

Real-Time HAP-Assisted Vehicular Edge Computing for Rural Areas

This is the peer reviewed version of the following article:

Original:

Traspadini, A., Giordani, M., Giambene, G., Zorzi, M. (2023). Real-Time HAP-Assisted Vehicular Edge Computing for Rural Areas. IEEE WIRELESS COMMUNICATIONS LETTERS, 12(4), 674-678 [10.1109/LWC.2023.3238851].

Availability:

This version is available <http://hdl.handle.net/11365/1228335> since 2023-05-13T12:42:31Z

Published:

DOI: <http://doi.org/10.1109/LWC.2023.3238851>

Terms of use:

Open Access

The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. Works made available under a Creative Commons license can be used according to the terms and conditions of said license.

For all terms of use and more information see the publisher's website.

(Article begins on next page)

Real-Time HAP-Assisted Vehicular Edge Computing for Rural Areas

Alessandro Traspadini, *Student Member, IEEE*, Marco Giordani, *Member, IEEE*,
Giovanni Giambene, *Senior Member, IEEE*, Michele Zorzi, *Fellow, IEEE*

Abstract—Non-Terrestrial Networks (NTNs) are expected to be a key component of 6th generation (6G) networks to support broadband seamless Internet connectivity and expand the coverage even in rural and remote areas. In this context, High Altitude Platforms (HAPs) can act as edge servers to process computational tasks offloaded by energy-constrained terrestrial devices such as Internet of Things (IoT) sensors and ground vehicles (GVs). In this paper, we analyze the opportunity to support Vehicular Edge Computing (VEC) via HAP in a rural scenario where GV can decide whether to process data onboard or offload them to a HAP. We characterize the system as a set of queues in which computational tasks arrive according to a Poisson arrival process. Then, we assess the optimal VEC offloading factor to maximize the probability of real-time service, given latency and computational capacity constraints.

Index Terms—6G, non-terrestrial networks (NTNs), HAP, vehicular edge computing (VEC), optimization.

I. INTRODUCTION

One of the grand objectives of 6th generation (6G) wireless networks [1] is to enhance broadband Internet coverage in remote areas [2]. The main issues towards this goal are the high costs of terrestrial deployments (especially for fiber-optic backhaul roll-out) as well as the lack of power sources and transport connectivity in rural regions. To overcome these problems, 6G will promote Non-Terrestrial Networks (NTNs) with Unmanned Aerial Vehicles (UAVs), High Altitude Platforms (HAPs), and satellites to provide fast connectivity thanks to their flexibility and inherent support for large coverage [3]. For example, in 2020 a network of HAPs was deployed by Loon to bring Internet connectivity to unserved regions in Kenya. Similarly, in 2021, Airbus has proved the suitability of the solar-powered Zephyr HAP to provide direct-to-device connectivity in the rural areas of Arizona, US.

Besides providing connectivity, NTNs can also host edge servers for processing, caching, and/or storing data generated from power-constrained Internet of Things (IoT) devices, thereby supporting the transition of remote areas towards smart cities [4]. Similarly, air/space-borne platforms can support Vehicular Edge Computing (VEC), where ground vehicles (GVs) offload computationally-intensive autonomous driving tasks like object detection and recognition, tracking/trajectory prediction, and/or semantic segmentation [5]. Processing these tasks onboard certain vehicles may be infeasible due to the limited autonomy and to the small capacity of their batteries. While in urban areas GV can delegate the burden of data

processing to roadside units [6], NTNs may represent a valid alternative to serve computational requests in poorly connected rural areas. One of the main requirements of VEC systems is the support for real-time latency to guarantee safe autonomous driving as specified by the standards [7]. While satellites incur large round-trip times, UAVs offer limited computational capacity and short-lived connectivity [8]. Instead, HAPs provide large coverage, and are deployed in the stratosphere at an altitude of 20 km, which helps reduce the propagation delay to less than 1 ms. The effectiveness of HAPs in Mobile Edge Cloud (MEC) networks was studied in [9], where the authors addressed the problems of partial offloading and bandwidth allocation in remote areas without the coverage of ground infrastructure. However the benefit of HAPs for VEC is still an open question, that we will analyze in this work.

Based on the above introduction, in this paper we present an optimization problem for real-time VEC in a rural scenario, in which multiple GV can decide whether to process perception data (generated in the form of video frames) onboard or offload them to a HAP. We consider wireless transmissions in the millimeter wave (mmWave) band [10]. With respect to our previous work [11], we develop a more accurate model to characterize both GV's and HAP's queuing systems, and determine a closed-form expression for the average waiting time experienced by computing tasks. We show that it is convenient to offload some VEC processing tasks from GV to the HAP under common assumptions for the processing capacity and load of the nodes.

