# Mutual Modality Learning for Video Action Classification

Stepan Komkov, Maksim Dzabraev, Aleksandr Petiushko

Lomonosov Moscow State University

Huawei Moscow Research Center

stepan.komkov@intsys.msu.ru, dzabraev.maksim@intsys.msu.ru, petyushko.alexander1@huawei.com

*Abstract*—The construction of models for video action classification progresses rapidly. However, the performance of those models can still be easily improved by ensembling with the same models trained on different modalities (e.g. Optical flow). Unfortunately, it is computationally expensive to use several modalities during inference. Recent works examine the ways to integrate advantages of multi-modality into a single RGB-model. Yet, there is still a room for improvement. In this paper, we explore the various methods to embed the ensemble power into a single model. We show that proper initialization, as well as mutual modality learning, enhances single-modality models. As a result, we achieve state-of-the-art results in the Something-Something-v2 benchmark.

*Index Terms*—Video Recognition, Video Action Classification, Video Labeling, Mutual Learning

## I. Introduction

Video Recognition has progressed a lot during the last several years. Datasets have enlarged from thousands of clips [15], [23] to hundreds of thousands [1], [10], [14] and even to hundreds of millions [17]. Neural network-based approaches for video processing evolved from simple 3D-convolutions [25] to Parvo- and Magnocellular counterparts emulation [7] and absorbed developments of classical Image Recognition [2], [16].

Nevertheless, classical results in the domain of Video Processing are still useful: Optical Flow estimation for a video sequence can significantly improve the quality of video recognition [22]. However, the common ways to estimate Optical Flow require an amount of calculation that is comparable to the whole further neural network inference. That is why a number of works are devoted to the implicit Optical Flow estimation during the RGB-based neural network inference [3], [5], [19], [24].

In our work, we target not only the improvement of RGB-based models but also simultaneous improvement of different single-modality models. To this end, we utilize Mutual Learning [34] that enables us to share the knowledge between single-modality models in both directions. We combine it with proper initialization that develops the performance of trained models.

We show that our approach not only improves each single-modality model but also boosts RGB-based models better than existing methods. Additionally, we examine how to use Mutual Learning to achieve the best results of multi-modality ensemble. Thus, we achieve state-of-the-art (SOTA) results among

the ones reported previously in the Something-Something-v2 benchmark [10].

Our code is available as a fork from the code presented in [16] (https://github.com/papermsucode/mutual-modality-learning).

## II. Related Works

First, we briefly describe the most common approaches for video action classification from the historical perspective. They can be divided into two groups: models with 3d-convolutions and models with 2d-convolutions.

Second, we describe the methods that improve single-modality models using other modalities and highlight the differences from our work.

### A. 3D-approaches

A video sequence is a 4d-tensor with the following parameters: height of frames, width of frames, number of frames, and number of channels per frame (3 in case of RGB input). Therefore, we can process it using Convolutional Neural Networks (CNNs) where 3d-convolutions are applied instead of 2d convolutions (with the new temporal dimension). Tran *et al.* are the first to propose 3d convolutional networks based on this idea [25]. Thereby, they achieved the SOTA results in a number of tasks.

Although a video model has to obtain temporal and motion information from the sequence of frames, it still needs to recognize the spatial information contained in each frame. Carreira and Zisserman propose to inflate the trained weights of the Image Recognition network and to use them as initialization for 3d-CNN [2]. Nowadays, this is a common approach for video model initialization.

Wang *et al.* implement an attention mechanism that helps to find dependencies between far positions on different frames. That is meaningful for fast-moving objects or quick movements of the camera [29].

The disadvantage of 3d-CNNs is that they require to work with much more parameters in comparison to their 2d analogs. To address this problem, the first convolutions can be replaced by the per-frame 2d-convolutions (top-heavy models) since those convolutions are mostly responsible for the evaluation of spatial features [30], [35]. Also, 3d-convolutions can be decomposed as 2d spatial-convolutions plus 1d temporal-convolutions. This kind of decomposition reduces the number

of parameters and operations and increases non-linearity at the same time [26], [30].

Feichtenhofer *et al.* present a SlowFast Network architecture that emulates Parvo- and Magnocellular counterparts by sampling video frames with two different framerates and by feeding them to two branches with different computational power [7]. Thus, the lightweight Fast pathway captures motion and temporal dynamics while the Slow pathway captures the spatial semantics. This approach achieved the SOTA results on the Kinetics-400 Action Classification dataset [14] among models without additional data.

The Temporal Pyramid Networks (TPN) of Yang *et al.* can be viewed as an extension of SlowFast networks [31]. A thinned out frames sequence flows to the different branches from intermediate layers instead of entering from the input. This approach is an add-on to existing architectures and can be implemented for 3d-CNNs and 2d-CNNs.

### B. 2D-approaches

The early CNN-based models for video with 2d-convolutions consist of two streams. The first stream called Spatial takes RGB frames as an input. The second stream called Temporal takes a stack of consecutive Optical Flow estimations [22], [27]. The final prediction is an average of the predictions of both streams. Note that the authors use pretrained weights from Image Recognition models for the temporal stream as well as for the spatial stream.

Nowadays, the idea of features sharing between frames is used to simulate a 3d-inference using 2d-convolutions. The pioneering work in this scope is Temporal Shift Modules network (TSM) by Lin *et al.* that applies ordinary 2d-ResBlocks [11] to each input frame [16]. The single difference is that TSM replaces a one-eighth of channels with the same channels from the previous frame and another one-eighth of channels with the same channels from the future frame before each first convolution of the ResBlock. Thereby, the authors achieved the SOTA results on the Something-Something-v2 dataset [10] and provided a powerful and efficient baseline for the future research in this area.

Based on the idea of feature sharing, Shao *et al.* present Temporal Interlacing Network [20]. The authors add extra lightweight blocks that decide on distances and weights of channels sending within each ResBlock. This approach is used instead of a fixed replacement of channels.

### C. Optical Flow distillation

Despite all the aforementioned progress, most of works can be improved by averaging of their predictions with the predictions of the same network trained on the Optical Flow modality [2], [16], [26], [27], [30], [35].

Since the Optical Flow calculation is a time-consuming operation, a number of works is devoted to incorporation of the motion-estimation blocks inside the CNN architecture [5], [13], [19]. However, knowledge distillation from the Optical Flow modality to any RGB single-modality network seems to be of more interest.

Three basic works that should be mentioned are Knowledge Distillation (KD) [12], Mutual Learning (ML) [34], and Born-Again Networks (BAN) [8].

The first proposes to use soft-predictions of the model called Teacher network to train the smaller model called Student network. It turns out that this technique is helpful for video action classification task not as a neural network compression method but as a transfering of modality knowledge. Zhang *et al.* use KD to train a two-stream network with Motion Vector as the second modality [33]. Stroud *et al.* confirm by constructing Distilled 3D Networks (D3D) [24] that KD from the Optical Flow stream improves the quality of the RGB stream. In addition, the authors of D3D show that KD teaches implicit Optical Flow calculation inside the RGB stream. Motion-Augmented RGB Stream (MARS) of Crasto *et al.* distills the knowledge not from the prediction of the Optical Flow stream but from its feature maps before the global averaging operation [3].

