# CONGESTION-AWARE DISTRIBUTED TASK OFFLOADING IN WIRELESS MULTI-HOP NETWORKS USING GRAPH NEURAL NETWORKS

*Zhongyuan Zhao⋆, Jake Perazzone‡, Gunjan Verma‡, and Santiago Segarra⋆*

⋆Rice University, USA ‡US Army's DEVCOM Army Research Laboratory, USA

## ABSTRACT

Computational offloading has become an enabling component for edge intelligence in mobile and smart devices. Existing offloading schemes mainly focus on mobile devices and servers, while ignoring the potential network congestion caused by tasks from multiple mobile devices, especially in wireless multi-hop networks. To fill this gap, we propose a low-overhead, congestion-aware distributed task offloading scheme by augmenting a distributed greedy framework with graph-based machine learning. In simulated wireless multi-hop networks with 20-110 nodes and a resource allocation scheme based on shortest path routing and contention-based link scheduling, our approach is demonstrated to be effective in reducing congestion or unstable queues under the context-agnostic baseline, while improving the execution latency over local computing.
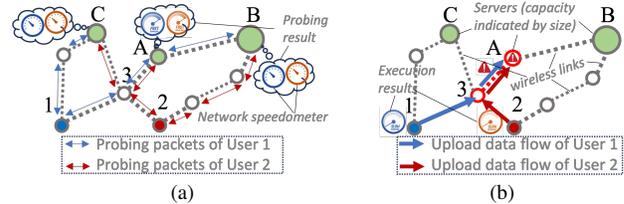
*Index Terms*— Computational offloading, queueing networks, wireless multi-hop networks, graph neural networks, shortest path.

## 1. INTRODUCTION

The proliferation of mobile and smart devices enables the collection of rich sensory data from both physical and cyber spaces, leading to many exciting applications, such as connected vehicles, drone/robot swarms, software-defined networks (SDN), and Internet-of-Things (IoT) [1–4]. To support these applications, wireless multi-hop networks, which have been traditionally used for military communications, disaster relief, and sensor networks, are now envisioned in the next-generation wireless networks, including xG (device-to-device, wireless backhaul, and non-terrestrial coverage), vehicle-to-everything (V2X), and machine-to-machine (M2M) communications [5–7]. This can be partially attributed to the self-organizing capability of wireless multi-hop networks, enabled by distributed resource allocation without relying on infrastructure [8–14]. Additionally, computational offloading [15–23] promises to improve the performance and energy efficiency of resource-limited mobile devices by moving their resource-intensive computational tasks to external servers, including remote computing clusters and edge servers located near the mobile devices [16]. In particular, computational offloading has become an enabling technology for edge intelligence, as it is often impractical to equip mobile devices with hardware accelerators due to economic or energy constraints.

**Fig. 1**: Challenges in distributed multi-hop offloading: (a) probing: nodes 1 and 2 query the communication and computing bandwidth of three servers. (b) offloading: nodes 1 and 2 both select server A based on collected information, however, such decisions lead to congestion at both the server A and link (3,A) in execution.

Current studies in computational offloading (also referred as mobile edge computing, fog computing, cloudlets, etc.) mostly focus on the performance and energy consumption of individual devices [15, 20, 23] under simplified networking assumptions, e.g., servers are within a single-hop [15, 23]. For offloading in wireless multi-hop networks [17–22], a centralized scheduler with global knowledge of the network is often assumed [17, 20]. However, centralized multi-hop offloading has the drawbacks of single-point-of-failure and poor scalability, due to the high communication overhead of collecting the full network state to a dedicated scheduler. Distributed multi-hop offloading based on pricing [18,21] and learning [22] only focus on the capacity of servers, while ignoring the potential network congestion caused by offloading [19], as illustrated by the motivating example in Fig. 1. In phase 1 (Fig. 1(a)), nodes 1 and 2 query the capacities of three servers and the corresponding links with probing packets, leading them to both conclude that server A offers the fastest execution. In phase 2 (Fig. 1(b)), nodes 1 and 2 start offloading their tasks to server A via the blue and red routes decided in phase 1, causing congestion at link $(3, A)$ and server $A$ since their capacities cannot simultaneously support the traffic of the two tasks. This problem is more pronounced for recurrent computational tasks, such as video surveillance and network traffic analysis. Moreover, the complexity of managing this issue by negotiation between devices or trial-and-error can increase exponentially with the number of mobile nodes.

In this work, we develop an intelligent distributed task offloading scheme that can exploit the network context via graph-based machine learning. Specifically, we build a distributable graph neural network (GNN) that can encode the network topology and information from all links, servers, and tasks in the network into congestion-aware link weights. Such link weights can mitigate network congestion (unstable queues) by enabling mobile devices to better estimate the costs of their offloading options in the presence of tasks on other mobile devices, leading to improved task execution latency.

The contributions of this paper are twofold:
1) To our best knowledge, we present the first low-complexity approach to integrate network context into fully distributed and near-simultaneous offloading decisions in wireless multi-hop networks.

2) Our approach is proved to be able to mitigate congestion while offloading tasks in simulated wireless multi-hop networks.

## 2. SYSTEM MODEL AND PROBLEM FORMULATION

We model a wireless multi-hop network as a *connectivity graph* $\mathcal{G}^n$ and a *conflict graph* $\mathcal{G}^c$. The connectivity graph is an undirected graph $\mathcal{G}^n = (\mathcal{V}, \mathcal{E})$, where $\mathcal{V}$ is a set of nodes representing wireless devices in the network and $\mathcal{E}$ is a set of links in which $e = (v_1, v_2) \in \mathcal{E}$ for $v_1, v_2 \in \mathcal{V}$ indicates that nodes $v_1$ and $v_2$ can communicate directly. $\mathcal{G}^n$ is assumed to be a connected graph, i.e., two arbitrary nodes in the network can always reach each other. The conflict graph, $\mathcal{G}^c = (\mathcal{E}, \mathcal{C})$, describes the conflict relationship between links and is defined as follows: each vertex $e \in \mathcal{E}$ corresponds to a link in $\mathcal{G}^n$ and each undirected edge $(e_1, e_2) \in \mathcal{C}$ indicates a conflict between links $e_1, e_2 \in \mathcal{E}$ in $\mathcal{G}^n$. Two links could be in conflict due to either 1) *interface conflict*, i.e., two links share a wireless device with only one wireless transceiver; or 2) *wireless interference*, i.e., their incident devices are within a certain distance such that their simultaneous transmission will cause the outage probability to exceed a prescribed level [24]. In this paper, we assume the conflict graph $\mathcal{G}^c$ to be known, possibly by each link monitoring the wireless channel [12], or through more sophisticated estimation as in [25].

