

Modality Mixer for Multi-modal Action Recognition

Sumin Lee Sangmin Woo Yeonju Park Muhammad Adi Nugroho Changick Kim
KAIST

{suminlee94, smwoo95, yeonju29, madin, changick}@kaist.ac.kr

Abstract

In multi-modal action recognition, it is important to consider not only the complementary nature of different modalities but also global action content. In this paper, we propose a novel network, named Modality Mixer (M-Mixer) network, to leverage complementary information across modalities and temporal context of an action for multi-modal action recognition. We also introduce a simple yet effective recurrent unit, called Multi-modal Contextualization Unit (MCU), which is a core component of M-Mixer. Our MCU temporally encodes a sequence of one modality (e.g., RGB) with action content features of other modalities (e.g., depth, IR). This process encourages M-Mixer to exploit global action content and also to supplement complementary information of other modalities. As a result, our proposed method outperforms state-of-the-art methods on NTU RGB+D 60, NTU RGB+D 120, and NW-UCLA datasets. Moreover, we demonstrate the effectiveness of M-Mixer by conducting comprehensive ablation studies.

1. Introduction

Humans experience their surroundings through a combination of various modality data, such as audio, sight, and touch. According to the recent advancements in sensor technology, multi-modal learning has attracted much research interest in the field of computer vision. Toward this direction, for video action recognition, many multi-modal methods have been developed, which achieved higher performance than other methods based on a single modality.

Earlier studies on action recognition mostly relied on a single RGB modality [4, 11, 34, 40], which are focused on spatio-temporal modeling. Lately, many models based on multi-modality have been developed to integrate information of different modalities such as RGB, optical flow, and depth [8, 9, 13, 14, 15, 21, 25, 29, 35]. Due to the different properties of sensors, each modality possesses different key characteristics that contribute to the overall action recognition. Specifically, while RGB images provide visual appearances, depth data contains the 3D structure of

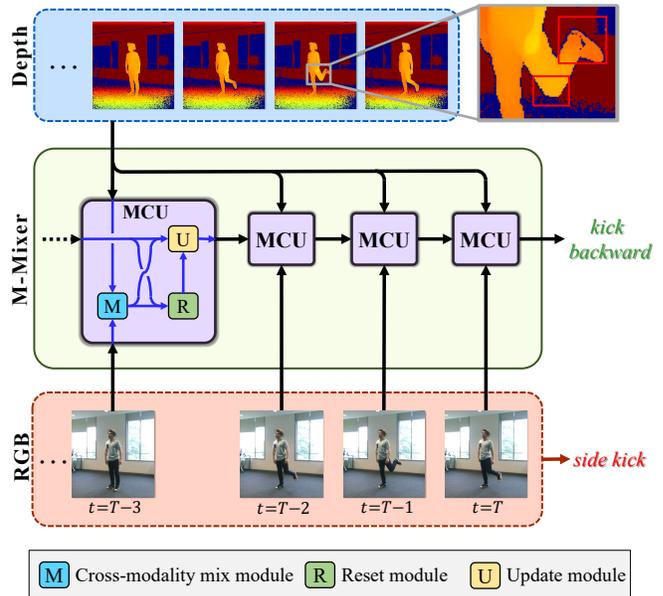


Figure 1. **Multi-modal Action Recognition with Modality Mixer (M-Mixer) network.** Solely relying on appearance data of the RGB modality, the action ‘kick backward’ is easily misclassified to ‘side kick’. On the other hand, our MCU supplements the information on foot orientation from the depth image: the foot is going behind the knee (the more yellowish the color is, the closer it is to the camera). As such, M-Mixer correctly classifies the action as ‘kick backward.’ Overall, MCU performs temporal encoding of the RGB sequence, while augmenting it with complementary information of depth content feature. Here, we assume to use RGB and depth input and depict only RGB stream for the sake of clarity.

2D frames, which is complementary to RGB modality. For example, as illustrated in Fig. 1, the action ‘kick backward’ may be incorrectly predicted as ‘side kick’ when employing only the RGB modality [26]. However, the depth data can indicate that the foot is going behind the knee, leading to the correct action class ‘kick backward’. Therefore, in order to distinguish a correct action class, it is necessary to integrate complementary information from multi-modal data as well as temporal encoding of given videos.

In this paper, we propose a novel network, Modality

Mixer (M-Mixer), which leverages two important factors for multi-modal action recognition: 1) complementary information across modalities and 2) temporal context of action. Taking feature sequences of multiple modalities as inputs, our M-Mixer temporally encodes each feature sequence with action content features of other modalities. The action content features include modality-specific information and the overall activity of videos. We also introduce a simple yet effective recurrent unit, called Multi-modal Contextualization Unit (MCU), which adaptively integrates a modality sequence and action content features. Our M-Mixer network employs a distinct MCU for each modality. As each MCU is dedicated to a specific modality, we describe our MCU in detail from an RGB perspective, as illustrated in Fig. 1. Our MCU consists of three modules: cross-modality mix module, reset module, and update module. Concretely, given an RGB feature at certain timestep and context features of other modalities, our cross-modality mix module models their relationship and adaptively integrates them by weighted summation. By doing so, MCU enables the network to exploit complementary information across modalities and global action content during temporal encoding. Then, reset and update modules learn the relationships between the integrated feature of the current timestep and the hidden state of the previous timestep. Based on MCU, our M-Mixer network assimilates more richer and discriminative information from multi-modal sequences for action recognition. Note that, our M-Mixer is not limited to only two modalities and is applicable to more modalities.

