

(Supplementary) Data-Free Class-Incremental Hand Gesture Recognition

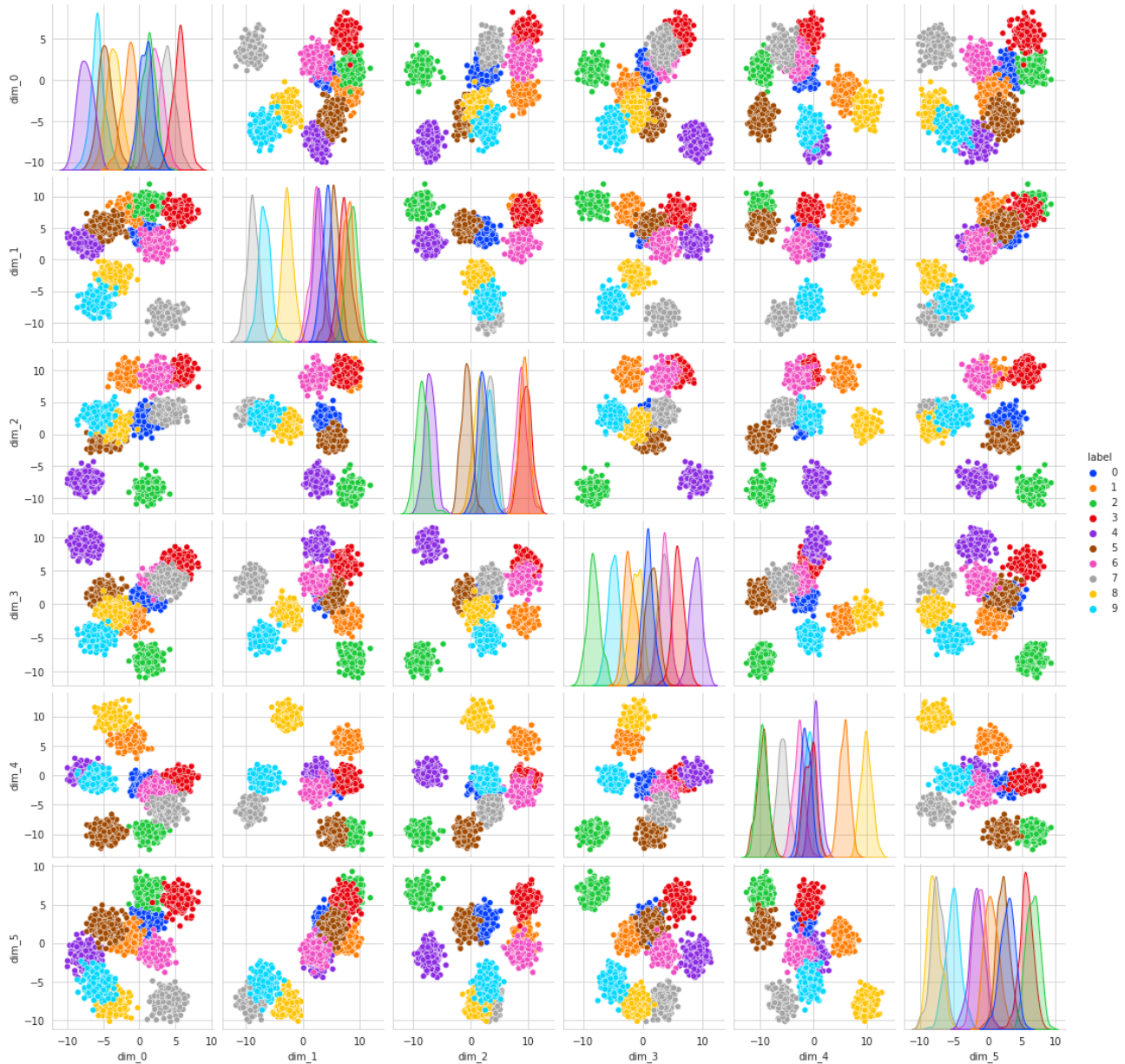


Figure 1: The pair plot of the synthetic dataset (6D Gaussian blobs) used for the preliminary proof of concept in Section 3 (BOAT-MI: Our Methodology) of the main paper. Different colors indicate different class indices (labels on the right). For the DFCIL setup, the first 5 classes (label 0 \rightarrow 4) are considered the old set \mathcal{C}_O , and the other 5 (label 5 \rightarrow 9) new \mathcal{C}_N .

