
SMPLy Private: From Masks to Meshes in Action Recognition

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this paper, we introduce MASK2MESH (M2M), a novel privacy-preserving
2 data augmentation framework that effectively bridges the realism gap commonly
3 associated with synthetic-based action recognition methods. Traditional privacy-
4 enhancing techniques, such as feature masking and synthetic data supplementation,
5 typically degrade data quality, thereby diminishing the performance of action
6 recognition models. In contrast, our approach leverages the SMPL-X model to
7 substitute real humans with 3D meshes within video data, preserving the intricate
8 details of human movements and expressions essential for accurate action recog-
9 nition. Our framework leverages a single dataset, augmented with superimposed
10 meshes, to streamline pre-training and fine-tuning. Unlike recent methods that
11 rely on synthetic data, our approach augments real data, eliminating associated
12 overheads and potential biases. Empirical results demonstrate that our approach
13 achieves performance within 0.5% of models trained on real video data, indicating
14 a promising direction for privacy-conscious, scalably efficient, and high-fidelity
15 video data processing.

16 1 Introduction

17 Action recognition, the process of classifying human activities based on video sequences, is crucial
18 for applications such as surveillance, human-computer interaction, and video analytics [27]. Tradi-
19 tional action recognition systems rely heavily on extensively annotated datasets to achieve optimal
20 performance. With advancements in deep learning and the emergence of vision transformers (ViTs)
21 [16], pre-training models on large datasets has become standard practice to enhance accuracy and
22 generalization [48]. However, these datasets often include identifiable individuals, raising significant
23 privacy and ethical concerns [61].

24 Data sharing, particularly without obtaining explicit consent from individuals, necessitates robust
25 de-identification methods. Conventional anonymization techniques, such as blurring and pixelation,
26 often degrade data quality, thereby reducing its efficacy for action recognition tasks. Moreover,
27 these methods rely on heuristics and may not effectively balance privacy protection with data utility
28 [51]. Ensuring individual privacy while sharing video data can significantly advance research and
29 development in action recognition and computer vision, where large datasets are imperative. Privacy-
30 preserving techniques that maintain data quality can enhance the accuracy and reliability of machine
31 learning models, facilitating more robust and fair applications. Nonetheless, these methods do not
32 fully address the realism gap introduced by synthetic methods, nor do they completely safeguard
33 privacy, as variable visuals like skin tone and gender can still be discerned [65, 15, 35].

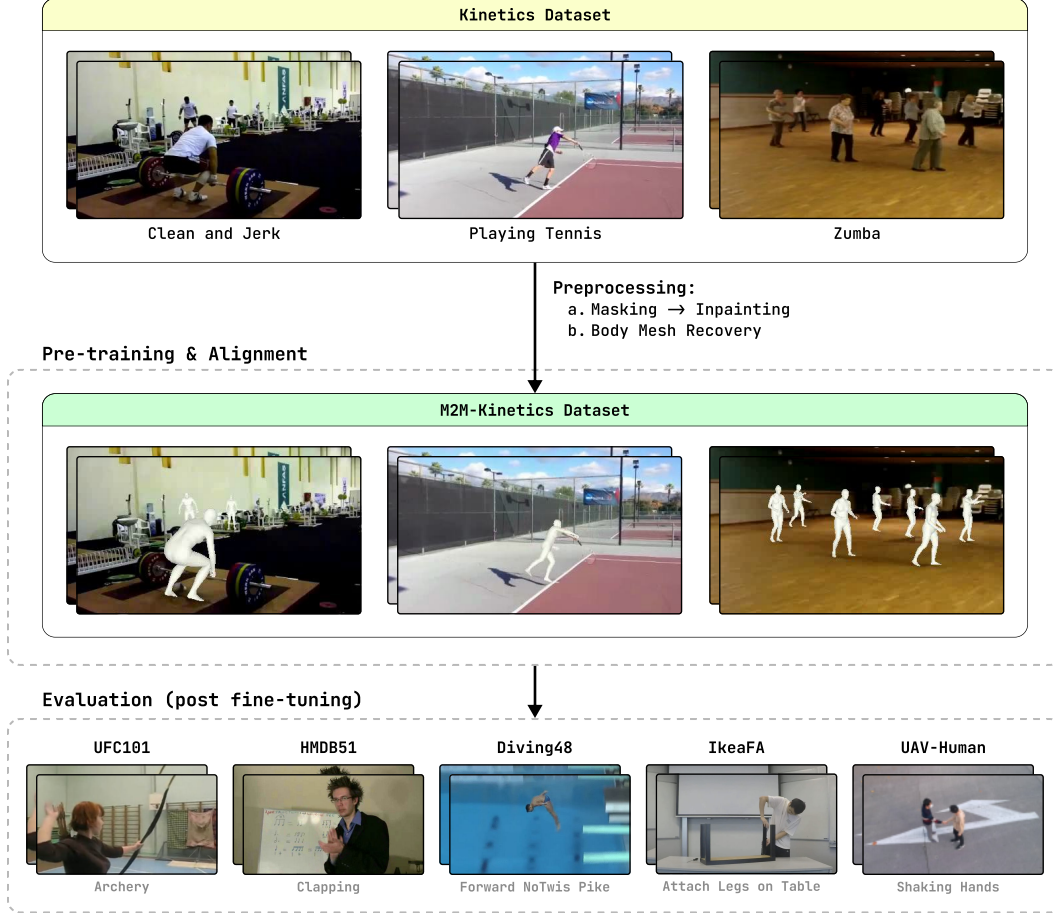


Figure 1: **SMPLY Private in Action**. Transforming human actions into privacy-preserving 3D meshes: videos from the Kinetics dataset are preprocessed using masking, inpainting, and body mesh recovery to replace humans with 3D mesh bodies. The M2M-augmented videos are then used for pre-training and alignment, with their performance evaluated across various human action recognition tasks, demonstrating the framework’s ability to maintain high fidelity and ensure ethical data usage.

Our research aims to address these challenges by demonstrating that using mesh bodies can effectively remove biases and close the realism gap in action video recognition while preserving privacy. In particular, we: (1) investigate whether meshes can preserve privacy while approximating real-world data and reducing the realism gap; (2) explore the potential of using meshes to mitigate biases related to background, scene-object interactions, race, and gender, examining whether they can serve as a standardization technique for the human form.

Subsequently, this approach is motivated by the premise that mesh bodies can replace humans in video data, maintaining privacy. However, to illustrate our interest in exploring the impact of mesh bodies on mitigating gender and race biases, consider an example. Video footage from the National Basketball Association (NBA) in the 1980s-1990s, predominantly features African-American players. Training models on these videos to predict contemporary players, who are more diverse, could introduce bias by predominantly associating basketball with African-American individuals. This concern extends to gender bias as well [4]. Employing meshes could standardize player representation across different demographics, thereby reducing biases and promoting ethical model development.