We perform extensive experiments on three benchmark datasets (*i.e.*, NTU RGB+D 60 [28], NTU RGB+D 120 [26], and Northwestern-UCLA (NW-UCLA) [36]). Our M-Mixer network achieves the state-of-the-art performance of 90.77%, 90.12%, and 94.43% on NTU RGB+D 60, NTU RGB+D 120, and NW-UCLA datasets, respectively. Also, we empirically show that our M-Mixer model can be extended to more than two modalities. Through extensive ablation experiments, we demonstrate the effectiveness of the proposed M-Mixer network.

Our main contributions are summarized as follows:

- We investigate how to take two important factors into account for multi-modal action recognition: 1) complementary information across modality, and 2) temporal context of an action.
- We introduce a novel network, named M-Mixer, with a new recurrent unit, called MCU. By effectively modeling the relation between a sequence of one modality and action contents of other modalities, our MCU facilitates M-Mixer to exploit rich and discriminative features.
- We evaluate the performance of multi-modal action recognition on three benchmark datasets. Experiments

show that M-Mixer outperforms the state-of-the-art methods. Moreover, we demonstrate the effectiveness of the proposed method by conducting comprehensive ablation studies.

2. Related Work

Video action recognition is one of the most representative tasks in the field of video understanding. Thanks to the emergence of deep learning, video action recognition has made significant progress over the last decade. Early deep learning models [12, 31, 37] were developed with a two-stream structure in which each stream captures the appearance and motion data of videos, respectively. Due to the high cost of computing accurate optical flow, other works [6, 24, 27, 32, 33, 41] studied to learn to mimic motion features from only RGB sequences. Stroud *et al.* [32] and Crasto *et al.* [6] suggested learning algorithms that distill knowledge from the temporal stream to the spatial stream in order to reduce the two-stream architecture into a single-stream model. Other studies [24, 27, 41] introduced modules to explore motion information in a unified network. After then, 3D convolution networks [4, 11, 34, 40] were proposed for action recognition, which led to significant performance improvements. Among them, SlowFast network [11] consists of two pathways for dealing with two different frame rates to capture spatial semantics and motion.

Due to the advance of sensor technologies, action recognition in multi-modal setting has attracted research interest [1, 2, 3, 7, 10, 16, 20, 39] RGB and depth are one of the most common combinations of modalities [8, 19, 21, 42]. Shahroudy *et al.* [29] introduced a shared-specific feature factorization network based on autoencoder structure for RGB and depth inputs. Liu *et al.* [25] presented a method of learning action features that are insensitive to camera viewpoint variation. Wang [38] proposed a cooperatively trained Convolutional neural Network (c-ConvNet), which enhances the discriminative information of RGB and depth modalities. In [9], a two-stream view-invariant framework was proposed with motion stream and Spatial-Temporal Dynamic (STD) stream with RGB and depth, where late fusion technique is employed to combine outputs of these RGB and depth streams. In [14, 15], frameworks of distillation and privileged information were suggested. Although these methods are trained with both RGB and depth data, a hallucination network of depth enables classifying actions with only RGB data. Garcia *et al.* [13] introduced an ensemble of three specialist networks, called the Distillation Multiple Choice Learning (DMCL) network, that works with missing modalities at inference time. DMCL includes each specialist network for RGB, depth, and optical flow videos that collaborates and strengthens each other and employ the late fusion scheme. Wang *et al.* [35] pro-

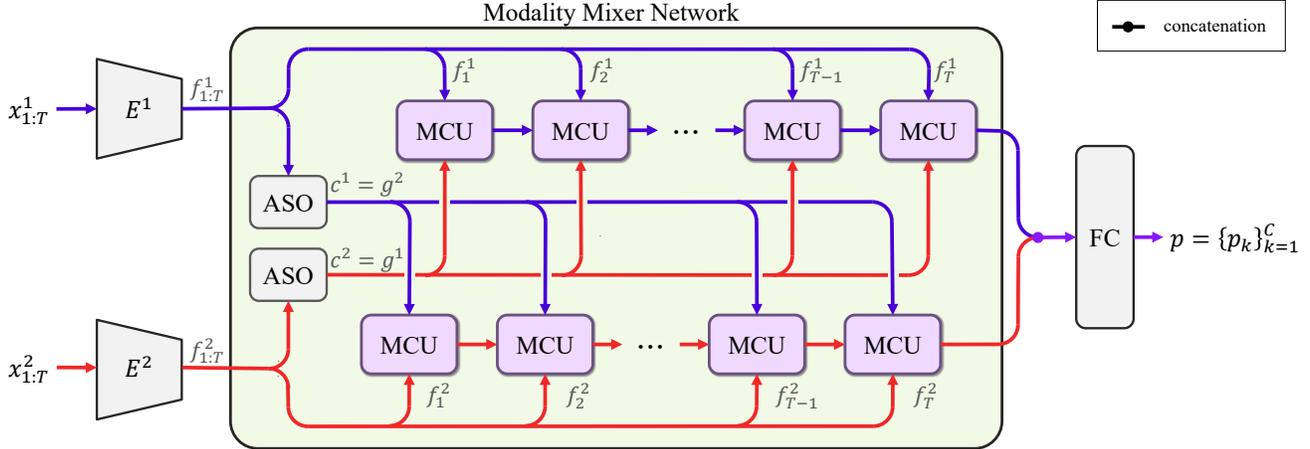


Figure 2. **Modality Mixer network.** We illustrate an example of using two modalities in this figure. M-Mixer network consists of Multi-modal Contextualization Unit (MCU) and Action Summarizing Operator (ASO). Our M-Mixer takes a feature sequence $f_{1:T}^i$, derived from a feature extractor E^i by a frame sequence $x_{1:T}^i$, as an input. First, ASO computes an action content feature c^i from $f_{1:T}^i$. Then, our MCU temporally encodes $f_{1:T}^i$ with a cross-modality action content feature g^i , which is c^j in this case, where $j \neq i$. By comparing $f_{1:T}^i$ with g^i during temporal encoding, MCU considers complementary information across modalities and overall action contents of videos. The final probability distribution over K action classes is calculated by using h_T^1 and h_T^2 . Here, the blue and red lines indicate streams of modality 1 and 2, respectively, and the purple line represents the fusion of modalities.

posed a hybrid network based on CNN (*e.g.*, ResNet50 [17] and 3D convolution) and RNN (*e.g.*, ConvLSTM [30]) to fuse RGB, depth, and optical flow modalities.

