# Application-Aware Resource Provisioning in a Heterogeneous Internet of Things

Eric Sturzinger*, Massimo Tornatore†, and Biswanath Mukherjee‡

*‡University of California, Davis †Politecnico di Milano

Email: *emsturzinger@ucdavis.edu, †tornator@polimi.it ‡bmukherjee@ucdavis.edu

*Abstract*—Internet of Things (IoT) traffic will become increasingly heterogeneous not only in terms of traditional metrics as required bandwidth and maximum latency, but also in terms of functional requirements such as compute power and temporary storage. Sophisticated planning and engineering approaches must be adopted by service providers to account for this heterogeneity, inherent in IoT applications. Metropolitan Area Networks (MANs) are ideally suited to manage and implement resource provisioning of heterogeneous IoT application traffic and, as a result, possess a unique ability to conserve MAN and Wide Area Network (WAN) bandwidth costs. We propose a novel comprehensive MAN resource provisioning model in a hybrid fog-cloud architecture which decouples compute and storage functions while accounting for traffic of a set of heterogeneous parameterized application profiles. This is intended to assist the MAN service provider to minimize the total operational cost of provisioning IoT traffic demands as well as provide a framework for dynamic lightpath reallocation within the MAN. The model demonstrates which application profile and topological parameters have the most significant effect on the individual cost components. As a result of the model, we demonstrate that optimal resource provisioning, i.e. whether functions are placed in the fog or cloud, depends heavily on application computational complexity, compression factor, and latency budget, as well as proportions of local and global traffic.

## I. INTRODUCTION

Internet of Things (IoT) traffic is expected to consume an increasing share of the total Internet bandwidth for the foreseeable future. Cisco predicts more than 50 billion individual devices will be connected by 2020 [1]. Traffic generated and consumed by these devices and networks, due to their anticipated rapid deployment, will require more intelligent and sophisticated network engineering and traffic provisioning approaches in Metropolitan Area Networks (MANs). MANs typically perform traffic distribution and aggregation (at the distribution layer) over a traditional three-tiered hierarchical network, between the core and access layers. With ever-increasing and heterogeneous IoT (and machine to machine - M2M) traffic requiring a more diverse set of network resources under more restrictive constraints, "Smart MANs" become a crucial network segment where complex resource-allocation decisions should be performed to minimize the total provisioning cost of heterogeneous IoT traffic. Specifically, a smart MAN must support application, destination, and resource-aware compute, storage, and routing solutions to accomplish this objective.

The Central Office Re-architectured as a Data Center (CORD) concept [2] seeks to "bring data center economies and cloud agility to service providers for their residential, enterprise, and mobile customers." It will allow COs to perform functions similar to large data centers at a much smaller scale. A CORD implementation in MANs would allow for an intelligent hybrid fog-cloud system, facilitated by our functional provisioning model. To maintain minimal operational cost, it will be necessary to adjust processing and storage locations as traffic fluctuates in application profile mixture and overall demand on any time scale. This will necessitate lightpath reconfiguration on a somewhat periodic basis or in response to less predictable traffic fluidity. Reconfigurable Optical Add-Drop Multiplexers (ROADMs) could perform this function in response to changes in traffic demand to conserve costs and provide enough bandwidth on each respective lightpath.

Cloud computing and storage are the catalysts that will allow compute-intensive IoT applications to flourish. Data center economies of scale allow large cloud service providers, such as Google Cloud Platform [3] and Amazon Web Services [4], to supply inexpensive compute and storage capabilities to customers. For real-time applications, cloud computing may become impractical, due to the significant propagation delays of the core (backbone) network, compared to maximum latency thresholds. Thus, the recent exploration into fog computing, which, according to [5], "extends the cloud to be closer to the things that produce and act on IoT data." It also documents several goals of an IoT-focused computing model, including, but not limited to: minimizing latency, conserving network bandwidth, and moving data to the best place for processing. Latency-sensitive, high-bandwidth applications, such as augmented/virtual reality, are prime examples of when intelligent fog-cloud resource allocation decisions would have a significant impact on total operational costs.

Our primary contributions in this study are: a unique application profile and how its individual parameters affect optimal resource-provisioning solutions and the decoupling of processing and storage functions - i.e. how profile subsets are modeled according to the necessity of each. The rest of the paper is organized as follows: Section II details related work, Section III describes the application profile and hybrid MAN topological parameters, Section IV shows the mathematical formulation, Section V discusses simulation results and analysis, and Section VI concludes the study.