Algorithm: Boundary-Aware Prototypical Model Inversion (Part 1)

```

1 Function Main ( $\mathcal{F}$ ,  $\mathcal{H}$ ,  $c$ ,  $\mu$ ,  $\Sigma$ ,  $m$ ,  $\alpha_f$ ,  $\alpha$ ,  $it$ ,  $\varepsilon$ ,  $n$ ,  $\delta$ ) :
2   #  $\mathcal{F}$ ;  $\mathcal{H}$ : feature extractor; SVM classifier
3   #  $c$ : index of class to invert
4   #  $\mu$ ;  $\Sigma$ : prototypical mean, covariance for class  $c$ 
5   #  $m$ ;  $\alpha_f$ ,  $\alpha$ : momentum; forward/reverse LR
6   #  $it$ ;  $\varepsilon$ : max iterations; tolerance threshold
7   #  $n$ ,  $\delta$ : max # of samples per class; margin
8   # Use inverted mean for initialization later
9    $\mathbf{x}_\mu \leftarrow \text{Invert}(\mathcal{F}, \mu, m, \alpha, it, \varepsilon, \text{NULL})$ 
10   $\mathbf{y}_s \leftarrow []$  # initialize feature list for inversion
11  # Extend the list with support vectors for class  $c$ 
12   $\mathbf{y}_s \text{ += } \mathcal{H}.\text{support\_vectors}[c]$ 
13  # Get principal directions of prototypical
14  # covariance and their unique linear combinations
15  # as prototypical features for inversion
16   $\mathbf{p}_s \leftarrow \text{GetProtoDirections}(\Sigma, n)$ 
17   $\mathbf{y}_s \text{ += } \text{GetProtoFeatures}(\mathbf{p}_s, \mu, c, \alpha_f, \delta)$ 
18   $n \leftarrow \text{len}(\mathbf{y}_s)$  # Final # of samples to invert
19   $\mathbf{x}_s \leftarrow [\mathbf{x}_\mu] * n$  # initialize input list for inversion
20  for  $i \leftarrow 1 : n$  do
21     $\mathbf{x}_s[i] \leftarrow \text{Invert}(\mathcal{F}, \mathbf{y}_s[i], m, \alpha, it, \varepsilon, \mathbf{x}_\mu)$ 
22  end
23  return  $\mathbf{x}_s$ 
24 End Function
25 Function GetProtoDirections ( $\Sigma$ ,  $n$ ) :
26   #  $\Sigma$ : class prototypical covariance
27   #  $n$ : maximum # of samples allowed
28   # Get principal directions retaining 95% variance
29    $\mathbf{p}_s \leftarrow \text{PCA}(\Sigma, 0.95)$ 
30    $\mathbf{p}_s \text{ += } (-\mathbf{p}_s)$  # extend with negative directions
31   if  $\text{ClampNumSamples}(\mathbf{p}_s, n)$  then
32     return  $\mathbf{p}_s$ 
33   end
34    $n \leftarrow n - \text{len}(\mathbf{p}_s)$  # update max # of samples
35    $ord = 2$  # start from adding 2nd order interactions
36   while TRUE do
37     # add unique linear interactions of order  $ord$ 
38     #  $p_{new} = \mathbf{p}_s[i] + \mathbf{p}_s[j]$  for  $ord = 2$ 
39     #  $p_{new} = \mathbf{p}_s[i] + \mathbf{p}_s[j] + \mathbf{p}_s[k]$  for  $ord = 3$ 
40      $\mathbf{p}_s.\mathbf{l} \leftarrow \text{LinearInteractions}(\mathbf{p}_s, ord)$ 
41     # normalize to get the unit vectors
42      $\mathbf{p}_s.\mathbf{l} \leftarrow \text{Normalize}(\mathbf{p}_s.\mathbf{l})$ 
43     if  $\text{ClampNumSamples}(\mathbf{p}_s.\mathbf{l}, n)$  then
44        $\mathbf{p}_s \text{ += } \mathbf{p}_s.\mathbf{l}$  # list extension
45       return  $\mathbf{p}_s$ 
46     end
47      $\mathbf{p}_s.\text{append}(\mathbf{p}_s.\mathbf{l})$ 
48      $n \leftarrow n - \text{len}(\mathbf{p}_s.\mathbf{l})$  # update max # of samples
49      $ord \leftarrow ord + 1$  # increment order
50   end
51 End Function