In this paper, we explore how to fuse multi-modal data with recurrent units. To this end, we propose a novel recurrent unit, Multi-modal Contextualization Unit (MCU), and network, Multi-modal Mixer (M-Mixer), for multi-modal action recognition. By encoding a feature sequence with content features of other modalities, our M-Mixer network facilitates to exploring complementary information across modalities and temporal action contents. Note that our M-Mixer network is not limited to types and the number of video modalities.

3. Proposed Method

In this section, we first describe the overall architecture of our proposed Modality Mixer (M-Mixer) network and then explain the proposed Modality Contextualization Unit (MCU) in detail. In Fig. 2, the framework of our M-Mixer network is illustrated, assuming the use of two modalities.

3.1. Modality Mixer Network

The goal of our M-Mixer network is to generate rich and discriminative features for action recognition from videos of N different modalities. Given a video of length T for the i -th modality, a feature extractor E^i converts a sequence of frames, $x_{1:T}^i \in \mathbb{R}^{3 \times T \times H \times W}$, to a sequence of features, $f_{1:T}^i \in \mathbb{R}^{d_f \times T}$, as follows:

$$f_{1:T}^i = E^i(x_{1:T}^i), \quad (1)$$

where H and W denote the height and width of a video, and $i = 1, 2, \dots, N$. Then, the proposed M-Mixer network takes the extracted feature sequences $f_{1:T}^i$ as inputs.

In our M-Mixer, firstly, Action Summarizing Operator (ASO) condenses an action content information $c^i \in \mathbb{R}^{d_c}$ from $f_{1:T}^i$, as follows:

$$c^i = \text{ASO}(f_{1:T}^i). \quad (2)$$

ASO can be any operation that can concentrate action content information from a feature sequence into a vector, such as max pooling, average pooling, and recurrent units. Among several instantiations, we empirically find that average pooling performs the best for ASO.

Our M-Mixer contains N MCUs that are responsible for each modality. Each MCU encodes a feature sequence of a designated modality with content features of other modalities in a temporal manner. In other words, an MCU for the i -th modality contextualizes $f_{1:T}^i$ with the j -th action content for all $j \in \{1, \dots, N\}$, where $j \neq i$. For these action contents, we define a cross-modality action content feature $g^i \in \mathbb{R}^{(N-1) \times d_c}$, as follows:

$$g^i = \left\| \left\| c^j, \text{ where } j \neq i. \right. \right\|_{\forall j} \quad (3)$$

Here, $\|$ indicates a vector concatenation. If only two modalities are used, g^i is equal to c^j . Then, MCU takes f_t^i and g^i to generate hidden state $h_t^i \in \mathbb{R}^{d_h}$ as follows:

$$h_t^i = \text{MCU}^i(f_t^i, g^i), \quad (4)$$

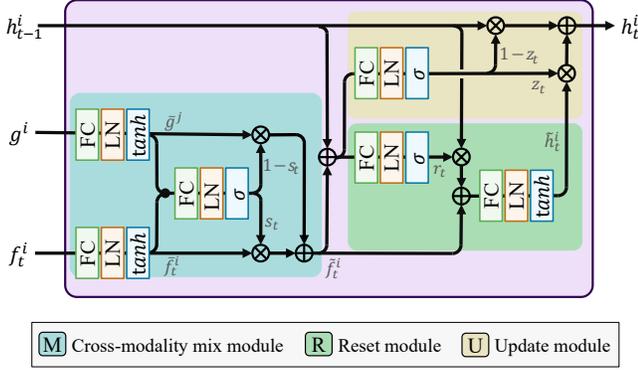


Figure 3. **Multi-modal Contextualization Unit (MCU)**. Our MCU consists of three modules: cross-modality mix module, reset module, and update module. In a cross-modality mix module, a cross-action content g^i is adaptively integrated with f_t^i , providing complementary information and the overall action content. Our reset module decides a reset gate r_t to effectively drop and take information of previous hidden state h_{t-1}^i and an integrated feature \tilde{f}_t^i . In an update module, an update gate z_t is computed to update previous hidden state h_{t-1}^i .

where MCU^i denotes an MCU for the i -th modality. By augmenting a cross-modality action content feature g^i , our MCU exploits complementary information as well as global action content,

To obtain final probability distribution $p = \{p_k\}_{k=1}^K$ over K action classes, we employ a fully connected layer to h_T^i for $i \in \mathcal{V}$, as follows:

$$p = \text{softmax} \left(\mathbf{W}_p \left(\left\| h_T^i \right\|_{i \in \mathcal{V}} \right) + b_p \right), \quad (5)$$

where p_k is a probability of the k -th action class, $\mathbf{W}_p \in \mathbb{R}^{(N \times d_h) \times K}$ is a learnable matrix, and $b_p \in \mathbb{R}^K$ is a bias term.