In contrast to the mentioned works, we utilize the idea of ML to train jointly several single-modality networks and improve the quality of each of them. Motivated by BAN, we show that the relaunch of training procedure can further boost the performance of models. Additionally, we show that proper initialization improves our results as well as results for MARS and D3D works.

Note that we target on the single-modality model quality. The improvement of the average predictions of several streams is a different branch of research. An example of an approach that addresses this problem is Gradient-Blending [28]. Nevertheless, we examine the ability of ML to improve the average prediction the multi-modality ensemble. It turns out that proposed initialization with relaunches of single-modality ML provides the best result for the ensemble.

### III. PROPOSED SOLUTION

The proposition of the best single-modality model training pipeline is depicted in Figure 1. The pipeline for the best ensemble training is described in section V.

The pipeline consists of three parts: initialization preparation, ML implantation and Mutual Modality Learning (MML).

The importance of each part is confirmed in section IV.

### A. Initialization preparation

The standard starting point for the Video Action Classification models training is an ImageNet [4] pretrained model. Inflating of 2d-convolutions proposed in [2] makes that possible for both 3d-models and 2d-models.

If we use the input modality different from RGB then we have to change the shape of the first convolution from $(C, 3, K, K)$ to $(C, N, K, K)$. Here, $C$ is a number of output channels of the first convolution, $K$ is a kernel size and $N$ is a number of channels of the new input. The pseudocode for the weights of the new convolution is as follows:

```
for i in 1:N do
    W_new[:,i] = (W[:,1]+W[:,2]+W[:,3])/3
end
```
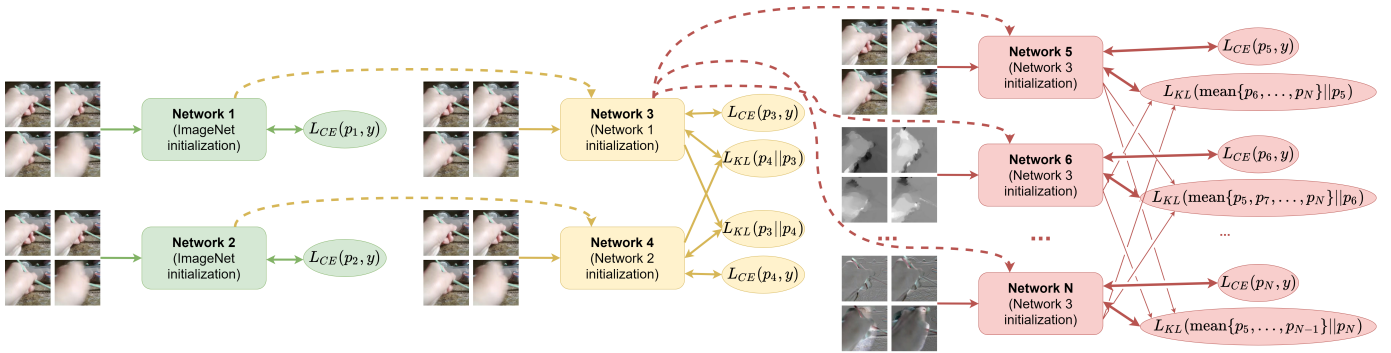
Fig. 1. Best viewed in color. Solid arrows denote flows of data. Dashed arrows denote weights transferring for initialization. Green part: first, we train two networks with RGB input initialized by ImageNet weights using cross-entropy loss. Yellow part: next, we use weights from the first step as initialization for two networks with RGB input that are trained jointly using Mutual Learning. Red part: finally, we apply Mutual Modality Learning to obtain the best single-modality model for each modality. We use weights of the network from the second step as initialization for each model in the third part.

In the proposed pipeline, we use ImageNet initialization only for the first step. The next two steps use weights from the previous step (with a change in the first convolution shape if it is needed).

### B. Mutual Learning implantation

ML is a technique of training two models together in a way that they help each other to reach better convergence. To achieve that, we modify the loss functions of the networks as follows:

$$L_1 = L_{CE}(p_1, y) + L_{KL}(p_2||p_1), \tag{1}$$

$$L_2 = L_{CE}(p_2, y) + L_{KL}(p_1||p_2). \tag{2}$$

Here, $L_i$ is a loss of the $i$-th network, $p_i$ is a vector of the predicted class probabilities by the $i$-th network, $y$ is a ground-true class label, $L_{CE}$ is a cross-entropy loss and $L_{KL}$ is the Kullback Leibler (KL) Divergence loss given by the formula

$$L_{KL}(p_i||p_j) = \sum_{n=1}^{N} p_i^n \cdot \log \frac{p_i^n}{p_j^n}. \tag{3}$$

In this formula $p_i^n$ stands for a probability for the $n$-th class predicted by the $i$-th model. Thus, models teach each other using dependencies that they found during training and thereby improve their performance.

If there are more than two models involved in ML then the loss function is

$$L_i = L_{CE}(p_i, y) + L_{KL}\left(\frac{\sum_{j \neq i} p_j}{M-1}||p_i\right) \tag{4}$$

where $M$ is a number of models.

### C. Mutual Modality Learning

In the original ML, both models use the same modality as an input. We propose to use different modalities of the video obtained from the same frames as inputs for different models. Thus, we share the knowledge obtained from one modality to other modalities.

Note that we need two consecutive frames to calculate the Optical Flow. Thus, if there are $N$ RGB frames in total then there are only $N-1$ Optical Flow frames in total.

So, suppose that the model requires $T$ input-frames for the prediction and we have two representations of the video by different modalities: one representation with $n$ frames and another with $N$ frames ($N > n$).

For this and similar cases in our work, we first sample frames with numbers $(i_1, \ldots, i_T)$ for the modality with the least number of frames, and then we use frames with numbers $(i_1 + \xi, \ldots, i_T + \xi)$ for the modality with the biggest number of frames. Here $\xi \sim \text{unif}\{0, \ldots, N-n\}$.

## IV. ABLATION STUDIES

There are several conclusions that we make:

- Initialization with the RGB model trained on the same video dataset significantly enhances the performance for various modalities and training scenarios (not only ML but MARS and D3D also).
- MML performs better than MARS or D3D approaches.
- Two iterations of ML are better than one and there is no need for more.
- MML performs better than ML as a final step.
- The behavior described above preserves when we use modalities different from the Optical Flow.

### A. Experiments setup

For the ablation studies, we use several models and benchmarks: TSM [20] on Something-Something-v2 [10] with the code provided by the authors (the main setup, we use it unless otherwise specified) and I3D [2] on Charades [21] with the code provided in https://github.com/facebookresearch/SlowFast.

We obtain the Optical Flow using TV-L1 algorithm [32] and combine 5 consecutive evaluations of the Optical Flow by the x- and y- axes as one input-frame.

For the RGBDiff modality, we take 6 consecutive RGB frames to obtain 5 consecutive differences between them. Obtained differences are concatenated and considered as one input-frame.

