062

063

064

065

066

067

068

069

070

071

072

073

Latent Video Dataset Distillation

Anonymous CVPR submission

Paper ID 11

Abstract

001 Dataset distillation has demonstrated remarkable effectiveness in high-compression scenarios for image datasets. 002 003 While video datasets inherently contain greater redundancy, existing video dataset distillation methods primarily 004 focus on compression in the pixel space, overlooking ad-005 vances in the latent space that have been widely adopted 006 in modern text-to-image and text-to-video models. In this 007 008 work, we bridge this gap by introducing a novel video dataset distillation approach that operates in the latent 009 space using a state-of-the-art variational encoder. Fur-010 thermore, we employ a diversity-aware data selection strat-011 egy to select both representative and diverse samples. Ad-012 013 ditionally, we introduce a simple, training-free method to 014 further compress the distilled latent dataset. By combining these techniques, our approach achieves a new state-of-015 the-art performance in dataset distillation, outperforming 016 prior methods on all datasets, e.g. on HMDB51 IPC 1, we 017 018 achieve a 2.6% performance increase; on MiniUCF IPC 5, 019 we achieve a 7.8% performance increase.

1. Introduction

Dataset distillation has emerged as a pivotal technique for 021 022 compressing large-scale datasets into computationally efficient representations that retain their essential characteris-023 tics [38]. While this technique has seen remarkable success 024 in compressing image datasets [4, 5, 22, 27, 36, 44], appli-025 026 cations onto video datasets remain an underexplored chal-027 lenge. Videos inherently possess temporal redundancy, as characterized by consecutive frames often sharing substan-028 tial similarity, presenting the potential for optimization via 029 030 dataset distillation.

Existing video distillation methods predominantly focus on pixel-space compression. VDSD [39] addresses the temporal information redundancy by disentangling static and dynamic information. Method IDTD [48] tackles the within-sample and inter-sample redundancies by leveraging a joint-optimization framework. However, these frameworks overlook the potential of latent-space compressions, which have proven transformative in generative models for
images and videos [34, 47]. Modern variational autoen-
coders (VAEs) [29, 40] offer a pathway to address this gap
by encoding videos into compact, disentangled representa-
tions in latent space.038
039040
041040

In this work, we improve video distillation by operating 043 entirely in the latent space of a VAE. Our framework dis-044 tills videos into low-dimensional latent codes, leveraging 045 the VAE's ability to model temporal dynamics [47]. Un-046 like previous methods, our approach encodes entire video 047 sequences into coherent latent trajectories to model tem-048 poral dynamics through its hierarchical architecture. We 049 compress the VAE itself through post-training quantiza-050 tion, largely reducing the model size, while retaining accu-051 racy [5]. After distillation, we apply Diversity-Aware Data 052 Selection using Determinantal Point Processes (DPPs) [17] 053 to select both representative and diverse instances. Unlike 054 clustering-based or random sampling methods, DPPs in-055 herently favor diversity by selecting samples that are well-056 spread in the latent space, reducing redundancy while en-057 suring comprehensive feature coverage [26]. This leads to 058 a more informative distilled dataset that enhances down-059 stream model generalization. 060

Our method further introduces a training-free latent compression strategy, which uses high-order singular value decomposition (HOSVD) to decompose spatiotemporal features into orthogonal subspaces [39]. This isolates dominant motion patterns and spatial structures, enabling further compression while preserving essential dynamics [34]. By factorizing latent tensors, we dynamically adjust the rank of the distilled representations, allowing denser instance packing under fixed storage limits. Experiments on the Mini-UCF dataset demonstrate that our method outperforms prior pixel-space approaches by 11.5% in absolute accuracy for IPC 1 and 7.8% for IPC 5.

Overall, our contributions are:

- We propose the first video dataset distillation framework operating in the latent space, leveraging a state-of-the-art VAE to efficiently encode spatiotemporal dynamics.
 074 075 076
- We address the challenge of sparsity in the video latent space by integrating Diversity-Aware Data Selection us-078

160

161

162

163

164

165

166

167

168

169

170

171

ing DPPs and High-Order Singular Value Decomposi-tion (HOSVD) for structured compression.

Our method generalizes to both small-scale and large scale video datasets, achieving a new state-of-the-art per formance on all settings compared to existing methods.

084 2. Related Work

Coreset Selection Coreset selection aims to identify a small 085 but representative subset of data that preserves the essential 086 properties of the full dataset, reducing computational com-087 plexity while maintaining model performance. One of the 088 foundational approaches utilizes k-center clustering [30] to 089 090 formulate coreset selection as a geometric covering problem, where a subset of data points is chosen to maxi-091 092 mize the minimum distance to previously selected points. By iteratively selecting the most distant samples in feature 093 space, this method ensures that the coreset provides broad 094 coverage of the dataset's distribution, making it a strong 095 096 candidate for reducing redundancy in large-scale datasets. Herding methods [40] take an optimization-driven approach 097 to coreset selection by sequentially choosing samples that 098 best approximate the mean feature representation of the 099 dataset. Probabilistic techniques leverage Bayesian infer-100 101 ence [24] and divergence minimization [33] to construct coresets that balance diversity and statistical representative-102 ness. Influence-based selection methods [41] instead fo-103 104 cus on quantifying the contribution of individual samples to generalization performance, retaining only the most im-105 pactful data points. 106

Image Dataset Distillation Dataset distillation [38] has 107 emerged as a powerful paradigm for compressing large-108 109 scale image datasets while preserving downstream task performance. Early gradient-based methods like Dataset Dis-110 tillation (DD) [38] optimized synthetic images by match-111 ing gradients between training trajectories on original and 112 113 distilled datasets. Later works introduced dataset condensation with gradient matching [46]. Further, Meta-learning 114 115 frameworks Like Matching Training Trajectories (MTT) [3] and Kernel Inducing Points (KIP) [28] advances perfor-116 mance by distilling datasets through bi-level optimization 117 over neural architectures. Dataset condensation with Distri-118 119 bution Matching (DM) [45] synthesizes condensed datasets 120 by aligning feature distributions between original and syn-121 thetic data across various embedding spaces.

