
Pack and Detect: Fast Object Detection in Videos Using Region-of-Interest Packing

Athindran Ramesh Kumar
Department of Computer Science
Indian Institute of Technology Madras
Chennai, Tamil Nadu, India
r.athindran@gmail.com

Balaraman Ravindran
Department of Computer Science
Robert Bosch Centre for Data Science and AI
Indian Institute of Technology Madras
Chennai, Tamil Nadu, India
ravi@cse.iitm.ac.in

Anand Raghunathan
Department of Electrical Engineering
Purdue University
West Lafayette, IN, USA
raghunathan@purdue.edu

Abstract

Object detection in videos is an important task in computer vision for various applications such as object tracking, video summarization and video search. Although great progress has been made in improving the accuracy of object detection in recent years due to the rise of deep neural networks, the state-of-the-art algorithms are highly computationally intensive. In order to address this challenge, we make two important observations in the context of videos: (i) Objects often occupy only a small fraction of the area in each video frame, and (ii) There is a high likelihood of strong temporal correlation between consecutive frames. Based on these observations, we propose Pack and Detect (PaD), an approach to reduce the computational requirements of object detection in videos. In PaD, only selected video frames called anchor frames are processed at full size. In the frames that lie between anchor frames (inter-anchor frames), regions of interest (ROIs) are identified based on the detections in the previous frame. We propose an algorithm to pack the ROIs of each inter-anchor frame together into a reduced-size frame. The computational requirements of the detector are reduced due to the lower size of the input. In order to maintain the accuracy of object detection, the proposed algorithm expands the ROIs greedily to provide additional background around each object to the detector. PaD can use any underlying neural network architecture to process the full-size and reduced-size frames. Experiments using the ImageNet video object detection dataset indicate that PaD can potentially reduce the number of FLOPS required for a frame by $4\times$. This leads to an overall increase in throughput of $1.25\times$ on a 2.1 GHz Intel Xeon server with a NVIDIA Titan X GPU at the cost of 1.1% drop in accuracy.

Keywords Object Detection · Neural Network · Temporal Correlation · Object occupancy · Region-of-Interest packing

1 Introduction

The task of object detection in videos (17; 47; 18; 16; 46; 45; 3; 31; 24; 39) has been gaining attention in recent years. It serves as an important preprocessing task for object tracking and for several other video processing tasks such as video summarization and video search. Many object detection applications require video frames to be processed in real-time in resource-constrained environments. It is thus imperative to design systems that can detect objects in videos accurately, but also in a computationally efficient manner.

Still-image object detection has been studied extensively in the past. The accuracy and speed of still-image object detection have improved by leaps and bounds in recent years due to advances in deep convolutional neural networks (CNN). Recent CNN-based object detectors include Faster-RCNN (37), SSD (25), YOLO (34; 35; 36) and RFCN(4). These still-image object detectors can be extended for object detection in videos by using them on a per-frame basis. However, this is inefficient as there is a strong temporal correlation between frames in a video. This temporal redundancy can be leveraged either to improve the accuracy or speed of object detection. In the recent past, there have been several attempts to improve the accuracy of object detection in videos by either integrating the bounding boxes (17; 18; 10; 40) or features (16; 46; 45; 12) across frames. However, there has not been enough attention on leveraging this temporal redundancy to improve speed. Some exceptions to this norm are (47; 3; 31; 24; 39). In this work, we propose a method, Pack and Detect (PaD), for fast object detection in videos that can work with any underlying object detector.

PaD leverages two key opportunities in the context of object detection in videos. First, the objects of interest often occupy only a small fraction of an image. Second, there is a strong correlation between successive frames in a video. In PaD, only selected frames called anchor frames are passed in their entirety to the underlying object detector. In frames that lie between anchor frames (inter-anchor frames), the detections from the previous frame are used to identify ROIs in the image where an object could potentially be located. The ROIs are packed in a reduced-size image that is fed into the detector, resulting in lower computational requirements.

We propose a ROI packing algorithm based on the following criteria:

1. Each ROI is expanded to provide as much background context as possible to maintain the accuracy of the detector.
2. There is minimal loss of resolution and no change in aspect ratio to maintain the accuracy of the detector.
3. Each object is present in a unique ROI.
4. The space in the reduced-size frame is used as efficiently as possible.

We evaluate PaD by implementing it on top of the SSD300 object detector and evaluating it with the ImageNet video object detection dataset. Our results indicate that PaD reduces the FLOP count for reduced-size frames by around $4\times$. Overall, PaD achieves $1.25\times$ increase in throughput with only a 1.1% drop in accuracy.

2 Related Work

Object Detection in Videos

The temporal redundancy present in videos has been exploited before to improve the accuracy and speed of object detection. In (17; 18; 10; 40), the aggregation of information from neighbouring frames is done at the bounding box level to improve accuracy. In (17), per-frame object detection is combined with multi-context suppression, motion-guided propagation and object tracking to improve detection accuracy. In (10; 40), non-maximum suppression is done over bags of frames. In (16; 46; 45; 12), integration of the CNN features across neighbouring frames is used to improve accuracy. In (16), a CNN is combined with a Long Short Term Memory (LSTM) to obtain temporal features for object detection. In (46; 45; 12), the features from neighbouring frames are aggregated together using optical flow information to improve feature quality. These methods (17; 18; 10; 40; 16) pose large computation requirements, making them often unsuitable for real-time processing.

