

Focus: Querying Large Video Datasets with Low Latency and Low Cost

Kevin Hsieh^{†§} Ganesh Ananthanarayanan[§] Peter Bodik[§] Paramvir Bahl[§] Matthai Philipose[§]
 Phillip B. Gibbons[†] Onur Mutlu^{*†}
[†]Carnegie Mellon University [§]Microsoft ^{*}ETH Zürich

Abstract

Large volumes of videos are continuously recorded from cameras deployed for traffic control and surveillance with the goal of answering “after the fact” queries: *identify video frames with objects of certain classes (cars, bags)* from many days of recorded video. While advancements in convolutional neural networks (CNNs) have enabled answering such queries with high accuracy, they are too expensive and slow. We build Focus, a system for low-latency and low-cost querying on large video datasets. Focus uses cheap ingestion techniques to index the videos by the objects occurring in them. At ingest-time, it uses compression and video-specific specialization of CNNs. Focus handles the lower accuracy of the cheap CNNs by judiciously leveraging expensive CNNs at query-time. To reduce query time latency, we cluster similar objects and hence avoid redundant processing. Using experiments on video streams from traffic, surveillance and news channels, we see that Focus uses 58× fewer GPU cycles than running expensive ingest processors and is 37× faster than processing all the video at query time.

1. Introduction

Cameras are ubiquitous, with millions of them deployed by government and private entities at traffic intersections, enterprise offices, and retail stores. Videos from these cameras are continuously recorded [2, 7]. One of the main purposes for recording the videos is answering “after-the-fact” queries: *identify video frames with objects of certain classes (like cars or bags)* over many days of recorded video. As results from these queries are used by analysts and investigators, achieving low query latencies is crucial.

Advances in convolutional neural networks (CNNs) backed by copious training data and hardware accelerators (e.g., GPUs [13]) have led to high accuracy in the computer vision tasks like object detection and object classification. For instance, the ResNet152 object classifier CNN [38] won the ImageNet challenge that evaluates classification accuracy on 1,000 classes using a public image dataset with labeled ground truths [63]. For each image, these classifiers return a ranked list of 1,000 classes in decreasing order of confidence.

Despite the accuracy of image classifier CNNs (like

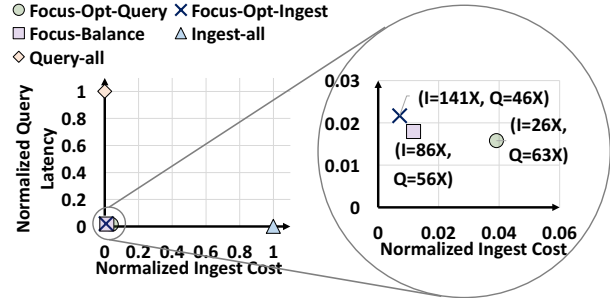


Figure 1: Effectiveness of Focus at reducing both ingest cost and query latency, for an example traffic video. We compare against two baselines: “Ingest-all” that runs ResNet152 on all video frames during ingest, and “Query-all” that runs ResNet152 on all the video frames at query time. By zooming in, we see that Focus (the Focus-Balance point) is simultaneously 86× cheaper than Ingest-all in its GPU consumption and 56× faster than Query-all in query latency, all the while achieving at least 95% precision and recall. (Also shown are two alternatives offering slightly different trade-offs.)

ResNet152), using them for video analytics queries is both expensive and slow. Using the ResNet152 classifier at *query-time* to identify video frames with cars on a month-long traffic video requires 280 GPU hours and costs \$250 in the Azure cloud. The latency for running queries is also high. To achieve a query latency of one minute on 280 GPU hours of work would require tens of thousands of GPUs classifying the frames of the video in parallel, which is many orders of magnitude more than what is typically provisioned (few tens or hundreds) by traffic jurisdictions or retail stores. Note that the above cost and latency values are *after* using motion detection techniques to exclude frames with no moving objects.

We believe that enabling *low-latency and low-cost querying over large video datasets* will make video analytics more useful and open up many new opportunities.

A natural approach to enabling low latency querying is doing all classifications with ResNet152 at *ingest-time*, i.e., on the *live* videos, and store the results in an index of object classes to video frames. Any queries for specific classes (e.g., cars) will thus involve only a simple index lookup at *query-time*. There are, however, at least two problems with this approach. First, the cost to index all the video at ingest-time, e.g., \$250/month/stream in the

above example, is prohibitively high. Second, most of this ingest-time cost is wasteful because typically only a small fraction of recorded videos get queried [16]. Following a theft, the police would query a few days of video from a handful of surveillance cameras, but not all the videos.

We present Focus, a system to support low-latency low-cost querying on large video datasets. To address the above drawbacks, Focus has the following goals: (1) low cost indexing of video at ingest-time, (2) providing high accuracy and low latency for queries, and (3) allowing trade offs between the cost at ingest-time against the latency at query-time. As input, the user specifies the *ground-truth CNN* (or “GT-CNN”, e.g., the ResNet152 classifier) and the desired accuracy of results that Focus needs to achieve relative to the GT-CNN.

Focus uses four key techniques – cheap CNNs for ingest, using top-K results from the ingest-time CNN, clustering similar objects, and judicious selection of system and model parameters.

First, to make video ingestion cheap, Focus uses *compressed* and *specialized* versions of CNNs, to create an ingest-time index of object classes to frames. CNN compression (e.g., [66]) creates new CNNs with fewer convolutional layers and smaller input images. Specialization [35, 65] trains those CNNs on a smaller set of object classes specific to each video stream so that those cheaper CNNs can classify these video-specific objects more accurately. Together, these techniques result in highly efficient CNNs for video indexing.

Second, the cheap ingest CNNs, however, are also less accurate than the expensive GT-CNN (like ResNet152), measured in terms of *recall* and *precision*. Recall is the fraction of frames in the video that contained objects of the queried class that were *actually* returned in the query’s results. Precision, on the other hand, is the fraction of frames in the query’s results that contained objects of the queried class. To increase recall, Focus relies on an empirical observation: while the top-most (i.e., most confident) classification results of the cheap and expensive CNNs may not always match, the top-most result of the expensive CNN falls within the *top-K* results of the cheap CNN. Therefore, at ingest-time, Focus indexes each object with the “top-K” results of the cheap CNN (instead of just the top-most). To increase precision, at query-time, we first filter the objects from the top-K index and then classify the filtered objects with the expensive GT-CNN.

Third, to reduce the query-time latency of using the expensive GT-CNN, Focus relies on the significant similarity between objects in videos. For example, a car moving across an intersection will look very similar in consecutive frames. Focus leverages this similarity by clustering the objects at ingest-time, classifying *only* the cluster centroids with the expensive GT-CNN at query-time, and assigning the same class to all objects in the cluster, thus

considerably reducing query latency.

In a nutshell, Focus’s ingest-time and query-time operations are as follows. At ingest-time, it classifies the detected objects using a cheap CNN, clusters similar objects, and indexes each cluster centroid using the top-K classification results. At query-time, when the user queries for class X, Focus looks up the ingest index for centroids that match class X and classifies them using the GT-CNN. For centroids that were classified as class X, it returns all objects from the corresponding clusters to the user.

Finally, Focus smartly chooses the ingest-time CNN and its parameters to meet user-specified targets on precision and recall. Among the choices that meet the accuracy targets, it allows the user to trade off between the ingest cost and query latency. For example, using a cheaper ingest CNN reduces the ingest cost but increases the query latency as Focus needs to use a larger K for the top-K index to retain the accuracy targets. Focus identifies the “sweet spot” in parameters that sharply improve one of ingest cost or query latency for a small worsening of the other.

We built Focus and evaluated it on thirteen 12-hour videos from three domains – traffic cameras, surveillance cameras, and news channels. We compare against two baselines: “Ingest-all” that runs GT-CNN on all video frames during ingest, and “Query-all” that runs GT-CNN on all the video frames at query time. We use ResNet152 as GT-CNN and augment both baselines with motion detection to remove frames with no objects, which is one of the core techniques in a recent prior work, NoScope [44]. Figure 1 shows a representative result, for a traffic video from a commercial intersection. On average, Focus is $58\times$ (up to $98\times$) cheaper than Ingest-all and $37\times$ (up to $57\times$) faster than Query-all. This leads to the cost of ingestion coming down from \$250/month/stream to \$4/month/stream, and the latency to query a 24 hour video dropping from 1 hour to under 2 minutes. See §6 for the full details.

We make the following contributions.

1. We formulate the problem of querying video datasets by showing the trade-offs between query latency, ingest cost, and accuracy (precision and recall) of results.
2. We propose techniques to ingest videos with low cost by leveraging compressed and video-specific specialization of CNNs, while retaining high accuracy targets by creating approximate (top-K) indexes.
3. We identify and leverage similarity between objects in a video to cluster them using CNN features and significantly speeding up queries.
4. We propose and build a new end-to-end system to support low-latency, low-cost querying on large video datasets. We show that our system offers new trade-off options between ingestion cost and query latency, as it is significantly cheaper than analyzing all videos frames at ingest time and significantly faster than analyzing queried video frames at query time.

