

# BACKGROUND MIXUP DATA AUGMENTATION FOR HAND AND OBJECT-IN-CONTACT DETECTION

*Koya Tango, Takehiko Ohkawa, Ryosuke Furuta, Yoichi Sato*

The University of Tokyo, Tokyo, Japan

## ABSTRACT

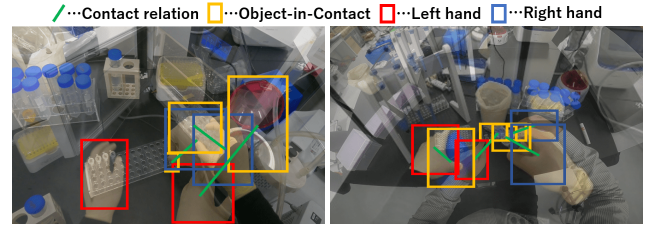
Detecting the positions of human hands and objects-in-contact (hand-object detection) in each video frame is vital for understanding human activities from videos. For training an object detector, a method called Mixup, which overlays two training images to mitigate data bias, has been empirically shown to be effective for data augmentation. However, in hand-object detection, mixing two hand-manipulation images produces unintended biases, e.g., the concentration of hands and objects in a specific region degrades the ability of the hand-object detector to identify object boundaries. We propose a data-augmentation method called Background Mixup that leverages data-mixing regularization while reducing the unintended effects in hand-object detection. Instead of mixing two images where a hand and an object in contact appear, we mix a target training image with background images without hands and objects-in-contact extracted from external image sources, and use the mixed images for training the detector. Our experiments demonstrated that the proposed method can effectively reduce false positives and improve the performance of hand-object detection in both supervised and semi-supervised learning settings.

**Index Terms**— Data Augmentation, Hand and Object-in-Contact Detection, Mixup

## 1. INTRODUCTION

Detecting the positions of a person’s hands and an object-in-contact (hand-object detection) from an image provides an important clue for understanding how the person interacts with the physical world. This hand-object detection is applicable to recognizing a person’s primitive actions, such as “taking” or “pushing”, and logging the person’s activity of interacting with the environment [1]. Shan et al. [2] built a hand-object detector for localizing hands and interacting objects on a large-scale dataset collected in naturalistic house-holding situations, such as in kitchen [3, 4, 5], DIY [2], and craft work [2, 5].

However, a hand-object detector trained on such house-holding images may not be well generalized to other hand-



(a) Ambiguous contact states (b) Ambiguous object boundaries

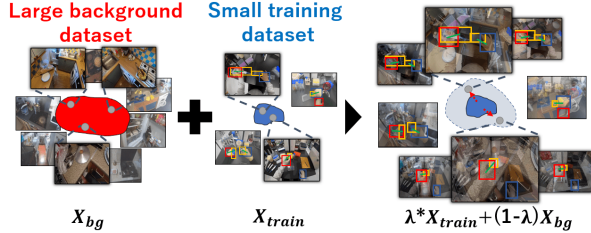
**Fig. 1: Problems with Mixup.** Naively mixing two images causes ambiguity in (a) contact states and (b) object boundaries.

manipulation images. For instance, the images in biological laboratories or factories have significantly different data distribution from the daily scenes used in training. To build an accurate hand-object detector for such unique application domains, a large amount of data and labels must be collected from scratch. However, data collection and annotation can be difficult due to various reasons, such as cost or privacy issues. In particular, expert knowledge is required to annotate the data in such specific application domains. Under these limitations, a hand-object detector may overfit the training data and lack the generalization ability due to the small amount of training data.

To improve the generalization ability of the detector trained on a small dataset, data augmentation is a key component in training. Recently, Mixup [6], a method that overlays two different images, has been used as an empirically strong augmentation for object detection [7]. Nevertheless, naively applying Mixup induces unintended biases in hand-object detection. As shown in Figure 1, (a) contact states become ambiguous when hand-object pairs from different images overlap, and (b) the concentration of hands and objects in a specific local region makes identifying object boundaries difficult. These unintended mixtures will degrade the performance of a hand-object detector.

