

On the Delivery of Augmented Information Services over Wireless Computing Networks

Hao Feng^{*†}, Jaime Llorca[†], Antonia M. Tulino^{†‡}, Andreas F. Molisch^{*}

^{*}University of Southern California, Email: {haofeng, molisch}@usc.edu

[†]Nokia Bell Labs, Email: {jaime.llorca, a.tulino}@nokia-bell-labs.com

[‡]University of Naples Federico II, Italy. Email: {antoniamaria.tulino}@unina.it

Abstract—In an *augmented information (AgI) service*, users consume information that results from the execution of a chain of service functions that process source information to create real-time augmented value. Applications may include real-time analysis of remote sensing data, real-time computer vision, personalized video streaming, and augmented reality, among others. We consider the problem of optimal distribution of AgI services over a wireless computing network, in which nodes are equipped with both communication and computing resources. We characterize the wireless computing network capacity region and design a joint flow scheduling and resource allocation algorithm that stabilizes the underlying queuing system while achieving arbitrarily close to minimum network cost, with a tradeoff in network delay. Our solution captures the unique chaining and flow scaling aspects of AgI services, while exploiting the use of the broadcast approach over the wireless channel.

I. INTRODUCTION

Internet traffic will soon be dominated by the consumption of what we refer to as *augmented information (AgI) services*. Unlike traditional information services, in which users consume information that is produced or stored at a given source and it is delivered via a communications network, AgI services provide end users with information that results from the *processing* of source information via possibly multiple service functions that can be hosted anywhere in the network. Examples include real-time analysis of remote sensing data, real-time computer vision, personalized video streaming, and augmented reality services, among others.

While today's AgI services are mostly implemented in the form of virtual functions instantiated over general purpose servers at centralized cloud data centers, the increasingly low latency requirements of next generation real-time AgI services is driving cloud resources closer to the end users in the form of small cloud nodes at the edge of the network, resulting in what is referred to as a distributed cloud network [1]. This naturally raises the question of where to execute each service function, a question that is impacted both by the computation and the communication resources of the cloud network infrastructure. Ref. [1] addressed the cloud service distribution problem as a global optimization problem, where the goal is to find the placement of service functions and the routing of network flows that meets a given set of demands with minimum total cloud network cost. The capacity of *wireline* cloud networks was recently addressed by the present authors in [2] and [3]. These works provided the first characterization of a cloud

network capacity region in terms of the closure of input rates that can be stabilized by any control algorithm and designed a throughput-optimal dynamic control policy that achieves arbitrarily close to minimum average network cost.

A key missing aspect in all these works is the wireless network. AgI services are increasingly sourced and accessed from wireless devices, and with the advent of mobile computing, service functions can also be hosted at wireless computing nodes (i.e., computing devices with wireless networking capabilities). When introducing the wireless network into the computing infrastructure, the often unpredictable nature of the wireless channel further complicates flow scheduling, routing, and resource allocation. In the context of traditional wireless multi-hop networks, the Lyapunov drift plus penalty (LPP) control methodology (see [4] and references therein) has shown to be a promising approach to tackle these intricate stochastic network optimization problems. Ref. [5] extends the LPP approach to multi-hop, multi-commodity wireless ad-hoc networks, leading to the Diversity Backpressure (DIVBAR) algorithm. DIVBAR exploits the broadcast nature of the wireless channel without the need of instantaneous channel state information (CSI), and it is shown to be throughput-optimal under the assumption that at most one packet can be transmitted in each transmission attempt and that no advanced coding scheme is used. Ref. [6] further incorporates rateless coding in the transmission of a single packet.

Motivated by the important role of wireless networks in the delivery of AgI services, in this paper, we address the problem of optimal distribution of AgI services over a multi-hop wireless computing network. The network is composed of nodes with communication and computing capabilities. We extend the *multi-commodity-chain* flow model in [2], [3] for the delivery of AgI services over wireless multi-hop computing networks, enabling the characterization of the unique flow chaining and scaling aspects of AgI services. We adopt the *broadcast approach* [7], [8], where information is encoded into superposition layers according to the channel conditions, in order to exploit routing diversity with enhanced transmission efficiency. We characterize the capacity region of a wireless computing network and design a fully distributed flow scheduling and resource allocation algorithm that adaptively stabilizes the underlying queuing system while achieving arbitrarily close to minimum network cost, with a tradeoff in network delay.

The remainder of the paper is organized as follows: Section II presents the system model. Section III characterizes the network capacity region of a wireless computing network. Section IV describes the proposed dynamic wireless computing network control (DWCNC) algorithm. Section V provides the performance analysis of the proposed algorithm. The paper is concluded in Section VI.

II. SYSTEM MODEL

A. Network model

We consider a wireless computing network composed of $N = |\mathcal{N}|$ distributed computing nodes that communicate over wireless links labeled according to node pairs (i, j) for $i, j \in \mathcal{N}$. Node $i \in \mathcal{N}$ is equipped with K_i^{tr} transmission resource units (e.g., transmit power) that it can use to transmit information over the wireless channel. In addition, node i is equipped with K_i^{pr} processing resource units (e.g., central processing units or CPUs) that it can use to process information as part of an AgI service (see Sec. II-B).

Time is slotted with slots normalized to integer units $t \in \{0, 1, 2, \dots\}$. We use the binary variable $y_{i,k}^{tr}(t) \in \{0, 1\}$ to indicate the allocation or activation of $k \in \{0, \dots, K_i^{tr}\}$ transmission resource units at node i at time t , which incurs $w_{i,k}^{tr}$ cost units. Analogously, $y_{i,k}^{pr}(t) \in \{0, 1\}$ indicates the allocation of $k \in \{0, \dots, K_i^{pr}\}$ processing resource units at node i at time t , which incurs $w_{i,k}^{pr}$ cost units. Notice that binary resource allocation variables $y_{i,k}^{tr}(t), y_{i,k}^{pr}(t)$ must satisfy $\sum_{k \in \mathcal{K}_i^{tr}} y_{i,k}^{tr}(t) \leq 1, \sum_{k \in \mathcal{K}_i^{pr}} y_{i,k}^{pr}(t) \leq 1$.