To train our M-Mixer network, we define a loss function L by utilizing the standard cross-entropy loss as follows:

$$L = \sum_{k=1}^K y_k \log(p_k), \quad (6)$$

where y_k is the ground-truth label for the k -th action class.

3.2. Multi-modal Contextualization Unit

We describe our new recurrent unit, MCU, which is the core component of the proposed M-Mixer network. Our MCU consists of three submodules: cross-modality mix module, reset module, and update module. At the t -th timestep, the proposed MCU takes f_t^i and g^i to contextualize a modality-specific feature with a cross-modality action content feature. This strategy enables MCU to supplement with complementary information of other modalities

in terms of global action content. As a result, the proposed MCU exploits rich and well-contextualized features for action recognition.

Cross-modality Mix Module. First, f_t^i and g^i are projected to the same embedding space, as follows:

$$\tilde{f}_t^i = \tanh(\text{LN}(\mathbf{W}_f f_t^i)), \quad (7)$$

$$\tilde{g}^i = \tanh(\text{LN}(\mathbf{W}_g g^i)), \quad (8)$$

where $\mathbf{W}_g \in \mathbb{R}^{((N-1) \times d_c) \times d_h}$ and $\mathbf{W}_f \in \mathbb{R}^{d_f \times d_h}$ are trainable matrices, and LN and tanh denote the layer normalization and the hyperbolic tangent function, respectively. Note that we exclude a bias term for simplicity.

Next, an integration score s_t is computed to determine how much representations of target modality and other modalities are activated, as follows:

$$s_t = \sigma(\text{LN}(\mathbf{W}_s [\tilde{f}_t^i \parallel \tilde{g}^j])), \quad (9)$$

where σ indicates the sigmoid function and $\mathbf{W}_s \in \mathbb{R}^{2d_h \times d_h}$ is a weight matrix. Then, \tilde{f}_t^i and \tilde{g}^i are combined to the supplemented feature \tilde{f}_t^i , as follows:

$$\tilde{f}_t^i = s_t \otimes \tilde{f}_t^i + (1 - s_t) \otimes \tilde{g}^j, \quad (10)$$

where \otimes denotes the element-wise multiplication.

Reset and Update Module. Our reset and update modules learn relationships between the supplemented feature \tilde{f}_t^i and previous hidden state h_{t-1}^i . In a reset module, a reset gate r_t effectively drops and takes information of h_{t-1}^i and \tilde{f}_t^i . And an update module measures an update gate z_t to amend previous hidden state h_{t-1}^i to current hidden state h_t^i . We compute r_t and z_t , as follows:

$$r_t = \sigma(\text{LN}(\mathbf{W}_r (\tilde{f}_t^i + h_{t-1}^i))), \quad (11)$$

$$z_t = \sigma(\text{LN}(\mathbf{W}_z (\tilde{f}_t^i + h_{t-1}^i))), \quad (12)$$

where $\mathbf{W}_r \in \mathbb{R}^{d_h \times d_h}$ and $\mathbf{W}_z \in \mathbb{R}^{d_h \times d_h}$ are learnable parameters. Then, the hidden state h_{t-1}^i is updated with z_t , as follows:

$$h_t^i = z_t \otimes \tilde{h}_t^i + (1 - z_t) \otimes h_{t-1}^i, \quad (13)$$

where \tilde{h} is defined as:

$$\tilde{h}_t^i = \tanh(\text{LN}(\mathbf{W}_h (r_t \otimes h_{t-1}^i + \tilde{f}_t^i))). \quad (14)$$

Here, $\mathbf{W}_h \in \mathbb{R}^{d_h \times d_h}$ is a trainable matrix.

4. Experiments

4.1. Dataset

NTU RGB+D 60. NTU RGB+D 60 [28] is a large-scale human action recognition dataset, consisting of 56,880 videos. It includes 40 subjects performing 60 action classes in 80 different viewpoints. As suggested in [28], we follow the cross-subject evaluation protocol. For this evaluation, this dataset is split into 40,320 samples for training and 16,560 samples for testing.

NTU RGB+D 120. As an extended version of NTU RGB+D 60, NTU RGB+D 120 [26] is one of the large-scale multi-modal dataset for video action recognition. It contains 114,480 video clips of 106 subjects performing 120 classes from 155 different viewpoints. We follow the cross-subject evaluation protocol as proposed in [26]. For the cross-subject evaluation, the 106 subjects are divided into 53 subjects for training and the remaining 53 subjects for testing.

Northwestern-UCLA (NW-UCLA). NW-UCLA [36] is composed of 1475 video clips with 10 subjects performing 10 actions. Each scenario is captured by three Kinect cameras at the same time from three different viewpoints. As suggested in [36], we follow the cross-view evaluation protocol, using two views for training and the other one for testing.

4.2. Implementation Details

For the feature extractor for each modality, we use ResNet-18 [17] for NTU RGB+D 60 and NW-UCLA and ResNet-34 for NTU RGB+D 120, which are initialized with ImageNet [22] pretrained weights. We set the size of the hidden dimension in MCU, d_h , to 512. The input of each modality is a video clip uniformly sampled with temporal stride 8. For the training procedure, we adopt random cropping and resize each frame to 224×224 . We also apply random horizontal flipping and random color jittering for RGB videos. For depth and IR frames, we use the same method to [15], a jet colormap, to convert those to color images.

To train our M-Mixer network, we use 4 GPUs of RTX 3090. We use the Adam [23] optimizer with the initial learning rate of 10^{-4} . A batch size per GPU is 8 on NTU RGB+D 60 and NTU RGB+D 120. Due to the small number of training samples, we use a single GPU with batch size 8 on NW-UCLA.