On the other hand, (47; 3; 31; 24; 39) are relatively faster methods aimed at object detection in videos. The methods in (3; 24; 39) are faster by virtue of using a faster still-image object detector or an efficient backbone network. In (47), the feature maps from selected anchor frames are transferred to neighbouring frames by warping them with optical flow information, leading to reduced computation. In (31), neighbouring frames are subtracted to give rise to a sparse input that is processed with a sparsity-aware hardware accelerator (8) to achieve computational savings.

PaD differs vastly from the prior methods proposed to speed-up video object detection. PaD can be used alongside previous methods such as (47) and on top of existing efficient object detectors that operate on a per-frame basis (3; 24; 39). Moreover, PaD does not require any specialized hardware accelerator like in (31) to obtain computational savings.

Efficient Neural Networks

Several efforts have attempted to reduce the computational requirements of neural networks. Quantization with retraining was shown to improve the efficiency of neural network implementations in (42). Deep compression (9) combined pruning, trained quantization and weight compression and demonstrated large speedups on a custom hardware accelerator (8). Subsequent efforts have explored structured sparsity (43) by pruning filters (26; 43; 7; 29; 22; 11) of a CNN. MobileNet (13) replaces the standard convolution with a combination of depth-wise and point-wise convolution to reduce computation. SqueezeNet (14) uses network architecture modifications to reduce the number of computations and memory. Scalable-effort classifiers reduce computational requirements by first using lower-complexity classifiers to process an input and subsequently using higher-accuracy classifiers only when needed (41). A similar approach is taken by Big-Little networks (32). Conditional computation (2) selectively activates certain parts of the network depending on the input. The policy for deciding which parts of the network to activate is learnt using reinforcement learning. Dynamic deep neural networks (D2NN) (23) work in a similar manner to conditional computation and turn on/off regions of the network using reinforcement learning. DyVEDeep (6) reduces computations in neural networks dynamically by using three strategies - saturation prediction and early termination, significance driven selective sampling and similarity-based feature map approximation.

The above methods focus on modifications to the network to reduce computations and achieve speedup. In this work, we take a complementary approach and compress the inputs that we feed into the network. Hence, PaD is orthogonal to most existing techniques and can be used in combination with them.

Visual Attention Mechanism

Inspired by human vision, there have been several attempts (20; 33; 30; 1; 19; 15) to reduce computation by processing an image as a sequence of glimpses rather than as a whole. The notion of a foveal glimpse is somewhat similar to the idea of ROI discussed here. However, there are several important differences. A foveal glimpse is a high resolution crop of an important region in the image that is crucial to the task at hand. In our work, we pack all the ROIs together in a single frame and do not process them sequentially. Further, the location of ROIs is inferred from the detections in the previous frame in a video and does not need an attention mechanism. Also, a foveal glimpse obtains crops by extracting pixels close to the location target at high resolution and pixels far from the location target at low resolution. We do not employ multi-resolution processing. Hence, our work, although inspired from the notion of foveal attention is considerably different.

Multiple Object Tracking

Object detection in video is a precursor to the problem of multiple object tracking. Once the objects are detected in the video, the detections are linked together to form a track. This problem is studied separately from the object detection problem in the literature. The ImageNet VID dataset used in this work does not have ground truth labels to measure the tracking metrics. In the MOT challenge (27), the detections that are input to the tracker are provided with the dataset. Several popular trackers such as (21; 44; 28; 38) have garnered attention through the challenge. While the use of detection and tracking to complement each other to improve accuracy or speed is possible, it is not well studied in the literature. In (5), the detection and tracking have been used in a complementary fashion to

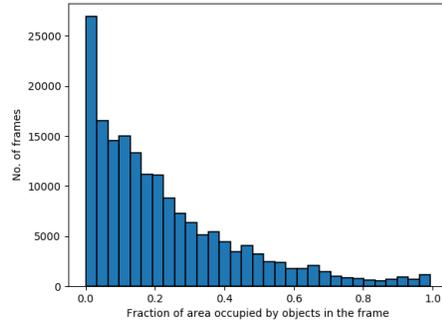


Figure 1: Histogram of object occupancy ratio

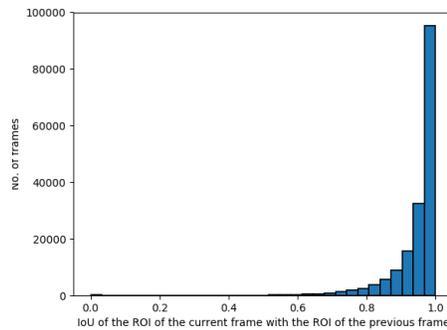


Figure 2: Histogram of Intersection over Union (IoU) of regions containing objects between consecutive frames

improve the accuracy. In future work, we will explore the potential of combining detection and tracking to improve speed.

