

Blazelt: Fast Exploratory Video Queries using Neural Networks

Daniel Kang, Peter Bailis, Matei Zaharia
Stanford InfoLab

ABSTRACT

As video volumes grow, analysts have increasingly turned to deep learning to process visual data. While these deep networks deliver impressive levels of accuracy, they execute as much as $10\times$ slower than real time (3 fps) on a \$8,000 GPU, which is infeasible at scale. In addition, deploying these methods requires writing complex, imperative code with many low-level libraries (e.g., OpenCV, MXNet), an often ad-hoc and time-consuming process that ignores opportunities for cross-operator optimization. To address the computational and usability challenges of video analytics at scale, we introduce BLAZEIT, a system that optimizes queries over video for spatiotemporal information of objects. BLAZEIT accepts queries via FRAMEQL, a declarative language for exploratory video analytics, that enables video-specific query optimization. We propose new query optimization techniques uniquely suited to video analytics that are not supported by prior work. First, we adapt control variates to video analytics and provide advances in specialization for aggregation queries. Second, we adapt importance-sampling using specialized NNs for cardinality-limited video search (i.e. scrubbing queries). Third, we show how to infer new classes of filters for content-based selection. By combining these optimizations, BLAZEIT can deliver over three order of magnitude speedups over the recent literature on video processing.

PVLDB Reference Format:

Daniel Kang, Peter Bailis, Matei Zaharia. Blazelt: Fast Exploratory Video Queries using Neural Networks. *PVLDB*, 11 (5): xxxx-yyyy, 2018.
DOI: <https://doi.org/TBD>

1 Introduction

Video is rich with semantic information and is a rapidly expanding source of data at scale. For example, London alone has over 500,000 CCTVs [2], and a single autonomous vehicle can generate terabytes of data per day [19]. This growing volume of video can provide answers to queries about the real world. Thus, analysts are increasingly interested in running *exploratory queries* to quickly understand higher-level information over video. For example, an urban planner working on traffic meter setting [61] or urban planning [9] may be interested in whether Mondays have notably different traffic volumes than Tuesdays, and thus counts the number of cars that pass through an intersection. An analyst at an autonomous car company may notice the car behaves strangely at yellow lights, with multiple pedestrians in the crosswalk and searches for

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. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Articles from this volume were invited to present their results at The 45th International Conference on Very Large Data Bases, August 2019, Los Angeles, California.

Proceedings of the VLDB Endowment, Vol. 11, No. 5

Copyright 2018 VLDB Endowment 2150-8097/18/10... \$ 10.00.

DOI: <https://doi.org/TBD>

events of a yellow light, a crosswalk, and at least three pedestrians [20]. However, it is not cost effective and is too time-consuming to manually watch these growing quantities of video, so automated methods of video analysis are increasingly important in answering such queries.

Modern computer vision techniques have made great strides in automating these tasks, with near human-levels of accuracy for some tasks [32]. In particular, a commonly used approach is to perform object detection [16], which returns a sequence of bounding boxes and object class information for each frame of video, over all the frames in a video [56]. This information and simple visual features (e.g., colors) can subsequently be used to answer queries regarding the time and location of objects.

Unfortunately, there are two significant challenges in deploying these vision techniques. First, from a usability perspective, using these methods requires complex, imperative programming across many low-level libraries, such as OpenCV, Caffe2, and Detectron [24]—an often ad-hoc, tedious process that ignores opportunity for cross-operator optimization. Second, from a computational perspective, the naive method of running object detection on every frame of video is infeasible at scale: state-of-the-art object detection (e.g., Mask R-CNN [31]) runs at ~ 3 frames per second (fps), which would take 8 decades of GPU time to process 100 cameras over a month of video.

Prior work has shown that certain video queries can be highly optimized [6, 36, 41]. For example, NOSCOPE [41] and FOCUS [36] optimize the task of binary detection (presence or absence of a target class). While promising, using these pipelines face the same usability challenges as above, requiring interfacing with low-level libraries. Additionally, these pipelines are typically limited in scope (e.g., only binary detection for NOSCOPE).

To address these usability and computational challenges, we present BLAZEIT, a video query system with a declarative query language and three novel optimizations for video analytics queries not supported by prior work. To our knowledge, BLAZEIT is the first system to combine a declarative query language and an optimizer to automatically generate query-specific pipelines for video analytics at scale (as opposed to prior works such as NOSCOPE, which implement fixed-function video processing pipelines).

System Architecture. BLAZEIT consists of two components: 1) a declarative query language called FRAMEQL, and 2) a query optimization and execution engine.

BLAZEIT's first component, called FRAMEQL, is its SQL-like query language, which lets users query the set of visible objects as a relation. Using a SQL-like language has two major benefits. First, as SQL is widely used, BLAZEIT can be quickly adopted by analysts and users. Second, a declarative language enables data independence, separating the specification of the system from the implementation, thus enabling new, video-specific query optimization. In Section 4, we demonstrate how, when combined with the standard relational algebra, FRAMEQL's

schema enables a range of queries for spatiotemporal information of objects in video.

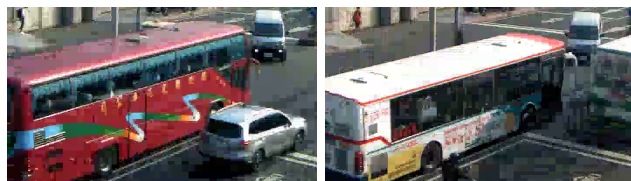
The second component of BLAZEIT is an end-to-end query optimizer and novel optimizations for several common classes of video queries: instead of simply running object detection to populate the rows in an FRAMEQL table, BLAZEIT leverages a range of techniques and a query optimizer to speed up query execution. BLAZEIT focuses on exploratory queries (Section 2), so BLAZEIT lazily populates rows and examines frames only when necessary. We show that using video-specific information enables several forms of query optimization beyond traditional relational databases.

Optimizations. BLAZEIT introduces three novel optimizations for common types video queries, which are not present in existing systems for querying video (Section 11):

Aggregation. First, we study the problem of optimizing aggregate queries, which can be used to provide higher-level statistics over video, (e.g., the average number of cars per hour). As in approximate query processing (AQP) [5, 33], BLAZEIT allows users to specify an error tolerance and therefore samples from the video, as opposed to performing object detection over every frame to compute the exact statistic (e.g., number of cars). However, as object detection remains the computational bottleneck, BLAZEIT further reduces the number of object detection calls by 1) rewriting queries using specialized neural networks (NNs) or 2) adapting the method of control variates [26] to video analytics. Specialized NNs [29, 41] are small NNs trained to predict the output of a larger network for a specific query. Prior work has used specialized NNs for binary classification, and in BLAZEIT we show that they can also learn to provide accurate statistics in some cases. However, in many cases, specialized NNs are not accurate enough. For these cases, we adapt the method of control variates [26], in which specialized NNs are used as a cheap to compute, correlated measure that is used to approximate an expensive statistic (i.e. object detection) to reduce sampling variance. Control variates are only beneficial when the cost of computing the control variate is significantly cheaper than the true statistic (e.g., specialized NNs vs object detection). In traditional relational DBs, the cost of materializing a row is not significantly higher than the cost of processing the row, so control variates do not reduce computation. We show that, when applied to video processing, these techniques can give up to three order-of-magnitude speedups over naive methods (e.g., using NOSCOPE to filter frames with no objects and applying object detection) and up to a $8.7\times$ speedups over AQP.

Scrubbing. Second, we study the problem of optimizing cardinality-limited scrubbing queries [49, 50], in which a fixed number of frames that match a target predicate are returned (e.g., Tesla’s autopilot is known to behave anomalously with lane dividers [45] so an analyst may look for such frames). Searching for these events sequentially or randomly is prohibitively slow when these events are rare (e.g., one event per hour). To address this problem, we adapt importance sampling from the rare-event simulation literature [40] to video analytics. Specifically, BLAZEIT uses the confidence score of a specialized NN as a proxy signal to choose which frames to perform full object detection. Intuitively, this biases the search for requested events towards sections that are more likely to contain them. We demonstrate this biased sampling can deliver up to $500\times$ speedups over naive methods (e.g., using NOSCOPE to filter out frames with no objects and subsequently applying object detection).

Selection. Third, we study the problem of optimizing selection queries, in which users perform content-based filtering of objects (e.g., searching for red tour buses, as in Figure 1). These queries must perform object detection to obtain bounding boxes, which is computationally expensive. Thus, to reduce this computational overhead, BLAZEIT learns a set of conservative filters from the FRAMEQL query to discard irrelevant



(a) Red tour bus.

(b) White transit bus.

Figure 1: Examples of buses in taipei.

frames or parts of frames before applying object detection. We show that BLAZEIT can infer and train for classes of filters: 1) NOSCOPE’s label-based filtering, 2) content-based filtering (using simple visual features), 3) temporal filtering (skipping parts of the video or subsampling frames), and 3) spatial filtering (cropping parts of the video). BLAZEIT inspects the query contents to determine which filters to apply and how to set the filter parameters. For example, if the FRAMEQL query were for red buses comprised of at least 512×512 pixels in the bottom right for at least one second, BLAZEIT can 1) filter frames for a certain redness content, 2) sample at a rate of 0.5s, and 3) crop the video to the bottom right. We demonstrate that these filters can achieve up to a $50\times$ speedup over naive methods.

In summary, we make the following contributions:

1. We introduce FRAMEQL, a query language for spatiotemporal information of objects in videos, and show it can answer a variety of real-world queries (Section 4).
2. We introduce an aggregation algorithm that uses the method of control variates to leverage imprecise specialized NNs for more efficient aggregation than existing AQP methods (Section 6).
3. We introduce a scrubbing algorithm that leverages specialized NNs in importance sampling for finding rare events (Section 7).
4. We introduce a selection algorithm that can infer filters for discarding irrelevant frames can be inferred from FRAMEQL queries and can be applied to content-based selection for up to $50\times$ speedups (Section 8).
5. We demonstrate that these optimizations can give up to three orders-of-magnitude speedups over naively applying object detection or NOSCOPE (Section 10).