3

#### TABLE I
#### ONE MODEL TRAINING

| Model | Top-1 / Top-5 | Model | Top-1 / Top-5 |
|---|---|---|---|
| RGB from Ima-geNet | **58.10 / 84.61** | RGB from Flow | 57.53 / 84.42 |
| Flow from Ima-geNet | 52.32 / 81.84 | Flow from RGB | **55.19 / 84.14** |
| Diff from Ima-geNet | 58.74 / 84.39 | Diff from RGB | **58.98 / 86.33** |

#### TABLE II
#### MARS AND D3D TRAINING OF TSM

| Model | Teacher modal-ity | MARS training Top-1 / Top-5 | D3D training Top-1 / Top-5 |
|---|---|---|---|
| RGB from Ima-geNet | Flow | 57.56 / 84.39 | 58.99 / 85.18 |
| RGB from RGB | Flow | **59.11 / 85.24** | **59.95 / 85.86** |
| Flow from Ima-geNet | RGB | 57.46 / 85.01 | 55.04 / 83.36 |
| Flow from RGB | RGB | **58.23 / 85.37** | **56.41 / 83.98** |

#### TABLE III
#### RGB RESULTS OF MML

| RGB results | Flow from Ima-geNet | Flow from Flow | Flow from RGB |
|---|---|---|---|
| RGB from Ima-geNet | 56.25 / 84.07 | 60.02 / 86.08 | 58.70 / 85.18 |
| RGB from RGB | **60.80 / 86.47** | **60.94 / 86.67** | **60.82 / 86.75** |
| RGB from Flow | 58.37 / 84.82 | 58.56 / 85.28 | 58.62 / 85.36 |

#### TABLE IV
#### FLOW RESULTS OF MML

| Flow results | Flow from Ima-geNet | Flow from Flow | Flow from RGB |
|---|---|---|---|
| RGB from Ima-geNet | 54.94 / 83.82 | 57.06 / 84,87 | **57.84 / 85.15** |
| RGB from RGB | 54.76 / 83.61 | 56.74 / 84.78 | **57.95 / 85.44** |
| RGB from Flow | 55.85 / 84.33 | 56.86 / 84.80 | **57.79 / 85.34** |

*1) TSM on Something-Something-v2:* We use the standard setup for the TSM+ResNet-50 [11] training proposed by the authors with batch size 64, ImageNet pretrain, 0.025 initial learning rate. The only difference is the frames sampling strategy. Instead of using one sampling strategy, we use both uniform sampling and dense sampling. The first one works as follows: we split video into $T$ equal parts and take a random frame from each of them. Dense sampling requires taking each $\tau$-th frame starting from a random position. We apply each of the two sampling strategies with $50\%$ probability. See Appendix A as an explanation for this strategy.

We use single uniform sampling with one spatial 224x224 center crop during testing for the ablation studies. That is why the baseline result is worse than the same in [16] where 256x256 central crop is used during testing.

*2) I3D on Charades:* This setup is used to show the advantages of MML regarding other approaches. Both D3D [24] and MARS [3] deal with 3d-models, that is why we use the I3D ResNet-50 model [2] to make a fair comparison with the mentioned methods.

Besides, Charades is the dataset with multiple corresponding classes per one clip, so we show how to extend the proposed MML to the multi-label task.

Optimizer, the number of epochs and other hyperparameters are taken from the standard config-file for the Charades training in https://github.com/facebookresearch/SlowFast without any changes. We use model trained on Kinetics-400 [2] as a standard initialization instead of ImageNet initialization.

### B. Initialization

An abbreviation "Flow from ImageNet" means that we initialize a model that takes Optical Flow as an input with the weights of the model trained on ImageNet. An abbreviation "Diff from RGB" means that we initialize a model that takes differences between RGB frames as an input with the weights of the model with RGB input trained on the current dataset using the cross-entropy loss and initialized by a model trained on ImageNet. We make other abbreviations in a similar way.

We do not include training from scratch into the ablation studies since this is a well-known fact that ImageNet initialization outperforms random initialization for the training of one-stream video models [2], [16].

We can see from Table I that RGB initialization significantly outperforms ImageNet initialization in the case of ordinary cross-entropy training of the Flow and Diff models. At the same time, Flow initialization is useless for RGB models.

We apply MARS [3] and D3D [24] approaches in both directions for RGB and Flow models. Table II shows that RGB initialization improves results in each scenario. It should be noted that both MARS and D3D approaches mainly target 3d-models. That is why the results of MARS training of "RGB from ImageNet" may be worse than the baseline ("RGB from ImageNet" using cross-entropy) since we use 2d-models.

Table III and IV are more representative. First, we train "RGB from ImageNet" and "Flow from ImageNet" models using cross-entropy. Then we train Flow and RGB models together using MML with all possible pairs of the initialization. The results of the RGB models trained using MML are presented in Table III. The results of the Flow models trained using MML are presented in Table IV.

As we can see, the middle values of each column in Table III are the best as well as the right values of each row in Table IV. Thus, the consistency of better initialization is preserved in the case of MML.

Finally, even if we train models on one modality using ML then RGB initialization is still the best. The first three

| First model | Top-1 / Top-5 | Second model | Top-1 / Top-5 |
|---|---|---|---|
| RGB from ImageNet | 57.76 / 84.42 | RGB from ImageNet2 | 58.15 / 84.64 |
| RGB from RGB | 57.84 / 84.55 | RGB from ImageNet | 60.20 / 86.33 |
| RGB from RGB | **60.54 / 86.23** | RGB from RGB2 | 60.47 / 86.08 |
| Flow from ImageNet | 52.94 / 82.21 | Flow from ImageNet2 | 53.44 / 82.50 |
| Flow from RGB | 57.58 / 85.17 | Flow from RGB2 | **57.71 / 85.26** |

| Training pipeline | RGB model mAP | Flow model mAP |
|---|---|---|
| Ordinary training from Kinetics | 33.72 | 15.81 |
| MARS training from Kinetics | 28.74 | |
| MARS training from RGB | 34.40 | |
| D3D training from Kinetics | 33.03 | |
| D3D training from RGB | 35.48 | |
| MML training from Kinetics | 33.84 | 17.34 |
| MML training from RGB | **35.96** | **29.12** |

| | First model | Top-1 / Top-5 | Second model | Top-1 / Top-5 | Tag |
|---|---|---|---|---|---|
| 1 | RGB from RGB | 60.82 / 86.75 | Flow from RGB | 57.95 / 85.44 | **A** |
| 2 | RGB from RGB | 60.88 / 86.86 | Flow from RGB2 | 57,87 / **85.53** | |
| 3 | RGB from **A**(RGB) | 61.18 / 86.81 | Flow from **A**(RGB) | 58.02 / 85.49 | **B** |
| 4 | RGB from **B**(RGB) | 61.15 / 86.81 | Flow from **B**(RGB) | 57.96 / 85.30 | |
| 5 | RGB from RGB | 60.54 / 86.23 | RGB from RGB2 | 60.47 / 86.08 | **C** |
| 6 | RGB from **C**(RGB) | 60.68 / 86.35 | RGB from **C**(RGB2) | 60.88 / 86.44 | |
| 7 | RGB from **C**(RGB) | **61.30 / 86.99** | Flow from **C**(RGB) | **58.36** / 85.49 | |

rows and the next two rows of Table V confirm that. We use abbreviations ImageNet2 and RGB2 to point out that we use different initialization obtained in the same way (KL loss is equal to zero otherwise).