4.3. Ablation Study

In this section, we conduct extensive experiments to show the effectiveness of the proposed M-Mixer network and MCU. All experiments in this section are conducted on NTU RGB+D 60 [28] with RGB and depth modalities.

ASO	Accuracy(%)
GRU [5]	90.66
Max pooling	90.65
Average pooling	90.77

Table 1. **Ablation study of Action Summarizing Operator (ASO).** We use three instantiations of ASO: GRU, max pooling, and average pooling.

Method	Accuracy (%)
LSTM [18]	84.58
GRU [5]	84.87
MCU	90.77

Table 2. **Comparisons with LSTM, GRU, and MCU.** The best score is marked in **bold**.

4.3.1 Action Summarizing Operator (ASO)

In Table 1, we test three instantiations of Action Summarizing Operator (ASO). For the average and max pooling, we apply mean and max operation across the temporal axis, respectively. Also, we employ a GRU [5] as ASO, which encodes a sequence into a vector on the time axis. Our M-Mixer network achieves 90.66% with GRU, 90.65% with max pooling, and 90.77% with average pooling. We observe that there are very small performance differences between the three operations. Among them, we empirically find that average pooling works slightly better than others. Therefore, we utilize the average pooling as ASO in all experiments of this paper.

4.3.2 Comparison with RNNs

To solely see the effectiveness of MCU, we replace our MCU in M-Mixer with LSTM [18] or GRU [5]. For these experiments, we use one LSTM or GRU for each modality. We input a feature sequence $f_{1:T}^i$ to each LSTM or GRU and do not use the cross-modality action content. The final predictions are calculated with concatenated output features of each LSTM or GRU, as same as our M-Mixer. Table 2 presents the performances of three networks. We analyze the effect of the action content feature by comparing the results of LSTM, GRU, and MCU. Compared to LSTM and GRU, our MCU learns relations between a current feature and the cross-modality action content to explore discriminative action information. Thanks to the global video content and complementary information, the proposed MCU achieves performance gains of 6.19% and 5.90% over LSTM and GRU, respectively.

Exp.	Cross-modality mix module	Cross-modality action content	LN	Acc. (%)
I	✓			86.82
II		✓		87.76
III	✓	✓		89.97
IV	✓	✓	✓	90.77

Table 3. **Ablation study of MCU.** MCU has three important components: a cross-modality mix module, cross-modality action content, and layer normalization indicated LN. The best scores are marked in **bold**.

Modality	Accuracy(%)		$\Delta(\%p)$
	MCU-self	MCU	
RGB	56.61	79.42	+ 22.81
Depth	84.31	88.59	+ 4.25
RGB+Depth	88.17	90.77	+ 2.60

Table 4. **Experiments on the effectiveness of the cross-modality action content.** For MCU-self, we use c^i instead of g^i to MCU. A single modality represents the performance of each modality stream. Δ indicates performance differences between MCU-self and MCU. The best scores are marked in **bold**.

4.3.3 Modality Contextualization Unit (MCU)

In this section, we validate the effects of three important components of our MCU: a cross-modality mix-module, the cross-modality action content, and the layer normalization. To observe the abilities of each model component, we conduct ablation experiments on these three components and report the performances in Table 3. In experiment I, we replace the cross-modality action content g^i to the self-modality action content c^i and turn off the layer normalization. Experiment II is conducted to investigate the effect of the cross-modality mix module, where we change it to simple concatenation and disable the layer normalization. Lastly, in experiment III, we only turn off the layer normalization of our MCU. Comparing the result of experiment I, utilizing cross-modality action content achieves 89.97% in experiment III, which is 3.15%p higher than experiment I. By comparing experiment II and III, we observe that using cross-modality mix module improves the performance from 87.76% to 89.97%. Finally, we obtain the best performance of 90.77% with all three components in experiment IV.

4.3.4 Cross-modality Action Content in MCU

To validate the effectiveness of the cross-modality action content, we strategically replace the cross-modality action content of MCU (see Eq. 3) to the self-modality action content (*i.e.*, c^i). Since the self-modality action content c^i is from the same modality as a sequence $f_{1:T}^i$ to be encoded,

the self-modality action content comprises global action information. On the other hand, the cross-modality action content g^i includes not only global action information but also complementary information of other modalities. We name MCU with the self-modality action content MCU-self.

Furthermore, in order to closely analyze the effect of the cross-modality action content, we evaluate action recognition performances of a single modality in our M-Mixer with MCU and MCU-self. In other words, we test the performance of RGB and depth features for employing MCU and MCU-self. To this end, we train two additional fully-connected layers to classify an action class with h_T^i , as follows:

$$p^i = \text{softmax}(\mathbf{W}_{p^i} h_T^i), \quad (15)$$

where p^i is a probability distribution of i -th modality, and $\mathbf{W}_{p^i} \in \mathbb{R}^{d_h \times K}$ is a learnable matrix. To train two classifiers, we use a loss function L_h with the standard cross-entropy loss, as follows:

$$L_h = \sum_{i=1}^2 \sum_{k=1}^K y_k \log(p_k^i), \quad (16)$$

where p_k^i is a probability of k -th action class for i -th modality. Note that the whole weights of M-Mixer network are fixed during training of the classifiers.

Comparative analysis in respect of modalities. Table 4 presents the results of comparative experiments about MCU and MCU-self. With RGB and depth modalities, MCU-self obtains 88.17%. Meanwhile, our MCU achieves 90.77%, which is 2.6%p higher than the performance of MCU-self. These results demonstrate the effectiveness of the cross-modality action content.