Based on the role of each wireless device, $\mathcal{V}$ can be partitioned into three subsets: $\mathcal{M}$ for edge nodes, $\mathcal{R}$ for relay nodes, and $\mathcal{S}$ for server nodes. An edge node $v \in \mathcal{M}$ is a sensor with limited computing capability and power supply, such as IoT sensors, SDN routers, phones, wearable devices, drones, or robots. A relay node $v \in \mathcal{R}$ is dedicated to communications, such as satellite, fixed, or mobile relay stations. A server node $v \in \mathcal{S}$ has richer computing resources and power supply than edge nodes, but may be located multiple hops away from the requesting edge nodes.

In this study, we consider a simplified task offloading scenario, in which each edge node may generate a non-trivial computational task for processing sensory data. A task $j_m$ is a series of similar jobs generated by an edge node $m \in \mathcal{M}$ at certain rate, and each job must be individually processed. In particular, task $j_m$ encompasses the information of the source $m$, the set of available external servers $\mathcal{S}_m \subseteq \mathcal{S}$, the job arrival rate $\lambda^j(m)$, and the numbers of upload and download data packets per job, respectively denoted as $\eta^u(j_m)$ and $\eta^d(j_m)$, where $\eta(j_m) = \eta^u(j_m) + \eta^d(j_m)$. A task $j_m$ can be executed locally at its source $m$ or remotely on an external server $v \in \mathcal{S}_m$ by uploading the data to the server $v$, and if required, sending back the results. We define set $\mathcal{S}_m^+ = \{m\} \cup \mathcal{S}_m$ as the action space of offloading task $j_m$. We assume that all jobs are atomic, $\eta^u(j_m) \gg \eta^d(j_m)$, jobs of different tasks arrive independently, and the execution time of a job grows by its description length. Finally, we denote $\mathcal{J}$ as the set of all tasks in the network.

We consider a time-slotted medium access control in the network. We model each wireless link $e \in \mathcal{E}$ as a $G/G/1$ queueing system with a first-in-first-out (FIFO) queue, a single server, and general packet arrival and service processes, where the arrival and service rates are $\lambda(e) \geq 0$ and $\mu(e) \geq 0$, respectively. Under this model, a packet leaves the queueing system when it reaches the other end of the link. Based on Little's law [26], if $\mu(e) > \lambda(e)$, the probability of a link having a non-empty queue is $\lambda(e)/\mu(e)$ and the expected time of a packet passing through a link is the response time of the queueing system, $1/(\mu(e) - \lambda(e))$. If $\mu(e) \leq \lambda(e)$, the queue becomes unstable and will grow infinitely, which is referred as congestion. However, to quantify the level of congestion, we use the expected time to deplete the queue of a link, $T\lambda(e)/\mu(e)$, assuming that new jobs only arrive in the first $T$ time slots. Notice that $\mu(e)$ of a link is generally unknown as it depends on both the average link

rate $r(e)$ and the link scheduling policy. A similar queueing model also applies to a computing node $v \in \mathcal{M} \cup \mathcal{S}$, i.e., an edge or server node, except that the service rate $\mu(v)$ is known in advance and not affected by link scheduling. Here, the arrival rate, link rate, and service rate are all defined in terms of number of packets per time slot.

We denote the features of links, nodes, and tasks under the following convention unless otherwise specified. For example, vector $\mathbf{r}^c = [r(e)|e \in \mathcal{E}]$ collects the link rates of all the wireless links, where superscript $c$ indicates that the dimension of $\mathbf{r}^c$ equals the number of nodes in graph $\mathcal{G}^c$, i.e., $\mathbf{r}^c \in \mathbb{R}^{|\mathcal{E}|}$. Similarly, vector $\boldsymbol{\mu}^n \in \mathbb{R}^{|\mathcal{V}|}$ describes the service rates of the node set $\mathcal{V}$ in graph $\mathcal{G}^n$ and vector $\mathbf{u}^j \in \mathbb{R}^{|\mathcal{J}|}$ represents the execution latency (in time slots) of all tasks. Matrices are denoted by upright bold upper-case symbols, where $\mathbf{X}_{ab}$ is the element at row $a$ and column $b$ of matrix $\mathbf{X}$, and $\mathbf{X}_{a*}$ is the $a$th row vector and $\mathbf{X}_{*b}$ is the $b$th column vector.

