

# Towards Systematic Design of Enterprise Networks

Yu-Wei Eric Sung, Xin Sun, Sanjay G.Rao, Geoffrey G. Xie, and David A. Maltz

**Abstract**—Enterprise networks are important, with size and complexity even surpassing carrier networks. Yet, the design of enterprise networks remains ad-hoc and poorly understood. In this paper, we show how a systematic design approach can handle two key areas of enterprise design: virtual local area networks (VLANs) and reachability control. We focus on these tasks given their complexity, prevalence, and time-consuming nature. Our contributions are three-fold. First, we show how these design tasks may be formulated in terms of network-wide performance, security, and resilience requirements. Our formulations capture the correctness and feasibility constraints on the design, and they model each task as one of optimizing desired criteria subject to the constraints. The optimization criteria may further be customized to meet operator-preferred design strategies. Second, we develop a set of algorithms to solve the problems that we formulate. Third, we demonstrate the feasibility and value of our systematic design approach through validation on a large-scale campus network with hundreds of routers and VLANs.

## I. INTRODUCTION

Recent empirical studies reveal that the size of some enterprise networks and the complexity of their routing design rival or even surpass those of carrier networks [1], [2]. Far more enterprise networks than carrier networks are in operation today, and their designs are highly customized to the needs of individual companies, universities, government agencies, or other types of organizations. However, despite their complexity, prevalence, and diversity, enterprise networks have received little attention from the research community.

Managers of enterprise networks face unique design challenges. They need to meet a wider range of security, resilience, and performance requirements than their counterparts at carrier networks. Examples of such challenges include the configuration of virtual local area networks (VLANs) to ease the management of different user groups [3], the integration of multiple routing domains to support company mergers [2], and the installation of packet filters to perform ingress filtering and to control access to privileged databases[4].

The unique challenges of enterprise network design have further exposed the limitations of the existing ad-hoc approach to network design and management. On the one hand, a manager faces high-level constraints such as performance, ease of manageability, security, and resilience to failures. On the other hand, to realize a network design, a manager must manually choose from a slew of protocols, low-level mechanisms, and options. Many of these protocols and mechanisms have profound interactions. However, the current “protocol by protocol” method of network configuration does not allow the network operator to see and control these interactions in a systematic manner. Design faults and configuration errors account for a substantial number of network problems [5], and are exploited by over 65% of cyber-attacks according to

recent statistics [6].

In this paper, we explore the feasibility of adopting a systematic approach to enterprise network design. The key elements include (i) identifying the network-wide performance, security, and resilience requirements of a task; (ii) formulating the requirements as one of optimizing desired (operator-customized) criteria subject to correctness and feasibility constraints on the design; and (iii) developing algorithms and heuristics to solve the formulated problems.

We show that two critical enterprise network design tasks lend themselves to such a systematic approach. These include (i) VLAN design; and (ii) reachability control through placement of packet filters. We model the objectives of VLAN design as achieving low costs associated with broadcast and data traffic, given constraints such as a categorization of hosts into distinct logical groups and a limit on the number of VLANs used. We model the objectives of packet filter placement as optimizing for operator-specified placement criteria such as balancing processing needs across routers, while correctly realizing desired security policies, and meeting feasibility constraints on the processing capacities of routers.

We evaluate the benefits of a systematic design approach in the context of algorithms we developed to solve our formulated problems. Our validations are conducted on a large-scale campus network dataset involving hundreds of routers and VLANs, and a few thousand switches. Beyond the general time savings in realizing a correct and easily customizable design, our results show that through systematic VLAN design, broadcast and data traffic can be reduced by over 24% and 55%, respectively. Our results also highlight the importance of a systematic approach to placing packet filters by identifying inconsistencies in the realization of operator security objectives in the campus network dataset. Overall, these results show the promise of a systematic design approach in these key areas, and are a first but key step towards the top-down design of enterprise networks in general.

## II. FRAMING ENTERPRISE DESIGN TASKS

The nature of the enterprise design problem is little known outside the operational community. For example, there is almost no coverage of this topic in college textbooks. Only through repeated inspections of router configuration files and close interactions with network managers have we obtained a basic understanding of what technical challenges it entails.

We observe that enterprise design can be decomposed into a sequence of distinct stages or tasks. The major tasks in order of execution are: (1) plan physical topology and wiring, (ii) create VLANs and layer-2 topology, (iii) select and configure routing protocols, and (iv) control reachability with packet filters or firewalls. This work focuses on tasks (ii) and (iv) because these tasks have been identified by network managers as challenging

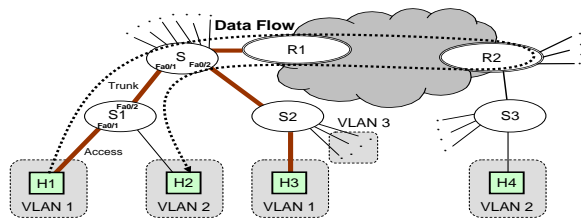


Fig. 1. Communication between different VLANs is routed through designated routers. When H1 talks to H2, R1 acts as a router in the outgoing direction, but as a switch in the return direction.

and time-consuming, and they have been relatively unexplored by the research community. In the rest of this section, we give a high level description of the technical challenges facing VLAN design and reachability control.

#### A. VLAN Design

Operators reduce the complexity of their configuration tasks by thinking about users as collective groups based on the role of each user in the organization (e.g., what resources they should be able to access). Today, these groupings are most commonly implemented by VLANs, which take a set of users in physically disparate locations and place them into a single logical subnet, even if the users are connected to different switches. For instance, an enterprise policy may permit access only for all sales personnel, and it may be desirable to ensure these users receive IP addresses from the same subnet so that routing policies and packet filters can be applied to them as a group. Consider Fig. 1. S, S1~S3 are switches, and R1~R2 are routers. Notice that even though hosts H1 and H3 are physically separated, they are both part of VLAN 1. Likewise, hosts H2 and H4 belong to VLAN 2.

Each VLAN constitutes a separate broadcast domain. Therefore, it is important to ensure that broadcast traffic is properly constrained to reduce unnecessary traffic for increased performance and security. To achieve this, every link is configured to permit only traffic for appropriate VLANs. In Fig. 1, the link S1-H1 is configured as an *access* link and forwards only VLAN 1 traffic. The link S1-S is configured as a *trunk* link and permits traffic for multiple explicitly specified VLANs (in this case, VLANs 1 and 2). Typically, a separate spanning tree rooted at a *root bridge* is constructed per VLAN. For example, the collection of bold links forms the spanning tree of VLAN 1, with S being its root bridge.

Each publicly accessible VLAN is assigned with what we term a *designated (gateway) router* for that VLAN. When a host inside a VLAN communicates with a host outside, the designated router is the first (last) router for outgoing (incoming) packets. In Fig. 1, R1 and R2 are respectively the designated routers for VLAN 1 and VLAN 2. The IP level path between H1 and H2 is:  $H1 - R1 \dots R2 - H2$ , with  $R1 \dots R2$  denoting there could be other routers in the path. The path of data flow is also highlighted in the figure.

In VLAN design, an operator faces two key tasks with unique technical challenges:

**(1) Grouping hosts into VLANs:** The operator must decide the appropriate number of VLANs in the design, and determine which hosts must belong to each VLAN. In doing

so, three factors must be considered. First, security policies and management objectives may influence the decision. For example, in a campus network, the manager may desire to separate faculty and student machines into different VLANs in order to provide faculty with greater access to servers hosting confidential documents. Second, hosts in a VLAN belong to the same broadcast domain, and it is important to keep the cost of broadcast traffic small. The cost depends both on (i) the number of hosts in the VLAN, and (ii) the span of the VLAN, i.e., how spread out the hosts of the VLAN are in the underlying network topology. Finally, the total number of VLANs in the network must be kept limited, as the demand on network hardware grows with the number of VLANs. For instance, a separate spanning tree is typically constructed and maintained for every VLAN in the network, and this increases the memory and processing requirements of individual switches.

**(2) Placement of router and bridge:** For each VLAN with the host assignment decided, the operator must determine the best locations of the designated router, and the root bridge of the spanning tree. A key consideration is the potential inefficiencies in data communication with VLANs. Consider Fig. 1. Even though H1 and H2 are physically connected to the same switch, the path along which data flows is substantially longer. Having longer paths not only leads to longer delays, but also increases the likelihood of failures, and complicates performance and failure diagnosis. For example, if H1 and H2 were in building X, and R2 were in building Y, communication could be disrupted by a power failure in a building located between X and Y.

The inefficiencies of communication between H1 and H2 would be reduced if R1 were chosen as the designated router of VLAN 2 instead of R2. An ideal placement strategy must consider both the location of all the hosts in the VLAN, and the traffic patterns of the hosts. For instance, if hosts in a VLAN tend to communicate with certain servers, it is more critical to limit the performance inefficiencies associated with communication involving those servers.

The placement of root bridge directly impacts the spanning tree constructed for a VLAN. This in turn determines (i) the network links that see broadcast traffic of the VLAN, and (ii) the hops traversed when a host in the VLAN communicates with its gateway router. Thus, it is important to place the root-bridge judiciously to lower broadcast traffic in the network and reduce inefficiencies in data communication.

