

# ApproxDet: Content and Contention-Aware Approximate Object Detection for Mobiles

Ran Xu  
Purdue University  
xu943@purdue.edu

Chen-lin Zhang  
Nanjing University  
zhangcl@lamda.nju.edu.cn

Pengcheng Wang  
Purdue University  
wang4495@purdue.edu

Jayoung Lee  
Purdue University  
lee3716@purdue.edu

Subrata Mitra  
Adobe Research  
subrata.mitra@adobe.com

Somali Chaterji  
Purdue University  
schaterji@purdue.edu

Yin Li  
University of Wisconsin-Madison  
yin.li@wisc.edu

Saurabh Bagchi  
Purdue University  
sbagchi@purdue.edu

## ABSTRACT

Advanced video analytic systems, including scene classification and object detection, have seen widespread success in various domains such as smart cities and autonomous transportation. With an ever-growing number of powerful client devices, there is incentive to move these heavy video analytics workloads from the cloud to mobile devices to achieve low latency and real-time processing and to preserve user privacy. However, most video analytic systems are heavyweight and are trained offline with some *pre-defined* latency or accuracy requirements. This makes them unable to adapt at runtime in the face of three types of dynamism — the input video characteristics change, the amount of compute resources available on the node changes due to co-located applications, and the user’s latency-accuracy requirements change. In this paper we introduce ApproxDet, an *adaptive* video object detection framework for mobile devices to meet accuracy-latency requirements in the face of changing content and resource contention scenarios. To achieve this, we introduce a multi-branch object detection kernel (layered on Faster R-CNN), which incorporates a data-driven modeling approach on the performance metrics, and a latency SLA-driven scheduler to pick the best execution branch at runtime. We couple this kernel with approximable video object tracking algorithms to create an end-to-end video object detection system. We evaluate ApproxDet on a large benchmark video dataset and compare quantitatively to AdaScale and YOLOv3. We find that ApproxDet is able to adapt to a wide variety of contention and content characteristics and achieves 52% lower latency and 11.1% higher accuracy over YOLOv3, outshining all baselines.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](https://www.acm.org/permissions).

SenSys’20, November 16–19, 2020, Virtual Event, Japan

© 2020 Association for Computing Machinery.

ACM ISBN 978-1-4503-9999-9/18/06...\$15.00

<https://doi.org/10.1145/1122445.1122456>

## CCS CONCEPTS

• **Human-centered computing** → Ubiquitous and mobile computing systems and tools; • **Computing methodologies** → Tracking; **Object detection**; *Classification and regression trees*; **Neural networks**; *Learning linear models*.

## KEYWORDS

Object Detection, Mobile Vision, Resource Contention, Content-aware Approximate Computing, Machine Learning

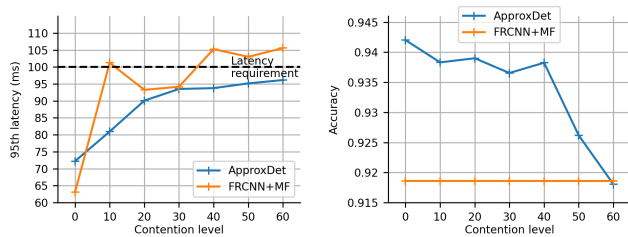
## ACM Reference Format:

Ran Xu, Chen-lin Zhang, Pengcheng Wang, Jayoung Lee, Subrata Mitra, Somali Chaterji, Yin Li, and Saurabh Bagchi. 2020. ApproxDet: Content and Contention-Aware Approximate Object Detection for Mobiles. In *Proceedings of The 18th ACM International Conference on Embedded Networked Sensor Systems (SenSys)*, Nov 16–19, 2020, Virtual Event, Japan. ACM, New York, NY, USA, 15 pages. <https://doi.org/10.1145/1122445.1122456>

## 1 INTRODUCTION

Mobile devices with integrated cameras have seen tremendous success in various domains. Equipped with increasingly powerful System-on-Chips (SoCs), mobile augmented reality (AR) devices such as the Microsoft HoloLens and Magic Leap One, along with top-of-the-line smartphones like iPhone 11 Pro and Samsung Galaxy S20, are opening up a plethora of new continuous mobile vision applications that were previously deemed impossible. These applications range from detection of objects around the environment for immersive experience in AR games such as Pokemon-Go [7], to recognition of road signs for providing directions in real-time [5], to identification of people for interactive photo editing [55], and to Manchester City’s AR-driven stadium tour. A fundamental vision task that all of these applications must perform, is object detection on the live video stream that the camera is capturing. To maintain the immersive experience of the user (e.g., for AR games) or to give usable output on time (e.g., for road sign recognition), such tasks must be performed in near real-time with very low latency.

Computer vision and computer systems research working together has made significant progress in lightweight object detection applicable to mobile settings for *still images* in recent years [39, 51, 59, 63], thanks to development of efficient deep neural networks (DNNs) [20, 24, 29, 70]. However, directly applying image-based



**Figure 1: ApproxDet: The first system of mobile video object detection that takes both video content-awareness and resource contention-awareness within its ambit. Compared to a widely used object detector [52] optimized for a target latency requirement (orange curve), ApproxDet (blue curve) keeps its runtime latency (left) below the requirement and achieves better accuracy (right).**

object detectors to video streams does not work well [71], especially in a mobile setting. *First*, applying a detector on all video frames introduces excessive computational cost and would often violate the latency requirements of our target continuous vision applications. *Second*, image-based object detectors are not cognizant of the significant temporal continuity that exists in successive video frames (e.g., a static scene with a slowly moving object) and therefore cannot leverage them for fitting in the latency budget. To overcome these algorithmic challenges, the computer vision community has proposed some DNN models [17, 34, 73] for video object detection and tracking. More recently, several lightweight DNN models [38, 72] that are suitable for mobile devices were developed.

Despite these efforts, we argue that the system challenges of video object detection for continuous vision applications on resource-constrained devices remain largely unsolved. A major shortcoming is that none of the existing approaches can adapt to runtime condition changes, such as *the content characteristics of the input videos*, and *the level of contention on the edge device*. Modern mobile devices<sup>1</sup> come with increasingly powerful SoCs having multiple heterogeneous processing units, and no longer process just a single application at a time. For example, both iOS and Android support multiple background tasks [3, 11, 12], such as an always-on personal assistant like Siri running a DNN for speech recognition (GPU contention), or a firewall constantly inspecting packets (memory bandwidth contention). These tasks can run simultaneously with a continuous vision application that requires a video object detector, leading to unpredictable resource contention on mobile devices similar to a traditional server setting [10, 42, 43, 67]. Such concurrent applications or background tasks can compete for resources critical to object detection, thus drastically increasing the latency of the object detector. Consider the example of a widely used DNN-based object detector – Faster R-CNN (FRCNN) [52], integrated with MedianFlow (MF) object tracking [33] and optimized for a latency requirement of 100 milliseconds (ms)<sup>2</sup>. The orange curve

in Figure 1 shows the latency (left) and accuracy (right) of this FRCNN+MF processing an input video at different GPU contention levels on an embedded device (NVIDIA TX2). Without contention, the detector has a latency of  $\approx 64$  ms. However, as the GPU contention level increases, we observe a severe increase in the detection latency. While the accuracy remains the same, the latency of the detector fluctuates significantly and violates our latency requirement of 100 ms. Different from server-class devices [67], our target mobile devices have less ability to isolate co-located applications from interference from one another. This happens due to the paucity of isolation mechanisms like VMs with resource reservation on our target class of mobile devices.

To address this issue, we propose **ApproxDet**, a novel system that takes both video content-awareness and resource contention-awareness within its ambit. In contrast to the static FRCNN+MF baseline in Figure 1, ApproxDet manages to keep a latency below the requirement with increased level of contention while achieving a better accuracy. To this end, ApproxDet uses a *single* model with multiple approximation knobs that are dynamically tuned at runtime to stay on the Pareto optimal frontier of the latency/accuracy tradeoff curve. We refer to the execution branch with a particular configuration of the approximation knob as an **approximation branch (AB)**.

This overall functionality is supported by *three core technical contributions* of this work. *First*, ApproxDet models the impacts of the contention level to the latency of the ABs. *Second*, our model combines an offline trained latency prediction model and an online contention sensor to precisely predict the latency of each AB in our system. Thus ApproxDet can adapt to resource contention at a given latency budget at runtime, an ability especially critical for the deployment on edge devices as their resources are limited and shared. *Third*, ApproxDet further considers how the video content influences both accuracy and latency. ApproxDet leverages video characteristics such as the object motion (fast vs. slow) and the sizes and the number of objects, to better predict the accuracy and latency of the ABs, and to select the best AB with reduced latency and increased accuracy.

Figure 2 presents an overview of ApproxDet. The object detection pipeline comprises of an object detector DNN based on our modified version of FRCNN [52], and a video object tracker that can be selected from among a set of choices. Both the detector and tracker are *jointly* approximated to achieve the required point in the accuracy-latency trade-off curve. Importantly, with the joint modeling of resource contention and content characteristics, ApproxDet can dynamically tune the approximation knobs, including the interval of performing object detection on video frames, the input shape of the frames and the number of proposals in the object detector DNN, the choice of object trackers, and the down-sampling ratio of the tracker. This means, in response to a resource contention from the GPU, ApproxDet can move to a more aggressive approximation setting of the detector DNN to bring down the latency since the detector DNN is more sensitive to the GPU contention. Moreover, in case of a content change in the video frame, e.g., a rapidly moving object, ApproxDet is able to switch over to a different AB where detector DNN is triggered more frequently to mitigate the tracker failure due to the fast-moving objects.

<sup>1</sup>In this paper, we use the term mobiles for the target platform, though without loss of generality, it applies to mobile and embedded platforms. The commonality is that both are using increasingly powerful SoCs, but are still resource constrained relative to the servers where streaming video analytics are typically run. Further, both are used to run co-located applications (mobiles more so than embedded), which can interfere with the video analytics.

<sup>2</sup>We pick a combination of detector and tracker configurations that satisfies the latency requirement in 95% of the video frames in a validation dataset.

