Kestrel: Using Mobile Devices to Augment Multi-camera Vehicle Tracking

Anonymous Author(s)

ABSTRACT

As mobile devices have become more powerful and GPU-equipped, vision-based applications are becoming feasible on these devices. Mobile devices can augment fixed camera surveillance systems to improve coverage and accuracy, but it is unclear how to leverage these mobile devices while respecting their processing, communication, and energy constraints. As a first step towards understanding this question, this paper discusses Kestrel, a system that tracks vehicles across a hybrid camera network. In Kestrel, fixed camera feeds are processed on the cloud, and mobile devices are invoked to resolve ambiguities in vehicle tracks. Kestrel's mobile pipeline detects objects using a deep neural network, extracts attributes using cheap visual features, and resolves path ambiguities by careful association of vehicle visual descriptors, while using several optimizations to conserve energy and reduce latency. We evaluate Kestrel on a heterogenous dataset including both mobile and static cameras from a campus surveillance network. The results demonstrate that Kestrel on a hybrid network can achieve precision and recall comparable to a fixed camera network of the same size and topology.

1 INTRODUCTION

Over the past five years, mobile device cameras have improved to the point where they can capture high-quality full-motion videos. Until recently, however, automated detection and classification has been beyond the capabilities of these devices. Moreover, given the bandwidth constraints of cellular and WiFi connections, it has also not been feasible to upload videos for processing in the cloud, so it has not been possible to exploit the large corpus of videos captured by mobile device users.

With the advent of GPUs on mobile devices, complex vision algorithms can now be executed on mobile devices. But a single camera on a single mobile device provides only a limited perspective, and most applications will require fusing the results of vision algorithms across multiple cameras, continuously and in near real-time. Such applications can stress the computational and network capabilities of mobile devices, and significantly impact their energy consumption.

To understand what role mobile devices can play in such applications, and how to architect these applications, we consider the challenging problem of *tracking an object through a network of heterogeneous cameras*. Vehicle tracking arises in the context of surveillance. Many large enterprises, such as university campuses, employ a network of fixed cameras to visually inspect vehicles (and people) traversing their streets. This camera network is used both for real-time surveillance, and for forensic purposes (to determine, after the fact, suspects or suspicious vehicles). With the prevalence of mobile devices, it is natural to consider *augmenting* these fixed cameras with videos captured by mobile devices, such as those carried by security officers or neighborhood watch personnel. Mobile devices are now in widespread use in police and security work [44]. In particular, mobile video surveillance is on the rise: cities like Chicago [27], and New York [53] have already equipped police officers with body cameras, and police cruisers [17] also have dashcams for recording traffic stops.

The Vehicle Tracking Problem and Challenges. The problem we consider in this paper is *to automatically detect a path taken by a vehicle through this hybrid fixed/mobile camera network.* The vehicle tracking problem stresses mobile device capabilities because it requires associating vehicles across multiple cameras, which itself requires significant network, processing and energy resources.

In our problem setting, each mobile or fixed camera either continuously or intermittently captures video, together with associated metadata (camera location, orientation, resolution, *etc.*). Conceptually, this creates a large corpus of videos over which users may pose several kinds of queries (either in near real-time, or after the fact), the most general of which is: What path did a given car, seen at a given camera (or at k distinct cameras), take through the camera network?

One approach to designing a system that can support such queries is to centrally collect this corpus in the cloud. This approach works well for fixed cameras, which are engineered with significant network bandwidth to the (private) cloud. Indeed, our approach assumes that fixed cameras in our hybrid network can continuously upload videos to the cloud.

However, for mobile devices with bandwidth-constrained network interfaces, such an approach is impractical. Our paper seeks to understand how to architect mobile device participation in this hybrid camera network. Specifically, in this network, resources are asymmetric, with the cloud having resource elasticity but with mobile devices being constrained along many dimensions. Given this asymmetry, we seek to design an architecture that must address three challenges. First, it must trigger video processing at a mobile device only when absolutely necessary, so that mobile devices are not required to continuously process every frame of the video in order to track vehicles. Second, it must use low-complexity descriptors for vehicles and their trajectories that, while respecting mobile device computational constraints, can also minimize bandwidth consumed without significantly impacting accuracy. Third, it must permit robust association of vehicles across different cameras using vehicle descriptors and their trajectory, even when cameras can have different perspectives (fixed cameras on light poles, mobile device cameras at the street level, with all cameras pointing in potentially different directions) and can be subject to different lighting conditions.

Existing commercial surveillance systems do not support automated vehicle tracking across multiple cameras, nor do they support a hybrid camera network. They either require a centralized collection

SenSys, Delft, The Netherlands 2016, 978-x-xxxx-x/YY/MM...\$15.00 DOI: 10.1145/nnnnnn.nnnnnn

of videos [2][3], or perform specific object detection on an individual camera, leaving it to a human operator to perform association by manual inspection [1].

Contributions. This paper describes the design of Kestrel (§2) for addressing the vehicle tracking problem. Users of Kestrel provide an image of a vehicle (captured, for example, by the user inspecting video from a static camera), and the system returns the sequence of cameras (i.e., the path through the camera network) at which this vehicle was seen. This system carefully selects, and appropriately adapts, the right combination of vision algorithms to enable accurate end-to-end tracking, even when individual components can have less than perfect accuracy, while respecting mobile device constraints. Kestrel addresses the challenge described above using one key architectural idea. Its cloud pipeline continuously processes videos streamed from fixed cameras to extract vehicles and their attributes, and when a query is posed, computes vehicle trajectories across these fixed cameras. When there is ambiguity in these trajectories, Kestrel queries one or more mobile device (which runs a *mobile pipeline* optimized for video processing on mobile devices). Thus, mobile devices are invoked only when absolutely necessary, significantly reducing resource consumption.

Our first contribution is techniques for fast execution of deep Convolutional Neural Networks (CNNs) on mobile device embedded GPUs (§3.1). The context for this contribution is that object detection and localization based on deep CNNs (those with 20 or more layers) has become significantly more accurate than earlier generations of detection algorithms. However, mobile GPUs cannot execute these deep CNNs because of insufficient memory. By quantifying the memory consumption of each layer, we have designed a novel approach that offloads the bottleneck layers to the mobile device's CPU and pipelines frame processing without impacting the accuracy. Kestrel leverages these optimizations as an essential building block to run a deep CNN (YOLO [57]) on mobile GPUs for *detecting* objects and drawing bounding boxes around them on a frame.

Our second contribution is the design of an accurate and efficient object *tracker* for objects within a single video (§3.2). This tracker runs on the mobile device, and we use it to extract object attributes, like the direction of motion and low-complexity object descriptors. Existing general purpose trackers (§6) are computationally expensive, or can only track single objects. Kestrel leverages the fact that object detection is relatively cheap with mobile GPUs and periodically detects objects (thereby reducing the energy cost of object detection), say every k video frames, so that the tracking algorithm has to only be accurate in between frames at which detection is applied. This enables even simpler tracking techniques, based on keypoint tracking, to be very effective in our setting, and an order of magnitude more efficient.

Our third contribution is accurate vehicle *association* (§4.1). Given an object in a camera's video, the association can determine which object in a second camera best matches the given object. Association is performed both on the cloud and in mobile devices, and uses the attributes determined from the tracking algorithms. Kestrel achieves highly accurate, low-complexity association using three simple ways to winnow candidate matching objects: the travel time between cameras, the direction of motion of the object in the first camera, and the color distribution of the vehicle. Association is used



Figure 1-Kestrel System Architecture

by a *path inference* algorithm that uses a dynamic programming approach to find the most likely *paths* through the camera network that a vehicle may have traversed.

The evaluation (\$5) of Kestrel uses a dataset of videos collected from a heterogeneous camera network. Specifically, our dataset consists of a total of nearly 4 hours of video footage collected on a university campus comprising of 11 static and 6 mobile cameras. In this dataset, we have manually annotated a ground truth dataset of ~120 cars. Kestrel is able to achieve over 90% precision and recall for the path inference problem across multiple hops in the camera network, with minimal degradation as the number of hops increases, even though, with each hop, the number of potential candidates increases exponentially. Its overall performance on a hybrid network is comparable to an identical fixed camera network: the use of mobile devices only minimally impacts performance. Finally, the association has nearly 97% precision and recall with each winnowing approach contributing to the association performance.