Formally, the latency-optimal multi-hop offloading problem is formulated as finding the optimal offload locations $\mathbf{s}^{j*}$ to minimize the total expected execution latency across all tasks, $\mathbf{1}^\top \mathbf{u}^j$,

$$\mathbf{s}^{j*} = \underset{\mathbf{s}^j}{\operatorname{argmin}} \sum_{j_m \in \mathcal{J}} u(j_m, m, s_m) \tag{1a}$$

$$\text{s.t. } \mathbf{s}^j \coloneqq [s_m | j_m \in \mathcal{J}], \ \mathbf{u}^j \coloneqq [u(j_m, m, s_m) | j_m \in \mathcal{J}], \tag{1b}$$

$$s_m \in \mathcal{S}_m^+, \text{for all } m \in \mathcal{M}, \tag{1c}$$

$$\mathbf{u}^j = f_r\left(\mathcal{G}^n, \mathcal{G}^c, \boldsymbol{\mu}^n, \mathbf{r}^c, \mathcal{J}, \mathbf{s}^j\right), \tag{1d}$$

where $u(j_m, m, s_m) \geq 0$ is the expected latency for task $j_m$ being executed on node $s_m$, and $f_r(\cdot)$ is a deterministic function that maps the network state $(\mathcal{G}^n, \mathcal{G}^c, \boldsymbol{\mu}^n, \mathbf{r}^c, \mathcal{J})$ and decision variables $\mathbf{s}^j$ to the expected execution latency of all tasks $\mathbf{u}^j$, defined in (1b). $f_r(\cdot)$ captures the effect of the resource allocation policy of the wireless network, such as the routing and scheduling protocols, on the expected execution latency of all the tasks under given offloading decisions $\mathbf{s}^j$. Problem (1) belongs to the class of generalized assignment problems (GAPs) [27], which is known to be NP-hard.

We specify the constraint in (1d) with the following resource allocation schemes commonly found in practice: 1) a contention-based scheduling policy that offers conflicting links (neighboring nodes on the conflict graph) with non-empty queues an equal chance of transmission, e.g., CSMA [11]; and 2) a shortest path routing scheme that determines the route between source $m$ and an external server $s \in \mathcal{S}_m$ for task $j_m$, based on the link weights; an example of such weights could be the expected time a data packet takes to pass through a link under the current traffic condition. Note that our approach can be easily adapted to other resource allocation schemes.

## 3. DISTRIBUTED TASK OFFLOADING WITH GNNS

We propose to approximately solve Problem (1) by augmenting distributed greedy decision-making with a trainable GNN. Our distributed decision framework is based on an extended connectivity graph, $\mathcal{G}^e = (\mathcal{V}^e, \mathcal{E}^e)$, built by adding virtual nodes and links to the original connectivity graph $\mathcal{G}^n$ as $\mathcal{V}^e = \mathcal{V} \cup \tilde{\mathcal{V}}$ and $\mathcal{E}^e = \mathcal{E} \cup \tilde{\mathcal{E}}$. For each edge or server node $v \in \mathcal{M} \cup \mathcal{S}$, there is a corresponding virtual node $\tilde{v} \in \tilde{\mathcal{V}}$ connected to $v$ by a virtual link $(v, \tilde{v}) \in \tilde{\mathcal{E}}$. The link rate of a virtual link equals to the service rate of node $v$, i.e., $r((v, \tilde{v})) = \mu(v)$, and the link rate of a physical link remains the same. We further introduce the line graph of $\mathcal{G}^e$ as $\mathcal{G}^\ell$: each vertex in $\mathcal{G}^\ell$ corresponds to an edge in $\mathcal{G}^e$, and an edge in $\mathcal{G}^\ell$ indicates that the two corresponding edges in $\mathcal{G}^e$ share a common vertex [28]. The vector of the extended link rates $\mathbf{r}^\ell \in \mathbb{R}^{|\mathcal{E}^e|}$ thus captures both the original link rates $\mathbf{r}^c$ and original service rates $\boldsymbol{\mu}^n$.

The distributed task offloading decision, denoted as $\mathbf{s}^j = h(\mathcal{G}^\ell, \boldsymbol{\delta}^\ell, \mathcal{J})$, lets each edge node select the offloading location of minimal cost based on given weights of the extended links $\boldsymbol{\delta}^\ell$,

$$\mathbf{s}^j = [s_m | j_m \in \mathcal{J}], \quad \text{where } s_m = \underset{v \in \mathcal{S}_m^+}{\arg\min} \, c(m, v), \qquad (2a)$$

$$c(m,v) = \max\left[\eta^u(j_m)\beta(m, \tilde{v}) + \eta^d(j_m)\beta(v, m), 2\zeta(v, m)\right]. \quad (2b)$$

In (2), the cost of offloading action $v \in \mathcal{S}_m^+$, denoted as $c(m, v)$, is the expected round-trip delay of a job on graph $\mathcal{G}^e$, while $\beta(v_1, v_2)$ and $\zeta(v_1, v_2)$ are the shortest path distance and hop distance between nodes $v_1, v_2$ on the edge-weighted graph $(\mathcal{V}^e, \mathcal{E}^e, \boldsymbol{\delta}^\ell)$, respectively. Eq. (2b) says that a job will take at least one time slot to travel through a physical link, e.g., even if all the data packets of a job can go across a link within one time slot, they can only be sent over the next link until the next time slot.

Next, we need to find a link weight vector, $\boldsymbol{\delta}^\ell = [\delta(e)|e \in \mathcal{E}^e]$, that can be translated to good offloading decisions by the distributed decision framework $h(\cdot)$ in (2). The baseline approach is to use the per-packet delay under a contention-free assumption, e.g., $\bar{\boldsymbol{\delta}}^\ell = 1/\mathbf{r}^\ell$, which, however, only holds for scenarios with only a few tasks and lightly-loaded task traffic.

