

# Recurrent Neural Networks for object detection

1<sup>st</sup> Bin Qasim Ahmad  
Technical University of Munich  
Department of Informatics  
Munich, Germany  
ahmad.qasim@tum.de

2<sup>nd</sup> Pettirsch Arnd  
Technical University of Munich  
Department of Informatics  
Munich, Germany  
a.pettirsch@outlook.de

**Abstract—ToDo**  
**Index Terms—TBD.**

## I. INTRODUCTION

### A. Image and Video Object Detection in general

- Image object detection history.
  - Bayesian methods before deep learning
  - ImageNet challenge and VID [15]
  - Deep Learning and AlexNet [16]
- Single stage and 2-stage image object detectors.
  - A two-stage pipeline firstly generates region proposals, which are then classified and refined. [17]
  - A single-stage method is often more efficient but less accurate. Directly regress on bounding boxes and classes. [18], [19]
- Why is video object detection harder?
  - Large size
  - Motion blur
  - Quality of the dataset
  - Partial occlusion
  - Unconventional Poses

### B. Recurrent Neural Networks in general

ToDo

## II. FEATURE-BASED VIDEO OBJECT DETECTION

### A. Definition

ToDo

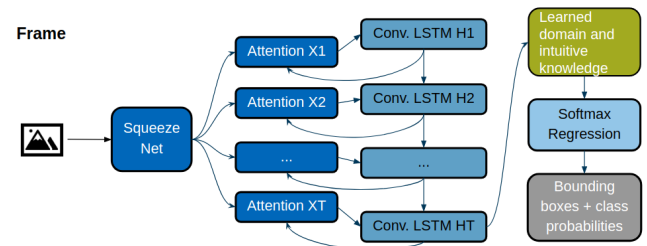
### B. Recurrent Multi-frame Single Shot Detector for Video Object Detection

ToDo

### C. Mobile Video Object Detection with Temporally Aware Feature Maps

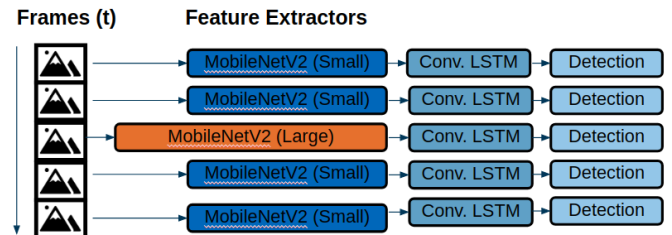
ToDo

### D. Feature Selective Small Object Detection via Knowledge-based recurrent attentive neural networks



- Compute feature maps using a modified SqueezeNet architecture.
- Propagate the features through a Recurrent Attentive Neural Network, comprised of:
  - Attention Mechanism to detect key areas within the feature maps.
  - Convolutional LSTM for temporal feature propagation.
- Reverse gaussian feature maps are combined with the maps obtained from Conv. LSTM.
  - These feature maps are based on learnable mean and covariance terms.
  - This prior knowledge is derived from the assumption that traffic signs are always located at the bias of the center.

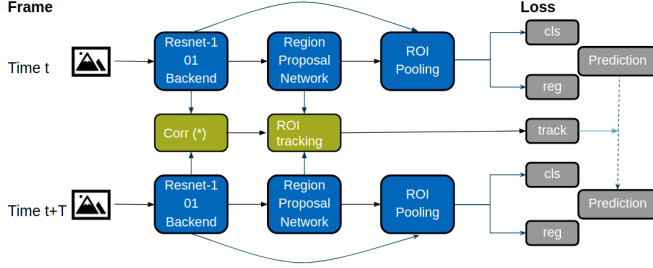
### E. Looking fast and slow: memory-guided mobile video object detection



- Run multiple feature extractors sequentially or concurrently to obtain feature maps.
  - The idea is to use small and large feature extractors to optimize performance.
- Aggregate and refine these feature maps using convolutional LSTM based memory network.

- To improve speed of LSTM network, add skip connections and LSTM state groups.
- Apply SSD-style detection on refined features to obtain classification and bounding boxes.
- Use a reinforcement learning based policy for selection of which feature extractor to run.
- Large and small frame extractors can run in parallel using asynchronous mode.

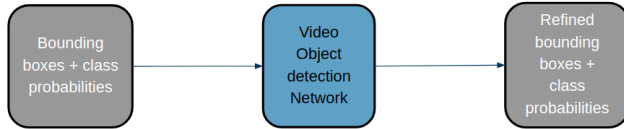
#### F. Detect to Track and track to detect



- Compute Convolutional feature maps using a Resnet-101 architecture.
- Use a RPN (region proposal network) to find candidate regions in the frame.
- ROI Pooling layer, to classify boxes and refine their coordinates (regression).
- Find correlation features between two frames' feature maps and do ROI tracking.
- Due to memory constraints, use tracklets, which are class-based optimal paths in video.

