

Applying Deep Learning Models to Mouse Behavior Recognition

Ngoc Giang Nguyen¹, Dau Phan¹, Favorisen Rosyking Lumbanraja¹, Mohammad Reza Faisal¹, Bahriddin Abapihi¹, Bedy Purnama¹, Mera Kartika Delimayanti¹, Kunti Robiatul Mahmudah¹, Mamoru Kubo², Kenji Satou²

¹Graduate School of Natural Science and Technology, Kanazawa University, Kanazawa, Japan; ²Institute of Science and Engineering, Kanazawa University, Kanazawa, Japan

Correspondence to: Ngoc Giang Nguyen, giangnn.bkace@gmail.com

Keywords: Mouse Behavior Recognition, Deep Learning, I3D Models, R(2 + 1)D Models

Received: November 28, 2018

Accepted: February 25, 2019

Published: February 28, 2019

Copyright © 2019 by authors and Scientific Research Publishing Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

ABSTRACT

In many animal-related studies, a high-performance animal behavior recognition system can help researchers reduce or get rid of the limitation of human assessments and make the experiments easier to reproduce. Recently, although deep learning models are holding state-of-the-art performances in human action recognition tasks, these models are not well-studied in applying to animal behavior recognition tasks. One reason is the lack of extensive datasets which are required to train these deep models for good performances. In this research, we investigated two current state-of-the-art deep learning models in human action recognition tasks, the I3D model and the R(2 + 1)D model, in solving a mouse behavior recognition task. We compared their performances with other models from previous researches and the results showed that the deep learning models that pre-trained using human action datasets then fine-tuned using the mouse behavior dataset can outperform other models from previous researches. It also shows promises of applying these deep learning models to other animal behavior recognition tasks without any significant modification in the models' architecture, all we need to do is collecting proper datasets for the tasks and fine-tuning the pre-trained models using the collected data.

1. INTRODUCTION

Researchers widely use many animals from fruit flies, mice to primates for studying biology, psychology or for developing new therapies or medicines. In many researches, observing the behaviors of the animals is a crucial step to get the data which is needed for answering research questions. Since watching

and annotating the behaviors of these animals in hours of video clips are hard works, it's necessary to have a reliable and automated behavior recognition system to delegate these works to computers. With a well-performed system, we could not only solve the problem of the limitation of human assessments but also make the experiments easier to reproduce.

Many studies reported works in creating such systems for animal behavior recognition tasks. In the paper of Jhuang H. *et al.* [1], they created a system to automatically analyze behaviors of mice in home-cages. The system contains two modules: a feature extraction module and a classification module. In the feature extraction module, for each frame, they calculated the mouse's position and velocity based features and combined them with motion features which are extracted from the adjacent frames using an algorithm in [2]. These features then fed into an SVMHMM (Support Vector Machine-Hidden Markov Models [3]) to assess the action in the frame. In another research [4], Jiang J. *et al.* also used a similar approach but with a different feature extractor and classifier. For the feature extraction module, they detected interest points using a modified version of the algorithm in [5], then they extracted contextual and visual features from these points. And they fed these extracted features into a shallow neural network that has only one hidden layer to assess the actions in the frames. The changes in the feature extraction method and the classification algorithm slightly improved the performance of the system in comparison to the previous paper. And it showed that the design of the feature extraction module can affect the performance of the whole system. However, creating good feature extractors is not an easy task. It requires much expert knowledge and carefulness and it is not always successful. And the abilities of these created systems are highly limited to the problems they were designed to solve. For example, an automated mouse behavior recognition system may not work well in a raccoon behavior recognition task, although the two animals are sharing many similarities in their appearance.

We could solve the above problem by using deep learning models which have the ability of automated learning to extract useful features from given data. Because of having this ability, deep learning models are widely used in many application fields from computer vision, speech recognition to natural language processing and often become state-of-the-art models in the field they applied to. And we can use the same models for tasks that have similarities without significant changes in the architecture of the models.

Though its high performance, it is not easy to apply deep learning models for whatever tasks we have because these models have too many parameters that it requires an extensive amount of data to train these parameters. And it is one of the reasons why deep learning models have very high performances in human action recognition tasks but not well-studied in applying to animal behavior recognition tasks.

To fill in this gap, in this research, we investigated two current state-of-the-art human action recognition deep learning models in applying to a mouse behavior recognition task. The first model we investigated is the I3D model [6], which implements an inflated version of the inception module architecture [7]. The most important features of inception module are the utilization of the combined effects of filters of different sizes and pooling kernels all in one layer; and the usage of 1×1 convolutional filters which not only help to reduce the number of parameters but also introduce new combinations of features to its next layers. The second model we investigated in this research is the R(2 + 1)D model [8] which implements a 3D version of the residual module architecture [9]. This architecture allows the model to go deeper by solving the vanishing of information when training deep models.

To deal with the scarcity of training data, we did not train the models from randomly-initialized parameters but we used the parameters that were pre-trained on human action recognition tasks. By doing so, we can transfer knowledge that related to action recognition from human's tasks to the new models [10].

In the next section, we show the dataset which we used to evaluate the performances of the two deep learning models in the mouse behavior recognition task. In Section 3, we describe in detail experiments and results of the evaluating process. Finally, we give some conclusions in Section 4.

2. THE MOUSE BEHAVIOR DATASET

In the research of H. Jhuang, *et al.* [9], they introduced a task of neurobehavioral analysis of mouse

phenotypes by monitoring the mouse's behaviors over long periods of time. In this experiment, each mouse is put in a transparent home cage, and these behaviors are recorded from a perpendicular angle to the side of the cages using consumer grade cameras.

In order to create a machine learning system to automatically analyze mouse's behaviors, Jhuang and his colleagues have created a mouse behavior dataset by annotating the mouse's behaviors in over 10 hours of recorded videos. In their dataset, they have annotated 8 types of behavior: drinking, eating, grooming, hanging, rearing, walking, resting and micro-movements of the head. Example scenes of these behaviors are shown in [Figure 1](#), and descriptions of these behaviors are shown in [Table 1](#).