To handle this, we propose a novel data-augmentation method, called Background Mixup, that utilizes data-mixing regularization while reducing the unintended effects in hand-object detection. As shown in Figure 2, we aim to augment

This work was supported by JST AIP Acceleration Research Grant Number JPMJCR20U1 and JST ACT-X Grant Number JPMJAX2007.



**Fig. 2: Overview of Background Mixup.** We aim to improve diversity in training data while preserving foreground’s semantics.

a training image by mixing it with the background of external image sources that does not contain the foreground (i.e., hands and objects-in-contact) and using the mixed images for training a hand-object detector. The contributions of this paper are summarized as follows.

- We propose a novel data-augmentation method, Background Mixup, that mixes a training image and a background image to improve the generalization ability of a hand-object detector in a small dataset.
- Compared with Mixup, our experiments showed that Background Mixup improves the performance of a hand-object detector in supervised and semi-supervised learning settings.
- Our method has also shown to be effective in reducing the number of false-positive predictions although Mixup has the disadvantage in this metric.

## 2. RELATED WORK

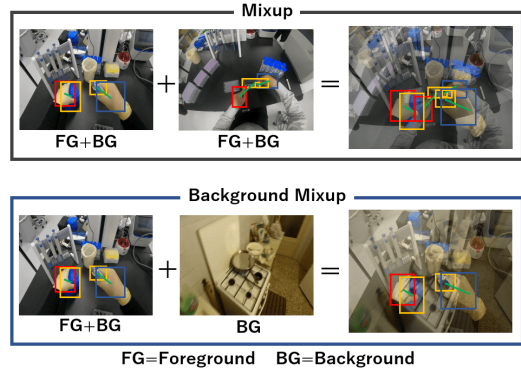
### 2.1. Hand and Object-in-Contact

Jointly analyzing hands and objects-in-contact serves to understand human behavior [8, 9]. While these studies have been conducted in a limited scale of data, Shan et al. [2] proposed a large-scale dataset for training a hand-object detector localizing hands and interacting objects, which is collected in daily situations such as EPIC-KITCHENS 2018 [3], EGTEA [4], CharadesEgo [5]. However, directly fine-tuning the hand-object detector on a small and specific dataset can lead to limited performance as discussed in Section 1.

To overcome this, we developed Background Mixup to improve the generalization ability of a hand-object detector on a small dataset in specific domains such as biomedical experiments and factory work.

### 2.2. Mixture-Based Data Augmentation

Mixture-based data augmentation mixes input data with other inputs to increase the diversity of data on a small dataset and



**Fig. 3: Comparison of Mixup [6] and Background Mixup.**

improve the generalization performance of the model. Several mixture-based methods, such as Mixup [6], CutMix [10], Mosaic [11], and Cutout [12], have been used in many downstream tasks.

These mixture-based methods are used for semi-supervised learning of object detection [13, 7]. Unbiased-Teacher [13] uses Cutout while Instant-Teaching [7] uses Mixup and Mosaic showing that Mixup particularly contributes to improving the performance of object detection. However, applying Mixup leads to unintended biases in the hand-object detection, as discussed in Section 1. These unintended mixtures will degrade the performance of a hand-object detector.

## 3. PROPOSED METHOD

In this section, we introduce our proposed training of a hand-object detector with Background Mixup data augmentation. Let  $\mathcal{X}_{train}$  and  $\mathcal{X}_{test}$  be sets of training and testing images, respectively. When the size of  $\mathcal{X}_{train}$  is small, a hand-object detector trained on  $\mathcal{X}_{train}$  may not generalize well to  $\mathcal{X}_{test}$  due to over-fitting to the training data. To solve this problem, we propose Background Mixup that uses a background image without foreground entities, i.e., hands and objects-in-contact, for increasing the diversity of the training data.