#### B. Reachability Control

From an operator's point of view, a primary objective of network security is to control packet level reachability, that is, what packets sent by a traffic source are permitted to reach a destination. Common security policies, such as restricting the types of external applications a host can access, limiting the scope of multicast traffic to specific subnets, and blocking unauthorized ICMP and SNMP probes, are essentially about permitting packets with particular header field combinations to be exchanged between hosts. Current design approaches are ad-hoc and error-prone, and current best practices for validating if a network configuration meets given reachability

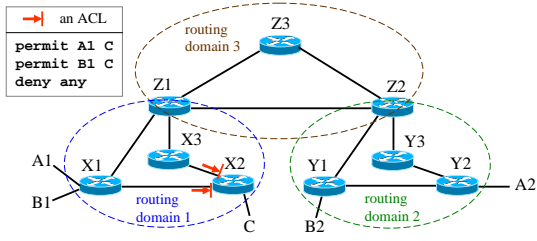


Fig. 2. Reachability control at data plane and control plane.

control objectives involve in-situ testing [4].

Today, operators realize reachability control objectives using two configuration options. The first is a data plane solution, which installs access control lists (ACLs), also commonly referred to as packet filters, on router interfaces. An ACL is a sequential collection of permit and deny conditions, called ACL rules. A packet’s header fields are matched against each rule successively. The order of rules is critical because testing stops with the *first match*. If no match is found, an implicit “deny any” rule is assumed and the packet is rejected.

The second approach to achieving reachability control objectives is a control plane solution. In particular, by either depriving some routers of certain routes, or creating black-hole routes in their forwarding tables, unwanted packets may be dropped by the routing logic. For example, one may partition a network into multiple routing domains and restrict the flow of routing information between the domains so that not all routers have routes to all destinations in the network.

Controlling reachability through the routing design has a much smaller CPU overhead because the execution of routing logic, particularly the lookup of the forwarding table, is mostly performed by forwarding hardware and requires little router CPU time. However, the routing-oriented solution is not always applicable because of its relatively limited range of conditions for matching packets. Unlike an ACL rule, which may simultaneously refer to multiple header fields, the routing logic matches packets either entirely based on source address or entirely on destination address.

Fig. 2 shows an example scenario where either configuration option can be used to meet a security policy. A1, A2, B1, B2, and C are subnets. Suppose the security policy does not permit any host in A2 and B2 to talk to C, but permits every host in A1 and B1 to talk to C. To realize this policy, the operator may configure an ACL, as shown in Fig. 2, in the inbound direction of both interfaces of router X2. Alternatively, the operator may block traffic between A2 and C, and between B2 and C, through routing design – one possible option is to install two source address based blackhole routes for traffic originated from A2 or B2 at router X2.

While routing design has been extensively studied (e.g., [7], [8], [9]), ACL placement has received little attention to date. In this paper, we focus on ACL placement. We assume that routing design is already completed, and routing domains are successfully configured before the operators proceed to determine the placement of ACLs in the network.

The key task with ACL placement is that operators need to construct a set of ACLs based on the security objectives and determine suitable locations, i.e., combinations of router

interface and traffic direction, to place them. In coming up with an ACL placement, the primary criterion is **correctness** of the design. The ACL and routing configurations must guarantee the delivery of all authorized packets while preventing all unauthorized traffic from reaching the destination. The solution should also be resilient to certain link or router failure scenarios - in particular, the alternate paths that may be taken when failures occur must also be correctly configured to ensure the reachability constraints are met.

Another consideration in ACL placement is the CPU overhead that routers incur from processing ACL rules *packet by packet*. There is a limit on the total number of ACL rules that a router can process consistently per packet. The limit varies from model to model. A low-end router may only be able to process dozens of ACL rules per packet without a noticeable reduction in link utilization. Therefore in some scenarios, it may be necessary to place ACLs throughout the network to distribute the computation cost. A recent study [1] reveals that some operational networks indeed have many ACLs placed at core routers, in addition to ACLs placed at access and distribution routers.

### III. SYSTEMATIC VLAN DESIGN

In this section, we present our approach for systematic VLAN design. We first describe the network-wide abstractions that we have developed to capture the most important factors of VLAN design. We then formulate the operator design tasks into optimization problems with general cost models. Finally, we present a set of heuristics for solving the optimization problems with particular cost models.

#### A. Network-Wide Abstractions

We model the VLAN design problem using the following abstractions:

- **Host Category:** This is a mapping  $\mathcal{P}$  that associates each host in the network with the logical category to which it belongs, such as engineering, sales, payroll, student cluster, faculty cluster, etc. While hosts in the same category need not belong to the same VLAN, *hosts in two different categories must belong to two different VLANs*. This is the correctness criterion for VLAN design.

- **Traffic Matrix:** A traffic matrix  $M_T$  which specifies expected traffic patterns between hosts in 2 different categories (or same category, or a given category and Internet). We assume information is provided about the *average* traffic between all host pairs in two categories. That is,  $M_T(i, j)$  specifies the average data traffic (in Kbps) sent by a host in category  $i$  to a host in category  $j$ . While a precise traffic matrix might be hard to obtain, we discuss in §III-D3 how to work with coarse traffic patterns if accurate information is unavailable.

#### B. Formulation of Operator Tasks

Given a complete network topology with hosts, switches, and routers, the goal of the operator is to put together a VLAN design with the above considerations. To make the problem more tractable, we model the VLAN design problem as a two-phase process:

**(i) Grouping hosts into VLANs:** The operator must decide the appropriate number of VLANs, denoted by  $x$ , in the design and which hosts must belong to each VLAN. More formally, the problem may be expressed as:

$$\begin{aligned} & \text{Minimize } [C(x) + \max_{1 \leq i \leq x} \{BroadcastCost_i\}] \\ & \text{subject to the correctness criterion defined by } \mathcal{P} \end{aligned}$$

Here,  $C(x)$  denotes the costs associated with having  $x$  VLANs in the design.  $BroadcastCost_i$  represents the cost of broadcast traffic associated with a given VLAN  $i$ .

**(ii) Placement of router and bridge:** For each created VLAN  $i$  with the host assignment decided, the operator wishes to determine the best location of the designated router, and the root of the spanning tree. The key objective is to minimize the combined costs of data traffic and broadcast traffic associated with the placement decisions. More formally, the operator task may be formulated as:

*Minimize TrafficCost, where*

$$TrafficCost = DataTrafficCost + \sum_i BroadcastCost_i \quad (1)$$

Here,  $DataTrafficCost$  represents the total cost of data traffic associated with all VLANs in the network, for a given design. In the future, it may be interesting to also constrain the number of VLANs that may be assigned to a given router, or root bridge.

Our formulation assumes that the two tasks are addressed sequentially to make the problem more tractable. In the future, it may be interesting to explore formulations that jointly optimize both design tasks.

### C. Phase 1: Grouping Hosts into VLANs

There are three key components in the design of a solver for grouping hosts into VLANs. These include (i) a model of the costs associated with a given number of VLANs; (ii) a model of the costs associated with broadcast traffic for a given VLAN; and (iii) an algorithm to realize the actual grouping. We present them in the rest of the section.

#### 1) Cost Models

**Costs associated with adding VLANs:** Our solver focuses on a particular cost function, where the manager specifies an acceptable bound on the total number of VLANs. In particular, if  $x$  VLANs are employed in the design, and  $MAX-VLANs$  is the maximum number of VLANs acceptable in the design (a constraint provided by the manager, and probably derived from the number of VLANs supported by the routers and switches being used), then:

$$\begin{aligned} C(x) &= 0, \text{ if } x \leq MAX-VLANs \\ C(x) &= \infty, \text{ if } x > MAX-VLANs \end{aligned}$$

We believe this is a natural cost function that is easy to express to the operator, and translates to many real-world design scenarios. While our current model may also be viewed as a feasibility criterion, it may be interesting to consider other kinds of cost functions in the future.

**Broadcast traffic costs:** Several applications may result in broadcast traffic in a network such as ARP, IPX, NetBIOS, SUNRPC, DHCP, and MS-SQL. We model the broadcast

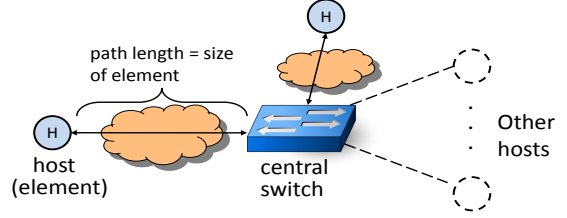


Fig. 3. Reduction from 3-partition problem to VLAN grouping problem.

traffic cost based on (i) the rate of broadcast traffic generated; and (ii) the number of links traversed as part of the broadcast. The links traversed by the broadcast traffic in a VLAN are simply the links present in the spanning tree for that VLAN. This may be easily generalized to a weighted sum of links, where weights are assigned to individual links to capture the cost of traversing that link.