2 KESTREL ARCHITECTURE

To address the challenges discussed above, Kestrel leverages several techniques. First, mobile GPUs can run deep neural nets trained for object detection, but these GPUs are energy hungry and fundamentally constrained with respect to their desktop counterparts. So, Kestrel features a novel vehicle detection and localization pipeline to invoke the neural nets on mobile devices only when necessary, and uses novel compute offload techniques to achieve low latency and energy. Next, Kestrel uses a fast tracking algorithm that can track and associate vehicles across successive frames in the same camera, while compensating for mobile camera movement. This tracking helps determine the vehicle's trajectory, which, together with other attributes for a vehicle (such as its color, or its direction of motion), can help associate (or re-identify) cars across multiple nonoverlapping cameras with different perspectives. Finally, Kestrel exploits a novel dynamic programming estimation framework to determine the most likely path of the vehicle across the hybrid camera network. Our novelty lies in determining the right set of object attributes (descriptors, trajectories, etc.) that can most effectively disambiguate between multiple possibilities while incurring minimal energy consumption and latency. This path inference framework runs on the cloud and provides the basis for identifying paths taken by specific vehicles, and for determining when to query a mobile device to increase confidence in the result.

Figure 1 shows the system architecture of Kestrel, which consists of a *mobile device pipeline* and a *cloud pipeline*. Videos from fixed

cameras are streamed and stored on a cloud. Mobile devices *store* videos locally and only send metadata to the cloud specifying when and where the video is taken. Given a vehicle to track through the camera network, the cloud performs object detection, tracking, attribute extraction, cross-camera association, and path inference only on videos from the fixed cameras. When the cloud determines that it has low confidence in the result, it uses metadata from the mobile devices to query relevant mobile devices to increase tracking accuracy. Once the mobile device receives the query, it executes a pipeline, optimized for energy and computation, to augment the tracks determined by the fixed cameras. In the following sections, we describe Kestrel's mobile and cloud components.

3 KESTREL MOBILE DEVICE PIPELINE

Kestrel's cloud component, discussed in §4, produces a rank-ordered set of vehicle paths in response to a query. These vehicle paths represent the cloud's best estimates of when, and which cameras the query vehicle traversed (we define paths more formally later). The cloud can determine which mobile devices could potentially have images of the query vehicle; the cloud can do this using metadata from mobile devices, which describe the locations and times at which each video was taken. To each such candidate mobile device, the cloud sends attributes of the queried vehicle (discussed below), a time window when the vehicle is likely to have been seen at the mobile device, and the set of ranked paths. Using this information, the mobile device re-ranks the vehicle paths, and returns this to the cloud, which can then iteratively refine the ranking by querying other mobile devices. This architectural choice, where the cloud only invokes a mobile device if it could potentially refine the result, is crucial for the feasibility of Kestrel on hybrid camera networks.

The primary challenge in Kestrel's mobile pipeline (Kestrel-Mobile, Figure 1) is to optimize energy, computation, and communication overhead. Kestrel-Mobile uses modern CNNs on mobile GPUs to detect vehicles in frames, but because running CNNs on GPUs can be both power hungry and time consuming, Kestrel-Mobile uses the time window to narrow down the set of frames to detect objects on, and a computationally efficient optical flow filter to determine which of these frames to run the object detector on. Moreover, Kestrel-Mobile features a lightweight state-of-the-art multi-object tracker to reduce object detector invocation times and further limit latency and energy consumption, as well as to extract crucial attributes for instance association (§3.2). Also, Kestrel-Mobile caches all instances and attributes extracted for reuse across subsequent queries. Using these components, the mobile pipeline associates cars seen in its videos with cars seen in vehicle paths received from the cloud, then re-ranks those paths, and returns them to the cloud. In this section, we discuss most of these components, and defer the association and re-ranking until §4.

3.1 Object Detection: Deep CNNs

Kestrel-Mobile uses GPUs on mobile devices to run Convolutional Neural Networks (CNNs) for vehicle detection. Driven by AR/VR applications, GPUs have started appearing in commercial off-theshelf mobile processors. An example mobile processor with a GPU is the Tegra K1, manufactured by nVidia, which contains 192 GPU cores and conforms to the Compute Unified Device Architecture (CUDA) specification. The Tegra K1 has already been incorporated into tablet-class devices such as the Nexus 9 [10] and the Google Tango [7]. Its successor, nVidia Tegra X1 [12], has been incorporated into Google's Pixel C [8]. Beyond tablets, mobile phones such as Lenovo Phab 2 Pro [42] and Asus ZenFone AR [18] are equipped with the latest Qualcomm Snapdragon 652 and 820 chips. In this paper, we use, for our evaluations, the Jetson TK1 board from nVidia [11], which contains the Tegra K1 mobile processor and an associated CUDA-capable development environment. Several benchmarks [14, 15] have already demonstrated that TK1's performance is comparable to that of other GPUs like the TX1 or the 820.

Traditional GPU architectures separate GPU and CPU memory, but starting from CUDA 6, platforms (like the TK1) have unified memory, so the CPU and GPU can share memory, which can result in programming simplicity [16]. However, relative to the run-time memory requirements of deep convolutional neural networks, devices like the TK1 have limited memory (2GB on the TK1).

YOLO. Kestrel-Mobile uses YOLO[57], a deep CNN for performing object detection. As shown in Figure 2, YOLO not only classifies the object but also draws a bounding box around it. YOLO is structured as a 27 layer CNN, with 24 convolutional layers, followed by 2 fully-connected (FC) layers and a detection layer. The convolutional layers extract the features that best suits the vision task, while the FC layers use these features to predict the output probabilities and the bounding box coordinates.

YOLO requires 2.8GB memory on TK1 using the regular GPU programming workflow described below (after removing variables that are only required in the training phase). Almost 80% of the memory allocated to the network parameters is taken by the first FC layer (also observed by other works [36]).

A normal GPU programming workflow involves loading the data to be operated upon (in our case the trained CNN weights) first in CPU RAM, then copying them to GPU memory and starting the kernel operation. If we had enough resources to hold the entire CNN model in memory, then multiple video frames can be dispatched for processing in sequence and the computation is very fast. However, when we ran YOLO on the TK1 using the workflow described above, it *failed to execute* due to insufficient memory resources. This motivated us to explore several methods to manage the memory constraint on mobile GPUs.

The CPU offload Optimization. Prior work has considered offloading computation to overcome processor limitations [25, 26]. We explore offloading only the FC layer computation to the CPU, not motivated by processor speed considerations, but by the fact that, because the CPU supports virtual memory it can more effectively manage memory over-subscription required for the FC layer. With CPU offloading, we observe that when the CPU is running the FC layer on the first video frame, the GPU cores are idle. So we adopt a *pipelining* strategy to start running the second video frame on the GPU cores rather than letting them idle. To achieve this kind of pipelining, we run the FC layers on the CPU on a separate thread, as GPU computation is managed by the main thread.

Other memory optimizations. As mentioned above, the FC layers are the ones that require most memory in a CNN. Thus we also explore *reducing the size* of the FC layer (other work [20, 22] has

Anon.



Figure 2—YOLO Detection





Figure 4-Camera Movement Subtraction

explored similar techniques in different settings), which can potentially reduce detection accuracy. To reduce the size of the FC layer, we reduce the number of filters and the number of outputs in the network configuration and re-train.

To reduce the memory footprint of the FC layer, we can also *split the computation* in some layers [51]. Specifically, in the FC layer, the input vector of previous layer is multiplied by the weight matrix to obtain the output vector. Splitting this matrix multiplication and reading the weight matrix in parts allows us to reuse memory. However, re-using memory and overwriting the used chunks incurs additional file access overheads (§5.5).

Speed optimizations. In addition to memory optimizations, our results incorporate two optimizations for speeding up the computation. First, we set the GPU clock-rates to the highest supported to ensure best performance. This approach trades off power for increased computation speed. Second, a few CUDA operations (matrix operation, random number state initialization) incur significant latency when they are invoked for the first time. We perform these operations before running the CNN to reduce latency.

Generality. Our optimizations are applicable to other complex vision tasks using CNNs as well. For example, FCN [46] (for semantic segmentation) requires 4.8GB memory, VGG 19 and 16 [61] (for object recognition) both require about 1.3 GB memory to load the base neural network. Although newer mobile GPU hardware (*e.g.*, TX1 [12]) has more memory, mobile GPU memory increases are limited by energy budgets: as [45] points out, DRAM power can account for up to 30% of the total power consumption of smartphones. Further, designing unified virtual memory across heterogeneous processors is still an active area of research [54, 55]. Virtual memories, however, incur performance overheads (*e.g.*, TLB and cache validation, memory coherence and address resolution). For these reasons, application-specific memory management techniques for deep CNNs will continue to be useful.

3.2 Attribute Extraction

After detecting objects, Kestrel-Mobile's pipeline extracts several *attributes* of the object from a video on a mobile device, including its (i) direction of motion, and (ii) a low-complexity visual descriptor. These attributes are used to associate objects across cameras (§4.1). To estimate the direction of motion, it is important to *track* the object across successive video frames.

