

STUDENT DANGEROUS BEHAVIOR DETECTION IN SCHOOL

Huayi Zhou* Fei Jiang† Hongtao Lu*

* Shanghai Jiao Tong University, sjtu_zhy@sjtu.edu.cn, htlu@sjtu.edu.cn

† East China Normal University, fjiang@mail.ecnu.edu.cn

ABSTRACT

Video surveillance systems have been installed to ensure the student safety in schools. However, discovering dangerous behaviors, such as fighting and falling down, usually depends on untimely human observations. In this paper, we focus on detecting dangerous behaviors of students automatically, which faces numerous challenges, such as insufficient datasets, confusing postures, keyframes detection and prompt response. To address these challenges, we first build a danger behavior dataset with locations and labels from surveillance videos, and transform action recognition of long videos to an object detection task that avoids keyframes detection. Then, we propose a novel end-to-end dangerous behavior detection method, named DangerDet, that combines multi-scale body features and keypoints-based pose features. We could improve the accuracy of behavior classification due to the highly correlation between pose and behavior. On our dataset, DangerDet achieves 71.0% mAP with about 11 FPS. It keeps a better balance between the accuracy and time cost.

Index Terms— dangerous behavior detection, keyframe, feature pyramid, keypoint detection, feature aggregation

1. INTRODUCTION

Recently, video surveillance systems are common to discover dangerous behaviors timely in school and shelter students from potential harm. Therefore, an effective and efficient monitoring alarm method for student dangerous behaviors is vital and urgent. However, student dangerous behaviors recognition has numerous challenges due to insufficient datasets, confusing postures, keyframes detection, and prompt response. In this paper, we focus on the primary school scene which is full of pupils who might behave dangerously because of their immaturity. Three typical potential dangerous behaviors, including fighting, tumbling and squatting, are automatically detected, as shown in Fig. 1.

Generally, dangerous behavior detection belongs to action recognition. The representative methods in this area include Two-stream convolutional networks [1], Temporal Segment Networks (TSN) [2] and YOWO [3]. However, the common disadvantages of these approaches include the expensive computation and consecutive frames analysis, which

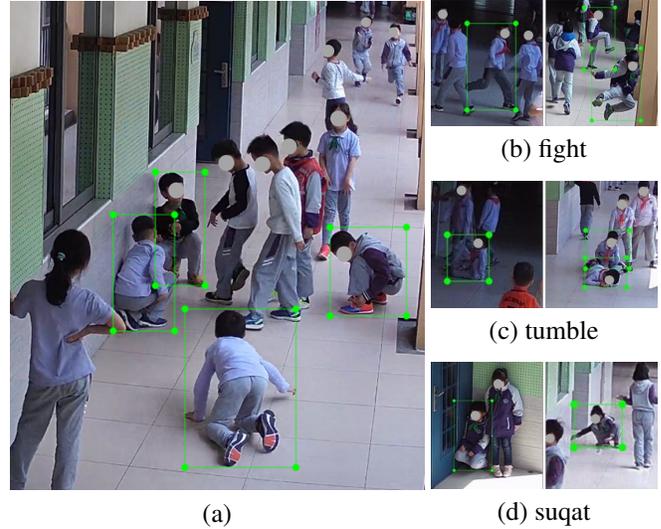


Fig. 1. Figure (a) shows the status of students in one corridor after class. Figure (b), (c) and (d) are annotation instances of three dangerous behaviors.

are unfriendly to the practical applications. To satisfy the real applications, instead of applying SOTA yet cumbersome methods in the action recognition area, we propose a novel idea that combines lightweight keyframe feature extraction network and 2D human pose features to fast localize and more accurately recognize dangerous behaviors. Both of object detection and pose estimation have been used for action recognition independently, but rarely together. Several works with similar opinions to ours were proposed, such as UAV surveillance system [4], Drone Surveillance System (DSS) [5] and student behavior recognition system [6], which cannot be trained end-to-end.