3 Motivation

3.1 Occupancy of objects in frames

PaD leverages the hypothesis that the objects of interest occupy only a small fraction of the area in the frame. We support this hypothesis using statistics from a popular video dataset. Figure 1 is a histogram of the object occupancy ratio in the ImageNet VID validation set containing 555 videos with 176126 frames. From the figure, we see that the objects occupy only 22.7% of the frame on average. In a vast majority of the frames, the objects occupy less than 30% of the frame.

3.2 Temporal correlation of object locations across frames

It is well known that successive frames in a video are likely to be highly correlated. We illustrate this through a statistical analysis of the ImageNet VID validation set. Figure 2 presents a histogram of the object occupancy area Intersection over Union (IoU) statistics between consecutive frames in the dataset. In the figure we can clearly see a sharp peak close to 1. On average, the IoU of areas containing objects between consecutive frames is 94.4%.

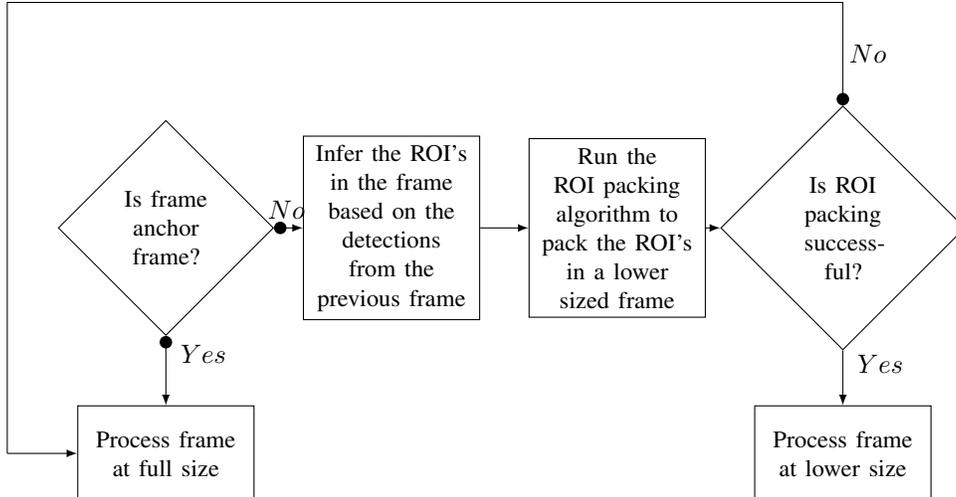


Figure 3: Overall Architecture of PaD

4 Pack-and-Detect: Approach and Algorithms

4.1 Overview

An overview of the PaD approach is presented in Figure 3. Full-sized video frames are processed at regular intervals (by designating the first of every d frames as an anchor frame). In other frames, ROIs are identified based on the locations of the detections from the previous frame. Only detections with a minimum confidence threshold τ are taken into consideration. An ROI packing algorithm attempts to pack the ROIs into a reduced-size frame. If the packing is successful, then the reduced-size frame is processed instead, giving rise to computational savings. Once the reduced-size frame is processed using the CNN detector, the object locations are mapped back to the original frame. However, if the packing is not successful, then the frame is processed at full size, incurring an overhead due to checking whether ROI packing is possible. We demonstrate that this tradeoff is often favorable, resulting in a net improvement in the speed of object detection.

4.2 ROI packing algorithm

Figure 4 describes the ROI packing algorithm. As a first step in the algorithm, we construct a graph where nodes represent ROIs and an edge connects two nodes if the corresponding ROIs intersect. We find all connected components of this graph. We then find the enclosing bounding box over the union of ROIs in each connected component. We iterate the connected components algorithm until the final bounding boxes do not overlap. This constraint is important because if two bounding boxes overlap, then parts of the same object could be present two or more times in the packed frame. Once the number and size of the bounding boxes are decided, the layout of the bounding boxes is determined by using the algorithm presented in Figure 5. Once the layout is decided, a check is done to see whether the bounding boxes can fit in the layout. If it is not possible to fit the bounding boxes in the layout, the image is processed at full size.

If the bounding boxes can be fit in the layout, a post-processing step is performed as described below. Our experiments indicated that neural network based object detectors are often overfit to the background context of the object to be detected. Consequently, the accuracy of the object detector degrades if there is no background context. To address this challenge, we extend each bounding box to provide as much context as possible to the detector. The algorithm for extending the bounding boxes works as follows. We decide whether to first extend the boxes horizontally or vertically. For the sake of discussion, let us assume that the choice is to first extend all the bounding boxes horizontally. We find all the bounding boxes that could potentially intersect when extended horizontally. We extend all bounding boxes horizontally until the layout size is reached or the bounding boxes start intersecting with each other. Then, we repeat the same procedure in the other dimension. Once