Light-weight Multi-Object Tracking. To be able to track objects¹ in a video captured at a single camera, Kestrel-Mobile needs to take multiple objects detected in the individual frames above and associate them across multiple frames of the video. Our tracking algorithm has two functions: (i) it can reduce the number of times YOLO is invoked, which in turn reduces the latency, and (ii) it can be used to estimate direction of motion. The state-of-the-art object tracking techniques [50, 52, 66] take as input the pre-computed bounding box of the object, then extract sophisticated features like SIFT [47], SURF [19], BRISK [43], *etc.* from the bounding box, and then iteratively infer the position of these key-points in subsequent frames using key point matching algorithms like RANSAC [29].

However, we observe that designing the most robust tracker for a single object is different from effectively tracking all the objects that appear within a given time window in a video. This is because in dynamic mobile camera scenarios, objects (especially fast moving vehicles) can enter and exit the scene frequently. This means that Kestrel-Mobile needs to run object detection algorithms like YOLO frequently in order not to miss out on tracking new objects. To avoid this, Kestrel-Mobile uses two optimizations. First, it performs optical flow 2 based scene change detection and whenever a scene change is detected, it runs YOLO to detect new objects. Second, to avoid running YOLO on every subsequent frame, we run YOLO every kframes (for small k) and stitch the detected objects together using a tracker that tracks light-weight features such as Good Features To Track (GFTT [60]). In §5, we compare our tracker against other state-of-the-art trackers to quantify the latency vs. accuracy tradeoffs between the two approaches.

Specifically, given a video stream, Kestrel-Mobile's tracker detects *keypoints* of each frame and uses a keypoint tracker, Kanade-Lucas-Tomasi (KLT [48]), to track keypoints across frames. For every k frames (called YOLO frames), Kestrel-Mobile runs YOLO to detect the objects of interest in the frame. Each detected object is stored in Kestrel-Mobile with a series of attributes: object ID (OID), the bounding box coordinates, the coordinates and features description of the keypoints in the bounding box, *etc.* (Kestrel-Mobile computes other attributes during this process, as discussed below). Kestrel-Mobile tries to associate each detected object with its previous appearance in the last YOLO frame by tracking the object over k frames. In the new YOLO frame, Kestrel-Mobile finds the match

¹While our paper is about tracking vehicles, many of the components we use are generic and can be used to detect and track objects in general (*e.g.*, people). Where a component is specific to vehicles, we will indicate as such.

 $^{^{2}}$ An optical flow in a video is the pattern of motion of points, edges, surfaces, objects in the scene.

by comparing the keypoints of tracked objects with those in the bounding box of the newly detected one. The new objects that are successfully associated with a tracked object will inherit the existing object ID, otherwise a new unique ID is assigned. Also, for each matching existing box, Kestrel-Mobile resets the tracking attributes, *i.e.*, the box coordinates are updated with the new box coordinates, the feature keypoints are updated with those from the YOLO frame inside the new box. Between two YOLO frames, Kestrel-Mobile tracks the object by updating the box attributes frame by frame. Specifically, it calculates the average displacement of all the keypoints associated with one object, and updates the box coordinates accordingly. By integrating the optimized YOLO with KLT, we have built a fast and effective mobile video processing toolkit that achieves real-time robust object detection, localization and tracking on mobile devices. Figure 3 shows Kestrel-Mobile tracking 2 vehicles.

Object State Extraction. An object seen in a single camera can go through different states: it can enter the scene, leave the scene, move or be stationary. Determining these states is important to filter out uninteresting objects like permanently parked cars, etc. Kestrel-Mobile differentiates moving objects from stationary ones by detecting the difference of the optical flow in the bounding boxes and in the background. For moving objects, the state of the object changes from entering, moving, to exit, as it moves through the camera scene. A moving object can become stationary in the scene, a stationary object can also start moving at any time. For example, a car may come to a stop sign or a traffic light at an intersection, or park temporarily for loading passengers. If the object is detected as exited from the scene, then the object can be evicted. Kestrel-Mobile uses fast eviction: after an object exits the scene, all the feature points and information for tracking are evicted from memory. This, in our experience, not only helps save memory, but also improves accuracy because evicting the keypoints from previous objects can reduce the search space for matching candidates (§4.1).

Extracting the Direction of Motion. Kestrel-Mobile estimates the direction of motion of a vehicle to re-rank the cloud-supplied paths. Suppose a vehicle is moving towards camera A as shown in Figure 5, but one of the cloud-supplied paths shows that it went through a different camera, say B located away from the vehicle's direction of motion. Since it could not possibly have gone through B, the mobile pipeline can lower that path's rank. Extracting the direction poses two challenges unique to mobile devices: how to deal with the *movement of the mobile device*, and how to *transform the direction of motion* to a global coordinate frame.

Kestrel-Mobile addresses the first challenge by compensating for camera movement. In general, to extract the direction of motion in the frame coordinates, we can use the difference between the bounding box positions from consecutive frames. In our experience, however, using the optical flow in the bounding box gives us more fine-grained direction estimation that is robust to the errors in YOLO's bounding boxes. To compensate for camera movement, we subtract the optical flow of the background from the optical flow of the object bounding box. The result is the direction of the object under consideration (Figure 4).

Kestrel-Mobile needs to transform the direction of motion to a global coordinate frame in order to reason about, and re-rank,

SenSys, November 6-8, 2017, Delft, The Netherlands



Figure 5-Extracting the Direction of Motion and Coordinate Transformation

other vehicle paths provided by the cloud. For an arbitrary camera orientation, the exact moving direction in global coordinates can be calculated using a homography transformation [68]. However, in most practical scenarios, cameras have a small pitch and roll (*i.e.*, no tilt towards ground or sky and no rotation, videos are often recorded in either portrait or landscape) especially across a small number of frames. So, we only use the azimuth of the camera, and infer only the horizontal direction of moving objects (Figure 5). As per the Android convention (that we use to obtain our dataset), right (east) is 0° and counterclockwise is positive in camera (global) coordinates. The transformation is simply $\theta_g = \gamma + (\theta_c - 90^{\circ})$ where θ_c (θ_g) is the direction of motion in the camera (global) coordinates, and γ is the azimuth of the camera orientation.

To obtain the azimuth for mobile cameras we use the orientation sensor on off-the-shelf mobile devices. We have a custom Android application that records meta-data of the video (GPS, orientation sensor, accelerometer data, *etc.*) along with the video. However, to use this information directly is not feasible as we cannot align the orientation sensor readings to a particular frame of the video. The timestamp of the video being taken on a mobile device may or may not be available.

To find the correct time alignment between frames and mobile sensors, Kestrel-Mobile takes the first few seconds optical flow pattern as a sliding window, and calculates its correlation with orientation sensors, as it slides through the beginning sequence of the sensor readings. By selecting the correlation peak, Kestrel-Mobile can find the time shift and calibrate the timestamp of each video frame (Figure 6). In this experiment, the camera is horizontally panned towards the left first and then towards the right. The accumulated optical flow reflects the total movement of the scene in the video. Both the orientation sensor and accumulated optical flow show the correct movement of the camera. We see that although there is correlation in the curves, they are shifted along the time axis and, after calibration, the optical flow is better aligned with the azimuth.

Extracting Object Descriptors. Kestrel-Mobile's pipeline requires a method to *re-identify* cars seen in cloud-provided vehicle paths with vehicles detected by its object detector. For this, Kestrel-Mobile extracts *descriptors* of objects. A good descriptor has the property that is *low-complexity* (can be computed easily and requires minimal network bandwidth to transmit), yet can effectively distinguish different objects. The simplest descriptor, a thumbnail image, can distinguish objects but can consume significant bandwidth (§5).

Kestrel-Mobile extracts *visual features* for use as descriptors. Visual features fall into two groups, local descriptive keypoints (SIFT, BRISK, ORB [59], *etc.*), and global features (color histogram

Anon.





Figure 6—Time Alignment of Motion Sensor and Optical Flow (OF)

Figure 7—Matching The Same Object From Different Angle Using SURF Features

Figure 8—Object Thumbnail and Corresponding Color Histogram Before and After GrabCut.

(CH), edge histogram, color layout descriptor (CLD), *etc.*). In multicamera scenarios, local features do not work well since objects can be captured in different angles with varying degrees of occlusion, illumination, *etc.* Figure 7 shows two cameras capturing the same vehicle at different angles and scales. Using SURF features and RANSAC [29] for keypoint matching demonstrates how slightly different views can lead to mismatches: the four matching keypoints (connected by green lines) are all false positives. In §5, we quantify the performance tradeoffs of these approaches.

For these reasons, Kestrel-Mobile extracts the color histogram. Before doing so, it preprocesses the bounding box to exclude background pixels. GrabCut [58] efficiently removes the background in a bounding box, leaving only the pixels of the object of interest to be counted. As we see in Figure 8 (left), with so many gray background pixels the color histogram can look very different from the color histogram of the car seen in Figure 8 (right) after removing the background pixels. Background removal helps other visual features as well, since keypoints from the background can confound matching during association.