### C. MML versus MARS and D3D

Since MARS [3] and D3D [24] works target mainly 3d-models, we use the I3D on Charades setup in this subsection.

It should be noted that ordinary KL loss implementation uses the "batchmean" regime of averaging, i.e. we divide the sum of losses by the number of instances in one batch. However, we have to use the "mean" regime of averaging when we train the multi-label model using Binary Cross-Entropy losses (BCE), i.e. we divide the sum of losses by the multiplication of two factors: the number of instances in one batch and the number of classes. See Appendix B as an explanation of this point.

By similar reasoning, we divide additional loss functions of MARS and D3D by the number of classes.

The mean Average Precision (mAP) results of all approaches are presented in Table VI. We can see again that RGB initialization improves the performance of each method. D3D is still better than MARS and MML is the best.

The right column in Table VI is empty for MARS and D3D approaches since these approaches do not modify the Optical Flow model during training.

We assume that performance correlates negatively with the strength of the supervision signal. Since we apply KL loss to probits, then any $l_2$ distance between logits is possible during MML. Thus, we weakly bound the feature extraction strategy of a network. In the case of D3D training, we minimize $l_2$ distance between logits only. Thus, D3D does not force a network to estimate the same features in contrast to MARS.

We want to stress that we can significantly improve a single-modality model different from RGB, e.g. MML improves mAP of the Flow model by about 2 times. With some further research these findings may be very helpful for video recognition by event cameras [9].

### D. Relaunch of the ML and MML versus ML

An abbreviation "RGB from **A**(RGB)" in Table VII means that we initialize an RGB model with the weights of the RGB model that was trained by ML tagged as **A**.

Rows number 1 and number 3 from Table VII demonstrate that relaunch of MML can improve the performance. At the same time, row number 4 demonstrates that the second relaunch is probably useless. Rows number 5 and number 6 demonstrate that the consistency is preserved for single-modality ML.

As we can see, the results of the RGB model trained using MML are better than the results of both RGB models trained using ML: row number 1 versus row number 5 from Table VII. This consistency is preserved for the relaunch of ML: row number 6 versus row number 7.

Rows number 1 and number 2 demonstrate that initialization with different RGB weights does not significantly affect the performance of MML.

Finally, row number 7 compared to row number 3 demonstrates that it is better to use MML only as a second step. We believe that the reason for that the separation of advantages of ML itself and additional information from another modality. We extensively examine this effect in section V.

5

TABLE VIII
RGBDIFF MODALITY

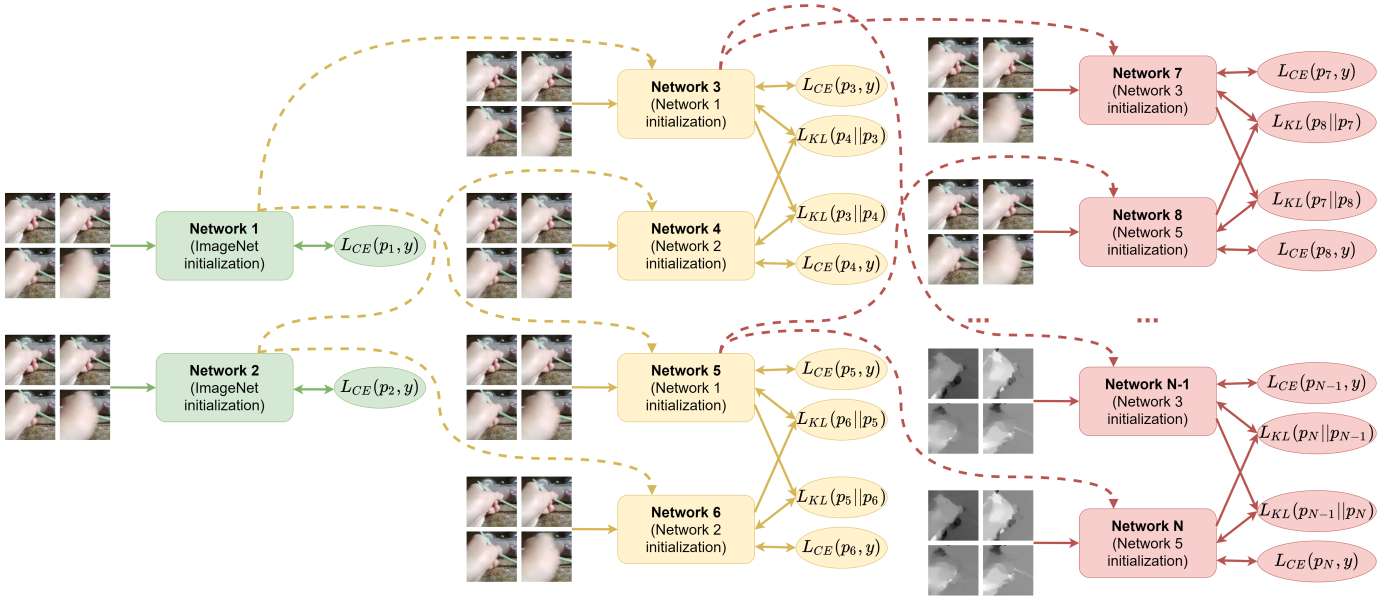| Row number | First model | Top-1 / Top-5 | Second model | Top-1 / Top-5 | Third model | Top-1 / Top-5 |
|---|---|---|---|---|---|---|
| 1 | RGB from RGB | 60.54 / 86.23 | RGB from RGB2 | 60.47 / 86.08 | | |
| 2 | RGB from RGB | 60.82 / **86.75** | Flow from RGB | 57.95 / 85.44 | | |
| 3 | Flow from RGB | 57.58 / 85.17 | Flow from RGB | 57.71 / 85.26 | | |
| 4 | Diff from RGB | 60.66 / 87.65 | Diff from RGB2 | 61.07 / 87.73 | | |
| 5 | RGB from RGB | 60.52 / 86.52 | Diff from RGB | 62.13 / 87.57 | | |
| 6 | RGB from RGB | **61.03** / 86.71 | Flow from RGB | **58.03 / 85.61** | Diff from RGB | **62.51 / 87.95** |



Fig. 2. Best viewed in color. Solid arrows denote flows of data. Dashed arrows denote weights transferring for initialization. Green part: first, we train two networks with RGB input initialized by ImageNet weights using cross-entropy loss. Yellow part: next, we launch RGB-only Mutual Learning for two times. We use the weights from the first step as initialization for each launch of Mutual Learning. We have to use two launches for the second step because we need to obtain two models for which KL loss has not been optimized yet. Red part: finally, we apply single-modality Mutual Learning to each modality that we want to use in the ensemble. We use the weight from one model from each pair from the previous step as the initialization.

### E. Other modalities

We expand our experiments to the Diff modality to examine the preservation of the consistency.

The last row from Table I confirms that RGB initialization is also useful for Diff model.