2 Use Cases

BLAZEIT focuses on *exploratory queries*: queries that can help a user understand a video quickly, e.g., queries for aggregate statistics (e.g., number of cars) or relatively rare events (e.g., events of many birds at a feeder) in videos. We assume a large amount of archival video (i.e. the batch analytics setting) and that the full object detector has not been run over the whole video, as this would be prohibitively expensive. However, we assume that a small representative sample of the video is annotated with an object detector: this data is used as training data for filters and specialized NNs. We denote this data as the *labeled set* (which can further be split into training data and held-out data). This labeled set can be constructed once, offline, and shared for multiple queries later.

We give several scenarios where BLAZEIT could apply:

Urban planning. Given a set of traffic cameras at various locations, an urban planner performs traffic metering based on the number of cars that pass by different intersections, and determine which days and times are the busiest [61]. The urban planner is interested in how public transit interacts with congestion [14] and look for times with at least one bus and at least five cars. Then, the planner seeks to understand how tourism affects traffic and looks for red buses as a proxy for tour buses (shown in Figure 1).

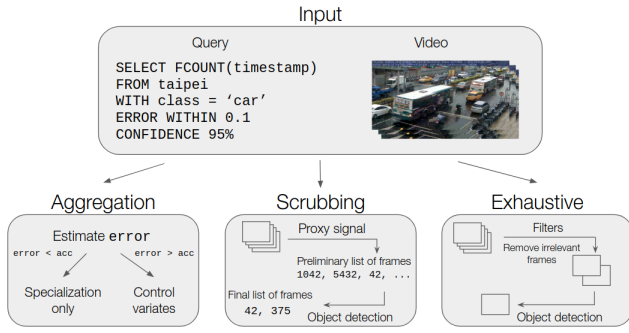


Figure 2: System diagram for BLAZEIT. BLAZEIT accepts FRAMEQL queries and chooses a query plan depending on the type of query.

Autonomous vehicle analysis. An analyst at an autonomous vehicle company notices anomalous behavior of the driving software when the car is in front of a yellow light and there are multiple pedestrians in the crosswalk [20]. Then, the analyst looks for events of a yellow light, a crosswalk, and at least three pedestrians.

Store planning. A retail store owner places a CCTV in the store [59]. The owner segments the video into aisles and counts the number of people that walk through each aisle to understand which products are popular and which ones are not. This information can be used for planning store layout, aisle layout, and product placement.

Ornithology. An ornithologist (a scientist who studies birds) may be interested in understanding bird feeding patterns. The ornithologist might place a webcam in front of a bird feeder [1]. Then, the ornithologist might put different bird feed on the left and right side of the feeder. Finally, the ornithologist can count the number of birds that visit the left and right side of the feeder. As a proxy for species, the ornithologist might then select red or blue birds.

These queries can be answered using spatiotemporal information of objects in the video, along with simple functions over the content of the boxes. Thus, these applications illustrate a need for a unified method of expressing such queries.

3 BLAZEIT System Overview

In this section, we give a brief overview of BLAZEIT’s system architecture: BLAZEIT’s query language, FRAMEQL, its query optimizer, and its execution engine.

FRAMEQL is a SQL-like declarative query language that allows users to express queries for spatiotemporal information of objects in video (further details are given in Section 4). The key challenge we address in this work is efficient execution of FRAMEQL queries: while performing object detection, entity resolution, and UDFs over each frame can answer FRAMEQL queries, this procedure is prohibitively slow. Therefore, we present optimizations for three common classes of queries: aggregation (Section 6), scrubbing (Section 7), and content-based selection (Section 8). Implementation details are given in Section 9. A system diagram of BLAZEIT is shown in Figure 2.

Configuration. As user needs are distinct, BLAZEIT contains several configurable components: the object detection method, the entity resolution method (resolving object across frames), and optional UDFs. While we provide defaults, depending on the users needs, these components can be changed. For example, a license plate reader could be used for resolving the identity of cars. The UDFs can be used to answer more complex queries, such as determining color, filtering by object size or location, or fine-grained classification. UDFs are functions that accept a timestamp, mask, and rectangular set of pixels. As an example, to

compute the “redness” of an object, the UDF could use OpenCV to average the red channel of the pixels.

Filters. Many filters BLAZEIT uses are statistical in nature, so the optimizer must account for their error rates (Section 8). For example, consider a content-based selection for red buses. BLAZEIT will train a specialized NN to filter for frames with buses, but the specialized NN may not be accurate on every frame (in this case, BLAZEIT will call the object detection method on uncertain frames). To account for this error rate, BLAZEIT uses a held-out set of frames to estimate the selectivity and error rate. Given an error budget, BLAZEIT’s query optimizer selects between the filters and uses rule-based optimization to select the fastest query plan. Finally, BLAZEIT can always ensure no false positives by running the most accurate method on the relevant frames.

Specialized neural networks. Throughout, we use specialized NNs as a core primitive [41]. A specialized NN is a neural network that has been trained to mimic a larger NN (e.g., Mask R-CNN on a simplified task, i.e. on a marginal distribution of the larger NN. For example, NOSCOPE’s simplified task (i.e. marginal) is the binary detection task. As the specialized NN is predicting a simpler output, they can run dramatically faster. Prior work has specialized NNs for binary detection [29, 41], but we extend specialization to count and perform multi-class classification. Additionally, we apply various statistical methods over the results of specialized NNs to accurately answer queries such as aggregation or scrubbing.

Bootstrapping filters. To bootstrap and train the filters and specialized NNs, we assume the presence of a labeled set: a small, representative sample that the object detector was run over. Notably, this procedure can be done automatically.

3.1 Limitations

While BLAZEIT can answer significantly more classes of video queries than prior work, we highlight several limitations.

Model Drift. Our current implementation assumes the labeled set for the filters obtained in BLAZEIT’s optimizer is from the same distribution as the remaining video to be analyzed. If the distribution changes dramatically in the new video (e.g., a sunny day vs an extremely foggy day), BLAZEIT will need to re-train the filters. To our knowledge, tracking model drift in visual data has not been well characterized and existing systems such as NOSCOPE [41] do not handle model drift. Thus we view the automatic detection and mitigation of model drift as an exciting area of future research.

Labeled set. Currently, BLAZEIT requires the object detection method to be run over a portion of the data to obtain data for training specialized NNs and filters. We view the problem of warm-starting filters and specialized NNs as an exciting area of future work.

Object detection. BLAZEIT is dependent on the user-defined object detection method and does not support object classes outside what the method returns. For example, the pretrained version of Mask R-CNN [24, 31] can detect cars, but cannot distinguish between sedans and SUVs. However, users could supply UDFs for these queries.

4 FrameQL: A Query Language for Complex Visual Queries over Video

To address the need for a unifying query language over video analytics, we introduce FRAMEQL, a SQL-like language for querying spatiotemporal information of objects in video. We choose a declarative language interface for two reasons. First, encoding queries via a declarative language interface separates the specification and implementation of the system, which

Field	Type	Description
timestamp	float	Time stamp
class	string	Object class (e.g., bus, car, person)
mask	(float, float)*	Polygon containing the object of interest, typically a rectangle
trackid	int	Unique identifier for a continuous time segment when the object is visible
features	float*	The feature vector output by the object detection method.

Table 1: FRAMEQL’s data schema contains spatiotemporal and content information related to objects of interest, as well as metadata (class, identifiers). Each record represents an object appearing in one frame; thus a frame may have many or no records. The features can be used for downstream tasks.

enables query optimization (discussed later). Second, as SQL is the lingua franca of data analytics, FRAMEQL can be easily learned by users familiar with SQL and enables interoperability with relational algebra.

FRAMEQL allows users to query the frame-level contents of a given video feed, specifically the objects appearing in the video over space and time by content and location. FRAMEQL represents videos (stored and possibly compressed in formats such as H.264) as relations, with one relation per video. As in SQL, FRAMEQL allows selection, projection, and aggregation of objects, and, by returning relations, can be composed with standard relational operators. By providing a table-like schema using the standard relational algebra, we enable users with only familiarity with SQL to query videos, whereas implementing these queries manually would require expertise in deep learning, computer vision, and programming.

We show FRAMEQL’s data schema in Table 1. FRAMEQL’s data schema contains fields relating to the time, location, and class of objects, scene and global identifiers, the box contents, and the features from the object detection method (described below). While BLAZEIT provides a default method of populating the schema, the user can specify the object detection method (which will populate `mask`, `class`, and `features`) and the entity resolution method (which will populate `trackid`). For example, an ornithologist may use an object detector that can detect different species of birds, but an autonomous vehicle analyst may not need to detect birds at all. Given these methods, BLAZEIT can automatically populate the data schema. BLAZEIT aims to be as accurate as the configured methods, specifically BLAZEIT does not aim to be more accurate than the configured methods.

Prior visual query engines have proposed similar schemas, *but assume that the schema is already populated* [42,46], i.e. that the data has been created through external means (typically by humans). In contrast, FRAMEQL’s schema can be automatically populated by BLAZEIT. However, as we focus on exploratory queries in this work, FRAMEQL’s schema is *virtual* and rows are only populated as necessary for the query at hand. This is similar to an unmaterialized view. This form of laziness enables a variety of optimizations via query planning.

FRAMEQL’s schema. We briefly describe the fields in FRAMEQL. Recall that a record corresponds to an object that appears in a frame.

- `timestamp` is the timestamp of the object. There is a one-to-one correspondence between timestamps and frames of the video.
- `class` is the object type (e.g., car or bus). The specificity of `class` is determined by the object detection method (e.g., by default Mask R-CNN on MS-COCO does not support sedan vs SUV).
- `mask` is the polygon that encloses the object. We only consider `mask` in the form of bounding boxes, but semantic segmentation [23] could be used for finer grained masks.

Syntactic element	Description
FCOUNT	Frame-averaged count. Equivalent to $COUNT(*) / MAX(timestamp)$. Also equivalent to a time-averaged count.
ERROR WITHIN	Absolute error tolerance
FPR WITHIN	Allowed false positive rate
FNR WITHIN	Allowed false negative rate
CONFIDENCE	Confidence interval
GAP	Minimum distance between returned frames

Table 2: Additional syntactic elements in FRAMEQL. Some of these were taken from BlinkDB.