The choice of color histogram does have drawbacks: vehicles with similar colors can have high correlation scores. Figure 9 illustrates, in our dataset, both a easy case, where Kestrel-Mobile searches for a red sedan among other vehicles of different colors, and one of the most challenging cases, where three white SUVs have very similar color layout that would require careful human inspection. But Kestrel does not rely on color alone: it uses other filters (direction, travel time estimation) to further narrow the candidate space and achieve high object association accuracy. Indeed, using these techniques, Kestrel was able to correctly distinguish between the three SUVs. We present the details of the association pipeline in the following sections.

4 KESTREL CLOUD PIPELINE

In Kestrel's architecture, video from fixed cameras are streamed to the cloud (Figure 1). The cloud *continuously* extracts objects and attributes from the streamed video. Unlike the mobile pipeline, which is only invoked on demand, the cloud pipeline processes every frame (and can afford to do so because of the cloud's resource elasticity). Unlike the mobile pipeline, the cloud pipeline does not need to compensate for camera motion, and, because the camera positions are known ahead of time, transformation to a global coordinate frame is straightforward. Kestrel's cloud pipeline (Kestrel-Cloud) re-uses the classifier and the matching descriptor from its mobile pipeline. It could have used a more heavyweight classifier, but YOLO has good accuracy while enabling full frame rate processing: more recent classifiers, like MSCNN [21] achieve higher accuracy, but with an order of magnitude higher execution time. Finally, the cloud pipeline would have the resources to use local descriptors for matching, but as discussed before, these can introduce matching error because of perspective differences (Figure 7). Beyond these similarities, the cloud pipeline performs several tasks that the mobile pipeline does not; we discuss these below.

To search for a vehicle's path through the camera network, a Kestrel user (e.g., a surveillance officer) selects a car in a specific frame of a fixed camera's feed. Given this input, the cloud pipeline runs a path inference algorithm, an important component of which is pair-wise instance association.

4.1 Pair-Wise Instance Association

An *instance* corresponds to a specific object captured at a specific camera, together with the associated attributes (§3.2). *Instance association* determines if object instances at two cameras represent the same object or not. Without loss of generality, a camera network can be modeled as an arbitrary topology shown in Figure 10. Assume that each intersection has a camera (*Cam* c), and each camera extracts several instances. We define the *n*th object instance of *Cam* c as $o_{c,n}$. Kestrel-Cloud infers association between any pair of object instances ($o_{x,i}, o_{y,j}$), using three key techniques: visual features, travel time estimation and the direction of motion.

Preprocessing. Before associating instances, Kestrel-Cloud preprocesses each instance to narrow the search space. Since YOLO can detect several types of objects, the preprocessing stage filters all object types other than vehicles. Also, it filters all stationary objects, such as parked cars on the side of the streets. Even so, there could still be hundreds of instances to search from multiple neighboring cameras. Kestrel-Cloud estimates the travel time ET(x, y) between a pair of camera x, y, and sets a much smaller but coarse time window (which is refined in the next step) to limit the number of candidate instances at this preprocessing stage.

Spatio-temporal Association. To associate between two object instances $(o_{x,i}, o_{y,i})$, Kestrel-Cloud uses the timestamp of the object's first scene entry TIN and its last appearance exiting the scene TOUT. Kestrel-Cloud calculates the total travel time from Cam x to Cam y as $\Delta T(o_{x,i}, o_{y,j}) = TIN_{y,j} - TOUT_{x,i}$, and compares this to the estimated time ET(x, y) needed to travel from Cam x to Cam y. Taking the exit timestamp from the first camera and the entry timestamp in the second camera filters out any variance due to a temporary stop (*e.g.*, at a stop sign) while in the camera view.

SenSys, November 6-8, 2017, Delft, The Netherlands



Figure 9—Object Descriptor Correlation

Figure 10—A Network of Camera Captured Object Instances



Figure 11-Heterogeneous Camera Network

One way to estimate travel time between two cameras is to use the Google Directions API [33]. For dense camera deployments typical of campuses, travel times between cameras are small (~30 secs) and can be inflated by noise. To avoid this, we estimate the travel time between each pair of cameras using a small annotated ground truth dataset. Further, since cars have different travel speeds and may stop in between, there is some variance in the travel time between cameras. We also estimate this variance from our training dataset. To filter out unlikely associations (*e.g.*, a car cannot appear at two cameras that are one block away from each other within 1 sec), we apply temporal constraints derived from the dataset: specifically, if the estimated travel time is *x*, over 95% of the vehicles take a travel time of $x \pm 0.6x$.

Visual Similarity. The visual similarity of object instances is also used in object association. Given two color histograms H_1 , H_2 from two instances, Kestrel-Cloud measures their similarity using a *correlation offset*, which is 1 minus the correlation between the two histograms. The lower the offset is, the more likely the two instances are the same. We find that this metric can discriminate well between similar and dissimilar objects. From a sample of 150 images (figure omitted for space reasons) we find that over 95% of identical objects have an offset of less than 0.4, while over 90% of different objects have an offset larger than 0.4. In a more general setting, this threshold can be learned from data.

Direction Filter. We also use the direction of motion (§3.2) to prune the search space for association. For example, in Figure 10, assume that the moving direction of instance $o_{1,1}$ from Cam 1 is from east to west, then Kestrel-Cloud ignores Cam 4, and searches over all object instances of Cam 2 within the temporal window around the expected travel time between Cam 1 and Cam 2, and considers only those instances exiting Cam 2 towards Cam 1.

Instance Association in the Mobile Pipeline. The mobile pipeline also contains some of these elements of instance association, which are necessary for re-ranking paths. Specifically, the mobile pipeline must eliminate stationary objects, but does not need spatio-temporal filtering and temporal filtering in the preprocessing stage (because Kestrel-Cloud selects the mobile device at the relevant location, and supplies the corresponding Kestrel-Mobile a time window over which to search for associations). The mobile pipeline computes visual similarity and applies the direction filter.

4.2 Path Inference

Building upon its object detection, attributes and association components, Kestrel-Cloud can infer, *in the cloud*, the path of a target object through a multi-camera network. Our pair-wise association simplifies the object instance network and retains only the valid paths *i.e.*, filters out objects moving in incorrect direction, those outside the temporal window and those whose visual correlation is too low. In this *pruned network*, the *weight* $w(o_x, o_y)$ of a *link* $l(o_x, o_y)$ is defined as the color histogram correlation offset (lower the offset, better is the correlation). A path from a source instance o_s to destination o_d in an instance network is defined as an *instance* path $p(o_s, o_d) = o_s, o_1, o_2, ..., o_d$. A physical route in the real world, termed a *camera path*, is defined by the sequence of cameras traversed by a vehicle: multiple instance paths can traverse the same camera path.

Consider one object of interest $o_{x,i}$, captured by either a mobile user or a camera operator from Cam x. Kestrel-Cloud seeks to infer the instance path it takes and answers this query in near-real time, generating the path and instances at each camera as a feedback to the user.

One straightforward approach is to assign the correlation offset as a weight to the links in the pruned instance network, and use maximum likelihood estimation (MLE) to find the minimum weighted path from any given instance. However, this approach fails to capture the correlation between non-neighboring instances, as it only relies on the link weight. One can associate between every pair of cameras, producing a clique, but our problem then reduces to finding the minimum weighted k-clique which is NP complete.

Instead, Kestrel-Cloud takes a hop-by-hop iterative approximation approach to find the best matching last hop instance each time. At every hop with multiple candidate paths joining at one instance, Kestrel-Cloud evaluates the weight of each path from the candidate set and uses the Viterbi decoding algorithm [65] to eliminate paths with heavy weight to keep the algorithm at polynomial complexity. Given an instance in the instance network (Figure 10), we define a valid neighbor as an instance that has a direct link and passes all the filters defined in §4.1. Every time a new instance is added to a path, the path weight is multiplied by an amount equal to the average correlation offset of this instance from every other instance in the path. The worst case complexity is $O(kn^2)$, assuming every camera is right next to each other, and every instance has a valid link to any instances. In practice, each camera has an average of 2.2 neighbors, and vehicles usually traverse each node in the network only once except for looping scenarios.

Re-ranking Paths in the Mobile Pipeline. The output of the path inference algorithm can be ambiguous: there may be multiple instance paths with high path weights. In this case, the cloud pipeline attempts to find (from previously uploaded metadata) a set of mobile



Figure 12-Example Ground Truth Instance Paths Across the Heterogeneous Camera Network

devices which might have videos that could help resolve this ambiguity. The cloud pipeline sends to the mobile device its set of ranked instance paths. Given this set, the mobile camera does *pair-wise association* with the instances from both the last hop and the next hop static cameras of each instance paths. If the locally extracted instance matches the path's direction, a new instance path is generated whose instance path weight is augmented with a *pair-wise association* score. Note that a single instance path could generate multiple new instance paths with different intermediate mobile instances and path weights. The mobile pipeline can re-rank the paths and transmit these to the cloud. For each new mobile instance, it only needs to transmit the color histogram and the instance's direction of motion in case the cloud needs to query another mobile device.

