

# Maximum Lifetime Analytics in IoT Networks

Víctor Valls\*, George Iosifidis\*, Theodoros Salonidis†

\*School of Computer Science and Statistics, Trinity College Dublin, Ireland

†IBM T.J. Watson Research Center, Yorktown Heights, NY, USA

**Abstract**—This paper studies the problem of allocating bandwidth and computation resources to data analytics tasks in Internet of Things (IoT) networks. IoT nodes are powered by batteries, can process (some of) the data locally, and the quality grade or performance of how data analytics tasks are carried out depends on where these are executed. The goal is to design a resource allocation algorithm that jointly maximizes the network lifetime and the performance of the data analytics tasks subject to energy constraints. This joint maximization problem is challenging with coupled resource constraints that induce non-convexity. We first show that the problem can be mapped to an equivalent convex problem, and then propose an online algorithm that provably solves the problem and does not require any a priori knowledge of the time-varying wireless link capacities and data analytics arrival process statistics. The algorithm’s optimality properties are derived using an analysis which, to the best of our knowledge, proves for the first time the convergence of the dual subgradient method with time-varying sets. Our simulations seeded by real IoT device energy measurements, show that the network connectivity plays a crucial role in network lifetime maximization, that the algorithm can obtain both maximum network lifetime and maximum data analytics performance in addition to maximizing the joint objective, and that the algorithm increases the network lifetime by approximately 50% compared to an algorithm that minimizes the total energy consumption.

## I. INTRODUCTION

We consider an Internet of Things (IoT) network where a set of nodes collect measurements that have to be later on analyzed by a data analytics or machine learning (ML) algorithm. For example, an algorithm for classifying, filtering or summarizing data. This class of services is one of the most important envisioned applications of the emerging IoT networks [1] and poses many technical challenges [2]; especially, when IoT networks operate subject to bandwidth, processing, and energy constraints.

Unlike previous generations of sensor networks, it is expected that IoT applications collect data at an unprecedented rate and that only a fraction of these will be non-ephemeral [3], [4]. Hence, the usual approach of transferring all the data to an IoT node gateway for processing, even if possible, may consume network resources unnecessarily. A promising solution to overcome this issue is to leverage the processing capacity of the IoT nodes and execute (some of) the data analytics tasks at the edge. That is, by carrying out some of the processing in situ it is possible to reduce the amount of data that needs to be transmitted over the network. However, executing data analytics at the IoT nodes has some costs. The two most significant ones are perhaps that processing is expensive in terms of energy, and that the grade or performance of how a

task is carried out depends on the algorithms that IoT nodes are able to run. For example, nodes with limited memory can only classify data with low complexity models, which makes their predictions less accurate in general.

In this paper, we introduce a resource allocation framework that allows us to design online execution and routing policies for data analytics tasks in IoT networks. That is, decide whether a data analytics task should be processed (i) locally, (ii) at the gateway, or (iii) elsewhere in the IoT network (*e.g.*, at a neighboring IoT node) depending on the available network resources (bandwidth, processing capacity and energy). An important feature of the framework is that it also allows us to maximize a combined criterion of *network lifetime*<sup>1</sup> and *data analytics performance*. The first objective is important because we would like the network to operate for as long as possible, and the second because the network has to fulfill its purpose besides “staying alive”. We call this problem Maximum Lifetime IoT Analytics (MLIA).

The problem of maximizing the time a network can operate has been addressed before in sensor networks (see Section II). However, unlike sensor networks, IoT encompasses more sophisticated scenarios where a variety of heterogeneous devices and applications coexist [2], and brings data analytics processing into play. The latter adds, technically, a new dimension to the existing routing algorithms (*e.g.*, [5], [6], [7], [8]) and raises technical challenges that cannot be addressed with previous solutions directly. In particular, to capture in-network processing we need to use a network model with “gains” (Section III-B), transform a non-convex problem into an equivalent convex one (Lemma 1), and develop a new online algorithm (Theorem 1) that can handle non-linearities in the utility as well as randomness in the actions (*i.e.*, the routing and processing decisions that can be made in each time slot). Furthermore, the fact that nodes can only operate during a limited time span<sup>2</sup> adds the difficulty of selecting the algorithm parameters so as to obtain the desired performance (see Section VI-B2). To this end, the main contributions of the paper are:

- (i) **MLIA problem**: we introduce the problem of jointly maximizing the lifetime of an IoT network and the performance of how data analytics tasks are carried out. This is an open problem arising in many IoT applications.
- (ii) **Problem model**: we formulate the MLIA problem as a convex optimization program. The model captures aspects such as the coexistence of different types of data analytics tasks, that data analytics may be carried out by different

This work was supported by Science Foundation Ireland under Grant No. 17/CDA/4760.

<sup>1</sup>Time the network can operate without any node running out of battery.

<sup>2</sup>Due to energy constraints.

IoT nodes, and that their computation cost may vary across nodes, among others.

- (iii) Online algorithm: we propose an online algorithm that solves the underlying convex problem and has non-asymptotic convergence guarantees. The algorithm determines the joint routing and processing policy for the IoT nodes in a myopic manner by only looking at the system's current state, *i.e.*, it does not require statistical knowledge of the underlying random processes such as the time-varying link capacities. Also, and to the best of our knowledge, this is the first paper that presents the convergence of the dual subgradient method with time-varying sets, which is another contribution.
- (iv) Performance evaluation: we perform a set of extensive experiments to (i) evaluate the performance of the proposed solution and to (ii) understand how this is affected by the network connectivity and system parameters. We also compare our algorithm to two benchmark policies, and show that our approach can increase the network lifetime by approximately 50% compared to an algorithm that minimizes the total energy consumption.

The rest of the paper is organized as follows. Section II presents the related work. In Section III, we introduce the system model, the problem, and the arising trade-offs. In Section IV, we formulate the MLIA problem as a convex program, and in Section V, we present the online algorithm. Section VI contains the numerical experiments and discussion.

## II. RELATED WORK

The problem of deciding how to transmit and process data to prolong the network lifetime has been studied before in Wireless Sensor Networks (WSNs). For instance, the work in [9] proposes load-balancing techniques to spread the energy consumption across nodes. In [10], the authors study the problem of designing a medium access protocol that takes into account the channel state information and the available energy. And in [11], it is shown that preprocessing data at the sensor nodes can help to reduce the network load and so the energy required to transmit data to the fusion center (*i.e.*, the gateway). From a problem formulation perspective, our approach differs from previous works in the literature of WSNs and IoT because we consider both the energy spent in routing and processing the data, and the performance of the analytics. The last point is crucial in heterogeneous IoT networks where the ability of nodes to run algorithms depends on their hardware.

The maximum lifetime objective has been extensively studied for routing in multi-hop sensor networks. The seminal work in [6] considered a static maximum lifetime routing problem and formulated it as a linear program. This work was extended in a sequel of papers [12], [5] to different types of wireless networks, and a *flow augmentation* algorithm was proposed to support fixed and arbitrary generation rates. The approach in [6] has also been adopted by other authors (see survey [13]); for example, [14] combines network lifetime with congestion control, and [7] proposes an algorithm to solve the maximum lifetime problem in a distributed fashion. Finally,

[15] proposed a static optimization solution for maximizing analytics performance in IoT networks with average power constraints, which is different than the lifetime criterion. Our work differs technically from all previous work on maximum lifetime routing in sensor networks because (i) it considers IoT sensing nodes that are heterogeneous and can perform data analytics computations in addition to routing data; (ii) it uses an objective that incorporates performance in addition to lifetime; and (iii) the proposed dynamic algorithm solves, provably, the underlying convex problem without knowledge of the arrivals or channel statistics. Regarding the algorithm, we use time-varying sets to handle the instantaneous routing and processing constraints, which is in marked contrast to stochastic Lagrange dual approaches where the stochasticity does not affect the decision variables.