- `trackid` is a unique identifier for the object as it is visible through a continuous segment in time. If the object exists and re-enters the scene, it will be assigned a new `trackid`.
- `content` is the pixels in as contained by `mask`.
- `features` are the features output by the object detection method. The features can be used for downstream tasks, such as finer-grained classification.

Syntactic sugar. FRAMEQL provides additional syntactic sugar beyond standard SQL as shown in Table 2; several were taken from BlinkDB [5]. We briefly provide the motivation behind each additional piece of syntax.

First, as in BlinkDB [5], users may wish to have fast response to queries and may tolerate some error. Thus, we allow the user to specify error bounds in the form of absolute error, false positive rate, and false negative rate, along with a specified confidence (see below for examples). NOSCOPE’s pipeline can be replicated with FRAMEQL using these constructs. We choose absolute error bounds in this work, as they allow for adaptive sampling procedures (Section 6).

Second, as we provide absolute error bounds, we provide a short-hand for returning a frame-averaged count, which we denote as `FCOUNT`. For example, consider two videos: 1) a 10,000 frame video with one car in every frame, 2) a 10 frame video with a car only in the first frame. Then, `FCOUNT` for the average number of cars per frame would return 1 in the first video and 0.1 in the second video. As videos vary in length, this allows for a normalized way of computing errors. `FCOUNT` can easily be transformed into a time-averaged count.

Finally, when the user selects timestamps, the `GAP` keyword provides a way to ensure that the returned frames are at least `GAP` frames apart. For example, if 10 consecutive frames contains a car and `GAP` = 100, only one frame of the 10 would be returned.

FRAMEQL examples. We first describe how the examples from Section 2 can be written in FRAMEQL. In the following examples, we assume the video is recorded at 30 fps.

Figure 3a shows how to count the average number of cars in a frame. Here, the query uses `FCOUNT` as the error bound are computed per-frame.

Figure 3b shows how to select frames with at least one bus and at least five cars. This query uses the `GAP` keyword to ensure the events found are a certain time apart. As the video is 30 fps, `GAP 300` corresponds to 10 seconds.

Figure 3c shows how to exhaustively select frames with red buses. Here, `redness` and `area` are UDFs, as described in Section 3. Here, 15 frames corresponds to 0.5 seconds.

The other example use-cases can be answered in a similar manner. We give further examples of queries to illustrate FRAMEQL’s syntactic elements.

First, counting the number of distinct cars can be written as:

```

SELECT FCOUNT(*)
FROM taipei
WHERE class = 'car'
ERROR WITHIN 0.1
AT CONFIDENCE 95%

SELECT timestamp
FROM taipei
GROUP BY timestamp
HAVING SUM(class='bus') >=1
AND SUM(class='car') >=5
LIMIT 10 GAP 300

```

(a) The FRAMEQL query for counting the frame-averaged number of cars within a specified error and confidence.

(b) The FRAMEQL query for selecting 10 frames of at least one bus and five cars, with each frame at least 10 seconds apart (at 30 fps, 300 frames corresponds to 10s).

```

SELECT *
FROM taipei
WHERE class = 'bus'
AND redness(content) >= 17.5
AND area(mask) > 100000
GROUP BY trackid
HAVING COUNT(*) > 15

```

(c) The FRAMEQL query for selecting all the information of red buses at least 100,000 pixels large, in the scene for at least 0.5s (at 30 fps, 0.5s is 15 frames). The last constraint is for noise reduction.

Figure 3: Three FRAMEQL example queries.

```

SELECT COUNT (DISTINCT trackid)
FROM taipei
WHERE class = 'car'

```

which is not the same as counting the average number of cars in a frame, as this query looks for distinct instances of cars using `trackid` (cf. Figure 3a).

Second, error rates can be set using syntax similar to BlinkDB [5]:

```

SELECT COUNT(*)
FROM taipei
WHERE class = 'car'
ERROR WITHIN 0.1 CONFIDENCE 95%

```

Third, NOSCOPE can be replicated as FRAMEQL queries of the form:

```

SELECT timestamp
FROM taipei
WHERE class = 'car'
FNR WITHIN 0.01
FPR WITHIN 0.01

```

Finally, a UDF could be used to classify cars:

```

SELECT *
FROM taipei
WHERE class = 'car'
AND classify(content) = 'sedan'

```

5 BLAZEIT Query Execution Overview

When the user issues an FRAMEQL query, BLAZEIT’s query engine optimizes and executes the query. BLAZEIT’s primary challenge is executing the query *efficiently*: naive methods, such as performing object detection on every frame or using NOSCOPE [41] as a filter, are often prohibitively slow. To optimize and execute the query, BLAZEIT inspects the query contents to see if optimizations can be applied. For example, BLAZEIT cannot optimize `SELECT *`, but can optimize aggregation queries with a user-specified error tolerance. Currently, BLAZEIT uses a rule-based optimizer.

Because object detection is the major computational bottleneck, BLAZEIT’s optimizer primarily attempts to reduce the number of object

detection calls while achieving the target accuracy. As object detection methods have increased in accuracy, they have similarly increased in computational complexity. For example, YOLOv2 [56] runs at approximately 80 fps, with an mAP score of 25.4 on MS-COCO [47] (mAP is an object detection metric, with values between 0 and 100, higher being better), but the most accurate version of Mask R-CNN [31] provided by the Detectron framework [24] run at 3 fps with a mAP of 45.2. As a result, object detection is, by far, the most computationally expensive part of a FRAMEQL query; for reference, the specialized NNs we use in this work run at 10,000 fps and some of our simple filters run at 100,000 fps.

BLAZEIT leverages existing techniques from NOSCOPE and three novel optimizations to reduce the computational cost of object detection, targeting aggregation (Section 6), scrubbing (Section 7), and content-based selection (Section 8). As the filters and specialized NNs we consider in these optimizations are cheap compared to the object detection methods, they are almost always worth calling: a filter that runs at 100,000 fps would need to filter 0.003% of the frames to be effective. Thus, we have found a rule-based optimizer to be sufficient in optimizing FRAMEQL queries.

We describe BLAZEIT’s novel optimizations in turn.

6 Optimizing Aggregates

In an aggregation query, the user is interested in some statistic over the data, such as the average number of cars per frame. To exactly answer these queries, BLAZEIT must call object detection on every frame, which is prohibitively slow. However, if the user specifies an error tolerance (Section 4), BLAZEIT can leverage a range of optimizations for accelerated query execution.

When the user issues an aggregation query with an error tolerance, BLAZEIT can efficiently execute the query using a range of techniques: 1) traditional AQP (Section 6.1), 2) query rewriting using specialized NNs (Section 6.2), 3) the method of control variates using specialized NNs (Section 6.3).

The overall procedure to optimize an aggregation query is shown in Algorithm 1.

The first step is determining whether a specialized NN can be trained for the query. Specifically, there must be sufficient training data. In cases where the training data does not contain a sufficient number of examples of interest (e.g., in a video of a street intersection, there are likely to be no examples of bears), BLAZEIT will default to traditional AQP. We present a slightly modified adaptive sampling algorithm that respects the user’s error bound (Section 6.1), which is inspired by Online Aggregation [33] and BlinkDB [5]. Notably, this adaptive sampling algorithm will terminate based on the sample variance, which allows for variance reduction methods (i.e. control variates) to execute faster.

When there is sufficient training data, BLAZEIT will use a specialized NN and estimate its error rate on a held-out set. If the error is within the user-specified error and confidence level, it will then execute the specialized NN on the unseen data and return the answer directly, foregoing the object detection method entirely. As specialized NNs are significantly faster than object detection, this results in much faster execution.

When the specialized NN is not accurate enough, it is used as a control variate: an auxiliary variable that is cheap to compute but highly correlated with the true statistic. We give full details below.

6.1 Sampling

When the query contains a tolerated error rate, BLAZEIT can sample from the video and only populate a small number of rows (or not populate them at all) for dramatically faster execution. Similar to online aggregation [33], we provide absolute error bounds. However, we present a sampling procedure that terminates based on the sampling variance and a CLT bound [48], so that variance reduction methods can terminate early.

Data: Training data, held-out data, unseen video, $uerr \leftarrow$ user’s requested error rate, $conf \leftarrow$ user’s confidence level

Result: Estimate of requested quantity
train specialized NN on training data;
 $err \leftarrow$ specialized NN error rate on held-out data;
 $\tau \leftarrow$ average of specialized NN over unseen video;
if $P(err < uerr) < conf$ **then**

 return τ ;
else
 $\hat{m} \leftarrow$ result of control variates;
 return \hat{m} ;

end

Algorithm 1: BLAZEIT’s procedure to return the results of an aggregate query. This is performed when there are enough examples to train a specialized NN.

Regardless of the confidence level, for an absolute error bound of ϵ , we require at least $\frac{K}{\epsilon}$ samples, where K is the range of the estimated quantity (derived from an ϵ -net argument [30]). For example, if the user queries for the average number of cars per frame, K would be the maximum number of cars over all frames plus one.

Thus, in BLAZEIT’s adaptive sampling procedure, we begin with $\frac{K}{\epsilon}$ samples. At every step, we linearly increase the number of samples. We terminate when the CLT bound gives that the error rate is satisfied at the given confidence level, namely,

$$Q\left(1 - \frac{\delta}{2}\right) \cdot \sigma_N < \epsilon$$

where δ is the confidence interval, σ_N is the sample standard deviation at round N , and Q is the percent point function (i.e. the inverse of the cumulative distribution function) for the normal distribution [33]. We use the finite sample correction to compute the sample standard deviation.

6.2 Specialized Networks for Query Rewriting

In cases where the specialized NN is accurate enough (the accuracy of the specialized NN depends on the noisiness of the video and object detection method), BLAZEIT can return the answer directly from the specialized NN for dramatically faster execution. In this work, we study counting the average number of an object in a frame, which is accomplished using multi-class classification.

To train the specialized NN, BLAZEIT selects the number of classes equal to the highest count that is at least 1% of the video plus one (e.g., if 1% of the video contains 3 cars, BLAZEIT will train a specialized NN with 4 classes, corresponding to 0, 1, 2, and 3 cars in a frame). BLAZEIT uses 150,000 frames for training and uses a standard training procedure for NNs (SGD with momentum [32]) for one epoch.

