Smart City Intersections: Intelligence Nodes for Future Metropolises

Zoran Kostić*, Alex Angus*, Zhengye Yang**, Zhuoxu Duan*,
Ivan Seskar***, Gil Zussman*, Dipankar Raychaudhuri***
*Dept. of Electrical Engineering, Columbia University, New York City
**Rensselaer Polytechnic Institute

***Winlab, Rutgers

Abstract—Traffic intersections are the most suitable locations for the deployment of computing, communications, and intelligence services for smart cities of the future. The abundance of data to be collected and processed, in combination with privacy and security concerns, motivates the use of the edgecomputing paradigm which aligns well with physical intersections in metropolises. This paper focuses on high-bandwidth, lowlatency applications, and in that context it describes: (i) system design considerations for smart city intersection intelligence nodes; (ii) key technological components including sensors, networking, edge computing, low latency design, and AI-based intelligence; and (iii) applications such as privacy preservation, cloud-connected vehicles, a real-time "radar-screen", traffic management, and monitoring of pedestrian behavior during pandemics. The results of the experimental studies performed on the COSMOS testbed located in New York City are illustrated. Future challenges in designing human-centered smart city intersections are summarized.

I. INTRODUCTION

Smart cities should be built with the primary goal of providing social good as defined by local communities [1], [2]. Contemporary technologies provide a plethora of components to support human-centered design of future metropolises. Issues of privacy, security, and local data governance on one hand, and optimization of bandwidth, computational resources, and latency on the other hand, implicate traffic intersections as the best locations for smart city intelligence nodes.

Smart-city intersections are the key locations for emerging smart cities, since city dynamics can be supported by the interconnection and collaboration between neighboring intersection intelligence nodes. The nodes will be equipped with artificial intelligence (AI)-enabled edge-computing [3] and communications equipment to facilitate automated low-latency data harvesting, inference, and decision making. This will enable the development of technologies like cloud connected vehicles, vehicle to infrastructure communications, and advanced sensory-based tools for alerting pedestrians and assisting handicapped individuals. Future applications will require intense AI-enabled computation, very high communication bandwidths, and ultra-low latencies.

We report the results of research on low-latency real-time applications for smart city intersections in metropolises and architectures, components, and methods for building intelli-



Fig. 1: COSMOS pilot site with cameras and edge-cloud nodes.

gent intersection nodes. The research utilizes COSMOS, an experimental testbed located in New York City.

II. SMART CITY INTERSECTIONS

The focus of this paper is low-latency high-bandwidth applications for smart city intersections. We explore how to support privacy-preserving real-time applications such as collaborative control of cloud connected vehicles and active pedestrian alert and assistance; both require the use of a number of sensors including multiple high-resolution video cameras. One of the key tasks for video-based applications is to detect and track objects in an intersection with high accuracy. We explore methods to achieve real time in smart city intersection applications defined by end-to-end latencies under 33.3 milliseconds. This includes (i) sensor data acquisition; (ii) communication between end-users, sensors, and edge cloud; (iii) AI-based inference computation; and (iv) providing feedback to participants in the intersection. The advanced "radarscreen" application is intended to broadcast the positions and velocities of objects to intersection participants in real time.

A. Privacy

Smart-city implementations prior to 2022 indicate that privacy and data security are the key concerns impeding successful large-scale deployments. Privacy concerns are further

amplified when video recordings are a part of data acquisition and processing. The COSMOS research program has a strong community outreach component. This is exemplified by multi-year activities on running NSF REM and RET programs where teachers from Harlem and other New York City schools get training and participate in developing STEM educational material for students in underprivileged schools (https://www.cosmos-lab.org/outreach/, [4]). Our approach to privacy is to integrate local communities into the data governance process. We will develop technologies that enable the communities to define and control data acquisition and processing supported by edge computing and temporary data storage paradigms.

B. Real-Time Interactions

The most important goal of smart city deployments is to improve the safety of pedestrians and other participants. Even in the most congested cities it is desirable to replace human drivers with safer self-driven vehicles. This motivates the concept of cloud-connected vehicles that interact with city infrastructure to improve their ability to navigate, and requires exceptionally low closed loop latencies associated with security-critical real-time actions.