5 EVALUATION

In this section, we evaluate the end-to-end performance of Kestrel, and the performance and overhead of each submodule.

5.1 Methodology

Dataset. Our campus surveillance network contains 98 cameras, 64 of which monitor residential areas and the remaining monitor main streets. We collected a video dataset from 17 distributed nonoverlapping cameras (6 mobile cameras and 11 surveillance cameras), deployed in a whole block of a residential area (Figure 11). Kestrel builds a hybrid camera network according to the topology shown in the figure. The surveillance cameras (in yellow in the figure) are commercial Axis Q-6035 and Q-6045 cameras, with a resolution of 1280×720 at 10 frames per second, and mounted on light poles. The mobile cameras (in pink in the figure) are from off-the-shelf smartphones held by human users at street level, and these capture 1920×1080 video at 30fps using a customized app which also captures inertial sensor readings from the phone. Users collect video by sweeping the visible area for interesting events: thus, unlike fixed cameras, the orientation (yaw and pitch) of the cameras are continuously changing.

Ground Truth. To evaluate the association, we manually annotate a ground truth list of instances that match across cameras (Figure 12). Specifically, each instance is labeled with a camera id (CID) and an object id (OID), *i.e.*, the first car that appears in camera 20 would be labeled as CID(OID) = 20(1). We visually track each vehicle from camera to camera, and label their corresponding instances in each camera with the same global object ID (GID). For example, if a vehicle travels through cameras 101, 20, 106, 19, appearing as the 8th, 12th, 15th, 4th object of each camera respectively, will form an *instance path* which can be represented as a set of CID(OID), {101(8), 20(12), 106(15), 19(4)}, and can be assigned a GID 2. In this section, we also evaluate Kestrel's accuracy in detecting the *camera path*, the sequence of camera traversed by the car. The camera path ground truth is represented as the corresponding set of CIDs.

From the dataset of 17 cameras, we have annotated, over a total of 235 minutes of footage, 120 global objects / vehicles. The longest path contains 7 cameras. In most cases, one vehicle going through the camera network can be seen at three to four cameras. The number can vary if the car goes into a plaza or parks in a private garage in the residential block. To form as many valid ground truth paths as possible, we issue a query with each instance of the path, except the last instance before it disappears from the camera network. Kestrel runs those queries to infer backwards as many hops as possible. For example, if one car goes through 3 cameras, a, b, c, Kestrel can infer from the instance at c to get b and a, or infer from b to get a. That gives both a two-hop path and a one-hop path for the evaluation. Collectively, we are able to establish 311 one-hop paths, 197 twohop paths, 125 three-hop paths, 55 four-hop paths, and 23 five-hop paths. There are a small number of longer paths, which we ignore because their number is too small to draw valid conclusions.

Metrics. The primary performance metric for Kestrel is accuracy. We use *recall* and *precision* to evaluate the accuracy of Kestrel's association and path inference. Among all the given paths, recall $(\frac{TP}{TP+FN})$ measures the fraction of the paths for which Kestrel can make a successful inference. Given one instance of an object, precision $(\frac{TP}{TP+FP})$ measures how often Kestrel correctly identifies all the instances on the path within its top k choices. Ideally, Kestrel should exhibit high recall and precision for as small a k as possible.

We explore two aspects of Kestrel's accuracy. Its *camera path* accuracy measures how accurately Kestrel can identify the sequence of cameras traversed by a vehicle. We measure this by majority voting across the top five instance paths. We also measure the accuracy of determining the *instance path*: such a path identifies not just the right sequence of cameras but also the right instances at each camera matching the queried vehicle. For this, we present the top-k accuracy: how often its top k choices contain the ground truth path. Specifically, if the annotated ground truth is among those top k choices, Kestrel is said to generate a true positive (TP), otherwise a false positive (FP). Meanwhile, all the other unselected choices are considered as negatives (N). The negatives are true (TN) when they belong to different objects, otherwise false negatives (FN). Ideally, Kestrel should exhibit high recall and precision for as small a k as possible.



Figure 13—Augmenting Sparse Camera Network with Mobile Cameras to Achieve Full Performance

We evaluate several other metrics as well. We measure mobile pipeline energy consumption in Joules per frame, by using a CA 60 current clamp [5] attached to the positive wire of the TK1's power supply. By inducing a magnetic field on the conductor and making use of a Hall effect sensor, we can compute the current passing through. We use a DataQ DI-149 Data Acquisition Kit [6] that samples at 80 Hz and allows us to save the readings in a file for post-processing. We measure detection latency of the object detector by CPU elapsed time, and detection accuracy by mean average precision (mAP [9]), which captures false positives as well as the false negatives.

5.2 Camera Path and Instance Path Accuracy

In this section, we evaluate Kestrel's accuracy. Kestrel supports a hybrid camera network, where videos captured on mobile devices can augment an existing fixed camera deployment. We would like to understand (a) how much the mobile cameras improve accuracy over the existing fixed camera deployment, and (b) how well a hybrid camera network performs relative to a fixed camera deployment *of the same size.* To this end, we evaluate accuracy over our 17-camera dataset by forming three topologies: a *Mobile* topology where the mobile cameras run the Kestrel mobile pipeline, a *Full* topology where all cameras are treated as fixed and run the cloud pipeline, and a *Partial* topology, which consists only of the 11 fixed cameras.

Figure 13 shows the average camera path inference precision and recall on these three different networks. To start with, with a Full topology where every intersection is monitored, Kestrel achieves a camera path precision of 99.2%, while Partial only infers two thirds of the paths correctly. Interestingly, the presence of a mobile camera enables Kestrel's lightweight and efficient mobile pipeline to effectively disambiguate path choices, achieving a precision of 97.7%. Recall results are qualitatively similar, but average recall numbers are slightly lower (explained in the next paragraph), with about 80% recall for Mobile vs. 90% for Full. Recall for partial is significantly lower. These results suggest that a hybrid camera network can approach the accuracy of an equivalently sized fixed network, and the addition of mobile devices to an existing fixed camera network can significantly improve performance.

For instance paths, Kestrel achieves good top-3 precision and recall for the Mobile topology, whose performance is close to the Full topology, and significantly better than Partial. In general, the recall performance for Mobile is lower than the Full topology (both for camera path and instance path accuracy) because of Kestrel-Mobile's energy optimization to avoid invoking the object detector sometimes leads to missed detections. We explore the trade-off between energy and recall below. Moreover, camera path accuracy is higher than instance path accuracy. This is because Kestrel can precisely infer direction of motion, but its association accuracy has room for improvement.

Figure 14 shows Kestrel accuracy as a function of path length. Its top 3 choices (k = 3) almost always (~90% precision and recall) find the correct paths for up to three camera hops. Increasing k to 4 or 5 provides marginal improvements, while reducing it to 1 gives a precision and recall of slightly less than 80% for up to 3 hops. Mobile has comparable performance to Full with less than 10% degradation. In subsequent sections, we discuss how much of this performance can be attributable to object detection, attribute extraction, and association.

This accuracy is significant, considering that, statistically, for each hop, Kestrel has to find the correct instance among an average of 20.38 candidates from 2.2 cameras. More important, the distribution of the number of candidate paths increases exponentially as the number of hops increases; at 5 hops, the number of candidate instance paths in our data set ranges from 10,000 to over a million. In this regime, a human operator can not manually identify the correct path just by visually inspecting pictures, but Kestrel can achieve nearly 70% recall and 80% precision at 5 hops, despite only using the color histogram descriptor, as a result of our techniques.

At 5 hops, the precision and recall might appear low, but Kestrel could be practicable even in this range, with a little operator input. Suppose that a user issues a path inference query and does not see the right 5-hop answer. If there is a correct 3-hop instance, she can re-issue a path inference query using as the starting hop instance the last correct instance, and stitch together the returned results.

5.3 Energy and Latency

Kestrel significantly reduces the energy consumption on the mobile device by only invoking the mobile device when presented with a query (unlike Kestrel-Cloud, which processes every frame), and only when Kestrel-Cloud has several high-ranked path candidates which a mobile device can help disambiguate. In our dataset, 84.6% of the queries require invoking a mobile device to resolve a path ambiguity.

