

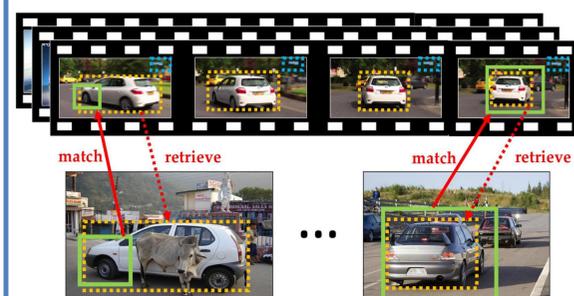
Krishna Kumar Singh

Fanyi Xiao

Yong Jae Lee

## Motivation and Idea:

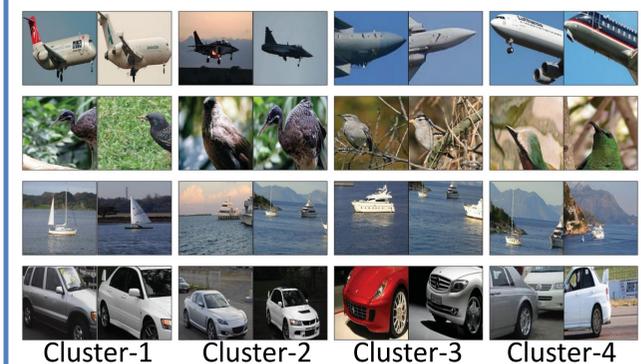
Train detectors with inexpensive image and video level tags, without using any localization information



- Existing weakly-supervised methods mine discriminative visual patterns, but they correspond to object-parts or include co-occurring background
- Key Idea:** Transfer tracked object boxes from videos as a substitute for strong human supervision to obtain more precise boxes

## Approach:

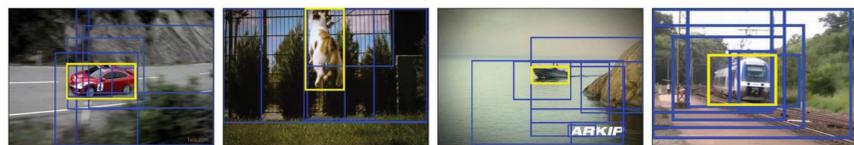
### Mining discriminative positive visual regions:



- Find discriminative object proposals by clustering in pre-trained Alexnet *pool-5* feature space and ranking clusters based on class consistency
- Most clusters do not tightly fit object due to scene clutter, occlusion, and intra class appearance variation

### Unsupervised video object tracking:

- Apply appearance and motion based unsupervised tracking method [1]
- For each video choose 9 highest ranked tracks
- Select highest scoring box:  $score(t_j^f) = \sum_i IOU(v_i^f, t_j^f) \times sim(r_i, v_i^f)$



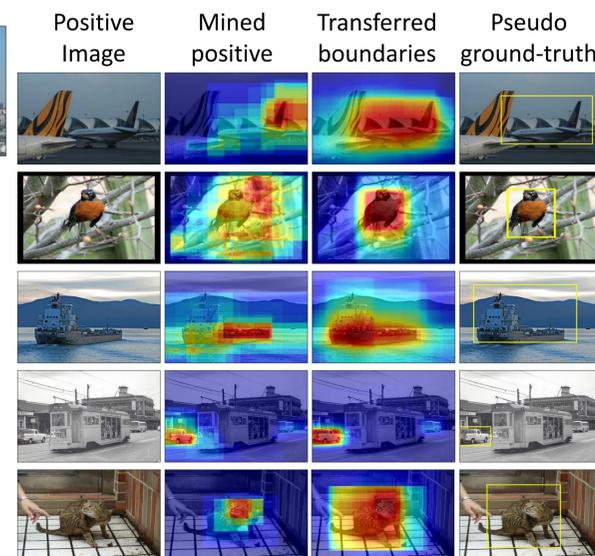
### Transferring tracked object boxes:



- For each positive region, find  $n$  best matching video regions and return corresponding tracked object boxes to the image
- Create 4-dimensional hough space where each box casts a vote for its coordinates

### Training an object detector:

- Train R-CNN using discovered pseudo ground-truth boxes
- Perform latent SVM update to improve the pseudo GT boxes
- Fine-tune the R-CNN model to update the features using pseudo GT boxes