Row number 5 compared to rows number 4 and number 1 from Table VIII confirms that MML is not worse than or even better than single-modality ML in case of RGB and Diff modalities.

Finally, the comparison of row number 6 to rows 1–5 demonstrates that MML with all three modalities outperforms or is not worse than any other ML in terms of individual results for each modality.

## V. ENSEMBLE PERFORMANCE

The predictions of RGB and Flow models can be highly correlated since we train them using KL loss. Thus, an averaging of the predictions may perform worse than the averaging of ordinary RGB and Flow models trained using cross-entropy. The same logic is applicable to MARS or D3D training.

We show results of ensembles of two models in Appendix C-A and some results of ensembles of three different models with RGB, Flow and Diff input modalities in Appendix C-B.

The main conclusions are as follows:

- RGB models that do not use Optical Flow during training perform the best in ensemble with Flow models. Models trained using ML with RGB only are the first, RGB models trained using MML with RGBDiff are the second.
- RGB models that use Optical Flow during training are the worst in the ensemble with Flow models.
  Performance in the ensemble with Flow models from better to worse: MML, D3D, MARS. We believe that this order is caused by the same reasons that are mentioned in the subsection IV-C.
- The same behavior preserves when we combine Flow models with/without RGB signals in loss function during

TABLE IX
SOTA on Something-Something-v2

| Solution | Ensemble | Base architecture | Number of input frames | Spatial crops × Temporal clips for prediction | Top-1 on validation | Top-5 on validation | Top-1 on test | Top-5 on test |
|---|---|---|---|---|---|---|---|---|
| TSM [16] | No | ResNet-50 | 8 | $1 \times 1$ | 59.1 | 85.6 | – | – |
| TIN [20] | No | ResNet-50 | 8 | $1 \times 1$ | 60.0 | 85.5 | – | – |
| TPN [31] | No | ResNet-50 | 8 | $1 \times 1$ | **62.0** | – | – | – |
| **MML (ours)** | No | ResNet-50 | 8 | $1 \times 1$ | 61.87 | **87.32** | – | – |
| STM [13] | No | ResNet-50 | 8 | $3 \times ?$ | 62.3 | 88.8 | 61.3 | 88.4 |
| W3 [18] | No | ResNet-50 | 16 | $? \times 2$ | **66.5** | **90.4** | – | – |
| STM [13] | No | ResNet-50 | 16 | $3 \times ?$ | 64.2 | 89.8 | 63.5 | 89.6 |
| TPN [31] | No | ResNet-101 | 16 | $3 \times 2$ | – | – | **67.72** | 91.28 |
| **MML (ours)** | No | ResNet-101 | 16 | $1 \times 3$ | 65.9 | 90.15 | 66.83 | **91.30** |
| bLVNet-TAM RGB+Flow [6] | Yes | ResNet-101 | 32+32 | $3 \times 10$ | 68.5 | 91.4 | 67.1 | 91.4 |
| TSM RGB+Flow [16] | Yes | ResNet-50 | 16+16 | $? \times ?$ | 66.0 | 90.5 | 66.55 | 91.25 |
| RGB-only ensemble (9702_10347) by Anonymous | Yes | – | – | $? \times ?$ | – | – | 68.18 | 91.26 |
| TSM ResNet-101, RGB+Flow by Anonymous | Yes | ResNet-101 | – | $? \times ?$ | – | – | 67.71 | 91.95 |
| **ML RGB+Flow (ours)** | Yes | ResNet-101 | 16+16 | $1 \times 3$ | 68.16 | 91.69 | – | – |
| **ML RGB+Flow+Diff (ours)** | Yes | ResNet-101 | 16+16+16 | $1 \times 3$ | **69.07** | **92.07** | **69.02** | **92.70** |

training with RGB models. The only point we want to stress is that "Flow from RGB" models still perform better than "Flow from ImageNet" models in ensembles with RGB models.

- It is also better to combine models trained using single-modality ML when we average the predictions of the RGB and Diff models.
- An ensemble of RGB and Diff models can achieve results that are similar to the results of the RGB and Flow ensemble.
- Models trained using single-modality ML achieve the best results in the ensemble of three different modalities in our experiments. See Appendix C-B for more details.

Thus, although MML provides the best single-modality models, ordinary ML performs better for ensembles. Considering the aforementioned observations, we propose a pipeline for the best ensemble training that is depicted in Figure 2.

First, we train two "RGB from ImageNet" models using cross-entropy. Second, we launch two single-modality ML procedures for the RGB models from the previous step. Finally, we train models using single-modality ML for each of three modalities (RGB, Flow, Diff) that we want to use in the ensemble. We use weights of the RGB models from the second step as an initialization for the third step. This is the reason why we have to launch two training procedures on the second step. KL loss is already optimized otherwise.

## VI. COMPARISON TO STATE-OF-THE-ART

Followed by the observations made above, we train TSM+ResNet-101 with 16 input frames per clip on Something-Something-v2 using the described pipelines. Thus,

we obtain an enhanced RGB model trained by MML and three models with RGB, Flow and Diff inputs for the best ensemble according to section V. The results are available in Table IX.

Our pipeline for the best ensemble achieves the SOTA results among the ones reported previously in the Something-Something-v2 benchmark.

We also make a comparison with other single-model solutions. There is only one single-model solution that outperforms our solution in one of two testing metrics in the Something-Something-v2 benchmark. This is a Temporal Pyramid Network [31] that is several times heavier than TSM and uses more launches per one prediction.

For the simplest scenario, when we use ResNet-50 as a base architecture with 8 input frames and one launch per prediction, we achieve +2.77% improvement of the top-1 performance without adding complexity for the inference.

We exclude STM [13] model from the comparison since it uses the average of predictions for three spatial crops. That is a more accurate but also a more computationally expensive approach.

## VII. CONCLUSION

We present Mutual Modality Learning, the approach that enhances the performance of single-modality model by joint training with models based on other modalities. In addition, we show that the proper initialization of network weights boosts the performance of various training scenarios. We check that our proposal works for different models and datasets, even for multi-label tasks. Our experiments lead to state-of-the-art results in the Something-Something-v2 benchmark.

## REFERENCES

[1] Joao Carreira, Eric Noland, Chloe Hillier, and Andrew Zisserman. A short note on the kinetics-700 human action dataset. *arXiv preprint arXiv:1907.06987*, 2019.

[2] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6299–6308, 2017.

[3] Nieves Crasto, Philippe Weinzaepfel, Karteek Alahari, and Cordelia Schmid. Mars: Motion-augmented rgb stream for action recognition. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

[4] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.

[5] Lijie Fan, Wenbing Huang, Chuang Gan, Stefano Ermon, Boqing Gong, and Junzhou Huang. End-to-end learning of motion representation for video understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6016–6025, 2018.

[6] Quanfu Fan, Chun-Fu Richard Chen, Hilde Kuehne, Marco Pistoia, and David Cox. More is less: Learning efficient video representations by big-little network and depthwise temporal aggregation. In *Advances in Neural Information Processing Systems*, pages 2264–2273, 2019.

[7] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6202–6211, 2019.