Kestrel-Mobile also reduces energy usage by careful design. The bottleneck of the mobile pipeline is running the neural net for object



Figure 14—Recall, Precision of Path Inference



Figure 16-Feature Size vs Computation Latency

detection: on TK1, YOLO alone consumes 2.25 J/Frame whereas attribute extraction only drains 0.42 J/Frame. In terms of latency, YOLO takes an average of 0.259 Sec/Frame, while tracking only takes 0.081 seconds. Therefore, the fact that the optical flow motion filter can avoid running YOLO on every frame significantly conserves energy and reduces latency on mobile devices at a cost of inference accuracy.

Energy and latency can be traded-off for higher accuracy by adjusting the threshold for the optical flow filter. Figure 15 quantifies this tradeoff between energy / latency and the inference error. Intuitively, the higher the motion filter threshold, the fewer the YOLO invocations, which result in less energy and lower latency. Compared to YOLO running on every frame (motion threshold being 0), a small motion filter of 4 pixels can almost cut the energy usage in half, which can double mobile battery life, without significantly sacrificing path inference performance. On the contrary, if the motion filter is too insensitive (larger than 6 pixels), Kestrel will miss a lot of vehicles, which results in much lower recall.

Choice of Descriptor. To validate our choice of color histogram, Figure 16 compares different features in terms of the average data size and computation latency per object. Local features like SIFT and SURF incur both high computation overhead and large size (and we have earlier shown that they are sensitive to perspective differences, and so not a good choice for Kestrel for that reason), while the lightweight color histogram incurs minimal latency and small size with reasonable performance.

5.4 Association Performance

In this section, we evaluate the performance of one-hop association. Given one annotated instance at any camera, this component tries to pick out the associated instance among all the instances from all the neighboring cameras.



Figure 15-Energy / Latency and Performance Trade-off

When conducting this evaluation, we also evaluate the efficacy of each of the components of the association algorithm (§4.1). We start with ranking the visual association score among all candidates without any additional processing at all (NONE). Then we add each component one by one in this order: preprocessing (PRE), direction filter (DIR), spatio-temporal association (ET), background subtraction (GrabCut). This process is cumulative: for example, DIR also includes PRE and the visual association score components.

Figure 17 shows the precision and recall of each successive combination. Generally, given an instance at any location at any time, Kestrel can almost always (> 97%) find the same object in a neighbor camera, if there were one to be found, within the top 3 returned instances. Comparing the performance of each component combination, PRE effectively narrows down the search space and brings precision and recall to nearly 90%. DIR further rules out false positives, as cars moving in wrong direction can confound visual association. Both DIR and ET increase recall and precision by moving true positives higher in the rank (which increases Top 1 precision). Grabcut also increases precision and recall noticeably.

Figure 18 shows an example of how various steps of the association algorithm can reduce the number of candidates for the association. In this specific query, Kestrel is able to filter down to 3 candidates from an initial set of 342, and finds the correct target vehicle by ranking the histogram correlation offset from among these three, even though the two cameras capture the vehicle from completely different perspectives. When averaged over all associations, each target has an average of 385 candidates to match, PRE prunes the search space to about 20, DIR removes about 8 cars on average going in the wrong direction, ET is able to further narrow down to 6 candidates. Finally, Kestrel ranks the correlation offset of the remaining candidates.

5.5 Object Detection Performance

Performance of CPU offload. In Table 1 we summarize the timing results of running the various strategies for GPU memory optimization. The *Original network* is already optimized (compared to the original network that runs on server class GPUs) in that it does not allocate memory for redundant variables not used in the testing phase. Running the computation on the CPU takes close to *half a minute* per frame, which is extremely slow, so it is imperative to leverage the GPU cores on the board for accelerated computation. Moreover, we see that for the original network, memory constraints



Figure 17—Precision and Recall Performance of the Semantics Association between Neighboring Cameras using Different Kestrel Component Combinations

Target					_				Avg Cand. #
Candidate				Į,		E	1		385.75
Search Space		8		Į (į	-	1.12			
PRE					JU				20.38
DIR	Wrong Direction			Ð	Wrong Direction		5		12.23
ET				S		Out of Bound	Out of Bound		6.77
Histogram Corr Offset		0.45329	0.31745	0.22681					-

Figure 18—Shrinking the Candidate Size at Each Step Matching a Target

							(
	Memory	FC layer: resulting size of	CDU	GPU	Split	CPU Offload	Dinalina		
	Requirement	matrix multiplication	CIU	010	Spin	CI O Ollioad	1 ipenne		
Original	2.8 GB	$1 \times 50176 * 50176 \times 4096$	25.026524	N/A	10.249577	0.703449	0.416910		
Medium	1.6 GB	$1 \times 25088 * 25088 \times 2048$	24.315950	0.272191	0.573386	0.400366	0.261985		
Small	1.3 GB	$1 \times 12544 * 12544 \times 1024$	24.003447	0.259000	0.394960	0.299940	0.261144		
T-bl-1 Annal Time T-law to Day Data time and Laws (area 1-)									

Table 1—Average Time Taken to Run Detection per Image (seconds)

do not allow running it on GPU at all; our subsequent optimizations to reduce the memory footprint make this possible.

CPU offload brings the computation time to under 1s. Using the CPU for the FC layer allows us to read the weight file only *once* for all the frames and store the weights in CPU memory (which supports demand paging). Finally, CPU offload along with pipelining gives the best results; it brings down the computation time to about 0.42s or almost *60 times faster* than running it on the CPU. We see this speed up because although the FC layer is memory intensive, the running time even on the CPU is not a bottleneck—in other words, the CPU processing completes within the time that the GPU is processing the convolutional layers of the next frame. This is interesting because it is achieved without compromising accuracy.

Other optimizations. Other optimizations are less effective. *Split* allows running original network on GPUs but it is still slow as it incurs a weight-file read overhead every frame. Reducing the size of the network, by reducing the size of the weight matrix in the FC layer to one fourth of the original size (*Medium network*) and 1/16 the original size (*Small network*), enables YOLO to run faster and follow similar trend as *Original network*, with one exception: for *Small network* the GPU is the fastest alternative as the FC layer is small enough that overhead of pipelining negates its benefit. However, these networks incur accuracy loss. On the Pascal VOC 2012 test dataset [13], we lose about 2-3% in mAP at each step of the reduction from Original to Medium to Small.

Power measurement. Interestingly, we found that CPU offloading based schemes are also the most energy-efficient. This is because these schemes optimize speed and so the circuitry is used for much shorter duration. They reduce the energy requirement by more than 95% for Original network (107J/frame for running on CPU, 3.78J with CPU offload, 3.95J with CPU Offload + Pipelining). Smaller CNNs do save energy (2.25 J/frame on Medium network for CPU Offload + Pipelining) but at the expense of accuracy.

5.6 Attribute Extraction

Tracking Accuracy vs Processing Latency. Most existing stateof-the-art multi-object tracking techniques assume offline processing for stored video and can incur significant overhead. To demonstrate this, we pick a representative of this class, MDP [66], which is regarded as the algorithm with the best accuracy with reasonable processing overhead. Also, we extend a robust single-object tracking tool, OpenCMT [52] to perform multi-object tracking for the comparison with Kestrel.

We use a standard multi-object tracking benchmark [41] to compare our tracker against these approaches. To level the playing field, we use YOLO for detection for all three approaches. In addition to precision, recall, and false negatives, we use a metric called MOTP [41] that measures the tightness of the bounding boxes generated by the tracking algorithm.

Method	Rcll	Prcn	FN	MOTP	TM	TD	
Kestrel	76.4	88.2	152	76.4	0.081	0.057	
MDP	70.1	87.1	192	75.4	N/A	0.146	
CMT	75.9	88.6	155	76.3	1.599	0.907	
Table 2—Tracking Accuracy and Latency on TK1 and Desktop							

We run the trackers on our vehicle dataset to compare the tracking performance as well as the average processing latency per frame, both on the TK1 (TM) and a desktop server (TD). Table 2 shows that Kestrel can achieve a slightly higher recall, comparable precision, lower false negatives and a comparable MOTP score compared to more sophisticated algorithms, while incurring an order of magnitude lower processing latency on the TK1 (one of the algorithms, MDP, cannot even be executed on TK1). The primary reason for the improved tracker performance is our design choice to periodically invoke YOLO, which can refine the box generated by the tracker. Between two YOLO frames, tracking and association is easier than the kinds of continuous tracking performed by modern trackers.

We also measure the energy consumption on TK1 using the same setup, and we find that our tracking requires 0.42 J/frame, but OpenCMT requires energy as high as 8.4 J/frame.