2. Background and Motivation

We first provide a brief overview of convolutional Neural Networks (CNN), the state-of-the-art approach to detecting and classifying objects in images (§2.1). We then discuss new observations we made about real-world videos, which motivate the design of our techniques (§2.2).

2.1. Convolutional Neural Networks

A Convolution Neural Network (CNN) [47] is a specific class of neural networks that works by extracting the visual features in images. During image classification, or “inference”, a CNN takes an input image and outputs the probability of each *class* (e.g., dog, flower, or car). CNNs are the state-of-the-art method used for many computer vision tasks, such as image classification (e.g., [38,45,71]) and face recognition (e.g., [46,64]).

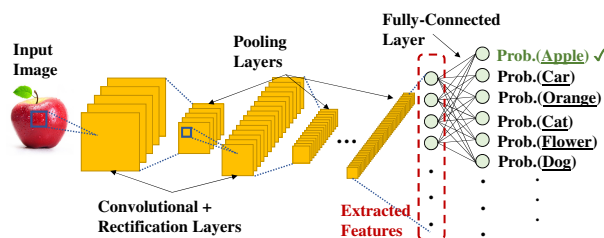


Figure 2: Architecture of an image classification CNN.

Figure 2 illustrates the architecture of an image classification CNN. Broadly, almost all CNNs consist of three key types of network layers: (1) *convolutional and rectification layers*, which detect visual features from input pixels, (2) *pooling layers*, which down-sample the input by merging neighboring pixel values, and (3) *fully-connected layers*, which provide the reasoning to classify the input object based on the outputs from previous layers. The outputs of an image classification CNN are the the probabilities of all object classes, and the class with the highest probability is the predicted class for the input image.

The output of the penultimate (i.e., previous-to-last) layer can be considered as “representative features” of the input image [45]. The features are a real-valued vector, with lengths between 512 and 4096 in state-of-the-art classifier CNNs (e.g., [38,45,66,71]). It has been shown that images with similar feature vectors (i.e., small Euclidean distances) are visually similar [22,23,45,58].

The high accuracy of CNNs comes at a cost: *inferring* (or classifying) using state-of-the-art CNNs to classify objects in images requires significant computational resources. This is because the higher accuracy of CNNs comes from using *deeper* architectures (i.e., more layers) to obtain better visual features. For instance, ResNet152 [38], the winner of the ImageNet competition [63] in 2015, has been trained to classify across 1000 classes from the ImageNet dataset using 152 layers, but

can only process 77 images/second even with a high-end GPU (NVIDIA K80 [13]). This makes querying on large video datasets using these CNNs to be *slow* and *costly*.

There are at least two recent techniques designed to reduce the cost of CNNs. First, *compression* is a set of techniques aiming to reduce the cost of CNN inference (classification) at the expense of reduced accuracy. Such techniques include removing some expensive convolutional layers [66], matrix pruning [31,37], and others [42,62] and can dramatically reduce the classification cost of a CNN. For example, ResNet18, which is a ResNet152 variant with only 18 layers is $8\times$ cheaper. Second, a more recent technique is CNN *specialization* [35], where the CNNs are trained on a subset of a dataset specific to a particular context, also making them much cheaper. Using the combination of cheap and expensive CNNs is a key facet of our solution, described in §4.

2.2. Characterizing Real-world Videos

We aim to support queries of the form, *find all frames in the video that contain objects of class X*. We identify some key characteristics of real-world videos towards supporting these queries: (1) large portions of videos can be excluded (§2.2.1), (2) only a limited set of object classes occur in each video (§2.2.2), and (3) objects of the same class have similar feature vectors (§2.2.3). The design of Focus is based on these characteristics.

We have analyzed 12 hours of video from six video streams each. The six video stream span across traffic cameras, surveillance cameras, and news channels. (§6.1 provides the details.) We detect the objects in each frame of these videos (using background subtraction [43]), and classify each object with the ResNet152 CNN [38] among the supported 1,000 object classes. In this paper, we use results from the costly ResNet152 CNN as ground truth.

2.2.1. Excluding large portions of videos. We find considerable potential to avoid processing large portions of videos at query-time. Significant portions of video streams either have *no objects* at all (as in a garage camera at night) or the objects are *stationary* (like parked cars). We find that in our video sets, one-third to one-half of the frames fall in these categories. Therefore, queries to *any* object class would benefit from pre-processing filters applied to exclude these portions of the videos.

Even among the frames that do contain objects, not all of them are relevant to a query because each query only looks for a *specific class* of objects. In our video sets, an object class occurs in only 0.01% of the frames on average, and even the most frequent object classes occur in no more than 16% – 43% of the frames in the different videos. This is because while there are usually some dominant classes (e.g., cars in a traffic camera, people in a news channel), most other classes are rare. Since queries are for specific object classes, there is considerable poten-

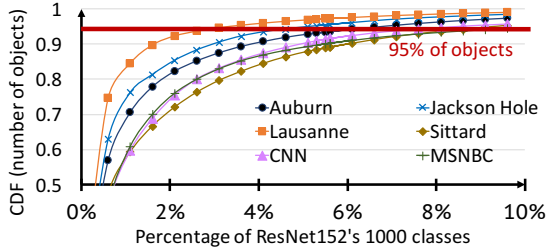


Figure 3: CDF of frequency of object classes. The x-axis is the fraction of classes out of the 1000 recognized by ResNet152 (truncated to 10%).

tial in *indexing* frames by the classes of objects.

2.2.2. Limited set of object classes in each video. We next focus on the classes of objects that occur in each of the videos and the disparity in frequency among them.

Most video streams have a limited set of objects because each video has its own context (e.g., traffic cameras can have automobiles, pedestrians or bikes but not airplanes). It is rare that a video stream contains objects of *all* the classes recognized by state-of-the-art classifier CNNs. Figure 3 shows the cumulative distribution function (CDF) of the frequency of object classes in our videos (as classified by ResNet152). We make two observations.

First, objects of only 22% – 33% (not graphed) of the 1,000 object classes occur in the less busy videos (Auburn, Jackson Hole, Lausanne, and Sittard). Even in the busier videos (CNN, and MSNBC), objects of only 50% – 69% of the classes appear. Also, there is little overlap between the classes of objects among the different videos. On average, the Jaccard indexes [72] (i.e., intersection over union) between the videos based on their object classes is only 0.46. Second, even among the object classes that do occur, a small fraction of classes disproportionately dominate. Figure 3 shows that 3% – 10% of the most frequent object classes cover $\geq 95\%$ of the objects in each video stream. This suggests that for each video stream, we can *automatically* (i) determine its most frequently occurring classes and (ii) train efficient CNNs *specialized* for classifying these classes (§2.1).

2.2.3. Feature vectors for finding duplicate objects.

Objects moving in the video often stay in the frame for several seconds; for example, a pedestrian might take a minute to cross a street. Instead of classifying *each instance* of the same object across the frames, we would like to *inexpensively* find duplicate objects and only classify one of them using a CNN (and apply the same label to all duplicates). Thus, given n duplicate objects, this requires only one CNN classification operation instead of n .

Comparing pixel values across frames is an obvious choice to identify duplicate objects, however, they turn out to be highly sensitive to even small changes in the camera’s real-time view of an object. Instead, feature vectors extracted from the CNNs are much more robust

since they are specifically trained to extract visual features for classification. We verify the robustness of feature vectors using the following analysis. In each video, for each object i , we find its *nearest* neighbor j using feature vectors from the cheap ResNet18 CNN and compute the fraction of object pairs that belong to the same class. This fraction is over 99% in each of our videos, which shows using feature vectors from cheap CNNs can potentially help identify duplicate objects.

3. Overview of Focus

The goal of Focus is to *index live video streams* by the object classes occurring in them and enable answering “after-the-fact” queries later on the stored videos of the form *find all frames that contain objects of class X*. Optionally, the query can be restricted to a subset of cameras and a time range. Such a query formulation is the basis for many widespread applications and could be used either on its own (such as for detecting all cars or bicycles in the video) or used as a basis for further processing (e.g., finding all collisions between cars and bicycles).

Focus is designed to work with a wide variety of current and future CNNs. At system configuration time, the user (system administrator) provides a *ground-truth CNN* (GT-CNN), which serves as the accuracy baseline for Focus, but is far too costly to run on every video frame. Through a sequence of techniques, Focus provides nearly-comparable accuracy but at greatly reduced cost. By default, and throughout this paper, we use the ResNet152 image classifier as the GT-CNN.

Because the acceptable target accuracy is application-dependent, Focus permits the user to specify the target, while providing reasonable defaults. Accuracy is specified in terms of *precision*, i.e., fraction of frames output by the query that actually contain an object of class X according to GT-CNN, and *recall*, i.e., fraction of frames that contain objects of class X according to GT-CNN that were actually returned by the query. The lower the target, the greater the cost-savings provided by Focus. Even for high targets such as 95%–99%, Focus is able to achieve order-of-magnitude or more cost savings.

Figure 4 presents the design of Focus.

