, significantly accelerating data processing while maintaining 530 compatibility with pretrained models. We demonstrated scalability across model sizes and initializa-531 tion strategies, including successful training of ViT-Huge models from scratch using limited computational resources. ViDROP's remarkable data efficiency in low-shot learning scenarios and its 532 effectiveness in both action classification and temporal action detection tasks highlight its versatility 533 and potential for real-world applications. 534

Looking forward, our work opens up possibilities for training state-of-the-art models with limited
 computational resources, making advanced video understanding more accessible. However, our re sults with VJEPA-initialized models highlight the importance of training on more diverse datasets to
 enhance generalization capabilities. Additionally, while our LEMON experiments for training ViT Huge models from scratch show promising results, they are not yet on par with models initialized
 from pretrained weights, indicating a need for further exploration of optimal scaling strategies.

540 REFERENCES 541

567

568

569

570

577

579

- Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lucic, and Cordelia Schmid. 542 Vivit: A video vision transformer. 2021 IEEE/CVF International Conference on Computer 543 Vision (ICCV), pp. 6816-6826, 2021. URL https://api.semanticscholar.org/ 544 CorpusID:232417054. 545
- 546 Mahmoud Assran, Mathilde Caron, Ishan Misra, Piotr Bojanowski, Florian Bordes, Pascal Vincent, 547 Armand Joulin, Michael G. Rabbat, and Nicolas Ballas. Masked siamese networks for labelefficient learning. In European Conference on Computer Vision, 2022. URL https://api. 548 semanticscholar.org/CorpusID:248178208. 549
- 550 Mahmoud Assran, Quentin Duval, Ishan Misra, Piotr Bojanowski, Pascal Vincent, Michael G. 551 Rabbat, Yann LeCun, and Nicolas Ballas. Self-supervised learning from images with a joint-552 embedding predictive architecture. 2023 IEEE/CVF Conference on Computer Vision and Pattern 553 Recognition (CVPR), pp. 15619-15629, 2023. URL https://api.semanticscholar. 554 org/CorpusID:255999752.
- 555 Alexei Baevski, Arun Babu, Wei-Ning Hsu, and Michael Auli. Efficient self-supervised learn-556 ing with contextualized target representations for vision, speech and language. In International Conference on Machine Learning, 2022a. URL https://api.semanticscholar.org/ 558 CorpusID:254685875. 559
- Alexei Baevski, Wei-Ning Hsu, Qiantong Xu, Arun Babu, Jiatao Gu, and Michael Auli. data2vec: 560 A general framework for self-supervised learning in speech, vision and language. ArXiv, 561 abs/2202.03555, 2022b. URL https://api.semanticscholar.org/CorpusID: 562 246652264. 563
- 564 Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. ArXiv, 565 abs/2106.08254, 2021. URL https://api.semanticscholar.org/CorpusID: 235436185. 566
 - Adrien Bardes, Quentin Garrido, Jean Ponce, Michael Rabbat, Yann LeCun, Mahmoud Assran, and Nicolas Ballas. Revisiting feature prediction for learning visual representations from video. arXiv:2404.08471, 2024.
- Florian Bordes, Randall Balestriero, and Pascal Vincent. Towards democratizing joint-571 embedding self-supervised learning. ArXiv, abs/2303.01986, 2023. URL https://api. 572 semanticscholar.org/CorpusID:257353289. 573
- 574 Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsu-575 pervised learning of visual features. In European Conference on Computer Vision, 2018. URL 576 https://api.semanticscholar.org/CorpusID:263891125.
- Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Un-578 supervised learning of visual features by contrasting cluster assignments. ArXiv, abs/2006.09882, 2020. URL https://api.semanticscholar.org/CorpusID:219721240. 580
- Mathilde Caron, Hugo Touvron, Ishan Misra, Herv'e J'egou, Julien Mairal, Piotr Bojanowski, and 581 Armand Joulin. Emerging properties in self-supervised vision transformers. 2021 IEEE/CVF 582 International Conference on Computer Vision (ICCV), pp. 9630-9640, 2021. URL https: 583 //api.semanticscholar.org/CorpusID:233444273. 584
- 585 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. A simple framework 586 for contrastive learning of visual representations. ArXiv, abs/2002.05709, 2020a. URL https: //api.semanticscholar.org/CorpusID:211096730.
- 588 Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 15745–15753, 2020. URL 590 https://api.semanticscholar.org/CorpusID:227118869.
- Xinlei Chen, Haoqi Fan, Ross B. Girshick, and Kaiming He. Improved baselines with 592 momentum contrastive learning. ArXiv, abs/2003.04297, 2020b. URL https://api. 593 semanticscholar.org/CorpusID:212633993.