In general, let  $N_i$  denote the number of hosts in VLAN  $i$ ,  $B_i$  denote the average broadcast traffic (in Kbps) generated by a host in VLAN  $i$ , and also let  $W_{ik}$  denote the number of links in the spanning tree for VLAN  $i$  when it chooses router  $S_k$  as its root bridge and designated router. (Note that our formulation focuses on the scenario where the root bridge and the designated router for a VLAN are always coupled. It is straightforward to extend our formulation to scenarios where the two are decoupled.) Then, we model the broadcast cost for VLAN  $i$  as

$$BroadcastCost_i = N_i \times B_i \times W_{ik} \quad (2)$$

We believe a linear dependence on the number of hosts in the network is a reasonable model. For instance, consider ARP queries, a key component of broadcast traffic. In typical scenarios, most ARP queries are sent by hosts in the VLAN for its designated router, or by the designated router for hosts in the VLAN, and a linear model fits well. Other models may be more appropriate in certain scenarios. For example, the entire IP address space of the VLAN may need to be considered for ARP broadcast storms due to port scans to non-existent hosts in the VLAN. As another example, a quadratic model is more appropriate if there is significant intra-VLAN ARP traffic. These scenarios are less typical, but we believe it is easy to extend our model to consider them.

Computing the number of links  $W_{ik}$  in the spanning tree of the VLAN depends on where the router and root bridge are located, which are themselves unknowns, and a degree of freedom the manager enjoys. When partitioning hosts into VLANs, our solver assumes the router and root bridge are placed in a manner that would result in the smallest number of links in the spanning tree. Thus, host grouping indicates the feasibility of keeping the broadcast costs small subject to appropriate router and bridge placement. The second phase of the solver (§III-D) determines router and bridge placement, with the broadcast traffic costs being one of the criteria.

## 2) Complexity of Problem

**Theorem III-1.** *The VLAN grouping problem is NP-hard with respect to the number of hosts to be grouped.*

*Proof:* Given our cost models, the problem of grouping hosts into VLANs involves minimizing the maximum broadcast cost across all VLANs subject to category constraints, with the broadcast cost defined as in Equation (2). We consider a decision version of the problem, where the goal is determine if a grouping exists such that all VLANs have a broadcast cost less than  $X$ , for a given  $X$ . We show this problem is NP-hard using a reduction from the well-known 3-partition problem, which is known to be NP-hard [10]. In the 3-partition problem, we are given a set  $A$  of  $3m$  elements, with each  $a \in A$  associated with an integer size  $s(a)$ . Further,  $\sum_{a \in A} s(a) = m * B$ , and  $\forall a \in A, B/4 < s(a) < B/2$ . The problem is to decide if the set of objects can be partitioned into  $m$  subsets such that the sum of the size of the objects in each subset is identical (or exactly  $B$ ). Note that since  $\forall a \in A, B/4 < s(a) < B/2$ , each subset is forced to consist of exactly three elements.

To show the reduction, we consider a special version of the VLAN grouping problem for each instance of the 3-partition problem. In particular, we consider a topology as shown in 3. In this topology there is a single central switch, and for each element in the 3-partition problem, we introduce a host which is connected to the central switch using a path of length equal to the size of the element. Further, all hosts are assumed to belong to the same category, and all hosts are assumed to produce the same amount of broadcast traffic corresponding to 1 unit (1Kbps). We note that for any VLAN involving two or more hosts in this topology, the spanning tree must consist of the central switch, and all switches on the path from the central switch to each of the hosts. Thus, the number of links in the spanning tree is simply the sum of the path lengths of each to the central switch, or equivalently, the sum of the sizes of the corresponding elements in the original 3-partition problem.

We claim that a feasible 3-partition exists in the original problem, if and only if the decision version of the VLAN grouping problem returns true for  $m$  permitted VLANs with a bound on broadcast cost of  $3 * B$ .

The proof of this claim has two parts:

- Lets assume a feasible 3-partition exists. Then for all the elements mapped to one subset, we take the corresponding hosts and group them in one VLAN. Since the sum of elements in each subset is exactly  $B$ , and each subset has exactly three elements, the number of spanning tree links in each VLAN is  $B$ , and each VLAN has 3 hosts. Hence, the broadcast traffic is  $3 * B$  for each VLAN, and the decision version of the VLAN grouping problem returns true for bound  $3 * B$ .
- Next, assume that the decision version of the VLAN grouping problem returns true for bound  $3 * B$ , that is we can group hosts into VLANs, so each VLAN has broadcast traffic at most  $3 * B$ . We first show all VLANs must have exactly 3 hosts. If this were not the case, there must be some VLAN with 4 or more hosts (as there are  $3m$  hosts to be partitioned into  $m$  VLANs). But, each host has a distance  $> B/4$  from

the central switch (as each element in original problem has size  $> B/4$ ). Hence, for this VLAN, the number of links in the spanning tree is  $> B$ . Since the number of hosts  $\geq 4$ , the broadcast cost for the VLAN is  $> 4B$  This is a contradiction to our assumption and is not feasible.

Given all  $m$  VLANs have exactly 3 hosts. Hence the number of links in the spanning tree of each VLAN must be  $\leq B$ , as the maximum broadcast cost is  $3 * B$  across all VLANs. But, the sum of the number of links in spanning trees of all VLANs must be equal to the sum of the sizes of all elements in original problem =  $m * B$ . This is only possible if the number of links in spanning tree of every VLAN is exactly  $B$ . This means an algorithm for solving the VLAN grouping problem can also be used to solve the 3-partition problem, where we create  $m$  subsets, with each subset corresponding to a VLAN, and elements in that subset corresponding to hosts in that VLAN. ■

## 3) Heuristic for Creating Host Groupings

Given the complexity of the problem and the scale of enterprise networks, it is impractical for any algorithm to find out the optimum grouping. Instead, our solver employs a greedy heuristic to determine grouping of hosts into VLANs. Initially, each category of hosts provided by the operator is assumed to constitute one VLAN. The solver then computes the minimum broadcast traffic costs for each VLAN. The VLAN with the largest broadcast traffic cost is taken, and is split into two VLANs if the total number of VLANs in the design is no more than  $MAX-VLANs$ . The process continues iteratively until the condition is violated.

When a VLAN  $i$  is chosen to be split, then, the goal is to split it in a manner that hosts close to one another in the underlying topology are placed in one VLAN to minimize the span. The solver employs the following 2-step algorithm:

- (i) For each host  $k$  in VLAN  $i$ ,  $H_{i,k}$ , we compute the shortest distances from  $H_{i,k}$  to all  $N_i$  hosts in VLAN  $i$ , including itself, to form a vector  $\{d(H_{i,k}, H_{i,h}) | h = 1..N_i\}$  of  $N_i$  values, where  $d(H_{i,k}, H_{i,h})$  denotes the shortest distance (i.e., number of layer-2 hops) from host  $k$  to host  $h$  in VLAN  $i$ .
- (ii) Using the vector of a host as its coordinate (or location) in the topology, we perform the k-means algorithm to cluster all hosts in VLAN  $i$  into two separate VLANs.

## D. Phase 2: Router and Bridge Placement

Once the solver groups hosts into VLANs, it then determines the recommended placement of the designated router, and the root bridge, for each VLAN. In doing so, the key objective is minimizing the combined costs of data and broadcast traffic. The broadcast traffic cost was formulated in Equation (2). In the rest of the section, we present a model for capturing data traffic communication costs and then describe the placement heuristics. We focus on the scenario where the designated router and the root bridge for a VLAN are always coupled. This is often the case in enterprise networks, as it simplifies the management tasks. We believe that our models and heuristics can be easily extended to scenarios where the two are decoupled.

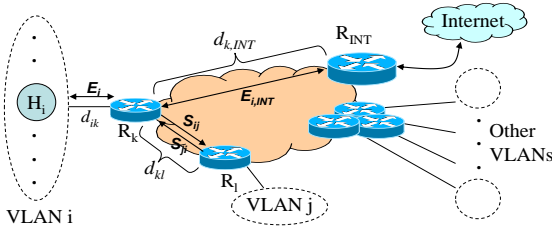


Fig. 4. Inter-VLAN traffic sent by a host in VLAN  $i$ .

### 1) Data Traffic Cost Model

The cost of data traffic communication depends on two factors (i) the amount of data traffic exchanged between a pair of hosts; and (ii) the number of hops (switches and routers) traversed as part of the communication. In modeling the data traffic, we separately consider the inter-VLAN and intra-VLAN traffic of every VLAN  $i$ . Thus,

$$\text{DataTrafficCost} = \text{InterVLAN} + \text{IntraVLAN} \quad (3)$$

- **Inter-VLAN traffic:** To model the costs associated with inter-VLAN traffic involving VLAN  $i$ , consider Fig. 4.  $H_i$  is a host in VLAN  $i$  that has designated router  $R_k$ . All inter-VLAN traffic sent, or received by  $H_i$  must traverse the path between  $H_i$  and router  $R_k$ . In addition, the portion of the traffic exchanged with a given VLAN  $j$  must traverse the path between  $R_k$  and  $R_l$ , where  $R_l$  is the designated router of VLAN  $j$ . Finally, the portion of the traffic exchanged with the Internet must traverse the path between  $R_k$  and  $R_{INT}$ , where  $R_{INT}$  is the gateway router to the Internet. Consider the following notations: (For clarity, in this section we always use  $i$  and  $j$  to refer to VLANs, and use  $k$  and  $l$  to refer to routers.)