## II. RELATED WORK

Recently, the technical literature is migrating from studies describing general architectural frameworks for IoT applica-

tions will operate to modeling anticipated IoT traffic characteristics and investigating resource allocation and analysis. It is widely agreed that cloud (and eventually fog) computing and storage will enable the vast increase in usable and archived information generated by IoT devices and networks. A distributed cloud-computing architecture is proposed in [6] where applications with various latency requirements are treated according to their source location by the physically-distributed virtual cloud. Ref. [7] puts forth a cloud-of-things architecture which horizontally integrates heterogeneous application domains within their respective vertical silos, i.e., smart city, health care, traffic monitoring, etc. via multiple abstraction layers. A description of fog computing's general role in IoT is described in [8]. Ref. [9] discusses the potential effects of mobility, reliable control and activation, and data aggregation and analytics on fog computing. A dynamic, fog-based resource management model is presented in [10] where three different types of end devices - static, small mobile, and large mobile - require different sets of resources.

Recent works have studied performance tradeoffs within a hybrid fog-cloud environment. Ref. [11] stresses that "the Fog complements the Cloud, does not substitute it," and presents a distributed fog infrastructure with Embedded Systems and Sensors and the Data Center Cloud as the bottom and top layers, respectively. Ref. [12] presents a model of four subsystems: LAN, Fog, WAN, and Cloud. The authors model the power consumption-delay tradeoff; however, offered traffic is not heterogeneous in nature, i.e., no consideration is given to respective application requirements. A performance comparison between a two-tier and three-tier cloud-of-things architecture is conducted in [13], the latter outperforming the former. The two-tier system consists of the cloud tier and physical end devices while the three-tier system includes the fog as the middle tier.

## III. MODELING OF APPLICATION PROFILES AND TOPOLOGICAL PARAMETERS

Our model consists of two primary components: a set of application profiles $\mathcal{A}$ and a MAN topology $G(\mathcal{N}, \mathcal{L})$ with associated node and link attributes. This approach allows us to understand which profile and network characteristics have a more significant impact on the comprehensive provisioning solution which minimizes total cost.

### A. Application Profile

We divide the set of application profiles $\mathcal{A}$ into four subsets such that:

$$\mathcal{A} = \mathcal{A}_p \cup \mathcal{A}_{sp} \cup \mathcal{A}_s \cup \mathcal{A}_n$$

where $\mathcal{A}_p$ represents the subset of profiles requiring processing only (cloud gaming, virtual reality), $\mathcal{A}_{sp}$ the subset of those requiring processing and storage (smart grid management [14]), $\mathcal{A}_s$ the subset of those requiring storage only (data used for future analytics in its original form), and $\mathcal{A}_n$ the subset of those requiring neither processing nor storage (simple machine

TABLE I: APPLICATION PROFILE

| Parameter | Description | Units |
|---|---|---|
| $\alpha$ | Computational complexity | $\frac{CPU}{Mbit/s}$ |
| $\beta$ | Compression factor | $N/A$ |
| $\kappa$ | Average flow size | $Mbit$ |
| $\Delta$ | Minimum data-availability time | $hrs$ |
| $\Theta$ | Maximum one-way latency | $ms$ |

TABLE II: TOPOLOGICAL PROPERTIES

| Name | Description | Units |
|---|---|---|
| $\mu_m$ | Unit compute cost at node $m$ | $\$/CPU$ |
| $\nu_f$ | Unit storage cost at node $f$ | $\$/GB$ |
| $C_m$ | Computation capacity of node $m$ | $CPU$ |
| $S_f$ | Storage Capacity of node $f$ | $GB$ |
| $\Lambda$ | Unit cost MAN BW | $\$/Mbit/s$ |
| $\epsilon_{up}$ | Unit cost upstream WAN BW | $\$/Mbit/s$ |
| $\epsilon_{down}$ | Unit cost downstream WAN BW | $\$/Mbit/s$ |
| $\tau_m$ | Computation time factor | $/CPU$ |
| $P_k^{s,m}$ | Total propagation delay of $k^{th}$ path from $s \to m$ | $ms$ |
| $T_k^{s,m}$ | Total transmission delay of $k^{th}$ path from $s \to m$ | $ms$ |
| $D_k^{s,m}$ | Total delay of $k^{th}$ path from $s \to m$ | $ms$ |

to machine traffic). Each individual application profile consists of a unique combination of parameters given by Table I.

Computational complexity, $\alpha$, refers to the number of CPUs, or percentage of a single CPU, required to process 1 Mbps of traffic of the respective application. Compression factor, $\beta$, is the ratio of processed output traffic to pre-processed input traffic at the processing node, which generally has the range: $0 < \beta \leq 1$. Average flow size, $\kappa$, represents the average amount of data that must be processed as a single entity. Applications that generate a large amount of data per flow will incur a much larger processing delay as their flow size can be several orders of magnitude greater than a single packet, which must be processed as a single entity. Minimum data-availability time, $\Delta$, refers to the minimum amount of time data of a particular application must remain accessible to a customer (human or machine). Maximum one-way latency, $\Theta$, is the most constraining element of this model as a significant proportion of next-generation application profiles. In this model, we do not define latency as completely end to end. For traffic with a destination in the local MAN, it is simply the delay from the source node (where the traffic entered the network) to the destination node, which includes delay across the core (for cloud processing and/or storage) and processing delay, if applicable.

