

Multimodal Cross-Domain Few-Shot Learning for Egocentric Action Recognition

– Supplementary Materials –

Masashi Hatano¹, Ryo Hachiuma², Ryo Fujii¹, and Hideo Saito¹

¹ Keio University

² NVIDIA



Fig. 4: Samples from each dataset. A curated selection of RGB images from each dataset showcases the domain gap between the source and target datasets.

A Datasets

In this section, we delve into the datasets used in our study, highlighting the significant domain gap observed in RGB images between the source (Ego4D [3]) and target (EPIC-Kitchens [2], MECCANO [9], and WEAR [1]) datasets. Fig. 4 presents a curated selection of RGB images drawn from each dataset. These visual examples underscore the diversity and complexity of cross-domain few-shot learning (CD-FSL) tasks.

- **Ego4D.** The Ego4D dataset comprises a diverse range of activities, including domestic chores in houses, gardening, cleaning, building, cooking, and arts/crafting. These actions contain those performed using a single hand, both hands, and various tools, demonstrating the diversity of human activities. Thus, this dataset is well-suited for learning generalizable features required for cross-domain and few-shot settings.
- **EPIC-Kitchens.** The EPIC-Kitchens dataset, serving as the target, is centered around kitchen activities, encompassing actions like “pouring flour”, “opening the refrigerator”, “moving a pizza”, and “pouring a smoothie”, each composed of a verb-noun pair. Such activities in the kitchens are also present within the Ego4D dataset, our source. However, the EPIC-Kitchens dataset distinguishes itself by its fine-grained action categories, with some actions

sharing verbs but differing in nouns and others vice versa. This granular differentiation of actions introduces a challenge for few-shot learning models, necessitating nuanced discernment in action categorization.

- **MECCANO.** The MECCANO dataset comprises detailed recordings of the assembly process for a toy bicycle in an industrial-like scenario. It contains fine-grained actions, including “aligning a screwdriver to screw”, “aligning objects”, “putting a tire”, and “putting a screw”. The diminutive size of the components poses a substantial challenge for action recognition, requiring exceptional precision to identify and understand the intricate interactions.
- **WEAR.** This dataset captures outdoor workout actions such as “stretching hamstrings”, “jogging”, “pushing-ups”, and “sitting-ups”. Unlike the other datasets, WEAR focuses on activities not involving hand-object interactions but body movements.

B Implementation Details

Experimental Setup. We utilize the three modalities: RGB, optical flow, and hand pose. For each modality, we select an input sequence comprising $T = 16$ frames sampled at a frequency of 8 FPS (frames per second). Spatial dimensions are standardized at 224×224 for RGB and optical flow inputs, while hand pose inputs, represented by the heatmap, are resized to 56×56 . The number of channels C_m are 3, 2, and 21 for RGB, optical flow, and hand pose, respectively. We employ FlowFormer [4] for estimating the optical flow between consecutive frames. For the prediction of 2D hand keypoints, which include 21 joints, we use RTMPose [5] trained on five public hand pose datasets available through MMPOSE³. Subsequently, the Gaussian heatmap of size 56×56 is produced, with a standard deviation σ set to 10. Following the VideoMAE [10] experimental setup, we adopt the mask ratio ρ_{pretrain} of 0.9 during our pretraining stage. In addition, we use the mask ratio $\rho_{\text{distill}} = 0.75$ during the multimodal distillation for all experiments. We use a machine equipped with Intel Xeon W-3235 CPU, 128GB RAM, and the NVIDIA Titan RTX GPU to compute the inference speed.