SenSys, November 6-8, 2017, Delft, The Netherlands

Period	Rcll	Prcn	TTRK	TYOLO	TTOT	
1	76.4	88.2	0.081	0.259	0.340	
3	74.7	82.5	0.081	0.086	0.167	
5	64.4	74.2	0.081	0.051	0.132	
Table 3—Tracking Performance and Latency Tradeoff on TK1						

Next, we explore the tradeoff between tracking accuracy and total latency by calling YOLO every 1, 3, 5 frames, and tracking only on non-YOLO frames. This evaluation uses a series of consecutive frames that have moving objects and pass the motion filter. In other words, we try to examine the best performance when every frame has objects to detect and track. Intuitively, the more frames between YOLO detection, the larger chance that the tracker may be led astray. As shown in Table 3, less frequent YOLO detection effectively reduces the total average processing latency (TTOT) per frame, while maintaining a reasonable precision and recall (in some cases, even better than MDP and CMT, not shown) for attribute extraction. Specifically, when YOLO is invoked every 3 frames, Kestrel's mobile pipeline can sustain about 6 fps without noticeable loss in tracking performance. With further optimization, and on more recent GPUs [12], we believe we can get to 10 fps, thereby reducing query latency. Finally, our camera movement compensation algorithm has nearly identical tracking precision and recall at walking speeds, when compared with a stationary camera.

Sensitivity to Camera Motion. To examine the tracking robustness of Kestrel-Mobile for mobile cameras, we collected a set of videos with a mobile camera while walking, biking, and driving on our campus. In these videos, we manually annotated the bounding boxes of vehicles appearing in every frame over nearly 6000 frames. Then, we ran our camera motion compensation algorithm, and our tracker, to evaluate its accuracy. Our results show a similar performance to Table 2, with an average precision of 74.9%, an average recall of 70.9%, and an MOTP of 78.7%. This indicates Kestrel-Mobile's algorithms can accurately track vehicles even when there is significant camera motion, with an accuracy comparable to a fixed camera.

Bandwidth. Kestrel architects the mobile and cloud pipeline to be bandwidth-efficient. To quantify bandwidth-efficiency, in our dataset a typical 10-minute video file recorded at 1920×1080 (30 fps) is around 1.4GB. Streaming every frame in real time requires ~20Mbps. Instead, Kestrel only sends attributes. In our dataset, after the cloud pipeline, the average size of a query sent to the mobile is only 1.52KB. The re-ranked result and its corresponding instance metadata average 0.92KB.

6 RELATED WORK

DNNs for mobile devices. Recent work has explored neural nets on mobile devices for audio sensing activity detection, emotion recognition and speaker identification using an accelerometer and a microphone [39, 40]. Their networks use only a small number of layers and are much simpler than the networks required for image recognition tasks. Another work [4] has benchmarked smaller Caffegenerated CNN (AlexNet) [37], with 5 convolutional layers and 3 FC layers, for vision-based classification-only (no localization) on

the TK1. DeepX [38] is able to achieve reduced memory footprint of the deep models via using compression, at the cost of a small loss in accuracy. LEO [30] schedules multiple sensing applications on mobile platforms efficiently. MCDNN [34] explores cloud offloading and on-device versus cloud execution tradeoff. However, the models they run are smaller than the ones required for our work, hence requiring us to pursue the optimizations explained in the paper.

Object Detection. Early object detectors use deformable part models [28] or cascade classifiers [64], but perform relatively poorly compared to recent CNN based classification-only schemes which achieve high accuracy at real time speeds [37]. However, top detection systems like R-CNN [32], Fast R-CNN [31], and MSCNN [21] exhibit less than real-time performance even on server-class machines. YOLO [57] is a one shot CNN-based detection algorithm that predicts bounding boxes and classifies objects, and we have used a mobile GPU on a deep CNN like YOLO.

Object Tracking. Tracking objects in videos is a well researched area. However, some tracking approaches like blob tracking [35, 62] work well in static camera networks, but not for mobile cameras. Many other trackers are targeted to single object tracking, e.g., KLT [48], OpenCMT [52]. Some work has gone into tracking multiple objects [24, 50, 56, 66]. However, most of these trackers are not targeted at execution on resource constrained devices. A recent work, Glimpse [23], achieves tracking on mobile devices using a combination of offloading and keeping an active cache of frames to work on, once a stale result is received from the server. By contrast, our tracking does not rely on offloading. We have compared the performance of our tracker with OpenCMT [52] and MDP [66] and the sensitivity to camera motion in §5.

Video Surveillance Systems and Object Association across Cameras. Vigil [67] is a surveillance system for wireless cameras that uses powerful edge devices at the cameras and performs simpler tasks on the edge device while offloading the more complex computations to the cloud. However, it does not specifically address the re-identification problem we consider. Prior work [63] tries to associate people *etc.* across different cameras using a query retrieval framework by ranking nodes in their camera network. Other work [49] proposes a centralized multi-hypothesis model to track a vehicle through a multi-camera network. While Kestrel can support such applications, our focus is to enable mobile camera based surveillance, so our architectural and design choices are different from this line of work.

7 CONCLUSION

This paper explores whether it is possible to perform complex visual detection and tracking by leveraging recent improvements in mobile device capabilities. Our system, Kestrel, tracks vehicles across a hybrid multi-camera network by carefully designing a combination of vision and sensor processing algorithms to detect, track, and associate vehicles across multiple cameras. Kestrel achieves > 90% precision and recall on vehicle path inference, and can do so while significantly reducing energy consumption on the mobile device. Future work includes experimenting with Kestrel at scale, extending Kestrel to support more queries, and different types of objects, so it can be used as a general visual analytics platform.

SenSys, November 6-8, 2017, Delft, The Netherlands

BIBLIOGRAPHY

- [1] Agent VI. http://www.agentvi.com/.
- [2] Avigilon Video Analytics. http://avigilon.com/products/video-analytics/solutions/.
- [3] BriefCam. http://briefcam.com/.[4] Caffe on TK1. http://goo.gl/6hgbM6.
- [5] Current Clamp. http://pdimeters.com/products/Accessories-Parts/PDI-CA60.php.
- [6] Data Acquisition Kit. http://www.dataq.com/products/di-149/.
- [7] Google Tango. https://www.google.com/atap/project-tango/.
- [8] Google's Pixel C. https://store.google.com/product/pixel_c.
- [9] Mean Average Precision. http://homepages.inf.ed.ac.uk/ckiw/postscript/ijcv_ voc09.pdf.
- [10] Nexus 9. https://www.google.com/nexus/9/.
- [11] nVidia Jetson TK1. https://developer.nvidia.com/jetson-tk1.
- [12] nVidia Jetson TX1. http://www.nvidia.com/object/jetson-tx1-module.html.
- [13] Pascal VOC. http://host.robots.ox.ac.uk/pascal/VOC/.
- [14] Snapdragon 820 Benchmarks Match The Tegra X1. http://wccftech.com/ snapdragon-820-benchmarks/.
- [15] Snapdragon 820 vs. Tegra K1. https://versus.com/en/ nvidia-tegra-k1-32-bit-vs-qualcomm-snapdragon-820-msm8996/.
- [16] Unified Memory Architecture. http://devblogs.nvidia.com/parallelforall/ unified-memory-in-cuda-6/.
- [17] WatchGuard: Law Enforcement Video Systems. http://watchguardvideo.com/.
- [18] Asus ZenFone AR. Asus ZenFone AR.
- [19] Herbert Bay, Tinne Tuytelaars, and Luc Van Gool. 2006. SURF: Speeded Up Robust Features. In Proceedings of the 9th European Conference on Computer Vision - Volume Part I (ECCV'06).
- [20] Sourav Bhattacharya and Nicholas D. Lane. 2016. Sparsification and Separation of Deep Learning Layers for Constrained Resource Inference on Wearables. In Proceedings of the 14th ACM Conference on Embedded Network Sensor Systems CD-ROM (SenSys '16).
- [21] Zhaowei Cai, Quanfu Fan, Rogerio Feris, and Nuno Vasconcelos. 2016. A Unified Multi-scale Deep Convolutional Neural Network for Fast Object Detection. In ECCV.
- [22] Guoguo Chen, C. Parada, and G. Heigold. Small Footprint keyword spotting using deep neural networks. In Proc. of IEEE ICASSP '14.
- [23] Tiffany Yu-Han Chen, Lenin S. Ravindranath, Shuo Deng, Paramvir Victor Bahl, and Hari Balakrishnan. 2015. Glimpse: Continuous, Real-Time Object Recognition on Mobile Devices. In 13th ACM Conference on Embedded Networked Sensor Systems (SenSys). Seoul, South Korea.
- [24] Wongun Choi, Caroline Pantofaru, and Silvio Savarese. 2012. A general framework for tracking multiple people from a moving camera. *Pattern Analysis and Machine Intelligence* 99 (2012), 1–1.
- [25] David Chu, Nicholas D. Lane, Ted Tsung-Te Lai, Cong Pang, Xiangying Meng, Qing Guo, Fan Li, and Feng Zhao. Balancing Energy, Latency and Accuracy for Mobile Sensor Data Classification. In *Proc. of Sensys* '11.
- [26] Heungsik Eom, P. St.Juste, R. Figueiredo, O. Tickoo, R. Illikkal, and R. Iyer. Machine Learning-Based Runtime Scheduler for Mobile Offloading Framework. In Proc. of UCC '13.
- [27] Every Chicago patrol officer to wear a body camera by 2018. Every Chicago patrol officer to wear a body camera by 2018.
- [28] Pedro F Felzenszwalb, Ross B Girshick, David McAllester, and Deva Ramanan. 2010. Object detection with discriminatively trained part-based models. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 32, 9 (2010), 1627–1645.
- [29] Martin A. Fischler and Robert C. Bolles. 1981. Random Sample Consensus: A Paradigm for Model Fitting with Applications to Image Analysis and Automated Cartography. *Commun. ACM* 24, 6 (June 1981).
- [30] Petko Georgiev, Nicholas D. Lane, Kiran K. Rachuri, and Cecilia Mascolo. 2016. LEO: Scheduling Sensor Inference Algorithms Across Heterogeneous Mobile Processors and Network Resources. In *Proceedings of the 22Nd Annual International Conference on Mobile Computing and Networking (MobiCom '16)*. 320–333.
- [31] Ross Girshick. 2015. Fast r-cnn. In Proceedings of the IEEE International Conference on Computer Vision. 1440–1448.
- [32] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. In *Proc. of IEEE CVPR* '14.
- [33] Google Maps Directions API. Google Maps Directions API.
- [34] Seungyeop Han, Haichen Shen, Matthai Philipose, Sharad Agarwal, Alec Wolman, and Arvind Krishnamurthy. 2016. MCDNN: An Approximation-Based Execution Framework for Deep Stream Processing Under Resource Constraints. In Proceedings of the 14th Annual International Conference on Mobile Systems, Applications, and Services (MobiSys '16).
- [35] Marko Heikkila and Matti Pietikainen. 2006. A texture-based method for modeling the background and detecting moving objects. *IEEE transactions on pattern* analysis and machine intelligence 28, 4 (2006), 657–662.
- [36] Alex Krizhevsky. 2014. One weird trick for parallelizing convolutional neural networks. CoRR abs/1404.5997 (2014). http://arxiv.org/abs/1404.5997