Finally, we note that data analytics optimization with routing costs has been considered in the cloud context [16], [17]. However, there the costs are not related to energy expenditure, neither the resource allocation decisions affect the time the network will be able to operate. Furthermore, in IoT, we have the additional inherent difficulties of routing data over wireless networks which do not appear, of course, in the cloud.

## III. SYSTEM MODEL AND PROBLEM STATEMENT

### A. System model

1) *Network*: We model an IoT network as a directed graph  $\mathcal{G} = (\mathcal{N}, \mathcal{E})$  consisting of  $n = |\mathcal{N}|$  nodes and  $l = |\mathcal{E}|$  links. We use  $\mathcal{C}$  to denote the set of applications in the network. For each application  $c \in \mathcal{C}$ , the nodes collect and send data to a gateway node—possibly over a multi-hop path. Parameter  $\lambda_i^{(c)}$ ,  $i \in \mathcal{N}$ ,  $c \in \mathcal{C}$  indicates the rate at which node  $i$  collects data of application  $c$ . These are exogenously given and injected into the network directly. The arrival rate matrix is given by

$$\boldsymbol{\lambda} = (\lambda_i^{(c)} : i \in \mathcal{N}, c \in \mathcal{C}). \quad (1)$$

Each link  $(i, j) \in \mathcal{E}$  has an average capacity of  $\mu_{ij}$  bits/s; we collect these in matrix  $\boldsymbol{\mu} = (\mu_{ij} : i, j \in \mathcal{N})$ . Links  $(i, j)$  and  $(j, i)$  can have different capacities. The transmission over link  $(i, j)$  induces an energy consumption  $e_{ij}^{\text{tx}}$  for the sender node  $i$ , and  $e_{ij}^{\text{rx}}$  for the receiver node  $j$ . Both are measured in J/bit. We model the wireless interference using the protocol interference model [18], according to which a transmission over link  $(i, j) \in \mathcal{E}$  is successful if and only if all nodes in range with  $i$  or  $j$  are idle. This requirement is based on the CSMA/CA protocol adopted by IEEE 802.11 standards.<sup>3</sup> The set of interfering links for each link  $(i, j) \in \mathcal{E}$  is defined as:

$$I(i, j) := \{(a, d), (d, b) : d \in \mathcal{N}_i \cup \mathcal{N}_j, a, b \in \mathcal{N}_d\}, \quad (2)$$

where  $\mathcal{N}_i := \{j \in \mathcal{N} : (i, j) \in \mathcal{E}\}$  is the set that contains the neighbors of a node  $i \in \mathcal{N}$ . Hence, when a link  $(i, j)$  is active none of the links in  $I(i, j)$  can be used.

<sup>3</sup>In detail, this interference model complies with the communication sequence RTS-CTS-Data-ACK, where each sender is also a receiver of the ACK packets; which is the strictest model. The set of interfering nodes can be reduced if the ACK operation is not used.

2) *Nodes*: IoT nodes may be heterogenous in terms of hardware. We use  $\rho_i$  to denote the processing capacity in FLOPS of a node  $i \in \mathcal{N}$ , and  $\gamma_i^{(c)}$  to denote the processing requirements (in FLOPS/bit) of an application in each of the nodes. The processing in each node may reduce the volume of each flow by a factor of  $0 \leq \beta_i^{(c)} \leq 1$ . For example, in an object recognition application the flow volume reduction is large since an image gets reduced to a collection of bounding boxes and tags [19], *i.e.*, to few bytes. Of course, this flow reduction factor may depend on where a data analytics task is carried out since IoT nodes may run different algorithms. We collect the flow reduction and processing requirements in matrices  $\beta = (\beta_i^{(c)} : i \in \mathcal{N}, c \in \mathcal{C})$  and  $\gamma = (\gamma_i^{(c)} : i \in \mathcal{N}, c \in \mathcal{C})$ .

Each IoT node  $i \in \mathcal{N}$  has an energy budget of  $E_i$  Joules that can spend transmitting, receiving, and processing data. When a node has used its energy budget, it dies meaning that it cannot transmit or process more data. We assume that the network gateway does not have energy constraints.

### B. Decision variables and constraints

IoT nodes can make two types of decisions: process and forward data. Variable  $x_{ij}^{(c)} \geq 0$  indicates the rate in bits/s at which application  $c$  is transmitted over link  $(i, j)$ . Similarly, variable  $y_i^{(c)}$  indicates the rate in bits/s that node  $i$  processes data of application  $c$ . We collect the decision variables in matrices  $\mathbf{x} = (x_{ij}^{(c)} : (i, j) \in \mathcal{E}, c \in \mathcal{C})$  and  $\mathbf{y} = (y_i^{(c)} : i \in \mathcal{N}, c \in \mathcal{C})$ . The transmission and processing rates must satisfy the link and node capacity constraints:

$$\sum_{c \in \mathcal{C}} x_{ij}^{(c)} \leq \mu_{ij}, \quad \sum_{c \in \mathcal{C}} \gamma_i^{(c)} y_i^{(c)} \leq \rho_i \quad \forall i \in \mathcal{N}, (i, j) \in \mathcal{E}. \quad (3)$$

The interference constraints affect how the links in  $\mathcal{E}$  can be activated and consequently the total amount of data that can be transferred over the network; see, for instance, Lemma 1 in [20]. These are formally given by:

$$\sum_{c \in \mathcal{C}} \frac{x_{ij}^{(c)}}{\mu_{ij}} + \sum_{(k,m) \in I(i,j)} \sum_{c \in \mathcal{C}} \frac{x_{km}^{(c)}}{\mu_{km}} \leq 1 \quad \forall (i, j) \in \mathcal{E}, \quad (4)$$

where  $I(i, j)$  is defined in (2).

Unlike classic max-flow type models [21], in processing-capable networks the amount of flow that arrives at a node may not be the same as the amount of flow that departs. Now, the flow conservation constraints are given by

$$\sum_{j \in \mathcal{N}_i} x_{ji}^{(c)} + \beta_i^{(c)} y_i^{(c)} + \lambda_i^{(c)} = \sum_{j \in \mathcal{N}_i} x_{ij}^{(c)} + y_i^{(c)}. \quad (5)$$

Equation (5) says that the data received from the other nodes in the network, plus the data received after the local processing, plus the exogenous arrivals must be equal to the traffic sent to other nodes in the network and for processing. Figure 1 shows an example of a network with five IoT nodes and a gateway.

Finally, the decisions  $\mathbf{x}$  and  $\mathbf{y}$  need to respect the energy budget of each node. These must satisfy:

$$T_s p_i(\mathbf{x}, \mathbf{y}) \leq E_i \quad i \in \mathcal{N} \quad (6)$$

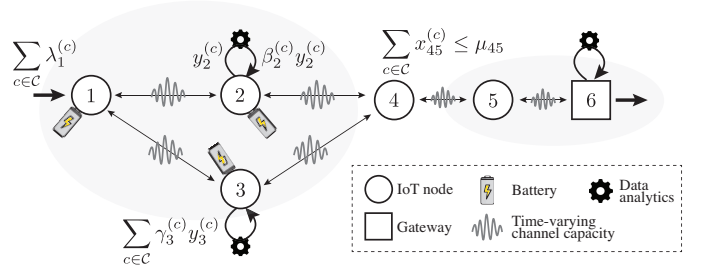


Figure 1. Illustrating an example of an IoT network with energy constraints. Node 1 receives data analytics tasks that must be processed either at node 2, node 3, or at the gateway. Nodes 1, 2 and 3 have energy constraints and the transmission of data over the network is affected by the interference (large shaded areas) and link capacity constraints.

where

$$p_i(\mathbf{x}, \mathbf{y}) := \sum_{j \in \mathcal{N}_i} \sum_{c \in \mathcal{C}} x_{ij}^{(c)} e_{ij}^{\text{tx}} + \sum_{j \in \mathcal{N}_i} \sum_{c \in \mathcal{C}} x_{ji}^{(c)} e_{ji}^{\text{rx}} + \sum_{c \in \mathcal{C}} y_i^{(c)} e_i^{\text{pr}}$$

and  $T_s$  is the network life time:

**Definition 1** (Network lifetime). *Time span where all the IoT nodes are alive, i.e., none of them runs out of energy.*

Hence, (6) must be satisfied by *all* the nodes in the network. We later extend this definition to scenarios where a subset of nodes are allowed to die before the network lifetime expires.