BLAZEIT estimates the error of the specialized NN on a held-out set using the bootstrap [15]. In this work, we assume no model drift (Section 3.1), thus we assume that the held-out set is representative of the unseen data. If the error is low enough at the given confidence level, BLAZEIT will process the unseen data using the specialized NN and return the result.

6.3 Control Variates

Unfortunately, specialized NNs are not always accurate enough to answer a query on their own. In these cases, BLAZEIT introduces a novel method to take advantage of the specialized NNs while still achieving high accuracy, by combining specialized NNs with AQP-like sampling. In particular, we adopt the method of control variates [26] to video analytics (to our knowledge, control variates have not been applied to database query optimization or video analytics). Specifically, control variates is a method of variance reduction (variance reduction is a

standard technique in Monte Carlo sampling [58] and stochastic optimization [39]) in which specialized NNs are used as a proxy for the statistic of interest. Intuitively, by reducing the variance of sampling, we can reduce the number of frames that have to be sampled and processed by the full object detector.

To formalize this intuition, suppose we wish to estimate the expectation of a quantity m and we have access to an auxiliary variable t . The desiderata for t are that: 1) t is cheaply computable, 2) t is correlated with m . We further assume we can compute $\mathbb{E}[t] = \tau$ and $Var(t)$ exactly. Then,

$$\hat{m} = m + c(t - \tau)$$

is an unbiased estimator of m for any choice of c . Standard analysis [26] shows that the optimal choice of c is

$$c = -\frac{Cov(m, t)}{Var(t)}$$

and using this choice of c gives that

$$\begin{aligned} Var(\hat{m}) &= Var(m) - \frac{Cov(m, t)^2}{Var(t)} \\ &= (1 - Corr(m, t)^2) Var(m). \end{aligned}$$

As an example, suppose $t = m$. Then, $\hat{m} = m + c(m - \mathbb{E}[m]) = \mathbb{E}[m]$ and $Var(\hat{m}) = 0$.

This formulation works for any choice of t , but choices where t is correlated with m give the best results. As we demonstrate in Section 10.2, specialized networks can provide a correlated signal to the ground-truth object detection method for several queries.

As an example, suppose we wish to count the number of cars per frame. Then, m is the random variable denoting the number of cars the object detection method returns. In BLAZEIT, we train a specialized NN to count the number of cars per frame. Ideally, the specialized model would exactly mimic the object detection counts, but this is typically not the case. However, the specialized NNs are typically correlated with the true counts. Thus, the random variable t would be the output of the specialized NN. As our choice of specialized NNs are extremely cheap to compute, we can calculate its mean and variance exactly on all the frames. In BLAZEIT’s adaptive sampling procedure, the covariance is estimated at every round.

7 Optimizing Scrubbing Queries

In cardinality-limited scrubbing queries, the user is typically interested in a rare event, such as a clip of a bus and five cars (if the event is common, the user can simply watch the video). To answer this query, BLAZEIT could run the object detection method over every frame to search for the event. However, if the event occurs infrequently, naive methods of random sampling or sequentially processing the video can be prohibitively slow (e.g., at a frame rate of 30 fps, an event that occurs, on average, once every 30 minutes corresponds to a rate of 1.9×10^{-5}).

Our key intuition is to bias the search towards regions of the video that likely contain the event. To bias the search, we use specialized NNs, and combine them with techniques from the rare-event simulation literature [40]. As an example of rare-event simulation, consider the probability of flipping 80 heads out of 100 coin flips. Using a fair coin, the probability of encountering this event is astronomically low (rate of 5.6×10^{-10}), but using a biased coin with $p = 0.8$ can be orders of magnitude more efficient (rate of 1.2×10^{-4}) [40].

In BLAZEIT, we use specialized NNs to bias which frames to sample. For a given query, BLAZEIT trains a specialized NN to recognize frames that satisfy the query. While we could train a specialized NN as a binary classifier of the frames that satisfy the predicate and that do not, we have found that rare queries have extreme class imbalance. Thus, we train the

specialized NN to incorporate as much information as possible; this procedure has the additional benefit of allowing the trained specialized NN to be reused for other queries such as aggregation. For example, suppose the user wants to find frames with at least one bus and at least five cars. Then, BLAZEIT trains a specialized NN to simultaneously count buses and cars. The signal BLAZEIT uses is the sum of the probability of the frame having at least one bus and at least five cars. BLAZEIT takes the most confident frames until the requested number of frames is found.

7.1 Planning Algorithm

BLAZEIT currently supports scrubbing queries searching for at least N of an object class (e.g., at least one bus and at least five cars). If there are no instances of the query in the training set, BLAZEIT will default to running the object detection method over every frame and applying applicable filters as described in Section 8. If there are instances of the query in the training set, BLAZEIT trains a specialized NN to count instances as above. In the case of multiple object classes, BLAZEIT trains a single NN to detect each object class separately (e.g., instead of jointly detecting “car” and “bus”, the specialized NN would return a separate confidence for “car” and “bus”). As with the counting case, we choose this procedure for class imbalance reasons.

Once the specialized NN is trained, the unseen data is labeled using the specialized NN. BLAZEIT rank-orders the frames by confidence and runs the object detection method over the frames in this order, until the requested number of frames is found. As a result, the frames are not returned in a particular order and may be in a different order compared to a brute-force sequential scan.

8 Optimizing Content-based Selection

In a content-based selection, the user is interested information about the mask or content of every instance of an event, e.g., finding red buses (Figure 3c). In these queries, the object detection method must be called to obtain the mask. As object detection is the overwhelming computational bottleneck, BLAZEIT aims to call object detection as few times as possible.

To achieve this, BLAZEIT infers filters to discard frames irrelevant to the query *before* running object detection on them. BLAZEIT currently supports four classes of filters: 1) label-based filtering, 2) content-based filtering, 3) temporal filtering, and 4) spatial filtering (described in detail below). Importantly, *these filter types and parameters are automatically selected from the query and training data.*

While some filters can be applied with no false positives, others filters are statistical in nature and may have some error rate. The error rate of these filters can be estimated on a held-out set, as in cross-validation [21]. However, as prior work, such as NOSCOPE [41], has considered how to set these error rates, we only consider the case where the filters are set to have no false negatives on the held-out set. Assuming the held-out set is representative of the unseen data (i.e. no model drift, see Section 3.1), this will incur few false negatives on the unseen data.

In BLAZEIT, we present instantiations of each class of filter to demonstrate their effectiveness. We describe each class of filter and BLAZEIT’s instantiations of the filter class.

Label-based filtering. In label-based filtering, the video is filtered based on the desired labels. We leverage similar techniques to NOSCOPE [41] for this type of filter.

Content-based filtering. In content-based filtering, the video is filtered based on fast to compute, low-level visual features, such as average color. If an analyst were to query for “red buses”, we could filter the video to have a certain number of red pixels, or a certain level of red.

For certain classes of filters, BLAZEIT can infer a filter to apply on the whole image from the filter that the user applies on the mask. For example, we define a UDF `redness` that returns a measure of

redness of an image, or portion of an image. In searching for red objects, we can filter frames that are not a certain level of red.

BLAZEIT currently only supports UDFs that return continuous values.

Temporal filtering. In temporal filtering, the video is filtered based on temporal cues. For example, the analyst may want to find buses in the scene for at least K frames. In this case, BLAZEIT subsamples the video at a rate of $\frac{K-1}{2}$. We additionally support basic forms of filtering such as “query the video from 10AM to 11AM.”

Spatial filtering. In spatial filtering, only regions of interest (ROIs) of the scene are considered. For example, a street may have cars parked on the side but the analyst may only be interested in vehicles in transit, so the analyst specifies in the query which parts of the scene contain moving vehicles. The ROI is specified by the user and can be used in smaller models for faster inference, and activity outside the ROI can be ignored, which can increase the selectivity of other filters.

Finally, standard object detectors run faster when the input is more square: in most existing detectors, the input image is resized so that the short-edge is a specific size and the aspect ratio is held constant [31, 57] (for a fixed short-edge size, reducing the long-edge size will make the image smaller). As the computation scales with the resolution, square images result in the least computation. Thus, BLAZEIT makes images more square if the given ROI allows such an operation. For example, if the query only looks for objects with `xmax(mask) < 720` in a 1280×720 video, BLAZEIT will resize the frames to be 720×720 .

8.1 Plan Selection

BLAZEIT will infer which filters can be applied from the user’s query. We describe how each class of filter can be inferred from the query.

First, if the user selects an area of the video, BLAZEIT resizes the frame to be as square as possible, while keeping the area (along with some padding) visible. This resizing is done because object detection methods typically run faster on square images.

Second, BLAZEIT infers the times in the video and the subsampling rate from the query to achieve exact results. For example, if the user queries for objects in the frame at least 30 frames (1 second), BLAZEIT can sample once every 14 frames.

Third, if the user selects a class or set of classes, BLAZEIT trains a specialized NN to detect these classes, as in NoScope. Then, BLAZEIT estimates the threshold on unseen data to ensure no false negatives.

Fourth, if the user selects a UDF over the content (e.g., determining the color of the object), BLAZEIT can apply the UDF over the entire frame (as opposed to the box), and filter frames that do not satisfy the UDF at the frame level. For this procedure to be effective, the UDF must return a continuous value (which can be scaled to a confidence) and return meaningful results at the frame level. Consider two possible UDFs for redness: 1) a UDF which returns true if the over 80% of the pixels have a red-channel value of at least 200 (out of 256), 2) a UDF that returns the average of the red-channel values. In estimating thresholds at the frame-level, BLAZEIT will learn that the first UDF can filter no frames, but the second that filter a large fraction of the frames (based on data from the held-out set). Thus, BLAZEIT can learn which filters can be used effectively as filters. To save computation, we allow users to specify which UDFs will likely be effective, and thus not compute the thresholds for UDFs that will likely not be effective. BLAZEIT sets UDF filter thresholds similar to how it sets thresholds for specialized NNs.

