

Motivation

Overview: *DFD* applies different treatments to regions with varying learning difficulties, simultaneously incorporating leniency and strictness, which enables better distillation efforts.

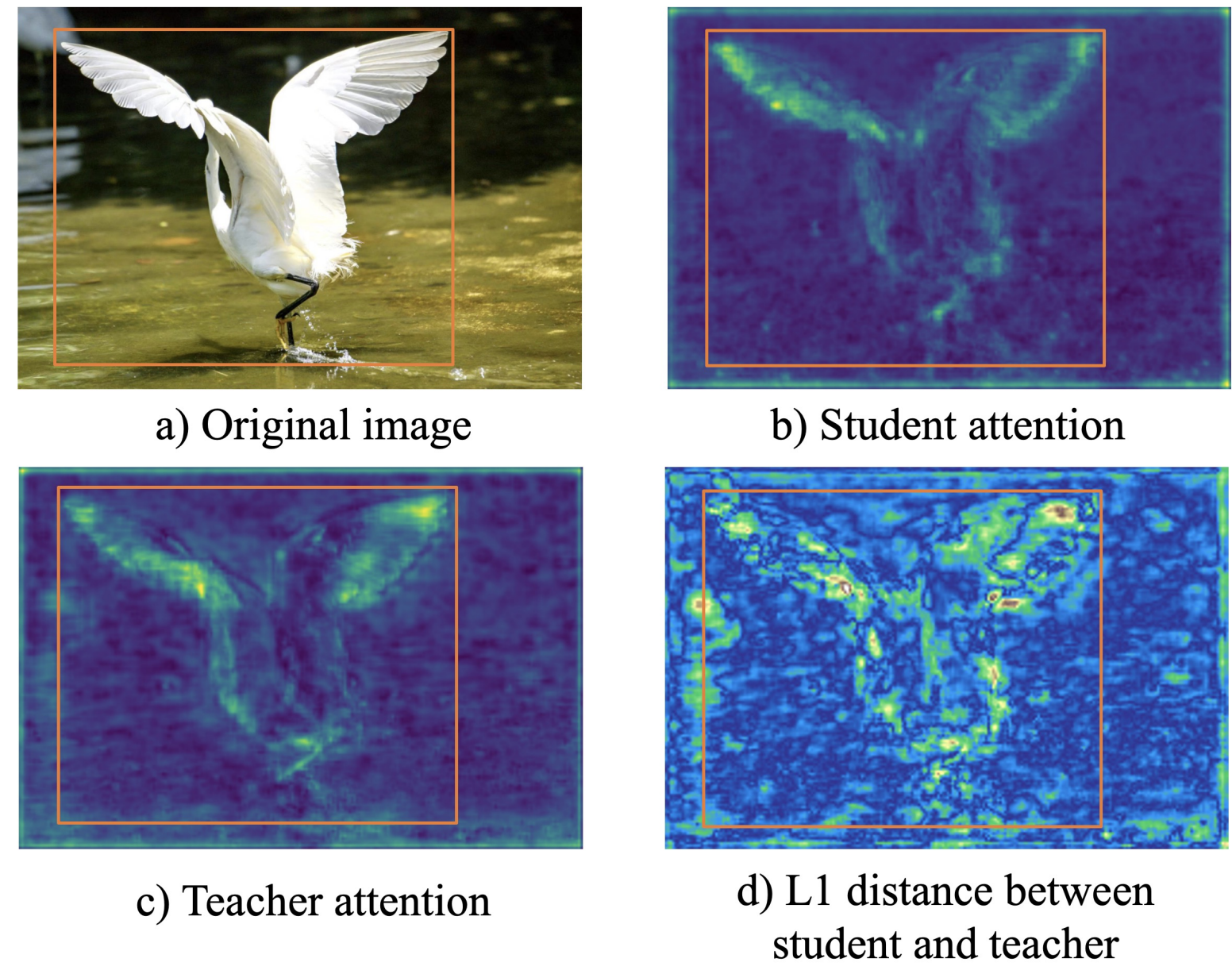


Figure 1. Visualization of spatial attention.

- **Region:** The feature responses do not fully conform to the distinction between foreground and background.
- **Disparity:** There are complex differences in both the regions and intensities of the features between the teacher and student models.

Table 1. Comparison of distillation on different regions.

Model	HD	LD	Split	DC	mAP
RetinaNet RX101-R50	-	-	-	-	36.4
	✓	-	-	-	39.8
	-	✓	-	-	38.2
	✓	✓	-	-	39.6
	✓	✓	✓	-	40.0
	✓	✓	✓	✓	40.4 (Ours)

- **HD:** High disparity region
- **LD:** Low disparity region.
- **Split:** Split these regions, and use different weights for the distillation losses of different parts.
- **DC:** Using different constraints for different regions.