### B. Topological Properties

To model the general characteristics of a MAN topology consisting of two hierarchical levels of interconnected rings, we begin with a graph $G(\mathcal{N}, \mathcal{L})$ where $\mathcal{N}$ is the set of nodes

Fig. 1: MAN topology with cloud extension

to be included in the analysis, ranging from data center nodes to edge nodes. Accordingly, $\mathcal{L}$ represents the set of all MAN links. Analogous to the subset of application profiles, each node in the MAN belongs to a specific subset of $\mathcal{N}$ such that

$$\mathcal{N} = \mathcal{N}_g \cup \mathcal{N}_c \cup \mathcal{N}_{DC}$$

where $\mathcal{N}_g = \{1, 2, ..., 16\}$ in Fig. 1 is the set of nodes that are ingress/source nodes, $\mathcal{N}_c = \{17, 18\}$ is the set of nodes that, combined, act as the Core CO (interface to WAN), and $\mathcal{N}_{DC} = \{19, 20\}$ is the set of all data center nodes (for cloud processing and/or storage). Other important subsets of nodes are those which contain processing and storage capabilities, $\mathcal{N}_p$ and $\mathcal{N}_s$, respectively. Subsets $\mathcal{N}_{pl}$ and $\mathcal{N}_{sl}$ contain only those physically located within the local MAN (excluding DCs). $\mathcal{F}_c = \{C1, C2, ..., C24\} \notin \mathcal{N}$ is the set of core nodes which are possible destinations of global traffic. Another parameter could identify those nodes with ROADMs available. Additionally, in older, more remote, rural networks, other constraints would be included related to the capabilities of installed optical equipment that can operate only at lower data rates or shorter distances than more advanced equipment at other locations.

Each node is characterized by a given unit compute and storage cost ($\mu_m$ & $\nu_f$) and a total compute and storage capacity ($C_m$ & $S_f$). We assign each node to specific processing and storage tiers (within the hybrid architecture), which dictate its unit costs and total capacities, depicted by subscripts of $P$ and $S$ in Fig. 1, respectively. DC nodes are assigned the smallest units costs and largest capacities, tier 4, while each subordinate tier is assigned smaller capacities and larger unit costs. Upstream and downstream WAN bandwidth costs, defined by $\epsilon_{up}$ and $\epsilon_{down}$, respectively, are the unit costs a backbone service provider charges for a WAN connection. MAN bandwidth unit cost, $\Delta$, is the approximate cost to operate an internal MAN link. One component of the model solution is the total amount of traffic each link would be required to support. As each path solution is determined by its source destination pair, and traffic demands fluctuate in short

and long-term, lightpaths can be reconfigured within the MAN as a result of changing optimal resource allocation solutions. This will require flexible systems, such as ROADMs, to implement this scheme, especially in more fluid scenarios where energy costs change throughout the day or on a longer-term seasonal basis. We only consider propagation ($P_k^{s,m}$), transmission ($T_k^{s,m}$), and processing delays, defined later, as individual components of total latency. We assume a 25-ms WAN delay (one-way) and a round-trip time (RTT) of 50 ms [15], affecting residual latency budgets of global traffic. Link transmission rates are a result of routing solutions and thus would render the model non-linear.

## IV. Mathematical Formulation

The objective of this resource assignment optimization problem is to minimize the combined operational cost while satisfying application and topological constraints. This is a network planning problem, and we model it using an Integer Linear Program (ILP).

### A. Variable Definition

We calculated k-shortest paths with respect to total path latency ($D_k^{s,m}$), setting $k = 4$ to ensure each source node has at least one path to each DC node through each core CO node to provide limited path diversity to each DC. Offered traffic from $s$ to $f$ of application $a$, processed at $m$ is defined as $v_{a,f}^{s,m}$. For any applications requiring storage, there is no initially defined destination node and thus offered traffic must be mirrored across all possible storage nodes, visually depicted in Fig. 2. $\mathcal{R}_{i,j}$ represents the set of all admissible paths containing link $i, j$.

$\mathbf{x}_{a,f}^{s,m} = 1$ if traffic of application profile $a \in \mathcal{A}$, of source $s \in \mathcal{N}_g$, is processed at node $m \in \mathcal{N}_p$, stored at node $f \in \mathcal{N}_s$, or both. If $a \in \mathcal{A}_s$, $m = f$ (Fig. 2d). If $a \in \mathcal{A}_p$, $f$ represents the final destination, whether in the local MAN (Fig. 2a) or in the core (Figs. 2b and 2c).