```

Algorithm: Boundary-Aware Prototypical Model Inversion (Part 2)

```

1 Function Invert ( $\mathcal{F}$ ,  $\mathbf{y}$ ,  $m$ ,  $\alpha$ ,  $iter$ ,  $\varepsilon$ ,  $\mathbf{x}_0$ ) :
2   #  $\mathcal{F}$ ;  $\mathbf{y}$ : feature extractor; target feature
3   #  $m$ ;  $\alpha$ : momentum; learning rate
4   #  $iter$ ;  $\varepsilon$ : max # of iterations; tolerance threshold
5   #  $\mathbf{x}_0$ : input initialization, can be NULL
6   if  $\mathbf{x}_0 == \text{NULL}$  then
7      $\mathbf{x} \leftarrow$  input prediction tensor (zero-initialized)
8   else
9      $\mathbf{x} \leftarrow \mathbf{x}_0$ 
10  end
11   $\mathbf{v} \leftarrow$  gradient update tensor (zero-initialized)
12  for  $i \leftarrow 1 : iter$  do
13     $\hat{\mathbf{y}} \leftarrow \mathcal{F}(\mathbf{x})$  # predict feature
14     $\mathcal{L} \leftarrow \|\mathbf{y} - \hat{\mathbf{y}}\|_2 / \|\mathbf{y}\|_2$  # distance norm
15    if  $\mathcal{L} < \varepsilon$  then
16      break
17    end
18     $\mathbf{v} \leftarrow m \mathbf{v} + \nabla_{\mathbf{x}} \mathcal{L}$  # gradient tensor update
19     $\mathbf{x} \leftarrow \mathbf{x} - \alpha \mathbf{v}$  # gradient descent
20  end
21  return  $\mathbf{x}$ 
22 End Function
23 Function GetProtoFeatures ( $\mathbf{p}_s$ ,  $\mu$ ,  $c$ ,  $\alpha_f$ ,  $\delta$ ) :
24   #  $\mathbf{p}_s$ : list of principal directions and combinations
25   #  $c$ ,  $\mu$ : class index and class prototypical mean
26   #  $\alpha_f$ ,  $\delta$ : forward learning rate and margin
27    $\mathbf{y}_s \leftarrow []$  # initialize feature list for inversion
28   for  $i \leftarrow 1 : \text{len}(\mathbf{p}_s)$  do
29      $\mathbf{f} \leftarrow \mu$ 
30     for  $k \leftarrow 1 : iter$  do
31        $\mathbf{f} \leftarrow \mathbf{f} + \alpha_f \mathbf{p}_s[i]$ 
32       if  $\text{Classify}(\mathbf{f}) \neq c$  then
33          $\mathbf{f} \leftarrow (1 - \delta) \mathbf{f} + \delta \mu$ 
34         ASSERT  $\text{Classify}(\mathbf{f}) == c$ 
35       end
36     end
37      $\mathbf{y}_s.\text{append}(\mathbf{f})$ 
38   end
39   return  $\mathbf{y}_s$ 
40 End Function
41 Function ClampNumSamples ( $\mathbf{p}_s$ ,  $n$ ) :
42   #  $\mathbf{p}_s$ : list of samples to clamp
43   #  $n$ : maximum # of samples allowed
44   if  $n > \text{len}(\mathbf{p}_s)$  then
45     return FALSE # not clamped
46   end
47   # randomly select  $n$  samples in-place
48    $\mathbf{p}_s \leftarrow \text{RandomSelect}(\mathbf{p}_s)$ 
49   return TRUE # clamped
50 End Function

```

Table 1: Results with BatchNorm (average of 3 runs with different class order) for DFCIL over six tasks in SHREC-2017.