Our contributions are summarized as follows:

1. We introduce **MASK2MESH (M2M)**, a framework that replaces real humans in video datasets with detailed 3D meshes using the **SMPL-X** [47] model, effectively preserving

privacy without compromising the integrity of action recognition. Our method outperforms existing privacy-preserving benchmarks and rivals non-privacy-preserving methods, closing the realism gap.

2. By strategically modifying the k-means clustering algorithm [42], we introduce K-NEXUS, a dataset sampling strategy designed to eliminate action class bias. This enhancement significantly improves the performance of M2M.
3. Our targeted study on gender representation in 3D meshes reveals that gender-neutral meshes improve average performance in action recognition tasks within gender-biased classes, suggesting that neutral representations can effectively mitigate gender biases while maintaining privacy and ensuring fairness in automated video analysis.
4. We demonstrate that models fine-tuned on M2M-augmented data learn representations quicker due to consistent mesh depictions, offering a new perspective on bridging the realism gap and serving as a standardization technique for ethical model development across diverse demographic groups.

2 Related Works

Preservation of Privacy in Action Recognition. Anonymization techniques, such as blurring, downsampling, adversarial augmentations, masking faces and other identifiable features, are commonly adapted strategies to maintain human confidentiality [14, 5, 49, 64, 15, 58, 50] but can be left susceptible to revealing an individual’s identity based on characteristics such as color and size [60, 46]. We eliminate the human form and replace it with 3D meshes, thus curating training data that does not have explicitly identifiable forms or features of the original human.

Body Mesh Recovery. Estimating 3D human body poses and shapes is a complex challenge addressed by various methods. SMPL-X uses a unified 3D model trained on extensive 3D scans, providing detailed and realistic representations. HybriK-X [34] effectively combines 3D keypoint estimation with body mesh reconstruction by converting precise 3D joint locations into relative body-part rotations. However, its application is limited to single individuals, which is not suitable for videos featuring multiple people. In contrast, our approach utilizes OSX [39], which excels in multi-human mesh recovery. OSX employs a unified encoder-decoder architecture integrated with a component-aware transformer. This setup not only predicts body parameters but also enhances segmentation, crucial for accurate face and hand estimation. By eliminating the need for separate networks and manual post-processing, OSX provides more natural and plausible 3D meshes. Given its simplicity and effectiveness in handling complex scenes with multiple individuals, OSX is our chosen method for accurate human body mesh recovery for M2M.

Biases and Synthetic Data. Object-scene bias in video action recognition refers to the tendency of models to rely on static objects or backgrounds rather than the dynamic actions themselves for classification [62]. To address this issue, various augmentation strategies such as loss augmentation [12], action-scene swapping [66], and video compositing [20] are prevalent but not privacy-preserving. One method mitigates this bias by first learning background information from real data and then temporal information from entirely synthetic data rather than augmenting components of the actual video [65]. Although this approach is privacy-preserving and focuses on learning background and actions, thereby addressing object-scene bias, the use of synthetic videos still leaves a desire for the realism gap to be bridged [19]. Instead, our approach returns to augmenting real videos by masking, in-painting, and overlaying appropriate mesh bodies in place of the original human. The plain “mannequin-like” mesh bodies remove any discriminatory bias, unlike the synthetic “video-game-like” humans, which tend to have features such as hair color, gender, skin tone, etc.

Self-Supervised Pretraining in Action Recognition. The training scheme for action recognition models is crucial to the performance on downstream tasks, given that most data in nature is unlabeled. Self-supervised learning (SSL) has proven to be a powerful pretraining mechanism in such schemes [1]. Furthermore, the default choice of encoder has shifted from convolutional neural networks

(CNNs) like temporal segment networks (TSNs) [56] and inflated 3D convolutional networks (I3Ds) [8] to vision transformers (ViTs) [16], as they effectively process frames as patch sequences, enabling the capture of long-range dependencies and patterns in videos. SSL is typically categorized into four paradigms: deep metric learning [10, 17, 32, 30], self-distillation [21, 11, 7], canonical correlation analysis [2, 63, 6, 18], and masked image modeling [24, 59, 9]. We chose the latter paradigm, employing a masked autoencoder (MAE) training scheme, specifically using VideoMAE [57], which typically incorporates the base vision transformer architecture (ViT-B) and has previously achieved state-of-the-art performance on benchmarks like UCF101 [53] and Kinetics [29].

3 Methodology

3.1 Dataset Curation: K-NEXUS

We use a subset of the Kinetics-400 [29] video dataset as our training dataset where we select 150 classes amongst the 400 along with at most 1,000 videos per class. Previous works [31, 65] have used random splits to curate their custom dataset, however, we consider the action class bias in the Kinetics dataset (e.g., actions like playing violin and playing guitar are visually closer than playing volleyball). To obtain discrete classes and reduce the bias, we uniquely deploy a k-means clustering algorithm [42] to obtain the final set of classes. Our approach aims to assemble a K -class dataset D^* of minimal bias from an existing dataset D with C classes. The class labels of D are denoted by $L = \{l_1, l_2, \dots, l_C\}$, where l_i represents the i -th class of D . Our objective is to find a subset $L^* = \{l_1^*, \dots, l_K^*\}$ such that: $l_i^* \in L$, $l_i^* \neq l_j^*$ for $i \neq j$, and D^* has minimal bias. Specifically, we perform the following steps for label sampling:

1. **Encoding action image-label pairs.** Let V be the set of all videos and L , previously defined, be the set of all class labels. $V_{l_j} \subseteq V$ is then the set of all videos with class label $l_j \in L$. For each video $v_i \in V_{l_j}$, we sample a random frame I_{v_i} . We then construct a set of tuples $\Theta = \{\theta_{v_i, l_j} \mid v_i \in V, l_j \in L\}$ where $\theta_{v_i, l_j} = (I_{v_i}, l_j)$. An embedding function is then defined as $f: \theta \rightarrow \mathbb{R}^d$, which maps a tuple θ to a d -dimensional embedding space using a LLaVA image encoder [41]. The embedding of a tuple θ_{v_i, l_j} is given by: $\mathbf{e}_{v_i, l_j} = f(\theta_{v_i, l_j})$. To compute the average embedding for class l_j , we aggregate the embeddings for the frame-label pairs from all videos in V_{l_j} : $E_{l_j} = \{\mathbf{e}_{v_i, l_j} \mid v_i \in V_{l_j}\}$. Then, the average embedding $\bar{\mathbf{e}}_{l_j}$ for class l_j is given by: $\bar{\mathbf{e}}_{l_j} = \frac{1}{|V_{l_j}|} \sum_{\mathbf{e} \in E_{l_j}} \mathbf{e}$. The set of average embeddings for all class labels is then given by: $\bar{E} = \{\bar{\mathbf{e}}_{l_j} \mid l_j \in L\}$.
2. **K-means clustering.** We then apply a modified k-means clustering algorithm (see Appendix B.1), that minimizes dataset bias instead of the within-cluster sum of squares, to partition the embeddings in \bar{E} into K clusters. Let $\kappa(\bar{E}, K)$ denote the modified clustering operation, resulting in cluster assignments $\Omega = \{\omega_1, \omega_2, \dots, \omega_C\}$, where each $\omega_i \in \{1, 2, \dots, K\}$.
3. **Selecting representative labels.** For each cluster $k \in \{1, 2, \dots, K\}$, we identify the labels that belong to cluster k . Let $L_k = \{l_i \mid \omega_i = k\}$ be the set of labels in cluster k . From each cluster k , we select the representative label l_k^* with the highest number of video samples within the cluster: $l_k^* = \arg \max_{l_i \in L_k} |V_{l_i}|$ where V_{l_i} is the set of all videos with class label l_i and $|V_{l_i}|$ denotes the number of videos in V_{l_i} . The final set of labels L^* is then given by: $L^* = \{l_k^* \mid k \in \{1, 2, \dots, K\}\}$. Thereby obtaining D^* (Kinetics-150).

Our method, termed k-Means Neural Embeddings eXploited for Unified Sampling (K-NEXUS), groups similar actions together and selects representative labels to reduce action class bias in Kinetics¹. By doing so, K-NEXUS ensures a more balanced set of labels and minimizes representation bias across different action categories.

¹We can compute the relative entropy [22] between adjacent frames and select a frame from a subset with the least change in entropy. However, Kinetics consists of short videos showcasing a single action class, resulting in minimal entropy change between adjacent frames. Hence, selecting a random frame per video is justified.

3.2 MASK2MESH Augmentation

Our proposed M2M augmentation framework is designed to achieve privacy-preserving video data by replacing real humans in videos with 3D mesh representations while preserving essential motion details. The framework consists of two main modules: (1) the masking and inpainting module; (2) the body mesh recovery module. We leverage the Kinetics-150 dataset curated in the aforementioned Section 3.1, resizing all video clips to 432×240 to standardize the input data. The first step in our framework involves detecting and removing human figures from the video frames. This process is divided into two sub-tasks: human detection and in-painting. We utilize Mask R-CNN [25] with a ResNet-101 [26] backbone, based on COCO [40] instance segmentation weights, to generate masks for human figures in each video frame. The generated masks are passed on to the subsequent inpainting module, which involves filling the regions occupied by human figures with plausible background content. We employ E²FGVI [38], an optical flow-based inpainting method. E²FGVI leverages temporal coherence and spatial context to generate high-quality inpainted frames, effectively removing humans while maintaining the integrity of the background.

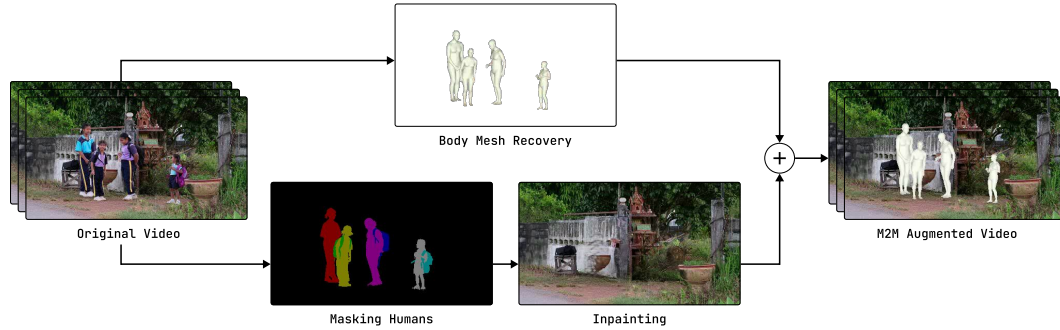


Figure 2: **MASK2MESH (M2M) Framework.** A visual representation of transforming real human actions into 3D mesh models across various activities for enhanced privacy and bias mitigation in action recognition. The figure details the flow from the initial video data of the Kinetics dataset, through masking, inpainting, and body mesh recovery to the final stage of mesh superimposition.

Once the humans are removed from the video frames, we focus on reconstructing the human motion using 3D mesh models. The body mesh recovery module processes the original (resized) videos to extract detailed 3D mesh models of the human figures. We utilize the SMPL-X model for this purpose and the OSX algorithm for mesh recovery. The SMPL-X model is chosen for its comprehensive representation of the human body, hands, and facial expressions, making it suitable for capturing intricate action details. Similarly, OSX is utilized due to its support for multi-human mesh recovery. The 3D meshes are then superimposed onto the inpainted video frames, replacing the real humans with their corresponding 3D mesh models. This step ensures that the augmented videos retain the essential motion cues required for accurate action recognition, while effectively anonymizing the individuals. A summarized overview of the M2M augmentation is provided in Figure 2.

3.3 Training Procedure

Pre-training is vital in video action recognition, enabling models to learn generalized spatial and temporal features from large, diverse datasets. This foundation enhances the model’s ability to recognize complex actions with less labeled data and improves performance and efficiency in downstream tasks [52]. Our training procedure consists of two stages: (1) self-supervised pre-training utilizing VideoMAE; (2) supervised pre-training to ensure label alignment.

Step 1: Video Masked Autoencoder for Self-Supervised Pre-training. We employ the traditional training methodology for MAEs [24] tailored explicitly for video data. This entails a configuration consisting of an encoder and a decoder, in which the model acquires the ability to approximate masked pixel values within video frames. During this stage, the encoder and decoder undergo joint training. Once the model is sufficiently trained, the decoder is removed, leaving only the encoder.

Step 2: Supervised Pre-training for Alignment of Labels. The pre-trained VideoMAE encoder is then augmented with a linear classification head. The encoder and the linear classifier are trained in tandem, using the action labels for supervision.