$\mathbf{r}_{a,k,f}^{s,m} = 1$ if pre-processed or direct traffic is routed over the $k^{th}$ path between node $s$ and processing (or destination - Figs. 2d and 2f) node $m$.

$\mathbf{r'}_{a,k,f}^{s,m} = 1$ if post-processed traffic is routed over the $k^{th}$ path between node $m$ and node $f$. It represents post-processed traffic where $f$ is either the destination (Fig. 2a) or storage node (Fig. 2e).

$\mathbf{r''}_{a,k,f}^{s,m,d} = 1$ if post-processed traffic of source $s$ and final destination $f$, is routed over the $k^{th}$ path between processing node $m$ and core CO node $d$ - Fig. 2b.

### B. Objective Function

We define the objective function as the total resource provisioning cost necessary to support the given offered traffic over the given network topology while satisfying all application and topological constraints. Total cost consists of the following five components and are defined below.

$$Cost_t = Cost_p + Cost_s + Cost_u + Cost_d + Cost_c \quad (1)$$

(a) $a \in \mathcal{A}_p$ - Local     (b) $a \in \mathcal{A}_p$ - Fog Processing - Global     (c) $a \in \mathcal{A}_p$ - Cloud Processing - Global

(d) $a \in \mathcal{A}_s$     (e) $a \in \mathcal{A}_{sp}$     (f) $a \in \mathcal{A}_n$ Pt to Pt - Global
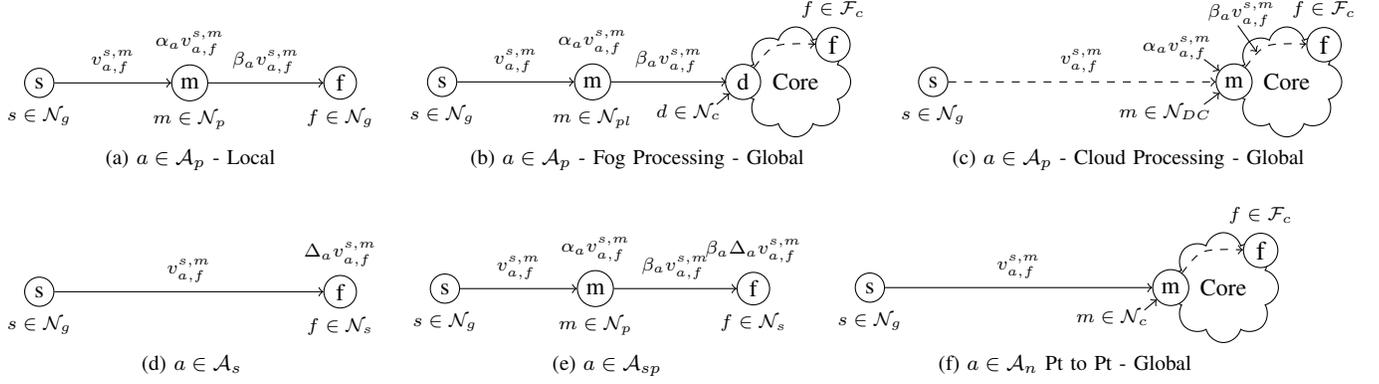
Fig. 2: Flow scenarios as a function of application profile subset and destination

We now describe the mathematical formulation of each cost component. Total processing cost is defined as:

$$Cost_p = \sum_{m \in \mathcal{N}_p} \mu_m \sum_{a \in \mathcal{A}_p \cup \mathcal{A}_{sp}} \alpha_a \sum_{s \in \mathcal{N}_g} \sum_{f \in \mathcal{N}_g \cup \mathcal{F}_c} x_{a,f}^{s,m} v_{a,f}^{s,m} \quad (2)$$

Eq. (2) represents all possible scenarios in which computation is performed. Storage cost is defined in Eq. (3), where the 2nd term is multiplied by compression factor, $\beta$ (see Fig. 2e).

$$Cost_s = \sum_{a \in \mathcal{A}_s} \Delta_a \sum_{f \in \mathcal{N}_s} \nu_f \sum_{s \in \mathcal{N}_g} x_{a,f}^{s,m} v_{a,f}^{s,m}$$
$$+ \sum_{a \in \mathcal{A}_{sp}} \beta_a \Delta_a \sum_{f \in \mathcal{N}_s} \nu_f \sum_{s \in \mathcal{N}_g} \sum_{m \in \mathcal{N}_p} x_{a,f}^{s,m} v_{a,f}^{s,m} \quad (3)$$

Eq. (4) represents total upstream WAN BW cost. First term is visualized in Fig. 2b, second in Fig. 2c, and third in Fig. 2f.