Real-Time for Safety-Critical Applications

Extracting intelligence that indicates a potential collision and providing feedback to vehicles or pedestrians presents computational and latency challenges. City street dynamics are determined by vehicles travelling at velocities between 0 and 100 kilometers per hour (km/h). If we consider a vehicle travelling at 10km/h, a typical speed of vehicles in congested intersections, the vehicle is moving at approximately 3 meters per second (m/s). If we divide 3m/s by the standard frame rate of conventional video, 30 frames per second, the result is a movement of 10cm, or the distance travelled by the vehicle in 33.3 milliseconds. If a vehicle's breaks could be activated in that time, it is conceivable that numerous life-threatening traffic accidents can ultimately be avoided. This approximate calculation leads us to target latencies below 33 milliseconds.

Sensor Latencies

Smart city sensors will have a wide range of operational frequencies and data acquisition bandwidths. CO_2 sensors may collect several bytes per hour, whereas high resolution cameras may stream data in compressed form at tens of Megabits per second, or in uncompressed form at several Gigabytes per second. Low-cost CMOS imaging sensors have latencies of several milliseconds, which are low enough not to obstruct the closed-loop target of 1/30 second. IP cameras use video encoding and streaming protocols that, because of inter-frame coding, may have buffers requiring hundreds of milliseconds to decode; this process severely impedes the ability to provide closed-loop services with less than 33.3 ms latencies.

Communications Latencies

Communications and networking latencies are determined as much by speed of physical media as they are driven by protocols at the application layer. The COSMOS optical network can provide up to 100 Gb/s, offering almost unlimited raw speed. On the other hand, conventional streaming of high resolution videos can create hundreds of milliseconds of latency. This suggests that video processing and inference is best done at the "extreme" edge - right next to the video sensor. More interestingly, this motivates research on integrated coding and video transmission protocols optimized for ultralow latency transmission of videos over high bandwidth edge communications infrastructure.

Inference and Decision Latencies

Inference latencies come from video preprocessing and deep learning algorithms for multiple object detection and tracking. The training of DL models is done offline and does not impact latencies for real-time interactions. Both published work and our own studies indicate that contemporary GPUs within specialized pipelines such as NVIDIA TensorRT and DeepStream can deliver speeds above 30 fps for object detection and tracking. We previously showed that inference speed varies as a function of input resolution and actual device capabilities, but we assess that inference computation will not be a bottleneck in meeting our real-time latency target.

The decision process is defined as a higher level of intelligence built on top of object detection and tracking. For example, this process would deduce the implications of a pedestrian being on a trajectory to intersect with a speedy vehicle and create a warning (or even a command) for the pedestrian or vehicle. Computational needs for this type of processes are subject to ongoing studies, but it is expected that the latencies will be less than a millisecond.

C. COSMOS Experimental Testbed

New York City (NYC) is an excellent example of a busy metropolis which provides formidable challenges for the deployment of smart city technologies. Busy urban traffic intersections have a large number of vehicles and pedestrians moving in many directions at various speeds, often with chaotic or unpredictable behavior. Furthermore, obstructions like building corners, parked vehicles, and construction equipment present difficulty to autonomous vehicle sensors requiring further advancements in traffic intersection based automation of monitoring, measuring, learning, and feedback.

The COSMOS testbed, NSF-funded Cloud Enhanced Open Software Defined Mobile Wireless Testbed for City-Scale Deployment [5], provides an experimentation platform for applications and architectures to support intelligence nodes of future metropolises. For our research, we use the COSMOS pilot site located at Columbia University, in New York City, at the intersection of the 120th Street and Amsterdam Avenue. The pilot node includes two street level and two bird's eye cameras, as illustrated in Fig.1. The COSMOS edge cloud servers can run real-time algorithms for detection and tracking of objects in the intersection to monitor and manage traffic flow and pedestrian safety. The node is equipped with an optical x-haul transport system that connects AI-enabled edge computing clusters. This allows for baseband processing with massively scalable CPU and GPU resources with FPGA



Fig. 2: COSMOS testbed camera views: (A) 1st-floor camera, 120th St; (B) 2nd-floor camera, Amsterdam Ave.; (C) 12th-floor camera, Amsterdam Ave.; (D) Calibrated 12th-floor camera.

assist, which can also support software defined radios. Four technology layers are provided for experimentation: the user device layer, radio hardware and front-haul network resources, radio cloud, and general purpose cloud.