Method	Task 0	Task 1	Task 2	Task 3	Task 4	Task 5	Task 6	Mean (Task 1 → 6)
Oracle (DGSTA-BN)	90.3							
		G↑ IFM↓	G↑ IFM↓	G↑ IFM↓	G↑ IFM↓	G↑ IFM↓	G↑ IFM↓	G↑ IFM↓
DeepInversion-BN [6]	92.8	81.4 7.0	62.6 18.2	52.9 29.1	46.3 35.6	34.8 48.1	31.5 51.3	51.6 31.6
ABD-BN [4]		81.7 1.9	72.2 2.0	65.9 10.7	57.3 14.8	51.7 16.3	42.7 25.2	61.9 11.8
R-DFCIL-BN [3]		72.0 5.1	61.5 7.3	53.2 9.6	50.7 5.1	46.6 8.9	39.9 0.6	54.0 6.1
BOAT-MI-BN (Ours)		86.5 7.5	82.2 9.3	74.7 1.2	72.6 2.6	65.5 8.7	62.3 1.7	74.0 5.2

Table 2: Results with BatchNorm (average of 3 runs with different class order) for DFCIL over six task in EgoGesture3D.

Method	Task 0	Task 1	Task 2	Task 3	Task 4	Task 5	Task 6	Mean (Task 1 → 6)
Oracle (DGSTA-BN)	75.0							
		G↑ IFM↓	G↑ IFM↓	G↑ IFM↓	G↑ IFM↓	G↑ IFM↓	G↑ IFM↓	G↑ IFM↓
DeepInversion-BN [6]	77.4	66.4 16.1	52.3 24.6	35.0 46.5	24.5 45.1	59.7 20.5	65.0 15.1	72.8 35.6
ABD-BN [4]		65.7 17.4	58.1 19.4	49.5 31.1	45.1 35.2	42.7 37.2	39.5 40.4	50.1 30.1
R-DFCIL-BN [3]		68.3 5.9	56.0 11.9	49.0 17.9	44.4 20.1	39.9 29.4	36.8 30.6	49.1 19.3
BOAT-MI-BN (Ours)		73.0 4.6	66.9 7.1	60.8 15.9	57.9 16.3	52.1 22.2	49.0 27.6	60.0 15.6

1. Experimental Details

We follow the original implementations of all DFCIL methods that we used for benchmarking in this work. However, these methods were originally designed for the image classification problem, as we discussed in Section 4.4 of the main paper. Therefore, we had to adjust them to our keypoint-based gesture recognition domain and tune their hyperparameters individually for our setup. All hyperparameter values and configuration files will be made available along with the codebase.

Base model training: We use the original DG-STA [1] architecture as our backbone to process the 3D keypoint sequences. We train the Oracle (upper bound, all classes in one task) and Task 0 models (base classes for each setup) with Adam optimizer (learning rate of $1e-3$, $\beta_1 = 0.9$ and $\beta_2 = 0.999$) for 150 epochs.

Learning rate for incremental tasks: All the methods start from the same Task 0 weights for Task 1. We modify the initial learning rate for the incremental tasks depending on the method since a stronger regularization requires a higher learning rate to learn the new classes. In this regard, we select the optimal learning rate for each method with a grid search in the range $\{1e-5 : 1e-3\}$.

Synthetic data generation: Following ABD [4], for a fair comparison among the model inversion-based approaches (DeepInversion [6], ABD [4] and R-DFCIL [3]), we use the same model inversion strategy for the synthetic data generation. However, we optimize the randomly initialized inputs directly instead of training a model to generate the samples. After a grid search we find $\{\alpha_{lr}, \alpha_{con}, \alpha_{stat}, \alpha_{temp}\}$ as $\{1e-1, 1, 2e2, 1e1\}$.

R-DFCIL hyperparameters: R-DFCIL [3] includes three hyperparameters to weight each of the loss terms, the **local**

classification loss (λ_{lce}), the **hard knowledge distillation loss** (λ_{hkd}) and the **relational knowledge distillation loss** (λ_{rkd}). For our experiments, we found $\{\lambda_{lce}, \lambda_{hkd}, \lambda_{rkd}\}$ as $\{1, 1.5e-1, 7.5e-1\}$ as the optimal values for SHREC-2017 and $\{1, 1e-1, 1e-1\}$ for EgoGesture3D.