II. SYSTEM MODEL

A. Problem Formulation

We analyze a vast remote region where n GV are distributed within an Area of Interest (AoI) of size A . Each GV generates sensory perceptions in the form of video frames of size n_{UL} according to a Poisson process with arrival rate r , as considered in some reference papers to model autonomous driving data [12]. Each frame requires a constant computational load C (e.g., for object detection [13]), which has to be executed fast enough to provide real-time services, i.e., at least at the frame rate of the sensor.

We consider that the GV can offload a subset η of their computational load to a HAP providing VEC functionalities. According to the split property of the Poisson processes, we can distinguish two independent Poisson processes, namely for frames that are offloaded to the HAP at a rate ηr , and for frames that are processed onboard at a rate $(1 - \eta)r$. On one side, local processing may involve long delays, given the limited computational capacity C_{GV} of low-budget car models. On the other side, the HAP can count on a higher computing

Alessandro Traspadini, Marco Giordani and Michele Zorzi are with the Department of Information Engineering, University of Padova, Padova, Italy (emails: {traspadini, giordani, zorzi}@dei.unipd.it). Giovanni Giambene is with the Department of Information Engineering and Mathematical Sciences, University of Siena, Siena, Italy (email: giovanni.giambene@unisi.it).

power C_{HAP} than available onboard a vehicle, i.e., $C_{\text{HAP}} \geq C_{\text{GV}}$, thus processing data faster, at the expense of a non-negligible communication delay for offloading data to the HAP server and for the processed data to be returned to the end users. We will investigate the optimized selection of $\eta \in [0, 1]$.

B. Delay Model

This section introduces our delay model for each data frame for both onboard and HAP-assisted processing.

a) *Onboard processing*: Each GV is modeled as an M/D/1 queue, and the average delay is equal to

$$\bar{t}_{\text{GV}} = \bar{W}_q(\text{M/D/1}) + (C/C_{\text{GV}}), \quad (1)$$

where $\bar{W}_q(\text{M/D/1})$ is the average waiting time (see Sec. II-D), while C/C_{GV} is the frame processing time onboard the GV.

b) *HAP-assisted processing*: Since the HAP has less severe space and energy constraints than a GV, it is modeled as an M/D/ c queue, where c servers can process up to c frames in parallel. The average delay for processing tasks can be expressed as

$$\bar{t}_{\text{HAP}} = 2\tau_p + t_{\text{UL}} + t_{\text{DL}} + \bar{W}_q(\text{M/D}/c) + (C/C_{\text{HAP}}), \quad (2)$$

where $\bar{W}_q(\text{M/D}/c)$ is the average waiting time (see Sec. II-D), $\tau_p = d/c_l$ is the propagation delay (where d is the distance between the generic GV and the HAP, and c_l is the speed of light), and t_{UL} and t_{DL} are the times to transmit data to and from the HAP, respectively (see Sec. II-C).

C. HAP-to-Ground Channel Model

According to the 3GPP specifications [14], HAPs may operate through highly directional links at mmWaves, where the large bandwidth available at these frequencies offers the potential for ultra-fast connectivity [10]. With the assumption of an interference-free environment considering transmissions over orthogonal bands, the Signal to Noise Ratio (SNR) between transmitter i and receiver j can be calculated as:

$$\gamma_{ij} = \frac{\text{EIRP}_i \cdot (G/T)_j}{\text{PL}_{ij} \cdot k \cdot B} \quad (3)$$

where EIRP_i is the effective isotropic radiated power, $(G/T)_j$ is the receiver antenna-gain-to-noise-temperature, PL is the path loss, k is the Boltzmann constant, and B is the bandwidth. The path loss depends on the frequency (here mmWaves) and on the distance between i and j , and accounts for additional atmospheric attenuations as described in [15].

Based on the SNR in Eq. (3), we introduce the median ergodic capacity R between the generic GV and the HAP as $R = B \log_2(1 + \gamma_{ij})$. Then, the transmission times t_{UL} and t_{DL} in Eq. (2), for each data frame, are given by $t_\ell = n_\ell/R_\ell$, $\ell \in \{\text{UL}, \text{DL}\}$, where n_ℓ is the size of the transmitted data.