We propose to use a distributable GNN to predict a congestion-aware edge weight vector as $\hat{\boldsymbol{\delta}}^\ell = f(\mathcal{G}^\ell, \boldsymbol{\lambda}^\ell, \mathbf{r}^\ell, \mathcal{G}^c; \boldsymbol{\omega})$, where $\boldsymbol{\omega}$ is the collection of trainable parameters of the GNN, and $\boldsymbol{\lambda}^\ell = [\lambda(e)|e \in \mathcal{E}^e]$ assigns the packet arrival rates of tasks to corresponding virtual links, e.g., $\lambda(e) = \lambda(v) = \lambda^j(v)\eta(j_v)$ if $e = (v, \tilde{v}) \in \tilde{\mathcal{E}}$ and $j_v \in \mathcal{J}$, otherwise, $\lambda(v) = 0$. In step 1, the GNN predicts the packet arrival rates on the extended links $\mathbf{x}^\ell \in \mathbb{R}^{|\mathcal{E}^e|}$, as $\mathbf{x}^\ell = \Psi_{\mathcal{G}^\ell}(\mathbf{X}^0; \boldsymbol{\omega})$, where $\Psi_{\mathcal{G}^\ell}$ is an $L$-layered graph convolutional neural network (GCNN) defined on graph $\mathcal{G}^\ell$, and $\boldsymbol{\omega}$ is the collection of trainable parameters. We define the output of an intermediate $l$-th layer of the GCNN as $\mathbf{X}^l \in \mathbb{R}^{|\mathcal{E}^e| \times g_l}$, and the input and output dimensions of the GCNN are set as $g_0 = 4$ and $g_L = 1$. The input node features are defined as $\mathbf{X}^0 = [\mathbf{q}^\ell, \mathbf{w}^\ell, \boldsymbol{\lambda}^\ell, \mathbf{r}^\ell]$, where $\mathbf{q}^\ell$ is an indicator vector of virtual links, i.e., $\mathbf{q}^\ell = [\mathbb{1}_{\tilde{\mathcal{E}}}(i)|i \in \mathcal{E}^e]$, and $\mathbf{w}^\ell$ is an indicator vector of virtual links for server nodes.

The $l$-th layer of the GCNN $\Psi_{\mathcal{G}^\ell}$ can be implemented in a fully distributed manner through neighborhood aggregation as the following local operation on an extended link $i \in \mathcal{E}^e$ (a vertex in $\mathcal{G}^\ell$)

$$\mathbf{X}_{i*}^l = \sigma_l\left(\mathbf{X}_{i*}^{l-1}\boldsymbol{\Theta}_0^l + \left[\mathbf{X}_{i*}^{l-1} - \sum_{e \in \mathcal{N}^\ell(i)} \frac{\mathbf{X}_{i*}^{l-1}}{\sqrt{d^\ell(e)d^\ell(i)}}\right]\boldsymbol{\Theta}_1^l\right), \quad (3)$$

where $\mathbf{X}_{i*}^l \in \mathbb{R}^{1 \times g_l}$ captures the $l$th-layer features on link $i$, $\mathcal{N}^\ell(i)$ denotes the set of neighbors of $i$ in $\mathcal{G}^\ell$, $d^\ell(\cdot)$ is the degree of a vertex in $\mathcal{G}^\ell$, $\boldsymbol{\Theta}_0^l, \boldsymbol{\Theta}_1^l \in \mathbb{R}^{g_{l-1} \times g_l}$ are trainable parameters (collected in $\boldsymbol{\omega}$), and $\sigma_l(\cdot)$ is an element-wise activation function of the $l$-th layer. Based on (3), the link arrival rate vector $\mathbf{x}^\ell$ can be computed in a fully distributed manner through $L$ rounds of local message exchanges between $i \in \mathcal{E}^e$ and its neighbors on $\mathcal{G}^\ell$.

In step 2, the congestion-aware weights for the extended links are estimated as $\hat{\boldsymbol{\delta}}^\ell = \varphi(\mathcal{G}^c, \mathcal{G}^\ell, \mathbf{r}^\ell, \mathbf{x}^\ell, T, K)$, where function $\varphi(\cdot)$ is described by Algorithm 1 and can be implemented in a fully distributed manner. In Algorithm 1, the service rate of each physical link, $\boldsymbol{\mu}_e^c$, is initialized to the worst-case scenario that every neighboring links always contend for transmission. Then, for each link $e \in \mathcal{E}$, the algorithm iteratively updates: $\mathbf{b}_e^c$, the probability of link $e$ contending for transmission due to non-empty queues, based on its input packet arrival rate $\mathbf{x}_e^c$ and service rate; $\mathbf{p}_e^c$, the expected number of contending neighbors of link $e$; and $\boldsymbol{\mu}_e^c$ the service rate of link $e$ based on the scheduling policy that contending links have the

---

**Algorithm 1** Iterative execution latency estimation algorithm $\varphi(\cdot)$

**Input**: $\mathcal{G}^c, \mathcal{G}^\ell, \mathbf{r}^\ell, \mathbf{x}^\ell, T, K$
**Output**: $\boldsymbol{\delta}^\ell$

1: Get $\mathbf{A}^c$ as the adjacency matrix of $\mathcal{G}^c$
2: Get conflict degree vector $\mathbf{d}^c = \mathbf{A}^c\mathbf{1}^c$
3: $\mathbf{x}^c = [\mathbf{x}_e^\ell|e \in \mathcal{E}], \boldsymbol{\mu}^c(0) = \mathbf{r}^c/(\mathbf{1}^c + \mathbf{d}^c)$
4: **for** $k \in \{1, \dots, K\}$ **do**
5: $\quad \mathbf{b}^c(k) = \min[\mathbf{x}^c/\boldsymbol{\mu}^c(k-1), \mathbf{1}]$
6: $\quad \mathbf{p}^c(k) = \mathbf{b}^c(k)\mathbf{A}^c$
7: $\quad \boldsymbol{\mu}^c(k) = \mathbf{r}^c/[\mathbf{1}^c + \mathbf{p}^c(k)]$
8: **end for**
9: $\hat{\boldsymbol{\mu}}^\ell = [\hat{\mu}(e)|e \in \mathcal{E}^e]$, where $\hat{\mu}(e) = \boldsymbol{\mu}_e^c(K)$ for a physical link $e \in \mathcal{E}$ and $\hat{\mu}(e) = \mathbf{r}_e^\ell$ for a virtual link $e \in \tilde{\mathcal{E}}$.
10: $\boldsymbol{\delta}^\ell = [\delta(e)|e \in \mathcal{E}^e]$, where $\delta(e) = (\hat{\boldsymbol{\mu}}_e^\ell - \mathbf{x}_e^\ell)^{-1}$ if $\hat{\boldsymbol{\mu}}_e^\ell > \mathbf{x}_e^\ell$, otherwise $\delta(e) = T\mathbf{x}_e^\ell/\hat{\boldsymbol{\mu}}_e^\ell$

---

same chance of transmission. Lastly, the per-packet latency of the extended links are updated based on their arrival and service rates.