- Ishan Rajendrakumar Dave, Simon Jenni, and Mubarak Shah. No more shortcuts: Realizing the potential of temporal self-supervision. ArXiv, abs/2312.13008, 2023. URL https://api.semanticscholar.org/CorpusID:266374984.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In North American Chapter of the Association for Computational Linguistics, 2019. URL https://api.semanticscholar.org/ CorpusID:52967399.
- Michael Dorkenwald, Fanyi Xiao, Biagio Brattoli, Joseph Tighe, and Davide Modolo. Scvrl:
 Shuffled contrastive video representation learning. 2022 IEEE/CVF Conference on Computer
 Vision and Pattern Recognition Workshops (CVPRW), pp. 4131–4140, 2022. URL https:
 //api.semanticscholar.org/CorpusID:249017621.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszko reit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition
 at scale. ArXiv, abs/2010.11929, 2020. URL https://api.semanticscholar.org/
 CorpusID:225039882.
- Alaaeldin El-Nouby, Michal Klein, Shuangfei Zhai, Miguel Angel Bautista, Alexander To shev, Vaishaal Shankar, Joshua M. Susskind, and Armand Joulin. Scalable pre-training of
 large autoregressive image models. ArXiv, abs/2401.08541, 2024. URL https://api.
 semanticscholar.org/CorpusID:267028705.
- Patrick Esser, Robin Rombach, and Björn Ommer. Taming transformers for high-resolution image synthesis. 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 12868–12878, 2020. URL https://api.semanticscholar.org/CorpusID: 229297973.
- Christoph Feichtenhofer, Haoqi Fan, Bo Xiong, Ross B. Girshick, and Kaiming He. A large-scale study on unsupervised spatiotemporal representation learning. 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 3298–3308, 2021. URL https://api.semanticscholar.org/CorpusID:233444206.
- Christoph Feichtenhofer, Haoqi Fan, Yanghao Li, and Kaiming He. Masked autoen coders as spatiotemporal learners. ArXiv, abs/2205.09113, 2022. URL https://api.
 semanticscholar.org/CorpusID:248863181.
- Evelyn Fix and Joseph L. Hodges. Discriminatory analysis nonparametric discrimination: Consistency properties. International Statistical Review, 57:238, 1989. URL https://api.semanticscholar.org/CorpusID:120323383.
- Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fründ, Peter N. Yianilos, Moritz Mueller-Freitag, Florian Hoppe, Christian Thurau, Ingo Bax, and Roland Memisevic. The "something something" video database for learning and evaluating visual common sense. 2017 IEEE International Conference on Computer Vision (ICCV), pp. 5843–5851, 2017. URL https://api. semanticscholar.org/CorpusID:834612.
- Jean-Bastien Grill, Florian Strub, Florent Altch'e, Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Ávila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own latent: A new approach to self-supervised learning. *ArXiv*, abs/2006.07733, 2020. URL https: //api.semanticscholar.org/CorpusID:219687798.
- Mark Hamilton, Zhoutong Zhang, Bharath Hariharan, Noah Snavely, and William T. Freeman. Un supervised semantic segmentation by distilling feature correspondences. *ArXiv*, abs/2203.08414, 2022. URL https://api.semanticscholar.org/CorpusID:247476291.
- Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Doll'ar, and Ross B. Girshick.
 Masked autoencoders are scalable vision learners. 2022 IEEE/CVF Conference on Computer
 Vision and Pattern Recognition (CVPR), pp. 15979–15988, 2021. URL https://api.semanticscholar.org/CorpusID:243985980.