Compared to the self-modality action content, the cross-modality action content contains complementary information of other modalities as well as global action content. Specifically, the RGB feature is strengthened with depth information, and the depth feature is augmented by RGB information in the setting of this experiment. As a result, our MCU achieves 79.42% in RGB stream and 88.59% in depth stream, which are 22.81%p and 4.28%p higher than RGB and depth streams of MCU-self, respectively. From these results, we demonstrate that the cross-modality action content effectively provides additional information across modalities and our MCU successfully utilizes complementary information in temporal encoding.

Comparison of class-wise performance. In Fig. 4, we report class-wise performance of the proposed M-Mixer and M-Mixer with MCU-self. 60 action classes on NTU

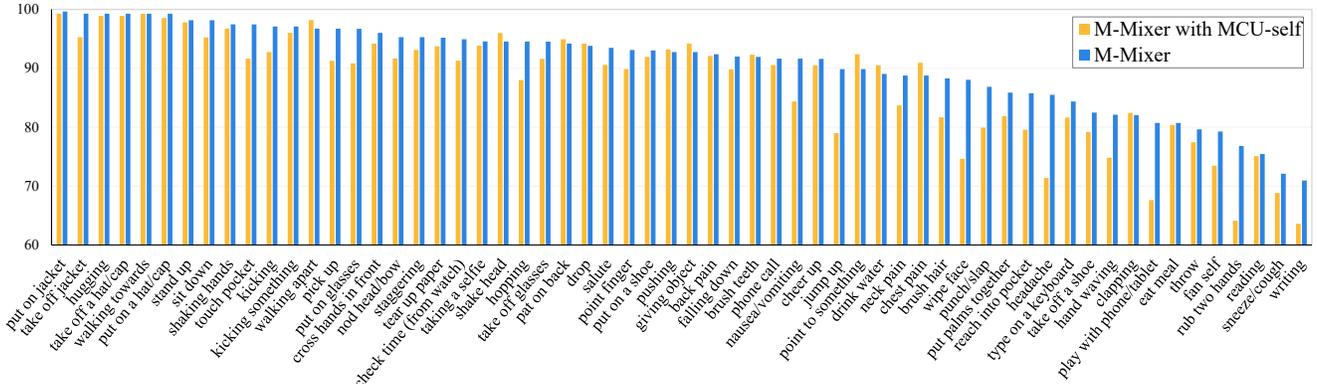


Figure 4. **Class-wise performance of M-Mixer network on NTU RGB+D 60 [28]**. 60 action classes are listed in descending order according to the performance of our M-Mixer.

Method	Modality	Accuracy(%)
Sharoudy <i>et al.</i> [29]	RGB + Depth	74.86
Liu <i>et al.</i> [25]	RGB + Depth	77.5
ADMD [15]	RGB + Depth	77.74
Dhiman <i>et al.</i> [9]	RGB + Depth	79.4
Garcia <i>et al.</i> [14]	RGB + Depth	79.73
c-ConvNet [38]	RGB + Depth	86.42
DMCL [13]	RGB + Depth + F	87.25
Wang <i>et al.</i> [35]	RGB + Depth + F	89.51
Ours (R18)	RGB + Depth	90.77
Ours (R18)	RGB + Depth + IR	91.13

Table 5. **Performance Comparison on NTU RGB+D 60 [28]**. ‘F’ denotes the optical flow. ‘R18’ indicates our M-Mixer with ResNet18 [17] feature extractor. The best scores are marked in **bold**.

RGB+D 60 [28] are depicted in descending order based on the performances of our M-Mixer. In most of the action classes, our M-Mixer achieves higher performances than using MCU-self. Especially, the proposed M-Mixer has significant performance improvements in the action classes, which is lower in using MCU-self (*e.g.*, ‘rub two hands’, ‘headache’, and ‘writing’).

4.4. Comparisons with state-of-the-arts

We compare our M-Mixer network with state-of-the-art methods on NTU RGB+D 60 [28], NTU RGB+D 120 [26], and NW-UCLA [36] for multi-modal action recognition.

NTU RGB+D 60. In Table 5, we report performances of our M-Mixer and state-of-the-art approaches. With RGB and depth modality, our M-Mixer achieves the best performance of 90.77%, which is 1.18%p and 3.44%p higher than DMCL [13] and the method proposed by Wang *et al.* [35],

Method	Modality	Accuracy(%)
VGG [26]	RGB+Depth	61.9
DMCL [13]	RGB + Depth + Flow	89.74
Ours (R34)	RGB+Depth	90.12

Table 6. **Performance Comparison on NTU RGB+D 120 [26]**. ‘R34’ indicates our M-Mixer with ResNet34 [17] feature extractor. The best scores are marked in **bold**.

Method	Modality	Accuracy(%)
Garcia <i>et al.</i> [14]	RGB + Depth	88.87
ADMD [15]	RGB + Depth	89.93
Dhiman <i>et al.</i> [9]	RGB+Depth	84.58
DMCL [13]	RGB + Depth + Flow	93.79
Ours (R18)	RGB+Depth	94.43

Table 7. **Performance Comparison on NW-UCLA [36]**. ‘R18’ indicates our M-Mixer with ResNet18 [17] feature extractor. The best scores are marked in **bold**.

respectively. Note that those two methods use additional information of optical flow. With RGB, depth, and IR videos, our M-Mixer obtains 91.13%. This result shows that the proposed M-Mixer can be extended to more than two modalities.