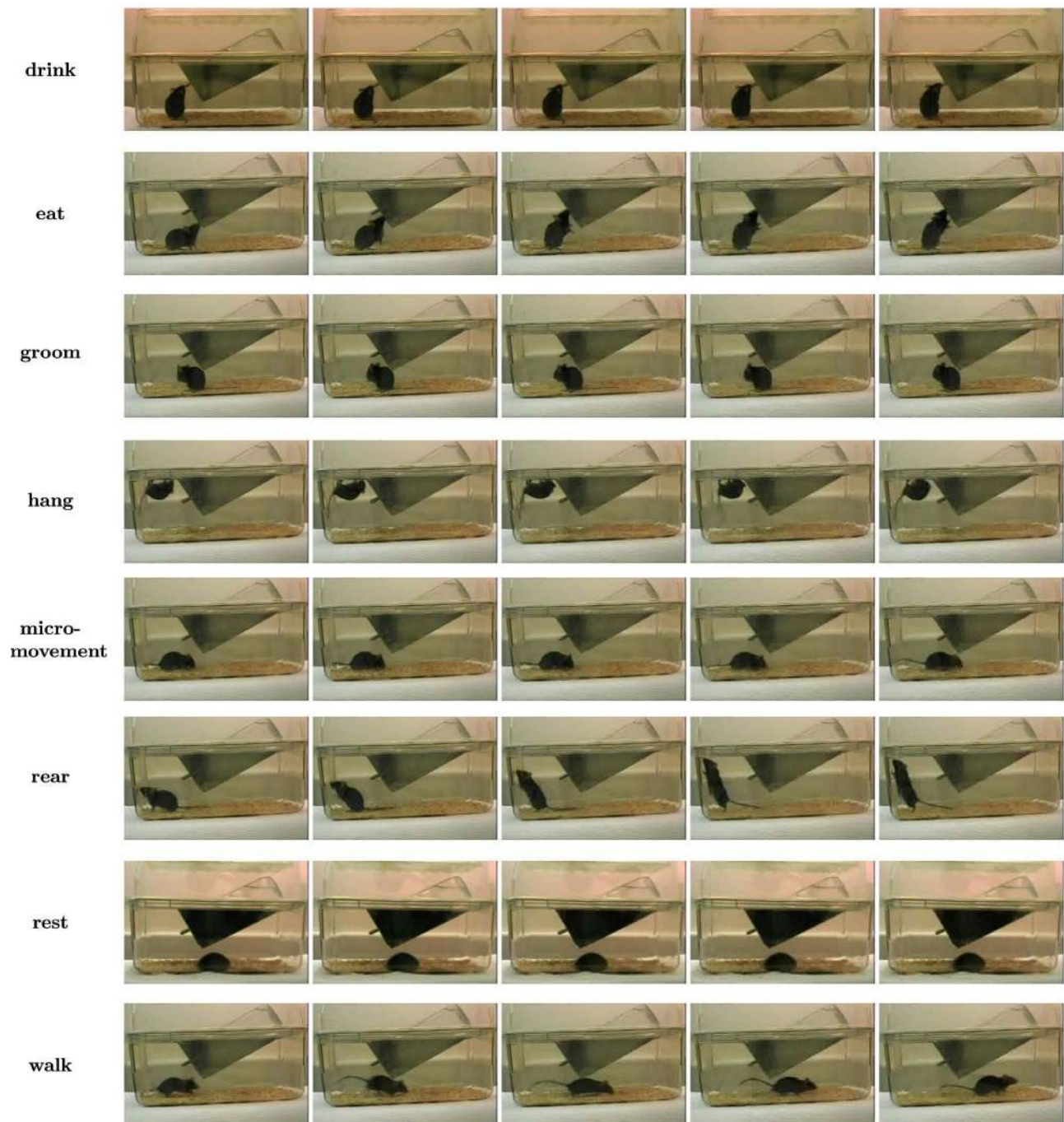


Figure 1. Example scenes of mouse's behaviors.

In the created mouse behavior dataset, among totally more than 9000 short clips, only 4200 clips that are most unambiguous were selected to create the “clipped database”. It includes about 285,000 frames and corresponds to about 2.5 hours of recorded videos. In this research, in order to properly evaluate the performance of the deep learning models, we decided to use only this subset to eliminate the ambiguous in the data that even human cannot declare. The distribution of number of frames of each behavior in the “clipped database” is shown in [Figure 2](#).

3. EXPERIMENTS AND RESULTS

3.1. Data Preparation

To generate optical-flow data from RGB data, we used the implementation of the TV-L1 algorithm from the research of [11] in OpenCV library. For each RGB frame, we input its previous frame and itself to the algorithm, and the algorithm outputs one optical-flow frame that has the same size as the inputs and contains two channels for horizontal and vertical movements respectively.

For data augmentation, we used the same method that used in the research of Carreira, J. and Zisserman, A. [6]. Each video frame in the dataset has a size of 320×240 pixels. For the I3D model, we resized

Table 1. Behaviors description.

No.	Behavior name	Description
1	<i>drink</i>	The mouse drinks water from the water-feed nipple.
2	<i>eat</i>	The mouse eats food from the food-feed door.
3	<i>groom</i>	The mouse grooms its coat.
4	<i>hang</i>	The mouse hangs on the top of the cage.
5	<i>micro-movement</i>	The mouse slightly moves its head around.
6	<i>rear</i>	The mouse rears on the side of the cage.
7	<i>rest</i>	The mouse stays stable or sleeps. There is no movement at all.
8	<i>walk</i>	The mouse walks or runs inside the cage.

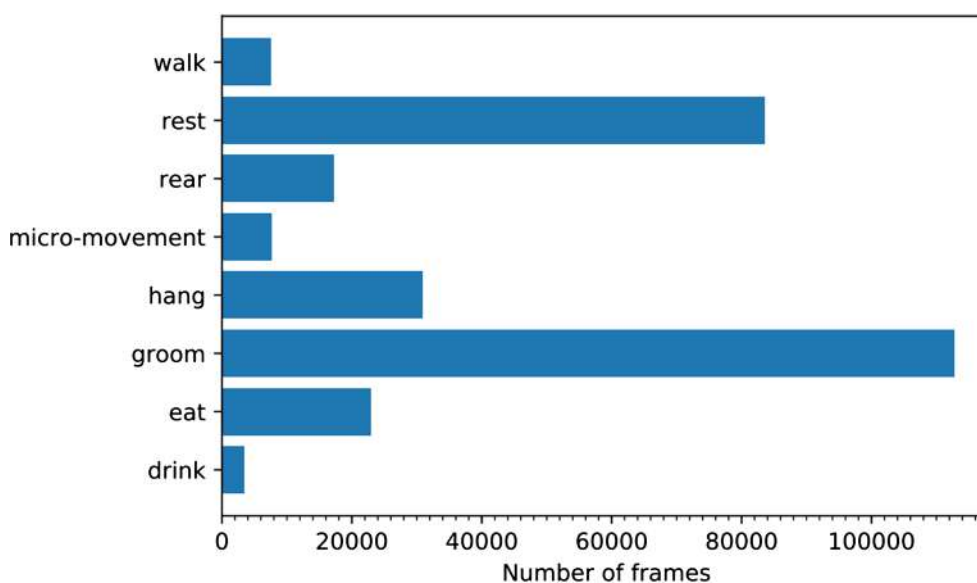


Figure 2. Distribution of number of frames of each behavior in the “clipped database”.

them to size 255×255 pixels. Then we randomly cropped them to a size of 224×224 pixels and randomly horizontal flip them to create input frames. For R(2 + 1)D model, we resized the images to 128×128 pixels then randomly cropped them to a size of 112×112 pixels. This data augmentation method helps us to increase the accuracy of prediction of about 3%.