D. Queuing Model

Sensory frames are generated at rate r , and are queued and eventually processed at the GVs and HAP according to a Poisson process of rate λ , equal to $(1 - \eta)r$ and ηrn , respectively; we follow a first-come-first-served (FCFS) discipline. Notably,

GVs operate via an M/D/1 queue, while HAPs via an M/D/ c queue. Following the analysis made in [16], the M/D/ c state probability p_j for state j , i.e., the probability that the queue at the HAP has j frames, can be computed by observing that any frame in the system at time $t + 1/\mu$ (where $1/\mu$ is the constant processing time) is either waiting in the queue at time t or has arrived in $[t, t + 1/\mu]$. Thus, given that the number of arrivals in an interval of duration $1/\mu$ is Poisson distributed with mean λ/μ , the state probability p_j is:

$$p_j = e^{-G} \frac{G^j}{j!} \sum_{k=0}^c p_k + e^{-G} \sum_{k=c+1}^{c+j} p_k \frac{G^{j-k+c}}{(j-k+c)!}, \quad (4)$$

where $G = \lambda/\mu$ is the offered traffic, and μ is the service rate at the HAP and GVs, i.e., C_{HAP}/C and C_{GV}/C , respectively. Eq. (4), combined with the normalization condition $(\sum_{j=0}^{+\infty} p_j = 1)$, generates an infinite system of linear equations that describe the state probabilities of the queue. We adopted the geometric tail approximation proposed in [16] to find a closed-form expression for this system of equations. Therefore, we can write

$$p_j = p_M \tau^{-(j-M)}, \quad j \geq M, \quad (5)$$

where τ is a constant depending on the server utilization, and M is a threshold state; we refer the interested reader to [16] to clarify how to compute these values. Based on Eq. (5), the infinite system in Eq. (4) can be reduced to a finite system of $M + 1$ linear equations that can be solved easily.

Lemma 1. *Based on Eq. (5), an empirical closed-form expression for the average waiting time of an M/D/ c queue with arrival rate λ and service rate μ is given by*

$$\bar{W}_q(\text{M/D}/c) = \frac{\sum_{k=c}^{M-1} p_k (k - c) + p_M \left[\frac{M - c + \frac{1}{1 - \frac{1}{\tau}} - 1}{1 - \frac{1}{\tau}} \right]}{\lambda}, \quad (6)$$

where a relatively small value of M can already provide an accurate approximation of $\bar{W}_q(\text{M/D}/c)$.

Proof. See Appendix. ■

III. OPTIMAL OFFLOADING FOR HAP-ASSISTED VEC

In our scenario, we require that autonomous driving tasks, especially object detection on the video frames, be executed at least at the same rate at which new frames are generated, i.e., $t_{\text{max}} = 1/r$. We model this aspect through the probability of real-time service P_{RT} , given by

$$P_{\text{RT}}(\eta) = \eta P(t_{\text{HAP}} \leq t_{\text{max}}) + (1 - \eta) P(t_{\text{GV}} \leq t_{\text{max}}), \quad (7)$$

where t_{max} is the maximum tolerable latency, and t_{GV} and t_{HAP} have been defined in Eqs. (1) and (2), respectively. Given the rural scenario, we assume that all GVs experience the same channel condition. Thus, we assume they follow the same policy, and offload data frames with the same probability η (*shared offloading factor*). The value of $P(t_{\text{HAP}} \leq t_{\text{max}})$ in Eq. (7) can be expressed by computing the maximum number of frames in the system that ensure real-time processing, that is

$$f_{\text{max}}^{\text{HAP}} = c \left\lfloor \frac{t_{\text{max}} - t_{\text{UL}} - t_{\text{DL}} - 2\tau_p}{C/C_{\text{HAP}}} \right\rfloor. \quad (8)$$

However, real-time processing is also possible when new offloading requests find the queue at full capacity, but some processing tasks are already in service and will leave the system soon. We define as Δ_t^{HAP} the percentage of in-service tasks that ensure real-time processing even if more than f_{\max}^{HAP} requests are in the system. We have

$$\Delta_t^{\text{HAP}} = \Gamma\left(\frac{t_{\max} - t_{\text{UL}} - t_{\text{DL}} - 2\tau_p}{C/C_{\text{HAP}}}\right), \quad (9)$$

where $\Gamma(x) = x - \lfloor x \rfloor$. Similarly, we can solve $P(t_{\text{GV}} \leq t_{\max})$ in Eq. (7) by introducing f_{\max}^{GV} and Δ_t^{GV} for the GV:

$$f_{\max}^{\text{GV}} = \left\lfloor \frac{t_{\max}}{C/C_{\text{GV}}} \right\rfloor, \quad \Delta_t^{\text{GV}} = \Gamma\left(\frac{t_{\max}}{C/C_{\text{GV}}}\right). \quad (10)$$