B. Augmented information service model

We consider the distribution of an *augmented information service* described by a chain of functions $\mathcal{M} = \{1, 2, \dots, M\}$ for a subset of destination nodes $\mathcal{D} \subseteq \mathcal{N}$. We adopt a multi-commodity-chain flow model [1], in which commodity $(d, m) \in \mathcal{D} \times \{\mathcal{M}, 0\}$ identifies the flow of information units (of arbitrarily fine granularity, e.g., bits) output of function $m \in \mathcal{M}$ for destination $d \in \mathcal{D}$. Commodity $(d, 0)$ denotes the source commodity for destination d , which identifies the flow of information units that arrive exogenously at the origin nodes $\mathcal{O} \subseteq \mathcal{N}$ for destination d (see Fig. 1).

Each service function has (possibly) different processing requirements. We use $r^{(m)}$ to denote the number of operations per information unit performed by function m . Another key aspect of AgI services is the fact that information flows can change size as they go through service functions. Let $\xi^{(m)} > 0$ denote the *scaling factor* of function m . Then, the size of the function's output flow is $\xi^{(m)}$ times as large as its input flow. Moreover, a processing delay $D_i^{(m)}$ (in timeslots) is incurred in executing function m at node i , as long as the processing flow satisfies the node's processing rate constraint.

C. Wireless transmission model

Due to the broadcast nature of the wireless medium, multiple receivers (RXs) may overhear the transmission of a given transmitter (TX). Multiple TXs may transmit simultaneously to overlapping RXs, due to the use of orthogonal broadcast

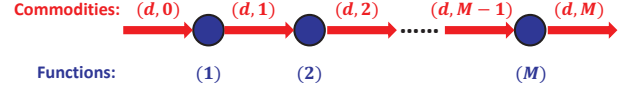


Fig. 1. Illustration of the AgI service chain for destination $d \in \mathcal{D}$. There are M functions and $M + 1$ commodities. The AgI service takes source commodity $(d, 0)$ and delivers final commodity (d, M) after going through the sequence of functions $\{1, 2, \dots, M\}$. Function m takes commodity $(d, m - 1)$ and generates commodity (d, m) .

channels of fixed bandwidth, a priori allocated by a given policy, whose design is out of the scope of this paper. We model the channel between node i and all other nodes in the network as a physically degraded Gaussian broadcast channel, where the *network state process* (the vector of all channel gains), denoted by $\mathbf{S}(t) \equiv \{s_{ij}(t), \forall i, j \in \mathcal{N}\}$, evolves according to a Markov process with state set \mathcal{S} and whose steady-state probability exists. We assume that the statistical channel state information (CSI) is known at the TX, while the instantaneous CSI can only be learned after the transmission has taken place and is thereby outdated (delayed).

It is well-known that superposition coding is optimal (capacity achieving) for the physically degraded broadcast channel with independent messages [9]. In particular, in this work we adopt the *broadcast approach* (see [7], [8] and references therein), which consists of sending incremental information using superposition layers, such that the number of decoded layers at any RX depends on its own channel state, and the information decoded by a given RX is a subset of the information decoded by any other RX with no worse channel gain. That is, for a given transmitting node i , if we sort the $N - 1$ potential receiving nodes in non-decreasing order of their channel gains $\{g_{i,1}, \dots, g_{i,N-1}\}$, such that $g_{i,n}$ with $n \in \{1, \dots, N - 1\}$ denotes the receiver with the n -th lowest channel gain, then the information decoded by receiver $g_{i,n}$ is also decoded by receiver $g_{i,l}$, for $l > n$. Moreover, let $\Omega_{i,n} \triangleq \{g_{i,n}, \dots, g_{i,N-1}\}$ be the set of receivers with the $N - n$ highest channel gains. Then, we can partition the information transmitted by node i during a given timeslot into $N - 1$ disjoint groups, with the n -th partition being the information decoded by every node in $\Omega_{i,n}$ and not by any node in $\Omega_{i,l}$ with $l < n$.

Let $\tilde{\gamma}_{i,k}(a)$ denote the optimal power density function of k transmission resource units at node i . Then, based on the *broadcast approach* [8], the maximum achievable rate on link (i, j) at time t is given by

$$R_{ij,k}(t) = \int_0^{s_{ij}(t)} \frac{a \tilde{\gamma}_{i,k}(a)}{1 + a \int_a^\infty \tilde{\gamma}_{i,k}(s) ds} da. \quad (1)$$

D. Communication protocol

The communication protocol between each TX-RX pair within one timeslot is illustrated in Fig. 2. At the beginning of each timeslot, TX and RX exchange all necessary control signals, including queue backlog state information (see Subsection II-F). Then, the TX decides how many transmission resource units to allocate for the given timeslot and how

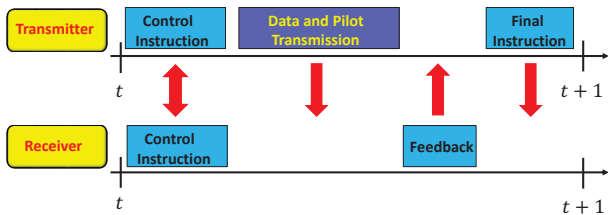


Fig. 2. Timing diagram of the communication protocol over a wireless link for one timeslot

to allocate the given bandwidth among the commodities for transmission.¹ Afterwards, the transmission starts and lasts for a fixed time period (within the timeslot); during that time both data and pilot tones (whose overhead is neglected) are transmitted. After the transmission ends, every potential RX provides immediate feedback, which may include decoded information and/or experienced CSI.² The TX then makes a *forwarding decision* that determines the information units for which each RX gets forwarding responsibility and sends the decision as a final instruction to all the RXs. Only the RX that gets the forwarding responsibility keeps the corresponding information units, while others discard their copies.³

We use $\mu_{ij}^{(d,m)}(t)$ to denote the information units of commodity (d, m) kept by node j after the transmission from node i during timeslot t . The total amount of information kept by node j from the transmission from node i can be no larger than the maximum achievable rate on link (i, j) :

$$\sum_{(d,m)} \mu_{ij}^{(d,m)}(t) \leq \sum_{k=0}^{K_i^{tr}} R_{ij,k}(t) y_{i,k}^{tr}(t), \quad \forall i, j, t, \quad (2)$$

with $R_{ij,k}(t)$ given by (1).