Representative Matching for Dataset Condensation 122 (DREAM) [21] improved sample efficiency by selecting 123 representative instances that retained the most informative 124 125 patterns from the original dataset, reducing redundancy in 126 synthetic samples. Generative modeling techniques have also been explored, with Distilling Datasets into Generative 127 Models (DiM) [37] encoding datasets into latent generative 128 spaces, allowing for smooth interpolation and novel sample 129 130 generation. Similarly, Hybrid Generative-Discriminative

Dataset Distillation (GDD) [19] balanced global structural131coherence with fine-grained detail preservation by combin-
ing adversarial generative models with traditional distilla-
tion objectives. However, temporal redundancy and frame
sampling complexities, as noted in [11, 20], highlight the
unique difficulties of extending image-focused distillation
to video datasets.131

Video Dataset Distillation While dataset distillation has 138 achieved significant success in static image datasets, di-139 rect application to videos presents unique challenges due 140 to temporal redundancy and the need for efficient frame 141 selection [34]. Recent attempts to address video dataset 142 distillation have primarily focused on pixel-space com-143 pression. Video Distillation via Static-Dynamic Disentan-144 glement (VDSD) [39] tackles temporal redundancies be-145 tween frames by separating static and dynamic components. 146 VDSD partitions videos into smaller segments and employs 147 learnable dynamic memory block that captures and syn-148 thesizes motion patterns, improving information retention 149 while reducing redundancy. IDTD [48] addresses the chal-150 lenges of within-sample redundancy and inter-sample re-151 dundancy simultaneously. IDTD employs an architecture 152 represented by a shared feature pool alongside multiple fea-153 ture selectors to selectively condense video sequences while 154 ensuring sufficient motion diversity. To retain the temporal 155 information of synthesized videos, IDTD introduces a Tem-156 poral Fusor that integrates diverse features into the temporal 157 dimension. 158

Text-to-Video Models and Their Role in Latent Space Learning Latent-space representations have become a cornerstone of modern video modeling, offering structured compression while maintaining high-level semantic integrity [34, 47]. Variational autoencoders provide to enable efficient storage and reconstruction [13]. Extending this concept, hierarchical autoregressive latent prediction [31] introduces an autoregressive component that improves temporal coherence, leading to high-fidelity video reconstructions. Further enhancing latent representations, latent video diffusion transformers [23] incorporate diffusion-based priors to refine video quality while minimizing storage demands.

Building upon these latent space techniques, recent text-172 to-video models have demonstrated their capability to gen-173 erate high-resolution video content from textual descrip-174 tions. These methods employ a combination of transformer-175 based encoders and diffusion models to synthesize realistic 176 video sequences. Imagen Video leverages cascaded video 177 diffusion models to progressively upsample spatial and tem-178 poral dimensions, ensuring high-quality output [9]. Mean-179 while, zero-shot generation approaches utilize decoder-180 only transformer architectures to process multimodal in-181 puts, such as text and images, without requiring explicit 182 video-text training data [15]. Hybrid techniques combin-183



Figure 1. Our training-free latent video distillation pipeline. The entire video dataset is encoded into latent space with a VAE. We further employ the DPPs to select both representative and diverse samples, followed by latent space compression with HOSVD for efficient storage.

ing pixel-space and latent-space diffusion modeling further
enhance computational efficiency while maintaining visual
fidelity by leveraging learned latent representations during
synthesis [43]. These advancements in latent space learning not only improve video compression but also drive the
development of scalable and high-quality text-driven video
generation.

191 3. Methodology

In this section, we first introduce the variational autoen-192 193 coder (VAE) used to encode video sequences into a com-194 pact latent space. We then discuss our Diversity-Aware Data Selection method. Next, we present our training-free la-195 tent space compression approach using High-Order Singu-196 lar Value Decomposition (HOSVD). Finally, we describe 197 198 our two-stage dynamic quantization strategy. The entire pipeline of our framework is shown in Fig. 1. 199

3.1. Preliminary

Problem Definition In video dataset distillation, given a 201 large dataset $\mathcal{T} = \{(x_i, y_i)\}_{i=1}^{|\mathcal{T}|}$ consisting of video sam-202 ples x_i and their corresponding class labels y_i , the objec-203 tive is to construct a significantly smaller distilled dataset 204 $\mathcal{S} = \{\tilde{x}_i, \tilde{y}_i\}_{i=1}^{|\mathcal{S}|}$, where $|\mathcal{S}| \ll |\mathcal{T}|$. The distilled dataset 205 is expected to achieve comparable performance to the orig-206 inal dataset on action classification tasks while significantly 207 208 reducing storage and computational requirements.

Latent image distillation has emerged as an effective alternative to traditional dataset distillation methods. Instead of distilling datasets at the pixel level, latent distillation leverages pre-trained autoencoders or generative models to encode images into a compact latent space. Latent Dataset Distillation with Diffusion Models [25], have

demonstrated that distilling image datasets in the latent 215 space of a pre-trained diffusion model improves generaliza-216 tion and enables higher compression ratios while maintain-217 ing image quality. Similarly, Dataset Distillation in Latent 218 Space [6] adapts conventional distillation methods like Gra-219 dient Matching, Feature Matching, and Parameter Match-220 ing to the latent space, significantly reducing computational 221 overhead while achieving competitive performance. Differ-222 ent from these methods, we extend latent space distillation 223 to video datasets by encoding both spatial and temporal in-224 formation into the latent space. 225

3.2. Variational Autoencoder

Variational Autoencoders (VAEs) are a class of generative227models that encode input data into a compact latent space228while maintaining the ability to reconstruct the original229data [13]. Unlike traditional autoencoders, VAEs enforce a230probabilistic structure on the latent space by learning a dis-231tribution rather than a fixed mapping. This allows for better232generalization and meaningful latent representations.233

A VAE consists of an encoder and a decoder. The en-234 coder maps the input x to a latent distribution $q_{\phi}(z|x)$, 235 where z is the latent variable. Instead of producing a deter-236 ministic latent representation, the encoder outputs the mean 237 and variance of a Gaussian distribution, from which sam-238 ples are drawn. This ensures that the latent space remains 239 continuous, facilitating smooth interpolation between data 240 points [14]. The decoder then reconstructs the input x from 241 a sampled latent variable z, following the learned distribu-242 tion $p_{\theta}(x|z)$. 243

To ensure a structured latent space, VAEs introduce a244regularization term that aligns the learned distribution with
a prior distribution, typically a standard normal distribution245

280

282

313

314

315

316

317

247 $p(z) = \mathcal{N}(0, 1)$. This prevents the model from collapsing 248 into a purely memorized representation of the data, encour-249 aging better generalization.

The training objective of a VAE is to maximize the Ev-250 251 idence Lower Bound (ELBO) [18], which consists of two terms: reconstruction loss and Kullback-Leibler (KL) di-252 vergence regularization [13]. The reconstruction loss en-253 sures that the decoded output remains similar to the original 254 255 input, while the KL divergence forces the learned latent distribution to be close to the prior, preventing overfitting and 256 257 promoting smoothness in the latent space. The overall loss function is formulated as follows. 258

$$\mathcal{L}_{\text{VAE}} = \mathbb{E}_{q_{\phi}(z|x)} [-\log p_{\theta}(x|z)] + \beta \cdot D_{\text{KL}}(q_{\phi}(z|x) \parallel p_{\theta}(z))$$
(1)

260 3.3. Diversity-Aware Data Selection

After encoding the entire video dataset into the latent space 261 using a state-of-the-art VAE, an effective data selection 262 strategy is crucial to maximize the diversity and represen-263 264 tativeness of the distilled dataset. To this end, we employ Diversity-Aware Data Selection using Determinantal Point 265 266 Processes (DPPs) [17], a principled probabilistic framework that promotes diversity by favoring sets of samples 267 that are well-spread in the latent space. 268

DPPs [17] provide a natural mechanism for selecting a 269 subset of latent embeddings that balance coverage and in-270 271 formativeness while reducing redundancy. Given the en-272 coded latent representations of the dataset, we construct a similarity kernel matrix L, where each entry L_{ij} quantifies 273 the pairwise similarity between latent samples z_i and z_j . 274 275 The selection process then involves sampling from a determinantal distribution parameterized by L, ensuring that the 276 277 chosen subset is both diverse and representative of the full 278 latent dataset. We define a kernel matrix L using the fol-279 lowing function:

$$L_{ij} = \exp(-\frac{\|z_i - z_j\|^2}{2\sigma^2})$$
(2)

281 Then subset *S* is sampled according to:

$$P(S) = \frac{\det(L_S)}{\det(L+I)} \tag{3}$$

here L_S is the submatrix of L that corresponds to the rows and columns indexed by S. The denominator det(L + I)serves as a normalization factor, ensuring that the probabilities across all possible subsets sum to 1. This normalization stabilizes the sampling process by incorporating an identity matrix I, which prevents numerical instability in cases where L is near-singular.