Several researches have demonstrated the effectiveness and efficiency of specific behaviors recognition using object detection methods. Zheng et. al [10] successfully applied improved Faster R-CNN [11] to the student behavior analysis, including hand-raising, standing and sleeping. After that, Zheng et. al [12] utilized lightweight network MobileNetV2-SSD [13] with multi-dimensional attention mechanisms to significantly reduce the computational cost. However, they

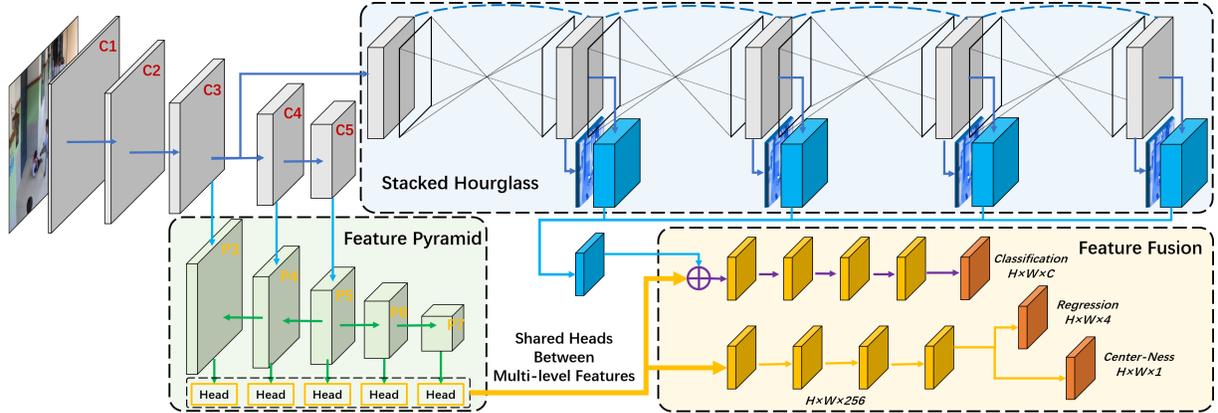


Fig. 2. The network architecture of DangerDet. We use $C_1 \sim C_5$ to denote the feature maps from the FPN [7] based on ResNet [8]. After C_3 are four stacked hourglasses [9] to extract pose features for better behavior classification. Feature levels $P_3 \sim P_7$ in feature pyramid are used for final prediction. Best viewed on screen.

ignore the pose features that are highly related to the behaviors. Meanwhile, EduSense [14] adopted OpenPose [15] for the behavior recognition of teachers and students in classrooms. Due to the highly occlusion among students, the real performances based on the 2D keypoints are obviously reduced. Compared with pose estimation methods, object detection algorithms for behavior recognition are less sensitive to the occlusion. Therefore, we propose a novel framework for the dangerous behavior recognition with object detection architecture appending pose estimation as the backbone.

Specifically, we firstly built a particular danger behavior dataset under school corridor scene. Then, we proposed the DangerDet for tackling dangerous behavior detection. The network architecture of DangerDet is shown in Fig. 2. We transform action recognition problem into behavior detection. The detection framework draws on the idea from FPN [7] and the single-stage anchor-free method FCOS [16] for its good balance between accuracy and time-consuming. And in order to use the information of keypoints to improve the accuracy of behavior classification, the backbone network introduces the versatile stacked hourglass [9] structure. In experiments, we verified the robustness and superiority of DangerDet compared with baseline on our self-made dataset.

2. RELATED WORK

In this section, we will briefly introduce the works related to deep feature extraction and keypoint heatmap generation involved in our DangerDet.

2.1. Deep Feature Extraction

The deep feature extraction based on convolutional neural network (CNN) and correlation strategies is an essential part of advanced object recognition algorithms [11, 17, 13, 7, 16].