3.2. The Models

The I3D models are derived from Inception-V1 models [7]. To benefit from the 2D architecture, all filters and pooling kernels of 2D models were inflated to 3D by endowing them with an additional temporal dimension, *i.e.* $N \times N$ filters become $N \times N \times N$ filters. Then, the weights of 2D filters are repeated N times along the temporal dimension to bootstrap parameters from pre-trained 2D models to the 3D models. We showed the architecture of an inflated inception module used in I3D models in Figure 3 and the detail of the architecture of the I3D model we used in this research in Figure 5.

The R(2 + 1)D models are derived from 2D versions [9] by replacing each 2D convolutional layer with two 3D convolutional layers, one for 2D image dimensions which have filters with size $1 \times N \times N$ and one for the time dimension which has filters with size $M \times 1 \times 1$. In some layers of the R(2 + 1)D models, to keep the total number of parameter to be the same as the 2D versions, the number of filter in these layers are calculated using the formula shown in Figure 4. The detail of the architecture of the R(2 + 1)D model we used in this research is shown in Figure 5.

For both models, we used 16 successive frames as an input (current frame, its 8 previous frames and its 7 next frames). To initialize parameters of the model, for the I3D models, we used weights from model-checkpoints that were pre-trained on ImageNet data [12]; and for the R(2 + 1)D models, we used weights from model-checkpoints that were pre-trained on Sport1M [13] and Kinetics data [14]. To fine-tune the models, we used momentum optimizer from the TensorFlow framework with momentum value equal to 0.9 and a learning rate start from $1e-3$ and decay to $5e-5$ after several thousands of iterations. We also

Inflated Inception Module

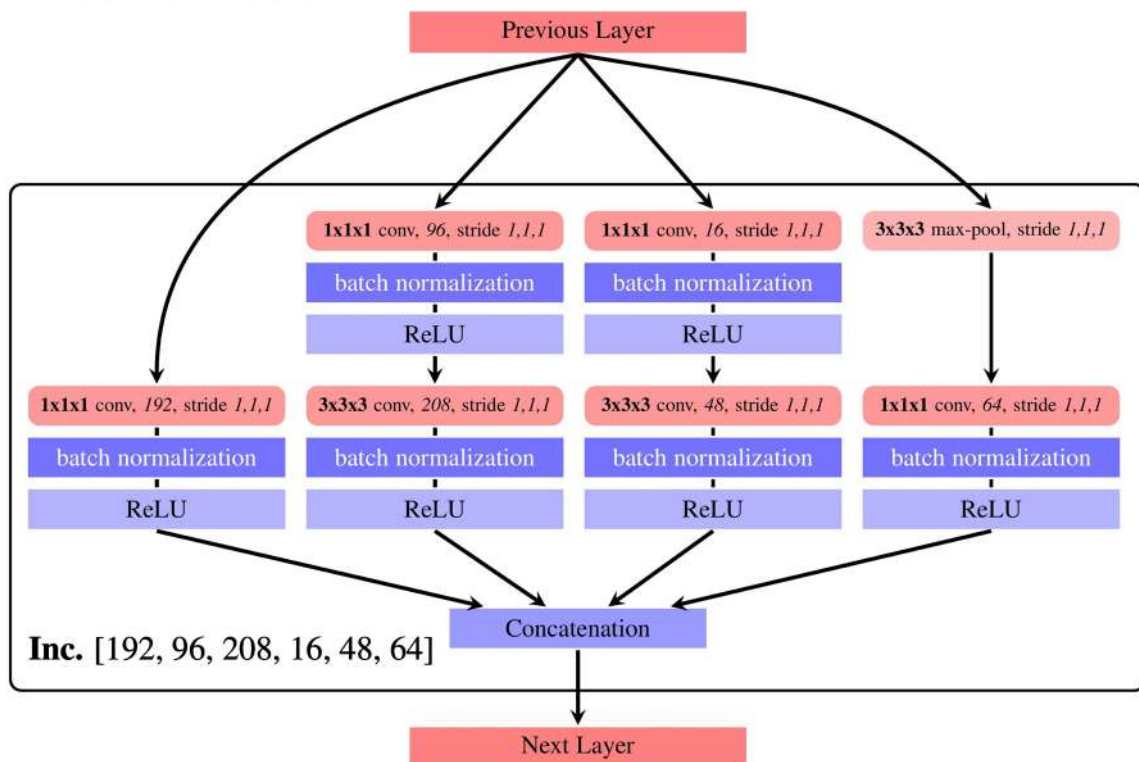
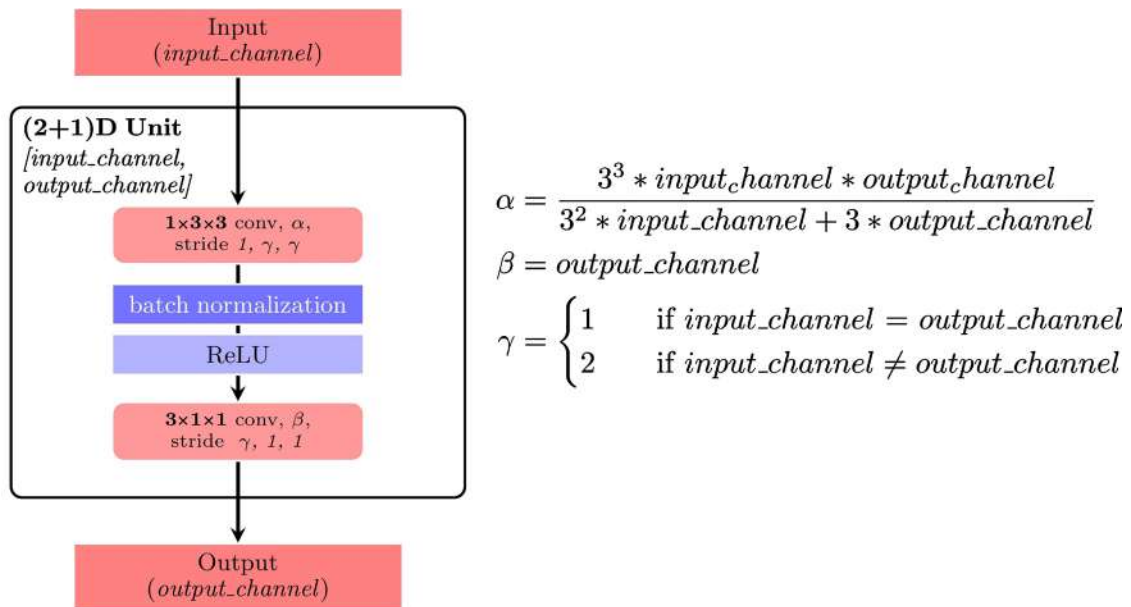


Figure 3. Architecture of an inflated inception module.