- **Knowledge Gap:** The feature disparity represents the varying capabilities of the models.
- **Distillation effort:** Using different distillation methods for different disparity regions can significantly enhance the distillation effect.

Methodology

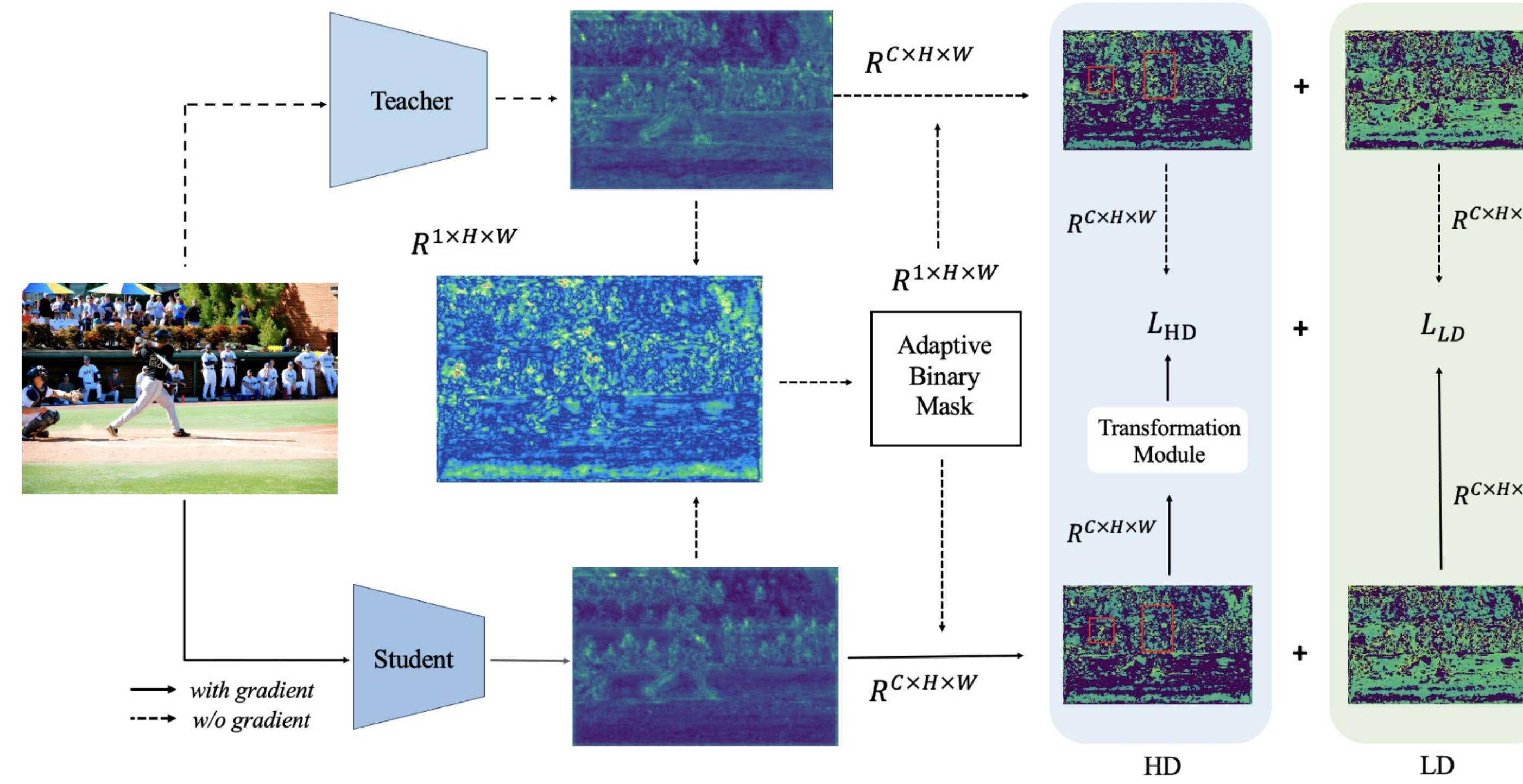


Figure 2: Overall architecture of DFD. We partition the feature into high disparity regions (HD) and low disparity regions (LD) by utilizing the mean L1 distance of spatial attention between the student and teacher as the threshold.