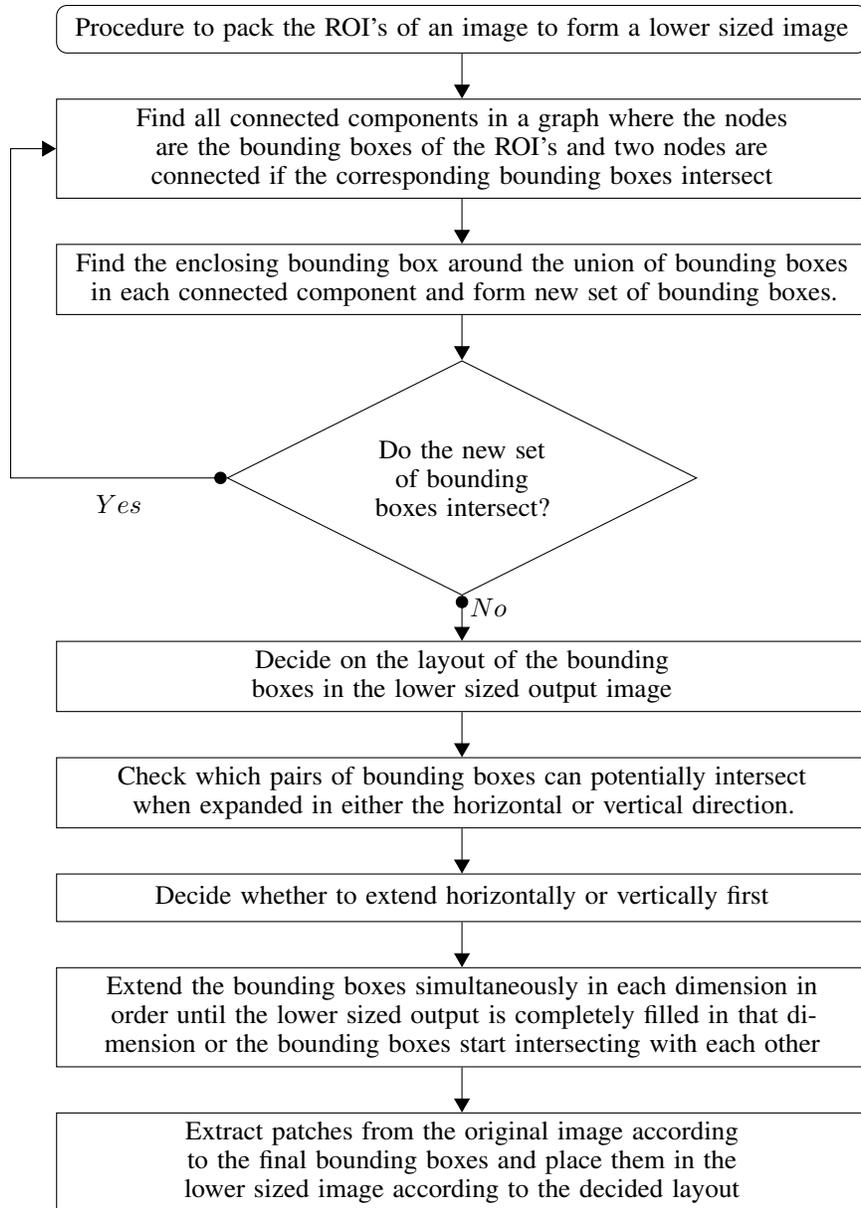


Figure 4: Flowchart illustrating the overall flow of the ROI packing algorithm. Please refer to figure 7 to observe the results of the ROI packing algorithm

the final bounding boxes are decided, the corresponding regions in the image are extracted and the reduced-size frame is composed according to the determined layout.

5 Experimental Methodology

The ImageNet object detection dataset (DET) is a dataset comprising 200 classes of objects that form a subset of the ImageNet 1000 classes. Further, the ImageNet video object detection dataset (VID) comprises of 30 classes of objects from among the DET 200 classes. The ImageNet video object detection (VID) dataset was the most appropriate choice for illustrating the results of our work. The ImageNet VID training set has 3862 video snippets and the ImageNet VID validation set has 555 video snippets. 53539 frames from the DET dataset comprising only of the classes from the VID