**Empirical task execution latency.** Based on Algorithm 1, the empirical per-packet latency of the extended links can also be estimated as $\boldsymbol{\tau}^\ell = \varphi(\mathcal{G}^c, \mathcal{G}^\ell, \mathbf{r}^\ell, \boldsymbol{\rho}^\ell, T, K)$. Here, vector $\boldsymbol{\rho}^\ell = \boldsymbol{\Gamma}\boldsymbol{\lambda}^j$ captures the traffic intensities on the extended links, which is determined by the packet arrival rates of all tasks $\boldsymbol{\lambda}^j = [\lambda(v)|(v, \tilde{v}) \in \tilde{\mathcal{E}} \ \& \ j_v \in \mathcal{J}]$, and the uploading route indicator matrix $\boldsymbol{\Gamma} \in \{0, 1\}^{|\mathcal{E}^e| \times |\mathcal{J}|}$. $\boldsymbol{\Gamma}$ is defined as: $\boldsymbol{\Gamma}_{e,j_m} = 1$ if link $e \in \mathcal{E}^e$ is on the route from $m$ to the virtual node $\tilde{s}_m$ of the corresponding offloading action $s_m = \mathbf{s}_m^j$ of task $j_m$, otherwise $\boldsymbol{\Gamma}_{e,j_m} = 0$. The downloading route matrix $\boldsymbol{\Gamma}^-$ is define as $\boldsymbol{\Gamma}_{e*}^- = \boldsymbol{\Gamma}_{e*}$ for all $e \in \mathcal{E}$, and $\boldsymbol{\Gamma}_{e*}^- = \mathbf{0}^\top$ for all $e \notin \mathcal{E}$. The empirical execution latency $\mathbf{u}^j$ can be found as

$$\mathbf{u}^j = f_u(\mathbf{s}^j) = \max\left[\boldsymbol{\tau}^{\ell\top}\boldsymbol{\Gamma} \odot \boldsymbol{\eta}^u + \boldsymbol{\tau}^{\ell\top}\boldsymbol{\Gamma}^- \odot \boldsymbol{\eta}^d, \ 2\mathbf{1}^\top\boldsymbol{\Gamma}^-\right]. \quad (4)$$

The empirical task execution latency vectors under the baseline and GNN-based policies are $\bar{\mathbf{u}}^j = f_u(\bar{\mathbf{s}}^\ell)$ and $\hat{\mathbf{u}}^j = f_u(\hat{\mathbf{s}}^\ell)$, respectively, where $\bar{\mathbf{s}}^j = h(\mathcal{G}^\ell, \bar{\boldsymbol{\delta}}^\ell, \mathcal{J})$ and $\hat{\mathbf{s}}^j = h(\mathcal{G}^\ell, \hat{\boldsymbol{\delta}}^\ell, \mathcal{J})$.

**Complexity.** For distributed execution, the local communication complexity (defined as the rounds of local exchanges between a node and its neighbors) of the GNN is $\mathcal{O}(L + K)$. For the distributed decision framework in (2), the distributed weighted single-source-shortest-path (SSSP) with the Bellman-Ford algorithm [29, 30] and all-pairs-shortest-path (APSP) with a very recent algorithm in [31] both take $\mathcal{O}(|\mathcal{V}|)$ rounds of message exchanges.

**Training.** The parameters $\boldsymbol{\omega}$ (collecting $\boldsymbol{\Theta}_0^l$ and $\boldsymbol{\Theta}_1^l$ across all layers $l$) of our GNN are trained on a set of random instances drawn from a target distribution $\Omega$. Based on each instance of the form $\varepsilon = (\mathcal{G}^n, \mathcal{G}^c, \mathcal{M}, \mathcal{R}, \mathcal{S}, \mathcal{J}, \mathbf{r}^c, \boldsymbol{\mu}^n) \sim \Omega$, we create the corresponding extended graph $\mathcal{G}^e$, its line graph $\mathcal{G}^\ell$, $\boldsymbol{\lambda}^\ell$, and $\mathbf{r}^\ell$, and run the full pipeline to get $\hat{\mathbf{u}}^j$. Since $\varphi(\cdot)$ and (4) are differentiable, we can estimate the gradient of the objective $o(\boldsymbol{\omega}) = \mathbf{1}^\top\hat{\mathbf{u}}^j$ in (1a) w.r.t. $\boldsymbol{\Gamma}$,
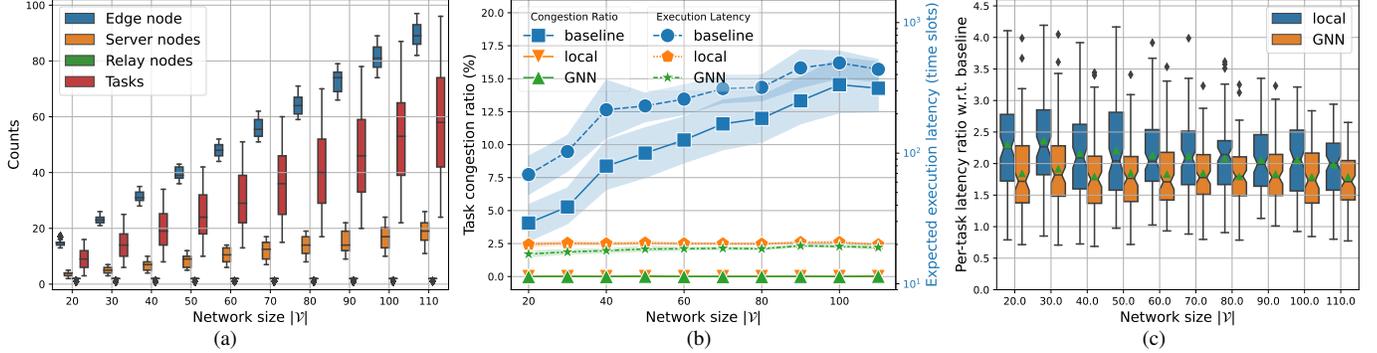
$$\nabla_{\boldsymbol{\Gamma}}o(\boldsymbol{\omega}) = \mathbb{E}_{\varepsilon \sim \Omega}\left[\frac{\partial\mathbf{1}^\top\hat{\mathbf{u}}^j(\varepsilon)}{\partial\boldsymbol{\Gamma}}\right], \ \widehat{\nabla_{\boldsymbol{\Gamma}}o(\boldsymbol{\omega})} = \frac{\partial\mathbf{1}^\top\hat{\mathbf{u}}^j(\varepsilon)}{\partial\boldsymbol{\Gamma}}. \quad (5)$$