III. BUILDING BLOCKS OF INTELLIGENT NODES

As of 2022, individual technological modules for implementing the vision of smart cities exist in the form of low power chips, high bandwidth modems, wired and wireless networks, and GPUs for machine learning (ML) and deep learning (DL). However, major challenges exist in the domains of privacy preservation, security, intelligent decision making, system integration, and in the interactions between technology and social good.

A. Sensors

Sensors range from dozens of low rate IoT-based devices collecting data about pollution to several high resolution lidars and cameras providing real-time feeds. Multi-modal data aggregation and collaborative intelligence are research topics of notable importance to smart intersection nodes [6].

B. Networking

For high bandwidth applications, networking at one intersection has to support wireless and wired connectivity from half a dozen infrastructure-installed cameras. Whereas coded video from a conventional IP-camera may require sub-hundred Mb/sec, experimentation with ultra low latency provides motivation to send raw video at several Gb/sec per camera. Support for cloud-connected vehicles could require harvesting videos and other data from each vehicle wirelessly, in either raw or meta format. Conventional video streaming protocols may be inadequate for accomplishing very low latencies, so research into edge-streaming protocols is an appealing topic.

C. Edge Computing

Smart city intersection applications require substantial computational resources, demand minimal latencies, and their functionality can be constrained to a limited geographical area. Furthermore, data privacy, security, and local data governance

are of utmost importance. This strongly implicates edge computing as the right modality. Two forms of edge computing can be used. In the extreme, AI-based computing can be done on devices located at the sensors such as Nvidia Jetson Nanos or ML-enabled ARM M1-M4 processors integrated into IoT chips. On the other hand, a more powerful computing node can be located in a facilities room of a building at the intersection. The node is then connected to sensors by high speed wireless, wired, or optical infrastructure. To support low latencies from sensors to actuators via AI computing, an edge computing node has to be integrated tightly with the network communications infrastructure.

D. AI-Enabled Data Processing

Intelligent tasks supporting smart city intersections are varied in complexity: CO_2 sensors generate several bytes once per hour, whereas high resolution cameras in our studies generate Megabits per second to be analyzed by visual deep learning models for object detection, tracking, and intelligent decisions for actuators. Automation and AI are crucial to scale systems for highly congested traffic intersections. Off the shelf AI models must be modified and retrained to accommodate the peculiarities of smart city intersection applications - one example being the detection of tiny pedestrians when viewed from bird's eye cameras.

Data Preprocessing

Visual deep learning tools require data preparation, labeling, and augmentation. The COSMOS pilot node contains low-elevation cameras and high-elevation bird's eye view cameras, each requiring different type of preprocessing (Fig. 2). The variation in angles and distances to the intersection, scale of objects, and overlapping field-of-views allow experimentation with the best view for a given application. For example, ground floor cameras are closer to traffic objects. They consequently provide more visual details for applications such as multicamera object reidentification, but are not as well suited to analyze large scale traffic patterns due to the scale distortion between objects at varying distances to the camera – the bird's eye view cameras offer a better perspective for this type of application.

High-elevation cameras allow us to perform calibration transforms to improve the effectiveness of deep learning models. See in Fig. 2 and Fig. 5 that the high-elevation camera view can be adjusted to appear perpendicular to the road by applying a homography transformation, after which resizing and cropping of the frame create the square aspect ratio required by many DL models. In our traffic intersection use case there are locations in the frame where relevant objects do not appear (i.e. no cars on building walls or pedestrians flying in the air). This motivates the creation of (black) masks overlayed on top of the frames, as seen in Fig. 3 and Fig. 5.

