

## Motivation

Recent single-stage detectors suffer from severe inconsistency problems:

- Classifier is confused by **misaligned classification and localization** due to the unreasonable IoU threshold.
- **Feature inconsistency**: the refined anchors are associated with the feature extracted from the previous location.

## Analysis in Inconsistency

Misaligned Classification and Localization:

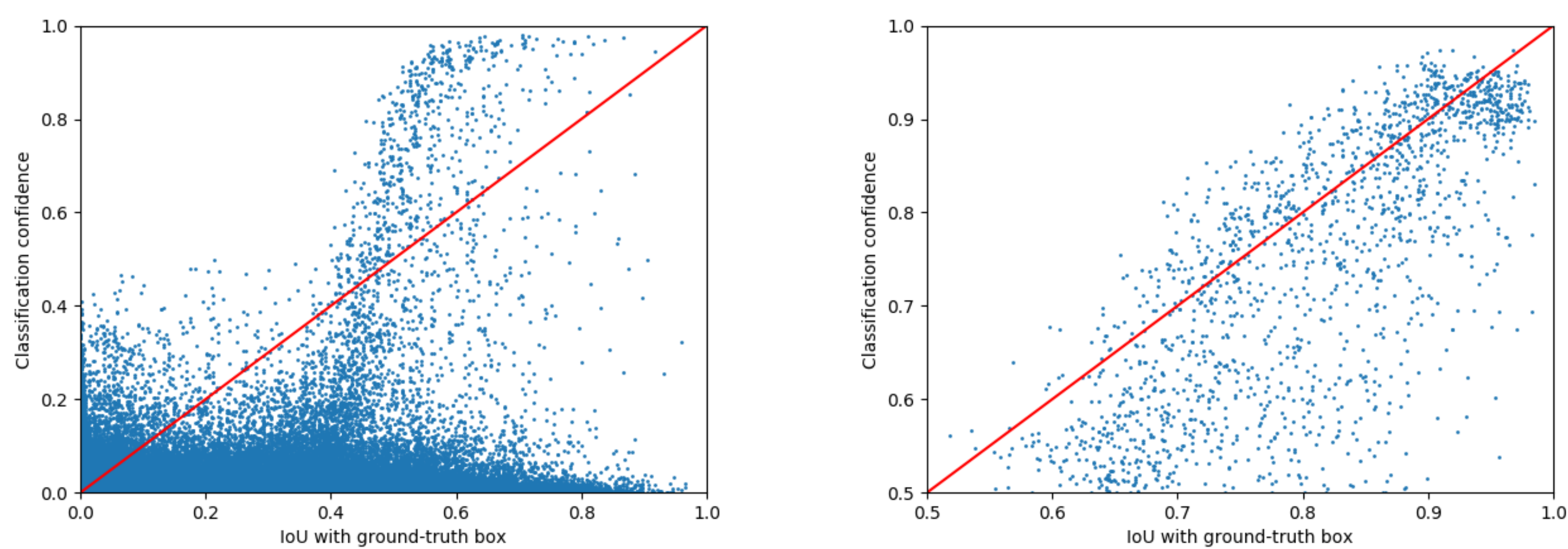


Fig 1: Classification confidence vs. IoU in different cascade stages.

- (left) Classification scores are **not well aligned** with the IoU for the first stage, especially for the confidences near IoU@0.5.
- (right) Improved consistency between classification and regression in the second stage using increased IoU threshold.