Eq. (7) depends on the queuing state probability of both GVs and HAP, based on the model in Sec. II-D. Using the Poisson Arrivals See Time Averages (PASTA) property, and the probability distribution described in Eqs. (4) and (5), the real-time probability for the HAP can be expressed as

$$P(t_{\text{HAP}} \leq t_{\max}) = \sum_{k=1}^c p_{\text{bin}}(k, c, \Delta_t^{\text{HAP}}) \sum_{j=0}^{k-1} p_{f_{\max}^{\text{HAP}}+j} + \sum_{i=0}^{f_{\max}^{\text{HAP}}-1} p_i, \quad (11)$$

where $p_{\text{bin}}(k, c, \Delta_t^{\text{HAP}}) = \binom{c}{k} (\Delta_t^{\text{HAP}})^k (1 - \Delta_t^{\text{HAP}})^{c-k}$ describes the probability that k offloading requests have been processed by at least a fraction Δ_t^{HAP} , and probabilities p_i are as in Eq. (4). Similarly, for the GV we have

$$P(t_{\text{GV}} \leq t_{\max}) = \sum_{i=0}^{f_{\max}^{\text{GV}}-1} p_i + \Delta_t^{\text{GV}} p_{f_{\max}^{\text{GV}}}. \quad (12)$$

The objective of our optimization problem is to choose the optimal offloading factor η , i.e., η^* , that maximizes P_{RT} , i.e.,

$$\arg \max_{\eta} P_{\text{RT}}(\eta), \quad (13a)$$

$$\text{subject to } G_{\text{HAP}} < c, \quad G_{\text{GV}} < 1 \quad (13b)$$

$$\eta \in [0, 1], \quad (13c)$$

where $G_{\text{HAP}} = \eta r C n / C_{\text{HAP}}$ and $G_{\text{GV}} = (1 - \eta) r C / C_{\text{GV}}$ are the offered traffic at the HAP and the GVs, respectively, while Eq. (13b) sets a constraint on η to keep the system stable. Therefore, the optimization problem can be rewritten as

$$\arg \max_{\eta} P_{\text{RT}}(\eta), \quad (14a)$$

$$\text{subject to } \eta \in [\eta_{\min}, \eta_{\max}], \quad (14b)$$

where $\eta_{\min} = \max(0, 1 - C_{\text{GV}}/(rC))$ and $\eta_{\max} = \min(1, cC_{\text{HAP}}/(rnC))$. Notably, C_{GV} (C_{HAP}) set a limit to the lower (upper) bound of η . Given that $P_{\text{RT}}(\eta)$ is a continuous function defined on a closed interval $[\eta_{\min}, \eta_{\max}]$, there exists a solution to the optimization problem in Eq. (14) according to Weierstrass' theorem. Then, the problem is solved numerically using the Brent solver [17].

We can analyze the average time it takes for a GV, which may offload computational tasks to the HAP, to process each video frame. Since, on average, each GV processes a frame in

\bar{t}_{GV} , and the latency for an offloaded frame is \bar{t}_{HAP} , then the average latency is evaluated as

$$\bar{t}(\eta) = \eta^* \bar{t}_{\text{HAP}} + (1 - \eta^*) \bar{t}_{\text{GV}}, \quad (15)$$

where η^* is the optimal offloading probability from Eq. (14).

IV. PERFORMANCE EVALUATION

In Sec. IV-A below, we introduce our simulator and its parameters, while numerical results are provided in Sec. IV-B.

A. Simulation Setup and Parameters

We analyze a rural/remote scenario with $n \in \{1, \dots, 200\}$ GVs, uniformly distributed over a large AoI of 1000 km². The coverage area of the HAP is large enough to ensure that GVs are always and continuously under coverage in the AoI, even in case of mobility. Each GV generates video frames of size $n_{\text{UL}} \in [1, 3]$ Mb from its camera sensor at an average rate $r = 10$ fps. Object detection on these frames requires a constant computational load of $C = 60$ GFLOP per frame, which is computed as the average between the computational performance of two popular object detectors, namely Gaussian YOLO and SqueezeDet+ [18]. If frames are offloaded to the HAP (with probability η^*), eventually the processed output (i.e., the bounding boxes of the detected objects) is returned to the GVs in a packet of a much smaller size than the original frame, i.e., $n_{\text{DL}} = 100$ kb, which implies that $t_{\text{DL}} \ll t_{\text{UL}}$.