290 Our approach is motivated by the observation that naive 291 random sampling or traditional clustering-based selection strategies [12] tend to underperform in high-dimensional292latent spaces [7], where redundancy is prevalent. By lever-293aging DPPs, we effectively capture a more comprehensive294distribution of video features, thereby improving the qual-295ity of the distilled dataset. Furthermore, the computational296efficiency of DPPs sampling allows us to scale our selection297process to large datasets without significant overhead.298

Applying DPPs in the latent space instead of the pixel 299 space offers several key advantages. First, latent representa-300 tions encode high-level semantic features, making it possi-301 ble to directly select samples that preserve meaningful vari-302 ations in motion and structure, rather than relying on pixel-303 wise differences that may be redundant or noisy. Second, 304 the latent space is significantly more compact and disen-305 tangled, allowing DPPs to operate more effectively with re-306 duced computational complexity compared to pixel-space 307 selection [39], which often involves large-scale feature ex-308 traction. Finally, in the latent space, similarity measures are 309 inherently more structured, which makes DPPs better suited 310 for ensuring diverse and representative selections that gen-311 eralize well to downstream tasks. 312

3.4. Training-free Latent Space Compression

While our Diversity-Aware Data Selection effectively distills a compact subset of the latent dataset, we observe that the selected latent representations remain sparse, leading to inefficiencies in storage and downstream processing.

Singular Value Decomposition (SVD) is a fundamental318matrix factorization technique widely used in dimensional-
ity reduction, data compression, and noise filtering. Given
a matrix $X \in \mathbb{R}^{m \times n}$, SVD decomposes it into three com-
ponents:320321

$$X = U\Sigma V^T \tag{4} 323$$

where U is an orthogonal matrix whose columns represent 324 the left singular vectors, Σ is a diagonal matrix containing 325 the singular values that indicate the importance of each cor-326 responding singular vector, and V is an orthogonal matrix 327 whose columns represent the right singular vectors. A key 328 property of SVD is that truncating the smaller singular val-329 ues allows for an effective low-rank approximation of the 330 original matrix, reducing storage requirements while pre-331 serving essential information. This property makes SVD 332 particularly useful in data compression and feature selec-333 tion. 334

However, when applied to higher-dimensional data, such335as video representations in latent space, SVD requires flat-336tening the tensor into a 2D matrix, which disrupts spa-337tial and temporal correlations. This limitation motivates338our adoption of High-Order Singular Value Decomposition339(HOSVD), which extends SVD to multi-dimensional tensors while preserving their inherent structure.341

| Dataset | | Mini | UCF | HMDB51 Kinetics-400 | | SSv2 | | | |
|----------------------|-----------------|----------------|----------------|---------------------|---------------|----------------|----------------|----------------|---------------|
| IF | C | 1 | 5 | 1 | 5 | 1 | 5 | 1 | 5 |
| Full Dataset | | 57.2 ± 0.1 | | 28.6 ± 0.7 | | 34.6 ± 0.5 | | 29.0 ± 0.6 | |
| | Random | 9.9 ± 0.8 | 22.9 ± 1.1 | 4.6 ± 0.5 | 6.6 ± 0.7 | 3.0 ± 0.1 | 5.6 ± 0.0 | 3.2 ± 0.1 | 3.7 ± 0.0 |
| Coreset Selection | Herding [40] | 12.7 ± 1.6 | 25.8 ± 0.3 | 3.8 ± 0.2 | 8.5 ± 0.4 | 4.3 ± 0.3 | 8.0 ± 0.1 | 4.6 ± 0.3 | 6.8 ± 0.2 |
| | K-Center [30] | 11.5 ± 0.7 | 23.0 ± 1.3 | 3.1 ± 0.1 | 5.2 ± 0.3 | 3.9 ± 0.2 | 5.9 ± 0.4 | 3.8 ± 0.5 | 4.0 ± 0.1 |
| Dataset Distillation | DM [45] | 15.3 ± 1.1 | 25.7 ± 0.2 | 6.1 ± 0.2 | 8.0 ± 0.2 | 6.3 ± 0.0 | 9.1 ± 0.9 | 4.1 ± 0.4 | 4.5 ± 0.3 |
| | MTT [3] | 19.0 ± 0.1 | 28.4 ± 0.7 | 6.6 ± 0.5 | 8.4 ± 0.6 | 3.8 ± 0.2 | 9.1 ± 0.3 | 3.9 ± 0.2 | 6.5 ± 0.2 |
| | FRePo [49] | 20.3 ± 0.5 | 30.2 ± 1.7 | 7.2 ± 0.8 | 9.6 ± 0.7 | _ | _ | _ | _ |
| | DM+VDSD [39] | 17.5 ± 0.1 | 27.2 ± 0.4 | 6.0 ± 0.4 | 8.2 ± 0.1 | 6.3 ± 0.2 | 7.0 ± 0.1 | 4.3 ± 0.3 | 4.0 ± 0.3 |
| | MTT+VDSD [39] | 23.3 ± 0.6 | 28.3 ± 0.0 | 6.5 ± 0.1 | 8.9 ± 0.6 | 6.3 ± 0.1 | 11.5 ± 0.5 | 5.7 ± 0.2 | 8.4 ± 0.1 |
| | FRePo+VDSD [39] | 22.0 ± 1.0 | 31.2 ± 0.7 | 8.6 ± 0.5 | 10.3 ± 0.6 | _ | _ | _ | _ |
| | IDTD [48] | 22.5 ± 0.1 | 33.3 ± 0.5 | 9.5 ± 0.3 | 16.2 ± 0.9 | 6.1 ± 0.1 | 12.1 ± 0.2 | _ | _ |
| | Ours | 34.8 ± 0.5 | 41.1 ± 0.6 | 12.1 ± 0.3 | 17.6 ± 0.4 | 9.0 ± 0.1 | 13.8 ± 0.1 | 6.9 ± 0.6 | 10.5 ± 0.4 |

Table 1. Performance comparison between our method and existing baselines on both small-scale and large-scale datasets. Follow previous works, we report Top-1 test accuracies (%) for small-scale datasets and Top-5 test accuracies (%) for large-scale datasets.