Feature Inconsistency:

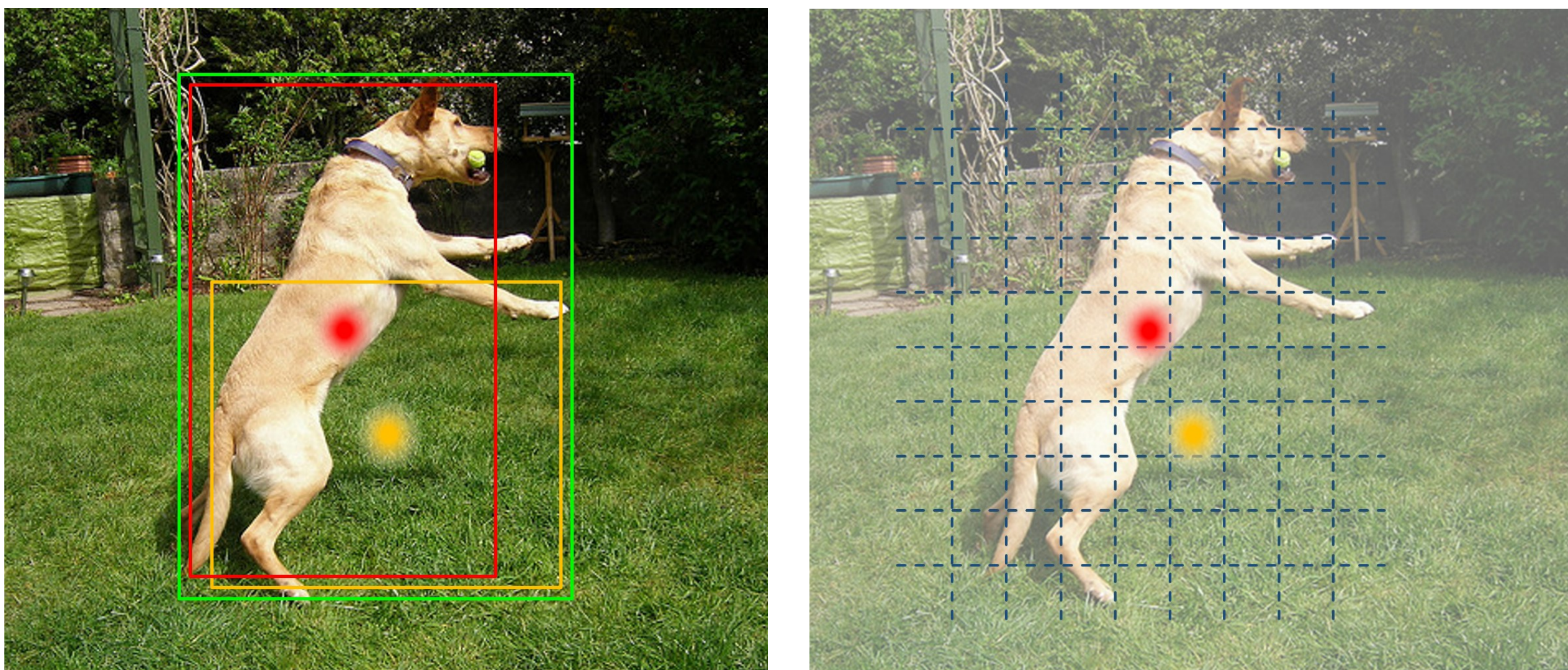


Fig 2: Feature misalignment between original and refined anchor.

- (left) Original image. The green box stands for the ground truth and the orange one represents the original anchor. The refined anchor is shown as the red bounding box.
- (right) Feature grid. Locations of center points for original and refined anchors are plot. It indicates that **simply extracting features from the previous location (orange point) is inaccurate**.

## Proposed Method

Designing rules for the cascade manner:

- Improving consistency between classification confidence and localization performance.
- Maintaining feature consistency between different stages.

Cascade RetinaNet:

- **Gradually increase the foreground IoU thresholds** to maintain the consistency between classification and localization.
- Encode the current localization information into the features of next stage by **Feature Consistency Module (FCM)**.