- $d_{ik}$ : the number of hops between a host in VLAN  $i$ , and a router  $R_k$ , averaged across all hosts in VLAN  $i$ .

- $d_{kl}$ : the number of hops on the path between routers  $R_k$  and  $R_l$ .

- $d_{k,INT}$ : the number of hops between router  $R_k$  and the gateway router to the Internet (i.e.,  $R_{INT}$ ).

- $S_{ij}$ : the amount of traffic that VLAN  $i$  sends to VLAN  $j$ .

- $E_{i,INT}$ : the amount of traffic that VLAN  $i$  exchanges with (i.e., sends to or receives from) the Internet.

- $E_i$ : the total amount of inter-VLAN traffic associated with VLAN  $i$ . Note that the inter-VLAN traffic associated with VLAN  $i$  consists of the traffic  $i$  sends to all the other VLANs, the traffic all the other VLANs send to  $i$ , and the traffic  $i$  exchanges with the Internet, i.e.,  $\sum_j S_{ij} + \sum_j S_{ji} + E_{i,INT} = E_i$ . These values can be easily obtained from  $\mathbf{M}_T$ .

- $x_{ik}$ : 1 if VLAN  $i$  chooses router  $R_k$  as its designated router; 0 otherwise. Note that:  $\sum_k x_{ik} = 1$

Then, the total inter-VLAN traffic costs *InterVLAN* can be formulated as follows:

$$\begin{aligned} & \sum_i \sum_j \sum_k \sum_l S_{ij} d_{kl} x_{ik} x_{jl} \\ & + \sum_i \sum_k E_i d_{ik} x_{ik} + \sum_i \sum_k E_{i,INT} d_{k,INT} x_{ik} \end{aligned} \quad (4)$$

Intuitively, the first term of (4) represents the inter-VLAN traffic that is routed from a VLAN  $i$ 's designated router  $R_k$  to another VLAN  $j$ 's designated router  $R_l$ . This includes all the traffic  $i$  sends to  $j$ . (Note that the traffic  $j$  sends to  $i$  is counted separately.) This term is summed over every pair of VLANs, and for each pair  $i$  and  $j$ , traffic  $i$  sends to  $j$  and traffic  $j$  sends to  $i$  are both counted. The second term of (4) represents the inter-VLAN traffic that traverses the path between a VLAN  $i$  and its designated router  $R_k$ . All inter-VLAN traffic sent or received by VLAN  $i$  must traverse this path. This term is summed over all VLANs. The last term of (4) represents the traffic that traverses the path between a VLAN  $i$ 's designated router  $R_k$ , and the gateway router to the Internet. Any traffic  $i$  exchanged with the Internet must traverse this path. This term is summed over all VLANs.

- **Intra-VLAN traffic:** When two hosts in the same VLAN communicate, the number of hops between them depends on the spanning tree of that VLAN, and is bounded by two times the total number of hops between each host and the root bridge of that VLAN. If router  $R_k$  is the root bridge for VLAN  $i$ , then the average number of hops between a host in VLAN  $i$  and  $R_k$  is  $d_{ik}$ . Hence the average number of hops traversed by intra-VLAN traffic is at most  $2d_{ik}$ . Let  $I_i$  denote the intra-VLAN traffic (in Kbps) associated with VLAN  $i$ , the total intra-VLAN traffic costs  $\sum_i \text{IntraVLAN}_i$  can be formulated as follows:

$$\sum_i \sum_k I_i 2d_{ik} x_{ik} \quad (5)$$

Now that we have formulated data traffic cost, we can formulate the total traffic cost as shown in (1). According to (1), to minimize total traffic cost, we must jointly minimize both data and broadcast traffic costs:

$$\begin{aligned} \min & \sum_i \sum_j \sum_k \sum_l S_{ij} d_{kl} x_{ik} x_{jl} \\ & + \sum_i \sum_k (E_i d_{ik} + E_{i,INT} d_{k,INT} + 2I_i d_{ik} + N_i B_i W_{ik}) x_{ik} \\ \text{s.t.} & \sum_k x_{ik} = 1, \quad i = 1, 2, \dots, n \end{aligned} \quad (6)$$

Here the term  $N_i B_i W_{ik}$  represents the broadcast traffic cost associated with VLAN  $i$  when it chooses router  $R_k$  as its root-bridge and designated router, which has been formulated in (2).

### 2) Complexity

**Theorem III-2.** *The router and bridge placement problem is NP-hard with respect to the number of routers and root bridges to choose from.*

*Proof:* The above router and bridge placement problem falls into a category of nonlinear assignment problems, namely quadratic semi-assignment problems (QSAP) [11]. QSAP models the problem of allocating a set of  $n$  facilities to a set of  $m$  locations, with the costs being the cumulative product of flow between any two facilities and the distance between any two locations, plus the costs associated with a facility being placed at a certain location. The objective is to assign each facility to a location such that the total cost is

minimized. QSAP is a variant of the well known quadratic assignment problem (QAP) [11]. The only difference between QSAP and QAP is that in the former each location may take none, one or more than one facilities, whereas in QAP each location has to obtain exactly one facility, and vice versa. Both problems are known to be NP-hard [11], [12].

Formally, we are given three matrices with real elements  $F = (f_{ij})$ ,  $D = (d_{kl})$  and  $B = (b_{ik})$ , where  $f_{ij}$  is the flow between facility  $i$  and facility  $j$ ,  $d_{kl}$  is the distance between location  $k$  and location  $l$ , and  $b_{ik}$  is the cost of placing facility  $i$  at location  $k$ . Note that  $F$  and  $D$  matrices can be either symmetric or not. The QSAP can be formulated as follows:

$$\begin{aligned} \min \quad & \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^m \sum_{l=1}^m f_{ij} d_{kl} x_{ik} x_{jl} + \sum_{i=1}^n \sum_{k=1}^m b_{ik} x_{ik} \\ \text{s.t.} \quad & \sum_j x_{ij} = 1, i = 1, 2, \dots, n \\ & x_{ij} \in \{0, 1\}, i = 1, 2, \dots, n, j = 1, 2, \dots, m. \end{aligned} \quad (7)$$

It is easy to see that our placement problem has the same structure of QSAP. Consider every VLAN as a facility, and every router as a location. Also consider the amount of traffic between VLANs in our problem to be the flow between facilities in QSAP, and the number of hops between routers in our problem to be the distance between locations. Then the first term of (6) can be viewed as the cost of moving flows between different facilities, i.e., the first term of (7), and the second term of (6) can be viewed as the cost of placing facilities at certain locations, i.e., the second term of (7). Hence our router and root bridge placement problem can be formulated as QSAP, and thus is NP-hard. ■

### 3) Heuristic for Router and Bridge Placement

Given the complexity of the problem and the scale of enterprise networks, it is practically impossible for any algorithm to find out the optimum placement. Further, obtaining an accurate estimate of  $M_T$  might be difficult, especially for a network that is yet in operation. We instead design a heuristic which is guided by observations of typical traffic patterns in enterprises. Many enterprises today dedicate a small number of VLANs to house important server machines, such as network file servers, DNS and DHCP servers. These VLANs are likely to be extremely popular in that most hosts in the enterprise communicate with these VLANs. For the vast majority of other non-server VLANs, however, most traffic exchanged is with the server VLANs, and with the Internet. We refer to these non-server VLANs as client VLANs.

Our solver requires an operator to indicate the set of server VLANs in the design. For every client VLAN, information is provided regarding what fraction of its traffic is exchanged with the Internet, and each server VLAN. If this information is unavailable to operators, it is assumed an equal amount of traffic is exchanged with each of the server VLANs.

Consider the terms in Equations (2), (4), and (5). The costs associated with broadcast and intra-VLAN traffic depend entirely on the placement choices of router and root bridge associated with that VLAN alone. The cost associated with

inter-VLAN traffic however has components that depend on the placement choices of other VLANs. The extent of this dependency on remote VLAN placement is likely higher if there is a strong bias in traffic to the remote VLAN.

The solver proceeds in two steps:

(i) Placement decisions are made for all server VLANs. In doing so, terms dependent on placement decisions of other VLANs are not considered.

(ii) The optimization is conducted for all client VLANs. Given that they primarily communicate with server VLANs, terms involving placement decisions of server VLANs alone are considered, and terms involving placement decisions of other client VLANs are neglected.

With this approach, solving each step above requires minimizing  $TrafficCost_i$  (i.e., sum of Equations (2) and(3)) individually for each VLAN, with the only unknowns being the router and bridge choices for that VLAN. A simple iterative algorithm that tries all possible choices of network elements as designated router and root bridge suffices to ensure the best placement can be found for each VLAN.