660

661

662

663

667

668

677

685

686

687

688

689

693

694

- 648 Sun-Kyoo Hwang, Jaehong Yoon, Youngwan Lee, and Sung Ju Hwang. Everest: Efficient masked 649 video autoencoder by removing redundant spatiotemporal tokens. ArXiv, abs/2211.10636, 2022. 650 URL https://api.semanticscholar.org/CorpusID:259188150.
- Andrew Jaegle, Sebastian Borgeaud, Jean-Baptiste Alayrac, Carl Doersch, Catalin Ionescu, David 652 Ding, Skanda Koppula, Andrew Brock, Evan Shelhamer, Olivier J. H'enaff, Matthew M. 653 Botvinick, Andrew Zisserman, Oriol Vinyals, and João Carreira. Perceiver io: A general ar-654 chitecture for structured inputs & outputs. ArXiv, abs/2107.14795, 2021. URL https: 655 //api.semanticscholar.org/CorpusID:236635379. 656
- 657 S. Jenni, Givi Meishvili, and Paolo Favaro. Video representation learning by recognizing temporal 658 transformations. In European Conference on Computer Vision, 2020. URL https://api. 659 semanticscholar.org/CorpusID:220665754.
 - Y.-G. Jiang, J. Liu, A. Roshan Zamir, G. Toderici, I. Laptev, M. Shah, and R. Sukthankar. THUMOS challenge: Action recognition with a large number of classes. http://crcv.ucf.edu/ THUMOS14/, 2014.
- Neha Mukund Kalibhat, Warren Morningstar, Alex Bijamov, Luyang Liu, Karan Singhal, and 665 P. A. Mansfield. Disentangling the effects of data augmentation and format transform in self-666 supervised learning of image representations. ArXiv, abs/2312.02205, 2023. URL https: //api.semanticscholar.org/CorpusID:265659216.
- Will Kay, João Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijaya-669 narasimhan, Fabio Viola, Tim Green, Trevor Back, Apostol Natsev, Mustafa Suleyman, and An-670 drew Zisserman. The kinetics human action video dataset. ArXiv, abs/1705.06950, 2017. URL 671 https://api.semanticscholar.org/CorpusID:27300853. 672
- 673 Hilde Kuehne, Hueihan Jhuang, Estíbaliz Garrote, Tomaso A. Poggio, and Thomas Serre. Hmdb: A 674 large video database for human motion recognition. 2011 International Conference on Computer 675 *Vision*, pp. 2556–2563, 2011. URL https://api.semanticscholar.org/CorpusID: 676 206769852.
- Guillaume Leclerc, Andrew Ilyas, Logan Engstrom, Sung Min Park, Hadi Salman, and Aleksander 678 Madry. Ffcv: Accelerating training by removing data bottlenecks. 2023 IEEE/CVF Conference 679 on Computer Vision and Pattern Recognition (CVPR), pp. 12011-12020, 2023. URL https: 680 //api.semanticscholar.org/CorpusID:259224879. 681
- 682 Min-Seob Lee, Song Park, Byeongho Heo, Dongyoon Han, and Hyunjung Shim. Seit++: Masked 683 token modeling improves storage-efficient training. ArXiv, abs/2312.10105, 2023. URL https: 684 //api.semanticscholar.org/CorpusID:266348450.
 - Jie Lei, Linjie Li, Luowei Zhou, Zhe Gan, Tamara L. Berg, Mohit Bansal, and Jingjing Liu. Less is more: Clipbert for video-and-language learning via sparse sampling. 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 7327–7337, 2021. URL https: //api.semanticscholar.org/CorpusID:231880022.
- 690 Xingbin Liu, Jinghao Zhou, Tao Kong, Xianming Lin, and Rongrong Ji. Exploring target rep-691 resentations for masked autoencoders. ArXiv, abs/2209.03917, 2022. URL https://api. 692 semanticscholar.org/CorpusID:252118863.
 - Stuart P. Lloyd. Least squares quantization in pcm. IEEE Trans. Inf. Theory, 28:129–136, 1982. URL https://api.semanticscholar.org/CorpusID:10833328.
- 696 Ishan Misra, C. Lawrence Zitnick, and Martial Hebert. Shuffle and learn: Unsupervised learning 697 using temporal order verification. In European Conference on Computer Vision, 2016. URL https://api.semanticscholar.org/CorpusID:9348728.
- Théo Moutakanni, Maxime Oquab, Marc Szafraniec, Maria Vakalopoulou, and Piotr Bojanowski. 700 You don't need data-augmentation in self-supervised learning. ArXiv, abs/2406.09294, 2024. 701 URL https://api.semanticscholar.org/CorpusID:270440405.