HOSVD is a tensor decomposition technique that gen-eralizes traditional SVD to higher-dimensional data. By treating the selected latent embeddings as a structured ten-sor rather than independent vectors, we exploit multi-modal correlations across feature dimensions to achieve more effi-cient compression. Specifically, given a set of selected la-tent embeddings $Z \in R^{d_1 \times d_2 \times \cdots \times d_n}$, we decompose it into a core tensor \mathcal{G} and a set of orthonormal factor matrices U_i , such that

$$Z = \mathcal{G} \times_1 U_1 \times_2 U_2 \times \dots \times_n U_n \tag{5}$$

where \times_i denotes the mode- *i* tensor-matrix product. By truncating the singular values in each mode with a rank compression ratio, we discard low-energy components while preserving the most informative structures in the latent space.

A key advantage of HOSVD over traditional SVD is its ability to retain the original tensor structure, rather than re-quiring flattening into a 2D matrix. More importantly, trun-cating the singular values in the temporal mode directly reduces temporal redundancy, ensuring that only the most rep-resentative motion patterns are retained. This enables more efficient storage and reconstruction, while minimizing the loss of critical temporal information.

Unlike conventional post-hoc compression techniques that require fine-tuning or retraining, HOSVD operates in a completely training-free manner, making it highly efficient and scalable. Furthermore, our empirical analysis shows that applying HOSVD after DPPs-based selection leads to a substantial reduction in storage and computational re-quirements while maintaining near-optimal performance in downstream tasks.

By integrating HOSVD into our dataset distillation
pipeline, we achieve an additional compression gain with
minimal loss of information, further pushing the boundaries
of efficiency in video dataset distillation.

3.5. VAE Quantization

To further improve storage efficiency, we apply a two-stage quantization process to the 3D-VAE [47], combining dynamic quantization for fully connected layers and mixedprecision optimization for all other layers.

The first stage involves dynamic quantization, where all382fully connected layers are reduced from 32-bit floating-383point to 8-bit integer representations. Dynamic quantization384works by scaling activations and weights dynamically dur-385ing inference. Formally, given an activation x and weight386matrix W, the quantized representation is computed as:387

$$W_q = \operatorname{round}\left(\frac{W}{s_W}\right) + z_W, \quad x_q = \operatorname{round}\left(\frac{x}{s_x}\right) + z_x$$
 (6) 388

where s_W and s_x are learned scaling factors, and z_W and z_x are zero points for weight and activation quantization, respectively. The dynamically scaled operation ensures that numerical stability is preserved while reducing the model size. This quantization is applied to all fully connected layers in the encoder and decoder, allowing for efficient memory compression without requiring retraining.

Unlike convolutional layers, fully connected layers pri-marily perform matrix multiplications, which exhibit high redundancy and are well-suited for integer quantization (INT8). Quantizing these layers from FP32 to INT8 signif-icantly reduces memory consumption and improves com-putational efficiency while maintaining inference stabil-ity [10]. Since fully connected layers do not require the high dynamic range of floating-point precision, INT8 quan-tization achieves optimal storage and performance benefits.

In the second stage, we employ mixed-precision optimization, where all remaining convolutional and batch normalization layers undergo reduced-precision floating-point compression, scaling them from FP32 to FP16. Unlike integer quantization, FP16 maintains a wider dynamic range, preventing significant loss of information in convolutional layers, which are more sensitive to precision reduction [42].

462

463

464

465

466

467

468

469

470

471

472

473

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

This hybrid quantization approach balances storage efficiency and numerical precision, ensuring that the 3D-VAE remains compact while preserving its ability to model spatiotemporal dependencies in video sequences. While applying post-training dynamic quantization on CV-VAE[47], we achieve a more than 2.6× compression ratio while maintaining high reconstruction fidelity.

419 **4.** Experiments

420 4.1. Datasets and Metrics

Following previous works VDSD [39] and IDTD [48], we 421 422 evaluate our proposed video dataset distillation approach 423 on both small-scale and large-scale benchmark datasets. For small-scale datasets, we utilize MiniUCF [39] and 424 HMDB51 [16], while for large-scale datasets, we con-425 duct experiments on Kinetics [2] and Something-Something 426 V2 (SSv2) [8]. MiniUCF is a miniaturized version of 427 428 UCF101 [32], consisting of the 50 most common action classes selected from the original UCF101 dataset. 429 HMDB51 is a widely used human action recognition dataset 430 containing 6,849 video clips across 51 action categories. 431 432 Kinetics is a large-scale video action recognition dataset, 433 available in different versions covering 400, 600, or 700 human action classes. SSv2 is a motion-centric video dataset 434 comprising 174 action categories. 435

436 4.2. Baselines

Based on previous work, we include the following base-437 438 line: (1) coreset selection methods such as random selection, Herding [40], and K-Center [30], and (2) dataset dis-439 tillation methods including DM [45], MTT [3], FRePo [49], 440 441 VDSD [39], and IDTD [48]. DM [45] ensures that the models trained on the distilled dataset produce gradient updates 442 similar to those trained on the full dataset. MTT [3] im-443 444 proves distillation by aligning model parameter trajectories between the synthetic and original datasets. FRePo [49] 445 focuses on generating compact datasets that allow pre-446 trained models to quickly recover their original perfor-447 mance with minimal training. VDSD [39] introduces a 448 449 static-dynamic disentanglement approach for video dataset distillation. IDTD [48] enhances video dataset distillation 450 by increasing feature diversity across samples while densi-451 fying temporal information within instances. 452

453 4.3. Implementation Details

Dataset Details For small-scale datasets, MiniUCF and HMDB51, we follow the settings from previous work [39, 48], where videos are dynamically sampled to 16 frames with a sampling interval of 4. Each sampled frame is then cropped and resized to 112×112 resolution. We adopt the same settings as prior work [39, 48] for Kinetics-400, each video is sampled to 8 frames and resized to 64×64 , maintaining a compact representation suitable for large-scale dataset distillation. In Something-Something V2 (SSv2), which is relatively smaller among the two largescale datasets, we sample 16 frames per video and resize them to 112×112, demonstrating the scalability of our method across datasets of varying sizes.

Evaluation Network Following the previous works, we use a 3D convolutional network, C3D [35] as the evaluation network. C3D [35] is trained on the distilled datasets generated by our method. Similar to previous works, we assess the performance of our distilled datasets by measuring the top-1 accuracy on small-scale datasets and top-5 accuracy on large-scale datasets.

Fair Comparison Throughout our experiments, we rig-474 orously ensure that the total storage space occupied by the 475 quantized VAE model and the decomposed matrices remain 476 within the constraints of the corresponding Instance Per 477 Class (IPC) budget. Specifically, on SSv2, our method uti-478 lizes no more than 68% of the storage space allocated to 479 the baseline methods DM and MTT, guaranteeing a fair and 480 consistent comparison. A comprehensive analysis of fair 481 comparisons across all four video datasets is provided in 482 the supplementary material. 483

4.4. Experimental Results

In Tab. 1, we present the performance of our method across MiniUCF [39], HMDB51 [16], Kinetics-400 [2], and SSv2 [8] under both IPC 1 and IPC 5 settings.