Training. In the pretraining stage, we train the model, which consists of ViT-S models, pretrained on the Kinetics-400 dataset [6] and a classifier head, which is attached to the ViT-S backbone, for all modalities for 100 epochs. For training settings, we generally follow the VideoMAE [10]. During the multimodal distillation stage, we train the student RGB model for 100 epochs using the AdamW optimizer [8], with a peak learning rate of $2e - 3$, linearly increased for the first 10 epochs of the training and decreased to $1e - 6$ until the end of training with cosine decay [7]. Note that we linearly scaled the peak learning rate with respect to the overall batch size. Regarding the parameters for the loss function, we empirically adapt the balancing hyperparameters λ_{ce_m} to $5e - 2$ for RGB, and $1e - 2$ for optical flow and hand pose input modality.

Evaluation Metrics. Following the existing CD-FSL work, we report the top-1 accuracy on the query set Q in the target validation set over 600 runs to

³ <https://github.com/open-mmlab/mmpose>

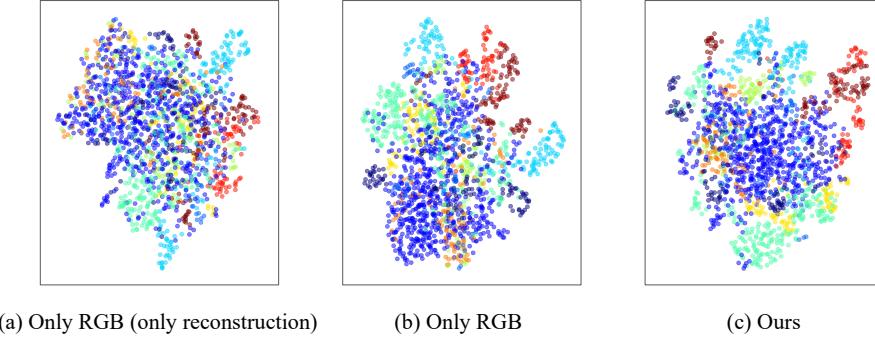


Fig. 5: Comparative UMAP visualization of feature representations. UMAP plot of 10 classes from EPIC-Kitchens validation set with features obtained from (a) Only RGB (only reconstruction), (b) Only RGB, and (c) Ours.

measure action recognition performance. To benchmark efficiency, we quantify the model’s performance by measuring the forward pass time during inference. We report the inference time averaged over 600 iterations.

C Visualization of Feature Representations

We compare the class-discriminativeness of embeddings extracted from three encoders: Only RGB (only reconstruction), Only RGB, and Ours. The only RGB (only reconstruction) model is trained without multimodal distillation, and $\lambda_{ce_{RGB}} = 0$ is used during the pretraining stage. The only RGB model is trained without multimodal distillation, and $\lambda_{ce_{RGB}} = 0.05$ is used during the pretraining stage. Fig. 5 shows the UMAP plot of 10 classes from EPIC-Kitchens datasets. We see that only RGB model creates better grouping on the embeddings of the target datasets than only RGB (only reconstruction) model. This result supports that using the cross-entropy loss helps learn the class-discriminative features during the pretraining stage. We further see that the multimodal distillation also helps learn discriminative features compared to the only RGB model.

D Loss Weight Ablation

We present an ablation study focused on the impact of adjusting the loss weight for the cross-entropy loss on the RGB modality $\lambda_{ce_{RGB}}$ during the pretraining stage. The loss weight for the cross-entropy loss $\lambda_{ce_{RGB}}$ serves as a critical hyperparameter that balances the contribution of the cross-entropy loss to the total loss function. For the ablation study, we varied the value of the loss weight

Table 6: Loss weight ablation. We conduct an ablation study on the loss weight for cross-entropy loss on RGB modality pertaining.

$\lambda_{\text{ce}_{\text{RGB}}}$	1	0.1	0.05	0.01
5-shot	54.58	56.68	57.07	52.40

across a predefined range $\lambda_{\text{ce}_{\text{RGB}}} \in \{1, 0.1, 0.05, 0.01\}$. We report the 5-way 5-shot action recognition accuracy on the EPIC-Kitchens dataset of the only RGB model with varied hyperparameter $\lambda_{\text{ce}_{\text{RGB}}}$ in Tab. 6. Our ablation analysis reveals a critical insight: a high or low loss weight for the RGB modality, $\lambda_{\text{ce}_{\text{RGB}}}$, during the pretraining stage can detrimentally affect the acquisition of class-discriminativeness on the target data. Assigning a high loss weight to the cross-entropy reduces the relative contribution of learning from unlabeled target data. Conversely, a low loss weight fails to capture class discriminativeness adequately.