[8] Tommaso Furlanello, Zachary C Lipton, Michael Tschannen, Laurent Itti, and Anima Anandkumar. Born again neural networks. *arXiv preprint arXiv:1805.04770*, 2018.

[9] Daniel Gehrig, Mathias Gehrig, Javier Hidalgo-Carrió, and Davide Scaramuzza. Video to events: Recycling video datasets for event cameras. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3586–3595, 2020.

[10] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The "something something" video database for learning and evaluating visual common sense. In *ICCV*, volume 1, page 5, 2017.

[11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[12] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.

[13] Boyuan Jiang, MengMeng Wang, Weihao Gan, Wei Wu, and Junjie Yan. Stm: Spatiotemporal and motion encoding for action recognition. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2000–2009, 2019.

[14] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017.

[15] Hildegard Kuehne, Hueihan Jhuang, Estíbaliz Garrote, Tomaso Poggio, and Thomas Serre. Hmdb: a large video database for human motion recognition. In *2011 International Conference on Computer Vision*, pages 2556–2563. IEEE, 2011.

[16] Ji Lin, Chuang Gan, and Song Han. Tsm: Temporal shift module for efficient video understanding. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7083–7093, 2019.

[17] Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and Josef Sivic. Howto100m: Learning a text-video embedding by watching hundred million narrated video clips. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2630–2640, 2019.

[18] Juan-Manuel Perez-Rua, Brais Martinez, Xiatian Zhu, Antoine Toisoul, Victor Escorcia, and Tao Xiang. Knowing what, where and when to look: Efficient video action modeling with attention. *arXiv preprint arXiv:2004.01278*, 2020.

[19] AJ Piergiovanni and Michael S Ryoo. Representation flow for action recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9945–9953, 2019.

[20] Hao Shao, Shengju Qian, and Yu Liu. Temporal interlacing network. *arXiv preprint arXiv:2001.06499*, 2020.

[21] Gunnar A. Sigurdsson, Gül Varol, Xiaolong Wang, Ali Farhadi, Ivan Laptev, and Abhinav Gupta. Hollywood in homes: Crowdsourcing data collection for activity understanding. In *ECCV*, 2016.

[22] Karen Simonyan and Andrew Zisserman. Two-stream convolutional networks for action recognition in videos. In *Advances in neural information processing systems*, pages 568–576, 2014.

[23] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.

[24] Jonathan Stroud, David Ross, Chen Sun, Jia Deng, and Rahul Sukthankar. D3d: Distilled 3d networks for video action recognition. In *The IEEE Winter Conference on Applications of Computer Vision (WACV)*, March 2020.

[25] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 4489–4497, 2015.

[26] Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A closer look at spatiotemporal convolutions for action recognition. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 6450–6459, 2018.

[27] Limin Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, Dahua Lin, Xiaoou Tang, and Luc Van Gool. Temporal segment networks: Towards good practices for deep action recognition. In *European conference on computer vision*, pages 20–36. Springer, 2016.

[28] Weiyao Wang, Du Tran, and Matt Feiszli. What makes training multi-modal networks hard? *arXiv preprint arXiv:1905.12681*, 2019.

[29] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7794–7803, 2018.

[30] Saining Xie, Chen Sun, Jonathan Huang, Zhuowen Tu, and Kevin Murphy. Rethinking spatiotemporal feature learning: Speed-accuracy trade-offs in video classification. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 305–321, 2018.

[31] Ceyuan Yang, Yinghao Xu, Jianping Shi, Bo Dai, and Bolei Zhou. Temporal pyramid network for action recognition. *arXiv preprint arXiv:2004.03548*, 2020.

[32] Christopher Zach, Thomas Pock, and Horst Bischof. A duality based approach for realtime tv-l 1 optical flow. In *Joint pattern recognition symposium*, pages 214–223. Springer, 2007.

[33] Bowen Zhang, Limin Wang, Zhe Wang, Yu Qiao, and Hanli Wang. Real-time action recognition with enhanced motion vector cnns. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2718–2726, 2016.

[34] Ying Zhang, Tao Xiang, Timothy M Hospedales, and Huchuan Lu. Deep mutual learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4320–4328, 2018.

[35] Mohammadreza Zolfaghari, Kamaljeet Singh, and Thomas Brox. Eco: Efficient convolutional network for online video understanding. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 695–712, 2018.

TABLE X
TESTING WITH DIFFERENT SAMPLING STRATEGIES

| Sampling during training | Dense 0 + Uniform 1 | Dense 0 + Uniform 2 | Dense 1 + Uniform 0 | Dense 2 + Uniform 0 | Dense 1 + Uniform 1 | Dense 2 + Uniform 1 | Dense 1 + Uniform 2 | Dense 2 + Uniform 2 |
|---|---|---|---|---|---|---|---|---|
| Dense sampling | 57.33 / 84.51 | 58.79 / 85.68 | 57.61 / 84.98 | **58.70 / 85.66** | 59.71 / 86.57 | 60.07 / 86.47 | 60.11 / 86.56 | 60.27 / 86.61 |
| Uniform sampling | 59.86 / **86.14** | 61.16 / **87.03** | 56.25 / 83.72 | 56.20 / 84.18 | 60.64 / 85.58 | 59.69 / 86.38 | 61.50 / 87.32 | 61.03 / **87.10** |
| Both samplings | **60.11** / 85.79 | **61.38** / 86.82 | **57.80** / 85.05 | 58.66 / 85.32 | **61.10 / 86.66** | **61.01 / 86.57** | <u>**61.71 / 87.40**</u> | 61.59 / 86.97 |

## APPENDIX A
### SAMPLING STRATEGY

The common procedure for the Something-Something-v2 final testing is an averaging of two predictions for each video. For each prediction, we use central full-resolution crop and uniform sampling: we use frames with numbers $\left\{ \left\lfloor \frac{0 \cdot T}{N} \right\rfloor, \ldots, \left\lfloor \frac{(N-1) \cdot T}{N} \right\rfloor \right\}$ for the first prediction and frames with numbers $\left\{ \left\lfloor \frac{0.5 \cdot T}{N} \right\rfloor, \ldots, \left\lfloor \frac{(N-0.5) \cdot T}{N} \right\rfloor \right\}$ for the second prediction, where $T$ is a total number of frames for current modality and $N$ is the shape of the temporal dimension of the input. Note that uniform sampling, unlike dense sampling, allows any period between input frames and depends on the total length of the video.