One of the most important ideas to deal with the object scale variation is the feature pyramid strategy. This method is firstly proposed in FPN [7] which excavated hierarchical features from discretized image pyramid and made significant improvement in object recognition. Then, based on this foundation work, many similar pyramid structures with multi-scale lightweight feature maps and subtle skip connections architecture are invented. The effective feature pyramid in FCOS [16] is a commonly used framework. For our DangerDet, we will adopt a homologous pyramid architecture for feature extraction to better and faster detect various dangerous behaviors. The legend is depicted in Fig. 2.

2.2. Keypoint Heatmap Generation

Here we introduce two representative pose estimation algorithms Stacked Hourglass [9] and PifPaf [18] which are conducive to our work. The hourglass network structure is firstly proposed in [9] to utilize multi-scale features produced by residual module in ResNet [8] to recognize person pose. By stacking multiple such hourglasses blocks, we could reuse the information of whole-body joints to further improve the recognition accuracy. Generally, the four stacked hourglasses configuration is the most cost-effective and thus embedded into our architecture to generate keypoint heatmaps for enhancing our behavior detection task. As for PifPaf [18], it is an excellent open-source bottom-up multi-person pose estimation algorithm that outperforms previous methods under various scenes including low resolution, crowd and occlusion thanks to (i) the new composite field PAF encoding fine-grained information and (ii) the choice of Laplace loss for regressions which incorporates a notion of uncertainty. We directly use PifPaf to estimate the human pose of all the keyframes in our self-made dataset. The obtained keypoints will be used to generate keypoints heatmaps as the weak ground-truths to support the supervised training.

3. OUR METHOD

In this section, we will first explain the network structure of our proposed method in detail. Then, the proposed feature aggregation strategy about multiple branches is illustrated. Finally, we will present the design of loss function about outputted features by multi-channel.

3.1. Network Architecture

The overall architecture of our DangerDet is the combination of both object detection and pose estimation. In Fig. 2, the input image shape is $(1152, 768, 3)$. After passing \mathcal{C}_1 and \mathcal{C}_2 , the size of feature map is down-sampled $4(2^2)$ times. Finally, the size of the feature map outputted by \mathcal{C}_5 is reduced to $32(2^5)$ times. Based on blocks $\mathcal{C}_3 \sim \mathcal{C}_5$, we construct the feature pyramid $\mathcal{P}_3 \sim \mathcal{P}_7$ following FPN [7]. We share the heads between different feature levels of $\mathcal{P}_i (i \in [3, 4, 5, 6, 7])$ for multi-level prediction. To regress different size range among \mathcal{P}_i , we increase a trainable scalar s_i to automatically adjust the base of the exponential function $exp(s_i x)$ for feature level \mathcal{P}_i . Besides, right after \mathcal{C}_3 , we add four stacked hourglass modules connected in series. The intermediate supervision and symmetric distribution of capacity strategy in [9] are maintained. These modules generate multiple keypoints heatmaps and pose features which are consistent with the output of \mathcal{C}_3 in size. Among them, the keypoints heatmap is used for supervised learning, and its dimension is K which denotes the keypoints number. And pose features are abstracted for subsequent behavior classification.

After shared heads, the final outputs have two pipelines and three branches similar to FCOS [16]. The difference is that we have only three categories ($C = 3$) in the *Classification* branch. Moreover, we fuse all feature maps outputted by \mathcal{P}_i with the combined pose feature map from the stacked hourglass module to obtain new augmented feature maps as inputs separately. For *Regression* branch, it outputs a 4D vector (l, t, r, b) encoding the location of a bounding box at each foreground pixel. The *Center-Ness* branch depicts the normalized distance from the location to the center of the object that the location is responsible for. It is invented to down-weight the scores of bounding boxes far from the center of an object. Given the regression targets l^*, t^*, r^* , and b^* for a location, it is defined in Eqn. 1 as below:

$$CenterNess^* = \sqrt{\frac{\min(l^*, r^*)}{\max(l^*, r^*)} \times \frac{\min(t^*, b^*)}{\max(t^*, b^*)}} \quad (1)$$

3.2. Feature Aggregation

In this part, we focus on the aggregation of multi-level features extracted from FPN and pose features predicted by stacked hourglasses. In the *Classification* prediction branch,