3.4 Downstream Evaluation

The evaluation of our SMPLy Private models is conducted on six distinct action-recognition tasks. The UCF101 dataset [53] comprises 13,320 YouTube videos spanning 101 action classes, showcasing notable diversity in activities performed and camera movement. The HMDB51 [33] dataset consists of 6,849 movie clips, categorized into 51 distinct action classes. The Diving48 dataset [37] is a highly specialized data collection designed for competition diving. It consists of 18,000 video clips that cover 48 different categories. This dataset aims to evaluate our models’ ability to handle the challenges posed by the identical background and object properties commonly found in competitive diving scenarios. Ikea Furniture [23, 54] Assembly provides a collection of 111 movies with 14 actors demonstrating assembling and disassembling furniture. These videos are filmed using the same camera and scenario settings, ensuring consistency. The videos are categorized into 12 different action categories. The UAV-Human dataset [36] comprises 22,476 films recorded by unmanned aerial vehicles, such as drones, featuring 155 distinct action categories and 119 individuals. Note that from UCF101 to UAV-Human, the scene-object bias generally decreases.

4 Experiments

For detailed technical specifications and additional training information related to our experimental setup, please refer to Appendix A.

4.1 MASK2MESH Performance

Pre-trained Model	Privacy	Step 1: MAE	Step 2: Align	UCF101		HMDB51		Diving48		IkeaFA		UAV-Human		Mean	
				FT	LP	FT	LP	FT	LP	FT	LP	FT	LP	FT	LP
VideoMAE-Align w/ Real	✗	Kinetics	Kinetics	93.4	91.5	73.5	69.8	66.3	19.9	72.2	58.4	34.8	13.8	68.0	50.7
TimeSformer w/ Kinetics	✗	-	Kinetics	92.1	89.4	59.5	55.4	46.4	17	61.9	47.7	23.3	8.4	56.6	43.6
TimeSformer w/ Synthetic	✗	-	Synthetic	89	82.1	54.4	49.2	44.9	19.2	63.6	45.5	25	13.8	55.4	42.0
TSN w/ RN50	✓	-	Synthetic	83.4	28	54.4	20.9	63.5	10.9	42.7	36	35.6	5.7	55.9	20.3
I3D w/ RN50	✓	-	Synthetic	82.1	27.6	55.7	22.6	55.0	10.1	42.7	33.2	35.1	5.8	54.2	19.9
OmniMAE-Align w/ Synthetic	✓	Synthetic	Synthetic	80	26.4	53.3	22.2	57.3	10	41.5	35.7	31.8	5.5	52.8	20.0
PPMA	✓	No Human Kinetics	No Human Kinetics + Synthetic	92.5	88.4	71.2	64.9	64.0	21.9	67.9	57.7	38.5	19.3	66.8	50.4
SMPLy Priv. (Ours)	✓	M2M Kinetics	M2M Kinetics	93.2	90.9	72.6	69.2	66.0	19.7	71.3	58.2	34.6	14.3	67.5	50.5
SMPLy Priv. w/ K-NEXUS (Ours)	✓	M2M Kinetics	M2M Kinetics	94.2	91.6	74.3	70.8	69.0	21.6	72.9	59.5	36.4	15.3	69.6	51.8

Table 1: M2M Performance Evaluation. Top-1 downstream task accuracy for linear probing (LP) and finetuning (FT) is reported. The mean FT and LP accuracy over all downstream tasks across datasets is represented in the final column. The results show that our SMPLy Private model outperforms, on average, prior benchmarks by at least 0.7% in FT and 0.1% in LP. If K-NEXUS is used (in teal), the improvement increases to at least 2.8% in FT and 1.4% in LP due to enhanced feature learning from class discretization. SMPLy Private rivals the VideoMAE trained with real human data (baseline in violet), reducing the realism gap by 0.5% in FT and 0.2% in LP. With K-NEXUS, SMPLy Private surpasses this baseline. All other scores are from [31, 65]. Our choice of alignment was selected based on results from Table 4.

Table 1 shows the average downstream performance on various classification tasks with multiple pre-trained models, including our proposed SMPLy Private and SMPLy Private with K-NEXUS. Models pre-trained on conventional large-scale real video data with humans typically have a performance edge over models trained with synthetic data due to a realism gap. We establish such a baseline by pre-training and then aligning VideoMAE with Kinetics-150 (first row in violet). Other privacy-preserving baselines present a significant realism gap (non-bolded and non-colored). However, with

SMPLY Private and the use of our M2M-augmented dataset (M2M Kinetics), the average downstream performance gap from the human baseline is reduced to 0.5% with FT and 0.2% with LP.

The performance gap is attributed to SMPLY Private performing slightly worse than the human baseline on tasks with high scene-object bias, such as UCF101 and HMDB51. This is likely because the inclusion of humans in the Kinetics videos helps the model better learn both scene-object cues and action features. Compared to the performance of “OmniMAE-Align w/ Synthetic,” SMPLY Private narrows the realism gap as it achieves a level of performance comparable to the “VideoMAE-Align w/ Real” baseline. Both “TSN with RN50” and “PPMA” utilize synthetic data for model training, yet they fall short of “VideoMAE-Align with Real” in downstream performance. With K-NEXUS, SMPLY Private further improves over “VideoMAE-Align w/ Synthetic” by 1.6% with FT and 1.1% with LP. Overall, SMPLY Private with K-NEXUS, which uses M2M Kinetics in both pre-training and alignment steps, achieves the best performance among privacy-preserving models, reducing the performance gap with the human-baseline model to minimal levels. This shows the effectiveness of our approach in achieving high performance through pre-trained representations for privacy-preserving action recognition without the need for synthetic data.

4.2 The Inpainting Influence

Mesh Recovery	Privacy	Inpainting	Downstream Accuracy (LP Only)					
			UCF101	HMDB51	Diving48	IkeaFA	UAV-Human	Mean
OSX	✗	None	91.1	69.4	20.0	58.3	14.6	50.7
	✓	E ² FGVI	90.9	69.2	19.7	58.2	14.3	50.5

Table 2: **Inpainting Effect Analysis.** We report the downstream task accuracy for OSX mesh recovery and inpainting. The downstream performance remains comparable, while the inpainting step ensures privacy.

We incorporate the E²FGVI inpainting technique to remove humans from video streams². While direct application of mesh recovery to videos is feasible, integrating an inpainting step markedly improves the privacy-preserving capabilities of our framework. Absent this inpainting process, residual demographic information can still be discerned, compromising both the privacy and the unbiased nature. Although omitting inpainting yields a higher performance (as demonstrated in Table 2), applying inpainting ensures that our pipeline remains fully privacy-preserving. We hypothesize that the observed higher performance may stem from the retention of features, which, although insufficiently masked, provide additional discriminative features that aid the learning process of the model. This highlights a trade-off between performance and privacy.