$$Cost_u = \epsilon_{up} \Big[ \sum_{a \in \mathcal{A}_p \cup \mathcal{A}_{sp}} \beta_a \sum_{s \in \mathcal{N}_g} \sum_{m \in \mathcal{N}_{pl}} \sum_{f \in \mathcal{N}_{DC} \cup \mathcal{F}_c} x_{a,f}^{s,m} v_{a,f}^{s,m}$$
$$+ \sum_{a \in \mathcal{A}_p \cup \mathcal{A}_s \cup \mathcal{A}_{sp}} \sum_{s \in \mathcal{N}_g} \sum_{m \in \mathcal{N}_{DC}} \sum_{f \in \mathcal{N}_g \cup \mathcal{N}_s \cup \mathcal{F}_c} x_{a,f}^{s,m} v_{a,f}^{s,m}$$
$$+ \sum_{a \in \mathcal{A}_n} \sum_{s \in \mathcal{N}_g} \sum_{m \in \mathcal{N}_c} \sum_{f \in \mathcal{F}_c} v_{a,f}^{s,m} \Big] \quad (4)$$

We only consider downstream traffic originally generated in the local MAN when it is processed in the cloud and sent back to the MAN to a destination or storage node.

$$Cost_d = \epsilon_{down} \sum_{a \in \mathcal{A}_p \cup \mathcal{A}_{sp}} \beta_a \sum_{s \in \mathcal{N}_g} \sum_{m \in \mathcal{N}_{DC}} \sum_{f \in \mathcal{N}_g \cup \mathcal{N}_{sl}} x_{a,f}^{s,m} v_{a,f}^{s,m} \quad (5)$$

We now define the internal MAN link capacity costs. Eq. (6) represents local and global traffic not requiring processing while Eq. (7) is the amount of pre-processed traffic destined for its processing node.

$$Cap_{1,i,j} = \sum_{a \in \mathcal{A}_n \cup \mathcal{A}_s} \sum_{s \in \mathcal{N}_g} \Big[ \sum_{f \in \mathcal{N}_g \cup \mathcal{N}_{sl}} \sum_{k \in \mathcal{R}_{i,j}} r_{a,k,f}^{s,m} v_{a,f}^{s,m}$$
$$+ \sum_{f \in \mathcal{F}_c} \sum_{m \in \mathcal{N}_c} \sum_{k \in \mathcal{R}_{i,j}} r_{a,k,f}^{s,m} v_{a,f}^{s,m} \Big] \quad (6)$$

$$Cap_{2,i,j} =$$
$$\sum_{a \in \mathcal{A}_p \cup \mathcal{A}_{sp}} \sum_{s \in \mathcal{N}_g} \sum_{m \in \mathcal{N}_p} \sum_{f \in \mathcal{N}_g \cup \mathcal{N}_s \cup \mathcal{F}_c} \sum_{k \in \mathcal{R}_{i,j}} r_{a,k,f}^{s,m} v_{a,f}^{s,m} \quad (7)$$

$$Cap_{3,i,j} = \sum_{a \in \mathcal{A}_p \cup \mathcal{A}_{sp}} \beta_a \sum_{s \in \mathcal{N}_g} \sum_{m \in \mathcal{N}_p} \sum_{f \in \mathcal{N}_g \cup \mathcal{N}_s \cup \mathcal{F}_c} \Big[ \sum_{k \in \mathcal{R}_{i,j}} r_{a,k,f}'^{s,m}$$
$$+ \sum_{d \in \mathcal{N}_c} \sum_{k \in \mathcal{R}_{i,j}} r_{a,k,f}''^{s,m,d} \Big] v_{a,f}^{s,m} \quad (8)$$

Eq. (8) includes the compression factor $\beta$, which converts pre-processed traffic to post-processed traffic. The total internal MAN link capacity is therefore:

$$Cost_c = \Lambda \sum_{i,j \in \mathcal{L}_l} \sum_h Cap_{h,i,j} \quad (9)$$

where $\mathcal{L}_l$ is the set of local links and $\Lambda$ is the unit cost per Mbps of MAN BW. We have now defined all components of the objective function and write the problem statement as:

$$min(Cost_t) \quad (10)$$

subject to the following constraints:

*C. Function Placement*

$$\sum_{m \in \mathcal{N}_p} x_{a,f}^{s,m} = 1, \forall\ (a \in \mathcal{A}_p, s \in \mathcal{N}_g, f \in \mathcal{N}_g \cup \mathcal{F}_c) \quad (11)$$

$$\sum_{f \in \mathcal{N}_s} x_{a,f}^{s,m} = 1, \forall\ (a \in \mathcal{A}_s, s \in \mathcal{N}_g, m = f) \quad (12)$$

$$\sum_{m \in \mathcal{N}_p} \sum_{f \in \mathcal{N}_s} x_{a,f}^{s,m} = 1, \forall\ (a \in \mathcal{A}_{sp}, s \in \mathcal{N}_g) \quad (13)$$
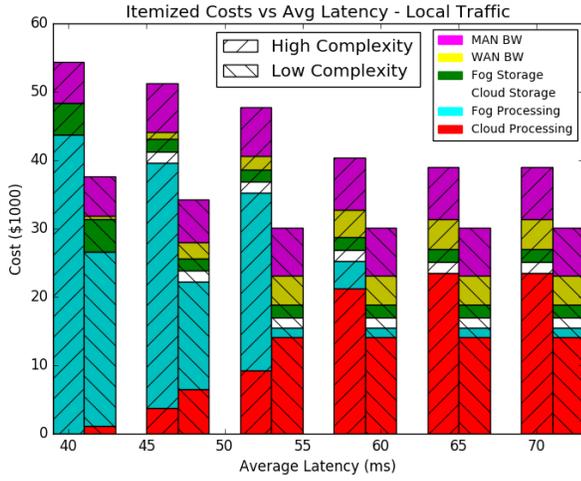
Fig. 3: Itemized costs as a function of average application latency budget

Eqs. (11)-(13) dictate that only one node processes traffic of a single application profile ($a \in \mathcal{A}_p$) and source/destination pair, a single storage node for each application profile ($a \in \mathcal{A}_s$) and source node, and a single node for each function for each application profile ($a \in \mathcal{A}_{sp}$) and source node.

### D. Compute Capacity

The following equation restricts total computation performed at each node $m \in \mathcal{N}_p$:

$$\sum_{a \in \mathcal{A}_p \cup \mathcal{A}_{sp}} \alpha_a \sum_{s \in \mathcal{N}_g} \sum_{f \in \mathcal{N}_g \cup \mathcal{N}_s \cup \mathcal{F}_c} x_{a,f}^{s,m} v_{a,f}^{s,m} \leq C_m, \forall m \in \mathcal{N}_p \tag{14}$$

### E. Storage Capacity

Storage capacity of each storage-capable node is defined in Eq. (15). If processed data is stored (2nd term), initial offered traffic must also be multiplied by compression factor, $\beta$:

$$\sum_{a \in \mathcal{A}_s} \Delta_a \sum_{s \in \mathcal{N}_g} x_{a,f}^{s,m} v_{a,f}^{s,m}$$
$$+ \sum_{a \in \mathcal{A}_{sp}} \beta_a \Delta_a \sum_{s \in \mathcal{N}_g} \sum_{m \in \mathcal{N}_p} x_{a,f}^{s,m} v_{a,f}^{s,m} \leq S_f, \; \forall f \in \mathcal{N}_s \tag{15}$$

### F. Solenoidality

Eqs. (16)-(18) dictate unsplittable traffic on all paths between any combination of source, destination, processing, and storage nodes, flow scenarios of which are shown in Fig. 2.

$$\sum_{k \in \mathcal{R}_{i,j}} r_{a,k,f}^{s,m} = x_{a,f}^{s,m},$$
$$\forall(a \in \mathcal{A}_p \cup \mathcal{A}_{sp} \cup \mathcal{A}_s, s \in \mathcal{N}_g, m \in \mathcal{N}_p, \tag{16}$$
$$f \in \mathcal{N}_g \cup \mathcal{N}_s \cup \mathcal{F}_c)$$

$$\sum_{k \in \mathcal{R}_{i,j}} r_{a,k,f}'^{s,m} = x_{a,f}^{s,m}, \tag{17}$$
$$\forall(a \in \mathcal{A}_p \cup \mathcal{A}_{sp}, s \in \mathcal{N}_g, m \in \mathcal{N}_p, f \in \mathcal{N}_g \cup \mathcal{N}_s)$$

$$\sum_{m \in \mathcal{N}_c} \sum_{k \in \mathcal{R}_{i,j}} r_{a,k,f}''^{s,m,d} = x_{a,f}^{s,m}, \tag{18}$$
$$\forall(a \in \mathcal{A}_p, s \in \mathcal{N}_g, f \in \mathcal{F}_c)$$

### G. Latency

To formulate latency constraints, we must first define a secondary variable, indicating the average processing delay at node $m$ of application $a$, $\gamma_{a,m}$:

$$\gamma_{a,m} = \alpha_a \kappa_a \tau_m \tag{19}$$

We approximate delay as a combination of computational complexity and average flow size of application $a$ as well as the computation factor of node $m$, which is normalized to tier four data center nodes ($\tau_{19}, \tau_{20} = 1$).

$$\sum_{k \in \mathcal{R}_{i,j}} r_{a,k,f}^{s,m} D_k^{s,m} \leq \theta_{a,f}, \forall(a \in A_n \cup \mathcal{A}_s, s \in \mathcal{N}_g, \tag{20}$$
$$m \in \mathcal{N}_l \cup \mathcal{N}_s, f \in \mathcal{N}_g \cup \mathcal{N}_s \cup \mathcal{F}_c)$$