We use a trained hand-object detector [2] to extract the background images from an external image source (e.g., kitchens), which are different from our target data of  $\mathcal{X}_{train}$  and  $\mathcal{X}_{test}$ . We extract the background images in which neither object-in-contact nor hand was detected by the hand-object detector, and construct a set of the background images  $\mathcal{X}_{bg}$ .

Figure 3 shows a comparison between Mixup and Background Mixup. With Mixup, the foreground and background are combined, causing unintended effects that make the contact state ambiguous or make it difficult to identify the boundaries of objects, as shown in Figure 1. In contrast, Background Mixup reduces such unintended effects by mixing the training image with the background image, which can retain the foreground of the training image.

We denote the training image as  $I_{train} \in \mathcal{X}_{train}$  and the

background image as  $I_{bg} \in \mathcal{X}_{bg}$ . We define Background Mixup as:

$$\hat{I} = \lambda I_{train} + (1 - \lambda) I_{bg} \quad (1)$$

$$\lambda \sim \text{Beta}(\alpha, \beta). \quad (2)$$

where  $I_{train}$  and  $I_{bg}$  are randomly sampled and  $\lambda \in [0, 1]$  is a parameter controlling the degree of the mixture. Following Mixup, the parameter  $\lambda$  is drawn from beta distribution  $\text{Beta}(\alpha, \beta)$  where  $\alpha$  and  $\beta$  indicate hyperparameters to determine the distribution shape.

We use  $\hat{I}$  for training the hand-object detector instead of  $I_{train}$ . This method can be implemented with a small computational cost at the training stage. Thus no additional computational cost is required in inference.

## 4. EXPERIMENTS

### 4.1. Experimental Setup

**Datasets.** We validated our method on various hand manipulation datasets including biomedical experiments, mock factories, and kitchens. We used a first-person video dataset of biomedical experiments and a mock factory environment dataset [14] as specific application domains where the data size and variety are limited. We also used a kitchen environment dataset [3, 2] including diverse cooking scenes.

For the dataset of biomedical experiments, we recorded 12 videos that contained basic actions such as preparing reagents in a biomedical lab. The bounding boxes of hands and objects-in-contact were annotated by an expert in the field. The total duration is 27 minutes, and the number of annotated frames is 3,093. We split the 12 videos into 6:3:3 for train:val:test.

For the set of background images in the experiments on biomedical and factory datasets, we used EPIC-KITCHENS-100 [15] of cooking scenes in kitchens as an external image source. For the experiments on the cooking dataset, we used Something-Something V2 [16] of daily scenes as an external image source to augment the background appearance.

**Training details.** We measure the performance of our method in supervised and semi-supervised learning settings. For supervised learning, we fine-tune the pre-trained hand-object detector proposed by Shan et al. [2]. For semi-supervised learning, we trained the hand-object detector in the training pipeline of Unbiased-Teacher (UB-Teacher) [13]. We evaluated the performance by the average precision (AP) of hands and objects-in-contact. Note that we do not provide hand AP in an experiment with mock factory environment dataset [14] because the hand bounding boxes are not annotated.

**Baselines.** We denote our proposed Background Mixup with EPIC-KITCHENS-100 [15] and Something-Something V2 [16] as **BG-Mix<sub>K</sub>** and **BG-Mix<sub>D</sub>**, respectively. We prepared two variants of Mixup as comparison methods.

**Mixup** is the original Mixup that combines two different images within a dataset, and **Mixup<sub>K</sub>** is Mixup that mixes a training image with a randomly selected image from EPIC-KITCHENS 2018 [2, 3].