## IV. SYSTEMATIC REACHABILITY CONTROL

In this section, we present our approach for systematic reachability control. We first describe the network-wide abstractions that we have developed to capture the ultimate requirements of reachability control. Next, we formulate the task of ACL placement into a set of optimization problems, each fashioning a different design strategy. We then show that finding the optimal placement is an NP-hard problem. Finally, we present greedy heuristics to approach these optimization problems.

### A. Network-Wide Abstractions

We consider the *Reachability Set (RS)* between two points in a network to be the subset of packets (from the universe of all IP packets) that the network may carry between those points. The RS notation has been shown to provide a unifying metric for determining the joint effect of packet filters and routing protocols on end-to-end reachability [4]. The RS metric provides the required building block towards a network-wide abstraction that can *completely* capture the operator intent in regard to reachability control. In addition, a network's reachability control policy is said to be resilient against an event if the network continues to uphold the reachability policy despite the occurrence of the event. We model the reachability requirement and the resiliency requirement of a reachability control policy at the granularity of VLANs (or subnets in general) using the following abstractions:

•**Reachability Matrix:** Consider a network with  $N$  VLANs. The network's reachability policy can be completely described by an  $N$  by  $N$  reachability matrix, denoted by  $M_R$ , where element  $M_R(i, j)$  denotes the maximum RS that will always reach an intended destination host in VLAN  $j$  if originated by a host of VLAN  $i$ .

•**Managed Event Set:** The resilience requirement of a network's reachability control policy can be completely described by a managed event set, denoted by  $E_m$ , with each element in the set specifying a topology-changing event to which the

network must respond without causing the reachability matrix to change.

### B. Formulation of Operator Tasks

The primary task of the operator is to place ACLs in a manner that meets the correctness and feasibility criteria below:

**(i) Correctness Criterion:** The network’s reachability matrix is invariant and as specified in  $\mathbf{M}_R$  under all events in  $\mathbf{E}_m$ .

**(ii) Feasibility Criterion:** Let  $c(r)$  represent the limit on the total number of ACL rules that can be configured on a router  $r$ , including all its interfaces and in both traffic directions, without overloading  $r$ . Let  $b(r)$  be the number of ACL rules that has been configured on router  $r$ . Then,  $\forall r, b(r) \leq c(r)$ .

In some network topologies, it may be possible to have multiple ACL placement strategies that meet the correctness and feasibility criteria. For instance, consider a cell of the reachability matrix,  $\mathbf{M}_R(i, j)$ . Consider the simplest case where only a single path of routers exists from VLAN  $i$  to VLAN  $j$ . The operator may place an ACL permitting only  $\mathbf{M}_R(i, j)$  at any of the routers to meet the criteria. We leverage this potential flexibility to allow operators to express their preference for an ACL placement design. In this paper, we consider the following four ACL placement strategies:

**Minimum Rules (MIN) Strategy.** The operator wishes to minimize the total number of filter rules installed on all routers in the network. More formally:

$$\text{Minimize } \sum_r b(r)$$

**Load Balancing (LB) Strategy:** The operator wishes to spread the ACL processing overhead across the network in order to avoid overburdening any router. Formally:

$$\text{Minimize } \max_r \{b(r)\}$$

The configuration derived from this strategy will not impose a need for costly super nodes. However, the operator may intentionally set  $c(r)$  to  $\infty$  when designing a new network (with no hardware purchased yet) or when it is feasible to upgrade existing router hardware.

**Capability Based (CB) Strategy:** The operator wishes to allocate the ACL processing overhead based on each router’s filtering capability. Formally:

$$\text{Maximize } \min_r \{c(r) - b(r)\}$$

Using this strategy, the derived configuration squeezes the most out of the capability of the current hardware.

**Security Centric (SEC) Strategy:** The operator wishes to minimize the security risk posed by unwanted traffic permitted in the network, by placing filters as close to the source as possible. For a filter  $f$ , let  $h(f)$  represent the hop count from the router on which  $f$  is installed to the gateway router of the traffic sources targeted by  $f$ , averaged across all traffic sources. Let  $H$  be the average  $h(f)$ , averaged across all filter rules installed in the network. Ideally,  $H$  should be 0. Formally, the goal of the strategy is:

$$\text{Minimize } H$$

### C. Complexity of ACL Placement

We model the problem of placing ACLs for the entire reachability matrix  $\mathbf{M}_R$  as processing each cell  $\mathbf{M}_R(i, j)$  of

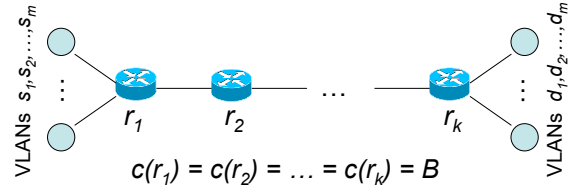


Fig. 5. Network setting used in reduction of the “bin packing” problem.

the matrix one by one until the reachability requirements of all cells are satisfied. The processing for each cell involves finding a correct and feasible placement to install the ACL that describes the reachability control requirement of that cell. Note that if a cell  $\mathbf{M}_R(i, j)$  contains “full-reachability” (i.e., any packet from VLAN  $i$  can reach VLAN  $j$ ), the processing for that cell is skipped since no ACL is required. The following theorem establishes that finding the optimal solution to the ACL placement problem is NP-hard.

**Theorem IV-1.** *The ACL placement problem is NP-hard with respect to the number of cells to be processed.*

*Proof:* We present a reduction of the well-known NP-complete “bin packing” decision problem [13] into the problem of ACL placement with the MIN strategy. The reduction holds for the other strategies too because they share the same decision problem as the MIN strategy.

The “bin packing” decision problem can be formally stated as follows. Given (i) a finite set  $U$  of  $m$  items, with each  $u \in U$  having a positive integer size  $s(u)$ , and (ii) positive integers  $B$  (called the bin capacity) and  $k \leq m$ , can  $U$  be partitioned into  $k$  disjoint sets  $U_1, \dots, U_k$  such that for each  $U_i$  the total sum of the sizes of the items in  $U_i$  does not exceed  $B$ ?

Next, we reduce this general problem to the question of whether it’s feasible to place ACLs for the special network setting illustrated in Fig. 5. First, we map each of the  $k$  bins into a router with  $c(r) = B$ . The routers form a linear topology that connects two groups of VLANs at the two ends, each with  $m$  VLANs. We then map each item  $u_i \in U; i = 1, 2, \dots, m$  to cell  $(s_i, d_i)$  of the network’s reachability matrix, i.e., one that affects packets originating from VLAN  $s_i$  on the left side and going to VLAN  $d_i$  on the right, such that the number of ACL rules required for that cell is  $s(u_i)$ . Finally, we set all the unmapped cells in the reachability matrix to “full reachability”, i.e., requiring no packet filter. Clearly, the answer to the “bin packing” problem is yes if and only if it’s feasible to place the ACLs for the network setting considered since the ACL for each cell must be placed in one of the routers with sufficient remaining capacity. ■

### D. Heuristics for ACL Placement

Since the ACL placement problem is NP-hard, we begin this section by presenting heuristics for processing individual cells (i.e.,  $\mathbf{M}_R(i, j)$ ) of the reachability matrix. These fine-grained heuristics provide insights on how our solvers ensure the correctness of placement and approximate various placement strategies. We then discuss placement strategies that involve processing  $\mathbf{M}_R$  one row or one column at a time.



### 1) Placement by Cell

Several polynomial-time heuristics exist for approximating an optimal solution to the “bin-packing” problem. Among them, the “first fit decreasing” strategy, whereby the items are first sorted from largest to smallest and then sequentially placed in the first feasible bin, strikes a good balance between the optimality of the solution and the time complexity. We have adopted the same strategy for ACL placement given a strong resemblance between the two problems. In particular, we first sort the cells in the decreasing order of the number of ACL rules they contain and then process them sequentially using the greedy per-cell placement heuristics presented in the remainder of the section.

To process a given cell  $\mathbf{M}_R(i, j)$  of the reachability matrix, we assume that the routing design stage is already completed so that a subgraph  $g(i, j)$  of the layer-3 network topology can be derived from the routing design which contains VLANs  $i$  and  $j$ , and satisfies the following conditions:

- The subgraph is sufficiently connected so that no event in  $\mathbf{E}_m$  will disconnect VLAN  $i$  from VLAN  $j$ . That is, we assume that the resilience is ensured by the routing design.
- For each path from VLAN  $i$  to VLAN  $j$  in the subgraph, either it is one of the default forwarding paths from VLAN  $i$  to VLAN  $j$  or there exists an event in  $\mathbf{E}_m$  under which it will be used to route traffic from VLAN  $i$  to VLAN  $j$ .

We note that obtaining  $g(i, j)$  may be nontrivial for some of the existing networks where route filters and route redistributions are configured in an ad-hoc fashion [2]. Here we assume that routing design has been accomplished systematically to ensure the predictability of  $g(i, j)$ . We also note that overestimating  $g(i, j)$ , i.e., including more nodes and edges than necessary, does not affect the correctness of the placement although the resulting solution may place more filter rules than necessary.