NTU RGB+D 120. Table 6 shows performance comparisons on NTU RGB+D 120. While NTU RGB+D 120 contains twice as many samples and classes as NTU RGB+D 60, our M-Mixer still obtains the state-of-the-art performance of 90.12%. Compared to the VGG architecture proposed in [26], our M-Mixer attains 28.2% higher performance. Also, M-Mixer achieves 0.38% higher performance than DMCL.

Videos	M-Mixer with MCU-self		M-Mixer		
	RGB	Depth	RGB	Depth	RGB+Depth
	Taking a selfie (20.77)	Brush teeth (43.48)	Taking a selfie (83.99)	Tear up paper (24.72)	Taking a selfie (73.75)
	Eat meal (52.25)	Sneeze/cough (49.88)	Wipe face (92.96)	Sneeze/cough (63.23)	Wipe face (88.42)
	Reading (48.21)	Writing (72.52)	Writing (85.80)	Writing (89.54)	Writing (95.69)
	Brush hair (77.46)	Brush hair (50.13)	Brush hair (86.05)	Brush hair (98.55)	Brush hair (99.92)

Figure 5. **Qualitative evaluation of M-Mixer network on NTU RGB+D 60 [28]**. Predicted results consistent with ground-truth are colored in green, otherwise in red. RGB, Depth, and RGB+Depth indicate prediction results from its respective stream. Also, confidence scores of predictions are presented in parentheses.

NW-UCLA. In Table 7, we summarize the results on NW-UCLA. Our M-Mixer network surpasses the state-of-the-art methods by achieving 94.43%. This performance is 0.64% higher than DMCL that utilizes three modalities, and 9.85% higher than the method proposed by Dhiman *et al.* [9] with RGB and depth. This demonstrates the superiority of our M-Mixer network in discriminating actions.

4.5. Qualitative Evaluation

Figure 5 shows the prediction results of M-Mixer on the sample videos of the NTU RGB+D 60 [28]. To clearly see the efficacy of cross-modality action content, we also report the prediction results of M-Mixer with MCU-self. In addition, we present the confidence score of each prediction under the predicted label. We observe that our M-Mixer significantly improves the prediction results of both RGB and depth streams in comparison to using MCU-self. For example, in the second row of 5, while RGB and depth of M-Mixer with MCU-self incorrectly predicts ‘wipe face’ to ‘eat meal’ and ‘sneeze/cough’, RGB and RGB+Depth of M-Mixer classify the video correctly to ‘wipe face’. Also, M-Mixer achieves higher confidence scores than M-Mixer with MCU-self for correctly predicted action classes (see the last row of 5). These results show that using cross-modality action content is more effective in leveraging complementary information from other modalities than self-modality action content. With these high performances of a single modality, our M-Mixer with all modalities successfully correctly predicts the action classes with high confidence scores. From

these results, we demonstrate the superiority of our proposed method. More qualitative results are in our supplementary material.

5. Conclusion

In this paper, we have proposed the Modality Mixer (M-Mixer) network for multi-modal action recognition. Also, we have introduced a novel recurrent unit, called Multi-modal Contextualization Unit (MCU), which is a key component of our M-Mixer network. Our MCU deals with two important factors for multi-modal action recognition: 1) complementary information across modalities and 2) temporal context of the action. Taking a feature sequence and a cross-modality content, MCU effectively learns the complementary relationships between modalities as well as the interactions between an action of the current timestep and the global action content. In comprehensive ablation studies, we demonstrate the effectiveness of our proposed methods. Moreover, we confirmed that our M-Mixer network outperforms state-of-the-art methods on NTU RGB+D 60 [28], NTU RGB+D 120 [26], and NW-UCLA [36] for multi-modal action recognition.

Acknowledgment

This work was supported by the Agency For Defense Development by the Korean Government (UD190031RD).