### C. Utilities

The execution of analytics at the IoT nodes and the gateway results in different performance. The latter can be related, for instance, to the precision or confidence that an image is classified using a machine learning algorithm. Formally, we define function  $\omega_i^{(c)} : \mathbf{R} \rightarrow \mathbf{R}$ ,  $i \in \mathcal{N}, c \in \mathcal{C}$  to capture the network benefit of processing an application at a node. We assume  $\omega_i^{(c)}$  is concave and non-decreasing, and that IoT nodes may have different performance for processing the same task. This performance diversity may arise due to hardware or software differences among the IoT nodes. See [19, Figure 1] for an example of how different object detection algorithms have different precisions.

### D. Problem statement

The IoT nodes collect data that must be processed locally, at another IoT node, or at the network gateway. Ideally, we would prefer to process all the data in the gateway since this has usually an “unlimited” energy budget and better hardware that allows it to run more sophisticated algorithms. However, that might not be always possible for several reasons. First, the amount of bandwidth available between the nodes and the gateway might not be enough to transport all the data. Second, nodes are subject to energy constraints, which in turn limit the total amount of data that can be transferred. And third, the central node may not be able to cope with the processing of all the data from all the IoT nodes.

Our goal is to design a routing and processing policy  $(\mathbf{x}, \mathbf{y})$  that maximizes (i) the network lifetime and (ii) the analytics performance. The lifetime criterion is crucial because when nodes run out of battery, the network operation is disrupted

e.g., nodes stop collecting information, monitoring an area, among others. On the other hand, the analytics performance metric is important because our goal is not just to keep the nodes in the network alive but also to maximize the service performance. Sometimes these two objectives may not be conflicting since processing data in the network reduces the network congestion and so the subsequent routing energy cost. However, in-network processing incurs an energy cost as well as analytics performance degradation since IoT nodes may run “lighter” algorithms that fit their hardware specifications. The overall balance depends on many parameters, such as the flow reduction due to processing, the energy cost of these tasks, and the network properties.

In sum, we would like to find a policy  $(\mathbf{x}, \mathbf{y})$  that maximizes a combined criterion of network lifetime and data analytics performance depending on (i) the capacity of network links; (ii) the nodes’ processing capacity; (iii) the applications data rates; (iv) the data analytics performance “quality”; and (v) the nodes’ energy budget. Next, we present the mathematical formulation of the problem and introduce an algorithm that is amenable to implementation in stochastic environments.

#### IV. MAXIMUM NETWORK LIFETIME AND ANALYTICS

The IoT operation constraints determine the transmission and processing policies that are implementable, and subsequently, the set of supportable data rates. First, we define sets:

$$\begin{aligned} X &:= \{\mathbf{x} \mid \text{constraints (3) and (4) are satisfied}\}, \\ Y &:= \{\mathbf{y} \mid \text{constraint (3) is satisfied}\}. \end{aligned}$$

These sets are bounded polytopes (so convex sets [22, Section 2.2.4]) because they are the intersection of inequality constraints and all links and nodes have bounded capacity. Using these sets, we can define the capacity region of the IoT system.

**Definition 2** (IoT network capacity region). *The IoT network capacity region is the set*

$$\Gamma(\lambda) := \{\mathbf{x} \in X, \mathbf{y} \in Y \mid \text{constraints (5), (6) are satisfied}\}$$

where  $\lambda$  is given in (1).

We will always assume that  $\Gamma(\lambda)$  is non-empty. We are now in position to formulate our optimization problem.

##### A. Maximum lifetime data analytics problem

Recall that  $T_s$  is the network lifetime. The optimization problem is given by

$$\begin{aligned} &\underset{T_s, \mathbf{x}, \mathbf{y}}{\text{maximize}} && (1 - \eta)T_s + \eta \sum_{i \in \mathcal{N}, c \in \mathcal{C}} \omega_i^{(c)}(y_i^{(c)}) \\ &\text{subject to} && (5), (6) \quad \forall i \in \mathcal{N} \\ &&& \mathbf{x} \in X, \mathbf{y} \in Y, T_s \geq 0 \end{aligned} \quad (7)$$

Parameter  $\eta \in [0, 1]$  is used to balance how much we prioritize the network lifetime over the analytics performance metric. If  $\eta = 0$ , then (7) aims only to maximize the network lifetime; if  $\eta = 1$ , the optimization only takes into account the data analytics performance; and for any value of  $\eta \in (0, 1)$ , it balances the two terms in the objective.

Problem (7) is non-convex because the inequality constraint (6) involves the product of variables  $\mathbf{x}$  and  $T_s$ .<sup>4</sup> Moreover, variables  $(\mathbf{x}, \mathbf{y})$  and  $T_s$  are *not* independent since  $(\mathbf{x}, \mathbf{y})$  affect the energy consumption and so indirectly the time the network will be able to operate. Nonetheless, and as we show in the following lemma, it is possible to transform the non-convex problem (7) into an “equivalent” convex problem.

**Lemma 1.** *The optimization problem*

$$\begin{aligned} &\underset{\mathbf{x}, \mathbf{y}}{\text{minimize}} && (1 - \theta) \max_{i \in \mathcal{N}} \left\{ \frac{p_i(\mathbf{x}, \mathbf{y})}{E_i} \right\} + \theta \sum_{i \in \mathcal{N}, c \in \mathcal{C}} -\omega_i^{(c)}(y_i^{(c)}) \\ &\text{subject to} && (5), \mathbf{x} \in X, \mathbf{y} \in Y \end{aligned} \quad (8)$$

is convex and “equivalent” to the non-convex problem in (7). It allows us to balance smoothly between maximum network lifetime and maximum data analytics performance.

*Proof:* Let  $(T_s^*, \mathbf{x}^*, \mathbf{y}^*)$  be a solution to (7). By Weierstrass’ theorem [23, Proposition 2.1.1; condition 1], one can show that a solution always exists and is finite.<sup>5</sup> The key part of the proof relies on showing that one of the energy constraints must be tight at the optimum. We proceed to show this by contradiction. Suppose there exists a  $T' > T_s^*$  such that  $T_s^* p_i(\mathbf{x}, \mathbf{y}) < T' p_i(\mathbf{x}^*, \mathbf{y}^*) \leq E_i$  for all  $i \in \mathcal{N}$ . Then, we must also have that  $(1 - \eta)T_s^* < (1 - \eta)T'$  (i.e., the objective value increases); however, this is not possible since by assumption  $T_s^*$  is an optimal value. Hence, we have that at least one constraint must be tight at the optimum, i.e.,  $T_s^* p_i(\mathbf{x}^*, \mathbf{y}^*) = E_i$ . Next, rearrange terms in the last equation and rewrite the energy constraints as

$$\max_{i \in \mathcal{N}} \{p_i(\mathbf{x}^*, \mathbf{y}^*)/E_i\} = 1/T_s^*, \quad (9)$$

where the max follows because the constraint must be satisfied by all the nodes. From (9), we can see that maximizing  $T_s^*$  is equivalent to minimizing the LHS of (9)—which is a convex function since  $p_i(\mathbf{x}, \mathbf{y})$  is linear [22, pp. 72] and  $E_i$  is a constant. Finally, considering that  $\sum_{i \in \mathcal{N}, c \in \mathcal{C}} -\omega_i^{(c)}(y_i^{(c)})$  is convex, we can use scalarization with  $\theta \in [0, 1]$  to obtain a convex problem. ■

It is important to emphasize that problems (7) and (8) are equivalent but not the same. Namely, if  $\theta = 0$ , then the optimization maximizes the network life time; if  $\theta = 1$ , it only takes into account the analytics performance; and if  $\theta \in (0, 1)$  it balances the two objectives smoothly. However, the solutions to problems (7) and (8) do *not* need to be the same for  $\theta = \eta$  when  $\theta \in (0, 1)$  since the network lifetime term in (7) is linear, whereas in (8), the term must be regarded as  $1/T_s$ . As we will show in detail in Section VI, parameter  $\theta$  has a huge impact on the system’s performance.