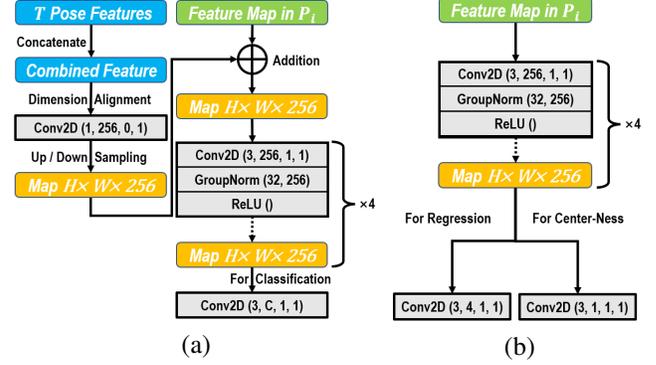


Fig. 3. Illustrations of feature extraction and aggregation in DangerDet. a) *Classification*. b) *Regression* and *Center-Ness* [16]. Please note that the format of Conv2D is (*kernel_size, out_channels, padding, stride*).

as shown in Fig. 3(a), the input map is obtained by element-wise addition of two kinds of features. One is feature map in \mathcal{P}_i . The another one is transformed from multiple pose features by a series of sequential operations including concatenation, dimension alignment and up(down)-sampling. Then, we extract the high-level semantics through four 3×3 Conv2D layers with 256 filters. Layers GroupNorm and ReLU follow closely each convolution layer. Finally, through a 3×3 Conv2D layer with C ($C = 3$) filters, this branch outputs a result map with behavior category in pixel level. In Fig. 3(b), we maintain the same feature extraction structure as the classification branch. However, pose features are removed. The number of output channels of *Regression* and *Center-Ness* branched are changed into 4 and 1 accordingly. The outputs of these two branches are also heatmaps in pixel level. In the post-processing, these result heatmaps will be decoded to obtain behavior types and bounding-boxes.

3.3. Loss Function

Our DangerDet has two major losses: (i) For object detection part, the losses are L_{cls} , L_{reg} , and L_{cen} for *Classification*, *Regression* and *Center-Ness* respectively. (ii) For pose estimation part, the loss is the MSE for keypoints heatmap named L_{kpt} . The final training loss function is defined as Eqn. 2:

$$L(\{p_{x,y}\}, \{t_{x,y}\}) = \frac{1}{N_{pos}} \sum_{x,y} L_{cls}(p_{x,y}, c_{x,y}^*) + \frac{\lambda}{N_{pos}} \sum_{x,y} 1_{c_{x,y}^* > 0} L_{reg}(t_{x,y}, t_{x,y}^*) + \sum_{t=1}^T L_{kpt}(C_t^r, C^r) \quad (2)$$

where L_{cls} is focal loss [19] and L_{reg} is the IOU loss as in UnitBox [20]. N_{pos} denotes the number of positive samples. λ is set to 1. C_t^r in L_{kpt} is the predicted joints confidence map. C^r is the ground-truth keypoint heatmap. T is the number of hourglass modules which is set as 4.

Table 1. Comparison between DangerDet and other methods.

Method	resnet	pose	T	AP_{50}	AP_{75}	mAP
FCOS	50	✗	0	83.7	69.5	59.5
DangerDet	50	✓	4	89.6	81.3	69.2
FCOS*	101	✗	0	89.4	78.4	67.6
DangerDet*	101	✓	4	85.3	82.6	71.0
DangerDet1	50	✓	1	86.1	69.2	60.5
DangerDet2	50	✓	2	87.8	72.6	62.9

4. EXPERIMENTS

In this section, we will describe the construction process of our danger behavior dataset, and then report the results of ablation experiments on this dataset.