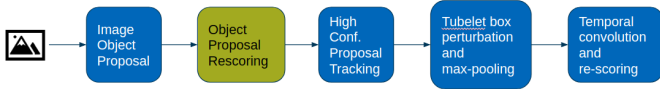
### III. BOX-LEVEL-BASED VIDEO OBJECT DETECTION

#### A. Definition



Bounding Boxes and Class probabilities are fed into the network and are refined temporally and/or spatially.

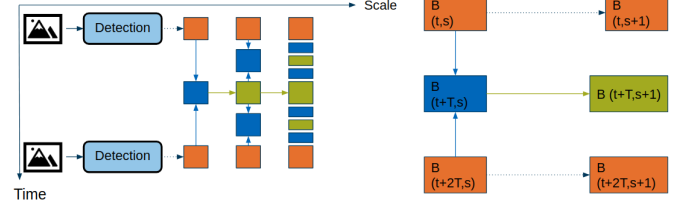
#### B. Object Detection from Video Tubelets with Convolutional Neural Networks



- Use selective search algorithm to generate around 2000 object proposals on each frame.
- Use GoogleNet for feature extraction and then 30 SVMs for 30 VID classes to generate object proposal scores for each object proposal.
- Track high confidence targets bi-directionally.
- Two kinds of Perturbations:

- The first method is to generate new boxes around each tubelet box on each frame by randomly perturbing the boundaries of the tubelet box.
- The second perturbation method is to replace each tubelet box with original object detections that have overlaps with the tubelet box beyond a threshold.
- Train a class-specific TCN using the tubelet features as input. The inputs are time series including detection scores, tracking scores and anchor offsets. The output values are probabilities whether each tubelet box contains objects of the class

#### C. Optimizing Video Object Detection via Scale-Time Lattice



- Apply object detection on keyframes extracted adaptively.
  - The extraction policy is based on number of objects and amount of movement in frames.
  - If higher number/movement of objects in frames then higher extraction rate.
- Propagation and refinement unit, propagates the frames temporally and refines spatially.
- For temporal propagation, use a small network such as resnet-18 to extract box features and a regressor to predict object movement from  $t$  to  $t + T$ .
- For spatial refinement, use a regressor to refine the bounding boxes over increasing scale.

#### D. Context Matters: Refining Object Detection in Video with Recurrent Neural Networks

ToDo

#### E. Spatially Supervised Recurrent Convolutional Neural Networks for Visual Object Tracking

ToDo

### IV. FLOW-BASED OBJECT DETECTION

#### A. Definition

ToDo

#### B. Deep Feature Flow for Video Recognition

ToDo

### V. COMPARISON OF DIFFERENT APPROACHES

#### A. General

TABLE I: Results on KITTI Dataset

Model	MAP	FPS	Machine	Architecture
Recurrent [1]	86.0	50	Nvidia TITAN X	Feature-Level
Feature Selective [6]	81.3	30.8	Nvidia TITAN X	Feature-Level

TABLE II: Results on ImageNet Dataset

Model	MAP	FPS	Machine	Architecture
DT [8]	82.0	7	Nvidia TITAN X	Feature-Level
DT [8]	78.5	55	Nvidia TITAN X	Feature-Level
Scale-Time Lattice [10]	79.6	20	Nvidia TITAN X	Box-Level
Scale-Time Lattice [10]	79	62	Nvidia TITAN X	Box-Level
DeepFeature Flow [3]	73.9	3	-	Flow-Based
DeepFeature Flow [3]	73.1	20.5	-	Flow-Based
Looking Fast and Slow [7]	60.7	48.8	Pixel Phone	Feature-Level
Object Detection with Temporally-Aware [2]	54.4	15	Pixel Phone	Feature-Level

TABLE III: Results on COCO Dataset

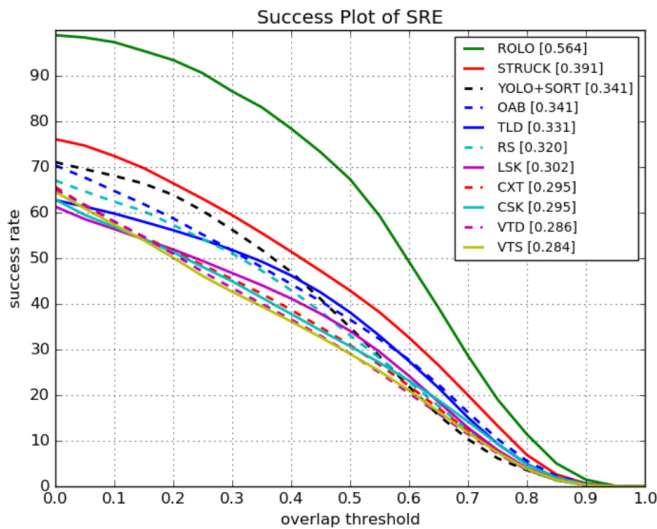
Model	MAP	FPS	Machine	Architecture
Feature Selective [6]	57.8	37.5	Nvidia TITAN X	Feature-Level