##### B. Generalized maximum network lifetime problem

Instead of considering the lifetime of individual nodes, we can consider the lifetime of groups of nodes  $\mathcal{M}_k \subset \mathcal{N}$ ,  $k \in$

<sup>4</sup>The product of two variables is generally not convex.

<sup>5</sup>Note that by construction  $X$  and  $Y$  are bounded sets. Variable  $T_s$  belongs also to a bounded set since it is nonnegative and nodes have a finite energy.

$\{1, 2, \dots, K\}$ . For instance, the type of data collected by a subset of nodes may be more important to the global system objective than the data collected by another subset of nodes. To consider groups of nodes in the optimization problem, we can replace the network lifetime term in the objective of problem (8) with

$$\sum_{k=1}^K \pi_k \max_{i \in \mathcal{M}_k} \{p_i(\mathbf{x}, \mathbf{y})/E_i\}.$$

Parameter  $\pi_k \geq 0$  is used to emphasize how much we would like a subset of nodes to “remain” alive with respect to another subset. Note that the formulation in (8) is a special case where  $|\mathcal{M}_k| = 1$  for all  $k \in \{1, \dots, K\}$  and  $\cup_{k=1}^K \mathcal{M}_k = \mathcal{N}$ . Hereafter, and to streamline exposition, we will use the formulation in (8) but our results apply to more general cases directly.

### C. Practical limitations

The resulting optimization problem is readily solvable by convex solvers such as SCS [24]. However, for this, one needs to know all the system parameters, which is rarely the case in practical scenarios. For example, the average capacity of a link connecting two IoT nodes is usually not known. Furthermore, there are generally instantaneous constraints, such as the level of noise in the system or interference, which affect the decisions that the IoT nodes can make. In the next section, we present a dynamic algorithm that learns and adapts to the instantaneous network/system characteristics.

## V. DYNAMIC MLIA ALGORITHM

In this section, we present a dynamic algorithm that solves (8) and does not require previous knowledge on (i) the average arrival rate of the data analytics; (ii) the capacity of the network links  $\mu$  in each time instance; and (iii) the average performance or reward obtained from carrying out data analytics at the IoT nodes.

### A. Dynamic problem formulation

We divide the time in slots  $t \in \mathbb{N}$  of normalized duration and parameterize the variables in the static model (8) with  $[t]$  to indicate their value in a time slot. For instance,  $\lambda_i^{(c)}[t]$  indicates the new arrivals of application  $c \in \mathcal{C}$  at node  $i$  and time  $t \in \mathbb{N}$ , and  $x_{ij}^{(c)}[t]$  is the amount of commodity transmitted over link  $(i, j) \in \mathcal{E}$ . We also need to capture the network constraints of the dynamic problem. The instantaneous link capacity constraints are

$$\sum_{c \in \mathcal{C}} x_{ij}^{(c)}[t] \leq \mu_{ij}[t], \quad \sum_{c \in \mathcal{C}} \gamma_i^{(c)} y_i^{(c)}[t] \leq \rho_i[t] \quad (10)$$

$\forall (i, j) \in \mathcal{E}, i \in \mathcal{N}$ . The (binary) interference constraints are

$$\mathbb{I} \left( \sum_{c \in \mathcal{C}} x_{ij}^{(c)}[t] \right) + \sum_{(k, m) \in I(i, j)} \mathbb{I} \left( \sum_{c \in \mathcal{C}} x_{km}^{(c)}[t] \right) \leq 1, \quad (11)$$

for all  $(i, j) \in \mathcal{E}$  where  $\mathbb{I}$  is the indicator function, i.e.,  $\mathbb{I}(x) = 1$  if  $x > 0$  and  $\mathbb{I}(x) = 0$  otherwise. Hence, the binary interference constraints only allow one node to transmit in the

interference range regardless of the transmission rate. We are now in position to define sets

$$\begin{aligned} X[t] &:= \{\mathbf{x} \mid \text{constraints (10) and (11) are satisfied}\}, \\ Y[t] &:= \{\mathbf{y} \mid \text{constraint (10) is satisfied}\}, \end{aligned}$$

which contain the admissible routing/processing policies that can be implemented in the system in each time slot. Note that  $X[t]$  may not be convex since (11) is not convex.

The operation of the network consists of selecting actions (or values) from the instantaneous action sets  $X[t]$  and  $Y[t]$  while taking into account the properties of the underlying convex problem. We explain next how this can be achieved.

### B. Algorithm overview

The key idea for solving (8) in the dynamic setting is to relax the flow conservation constraints and formulate the Lagrange dual problem. The Lagrangian is

$$\begin{aligned} L(\mathbf{x}, \mathbf{y}, \boldsymbol{\nu}) &= (1 - \theta) \max_{i \in \mathcal{N}} \left\{ \frac{p_i(\mathbf{x}, \mathbf{y})}{E_i} \right\} + \theta \sum_{i \in \mathcal{N}, c \in \mathcal{C}} -\omega_i^{(c)}(y_i^{(c)}) \\ &+ \sum_{i \in \mathcal{N}} \sum_{c \in \mathcal{C}} \nu_i^{(c)} \left( \sum_{j \in \mathcal{N}_i} (x_{ji}^{(c)} - x_{ij}^{(c)}) + (\beta_i^{(c)} - 1)y_i^{(c)} + \lambda_i^{(c)} \right) \end{aligned}$$

where  $\nu_i^{(c)} \in \mathbf{R}$  is the Lagrange multiplier associated to each of the flow conservation constraints, and we define  $\boldsymbol{\nu} = (\nu^{(c)} : i \in \mathcal{N}, c \in \mathcal{C})$ . The Lagrange dual function is defined as

$$h(\boldsymbol{\nu}) := \min_{\mathbf{x} \in X, \mathbf{y} \in Y} L(\mathbf{x}, \mathbf{y}, \boldsymbol{\nu}), \quad (12)$$

and recall that it is concave [22, Section 5.2]. Hence, it can be maximized with the subgradient method applied to the Lagrange dual problem.<sup>6</sup> Specifically, with update:

$$\begin{aligned} \nu_i^{(c)}[t+1] &= \nu_i^{(c)}[t] + \alpha \left( \sum_{j \in \mathcal{N}_i} (x_{ji}^{(c)}[t] - x_{ij}^{(c)}[t]) \right. \\ &\quad \left. + (\beta_i^{(c)} - 1)y_i^{(c)}[t] + \lambda_i^{(c)} \right) \quad \forall c \in \mathcal{C}, i \in \mathcal{N}, \end{aligned} \quad (13)$$

where the term in parenthesis is a subgradient of  $h$  at  $\boldsymbol{\nu}[t]$ ,

$$(\mathbf{x}[t], \mathbf{y}[t]) \in \arg \min_{\mathbf{x} \in X, \mathbf{y} \in Y} L(\mathbf{x}, \mathbf{y}, \boldsymbol{\nu}), \quad (14)$$

and  $\alpha > 0$  the step size or parameter that controls the accuracy of solution in the optimization. Recall that the solution of the primal and dual problem coincides when strong duality holds [22, Section 5.2.3]; which is the case in our problem since the constraints in (8) are linear and the network capacity or feasible set  $\Gamma(\boldsymbol{\lambda})$  is non-empty by assumption.

Next, we proceed to explain how to include constraints (10)-(11) so that our dynamic algorithm can be implemented in a real system. We do this incrementally from the static problem.

<sup>6</sup>We use subgradient instead of gradient because the Lagrange dual function may not be differentiable.