**Table 1: A comparison of the key features of our ApproxDet solution to previous approaches. ApproxDet provides the most flexible framework for adaptive video object detection. “Single/multi model”: if a method uses shared execution branch for different control parameters. “Switching cost considered”: a technique takes into account switching cost while making its decision.**

Solution	Single/Multi Model (S/M)	Tuning Knobs	Switching cost considered	Dynamic Scenario	Open Source	Mobile/Server (M/S)	Video/Image
ApproxDet	S	si, shape, nprop, tracker, and ds	✔	✔	✔	M	VID
Faster R-CNN	S	shape, nprop	✘	✘	✔	S	VID
YOLO, SSD	S	shape	✘	✘	✔	S	VID
AdaScale	S	scale	✘	⚠	✔	S	VID
MCDNN [MobiSys16]	M	models	✘	✘	✘	M+S	IMG
NestDNN [MobiCom18]	S	# filters	⚠	✘	✘	M	IMG
RANet [CVPR20]	M	resource budget	✘	✘	✔	S	IMG

✔ Supported    ⚠ Partially Supported    ✘ Not Supported

To our best knowledge, ApproxDet is the first system that accounts for *the joint adaptation to video content and resource contention* for mobile vision applications. Distinct from prior work that optimizes multiple concurrent DNN applications [15, 21, 31], our system treats the contention as a black box. Our principle is that we neither know the context nor have control over the contention in real-world systems, as these video-analytic systems are typically user-space processes without any OS privilege. To estimate the contention, ApproxDet uses the current observed latency to map to the contention level, and adapts to use the AB that can satisfy the latency requirement from the user. Our evaluation bears out that this design choice is particularly effective for handling varied forms of contention under one simple algorithm. Table 1 further contrasts the features of ApproxDet with existing works.

To evaluate our model, we conduct extensive experiments on ImageNet Video Object Detection (VID) dataset and compare our ApproxDet to a number of baselines, including AdaScale [6], Faster R-CNN [52], Faster R-CNN with tracking [33], and YOLOv3 [51]. Our results suggest that ApproxDet is able to adapt to a wide variety of contention and content characteristics and achieves 52% lower latency and 11.1% higher accuracy over the latest YOLOv3 optimized for efficiency and accuracy, and outshining all other baselines.

In summary, our work makes the following contributions:

- (1) We show that contention in mobile/embedded devices can significantly degrade the latency requirements of continuous vision applications.
- (2) We propose ApproxDet— an adaptive object detection framework that takes the runtime content characteristics *and* resource availability into consideration to dynamically optimize for the best accuracy-vs-latency tradeoff. This optimization is done using a single model with different approximation branches, rather than

using an ensemble of models, reducing switching overhead and the memory footprint.

- (3) ApproxDet makes use of video-specific features, i.e., it does not consider video to be simply a set of discrete image frames. For example, it uses near past video content characteristics (such as, size of the objects) to guide its choice of the optimal approximation branch.

## 2 BACKGROUND

This section introduces the background of ApproxDet, including object detection and tracking models used in ApproxDet and the context of approximate computing on edge devices.

### 2.1 Object Detection

Given an input image or video frame, an object detector aims at locating tight bounding boxes of object instances from target categories. In terms of network architecture, a CNN-based object detector can be divided into the backbone part that extracts image features, and the detection part that classifies object regions based on the extracted features. The detection part can be further divided into two-stage [9, 52] and single-stage detectors [37, 39, 51]. Two-stage detectors usually make use of Region Proposal Networks (RPN) for generating regions-of-interest (RoIs), which are further refined through the detection head and thus more accurate.

Our work builds on Faster-RCNN [52] — an accurate and flexible framework for object detection and a canonical example of a two-stage object detector. An input image or video frame is first resized to a specific *input shape* and fed into a DNN, where image features are extracted. Based on the features, a RPN identifies a *pre-defined number* of candidate object regions, known as region proposals. Image features are further aggregated within the proposed regions, followed by another DNN to classify the proposals into either background or one of target object categories and to refine the location of the proposals. Our key observation is that the input shape and the number of proposals have significant impact to the accuracy and latency. Therefore, we propose to expose *input shape* and *number of region proposals* as tuning knobs in ApproxDet.

Another line of research develops single-stage object detection [39, 51]. Without using region proposals, these models are optimized for efficiency and oftentimes less flexible. We consider one of the representative single-stage detectors as a baseline [51] in our experiments. YOLO simplifies object detection as a regression problem by directly predicting bounding boxes and class probabilities without the generation of region proposals.

### 2.2 Object Tracking

Object tracking seeks to locate moving objects over time within a video. We focused on motion-based visual tracking due to its simplicity and efficiency. A motion-based tracker assumes the initial position of each object is given in a starting frame, and makes use of local motion cues to predict the object’s position in the next batch of frames. Our system considers a set of existing motion-based object trackers — MedianFlow [33], KCF [22], and CSRT [40]. The key difference lies in the extraction of motion cues, via e.g., optical flow or correlation filters, leading to varying accuracy and efficiency under different application scenarios. We thus propose to enable the adaptive *choice of the trackers* as one of our tuning knobs.

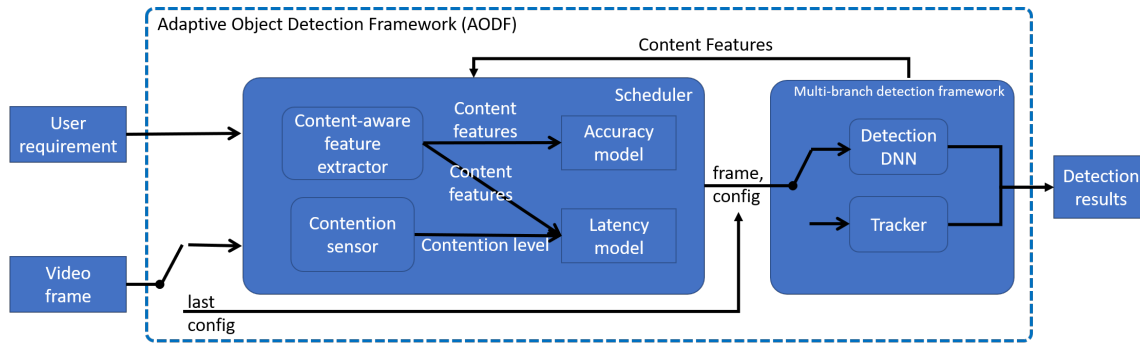


Figure 2: The workflow of the adaptive object detection framework.

Another important factor of object tracking performance is the input resolution to a motion-based tracker. A downsampled version of the input image allows to better capture large motion and thus to track fast-moving objects, while a high-resolution input image facilitates the accurate tracking of objects that move slowly. Therefore, we further expose *the downsampling ratio of the input image* as another tuning knob for tracking.

### 2.3 Approximate Computing and Adaptation

Many computations are inherently approximate—they trade off quality of results for lower execution time or lower energy. Approximate computing has emerged as an area that exposes additional sources of approximation at the computer system level, including resource-constrained mobile and embedded platforms. One challenge in approximate computing is that the accuracy and performance of applying approximate techniques to a specific application and input sets are hard to predict and control [1, 13, 47]. This may lead to missed optimization opportunities, unacceptable quality outputs, and even incorrect executions. The two fundamental causes is that approximation techniques are not content-aware and contention-aware. Some recent work has started to address these issues. For example, Input Responsive Approximation (IRA) [35] and VideoChef [66] have brought in content-aware approximation for image processing and video processing pipelines respectively.

To the best of our knowledge there is no solution that makes video analytics on mobile platforms adaptive to resource contention. Recently, Min *et al.* [46] assess the runtime quality of sensing models in the multi-mobile-device environment so that the best device is selected as a function of model accuracy. However, we have not seen similar work on the video object detection task. There are several works that provide tunable knobs to trade off accuracy-versus-latency, primarily in the image processing context [15, 27, 68], and some in video processing context [6]. It is conceivable that these knobs can be reconfigured to an optimal setting continuously as contention varies. However, there are key systems challenges that have to be solved before that end goal can be achieved. Such challenges include how to sense contention, how to change the knob in response to a specific level of contention, and how to optimize for the switching overhead from one approximation level to another.

## 3 OVERVIEW

Figure 2 presents the overall workflow of our system ApproxDet. Our system consists of two modules — a scheduler and a multi-branch object detection framework. The detection framework takes a video frame and a configuration as an input and produces the detection results while the scheduler decides which configuration the detection framework should use.

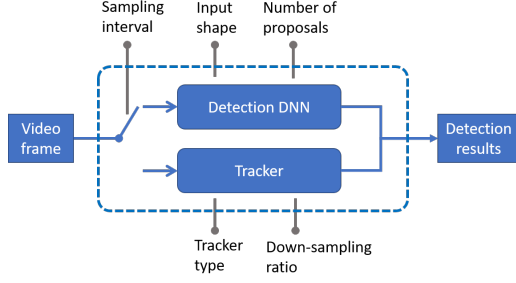
The detection framework includes two kernels: a detection kernel and a tracking kernel. This follows the common practice for video object detection that combines the heavy-weight detection and the light-weight tracker [38, 74]. At a high-level, the detection framework exposes five tuning knobs. With each tuning knob varying in a dynamic range, we construct a multi-dimensional configuration space and call the execution path of each configuration **an approximation branch (AB)**. The accuracy and the latency (execution time) are different for each AB and the values depend upon the video content characteristics (e.g., still versus fast-moving) and the compute resources available (e.g., lightly-loaded versus heavily-loaded mobile). To efficiently select an AB at runtime according to the given (and possibly changing) user requirement, the scheduler estimates the current latency and accuracy of each branch. The scheduler then selects the most accurate/fastest branch according to the specific user requirement.

We train an accuracy model and a latency model offline to support such estimation online. To better predict such online performance metric, we build two lightweight online modules – (1) a *content-aware feature extractor*, which extracts the height, width, tracks the object information of the last frame, and calculates the object movements of the past few frames, and (2) a *contention sensor*, which senses the current resource contention level. The scheduler is designed to run occasionally to re-calibrate the best approximation branch based on a learnable interval called “scheduler interval”, which represents the number of frames that the configuration of the detection framework can be maintained.

## 4 DESIGN ELEMENTS OF APPROXDET

### 4.1 Multi-branch Object Detection Framework

To support the runtime adaptive object detection framework on videos, we first design a multi-branch object detection framework, with light switching overheads between different branches so that it can quickly adapt to runtime changes. Different from object detection on still images, videos have temporal similarities and an



**Figure 3: Our multi-branch object detection framework ApproxDet with the five runtime tuning knobs.**

object tracker is widely used to reduce the runtime cost with minor accuracy drop. Similarly, we compose an object detection DNN and an object tracker. The detection DNN produces initial bounding boxes for each object in the input image, while the tracker tracks objects between successive frames.

The overwhelming majority of work on lightweight object detection (such as that suited for our target class of devices, mobiles or embedded devices) is for images, such as YOLOv3 [51] and SSD [39] and thus, does not leverage the video characteristics inherent in relation across image frames. For us, this leads to several insights and guides several design decisions.