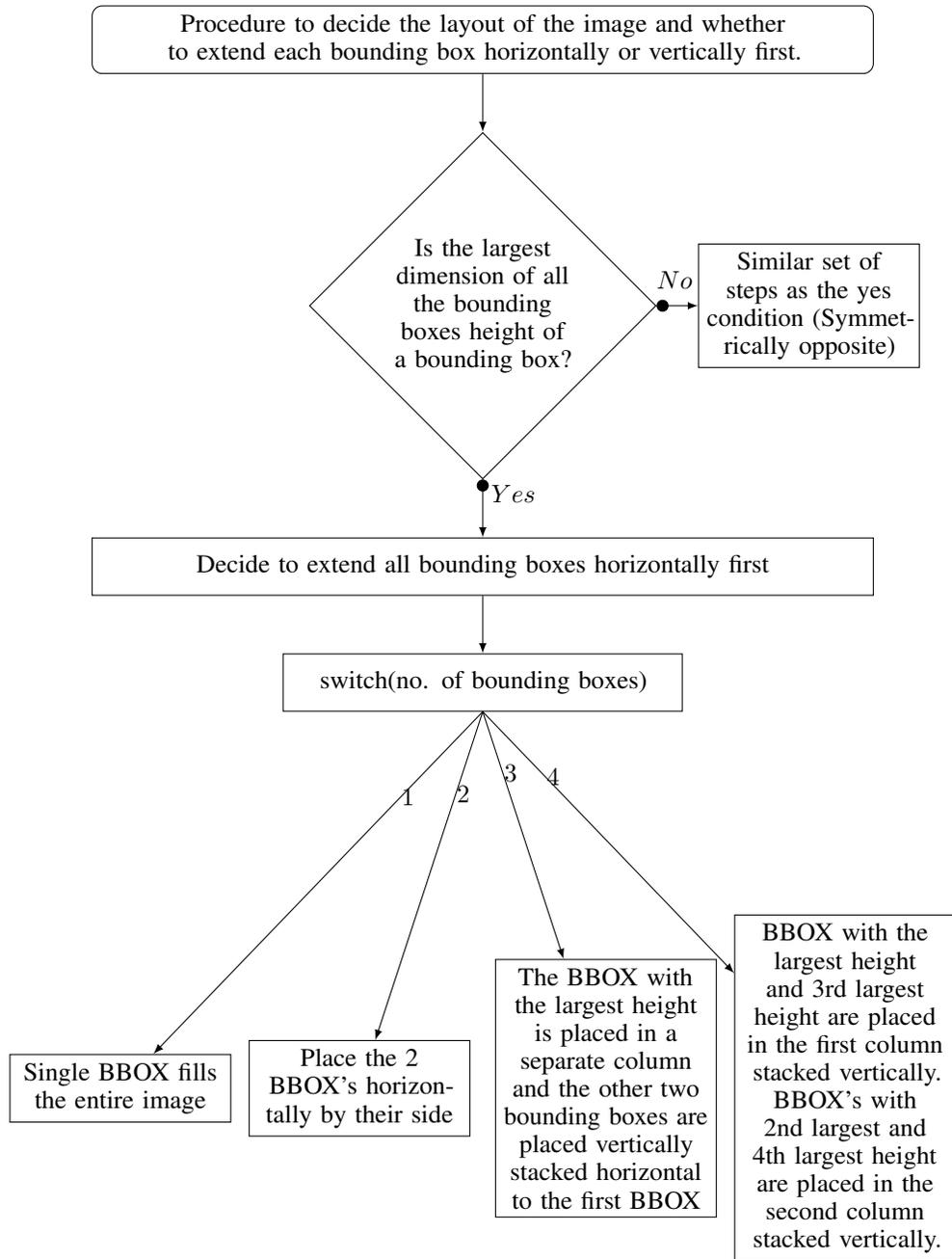


Figure 5: Flowchart detailing the layout procedure. Please refer to figure 7 to observe the results of the layout procedure. The cases with 1,2 and 4 non-overlapping bounding boxes have been illustrated in figure 7

Dataset	No. of classes	No. of video snippets	No. of frames selected
DET training set	200	N/A	53539
VID training set	30	3862	57834
VID validation set	30	555	176126

Figure 6: Details of the ImageNet DET and ImageNet VID dataset

dataset and 57834 frames from the VID training set were combined to form the final training set in our experiments.

The SSD300 (25) object detector operated on a per-frame basis was used as the baseline for our work. The SSD300 object detector uses VGG16 as feature extractor. The SSD300 pretrained model on the DET dataset was further trained on our training set for 210k iterations with a learning rate of 10^{-3} for the first 80000 iterations, 10^{-4} for the next 40000 iterations and 10^{-5} for the rest of the training. This SSD300 trained model gave a mAP score of 70.6 on the VID validation set. Further, this model has a network throughput of 47 fps and a overall throughput (including standard pre-processing time) of 18 fps.

The SSD300 network processes images at 300×300 as the name suggests. However, closer observation of the network suggested that the same network can process 150×150 images as well by stopping processing at the penultimate layer. Hence, we use the same SSD300 network to process both full-size and reduced-size images. In all our experiments, the full size s_1 is 300 and the reduced size s_2 is 150. When a 150×150 sized image is passed on to the SSD300 network, processing is configured to stop at the penultimate layer. All the experiments were performed using the SSD Caffe framework running on a 2.1 GHz Intel Xeon CPU with a Nvidia TITAN X GPU. In all the experiments, the batch size was 1 to emulate a real-time processing scenario. The detection threshold τ used to select ROIs was fixed at 0.2 in our experiments unless explicitly specified otherwise.

6 Experimental Results

Results from sample videos

We show results on processing some sample videos with PaD. Figure 7 provides sample detections with our ROI-packing algorithm. The first column shows frame i . The second column shows the ROI-packed reduced-size frame $i + 1$ with the detections. The third column shows the original frame $i + 1$ with detections mapped from ROI-packed frame $i + 1$. For this experiment, the confidence threshold τ for selecting a detection as an ROI for the next frame was set to 0.3 for the sake of illustration. All bounding boxes with a minimum threshold of 0.2 are shown in the figure.

In Figure 8, we plot the per-frame time as well as the cumulative time for processing a sample video using PaD and the baseline. It can be seen that processing the lower sized frame of 150×150 is almost $3 \times$ faster. When ROI packing fails, there is a slight overhead incurred which is visible towards the end of the video in Figure 8(a). Also, we see that some frames require almost no time for processing. This is because the previous frame had no detections. Overall, from Figure 8(b), we note that processing the video using PaD requires almost 8s lesser time than the baseline.