1) *Unknown arrivals*: Suppose now that  $X[t] = X$  and  $Y[t] = Y$  for all  $t \in \mathbf{N}$  (the sets of actions do not change over time). We want to relax the fact that the data analytics arrival  $\lambda$  is not known a priori. For this, we need to solve the Lagrange dual problem (*i.e.*, maximize  $h(\nu)$ ) with the subgradient method replacing  $\lambda_i^{(c)}$  with the random variable  $\lambda_i^{(c)}[t]$  in the dual update (13). This change amounts to making the subgradient of  $h(\nu)$  stochastic, or equivalently, to using the stochastic dual subgradient method to solve the underlying convex problem (8). For the algorithm to converge we need to make the following assumption.

**Assumption 1.**  $\{\lambda_i^{(c)}[t]\}$  is an i.i.d. process with expected value  $\lambda_i^{(c)}$  for all  $i \in \mathcal{N}$ ,  $c \in \mathcal{C}$ ,  $t \in \mathbf{N}$  and  $\lambda_i^{(c)}[t]$  is uniformly bounded for all  $t \in \mathbf{N}$ .

2) *Time-varying sets*: We can incorporate time-varying sets in our algorithm by replacing  $X$  and  $Y$  with  $X[t]$  and  $Y[t]$  in update (14). In order to solve the underlying problem (8), we need to make the following assumption, which can be regarded as if we used stochastic subgradients in update (13).

**Assumption 2.**  $\mathbf{E}(X[t]) = X$  and  $\mathbf{E}(Y[t]) = Y$  for all  $t \in \mathbf{N}$  and  $X[t]$  and  $Y[t]$  are bounded sets. Here, the expectations are defined using (Minkowski) set addition [25, pp. 32].

3) *Noisy utility function*: We may not have access to the utility function that measures the *exact* reward derived from processing analytics locally. To capture this, we define  $\tilde{\omega}_i^{(c)} : \mathbf{R} \rightarrow \mathbf{R}$ ,  $i \in \mathcal{N}$ ,  $c \in \mathcal{C}$  to be an estimate of utility function  $\omega_i^{(c)}$ . Also, let  $\tilde{L}(\mathbf{x}, \mathbf{y}, \nu)$  be the Lagrangian where  $\tilde{\omega}_i^{(c)}$  is used instead of  $\omega_i^{(c)}$ , and replace the Lagrangian in (14) with  $\tilde{L}(\mathbf{x}, \mathbf{y}, \nu)$ . Technically, minimizing  $\tilde{L}(\mathbf{x}, \mathbf{y}, \nu)$  instead of  $L(\mathbf{x}, \mathbf{y}, \nu)$  corresponds to computing (sub)gradients of the Lagrange dual function approximately, *i.e.*, having  $\epsilon$ -subgradient [26, pp. 625]. The convergence depends on the mild assumption that the errors are bounded.

**Assumption 3.** The maximum error between  $\omega_i^{(c)}$  and  $\tilde{\omega}_i^{(c)}$  is bounded, *i.e.*,  $\max_{y \leq p_i} |\omega_i^{(c)}(y) - \tilde{\omega}_i^{(c)}(y)| := \xi_i^{(c)} < \infty$ .

### C. Algorithm & convergence

Algorithm 1 consists of three steps. The first one is to obtain the network state or set of possible actions, *i.e.*,  $X[t]$  and  $Y[t]$ . We assume these are obtained by the IoT nodes and that the information is transmitted to the network gateway where the routing decisions are made. The second step is to minimize the Lagrangian using sets  $X[t]$  and  $Y[t]$ . The complexity of this step depends on the number of elements in the sets. When these are discrete and contain few elements, the minimization can be carried out by exhaustive search, and if  $X[t]$  and  $Y[t]$  are convex by standard convex optimization techniques [22], [23]. The third and final step is to update the dual variables, *i.e.*, carry out the (stochastic) subgradient update with  $\epsilon$ -subgradients. Parameter  $\alpha > 0$  is the step size that controls the accuracy/speed tradeoff of the algorithm [23].

We establish the convergence of the algorithm with respect to an optimal policy implemented in a random fashion. We need to make the following definitions.

**Definition 3** (Average policy space).

$$X := \lim_{T \rightarrow \infty} \frac{1}{T} \bigoplus_{t=1}^T X[t], \quad Y := \lim_{T \rightarrow \infty} \frac{1}{T} \bigoplus_{t=1}^T Y[t]$$

where  $\oplus$  denotes the (Minkowski) set addition [25, pp. 32].

**Definition 4** (Optimal policy). An optimal policy is a pair  $(x^*, y^*) \in (X, Y)$  that solves (8). By construction we have

$$x^* = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T x^*[t], \quad y^* = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T y^*[t]$$

where  $x^*[t] \in X[t]$  and  $y^*[t] \in Y[t]$  for all  $t \in \mathbf{N}$ .

**Theorem 1.** Let  $f$  be the objective function in the optimization problem (8). Suppose Assumptions 1, 2 and 3 are satisfied. Algorithm 1 ensures that

- (i)  $\frac{1}{T} \sum_{t=1}^T \mathbf{E} (f(x[t], y[t]) - f(x^*[t], y^*[t])) \leq \alpha \epsilon_1 + \epsilon_2$ ,
- (ii)  $\lim_{T \rightarrow \infty} \mathbf{E} \left( (\bar{x}_{ji}^{(c)} - \bar{x}_{ij}^{(c)}) - (\beta_i^{(c)} - 1) \bar{y}_i^{(c)} + \lambda_i^{(c)}[t] \right) = 0$ ,

$\forall i \in \mathcal{N}$ ,  $c \in \mathcal{C}$  where  $\bar{\mathbf{x}} := \frac{1}{T} \sum_{t=1}^T \bar{\mathbf{x}}[t]$ ,  $\bar{\mathbf{y}} = \frac{1}{T} \sum_{t=1}^T \bar{\mathbf{y}}[t]$ ;  $\epsilon_1$  is a bounded constant related to the variance of the arrival process and sets  $X[t]$  and  $Y[t]$ , and  $\epsilon_2$  a constant such that  $|\xi_i^{(c)}|_2 \leq \epsilon_2$ .

Theorem 1 establishes an upper bound on the objective function of our policy  $(\mathbf{x}[t], \mathbf{y}[t])$  compared to an optimal policy  $(\mathbf{x}^*, \mathbf{y}^*)$  implemented in a “randomized” fashion—that is, the policy that we would implement if we knew in advance all the parameters in the system. Note that the bound on the optimality gap depends on  $\epsilon_1$  and  $\epsilon_2$ . The first term is related to the variation of the arrivals and link capacities, and the second on how good the estimate of the reward functions are. Importantly, note that the first term depends on the step size  $\alpha$ , which means that we can make it arbitrarily small. Theorem 1 establishes also that the flow conservation constraints are satisfied on expectation as  $T \rightarrow \infty$ , which means that we will recover, asymptotically, a policy in the network capacity region  $\Gamma(\lambda)$ .

### D. Practical aspects

1) *Non-asymptotic analysis*: Differently from previous stochastic network optimization analyses, *e.g.*, [27], [28], our results do *not* compare the performance of the average policy. Furthermore, and unlike these previous works that give optimality bounds only asymptotically (*i.e.*, as  $T \rightarrow \infty$ ), here we provide guarantees on how the policy performs on expectation per time slot. This is very important for this problem because IoT nodes have a *finite* energy budget that may not allow them to reach a “steady” state—this will be illustrated in the experiments in Section VI-B2.