In addition, it shall be useful to denote by $\mu_{ig_{i,l}}^{(d,m)}(t)$ the information kept by node $g_{i,l}$ from the l -th partition of node i 's transmitted information at time t . We then have

$$\mu_{ig_{i,l}}^{(d,m)}(t) = \sum_{n=1}^l \mu_{ig_{i,l,n}}^{(d,m)}(t), \quad \forall i, l, d, m, t, \quad (3)$$

$$\sum_{(d,m)} \sum_{l=n}^{N-1} \mu_{ig_{i,l,n}}^{(d,m)}(t) \leq \sum_{k=0}^{K_i^{tr}} [R_{ig_{i,n},k}(t) - R_{ig_{i,n-1},k}(t)] y_{i,k}^{tr}(t), \quad \forall i, l, n, t, \quad (4)$$

where $R_{ig_{i,0},k}(t) = 0$, $\forall i, k, t$. Note that (4) is consistent with the fact that only one copy of the transmitted information units is kept for forwarding.

E. Computing model

We assume a static dedicated *computing channel* model, where the achievable processing rate at node i with the allocation of k processing resource units is given by $R_{i,k}$. We use $\mu_{i,pr}^{(d,m)}(t)$ to denote the assigned flow rate from node

¹Note that we require by definition uniform power spectral density.

²Note that with the broadcast approach, the TX can identify the information units decoded by each RX directly from the CSI feedback.

³The control information, feedbacks and final instruction are sent though a stable control channel.

i to its processing unit for commodity (d, m) at time t , and $\mu_{pr,i}^{(d,m)}(t)$ for the flow rate from the processing unit back to node i . Then,

$$\mu_{pr,i}^{(d,m)}(t) = \xi^{(m)} \mu_{i,pr}^{(d,m-1)}(t - D_i^{(m)}), \quad \forall i, d, m > 0, t, \quad (5)$$

$$\sum_{(d,m>0)} \mu_{i,pr}^{(d,m-1)}(t) r^{(m)} \leq \sum_{k=0}^{K_i^{pr}} R_{i,k} y_{i,k}^{pr}(t), \quad \forall i, t, \quad (6)$$

where (5) and (6) are multi-commodity-chain and maximum rate constraints, respectively [1].

F. Queuing model

We denote by $a_i^{(d,m)}(t)$ the exogenous arrival rate of commodity $(d, m) \in \mathcal{D} \times \{\mathcal{M}, 0\}$ at node i during timeslot t , and by $\lambda_i^{(d,m)}$ its expected value. We assume that $a_i^{(d,m)}(t)$ is independently and identically distributed (i.i.d.) across timeslots and that its second moment is upper bounded by constant A_{\max}^2 . Recall that, in an AgI service, only the source commodity $(d, 0)$ enters the network exogenously, while all other commodities are created inside the network as the output of a service function. Hence, $a_i^{(d,m)}(t) = 0$, for all t when $m > 0$ or $i \notin \mathcal{O}$.

During network evolution, internal network queues buffer the data according to their commodities. We define the *queue backlog* of commodity (d, m) at node i , $Q_i^{(d,m)}(t)$, as the amount of commodity (d, m) in the queue of node i at the beginning of timeslot t , which evolves over time as follows:

$$Q_i^{(d,m)}(t+1) \leq \left[Q_i^{(d,m)}(t) - \sum_{j:j \neq i} \mu_{ij}^{(d,m)}(t) - \mu_{i,pr}^{(d,m)}(t) \right]^+ + \sum_{j:j \neq i} \mu_{ji}^{(d,m)}(t) + \mu_{pr,i}^{(d,m)}(t) + a_i^{(d,m)}(t). \quad (7)$$

Note that, in an AgI service, only the final commodity (d, M) is allowed to exit the network once it arrives to its destination $d \in \mathcal{D}$, while any other commodity (d, m) , $m < M$, can only get consumed by being processed into the next commodity $(d, m+1)$ in the service chain. Final commodity (d, M) is assumed to leave the network immediately upon arrival/decoding, i.e., $Q_d^{(d,M)}(t) = 0$, for all d, t .

G. Network Objective

The goal is to develop a control algorithm that dynamically distributes the service and schedules the flow in the network to minimize the average resource cost

$$\limsup_{t \rightarrow \infty} \frac{1}{t} \sum_{\tau=0}^{t-1} \mathbb{E} \{h(\tau)\}, \quad (8)$$

where $h(\tau)$ is the total cost of the network at time τ :

$$h(\tau) \triangleq \sum_{i \in \mathcal{N}} \left[\sum_{k=0}^{K_i^{pr}} w_{i,k}^{pr} y_{i,k}^{pr}(\tau) + \sum_{k=0}^{K_i^{tr}} w_{i,k}^{tr} y_{i,k}^{tr}(\tau) \right], \quad (9)$$

while ensuring the network is *rate stable* [4], i.e.,

$$\lim_{t \rightarrow \infty} \frac{1}{t} Q_i^{(d,m)}(t) = 0 \quad \text{with prob. 1, } \forall i, d, m. \quad (10)$$

III. WIRELESS COMPUTING NETWORK CAPACITY REGION

The wireless computing network capacity region Λ is defined as the closure of all input rate matrices $\{\lambda_i^{(d,m)}\}$ that can be stabilized by the network according to some control algorithm conforming to the network and service structure specified in Section II.

Let $\pi_{\mathbf{s}}$ denote the steady state probability distribution of the network state process $\mathbf{S}(t)$.

Theorem 1. *The wireless computing network capacity region Λ consists of all average exogenous input rates $\{\lambda_i^{(d,m)}\}$ for which there exist multi-commodity flow variables $f_{ij}^{(d,m)}$, $f_{pr,i}^{(d,m)}$, $f_{i,pr}^{(d,m)}$, together with probability values $\alpha_{i,k}^{pr}$, $\alpha_{i,k}^{tr}(\mathbf{s})$, $\beta_{i,pr}^{(d,m)}(k)$, $\beta_{i,tr}^{(d,m)}(\mathbf{s}, k)$, $\eta_{ij}^{(d,m)}(\mathbf{s}, k, n)$, for all $i, j \neq i, k, d, m$, and all network states $\mathbf{s} \in \mathcal{S}$, such that:*

$$\sum_j f_{ji}^{(d,m)} + f_{pr,i}^{(d,m)} + \lambda_i^{(d,m)} \leq \sum_j f_{ij}^{(d,m)} + f_{i,pr}^{(d,m)}, \quad (11)$$

$$f_{pr,i}^{(d,m)} = \xi_i^{(m)} f_{i,pr}^{(d,m-1)}, \quad m > 0 \quad (12)$$

$$f_{i,pr}^{(d,m)} \leq \frac{1}{r^{(m+1)}} \sum_{k=0}^{K_i^{pr}} \alpha_{i,k}^{pr} \beta_{i,pr}^{(d,m)}(k) R_{i,k}, \quad (13)$$