References

- Bock, M., Moeller, M., Van Laerhoven, K., Kuehne, H.: Wear: A multimodal dataset for wearable and egocentric video activity recognition. arXiv preprint arXiv:2304.05088 (2023)
- Damen, D., Doughty, H., Farinella, G.M., Fidler, S., Furnari, A., Kazakos, E., Moltisanti, D., Munro, J., Perrett, T., Price, W., Wray, M.: Scaling egocentric vision: The epic-kitchens dataset. In: ECCV (2018)
- Grauman, K., Westbury, A., Byrne, E., Chavis, Z., Furnari, A., Girdhar, R., Hamburger, J., Jiang, H., Liu, M., Liu, X., Martin, M., Nagarajan, T., Radosavovic, I., Ramakrishnan, S.K., Ryan, F., Sharma, J., Wray, M., Xu, M., Xu, E.Z., Zhao, C., Bansal, S., Batra, D., Cartillier, V., Crane, S., Do, T., Doulaty, M., Erapalli, A., Feichtenhofer, C., Fragomeni, A., Fu, Q., Gebreselasie, A., González, C., Hillis, J., Huang, X., Huang, Y., Jia, W., Khoo, W., Kolář, J., Kottur, S., Kumar, A., Landini, F., Li, C., Li, Y., Li, Z., Mangalam, K., Modhuguri, R., Munro, J., Murrell, T., Nishiyasu, T., Price, W., Ruiz, P., Ramazanova, M., Sari, L., Somasundaram, K., Southerland, A., Sugano, Y., Tao, R., Vo, M., Wang, Y., Wu, X., Yagi, T., Zhao, Z., Zhu, Y., Arbeláez, P., Crandall, D., Damen, D., Farinella, G.M., Fuegen, C., Ghanem, B., Ithapu, V.K., Jawahar, C.V., Joo, H., Kitani, K., Li, H., Newcombe, R., Oliva, A., Park, H.S., Rehg, J.M., Sato, Y., Shi, J., Shou, M.Z., Torralba, A., Torresani, L., Yan, M., Malik, J.: Ego4d: Around the world in 3,000 hours of egocentric video. In: CVPR (2022)
- Huang, Z., Shi, X., Zhang, C., Wang, Q., Cheung, K.C., Qin, H., Dai, J., Li, H.: FlowFormer: A transformer architecture for optical flow. In: ECCV (2022)
- Jiang, T., Lu, P., Zhang, L., Ma, N., Han, R., Lyu, C., Li, Y., Chen, K.: Rtm-pose: Real-time multi-person pose estimation based on mmpose. arXiv preprint arXiv:2303.07399 (2023)
- Kay, W., Carreira, J., Simonyan, K., Zhang, B., Hillier, C., Vijayanarasimhan, S., Viola, F., Green, T., Back, T., Natsev, P., et al.: The kinetics human action video dataset. arXiv preprint arXiv:1705.06950 (2017)
- Loshchilov, I., Hutter, F.: SGDR: Stochastic gradient descent with warm restarts. In: ICLR (2017)

8. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. In: ICLR (2019)
9. Ragusa, F., Furnari, A., Farinella, G.M.: Meccano: A multimodal egocentric dataset for humans behavior understanding in the industrial-like domain. Computer Vision and Image Understanding (CVIU) (2023)
10. Tong, Z., Song, Y., Wang, J., Wang, L.: VideoMAE: Masked autoencoders are data-efficient learners for self-supervised video pre-training. In: NeurIPS (2022)