Supervised object detection and tracking models require a large number of precisely annotated ground truth labels to train the algorithms "by example". Producing accurate and consistent sets of labeled videos is difficult as both domain knowledge and significant amounts of time are needed. To

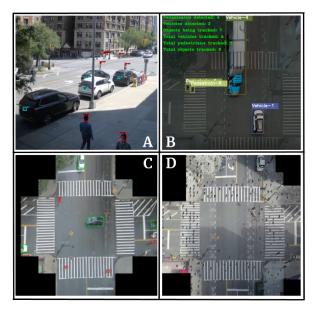


Fig. 3: (A) YOLOv4 detections of faces and license plates in ground floor video; (B) SORT tracking of vehicles and pedestrians in bird's eye video; (C) Bird's eye ground truth bounding box labels of intersection objects; (D) Pedestrian copy-paste data augmentation for improving detection of small objects

TABLE I: Object Detection Performance

Model	Pedestrian AP (%)	Vehicle AP (%)	mAP (%)	Inference Speed*
YOLOv4	66.31	97.58	81.95	34.99
SSD	57.04	94.81	75.93	11.31
RetinaNet	20.83	95.59	58.21	22.97

^{*}Inference speed (FPS) on NVIDIA T4 GPU.

support our experiments we annotated thousands of frames capturing the intersection in various weather, lighting, and congestion conditions.

Object detection models typically struggle with small object detection. Tiny pedestrians in the bird's eye camera view, as well as far-away license plates in the ground-floor camera view, convey very little information. This results in relatively poor detection and tracking accuracies. To improve the performance, we have deployed techniques of pretraining the DL models with a small-object dataset [7] and applying data augmentation techniques such as the copy/paste method illustrated in Fig. 3 (D).

Object Detection and Tracking

In smart traffic intersections, detecting pedestrians and vehicles and tracking their trajectories are the prerequisites for all downstream applications, Fig. 4. This involves two computer vision tasks: Object Detection and Multiple Object Tracking (MOT). The objective of object detection is to localize and classify objects within the frame. MOT aims to associate object identities across successive frames. State-of-the-art methods rely on deep learning blocks such as Convolutional Neural



Fig. 4: Pedestrian and Vehicle detection on 120th St. and Amsterdam Ave.

Networks (CNN) [8] and Vision Transformers [9]. These methods bring heavy computational cost, and the accuracyspeed trade-off - the budgeting between computational complexity and inference speed - is vital to the success of smart city applications. With this consideration in mind, we experimented with a series of algorithms for detecting and tracking objects to find the best approach [10] based on our custom annotated dataset for bird's eye videos. We choose YOLOv4 [11] as the base detector for all downstream applications since it is able to provide accurate results in real-time. Object detection performance is shown in Table I, where the Average Precision (AP) and mean Average precision (mAP) are used as the evaluation metrics. On our bird's eye view intersection data, YOLOv4 outperforms both RetinaNet [12] and SSD [13] in terms of AP and inference speed, where inference speed is measured as the average time for a forward pass through the model with batch size equal to 1. For MOT, different scenarios need to be considered separately. For bird's eye cameras, object occlusions barely occur, so re-identification (reID) calculation is not as necessary as for the ground level cameras. The reID calculation is often the computation bottleneck in MOT algorithms. "Simple Online and Realtime Tracking" (SORT) and "Simple Online and Realtime Tracking with a Deep Association Metric" (DeepSORT) suffice for the bird's eye view cameras. Illustrations for detection are shown in Fig. 3.

Image Resolution and Object Density

Highly elevated bird's eye cameras have a good view of the overall scene, shown in Fig. 2. Pedestrians, which appear small, become a problem for object detection and tracking. Intuitively, the higher the resolution of the input image, the more object features can be preserved. However, higher resolution leads to a larger computational cost, thus making the inference slow. We tested a dozen combinations of image input resolutions and aspect ratios to find the best balance between accuracy and speed, three of which are shown in Fig. 5. Some deep learning models, like YOLOv4 [11], perform better on input images with a fixed-sized, square aspect ratio. To maximize the preservation of important features of the intersection scene and to minimize the irrelevant components, the experiments indicate that the "squared cropped" 832×832

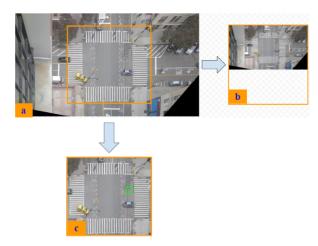


Fig. 5: (a) Calibrated 16:9 native frame; (b) 16:9 frame squared using zero-padding; (c) Square cropped frame.