We then estimate the gradient of $o(\boldsymbol{\omega})$ w.r.t. $\hat{\boldsymbol{\delta}}^\ell$ and $\boldsymbol{\omega}$ as

$$\widehat{\nabla_{\hat{\boldsymbol{\delta}}^\ell}o(\boldsymbol{\omega})} = -\mathbf{1}^\top\widehat{\nabla_{\boldsymbol{\Gamma}}o(\boldsymbol{\omega})} + \nabla_{\hat{\boldsymbol{\delta}}^\ell}\left(\hat{\boldsymbol{\tau}}^\ell - \hat{\boldsymbol{\delta}}^\ell\right)^2/|\mathcal{V}^e|, \quad (6a)$$

$$\widehat{\nabla_{\boldsymbol{\omega}}o(\boldsymbol{\omega})} = \widehat{\nabla_{\hat{\boldsymbol{\delta}}^\ell}o(\boldsymbol{\omega})}\left[\partial f(\mathcal{G}^\ell, \boldsymbol{\lambda}^\ell, \mathbf{r}^\ell, \mathcal{G}^c; \boldsymbol{\omega})/\partial\boldsymbol{\omega}\right]. \quad (6b)$$

The estimation in (6a) is based on the facts that a modification in $\hat{\boldsymbol{\delta}}^\ell$ would reduce the cost $o(\boldsymbol{\omega})$ if: i) It more faithfully captures the

**Fig. 2**: (a) Scale of instances: the numbers of edge nodes, server nodes, relay nodes, and tasks by network size. (b) Task congestion probability by network size with $95\%$ confidence interval. (c) Average per-task execution latency ratio w.r.t. the baseline policy across network sizes.

empirical per-packet latency $\hat{\boldsymbol{\tau}}^\ell$ (second term), and ii) Aggregated over jobs, the incorporation of the given link in the corresponding routes reduces $o(\boldsymbol{\omega})$ (first term). Furthermore, recalling that $\hat{\boldsymbol{\delta}}^\ell = f(\mathcal{G}^\ell, \boldsymbol{\lambda}^\ell, \mathbf{r}^\ell, \mathcal{G}^c; \boldsymbol{\omega})$, expression (6b) stems from applying the chain rule. For each sampled instance $\varepsilon \sim \Omega$, $\boldsymbol{\omega}$ is updated by stochastic gradient descent (SGD), $\boldsymbol{\omega} = \boldsymbol{\omega} - \alpha \widehat{\nabla_{\boldsymbol{\omega}} o(\boldsymbol{\omega})}$, where $\alpha > 0$ is the learning rate. The training ends based on an early stop mechanism.

## 4. NUMERICAL EXPERIMENTS

The GNN-enhanced distributed offloading is evaluated on simulated wireless ad-hoc networks. Each instance in the training and test sets of the form $\varepsilon = (\mathcal{G}^n, \mathcal{G}^c, \mathcal{M}, \mathcal{R}, \mathcal{S}, \mathcal{J}, \mathbf{r}^c, \boldsymbol{\mu}^n)$ is generated as follows. The connectivity graph $\mathcal{G}^n$ is drawn as a random graph from the Barabási–Albert (BA) [32] model, and the conflict graph $\mathcal{G}^c$ is the line graph of $\mathcal{G}^n$. The BA model has two parameters: the number of vertices $|\mathcal{V}|$ and the number of edges, $\nu$, that each new vertex forms during a preferential attachment process. We set $\nu = 2$, and $|\mathcal{V}| \in \{20, 30, \ldots, 110\}$. The relay node set $\mathcal{R}$ is selected as well-connected nodes by applying minimal node cut on $\mathcal{G}^n$. We then cut $\mathcal{V}$ into a larger partition $\mathcal{V}^b$ and a smaller partition $\mathcal{V}^s$ by solving the minimal cut problem on $\mathcal{G}^n$ with the Stoer-Wagner algorithm [33]. The server set $\mathcal{S}$ contains $10\% \sim 25\%$ of nodes, i.e., $|\mathcal{S}| = \lfloor \mathbb{U}(0.1, 0.25)|\mathcal{V}| \rfloor$, randomly selected from $\mathcal{V}^s - \mathcal{R}$, and if $|\mathcal{V}^s - \mathcal{R}| < |\mathcal{S}|$ then from $\mathcal{V}^b - \mathcal{R}$, where $\mathbb{U}$ stands for uniform distribution. The rest are edge nodes, $\mathcal{M} = \mathcal{V} \setminus (\mathcal{R} \cup \mathcal{S})$. In this way, we make sure that server nodes are likely multiple hops away from edge nodes. The link rate $\mathbf{r}^c_e \sim \mathbb{U}(30, 70)$, for all $e \in \mathcal{E}$, and the service rate $\boldsymbol{\mu}^n_v$ are drawn from Pareto distributions with shape $a = 2$ and mode $q = 100$ for $v \in \mathcal{S}$ and $q = 8$ for $v \in \mathcal{M}$, and $\boldsymbol{\mu}^n_v = 0$, for all $v \in \mathcal{R}$. The training and test sets, respectively, contain 2000 and 1000 network instances $(\mathcal{G}^n, \mathcal{G}^c, \mathcal{M}, \mathcal{R}, \mathcal{S}, \mathbf{r}^c, \boldsymbol{\mu}^n)$ generated under the same configuration but with different sets of pseudo-random seeds. The task set $\mathcal{J}$ is created randomly on-the-fly during training and testing, with job parameters $\eta^u(j_m) = 100, \eta^d(j_m) = 1, \lambda^j(j_m) \sim \mathbb{U}(0.015, 0.075)$, $\mathcal{S}_m = \mathcal{S}$, and $|\mathcal{J}| = \lfloor \mathbb{U}(0.3, 1)|\mathcal{M}| \rfloor$. For each network instance, we draw 10 random instances of task set $\mathcal{J}$. The scale of simulation instances by network size is illustrated in Fig. 2(a).