9 Implementation

We implement an open-source¹ BLAZEIT prototype that implements the above query optimizer (currently, the plans are hard-coded; we plan on creating a parser later). We implement our prototype in Python 3.5 as the

¹<https://github.com/stanford-futuredata/blazeit>

deep learning frameworks we use for object detection require Python. To interface with these libraries, we implement the control plane in Python. For efficiency purposes, we implement the non-NN filters in C++. We use PyTorch v0.4 for the training and evaluation of specialized models. For object detection, we use FGFA [64] using MXNet v1.2 and Mask R-CNN [31] using the Detectron framework [24] in Caffe v0.8. We modify the implementations to accept arbitrary parts of video. For FGFA, we use the provided pre-trained weights and for Mask R-CNN, we use the pretrained `X-152-32x8d-FPN-IN5k` weights. We ingest video via OpenCV.

Currently, BLAZEIT uses a rule-based optimizer for query rewriting [54], which supports the queries in Section 4. Most queries follow the following general steps: 1) train specialized neural networks and filters for the query at hand, 2) compute statistics on a held-out dataset to estimate the error or selectivity of the NNs and filters, 3) choose a plan for the unseen data, 4) execute the plan.

We briefly overview different parts of the implementation.

Video ingestion. BLAZEIT initially loads the video using OpenCV, resizes the frames to the appropriate size for each model (65×65 for specialized NNs, short side of 600 pixels for object detection methods), and normalizes the pixel values appropriately. Additionally, we can preprocess the video and directly store the result for faster ingestion.

Specialized NN training. We train the specialized NNs using PyTorch v0.4. Video are ingested and resized to 65×65 pixels and normalized using standard ImageNet normalization [32]. Standard cross-entropy loss is used for training, with a batch size of 16. We used SGD with a momentum of 0.9. Our specialized NNs use a “tiny ResNet” architecture, a modified version of the standard ResNet architecture [32], which has 10 layers and a starting filter size of 16.

Identifying objects across frames. Our default implementation for computing `trackid` use motion IOU [64], but is configurable. Given the set of objects in two consecutive frames, we compute the pairwise IOU of each object in the two frames. We use a cutoff of 0.7 to call an object the same across consecutive frames.

10 Evaluation

We evaluated BLAZEIT on a variety of FRAMEQL queries on real-world video stream in three scenarios: 1) aggregate queries, 2) scrubbing queries for rare events, and 3) accurate, spatiotemporal queries over a variety of object classes. We illustrate that:

1. BLAZEIT achieves up to $4000 \times$ increased throughput compared to a naive baseline, a $2500 \times$ speedup compared to NOSCOPE, and up to a $8.7 \times$ speedup over AQP (Section 10.2).
2. BLAZEIT achieves up to $1000 \times$ speedup compared to a naive baseline and a $500 \times$ speedup compared to NOSCOPE for video scrubbing queries (Section 10.3).
3. BLAZEIT achieves up to $50 \times$ speedup for content-based selection over naive methods by automatically inferring filters to apply before object detection (Section 10.4).

10.1 Experimental Setup

Evaluation queries and videos. We evaluated BLAZEIT on six videos shown in Table 3, which were scraped from YouTube. `taipei`, `night-street`, `amsterdam`, and `archie` were from the same cameras as in NOSCOPE (the other streams were removed from YouTube, so we were unable to use them) and we collected two other streams. We only consider times where the object detection method can perform well (due to lighting conditions), which resulted in 6-11 hours of video per day. These datasets vary in object class (car, bus, boat), occupancy (12% to 90%), and average duration of object appearances (1.4s to 10.7s).

For each webcam, we used three days of video: one day for training labels, one day for threshold computation, and one day for testing (as in [41]).

We evaluated on queries similar to Figure 3, in which the class and video were changed.

Choice of object detection method. We labeled a portion of each video using Mask R-CNN [31], FGFA [64], and YOLOv2 [56], and manually selected the object detection that was the most accurate for each video. As object detection methods have improved since NOSCOPE was published, we did not select YOLOv2 for any of the videos.

Data preprocessing. The literature reports that state-of-the-art object detection methods still suffer in performance for small objects [31, 64], which we have empirically observed even for newer detectors. Thus, we only consider regions where objects are large relative to the size of the frame (these regions are video dependent). Object detectors will return a set of boxes and confidences values. We manually selected confidence thresholds for each video and object class for when to consider an object present, shown in Table 3.

Evaluation metrics. We computed all accuracy measures with respect to the object detection method, in which we treat the object detection method as ground truth. For aggregate statistical queries, we report the absolute error. For scrubbing queries, we guarantee only true positives are returned, thus we only report throughput. Finally, for queries that require detailed information about object (i.e. queries that require performing object detection), all our errors are false negatives, because every frame chosen by our methods is passed to the object detector. Thus, we report the false negative rate for these queries.

In this work, we consider accuracy at the *frame-level*, as we have empirically found that modern object detection methods can return frame-level accurate results. This is in contrast to the one-second binning that is used in [41] to mitigate label flickering in NOSCOPE.

We measure throughput by timing the complete end-to-end system excluding the time taken to decode video, as is standard [41]. We additionally assume the labeled set is computed offline one, so we exclude the time to generate the labeled set (as in [41], we currently use a day of video for training and a day of video for the held-out set). Unlike in [41], we also show runtime numbers *when the training time of the specialized model is included* (excluded in [41]). We include this time as BLAZEIT focuses on exploratory queries, whereas NOSCOPE focuses on long-running streams of data. We additionally show numbers where the training time is excluded, which could be achieved if the specialized NNs were indexed ahead of time.

Hardware Environment. We perform our experiments on a server with an NVIDIA Tesla P100 GPU and two Intel Xeon E5-2690v4 CPUs (56 threads). The system has a total of 504 GB of RAM.

10.1.1 NoScope Baseline Configuration

To our knowledge, NOSCOPE is the closest system to BLAZEIT. NOSCOPE focuses on *binary detection*: the presence or absence of a particular object class. Namely, NOSCOPE cannot directly answer queries in the form of counting or scrubbing for multiple instances of an object or objects.

As NOSCOPE is not directly applicable to the tasks we consider, where relevant, we compare against a NOSCOPE *oracle*, namely a method that returns (on a frame-by-frame basis) whether or not an object class is present in the scene. We assume the oracle is free to query. Thus, this oracle is strictly more powerful—both in terms of accuracy and speed—than NOSCOPE. We describe how the NOSCOPE oracle can be used to answer each type of query.

Aggregates. As NOSCOPE cannot distinguish between one and several objects, whenever NOSCOPE detects an object class is present, it must

Video Name	Object	Occupancy	Average duration	Distinct count	Resol.	FPS	# Eval frames	Length (hrs)	Object detection method	Thresh
taipei	bus	11.9%	2.82s	1749	720p	30	1188k	33	FGFA	0.2
	car	64.4%	1.43s	32367						
night-street	car	28.1%	3.94s	3191	720p	30	973k	27	Mask	0.8
rialto	boat	89.9%	10.7s	5969	720p	30	866k	24	Mask	0.8
grand-canal	boat	57.7%	9.50s	1849	1080p	60	1300k	18	Mask	0.8
amsterdam	car	44.7%	7.88s	3096	720p	30	1188k	33	Mask	0.8
archie	car	51.8%	0.30s	90088	2160p	30	1188k	33	Mask	0.8

Table 3: Video streams and object labels queried in our evaluation. While there are three days of video total for each stream, we show the data from the test set, as the data from the test set will influence the runtime of the baselines and BLAZEIT.

call the object detection method to identify the individual objects. Thus, to count cars in `taipei` would require performing object detection on 64.4% of the frames (i.e. the occupancy rate of cars).

Cardinality-limited scrubbing. As above, NOSCOPE can be used to filter frames that do not contain the objects of interest. For example, if the query were searching for at least one bus and at least five cars in `taipei`, NOSCOPE can be used to remove frames that do not have a bus and a car. Object detection will then be performed on the remaining frames until the requested number of events is found. Thus, for finding rare events, NOSCOPE fares poorly.

Content-based selection. NOSCOPE can only use label-based filtering, but not the other filters classes.

10.2 Aggregate Queries

We evaluate BLAZEIT on six aggregate queries across six videos. The queries are similar to Query 3a (shown in Section 2), with the video and object class changed. We ran five variants of each query:

- Naive: object detection on every frame.
- NOSCOPE oracle: the object detection method on every frame with the object class present.
- Naive AQP: sample from the video.
- BLAZEIT: we use specialized NNs and control variates for efficient sampling.
- BLAZEIT (no train): we exclude the training time from BLAZEIT.

There are two qualitatively different modes in which BLAZEIT executes these queries: 1) where BLAZEIT rewrites the query using a specialized NN, and 2) when BLAZEIT samples using specialized NNs as control variates, and select between these methods as described in Section 6. We analyze these cases separately.

Query rewriting via specialized NNs. We evaluate the runtime and accuracy of specialized networks when the query can be entirely rewritten by running the specialized NN instead. We ran each query with a target error rate of 0.1 and a confidence interval of 95%. We show the average of three runs. The results are shown in Figures 4. The specialized NNs were unable to achieve this accuracy target for `archie`, so we exclude it. However, we show below that specialized NNs can be used as a control variate even in this case.

As shown, the BLAZEIT can achieve up to 8500 \times speedup if the model is cached and a 3200 \times speedup when including the training time and the time to compute thresholds. In contrast, [41] does not include this time in their evaluation. The NOSCOPE oracle baseline does not perform well when the video has many objects of interest (e.g., `rialto`).

In some cases, naive AQP outperform BLAZEIT when BLAZEIT trains the specialized NNs from scratch. However, in all cases, BLAZEIT outperform AQP when the models are cached.