Results over the entire dataset

PaD was run with a inter-anchor distance $d = 5$ and $s_2 = 150$. In Figure 9, we plot the histogram of average per-frame processing time on a video-by-video basis. In other words, the average time taken per frame was obtained for each video and is plotted as a histogram across videos. From the figure, we can clearly see that the average time taken to process a frame is lower using PaD for more videos than the baseline. The average per-frame speedup is around $1.25 \times$ and the FLOP reduction on the average is 32%. The average overhead incurred for ROI-packing is around 9% of the total time taken. The mAP score drops by 1.1% (from 70.6 to 69.5).



Figure 7: Sample set of consecutive frames processed with the ROI packing mechanism. The first column shows frame i . The second column shows the ROI packed frame $i + 1$ with the detections. The third column shows the original frame $i + 1$ with detections transformed from ROI packed frame $i + 1$. For this experiment, the confidence threshold for selecting a detection as an ROI for the next frame is chosen as 0.3 for the sake of illustration. All bounding boxes with a minimum threshold of 0.2 are shown in the figure.

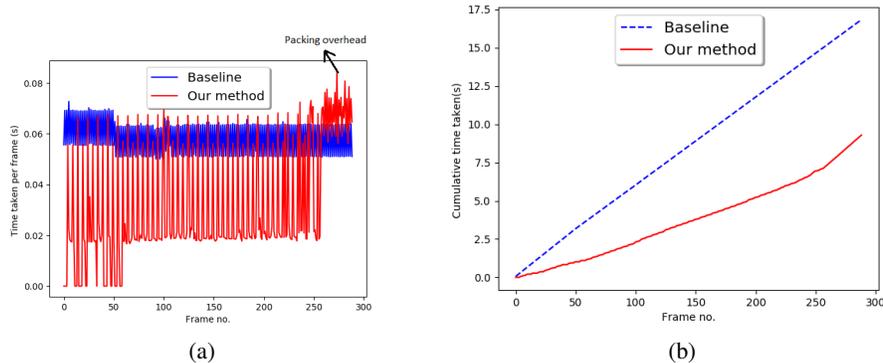


Figure 8: (a)Time taken per frame for processing a sample video (b)Cumulative time taken processing a sample video

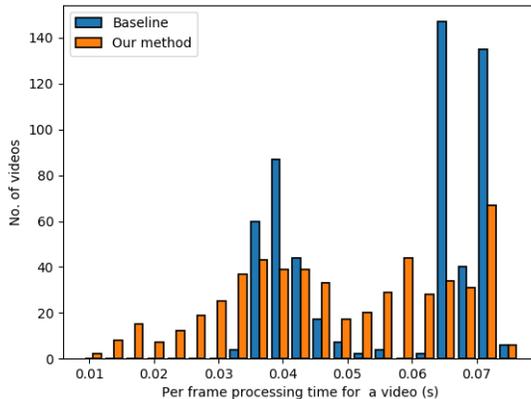


Figure 9: Histogram of average per-frame processing time on a video-by-video basis

Comparison with a naive ROI-packing algorithm

In order to illustrate the benefits of our ROI-packing algorithm discussed in section 4, we compare the accuracy drop when compared with a naive ROI-packing algorithm.

The naive ROI-packing algorithm can accommodate upto four ROIs just like the sophisticated method. If there are more than four objects in the frame, the frame is processed at full size. Otherwise, the bounding box surrounding each frame is extended by a factor of $1.2\times$ and is treated as an ROI. If there is only one object, the ROI surrounding the bounding box is rescaled to size $s_2 \times s_2$ and is processed by the detector. If there are two objects, the lower sized frame is divided into two columns of size $s_2 \times \frac{s_2}{2}$. The two ROIs are rescaled to the appropriate sizes and laid out on the lower sized frame. In the case of three or four objects, the lower sized frame is divided into four regions in two columns and two rows of size $\frac{s_2}{2} \times \frac{s_2}{2}$. In the case of three ROIs, the ROIs will be rescaled to occupy three of the four regions in the frame and the fourth region will be left blank. In the case four ROIs, the ROIs will be rescaled and fit to these four regions. We do not perform a greedy expansion of the RoIs to provide additional background context. Instead, the ROIs are just expanded by a constant factor of $1.2\times$ and rescaled to appropriate size.

PaD’s ROI packing method with inter-anchor distance $d = 5$ gave a mAP score of 69.5. With the same parameter setting, the naive ROI packing algorithm gave a mAP score of 56.8. This clearly illustrates the need for an ROI-packing algorithm that preserves the scale and aspect ratio of the ROIs and provides as much background context as possible.

7 Conclusion and Future Work

Still-image object detection has improved by leaps and bounds in recent years due to the success in training and deploying neural networks. However, the opportunities that are available in the context of videos have not been fully exploited. Neural networks are in general very compute-intensive. In this work, we use the opportunities available in the context of videos to speed up and reduce the amount of computation in neural network based object detectors. In the proposed method, called PaD, the full-sized input is only processed in selected anchor frames. In the inter-anchor frames, ROIs are identified based on the locations of objects in the previous frame. These ROIs are packed together in a reduced-size frame that is fed to the CNN object detector. The ROI packing algorithm needs to ensure that the scales and aspect ratios of the objects are preserved and enough background context is provided. With this setup, we observed $1.25\times$ speedup with 1.1% drop in accuracy on the ImageNet VID validation set. Further, the time taken to process a lower sized frame is almost $3\times$ lesser and the FLOP count reduces by $4\times$.