We found out that the use of more than two temporal crops with the same sampling strategy or more number of spatial crops insignificantly improves the validation results. At the same time, the use of different sampling strategies during testing significantly improves results regardless of the sampling strategy during training. That is why we incorporate both samplings into training. The median testing results for full-resolution central crops testing are shown in Table X. Label "Dense $k$ + Uniform $m$" means that we use $k + m$ predictions per video using frames with numbers $\left\{ \left\lfloor \frac{i \cdot T'}{k} \right\rfloor, \left\lfloor \frac{i \cdot T'}{k} + \tau \right\rfloor, \ldots, \left\lfloor \frac{i \cdot T'}{k} + \tau \cdot (N-1) \right\rfloor \right\}$, $i \in \{0, \ldots, k-1\}$ when $k > 1$ or frames with numbers $\left\{ \left\lfloor \frac{T'}{2} \right\rfloor, \left\lfloor \frac{T'}{2} + \tau \right\rfloor, \ldots, \left\lfloor \frac{T'}{2} + \tau \cdot (N-1) \right\rfloor \right\}$ when $k = 1$ and frames with numbers $\left\{ \left\lfloor \frac{i/m \cdot T}{N} \right\rfloor, \ldots, \left\lfloor \frac{(N-1+i/m) \cdot T}{N} \right\rfloor \right\}$, $i \in \{0, \ldots, m-1\}$. Here $T$ is a total number of frames for current modality, $N$ is the shape of the temporal dimension of the input, $\tau < \frac{T}{N-1}$ is a dense for the dense sampling and $T' = T - \tau \cdot (N-1)$. Note that there is no random nature in frame numbers during testing.

We make the next conclusions based on Table X:

- Dense sampling training is not suitable for the Something-Something-v2.
- Uniform sampling training and Both samplings training are nearly equal if we use prediction for Uniform sampling.
- Both samplings training outperforms Uniform sampling strategy by up to one percent when tested with both strategies.
- It is better to average predictions for two Uniform samplings and one Dense sampling during testing.

## APPENDIX B
### LOSS MODIFICATION FOR THE BCE TRAINING

The ordinary implementation of the KL loss divides the sum of $B \cdot N$ terms by the $B$, where $B$ is a batch size and $N$ is a number of classes. The reason for that is that ordinary Cross-Entropy loss also divides the sum of $B \cdot N$ terms by the $B$, which can be unobvious:

$$
L_{CE} = \frac{1}{B} \cdot \sum_{b=1}^{B} -\log \frac{e^{l_b^{\mathrm{gt}_b}}}{\sum_{j=1}^{N} e^{l_b^j}} =
$$
$$
= \frac{1}{B} \cdot \sum_{b=1}^{B} \sum_{i=1}^{N} -y_b^i \cdot \log \frac{e^{l_b^i}}{\sum_{j=1}^{N} e^{l_b^j}} = \quad (5)
$$
$$
= \frac{1}{B} \cdot \sum_{b=1}^{B} \sum_{i=1}^{N} -y_b^i \cdot \log p_b^i = \frac{1}{B} \cdot H(y, p).
$$

Here $l_b^i$ is a predicted logit for the class number $i$ for the instance number $b$, $\mathrm{gt}_b$ — ground truth class for the instance number $b$, $y_b^i = I_{\mathrm{gt}_b}(i)$ and $p_b^i = \frac{e^{l_b^i}}{\sum_{j=1}^{N} e^{l_b^j}}$.

So the magnitudes of the CE loss and KL loss are the same. Since the multi-label BCE loss is divided by the $B \cdot N$:

$$
L_{mlBCE} = \quad (6)
$$
$$
\frac{1}{B \cdot N} \cdot \sum_{b=1}^{B} \sum_{i=1}^{N} - \left( y_b^i \cdot \log \sigma(l_b^i) + (1 - y_b^i) \cdot \log \left( 1 - \sigma(l_b^i) \right) \right),
$$

then we divide the KL loss by the $B \cdot N$ to make the magnitudes the same again.

The authors of the MARS and D3D approaches found the best weights for their loss functions in the case of the Cross-Entropy training (50 for MARS and 1 for D3D). Our experiments confirm that additional division of the loss by the number of classes improves the performance of these two methods in the case of multi-label training according to the reasoning made above.

## APPENDIX C
### ENSEMBLES

#### A. Ensembles of two models

Results of the ensembles of RGB and Flow models are depicted in Table XI. Results of the ensembles of RGB and Diff models are depicted in Table XII.

TABLE XI
ENSEMBLE OF RGB AND FLOW MODELS

| RGB / Flow | CE from ImageNet | MARS from RGB | D3D from RGB | MML with Flow from RGB (A) | MML from RGB with Flow from Flow | ML from RGB 1 (B) | ML from RGB 2 | ML from B 1 | ML from B 2 | MML fwith Flow from A | MML with Flow from B | MML with Diff from RGB | MML with Flow and Diff from RGB |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| CE from ImageNet | 63.67 / 88.41 | 60.40 / 86.68 | 61.72 / 87.15 | 62.65 / 88.03 | 62.56 / 88.02 | 64.62 / 89.14 | 64.33 / 88.76 | 64.45 / 88.82 | 64.42 / 88.97 | 62.70 / 88.09 | 62.96 / 88.19 | 63.77 / 88.64 | 63.31 / 88.29 |
| CE from RGB | 63.74 / 88.73 | 61.18 / 87.42 | 62.12 / 87.91 | 62.89 / 88.59 | 62.79 / 88.41 | 64.34 / 89.44 | 64.39 / 89.21 | 64.52 / 89.45 | 64.45 / 89.35 | 62.94 / 88.68 | 63.14 / 88.68 | 64.00 / 88.94 | 63.18 / 88.73 |
| MARS from RGB | 62.26 / 87.81 | 61.61 / 87.24 | 61.92 / 87.66 | 62.37 / 88.25 | 62.50 / 88.19 | 63.57 / 88.84 | 63.42 / 88.57 | 63.62 / 88.75 | 63.58 / 88.71 | 62.78 / 88.42 | 62.88 / 88.49 | 62.98 / 88.44 | 62.62 / 88.37 |
| MARS from ImageNet | 62.57 / 87.67 | 61.28 / 87.21 | 61.92 / 87.71 | 62.55 / 88.08 | 62.61 / 88.02 | 63.73 / 88.80 | 63.61 / 88.62 | 63.76 / 88.76 | 63.70 / 88.81 | 62.92 / 88.26 | 62.96 / 88.38 | 63.28 / 88.47 | 62.72 / 88.31 |
| ML with Flow from ImageNet | 63.41 / 88.44 | 60.73 / 86.72 | 62.24 / 87.35 | 63.13 / 88.26 | 63.03 / 88.25 | 64.88 / 89.23 | 64.67 / 88.93 | 64.64 / 89.13 | 64.57 / 89.33 | 63.32 / 88.43 | 63.64 / 88.46 | 64.31 / 88.79 | 63.56 / 88.38 |
| ML Flow from RGB 1 | 63.97 / 88.93 | 61.91 / 87.78 | 62.78 / 88.31 | 63.74 / 89.01 | 63.60 / 88.91 | 65.19 / 89.77 | 64.98 / 89.45 | 65.20 / 89.74 | 65.06 / 89.72 | 63.74 / 89.09 | 64.08 / 89.20 | 64.66 / 89.22 | 64.16 / 89.16 |
| ML Flow from RGB 2 | 64.23 / 88.99 | 62.04 / 87.82 | 62.93 / 88.31 | 63.89 / 88.93 | 63.66 / 88.94 | 65.24 / 89.72 | 65.11 / 89.48 | 65.08 / 89.78 | 65.09 / 89.77 | 64.02 / 89.19 | 64.21 / 89.16 | 64.83 / 89.41 | 64.28 / 89.09 |
| MML with RGB from RGB | 63.47 / 88.41 | 61.86 / 87.62 | 62.62 / 88.03 | 63.22 / 88.53 | 63.38 / 88.55 | 64.79 / 89.35 | 64.45 / 88.91 | 64.75 / 89.32 | 64.55 / 89.21 | 63.53 / 88.75 | 63.91 / 88.72 | 64.32 / 88.81 | 63.82 / 88.78 |
| MML with RGB from A | 63.78 / 88.65 | 62.08 / 87.69 | 62.89 / 88.21 | 63.57 / 88.83 | 63.60 / 88.71 | 64.98 / 89.47 | 64.71 / 89.24 | 64.87 / 89.48 | 64.96 / 89.40 | 63.57 / 88.88 | 64.04 / 88.91 | 64.35 / 89.16 | 63.85 / 88.99 |
| MML with RGB from B | 63.81 / 88.73 | 62.20 / 87.71 | 62.95 / 88.04 | 63.70 / 88.76 | 63.74 / 88.70 | 64.90 / 89.37 | 64.79 / 89.20 | 64.88 / 89.27 | 64.81 / 89.40 | 63.98 / 88.78 | 64.00 / 88.84 | 64.28 / 88.95 | 64.08 / 88.97 |
| MML with RGB and Diff from RGB | 63.62 / 88.61 | 61.94 / 87.56 | 62.57 / 88.14 | 63.41 / 88.64 | 63.50 / 88.64 | 64.90 / 89.55 | 64.71 / 89.20 | 64.76 / 89.42 | 64.50 / 89.48 | 63.59 / 88.81 | 63.82 / 88.89 | 64.43 / 89.06 | 63.64 / 88.83 |