$$f_{ij}^{(d,m)} \leq \sum_{\mathbf{s} \in \mathcal{S}} \pi_{\mathbf{s}} \sum_{k=0}^{K_i^{tr}} \alpha_{i,k}^{tr}(\mathbf{s}) \beta_{i,tr}^{(d,m)}(\mathbf{s}, k) \times \sum_{n=1}^{g_{i,s}^{-1}(j)} [R_{ig_{i,n},k}(\mathbf{s}) - R_{ig_{i,n-1},k}(\mathbf{s})] \eta_{ij}^{(d,m)}(\mathbf{s}, k, n), \quad (14)$$

$$f_{i,pr}^{(d,M)} = 0, \quad f_{pr,i}^{(d,0)} = 0, \quad f_{dj}^{(d,M)} = 0, \quad (15)$$

$$f_{i,pr}^{(d,m)} \geq 0, \quad f_{ij}^{(d,m)} \geq 0, \quad (16)$$

$$\sum_{k=0}^{K_i^{pr}} \alpha_{i,k}^{pr} \leq 1, \quad \sum_{k=0}^{K_i^{tr}} \alpha_{i,k}^{tr}(\mathbf{s}) \leq 1, \quad (17)$$

$$\sum_{(d,m)} \beta_{i,pr}^{(d,m)}(k) \leq 1, \quad \sum_{(d,m)} \beta_{i,tr}^{(d,m)}(\mathbf{s}, k) \leq 1, \quad (18)$$

$$\sum_j \eta_{ij}^{(d,m)}(\mathbf{s}, k, n) \leq 1, \quad (19)$$

where \mathbf{s} denotes the network state, whose (i, j) -th element $(\mathbf{s})_{ij}$ indicates the channel state of link (i, j) and $g_{i,s}^{-1}(j)$ is the index of node j in the sequence $\{g_{i,1}, \dots, g_{i,N-1}\}$, given the network state \mathbf{s} . Finally, with an abuse of notation, in (14), $R_{ij,k}(\mathbf{s})$ denotes the rate given by (1) with $s_{ij}(t)$ equal to the (i, j) -th element of \mathbf{s} .

Furthermore, the minimum average network cost required for network stability is given by

$$\bar{h}^* = \min_{\underline{h}} \quad (20)$$

where

$$\underline{h} = \sum_{i \in \mathcal{N}} \left(\sum_{k=0}^{K_i^{pr}} \alpha_{i,k}^{pr} w_{i,k}^{pr} + \sum_{k=0}^{K_i^{tr}} w_{i,k}^{tr} \sum_{\mathbf{s} \in \mathcal{S}} \pi_{\mathbf{s}} \alpha_{i,k}^{tr}(\mathbf{s}) \right), \quad (21)$$

and the minimization is over all $f_{i,pr}^{(d,m)}$, $f_{ij}^{(d,m)}$, $\alpha_{i,k}^{pr}$, $\alpha_{i,k}^{tr}(\mathbf{s})$, $\beta_{i,pr}^{(d,m)}(k)$, $\beta_{i,tr}^{(d,m)}(\mathbf{s}, k)$, and $\eta_{ij}^{(d,m)}(\mathbf{s}, k, n)$ satisfying Eqs. (11)-(19). \square

Proof: The proof of Theorem 1 is omitted here due to space limitations and can be found in [10]. \blacksquare

In the above theorem, (11) are flow conservation constraints, (13) and (14) are rate constraints, and (15) and (16) are non-negativity and flow efficiency constraints. The probability values $\alpha_{i,k}^{pr}$, $\alpha_{i,k}^{tr}(\mathbf{s})$, $\beta_{i,pr}^{(d,m)}(k)$, $\beta_{i,tr}^{(d,m)}(\mathbf{s}, k)$ and $\eta_{ij}^{(d,m)}(\mathbf{s}, k, n)$ define an optimal *stationary randomized policy*, where:

- $\alpha_{i,k}^{pr}$: the probability that k processing resource units are allocated at node i ;
- $\alpha_{i,k}^{tr}(\mathbf{s})$: the probability that k transmission resource units are allocated at node i , conditioned on the network state \mathbf{s} ;
- $\beta_{i,pr}^{(d,m)}(k)$: the probability that node i processes commodity (d, m) , conditioned on the allocation of k processing resource units;
- $\beta_{i,tr}^{(d,m)}(\mathbf{s}, k)$: the probability that node i transmits commodity (d, m) , conditioned on the network state \mathbf{s} and on the allocation of k transmission resource units;
- $\eta_{ij}^{(d,m)}(\mathbf{s}, k, n)$: the probability that node i forwards the information of the n -th partition to node j , conditioned on the network state \mathbf{s} , and on the allocation of k transmission resource units.

It is important to note that this optimal stationary randomized policy is hard to obtain in practice, as it requires the knowledge of $\{\lambda_i^{(d,m)}\}$ and solving a complex nonlinear program. However, its existence is essential for proving the performance of our proposed dynamic control algorithm.

IV. THE DYNAMIC WIRELESS COMPUTING NETWORK CONTROL ALGORITHM

Defining a non-negative control parameter V representing the degree to which we emphasize resource cost minimization, we propose a dynamic wireless computing network control strategy that accounts for both transmission and processing-related flow and resource allocation decisions in a fully distributed manner.

Dynamic Wireless Computing Network Control (DWCNC):

Local wireless transmission decisions: At timeslot t , each node i observes its local queue backlogs, the queue backlogs of its potential RXs and the associated statistical CSI, and performs the following operations:

- 1) For each commodity (d, m) and each receiving node j , compute the *differential backlog weight*

$$W_{ij}^{(d,m)}(t) \triangleq [Q_i^{(d,m)}(t) - Q_j^{(d,m)}(t)]^+.$$

- 2) For each transmission resource allocation choice $k \in \{0, \dots, K_i^{tr}\}$, compute the *transmission utility weight* for each commodity (d, m) :

$$W_{i,k,tr}^{(d,m)}(t) \triangleq \sum_{\mathbf{s} \in \mathcal{S}} \Pr(\mathbf{S}(t) = \mathbf{s} | \mathbf{S}(t-1) = \tilde{\mathbf{s}}) \times \sum_{n=1}^{N-1} [R_{ig_{i,n},k}(\mathbf{s}) - R_{ig_{i,n-1},k}(\mathbf{s})] \max_{j \in \Omega_{i,n}(\mathbf{s})} \{W_{ij}^{(d,m)}(t)\},$$

where $\tilde{\mathbf{s}}$ denotes the CSI feedbacks at time $t-1$, and, with an abuse of notation, $\Omega_{i,n}(\mathbf{s})$ is used to indicate the dependence of $\Omega_{i,n}$ on the the network state.