In the simulations, we compare a fully-local scheme in which all data frames are processed onboard (i.e., $\eta = 0$) and the optimal offloading policy (where $\eta = \eta^*$) based on the solution of Eq. (14) with $t_{\max} = 1/r$. GVs have a computational capacity of $C_{\text{GV}} \in \{200, 600, 800, 1000\}$ GFLOPS, whereas the HAP operates via $c = 15$ parallel servers, each of which offers a computational capacity of $C_{\text{HAP}} \in \{3000, 4000, 5000\}$ GFLOPS, so as to simulate different computing conditions. These values are consistent with the capacity of off-the-shelf computing units: for example, the Nvidia GeForce RTX 3080 Mobile CPU (GeForce GTX 1080 GPU), which is compatible with the space/power constraints onboard GVs (HAP), offers a capacity of 300 (9000) GFLOPS.

All devices operate at a carrier frequency of $f_c = 38$ GHz (mmWave) and with a bandwidth $B = 400$ MHz, while for a complete description of the channel parameters in Sec. II-C we refer the interested reader to [11, Table 1].

B. Numerical Results

a) *Number of users (GVs):* In Fig. 1a, we evaluate the impact of the number of GVs in terms of VEC performance. We can see that, on average, fully local processing onboard the GVs, with $C_{\text{GV}} = 800$ GFLOPS, requires about 200 ms for each video frame, regardless of the value of n , which is not compatible with real-time services. As expected, the average latency for processing data via the HAP grows with n , due to the more frequent offloading requests and the resulting populated queues as the number of GVs increases. In any case, HAP-assisted VEC with powerful processing can reduce the average latency compared to the fully local scenario (up to almost 5 times when $n = 50$ GVs), despite the communication

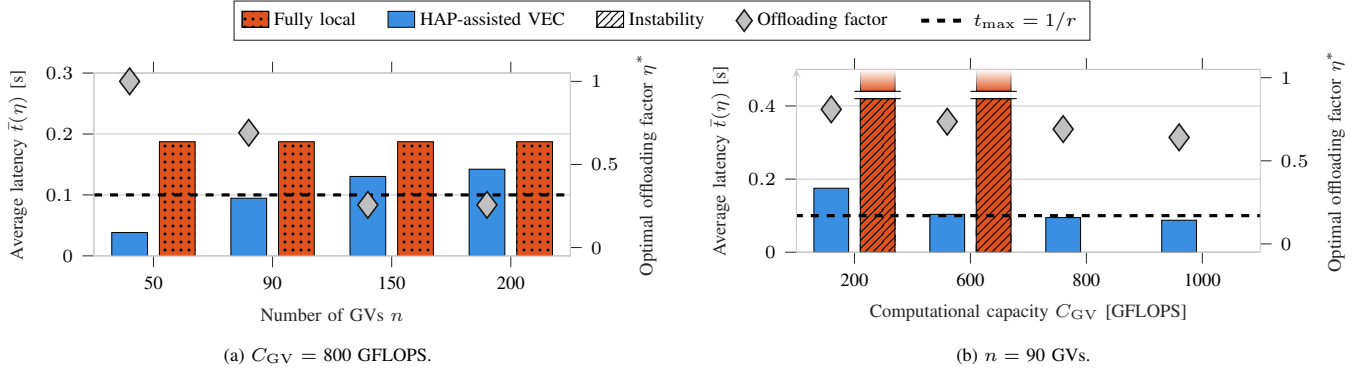


Fig. 1: Optimal offloading factor η^* (right axis) and average latency $\bar{t}(\eta)$ (left axis) vs. n and C_{GV} , for $r = 10$ fps and $n_{UL} = 3$ Mb. Striped bars are interrupted to represent the case of unstable queues where the latency increases indefinitely in the long term.