The foremost concern of reachability control is the correctness of the solution. The heuristics for all four optimization strategies use the same approach to ensure correctness. They guarantee that the ACL for each cell is placed along all members of an  $(i, j)$  edge-cut-set in  $g(i, j)$ . In other words, all packets that go from VLAN  $i$  to VLAN  $j$  will encounter an instance of the ACL no matter which physical path they take.

We assume that the address spaces of different VLANs don’t overlap and that an algorithm exists to convert  $\mathbf{M}_R(i, j)$  into a sequential set  $f(i, j)$  of ACL rules. If VLAN  $i$  and VLAN  $j$  are respectively assigned address blocks of  $A$  and  $B$ , each rule in  $f(i, j)$  looks like the following.

```
{permit or deny} a b [more fields]
```

where  $a \subseteq A$  and  $b \subseteq B$ . In addition, to avoid ambiguity,  $f(i, j)$  must end with

```
deny A B
```

Finally, the heuristics require a post-processing step be performed after the entire reachability matrix is processed. The post-processing step overrides the implicit “deny any” on each interface by adding an explicit “permit any” at the end of all rules placed on that interface. In addition,

*Input:* (1) Topology  $g(i, j) = (V, E)$  where nodes in  $V$  may be VLAN  $i$ , VLAN  $j$ , or intermediate routers and subnets connecting  $i$  and  $j$ . The set of all routers in  $V$  is denoted by  $R$ . (2) Sequential ACL rule set  $f(i, j)$  with  $n(i, j)$  members. *Output:* Set of 2-tuple  $D$ , where  $D[0]$  is a router interface and  $D[1]$  takes a value of either 0 or 1, representing the direction of the ACL with respect to traffic – 0 means inbound and 1 means outbound.

- 1: Label all routers with insufficient filter capacity left, i.e.,  $c(r) - b(r) < n(i, j)$  as ineligible for inclusion into  $S$ .
- 2: Sort  $R$  into array based on increasing  $b(r)$  values; i.e.,  $b(R[0]) \leq b(R[1]) \leq \dots$ ; choosing minimum router hop count from  $i$  or  $j$  as tie breaker
- 3:  $S = \emptyset$ ;
- 4: **for**  $k = 0$  to  $\|R\| - 1$  **do**
- 5:   Add  $R[k]$  to  $S$ ;
- 6:   Try finding the smallest edge-cut-set between  $i$  and  $j$  using only edges connecting a node in  $S$ ;
- 7:   **if** successful **then**
- 8:     {denote the minimum cut-set by  $CUT$ }
- 9:     **for** each edge  $e \in CUT$  **do**
- 10:      **if** both ends of  $e$  are routers **then**
- 11:        **if** starting end of  $e$  has smaller  $b(r)$  **then**
- 12:          Add (starting end, 1) to  $D$ ;
- 13:        **else**
- 14:          Add (the other end, 0) to  $D$ ;
- 15:        **end if**
- 16:      **else if** starting end of  $e$  is a router **then**
- 17:        Add (starting end, 1) to  $D$ ;
- 18:      **else if** ending end of  $e$  is a router **then**
- 19:        Add (ending end, 0) to  $D$ ;
- 20:      **end if**
- 21:    **end for**
- 22:    return  $D$ ;
- 23:   **end if**
- 24: **end for**

Fig. 6. ACL placement solver for the LB strategy.

the post-processing step may optionally apply compression algorithms [14], [15] to further reduce the number of rules placed on each interface.

Fig. 6 presents the algorithm for the LB strategy. Initially, routers with insufficient capacity to accept  $f(i, j)$  are eliminated. The remaining routers are sorted in ascending order of  $b(r)$ . The number of router hops from either the source or destination VLAN is used as the tie breaker because it is more likely to find small edge-cut-sets closer to the network edge which is generally less connected than the middle of the topology. The first  $k$  routers in the sorted list are considered in set  $S$ . The algorithm iterates over  $k$  until a minimum edge-cut-set between VLAN  $i$  and VLAN  $j$  can be found using only edges connecting a node in  $S$ . The remaining steps of the algorithm (line 8 onwards) identify the appropriate router interfaces on which the filters must be applied. The algorithm can be implemented in polynomial time with well known efficient polynomial algorithms for finding the minimum edge-cut-set in a network [16].

The heuristics for the other strategies follow the same

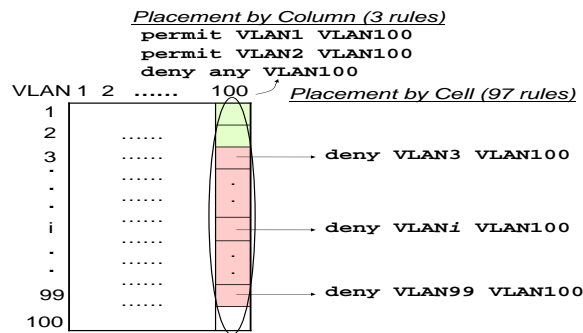


Fig. 7. Hypothetical reachability matrix highlighting the difference between fine-grained and column-based placement.

algorithm with minor variations. The CB strategy simply involves changing the sorting criterion in line 2 from “increasing  $b(r)$  values” to “decreasing  $(c(r) - b(r))$  values” while keeping the same tie breaker. The SEC strategy involves changing the sorting criterion to “increasing hop count from the gateway router of VLAN  $i$ ” and changing the tie breaker to “decreasing  $(c(r) - b(r))$  values”. Finally, the MIN strategy involves replacing lines 2-5 by including all routers in  $S$ , and then finding the minimum edge-cut-set.

## 2) Placement by Row or Column

Our discussion so far assumes a *fine-grained* strategy, where each cell of the reachability matrix is placed independently of other cells. Another degree of freedom for a placement scheme involves placing an entire row or column of the reachability matrix. For instance, security policies such as server access control by nature restrict traffic *to* one VLAN from all other VLANs. For such policies, one strategy is to place the entire column of the reachability matrix corresponding to the destination VLAN. Likewise, security policies like ingress filtering or blocking of unauthorized email servers by nature restrict traffic *from* one VLAN to all other VLANs. In such cases, a potential strategy is to place the entire row of the reachability matrix corresponding to the source VLAN. Note that placement by row/column does not reduce the inherent complexity of finding the optimal solution to the ACL placement problem, which can be shown to remain NP-hard using a similar proof as in §IV-C.

Placement by row/column offers interesting trade-offs compared to a fine-grained placement strategy. On the one hand, a fine-grained strategy may distribute rules over multiple routers, and require fewer rules on any given router than placement by row/column. In fact, in some scenarios, placement by row/column may not be feasible as the capacity of the router may be exceeded. On the other hand, placement by row/column may offer opportunities to compress the number of rules to be placed by using the wildcard “any” to represent any source or destination. For instance, Fig. 7 shows the reachability matrix for a hypothetical scenario where all hosts in VLANs 1 and 2 have full reachability to VLAN 100 (so no ACL rules are required for the corresponding cells), but all hosts in VLANs 3-99 are denied access to VLAN 100. If cells in the entire column for VLAN 100 are placed together, only 3 rules are required, as the deny rules from every other source VLAN 3 to VLAN 99 can be effectively compressed using the wildcard “any”. However, if a fine-grained strategy

is used, potentially 97 rules in total are required to be placed individually, and the rules may be distributed across many routers.

The algorithm in Fig. 6 can be easily extended to process one row or one column of the reachability matrix at a time. The key change is that the target edge-cut-set at line 6 needs to be enlarged to disconnect one source VLAN from many destination VLANs for row-based placement, or one destination VLAN from all source VLANs for column-based placement. Alternatively, the reachability matrix could be processed using a hybrid approach, where some entries are processed by row/column, and others are placed using a fine-grained approach. We omit further details for lack of space.

## V. EVALUATIONS AND VALIDATION

We evaluate our heuristics on a large-scale campus network with tens of thousands of hosts. The network consists of about 200 routers, 1300 switches, and hundreds of VLANs. Four routers form the core of the network. Typically, each building has a router with a link to one of the core routers. This link connects all hosts in the building to the rest of the network. Our data includes configuration files of all switches and routers, and the physical topology of the network.

**VLAN Usage:** While the campus IT operators provide routing services for the entire campus, each logical group such as the School of Engineering, the School of Liberal Arts, and the Libraries has its own administrators. Each administrative unit is given an IP address block and is free to assign addresses within that block to individual hosts. The operator policy requires that hosts in different administrative units must belong to different VLANs. VLANs are extensively used to meet this goal, as well as to constrain the size of broadcast domains. Most VLANs span a small section of the campus - about 50% of them span only one building. However, about 10% of the VLANs span 5+ buildings, and the largest VLAN spans over 60 buildings. VLANs with a large span correspond to administrative units that have hosts in most buildings on campus, e.g., hosts in all classrooms are administered together and are grouped into a VLAN.

**ACL Usage:** Prominent ACL policies used by the campus network include (i) ingress filtering to ensure that packets have a source IP address from the address space of their originating subnets; (ii) restricting communication involving dormitory hosts; (iii) restrictions involving wireless traffic; and (iv) restricting communication with data centers that house many key servers. Overall, ACL rules are placed in over 70 routers, with about 20% of the routers having 300+ rules, which may include rules from multiple ACLs.