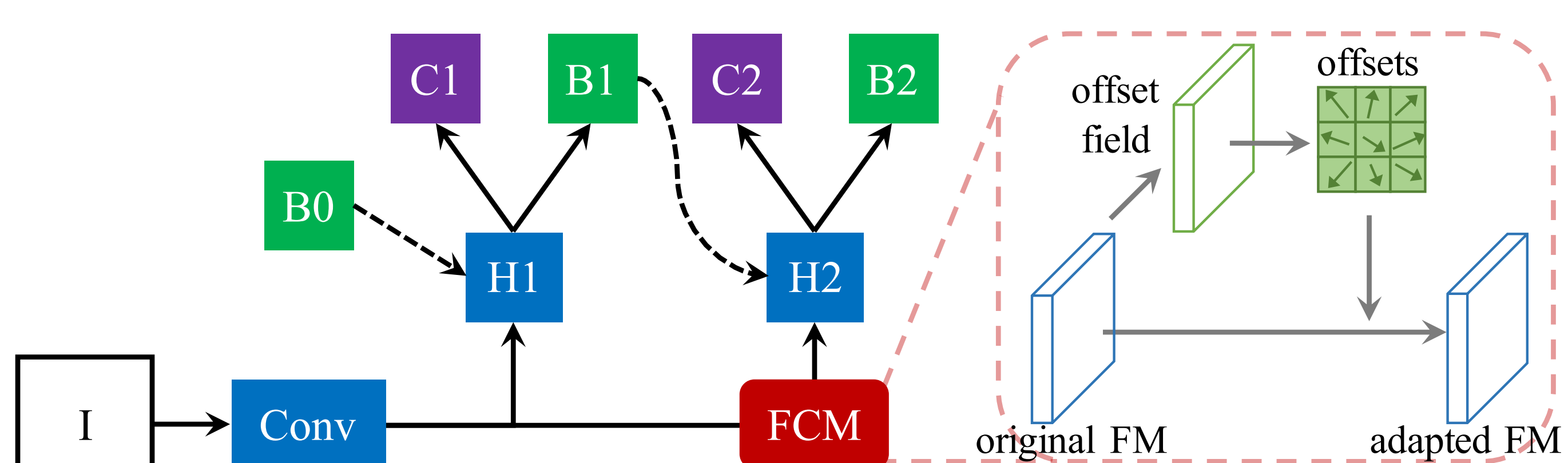


Fig 3: The overall architecture of Cascade RetinaNet.

## Experiments

Ablation study:

| Method         | Scale | IoU | AP          | AP <sub>50</sub> | AP <sub>60</sub> | AP <sub>70</sub> | AP <sub>80</sub> | AP <sub>90</sub> |
|----------------|-------|-----|-------------|------------------|------------------|------------------|------------------|------------------|
| RetinaNet [20] | 600   | -   | 34.0        | 52.5             | -                | -                | -                | -                |
| Cas-RetinaNet  | 600   | 0.5 | 33.8        | 52.3             | 48.1             | 41.5             | 29.8             | 11.2             |
| Cas-RetinaNet  | 600   | 0.6 | 34.4        | 52.5             | 48.5             | 41.9             | 30.5             | 11.7             |
| Cas-RetinaNet  | 600   | 0.7 | 34.4        | 52.0             | 48.1             | 41.7             | 31.3             | <b>12.5</b>      |
| RetinaNet [20] | 800   | -   | 35.4        | 53.9             | -                | -                | -                | -                |
| Cas-RetinaNet  | 800   | 0.5 | 35.4        | 54.6             | 50.4             | 43.0             | 31.4             | 11.8             |
| Cas-RetinaNet  | 800   | 0.6 | <b>36.1</b> | <b>55.0</b>      | <b>50.8</b>      | <b>43.9</b>      | <b>32.5</b>      | <b>12.5</b>      |

Tab 1: Ablation study on IoU thresholds.