- 3) Compute the optimal number of resource units k_{tr}^\dagger to allocate and the optimal commodity $(d, m)_{tr}^\dagger$ to transmit:

$$\left[k_{tr}^\dagger, (d, m)_{tr}^\dagger \right] = \arg \max_{k, (d, m)} \left\{ W_{i, k, tr}^{(d, m)}(t) - V w_{i, k}^{tr} \right\}. \quad (22)$$

Denote by $W_i^{tr\dagger}(t)$ the corresponding metric value.

- 4) If $W_i^{tr\dagger}(t) > 0$, node i transmits commodity $(d, m)_{tr}^\dagger$ by allocating k_{tr}^\dagger transmission resource units; otherwise, node i keeps silent ($k_{tr}^\dagger = 0$) at time t .
- 5) After receiving the CSI feedbacks, node i derives the information decoded by each RX and assigns the forwarding responsibility for the n -th partition of the transmitted information to the RX in $\Omega_{i, n}(\mathbf{S}(t))$ with the largest positive $W_{ij}^{(d, m)}(t)$, while all other RXs in $\Omega_{i, n}(\mathbf{S}(t))$ and node i discard the information. If no receiver in $\Omega_{i, n}(\mathbf{S}(t))$ has positive $W_{ij}^{(d, m)}(t)$, node i retains the information of partition n , while all the receivers in $\Omega_{i, n}(\mathbf{S}(t))$ discard it. Note that, with an abuse of notation, $\Omega_{i, n}(\mathbf{S}(t))$ is here used to indicate the dependence of $\Omega_{i, n}$ on the realization of the network state at time t .

Local processing decisions: At timeslot t , each node i observes its local queue backlogs and performs the following operations:

- 1) For each commodity (d, m) , compute the *processing utility weights*

$$W_i^{(d, m)}(t) \triangleq \frac{1}{r^{(m+1)}} \left[Q_i^{(d, m)}(t) - \xi^{(m+1)} Q_i^{(d, m+1)}(t) \right]^+,$$

Specifically, $W_i^{(d, m)}(t)$ indicates the benefit of executing function $(m+1)$ to process commodity (d, m) into commodity $(d, m+1)$ at time t , in terms of the local backlog reduction per processing unit cost.

- 2) Compute the optimal number of resource units k_{pr}^\dagger to allocate and the optimal commodity $(d, m)_{pr}^\dagger$ to process:

$$\left[k_{pr}^\dagger, (d, m)_{pr}^\dagger \right] = \arg \max_{k, (d, m)} \left\{ R_{i, k} W_i^{(d, m)}(t) - V w_{i, k}^{pr} \right\}. \quad (23)$$

- 3) Make the following flow rate assignment decisions:

$$\begin{aligned} \mu_{i, pr}^{(d, m)\dagger}(t) &= R_{i, k_{pr}^\dagger} / r^{(m_{pr}^\dagger + 1)}; \\ \mu_{i, pr}^{(d, m)}(t) &= 0, \quad \forall (d, m) \neq (d, m)_{pr}^\dagger. \end{aligned}$$

Note from the above description that the complexity associated to the transmission decisions at node i in each timeslot is $O(K_i^{tr} NM)$, dominated by the need to compute a transmission utility weight for each $(k, (d, m))$ pair in order to choose the pair that maximizes the metric in (22). In contrast, for the local processing decisions at node i in each timeslot, maximizing the metric in (23) over $(k, (d, m))$ can be decomposed into first maximizing $W_i^{(d, m)}(t)$ over (d, m) and then maximizing the metric over k given the maximized $W_i^{(d, m)}(t)$, requiring $O(K_i^{pr} + NM)$ operations. Hence, the overall complexity of the DWCNC decisions at node i in each timeslot is $O(K_i^{tr} NM + K_i^{pr})$.

V. PERFORMANCE ANALYSIS

Theorem 2. *If the rate matrix $\boldsymbol{\lambda} \triangleq \{\lambda_i^{(d, m)}\}$ is strictly interior to the capacity region Λ , then DWCNC stabilizes the wireless computing network, while achieving a (statistical and temporal) average total resource cost arbitrarily close to minimum average cost $\bar{h}^*(\boldsymbol{\lambda})$; i.e.,*

$$\limsup_{t \rightarrow \infty} \frac{1}{t} \sum_{\tau=0}^{t-1} \mathbb{E}[h(\tau)] \leq \bar{h}^*(\boldsymbol{\lambda}) + \frac{NB}{V}$$

$$\limsup_{t \rightarrow \infty} \frac{1}{t} \sum_{\tau, i, d, m} \mathbb{E}\left[Q_i^{(d, m)}(\tau)\right] \leq \frac{NB + V \left[\bar{h}^*(\boldsymbol{\lambda} + \varepsilon \mathbf{1}) - \bar{h}^*(\boldsymbol{\lambda})\right]}{\varepsilon}$$

where B is a constant that depends on the system parameters $R_{i, K_i^{tr}}(\mathbf{s})$, $R_{i, K_i^{pr}}$, A_{\max} , $\xi^{(m)}$, $D_i^{(m)}$ and $r^{(m)}$; ε is a positive constant satisfying $(\boldsymbol{\lambda} + \varepsilon \mathbf{1}) \in \Lambda$; and $\bar{h}^*(\boldsymbol{\lambda})$ denotes the average cost obtained by the optimal solution. \square

The proof of Theorem 2 is shown in Appendix A.

VI. NUMERICAL EXPERIMENTS

In this section, we simulate the performance of DWCNC over 5×10^5 timeslots in a wireless computing network composed of 10 nodes, as shown in Fig. 3a. We assume all nodes are equipped with 1 transmission resource unit, whose activation incurs 1 cost unit. All links can be in 3 channel states, $\{0, 1, 2\}$, with associated transmission rates $\{0, 1, 2\}$ information units per timeslot, respectively. We assume all *active* links, represented by an edge in Fig. 3a, have channel state probabilities $\{0, 1/2, 1/2\}$, while all *inactive* links, represented by the absence of an edge in Fig. 3a, have channel state probabilities $\{1, 0, 0\}$. In terms of processing resources, we assume all 10 nodes are equipped with 1 processing resource unit with processing rate 2 operations per timeslot. Activating the processing resource unit incurs 1 cost unit at every node, except at nodes 5 and 6, where it incurs 0.5 cost units.