For the detection DNN, we choose the popular Faster-RCNN with ResNet-50 as the backbone [52]. It shows state-of-the-art performance with medium speed when compared with other detection models. For the tracker part, we experiment with a set of 4 trackers — MedianFlow [33], KCF [22], CSRT [40], and Dense Optical Flow [16]. These trackers are open-sourced in OpenCV with reasonable performance. We then expose *five tuning knobs* for this object detection framework that our scheduler controls programmatically at runtime to achieve the right accuracy-latency tradeoff. We introduce them below and illustrate them in Figure 3.

- Sampling interval ( $si$ ): For every  $si$  frame, we run the heavyweight object detection DNN on the first frame and light-weight object tracker on the rest of the frames.
- Input shape ( $shape$ ): The resized shape of the video frame that is fed into the detection DNN.
- Number of proposals ( $nprop$ ): The number of proposals generated from the Region Proposal Networks (RPN) in our detection DNN.
- Tracker type ( $tracker$ ): Type of object tracker.
- Down-sampling ratio ( $ds$ ): The downsampling ratio of the frame used by the object tracker.

Generally, we have empirically observed that smaller  $si$ , larger  $shape$ , more  $nprop$ , and smaller  $ds$  will raise the accuracy and vice-versa. We will discuss the specifics of the knobs in Section 5.1.

## 4.2 Content Feature Extraction

As multi-branch object detection framework is designed, an important prerequisite is to *precisely estimate the accuracy and computation time (latency) of each approximation branch*. To start with, the content feature has great impact on both the accuracy and latency of each AB based on the following two observations – (1) tracker latency is affected by the number and area of the objects because tracker algorithms take the bounding boxes of the detection frames as inputs and calculate features inside each box; (2) both detection

and tracker accuracy are affected by the content in the video. For example, detection DNNs perform consistently poorly with small objects on MS COCO dataset, including Faster-RCNN [52], SSD [39], and YOLO [50]. Moreover, both detection DNN and tracker find it harder to deal with fast-moving objects. Some previous works [5] mention that movement between frames can be used as a feature to trigger the heavy detection process. This implies that for video object detection systems, we need to extract these content features to improve the accuracy and latency of our models. In this paper, we mainly consider two types of content features, described next.

**4.2.1 Object Basic Features.** We use *the number of objects* and *the summed area of the objects* as features for modeling the tracker latency. The intuition is that some light-weight trackers’ latency increases proportionally with the number of objects and the area of the objects since each object is tracked independently, and the larger the area, the more tracking-related features need computation. These features can be easily extracted with no extra cost from the output of the detection DNN.

We empirically verify that the latency of the object trackers is affected by both the number and sizes of the objects, as shown in Figure 4 and 5. Specifically, we use 10% of the ImageNet video object detection (VID) training dataset to generate the latency data samples. The detailed data split can be found in Section 6.2. We use all latency data in our validation set to plot Figure 4 and 5. Results have shown that the number of objects ( $n_{obj}$ ) and the average size of objects ( $avg\_size$ ) have a significant impact on the tracking latency, which we use as the object’s basic features.

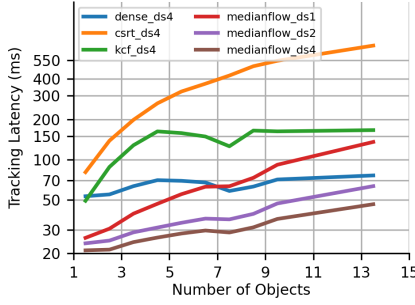
**4.2.2 Object Movement Features.** We use *the recent movement of objects* as a feature for modeling the framework accuracy. More rigorously, we define the movement as the Euclidean distance of the objects’ centers and we take the mean movement of all the objects in the recent frames. The intuition is that the faster the objects move in the video frame, the lower the accuracy, especially for the execution branches with higher sampling interval. We use the same data split as in Figure 4 and 5 to generate our accuracy data. We empirically show this in Figure 6, where we divide the validation dataset into three subsets — videos with slow, medium, and fast moving objects and show the accuracy reduction (compared to the detection-only branch) as we increase the sampling interval of the object detection kernel. For the detection kernel, we choose 100-proposal, 576-shape branch. The results show that the accuracy of high  $si$  branches ( $si = 100$ ) does not drop significantly ( $\approx 10\%$ ) on slow moving videos but reduces ( $> 30\%$ ) on fast moving videos.

## 4.3 Latency Modeling

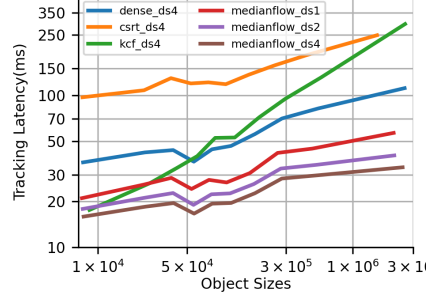
Latency prediction models aim to predict the frame-wise latency of each AB for future frames. Denote  $L_{fr}$  as the per-frame latency of our adaptive object detection framework.  $L_{fr}$  is a function of the DNN based detection latency  $L_{DNN}$  and the tracking latency  $L_{tracker}$ . If object detection DNN runs every  $si$  frames (sampling interval), the latency  $L_{fr}$  is given by

$$L_{fr} = \frac{L_{DNN}}{si} + L_{tracker}, \quad (1)$$

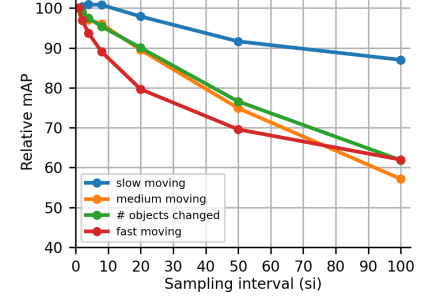
We now describe the models of the detection latency  $L_{DNN}$  and the tracking latency  $L_{tracker}$ , respectively.



**Figure 4: The latency curve of object trackers on different numbers of the objects on the validation dataset.**



**Figure 5: The latency curve of object trackers on different sizes of the objects on the validation dataset.**



**Figure 6: Detection plus Median-Flow tracker vis-à-vis detection-only branch on slow/medium/fast subsets.**

**4.3.1 Latency Prediction for Object Detection DNN.** The latency of the object detection DNN  $L_{DNN}$  is jointly determined by the two detector tuning knobs – the input image size *shape* and the number of proposals  $nprop$ . Moreover, considering the input shape of frames may vary in different videos, we add the *height* and *width* of the input image as additional features. These features could be ignored if the video source is a video camera (which outputs fixed sized frames). Besides the input shape of video frames, system *contention* (CPU/GPU usage and memory bandwidth, as detailed in Section 4.5) will also affect the DNN latency. Thus, the latency equation of the DNN is given by

$$L_{DNN} = f_{DNN}(nprop, shape, height, width, contention). \quad (2)$$

We fit a quadratic regression model for  $f_{DNN}$  to characterize the latency of the detection DNN. Once trained, the regression model is evaluated on a subset of the test set (sparsely sampled), where the mean squared error (MSE) between the prediction  $\hat{L}_{DNN}$  and the ground-truth  $L_{DNN}$  latency are reported.

**4.3.2 Latency Prediction for Object Trackers.** As discussed in Section 4.2.1, the number of objects and average sizes of objects play a major role for the tracking latency. We further construct a model  $f_{tracker}$  to characterize the latency of the object tracker under the system contention. Similar to the detection latency model, we also add the *height* and *width* of the input image as additional features. Thus,  $f_{tracker}$  is given by:

$$L_{tracker} = f_{tracker}(height, width, n\_obj, avg\_size, contention) \quad (3)$$

We fit quadratic regression models to the ground-truth  $L_{tracker}$ . Moreover, since the model depends on  $n\_obj$  and  $avg\_size$  of the previous frame, we use the previous frame’s  $n\_obj$  and  $avg\_size$  to train  $L_{tracker}$ . After the training process, we compute the predicted  $\hat{L}_{tracker}$  and measure the MSE across a subset of the test set.

## 4.4 Accuracy Modeling

Accuracy prediction models aim to predict the expectation of the accuracy of each AB for near future frames. Predicting a model’s accuracy in the future is a very challenging task. To this end, we start with a baseline, content-agnostic accuracy model that models

the average accuracy of approximation branches on an offline validation set. Improving on the baseline, we propose a content-aware model that estimates the accuracy based on the input video content.

The accuracy of an object detector is usually defined by the metric mean average precision (mAP). However, we find that predicting the absolute mAPs given a test video is difficult. To address this issue, we propose to convert the absolute mAP metric into a relative percentage metric. More precisely, a base branch is identified in the detection framework using the detection-only branch ( $si = 1$ ) with  $nprop = 100$  and  $shape = 576$ . This base branch sets the performance upper-bound for all approximation branches (62.3% mAP on the validation set). The mAP of each AB is normalized to its percentage value by dividing its mAP by the base branch’s mAP.

**4.4.1 Content-agnostic Accuracy Modeling.** We first present a baseline model that predicts the average accuracy of each AB on a target dataset, i.e., the validation set. Different from the latency models, the factors on the accuracy are coupled all together (i.e., no distinction between detection DNN and tracking). Thus, we have a single unified model, given by:

$$A = f_A(si, shape, nprop, tracker, ds) \quad (4)$$

where *tracker* is the tracker type and *ds* is the downsampling ratio of the input to the tracker. A decision tree model  $f_A$  was learned to predict the accuracy  $A$ , trained with the MSE loss across the whole training dataset.

**4.4.2 Content-aware Accuracy Model.** A content-agnostic accuracy model is independent of the video content, and thus often fails to predict the accuracy of our detection framework on individual videos that may differ significantly from the average accuracy on the validation set. We further design a content-aware accuracy model that predicts the accuracy of an AB taking the video characteristics into account. Our model is presented in Eq. 5 where both the configurations of the detection framework and content characteristics are taken into consideration.

$$A = f(si, shape, nprop, tracker, ds, movement) \quad (5)$$

where *movement* is the object movement features extracted from the video content. We find that linear regression models demonstrate 8% improvement (in MSE) vis-à-vis the decision tree models used in our content-agnostic models.

## 4.5 Synthetic Contention Generator

Synthetic Contention Generator (CG) is designed as a stand-in for any resource contention on the device that may affect ApproxDet. A detection framework may suffer from unpredictable levels of resource contention when it is running on mobile platforms due to the instantiation of other co-located applications, for which we will not have information. We focus on three important types of resources on mobile platforms – CPU, memory bandwidth (MB), and GPU. We control CPU contention by the number of CPU cores our CG occupies. We control MB contention by the amount of memory-to-cache bandwidth that it consumes. The code is modified from the widely used STREAM benchmark [44, 45] that is meant to measure the MB capacity of a chip. For the GPU contention, we control the number of GPU cores that are utilized. The three-dimensional CG is orthogonal, which means we can tune each dimension without affecting the other dimensions. The CG is representative because we executed and mapped the contention caused by some widely-used applications in the 3D contention space (Table 2). The first one is an anomaly detection program that uses Robust Random Cut Forest (RRCF) [19] to detect anomalies from a local temperature and humidity sensor data. We also used our two object detection DNNs, Faster R-CNN and YOLOv3, for checking how much contention they generate.

