Motivation

Detection in varying scale and sparsity scenarios ?

- Two-stage detectors (e.g. Faster-RCNN [1]) \rightarrow exhaustive but slow
- One-pass grid-based models (e.g. YOLO [2]) \rightarrow Input resolution must account for small objects

• Trade-off accuracy vs memory/computational efficiency (e.g., applications involving embedded systems)



High density and saliency \leftarrow - - - - - - - - - - - Sparse and small objects

Objective: Design a detection scheme that makes use of group structures in the image to curate few, meaningful, proposals

One-stage detectors

- In each cell in the (I, J) grid, output **K** boxes B'_{k}
- In empty cells, $c_k^{i,j} \rightarrow 0$

• In active cells, assignment $A(B_k^{i,j}, b^*) \in \{0, 1\}.$ $c_k^{i,j} \rightarrow iou(B_k^{i,j}, b^*)$ and match coordinates

[1] S. Ren et al., "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS, 2015. [2] J. Redmon et al, "You Only Look Once: Unified, Real-Time Object Detection". CVPR, 2016.

Localizing Grouped Instances for Efficient Detection in Low-Resource Scenarios Amélie Royer, Christoph H. Lampert

Model Overview



Based on **YOLO**: fast (one-pass) yet dense (grid-based)

Proposed Solution

How is the assignment A defined ? $A(B_k^{i,j}, b^*) = [|b^* \text{ in cell } (i,j)|] [|k = \arg \max_{k'} \operatorname{iou}(B_{k'}^{i,j}, b^*)|]$

The total loss gathers contributions from all empty cells (push confidence scores to 0) and active cells (match assigned bounding boxes to target coordinates and confidences to iou)

Shortcomings:

- Grid parameters (I, J, K) need to match the data resolution
- or too low (assignment collisions)

Proposed solution, ODGI:

Define groups of objects as intermediary structures that are easier to detect at low resolution and can be refined if needed. New group assignment

$$\mathbf{A} = \arg\max_{\mathbf{A}'} \sum_{i,j,k} \sum_{q \in \{1...|\mathbf{B}|\}} \sum_{\mathbf{b}^* \in \mathcal{P}^q}$$

B • Hard to solve (exhaustive search), inefficient training procedure. We choose to fix

 $K = 1 \rightarrow$ intuitive interpretation

• Ground-truth group flag is defined as $grp(\bigcup b^*) = [|b^*| > 1]$

Institute of Science and Technology Austria, Klosterneuburg, Austria

• Choosing K is non-trivial. Too high (imbalance between empty and active cells)

$$A'(B_k^{i,j}, \bigcup b^*)$$
 iou $(B_k^{i,j}, \bigcup b^*)$

Experiments

- <u>______</u>0.6





• Ablation experiment (1): No groups (\sim two stages YOLO). Achieves lower map, though compensated by the learned offsets. • Ablation experiment (2): No offsets or fixed offsets.

Conclusions



• Accuracy (map@0.5) vs efficiency (runtime) trade-off • Two aerial views datasets (VEDAI and SDD), 1024x1024 px • Three backbones: tiny-YOLOv2, YOLOv2, MobileNet

- [+] Focusing on meaningful patches defined by groups allows to start at lower resolution and avoid unnecessay regions
- [+] Stage transition parameters can be changed without retraining to match the application at hand
- [-] Sequential yet not perfectly end-to-end (joint training)