We consider 2 services, each composed of 2 functions. All 4 functions require 1 operation per information unit and take 10 timeslots per operation. The first and second functions of Service 1 have a scaling factor of 1 and 2, respectively, while the first and second functions of Service 2 have a scaling factor of 0.5 and 1, respectively. That is, the second function of Service 1 is an expansion function, while the first function of Service 2 is a compression function. We assume one source-destination pair for Service 1, with source in node 2 and destination in node 10, and another for Service 2, with source in node 1 and destination in node 8. Both source nodes receive exogenous arrivals with rate satisfying i.i.d. Poisson distribution across timeslots with mean value 1 information unit per timeslot. We assume every node in the network has the ability to implement all 4 functions.

Fig. 3b shows the processing flow rate distribution for the 4 functions across the network nodes. Observe that function (1, 1) (Service 1, Function 1) is mostly implemented at node 5, which is the node with lowest processing cost along the

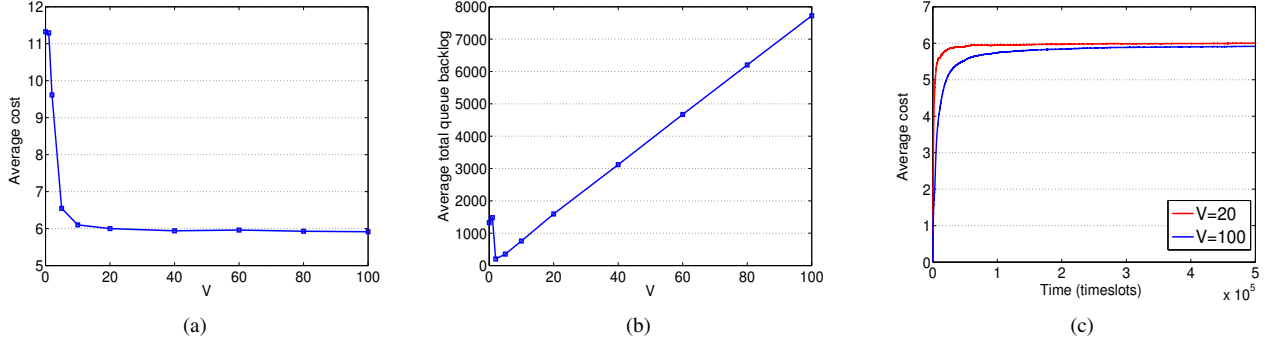


Fig. 4. AgI service distribution over a 10-node wireless computing network. a) Time average cost v.s. control parameter V ; b) Time average total backlog (occupancy) v.s. control parameter V ; c) Time average cost evolution over time.

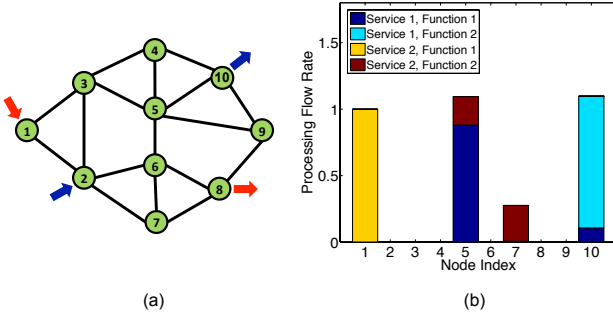


Fig. 3. a) 10-node wireless computing network supporting 2 services with source-destination pairs 2-10 (blue arrows) and 1-8 (red arrows), respectively; b) Processing flow rate distribution across wireless computing nodes.

shortest path from the source (node 2) to the destination (node 10) of Service 1. Note, however, that function (1,2), which is an expansion function, is entirely implemented at node 10, the final destination of Service 1, in order to minimize the transmission cost impact of the larger-size final commodity. It is interesting to see that function (1,1) gets also implemented at node 10 in order to process a small portion of its input commodity. While not shown in the figure, the portion of this commodity that gets processed at node 10 is routed through the path 2-3-4-10 in order to avoid congesting node 5, which is also implementing function (2,2) for Service 2. Looking at Service 2, note that function (2,1), which is a compression function, is entirely implemented at node 1, the source of Service 2, in order to reduce flow size before even entering the network. On the other hand, function (2,2) gets split between nodes 7 and 5 (cheapest processing nodes) in order to avoid congesting node 2, which is on the shortest path of both Service 1 and Service 2. While not shown in the figure, the portion of the input commodity to function (2,2) that gets implemented at node 5 goes around node 2 by following the path 1-3-5-6-8.

Figs. 4a and 4b illustrate the tradeoff between average cost and average total queue backlog as a function of the control parameter V . The average cost is shown to decrease inversely proportional to V until convergence to approximately 5.9 cost units, while the average queue length increases linearly with V , confirming the $[O(1/V), O(V)]$ cost-delay tradeoff of Theorem 2. Finally, Fig. 4c shows the time evolution of

the average cost for $V = 20$ and $V = 100$. Consistent with Theorem 2, a larger V provides a smaller deviation from the minimum average cost at the expense of slower convergence.

VII. CONCLUSIONS

We considered the problem of optimal distribution of augmented information services over wireless computing networks. We characterized the capacity region of a wireless computing network and designed a dynamic wireless computing network control (DWCNC) algorithm that drives local transmissions-plus-processing flow scheduling and resource allocation decisions, shown to achieve arbitrarily close to minimum average network cost, with a tradeoff in network delay. Our solution captures the unique chaining and flow scaling aspects of AgI services, while exploiting the use of the broadcast approach over the wireless channel.

ACKNOWLEDGEMENTS

This work was supported in part by US National Science Foundation.

APPENDIX A PROOF OF THEOREM 2

Let the network's *Lyapunov drift* [4] be defined as

$$\Delta(\mathcal{H}(t)) \triangleq \frac{1}{2} \sum_{i,(d,m)} \mathbb{E} \left[\left(Q_i^{(d,m)}(t+1) \right)^2 - \left(Q_i^{(d,m)}(t) \right)^2 \middle| \mathcal{H}(t) \right],$$

where $\mathcal{H}(t) \triangleq \{\mathbf{Q}(t), \mathbf{S}(t-1)\}$.