On MiniUCF, our approach outperforms the best baseline (IDTD) by 12.3% under IPC 1, achieving 34.8% accuracy compared to 22.5%, and by 7.8% under IPC 5, reaching 41.1% accuracy. Similarly, on HMDB51, our method achieves 12.1% accuracy under IPC 1, surpassing the strongest baseline by 2.6%, while under IPC 5, it reaches 17.6%, a 1.4% improvement. These results highlight the effectiveness of our latent-space distillation framework, which provides superior compression efficiency and classification performance compared to pixel-space-based approaches. The consistent performance gains across both IPC settings demonstrate the robustness of our method in preserving essential video representations while achieving high compression efficiency.

Furthermore, the results in Kinetics-400 and SSv2 re-502 inforce our findings, as our approach consistently outper-503 forms all baselines. Improvements in low-IPC regimes 504 (IPC 1) suggest that our training-free latent compression 505 and diversity-aware data selection are particularly effective 506 when dealing with extreme data reduction. Our method 507 achieves 9.0% accuracy on Kinetics-400 IPC 1, outperform-508 ing the strongest baseline (IDTD) by 2.9%, and 6.9% ac-509 curacy on SSv2 IPC 1, surpassing VDSD by 2.2%. The 510 trend continues in IPC 5, where our model achieves 13.8% 511 on Kinetics-400 and 10.5% on SSv2, both establishing new 512

559

560

561

562



Figure 2. Comparison between different dataset distillation methods and data sampling methods on the SSv2 when IPC is 1.

513 state-of-the-art results in video dataset distillation.

4.5. Ablation Study 514

In this section, we systematically analyze the key compo-515 nents of our method to understand their contributions to 516 517 overall performance. We evaluate cross-architecture gen-518 eralization, various sampling methods, different rank compression ratios in HOSVD, and different latent space com-519 520 pression techniques.

Cross Architecture Generalization To further evalu-521 ate the generalization capability of our method, we con-522 523 duct experiments on cross-architecture generalization, as presented in Tab. 2. The results demonstrate that datasets 524 525 distilled using our method consistently achieve superior performance across different evaluation models-ConvNet3D, 526 CNN+GRU, and CNN+LSTM-compared to previous 527 state-of-the-art methods. 528

Our approach achieves 34.8% accuracy with Con-529 vNet3D, significantly surpassing all baselines, including 530 531 MTT+VDSD (23.3%) and DM+VDSD (17.5%). Notably, our method also outperforms all baselines when eval-532 uated on recurrent-based architectures (CNN+GRU and 533 CNN+LSTM), obtaining 19.9% and 18.3% accuracy, re-534 spectively. This highlights the robustness of our distilled 535 536 dataset in preserving spatiotemporal coherence, which is crucial for models that leverage sequential dependencies. 537

Sampling Methods We evaluate the impact of different 538 sampling strategies on dataset distillation, comparing our 539 540 Diversity-Aware Data Selection using Determinantal Point 541 Processes (DPPs) against random sampling, Kmeans clustering [12], and prior dataset distillation methods (DM + 542 VDSD, MTT + VDSD, IDTD). As shown in Fig. 2, our 543 method achieves the highest performance, demonstrating 544 the effectiveness of DPP-based selection in video dataset 545 546 distillation.

Among sampling strategies, DPPs-only selection outper-547 forms Kmeans and random sampling, indicating that DPPs 548 promote a more diverse and representative subset of the la-549 tent space. Compared to Kmeans (7.2%), DPPs selection 550 achieves 9.3% accuracy, validating its ability to reduce re-551 dundancy and improve feature coverage. Furthermore, our 552 full method, which integrates DPPs-based selection with 553 HOSVD, achieves the best overall performance at 10.5%, 554 surpassing both previous dataset distillation methods and 555 other alternative sampling techniques. The complete evalu-556 ation accuracies are detailed in the supplementary material. 557

These results highlight the importance of an effective data selection strategy in video dataset distillation. Our approach leverages DPPs to maximize diversity while retaining representative samples, leading to superior generalization in downstream tasks.

| | Evaluation Model | | | |
|-----------------|------------------|----------------|--------------|--|
| | ConvNet3D | CNN+GRU | CNN+LSTM | |
| Random | 9.9 ± 0.8 | 6.2 ± 0.8 | 6.5 ± 0.3 | |
| DM [45] | 15.3 ± 1.1 | 9.9 ± 0.7 | 9.2 ± 0.3 | |
| DM + VDSD [39] | 17.5 ± 0.1 | 12.0 ± 0.7 | 10.3 ± 0.2 | |
| MTT [3] | 19.0 ± 0.1 | 8.4 ± 0.5 | 7.3 ± 0.4 | |
| MTT + VDSD [39] | 23.3 ± 0.6 | 14.8 ± 0.1 | 13.4 ± 0.2 | |
| Ours | 34.8 ± 0.5 | 19.9 ± 0.7 | 18.3 ± 0.7 | |

Table 2. Result of experiment on cross-architecture generalization for MiniUCF when IPC is 1.

Rank Compression Ratio

We evaluate the impact of different rank compression ratios in HOSVD on overall performance in Tab. 3. Empirical results show that a rank compression ratio of r=0.75 consistently provides a strong balance between storage efficiency and model accuracy across datasets. While increasing the compression ratio reduces storage requirements, overly aggressive compression can lead to significant information loss, negatively affecting downstream tasks. Notably, as shown in Fig. 3, when the rank compression ratio is set to r = 0.1, both datasets exhibit classification accuracy around 4.0%, suggesting that excessive compression leads to degraded latent representations, making the distilled dataset nearly indistinguishable from random noise.

| Dataset | Rank Compression Ratio | | | | | |
|---------|------------------------|--------------|--------------|----------------|--------------|--|
| | 0.10 | 0.25 | 0.50 | 0.75 | 1.00 | |
| MiniUCF | 4.1 ± 0.1 | 19.0 ± 1.3 | 31.5 ± 0.7 | 34.8 ± 0.5 | 28.9 ± 0.5 | |
| HMDB51 | 3.9 ± 0.6 | 7.6 ± 1.0 | 11.5 ± 0.1 | 12.1 ± 0.3 | 8.9 ± 0.5 | |

Table 3. Accuracies under different rank compression ratios. Both MiniUCF and HMDB51 datasets are evaluated under IPC 1.

HOSVD vs Classic SVD To evaluate the effectiveness 577 of our latent-space compression strategy, we compare truncated SVD with HOSVD under the same storage budget at 579 IPC 5. Truncated SVD is a matrix factorization technique that approximates a data matrix by keeping only its largest 581

563 564

565

575 576

578



Figure 3. Accuracies of HMDB51 (IPC 1) and MiniUCF (IPC 1) under different rank compression ratios utilized in HOSVD.

singular values, thereby reducing dimensionality while retaining the most informative components. However, SVD
operates on flattened data matrices, leading to a loss of
structural information, particularly in spatiotemporal representations.