While specialized NNs do not provide a guarantee for error on unseen data, we show that the absolute error stays within the 0.1 for the given videos in Table 4. Thus, we empirically demonstrate that specialized

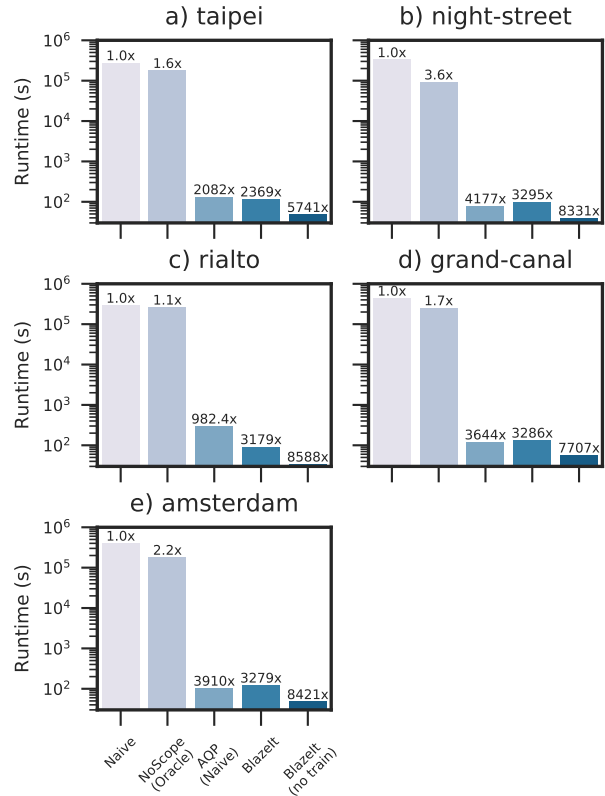


Figure 4: End-to-end runtime of baselines and BLAZEIT on aggregate queries where the query is rewritten with a specialized network, measured in seconds. Note the y-axis is on a log-scale. All queries targeted an error of 0.1.

Video Name	Error
taipei	0.043
night-street	0.022
rialto	-0.031
grand-canal	0.081
amsterdam	0.050

Table 4: Average error over 3 runs of query-rewriting using a specialized NN for counting. These videos stayed within the requested 0.1 error bound.

NNs can be used for query rewriting while respecting the user’s error bounds.

Sampling and control variates. We evaluate the runtime and accuracy

Video Name	Pred (day 1)	Actual (day 1)	Pred (day 2)	Actual (day 2)
taipei	0.86	0.85	1.21	1.17
night-street	0.76	0.84	0.40	0.38
rialto	2.25	2.15	2.34	2.37
grand-canal	0.95	0.99	0.87	0.81

Table 5: Estimated and true counts for specialized NNs run on two different days of video.

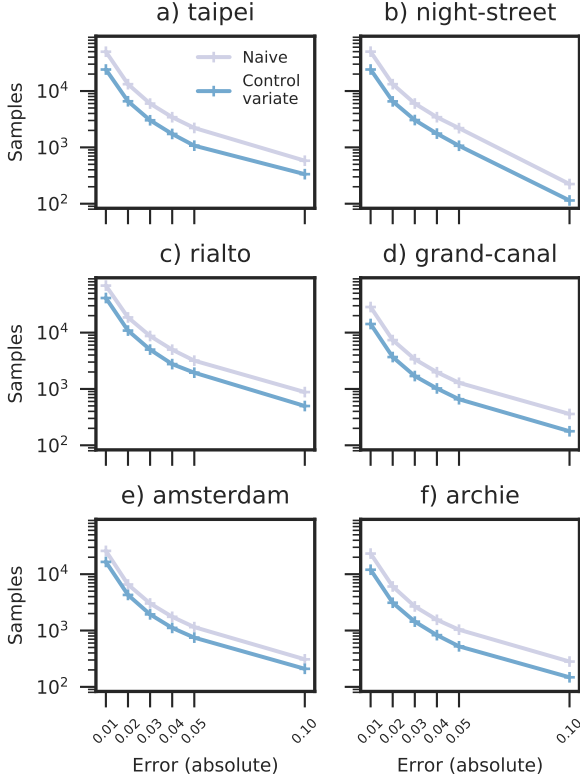


Figure 5: Sample complexity of naive AQP and AQP with control variates. Note the y-axis is on a log-scale.

of using sampling and sampling with control variates. Because of the extreme computational cost of running object detection, we ran the object detection method once and recorded the results. Thus, the run times in this section are estimated from the number of object detection calls.

We targeted error rates of 0.01, 0.02, 0.03, 0.04, 0.05, and 0.1. Each query was run with a confidence interval of 95%. We averaged the number of samples for each error level over 100 runs.

The results are shown in Figure 5. As shown, using specialized NNs as a control variate can deliver up to a $2\times$ reduction in sample complexity. As predicted by theory, the reduction in variance depends on the correlation coefficient between the specialized NNs and the object detection methods. Specifically, as the correlation coefficient increases, the sample complexity decreases.

Specialized NNs do not learn the average. A potential concern of specialized NNs is that they simply learn the average number of cars. To demonstrate that they do not, we swap the day of video for choosing thresholds and testing data. We show the true counts for each day and

Video name	Object	Number	Instances
taipei	car	6	70
night-street	car	5	29
rialto	boat	7	51
grand-canal	boat	5	23
amsterdam	car	4	86
archie	car	4	102

Table 6: Query details and number of instances. We selected rare events with at least 10 instances.

the average of 3 runs in Table 5. Notably, we see that the specialized NNs return different results for each day. This shows that the specialized NNs do not learn the average and return meaningful results.

10.3 Cardinality-limited Scrubbing Queries

We evaluate BLAZEIT on six scrubbing queries, in which frames of interest are returned to the user, up to the requested number of frames. The queries are similar to Query 3b, as shown in Section 2. We show in Table 6 query details and the number of instances of each query. If the user queries more than the maximum number of instances, BLAZEIT must query every frame. Thus, we chose queries with at least 10 instances.

In scrubbing queries, BLAZEIT will only return true positives (as it calls the full object detection method to verify frames of interest), thus we only report the runtime. Additionally, if we suppose that the videos are pre-indexed with the output of the specialized NNs, we can simply query the frames using information from the index. This scenario might occur if, for example, the user executed an aggregate query as above. Thus, we additionally report sample complexity as an objective metric across object detection methods.

We run the following variants:

- Naive: the object detection method is run until the requested number of frames is found.
- NOSCOPE: the object detection method is run over the frames containing the object class(es) of interest until the requested number of frames is found.
- BLAZEIT: specialized NNs are used as a proxy signal to rank the frames (Section 7).
- BLAZEIT (indexed): we assume the specialized NN has been trained and run over the remaining data, as might happen if a user runs queries about some class repeatedly.

Single object class. As shown in Figure 6, BLAZEIT can achieve over a $1000\times$ speedup compared to several baselines. We see that the non-specialized baselines do poorly in finding rare objects, where BLAZEIT’s specialized NNs can serve as a high-fidelity signal for the query at hand.

We additionally varied the number of cars in *taipei* to see if BLAZEIT could also search for common objects. The results are shown in Figure 7. For both the naive method and the NOSCOPE oracle, the same complexity increases as the number of cars increases. However, for up to 5 cars, BLAZEIT’s sample complexity remains nearly constant, which demonstrates the efficacy of biased sampling. While BLAZEIT shows degraded performance with 6 cars, there are only 70 such instances, and is thus significantly harder to find.

Multiple object classes. We additionally test BLAZEIT on multiple object classes by searching for at least one bus and at least five cars in *taipei*. There are 63 instances of such events in the test set.

The end-to-end speedups are shown in Figure 8. Searching for multiple object classes is favorable for the NOSCOPE oracle, as it becomes more selective. Nonetheless, BLAZEIT significantly outperforms the

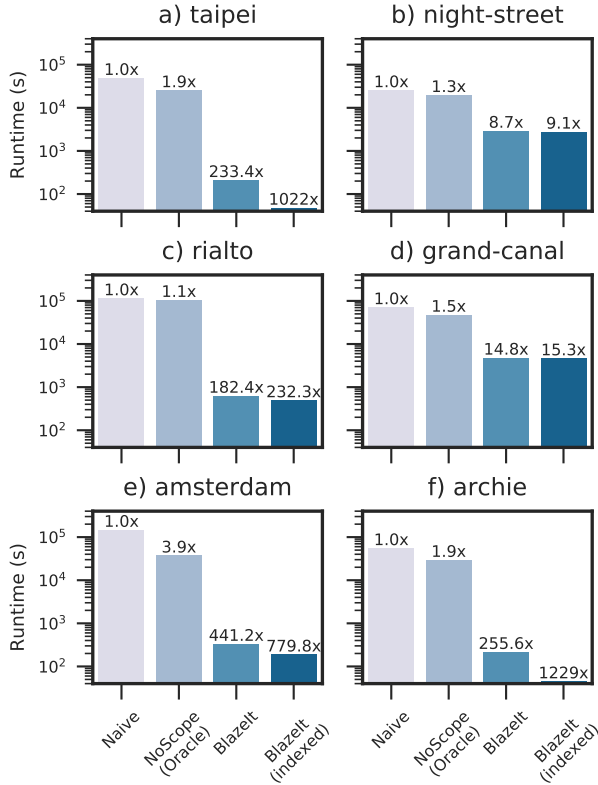


Figure 6: End-to-end runtime of baselines and BLAZEIT on scrubbing queries. Note the y-axis is on a log-scale. All queries looked for 10 events. The average over three runs is shown.

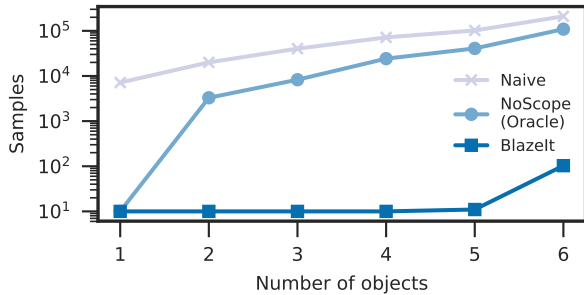


Figure 7: Sample complexity of baselines and BLAZEIT when searching for at least N cars in *taipei*. Note the y-axis is on a log-scale. All queries looked for 10 events.

NO SCOPE oracle, giving up to a $81\times$ performance increase. BLAZEIT also significantly outperforms the naive baseline, giving over a $966\times$ speedup.