### 4.2. Quantitative Evaluation

#### 4.2.1. Supervised Learning

Table 1 lists the results on the supervised learning settings. In the biomedical and mock factory environment datasets, BG-Mix<sub>K</sub> exhibited the highest performance in mAP and object AP while Mixup and Mixup<sub>K</sub> decreased in these metrics. This is because our method avoids the unintended effects shown in Figure 1, which degrades the performance of detecting an object-in-contact. In the cooking dataset, however, Mixup, Mixup<sub>K</sub> and BG-Mix<sub>D</sub> all obtained lower obj AP and hand AP than the Supervised baseline because the dataset already has diverse foreground and background appearances even without data augmentation. The hand APs of Background Mixup did not increase because the hand APs were already saturated.

#### 4.2.2. Semi-Supervised Learning

Table 2 shows the results for semi-supervised learning with 1% labeled data. BG-Mix<sub>K</sub> and BG-Mix<sub>D</sub> exhibited the highest performance on all datasets, except for the hand AP on the biomedical data. In the cooking dataset having a variety of objects and backgrounds, although the performance of supervised learning decreased with both BG-Mix<sub>D</sub> and Mixup as shown in Section 4.2.1, BG-Mix<sub>D</sub> improved the performance of semi-supervised learning. This indicates that Background Mixup is effective under limited labels where the fully-supervised model suffers from generalizing to unknown test data.

#### 4.2.3. Analysis of False-positive Predictions

Although mAP is a standard evaluation criterion in object detection, there is a technique that can improve the mAP score by allowing many false positives with low confidence [17]. However, detection results that contain many false positives are problematic in real scenarios. Therefore, we experimented with precision to measure the percentage of false positives in detection results of a hand-object detector. Precision indicates the percentage of true positives among the predictions detected as positive. In other words, the lower precision, the higher the percentage of false positives.

Table 3 shows a comparison of precision when the percentage of labeled data is 1% in semi-supervised learning, and the confidence threshold is 0.1. While Mixup<sub>K</sub> improves mAP, the precision is decreased. This indicates training the detector on mixed images with many overlapping bounding boxes in a specific area, as illustrated in Fig. 1(b), induces the

**Table 1:** Quantitative comparisons on supervised learning.

Model	Biomedical			Factory		Cooking		
	hand AP	obj AP	mAP	hand AP	obj AP	hand AP	obj AP	mAP
Supervised	<b>90.9</b> $\pm 0.0$	70.6 $\pm 0.3$	80.7 $\pm 0.2$	-	45.0 $\pm 0.1$	<b>90.6</b> $\pm 0.0$	<b>66.4</b> $\pm 0.1$	<b>78.5</b> $\pm 0.0$
+ Mixup	<b>90.9</b> $\pm 0.0$	70.4 $\pm 0.1$	80.6 $\pm 0.0$	-	44.6 $\pm 0.0$	<b>90.6</b> $\pm 0.0$	65.6 $\pm 0.2$	78.1 $\pm 0.1$
+ Mixup <sub>K</sub>	<b>90.9</b> $\pm 0.0$	69.8 $\pm 0.4$	80.4 $\pm 0.2$	-	44.6 $\pm 0.1$	-	-	-
+ BG-Mix <sub>K</sub>	<b>90.9</b> $\pm 0.1$	<b>72.2</b> $\pm 0.2$	<b>81.0</b> $\pm 0.1$	-	<b>45.2</b> $\pm 0.1$	-	-	-
+ BG-Mix <sub>D</sub>	-	-	-	-	-	<b>90.6</b> $\pm 0.0$	65.7 $\pm 0.2$	78.2 $\pm 0.1$