Residual Module

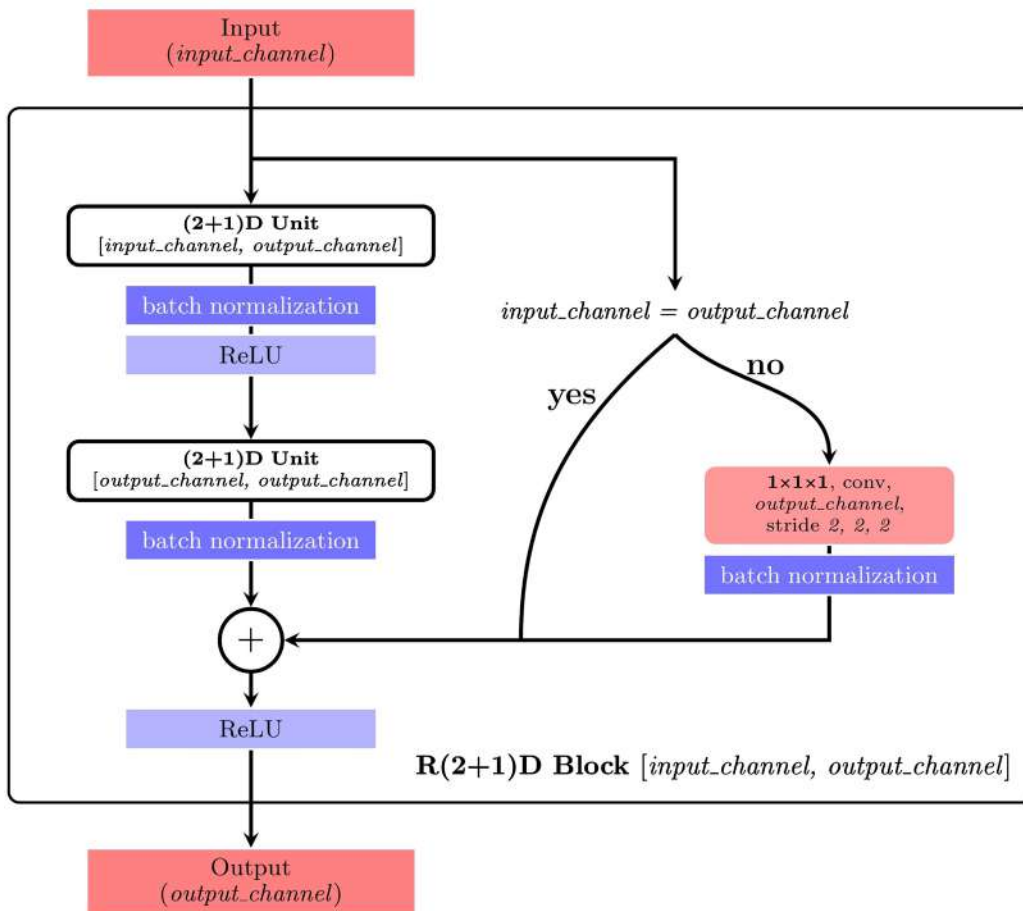


Figure 4. Architecture of a (2 + 1)D residual module.