Additionally, we show the sample complexity as a function of the LIMIT in Figure 9 of BLAZEIT and the baselines, for *taipei*. We see that BLAZEIT can be up to 5 orders of magnitude more sample efficient over both the naive baseline and NoScope.

10.4 Content-based Selection Queries

To illustrate the effectiveness of content-based filters, we evaluate BLAZEIT on the query shown in Figure 3c.

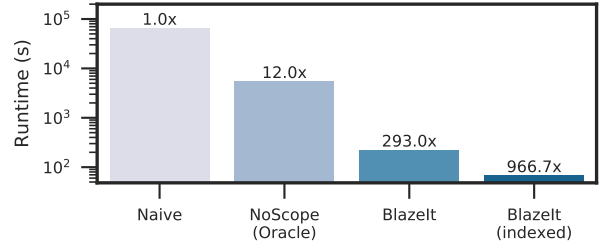


Figure 8: End-to-end runtime of baselines and BLAZEIT on finding at least one bus and at least five cars in *taipei*. Note the y-axis is on a log scale.

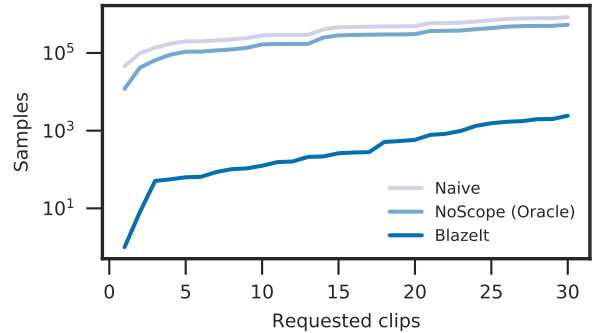


Figure 9: Sample complexity of BLAZEIT, NoScope and the naive method when searching for at least one bus and at least five cars in *taipei*. The x-axis is the number of requested frames. Note the y-axis is on a log scale.

We run the following variants:

- Naive: we run the object detection method on every frame.
- NO SCOPE oracle: we run the object detection method on the frames that contain the object class of interest.
- BLAZEIT: we apply the filters described in Section 8.

For each query, BLAZEIT’s CBO trained, estimated the selectivity, and computed the threshold for each filter applicable to the query (which was determined by BLAZEIT’s rule-based optimizer). We include the time to train the filters and select the thresholds in the runtime. Due to the large computational cost of running the object detector, we extrapolate its cost by multiplying the number of calls by the runtime of the object detector.

End-to-end performance. The results for the end-to-end runtime of the naive baseline, the NoScope oracle, and BLAZEIT are shown in Figure 10. As buses are relatively rare (12% occupancy, see Table 3), NO SCOPE performs well on this query, giving a $8.4\times$ performance improvement over the naive method. However, BLAZEIT outperforms the NO SCOPE oracle by $6.4\times$, due to its extended classes of filters. Furthermore, BLAZEIT delivers up to $54\times$ improved throughput over naive methods for this query.

Factor analysis. We performed a factor analysis and lesion study to understand the impact of each class of filter. In the factor analysis, we added the filters one at a time. In the lesion study, we individually removed the filters.

Results are shown in Figure 11. As shown in the factor analysis, every filter adds a non-trivial speedup. Additionally, removing any

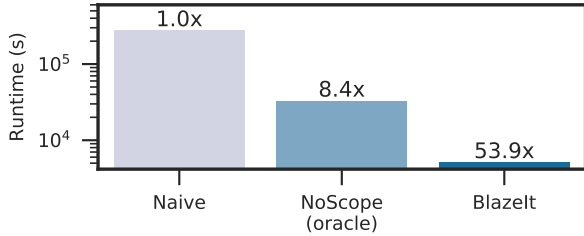


Figure 10: End-to-end throughput of baselines and BLAZEIT on the query in Figure 3c. Note the y-axis is on a log scale.

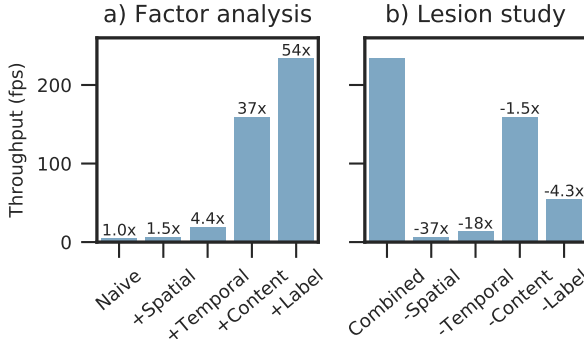


Figure 11: Factor analysis and lesion study of BLAZEIT’s filters on the query in Figure 3c. In the factor analysis, we add filters one at a time. In the lesion study, the filters are individually removed. The filters are described in Section 8.

class of filter reduces performance. Thus, every class of filter improves performance for this query.

11 Related Work

BLAZEIT builds on a long tradition of data management for multimedia and video, and on recent advances in computer vision. We outline some relevant parts of the literature below.

Approximate Query Processing. In AQP systems, the result of a query is returned significantly faster by subsampling the data [18]. Typically, the user specifies an error bound [5], or the error bound is refined over time [33]. Prior work has leveraged various sampling methods [4, 11], histograms [3, 13, 25, 53], and sketches [10, 34, 38].

We leverage ideas from this space and introduce a new form of variance reduction in the form of control variates [26] by using specialized networks. This form of variance reduction, and others involving auxiliary variables, does not make sense in a traditional relational database: the cost of materializing a tuple must be disproportionately large compared to computing the auxiliary variable.

Additionally, we use specialized NNs as a form of importance sampling to bias the search for cardinality-limited scrubbing queries.

Visual Data Management. Visual data management has aimed to organize and query visual data, starting from systems such as Chabot [51] and QBIC [17]. These systems were followed by a range of “multimedia” database for storing [8, 44], querying [7, 43, 52], and managing [22, 37, 62] video data. Many of these systems use classic computer vision techniques such as *low-level* image features (e.g colors, textures) and rely on textual annotations for semantic queries. However, recent

advances in computer vision allow the *automatic* population of semantic data and thus we believe it is critical to reinvestigate these systems.

Modern video analytics. Systems builders have created systems for analyzing video; perhaps the most related is NOSCOPE [41]. NOSCOPE is a highly tuned pipeline for binary detection: it returns the presence or absence of a particular object class in video. Similar to NOSCOPE, other systems, such as FOCUS [36] and TAHOMA [6] have optimized binary detection. However, these systems are inflexible and cannot adapt to user’s queries. Additionally, as NOSCOPE does not focus on the exploratory setting, it does not aim to minimize the training time of specialized NNs. In BLAZEIT, we leverage and extend specialization and present novel optimizations for aggregation, scrubbing, and content-based selection, which these systems do not support.

Other systems, such as VideoStorm [63], aim to reduce latency of live queries that are pre-defined as a *computation graph*. As the computation is specified as a black-box, VideoStorm cannot perform cross-operator optimization. In BLAZEIT, we introduce FRAMEQL and an optimizer that can infer optimizations from the given query. Additionally, we focus on the batch analytics setting. However, we believe BLAZEIT could be run on a system such as VideoStorm for live analytics.

In the batch setting, SCANNER [55] takes a pre-defined computation graph and executes the graph using all the hardware resources available. However, SCANNER does not do automatic cross-operator optimizations or use specialization. We believe BLAZEIT could be run on SCANNER for scale-out.

Speeding up deep networks. We briefly discuss two of the many forms of improving deep network efficiency.

First, a large body of work changes model architecture or weights for improved inference efficiency, that preserve the full generality of these models. Model compression uses a variety of techniques from pruning [28] to compressing [12] weights from the original model, which can be amenable to hardware acceleration [27]. Model distillation uses a large model to train a smaller model [35]. However, these methods aim to retain or nearly retain the accuracy of the original model. These methods do not allow for adapting to the task at hand, as BLAZEIT does. Additionally, these methods are largely orthogonal to BLAZEIT, and reducing the cost of object detection would also improve BLAZEIT’s runtime.

Second, model specialization [41, 60] aims to dramatically improve inference speeds by training a smaller model to mimic the larger model on a *reduced task*. However, specialization has typically been applied in *specific* pipelines, such as NOSCOPE’s binary detection. In BLAZEIT, we leverage and extend specialization to counting and multi-class classification.

12 Conclusions

Querying video for semantic information has become possible with recent advances in computer vision, but these models run as much as 10× slower than real-time. Additionally, deploying these models requires complex programming with low-level libraries. In response, we present a declarative language for video analytics, FRAMEQL, and BLAZEIT, a system that accepts, automatically optimizes, and executes FRAMEQL queries up to three orders of magnitude faster. We demonstrate that FRAMEQL can answer a range of real-world queries, of which we focus on exploratory queries in the form of aggregates and searching for rare events. BLAZEIT introduces new techniques based on AQP, Monte Carlo sampling, and rare-event simulation, and extends specialization to answer these exploratory queries up to three orders of magnitude faster. These results suggest that large video datasets can be explored with orders of magnitude lower computational cost.

Acknowledgments

This research was supported in part by affiliate members and other supporters of the Stanford DAWN project—Google, Intel, Microsoft, NEC, Teradata, and VMware—as well as DARPA under No. FA8750-17-2-0095 (D3M), industrial gifts and support from Toyota Research Institute, Juniper Networks, Keysight Technologies, Hitachi, Facebook, Northrop Grumman, NetApp, and the NSF under grants DGE-1656518 and CNS-1651570.