BOAT-MI implementation: The complete Python-style pseudocode is provided in Algorithm 1 and 2. In Algorithm 1 Line 1, `Function Main` is the driver for our proposed model inversion mechanism for a single class c . Similar to the other baselines, we employ DG-STA [1] backbone as our feature extractor \mathcal{F} here. We use the radial basis function as the kernel for our SVM classifier (\mathcal{H}) here. To generate proto-SVs for model inversion, it uses the principal directions of the prototypical variance Σ (Line 29, Algorithm 1). Ray casting from the class-prototypical mean μ is done iteratively with a learning rate α_f (Line 29-36, Algorithm 2) until the ray hits the boundary. Also, a margin δ , normalized with respect to the distance between the mean μ and the boundary vector, is used to avoid noise inherent in the boundary estimation process (Figure 2 main paper, Line 33, Algorithm 2). In other words, the proto-SV feature is taken to be the on the margin, which is δ inside the boundary.

Note that the number of principal directions (PDs) is pretty low compared to the dimensionality of the feature vector. Along with these raw PDs, we consider upto third-order interactions, i.e. simple unweighted linear combination of PDs ($p[i] + p[j]$ and $p[i] + p[j] + p[k]$, Line 37-40, Algorithm 1, $p[\cdot]$ is the list of PDs) ignoring duplicates, which we found to be empirically sufficient for all our studies. Thus, our augmentations are deterministic (not random), conditioned on the set of PDs.

The feature inversion procedure (feature \rightarrow input, Line 1-22, Algorithm 2) deviates from the vanilla implementation [2] in three aspects. First, we initialize the input tensor with

the inverted class-prototypical mean (Line 9, Algorithm 2). Second, we employ the normalized \mathcal{L}_2 function as the distance metric (Line 14, Algorithm 2). Third, we replace the vanilla gradient descent with its momentum-based counterpart [5] (Line 18-19, Algorithm 2). All these modifications are empirically found to expedite the convergence.

The forward learning rate (α_f), momentum (m), and reverse learning rate (α) are set to 0.05, 0.9, and 1.0 respectively. We use the value of 0.2 as our normalized margin (δ) based upon the ablation study presented in Table 4 of the main paper. The rest of the parameters are set similarly to the baselines. All these configurations will be open-sourced with the codebase.

2. Additional Experiments

The effect of stat alignment loss during model inversion: Model inversion-based methods like DeepInversion, ABD and R-DFCIL use the BatchNorm (BN) statistics to compute a regularization loss that aligns the synthetic and real data distributions during model inversion. This loss is based on the KL divergence between the synthetic data distribution and the BatchNorm distribution of the previous task model, which reflects the real data statistics from the previous training. However, our backbone architecture, DG-STA [1], does not have BatchNorm layers and thus this loss term is not included. To investigate how the **stat alignment loss** affects model inversion for these baselines, we replace the Layer-Norm layers with BatchNorm layers in DG-STA and run additional experiments.

Table 1 and 2 show the results for SHREC-2017 and EgoGesture3D setups, respectively. As can be seen, the incorporation of the BN layer does not change the overall ranking of the methods. With BN, our proposed approach achieves significantly higher global accuracy in each stage there with 12.1% and 9.9% improvements on average for SHREC-2017 (Table 1) and EgoGesture3D (Table 2), respectively, over the next best methods. Such improvements are also accompanied by the lower mean instantaneous forgetting in general – 0.9% and 3.7% lower *IFM* than the second best method for SHREC-2017 (Table 1) and EgoGesture3D (Table 2), respectively. Therefore, regardless of the choice of normalization layers, BOAT-MI excels the SOTA methods by a large margin.

Performance degradation/comparison to when all data is available, with large number of incremental stages: We first clarify that the upper bound of performance when all data is available – 75.8% – is reported on the “Oracle” row of Table 3 in our paper. To study this further with large number of incremental stages, we experimented with an extreme incremental learning setting on EgoGesture3D, where we added one class at a time to a model trained on the 59 base classes. Concretely, instead of adding 4 classes at a time over

6 incremental stages (as in paper), we added 1 new class at a time over 24 incremental stages. These results are shown in Figure 2 for both global accuracy and IFM. Our (BOAT-MI’s) performance drops by about half after 24 incremental stages – much lower compared to SOTA methods.