As part of future work, we plan to incorporate a motion model to obtain the ROIs in the current frame. Incorporating a motion model could also help extend this framework to larger batch sizes. Also, it is possible to use two different models or networks to process larger sized and smaller sized frames. This will help reduce the accuracy drop but will in turn increase the memory footprint. There is an overhead incurred in checking whether the ROIs can fit in the lower sized frame. Currently, we select anchor frames at regular intervals. However, information on whether ROIs were packed successfully in previous frames can help us decide how frequently we select anchor frames. Thus, another line of future work is a dynamic mechanism for selecting anchor frames in order to reduce the overhead. It would be interesting to test PaD in more resource constrained platforms like mobile GPUs and CPUs. We expect the benefits to be more pronounced in such platforms.

Acknowledgements We thank Intel India, Robert Bosch Centre for Data Science and AI (RBC-DSAI) and Center for Computational Brain Research (CCBR) for supporting this research.

References

- [1] Jimmy Ba, Volodymyr Mnih, and Koray Kavukcuoglu. Multiple object recognition with visual attention. *arXiv preprint arXiv:1412.7755*, 2014.
- [2] Emmanuel Bengio, Pierre-Luc Bacon, Joelle Pineau, and Doina Precup. Conditional computation in neural networks for faster models. *CoRR*, abs/1511.06297, 2015.
- [3] Xingyu Chen, Zhengxing Wu, and Junzhi Yu. TSSD: temporal single-shot object detection based on attention-aware LSTM. *CoRR*, abs/1803.00197, 2018.
- [4] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. R-fcn: Object detection via region-based fully convolutional networks. In *Advances in neural information processing systems*, pages 379–387, 2016.
- [5] Christoph Feichtenhofer, Axel Pinz, and Andrew Zisserman. Detect to track and track to detect. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3038–3046, 2017.
- [6] Sanjay Ganapathy, Swagath Venkataramani, Balaraman Ravindran, and Anand Raghunathan. Dyvedeep: Dynamic variable effort deep neural networks. *CoRR*, abs/1704.01137, 2017.
- [7] Jia Guo and Miodrag Potkonjak. Pruning filters and classes: Towards on-device customization of convolutional neural networks. In *Proceedings of the 1st International Workshop on Deep Learning for Mobile Systems and Applications*, EMDL '17, pages 13–17, New York, NY, USA, 2017. ACM.
- [8] Song Han, Xingyu Liu, Huizi Mao, Jing Pu, Ardavan Pedram, Mark A. Horowitz, and William J. Dally. Eie: Efficient inference engine on compressed deep neural network. In *Proceedings of the 43rd International Symposium on Computer Architecture*, ISCA '16, pages 243–254, Piscataway, NJ, USA, 2016. IEEE Press.

- [9] Song Han, Huizi Mao, and William J. Dally. Deep compression: Compressing deep neural network with pruning, trained quantization and Huffman coding. *CoRR*, abs/1510.00149, 2015.
- [10] Wei Han, Pooya Khorrami, Tom Le Paine, Prajit Ramachandran, Mohammad Babaeizadeh, Honghui Shi, Jianan Li, Shuicheng Yan, and Thomas S Huang. Seq-nms for video object detection. *arXiv preprint arXiv:1602.08465*, 2016.
- [11] Yihui He, Xiangyu Zhang, and Jian Sun. Channel pruning for accelerating very deep neural networks. *CoRR*, abs/1707.06168, 2017.
- [12] Congrui Hetang, Hongwei Qin, Shaohui Liu, and Junjie Yan. Impression network for video object detection. *CoRR*, abs/1712.05896, 2017.
- [13] Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *CoRR*, abs/1704.04861, 2017.
- [14] Forrest N. Iandola, Matthew W. Moskewicz, Khalid Ashraf, Song Han, William J. Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and <1mb model size. *CoRR*, abs/1602.07360, 2016.
- [15] Samira Ebrahimi Kahou, Vincent Michalski, and Roland Memisevic. Ratm: recurrent attentive tracking model. In *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops*, pages 1613–1622, 2015.
- [16] Kai Kang, Hongsheng Li, Tong Xiao, Wanli Ouyang, Junjie Yan, Xihui Liu, and Xiaogang Wang. Object detection in videos with tubelet proposal networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 2, page 7, 2017.
- [17] Kai Kang, Hongsheng Li, Junjie Yan, Xingyu Zeng, Bin Yang, Tong Xiao, Cong Zhang, Zhe Wang, Ruohui Wang, Xiaogang Wang, and Wanli Ouyang. T-CNN: tubelets with convolutional neural networks for object detection from videos. *CoRR*, abs/1604.02532, 2016.
- [18] Kai Kang, Wanli Ouyang, Hongsheng Li, and Xiaogang Wang. Object detection from video tubelets with convolutional neural networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [19] Adam Kosior, Alex Bewley, and Ingmar Posner. Hierarchical attentive recurrent tracking. In *Advances in Neural Information Processing Systems*, pages 3053–3061, 2017.
- [20] Hugo Larochelle and Geoffrey E Hinton. Learning to combine foveal glimpses with a third-order Boltzmann machine. In J. D. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. S. Zemel, and A. Culotta, editors, *Advances in Neural Information Processing Systems 23*, pages 1243–1251. Curran Associates, Inc., 2010.
- [21] Laura Leal-Taixé, Cristian Canton-Ferrer, and Konrad Schindler. Learning by tracking: Siamese CNN for robust target association. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 33–40, 2016.
- [22] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. *CoRR*, abs/1608.08710, 2016.
- [23] Lanlan Liu and Jia Deng. Dynamic deep neural networks: Optimizing accuracy-efficiency trade-offs by selective execution. *CoRR*, abs/1701.00299, 2017.
- [24] Mason Liu and Menglong Zhu. Mobile video object detection with temporally-aware feature maps. *CoRR*, abs/1711.06368, 2017.
- [25] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016.
- [26] J. Luo, J. Wu, and W. Lin. Thinet: A filter level pruning method for deep neural network compression. In *2017 IEEE International Conference on Computer Vision (ICCV)*, volume 00, pages 5068–5076, Oct. 2018.