- At *ingest-time* (left part of Figure 4), Focus classifies objects in the incoming video frames and extracts their feature vectors. To make this step cheap, it uses a highly compressed and specialized version of the GT-CNN model (IT_1 in Figure 4). Focus then clusters objects based on their feature vectors (IT_2) and assign to each cluster the *top K* most likely classes these objects belong to (based on classification confidence of the ingest CNN); (IT_3). It creates a *top-K index*, which maps each class to the set of object clusters (IT_4). The top-K index is the output of Focus’ ingest-time processing of videos.
- At *query-time* (right part of Figure 4), when the user

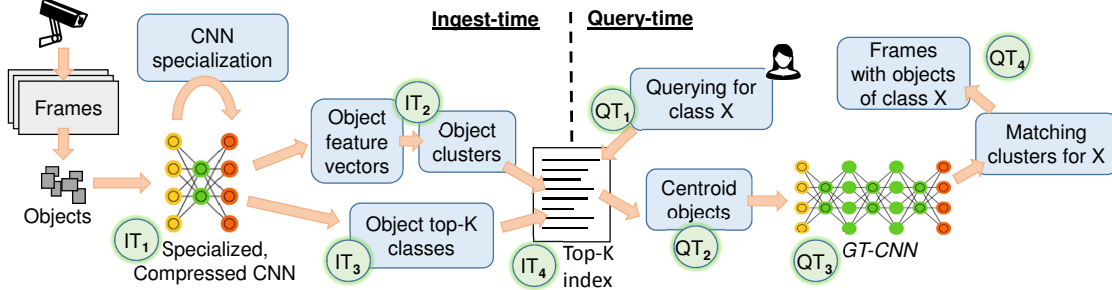


Figure 4: Overview of Focus.

queries for a certain class X (QT_1), Focus retrieves the matching clusters from the top-K index (QT_2), runs the *centroids* of the clusters through GT-CNN (QT_3), and returns all frames from the clusters whose centroids were classified by GT-CNN as class X (QT_4).

The top-K ingest index is a mapping between the object class to the clusters. Specifically,

object class \rightarrow \langle cluster ID \rangle
cluster ID \rightarrow [centroid object, \langle objects \rangle
in cluster, \langle frame IDs \rangle of objects]

We next explain how Focus’ key techniques keep ingest cost and query latency low while also meeting the user-specified accuracy targets.

1) Cheap Ingest-time CNN: Focus makes indexing at ingest-time cheap by compressing and specializing the GT-CNN model for *each* video stream. (i) *Compression* of CNN models [31, 37, 42, 62, 66] uses fewer convolutional layers and other approximation techniques (§2.1). (ii) *Specialization* of CNNs [35, 65] uses the observation that a specific video stream contains only a small number of object classes and their appearance is more constrained than in a generic video (§2.2.2). Both techniques are done automatically and together result in ingest-time CNN models that are up to $98\times$ cheaper than GT-CNN.

2) Top-K ingest index: The cheap ingest-time CNNs are less accurate, i.e., their top-most results do not often match the top-most classifications of GT-CNN. Therefore, to keep the recall high, Focus associates each object with the *top-K* classification results of the cheap CNN, instead of just its top-most result. Increasing the K increases recall because the top-most result of GT-CNN often falls within the ingest-time CNN’s top-K results. At query-time, Focus uses the GT-CNN to remove objects in this larger set that do not match the class, to regain precision lost by including all the top-K.

3) Clustering similar objects: A high value of K at ingest-time increases the work to do at query time, thereby increasing query latency. To reduce this overhead, Focus clusters similar objects at ingest-time using feature vectors from the ingest-time CNN. In each cluster, at query-time, we run only the cluster centroid through GT-CNN and apply the classified result from the GT-CNN to all

objects in the cluster. Thus, if the objects are not tightly clustered, clustering can reduce precision and recall.

4) Trading off ingest vs. query costs: Focus automatically chooses the cheap CNN, its K , and specialization and clustering parameters to achieve the desired precision and recall targets. These choices also help Focus trade off between the work done at ingest-time and query-time. For instance, to save ingest work, Focus can select a cheaper ingest-time CNN, and then counteract the resultant loss in accuracy by running the expensive GT-CNN on more objects at query time. Focus chooses its parameters so as to offer a sharp improvement in one of the two costs for a small degradation in the other cost. Because the desired trade-off point is application-dependent, Focus provides users with a choice of three options: ingest-optimized, query-optimized, and balanced (the default).

Note that while our explanation is anchored on image classification CNNs, the architecture of Focus is generally applicable to all existing CNNs (e.g., face recognition). Techniques that we use for CNN compression [42, 62] and specialization [35], and feature extraction from the CNNs are all broadly applicable to all CNNs.

4. Video Ingest & Querying Techniques

In this section, we describe the main techniques used in Focus: using cheap CNN models at ingest-time (§4.1), identifying similar objects and frames to save on redundant CNN processing (§4.2), and specializing the CNNs to the specific videos that are being analyzed (§4.3). §4.4 describes setting parameters in Focus.

4.1. Cheap Ingestion

Focus indexes the live videos at *ingest-time* to reduce the *query-time* latency. We perform object detection on each frame, typically an inexpensive operation, and then will classify the extracted objects using *ingest-time* CNNs that are far cheaper than the ground-truth GT-CNN. We use these classifications to index objects by class.

Cheap ingest-time CNN: As noted earlier, the user provides Focus with a GT-CNN. Optionally, the user can also provide other classifier architectures to be used in Focus’ search for cheap CNNs, such as AlexNet [45] and

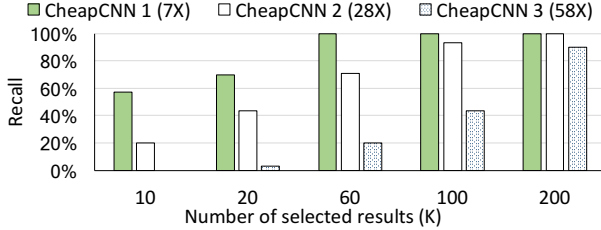


Figure 5: Effect of K on recall for three cheap CNNs. The number within the parenthesis indicates how much cheaper the model is compared to our GT-CNN, ResNet152.

VGG [66] (which vary in their resource costs and accuracies). Starting from these user-provided CNNs, Focus applies various levels of compression, such as removing convolutional layers and reducing the input image resolution (§2.1). This results in a large set of CNN options for ingestion, $\{\text{CheapCNN}_1, \dots, \text{CheapCNN}_n\}$, with a wide range of costs and accuracies.

Top-K Ingest Index: To keep recall high, Focus indexes each object using the *top* K object classes from CheapCNN_i ’s output, instead of using just the top-most class. Recall from §2.1 that the output of the CNN is a list of object classes in descending order of confidence. We empirically observe that the top-most output of the expensive GT-CNN is often contained within the top-K classes output by the cheap CNN (for a small value of K relative to the 1,000 classes recognized by the CNNs).

Figure 5 plots the effect of K on recall on one of our video streams, lausanne (see §6.1). The three models in the figure are ResNet18 [38], and ResNet18 with 3 and 5 layers removed; additionally, the input images were rescaled to 224, 112, and 56 pixels, respectively. All models were retrained on their original training data (ImageNet [63]). We make two observations.

First, we observe steady increase in recall with increasing K, for all three CheapCNNs. As the figure shows, CheapCNN₁, CheapCNN₂, and CheapCNN₃ reach 90% recall when $K \geq 60$, $K \geq 100$, and $K \geq 200$, respectively. Note that all these models recognize 1000 classes, so even $K = 200$ represents only 20% of the possible classes. Second, there is a *trade-off* between different models – the cheaper they are, the lower their recall with the same K. Overall, we conclude that by selecting the appropriate K, Focus can achieve the target recall.

Focus creates the *top-K index* of an object’s top-K classes output by CheapCNN_i at ingest-time. While filtering for objects of the queried class X using the top-K index (with the appropriate K) will have a high recall, it will have very low precision. Since we associate each object with K classes (while it has only one true class), the average precision is only $1/K$. Thus, at query time, to keep the precision high, Focus determines the *actual* class of objects from the top-K index using the expensive GT-CNN and only return objects that match the queried

class.

The selection of the cheap ingest-time CNN model (CheapCNN_i) and the K value (for the top-K results) has a significant influence on the recall of the outputs produced. Lower values of K reduce recall, i.e., Focus will miss returning frames that contain the queried objects. At the same time, higher values of K increase the number of objects to classify with GT-CNN at query time to keep precision high, and hence adds to the latency. We defer to §4.4 on how Focus sets these parameters as they have to be jointly set with other parameters in §4.2 and §4.3.

4.2. Redundancy Elimination

At query time, Focus retrieves the objects likely matching the user-specified class from the top-K index and infers their actual class using the GT-CNN. This would ensure precision of 100%, but could cause significant latency at query time. Even if this inference is parallelized across many GPUs, it would still incur a large cost. Focus uses the following observation to reduce this cost: if two objects are visually similar, their feature vectors would be closely aligned and they would likely be classified as the same class (e.g., “cars”) by the GT-CNN model (§2.2.3).