The hyperparameters of our GNN are $L = 5, K = 10, \alpha = 10^{-6}$, and $g_l = 32$ for $l \in \{1, 2, 3, 4\}$.[1] We limit the arrival of new jobs to the first $T = 1000$ time slots, so that the execution latency of a task is always finite, even if it is congested. A task $j_m \in \mathcal{J}$ is con-

gested if the queues of any link $e \in \mathcal{E}^e$ on its route in graph $\mathcal{G}^\ell$ is unstable, i.e., $\hat{\boldsymbol{\mu}}^\ell_e < \boldsymbol{\rho}^\ell_e$ while evaluating $\boldsymbol{\tau}^\ell = \varphi(\mathcal{G}^c, \mathcal{G}^\ell, \mathbf{r}^\ell, \boldsymbol{\rho}^\ell, T, K)$ (see line 10 of algorithm 1). If $j_m$ is congested, its latency $\mathbf{u}^j_{j_m} \geq T$. The GNN-based offloading policy is compared against the baseline and local (all tasks computed locally without offloading) policies, on a set of 10000 test instances, which is generated by creating 10 random task instances for each network instance in the test set.

The average execution latency and task congestion ratio as a function of the network size under the tested policies are presented in Fig 2(b). The task congestion ratio is the ratio between the total number of congested tasks and the total number of tasks for the 10 random task instances on each network instance. The average execution latency under the baseline policy is 288.7 compared to the GNN-based (18.4) and local (20.4) policies due to its high congestion ratio (up to $14.5\%$ in larger networks). Both the GNN-based and local policies can avoid task congestion under the test traffic configuration, whereas the GNN-based policy has the lowest execution latency across different network sizes, which on average is $9.8\%$ lower than the local policy.

In Fig. 2(c), we present the boxplot of the average per-task latency ratio of the GNN and local policies w.r.t. the baseline across network sizes, defined as $\mathbb{E}_{\mathcal{J}}(\hat{\mathbf{u}}^j_{j_m}/\bar{\mathbf{u}}^j_{j_m})$ per test instance. This metric can better describe the effects of a policy on free-flowing tasks under the baseline, as it is not dominated by the congested tasks like the average execution latency, e.g., $\hat{\mathbf{u}}^j_{j_m}/\bar{\mathbf{u}}^j_{j_m} \approx 0$ if task $j_m$ is congested under the baseline. The average per-task latency ratios of GNN-based and local policies are respectively 1.91 and 2.24. It shows that the GNN-based policy is in between the most conservative local policy that prevents congestion by avoiding offloading altogether, and the baseline policy that offloads tasks aggressively for faster execution but often causes congestion in the network.

## 5. CONCLUSIONS

In this paper, we develop a fully-distributed approach for computational offloading in wireless multi-hop networks. The key idea is to use graph neural networks to encode the information of computational tasks, links, and servers across the network into link weights, which can bring the awareness of other tasks across the network to distributed offloading and routing decisions that are otherwise agnostic to that information. In this way, we can lower the execution latency of tasks by reducing the congestion in the network. Compared with traditional methods, where each mobile device estimates the costs of its offloading options through probing packets, our approach can better understand the network context at lower communication overhead. Future work includes improving the training method and GNN design for better latency and energy performance.

---

[1] Training takes 0.5 hours on a workstation with a specification of 16GB memory, 8 cores, and Geforce GTX 1070 GPU. The source code is published at https://github.com/zhongyuanzhao/multihop-offload

# 6. REFERENCES

[1] S. K. Sarkar, T. G. Basavaraju, and C. Puttamadappa, *Ad hoc Mobile Wireless Networks: Principles, Protocols and Applications*. CRC Press, 2013.

[2] D. Kreutz, F. M. Ramos, P. E. Verissimo, C. E. Rothenberg, S. Azodolmolky, and S. Uhlig, "Software-defined networking: A comprehensive survey," *Proceedings of the IEEE*, vol. 103, no. 1, pp. 14–76, 2014.

[3] A. Kott, A. Swami, and B. J. West, "The internet of battle things," *Computer*, vol. 49, no. 12, pp. 70–75, 2016.

[4] "Cisco annual internet report (2018–2023)," white paper, Cisco Systems, Inc., Mar. 2020.

[5] I. F. Akyildiz, A. Kak, and S. Nie, "6G and beyond: The future of wireless communications systems," *IEEE Access*, vol. 8, pp. 133995–134030, 2020.

[6] X. Chen, D. W. K. Ng, W. Yu, E. G. Larsson, N. Al-Dhahir, and R. Schober, "Massive access for 5G and beyond," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 3, pp. 615–637, 2021.

[7] M. Noor-A-Rahim, Z. Liu, H. Lee, M. O. Khyam, J. He, D. Pesch, K. Moessner, W. Saad, and H. V. Poor, "6g for vehicle-to-everything (v2x) communications: Enabling technologies, challenges, and opportunities," *Proceedings of the IEEE*, vol. 110, no. 6, pp. 712–734, 2022.