13 References

- [1] Cornell lab bird cams. <http://cams.allaboutbirds.org/>.
- [2] Cctv: Too many cameras useless, warns surveillance watchdog tony porter, 2015.
- [3] S. Acharya, P. B. Gibbons, and V. Poosala. Aqua: A fast decision support systems using approximate query answers. In *Proceedings of the 25th International Conference on Very Large Data Bases*, pages 754–757. Morgan Kaufmann Publishers Inc., 1999.
- [4] S. Agarwal, H. Milner, A. Kleiner, A. Talwalkar, M. Jordan, S. Madden, B. Mozafari, and I. Stoica. Knowing when you're wrong: building fast and reliable approximate query processing systems. In *Proceedings of the 2014 ACM SIGMOD international conference on Management of data*, pages 481–492. ACM, 2014.
- [5] S. Agarwal, B. Mozafari, A. Panda, H. Milner, S. Madden, and I. Stoica. Blinkdb: queries with bounded errors and bounded response times on very large data. In *Proceedings of the 8th ACM European Conference on Computer Systems*, pages 29–42. ACM, 2013.
- [6] M. R. Anderson, M. Cafarella, T. F. Wenisch, and G. Ros. Predicate optimization for a visual analytics database. *SysML*, 2018.
- [7] W. Aref, M. Hammad, A. C. Catlin, I. Ilyas, T. Ghanem, A. Elmagarmid, and M. Marzouk. Video query processing in the vdbms testbed for video database research. In *Proceedings of the 1st ACM international workshop on Multimedia databases*, pages 25–32. ACM, 2003.
- [8] F. Arman, A. Hsu, and M.-Y. Chiu. Image processing on compressed data for large video databases. In *Proceedings of the first ACM international conference on Multimedia*, pages 267–272. ACM, 1993.
- [9] J. A. Burke, D. Estrin, M. Hansen, A. Parker, N. Ramanathan, S. Reddy, and M. B. Srivastava. Participatory sensing. In *WSW*, 2006.
- [10] M. Charikar, K. Chen, and M. Farach-Colton. Finding frequent items in data streams. In *International Colloquium on Automata, Languages, and Programming*, pages 693–703. Springer, 2002.
- [11] S. Chaudhuri, G. Das, and V. Narasayya. Optimized stratified sampling for approximate query processing. *ACM Transactions on Database Systems (TODS)*, 32(2):9, 2007.
- [12] W. Chen, J. Wilson, S. Tyree, K. Weinberger, and Y. Chen. Compressing neural networks with the hashing trick. In *International Conference on Machine Learning*, pages 2285–2294, 2015.
- [13] G. Cormode, F. Korn, S. Muthukrishnan, and D. Srivastava. Space-and time-efficient deterministic algorithms for biased quantiles over data streams. In *Proceedings of the twenty-fifth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 263–272. ACM, 2006.
- [14] J. De Cea and E. Fernández. Transit assignment for congested public transport systems: an equilibrium model. *Transportation science*, 27(2):133–147, 1993.
- [15] B. Efron and R. J. Tibshirani. *An introduction to the bootstrap*. CRC press, 1994.
- [16] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part-based models. *IEEE transactions on pattern analysis and machine intelligence*, 32(9):1627–1645, 2010.
- [17] M. Flickner, H. Sawhney, W. Niblack, J. Ashley, Q. Huang, B. Dom, M. Gorkani, J. Hafner, D. Lee, D. Petkovic, et al. Query by image and video content: The qbic system. *computer*, 28(9):23–32, 1995.
- [18] M. N. Garofalakis and P. B. Gibbons. Approximate query processing: Taming the terabytes. In *VLDB*, 2001.
- [19] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013.
- [20] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 3354–3361. IEEE, 2012.
- [21] S. Geisser. *Predictive inference*. Routledge, 2017.
- [22] S. Gibbs, C. Breiteneder, and D. Tsichritzis. Audio/video databases: An object-oriented approach. In *Data Engineering, 1993. Proceedings. Ninth International Conference on*, pages 381–390. IEEE, 1993.
- [23] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 580–587, 2014.
- [24] R. Girshick, I. Radosavovic, G. Gkioxari, P. Dollár, and K. He. Detectron. <https://github.com/facebookresearch/detectron>, 2018.
- [25] M. Greenwald and S. Khanna. Space-efficient online computation of quantile summaries. In *ACM SIGMOD Record*, volume 30, pages 58–66. ACM, 2001.
- [26] J. M. Hammersley and D. C. Handscomb. General principles of the monte carlo method. In *Monte Carlo Methods*, pages 50–75. Springer, 1964.
- [27] S. Han, X. Liu, H. Mao, J. Pu, A. Pedram, M. A. Horowitz, and W. J. Dally. Eie: efficient inference engine on compressed deep neural network. In *Computer Architecture (ISCA), 2016 ACM/IEEE 43rd Annual International Symposium on*, pages 243–254. IEEE, 2016.
- [28] S. Han, H. Mao, and W. J. Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding. *arXiv preprint arXiv:1510.00149*, 2015.
- [29] 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 *Proceedings of the 14th Annual International Conference on Mobile Systems, Applications, and Services*, pages 123–136. ACM, 2016.
- [30] D. Haussler and E. Welzl. -nets and simplex range queries. *Discrete & Computational Geometry*, 2(2):127–151, 1987.
- [31] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, pages 2980–2988. IEEE, 2017.
- [32] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [33] J. M. Hellerstein, P. J. Haas, and H. J. Wang. Online aggregation. In *Acm Sigmod Record*, volume 26, pages 171–182. ACM, 1997.
- [34] M. R. Henzinger, P. Raghavan, and S. Rajagopalan. Computing on data streams. *External memory algorithms*, 50:107–118, 1998.

- [35] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [36] K. Hsieh, G. Ananthanarayanan, P. Bodik, P. Bahl, M. Philipose, P. B. Gibbons, and O. Mutlu. Focus: Querying large video datasets with low latency and low cost. *arXiv preprint arXiv:1801.03493*, 2018.
- [37] R. Jain and A. Hampapur. Metadata in video databases. *ACM Sigmod Record*, 23(4):27–33, 1994.
- [38] C. Jin, W. Qian, C. Sha, J. X. Yu, and A. Zhou. Dynamically maintaining frequent items over a data stream. In *Proceedings of the twelfth international conference on Information and knowledge management*, pages 287–294. ACM, 2003.
- [39] R. Johnson and T. Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in neural information processing systems*, pages 315–323, 2013.
- [40] S. Juneja and P. Shahabuddin. Rare-event simulation techniques: an introduction and recent advances. *Handbooks in operations research and management science*, 13:291–350, 2006.
- [41] D. Kang, J. Emmons, F. Abuzaid, P. Bailis, and M. Zaharia. Noscope: optimizing neural network queries over video at scale. *Proceedings of the VLDB Endowment*, 10(11):1586–1597, 2017.
- [42] T. C. Kuo and A. L. Chen. A content-based query language for video databases. In *Multimedia Computing and Systems, 1996., Proceedings of the Third IEEE International Conference on*, pages 209–214. IEEE, 1996.
- [43] M. La Cascia and E. Ardizzone. Jacob: Just a content-based query system for video databases. In *Acoustics, Speech, and Signal Processing, 1996. ICASSP-96. Conference Proceedings., 1996 IEEE International Conference on*, volume 2, pages 1216–1219. IEEE, 1996.
- [44] J. Lee, J. Oh, and S. Hwang. Strg-index: Spatio-temporal region graph indexing for large video databases. In *Proceedings of the 2005 ACM SIGMOD international conference on Management of data*, pages 718–729. ACM, 2005.
- [45] T. Lee. Theres growing evidence teslas autopilot handles lane dividers poorly, 2018.
- [46] J. Z. Li, M. T. Ozsu, D. Szafron, and V. Oria. Moql: A multimedia object query language. In *Proceedings of the 3rd International Workshop on Multimedia Information Systems*, pages 19–28, 1997.
- [47] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.
- [48] S. Lohr. *Sampling: design and analysis*. Nelson Education, 2009.
- [49] J. Matejka, T. Grossman, and G. Fitzmaurice. Swift: reducing the effects of latency in online video scrubbing. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 637–646. ACM, 2012.
- [50] J. Matejka, T. Grossman, and G. Fitzmaurice. Swifter: improved online video scrubbing. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 1159–1168. ACM, 2013.
- [51] V. E. Ogle and M. Stonebraker. Chabot: Retrieval from a relational database of images. *Computer*, 28(9):40–48, 1995.
- [52] J. Oh and K. A. Hua. Efficient and cost-effective techniques for browsing and indexing large video databases. In *ACM SIGMOD Record*, volume 29, pages 415–426. ACM, 2000.
- [53] G. Piatetsky-Shapiro and C. Connell. Accurate estimation of the number of tuples satisfying a condition. *ACM Sigmod Record*, 14(2):256–276, 1984.
- [54] H. Pirahesh, J. M. Hellerstein, and W. Hasan. Extensible/rule based query rewrite optimization in starburst. In *ACM Sigmod Record*, volume 21, pages 39–48. ACM, 1992.
- [55] A. Poms, W. Crichton, P. Hanrahan, and K. Fatahalian. Scanner: Efficient video analysis at scale (to appear). 2018.
- [56] J. Redmon and A. Farhadi. Yolo9000: better, faster, stronger. *arXiv preprint*, 2017.
- [57] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015.
- [58] C. P. Robert. *Monte carlo methods*. Wiley Online Library, 2004.
- [59] A. W. Senior, L. Brown, A. Hampapur, C.-F. Shu, Y. Zhai, R. S. Feris, Y.-L. Tian, S. Borger, and C. Carlson. Video analytics for retail. In *Advanced Video and Signal Based Surveillance, 2007. AVSS 2007. IEEE Conference on*, pages 423–428. IEEE, 2007.
- [60] H. Shen, S. Han, M. Philipose, and A. Krishnamurthy. Fast video classification via adaptive cascading of deep models. *arXiv preprint*, 2016.
- [61] X. Sun, L. Muñoz, and R. Horowitz. Highway traffic state estimation using improved mixture kalman filters for effective ramp metering control. In *Decision and Control, 2003. Proceedings. 42nd IEEE Conference on*, volume 6, pages 6333–6338. IEEE, 2003.
- [62] A. Yoshitaka and T. Ichikawa. A survey on content-based retrieval for multimedia databases. *IEEE Transactions on Knowledge and Data Engineering*, 11(1):81–93, 1999.
- [63] H. Zhang, G. Ananthanarayanan, P. Bodik, M. Philipose, P. Bahl, and M. J. Freedman. Live video analytics at scale with approximation and delay-tolerance. In *NSDI*, volume 9, page 1, 2017.
- [64] X. Zhu, Y. Wang, J. Dai, L. Yuan, and Y. Wei. Flow-guided feature aggregation for video object detection. *arXiv preprint arXiv:1703.10025*, 2017.