Focus *clusters* objects that are similar, invokes the expensive GT-CNN only on the cluster centroids, and assigns the centroid’s label to all objects in each cluster. Doing so dramatically reduces the work done by the GT-CNN classifier at query time. Focus uses the feature vector output by the previous-to-last layer of the cheap ingest CNN (see §2.1) for clustering. Note that Focus clusters the *objects* in the frames and not the frames as a whole. The key questions regarding clustering are *how* do we cluster (algorithm) and *when* do we cluster (system). We discuss both these key questions below.

Clustering Heuristic: We require two properties in our clustering technique. First, given the high volume of video data, it should be a single-pass algorithm to keep the overhead low, as the complexities of most clustering algorithms are *quadratic*. Second, it should make no assumption on the number of clusters and adapt to outliers in data points on the fly. Given these requirements, we use the following simple approach for *incremental* clustering, which has been well-studied in the literature [27, 55].

We put the first object into the first cluster c_1 . To cluster a new object i with a feature vector f_i , we assign it to the closest cluster c_j if c_j is at most distance T away from f_i . However, if none of the clusters are within a distance T , we create a new cluster with centroid at f_i , where T is a distance threshold. We measure distance as the L_2 norm [10] between cluster centroid and object feature vector. We keep the number of clusters at a constant M by removing the smallest ones and storing their data in the top-K index. Using this algorithm, we can keep growing the popular clusters (such as similar cars), while keeping

the complexity as $O(Mn)$, which is linear to n , the total number of objects.

Clustering can reduce both precision and recall depending on parameter T . If the centroid is classified by GT-CNN as the queried class X but the cluster contains another object of a different class, it reduces precision. If the centroid is classified as a class different than X but the cluster has an object of class X , it reduces recall. We discuss setting T in §4.4.

Clustering at Ingest vs. Query Time: Focus clusters the objects at ingest-time rather than at query-time. Clustering at query-time would involve *storing* all feature vectors, *loading* them for objects filtered from the ingest index and then clustering them. Instead, clustering *at ingest time* creates clusters right when the feature vectors are created and only stores the cluster centroids in the top- K index. This makes the query-time latency much lower and also reduces the size of the top- K index. We observe that the ordering of indexing and clustering operations is mostly *commutative* in practice and has little impact on result accuracy (we do not present these results due to space constraints). We therefore use ingest-time clustering due to its latency and storage benefits.

Pixel Differencing of Objects: While clustering primarily reduces work done at query-time (number of objects to be classified by the GT-CNN), Focus also employs *pixel differencing* among objects in adjacent incoming frames to reduce ingest cost. Specifically, if two objects have very similar pixel values, it only runs the cheap CNN on one of them and assign them both to the same cluster in our top- K index.

4.3. Video-specific Specialization of CNNs

Recall from §4.1 that Focus uses a cheap ingest-time CNN, CheapCNN $_i$ to index object classes. Focus further reduces its cost by *specializing* the ingest-time CNN model to each video stream. Model specialization benefits from two properties of objects in each video stream. First, while object classification models are trained to differentiate between thousands of object classes, many video streams contain only a small number of classes (§2.2.2). Second, objects in a specific stream are often visually more constrained than objects in general (say, compared to the ImageNet [63] dataset). The cars and buses that occur in a specific traffic camera have much less variability, e.g., they have very similar angle, distortion and size, than a generic set of vehicles.

Instead of trying to differentiate among thousands of object classes, differentiating among just (say) fifty classes *and* in a specific camera’s video is a much simpler task, requiring simpler image features and smaller image resolutions. As a result, the specialized models are *smaller* and *more accurate* [35]. For example, by retraining a stream-specific CheapCNN $_i$, we can achieve similar ac-

curacy on video streams, while removing 1/3 of the convolutional layers and making the input image $4\times$ smaller in resolution. This leads to the specialized CheapCNN $_i$ being $10\times$ cheaper than even the generic CheapCNN $_i$.

Since the specialized CNN classifies across fewer classes, they are more accurate, which allows Focus to select a much smaller K (for the top- K ingest index) to meet the desired recall. We find that specialized models can use $K = 2$ or 4 , much smaller than the typical $K = 60 \sim 200$ for the generic cheap CNNs (Figure 5). Smaller K directly translates to fewer objects that have to be classified by GT-CNN at query time, thus reducing latency.

Model Retraining: On each video stream Focus periodically obtains a small sample of video frames and classifies their objects using GT-CNN to estimate the ground truth of distribution of object classes for the video (similar to Figure 3). From this distribution, Focus selects the most frequently occurring L_s object classes to retrain new specialized models. As we saw in §2.2.2, there is usually a “power law” in the distribution of classes – just a handful of classes account for a dominant majority of the objects – thus, low values of L_s usually suffice.¹

Specialization is also based off a family of CNN architectures (such as ResNet [38], AlexNet [45], and VGG [66]) with different number of convolution layers, similar to §4.1. Specialization adds to the set of options available for ingest CNNs ($\{\text{CheapCNN}_1, \dots, \text{CheapCNN}_n\}$ in §4.1), and Focus picks the best model (CheapCNN $_i$) and the corresponding K for the index.

“OTHER” class: While Focus specializes the CNN towards the most frequently occurring L_s classes, we also want to support querying of the *less* frequent classes. For this purpose, Focus includes an additional class called “OTHER” in the specialized model.² Being classified as OTHER simply means not being one of the L_s classes. At query time, if the queried class is among the OTHER classes of the ingest CNN’s index, Focus extracts all the clusters that match the OTHER class and classifies their centroids through the GT-CNN model.

The parameter L_s (for each stream) exposes the following trade-off. Using a small L_s allows us to train a simpler model with cheaper ingest cost and lower query-time latency *for the popular classes*, however, it also leads to a larger fraction of objects falling in the OTHER class; querying for them will be expensive because all those objects will have to be classified by the GT-CNN. Using a larger value of L_s , on the other hand, leads to a more expensive ingest and query-time models, but cheaper querying for the OTHER classes. We select L_s next in §4.4.

¹ Specialized CNNs can be retrained quickly on a small dataset. Retraining is relatively infrequent and done once every few days.

² Since there will be considerably fewer objects in the video belonging to the OTHER class, we proportionally re-weight the training data to contain equal number of objects of all the classes.

4.4. Balancing Accuracy, Latency, and Cost

Focus’ goals of high accuracy, low ingest cost and low query latency are impacted by the parameters in Focus’ techniques – K , the number of top results from the ingest-time CNN to index an object; L_s , the number of popular object classes we use to create a specialized model; CheapCNN $_i$, the specialized ingest-time cheap CNN; and T , the distance threshold for clustering objects.

The effect of these four parameters is intertwined. All the four parameters impact ingest cost, query latency, and recall, but only T impacts precision. This is because we apply the cluster centroid’s classification by GT-CNN to all the objects in its cluster. Thus, if the clustering is not tight (i.e., high value of T), we lose precision.

Parameter Selection: Focus selects parameter values *per video stream*. It samples a representative fraction of frames of the video stream and classifies them using GT-CNN for the ground truth. For each combination of parameter values, Focus computes the expected precision and recall (using the ground truths generated by GT-CNN) that would be achieved for each of the object classes. To navigate the combinatorial space of options for these parameters, we adopt a two-step approach. In the first step, Focus chooses CheapCNN $_i$, L_s and K using only the recall target. In the next step, Focus iterates through the values of T , the clustering distance threshold, and only select values that meet the precision target.

Trading off Ingest Cost and Query Latency: Among the combination of values that meet the precision and recall targets, the selection is based on *balancing* the ingest- and query-time costs. For example, picking a model CheapCNN $_i$ that is more accurate will have higher ingest cost, but lower query cost because we can use a lower K . Using a less accurate CheapCNN $_i$ will have the opposite effect. Focus identifies “intelligent defaults” that sharply improve one of the two costs for a small worsening of the other cost.

Figure 6 illustrates the parameter selection based on the ingest cost and query latency for one of our video streams (auburn_c). The figure plots all the viable “configurations” (i.e., set of parameters that meet the precision and recall target) based on their ingest cost (i.e., cost of CheapCNN $_i$) and query latency (i.e., the number of clusters according to K, L_s, T). We first draw the *Pareto boundary* [17], which is the set of configurations that cannot improve one metric without worsening the other. Focus can discard all the other configurations because at least one point on the Pareto boundary is better than them in both metrics. Focus balances between the ingest cost and query latency (Balance in Figure 6) by selecting the configuration that minimizes the *sum of ingest and query cost* (measured in total GPU cycles).

Focus also allows for other configurations based on the application’s preferences and query rates. Opt-Ingest

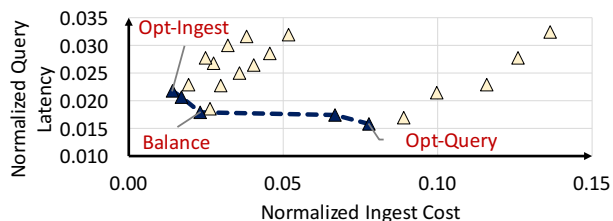


Figure 6: Parameter selection based on trading off ingest cost and query latency. The ingest cost is normalized to ingesting all objects with ResNet152, while the query latency is normalized to the time for querying all objects with ResNet152. The dashed line is the Pareto boundary.

minimizes the ingest cost and is applicable when the application expects most of the video streams to not get queried (such as a surveillance cameras), as this policy also minimizes the amount of wasted ingest work. On the other hand, Opt-Query minimizes query latency even if it incurs a heavy ingest cost. Such flexibility allows Focus to fit different applications.