"MML with Flow from RGB (**A**)" in the first row means that we use the model with RGB input that was jointly trained using Mutual Modality Learning with the model with Optical Flow input using RGB initialization for both models. Tag **A** means that we use the weights of this model as initialization for other models in the table.

"ML from **B** 2" in the first row means that we use the second model with RGB input that was jointly trained using Mutual Learning with the other (the first) model with RGB input. Both models were initialized by the weights obtained by the procedure with tag **B**

We color the cell on the intersection of the column and the row that are marked "CE from ImageNet" in Table XI as white since it is the baseline ensemble. The more intense red color is, the higher the top-1 value for the ensemble is. The more intense light blue color is, the lower the top-1 value for the ensemble is.

The analysis of the tables is in section V.

*B. Ensembles of three models*

We evaluate the validation results for each combination of three models with different input modalities. We sort all the results of the ensembles of three models by the descending order. We show the sum of all indexes of positions for each model in Table XIII. So, the smaller value stands in

TABLE XII
ENSEMBLE OF RGB AND DIFF MODELS

| RGB / Diff | CE from ImageNet | MARS from RGB | D3D from RGB | MML with Flow from RGB (A) | MML from RGB with Flow from Flow | ML from RGB 1 (B) | ML from RGB 2 | ML from B 1 | ML from B 2 | MML fwith Flow from A | MML with Flow from B | MML with Diff from RGB | MML with Flow and Diff from RGB |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| CE from ImageNet | 63.26 / 88.46 | 63.24 / 88.29 | 63.55 / 88.59 | 63.97 / 88.89 | 63.90 / 88.77 | 64.54 / 89.33 | 64.26 / 89.01 | 64.52 / 89.27 | 64.48 / 89.14 | 64.15 / 89.03 | 64.03 / 89.07 | 63.94 / 88.86 | 63.95 / 88.98 |
| CE from RGB | 63.10 / 88.37 | 63.24 / 88.44 | 63.57 / 88.66 | 63.64 / 88.61 | 63.86 / 88.59 | 64.37 / 88.91 | 64.03 / 88.70 | 64.55 / 88.95 | 64.16 / 88.90 | 63.87 / 88.91 | 64.00 / 88.80 | 63.85 / 88.55 | 63.83 / 88.59 |
| ML from RGB 1 | 63.68 / 88.45 | 64.13 / 88.83 | 64.11 / 88.88 | 64.49 / 89.13 | 64.52 / 89.06 | 64.87 / 89.38 | 64.51 / 88.96 | 65.05 / 89.22 | 64.86 / 89.22 | 64.63 / 89.19 | 64.74 / 89.32 | 64.30 / 88.92 | 64.41 / 88.95 |
| ML from RGB 2 | 63.87 / 88.59 | 64.29 / 89.03 | 64.33 / 89.05 | 64.56 / 89.33 | 64.75 / 89.25 | 64.93 / 89.33 | 64.86 / 89.04 | 65.08 / 89.26 | 64.87 / 89.36 | 65.02 / 89.43 | 64.92 / 89.36 | 64.64 / 89.13 | 64.77 / 89.29 |
| MML with RGB from RGB | 62.75 / 88.12 | 63.70 / 88.62 | 63.73 / 88.66 | 63.59 / 88.71 | 63.92 / 88.73 | 64.35 / 89.07 | 63.69 / 88.61 | 64.30 / 88.89 | 63.98 / 88.90 | 64.13 / 88.83 | 64.04 / 89.00 | 63.54 / 88.47 | 63.72 / 88.64 |
| MML with RGB and Flow from RGB | 63.48 / 88.40 | 63.38 / 88.40 | 63.61 / 88.61 | 64.13 / 88.97 | 64.15 / 88.91 | 64.89 / 89.42 | 64.46 / 88.99 | 64.91 / 89.35 | 64.52 / 89.32 | 64.30 / 89.12 | 64.48 / 89.20 | 64.18 / 88.91 | 64.18 / 89.01 |

TABLE XIII
THE RELEVANCE FOR THE ENSEMBLE WITH OTHER MODALITIES

| RGB model | Sum of positions | Flow model | Sum of position | Diff model | Sum of positions |
|---|---|---|---|---|---|
| ML from B 1 | 9486 | ML Flow from RGB 2 | 22434 | ML from RGB 2 | 42079 |
| ML from RGB 1 (B) | 10711 | ML Flow from RGB 1 | 25793 | ML from RGB 1 | 49533 |
| ML from B 2 | 13112 | MML with RGB from A | 28193 | CE from RGB | 67475 |
| ML from RGB 2 | 15670 | MML with RGB from B | 28954 | MML with RGB from RGB | 67528 |
| MML with Diff from RGB | 26918 | ML with Flow from ImageNet | 28972 | MML with RGB and Flow from RGB | 69895 |
| MML with Flow from B | 28069 | CE from ImageNet | 29238 | CE from ImageNet | 71179 |
| MML with Flow and Diff from RGB | 31281 | CE from RGB | 32372 | | |
| CE from ImageNet | 31793 | MML with RGB and Diff from RGB | 32659 | | |
| MML fwith Flow from A | 31968 | MML with RGB from RGB | 34290 | | |
| MML from RGB with Flow from Flow | 35475 | MARS from ImageNet | 50577 | | |
| MML with Flow from RGB (A) | 36167 | MARS from RGB | 54171 | | |
| D3D from RGB | 46044 | | | | |
| MARS from RGB | 50959 | | | | |

the table the better model is in ensemble with two other modalities. Note that the magnitude of sums vary across the input modalities since there are different numbers of models for each modality are tested.