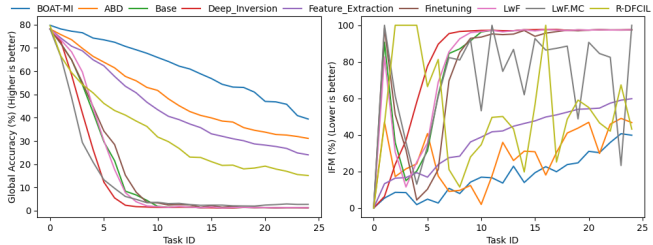


Figure 2: Global accuracy (%) and IFM (%) for 59 + 1 × 24 setup on EgoGesture3D dataset. Like other setups, our BOAT-MI excels the SOTA approaches here as well. Best viewed in digital format.

3. EgoGesture3D Processing

The new 3D skeleton dataset used in this paper is the derivative of the EgoGesture [8] comprising RGB-D images. Following the nomenclature of the parent dataset, we call it EgoGesture3D. Thanks to the MediaPipe [7] API that we use for 3D skeleton extraction from the images. The temporal association is turned on during the detection process with minimum detection and tracking confidences of 0.8 and 0.5, respectively. These numbers are chosen empirically for better extraction results.

However, merely using the extracted output from the temporally segmented gesture sequences is problematic for several reasons. First, the MediaPipe API fails to detect all the keypoints in all frames successfully. Next, it detects additional spurious hands (more than 2) in some of the frames. Also, in a few cases, for two hands detection, we find the handedness prediction to be the same (both left or both right hands). These issues are resolved with proper heuristics for each frame independently following the temporal detection with MediaPipe. Moreover, the baseline DG-STA model provides the optimal performance with 8 frames generally equally distant along the time dimension. Therefore, we save the samples with at least 8 frames of successful detections. In this regard, we experimented with interpolating with the less number of original frames. We find the oracle accuracy comparatively much lower that prevents us from saving samples with less than 8 frames. For the single-handed gestures, the values for the keypoints representing the absent hand are set to zeros. Very occasionally hand swapping occurs for just a few frames (i.e. 1-2 frames detect right hands with all others left). We keep these detections unchanged as some gestures may contain one hand mostly with the other in a cameo. Figure 3 and 4 show sample gestures in each column with the 2D projection of the skeleton on top of the corresponding RGB images.