702 703 704 705 706 707 708	Maxime Oquab, Timoth'ee Darcet, Théo Moutakanni, Huy Q. Vo, Marc Szafraniec, Vasil Khali- dov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Ass- ran, Nicolas Ballas, Wojciech Galuba, Russ Howes, Po-Yao (Bernie) Huang, Shang-Wen Li, Ishan Misra, Michael G. Rabbat, Vasu Sharma, Gabriel Synnaeve, Huijiao Xu, Hervé Jégou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning ro- bust visual features without supervision. <i>ArXiv</i> , abs/2304.07193, 2023. URL https://api. semanticscholar.org/CorpusID:258170077.
709 710 711 712	Song Park, Sanghyuk Chun, Byeongho Heo, Wonjae Kim, and Sangdoo Yun. Seit: Storage- efficient vision training with tokens using 1% of pixel storage. 2023 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 17202–17213, 2023. URL https://api. semanticscholar.org/CorpusID:257631701.
713 714 715 716 717	Ed Pizzi, Giorgos Kordopatis-Zilos, Hiral Patel, Gheorghe Postelnicu, Sugosh Nagavara Ravindra, Akshay Kumar Gupta, Symeon Papadopoulos, Giorgos Tolias, and Matthijs Douze. The 2023 video similarity dataset and challenge. <i>ArXiv</i> , abs/2306.09489, 2023. URL https://api.semanticscholar.org/CorpusID:259187976.
718 719 720	Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbeláez, Alexander Sorkine-Hornung, and Luc Van Gool. The 2017 davis challenge on video object segmentation. <i>ArXiv</i> , abs/1704.00675, 2017. URL https://api.semanticscholar.org/CorpusID:3619941.
721 722 723 724	Rui Qian, Tianjian Meng, Boqing Gong, Ming-Hsuan Yang, H. Wang, Serge J. Belongie, and Yin Cui. Spatiotemporal contrastive video representation learning. 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 6960–6970, 2020. URL https://api.semanticscholar.org/CorpusID:221090567.
725 726 727 728 729	Kanchana Ranasinghe, Muzammal Naseer, Salman Hameed Khan, Fahad Shahbaz Khan, and Michael S. Ryoo. Self-supervised video transformer. 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 2864–2874, 2021. URL https://api.semanticscholar.org/CorpusID:244800737.
730 731 732 733 734	Adrià Recasens, Pauline Luc, Jean-Baptiste Alayrac, Luyu Wang, Florian Strub, Corentin Tallec, Mateusz Malinowski, Viorica Patraucean, Florent Altch'e, Michael Valko, Jean-Bastien Grill, Aäron van den Oord, and Andrew Zisserman. Broaden your views for self-supervised video learning. 2021 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 1235–1245, 2021. URL https://api.semanticscholar.org/CorpusID:232417490.
735 736 737 738	Mohammadreza Salehi, Michael Dorkenwald, Fida Mohammad Thoker, Efstratios Gavves, Cees G. M. Snoek, and Yuki M. Asano. Sigma: Sinkhorn-guided masked video mod- eling. <i>ArXiv</i> , abs/2407.15447, 2024. URL https://api.semanticscholar.org/ CorpusID:271328916.
739 740 741 742 743	Sepehr Sameni, Simon Jenni, and Paolo Favaro. Spatio-temporal crop aggregation for video representation learning. 2023 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 5641–5651, 2022. URL https://api.semanticscholar.org/CorpusID: 254096149.
744 745 746	Khurram Soomro, Amir Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human ac- tions classes from videos in the wild. <i>ArXiv</i> , abs/1212.0402, 2012. URL https://api. semanticscholar.org/CorpusID:7197134.
747 748 749	Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are data- efficient learners for self-supervised video pre-training. <i>ArXiv</i> , abs/2203.12602, 2022. URL https://api.semanticscholar.org/CorpusID:247619234.
750 751 752 753	Aäron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. ArXiv, abs/1711.00937, 2017. URL https://api.semanticscholar.org/ CorpusID:20282961.
754 755	Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predic- tive coding. ArXiv, abs/1807.03748, 2018. URL https://api.semanticscholar.org/ CorpusID:49670925.