4.1. Our Danger Behavior Dataset

Our raw data is collected from a primary school in Shanghai, China. It consists of about 100 surveillance videos captured by six 4K (3840×2160) cameras installed in a corridor without overlapping. These cameras start recording student activities during a pre-set period of time, about ten minutes after class. Fig. 1(a) gives an example frame. Then we sampled these videos every 0.5 seconds (skipping about 15 frames with FPS=30) and produced latent keyframes of our danger behavior dataset. After removing irrelevant monitoring areas, we cropped all the frames to $2400 \times 1600(3 : 2)$. We used LabelImg [21], a popular open-source annotation tool in object detection, to finish the annotation work.

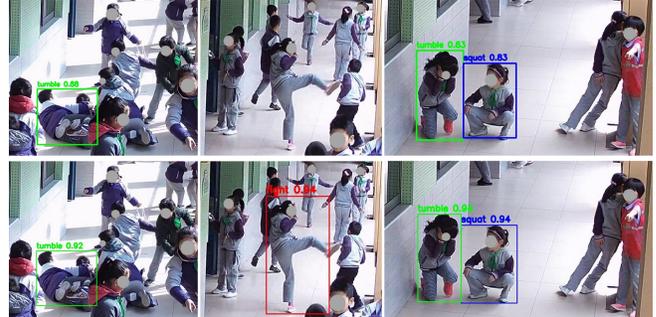
Specifically, we selected three representative dangerous behavior categories to label. They are **fight**, **tumble**, and **squat**. These behaviors are the precursor of potential danger, or themselves are dangerous actions, which need early warning in time once happening. Some examples of annotations are shown in Fig. 1(b)-(d). In addition, we have directly applied PifPaf [18] to estimate pose of all annotated frames and generate rough keypoints coordinates of all students in our dataset. Finally, the dataset includes **7161** static images, and the number of three behaviors (fight, tumble, squat) is **665**, **6962**, and **1788**, respectively.

4.2. Ablation Experiments

We divided our dataset into training and validation set in the ratio 8:2. On the validation set, the DangerDet was evaluated by measuring precision, recall, and mean average precision (mAP) at an intersection-over-union (IOU) threshold from 0.5 to 0.95. For comparison, we trained the original FCOS model as the baseline without the pose module or L_{kpt} part in loss function. The selected backbones include ResNet-50 and ResNet-101. Besides, we investigated the influence of different number T of stacked hourglass modules in DangerDet.

Table 2. Three behaviors' detection results of DangerDet*.

Behavior	fight	tumble	squat
Labels	132	1385	362
mAP	68.0	71.9	73.1

**Fig. 4.** Detection results of baseline FCOS (top column) and our DangerDet (bottom column) with IoU threshold = 0.8.

The validation results of DangerDet and other methods for comparison are shown in the table 1. With ResNet-50 and ResNet-101 as the backbone, our methods DangerDet(*) are 9.7% and 3.4% points higher than the baselines FCOS(*) in mAP, respectively. With ResNet-50 as the backbone, our methods DangerDet/1/2 with $T = 1$, $T = 2$ and $T = 4$ always perform better than the baseline FCOS. With lighter backbone, DangerDet1 obtains larger mAP than FCOS*. All these proved the superiority of our method. With $T = 4$ and the ResNet-101 as backbone, DangerDet* achieved the best result. Some qualitative results are shown in Fig. 4.

Table 2 shows the detection results of three behaviors using DangerDet*. Unsurprisingly, the action fight gets worse accuracy than the other two behaviors for its higher dependency on continuous spatiotemporal features. We expect that this situation might be improved by collecting more data or associating continuous frames.

5. CONCLUSION

We presented a novel method DangerDet which integrates both object detection and pose estimation for detecting three dangerous behaviors of students on campus. The proposed end-to-end network structure allows us to easily combine deep features and two-dimensional keypoints. By designing efficient feature aggregation methods and adjusting the loss function, the training robustness and final accuracy of DangerDet are improved. The integration of these strategies achieves an impressive result in dangerous behaviors detection of the real school scenario. Our method has potential value to be applied to campus security alarm system. We have released our codes for academic use in https://github.com/hnuzhy/fcos_pose.