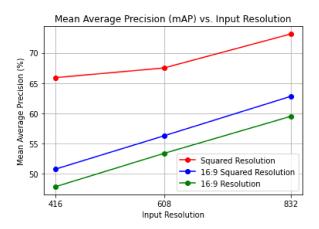


Fig. 6: mAP for pedestrians and vehicles, 9 cases of image resolution vs. aspect ratio

input produces the best results [14].

Object density refers to the number of objects in a scene, which may impact the speed of inference as the busyness of the streets change through the day. We explored the inference time for ten 90-second videos where the number of objects varied from 4,000 to 26,000. The results show approximately 40 percent increase in computational load from the lowest to the highest density case. This is important in that it shows that object density can be used to switch between computational resources to obtain the optimal power/accuracy balance.

IV. APPLICATIONS

Advances in video based object detection and tracking have enabled the deployment of a number of traffic intersection applications, where one can identify the locations of objects in the intersection and classify them by type of vehicle, pedestrian, bicycle, etc. They can be tracked as unique entities which persist through the duration of traffic cycles, different camera views, and times of the day, week, month, and year. The abundance of spatial, temporal, and visual data makes it



Fig. 7: Input (left) and output (right) of the face and license plate blurring pipeline.

possible to perform data anonymization, quantization of traffic trends, crowd behavior surveillance, real time intersection radar mapping, and more.

A. Privacy Protection - Face and License Plate Anonymization

Collecting real-time images and videos of public spaces from street level inadvertently involves capturing sensitive information such as faces and license plates. To avoid leaking private information with our datasets, we generated a pipeline to automatically blur these sensitive areas. We trained several object detection models on a custom labeled dataset to detect faces and licenses for subsequent anonymization. When training with sequential video datasets, it is important to leave entire videos out of the training process to use for validation. Stationary objects – parked cars, seated pedestrians, chained bicycles, etc. – occur identically in many frames, and model evaluation on these stationary objects yields biased results. This leads to model overfitting and poor generalization to new intersection scenes, which has to be addressed.

Fig. 7 shows an example input and output frame of the anonymization pipeline. For our face and license detection model, we chose YOLOv4 [11] for its compromise between detection accuracy and inference speed. For privacy critical applications, the most relevant performance measure is recall, the number of relevant faces and licenses that are detected out of the total number that pass through the frame. False positives are less of an issue than false negatives, as they result in an extra blurred area of the frame, but not a privacy leak. In our case, not all faces and licenses are "relevant" - some are too far away and too low resolution to be identifiable. We exclude these instances from the recall evaluation by defining pixel area thresholds below which the objects are ignored. We found that, below certain thresholds, facial features and license plate characters could not be reliably identified. While there exist information reconstruction techniques that could potentially recover these features, this is outside the scope of this project to consider them. Furthermore, we would need to reconsider our choice of anonymization as any form of blurring becomes ineffective. In the visible object evaluation our pipeline blurs over 99% of visible faces and licenses and in the total evaluation it blurs over 96% of objects greater than 100 pixels.

To increase our confidence in the anonymization pipeline, we performed manual evaluations by inspecting anonymized



Fig. 8: Normal blurring pipeline detections (top) vs. edge cases (bottom).

output videos for misses, where a miss is defined as an object with more than a quarter of the face or license plate exposed. The results of the manual evaluations confirmed the results of the programmatic evaluations and shed some light on edge cases where our models consistently missed, Fig. 8. Most edge cases were due to occlusions such as occluded borders of license plates, pedestrian body occlusion, and tree branch occlusion, resulting in consistent false negatives. More data collection and training is needed to rectify these edge cases.

B. Counting Objects

An important goal for smart intersections is to analyze traffic flow in real time. To this end, we use detection and tracking to classify and count vehicles and pedestrians and follow their paths through the intersection. Accumulation of the tracks provides sufficient data for traffic trend analyses that can be used to optimize traffic flow and improve pedestrian safety in the intersection.

To perform object tracking we use the detection based (MOT) algorithm DeepSORT. DeepSORT requires an object detection model to provide the locations and features of an object to be tracked. Given detections of vehicles and pedestrians, DeepSORT uses a Kalman filter to map detections with similar sizes and motions across frames of a video. In this way we can assign IDs to detected objects that persist throughout multiple video frames. Additionally, DeepSORT uses visual features of the object to increase the reliability of the tracking. Even if the object fails to be detected in consecutive frames, it can be assigned to the correct track by the re-identification model (reID) based on its visual features.