4.3 To Gender or Not To Gender?

Method	Women in Men-Biased (FT)	Men in Women-Biased (FT)	Mean
VideoMAE w/ real humans	81.4	78.8	80.1
SMPLY Private w/ male meshes	82.3	83.3	82.8
SMPLY Private w/ female meshes	83.4	82.5	82.9
SMPLY Private w/ neutral meshes	83.2	83.1	83.1

Table 3: **Gender-Action Bias Analysis.** We show the performance on gender-biased tasks using real human data vs. 3D meshes. Neutral meshes achieve the highest average accuracy, demonstrating effective mitigation of gender-action bias.

²In a resource-constrained setting, OpenCV’s implementation of Navier-Stokes inpainting [3, 28] can be used. We speculate that since inpainting is an intermediate step, any performance loss would be minimal.

In this section, we analyze the impact of using 3D meshes on gender-action bias in action recognition tasks. Specifically, we compare the performance of a model trained on real human data (Kinetics) with those trained on our augmented meshed data (M2M Kinetics). We conduct experiments on a specifically curated split of the Kinetics dataset in which women perform male-dominated tasks and men perform female-dominated tasks. The results are summarized in Table 3. Training on real human data revealed significant gender-action bias, with lower performance on “men in women-biased” tasks compared to “women in men-biased” tasks. In contrast, models trained on 3D meshes showed improved performance. Male meshes increased accuracy for “men in women-biased” tasks, while female meshes for “women in men-biased” tasks. Neutral meshes performed consistently well across both subsets. Overall, using 3D meshes outperformed the real human data approach, with higher average scores across all mesh-based methods. This indicates that 3D meshes help mitigate gender-action bias by offering a gender-agnostic representation.

4.4 The Tortoise & The Hare

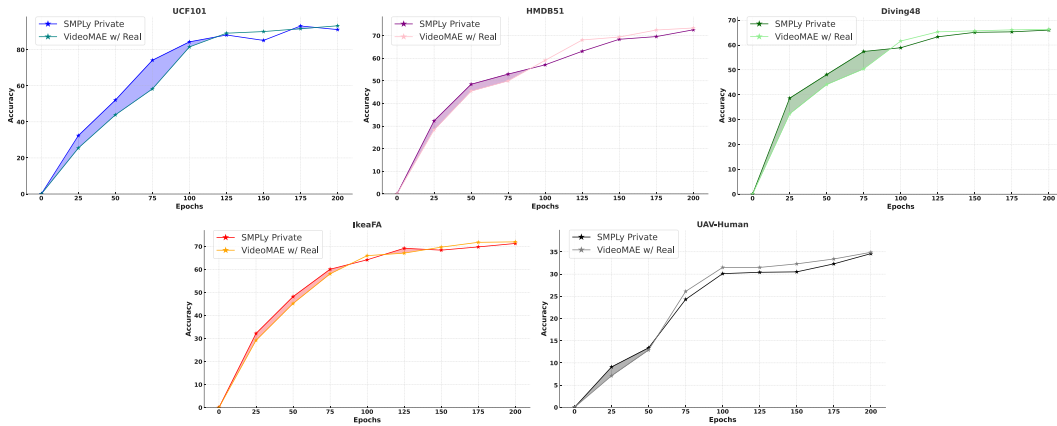


Figure 3: **Representation Learning Efficiency Comparison.** In the initial stages of training, the model trained with M2M Kinetics (teal in Table 1) demonstrates faster representation learning compared to the model trained on Kinetics (violet in Table 1), which catches up in the later stages.

As seen in Figure 3, we demonstrate that our model, fine-tuned from each pre-training checkpoint on M2M-augmented data, learns representations quicker in earlier stages because humans are consistently depicted as meshes. This consistency allows the model to isolate and understand the actions performed by the meshes, whereas, in real data, the varied depictions of humans make representation learning more challenging. It is observed that our model learns representations faster for longer with higher scene-object bias as indicated by the descending size of the shaded areas from UCF101 to UAV-Human (left to right). This showcases another perspective on further closing the realism gap; however, the “tortoise” (VideoMAE trained on real human data) eventually catches and overtakes the “hare” (SMPLy Private), typically in the latter stages of training.

4.5 In the Alignment Arena: Human vs. Mesh

In this experiment, we examine the effects of various alignment techniques following pre-training with the M2M Kinetics dataset. Our findings (see Table 4) indicate that when the dataset is exposed only to backgrounds (the no-human scenario, involving only inpainting without mesh recovery), performance is significantly lower compared to when actual humans are shown to the model. This approach notably falls short of the Kinetics baseline (in violet) when the model is exposed to real humans during the alignment phase. However, when our pretrained model is exposed to the M2M Kinetics dataset (with both inpainting and mesh recovery), we closely approximate the performance when compared to the Kinetics baseline ($|\Delta|$ is 0.3% and 0.2% for FT and LP respectively). This demonstrates that the mesh recovery technique in videos is effective in understanding real human actions, thus reinforcing our claim of bridging the realism gap.

Step 1: MAE	Step 2: Align	Privacy	UCF101		HMDB51		Diving48		IkeaFA		UAV-Human		Mean		Δ from "real"	
			FT	LP	FT	LP	FT	LP	FT	LP	FT	LP	FT	LP	FT	LP
M2M Kinetics	No Human Kinetics	✓	91.5	87.3	68.5	61.2	63.4	18.4	69.9	51.5	32.3	12.2	65.1	46.1	2.7	4.6
	Kinetics	✗	93.7	92.1	73.1	69.4	66.2	19.7	71.5	58.4	34.3	13.8	67.8	50.7	0.0	0.0
	M2M Kinetics	✓	93.2	90.9	72.6	69.2	66	19.7	71.3	58.2	34.6	14.3	67.5	50.5	0.3	0.2

Table 4: **Alignment Study for Humans and Meshes.** We present a comparative analysis of various alignment strategies following the pre-training of our model on the M2M Kinetics dataset. The alignment datasets under examination include the original Kinetics, Kinetics with humans removed, and M2M Kinetics. Our findings indicate that the M2M Kinetics dataset exhibits the smallest absolute difference, $|\Delta|$, from the Kinetics baseline (in violet), thereby further reducing the realism gap.

5 Conclusion

This study introduces M2M, a framework that replaces real humans in videos with detailed 3D meshes using the SMPL-X model, preserving privacy without compromising action recognition integrity. M2M avoids synthetic data biases by augmenting real data [31]. Experiments show models pre-trained and aligned on M2M-augmented data outperform existing privacy-preserving benchmarks and rival non-privacy-preserving ones. A targeted study on gender representation reveals that gender-neutral meshes produce better average performance in biased classes. Finally, we demonstrate that models trained on M2M-augmented data learn representations quicker, offering a new perspective on bridging the realism gap and standardizing action recognition data.