**Table 2: Applications running in the 3D contention space.**

Real Apps	CPU	MB (MB/s)	GPU
Anomaly detection	99.80%	500	0%
Faster R-CNN	69.75%	1000	99%
YOLOv3	65.85%	800	98.50%

## 4.6 Profiling Cost and Sub-sampling

The cost of collecting ground truth data with design features for performance prediction models is significant without proper sampling techniques. We measure our profiling cost for the accuracy, detection latency, and tracker latency models in Table 3.

To efficiently collect the profiling data, we use the master and worker model, where the master node manages a list of configurations of the detection framework and distributes the profiling work, while workers run the particular configuration to collect the training data for the modeling. As the feature space is huge, we sparsely sample the multi-dimensional space of (“number of proposals”, “re-sized shape”, “sampling interval”, “tracker”, “down-sampling ratio of the tracker”). We finally use 20% of the configurations to train our accuracy model.

Similar sub-sampling techniques are used for the latency models as well, and we sample data points on videos of various height and width, various numbers of objects and object sizes, under discrete 3D contention levels. We finally use 15 out of a million feature points (defined in Section 5.2 and 5.1) to train our detection latency model and 169 out of a million feature points to train our tracker latency model.

## 4.7 Scheduler

The scheduler is the core component of ApproxDet that makes the decision at runtime on which AB should be used to run the inference

**Table 3: Cost of profiling.**

Task	Cost
Framework accuracy	2,414 hr · core (20% of the configurations)
Detection latency	7 hr · machine w/ 15 out 1 million sampling
Tracker latency	1 hr · machine w/ 169 out 1 million sampling

on the input video frames. Formally, the scheduler maximizes the estimated detection accuracy of ApproxDet given a latency requirement  $L_{req}$ . This is done by identifying a feasible set of branches that satisfy the target latency requirements, and choosing the most accurate branch. In case of an empty feasible set, the fastest branch is returned. Thus, we formulate the optimal AB  $b_{opt}$  as follows,

$$b_{opt} = \begin{cases} \operatorname{argmax}_{b \in \hat{\mathcal{B}}} (A_b), & \text{if } \hat{\mathcal{B}} \neq \emptyset, \\ \operatorname{argmin}_{b \in \mathcal{B}} (L_{est,b}) & \text{otherwise} \end{cases} \quad (6)$$

where  $\hat{\mathcal{B}}$  is all ABs considered,  $\hat{\mathcal{B}}$  is the feasible set, i.e.,  $\hat{\mathcal{B}} = \{b \in \mathcal{B} \mid \text{iff } L_{est,b} < L_{req}, A_b \text{ and } L_{est,b} \text{ are the estimated accuracy and latency of the AB respectively. The search space } \mathcal{B}, \text{ composed of five orthogonal knobs, has millions of states. To reduce the scheduler overhead, we use a sampling technique in Section 5.1 and design light-weight online feature extractors.}$

To further reduce the scheduler overhead and enhance our system robustness, we restrict the scheduler to make decision at least every  $sw$  frames. The motivation of introducing  $sw$  is to prevent the scheduler to make very frequent decisions. Specifically, we set  $sw = \max(8, si)$ . The scheduler will thus make a decision at least every 8 frames. When the scheduler chooses a branch with a long  $si$ , it will make a following decision every  $si$  frames. In addition to the latency of the detection and tracking kernels, we add switching overhead  $L_{sw}$  and the scheduler overhead  $L_{sc}$  into the overall latency estimation of a AB  $b$ , i.e.,  $L_{est,b} = L_{b,fr} + (L_{sw} + L_{sc})/si$ . We also design the light-weight online feature extractors so that we can adapt seamlessly to the content and contention changes.

**Online content feature extractor** The online content feature extractor maintains the content features of the video by extracting *height*, *width* from current frame, memorizing  $n_{obj}$ , *avg\_size* of last frame and *movement* from past frames. It is lightweight in terms of the compute load it puts on the target platform and this is desirable since we have to extract the features at runtime on the target board for feeding into our models.

**Online contention sensor** The online contention sensor is designed to sense the contention level in the system so that we can refer to the modeling and make the right prediction on the latency of each AB. There are generally two approaches—one is to probe the system by directly measuring CPU, memory bandwidth and GPU usage by other processes, and the other is to record the latency of ApproxDet in the past few runs and match with the offline log. Although the first one can theoretically get the ground truth of the resource contention, it is not practical. As a normal application in the user space, it is difficult for ApproxDet to collect the exact resource information from other processes. The hardware is also lacking sufficient support for such fine-grained measurement on mobile or embedded devices [60]. In contrary, the offline latency log under various contention levels and the online latency log of the current branch in the past few runs are a natural observation of

the contention level. Thus, we proposed the log-based contention sensor.

The log-based contention sensor tries to find a contention level where the offline latency log matches the averaged online latency most closely. We use the nearest-neighbor principle to search for such contention levels in our pre-defined orthogonal 3D contention space. As multiple contention levels may cause the same impact on the latency of a given AB, we call it a cluster of contention levels and we pick one level out of it as the representative. In comparison to some previous work in the systems community [15], our contention sensor is lightweight, efficient, and does not require additional privileges at system level, making it a more practical offering in real-world systems.

## 5 IMPLEMENTATION

We implement ApproxDet in Python 3 and C with tensorflow-gpu, CUDA, and cuDNN libraries and release our code at <https://github.com/StarsThu2016/ApproxDet>. For detection kernel, we choose Faster R-CNN due to its high accuracy and moderate computational burden. For tracking kernel, we implement four variants and introduce the details in Section 5.1.

### 5.1 Configuration of the Tuning Knobs

Our five tuning knobs include the sampling interval ( $si$ ), the input image size ( $shape$ ) to the detection DNN, the number of proposals ( $nprop$ ) in the detection DNN, the type of object tracker ( $tracker$ ) and the downsampling ratio of the input to the tracker ( $ds$ ). We now describe the implementation details of these knobs, including their data types and value ranges.

**Sampling Interval ( $si$ ).**  $si$  defines the interval of running the object detector. The object tracker runs on the following  $si - 1$  frames. For example, our system runs object detection on every frame when  $si = 1$ . To reduce the search space of  $si$ , we constrain  $si$  in a preset set—{1, 2, 4, 8, 20, 50, 100}. These pre-defined  $si$  are chosen empirically to cover common video object detection scenarios. With the max value of  $si = 100$ , the detector runs at a large interval of 3-4 seconds and the tracker runs in-between.

**Input Video Frame Shape to Detector ( $shape$ ).**  $shape$  defines the shortest side of the input video frame to the object detector. The value of  $shape$  must be a multiple of 16 to make the precise alignment of the image pixels and the feature map [52]. We set the  $shape$  range from 224 to 576, since smaller shape than 224 significantly reduces the accuracy and larger shape than 576 will result in heavy computational burden and does not improve the accuracy based on results on the validation set.

**Number of Proposals ( $nprop$ ).**  $nprop$  controls the number of candidate regions considered for classification in the object detector. We limit the value of  $nprop$  (integer) between 1 and 100. With  $nprop = 1$ , only the top ranked proposal from RPN is used for detection. Increasing  $nprop$  will boost the detector's performance yet with increased computational cost and runtime.

**Type of Trackers ( $tracker$ ).**  $tracker$  defines which tracker to use from MedianFlow [33], KCF [22], CSRT [40], and dense optical flow trackers [16]. These trackers are selected based on their efficiency and accuracy. Different trackers have varying performance under different scenarios. For example, CSRT tracker is most accurate among these trackers, but is also most time consuming. MedianFlow

tracker is fast and accurate when a object move slowly in the video, yet have poor performance for a fast moving object. We use the implementation from OpenCV for all trackers.

**Downsampling ratio for the tracker ( $ds$ ).**  $ds$  controls the input image size to the tracker. The value of  $ds$  is limited to 1, 2, and 4, i.e., no downsampling, downsampling by a factor of 2 and 4, respectively. A larger  $ds$  reduces the computational cost, and favors the tracking of fast moving objects. A smaller  $ds$  increase the latency, yet provide more accurate tracking of slowly moving objects.

### 5.2 3D Contention Generator (CG)

Our 3D CG is lightweight in code size. It is configured to generate contention in CPU, memory bandwidth (MB), and GPU (as introduced in Section 4.5). For MB CG, we modify STREAM by having the code write continuously to a 152 MB memory space and controlling the interval of array elements to operate on the array data so as to control the MB occupied. We add a feedback loop to dynamically adjust the number of elements in the array to be skipped for the write operation, thus maintaining the MB contention at the experimentally given level. To increase the maximum contention that can be generated by the MB CG, we spread out the MB CG among the CPU cores on which contention is to be generated and in aggregate can generate an intense bandwidth contention of up to 18 GB/s when 6 CPU can be used. The maximum input for the MB CG depends on both the CPU and MB part, the more CPU we can occupy, the higher maximum MB contention we can achieve. For the GPU CG, we fixed the working frequency of TX2 board at 1300 MHz. Then our GPU CG performs add operation on a certain size of arrays by using a CUDA program. By changing the size of the arrays and the input to the CUDA kernel functions, we control the number of GPU cores that are kept busy. We generate contention from 1% to 99%, as measured by `tegrastats`. For the experiments, we use 11 discrete levels 1%, 99%, and 9 levels in increments of 10% starting at 10%. If there is no MB or GPU contention specified and we need CPU contention on a certain number of CPU cores, we use a MB CG with minimum input 1MB/s to serve as the CPU CG and pin it to the experimentally given number of cores. For the experiments, we use 6 discrete levels from 100% to 600% in increments of 100%.

### 5.3 Training of Latency and Accuracy Models

The latency and accuracy model in ApproxDet is trained using a subset of the validation set of ImageNet VID. Specifically, we sparsely sample the 5-D feature space ( $si$ ,  $shape$ ,  $nprop$ ,  $tracker$ ,  $ds$ ) and run ApproxDet using the sampled configuration on a subset of videos randomly sampled from the validation set. Standard gradient descent is then used for fitting the regression. And CART is used for the decision tree. To train content-aware accuracy model, during the training phase, the *movement* feature as the average motion across all objects and all frames is required in each video snippet. During the test phase, since movement of the objects are not available, we use average across all past frames as a substitute.