Though DeepSORT is a robust tracking system, it is still dependent on high quality object detection. If an object is not detected or misclassified for multiple consecutive frames, it will be regarded by the algorithm as a "new track" – the old track disappears and a new one is created upon redetection. For vehicles, we achieve consistent high accuracy detection,

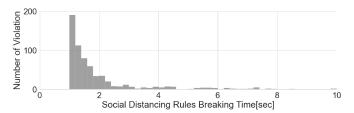


Fig. 9: B-SDA: Distribution of the duration of social distancing violations.

and corresponding high accuracy tracking, but for pedestrians, which have 4-5x smaller cross sections, high accuracy detection is a more difficult task. Pedestrian tracking accuracy suffers as a result of lower accuracy pedestrian detection. Data augmentation techniques such as the copy-paste pedestrian method shown in Fig. 3 (D) and pretraining object detectors on small-object datasets show improvement for small-object detection, but pedestrian detection and tracking accuracies are still lower than for vehicles, with multiple object tracking accuracies (MOTA) of 75.16% and 18.23% for vehicles and pedestrians, respectively.

The vehicle tracking performance is sufficient for applications that quantify traffic flow. For example, in an automatic counting task we record vehicles passing through the intersection as turning right, turning left, or going straight from all four directions with an accuracy of 95% evaluated over 21 minutes of video recording.

C. Social Distancing in Pandemics

Smart cities can assist in combating global pandemics, such as COVID-19, by providing means for monitoring, analyzing, and potentially controlling social distancing behavior. We proposed several techniques and applied them to video datasets collected at the COSMOS pilot intersection.

The fundamental idea is to estimate distances between pedestrians and compare them against the recommended minimal distance threshold. The first step is to detect the pedestrians. The real-world distance is then estimated by calculating the pixel-wise distance between pedestrians within one frame. The tracking of pedestrians between frames facilitates the calculation of higher order statistics, related to safe social distancing groups, which are more meaningful than an individual-to-individual social distancing violation rates. When acquaintances are walking together on the street as a "safe group", the intra-group distance is often smaller than the social distancing threshold, which triggers the indication of the violation. To solve this problem, we utilize the pedestrian trajectory similarity and stability, which can evaluate the motion dynamic between every pedestrian pair. This group validation approach is able to significantly reduce the number of false positive violations, achieving the F1 score of 0.92. Based on this approach, we built a social distancing analysis system B-SDA [15] for bird's eye view cameras, as well as a complementary method Auto-SDA [16], [17] with ground level cameras.

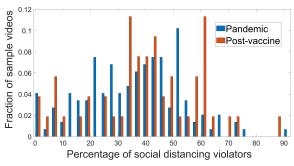


Fig. 10: Auto-SDA: Normalized histogram of the percentage of social distancing violations.

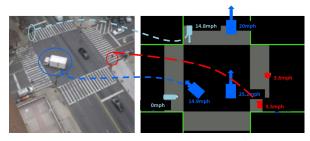


Fig. 11: The "radar screen": one frame of a video containing locations and velocities of objects within an intersection.

An example of the results obtained with the bird's eye video dataset, illustrated in Fig. 9, shows the distribution of the duration of social distancing violations during the Covid-19 pandemic. Fig. 10 shows the social distancing violation rates for the ground-floor camera dataset (i) during the pandemic and (ii) after the vaccine is widely available. Detailed analyses and comparisons of multiple statistics before the pandemic and during the pandemic demonstrate that the proposed systems can reliably identify social distancing violations.

D. Real Time "Radar-Screen"

The "Radar-screen" application aims to infer positions and velocities of objects within a traffic intersection and broadcast them to the participants in the intersection, as illustrated in Fig. 11. The information can be distributed in raw or coded/meta format. The application intends to provide a real-time service with latency of 1/30 seconds between the observation of objects and the wireless broadcast delivery. As described previously, this is motivated by the approximation of a 10 centimeter vehicle movement with speed of 10 km/h.