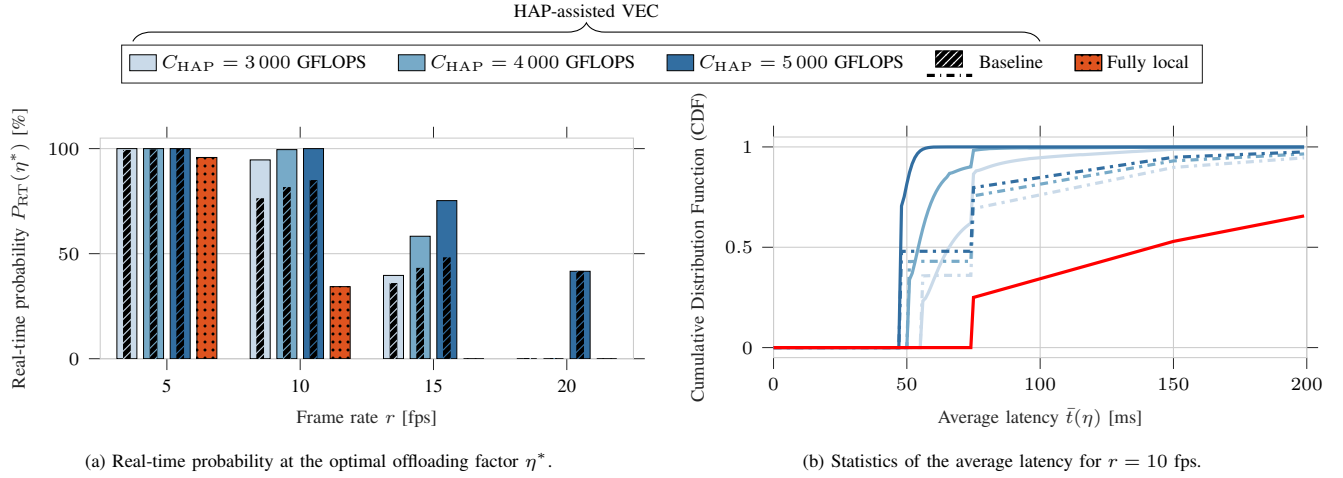


Fig. 2: Real-time probability (left) and average latency (right) vs. C_{HAP} and r , when $C_{GV} = 800$ GFLOPS, $n = 100$ GVs and $n_{UL} = 1$ Mb. Dash-dot lines and black striped bars refer to a baseline offloading scheme which balances the load of the HAP and of the GVs equally.

delay for uploading data frames to the HAP servers and for delivering the processed output to the end nodes.

Interestingly, the optimal offloading probability η^* decreases with the number of GVs. When $n < 150$, the best choice is to offload data with probability $\eta^* > 0.7$. This approach allows to achieve, on average, real-time data processing at the frame rate of the sensors, i.e., $\bar{t}(\eta) < t_{\max}$. On the other hand, when $n \geq 150$, the average latency $\bar{t}(\eta)$ to/from the HAP alone exceeds the latency constraint, i.e., $P(t_{HAP} \leq t_{\max}) = 0$. This is due to the fact that the more populated system may overload the available channel bandwidth, thus resulting in longer transmission delays, which makes fully local processing an increasingly desirable option: in these conditions, the optimal offloading factor is as low as $\eta^* = 0.25$.

b) GV's computational capacity: We let C_{GV} vary from 200 to 1000 GFLOPS, as considered in [11], in a scenario with $n = 90$ GVs, and plot the optimal offloading factor η^* in Fig. 1b. We see that η^* decreases as C_{GV} increases, specifically from 0.81 (200 GFLOPS) to 0.64 (1000 GFLOPS). Fully local processing becomes increasingly more attractive as vehicles incorporate more powerful processing hardware, to avoid additional delays for data offloading. Still, for the values of capacity in Fig. 1b, HAP-assisted VEC remains a more convenient choice to achieve real-time performance, even though the gap in terms of average latency compared to the fully local scenario is only 17% when $C_{GV} = 1000$ GFLOPS. Notably, Fig. 1b shows that fully local processing

leads to queue instability if $C_{GV} \leq 600$ GFLOPS, which further motivates the need for offloading to HAP.

c) Frame rate: In Fig. 2a, we study the real-time probability at the optimal offloading factor η^* , as a function of the frame rate r and the computational capacity of the HAP. In general, increasing r allows the sensor to capture data at better resolution, which leads to more accurate autonomous driving perception, at the cost of a higher data rate and complexity. As expected, $P_{RT}(\eta^*)$ is a decreasing function of r : for $C_{HAP} = 5000$ GFLOPS, P_{RT} ranges from 1 at $r = 10$ fps to 0.41 at $r = 20$ fps. For $C_{HAP} < 5000$ GFLOPS, the system is constrained by the capacity of the queues at the HAP servers; in particular, the system becomes unstable as $r \geq 20$ fps. Moreover, for the fully local configurations, $P_{RT} > 0.95$ only with $r = 5$ fps and $C_{GV} = 800$ GFLOPS, an indication that the limited computational capacity at the GVs is not enough to support automotive tasks at the frame rate of the sensors.