## 6 EVALUATION

### 6.1 Evaluation Platform

We evaluate ApproxDet on an NVIDIA Jetson TX2 board [8], which includes 256 NVIDIA Pascal CUDA cores, a dual-core Denver CPU,



a quad-core ARM CPU on a 8GB unified memory between CPU and GPU. The specification of this board is comparable to what is available in today's high-end smartphones such as Samsung Galaxy S20 and Apple iPhone 11 Pro. We train our neural network models on a server with NVIDIA Tesla K40c GPU with 12GB dedicated memory and an octa-core Intel i7-2600 CPU with 24GB RAM. For both the TX2 and the edge server, we install Ubuntu OS and Tensorflow v1.14, Pytorch v1.1, and MXNet v1.4.1.

## 6.2 Datasets, Task, and Metrics

We evaluate ApproxDet on the object detection task using ILSVRC 2015 VID dataset [53]. For the purpose of training, since ILSVRC 2015 VID training set is a video dataset, due to the redundant video dataset and limited resources, we follow the practice in [34] such that the VID training dataset is sub-sampled every 100 frames. We use 90% of this video dataset as training set to train ApproxDet's DNN model and keep aside another 10% as validation set to fine-tune ApproxDet (modeling). To evaluate ApproxDet's system performance, we use ILSVRC 2015 VID validation set – we refer to this as the “test set” throughout the paper. For all our baselines, we follow the data split of ApproxDet's DNN model to train and test.

We use latency and mean average precision (mAP) as the two metrics. We define the latency as the short-window averaged latency among one detection frame and its following tracking frames. The definition applies to ApproxDet's baselines with different trackers used. The latency also includes the overheads of the respective solutions, e.g., the switching overhead, the execution time of the online feature extractor, the online contention sensor, and the scheduler. For common detection methods, they will use low confidence threshold, e.g., 0.001 [39, 52] to achieve higher mAP. However, low confidence threshold will lead to many false positive outputs, which is not practical in real-world systems. To simulate the environment of real-world systems, we use a high confidence threshold of 0.3 to reduce the number of output objects of detection DNN, and make it the same for all baseline detection DNNs. Note that we always use the ground truth annotations available in the dataset to examine the true accuracy and never use the detection results of some state-of-the-art models as pseudo ground truth (fake annotations as in [25]).

## 6.3 Baselines

In this section, we will introduce baseline models used in our paper. **Faster R-CNN (FRCNN)** [52] is a popular two-stage detector in the computer vision community. In this paper, we vary the  $nprop$  and  $shape$  of ApproxDet's detection DNN to create different FRCNN baseline detectors. From this, we get a total of 28 ( $4shape \times 7nprop$ ) FRCNN baseline models. In these baseline models, we follow the multi-model design on the detection model, where model variants in different sizes (number of proposals and resized shape) achieve different latency and accuracy specifications. We also profile the latency of these model variants and evaluate FRCNN with our scheduler defined in Section 4.7.

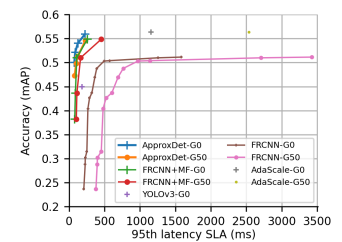
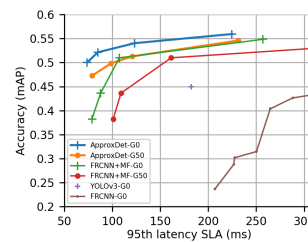
**FRCNN+MedianFlow (FRCNN+MF):** Since FRCNN only uses detection, we want to enhance this baseline for a fair comparison. For the enhanced baseline, we follow the mainstream “detection plus tracking” design in lightweight object detection tasks. We pick an unmodified detection variant with the highest accuracy ( $nprop =$

100,  $shape = 576$ ) and a fast object tracker—MedianFlow (without downsampling technique). Thus, the enhanced baseline models include FRCNN plus MedianFlow tracker with varying  $si$ . We also profile the latency of these model variants (each  $si$ ) and evaluate FRCNN+MF with our scheduler defined in Section 4.7. To reduce the scheduler cost, we perform the same sampling strategy in Section 5.1. Though we can pick different models based on the latency budget, these models are static and cannot respond to contention and context changes.

**AdaScale:** Among latest methods, we choose AdaScale [6] as it dynamically adjusts the input scale to improve accuracy and running speed simultaneously. AdaScale is thus most relevant to our work, with similar tuning knobs for resource-constrained devices. For conducting a fair comparison between AdaScale and our proposed ApproxDet, we re-implemented AdaScale and modified the following settings:

- Pretrained settings: In AdaScale, they use pretrained models on ImageNet VID + ImageNet Detection (DET). In this paper, ApproxDet and AdaScale both start from ImageNet pretrained models.
- Training sets: In AdaScale, they use ImageNet detection dataset (DET) joint with ImageNet VID dataset and subsample every 15 frames in each VID video. In this paper, we remove all training images in ImageNet DET dataset and we use same sampling techniques (sub-sampled every 100 frames) for training Adascale.

**YOLOv3:** YOLOv3 is a modern one-stage detector with a fast running speed. It is widely used in many mobile applications. In this paper, we use the same training/testing set and scheduler as ApproxDet to train and evaluate YOLOv3.



**Figure 7: Accuracy of the models vs 95-th latency SLA. G50 represents the 50% GPU contention and G0 represents no contention. (zoomed in)**

**Figure 8: Accuracy of the models vs 95-th latency SLA. G50 represents the 50% GPU contention and G0 represents no contention. (zoomed out)**

## 6.4 End-to-End Evaluation on Budgeted Latency

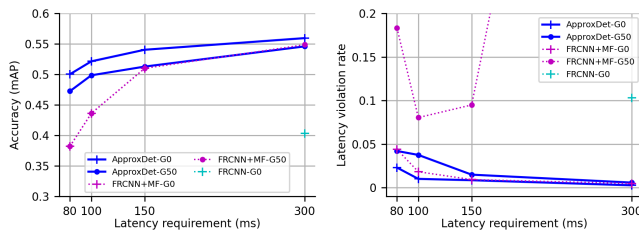
We first examine the end-to-end performance of ApproxDet vs. various baselines. Figure 7 shows a micro-view of the accuracy and 95-th percentile latency<sup>3</sup> under no contention (G0) and injected 50% GPU contention (G50) from the CG. The reason why our evaluation injects GPU contention only, is that the inference time of the DNN

<sup>3</sup>We choose the 95-th percentile latency service level agreement (SLA) because it is a promise to the users that in most cases the latency is below the number and shows much stronger latency guarantee than either mean or median latency.

has the most impact on the detection latency and such inferences mostly run on the device's GPU (when it is available). Hence, the impact of GPU contention would focus on the most interesting scenarios showing how resilient ApproxDet is with respect to changes in the dynamics within the mobile device.

The results show the superior accuracy-vs-latency trade-off (blue vs. green vs. purple) over a FRCNN+MF and YOLOv3 under no contention where particularly, ApproxDet's accuracy is 11.8%, 9.5% higher than FRCNN+MF given 80ms and 90ms latency SLA and also 9.1% higher than YOLOv3 given 170ms latency SLA (exactly following YOLOv3's posterior 95-th latency). The accuracy gain over FRCNN+MF reduces as latency SLA grows larger since both ApproxDet and FRCNN+MF will merge into detection-only (or detection in every alternate frame) branch.

Furthermore, when 50% GPU contention is injected, ApproxDet with strong adaptability on sensing and reacting to the contention, the accuracy drops by around 2% and preserves the 95th latency SLA with minor changes (max 16%). However, the best baseline FRCNN+MF cannot adapt to the contention and the latency increases significantly (25% - 50%). Figure 8 shows the performance of FRCNN and AdaScale. Both of their latency increase by 75% to 120% with increased contention levels. This is significantly worse than ApproxDet. Further, FRCNN is inferior to FRCNN+MF in terms of accuracy-latency tradeoff. This is because adding the light-weight MedianFlow tracker largely reduces the latency while preserving most of the accuracy.



**Figure 9: Given latency requirement, accuracy (mAP). (FRCNN+MF)'s accuracy does not change since it cannot adapt to contention and the execution is the same. FRCNN does not work for the latency SLA < 300 ms.**

Although accuracy vs latency SLA shows the posterior performance metric of each model, we also want to examine under a certain latency requirement, how each solution chooses a model variant to execute and what the accuracy (Figure 9) and the latency violation rates (Figure 10) are. We evaluate at 80ms, 100ms, 150ms, and 300ms. Latency requirement that is smaller than 80ms is not chosen because no branch will be returned from both baselines and larger latency requirement is not meaningful as discussed before. Figure 9 and 10 show that under no contention scenario (GO), ApproxDet and FRCNN+MF are equally good at controlling low latency violations (ApproxDet is slightly smaller), while the

accuracy of ApproxDet is 11.8%, 8.5%, 3.0%, and 1.1% higher than FRCNN+MF. Even under great reduction to 80ms latency requirement, ApproxDet is still able to maintain the accuracy above 50%.

Then, as 50% GPU contention is injected, the accuracy of ApproxDet is slightly reduced (by 2%) and is as good as FRCNN+MF under 150ms and 300ms latency requirements. However, ApproxDet is still able to control within the 5% violation rate while FRCNN+MF has a violation rate of 18.3%, 8.0%, 9.5%, and 100%. To summarize, ApproxDet is best at reducing the latency requirement to as low as 80ms with slightly reduced accuracy and is able to adapt well to the contention without violating the latency requirements.

## 6.5 Case Studies on Changing Conditions

Although the macro benchmark shows the overall performance of ApproxDet on a whole dataset, it is mostly the static behavior of the models since the latency requirement and resource contention levels do not change. To examine how ApproxDet adapts to the runtime changing conditions, we set up these case studies by injecting changing contention levels and latency requirements.

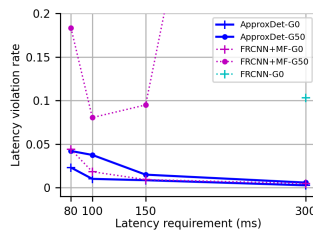
### Changing Contention Levels.

We first examine the performance of ApproxDet given a fixed latency requirement of 100 millisecond on one test video. We then manually inject the GPU contention through CG and gradually tune up the contention level by 20% for every 200 frames. Figure 11 shows that ApproxDet with quick sensing and adaptation, is always able to meet the latency requirements while a FRCNN+MF baseline will either exceed the latency requirement by 20% ( $si = 20$ ) or stay too conservative ( $si = 50$ ) and suffers from low accuracy (there is a 7% mAP difference between the two branches on the test dataset).