The application includes the acquisition of videos from surrounding buildings, potential harvesting of videos (or encoded data) from cameras within vehicles, harvesting of IoT sensor data, transmission via a high speed network to the inference computer, data aggregation and preprocessing, DL-based object detection and tracking, extraction of information at a higher abstraction level, and (in a more advanced version) deduction of commands that may be issued to individual vehicles after optimizing the traffic flow. The final step is the broadcasting of information. This is an aspirational applica-

tion in that achieving the cumulative latency of 33.3 ms is technologically challenging. Balancing between computational capabilities, power consumption, and latency minimization of the extreme edge compute units, or edge computing centers, requires rapid sensor data acquisition and dynamic network and resource control. This application motivates research to optimize each of the building blocks described in previous sections of this paper as well latency-focused cross-module system integration.

E. Traffic Management

Intelligent nodes located at individual intersections provide powerful data acquisition and intelligent edge-computing. On a larger scale, smart cities require the aggregation of data from multiple intersections and mutual coordination. In that vein, we have commenced collaborative studies with traffic engineering experts on the definition of key parameters such as timing resolution, sensor locations, and APIs for data exchange between intelligent smart intersection nodes and traffic optimization systems [18]. We are building simulators and defining digital twins that will play predictive roles in the behavior of individual traffic participants and in global optimization of traffic management.

V. CONCLUSION AND FUTURE CHALLENGES

A vision of the smart city intersection as the intelligence node for future metropolises has been presented. The proposed architecture is driven by societal needs to preserve privacy, which strongly implicate edge computing and intelligence as the key paradigm for data management and processing. Key technological components have been reviewed such as sensors, networks, and edge AI computing. Real time needs of future safety-critical systems have been examined, and design considerations for a "radar-screen" application, which closes the loop from sensors to actuators, have been summarized. The requirements for low latency, based on the 33.3 ms target, have been explored. System integration challenges have been illustrated using the examples from experiments performed on the pilot node of the COSMOS testbed in New York City.

Our research points to the following exploration topics: (i) State of the art DL-based object detection models are comprised of over 60 million parameters and require passing more than 100 convolutional layers, where each convolution has $O(n^4)$ complexity. Model optimization techniques like weight pruning, inference scheduling, and neural algorithmic search strategies [19] need to be incorporated into practical systems; (ii) Reliance on supervised datasets for video processing is not scalable due to the labeling cost and quality concerns. This necessitates research on unsupervised learning methodologies which should be based on continuous or active learning, and take advantage of the peculiarities of the fixed scene within a traffic intersection [20]; (iii) Data fusion from multiple cameras is expected to yield notable improvements in detection and tracking accuracies; (iv) Achieving low latency for low rate little-data applications is possible by using processing on the "extreme edge", but meeting the requirements of 1/30 second latency for high resolution videos is a challenge. New video coding methods and streaming protocols should be explored with focus on localized low-latency performance.

VI. ACKNOWLEDGMENT

This work was supported in part by NSF grants CNS-1827923, OAC-2029295, CNS-2038984, CNS-1910757, and AT&T VURI award.