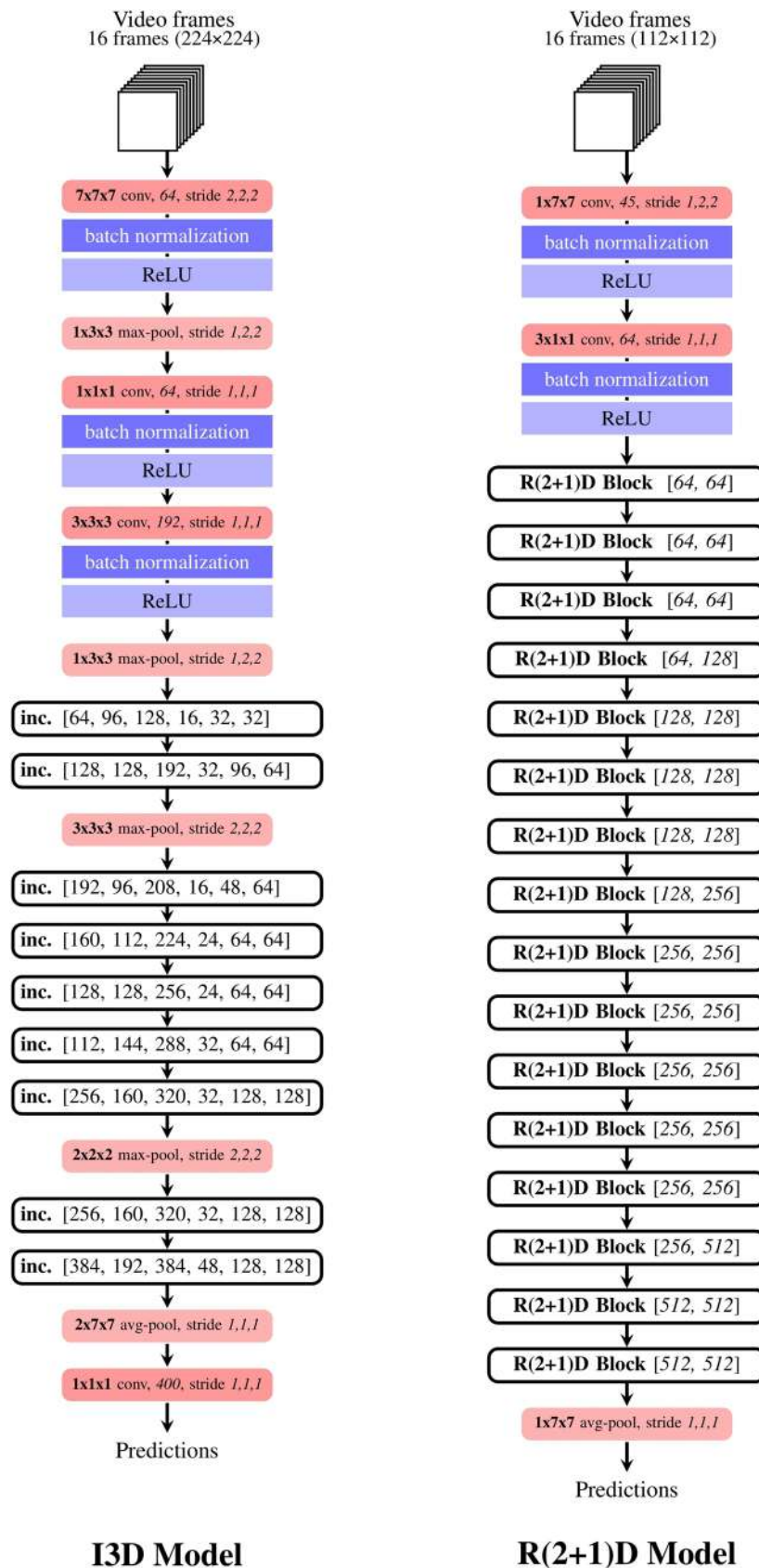


Figure 5. Architecture of the I3D model and the R(2 + 1)D model.

used dropout in fully connected layers with keep-probability of 36% to reduce the effect of overfitting when fine-tuning the models.

As discussed in the paper of Carreira, J. and Zisserman, A. [6], although I3D models can learn motion features from RGB input videos, using optical-flow as inputs can introduce some recurrent sense to the models and significantly improved the performances. In this research, we also use the same fusion method to combine output predictions of I3D models and R(2 + 1)D models. The two-stream fusion method is illustrated in Figure 6. To investigate the effects of different two-stream fusion ratios in prediction performances, we tested various fusion ratios of the two models by setting different values for *rgb_weight* and *flow_weight* in the Two-stream fusion module. For example, if only using 30% of RGB data fine-tuned model's output and 70% of optical-flow data fine-tuned model's output then *rgb_weight* is set to 0.3 and *flow_weight* is set to 0.7.

Because frames of the dataset come from 12 different videos, we used leave-one-videos-out cross-validation to properly evaluate the performance of the models. For each video, we used all the frames extracted from it as testing data and all the frames extracted from the other videos as training data. We used the training data to fine-tune the models and used the fine-tuned models to predict labels for testing data. Then we count the total number of correct and incorrect prediction and calculate the accuracy.

3.3. Results

Figure 7 shows the results of using different fusion ratio of RGB and optical-flow data fine-tuned models on accuracies of prediction of each behavior. And Figure 8 and Figure 9 show confusion matrixes of correct and incorrect prediction ratio of behaviors in combinations of *rgb_weight* and *flow_weight*.

In Figure 7, we can see that for “drink” behaviors, combinations with more portion of RGB fine-tuned models have better performance than combinations with more portion of optical-flow fine-tuned models for both I3D models and R(2 + 1)D models. And the performance of R(2 + 1)D models

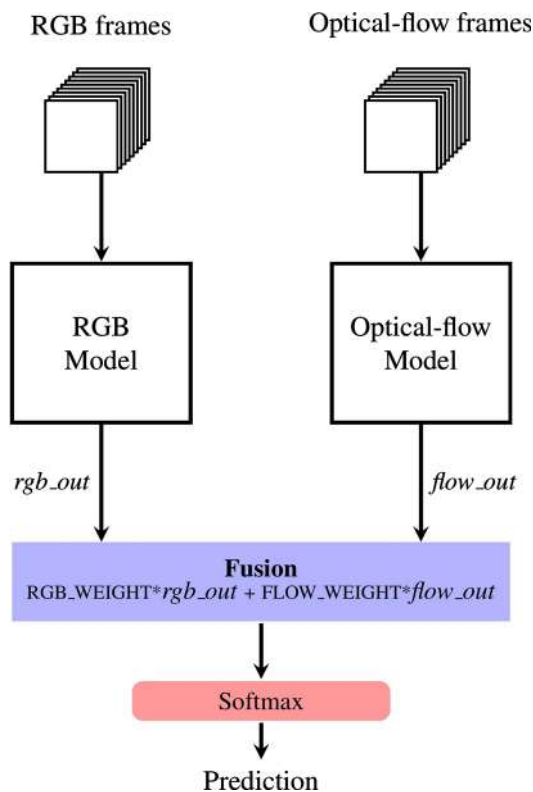


Figure 6. The two-stream fusion method.