5. Implementation Details

We describe the key aspects in Focus’s implementation.

Worker Processes. Focus’s ingest-time work is distributed across many machines, with each machine running one *worker* process for *each video stream’s* ingestion. The ingest worker receives the live video stream, and extracts the moving objects (using background subtraction [81]); it is extensible to plug in any other object detector. The detected objects are sent to the ingest-time CNN to infer the top- K classes and the feature vectors. The ingest worker uses the features to cluster objects in its video stream and stores the top- K index in MongoDB [12] for efficient retrieval at query-time.

Worker processes also serve queries by fetching the relevant frames off the top- K index database and classifying the objects with GT-CNN. We parallelize a query’s work across many worker processes if resources are idle.

GPUs for CNN classification. The cheap CNNs and GT-CNN execute on GPUs (or other hardware accelerators for CNNs) which could either be local on the same machine as the worker processes or “disaggregated” on a remote cluster. This detail is abstracted away from our worker process and it seamlessly works with both designs.

Dynamically adjusting K at query-time. As an enhanced technique, we can select a new $K_x \leq K$ at query time and only extract clusters where class X appears among the top- K_x classes; this will result in fewer clusters and thus also lower query-time latency. This technique is useful in two scenarios: 1) some classes might be very accurately classified by the cheap CNN; using a lower K_x will still meet the user-specified accuracy, yet will result in much lower latency; 2) if we want to retrieve only *some objects of class X* , we can use very low K_x to quickly retrieve them. If more objects are required, we can increase

K_x to extract a new batch of results.

6. Evaluation

We evaluate the Focus prototype with more than 150 hours of videos from 13 real video streams that span across traffic cameras, surveillance cameras, and news channels. Our highlights are:

1. On average, Focus is simultaneously $58\times$ (up to $98\times$) cheaper than the Ingest-all baseline in its GPU consumption and $37\times$ (up to $57\times$) faster than the Query-all baseline in query latency, all the while achieving at least 95% precision and recall (§6.2, §6.3).
2. Focus provides a rich trade-off space between ingest cost and query latency. Among the video streams, the ingest cost is up to $141\times$ cheaper than the Ingest-all baseline (and reduces query latency by $46\times$) if optimizing for low-cost ingest. The query latency is reduced by up to $66\times$ (with $11\times$ cheaper ingest) if optimizing for query latency (§6.4).
3. Focus is effective under broad conditions such as high accuracy targets (one order-of-magnitude savings even for 99% accuracy target, §6.5) and various frame sampling rates (30 fps-1 fps, §6.6).

6.1. Setup

Software Tools. We use OpenCV 3.2.0 [14] to decode the videos into frames, and then use the built-in background subtraction algorithm [43] in OpenCV to extract moving objects from video frames. We use background subtraction instead of object detector CNNs (e.g., YOLOv2 [59] or Faster R-CNN [60]) to detect objects because: (1) running background subtraction is orders of magnitude faster than running these CNNs, and (2) background subtraction can detect moving objects more reliably, while object detector CNNs usually have difficulties on small objects [52]. Nonetheless, our system can seamlessly use object detector CNNs as well. We run and train CNNs with Microsoft Cognitive Toolkit 2.1 [54], an open-source deep learning system.

Video Datasets. We evaluate 13 live video streams that span across traffic cameras, surveillance cameras, and news channels. We evaluate each video stream for 12 hours, which evenly cover day time and night time. Table 1 summarizes the video characteristics. By default, we evaluate each video at 30 fps and also evaluate the sensitivity to other frame rates (§6.6). In some figures we only show a representative sample of 9 cameras to improve legibility.

Accuracy Target. We use ResNet152, a state-of-the-art CNN, as our ground-truth CNN (GT-CNN). We evaluate all extracted objects with the GT-CNN and use the results as the correct answers. We define a class present

³The video streams are obtained from real and operational traffic cameras in a city. We mask the city name for anonymity.

Table 1: Video dataset characteristics

Type	Name	Location	Description
Traffic	auburn_c	AL, USA	A commercial area intersection in the City of Auburn [6]
	auburn_r	AL, USA	A residential area intersection in the City of Auburn [5]
	city_a_d	USA	A downtown intersection in City A ³
	city_a_r	USA	A residential area intersection in City A ³
	bend	OR, USA	A road-side camera in the City of Bend [8]
	jacksonh	WY, USA	A busy intersection (Town Square) in Jackson Hole [9]
Surveillance	church_st	VT, USA	A video stream rotates among cameras in a shopping mall (Church Street Marketplace) [3]
	lausanne	Switzerland	A pedestrian plazalatency (Place de la Palud) in Lausanne [11]
	oxford	England	A bookshop street in the University of Oxford [15]
	sittard	Netherlands	A market square in Sittard [4]
News	cnn	USA	News channel
	foxnews	USA	News channel
	msnbc	USA	News channel

in a one-second segment of video if the GT-CNN reports such class in 50% of the frames in that segment. We use this criteria as our ground truth because our GT-CNN (ResNet152) sometimes gives different answers to the exact same object in consecutive frames, and this criteria can effectively eliminate these random, erroneous results. We set our default accuracy target as 95% recall and 95% precision. We also evaluate the results with other accuracy targets such as 97%, 98% and 99% (§6.5). Note that in most practical cases, only one of the two metrics (recall or accuracy) needs to be high. For example, an investigator cares about high recall, and looking through some irrelevant results is an acceptable trade-off. By setting both targets high, we are lower bounding the performance improvements that Focus can achieve.

Baselines. We use two baselines for comparisons: (1) Ingest-all, the baseline system that uses GT-CNN to analyze all objects at ingest time, and stores the inverted index for query; and (2) Query-all, the baseline system that simply extracts objects at ingest time, and uses GT-CNN to analyze all the objects that fall into the query interval at query time. Note that we strengthen both baselines with basic motion detection (background subtraction). Therefore, the baselines *do not* run any GT-CNN on the frames that have no moving objects. Note that not running GT-CNN on frames with no moving objects is one of the core techniques in the recent NoScope work [44].

Metrics. We use two performance metrics. The first metric is *ingest cost*, which is the GPU time to ingest each video. The second metric is *query latency*, which is the latency for an object class query. Specifically, for each video stream, we evaluate all dominant object classes and take the average of their latencies. (Querying for non-dominant “OTHER” classes is much cheaper than querying popular classes, and would skew the results because there are far more such classes; thus, we focus on

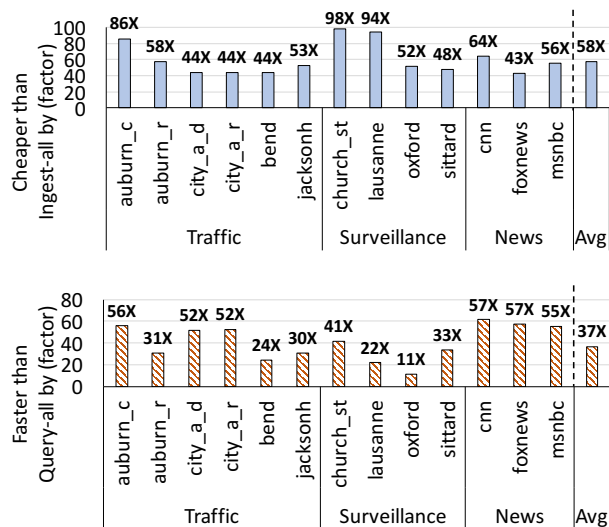


Figure 7: (Top) Focus ingest cost compared to Ingest-all. **(Bottom)** Focus query latency compared to Query-all.

the popular ones.) Both metrics include only GPU time spent classifying images and exclude other (CPU) time spent decoding video frames, detecting moving objects, recording and loading video, and reading and writing to the top-K index. We focus solely on GPU time because when the GPU is involved, it is the bottleneck resource. The query latency of Ingest-all is 0 and the ingest cost of Query-all is 0.

Experiment Platform. We run the experiments on our local cluster. Each machine in the cluster is equipped with a state-of-the-art GPU (NVIDIA GTX Titan X), 16-core Intel Xeon CPU (E5-2698), 64 GB RAM, a 40 GbE NIC, and runs 64-bit Ubuntu 16.04 LTS.

6.2. End-to-End Performance

We first show the end-to-end performance of Focus by showing its ingest cost and query latency when Focus aims to balance these two metrics (§4.4). Figure 7 compares the ingest cost of Focus with Ingest-all and the query latency of Focus with Query-all. We make two main observations.

First, Focus significantly improves query latency with a very small ingest cost. Focus makes queries by an average of 37× faster than Query-all with a very small cost at ingest time (an average of 58× cheaper than Ingest-all). With a 10-GPU cluster, the query latency on a 24-hour video goes down from one hour to less than two minutes. The processing cost of each video stream also goes down from \$250/month to \$4/month. This shows that Focus can strike a very good balance between these two competing goals very effectively.