For comparison, we considered a simple baseline offloading scheme that balances the load of the HAP and the GVs using an offloading factor η_{bl} given by

$$\eta_{bl} = \left(\frac{n \cdot C_{GV}}{c \cdot C_{HAP}} + 1 \right)^{-1}, \quad (16)$$

which guarantees $G_{HAP}(\eta_{bl})/c = G_{GV}(\eta_{bl})$. Results in Fig. 2a show that this baseline can support real-time processing with a probability higher than 0.85 as long as the frame rate is less than 10 fps, otherwise the system becomes

unstable. On the other hand, our proposed framework is able to regulate the offloading factor based on the actual capacity of the network, and can offer real-time performance even in more congested scenarios.

In Fig. 2b, we plot the Cumulative Distribution Function (CDF) of the average latency for $r = 10$ fps. HAP-assisted VEC, even with $C_{\text{HAP}} = 3000$ GFLOPs, achieves real-time processing (i.e., $\bar{t}(\eta) \leq 1/r = 100$ ms) with a probability higher than 0.94, vs. 0.34 for fully local processing. The CDF of the baseline shows a trend similar to the fully local scheme, but with a constant gain from 0.44 to 0.54 which depends on C_{HAP} . From Fig. 2b, it is clear that the CDF shows two steps: the first step (at $\bar{t}(\eta) \simeq 50$ ms) corresponds to the minimum processing time for HAP-assisted VEC, i.e., the time for downlink and uplink transmissions, and the processing time at the HAP; the second step (at $\bar{t}(\eta) \simeq 75$ ms) corresponds to the minimum processing time for fully-local processing, i.e., $C/C_{\text{GV}} \simeq 75$ ms.

V. CONCLUSIONS AND FUTURE WORKS

In this paper, we analyzed a rural scenario where autonomous GVs capture sensory data to perform real-time object detection, and addressed the following research question: “Is it feasible to offload and process self-driving-related tasks to HAPs?” For this purpose, we described both the GVs and the HAP as queues that receive sensory data, and determined a closed-form expression for the average waiting time of data in those queues. We then formalized a VEC optimization problem and compared the case in which data processing is performed onboard the GVs (fully local scenario) vs. the case in which a fraction of the processing load is offloaded to HAP servers. Simulation results showed that the limited computational capacity of GVs is generally not compatible with real-time operations, while an optimal offloading factor exists to minimize the processing time. In particular, real-time performance requires a computational capacity at the HAP higher than 3000 GFLOPs for a frame rate $r \leq 10$ fps.

In future work, we will extend our offloading optimization problem by incorporating the impact of power consumption.

APPENDIX: PROOF OF LEMMA 1

Let us consider an M/D/c queue with arrival rate λ and service rate μ . The state probability distribution fulfills Eq. (4). The average length of the queue $E[L_q]$ can be written as

$$E[L_q] = \sum_{k=c}^{M-1} p_k(k-c) + \sum_{k=M}^{+\infty} p_k(k-c). \quad (17)$$

Using the approximation in Eq. (5), the second summation in Eq. (17) can be rewritten as

$$\sum_{k=M}^{+\infty} p_k(k-c) = \sum_{k=M}^{+\infty} p_M \tau^{-(k-M)}(k-c) \quad (18)$$

$$= p_M \left[\sum_{j=0}^{+\infty} \tau^{-j} j + \sum_{j=0}^{+\infty} \tau^{-j} (M-c) \right] \quad (19)$$

$$= p_M \left[\sum_{j=0}^{+\infty} \tau^{-j} j + \frac{M-c}{1-\frac{1}{\tau}} \right]. \quad (20)$$

The summation in Eq. (20) is equal to

$$\sum_{j=0}^{+\infty} \tau^{-j} j = \sum_{k=1}^{+\infty} \sum_{j=k}^{+\infty} \tau^{-j} = \sum_{k=1}^{+\infty} \tau^{-k} \sum_{j=0}^{+\infty} \tau^{-j} \quad (21)$$

$$= \sum_{k=1}^{+\infty} \left[\sum_{j=0}^{+\infty} \tau^{-j} - \sum_{j=0}^{k-1} \tau^{-j} \right] \quad (22)$$

$$= \sum_{k=1}^{+\infty} \left[\left(\frac{1}{\tau} \right)^k \right] = \frac{\frac{1}{1-\frac{1}{\tau}} - 1}{1 - \frac{1}{\tau}}. \quad (23)$$