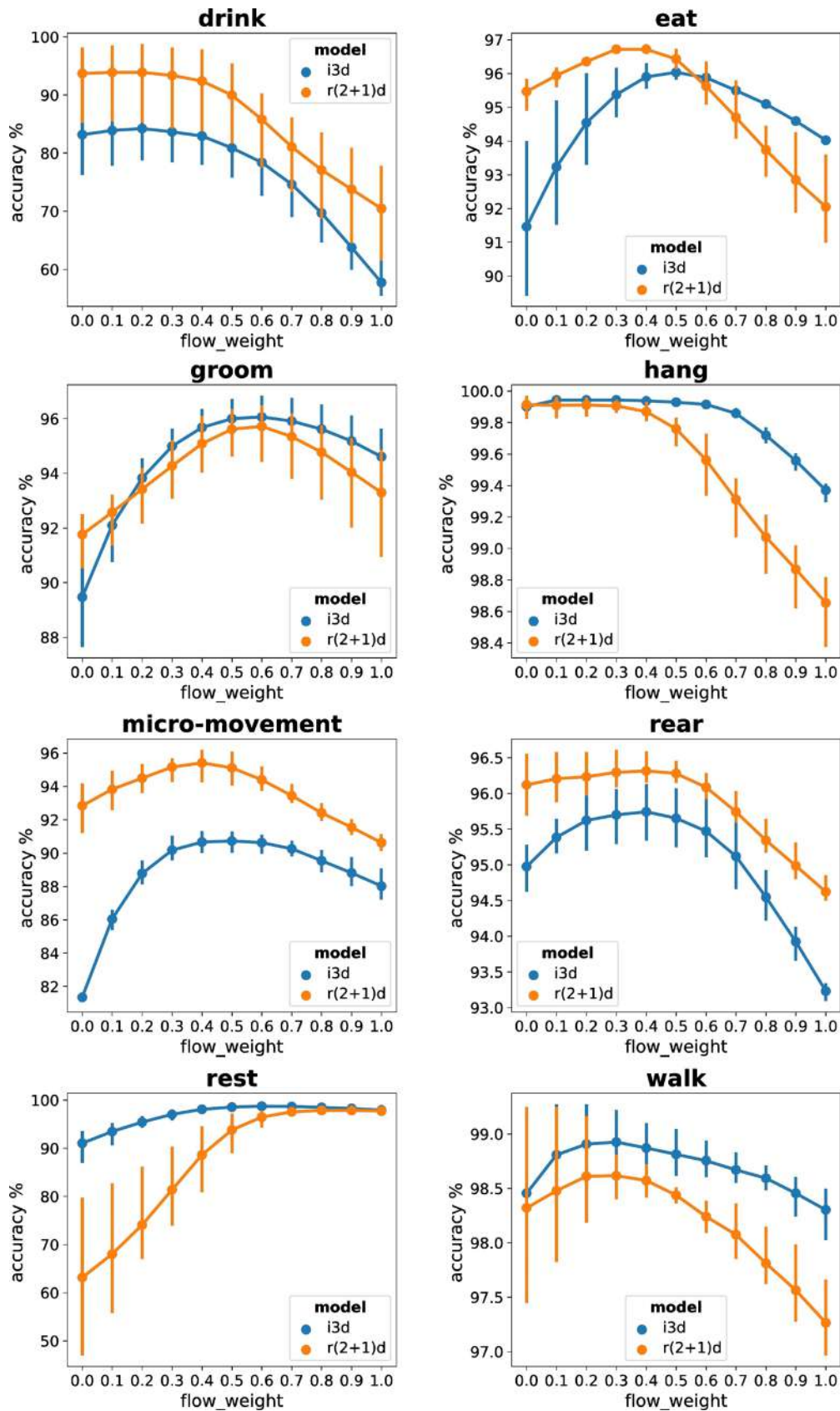


Figure 7. Accuracies of the prediction of each behavior with different two-stream fusion ratios.

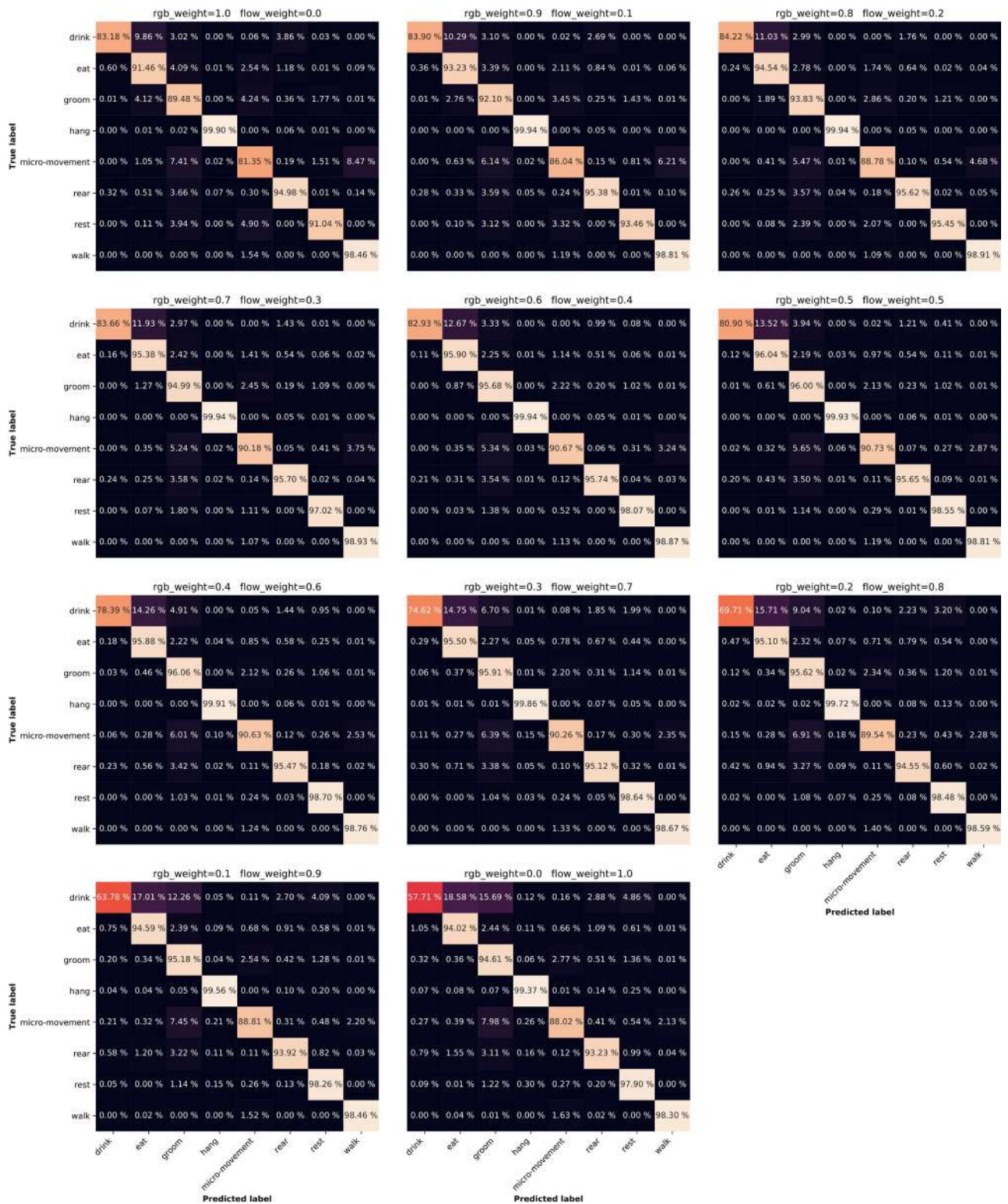


Figure 8. Confusion matrices of predictions of I3D models with different two-stream fusion ratios.

are better than the performance of the I3D models in this behaviors. In **Figure 8** and **Figure 9**, from confusion matrixes of both I3D model and R(2 + 1)D model, we can see that almost false predicted samples of “*drink*” behaviors are misclassified as “*eat*” behaviors. We can explain that as the water-feed nipple and