TABLE IV: Results on YT Dataset

Model	MAP	FPS	Machine	Architecture
Context Matters [4]	68.73	-	-	Box-Level

TABLE V: Results on OTB Challenge Dataset

Model	Success Rate	IoU	FPS	Machine
Spatially Supervised [5]	0.564	0.455	20/60	Nvidia TITAN X



## B. Conclusion Performance

Todo

## C. Conclusion Prediction Quality

Todo

## VI. OUTRO

### A. Conclusion

Todo

### B. Further work

Todo

## ACKNOWLEDGMENT

Todo

## REFERENCES

Please number citations consecutively within brackets [1]. The sentence punctuation follows the bracket [2]. Refer simply to the reference number, as in [3]—do not use “Ref. [3]” or “reference [3]” except at the beginning of a sentence: “Reference [3] was the first . . .”

Number footnotes separately in superscripts. Place the actual footnote at the bottom of the column in which it was cited. Do not put footnotes in the abstract or reference list. Use letters for table footnotes.

Unless there are six authors or more give all authors’ names; do not use “et al.”. Papers that have not been published, even if they have been submitted for publication, should be cited as “unpublished” [4]. Papers that have been accepted for publication should be cited as “in press” [5]. Capitalize only the first word in a paper title, except for proper nouns and element symbols.

For papers published in translation journals, please give the English citation first, followed by the original foreign-language citation [6].

## REFERENCES

- [1] Alexander Broad, Michael Jones, Teng-Yok Lee. Recurrent Multi-frame Single Shot Detector for Video Object Detection. 2018.
- [2] Mason Liu, Menglong Zhu. Mobile Video Object Detection with Temporally-Aware Feature Maps. 2018.
- [3] Xizhou Zhu, Yuwen Xiong, Jifeng Dai, Lu Yuan, Yichen Wei. Deep Feature Flow for Video Recognition. 2017.
- [4] Subarna Tripathi, Zachary C. Lipton, Serge Belongie, Truong Nguyen. Context Matters: Refining Object Detection in Video with Recurrent Neural Networks.
- [5] Guanghan Ning, Zhi Zhang, Chen Huang, Zhihai He, Xiaobo Ren, Haohong Wang. Spatially Supervised Recurrent Convolutional Neural Networks for Visual Object Tracking. 2016.
- [6] Kai Yi, Zhiqiang Jian, Shitao Chen, Nanning Zheng. Feature Selective Small Object Detection via Knowledge-based Recurrent Attentive Neural Network. 2019.
- [7] Mason Liu, Menglong Zhu, Marie White, Yinxiao Li, Dmitry Kalenichenko. Looking Fast and Slow: Memory-Guided Mobile Video Object Detection. 2019.
- [8] Christoph Feichtenhofer, Axel Pinz, Andrew Zisserman. Detect to Track and Track to Detect. 2017.
- [9] Kai Kang, Wanli Ouyang, Hongsheng Li, Xiaogang Wang. Object Detection from Video Tubelets with Convolutional Neural Networks. 2016.
- [10] Kai Chen, Jiaqi Wang, Shuo Yang, Xingcheng Zhang, Yuanjun Xiong, Chen Change Loy, Dahua Lin. Optimizing Video Object Detection via a Scale-Time Lattice. 2018.
- [11] Ian Goodfellow, Yoshua Bengio, Aaron Courville. Deep Learning (Adaptive Computation and Machine Learning). 2017.

- [12] Alexander Broad, Michael Jones, Teng-Yok Lee. Supplementary Material for Recurrent Multi-frame Single Shot Detector for Video Object Detection. 2018.
- [13] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li and Li Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database.
- [14] Alex Krizhevsky, Ilya Sutskever, Geoffrey E. Hinton. ImageNet Classification with Deep Convolutional Neural Networks. 2012.
- [15] Shaoqing Ren, Kaiming He, Ross Girshick, Jian Sun. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. 2016.
- [16] Joseph Redmon, Ali Farhadi. YOLOv3: An Incremental Improvement.
- [17] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, Alexander C. Berg. SSD: Single Shot MultiBox Detector. 2016.