- Naively adding a new stage with the same setting **brings no gains**.
- **Increasing the foreground IoU** for the second stage is beneficial since it leads to a more consistent distribution.

| Backbone   | Scale | FCM | AP          | AP <sub>50</sub> | AP <sub>60</sub> | AP <sub>70</sub> | AP <sub>80</sub> | AP <sub>90</sub> |
|------------|-------|-----|-------------|------------------|------------------|------------------|------------------|------------------|
| ResNet-50  | 600   |     | 34.4        | 52.5             | 48.5             | 41.9             | 30.5             | 11.7             |
| ResNet-50  | 600   | ✓   | 35.5        | 54.0             | 49.7             | 43.3             | 32.0             | 12.6             |
| ResNet-50  | 800   |     | 36.1        | 55.0             | 50.8             | 43.9             | 32.5             | 12.5             |
| ResNet-50  | 800   | ✓   | 37.1        | 56.3             | 52.2             | 45.3             | 33.5             | 12.8             |
| ResNet-101 | 800   |     | 37.9        | 56.8             | 52.8             | 46.0             | 34.9             | 13.9             |
| ResNet-101 | 800   | ✓   | <b>38.9</b> | <b>58.1</b>      | <b>53.9</b>      | <b>47.1</b>      | <b>36.2</b>      | <b>14.3</b>      |

Tab 2: Ablation study on FCM.

- Steadily improvements under different settings are achieved due to the effectiveness of the **adapted feature produced by FCM**.

| #Stages | Test stage | AP          | AP <sub>50</sub> | AP <sub>60</sub> | AP <sub>70</sub> | AP <sub>80</sub> | AP <sub>90</sub> |
|---------|------------|-------------|------------------|------------------|------------------|------------------|------------------|
| 1       | 1          | 34.0        | 52.5             | -                | -                | -                | -                |
| 2       | 1 ~ 2      | <b>35.5</b> | <b>54.0</b>      | <b>49.7</b>      | <b>43.3</b>      | <b>32.0</b>      | 12.6             |
| 3       | 1 ~ 2      | 35.0        | 53.1             | 49.1             | 42.5             | <b>32.0</b>      | 12.6             |
| 3       | 1 ~ 3      | 34.9        | 52.9             | 49.0             | 42.4             | 31.9             | <b>12.7</b>      |

Tab 3: Ablation study on the number of stages.

- Cascading three stages leads to a slight drop in the overall performance while achieves the best for high IoU. It is a **trade-off between sample quality and quantity** as mentioned in Cascade R-CNN [1].

Overall performances:

| Method                 | Backbone      | AP          | AP <sub>50</sub> | AP <sub>75</sub> | AP <sub>S</sub> | AP <sub>M</sub> | AP <sub>L</sub> |
|------------------------|---------------|-------------|------------------|------------------|-----------------|-----------------|-----------------|
| RetinaNet [20]         | ResNet-50     | 35.7        | 55.0             | 38.5             | 18.9            | 38.9            | 46.3            |
| RefineDet512 [35]      | ResNet-101    | 36.4        | 57.5             | 39.5             | 16.6            | 39.9            | 51.4            |
| GA-RetinaNet [3]       | ResNet-50     | 37.1        | 56.9             | 40.0             | 20.1            | 40.1            | 48.0            |
| RetinaNet [20]†        | ResNet-101    | 39.1        | 59.1             | 42.3             | 21.8            | 42.7            | 50.2            |
| ConRetinaNet [16]†     | ResNet-101    | 40.1        | 59.6             | 43.5             | 23.4            | 44.2            | 53.3            |
| CornerNet511 [17]      | Hourglass-104 | 40.5        | 56.5             | 43.1             | 19.4            | 42.7            | <b>53.9</b>     |
| <b>Cas-RetinaNet</b>   | ResNet-50     | 37.4        | 56.6             | 40.7             | 20.9            | 40.3            | 47.5            |
| <b>Cas-RetinaNet</b> † | ResNet-101    | <b>41.1</b> | <b>60.7</b>      | <b>45.0</b>      | <b>23.7</b>     | <b>44.4</b>     | 52.9            |

Tab 4: Overall performances on COCO minival set.

## Conclusion

- Analysis shows that **inconsistency in single-stage detectors** is the key factor limiting the detection performance.
- Two main designing rules for **maintaining consistency** are proposed: *improving consistency between classification and localization*, and *maintaining feature consistency between different stages*.
- **Cascade RetinaNet**, a simple but effective architecture, can maintain the consistency by increasing thresholds and adopting FCM, which leads to improved detection performance.