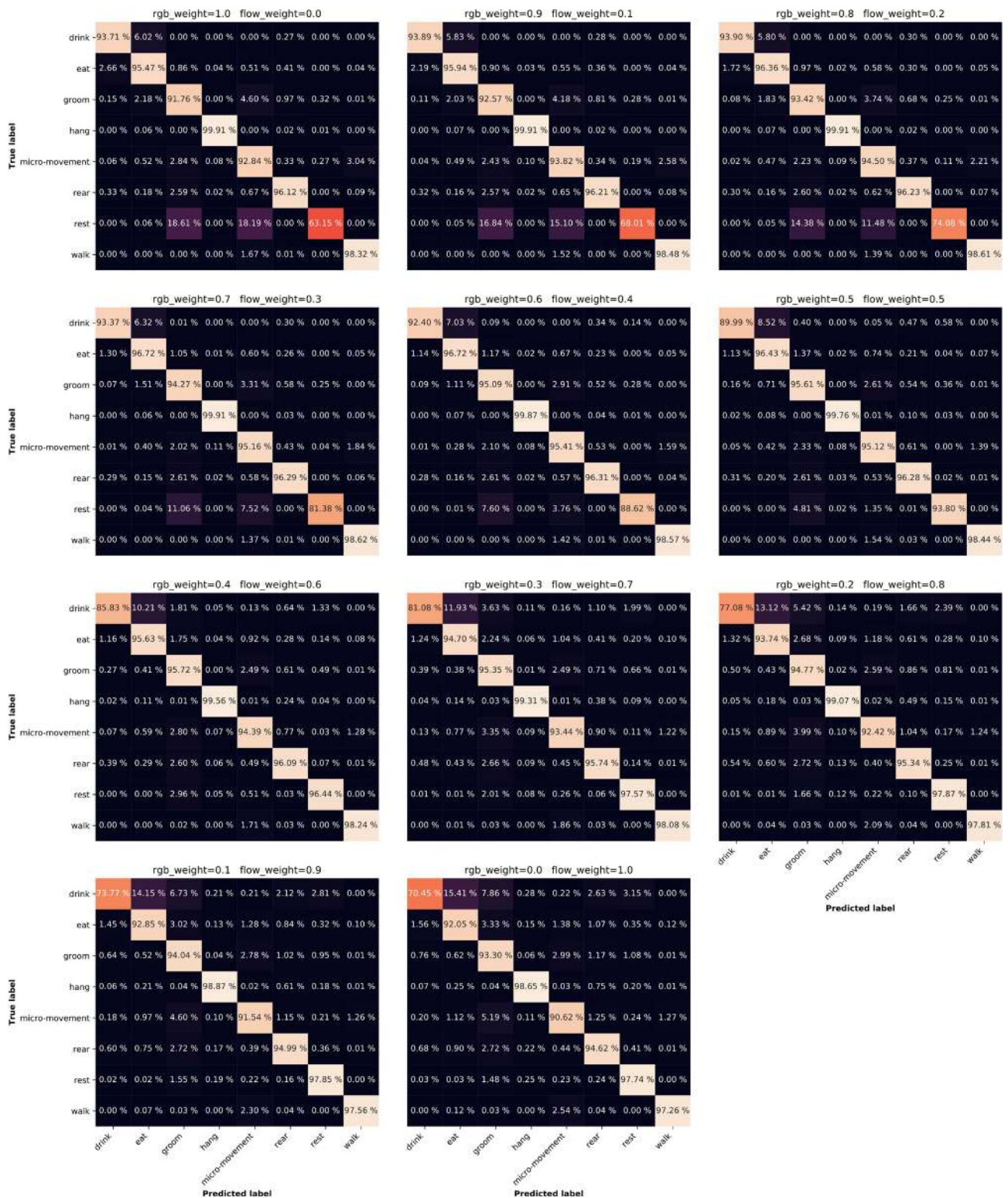


Figure 9. Confusion matrices of predictions of R(2 + 1)D models with different two-stream fusion ratios.

food-feed door are quite close to each other; Sometimes the two behaviors look very similar and the dataset is also imbalanced with the ratio of the number of “*drink*” frames to the number of “*eat*” frames is about 1:6.85. Therefore, the models tend to predict “*drink*” behaviors as “*eat*” behaviors in many cases. It

also explains why models fine-tuned using RGB data are more precise in distinguishing the two behaviors than models fine-tuned using optical-flow data. RGB data fine-tuned models can utilize the information of the mouse's mouth contact with water-feed nipple or food-feed door. However, this information is lost in optical-flow data because there is no motion of water-feed nipple or food-feed door in the scenes.

For “eat”, “groom”, “micro-movement”, “rear”, and “walk” behaviors, we can see that using a right combination of RGB data fine-tuned models and optical-flow data fine-tuned models give a better performance than using these models only. The R(2 + 1)D model outperforms I3D model in classifying “micro-movement” and “rear” behaviors but the I3D model is better in classifying “walk” behaviors.

The two models work very well on classifying “hang” behaviors and their performances just slightly reduce when we use a high portion of optical-flow data fine-tuned models because of the lack of the mouse surrounding environment in the optical-flow data.

And for the “rest” behaviors, it is easy to understand why using a high portion of optical-flow data fine-tuned models give better performance as “rest” behaviors are different from other behaviors that they have no movement in the scenes.

Overall, the two deep learning models we investigated in this research outperform the previous research model in the Mouse behavior dataset as shown in Table 2. The accuracies of the two models with different fusion ratio are shown in Figure 10. Both models have best performances at fusion ratio of

Table 2. Comparison of performance of models.

No	Model	Accuracy (%)
1	MF + SVMHMM [1]	93.0
2	FV + NN [4]	95.9
3	I3D	96.9
4	R(2 + 1)D	96.3

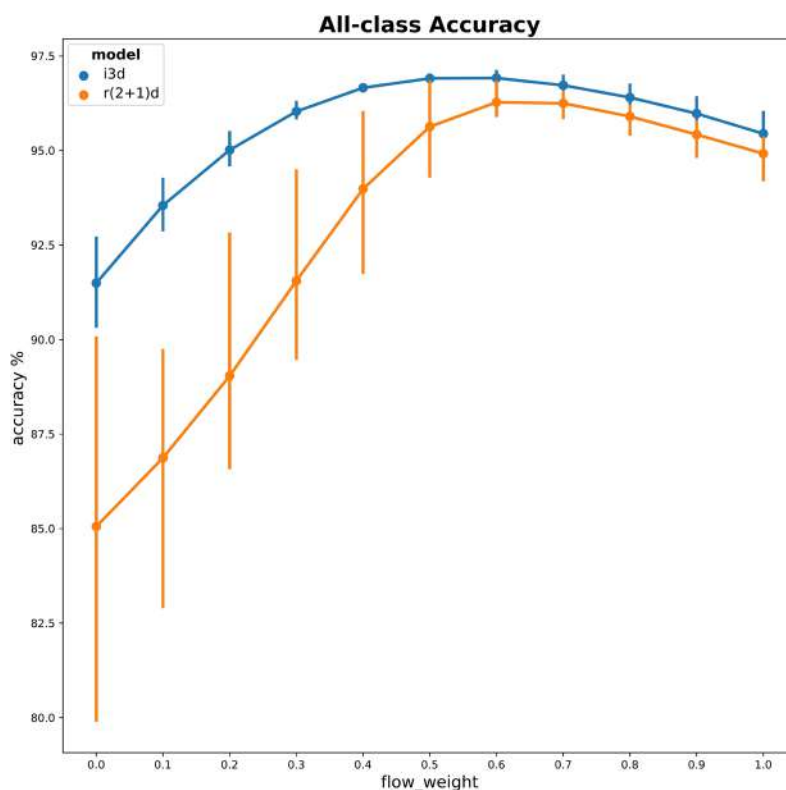


Figure 10. Accuracies of I3D models and R(2 + 1)D models with different two-stream fusion ratios.