References

- [1] Jean-Baptiste Alayrac, Adria Recasens, Rosalia Schneider, Relja Arandjelović, Jason Ramapuram, Jeffrey De Fauw, Lucas Smaira, Sander Dieleman, and Andrew Zisserman. Self-supervised multimodal versatile networks. *Advances in Neural Information Processing Systems*, 33:25–37, 2020.
- [2] Humam Alwassel, Dhruv Mahajan, Bruno Korbar, Lorenzo Torresani, Bernard Ghanem, and Du Tran. Self-supervised learning by cross-modal audio-video clustering. *Advances in Neural Information Processing Systems*, 33:9758–9770, 2020.
- [3] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 11621–11631, 2020.
- [4] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 6299–6308, 2017.
- [5] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*, pages 1724–1734, 2014.
- [6] Nieves Crasto, Philippe Weinzaepfel, Karteek Alahari, and Cordelia Schmid. Mars: Motion-augmented rgb stream for action recognition. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 7882–7891, 2019.
- [7] Srijan Das, Rui Dai, Di Yang, and Francois Bremond. Vpn++: Rethinking video-pose embeddings for understanding activities of daily living. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [8] Srijan Das, Saurav Sharma, Rui Dai, Francois Bremond, and Monique Thonnat. Vpn: Learning video-pose embedding for activities of daily living. In *Proceedings of European Conference on Computer Vision*, pages 72–90. Springer, 2020.
- [9] Chhavi Dhiman and Dinesh Kumar Vishwakarma. View-invariant deep architecture for human action recognition using two-stream motion and shape temporal dynamics. *IEEE Transactions on Image Process.*, 29:3835–3844, 2020.
- [10] Haodong Duan, Yue Zhao, Kai Chen, Dahua Lin, and Bo Dai. Revisiting skeleton-based action recognition. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 2969–2978, 2022.
- [11] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *Proceedings of IEEE International Conference on Computer Vision (ICCV)*, pages 6202–6211, 2019.
- [12] Christoph Feichtenhofer, Axel Pinz, and Andrew Zisserman. Convolutional two-stream network fusion for video action recognition. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 1933–1941, 2016.
- [13] Nuno Cruz Garcia, Sarah Adel Bargal, Vitaly Ablavsky, Pietro Morerio, Vittorio Murino, and Stan Sclaroff. Distillation multiple choice learning for multimodal action recognition. In *Proceedings of IEEE Winter Conference on Applications of Computer Vision*, pages 2755–2764, 2021.
- [14] Nuno C Garcia, Pietro Morerio, and Vittorio Murino. Modality distillation with multiple stream networks for action recognition. In *Proceedings of European Conference on Computer Vision*, pages 103–118, 2018.
- [15] Nuno C Garcia, Pietro Morerio, and Vittorio Murino. Learning with privileged information via adversarial discriminative modality distillation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(10):2581–2593, 2019.
- [16] Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R Channing Moore, Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for audio events. In *2017 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 776–780. IEEE, 2017.
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- [18] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [19] Jian-Fang Hu, Wei-Shi Zheng, Jiahui Pan, Jianhuang Lai, and Jianguo Zhang. Deep bilinear learning for rgb-d action recognition. In *Proceedings of European Conference on Computer Vision*, pages 335–351, 2018.
- [20] Yu Huang, Chenzhuang Du, Zihui Xue, Xuanyao Chen, Hang Zhao, and Longbo Huang. What makes multi-modal learning better than single (provably). *Advances in Neural Information Processing Systems*, 34:10944–10956, 2021.
- [21] Md Mofijul Islam and Tariq Iqbal. Hamlet: A hierarchical multimodal attention-based human activity recognition algorithm. In *Proceedings of IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 10285–10292. IEEE, 2020.
- [22] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. Large-scale video classification with convolutional neural networks. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 1725–1732, 2014.
- [23] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proceedings of 3rd International Conference on Learning Representations*, 2015.
- [24] Myunggi Lee, Seungeui Lee, Sungjoon Son, Gyutae Park, and Nojun Kwak. Motion feature network: Fixed motion filter for action recognition. In *Proceedings of European Conference on Computer Vision*, pages 387–403, 2018.
- [25] Jian Liu, Naveed Akhtar, and Ajmal Mian. Viewpoint invariant action recognition using rgb-d videos. *IEEE Access*, 6:70061–70071, 2018.
- [26] Jun Liu, Amir Shahrudy, Mauricio Perez, Gang Wang, Ling-Yu Duan, and Alex C Kot. Ntu rgb+ d 120: A large-scale benchmark for 3d human activity understanding. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(10):2684–2701, 2019.
- [27] AJ Piergiovanni and Michael S Ryoo. Representation flow

- for action recognition. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 9945–9953, 2019.
- [28] Amir Shahroudy, Jun Liu, Tian-Tsong Ng, and Gang Wang. Ntu rgb+d: A large scale dataset for 3d human activity analysis. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 1010–1019, 2016.
- [29] Amir Shahroudy, Tian-Tsong Ng, Yihong Gong, and Gang Wang. Deep multimodal feature analysis for action recognition in rgb+d videos. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(5):1045–1058, 2017.
- [30] Xingjian Shi, Zhoung Chen, Hao Wang, Dit-Yan Yeung, Wai-Kin Wong, and Wang-chun Woo. Convolutional lstm network: A machine learning approach for precipitation nowcasting. *Advances in Neural Information Processing Systems*, 28, 2015.
- [31] Karen Simonyan and Andrew Zisserman. Two-stream convolutional networks for action recognition in videos. *Advances in Neural Information Processing System*, 27, 2014.
- [32] Jonathan Stroud, David Ross, Chen Sun, Jia Deng, and Rahul Sukthankar. D3d: Distilled 3d networks for video action recognition. In *Proceedings of IEEE Winter Conference on Applications of Computer Vision*, pages 625–634, 2020.
- [33] Shuyang Sun, Zhanghui Kuang, Lu Sheng, Wanli Ouyang, and Wei Zhang. Optical flow guided feature: A fast and robust motion representation for video action recognition. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 1390–1399, 2018.
- [34] 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 IEEE Conference on Computer Vision and Pattern Recognition*, pages 6450–6459, 2018.
- [35] Huogen Wang, Zhanjie Song, Wanqing Li, and Pichao Wang. A hybrid network for large-scale action recognition from rgb and depth modalities. *Sensors*, 20(11):3305, 2020.
- [36] Jiang Wang, Xiaohan Nie, Yin Xia, Ying Wu, and Song-Chun Zhu. Cross-view action modeling, learning and recognition. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 2649–2656, 2014.
- [37] 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 *Proceedings of European Conference on Computer Vision*, pages 20–36. Springer, 2016.
- [38] Pichao Wang, Wanqing Li, Jun Wan, Philip Ogunbona, and Xinwang Liu. Cooperative training of deep aggregation networks for rgb-d action recognition. In *Proceedings of Association for the Advancement of Artificial Intelligence (AAAI) Conference on Artificial Intelligence*, volume 32, 2018.
- [39] Weiyao Wang, Du Tran, and Matt Feiszli. What makes training multi-modal classification networks hard? In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 12695–12705, 2020.
- [40] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 7794–7803, 2018.
- [41] Jiaojiao Zhao and Cees GM Snoek. Dance with flow: Two-in-one stream action detection. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 9935–9944, 2019.
- [42] Yi Zhu and Shawn Newsam. Random temporal skipping for multirate video analysis. In *ACCV*, pages 542–557. Springer, 2018.