- Calculate the spatial attention:

$$A(F) = H \cdot W \cdot \text{softmax} \left(\frac{1}{C} \cdot \sum_{c=1}^C |F| \right)$$

- Calculate the attention disparity map:

$$D = |A(F^T) - A(F^S)|$$

- Generate adaptive binary mask:

$$Mask_{i,j} = \begin{cases} 0, & \text{if } D_{i,j} < P \\ 1, & \text{Otherwise} \end{cases}$$

- Distillation losses:

$$L_{LD} = \sum_k^C \sum_i^H \sum_j^W (F_{i,j,k}^T - F_{i,j,k}^S)^2, i, j \in R_{LD}$$

$$L_{HD} = \sum_k^C \sum_i^H \sum_j^W (F_{i,j,k}^T - f_{trans}(F_{i,j,k}^S))^2, i, j \in R_{HD}$$

$$L_{total} = L_{task} + \alpha \cdot L_{HD} + \beta \cdot L_{LD}$$

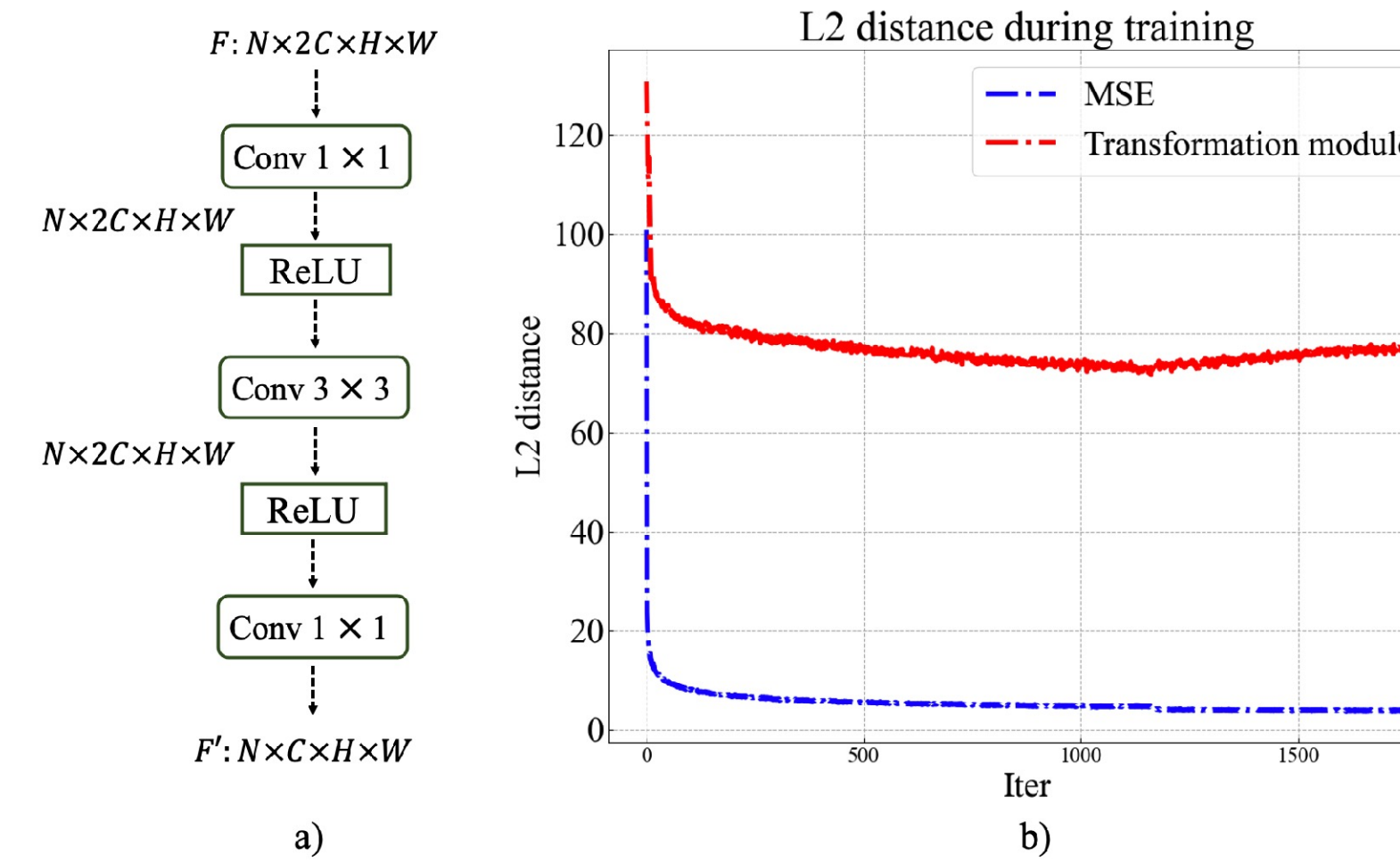


Figure 3. (a). The structure of the transformation module we used. (b). L2 distance during the training with transformation module.

Results and Ablation Study

Table 2. Main results on COCO dataset.