Ethical Implications. M2M has significant social and ethical implications, especially in privacy and data security. By replacing real individuals in video footage with 3D meshes, M2M addresses privacy concerns and aligns with regulations like GDPR [55] and ADPPA [13]. It mitigates risks of identity theft and privacy invasion, while reducing biases related to race and gender. However, its capability to generate highly realistic video data could also be repurposed for more invasive surveillance systems, potentially enhancing monitoring capabilities in workplaces or public areas and infringing on individual privacy and autonomy.

Limitations. The state-of-the-art mesh recovery methods responsible for overlaying the SMPL-X body on the human show degraded performance when the body joints of the person are occluded. This introduces a few incorrect representations as part of our dataset (see Appendix B.2). Since mesh recovery forms an integral part of our pipeline, our method suffers from the same plight.

Future Work. In our study, we perform image-instance segmentation without considering temporal relationships between frames. However, given that our data is video-based, an optimal approach would involve video-instance segmentation. This would likely yield more accurate masks and improved performance, mitigating issues such as transient mesh disappearances. Lastly, our present method focuses on body mesh recovery for individual video frames, whereas this process could be extended to entire videos to incorporate temporal relations. Most contemporary approaches to body mesh recovery leverage 3D pose estimation from videos, a complex problem [45]. Accurately recovering 3D meshes in videos is challenging but could significantly reduce glitches and enhance performance. We can better capture and use temporal information by replacing frame-by-frame processing with video-level analysis.

References

- [1] Randall Balestriero, Mark Ibrahim, Vlad Sobal, Ari Morcos, Shashank Shekhar, Tom Goldstein, Florian Bordes, Adrien Bardes, Gregoire Mialon, Yuandong Tian, Avi Schwarzschild, Andrew Gordon Wilson, Jonas Geiping, Quentin Garrido, Pierre Fernandez, Amir Bar, Hamed Pirsiavash, Yann LeCun, and Micah Goldblum. A cookbook of self-supervised learning, 2023.
- [2] Adrien Bardes, Jean Ponce, and Yann LeCun. Vicreg: Variance-invariance-covariance regularization for self-supervised learning. *arXiv preprint arXiv:2105.04906*, 2021.

- [3] M. Bertalmio, A.L. Bertozzi, and G. Sapiro. Navier-stokes, fluid dynamics, and image and video inpainting. In *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. CVPR 2001*, volume 1, pages I–I, 2001.
- [4] Kaylee Burns, Lisa Anne Hendricks, Kate Saenko, Trevor Darrell, and Anna Rohrbach. Women also snowboard: Overcoming bias in captioning models, 2019.
- [5] Daniel J Butler, Justin Huang, Franziska Roesner, and Maya Cakmak. The privacy-utility tradeoff for remotely teleoperated robots. In *Proceedings of the tenth annual ACM/IEEE international conference on human-robot interaction*, pages 27–34, 2015.
- [6] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [7] Mathilde Caron, Hugo Touvron, Ishan Misra, Herve Jegou, Piotr Bojanowski, Armand Joulin, Matthijs Douze, Matthieu Cord, and Ivan Laptev. Emerging properties in self-supervised vision transformers. *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.
- [8] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset, 2018.
- [9] Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T. Freeman, Michael Rubinstein, Yuanzhen Li, and Dilip Krishnan. Muse: Text-to-image generation via masked generative transformers, 2023.
- [10] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, 2020.
- [11] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [12] Jinwoo Choi, Chen Gao, Joseph C. E. Messou, and Jia-Bin Huang. Why can’t i dance in the mall? learning to mitigate scene bias in action recognition, 2019.
- [13] United States Congress. American data privacy and protection act, 2022.
- [14] Ji Dai, Behrouz Saghaei, Jonathan Wu, Janusz Konrad, and Prakash Ishwar. Towards privacy-preserving recognition of human activities. In *2015 IEEE International Conference on Image Processing (ICIP)*, pages 4238–4242. IEEE, 2015.
- [15] Ishan Rajendrakumar Dave, Chen Chen, and Mubarak Shah. Spact: Self-supervised privacy preservation for action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20164–20173, 2022.
- [16] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
- [17] Debidatta Dwibedi, Pavel Tokmakov, Ishan Misra, and Martial Hebert. Little help, big help: Simple unsupervised object detection. *arXiv preprint arXiv:2102.09084*, 2021.
- [18] Aleksandr Ermolov, Xiangyi Kong, Mikhail Petrov, and Cristian Sminchisescu. Whitening for self-supervised representation learning. *arXiv preprint arXiv:2007.06346*, 2021.

- [19] Eli Friedman, Assaf Lehr, Alexey Gruzdev, Vladimir Loginov, Max Kogan, Moran Rubin, and Orly Zvitia. Knowing the distance: Understanding the gap between synthetic and real data for face parsing, 2023.
- [20] Shreyank N Gowda, Marcus Rohrbach, Frank Keller, and Laura Sevilla-Lara. Learn2augment: Learning to composite videos for data augmentation in action recognition, 2022.
- [21] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Remi Munos, and Michal Valko. Bootstrap your own latent: A new approach to self-supervised learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [22] Yuejun Guo, Qing Xu, Shihua Sun, Xiaoxiao Luo, and Mateu Sbert. Selecting video key frames based on relative entropy and the extreme studentized deviate test. *Entropy*, 18(3):73, 2016.
- [23] Tengda Han, Jue Wang, Anoop Cherian, and Stephen Gould. Human action forecasting by learning task grammars. *arXiv:1709.06391*, 2017.
- [24] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [25] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn, 2018.
- [26] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015.
- [27] Sanjeewa Herath, Mehrtash T Harandi, and Fatih Porikli. Going deeper into action recognition: A survey. *Image and Vision Computing*, 60:4–21, 2017.
- [28] Itseez. Open source computer vision library. <https://github.com/itseez/opencv>, 2015.
- [29] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, Mustafa Suleyman, and Andrew Zisserman. The kinetics human action video dataset, 2017.
- [30] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. In *Proceedings of the 34th International Conference on Neural Information Processing Systems (NeurIPS)*, 2020.
- [31] Yo Whan Kim, Samarth Mishra, SouYoung Jin, Rameswar Panda, Hilde Kuehne, Leonid Karlinsky, Venkatesh Saligrama, Kate Saenko, Aude Oliva, and Rogerio Feris. How transferable are video representations based on synthetic data? In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022.
- [32] Sareh Koohpayegani, Tim Salimans, Ali Farhadi, and Hamed Javadi. Mean shift for self-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [33] H. Kuehne, H. Jhuang, E. Garrote, T. Poggio, and T. Serre. Hmdb: A large video database for human motion recognition. In *2011 International Conference on Computer Vision*, pages 2556–2563, 2011.
- [34] Jiefeng Li, Siyuan Bian, Chao Xu, Zhicun Chen, Lixin Yang, and Cewu Lu. Hybrik-x: Hybrid analytical-neural inverse kinematics for whole-body mesh recovery, 2023.
- [35] Ming Li, Xiangyu Xu, Hehe Fan, Pan Zhou, Jun Liu, Jia-Wei Liu, Jiahe Li, Jussi Keppo, Mike Zheng Shou, and Shuicheng Yan. Stprivacy: Spatio-temporal privacy-preserving action recognition. In *CVPR*, pages 1502–1511, 2022.