As shown in Tab. 4, HOSVD consistently outperforms 587 truncated SVD across all datasets, demonstrating its ability 588 to better preserve spatial and temporal dependencies in the 589 latent space. The performance gains are especially notable 590 on MiniUCF (+2.6%) and HMDB51 (+1.8%). Similarly, on 591 Kinetics-400 and SSv2, HOSVD achieves higher classifica-592 593 tion accuracy (+1.4% and +1.2%, respectively), highlighting its advantage in handling large-scale datasets. These 594 results confirm that HOSVD's tensor-based decomposition 595 596 provides a more compact yet expressive representation.

| Dataset | MiniUCF | HMDB51 | Kinetics-400 | SSv2 |
|---------|----------------|----------------|----------------|----------------|
| SVD | 38.5 ± 0.4 | 15.8 ± 0.2 | 12.4 ± 0.3 | 9.3 ± 0.2 |
| HOSVD | 41.1 ± 0.6 | 17.6 ± 0.4 | 13.8 ± 0.1 | 10.5 ± 0.4 |

Table 4. Classification accuracies comparison between different latent compression techniques under the same storage budget for each dataset at IPC 5.

597 4.6. Visualization

598 Following previous works, we provide an inter-frame contrast between DM and our method to illustrate the differ-599 ences in temporal consistency in Fig. 4. Specifically, we 600 sample three representative classes (CleanAndJerk, Playing 601 602 Violin, and Skiing) from the MiniUCF dataset and visualize the temporal evolution of distilled instances. The results 603 clearly demonstrate that our method retains more temporal 604 information, preserving smooth motion transitions across 605 frames. These visualizations further validate the effective-606 ness of our latent-space video distillation framework in pre-607 608 serving critical spatiotemporal dynamics.



(c) Skiing

Figure 4. Inter-frame comparison between DM and our method. Our frames are reconstructed from saved tensors and decoded by a 3D-VAE.

5. Conclusion

In this work, we introduce a novel latent-space video 610 dataset distillation framework that leverages VAE encod-611 ing, Diversity-Aware Data Selection, and High-Order Sin-612 gular Value Decomposition (HOSVD) to achieve state-of-613 the-art performance with efficient storage. By applying 614 training-free latent compression, our method preserves es-615 sential spatiotemporal dynamics while significantly reduc-616 ing redundancy. Extensive experiments demonstrate that 617 our approach outperforms prior pixel-space methods across 618 multiple datasets, achieving higher accuracy. We believe 619 our method provides an effective and scalable solution for 620 video dataset distillation, enabling improved efficiency in 621 training deep learning models. 622

Future Work Although we have achieved strong perfor-623 mance using selection-based, training-free methods, there 624 remains room for further improvement. In future work, 625 we aim to explore learning-based approaches in addition to 626 selection-based methods to further enhance dataset distilla-627 tion performance. By incorporating trainable mechanisms 628 for optimizing distilled representations, we expect to im-629 prove both efficiency and generalization. Additionally, we 630 plan to investigate non-linear decomposition techniques for 631 latent-space compression, moving beyond linear factoriza-632 tion methods such as HOSVD. Leveraging non-linear de-633 composition could lead to a more compact and expressive 634 latent space, enabling even greater storage efficiency while 635 preserving essential video dynamics. 636

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742 743

744

745

746

747

748

749

750

References 637

- 638 [1] stabilityai/sd-vae-ft-mse · Hugging Face — huggingface.co. 639 https://huggingface.co/stabilityai/sdvae-ft-mse. [Accessed 07-03-2025]. 1 640
- 641 [2] Joúo Carreira and Andrew Zisserman. Quo vadis, action 642 recognition? a new model and the kinetics dataset. In CVPR, 643 pages 4724-4733, 2017. 6
- 644 [3] George Cazenavette, Tongzhou Wang, Antonio Torralba, 645 Alexei A Efros, and Jun-Yan Zhu. Dataset distillation by 646 matching training trajectories. In CVPR, 2022. 2, 5, 6, 7
- 647 [4] Justin Cui, Ruochen Wang, Si Si, and Cho-Jui Hsieh. Dc-648 bench: Dataset condensation benchmark. arXiv preprint 649 arXiv:2207.09639, 2022. 1
- 650 [5] Justin Cui, Ruochen Wang, Si Si, and Cho-Jui Hsieh. Scaling up dataset distillation to imagenet-1k with constant memory. 652 In Proceedings of the International Conference on Machine 653 Learning (ICML), pages 6565-6590, 2023. 1
- 654 [6] Yuxuan Duan, Jianfu Zhang, and Liqing Zhang. Dataset distillation in latent space. arXiv preprint arXiv:2311.15547, 655 656 2023. 3
- [7] Lorenzo Ghilotti, Mario Beraha, and Alessandra Guglielmi. 657 658 Bayesian clustering of high-dimensional data via latent repulsive mixtures. arXiv preprint arXiv:2303.02438, 2023. 659 660 4
- 661 [8] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzyäska, Susanne Westphal, Heuna Kim, 662 663 Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz 664 Mueller-Freitag, Florian Hoppe, Christian Thurau, Ingo Bax, 665 and Roland Memisevic. The "something something" video 666 database for learning and evaluating visual common sense, 667 2017. 6
- 668 [9] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, 669 Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben 670 Poole, Mohammad Norouzi, David J Fleet, et al. Imagen 671 video: High definition video generation with diffusion mod-672 els. arXiv preprint arXiv:2210.02303, 2022. 2
- [10] Xing Hu, Yuan Cheng, Dawei Yang, Zhihang Yuan, Jiangy-673 674 ong Yu, Chen Xu, and Sifan Zhou. I-llm: Efficient integer-675 only inference for fully-quantized low-bit large language 676 models. arXiv preprint arXiv:2405.17849, 2024. 5
- 677 [11] De-An Huang, Vignesh Ramanathan, Dhruv Mahajan, Lorenzo Torresani, Manohar Paluri, Li Fei-Fei, and 678 679 Juan Carlos Niebles. What makes a video a video: Ana-680 lyzing temporal information in video understanding models 681 and datasets. In CVPR, pages 7366–7375, 2018. 2
- 682 [12] Abiodun M Ikotun, Absalom E Ezugwu, Laith Abualigah, Belal Abuhaija, and Jia Heming. K-means clustering algo-683 684 rithms: A comprehensive review, variants analysis, and ad-685 vances in the era of big data. Information Sciences, 622: 686 178-210, 2023. 4, 7
- 687 [13] Diederik P Kingma, Max Welling, et al. Auto-encoding vari-688 ational bayes, 2013. 2, 3, 4
- 689 [14] Diederik P Kingma, Max Welling, et al. An introduction to 690 variational autoencoders. Foundations and Trends® in Ma-691 chine Learning, 12(4):307-392, 2019. 3
- 692 [15] Dan Kondratyuk, Lijun Yu, Xiuye Gu, José Lezama, 693 Jonathan Huang, Rachel Hornung, Hartwig Adam, Hassan

Akbari, Yair Alon, Vighnesh Birodkar, et al. Videopoet: A large language model for zero-shot video generation. arXiv preprint arXiv:2312.14125, 2023. 2