2) *Distributed execution*: MLIA describes a centralized algorithm where the IoT gateway collects the nodes’ parameters to solve (15). However, the only necessary central calculation is that for devising an eligible link activation schedule (based on the interference constraints). If such a schedule is already



**Algorithm 1** Maximum Lifetime IoT Analytics (MLIA)**Parameters:** step size  $\alpha \geq 0$ Initialize  $t = 1$ ;  $\nu_i^{(c)}[t] = 0$ , for all  $c \in \mathcal{C}, i \in \mathcal{N}$ .**In each time slot**  $t = 1, 2, \dots$ 

- 1) **Obtain network state:** the network gateway collects the network connectivity information and computes sets  $X[t]$  and  $Y[t]$ .
- 2) **Compute action:** the network gateway obtains

$$(\mathbf{x}[t], \mathbf{y}[t]) \in \arg \min_{\mathbf{x} \in X[t], \mathbf{y} \in Y[t]} \tilde{L}(\mathbf{x}, \mathbf{y}, \alpha \nu[t]) \quad (15)$$

and broadcast the solution to the IoT nodes.

- 3) **Update the dual variables:** for all  $\nu_i^{(c)}$ ,  $i \in \mathcal{N}$ ,  $c \in \mathcal{C}$  perform update (13) with  $\lambda_i^{(c)}[t]$  instead of  $\lambda_i^{(c)}$ .

**end loop**

given or if the links are orthogonal (or, if any other interference control scheme is used) then each node can independently optimize its actions. Namely, each node  $i$  can calculate separately the quantity  $\partial \tilde{L}(\mathbf{x}, \mathbf{y}, \nu) / \partial x_{ij}$ ,  $\forall j \in \mathcal{N}_i$ , and similarly for  $y_{ij}$  variables; and then each pair of 1-hop neighbors exchange the respective dual variables. This is a standard approach for enabling a decentralized implementation of such protocols, see [25] for a survey, and applies directly to MLIA.

## VI. EXPERIMENTS AND EVALUATION

We evaluate the performance of our approach numerically using real hardware and application parameters. We investigate three points. First, how parameter  $\theta$  and the network connectivity affect the network lifetime and the analytics performance in the static or offline problem (8). Second, the convergence of the *dynamic* algorithm to the solution of the static problem depending on  $\alpha$ . And third, how the proposed algorithm compares to two benchmark algorithms. We fix the network size in the experiments to 20 nodes (including the gateway), but similar results are obtained with different network sizes.

## A. Experiments setup

1) **Network:** In all simulations, we use random geometric graphs (RGG) to model an IoT network—this is common practice in wireless networks [29]. Recall that in RGGs nodes are placed uniformly at random on an area of  $1 \times 1$  normalized units, and that their connectivity depends on the *radius* or *distance* each node covers. For simplicity, we assume that *all* links have an average capacity of 24 Mbits.

2) **Nodes:** The IoT nodes are Raspberry Pi(es) 3B equipped with an ARMv7 CPU, 1 GB RAM and a 802.11.b/g/n network card. According to our measurements, the power required to transmit and receive data is 0.4 W (with small variations depending on the channel quality), and 2.1 W for processing at full power. We assume that the energy spent to collect data and in idle mode is negligible.<sup>7</sup> All IoT nodes have batteries with a capacity of 2500 J;<sup>8</sup> the gateway is connected to a constant energy source and so does not have energy constraints.

<sup>7</sup>That is, IoT nodes can switch to low-power consumption mode.

<sup>8</sup>The batteries have a small capacity to keep simulations short. An IoT node can easily be equipped with a battery that has one hundred times more energy.

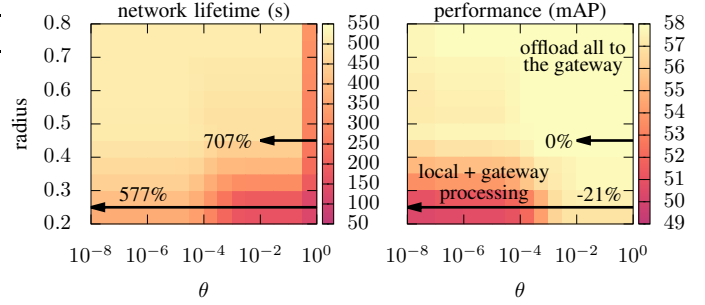


Figure 2. Illustrating the sensitivity of the solution depending on parameter  $\theta$  and the network radius.

3) **Application:** This consists of analyzing video streams of rate 1 frame/s with the object detection algorithm YOLO [19]. Frames have size 0.5 MB and the outcome of the processing (a collection of bounding boxes and tags that indicate where the objects are in a frame; see the video in [30] for an example) can be represented with at most 0.5 kB. Hence,  $\beta = 0.5\text{MB}/0.5\text{kB} = 10^{-3}$ . We consider that an IoT node is a “source” node (receives a video stream) with probability 1/2. According to our measurements, the IoT nodes and gateway (a desktop with a GPU) can process a frame in 3 s and 100 ms respectively. The analytics performance metric is given by the mean Average Precision (mAP [31]) that YOLO can detect objects correctly in a frame. This is 33.1 (YOLOv3-tiny) for an IoT node and 57.9 (YOLOv3-608) for the gateway [30].

## B. Evaluation

1) **Sensitivity analysis:** Figure 2 shows the network lifetime (left) and data analytics performance (right) for a range of values  $\theta$  and nodes’ radius. Recall that the radius affects the structure of the graph (e.g., number of links, shortest path to the gateway). The results are the average of 50 different networks generated as described in the previous section. Observe from the figure that independently of the radius, the network lifetime increases monotonically as  $\theta \rightarrow 0$ . Nonetheless, the radius plays an important role: when the network connectivity is high,<sup>9</sup> the best strategy is to offload all the processing to the gateway.<sup>10</sup> Otherwise, the solution balances between local and gateway processing. Also, observe that in this specific example selecting  $\theta = 1$  is generally not the “best” choice. When the radius is equal to 0.25, if we change from  $\theta = 1$  to  $\theta = 10^{-8}$  we can increase the network lifetime by 577% at the expense of just  $-21\%$  in the data analytics performance (see Figure 2). Similarly, when the radius is equal to 0.45 (the network is dense), by changing  $\theta$  from 1 to  $10^{-8}$  the network lifetime increases by 707% without affecting the analytics (0% change). **Conclusion:** the data analytics performance and network lifetime terms in the objective are very sensitive to parameter  $\theta$ . Also, these may not necessarily conflict: it is possible to obtain both, good analytics performance and network lifetime, by setting  $\theta$  properly.

<sup>9</sup>Nodes are on average less than two hops away from the gateway.

<sup>10</sup>That is because processing analytics at the nodes locally is more expensive than transmitting them, and the processing reward at the gateway is larger than the reward obtained as a result of processing the analytics at the nodes.

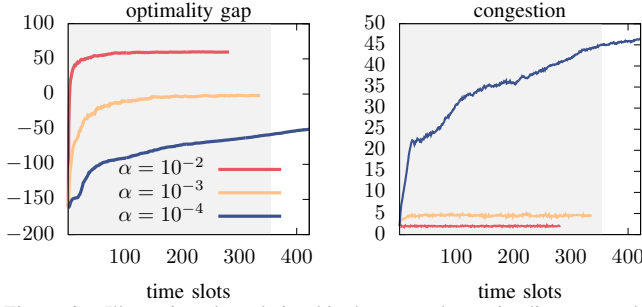


Figure 3. Illustrating the relationship between the optimality gap and the congestion in the network depending on the step size  $\alpha$ .

2) *Performance of the dynamic algorithm:* We now investigate how Algorithm 1 solves the static problem for a network with radius equal to 0.25. Parameter  $\theta$  is fixed to  $10^{-5}$ . The video frames arrive in the IoT nodes following a Poisson process. To capture the errors on the prediction function, we add noise of 20%. We run Algorithm 1 with  $\alpha \in \{10^{-2}, 10^{-3}, 10^{-4}\}$  and show the results in Figure 3. The results are the average of 50 samples. On the left, we have the normalized optimality gap<sup>11</sup> per iteration, and on the right, the system's congestion or sum of all the Lagrange multipliers  $\nu_i^{(c)}$ . The gray area indicates the lifetime of the system computed in the static problem.