- [27] Anton Milan, Laura Leal-Taixé, Ian Reid, Stefan Roth, and Konrad Schindler. Mot16: A benchmark for multi-object tracking. *arXiv preprint arXiv:1603.00831*, 2016.
- [28] Anton Milan, Seyed Hamid Rezaatofighi, Anthony R Dick, Ian D Reid, and Konrad Schindler. Online multi-target tracking using recurrent neural networks. In *AAAI*, volume 2, page 4, 2017.
- [29] Deepak Mittal, Shweta Bhardwaj, Mitesh M. Khapra, and Balaraman Ravindran. Recovering from random pruning: On the plasticity of deep convolutional neural networks. In *Eighteenth IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2018.
- [30] Volodymyr Mnih, Nicolas Heess, Alex Graves, and Koray Kavukcuoglu. Recurrent models of visual attention. In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2, NIPS'14*, pages 2204–2212, Cambridge, MA, USA, 2014. MIT Press.
- [31] Bowen Pan, Wuwei Lin, Xiaolin Fang, Chaoqin Huang, Bolei Zhou, and Cewu Lu. Recurrent residual module for fast inference in videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1536–1545, 2018.
- [32] E. Park, D. Kim, S. Kim, Y. D. Kim, G. Kim, S. Yoon, and S. Yoo. Big/little deep neural network for ultra low power inference. In *2015 International Conference on Hardware/Software Codesign and System Synthesis (CODES+ISSS)*, pages 124–132, Oct 2015.
- [33] Marc’Aurelio Ranzato. On learning where to look. *CoRR*, abs/1405.5488, 2014.
- [34] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 779–788, June 2016.
- [35] J. Redmon and A. Farhadi. Yolo9000: Better, faster, stronger. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6517–6525, July 2017.
- [36] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *CoRR*, abs/1804.02767, 2018.
- [37] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 91–99. Curran Associates, Inc., 2015.
- [38] Amir Sadeghian, Alexandre Alahi, and Silvio Savarese. Tracking the untrackable: Learning to track multiple cues with long-term dependencies. *arXiv preprint arXiv:1701.01909*, 4(5):6, 2017.
- [39] Mohammad Javad Shafiee, Brendan Chywl, Francis Li, and Alexander Wong. Fast YOLO: A fast you only look once system for real-time embedded object detection in video. *CoRR*, abs/1709.05943, 2017.
- [40] Peng Tang, Chunyu Wang, Xinggang Wang, Wenyu Liu, Wenjun Zeng, and Jingdong Wang. Object detection in videos by short and long range object linking. *CoRR*, abs/1801.09823, 2018.
- [41] Swagath Venkataramani, Anand Raghunathan, Jie Liu, and Mohammed Shoaib. Scalable-effort classifiers for energy-efficient machine learning. In *Proceedings of the 52nd Annual Design Automation Conference*, page 67. ACM, 2015.
- [42] Swagath Venkataramani, Ashish Ranjan, Kaushik Roy, and Anand Raghunathan. Axnn: energy-efficient neuromorphic systems using approximate computing. In *Proceedings of the 2014 international symposium on Low power electronics and design*, pages 27–32. ACM, 2014.
- [43] Wei Wen, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Learning structured sparsity in deep neural networks. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 2074–2082. Curran Associates, Inc., 2016.

- [44] Yu Xiang, Alexandre Alahi, and Silvio Savarese. Learning to track: Online multi-object tracking by decision making. In *Proceedings of the IEEE international conference on computer vision*, pages 4705–4713, 2015.
- [45] Xizhou Zhu, Jifeng Dai, Lu Yuan, and Yichen Wei. Towards high performance video object detection. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [46] Xizhou Zhu, Yujie Wang, Jifeng Dai, Lu Yuan, and Yichen Wei. Flow-guided feature aggregation for video object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, volume 3, 2017.
- [47] Xizhou Zhu, Yuwen Xiong, Jifeng Dai, Lu Yuan, and Yichen Wei. Deep feature flow for video recognition. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, page 3, 2017.