- [16] H. Kuehne, H. Jhuang, E. Garrote, T. Poggio, and T. Serre. Hmdb: A large video database for human motion recognition. In ICCV, pages 2556-2563, 2011. 6
- [17] Alex Kulesza and Ben Taskar. Determinantal point processes for machine learning. Foundations and Trends in Machine Learning, 5(2-3):123-286, 2012. 1, 4
- [18] Carlotta Langer, Yasmin Kim Georgie, Ilja Porohovoj, Verena Vanessa Hafner, and Nihat Ay. Analyzing multimodal integration in the variational autoencoder from an informationtheoretic perspective. arXiv preprint arXiv:2411.00522, 2024. Accessed: 2024-11-01. 4
- [19] Longzhen Li, Guang Li, Ren Togo, Keisuke Maeda, Takahiro Ogawa, and Miki Haseyama. Generative Dataset Distillation: Balancing global structure and local details. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Workshop, pages 7664-7671, 2024. 2
- [20] Xin Liu, Silvia L. Pintea, Fatemeh Karimi Nejadasl, Olaf Booij, and Jan C. van Gemert. No frame left behind: Full video action recognition. In CVPR, pages 14892-14901, 2021. 2
- [21] Yanqing Liu, Jianyang Gu, Kai Wang, Zheng Zhu, Wei Jiang, and Yang You. Dream: Efficient dataset distillation by representative matching. arXiv preprint arXiv:2302.14416, 2023.
- [22] Noel Loo, Ramin Hasani, Alexander Amini, and Daniela Rus. Efficient dataset distillation using random feature approximation. arXiv preprint arXiv:2210.12067, 2022. 1
- [23] Xin Ma, Yaohui Wang, Gengyun Jia, Xinyuan Chen, Ziwei Liu, Yuan-Fang Li, Cunjian Chen, and Yu Qiao. Latte: Latent diffusion transformer for video generation. arXiv preprint arXiv:2401.03048, 2024. 2
- [24] Dionysis Manousakas, Zuheng Xu, Cecilia Mascolo, and Trevor Campbell. Bayesian pseudocoresets. In NeurIPS, 2020. 2
- [25] Brian B Moser, Federico Raue, Sebastian Palacio, Stanislav Frolov, and Andreas Dengel. Latent dataset distillation with diffusion models. arXiv preprint arXiv:2403.03881, 2024. 3
- [26] Elvis Nava, Mojmir Mutny, and Andreas Krause. Diversified sampling for batched bayesian optimization with determinantal point processes. In International Conference on Artificial Intelligence and Statistics, pages 7031-7054. PMLR, 2022.1
- [27] Timothy Nguyen, Zhourong Chen, and Jaehoon Lee. Dataset meta-learning from kernel ridge-regression. arXiv preprint arXiv:2011.00050, 2020. 1
- [28] Timothy Nguyen, Roman Novak, Lechao Xiao, and Jaehoon Lee. Dataset distillation with infinitely wide convolutional networks. In NeurIPS, 2021. 2
- [29] Rajesh Ranganath, Sean Gerrish, and David Blei. Black box variational inference. In AISTATS, 2014. 1
- [30] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. In ICLR, 2018. 2, 5, 6

787

788

789

790

- [31] Younggyo Seo, Kimin Lee, Fangchen Liu, Stephen James, and Pieter Abbeel. Harp: Autoregressive latent video prediction with high-fidelity image generator. In 2022 IEEE International Conference on Image Processing (ICIP), pages 3943–3947. IEEE, 2022. 2
- [32] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah.
 UCF101: A dataset of 101 human actions classes from videos in the wild. *CoRR*, abs/1212.0402, 2012. 6
- [33] Piyush Tiwary, Kumar Shubham, Vivek Kashyap, et al. Constructing bayesian pseudo-coresets using contrastive divergence. arXiv preprint arXiv:2303.11278, 2023. 2
- [34] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang.
 Videomae: Masked autoencoders are data-efficient learners
 for self-supervised video pre-training. In *NeurIPS*, pages
 10078–10093. Curran Associates, Inc., 2022. 1, 2
- [35] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In 2015 IEEE International Conference on Computer Vision (ICCV), pages 4489–4497, 2015. 6
- [36] Kai Wang, Bo Zhao, Xiangyu Peng, Zheng Zhu, Shuo Yang,
 Shuo Wang, Guan Huang, Hakan Bilen, Xinchao Wang, and
 Yang You. Cafe: Learning to condense dataset by aligning
 features. In *CVPR*, 2022. 1
- [37] Kai Wang, Jianyang Gu, Daquan Zhou, Zheng Zhu, Wei
 Jiang, and Yang You. Dim: Distilling dataset into generative model. *arXiv preprint arXiv:2303.04707*, 2023. 2
- [38] Tongzhou Wang, Jun-Yan Zhu, Antonio Torralba, and
 Alexei A Efros. Dataset distillation. *arXiv preprint arXiv:1811.10959*, 2018. 1, 2
- [39] Ziyu Wang, Yue Xu, Cewu Lu, and Yong-Lu Li. Dancing
 with still images: Video distillation via static-dynamic disentanglement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages
 6296–6304, 2024. 1, 2, 4, 5, 6, 7
 - [40] Max Welling. Herding dynamical weights to learn. In *ICML*, 2009. 1, 2, 5, 6
 - [41] Shuo Yang, Zeke Xie, Hanyu Peng, Minjing Xu, Mingming Sun, and P. Li. Dataset pruning: Reducing training data by examining generalization influence. *ArXiv*, abs/2205.09329, 2022. 2
- [42] Juyoung Yun, Sol Choi, Francois Rameau, Byungkon Kang,
 and Zhoulai Fu. Standalone 16-bit training: Missing study
 for hardware-limited deep learning practitioners. *arXiv preprint arXiv:2305.10947*, 2023. 5
- [43] David Junhao Zhang, Jay Zhangjie Wu, Jia-Wei Liu,
 Rui Zhao, Lingmin Ran, Yuchao Gu, Difei Gao, and
 Mike Zheng Shou. Show-1: Marrying pixel and latent diffusion models for text-to-video generation. *arXiv preprint arXiv:2309.15818*, 2023. 3
- [44] Bo Zhao and Hakan Bilen. Dataset condensation with differentiable siamese augmentation. In *ICML*, 2021. 1
- [45] Bo Zhao and Hakan Bilen. Dataset condensation with distribution matching. In *WACV*, 2023. 2, 5, 6, 7
- [46] Bo Zhao, Konda Reddy Mopuri, and Hakan Bilen. Dataset
 condensation with gradient matching. *arXiv preprint arXiv:2006.05929*, 2020. 2

- [47] Sijie Zhao, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Muyao Niu, Xiaoyu Li, Wenbo Hu, and Ying Shan. Cvvae: A compatible video vae for latent generative video models. Advances in Neural Information Processing Systems, 37: 12847–12871, 2025. 1, 2, 5, 6
 812
- [48] Yinjie Zhao, Heng Zhao, Bihan Wen, Yew-Soon Ong, and Joey Tianyi Zhou. Video set distillation: Information diversification and temporal densification. *arXiv preprint arXiv:2412.00111*, 2024. 1, 2, 5, 6
 816
- [49] Yongchao Zhou, Ehsan Nezhadarya, and Jimmy Ba. Dataset distillation using neural feature regression. *arXiv preprint arXiv:2206.00719*, 2022. 5, 6
 819

Latent Video Dataset Distillation

Supplementary Material

820 6. VAE

821 6.1. 2D-VAE Quantization

Variational Autoencoders (VAEs) enable significant data
compression by encoding each image as a probability distribution in a learned latent space, having the architecture
like in Fig. 5. The 2D-VAE used in this paper optimizes the
following loss function:

$$\mathcal{L}_{\text{VAE}} = \mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{x})} \left[\log p_{\theta}(\mathbf{x} \mid \mathbf{z}) \right] - D_{\text{KL}} \left(q_{\phi}(\mathbf{z} \mid \mathbf{x}) \parallel p_{\theta}(\mathbf{z}) \right)$$
(7)

The first term minimizes the reconstruction loss when decoding the latent representation of an image, while the second term, the KL divergence, ensures each encoded distribution aligns with a normal prior distribution. Combined, the objective balances the quality of decoded images and the smoothness of the latent distribution.