6. REFERENCES

- [1] Karen Simonyan and Andrew Zisserman, “Two-stream convolutional networks for action recognition in videos,” *NIPS*, vol. 27, pp. 568–576, 2014.
- [2] 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 *ECCV*. Springer, 2016, pp. 20–36.
- [3] Okan Köpüklü, Xiangyu Wei, and Gerhard Rigoll, “You only watch once: A unified cnn architecture for real-time spatiotemporal action localization,” *arXiv preprint arXiv:1911.06644*, 2019.
- [4] Surya Penmetsa, Fatima Minhuj, Amarjot Singh, and SN Omkar, “Autonomous uav for suspicious action detection using pictorial human pose estimation and classification,” *ELCVIA: electronic letters on computer vision and image analysis*, vol. 13, no. 1, pp. 0018–32, 2014.
- [5] Amarjot Singh, Devendra Patil, and SN Omkar, “Eye in the sky: Real-time drone surveillance system (dss) for violent individuals identification using scatternet hybrid deep learning network,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2018, pp. 1629–1637.
- [6] Feng-Cheng Lin, Huu-Huy Ngo, Chyi-Ren Dow, Ka-Hou Lam, and Hung Linh Le, “Student behavior recognition system for the classroom environment based on skeleton pose estimation and person detection,” *Sensors*, vol. 21, no. 16, pp. 5314, 2021.
- [7] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie, “Feature pyramid networks for object detection,” in *CVPR*, 2017, pp. 2117–2125.
- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, “Deep residual learning for image recognition,” in *CVPR*, 2016, pp. 770–778.
- [9] Alejandro Newell, Kaiyu Yang, and Jia Deng, “Stacked hourglass networks for human pose estimation,” in *ECCV*. Springer, 2016, pp. 483–499.
- [10] Rui Zheng, Fei Jiang, and Ruimin Shen, “Intelligent student behavior analysis system for real classrooms,” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 9244–9248.
- [11] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” *NIPS*, vol. 28, pp. 91–99, 2015.
- [12] Rui Zheng, Fei Jiang, and Ruimin Shen, “Gesture-det: Real-time student gesture analysis with multi-dimensional attention-based detector,” in *IJCAI*, 2020, pp. 680–686.
- [13] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen, “Mobilenetv2: Inverted residuals and linear bottlenecks,” in *CVPR*, 2018, pp. 4510–4520.
- [14] Karan Ahuja, Dohyun Kim, Francesca Xhakaj, Virag Varga, Anne Xie, Stanley Zhang, Jay Eric Townsend, Chris Harrison, Amy Ogan, and Yuvraj Agarwal, “Edusense: Practical classroom sensing at scale,” *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 3, no. 3, pp. 1–26, 2019.
- [15] Zhe Cao, Gines Hidalgo, Tomas Simon, Shih-En Wei, and Yaser Sheikh, “Openpose: realtime multi-person 2d pose estimation using part affinity fields,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 43, no. 1, pp. 172–186, 2019.
- [16] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He, “Fcos: Fully convolutional one-stage object detection,” in *ICCV*, 2019, pp. 9627–9636.
- [17] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun, “R-fcn: Object detection via region-based fully convolutional networks,” in *NIPS*, 2016, pp. 379–387.
- [18] Sven Kreiss, Lorenzo Bertoni, and Alexandre Alahi, “Pifpaf: Composite fields for human pose estimation,” in *CVPR*, 2019, pp. 11977–11986.
- [19] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár, “Focal loss for dense object detection,” in *ICCV*, 2017, pp. 2980–2988.
- [20] Jiahui Yu, Yuning Jiang, Zhangyang Wang, Zhimin Cao, and Thomas Huang, “Unitbox: An advanced object detection network,” in *Proceedings of the 24th ACM international conference on Multimedia*, 2016, pp. 516–520.
- [21] D Tzutalin, “Labelimg (2015),” *GitHub repository <https://github.com/tzutalin/labelImg>*, vol. 6.