By squaring both sides of (7) and summing over $(d,m) \in \mathcal{D} \times \{\mathcal{M}, 0\}$, after further algebraic manipulations, we have

$$\begin{aligned} \Delta(\mathcal{H}(t)) + V \mathbb{E} \{ h(t) | \mathcal{H}(t) \} &\leq NB_0 + \sum_{i,(d,m)} \lambda_i^{(d,m)} Q_i^{(d,m)}(t) \\ &\quad - \sum_i \mathbb{E} \{ \Upsilon(t) + Z_i^{pr}(t) - V h_i^{pr}(t) + Z_i^{tr}(t) - V h_i^{tr}(t) | \mathcal{H}(t) \}, \end{aligned} \quad (24)$$

where NB_0 bounds the sum of the quadratic flow terms, and

$$\begin{aligned} \Upsilon(t) &= \sum_{i,(d,m)} Q_i^{(d,m)}(t) \left[\mu_{pr,i}^{(d,m)}(t) - \mu_{pr,i}^{(d,m)}(t + D_i^{(d,m)}) \right], \\ Z_i^{pr}(t) &= \sum_{(d,m)} \mu_{i,pr}^{(d,m)}(t) \left[Q_i^{(d,m)}(t) - \xi^{(m+1)} Q_i^{(d,m+1)}(t) \right], \\ Z_i^{tr}(t) &= \sum_{l=1}^{N-1} \sum_{(d,m)} \mu_{ig,i,l}^{(d,m)}(t) \left[Q_i^{(d,m)}(t) - Q_{g,i,l}^{(d,m)}(t) \right], \\ h_i^{pr}(t) &= \sum_{k=0}^{K_i^{pr}} w_{i,k}^{pr} y_{i,k}^{pr}(\tau), \quad h_i^{tr}(t) = \sum_{k=0}^{K_i^{tr}} w_{i,k}^{tr} y_{i,k}^{tr}(\tau). \end{aligned} \quad (25)$$

Lemma A.1. *The DWCNC algorithm, at each timeslot t , maximizes $\mathbb{E}\{Z_i^{tr}(t) - Vh_i^{tr}(t)|\mathcal{H}(t)\}$ subject to (2)-(4) and $\mathbb{E}\{Z_i^{pr}(t) - Vh_i^{pr}(t)|\mathcal{H}(t)\}$ subject to (5)-(6).*

Proof: See Appendix B. \blacksquare

Lemma A.1 implies that the right hand side of (24) under DWCNC is no larger than the corresponding expression under the stationary randomized policy of Theorem 1 that supports $(\lambda + \varepsilon\mathbf{1}) \in \Lambda$ and achieves average cost $\bar{h}^*(\lambda + \varepsilon\mathbf{1})$:

$$\begin{aligned} \Delta(\mathcal{H}(t)) + V\mathbb{E}\{h(t)|\mathcal{H}(t)\} &\leq NB + \sum_{i,(d,m)} \lambda_i^{(d,m)} Q_i^{(d,m)}(t) + \\ &\sum_i \mathbb{E}[\Upsilon(t) - Z_i^{*pr}(t) - Vh_i^{*pr}(t) + Z_i^{*tr}(t) - Vh_i^{*tr}(t)|\mathcal{H}(t)], \\ &\leq NB + V\bar{h}^*(\lambda + \varepsilon\mathbf{1}) + \mathbb{E}\{\Upsilon(t)|\mathcal{H}(t)\} - \varepsilon \sum_i \sum_{(d,m)} Q_i^{(d,m)}(t). \end{aligned} \quad (26)$$

Using the fact that $\frac{1}{t} \sum_{\tau=0}^{t-1} \mathbb{E}\{\Upsilon(\tau)\}$ can be upper bounded by a constant NB_{Υ} and denoting $B \triangleq B_0 + B_{\Upsilon}$, then, from inequality (26), standard Lyapunov drift manipulations [4] readily lead to the conclusions of Theorem 2.

APPENDIX B PROOF OF LEMMA A.1

Due to the deterministic nature of the computing channel, maximizing $\mathbb{E}\{Z_i^{pr}(t) - Vh_i^{pr}(t)|\mathcal{H}(t)\}$ is equivalent to maximizing $Z_i^{pr}(t) - Vh_i^{pr}(t)$. The maximization of $Z_i^{pr}(t) - Vh_i^{pr}(t)$ subject to (5)-(6) can be achieved by the greedy choice of commodity (d, m) , resource allocation k , and the greedy assignment of processing rate $\mu_{i,pr}^{(d,m)}(t)$, as described by the local processing decisions of DWCNC.

With respect to the transmission decisions, it follows by plugging (3) into (25) that

$$Z_i^{tr}(t) = \sum_{(d,m)} \sum_{n=1}^{N-1} \sum_{l=n}^{N-1} \mu_{ig_{i,l,n}}^{(d,m)}(t) [Q_i^{(d,m)}(t) - Q_{g_{i,l}}^{(d,m)}(t)]. \quad (27)$$

Let $\beta_{i,tr}^{(d,m)}(t)$ be the fraction of time-frequency resources allocated to the transmission of commodity (d, m) at timeslot t , and let $\eta_{ij,n}^{(d,m)}(t)$ be the fraction of commodity (d, m) decoded by the set of nodes $\Omega_{i,n}$ and forwarded to node j , with $n \leq g_{i,\mathbf{S}(t)}^{-1}(j)$. We then have

$$\mu_{ig_{i,l,n}}^{(d,m)}(t) = \beta_{i,tr}^{(d,m)}(t) \eta_{ig_{i,l,n}}^{(d,m)}(t) [R_{ig_{i,n,k}}(\mathbf{S}(t)) - R_{ig_{i,n-1,k}}(\mathbf{S}(t))], \quad \forall i, t, \quad (28)$$

$$\sum_{(d,m)} \beta_{i,tr}^{(d,m)}(t) \leq 1, \quad \forall i, t, \quad (29)$$

$$\sum_j \eta_{ij,n}^{(d,m)}(t) \leq 1, \quad \forall i, t, (d, m). \quad (30)$$

Plugging (28) into (27) and taking the expectation conditioned on $\mathcal{H}(t)$ and $\{y_{i,k}^{tr}(t) = 1\}$, it follows that