In addition, we also pick a real application—Gaussian Elimination from the Rodinia Benchmark Suite [4], a GPU-intensive linear algebra routine widely used by many applications. We examine the performance of ApproxDet and baselines without and with the background app. Figure 12 shows that ApproxDet can adapt to the resource contention produced by the background app. The latency goes up dramatically when the app starts running and quickly drops as ApproxDet senses such contention and schedule a more efficient AB. A repeated experiment during the 600–800 frames has further confirmed our adaptability. In contrast, the baseline FRCNN+MF with different  $si$  are either too conservative or too aggressive under varying contention levels in terms of latency. When the contention app starts running, there is a larger latency spike of ApproxDet than FRCNN+MF. This is because our system schedules the AB with smaller  $si$  as long as the latency of such AB is below the latency requirement. However, since the object detector takes a larger portion in latency and is more sensitive to GPU contention, the latency increase is much higher when the system has not responded yet.

**Changing Latency Requirements.** We then examine the performance of ApproxDet given more relaxed latency requirements of 80ms, 100ms, 150ms, and 300ms per frame in four equally chunked phases of 200 frames. Figure 13 shows that we can always keep up with the latency requirement and always run below the latency requirement while FRCNN+MF, although is choosing a branch that satisfy 95% of the video frames in the validation dataset, still violates the latency requirement by 20% under 80ms requirement and is slightly above the threshold given 100ms latency requirement.

**Figure 10: Given latency requirement, latency violations. The latency violation is higher with contention although FRCNN+MF is using a very conservative p95 scheduler.**



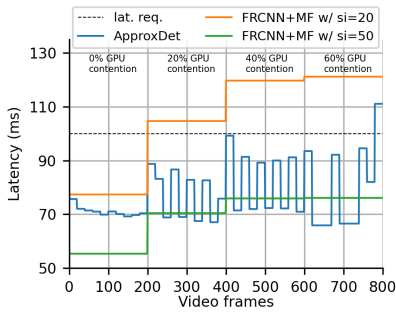


Figure 11: Models' latency under changing contention.

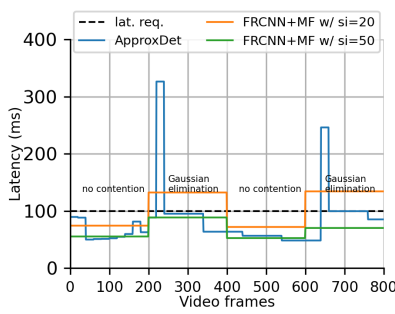


Figure 12: Models' latency without and with real app.

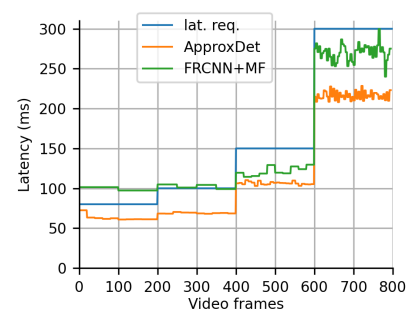


Figure 13: Models' latency under changing latency budget.

### 6.6 Micro-benchmark: Performance Prediction Models

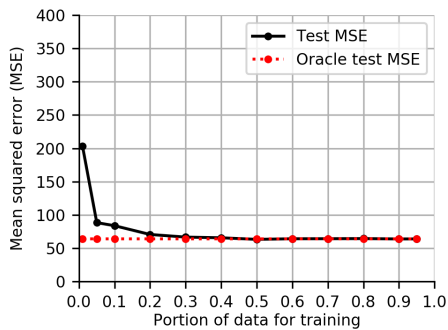


Figure 14: MSE of the accuracy prediction model given certain amount of training data in the validation dataset.

6.6.1 Accuracy prediction model. We set up the oracle method as using the ground truth validation data to predict on the test dataset. According to our belief that the accuracy reduction should be same in validation dataset and test dataset, the oracle approach should achieve zero MSE and if it is not zero, it represents the gap of accuracy reduction between the two datasets.

Figure 14 shows the training curve of the accuracy prediction model. We use different amounts of data to train the model and examine the mean squared error (MSE) on the rest of validation dataset for cross validation and the whole test dataset to report final performance. We can see that with only 20% of the training data, we are able to predict on the test dataset with 74.58 MSE. We can further reduce the MSE to 71.67 if 95% of the training data is available, while the oracle predicts with a comparable 71.65 MSE.

6.6.2 Tracker latency prediction model: We set up the baseline approach, which always predicts a constant latency for each tracker using the averaged latency across all frames under specific contention levels.

As seen from Table 4, we can largely reduce the prediction root mean squared error (RMSE) on the test dataset. Some trackers return high prediction RMSE (CSRT for example). The reason is some trackers may have unstable latency under high contention levels. We leave further exploration of this as our future work.

Table 4: Precision of the tracker latency prediction models on the validation dataset. Please note that we only show the results of tracking frames on test datasets. In this table, “no” means no contention, “g50” means 50% GPU contention, and “c6m3600” means contention with 6 CPU cores and 3600 memory bandwidth. RMSE means root squared mean squared error. We use “our model/baseline” formats to show the result.

	RMSE (no)	RMSE (g50)	RMSE (c6m3600)
Medianflow_ds1	11.98/17.86	6.37/19.83	14.81/44.41
Medianflow_ds2	6.80/12.37	3.14/12.90	9.91/26.96
Medianflow_ds4	7.32/12.30	3.72/12.43	11.14/26.18
KCF_ds4	32.23/41.24	23.81/43.38	41.46/81.37
CSRT_ds4	46.66/92.54	44.17/102.98	78.97/179.77
Dense_ds4	14.15/24.45	6.45/26.32	11.75/53.26

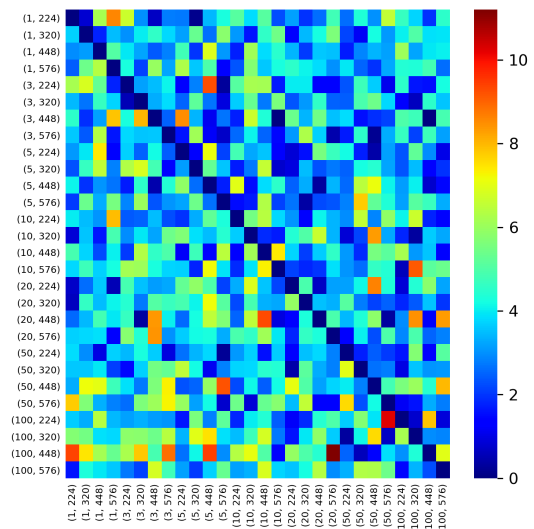


Figure 15: Heatmap view of switching latency among AB with varying “numbers of proposals” and input “shapes”.

## 6.7 Overhead: Switching, Scheduler and Online Components

Figure 15 shows a visualization of a heatmap of the switching latency from any approximation branch of the detection framework to another, especially inside the detection DNN. This is an average of 10 such switches and the switching latency is defined as the difference of the execution time of the first frame in the new branches over that of following frames. We can find that switching latency is bounded by 12 milliseconds.

**Table 5: Latency percentage of different parts in ApproxDet system. These are measured with zero contention.**

User requirement	Scheduler overhead (%)	Detection latency (%)	Tracking latency (%)
80ms	0.82%	59.65%	39.53%
100ms	0.62%	66.63%	32.75%
150ms	0.52%	69.10%	30.38%

The scheduler overhead comes from the accuracy and latency prediction models, as well as the contention sensor. Generally, the overhead of the scheduler is 11.09 milliseconds per execution. Since we will only execute the scheduler one time during all *si* frames, the average overhead of the scheduler is under 1 millisecond in most cases. This is supported by the profiling results in Table 5. The overhead of the scheduler occupies less than 1% of the total latency, suggesting that our scheduler is sufficiently fast to satisfy the requirement of an online system.

## 7 DISCUSSION AND FUTURE WORK

**Relation to the OS and Standing on the Contention.** ApproxDet is positioned as a user-level application and does not change the OS's orchestration. Thus, we design ApproxDet to treat contention in a black-box manner. The advantage of such standing makes ApproxDet easier to implement, free from the OS restrictions, and reduce the complexity of contention scenarios by observing the impact of contention on ApproxDet's latency. However due to the black-box method, ApproxDet has the limitation of not knowing the exact contention details to adapt accordingly. Different contention scenarios may have the same impact on ApproxDet's latency and however ApproxDet will not be able to differentiate them, leading to potential sub-optimal adaptation. Future work may add OS privilege, observe the true contention levels from the OS, and even control the contention as well. The marginal benefit over the cost and lost features is yet to be revealed.

**Contention Scenarios.** We evaluate ApproxDet mainly with GPU contention from the synthetic CG because the object detector mainly runs on the GPU and the detection latency is much higher than the tracker latency. ApproxDet has the limitation on the contention source and scenarios, though we also include a case study with a real App. Future work may evaluate with various mobile background workloads from the real trace data and study the propagation effect of ApproxDet to the OS or other Apps.

**Generalize to Other Detection Methods.** Our adaptive detection framework allows any architectures of detection DNNs beyond FRCNN as long as we can expose the tuning knobs from them. The

choice of the detection method can be another approximate knob as well.

**Features of Performance Prediction Models.** In ApproxDet, we consider the content features motivated by Figure 4, 5 and 6 with clear individual impact. We have considered the low-level content feature like the edges in the video frame but it does not have clear impact and thus has been excluded. More features can be introduced to improve the accuracy and latency prediction models, like object type, shape deformation, context, etc.

**Reinforcement Learning (RL) for the Scheduler.** RL based models can be an alternative method for implementing the scheduler. RL models learn to make schedule decisions based on training data and can potentially outperform the rule based policies.

**Using Neural Networks to Infer Configurations.** We leverage the sampling techniques to reduce the searching cost of configurations. Neural networks can be used to improve the searching efficiency if we disable the sampling and allow more flexible choices.

**Evaluation Dataset and Devices.** We plan to further evaluate ApproxDet on data-sets with high resolution videos (e.g., 1080p and 2K videos) and other mobile platforms, e.g., devices of lower computation power and smartphones.

## 8 RELATED WORK

**Object Detection.** Object detection is a well established topic in computer vision and has made recent progress, thanks to DNNs. DNN-based object detectors can be summarized into two categories: single-stage detectors, such as YOLO series [50, 51], SSD [39], RetinaNet [37], and two-stage detectors, such as Faster-RCNN [52], R-FCN [9], Cascade R-CNN [56]. Single-stage detectors classify a dense set of grids, while two-stage detectors focus on a sparse set of region proposals oftentimes producing a separate region proposal network (RPN). Two-stage detectors are usually more accurate with reduced efficiency. While these detectors operate on single images, several recent works seek to extend them to the task of video object detection [17, 34, 72, 73]. A key idea behind these methods is to explore the temporal continuity of videos, where motion tracking is used to enhance the video object detection performance. We used this similar idea for our system, ApproxDet.