- [37] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. ImageNet Classification with Deep Convolutional Neural Networks. In Proc. of NIPS '12.
- [38] Nicholas D. Lane, Sourav Bhattacharya, Petko Georgiev, Claudio Forlivesi, Lei Jiao, Lorena Qendro, and Fahim Kawsar. 2016. DeepX: A Software Accelerator for Low-power Deep Learning Inference on Mobile Devices. In Proceedings of the 15th International Conference on Information Processing in Sensor Networks (IPSN '16).
- [39] Nicholas D. Lane and Petko Georgiev. Can Deep Learning Revolutionize Mobile Sensing?. In Proc. of HotMobile '15.
- [40] Nicholas D. Lane, Petko Georgiev, and Lorena Qendro. DeepEar: robust smartphone audio sensing in unconstrained acoustic environments using deep learning. In Proc. of UbiComp '15.
- [41] L. Leal-Taixé, A. Milan, I. Reid, S. Roth, and K. Schindler. 2015. MOTChallenge 2015: Towards a Benchmark for Multi-Target Tracking. arXiv:1504.01942 [cs] (April 2015).
- [42] Lenovo Phab 2 Pro. Lenovo Phab 2 Pro.
- [43] Stefan Leutenegger, Margarita Chli, and Roland Y. Siegwart. 2011. BRISK: Binary Robust Invariant Scalable Keypoints. In Proceedings of the 2011 International Conference on Computer Vision (ICCV '11).
- [44] Rachael Lindsay, Louise Cooke, and Tom Jackson. 2009. The impact of mobile technology on a UK police force and their knowledge sharing. *Journal of Information & Knowledge Management* 8, 02 (2009), 101–112.
- [45] Song Liu, Karthik Pattabiraman, Thomas Moscibroda, and Ben Zorn. Flikker: Saving DRAM Refresh-power through Critical Data Partitioning. In Proc. of ACM ASPLOS '11.
- [46] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In Proc. of IEEE CVPR '15.
- [47] David G. Lowe. 2004. Distinctive Image Features from Scale-Invariant Keypoints. Int. J. Comput. Vision 60, 2 (Nov. 2004).
- [48] Bruce D. Lucas and Takeo Kanade. 1981. An Iterative Image Registration Technique with an Application to Stereo Vision. In *Proceedings of the 7th International Joint Conference on Artificial Intelligence - Volume 2 (IJCAI'81).*
- [49] B. C. Matei, H. S. Sawhney, and S. Samarasekera. Vehicle tracking across nonoverlapping cameras using joint kinematic and appearance features. In *Proc.* of *IEEE CVPR* '11.
- [50] Å. Milan, K. Schindler, and S. Roth. 2016. Multi-Target Tracking by Discrete-Continuous Energy Minimization. *IEEE TPAMI* (2016).
- [51] Mohammad-Mahdi Moazzami, Dennis E. Phillips, Rui Tan, and Guoliang Xing, ORBIT: A Smartphone-based Platform for Data-intensive Embedded Sensing Applications. In Proc. of ACM IPSN '15.
- [52] Georg Nebehay and Roman Pflugfelder. 2015. Clustering of Static-Adaptive Correspondences for Deformable Object Tracking. In *Computer Vision and Pattern Recognition*. IEEE.
- [53] NYPD Plans to Put Body Cameras on All 23,000 Patrol Officers by 2019. NYPD Plans to Put Body Cameras on All 23,000 Patrol Officers by 2019.
- [54] Bharath Pichai, Lisa Hsu, and Abhishek Bhattacharjee. Architectural Support for Address Translation on GPUs: Designing Memory Management Units for CPU/GPUs with Unified Address Spaces. In Proc. ACM ASPLOS '14.
- [55] Jason Power, Mark D. Hill, and David A. Wood. Supporting x86-64 address translation for 100s of GPU lanes. In Proc. of IEEE HPCA '14.
- [56] Vignesh Ramanathan, Jonathan Huang, Sami Abu-El-Haija, Alexander N. Gorban, Kevin Murphy, and Li Fei-Fei. 2015. Detecting events and key actors in multiperson videos. *CoRR* abs/1511.02917 (2015). http://arxiv.org/abs/1511.02917
- [57] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. 2016. You Only Look Once: Unified, Real-Time Object Detection. In *The IEEE Conference* on Computer Vision and Pattern Recognition (CVPR).
- [58] Carsten Rother, Vladimir Kolmogorov, and Andrew Blake. 2004. "GrabCut": Interactive Foreground Extraction Using Iterated Graph Cuts. In ACM SIGGRAPH 2004 Papers (SIGGRAPH '04).
- [59] Ethan Rublee, Vincent Rabaud, Kurt Konolige, and Gary Bradski. 2011. ORB: An Efficient Alternative to SIFT or SURF. In Proceedings of the 2011 International Conference on Computer Vision (ICCV '11).
- [60] Jianbo Shi and Carlo Tomasi. 1994. Good features to track. In CVPR. IEEE, 593–600.
- [61] Karen Simonyan and Andrew Zisserman. Very Deep Convolutional Networks for Large-Scale Image Recognition. In Proc. of ICLR '15.
- [62] Chris Stauffer and W Eric L Grimson. 2000. Learning patterns of activity using real-time tracking. Pattern Analysis and Machine Intelligence, IEEE Transactions on 22, 8 (2000), 747–757.
- [63] S. Sunderrajan, Jiejun Xu, and B. S. Manjunath. Context-aware graph modeling for object search and retrieval in a wide area camera network. In *Proc. of ICSDC* '13.
- [64] Paul Viola and Michael J. Jones. 2004. Robust Real-Time Face Detection. Int. J. Comput. Vision 57, 2 (May 2004).
- [65] A. Viterbi. 2006. Error Bounds for Convolutional Codes and an Asymptotically Optimum Decoding Algorithm. *IEEE Trans. Inf. Theor.* 13, 2 (Sept. 2006).
- [66] Yu Xiang, Alexandre Alahi, and Silvio Savarese. 2015. Learning to Track: Online Multi-Object Tracking by Decision Making. In Proceedings of the IEEE International Conference on Computer Vision. 4705–4713.

[67] Tan Zhang, Aakanksha Chowdhery, Paramvir (Victor) Bahl, Kyle Jamieson, and Suman Banerjee. 2015. The Design and Implementation of a Wireless Video Surveillance System. In Proceedings of the 21st Annual International Conference

on Mobile Computing and Networking (MobiCom '15).
[68] Zhengyou Zhang. 2000. A Flexible New Technique for Camera Calibration. IEEE Trans. Pattern Anal. Mach. Intell. 22, 11 (Nov. 2000).