Second, Focus is effective across different video streams with various characteristics. It makes queries 11× to 57× faster with a very small ingest time cost (48× to 98× cheaper) across busy intersections

(auburn_c, city_a_d and jacksonh), normal intersections or roads (auburn_r and city_a_r, bend), rotating cameras (church_st), busy plazas (lausanne and sittard), a university street (oxford), and different news channels (cnn, foxnews, and msnbc). Among these videos, the gains in query latency are smaller for relatively less busy videos (auburn_r, bend, lausanne, and oxford). This is because these videos are dominated by fewer object classes, and Focus has more work (i.e., analysis using GT-CNN) to do at query time for these classes. We conclude that the core techniques of Focus are general and effective on a variety of real-world videos.

6.3. Effect of Different Focus Components

Figure 8 shows the breakdown of ingest-time cost and query latency across different design points of Focus: (1) Compressed model, which applies a generic compressed model for indexing at ingest time, (2) Compressed + Specialized model, which uses a per-stream specialized and compressed model for indexing, and (3) Compressed + Specialized model + Clustering, which adds feature-based clustering at ingest time to reduce redundant work at query time. All of the above include the top-K index and using GT-CNN at query-time, and achieve the same accuracy of 95%. Three main observations are in order.

First, generic compressed models provide benefits for both ingest cost and query latency, but they are not the major source of improvement. This is because the accuracy of a generic compressed model degrades significantly when we remove convolutional layers. In order to retain the accuracy target, we need to choose relatively expensive compressed models (CheapCNN_i) and a larger K, which incur higher ingest cost and query latency.

Second, specializing the model (in addition to compressing it) greatly reduces ingest cost and query latency. Because of fewer convolutional layers and smaller input resolution, our specialized models are 7× to 71× cheaper than the GT-CNN, while retaining the accuracy target for each video streams. Running a specialized model at ingest time speeds up query latency by 5× to 25× (Figure 8b).

Third, clustering is a very effective technique to further reduce query latency with unnoticeable costs at ingest time. As Figure 8b shows, using clustering (on top of a specialized compressed model) reduces the query latency by up to 56×, significantly better than just running a specialized model at ingest time. This gain comes with a negligible cost (Figure 8a), because we run our clustering algorithm (§4.2) on the CPUs of the ingest machine, which is fully pipelined with the GPUs that run the specialized CNN model.

6.4. Ingest Cost vs. Query Latency Trade-off

One of the interesting features of Focus is the flexibility to tune its system parameters to achieve different application

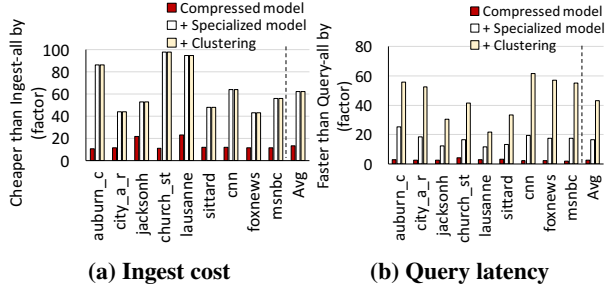


Figure 8: Effect of different Focus components

goals (§4.4). Figure 1 from §1 depicted three alternative settings for Focus that illustrate the trade-off space between ingest cost and query latency, using the auburn_c video stream: (1) Focus-Opt-Query, which optimizes for query latency by increasing ingest cost, (2) Focus-Balance, which is the default option that balances these two metrics (§4.4), and (3): Focus-Opt-Ingest, which is the opposite of Focus-Opt-Query. The results are shown relative to the two baselines. The chart at the right of the figure is the zoomed-in region that covers the three settings of Focus, and each data label (I, Q) indicates its ingest cost is $I\times$ cheaper than Ingest-all, while its query latency is $Q\times$ faster than Query-all.

As Figure 1 shows, Focus offers very good options in the trade-off space between ingest cost and query latency. Focus-Opt-Ingest achieves $141\times$ cheaper cost than Ingest-all to ingest the video stream, and makes the query $46\times$ faster than doing nothing at ingest (Query-all). On the other hand, Focus-Opt-Query reduces query latency by $63\times$ with a relatively higher ingest cost, but it is still $26\times$ cheaper than Ingest-all. As they are all good options compared to the baselines, such flexibility allows a user to tailor Focus for different contexts. For example, a traffic camera that requires fast turnaround time for queries can use Focus-Opt-Query, while a surveillance video stream that will be queried very rarely would choose Focus-Opt-Ingest to reduce the amount of wasted ingest cost.

Figure 9 shows the (I, Q) values for both Focus-Opt-Ingest (Opt-I) and Focus-Opt-Query (Opt-Q) for the representative videos. As the figure show, the trade-off flexibility exists among all the other videos. On average, Focus-Opt-Ingest spends only $95\times$ cheaper ingest cost to provide $35\times$ query latency reduction. On the other hand, Focus-Opt-Query makes queries $49\times$ faster with a higher ingest cost ($15\times$ cheaper than Ingest-all). We conclude that Focus provides good flexibility between ingest cost and query latency, and makes it a better fit in different contexts.

6.5. Sensitivity to Accuracy Target

Figures 10 and 11 illustrate the improvements of ingest cost and query latency of Focus compared to the baselines under different accuracy targets. Other than the default 95% accuracy target (recall and precision), we evaluate

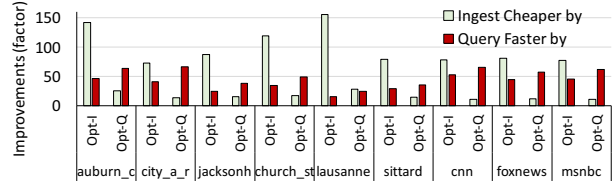


Figure 9: Trade-offs between ingest cost and query latency

three higher targets, 97%, 98%, and 99%.

As the figures show, with higher accuracy targets, the ingest costs are about the same, and the improvement of query latency decreases. Focus keeps the ingest cost similar ($62\times$ to $64\times$ cheaper than the baseline) because it still runs the specialized and compressed CNN at ingest time. However, when the accuracy targets are higher, Focus needs to select more top-K classification results, which increases the work at query time. On average, the query latency of Focus is faster than Query-all by $15\times$, $12\times$, and $8\times$ with respect to 97%, 98%, and 99% accuracy targets. We conclude that the techniques of Focus can achieve higher accuracy targets with significant improvements on both ingest cost and query latency.

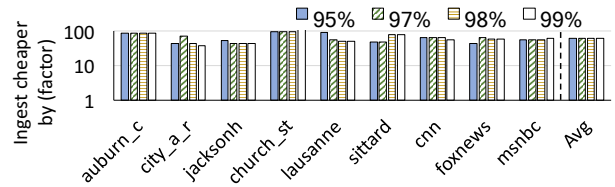


Figure 10: Ingest cost sensitivity to accuracy target

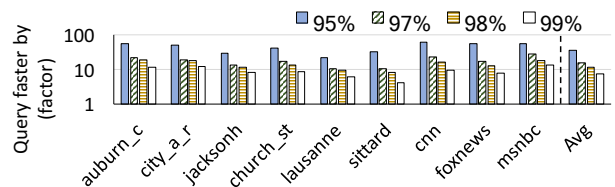


Figure 11: Query latency sensitivity to accuracy target

6.6. Sensitivity to Frame Sampling

A common approach to reduce the video processing time is to use frame sampling (i.e., periodically select a frame to process). However, not all applications can use frame sampling because it can miss objects that show up and disappear within a frame sampling window. As the frame sampling rate is an application dependent choice, we study the sensitivity of Focus’s performance to different frame rates. Figures 12 and 13 show the ingest cost and query latency of Focus at different frame rates (i.e., 30 fps, 10 fps, 5 fps, and 1 fps) compared to Ingest-all and Query-all, respectively. We make two observations.

First, the ingest cost reduction is roughly the same across the different frame rates. On average, the ingest

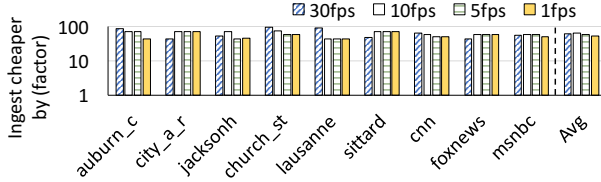


Figure 12: Ingest cost sensitivity to frame sampling

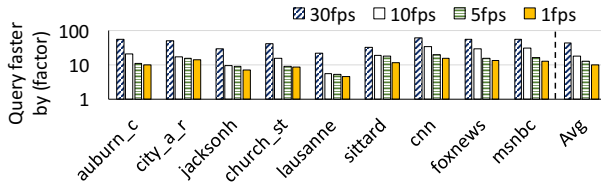


Figure 13: Query latency sensitivity to frame sampling

cost of Focus is $62\times$ cheaper than Ingest-all at 30 fps, and it is $64\times$ to $58\times$ cheaper at lower frame rates. This is because the major ingest cost saving comes from the specialized and compressed CNN models (§6.3), which are orthogonal to frame sampling rates.