**Efficient DNNs for Mobile Applications** There is an emerging interest to develop efficient DNNs for mobile vision applications. An important line of research is to identify lightweight network architectures with low computational cost. Examples include manually designed MobileNet family [24, 54], SqueezeNet [30], and ShuffleNet [70], as well as the more recent MNasNet [57], MobileNetV3 [23], FBNet [64], and EfficientNet [58] produced by neural architecture search [14]. Other techniques include the removal of redundant parameters (network pruning) [20, 26, 36, 41], the quantization of parameters (network quantization) [28, 29, 49]. While these architectures and techniques are primarily designed for image classification, they can be used as a building block for efficient object detectors [59, 63]. *In spite of the efficiency, none of these approaches can adapt to contention and content during runtime.*

**Adaptive Inference for DNNs** The early work on anytime prediction [75] presents the first set of methods that are adaptive to a computing budget at runtime. This idea was recently revisited in

the computer vision community using DNNs. For example, at inference time, a DNN can choose to drop certain operations [62, 65], or to select one of the many exits [27, 68] or branches [61], based on the image content or a given latency budget. Unfortunately, none of these approaches is designed for object detection. Besides, they do not consider the modeling of resource contention. In parallel to these developments, several recent work in the systems community also seek to build adaptive inference systems for DNNs. For example, NestDNN [15] uses network pruning to convert a static DNN into multiple DNNs, and dynamically selects from these to fit the resource requirement for image classification. AdaScale [6] learns to adaptively change the input shape of an object detection DNN, in order to achieve a latency-accuracy tradeoff. In comparison to NestDNN, our system ApproxDet uses a single adaptive DNN model for object detection. In contrast to AdaScale, ApproxDet considers a joint adaptation of video content and resource contention.

**Approximate Video Analytics.** Our work shares similar ideas to several recent works in video analytics in the systems community. For example, server-side solutions like VideoStorm [69], Chameleon [32], and Focus [25] exploit various configurations and DNN models to optimize video analytics queries, but they also require loading of multiple models at the same time, which are challenging in resource-constrained mobile devices. If memory constraints prevent multiple models being co-resident on a device, it is conceivable to send them over the network on demand, using efficient wireless reprogramming [48], but the times involved are such that the prediction of which model will be required will have to be done far in advance. Liu *et al.* [38] explore the offloading of object detection to an edge device in combination with fast on-device tracking for mobile AR. ExCamera [18] enables low-latency video processing on the cloud using serverless architecture. VideoChef [66] uses approximation knobs of traditional video preprocessing filters in a content-aware manner. It cannot handle object detection and is not applicable to DNN-based video processing. A recent work in this space called MARLIN [2] shows that for AR applications, instead of continuously running the detection DNN, they can decide when to run their specially designed lightweight trackers. We take the idea of running the tracker at a configurable interval  $si$ , but we use traditional trackers.

## 9 CONCLUSION

In this paper, we present ApproxDet, a single model adaptive system for video object detection in a video content-aware and resource contention-aware fashion, focusing on resource-constrained mobile/embedded devices. The joint context and content aware mechanisms make ApproxDet can adaptively respond to both content and context changes while most baseline models cannot respond to any changes or only one change. Case studies have shown our ApproxDet can adjust to different scenarios with best accuracy and least latency over previous models. Further, We contrast ApproxDet with multiple baselines, including AdaScale, a content-aware adaptive server-based video object detection system, and YOLOv3, a single-stage objection system known for its efficiency on the ImageNet VID dataset. We find that ApproxDet is 52.9% lower in latency and 11.1% higher in the accuracy metric over YOLOv3 and outperforms all baselines with significantly lower switching overhead,

stemming from its single model design. Results on ImageNet VID dataset proves the efficiency and effectiveness of our ApproxDet, highlighting ApproxDet's promise in enabling latency-sensitive applications, pushing the frontiers of ever-evolving AR/MR (mixed reality) experiences instantiated on embedded platforms.

## ACKNOWLEDGMENTS

We thank the anonymous reviewers and shepherd for their valuable comments to improve the quality of this paper. This material is based in part upon work supported by the National Science Foundation under Grant Number CNS-1527262, Army Research Lab under Contract number W911NF-20-2-0026, the Lilly Endowment (Wabash Heartland Innovation Network, WHIN-Purdue), and gift funding from Adobe Research. Yin Li acknowledges the support by the UW VCRGE with funding from WARF. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the sponsors.

## REFERENCES

- [1] Jason Ansel, Yee Lok Wong, Cy Chan, Marek Olszewski, Alan Edelman, and Saman Amarasinghe. 2011. Language and compiler support for auto-tuning variable-accuracy algorithms. In *International Symposium on Code Generation and Optimization (CGO 2011)*. IEEE, 85–96.
- [2] Kittipat Apicharttrisorn, Xukan Ran, Jiasi Chen, Srikanth V Krishnamurthy, and Amit K Roy-Chowdhury. 2019. Frugal following: Power thrifty object detection and tracking for mobile augmented reality. In *Proceedings of the Conference on Embedded Networked Sensor Systems (SenSys)*. 96–109.
- [3] Ross Bulat. 2020. React Native: Background Task Management in iOS. <https://medium.com/@rossbulat/react-native-background-task-management-in-ios-d0f05ae53cc5>
- [4] Shuai Che, Michael Boyer, Jiayuan Meng, David Tarjan, Jeremy W Sheaffer, Sang-Ha Lee, and Kevin Skadron. 2009. Rodinia: A benchmark suite for heterogeneous computing. In *2009 IEEE international symposium on workload characterization (IISWC)*. Ieee, 44–54.
- [5] Tiffany Yu-Han Chen, Lenin Ravindranath, Shuo Deng, Paramvir Bahl, and Hari Balakrishnan. 2015. Glimpse: Continuous, real-time object recognition on mobile devices. In *Proceedings of the ACM Conference on Embedded Networked Sensor Systems (SenSys)*. 155–168.
- [6] Ting-Wu Chin, Ruizhou Ding, and Diana Marculescu. 2019. AdaScale: Towards real-time video object detection using adaptive scaling. In *Proceedings of the Conference on Machine Learning and Systems (SysML)*.
- [7] The Pokemon Company. 2020. Pokémon GO | Augmented Reality Mobile Game. <https://pokemongolive.com/en/>
- [8] NVIDIA Corporation. 2018. *Jetson TX2 Module*. Retrieved May 5, 2020 from <https://developer.nvidia.com/embedded/buy/jetson-tx2>
- [9] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. 2016. R-FCN: Object detection via region-based fully convolutional networks. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*. 379–387.
- [10] Christina Delimitrou and Christos Kozyrakis. 2013. ibench: Quantifying interference for datacenter applications. In *2013 IEEE international symposium on workload characterization (IISWC)*. IEEE, 23–33.
- [11] Android Developer. 2019. Guide to background processing: Android. <https://developer.android.com/guide/background>
- [12] Apple Developer. 2019. Services provided by an app that require it to run in the background. [https://developer.apple.com/documentation/bundleresources/information\\_property\\_list/uibackgroundmodes](https://developer.apple.com/documentation/bundleresources/information_property_list/uibackgroundmodes)
- [13] Yufei Ding, Jason Ansel, Kalyan Veeramachaneni, Xipeng Shen, Una-May O'Reilly, and Saman Amarasinghe. 2015. Autotuning algorithmic choice for input sensitivity. *ACM SIGPLAN Notices* 50, 6 (2015), 379–390.
- [14] Thomas Elsken, Jan Hendrik Metzen, and Frank Hutter. 2019. Neural Architecture Search: A Survey. *Journal of Machine Learning Research* 20, 55 (2019), 1–21.
- [15] Biyi Fang, Xiao Zeng, and Mi Zhang. 2018. NestDNN: Resource-aware multi-tenant on-device deep learning for continuous mobile vision. In *Proceedings of the Annual International Conference on Mobile Computing and Networking (MobiCom)*. 115–127.
- [16] Gunnar Farneback. 2003. Two-frame motion estimation based on polynomial expansion. In *Proceedings of Scandinavian Conference on Image Analysis*. 363–370.
- [17] Christoph Feichtenhofer, Axel Pinz, and Andrew Zisserman. 2017. Detect to track and track to detect. In *Proceedings of the IEEE International Conference on*