Observe that the plots have different lengths (*i.e.*, different network lifetimes) and convergence behavior. Specifically, with  $\alpha = 10^{-2}$ , the transient phase is short,<sup>12</sup> the optimality gap is large, and the network lifetime shorter than in the static problem. With  $\alpha = 10^{-3}$ , the transient phase is slightly longer, but we obtain instead a much smaller optimality gap, and the network lifetime matches nearly the one of the static problem. Finally, with  $\alpha = 10^{-4}$ , the transient phase is so long that the optimality gap is always negative (*i.e.*, better than the static optimum) and the network lifetime longer than in the static problem. However, note that this is possible because congestion keeps building up in the system, *i.e.*, data analytics get accumulated in the nodes without being processed neither transmitted. **Conclusion:** parameter  $\alpha$  controls not only the optimality gap but also the duration of the transient phase. Selecting  $\alpha$  too small may result in generating congestion and not routing/processing data analytics.

3) *Comparison with other algorithms:* Given that this paper introduces a new problem, there are no other algorithms with which to compare. Hence, we use two intuitive alternatives: (i) a max-flow type algorithm [21] that makes routing/processing decisions based only on the network congestion (*i.e.*, the Lagrange multipliers); and (ii) an algorithm that balances analytics performance and energy consumption (instead of network lifetime). Specifically, the first term in the objective in problem (7) is replaced with  $\sum_{i \in \mathcal{N}} p_i(\mathbf{x}, \mathbf{y})$ .<sup>13</sup> Also, to compare the algorithms fairly, we fix the reward for processing analytics to 40 mAP and evaluate their performance in terms of network lifetime gain (we do this by selecting  $\eta$  and  $\theta$  accordingly). Figure 4 shows the average data analytics reward for 50 different realizations (time-varying arrivals and link

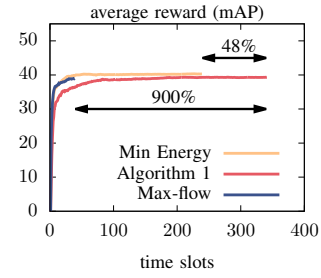


Figure 4. Illustrating how the Algorithm 1 compares to a min-energy and max-flow type algorithms.

capacities). Parameter  $\alpha$  is selected equal to  $10^{-3}$ . Observe that with the proposed algorithm, we obtain 48% longer lifetime than with the algorithm that minimizes only the total energy consumption (min-energy), and a 900% gain compared to the algorithm that only considers the network congestion (max-flow). **Conclusion:** minimizing the total energy consumption does not necessarily maximize the network lifetime. Algorithms that do not take into account the energy constraints can degrade the lifetime of the network significantly.

## VII. CONCLUSIONS

We have studied a crucial new problem in emerging IoT networks: how to allocate bandwidth and computation resources to data analytics tasks while considering the time the network can operate. The algorithm proposed can (i) balance between maximizing the network lifetime and the grade in which analytics are carried out; and (ii) operate without knowledge of the traffic arrivals or channel statistics. Our simulations seeded by real IoT device energy measurements, show that the network connectivity plays a crucial role in network lifetime maximization, that the algorithm can obtain both maximum network lifetime and maximum data analytics performance in addition to maximizing the joint objective, and that the algorithm increases the network lifetime by approximately 50% compared to an algorithm that minimizes the total energy consumption in the network.

## REFERENCES

- [1] N. J. Kaminski, I. Macaluso, E. D. Pascale, A. Nag, J. Brady, M. Y. Kelly, K. E. Nolan, W. Guibène, and L. Doyle, "A neural-network-based realization of in-network computation for the Internet of Things," in *IEEE ICC*, 2017.
- [2] O. Vermesan and J. Bacquet, "Cognitive hyperconnected digital transformation internet of things intelligence evolution," *Series in Communications*, 2017.
- [3] "Forecast and methodology," White paper, Cisco Global Cloud, 2015.
- [4] "WorldWide IoT 2016 Predictions," IDC FutureScape, 2016.
- [5] J.-H. Chang and L. Tassiulas, "Maximum lifetime routing in wireless sensor networks," *IEEE/ACM ToN*, vol. 12, no. 4, pp. 609–619, 2004.
- [6] —, "Routing for maximum system lifetime in wireless ad-hoc networks," in *Allerton Conference*, 1999.
- [7] R. Madan and S. Lall, "Distributed algorithms for maximum lifetime routing in wireless sensor networks," *IEEE Transactions on Wireless Communications*, vol. 5, no. 8, pp. 2185–2193, Aug 2006.
- [8] Y. Xue, Y. Cui, and K. Nahrstedt, "A utility-based distributed maximum lifetime routing algorithm for wireless networks," in *Quality of Service in Heterogeneous Wired/Wireless Networks*, 2005.
- [9] H. Zhang and H. Shen, "Balancing energy consumption to maximize network lifetime in data-gathering sensor networks," *IEEE Trans. on Parallel and Distributed Systems*, vol. 20, no. 10, pp. 1526–1539, 2009.
- [10] Y. Chen and Q. Zhao, "On the lifetime of wireless sensor networks," *IEEE Communications Letters*, vol. 9, no. 11, pp. 976–978, Nov 2005.

<sup>11</sup>The optimality gap divided by  $\theta$ .

<sup>12</sup>Time required to converge to a ball or steady value around the optimum.

<sup>13</sup>We assume that averages  $\lambda$  and  $\mu$  are known in the min-energy algorithm.



- [11] B. Chen, L. Tong, and P. K. Varshney, "Channel-aware distributed detection in wireless sensor networks," *IEEE Signal Processing Magazine*, vol. 23, no. 4, pp. 16–26, July 2006.
- [12] J.-H. Chang and L. Tassiulas, "Fast approximate algorithms for maximum lifetime routing in wireless ad-hoc networks," in *International Conference on Research in Networking*, 2000, pp. 702–713.
- [13] R. M. Curry and J. C. Smith, "A survey of optimization algorithms for wireless sensor network lifetime maximization," *Computers & Industrial Engineering*, vol. 101, pp. 145–166, 2016.
- [14] J. Zhu, S. Chen, B. Bensou, and K. L. Hung, "Tradeoff between lifetime and rate allocation in wireless sensor networks: A cross layer approach," in *IEEE INFOCOM*, 2007.
- [15] A. Galanopoulos, G. Iosifidis, and T. Salonidis, "Optimizing data analytics in energy constrained iot networks," in *WiOpt*, 2018.
- [16] H. Feng, J. Llorca, A. M. Tulino, and A. F. Molisch, "Dynamic network service optimization in distributed cloud networks," in *ICC*, 2016.
- [17] H. Feng, J. Llorca, A. M. Tulino, D. Raz, and A. F. Molisch, "Approximation algorithms for the nfv service distribution problem," in *IEEE INFOCOM*, 2017.
- [18] P. Gupta and P. R. Kumar, "The capacity of wireless networks," *IEEE Transactions on Information Theory*, vol. 46, no. 2, pp. 388–404, 2000.
- [19] J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," *arXiv:1804.02767*, 2018.
- [20] M. Kodialam and T. Nandagopal, "Characterizing the capacity region in multi-radio multi-channel wireless mesh networks," *Mobicom*, 2005.
- [21] G. B. Dantzig, *Linear Programming and Extensions*. Princeton University Press, 1963.
- [22] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge University Press, 2004.
- [23] D. P. Bertsekas, A. Nedić, and A. E. Ozdaglar, *Convex analysis and optimization*. Athena Scientific, 2003.
- [24] B. O'Donoghue, E. Chu, N. Parikh, and S. Boyd, "SCS: Splitting conic solver, version 2.0.2," <https://github.com/cvxgrp/scs>, 2017.
- [25] L. Georgiadis, M. J. Neely, and L. Tassiulas, *Resource Allocation and Cross-Layer Control in Wireless Networks*. Foundations and Trends in Optimization, 2006.
- [26] D. P. Bertsekas, *Nonlinear Programming: Second Edition*. Athena Scientific, 1999.
- [27] A. Eryilmaz and R. Srikant, "Fair resource allocation in wireless networks using queue-length-based scheduling and congestion control," *IEEE/ACM ToN*, vol. 15, no. 6, pp. 1333–1344, 2007.
- [28] M. J. Neely, "Stability and Capacity Regions or Discrete Time Queueing Networks," *ArXiv e-prints*, Mar. 2010.
- [29] M. Haenggi and et al., "Stochastic geometry and random graphs for the analysis and design of wireless networks," *IEEE JSAC*, vol. 27, no. 7, pp. 1029–1046, September 2009.
- [30] [Online]. Available: <https://pjreddie.com/darknet/yolo/>
- [31] [Online]. Available: "http://cocodataset.org/#detection-eval"