### A. VLAN Design

In this section, we present results evaluating our systematic design approach for each of the VLAN design tasks.

**Grouping Hosts into VLANs:** With help from the operators, we categorize the hosts on a large segment of the campus. Each category corresponds to a different administrative unit. In total, there are 119 categories and 15084 hosts. Many categories are small, and the median category has only 79 hosts. However, the largest category includes 2000+ hosts.

	# Hosts per VLAN (182 VLANs)	
	Current	Systematic
Mean	82.9	82.9
Std Dev	71.9	57.1
90%ile	193	167
Max	254	195

TABLE I  
NUMBER OF HOSTS PER  
VLAN WITH THE  
CURRENT AND THE  
SYSTEMATIC DESIGNS.

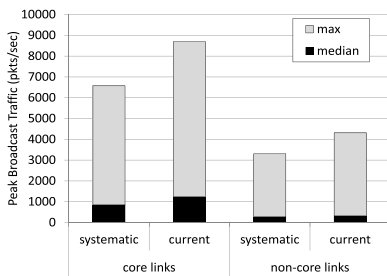


Fig. 8. Estimated peak broadcast traffic load per link.

We group hosts into VLANs using our systematic approach. Our algorithms are subject to two constraints. First, a maximum of 182 VLANs is permitted, as this is the number of VLANs used in the current design. Second, hosts from different categories are required to belong to different VLANs.

Table I shows the number of hosts per VLAN produced by our approach and compares the results to the current design. The results show the effectiveness of our approach in avoiding the creation of large VLANs with many hosts. The maximum number of hosts in any VLAN is reduced from 254 to 195, and the 90%ile is reduced from 193 to 167. This is achieved by a more equitable distribution of hosts across VLANs as indicated by the lower standard deviation. In addition, we also found (though not shown in the table) that our systematic approach also reduces the span of large VLANs by decreasing the number of links in their spanning trees. In particular, the maximum number of spanning tree links in any VLAN is reduced from 417 to 254.

We next study the potential benefit of our systematic grouping in reducing broadcast traffic, which is usually dominated by VLANs with a large size and span. To get a realistic estimate of broadcast traffic pattern, we measured the broadcast traffic sent by hosts in one of the VLANs over a 24-hour period. We observed an average and peak packet rate of 0.004 pkt/s/source and 2.12 pkt/s/source, respectively. We then estimated the peak broadcast traffic seen per link, assuming every host generates broadcast traffic at the peak rate.

Fig. 8 shows the median and maximum estimated peak broadcast packet rates per network link for the current grouping and our systematic grouping. Two types of links, *core links* and *non-core links*, are shown. The core links include links between core routers, and links connecting a core router to routers of various buildings in campus. All the remaining links are non-core links. Overall, there are about 500 core links and 41000 non-core links. Our systematic design results in similar median broadcast traffic to the current design, but significantly reduces the maximum broadcast traffic rate by around 1000 pkts/sec and 2000 pkts/sec for non-core links and core-links, respectively. The decrease of broadcast traffic in core links comes from both reducing the number of hosts in large VLANs as well as ensuring VLANs span as few links as possible. The drop in broadcast packet rate on core links allows core routers to potentially save their processing power for more important tasks, e.g., assuring critical traffic is quickly transported through the backbone.

#### Router and Bridge Placement:

The operators provided a set of six server VLANs that

housed many of the popular servers that other hosts would access. These include servers like campus web servers, DNS and DHCP servers, and other important data servers. The operators also confirmed that a large portion of traffic from the other VLANs (client VLANs) is either exchanged with these server VLANs, or with the Internet. We then compute the optimal placement of their routers using our algorithm in §III-D. We assume router and bridge placement are coupled, given this is true of the current design, and given the operator preference for such a choice. In addition, we assume that intra-VLAN data traffic is negligible, and 1% of inter-VLAN data traffic incurs broadcast traffic. Among the remaining 99% of inter-VLAN data traffic,  $f\%$  is exchanged with the Internet, and the rest is exchanged evenly with each server VLAN. We believe these models are realistic in many enterprise settings, and the operators confirmed these are reasonable traffic models.

Fig. 9 explores the effectiveness of our systematic router placement in reducing the number of hops traversed by data traffic when  $f$  is varied. There are two bars for each choice of  $f$ , one for the current placement and the other for our systematic placement. Each bar represents the 90%ile of the average weighted hop count for hosts in a client VLAN. The weighted hop count is the average number of hops from a client host to the gateway routers of the server VLANs and the Internet, weighted by the corresponding fraction of data traffic exchanged with them. For all scenarios, the average weighted hop count is decreased by 1-1.5 hops using our systematic placement, since our systematic approach takes traffic patterns into account. Reducing the number of hops traversed by data traffic not only results in lower delays, but also reduces the possibility of communication being disrupted by failures. Further, the data traffic carried by network links could also be reduced.

We next study the potential benefit of our systematic placement in reducing data traffic on network links. To model the traffic behavior of end hosts, we consider two models: a uniform model and a trace model. The uniform model assumes every host transmits data uniformly at 10Kbps. The trace model is based on traffic traces collected at LBNL [17]. The traces were recorded over a 22-hour period in December 2004, covering about 8000 internal addresses. We computed a list of average data rate sent/received by each internal address, which ranges from 0-8183Kbps with a mean of 14.6Kbps. We then randomly assigned a rate from this list to each host in our campus network and evaluated the traffic load on each link. Fig. 10 shows the median and 95%ile traffic load on the core links using both traffic models under the current and systematic designs. While the median core link load is similar for both designs using the two traffic models, our systematic placement improves the 95%ile load from 20.9Mbps to 6.4Mbps and from 27Mbps to 12.1Mbps for the uniform model and the trace model, respectively. The results show that shorter data paths may involve traversal of fewer core links, and the potential reductions in data traffic on these core links is significant.

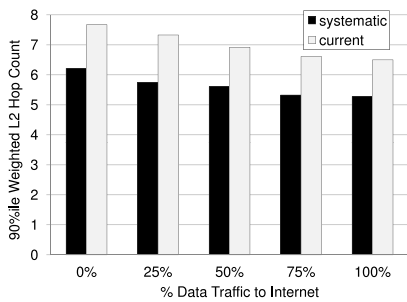


Fig. 9. Reduction of hops traversed by data traffic using our systematic router placement, with varying  $f$  (i.e., fraction of traffic to the Internet).

### B. Placement of ACL rules

The campus network we analyzed is well-run, and many hours of design time have been spent on its ACL rules. Using our systematic design algorithms, we were able to automatically create an ACL placement that mostly matches the current placements in this large-scale network using only an hour of CPU time. Beyond the general time savings in creating placements and adapting them as the network changes, we found two interesting examples that illustrate the importance and benefits of systematic placement of ACL rules.

**Correctness of Placement:** Our analysis discovered an inconsistency between operator intent and the current ACL placement. One operator policy is to prevent access from unregistered dormitory users to any host other than a small number of well-known registration servers. Fig. 11(a) illustrates the relevant segment of the network. Hosts in the dormitories are separated into a group of VLANs. These VLANs share the same gateway router. The gateway router and a core router are part of a broadcast subnet. In order to regulate the traffic, the operators applied an ACL on the outbound interface from each router to the broadcast subnet. However, this decision results in leakage of undesirable traffic from unregistered users in one VLAN to other VLANs that share the same designated router. Since some routers are the first-hop gateways for over twenty VLANs, undesired communication is being permitted between a large number of hosts. The operators confirmed that systematic design had identified a previously unknown error in their ACL placement, and thanked us for pointing it out.

Fig. 11(b) illustrates a correct placement. It involves duplicating and moving the ACL to each inbound VLAN interface, and could result in significantly more rules. We hypothesize that the inconsistency arose as the operators tried to cut the number of rules in an ad-hoc fashion. Such errors can be easily avoided by systematic design approaches.

**Customizing placement for operator objectives:** To illustrate our systematic approach for customizing ACL placement, we consider the largest ACL in the campus network. This ACL consists of 693 rules - in contrast, all other ACLs in the network have no more than 60 rules. The ACL policy permits a specified list of hosts across various client VLANs to access a server VLAN - all other hosts are denied access to the server VLAN.

In the current design, all rules are placed in the last-

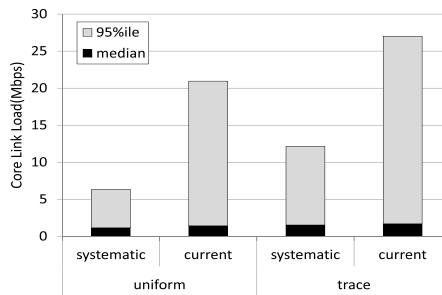


Fig. 10. Data traffic load on core links using the uniform and the trace traffic models, with  $f=50$ .