Second, the query latency improvement of Focus degrades with lower frame rates. This is expected because one of our key techniques to reduce query latency is redundancy elimination, especially clustering similar objects using CNN feature vectors. At lower frame rates, the benefit of this technique reduces because there are fewer redundancies. Nonetheless, on average, Focus is still one order of magnitude faster than Query-all at a very low frame rate (1 fps).

6.7. Applicability with Different Query Rate

There are two factors that can affect the applicability of Focus: 1) the number of classes that get queried over time and 2) the fraction of videos that get queried. In the first extreme case where all the classes and all the videos are queried, Ingest-all could be a good option because its cost is amortized among all the queries. In our study, even in such an extreme case, the overall cost of Focus is still $4\times$ cheaper than Ingest-all on average (up to $6\times$ cheaper) because we run a very cheap CNN at ingest time, and we run GT-CNN *per object cluster* only *once* (§5), so the overall cost is still cheaper than Ingest-all.

The second extreme case is only a tiny fraction of videos gets queried. While Focus can save the ingest cost by up to $141\times$ (§6.4), it can be more costly than Query-all if the fraction of videos gets queried is less than $\frac{1}{141} = 0.7\%$. In such a case, we can choose to do nothing at ingest time and run all the techniques of Focus only at query time when we know the fraction of videos that get queried. While this approach increases query latency, it still reduces the query latency by a average of $22\times$ (up to $34\times$) than Query-all in our evaluation. We conclude that Focus is still better than both baselines even under extreme query rates.

7. Related Work

To our best knowledge, Focus is the first system that offers low-cost, low-latency, and high-accuracy video queries by balancing between ingest-time cost and query latency. We now discuss work related to our key techniques.

1) Cascaded classification. Various works in vision research propose speeding up classification by cascading a series of classifiers. Viola et al. [75] is the earliest work which cascades a series of classifiers (from the simplest to the most complicated) to quickly disregard regions in an image. Many improvements followed (e.g., [50, 76, 77]). CNNs are also cascaded (e.g., [26, 35, 49, 70]) to reduce object detection latency. Our work is different in two major ways. First, we *decouple* the compressed CNN from the GT-CNN, which allows us to choose from a wider range for ingest-time CNNs and allows for better trade-offs between ingest cost and query latency, a key aspect of our work. Second, we cluster similar objects using CNN features to eliminate redundant work, which is a new and effective technique for video streams.

2) Neural network compression. Recent work proposes various techniques to reduce the running time of CNNs. These techniques include shallow models [21], predicting weights [33], matrix pruning [31, 37], model quantization [36], and others (e.g., [20, 34, 39, 41, 42, 57, 62]). Our work is largely orthogonal to these, in that our system is not tied to a specific model compression technique, and we can employ any of these techniques.

3) Context-specific model specialization. Context-specific specialization of models can improve accuracy [53] or reduce running time [35, 44, 65]. Among these, the closest to our work is Kang et al.’s proposal, NoScope [44], which aims to optimize CNN-based video queries. A few key differences stand out. First, NoScope applies all the optimizations at query-time, while Focus adopts a different architecture by splitting work between ingest- and query-time. Thus, Focus trades off higher ingest cost for *even lower* query latency. Second, NoScope optimizes CNNs for a single class, while we optimize ingest CNNs for all frequent classes in the stream and allow queries even for the rare – OTHER – classes. Finally, we use the object feature vectors to cluster similar objects and create an index to map classes to clusters; this allows us to efficiently query across *all* classes, while NoScope has to redo all query-time work, including training specialized CNNs, for each query.

4) Stream processing systems. Systems for general stream data processing (e.g., [1, 18, 19, 24, 28, 29, 51, 56, 73, 74, 79]) and specific to video analytics (e.g., [80]) mainly focus on the general stream processing challenges such as load shedding, fault tolerance, distributed execution, or limited network bandwidth. In contrast, our work is specific for querying on recorded video data with ingest and query trade-offs, thus it is mostly orthogonal to these.

We can integrate Focus with one of these general stream processing system to build a more fault tolerable system.

5) Video indexing and retrieval. A large body of works in multimedia and information retrieval research propose various content-based video indexing and retrieval techniques to facilitate queries on videos (e.g., [40,48,68,69]). Among them, most works focus on indexing videos for different types of queries such as shot boundary detection [78], semantic video search [30], video classification [25], or spatio-temporal information-based video retrieval [61]. Some works (e.g., [32, 67]) focus on the query interface to enable query by keywords, concepts, or examples. These works are largely orthogonal to our work because we focus on the *cost and latency* of video queries, not query types or interfaces. We believe our idea of splitting ingest-time and query-time work is generic for videos queries, and can be extended to different types of queries.

8. Conclusion

Answering queries of the form, *find me frames that contain objects of class X* is an important workload on recorded video datasets. Such queries are used by analysts and investigators, and it is crucial to answer them with low latency and low cost. We present Focus, a system that performs low cost ingest-time analytics on live video that later facilitates low-latency queries on the recorded videos. Focus uses compressed and specialized CNNs at ingest-time that substantially reduces cost. It also clusters similar objects to reduce the work done at query-time, and hence the latency. Focus selects the ingest-time CNN and its parameters to smartly trade-off between the ingest-time cost and query-time latency. Our evaluations using 150 hours of video from traffic, surveillance, and news domains show that Focus reduces GPU consumption by $58\times$ and makes queries $37\times$ faster compared to current baselines. We conclude that Focus is a promising approach to querying large video datasets. We hope that Focus will enable future works on better determining the ingest-time and query-time trade-offs in video querying systems. Our next steps include training a specialized and highly accurate query-time CNN for each stream and object to further reduce query latency.

References

- [1] “Apache Storm.” [Online]. Available: <http://storm.apache.org/index.html>
- [2] “Avigilon,” <http://avigilon.com/products/>.
- [3] “Church Street Market Place.” [Online]. Available: <https://www.youtube.com/watch?v=S3B18AuKPds>
- [4] “City Cam, WebcamSittard: Town Square Sittard (NL).” [Online]. Available: <https://www.youtube.com/watch?v=Zb9koIwo3Js>
- [5] “City of Auburn North Ross St and East Magnolia Ave.” [Online]. Available: <https://www.youtube.com/watch?v=cjuskMMYILA>
- [6] “City of Auburn Toomer’s Corner Webcam.” [Online]. Available: https://www.youtube.com/watch?v=yJAK_FozAmI
- [7] “Genetec,” <https://www.genetec.com/>.
- [8] “Greenwood Avenue Bend, Oregon.” [Online]. Available: <https://www.youtube.com/watch?v=SNz323Cyago>
- [9] “Jackson Hole Wyoming USA Town Square.” [Online]. Available: <https://www.youtube.com/watch?v=psfFJR3vZ78>
- [10] “L² Norm.” [Online]. Available: <http://mathworld.wolfram.com/L2-Norm.html>
- [11] “Lausanne, Place de la Palud.” [Online]. Available: <https://www.youtube.com/watch?v=GdhEsWcV4iE>
- [12] “MongoDB.” [Online]. Available: <https://www.mongodb.com/>
- [13] “Nvidia Tesla K80.” [Online]. Available: <http://www.nvidia.com/object/tesla-k80.html>
- [14] “OpenCV 3.2.” [Online]. Available: <http://opencv.org/opencv-3-2.html>
- [15] “Oxford Martin School Webcam - Broad Street, Oxford.” [Online]. Available: <https://www.youtube.com/watch?v=Qhq4vQdfrFw>
- [16] “Top Video Surveillance Trends for 2016.” [Online]. Available: <https://technology.ihc.com/api/binary/572252>
- [17] “Wikipedia: Pareto efficiency.” [Online]. Available: https://en.wikipedia.org/wiki/Pareto_efficiency
- [18] D. J. Abadi, Y. Ahmad, M. Balazinska, U. Çetintemel, M. Cherniack, J. Hwang, W. Lindner, A. Maskey, A. Rasin, E. Ryvkina, N. Tatbul, Y. Xing, and S. B. Zdonik, “The design of the Borealis stream processing engine,” in *CIDR*, 2005.
- [19] L. Amini, H. Andrade, R. Bhagwan, F. Eskesen, R. King, P. Selo, Y. Park, and C. Venkatramani, “SPC: A distributed, scalable platform for data mining,” in *DM-SSP*, 2006.
- [20] S. Anwar, K. Hwang, and W. Sung, “Fixed point optimization of deep convolutional neural networks for object recognition,” in *ICASSP*, 2015.
- [21] J. Ba and R. Caruana, “Do deep nets really need to be deep?” in *NIPS*, 2014.
- [22] A. Babenko and V. S. Lempitsky, “Aggregating deep convolutional features for image retrieval,” in *ICCV*, 2015.
- [23] A. Babenko, A. Slesarev, A. Chigorin, and V. S. Lempitsky, “Neural codes for image retrieval,” in *ECCV*, 2014.
- [24] P. Bailis, E. Gan, S. Madden, D. Narayanan, K. Rong, and S. Suri, “MacroBase: Prioritizing attention in fast data,” in *SIGMOD*, 2017.
- [25] D. Brezeale and D. J. Cook, “Automatic video classification: A survey of the literature,” *IEEE Trans. Systems, Man, and Cybernetics, Part C*.
- [26] Z. Cai, M. J. Saberian, and N. Vasconcelos, “Learning complexity-aware cascades for deep pedestrian detection,” in *ICCV*, 2015.
- [27] F. Cao, M. Ester, W. Qian, and A. Zhou, “Density-based clustering over an evolving data stream with noise,” in *SIAM International Conference on Data Mining*, 2006.
- [28] D. Carney, U. Çetintemel, M. Cherniack, C. Convey, S. Lee, G. Seidman, M. Stonebraker, N. Tatbul, and S. B.