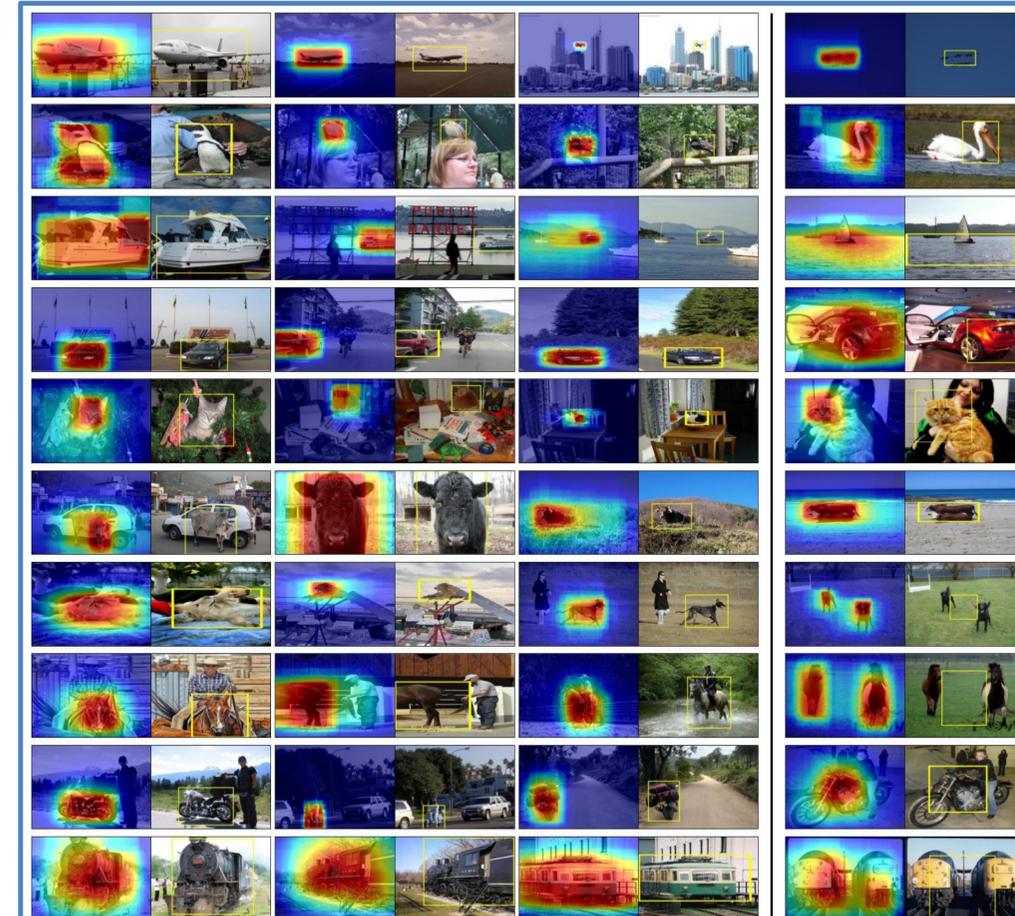
### Quantitative Results:

VOC 2007	Aero	Bird	Boat	Car	Cat	Cow	Dog	Horse	Mbike	Train	mAP
Wang et al., 2014	48.9	26.1	11.3	40.9	34.7	<b>34.7</b>	34.4	35.4	<b>52.7</b>	34.8	35.4
Cinbis et al., 2015	39.3	28.8	<b>20.4</b>	47.9	22.1	33.5	29.2	38.5	47.9	41.0	34.9
Ours	<b>53.9</b>	<b>37.7</b>	13.7	<b>56.6</b>	<b>51.3</b>	24.0	<b>38.5</b>	<b>47.9</b>	47.0	<b>48.4</b>	<b>41.9</b>

VOC 2010	Aero	Bird	Boat	Car	Cat	Cow	Dog	Horse	Mbike	Train	mAP
Cinbis et al., 2015	44.6	25.5	<b>14.1</b>	36.3	23.2	26.1	29.2	36.0	<b>54.3</b>	31.2	32.1
Ours	<b>53.5</b>	<b>37.5</b>	8.0	<b>44.2</b>	<b>49.4</b>	<b>33.7</b>	<b>43.8</b>	<b>42.5</b>	47.6	<b>40.6</b>	<b>40.1</b>

Ablation Study (VOC 2007)	Aero	Bird	Boat	Car	Cat	Cow	Dog	Horse	Mbike	Train	mAP
Initial pseudo GT	43.4	30.5	11.9	50.2	39.6	16.7	31.6	36.7	42.2	40.7	34.4
Updated pseudo GT (UGT)	48.0	34.2	12.2	51.3	43.0	21.9	33.4	39.1	43.8	42.2	36.9
UGT + bbox-reg	50.7	36.6	13.4	53.1	50.8	21.6	37.6	44.0	46.1	43.4	39.7
UGT + bbox-reg + finetune	<b>53.9</b>	<b>37.7</b>	<b>13.7</b>	<b>56.6</b>	<b>51.3</b>	<b>24.0</b>	<b>38.5</b>	<b>47.9</b>	<b>47.0</b>	<b>48.4</b>	<b>41.9</b>

## Qualitative Results:



- Each image pair consists of heatmap of transferred video object boxes and final selected pseudo GT box; last column shows some failure cases
- Our method accurately localizes the object in many images

## Conclusion:

- A novel weakly-supervised object detection framework that tracks and transfers object boxes from weakly-labeled videos to images
- State-of-the-art-results on PASCAL 2007 and 2010 datasets for the 10 categories of the YouTube-Objects dataset

**Acknowledgement:** This work was supported in part by an Amazon Web Services Education Research Grant and GPUs donated by NVIDIA

**Reference:** [1] F. Xiao and Y. J. Lee. Track and segment: An iterative unsupervised approach for video object proposals. In CVPR, 2016.