In Eq. (20), no processing is required and $f \in \mathcal{N}$, $\Theta_a = \theta_{a,f}$.

$$\sum_{k \in \mathcal{R}_{i,j}} r_{a,k,f}^{s,m} D_k^{s,m} + \gamma_{a,m} + \sum_k r_{a,k,f}'^{s,m} D_k^{m,f} \leq \theta_{a,f,m}, \tag{21}$$
$$\forall(a \in \mathcal{A}_p \cup \mathcal{A}_{sp}, s \in \mathcal{N}_g, m \in \mathcal{N}_p, f \in \mathcal{N}_g \cup \mathcal{N}_s)$$

Eq. (21) accounts for total delay of the pre-processed data path, processing delay, and processed data path, whether destined for a storage node (Fig. 2e) or a local destination (Fig. 2a).

$$\sum_{k \in \mathcal{R}_{i,j}} r_{a,k,f}^{s,m} D_k^{s,m} + \gamma_{a,m} + \sum_k r_{a,k,f}''^{s,m,d} D_k^{m,d} \leq \theta_{a,f}, \tag{22}$$
$$\forall(a \in \mathcal{A}_p, s \in \mathcal{N}_g, m \in \mathcal{N}_p, d \in \mathcal{N}_c, f \in \mathcal{F}_c)$$

Eq. (22) where $f \in \mathcal{F}_c$, $\Theta_{a,f} = \theta_a - 25$ ms accounts for the average one-way WAN delay, as shown in Fig. 2b.

## V. Simulation and Results

We used PuLP to generate the LP model, written in Python 3.4. Offered traffic is distributed uniformly across all source nodes ($s \in \mathcal{N}_g$) and all local ($f \in \mathcal{N}_g$) and global ($f \in \mathcal{F}_c$) destination nodes, the sum total of which is only 100 Gbps in all plots (thus the small total cost). Any traffic requiring storage is automatically classified as local as it does not have a destination outside of the local MAN and either data center.

Fig. 3 shows the effect of application's latency when comparing those with low and high complexities. Most interesting here is that, under more restrictive latencies, processing delay of higher-complexity applications demand more of the overall budget which results in computation restricted to the fog. As latency increases, processing delay demands less of overall latency and higher-complexity applications can be shifted to the cloud to reduce costs.

Fig. 4 depicts the effect average compression factor has on provisioning solutions, to a greater extent with global as compared to local traffic. Fig. 4a shows how the decoupling of processing and storage is an advantage as computation can remain in the fog while a much smaller volume of traffic is transmitted over the WAN and stored in the cloud at lower cost. MAN BW and cloud processing requirements increase at a higher rate with increasing compression factor at lower complexities as computation is a smaller proportion of total cost. Global traffic of applications with low computational
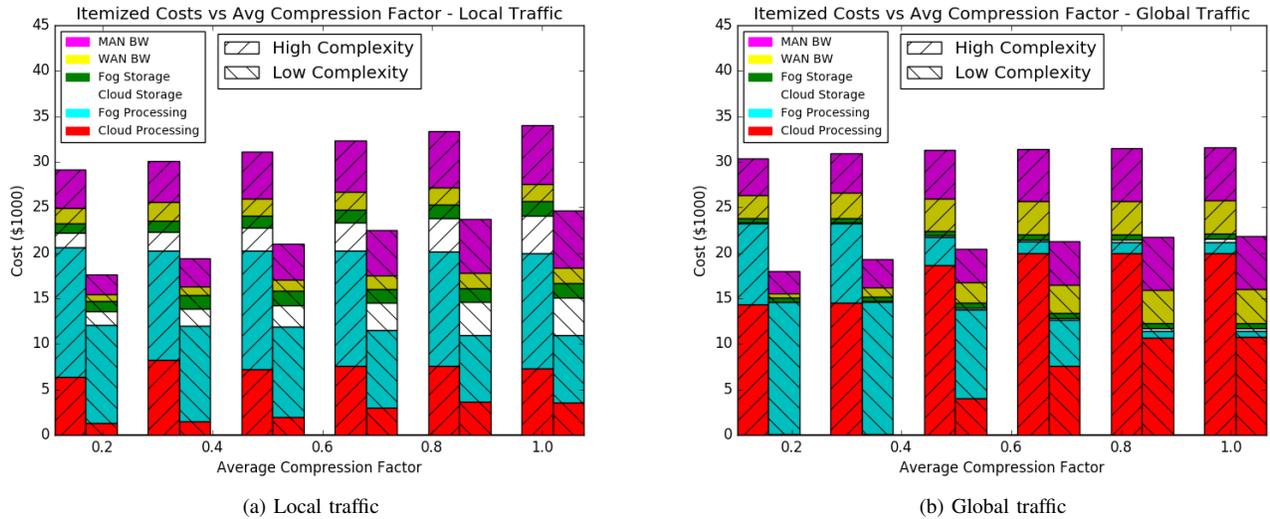
Fig. 4: Comparison of effect of compression factor on itemized costs with primarily (a) local and (b) global traffic.