## VIII. APPENDIX

### A. Preliminaries

To streamline exposition and due to space constraints, we assume in the analysis there is only one commodity. The extension to multiple commodities is nonetheless straightforward. Now, we write the Lagrangian as  $L(x, y, \nu) = f(x, y) + \nu^\top (g(x, y) + \lambda)$ , where  $f(x, y)$  is the objective function in (8) and  $g(x, y) = (g_1(x, y), \dots, g_n(x, y))$  with  $g_i(x, y) = (\beta_i - 1)y_i + \sum_{j \in \mathcal{N}_i} (x_{ji} - x_{ij})$ ,  $i \in \{1, \dots, n\}$ . That is,  $g_i(x, y) + \lambda_i$  is the flow conservation constraint of node  $i \in \mathcal{N}$  and  $g(x, y)$  a collection of  $n$  flow conservation constraints. Note we do not use bold notation and that  $\nu^\top (g(x, y) + \lambda)$  is the inner product of vectors  $\nu$  and  $(g(x, y) + \lambda)$ .

### B. Proof of Theorem 1

We have the dual (sub)gradient update  $\nu[t+1] = \nu[t] + \alpha(g(x[t], y[t]) + \lambda)$  where  $x[t] \in X[t]$  and  $y[t] \in Y[t]$  are obtained by minimizing the Lagrangian for a fixed  $\nu[t]$ . We start by showing an upper bound on the objective function.

**Lemma 2** (Objective upper bound). *It holds*

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^T (f(x[t], y[t]) - f(x^*[t], y^*[t])) \\ & \leq -\frac{1}{T} \sum_{t=1}^T (\nu[t]^\top (g(x[t], y[t]) + \lambda[t])) \\ & \quad + \frac{1}{T} \sum_{t=1}^T (\nu[t]^\top (g(x^*[t], y^*[t]) + \lambda[t])) + \epsilon_2 \end{aligned} \quad (16)$$

*Proof:* Since  $x[t]$  and  $y[t]$  are selected from  $X[t]$  and  $Y[t]$  to minimize the Lagrangian, we can write  $\tilde{f}(x[t], y[t]) + (\nu[t]^\top (g(x[t], y[t]) + \lambda[t])) = f(x[t], y[t]) + \xi(x[t], y[t]) + (\nu[t]^\top (g(x[t], y[t]) + \lambda[t])) \leq f(x^*[t], y^*[t]) + \xi(x^*[t], y^*[t]) + (\nu[t]^\top (g(x^*[t], y^*[t]) + \lambda[t]))$ . Rearranging terms, summing from  $t = 1, \dots, T$ , dividing by  $T$ , and using the fact that  $|\xi(x[t], y[t])| \leq \epsilon_2$  by assumption for all  $(x[t], y[t]) \in (X[t], Y[t])$  yields the stated result. ■

Next, we bound the first two terms in the RHS of (16).

**Lemma 3.** *The first term in the RHS of (16) is upper bounded by  $\frac{\alpha}{2T} \sum_{t=1}^T \sigma[t]$  where  $\sigma[t] := \|g(x[t], y[t]) + \lambda[t]\|_2^2$ .*

*Proof:* Let  $\theta \in \mathbf{R}^n$ . From the standard dual subgradient update, we have that  $\|\nu[t+1] - \theta\|_2^2 = \|\nu[t] + \alpha(g(x[t], y[t]) + \lambda[t]) - \theta\|_2^2 = \|\nu[t] - \theta\|_2^2 + \alpha^2 \|g(x[t], y[t]) + \lambda[t]\|_2^2 + 2\alpha(\nu[t] - \theta)^\top (g(x[t], y[t]) + \lambda[t]) \leq \|\nu[t] - \theta\|_2^2 + \alpha^2 \sigma[t] + 2\alpha(\nu[t] - \theta)^\top (g(x[t], y[t]) + \lambda[t])$ . Parameter  $\sigma[t]$  is bounded since  $X[t], Y[t], \lambda[t]$  are bounded for all  $t \in \mathbf{Z}_+$ . Next, rearrange terms, apply the expansion recursively for  $t = 1, \dots, T$  to obtain  $-2\alpha \sum_{t=1}^T (\nu[t] - \theta)^\top (g(x[t], y[t]) + \lambda[t]) \leq \alpha^2 \sum_{t=1}^T \sigma[t] + \|\nu[1] - \theta\|_2^2 - \|\nu[T] - \theta\|_2^2$ . Drop the third term in the RHS of the last equation since it is nonnegative, let  $\theta = 0$  and  $\nu[1] = 0$ , and divide across by  $2\alpha T$  to obtain the stated bound. ■

**Lemma 4.** *The second term in the RHS of (16) is equal to zero on expectation.*

*Proof:* Take expectations with respect to random variable  $X[t], Y[t], \lambda[t]$ , and write  $\mathbf{E}(\frac{1}{T} \sum_{t=1}^T (\nu[t]^\top (g(x^*[t], y^*[t]) + \lambda[t]))) \stackrel{(a)}{=} \frac{1}{T} \sum_{t=1}^T (\nu[t]^\top \mathbf{E}(g(x^*[t], y^*[t]) + \lambda[t])) \stackrel{(b)}{=} 0$  where (a) follows from the linearity of the expectation, and (b) since  $\mathbf{E}(g(x^*[t], y^*[t]) + \lambda[t]) = \mathbf{E}(g(x^*[t], y^*[t])) + \mathbf{E}(\lambda[t]) = g(x^*, y^*) + \lambda = 0$  (i.e., at the optimum the problem must be feasible). Note that this is always case since random variables  $X[t], Y[t], \lambda[t]$  do not depend on  $\nu[t]$  for all  $t \in \mathbf{Z}_+$ . ■

We are now in position to prove claim (i). Take expectations in the bound obtain in Lemma 2. The first term can be upper bounded by Lemma 3 and the second with Lemma 4. By letting  $\mathbf{E}(\sigma[t]) \leq 2\epsilon_1$  for all  $t \in \mathbf{Z}_+$  the stated result follows.

**Lemma 5** (Feasibility).  *$\mathbf{E}(g(\bar{x}, \bar{y}) + \lambda) = 0$  converges to a feasible point asymptotically as  $T \rightarrow \infty$ .*

*Sketch:* Recall  $\nu[t+1] = \nu[t] + \alpha(g(x[t], y[t]) + \lambda[t])$ . Rearrange terms and apply the iteration recursively to obtain  $\nu[t+1] - \nu[1] = \sum_{t=1}^T \alpha(g(x[t], y[t]) + \lambda[t])$ . Dividing by  $T$  and using the fact that  $\nu[1] = 0$  we have  $\frac{1}{T} \sum_{t=1}^T g(x[t], y[t]) + \lambda[t] = g(\bar{x}, \bar{y}) + \bar{\lambda} = \nu[t]/(\alpha T)$  which implies that  $(g(\bar{x}, \bar{y}) + \bar{\lambda}) \rightarrow 0$  as  $T \rightarrow \infty$  if  $\nu[t]$  is bounded. The latter will be the case, on expectation, when  $\nu[t]$  is a nonnegative process (nodes do not pre-serve/process data) and the flow conservation constraints are not tight: there exists a  $\chi \succ 0$  and  $(\hat{x}, \hat{y})$  such that  $g(\hat{x}, \hat{y}) + \lambda + \chi = 0$  holds. ■