Finally, $E[L_q]$ can be written from (17), (20), and (23) as

$$E[L_q] = \sum_{k=c}^{M-1} p_k(k-c) + p_M \left[\frac{M-c + \frac{1}{1-\frac{1}{\tau}} - 1}{1 - \frac{1}{\tau}} \right]. \quad (24)$$

where p_0, p_1, \dots, p_M are computed by a linear system, as explained in Section II-D. The average waiting time in the queue, $E[W_q]$, is derived from Little's Law, i.e., $E[W_q] = E[L_q]/\lambda$, so that we obtain the expression in Eq. (6).

REFERENCES

- [1] M. Giordani *et al.*, “Toward 6G networks: Use cases and technologies,” *IEEE Commun. Mag.*, vol. 58, no. 3, pp. 55–61, March 2020.
- [2] A. Chaoub *et al.*, “6G for Bridging the Digital Divide: Wireless Connectivity to Remote Areas,” *IEEE Wireless Commun.*, pp. 160–168, July 2021.
- [3] M. Giordani and M. Zorzi, “Non-Terrestrial Networks in the 6G Era: Challenges and Opportunities,” *IEEE Network*, vol. 35, no. 2, pp. 244–251, Dec. 2021.
- [4] D. C. Nguyen *et al.*, “6G Internet of Things: A Comprehensive Survey,” *IEEE Internet of Things Journal*, vol. 9, no. 1, pp. 359–383, Jan. 2022.
- [5] L. Liu *et al.*, “Vehicular edge computing and networking: A survey,” *Mob. Netw. Appl.*, vol. 26, no. 3, pp. 1145–1168, Jun 2021.
- [6] A. Belogae *et al.*, “Cost-Effective V2X Task Offloading in MEC-Assisted Intelligent Transportation Systems,” *IEEE Access*, vol. 8, pp. 169 010–169 023, Sept. 2020.
- [7] G. Velez *et al.*, “5G beyond 3GPP Release 15 for connected automated mobility in cross-border contexts,” *Sensors*, vol. 20, no. 22, p. 6622, Nov. 2020.
- [8] M. Ke *et al.*, “An Edge Computing Paradigm for Massive IoT Connectivity over High-Altitude Platform Networks,” *IEEE Wireless Commun.*, vol. 28, no. 5, pp. 102–109, Jun. 2021.
- [9] D. S. Lakew *et al.*, “Intelligent Offloading and Resource Allocation in HAP-Assisted MEC Networks,” in *ICTC*, 2021.
- [10] D. Wang *et al.*, “The Potential of Multi-Layered Hierarchical Non-Terrestrial Networks for 6G: A Comparative Analysis Among Networking Architectures,” *IEEE Veh. Technol. Mag.*, vol. 16, no. 3, pp. 99–107, Sep. 2021.
- [11] A. Traspadini *et al.*, “UAV/HAP-Assisted Vehicular Edge Computing in 6G: Where and What to Offload?” *EuCNC-6G Summit*, 2022.
- [12] M. Lundgren *et al.*, “Vehicle self-localization using off-the-shelf sensors and a detailed map,” in *IEEE Intelligent Vehicles Symposium*, 2014.
- [13] V. Rossi *et al.*, “On the Role of Sensor Fusion for Object Detection in Future Vehicular Networks,” in *EuCNC-6G Summit*, 2021.
- [14] 3GPP, “Solutions for NR to support Non-Terrestrial Networks (NTN),” TR 38.821 (Release 16), 2020.
- [15] M. Giordani and M. Zorzi, “Satellite Communication at Millimeter Waves: a Key Enabler of the 6G Era,” *IEEE ICNC*, Feb 2020.
- [16] H. Tijms, *A First Course in Stochastic Models*. Wiley Ed., 2003.
- [17] K. R. Gegenfurtner, “PRAXIS: Brent’s algorithm for function minimization,” *Behav. Res. Methods Instrum. Comput.*, vol. 24, no. 4, pp. 560–564, Dec. 1992.
- [18] W. Li and K. Liu, “Confidence-Aware Object Detection Based on MobileNetV2 for Autonomous Driving,” *Sensors*, vol. 21, no. 7, p. 2380, March 2021.