$$\begin{aligned} &\mathbb{E}\{Z_i^{tr}(t)|\mathcal{H}(t), y_{i,k}^{tr}(t) = 1\} \\ &\stackrel{(a)}{\leq} \sum_{(d,m)} \sum_{n=1}^{N-1} \mathbb{E}\left\{\beta_{i,tr}^{(d,m)}(t) \sum_{l=n}^{N-1} \eta_{ig_{i,l,n}}^{(d,m)}(t) W_{ig_{i,l}}^{(d,m)}(t) \right. \\ &\quad \left. \times [R_{ig_{i,n,k}}(\mathbf{S}(t)) - R_{ig_{i,n-1,k}}(\mathbf{S}(t))]| \mathcal{H}(t), y_{i,k}^{tr}(t) = 1\right\} \end{aligned}$$

$$\begin{aligned} &\stackrel{(b)}{\leq} \sum_{(d,m)} \sum_{n=1}^{N-1} \mathbb{E}\left\{\beta_{i,tr}^{(d,m)}(t) \max_{j \in \Omega_{i,n}(\mathbf{S}(t))} \{W_{ij}^{(d,m)}(t)\} \right. \\ &\quad \left. \times [R_{ig_{i,n,k}}(\mathbf{S}(t)) - R_{ig_{i,n-1,k}}(\mathbf{S}(t))]| \mathcal{H}(t), y_{i,k}^{tr}(t) = 1\right\}, \\ &\stackrel{(c)}{\equiv} \sum_{(d,m)} \mathbb{E}\left\{\beta_{i,tr}^{(d,m)}(t) | \mathcal{H}(t), y_{i,k}^{tr}(t) = 1\right\} \sum_{n=1}^{N-1} \mathbb{E}\left\{\max_{j \in \Omega_{i,n}(\mathbf{S}(t))} \{W_{ij}^{(d,m)}(t)\} \right. \\ &\quad \left. \times [R_{ig_{i,n,k}}(\mathbf{S}(t)) - R_{ig_{i,n-1,k}}(\mathbf{S}(t))]| \mathcal{H}(t), y_{i,k}^{tr}(t) = 1\right\} \\ &\stackrel{(d)}{\leq} \max_{(d,m)} \left\{ \sum_{n=1}^{N-1} \mathbb{E}\left\{\max_{j \in \Omega_{i,n}(\mathbf{S}(t))} \{W_{ij}^{(d,m)}(t)\} \right. \right. \\ &\quad \left. \left. \times [R_{ig_{i,n,k}}(\mathbf{S}(t)) - R_{ig_{i,n-1,k}}(\mathbf{S}(t))]| \mathcal{H}(t), y_{i,k}^{tr}(t) = 1\right\} \right\} \\ &= \max_{(d,m)} \{W_{i,k,tr}^{(d,m)}(t)\} \quad (31) \end{aligned}$$

In (31), inequality (a) follows from the definition of $W_{ij}^{(d,\phi,m)}(t)$; inequality (b) follows from (30); equality (c) holds because, given $\mathcal{H}(t)$ and $\{y_{i,k}^{tr}(t) = 1\}$, the values of $R_{ig_{i,n,k}}(\mathbf{S}(t)) - R_{ig_{i,n-1,k}}(\mathbf{S}(t))$ and $\max_{j \in \Omega_{i,n}(\mathbf{S}(t))} \{W_{ij}^{(d,m)}(t)\}$ are a function of $\mathbf{S}(t)$ and independent of $\beta_{i,tr}^{(d,m)}(t)$; inequality (d) follows from (29).

Finally, based on (31), we further have

$$\begin{aligned} &\mathbb{E}\{Z_i^{tr}(t) - Vh_i^{tr}(t)|\mathcal{H}(t)\} \\ &\leq \sum_{k=0}^{K_i^{tr}} \left[\max_{(d,m)} \{W_{i,k,tr}^{(d,m)}(t)\} - Vw_{i,k}^{tr} \right] \Pr[y_{i,k}^{tr}(t) = 1] \\ &\stackrel{(e)}{\leq} \max_{k,(d,m)} \{W_{i,k,tr}^{(d,m)}(t) - Vw_{i,k}^{tr}\}, \quad (32) \end{aligned}$$

In (31) and (32), the upper bounds (a) and (b) can be achieved by implementing step 5) of the local transmission decisions of DWCNC; the upper bound (d) and (e) can be achieved by implementing step 3) of the local transmission decisions of DWCNC, concluding the proof of Lemma A.1.

REFERENCES

- [1] M. Barcelo, J. Llorca, A. M. Tulino, N. Raman, "The Cloud Service Distribution Problem in Distributed Cloud Networks," *IEEE ICC*, Sept. 2015.
- [2] H. Feng, J. Llorca, A. M. Tulino and A. F. Molisch, "Dynamic Network Service Optimization in Distributed Cloud Networks," *IEEE INFOCOM SWFAN Workshop*, Sept. 2016.
- [3] H. Feng, J. Llorca, A. M. Tulino and A. F. Molisch, "Optimal Dynamic Resource Allocation in Distributed Cloud Networks," *IEEE ICC*, Sept. 2016.
- [4] M. J. Neely, "Stochastic network optimization with application to communication and queueing systems", *Synthesis Lectures on Communication Networks*, Morgan & Claypool, 2010.
- [5] M. J. Neely, "Optimal Backpressure Routing for Wireless Networks with Multi-Receiver Diversity", *Ad Hoc Networks*, vol. 7, pp. 862-881, July 2009.
- [6] H. Feng and A. F. Molisch, "Diversity Backpressure Scheduling and Routing with Mutual Information Accumulation in Wireless Ad-hoc Networks", *IEEE Transactions on Information Theory*, vol. 62, no. 12: pp. 7299-7323, Dec. 2016.
- [7] S. Shamai and A. Steiner, "A broadcast approach for a single-user slowly fading MIMO channel." *IEEE Transactions on Information Theory*, vol. 49, no. 10: pp. 2617-2635, Oct. 2003.
- [8] A. M. Tulino, G. Caire and S. Shamai, "Broadcast approach for the sparse-input random-sampled MIMO Gaussian channel." *IEEE ISIT*, Aug. 2014.
- [9] A. El Gamal and Y.-H. Kim, *Network information theory*. Cambridge University Press, 2011.
- [10] H. Feng, J. Llorca, A. M. Tulino, A. F. Molisch, "Optimal Delivery of Augmented Information Services over Wireless Computing Networks", *to be submitted*.