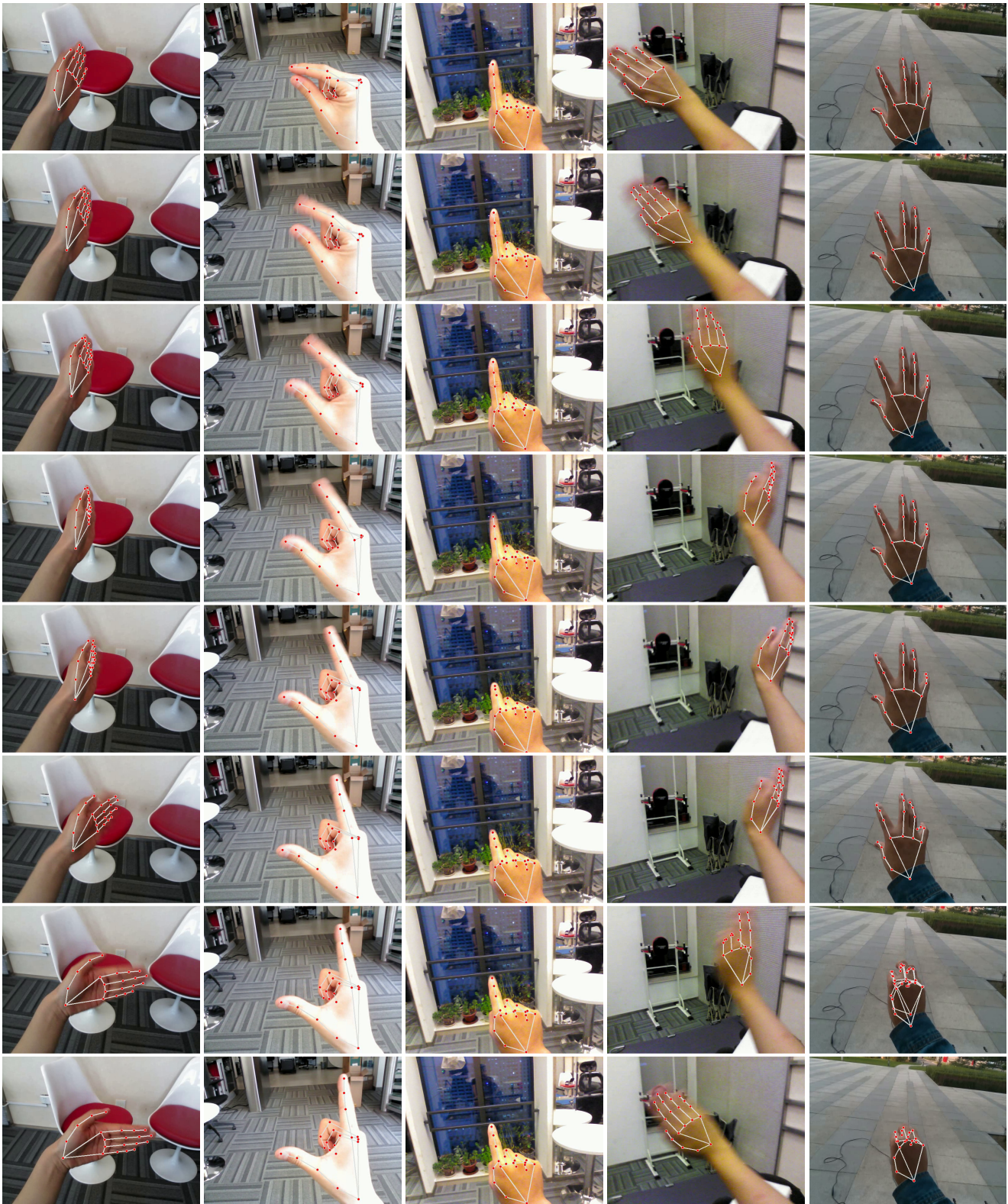


Figure 3: 2D projection of the EgoGesture3D skeletons on EgoGesture RGB images. (Top to bottom) Each column shows the temporal sequence for a single gesture (8 successfully detected frames selected in temporal order at random). (Left to right) [1] wave palm towards right; [12] zoom in with two fingers; [16] click with index finger; [39] wave hand; [48] grab.