In order to ensure a fair comparison with previous work,
the weights of the VAE are quantized through post-training
static quantization, reducing the bid-width from 32 to 8 bits:

$$x_q = \operatorname{round}\left(\frac{x}{s}\right) + z \tag{8}$$

838 Where s is the scaling factor, and z is the zero point.

By applying linear quantization, the size of the pre-839 trained model is reduced to one-fourth of its original size. 840 841 Empirically, the quantized VAE continues to yield high accuracy during experimentation. Compared to other methods 842 such as quantization-aware training, static quantization has 843 the advantage of retaining a high level of accuracy while 844 845 offering lower computational complexity during the quantization phase. 846

847 7. Implementation Details

In this section, we provide implementation details of our
experiments, including the selection of VAEs, the preprocessing steps applied to video datasets, and the measures
taken to ensure a fair comparison.

852 7.1. Additional VAE Selection

We have adopted and quantized SD-VAE-FT-MSE[1] and CV-VAE[47] in our experiments. The variational autoencoders are used to encode video sequences into a compact latent space, enabling efficient dataset distillation. When dealing with IPC 1, where storage constraints are particularly strict, we employ SD-VAE-FT-MSE, a 2D-VAE, which compresses videos as independent frames, allowing



Figure 5. Architecture of Variational Autoencoder(VAE).

for highly compact storage. In contrast, for IPC 5, we uti-860 lize CV-VAE, a 3D-VAE, which explicitly models temporal 861 dependencies in video sequences. Unlike 2D-VAEs, which 862 treat frames as separate entities, 3D-VAEs capture motion 863 continuity and temporal redundancy, effectively reducing 864 redundant information across consecutive frames. This re-865 sults in a more structured latent representation, ensuring 866 that only the most informative motion features are retained, 867 leading to improved efficiency in video dataset distillation. 868 This selective choice of VAE architectures ensures that our 869 distilled datasets achieve the optimal balance between com-870 pression efficiency and information retention across differ-871 ent IPC levels. 872

7.2. Quantized VAE Model Size

We apply post-training static quantization on SD-VAE-FT-MSE, compressing the model from original 335MB to 80MB, achieving around 76% compression rate.

7.3. Fair Comparison

Throughout our experiments across four video datasets un-878 der two IPC settings (1 and 5), we rigorously ensure that the 879 storage used by our method does not exceed predefined stor-880 age constraints. For example, in MiniUCF IPC 1, previous 881 methods allocate a storage limit of 115MB. Under the same 882 setting, we sample 24 instances per class and apply HOSVD 883 with a compression rate of 0.75, saving the core tensor 884 and factor matrices. The resulting distilled dataset occu-885 pies 27MB, while the quantized 2D-VAE requires 80MB, 886 leading to a total memory consumption of 107MB, which 887 remains within the 115MB storage budget. The detailed 888 storage consumption can be found in Tab. 5. 889

7.4. Sampling Methods

In Tab. 6, we have provided a detailed accuracies on different sampling and dataset distillation techniques evaluating on the dataset SSv2 when IPC is 5.

890

873

874

875

876

| Dataset | MiniUCF | HMDB51 | Kinetics-400 | SSv2 |
|---------|---------|--------|--------------|--------|
| IPC 1 | 107 MB | 107 MB | 148 MB | 223 MB |
| IPC 5 | 475 MB | 475 MB | 455 MB | 458 MB |

Table 5. Storage consumed by our method for each dataset. Storage represents the total size of the distilled tensors and the associated VAE model.

| Random | DM + VDSD | MTT + VDSD | IDTD | Kmeans | DPPs only | Ours |
|-------------|---------------|---------------|-------------|---------------|---------------|--------------|
| 3.9 ± 0.1 | 4.0 ± 0.1 | 8.3 ± 0.1 | 9.5 ± 0.3 | 7.2 ± 0.3 | 9.3 ± 0.1 | 10.5 ± 0.2 |

Table 6. Performance of different dataset distilation and data sampling methods on the SSv2 dataset under IPC 1.

894 8. Peak Memory Analysis

To assess the efficiency of our method in terms of memory consumption, we compare the peak GPU memory usage during dataset distillation with other methods: DM and VDSD. As shown in Tab. 7, our method achieves the lowest peak memory consumption at 11,085 MiB, significantly reducing memory usage compared to DM (20,457 MiB) and VDSD (12,545 MiB).

| Method | DM | VDSD | Ours |
|------------|------------|------------|------------|
| GPU Memory | 20,457 MiB | 12,545 MiB | 11,085 MiB |

Table 7. Peak memory comparison between different dataset distillation methods on MiniUCF when IPC is 5.

Our method minimizes peak memory usage by operating in the latent space and leveraging training-free compression via HOSVD, significantly reducing redundant memory
allocation during dataset distillation. This lower memory
footprint allows our approach to scale to larger datasets and
higher IPC settings while maintaining efficiency.

908 9. Runtime Analysis

To assess the computational efficiency of our method, we
compare its distillation runtime with VDSD across different datasets. All experiments are conducted on an NVIDIA
H100 SXM GPU. Our training-free method demonstrates
a significant speed advantage, particularly on large-scale
datasets, due to its latent-space processing and training-free
compression strategy.

On small-scale datasets, such as HMDB51 and MiniUCF, our method completes the dataset distillation process
in under 10 minutes, whereas VDSD requires 2.5 hours.
The efficiency gain is even more pronounced on large-scale
datasets, where our method finishes in approximately 1 hour
on Kinetics-400 and SSv2, while VDSD exceeds 5 hours.

922 These results confirm that our latent-space approach
923 significantly reduces computational overhead compared to
924 pixel-space distillation methods like VDSD. By leverag925 ing structured compression techniques such as HOSVD and

eliminating costly iterative optimization steps, our method
achieves faster dataset distillation without compromising
performance. This makes our approach highly scalable and
practical for real-world applications, especially in large-
scale video analysis scenarios.926
927
928

10. Visualization

We provide the reconstructed and decoded frames of our method for MiniUCF across 20 classes in Fig. 6. 933 CVPR 2025 Submission #11. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.



Figure 6. Reconstructed and decoded frames of our method for MiniUCF with a 3D-VAE.