complexity and low compression factor in Fig. 4b are processed in the fog as WAN BW cost savings overcome higher computational cost, aiding in minimizing the total traffic sent over the WAN. Naturally, processing of both high and low complexity applications shift toward the cloud exclusively as cost savings of a smaller volume of post-processed traffic sent over the WAN diminishes with increasing compression factor. MAN BW also increase at a higher rate with low complexity than high complexity, as with local traffic.

## VI. CONCLUSION

In this study, we divided IoT-specific application profiles into four primary subsets that have various functional and performance requirements and developed an ILP model which determines optimal functional location. We also demonstrated how specific application profile parameters affect individual MAN operating costs in a hybrid fog-cloud architecture. We have briefly discussed how the optical layer would be affected by our model, specifically in lightpath reconfiguration in short and long-term scenarios. Adjusting lightpath configuration is required by a newly computed optimal resource allocation solution which minimizes operational costs in response to constantly changing traffic demands, implemented within the CORD architecture. This would allow the necessary flexibility for MANs to adapt to unpredictable changes in traffic volume and characteristics. In future work, we will quantitatively show how, and in what circumstances, MAN bandwidth is affected by traffic growth, especially that which is nonuniform in source destination pair as well as application profile.

## REFERENCES

[1] "The Internet of Things: How the Next Evolution of the Internet is Changing Everything." White Paper, 2011.

[2] A. Al-Shabibi and L. Peterson, "CORD: Central Office Re-architectured as a Datacenter," *Open Stack Summit*, 2015.

[3] "Google Cloud Platform." https://cloud.google.com. Accessed: 2017.

[4] "Amazon Web Services." https://aws.amazon.com. Accessed: 2017.

[5] "Fog Computing and the Internet of Things: Extend the Cloud to Where the Things Are." White Paper, 2015.

[6] D. Mazmanov, C. Curescu, H. Olsson, A. Ton, and J. Kempf, "Handling Performance Sensitive Native Cloud Applications with Distributed Cloud Computing and SLA Management," pp. 470–475, 2013.

[7] R. Petrolo, V. Loscrì, and N. Mitton, "Towards a Smart City Based on Cloud of Things, a Survey on the Smart City Vision and Paradigms," *Transactions on Emerging Telecommunications Technologies*, 2015.

[8] F. Bonomi, R. Milito, J. Zhu, and S. Addepalli, "Fog Computing and its Role in the Internet of Things," in *Proc., First Edition of the MCC Workshop on Mobile Cloud Computing*, pp. 13–16, ACM, 2012.

[9] M. Yannuzzi, R. Milito, R. Serral-Gracià, D. Montero, and M. Nemirovsky, "Key Ingredients in an IoT Recipe: Fog Computing, Cloud Computing, and More Fog Computing," in *IEEE 19th International Workshop on Computer Aided Modeling and Design of Communication Links and Networks (CAMAD)*, pp. 325–329, 2014.

[10] M. Aazam and E.-N. Huh, "Fog Computing Micro Datacenter Based Dynamic Resource Estimation and Pricing Model for IoT," in *IEEE 29th International Conference on Advanced Information Networking and Applications*, pp. 687–694, 2015.

[11] F. Bonomi, R. Milito, P. Natarajan, and J. Zhu, "Fog Computing: A Platform for Internet of Ihings and Analytics," in *Big Data and Internet of Things: A Roadmap for Smart Environments*, pp. 169–186, Springer, 2014.

[12] R. Deng, R. Lu, C. Lai, and T. H. Luan, "Towards Power Consumption-Delay Tradeoff by Workload Allocation in Cloud-Fog Computing," in *IEEE International Conference on Communications (ICC)*, pp. 3909–3914, 2015.

[13] W. Li, I. Santos, F. C. Delicato, P. F. Pires, L. Pirmez, W. Wei, H. Song, A. Zomaya, and S. Khan, "System Modelling and Performance Evaluation of a Three-tier Cloud of Things," *Future Generation Computer Systems*, 2016.

[14] I. Lee and K. Lee, "The Internet of Things (IoT): Applications, Investments, and Challenges for Enterprises," *Business Horizons*, vol. 58, no. 4, pp. 431–440, 2015.

[15] "NTT America Backbone Network SLA Terms and Conditions." http://www.us.ntt.net/support/sla/terms.cfm. Accessed: 2017-01-02.