[8] L. Tassiulas, "Stability properties of constrained queueing systems and scheduling policies for maximum throughput in multihop radio networks," *IEEE Trans. on Automatic Control*, vol. 31, no. 12, 1992.

[9] C. Joo, X. Lin, and N. B. Shroff, "Understanding the capacity region of the greedy maximal scheduling algorithm in multihop wireless networks," *IEEE/ACM Trans. Netw.*, vol. 17, no. 4, pp. 1132–1145, 2009.

[10] X. Lin and S. B. Rasool, "Constant-time distributed scheduling policies for ad hoc wireless networks," *IEEE Trans. on Automatic Control*, vol. 54, no. 2, pp. 231–242, 2009.

[11] L. Jiang and J. Walrand, "A distributed CSMA algorithm for throughput and utility maximization in wireless networks," *IEEE/ACM Trans. Netw.*, vol. 18, no. 3, pp. 960–972, 2010.

[12] Z. Zhao, G. Verma, C. Rao, A. Swami, and S. Segarra, "Link scheduling using graph neural networks," *IEEE Trans. Wireless Commun.*, vol. 22, no. 6, pp. 3997–4012, 2023.

[13] Z. Zhao, B. Radojicic, G. Verma, A. Swami, and S. Segarra, "Delay-aware backpressure routing using graph neural networks," in *IEEE Int. Conf. on Acoustics, Speech and Signal Process. (ICASSP)*, pp. 4720–4724, 2023.

[14] Z. Zhao, A. Swami, and S. Segarra, "Graph-based deterministic policy gradient for repetitive combinatorial optimization problems," in *Intl. Conf. Learn. Repres. (ICLR)*, 2023.

[15] D. Van Le and C.-K. Tham, "Quality of service aware computation offloading in an ad-hoc mobile cloud," *IEEE Trans. Vehicular Tech.*, vol. 67, no. 9, pp. 8890–8904, 2018.

[16] A. J. Ferrer, J. M. Marquès, and J. Jorba, "Towards the decentralised cloud: Survey on approaches and challenges for mobile, ad hoc, and edge computing," *ACM Computing Surveys (CSUR)*, vol. 51, no. 6, pp. 1–36, 2019.

[17] C. Funai, C. Tapparello, and W. Heinzelman, "Computational offloading for energy constrained devices in multi-hop cooperative networks," *IEEE Transactions on Mobile Computing*, vol. 19, no. 1, pp. 60–73, 2019.

[18] R. Chattopadhyay and C.-K. Tham, "Fully and partially distributed incentive mechanism for a mobile edge computing network," *IEEE Transactions on Mobile Computing*, vol. 21, no. 1, pp. 139–153, 2020.

[19] Y. Cai, J. Llorca, A. M. Tulino, and A. F. Molisch, "Mobile edge computing network control: Tradeoff between delay and cost," in *IEEE Global Communications Conference (GLOBECOM)*, pp. 1–6, 2020.

[20] G. Feng, X. Li, Z. Gao, C. Wang, H. Lv, and Q. Zhao, "Multipath and multi-hop task offloading in mobile ad hoc networks," *IEEE Trans. Vehicular Tech.*, vol. 70, no. 6, pp. 5347–5361, 2021.

[21] L. Liu, M. Zhao, M. Yu, M. A. Jan, D. Lan, and A. Taherkordi, "Mobility-aware multi-hop task offloading for autonomous driving in vehicular edge computing and networks," *IEEE Transactions on Intelligent Transportation Systems*, vol. 24, no. 2, pp. 2169–2182, 2022.

[22] X. Dai, Z. Xiao, H. Jiang, H. Chen, G. Min, S. Dustdar, and J. Cao, "A learning-based approach for vehicle-to-vehicle computation offloading," *IEEE Internet of Things Journal*, vol. 10, no. 8, pp. 7244–7258, 2022.

[23] X. Li, T. Chen, D. Yuan, J. Xu, and X. Liu, "A novel graph-based computation offloading strategy for workflow applications in mobile edge computing," *IEEE Transactions on Services Computing*, vol. 16, no. 2, pp. 845–857, 2022.

[24] W. Cheng, X. Cheng, T. Znati, X. Lu, and Z. Lu, "The complexity of channel scheduling in multi-radio multi-channel wireless networks," in *IEEE Intl. Conf. on Computer Comms. (INFOCOM)*, pp. 1512–1520, 2009.

[25] J. Yang, S. C. Draper, and R. Nowak, "Learning the interference graph of a wireless network," *IEEE Trans. Signal Inf. Process. Netw.*, vol. 3, no. 3, pp. 631–646, 2016.

[26] J. D. Little, "A proof for the queuing formula: L= λ w," *Operations research*, vol. 9, no. 3, pp. 383–387, 1961.

[27] T. Öncan, "A survey of the generalized assignment problem and its applications," *INFOR: Information Systems and Operational Research*, vol. 45, no. 3, pp. 123–141, 2007.

[28] F. Harary and R. Z. Norman, "Some properties of line digraphs," *Rendiconti del circolo matematico di palermo*, vol. 9, pp. 161–168, 1960.

[29] R. Bellman, "On a routing problem," *Quarterly of Applied Mathematics*, vol. 16, no. 1, pp. 87–90, 1958.

[30] L. R. Ford Jr, "Network flow theory," tech. rep., Rand Corp Santa Monica CA, 1956.

[31] A. Bernstein and D. Nanongkai, "Distributed exact weighted all-pairs shortest paths in near-linear time," in *ACM SIGACT Symp. Theory of Comp.*, pp. 334–342, 2019.

[32] R. Albert and A.-L. Barabási, "Statistical mechanics of complex networks," *Rev. Mod. Phys.*, vol. 74, pp. 47–97, Jan 2002.

[33] M. Stoer and F. Wagner, "A simple min-cut algorithm," *Journal of the ACM (JACM)*, vol. 44, no. 4, pp. 585–591, 1997.