756	Diane Wagner, Fábio Ferreira, Daniel Stoll, Robin Tibor Schirrmeister, Samuel G, Müller, and
757	Frank Hutter. On the importance of hyperparameters and data augmentation for self-supervised
758	learning. ArXiv, abs/2207.07875, 2022. URL https://api.semanticscholar.org/
759	CorpusID:250627318.

- Limin Wang, Bingkun Huang, Zhiyu Zhao, Zhan Tong, Yinan He, Yi Wang, Yali Wang, and Yu Qiao. Videomae v2: Scaling video masked autoencoders with dual masking. 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 14549–14560, 2023a. URL https://api.semanticscholar.org/CorpusID:257805127.
- Yite Wang, Jiahao Su, Hanlin Lu, Cong Xie, Tianyi Liu, Jianbo Yuan, Haibin Lin, Ruoyu Sun, and Hongxia Yang. Lemon: Lossless model expansion. ArXiv, abs/2310.07999, 2023b. URL https://api.semanticscholar.org/CorpusID:263909329.
- Olivia Wiles, João F. M. Carreira, Iain Barr, Andrew Zisserman, and Mateusz Malinowski. Compressed vision for efficient video understanding. In *Asian Conference on Computer Vision*, 2022. URL https://api.semanticscholar.org/CorpusID:252735173.
- Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu. Simmim: a simple framework for masked image modeling. 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 9643–9653, 2021. URL https://api.semanticscholar.org/CorpusID:244346275.
- Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self supervised learning via redundancy reduction. ArXiv, abs/2103.03230, 2021. URL https:
 //api.semanticscholar.org/CorpusID:232110471.
- Runtian Zhai, Bing Liu, Andrej Risteski, Zico Kolter, and Pradeep Ravikumar. Understanding augmentation-based self-supervised representation learning via rkhs approximation. ArXiv, abs/2306.00788, 2023. URL https://api.semanticscholar.org/CorpusID: 258999873.
- Chen-Lin Zhang, Jianxin Wu, and Yin Li. Actionformer: Localizing moments of actions with transformers. In *European Conference on Computer Vision*, volume 13664 of *LNCS*, pp. 492–510, 2022.
- Yue Zhao and Philipp Krahenbuhl. Training a large video model on a single machine in a day. ArXiv, abs/2309.16669, 2023. URL https://api.semanticscholar.org/ CorpusID:263135660.
- Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Loddon Yuille, and Tao Kong.
 ibot: Image bert pre-training with online tokenizer. *ArXiv*, abs/2111.07832, 2021. URL https:
 //api.semanticscholar.org/CorpusID:244117494.

804 805

808 809