- [36] Tianjiao Li, Jun Liu, Wei Zhang, Yun Ni, Wenqian Wang, and Zhiheng Li. Uav-human: A large benchmark for human behavior understanding with unmanned aerial vehicles, 2021.
- [37] Yingwei Li, Yi Li, and Nuno Vasconcelos. Resound: Towards action recognition without representation bias. In Vittorio Ferrari, Martial Hebert, Cristian Sminchisescu, and Yair Weiss, editors, *Computer Vision – ECCV 2018*, volume 11210 of *Lecture Notes in Computer Science*. Springer, Cham, 2018.
- [38] Zhen Li, Cheng-Ze Lu, Jianhua Qin, Chun-Le Guo, and Ming-Ming Cheng. Towards an end-to-end framework for flow-guided video inpainting, 2022.
- [39] Jing Lin, Ailing Zeng, Haoqian Wang, Lei Zhang, and Yu Li. One-stage 3d whole-body mesh recovery with component aware transformer, 2023.
- [40] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco: Common objects in context, 2015.
- [41] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023.
- [42] Stuart P. Lloyd. Least squares quantization in pcm. *IEEE Transactions on Information Theory*, 28(2):129–137, 1982.
- [43] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts, 2017.
- [44] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019.
- [45] L Needham, M Evans, DP Cosker, L Wade, PM McGuigan, JL Bilzon, and SL Colyer. The accuracy of several pose estimation methods for 3d joint centre localisation. *Scientific Reports*, 11(1):20673, 2021.
- [46] Seong Joon Oh, Rodrigo Benenson, Mario Fritz, and Bernt Schiele. Faceless person recognition; privacy implications in social media, 2016.
- [47] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed AA Osman, Dimitrios Tzionas, and Michael J Black. Expressive body capture: 3d hands, face, and body from a single image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [48] Hieu H Pham, Louahdi Khoudour, Alain Crouzil, Pablo Zegers, and Sergio A Velastin. Video-based human action recognition using deep learning: A review. *arXiv preprint arXiv:2208.03775*, 2022.
- [49] AJ Piergiovanni and Michael S Ryoo. Avid dataset: Anonymized videos from diverse countries, 2020.
- [50] Francesco Pittaluga, Sanjeev Koppal, and Ayan Chakrabarti. Learning privacy preserving encodings through adversarial training. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 791–799. IEEE, 2019.
- [51] Zhongzheng Ren, Yong Jae Lee, and Michael S. Ryoo. Learning to anonymize faces for privacy preserving action detection. *arXiv preprint arXiv:1803.11556*, 2018.
- [52] Madeline Chantry Schiappa, Naman Biyani, Prudvi Kamtam, Shruti Vyas, Hamid Palangi, Vibhav Vineet, and Yogesh Rawat. Large-scale robustness analysis of video action recognition models, 2023.
- [53] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild, 2012.

- 430 [54] Sam Toyer, Anoop Cherian, Tengda Han, and Stephen Gould. Human pose forecasting via deep
431 Markov models. In *DICTA*, 2017.
- 432 [55] European Union. General data protection regulation, 2016.
- 433 [56] Limin Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, Dahua Lin, Xiaoou Tang, and Luc Van Gool.
434 Temporal segment networks: Towards good practices for deep action recognition, 2016.
- 435 [57] Zehao Wang, Christoph Feichtenhofer, Lorenzo Torresani, Du Tran, Joao Carreira, and Andrew
436 Zisserman. Videomae: Masked autoencoders are data-efficient learners for self-supervised
437 video pre-training. In *Advances in Neural Information Processing Systems*, 2022.
- 438 [58] Zhenyu Wu, Haotao Wang, Zhaowen Wang, Hailin Jin, and Zhangyang Wang. Privacy-
439 preserving deep action recognition: An adversarial learning framework and a new dataset. *IEEE*
440 *Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- 441 [59] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han
442 Hu. Simmim: A simple framework for masked image modeling, 2022.
- 443 [60] Kaiyu Yang, Jacqueline H Yau, Li Fei-Fei, Jia Deng, and Olga Russakovsky. A study of face
444 obfuscation in imagenet. In *International Conference on Machine Learning*, pages 25313–25330.
445 PMLR, 2022.
- 446 [61] Jang-Hee Yoo and Kyoung-Ho Choi. A review on video-based human activity recognition.
447 *Computers*, 2(2):88–131, 2013.
- 448 [62] Sangdoo Yun, Seong Joon Oh, Byeongho Heo, Dongyoon Han, and Jinhyung Kim. Videomix:
449 Rethinking data augmentation for video classification, 2020.
- 450 [63] Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-
451 supervised learning via redundancy reduction. *Proceedings of the 38th International Conference*
452 *on Machine Learning (ICML)*, 2021.
- 453 [64] Zhixiang Zhang, Thomas Cilloni, Charles Walter, and Charles Fleming. Multi-scale, class-
454 generic, privacy-preserving video. *Electronics*, 10(10):1172, 2021.
- 455 [65] Howard Zhong, Samarth Mishra, Donghyun Kim, SouYoung Jin, Rameswar Panda, Hilde
456 Kuehne, Leonid Karlinsky, Venkatesh Saligrama, Aude Oliva, and Rogerio Feris. Learning
457 human action recognition representations without real humans, 2023.
- 458 [66] Yuliang Zou, Jinwoo Choi, Qitong Wang, and Jia-Bin Huang. Learning representational
459 invariances for data-efficient action recognition, 2022.

Appendix

A Technical Specifications

General Details. All our model training is distributed over NVIDIA 8×A100-80GB (SXM). In all experiments, we employ the ViT-B backbone and to uphold rigorous privacy preservation, we develop conduct training without pre-trained weights from ImageNet. We pre-train VideoMAE for 200 epochs using a tube masking technique that masks 90% of image patches to enhance learned video representations. After pre-training, we remove the VideoMAE decoder, retaining only the encoder. For label alignment, we conduct supervised pre-training for 50 epochs using the same subset of 150 Kinetics classes as SynAPT [31]. For downstream evaluation, we fine-tune (FT) the entire network or train a linear probe (LP) for 30 epochs. Both steps use video inputs as 4D tensors (C, T, H, W) , with $C = 3$ (RGB channels), $T = 16$ frames, and spatial dimensions H and W as the video input is resized to 224×224 and normalized.