REFERENCES

- [1] L. Sánchez, L. Muñoz, J. A. Galache, P. Sotres, J. R. Santana, V. Gutiérrez, R. Ramdhany, A. D. Gluhak, S. Krco, E. Theodoridis, and D. Pfisterer, "Smartsantander: Iot experimentation over a smart city testbed," *Comput. Networks*, vol. 61, pp. 217–238, 2014.
- [2] L. Belli, A. Cilfone, L. Davoli, G. Ferrari, P. Adorni, F. Di Nocera, A. Dall'Olio, C. Pellegrini, M. Mordacci, and E. Bertolotti, "Iot-enabled smart sustainable cities: Challenges and approaches," *Smart Cities*, vol. 3, no. 3, pp. 1039–1071, 2020.
- [3] A. Y. Ding, E. Peltonen, T. Meuser, A. Aral, C. Becker, S. Dustdar, T. Hiessl, D. Kranzlmüller, M. Liyanage, S. Maghsudi, N. Mohan, J. Ott, J. S. Rellermeyer, S. Schulte, H. Schulzrinne, G. Solmaz, S. Tarkoma, B. Varghese, and L. Wolf, "Roadmap for edge ai: A dagstuhl perspective," SIGCOMM Comput. Commun. Rev., vol. 52, p. 28–33, mar 2022.
- [4] P. Skrimponis, N. Makris, S. B. Rajguru, K. Cheng, J. Ostrometzky, E. Ford, Z. Kostic, G. Zussman, and T. Korakis, "Cosmos educational toolkit: Using experimental wireless networking to enhance middle/high school stem education," SIGCOMM Comput. Commun. Rev., vol. 50, p. 58–65, oct 2020.
- [5] D. Raychaudhuri, I. Seskar, G. Zussman, T. Korakis, D. Kilper, T. Chen, J. Kolodziejski, M. Sherman, Z. Kostic, X. Gu, et al., "Challenge: Cosmos: A city-scale programmable testbed for experimentation with advanced wireless," in *Proc. ACM MobiCom*, 2020.
- [6] X. Xu, Q. Huang, X. Yin, M. Abbasi, M. R. Khosravi, and L. Qi, "Intelligent offloading for collaborative smart city services in edge computing," *IEEE Internet Things J.*, vol. 7, no. 9, pp. 7919–7927, 2020.
- [7] P. Zhu, L. Wen, D. Du, X. Bian, H. Ling, Q. Hu, Q. Nie, H. Cheng, C. Liu, X. Liu, et al., "Visdrone-det2018: The vision meets drone object detection in image challenge results," in *Proc. ECCV Workshops*, 2018.
- [8] L. Liu, W. Ouyang, X. Wang, P. Fieguth, J. Chen, X. Liu, and M. Pietikäinen, "Deep learning for generic object detection: A survey," *Int. J. Comput. Vis.*, vol. 128, no. 2, pp. 261–318, 2020.
- [9] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al., "An image is worth 16x16 words: Transformers for image recognition at scale," arXiv preprint:2010.11929, 2020.
- [10] S. Yang, E. Bailey, Z. Yang, J. Ostrometzky, G. Zussman, I. Seskar, and Z. Kostic, "COSMOS smart intersection: Edge compute and communications for bird's eye object tracking," in *Proc. SmartEdge*, 2020.
- [11] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "YOLOv4: Optimal speed and accuracy of object detection," arXiv preprint:2004.10934, 2020.
- [12] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proc.ICCV*, 2017.
- [13] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "SSD: Single shot multibox detector," in *Proc. ECCV*, 2016.
- [14] Z. Duan, Z. Yang, R. Samoilenko, D. S. Oza, A. Jagadeesan, M. Sun, H. Ye, Z. Xiong, G. Zussman, and Z. Kostic, "Smart city traffic intersection: Impact of video quality and scene complexity on precision and inference," in *Proc. IEEE Smart City* '21, 2021.
- [15] Z. Yang, M. Sun, H. Ye, Z. Xiong, G. Zussman, and Z. Kostic, "Birds eye view social distancing analysis system," arXiv preprint:2112.07159, 2021.
- [16] M. Ghasemi, Z. Yang, M. Sun, H. Ye, Z. Xiong, Z. Kostic, and G. Zussman, "Demo: Video-based social distancing evaluation in the COSMOS testbed pilot site," in *Proc. ACM MOBICOM'21*, 2021.
- [17] M. Ghasemi, Z. Kostic, J. Ghaderi, and G. Zussman, "Auto-SDA: Automated video-based social distancing analyzer," in *Proc. ACM Hot-EdgeVideo*, 2021.

- [18] G. Karagiannis, O. Altintas, E. Ekici, G. Heijenk, B. Jarupan, K. Lin, and T. Weil, "Vehicular networking: A survey and tutorial on requirements, architectures, challenges, standards and solutions," *IEEE Commun. Surveys Tuts.*, vol. 13, no. 4, pp. 584–616, 2011.
- [19] C. R. Banbury, V. J. Reddi, M. Lam, W. Fu, A. Fazel, J. Holleman, X. Huang, R. Hurtado, D. Kanter, A. Lokhmotov, et al., "Benchmarking tinyml systems: Challenges and direction," arXiv preprint:2003.04821, 2020.
- [20] Z. Dai, G. Wang, S. Zhu, W. Yuan, and P. Tan, "Cluster contrast for unsupervised person re-identification.," arXiv preprint:2103.11568, 2021