40% RGB data fine-tuned models and 60% optical-flow data fine-tuned models. The I3D model achieves 96.9% of accuracy and the R(2 + 1)D achieves 96.3% of accuracy.

4. CONCLUSIONS

We have investigated two current state-of-the-art deep learning models for human action recognition in a mouse behavior recognition task. Both models outperformed the models from previous researches. It proves that our approach of utilizing deep learning models that pre-trained on human action datasets and fine-tuning them for animal behavior recognition tasks is efficient despite the scarcity of training data. We also showed the effect of two-stream fusion ratios on the predictions.

The fine-tuned models can precisely recognize most of behaviors they learned from the mouse behavior dataset. But there are some difficulties in classifying behaviors that are ambiguous or similar to other behaviors. Our proposal to solve the problem is to collect more data on difficult-to-classify behaviors. And we can redesign experimental environment such as changing the camera position or the cage configuration in order to minimize the ambiguity between behaviors.

For further research, we will collect behavior data of other animals. Then we will use them to fine-tune the fine-tuned models we achieved from this research to experiment if we can really efficiently utilize deep learning models for animal behavior recognition tasks without any requirements of extensive data for training these models.

ACKNOWLEDGEMENTS

In this research, the super-computing resource was provided by Human Genome Center, the Institute of Medical Science, The University of Tokyo. Additional computation time was provided by the super computer system in Research Organization of Information and Systems (ROIS), National Institute of Genetics (NIG). This work was supported by JSPS KAKENHI Grant Number JP18K11525.

CONFLICTS OF INTEREST

The authors declare no conflicts of interest regarding the publication of this paper.

REFERENCES

1. Jhuang, H., Garrote, E., Yu, X., Khilnani, V., Poggio, T., Steele, A.D. and Sere, T. (2010) Automated Home-Cage Behavioural Phenotyping of Mice. *Nature communications*, 1, Article Number: 68. <https://doi.org/10.1038/ncomms1064>
2. Jhuang, H., Serre, T., Wolf, L. and Poggio, T. (2007) A Biologically Inspired System for Action Recognition. *2007 IEEE 11th International Conference of Computer Vision*, Rio de Janeiro, 14-21 October 2007, 716-725. <https://doi.org/10.1109/ICCV.2007.4408988>
3. Altun, Y., Tsochantaridis, I. and Hofmann, T. (2003) Hidden Markov Support Vector Machines. *International Conference on Machine Learning*, Washington DC, 21-24 August 2003, 3-10.
4. Jiang, Z., Crokes, D., Green, B.D., Zhang, S. and Zhou, H. (2017) Behaviour Recognition in Mouse Videos Using Contextual Features Encoded by Spatial-Temporal Stacked Fisher Vectors. *International Conference on Pattern Recognition Applications and Methods*, 259-269. <https://doi.org/10.5220/0006244602590269>
5. Dollar, P., Rabaud, V., Cottrell, G. and Belongie, S. (2005) Behavior Recognition via Sparse Spatio-Temporal Feature. *IEEE International Workshop on Visual Surveillance and Performance Evaluation of Tracking and Surveillance*, Beijing, 15-16 October 2005, 65-72. <https://doi.org/10.1109/VSPETS.2005.1570899>
6. Carreira, J. and Zisserman, A. (2018) Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset. *2017 IEEE Conference on Computer Vision and Pattern Recognition*, Honolulu, HI, 21-26 July 2017, 4724-4733.

<https://doi.org/10.1109/CVPR.2017.502>

7. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V. and Rabinovich, A. (2015) Going Deeper with Convolutions. 2015 *IEEE Conference on Computer Vision and Pattern Recognition*, Boston, MA, 7-12 June 2015, 1-9. <https://doi.org/10.1109/CVPR.2015.7298594>
8. Tran, D., Wang, H., Torresani, L., Ray, J., LeCun, Y. and Paluri, M. (2018) A Closer Look at Spatiotemporal Convolutions for Action Recognition. *Computer Vision and Pattern Recognition*. <https://arxiv.org/abs/1711.11248>
9. He, K., Zhang, X., Ren, S. and Sun, J. (2015) Deep Residual Learning for Image Recognition. 2016 *IEEE Conference on Computer Vision and Pattern Recognition*, Las Vegas, NV, 27-30 June 2016, 770-778. <https://doi.org/10.1109/CVPR.2016.90>
10. Torrey, L. and Shavlik, J. (2009) Transfer Learning. In: Soria, E., Martin, J., Magdalena, R., Martinez, M. and Serrano, A., Eds., *Handbook of Research on Machine Learning Applications*, IGI Global, 242-264.
11. Zach, C., Pock, T. and Bischof, H. (2007) A Duality Based Approach for Realtime TV- L^1 Optical Flow. *Proceeding of 29th DAGM Symposium of Pattern Recognition*, **4713**, 214-223. https://doi.org/10.1007/978-3-540-74936-3_22
12. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C. and Fei, L.F. (2015) ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, **115**, 211-252. <https://doi.org/10.1007/s11263-015-0816-y>
13. Karpathy, A., Toderici, G., Shetty, S., Leung, T., Sukthankar, R. and Fei, L.F. (2014) Large-Scale Video Classification with Convolutional Neural Networks. 2014 *IEEE Conference on Computer Vision and Pattern Recognition*, Columbus, OH, 23-28 June 2014, 1725-1732. <https://doi.org/10.1109/CVPR.2014.223>
14. Kay, W., Carreira, J., Simonyan, K., Zhang, B., Hillier, C., Vijayanarasimhan, S., Viola, F., Green, T., Back, T., Natsev, P., Suleyman, M. and Zisserman, A. (2017) The Kinetics Human Action Video Dataset. *Computer Vision and Pattern Recognition*. <https://arxiv.org/abs/1705.06950>