Table 5: Summary of Training Details.

General Specifications	
GPU Configuration	NVIDIA 8×A100-80GB (SXM)
Model Backbone	ViT-B (12 encoder blocks, 768 emb. dim.; 4 decoder blocks, 384 emb. dim.)
Input Tensor Shape	$3 \times 16 \times 224 \times 224$
Normalization	$\mu = [0.485, 0.456, 0.406], \sigma = [0.229, 0.224, 0.225]$
Self-Supervised Pre-Training	
Pre-training Method	VideoMAE
Epochs	200 (10 warm-up)
Masking Strategy	Tube, ratio = 0.9
Batch Size	128
Patch Size	$2 \times 16 \times 16$
Loss Function	MSE
Optimizer	AdamW
Learning Rate (Max)	0.0008
Learning Rate Scheduler	Cosine
Supervised Label Alignment	
Epochs	50 (6 warm-up)
Loss Function	Cross-Entropy
Optimizer	AdamW
Learning Rate (Max)	0.002
Downstream Evaluation	
Adjustment Epochs (FT or LP)	30

Step 1. During self-supervised pre-training, we use mean squared error (MSE) pixel reconstruction loss between the original and reconstructed frames. The batch size is 128, and patch sizes are $2 \times 16 \times 16$. We use the AdamW optimizer [44] with a cosine learning rate scheduler [43], a maximum learning rate of 0.0008, and a 10-epoch warm-up.

Step 2. For supervised label alignment, we add a final linear head to the ViT-B model for supervised training. The model shares an encoder with distinct linear classifiers for each dataset, using cross-entropy loss to measure discrepancies between the predicted and actual action categories. We train the model for 50 epochs with the AdamW optimizer and a cosine rate scheduler with a maximum learning rate of 0.002 and a 6-epoch warm-up.

B Dataset Considerations

B.1 Modified Clustering Operation

The K-NEXUS clustering operation, $\kappa(\bar{E}, K)$, aims to minimize the within-cluster bias. The clustering operation seeks to minimize the following objective from [37]: $\arg \min_{\omega} \sum_{k=1}^K \sum_{\bar{e} \in \omega_k} \mathcal{B}(D, \bar{e})$ where $\mathcal{B}(D, \bar{e})$ is the bias measurement for a dataset D using class embedding \bar{e} and is defined as:

$$\mathcal{B}(D, \bar{e}) = \log(\mathcal{M}(D, \bar{e})) - \log(\mathcal{M}_{\text{chance}})$$

Here, $\mathcal{M}(D, \bar{e})$ represents the performance of the representation \bar{e} on dataset D , and $\mathcal{M}_{\text{chance}}$ is the performance at the chance level, defined as:

$$\mathcal{M}_{\text{chance}} = \min_{\bar{e}} \mathcal{M}(D, \bar{e})$$

The centroid \bar{e}_k of cluster ω_k is defined as the mean of the bias measurements of all points in ω_k :

$$\bar{e}_k = \frac{1}{|\omega_k|} \sum_{\bar{e} \in \omega_k} \mathcal{B}(D, \bar{e})$$

This involves the following iterative steps:

1. **Assignment step.** Assign each point \bar{e}_i to the cluster with the nearest centroid based on the bias measurement:

$$\omega_k^{(t+1)} = \left\{ \bar{e}_i : \mathcal{B}(D, \bar{e}_i) \leq \mathcal{B}(D, \bar{e}_j^{(t)}), \forall j = 1, 2, \dots, K \right\}$$

2. **Update step.** Calculate the new centroids for each cluster:

$$\bar{e}_k^{(t+1)} = \frac{1}{|\omega_k^{(t+1)}|} \sum_{\bar{e}_i \in \omega_k^{(t+1)}} \mathcal{B}(D, \bar{e}_i)$$

Furthermore, the optimization problem to select a subset of classes from the original dataset, as laid out in [37], presents an exponential time complexity of $O(2^n)$. It is possible to converge to a solution for the selection of a small number of classes. However, it lacks feasibility for our case ($K = 150$ to obtain the Kinetics-150 dataset). Our K-NEXUS approach, converges while having the time complexity of $O(n \times k \times t)$, where n is the number of classes, k is the number of clusters, t is the number of update steps. Thus, we are able to perform class sampling for larger values with a linear time complexity.

B.2 SMPLY Failing



Figure 4: **Failure Cases.** Our augmentation framework suffers when the human joints are occluded. The pottery wheel, music stand, and drums are partially obscured by the superimposed mesh, demonstrating the challenges in handling occlusions within the scene.

In our investigation of the SMPLY Private framework’s M2M-augmentation, we identified challenges related to occlusion-based scenarios. To quantify this issue, we manually reviewed 5 randomly

497 selected videos per class from the Kinetics-150 dataset. Our findings indicated that occlusion-related
498 difficulties were present in only 2.8% of the videos. In these instances, the occlusions involved
499 informative objects or backgrounds that contributed to learned features and supervisory signals.
500 Additionally, there were a few cases outside this 2.8% where occlusions were present; however,
501 these did not involve the occlusion of significant objects or scenes essential to the video’s labels (i.e.,
502 potentially irrelevant parts of the video were occluded, not the major components). As advancements
503 in mesh models, particularly those capable of handling occlusions more effectively, continue to
504 emerge, we anticipate improvements in this aspect of the M2M framework.

Checklist

1. For all authors...

- (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [\[Yes\]](#)
- (b) Did you describe the limitations of your work? [\[Yes\]](#) See Section 5.
- (c) Did you discuss any potential negative societal impacts of your work? [\[Yes\]](#) See Section 5.
- (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)

2. If you are including theoretical results...

- (a) Did you state the full set of assumptions of all theoretical results? [\[N/A\]](#)
- (b) Did you include complete proofs of all theoretical results? [\[N/A\]](#)

3. If you ran experiments (e.g. for benchmarks)...

- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[Yes\]](#)
- (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#) See Appendix A.
- (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [\[No\]](#)
- (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [\[Yes\]](#) See Appendix A.

4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- (a) If your work uses existing assets, did you cite the creators? [\[Yes\]](#)
- (b) Did you mention the license of the assets? [\[N/A\]](#)
- (c) Did you include any new assets either in the supplemental material or as a URL? [\[Yes\]](#)
- (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [\[N/A\]](#)
- (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [\[N/A\]](#)

5. If you used crowdsourcing or conducted research with human subjects...

- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [\[N/A\]](#)
- (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [\[N/A\]](#)
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [\[N/A\]](#)