- Computer Vision (ICCV)*. 3038–3046.
- [18] Sadjad Fouladi, Riad S Wahby, Brennan Shacklett, Karthikeyan Balasubramaniam, William Zeng, Rahul Bhalerao, Anirudh Sivaraman, George Porter, and Keith Winstein. 2017. Encoding, Fast and Slow: Low-Latency Video Processing Using Thousands of Tiny Threads.. In *Proceedings of the Symposium on Networked Systems Design and Implementation (NSDI)*. 363–376.
  - [19] Sudipto Guha, Nina Mishra, Gourav Roy, and Okke Schrijvers. 2016. Robust random cut forest based anomaly detection on streams. In *Proceedings of the International Conference on Machine Learning (ICML)*. 2712–2721.
  - [20] Song Han, Huizi Mao, and William J Dally. 2016. Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding. (2016), 1–13.
  - [21] Seungyeop Han, Haichen Shen, Matthai Philipose, Sharad Agarwal, Alec Wolman, and Arvind Krishnamurthy. 2016. MCDNN: An approximation-based execution framework for deep stream processing under resource constraints. In *Proceedings of the Annual International Conference on Mobile Systems, Applications, and Services (MobiSys)*. 123–136.
  - [22] J. F. Henriques, R. Caseiro, P. Martins, and J. Batista. 2014. High-Speed Tracking with Kernelized Correlation Filters. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 37, 3 (2014), 583–596.
  - [23] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. 2019. Searching for MobileNetV3. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. 1314–1324.
  - [24] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Heyand, Marco Andreetto, and Hartwig Adam. 2017. MobileNets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861* (2017).
  - [25] Kevin Hsieh, Ganesh Ananthanarayanan, Peter Bodik, Shivaram Venkataraman, Paramvir Bahl, Matthai Philipose, Phillip B Gibbons, and Onur Mutlu. 2018. Focus: Querying large video datasets with low latency and low cost. In *13th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 18)*. 269–286.
  - [26] Hengyuan Hu, Rui Peng, Yu-Wing Tai, and Chi-Keung Tang. 2016. Network trimming: A data-driven neuron pruning approach towards efficient deep architectures. *arXiv preprint arXiv:1607.03250* (2016).
  - [27] Gao Huang, Danlu Chen, Tianhong Li, Felix Wu, Laurens van der Maaten, and Kilian Q Weinberger. 2018. Multi-scale dense networks for resource efficient image classification. In *Proceedings of International Conference on Learning Representations (ICLR)*.
  - [28] Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. 2016. Binarized neural networks. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*. 4107–4115.
  - [29] Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. 2017. Quantized Neural Networks: Training Neural Networks with Low Precision Weights and Activations. *Journal of Machine Learning Research* 18 (2017), 187–1.
  - [30] Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. 2016. SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and- 0.5 MB model size. In *Proceedings of International Conference on Learning Representations (ICLR)*. 1–13.
  - [31] Angela H Jiang, Daniel L-K Wong, Christopher Canel, Lilia Tang, Ishan Misra, Michael Kaminsky, Michael A Kozuch, Padmanabhan Pillai, David G Andersen, and Gregory R Ganger. 2018. Mainstream: Dynamic Stem-Sharing for Multi-Tenant Video Processing. In *Proceedings of USENIX Annual Technical Conference (USENIX ATC)*. 29–42.
  - [32] Junchen Jiang, Ganesh Ananthanarayanan, Peter Bodik, Siddhartha Sen, and Ion Stoica. 2018. Chameleon: scalable adaptation of video analytics. In *Proceedings of the Conference of the ACM Special Interest Group on Data Communication (SIGCOMM)*. 253–266.
  - [33] Zdenek Kalal, Krystian Mikolajczyk, and Jiri Matas. 2010. Forward-Backward error: Automatic detection of tracking failures. In *Proceedings of IEEE International Conference on Pattern Recognition (CVPR)*. 2756–2759.
  - [34] Kai Kang, Hongsheng Li, Junjie Yan, Xingyu Zeng, Bin Yang, Tong Xiao, Cong Zhang, Zhe Wang, Ruohui Wang, Xiaogang Wang, et al. 2017. T-CNN: Tubelets with convolutional neural networks for object detection from videos. *IEEE Transactions on Circuits and Systems for Video Technology* 28, 10 (2017), 2896–2907.
  - [35] Michael A Laurenzano, Parker Hill, Mehrzad Samadi, Scott Mahlke, Jason Mars, and Lingjia Tang. 2016. Input responsiveness: using canary inputs to dynamically steer approximation. In *Proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*. 161–176.
  - [36] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. 2017. Pruning filters for efficient ConvNets. In *Proceedings of International Conference on Learning Representations (ICLR)*. 1–13.
  - [37] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. 2017. Focal loss for dense object detection. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. 2980–2988.
  - [38] Luyang Liu, Hongyu Li, and Marco Gruteser. 2019. Edge assisted real-time object detection for mobile augmented reality. (2019), 1–16.
  - [39] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. 2016. SSD: Single shot multibox detector. In *Proceedings of the European conference on Computer Vision (ECCV)*, Vol. 9907. 21–37.
  - [40] Alan Lukežič, Tom'aš Voj'ir, Luka Čehovin Zajc, Jir'i Matas, and Matej Kristan. 2018. Discriminative Correlation Filter Tracker with Channel and Spatial Reliability. *International Journal of Computer Vision* 126 (2018), 671–688.
  - [41] Jian-Hao Luo, Hao Zhang, Hong-Yu Zhou, Chen-Wei Xie, Jianxin Wu, and Weiyao Lin. 2018. ThiNet: pruning CNN filters for a thinner net. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 41, 10 (2018), 2525–2538.
  - [42] Amiya K Maji, Subrata Mitra, Bowen Zhou, Saurabh Bagchi, and Akshat Verma. 2014. Mitigating interference in cloud services by middleware reconfiguration. In *Proceedings of the 15th International Middleware Conference*. 277–288.
  - [43] Jason Mars, Lingjia Tang, Robert Hundt, Kevin Skadron, and Mary Lou Soffa. 2011. Bubble-up: Increasing utilization in modern warehouse scale computers via sensible co-locations. In *Proceedings of the 44th annual IEEE/ACM International Symposium on Microarchitecture*. 248–259.
  - [44] John D. McCalpin. 1991-2007. *STREAM: Sustainable Memory Bandwidth in High Performance Computers*. Technical Report. University of Virginia, Charlottesville, Virginia. <http://www.cs.virginia.edu/stream/>
  - [45] John D. McCalpin. 1995. Memory Bandwidth and Machine Balance in Current High Performance Computers. *IEEE Computer Society Technical Committee on Computer Architecture (TCCA) Newsletter* (Dec. 1995), 19–25.
  - [46] Chulhong Min, Alessandro Montanari, Akhil Mathur, and Fahim Kawsar. 2019. A closer look at quality-aware runtime assessment of sensing models in multi-device environments. In *Proceedings of the 17th Conference on Embedded Networked Sensor Systems*. 271–284.
  - [47] Subrata Mitra, Manish K Gupta, Sasa Misailovic, and Saurabh Bagchi. 2017. Phase-aware optimization in approximate computing. In *2017 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*. IEEE, 185–196.
  - [48] Rajesh Krishna Panta, Saurabh Bagchi, and Samuel P Midkiff. 2011. Efficient incremental code update for sensor networks. *ACM Transactions on Sensor Networks (TOSN)* 7, 4 (2011), 1–32.
  - [49] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. 2016. XNOR-Net: ImageNet classification using binary convolutional neural networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, Vol. 9908. 525–542.
  - [50] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. 2016. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 779–788.
  - [51] Joseph Redmon and Ali Farhadi. 2018. YOLOv3: An incremental improvement. *arXiv preprint arXiv:1804.02767* (2018).
  - [52] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster R-CNN: Towards real-time object detection with region proposal networks. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*. 91–99.
  - [53] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. 2015. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision* 115, 3 (2015), 211–252.
  - [54] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. 2018. MobileNetV2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 4510–4520.
  - [55] Xiaoyong Shen, Aaron Hertzmann, Jiaya Jia, Sylvain Paris, Brian Price, Eli Shechtman, and Ian Sachs. 2016. Automatic portrait segmentation for image stylization. In *Computer Graphics Forum*, Vol. 35. Wiley Online Library, 93–102.
  - [56] Hui Shuai, Qingshan Liu, Kaihua Zhang, Jing Yang, and Jiankang Deng. 2018. Cascaded Regional Spatio-Temporal Feature-Routing Networks for Video Object Detection. *IEEE Access* 6 (2018), 3096–3106.
  - [57] Mingxing Tan, Bo Chen, Ruoming Pang, Vijay Vasudevan, Mark Sandler, Andrew Howard, and Quoc V Le. 2019. MNasNet: Platform-aware neural architecture search for mobile. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2820–2828.
  - [58] Mingxing Tan and Quoc V Le. 2019. EfficientNet: Rethinking model scaling for convolutional neural networks. In *Proceedings of the International Conference on Machine Learning (ICML)*.
  - [59] Mingxing Tan, Ruoming Pang, and Quoc V Le. 2020. EfficientDet: Scalable and efficient object detection. (2020), 10781–10790.
  - [60] Matthew Tancreti, Mohammad Sajjad Hossain, Saurabh Bagchi, and Vijay Raghunathan. 2011. Aveksha: A hardware-software approach for non-intrusive tracing and profiling of wireless embedded systems. In *Proceedings of the 9th ACM Conference on Embedded Networked Sensor Systems*. 288–301.
  - [61] Surat Teerapittayanon, Bradley McDanel, and HT Kung. 2016. BranchyNet: Fast inference via early exiting from deep neural networks. In *Proceedings of IEEE International Conference on Pattern Recognition (ICPR)*. 2464–2469.

- [62] Andreas Veit and Serge Belongie. 2019. Convolutional networks with adaptive inference graphs. *International Journal of Computer Vision* 128 (2019), 730–741.
- [63] Robert J Wang, Xiang Li, and Charles X Ling. 2018. PELEE: A real-time object detection system on mobile devices. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*. 1963–1972.
- [64] Bichen Wu, Xiaoliang Dai, Peizhao Zhang, Yanghan Wang, Fei Sun, Yiming Wu, Yuandong Tian, Peter Vajda, Yangqing Jia, and Kurt Keutzer. 2019. FBNet: Hardware-aware efficient convnet design via differentiable neural architecture search. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 10734–10742.
- [65] Zuxuan Wu, Tushar Nagarajan, Abhishek Kumar, Steven Rennie, Larry S Davis, Kristen Grauman, and Rogerio Feris. 2018. BlockDrop: Dynamic inference paths in residual networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 8817–8826.
- [66] Ran Xu, Jinkyu Koo, Rakesh Kumar, Peter Bai, Subrata Mitra, Sasa Misailovic, and Saurabh Bagchi. 2018. VideoChef: efficient approximation for streaming video processing pipelines. In *Proceedings of the USENIX Annual Technical Conference (USENIX ATC)*. 43–56.
- [67] Ran Xu, Subrata Mitra, Jason Rahman, Peter Bai, Bowen Zhou, Greg Bronevetsky, and Saurabh Bagchi. 2018. Pythia: Improving datacenter utilization via precise contention prediction for multiple co-located workloads. In *Proceedings of the International Middleware Conference (Middleware)*. 146–160.
- [68] Le Yang, Yizeng Han, Xi Chen, Shiji Song, Jifeng Dai, and Gao Huang. 2020. Resolution Adaptive Networks for Efficient Inference. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2369–2378.
- [69] Haoyu Zhang, Ganesh Ananthanarayanan, Peter Bodik, Matthai Philipose, Paramvir Bahl, and Michael J Freedman. 2017. Live Video Analytics at Scale with Approximation and Delay-Tolerance. In *Proceedings of the Symposium on Networked Systems Design and Implementation (NSDI)*, Vol. 9. 377–392.
- [70] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. 2018. ShuffleNet: An extremely efficient convolutional neural network for mobile devices. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 6848–6856.
- [71] Zhong-Qiu Zhao, Peng Zheng, Shou-tao Xu, and Xindong Wu. 2019. Object detection with deep learning: A review. *IEEE Transactions on Neural Networks and Learning Systems* 30, 11 (2019), 3212–3232.
- [72] Xizhou Zhu, Jifeng Dai, Lu Yuan, and Yichen Wei. 2018. Towards high performance video object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 7210–7218.
- [73] Xizhou Zhu, Yujie Wang, Jifeng Dai, Lu Yuan, and Yichen Wei. 2017. Flow-guided feature aggregation for video object detection. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. 408–417.
- [74] Xizhou Zhu, Yuwen Xiong, Jifeng Dai, Lu Yuan, and Yichen Wei. 2017. Deep feature flow for video recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2349–2358.
- [75] Shlomo Zilberstein. 1996. Using anytime algorithms in intelligent systems. *AI magazine* 17, 3 (1996), 73–73.