**Table 2:** Quantitative comparisons on semi-supervised learning at 1% labels.

Model	Biomedical			Factory		Cooking		
	hand AP	obj AP	mAP	hand AP	obj AP	hand AP	obj AP	mAP
UB-Teacher	90.6 $\pm 0.1$	64.0 $\pm 1.1$	77.3 $\pm 0.4$	-	27.0 $\pm 0.5$	<b>90.5</b> $\pm 0.0$	41.4 $\pm 0.5$	65.9 $\pm 0.3$
+ Mixup	<b>90.9</b> $\pm 0.0$	62.4 $\pm 3.4$	76.6 $\pm 1.7$	-	30.4 $\pm 1.4$	90.4 $\pm 0.0$	46.5 $\pm 0.1$	68.5 $\pm 0.1$
+ Mixup <sub>K</sub>	90.8 $\pm 0.1$	65.4 $\pm 0.6$	78.1 $\pm 0.3$	-	31.0 $\pm 1.4$	-	-	-
+ BG-Mix <sub>K</sub>	90.8 $\pm 0.1$	<b>66.4</b> $\pm 0.1$	<b>78.6</b> $\pm 0.1$	-	<b>32.6</b> $\pm 0.7$	-	-	-
+ BG-Mix <sub>D</sub>	-	-	-	-	-	<b>90.5</b> $\pm 0.0$	<b>47.2</b> $\pm 0.5$	<b>68.9</b> $\pm 0.2$

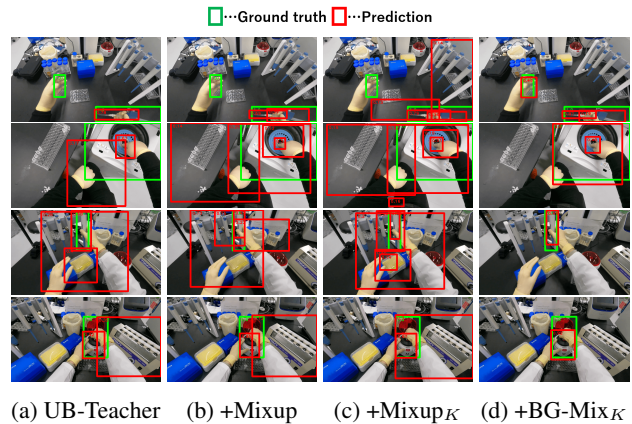
**Table 3:** Comparisons of false positive predictions on biomedical experiments dataset.

Model	mAP	Precision <sub>hand</sub>	Precision <sub>obj</sub>
UB-Teacher	77.3 $\pm 0.6$	87.1 $\pm 1.5$	<b>49.7</b> $\pm 0.1$
+ Mixup	76.6 $\pm 1.7$	76.8 $\pm 1.8$	42.0 $\pm 3.2$
+ Mixup <sub>K</sub>	78.1 $\pm 0.3$	75.6 $\pm 5.3$	38.7 $\pm 2.9$
+ BG-Mix <sub>K</sub>	<b>78.6</b> $\pm 0.1$	<b>89.1</b> $\pm 2.5$	48.8 $\pm 2.6$

bias of increasing the number of false positives. BG-Mix<sub>K</sub> can improve the mAP without increasing the number of false positives because it keeps the information of hand-object contact and avoids the concentration of target hands and objects-in-contact in a local region.

### 4.3. Qualitative Evaluation

Figure 4 shows the inference results for object-in-contact with a confidence threshold of 0.1. We observed that Figure 4 (b) Mixup and (c) Mixup<sub>K</sub> increased the number of false positives (e.g., red bounding boxes far from the ground truth bounding boxes). In contrast, the predictions of Figure 4 (d) Background Mixup is less noisy and accurately represent the location of the object-in-contact compared to ground truth. Our method of increasing the diversity of the background without changing the foreground semantics could improve the performance of a hand-object detector without increasing the number of false positives.

**Fig. 4:** Qualitative results in detecting object-in-contact.

## 5. CONCLUSION

We proposed Background Mixup, which mixes training images with background images that do not contain the hands and the objects-in-contact, whereas Mixup mixes both the foreground (i.e., the hand and the object-in-contact) and the background. Background Mixup can improve the performance of a hand-object detector in small datasets, such as biomedical experiments and mock factory environments, by increasing the diversity of the background appearances while inhibiting the unintended effects caused by Mixup. We have also shown that Background Mixup was effective in reducing the number of false positives.