Teacher	Method	schedule	mAP	AP s	AP m	AP L
RetinaNet ResNext101	RetinaNet ResNet50	2x	37.4	20.0	40.7	49.7
	FGD	2x	40.9	23.1	45.1	54.9
	MGD	2x	41.2	23.6	45.3	54.6
	PKD	2x	41.2	23.0	45.4	55.6
	Ours	2x	41.7(+4.0)	24.4	46.0	55.7
CascadeMaskRCNN ResNeXt 101	FasterRCNN ResNet 50	2x	38.4	21.5	42.1	50.3
	FGD	2x	42.0	23.8	46.4	55.5
	MGD	2x	42.1	23.7	46.4	56.1
	PKD	2x	41.4	22.7	45.1	56.0
	Ours	2x	42.4(+4.0)	23.9	46.8	56.3
Reppoints ResNeXt101	Reppoints ResNet 50	2x	38.6	22.5	42.2	50.4
	FGD	2x	42.0	24.0	45.7	55.6
	MGD	2x	42.3	24.4	46.2	55.9
	PKD	2x	42.4	24.3	46.7	56.4
	Ours	2x	42.7(+4.1)	24.8	46.7	56.2

Table 4. combining DFD with other distillation methods.

Model	Method	Schedule	mAP
RetinaNet RX101	MGD	1x	40.0
	MGD + Ours	1x	40.3
- RetinaNetR50	PKD	1x	39.9
	PKD + Ours	1x	40.3
	PKD + MGD	1x	39.9

Table 3. Experiments with progressively stronger teacher on COCO dataset.

Teacher	Method	Schedule	mAP
FasterRCNN R101 39.8	<i>student</i>	1x	37.4
	PKD	1x	39.6
	Ours	1x	39.7
FasterRCNN Rx101 41.2	PKD	1x	40.0
	Ours	1x	40.5
MaskRCNN Rx101 42.2	PKD	1x	40.5
	Ours	1x	41.1

Table 5. Using different constraints in different regions.

HD	LD	Schedule	mAP
MSE	MSE	1x	39.7
TF	TF	1x	40.2
TF	-	1x	39.9
-	TF	1x	39.4
MSE	TF	1x	39.9
TF	TF + MSE	1x	40.2
TF	MSE	1x	40.4

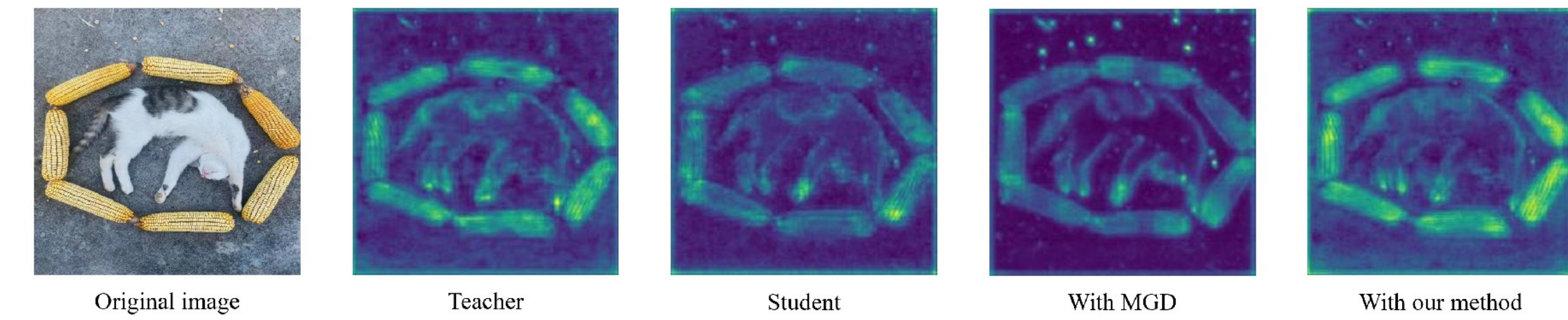


Figure 4. Visualization of spatial attention after distillation.

Table 6. Experiments of other detectors.

Method	Schedule	mAP
MaskRCNN-SwinS(Teacher)	1x	48.2
MaskRCNN-SwinT(Student)	1x	42.7
PKD	1x	43.9
Ours	1x	44.4
Method	Epoch	mAP
YOLOv6-Small(Teacher)	400	44.0
YOLOv6-Tiny(Student)	300	40.6
PKD	300	41.3
Ours	300	41.7

Table 7. Results of pose estimation on COCO-Body and segmentation on Cityscapes

Pose estimation	Method	Input Size	mAP
Heatmap Res50	Teacher	256 × 192	71.8
Heatmap MobileNetV2	<i>student</i>	256 × 192	62.0
	CWD	256 × 192	62.2
	Ours	256 × 192	62.6
Segmentation	Method	Input Size	mAP
PspNet Res101	Teacher	512 × 512	78.34
PspNet Res18	<i>student</i>	512 × 512	69.85
	CWD	512 × 512	73.53
	Ours	512 × 512	73.74