- Zdonik, "Monitoring streams - A new class of data management applications," in *VLDB*, 2002.
- [29] S. Chandrasekaran, O. Cooper, A. Deshpande, M. J. Franklin, J. M. Hellerstein, W. Hong, S. Krishnamurthy, S. Madden, F. Reiss, and M. A. Shah, "TelegraphCQ: Continuous dataflow processing," in *SIGMOD*, 2003.
- [30] S. Chang, W. Ma, and A. W. M. Smeulders, "Recent advances and challenges of semantic image/video search," in *ICASSP*, 2007.
- [31] W. Chen, J. T. Wilson, S. Tyree, K. Q. Weinberger, and Y. Chen, "Compressing neural networks with the hashing trick," *CoRR*, vol. abs/1504.04788, 2015.
- [32] M. G. Christel and R. M. Conescu, "Mining novice user activity with TRECVID interactive retrieval tasks," in *CIVR*.
- [33] M. Denil, B. Shakibi, L. Dinh, M. Ranzato, and N. de Freitas, "Predicting parameters in deep learning," in *NIPS*, 2013.
- [34] E. L. Denton, W. Zaremba, J. Bruna, Y. LeCun, and R. Fergus, "Exploiting linear structure within convolutional networks for efficient evaluation," in *NIPS*, 2014.
- [35] S. Han, H. Shen, M. Philipose, S. Agarwal, A. Wolman, and A. Krishnamurthy, "MCDNN: An approximation-based execution framework for deep stream processing under resource constraints," in *MobiSys*, 2016.
- [36] S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural network with pruning, trained quantization and huffman coding," in *ICLR*, 2016.
- [37] S. Han, J. Pool, J. Tran, and W. Dally, "Learning both weights and connections for efficient neural network," in *NIPS*, 2015.
- [38] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.
- [39] G. E. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *CoRR*, vol. abs/1503.02531, 2015.
- [40] W. Hu, N. Xie, L. Li, X. Zeng, and S. J. Maybank, "A survey on visual content-based video indexing and retrieval," *IEEE Trans. Systems, Man, and Cybernetics, Part C*.
- [41] K. Hwang and W. Sung, "Fixed-point feedforward deep neural network design using weights +1, 0, and -1," in *SiPS*, 2014.
- [42] M. Jaderberg, A. Vedaldi, and A. Zisserman, "Speeding up convolutional neural networks with low rank expansions," *CoRR*, vol. abs/1405.3866, 2014.
- [43] P. KaewTraKulPong and R. Bowden, "An improved adaptive background mixture model for real-time tracking with shadow detection," in *AVSS*, 2001.
- [44] D. Kang, J. Emmons, F. Abuzaid, P. Bailis, and M. Zaharia, "NoScope: Optimizing deep CNN-based queries over video streams at scale," *PVLDB*, 2017.
- [45] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in *NIPS*, 2012.
- [46] S. Lawrence, C. L. Giles, A. C. Tsoi, and A. D. Back, "Face recognition: a convolutional neural-network approach," *IEEE Trans. Neural Networks*, 1997.
- [47] Y. LeCun, B. E. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. E. Hubbard, and L. D. Jackel, "Backpropagation applied to handwritten zip code recognition," *Neural Computation*, 1989.
- [48] M. S. Lew, N. Sebe, C. Djeraba, and R. Jain, "Content-based multimedia information retrieval: State of the art and challenges," *TOMCCAP*.
- [49] H. Li, Z. Lin, X. Shen, J. Brandt, and G. Hua, "A convolutional neural network cascade for face detection," in *CVPR*, 2015.
- [50] R. Lienhart and J. Maydt, "An extended set of Haar-like features for rapid object detection," in *ICIP*, 2002.
- [51] W. Lin, H. Fan, Z. Qian, J. Xu, S. Yang, J. Zhou, and L. Zhou, "StreamScope: Continuous reliable distributed processing of big data streams," in *NSDI*, 2016.
- [52] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. E. Reed, C. Fu, and A. C. Berg, "SSD: Single shot multibox detector," in *ECCV*.
- [53] A. Mhalla, H. Maâmatou, T. Chateau, S. Gazzah, and N. E. B. Amara, "Faster R-CNN scene specialization with a sequential monte-carlo framework," in *DICTA*.
- [54] Microsoft, "Microsoft Cognitive Toolkit." [Online]. Available: <https://www.microsoft.com/en-us/cognitive-toolkit/>
- [55] L. O'Callaghan, N. Mishra, A. Meyerson, and S. Guha, "Streaming-data algorithms for high-quality clustering," in *ICDE*, 2002.
- [56] A. Rabkin, M. Arye, S. Sen, V. S. Pai, and M. J. Freedman, "Aggregation and degradation in JetStream: Streaming analytics in the wide area," in *NSDI*, 2014.
- [57] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi, "XNOR-Net: ImageNet classification using binary convolutional neural networks," in *ECCV*, 2016.
- [58] A. S. Razavian, H. Azizpour, J. Sullivan, and S. Carlsson, "CNN features off-the-shelf: An astounding baseline for recognition," in *CVPR Workshops*, 2014.
- [59] J. Redmon and A. Farhadi, "YOLO9000: Better, faster, stronger," *CoRR*, vol. abs/1612.08242, 2016.
- [60] S. Ren, K. He, R. B. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," in *NIPS*, 2015.
- [61] W. Ren, S. Singh, M. Singh, and Y. S. Zhu, "State-of-the-art on spatio-temporal information-based video retrieval," *Pattern Recognition*.
- [62] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio, "FitNets: Hints for thin deep nets," *CoRR*, vol. abs/1412.6550, 2014.
- [63] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "ImageNet large scale visual recognition challenge," *IJCV*, 2015.
- [64] F. Schroff, D. Kalenichenko, and J. Philbin, "FaceNet: A unified embedding for face recognition and clustering," in *CVPR*, 2015.
- [65] H. Shen, S. Han, M. Philipose, and A. Krishnamurthy, "Fast video classification via adaptive cascading of deep models," in *CVPR*, 2017.
- [66] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *ICLR*, 2015.
- [67] C. Snoek, K. E. A. van de Sande, O. de Rooij, B. Huurnink, E. Gavves, D. Odijk, M. de Rijke, T. Gevers, M. Worring, D. Koelma, and A. W. M. Smeulders, "The medi-amill TRECVID 2010 semantic video search engine," in *TRECVID 2010 workshop participants notebook papers*, 2010.

- [68] C. Snoek and M. Worring, "Multimodal video indexing: A review of the state-of-the-art," *Multimedia Tools Appl.*
- [69] C. G. M. Snoek and M. Worring, "Concept-based video retrieval," *Foundations and Trends in Information Retrieval.*
- [70] Y. Sun, X. Wang, and X. Tang, "Deep convolutional network cascade for facial point detection," in *CVPR*, 2013.
- [71] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. E. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *CVPR*, 2015.
- [72] P.-N. Tan, M. Steinbach, and V. Kumar, *Introduction to Data Mining, (First Edition)*. Boston, MA, USA: Addison-Wesley Longman Publishing Co., Inc., 2005.
- [73] N. Tatbul, U. Çetintemel, and S. B. Zdonik, "Staying FIT: efficient load shedding techniques for distributed stream processing," in *VLDB*, 2007.
- [74] Y. Tu, S. Liu, S. Prabhakar, and B. Yao, "Load shedding in stream databases: A control-based approach," in *VLDB*, 2006.
- [75] P. A. Viola and M. J. Jones, "Rapid object detection using a boosted cascade of simple features," in *CVPR*, 2001.
- [76] Z. E. Xu, M. J. Kusner, K. Q. Weinberger, and M. Chen, "Cost-sensitive tree of classifiers," in *ICML*, 2013.
- [77] Q. Yang, C. X. Ling, X. Chai, and R. Pan, "Test-cost sensitive classification on data with missing values," *IEEE Trans. Knowl. Data Eng.*, 2006.
- [78] J. Yuan, H. Wang, L. Xiao, W. Zheng, J. Li, F. Lin, and B. Zhang, "A formal study of shot boundary detection," *IEEE Trans. Circuits Syst. Video Techn.*
- [79] M. Zaharia, T. Das, H. Li, T. Hunter, S. Shenker, and I. Stoica, "Discretized streams: fault-tolerant streaming computation at scale," in *SOSP*, 2013.
- [80] H. Zhang, G. Ananthanarayanan, P. Bodík, M. Philipose, P. Bahl, and M. J. Freedman, "Live video analytics at scale with approximation and delay-tolerance," in *NSDI*, 2017.
- [81] Z. Zivkovic, "Improved adaptive gaussian mixture model for background subtraction," in *ICPR*, 2004.