## 6. REFERENCES

- [1] T. Yagi, T. Nishiyasu, K. Kawasaki, M. Matsuki, and Y. Sato, “GO-Finder: A Registration-Free Wearable System for Assisting Users in Finding Lost Objects via Hand-Held Object Discovery,” in *26th International Conference on Intelligent User Interfaces (IUI)*, 2021, pp. 139–149.
- [2] D. Shan, J. Geng, M. Shu, and D. F. Fouhey, “Understanding Human Hands in Contact at Internet Scale,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020, pp. 9869–9878.
- [3] D. Damen, H. Doughty, G. M. Farinella, S. Fidler, A. Furnari, E. Kazakos, D. Moltisanti, J. Munro, T. Perrett, W. Price, et al., “Scaling Egocentric Vision: The EPIC-KITCHENS Dataset,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 720–736.
- [4] Y. Li, M. Liu, and J. M. Rehg, “In The Eye of Beholder: Joint Learning of Gaze and Actions in First Person Video,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 619–635.
- [5] G. A. Sigurdsson, A. Gupta, C. Schmid, A. Farhadi, and K. Alahari, “Actor and Observer: Joint Modeling of First and Third-Person Videos,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 7396–7404.
- [6] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, “Mixup: Beyond Empirical Risk Minimization,” *arXiv preprint arXiv:1710.09412*, 2017.
- [7] Q. Zhou, C. Yu, Z. Wang, Q. Qian, and H. Li, “Instant-Teaching: An End-to-End Semi-Supervised Object Detection Framework,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 4081–4090.
- [8] Y. Li, Z. Ye, and J. M. Rehg, “Delving Into Egocentric Actions,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 287–295.
- [9] S. Bambach, S. Lee, D. J. Crandall, and C. Yu, “Lending a Hand: Detecting Hands and Recognizing Activities in Complex Egocentric Interactions,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2015, pp. 1949–1957.
- [10] S. Yun, D. Han, S. J. Oh, S. Chun, J. Choe, and Y. Yoo, “Cutmix: Regularization Strategy to Train Strong Classifiers with Localizable Features,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 6023–6032.
- [11] Z. Ge, S. Liu, F. Wang, Z. Li, and J. Sun, “YOLOX: Exceeding YOLO Series in 2021,” *arXiv preprint arXiv:2107.08430*, 2021.
- [12] T. DeVries and G. W. Taylor, “Improved Regularization of Convolutional Neural Networks with Cutout,” *arXiv preprint arXiv:1708.04552*, 2017.
- [13] Y.-C. Liu, C.-Y. Ma, Z. He, C.-W. Kuo, K. Chen, P. Zhang, B. Wu, Z. Kira, and P. Vajda, “Unbiased Teacher for Semi-Supervised Object Detection,” in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- [14] F. Ragusa, A. Furnari, S. Livatino, and G. M. Farinella, “The MECCANO Dataset: Understanding Human-Object Interactions from Egocentric Videos in an Industrial-Like Domain,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 2021, pp. 1569–1578.
- [15] D. Damen, H. Doughty, G. M. Farinella, A. Furnari, J. Ma, E. Kazakos, D. Moltisanti, J. Munro, T. Perrett, W. Price, and M. Wray, “Rescaling Egocentric Vision: Collection, Pipeline and Challenges for EPIC-KITCHENS-100,” *International Journal of Computer Vision (IJCV)*, vol. 130, no. 1, pp. 33–55, 2022.
- [16] R. Goyal, S. Ebrahimi Kahou, V. Michalski, J. Materzynska, S. Westphal, H. Kim, V. Haenel, I. Fruend, P. Yianilos, M. Mueller-Freitag, et al., “The” Something Something” Video Database for Learning and Evaluating Visual Common Sense,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2017, pp. 5842–5850.
- [17] T. Ito, “mAP understanding with code and its tips,” <https://www.kaggle.com/its7171/map-understanding-with-code-and-its-tips>, 2021.