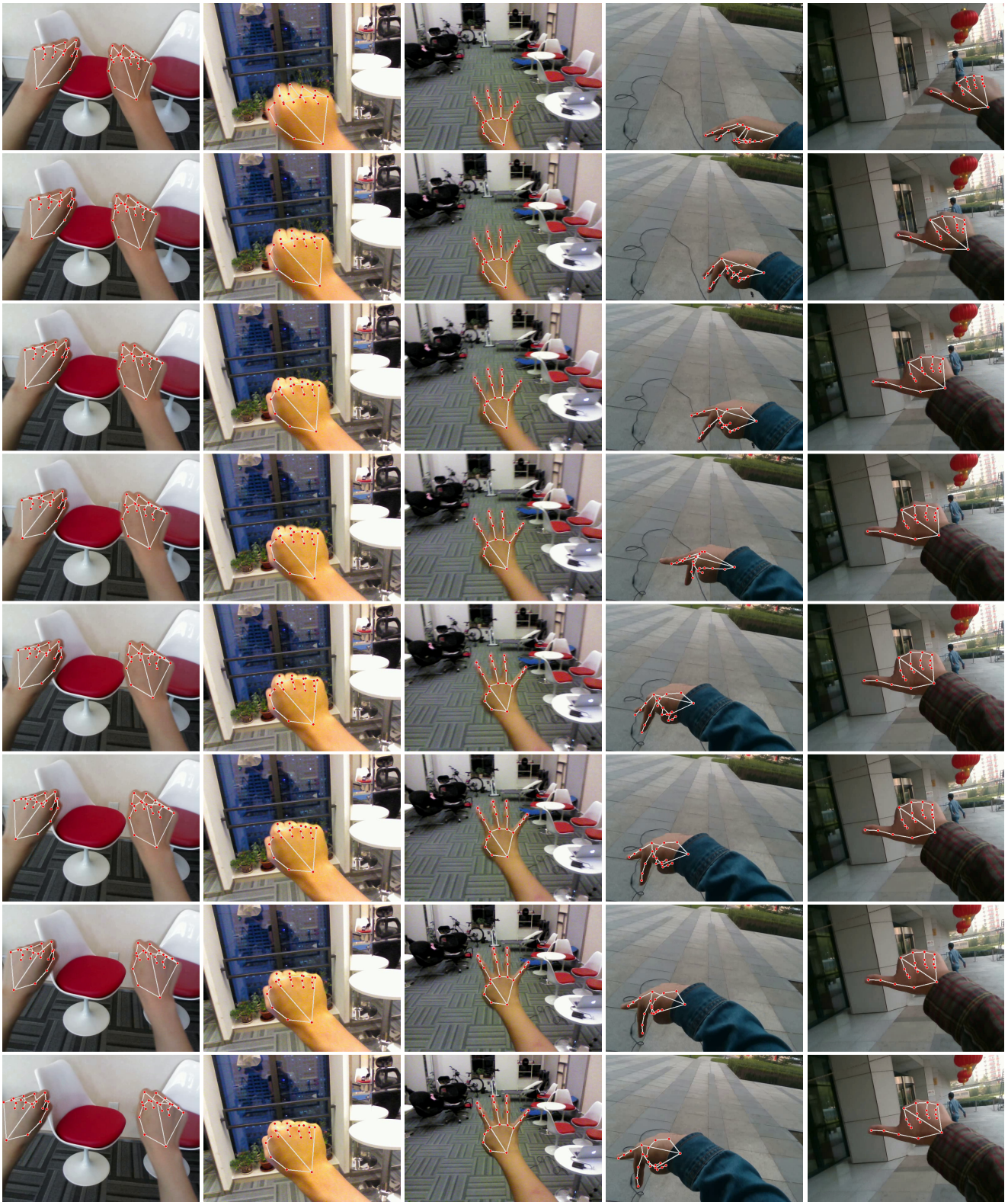


Figure 4: 2D projection of the EgoGesture3D skeletons on EgoGesture RGB images. (Top to bottom) Each column shows the temporal sequence for a single gesture (8 successfully detected frames selected in temporal order at random). (Left to right) [8] zoom in with two fingers; [21] static fist; [28] number 4; [49] walk; [66] thumb towards left.

References

- [1] Yuxiao Chen, Long Zhao, Xi Peng, Jianbo Yuan, and Dimitris Metaxas. Construct dynamic graphs for hand gesture recognition via spatial-temporal attention. In *BMVC*, 2019.
- [2] Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. Model inversion attacks that exploit confidence information and basic countermeasures. In *ACM Conference on Computer and Communications Security*, 2015.
- [3] Qiankun Gao, Chen Zhao, Bernard Ghanem, and Jian Zhang. R-dfcil: Relation-guided representation learning for data-free class incremental learning. In *ECCV*. Springer, 2022.
- [4] James Smith, Yen-Chang Hsu, Jonathan Balloch, Yilin Shen, Hongxia Jin, and Zsolt Kira. Always be dreaming: A new approach for data-free class-incremental learning. In *ICCV*, 2021.
- [5] Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In *ICML*, 2013.
- [6] Hongxu Yin, Pavlo Molchanov, Jose M. Alvarez, Zhizhong Li, Arun Mallya, Derek Hoiem, Niraj K. Jha, and Jan Kautz. Dreaming to distill: Data-free knowledge transfer via deepinversion. In *CVPR*, 2020.
- [7] Fan Zhang, Valentin Bazarevsky, Andrey Vakunov, Andrei Tkachenka, George Sung, Chuo-Ling Chang, and Matthias Grundmann. Mediapipe hands: On-device real-time hand tracking. *arXiv preprint arXiv:2006.10214*, 2020.
- [8] Yifan Zhang, Congqi Cao, Jian Cheng, and Hanqing Lu. Egogesture: A new dataset and benchmark for egocentric hand gesture recognition. *IEEE TMM*, 2018.