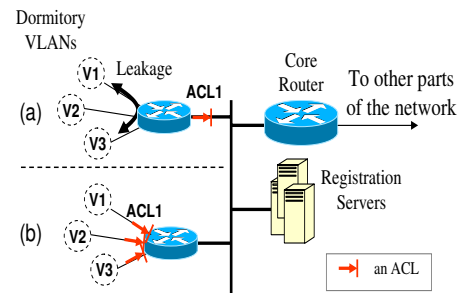


Fig. 11. (a) Scenario of ACL placement inconsistent with intent. (b) The corrected placement.

Metrics	$c(r) = \infty$					$c(r) \leq 300$					
	By	Fine-Grained				By	Fine-Grained				
	Col.	MIN	LB	CB	SEC	Col.	MIN	LB	CB	SEC	
$b(r)=\#$ rules on $r$	693	1169	2434	1169	1169	N/A	1369	2408	2389	1369	
$c(r)=\text{ACL capacity of } r$	693	418	280	1169	418	N/A	280	280	280	280	
$\sum_r b(r)$	693	418	280	1169	418	N/A	280	280	280	280	
$\max_r \{b(r)\}$	$\infty$	$\infty$	$\infty$	$\infty$	$\infty$	N/A	20	20	20	20	
$\min_r \{c(r) - b(r)\}$	$\infty$	$\infty$	0	0.1	2.06	0	N/A	0	0.09	0.08	0
$H$	1.69	0	0	0.1	2.06	0	N/A	0	0.09	0.08	0

TABLE II  
PLACEMENT OF ACL RULES BASED ON VARIOUS OPERATOR OBJECTIVES UNDER TWO EXTREME RESOURCE CONSTRAINTS.

hop router to the destination server VLAN. While this is a reasonable placement, there are alternative strategies that may be of interest to an operator. For instance, an operator may prefer to drop unwanted traffic closer to the source, or may wish to reduce the total rules placed on the router.

Table II illustrates how our approach can enable an operator to flexibly choose from a range of placement strategies based on the desired criteria of interest. Each column corresponds to a placement scheme, and each row corresponds to the metric used to rate a placement scheme.

The left half of the table presents results with these schemes assuming no constraints on the number of rules that may be placed on any router ( $c(r)=\infty$ ). One of our strategies (column-based placement) does match the design currently employed in the network. This strategy performs best in terms of keeping the total rules across the network small, for reasons elaborated in §IV-D2. However, other strategies offer benefits in alternate metrics of interest to the operator. For instance, the fine-grained SEC strategy pushes all rules to the first-hop router ( $H=0$ ), ensuring that traffic is filtered as early as possible, while the LB strategy ensures the maximum number of rules in any router is at most 280.

In networks built with low-end routers, it may not be feasible to place all rules in one router. To show the potential value of our systematic approach in such environments, we limit the processing capability of all routers in the network to be fewer than 300 rules ( $c(r)\leq 300$ ). The right half of Table II presents the results from systematic placement in this regime. Unlike column-based placement, all fine-grained strategies are able to produce a feasible placement despite the tight constraint. In addition, the various strategies offer benefits in metrics they target. For instance, the MIN strategy ensures the total number of rules is small (1369). Interestingly, the strategy also performs well in the other metrics.

Fig. 12 depicts how rules are distributed in the network after applying the fine-grained LB strategy in this setting.

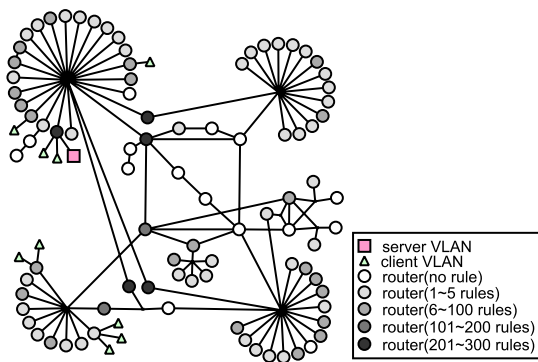


Fig. 12. Layer-3 topology showing systematic distribution of ACL rules after applying fine-gained, LB placement strategy.

Only routers and relevant VLANs (i.e., the server VLAN, and client VLANs with permitted hosts to the server VLAN) are shown. The number of rules varies per router, depending on the topology and the number of client VLANs attached to the router. Overall, the LB strategy spreads the load across the network, with no router having more than 280 rules. This exhibits the potential to systematically design the placement for the entire network with only lower-end hardware.

## VI. RELATED WORK

Many prior efforts on systematic network design focus on tasks encountered in carrier networks, such as configuring BGP policies [7], [8], [18], [9], optimizing OSPF weights, and redundancy planning [19]. In contrast, we focus on tasks in enterprise networks, which has received limited attention.

A few recent studies [20], [21], [22], [23], [24], [25] are partially motivated by enterprise networks. Most of them consider clean-slate designs by rearchitecting the control plane itself to contain the complexity of network design. In contrast, our work is relevant to both existing enterprise environments and clean slate designs.

Industry-driven efforts to simplifying enterprise network configuration involve template-based approaches [26], [27], [28], [29], [30], and abstract languages to specify configurations in a vendor-neutral fashion [31], [32], [33]. However, these approaches merely model the low-level mechanism and configuration, and do not abstract high-level operator intent.

A logic-based approach to configuration generation based on model-finding is presented in [6]. The focus is on the generation of correct configurations, and the system does not support optimization to meet desired performance objectives. Our previous papers [3], [4] have looked at bottom-up analysis of the VLAN design of an operational network, and reachability policies of existing networks. In contrast, our focus in this paper is on systematic design in these areas.

## VII. DISCUSSION AND OPEN ISSUES

In this paper, we have taken a first step towards the systematic design of enterprise networks. The contribution of this work is not only in providing the first set of heuristics for automating arguably two of the most complex tasks in enterprise network design, but also in the methodology that we have used to derive these heuristics.

Our methodology consists of three distinct steps. First,

we model operational goals with network-wide abstractions: e.g., the traffic matrix for the task of VLAN design, and the reachability matrix for the task of reachability control. Second, we formulate each task as a set of optimization problems, each modeling a different design strategy, and all subject to correctness and feasibility criteria associated with the task. Third, we develop heuristics to solve each of the optimization problems. While our goal is to devise practical heuristics that provide “good” solutions to these problems, it may be interesting to conduct an extensive study on the optimality of our heuristics by comparing our solutions with the “optimal” solutions. We leave this study for future work.

We recognize that this methodology is not without technical challenges when applied to a new enterprise network design task. The most challenging part is to find suitable network-wide abstractions to model the operational goals. While our experience suggests that it is very beneficial to study the configurations of existing operational networks [1], [3], whether there exists a general method for finding such abstractions remains an open research question. Another open question is how to best integrate the solutions for different design tasks into a complete network design. The design space of different tasks may overlap. For example, a particular choice of routing design may impact how optimal a solution our packet filter placement heuristics can achieve.

The ultimate goal for this area of research is to develop a system that enterprise network managers can use to produce, for a given topology of routers and switches, a complete set of configuration files *ready to be installed* into all the devices. While we view our work as an important step towards this goal, there is a semantic gap between the input and output we consider for the heuristics and the actual information network managers deal with. We envision the need for human-friendly languages (or GUIs) and associated interpreters to specify and translate operational goals into the network-wide abstractions proposed in this paper. When upgrading an existing network, the baseline data including the traffic matrix, reachability matrix, etc., can be obtained by measurements or static analysis of existing network configurations [4]. We also envision the need for tools similar to PRESTO [30] to convert systematic design solutions into device-vendor-specific configuration commands. All these requirements create a fertile ground for future research.

In this work, we mainly focus on the design of greenfield networks, or the networks to be deployed. One interesting problem is how we may optimize an *existing* network in the wild, taking into account the costs of making changes. That is, how we may not only optimize the design of the network, but also minimize the amount of changes needed to go from the current design to the near optimal one. Also, given that enterprise networks usually keep growing and evolving, another interesting problem is how we may deal with the evolution of the network after the initial systematic deployment. We are investigating these problems in our ongoing work.

One limitation of this work is that we have validated the performance of our heuristics only on a single network. Obtaining access to data not only takes significant effort, and extensive interactions with operators, but is sometimes infeasible given

the sensitive nature of such data-sets. Access to enterprise network data is a key challenge for the community, and in our parallel ongoing efforts, we are investigating the feasibility of creating enterprise data repositories that can be shared by the community.

### VIII. CONCLUSIONS

In this paper, we have shown the viability and importance of a systematic approach to two key design tasks in enterprise networks: VLAN design and reachability control. Our contributions include (i) a systematic formulation of these critical but poorly understood enterprise design tasks, (ii) a set of algorithms to solve the formulated problems, and (iii) a validation of the systematic approach on a unique large-scale campus network dataset.

Our evaluations show the promise of our approach. The campus network we analyzed is well-run, and many hours of human design time have been spent on it. Yet, our approach produces better results with less human effort. Beyond the general time savings in the design process, a systematic approach can ensure correctness and lead to significantly better designs. For example, through systematic VLAN design, broadcast and data traffic on the core links of the campus network can be reduced by over 24% and 55%, respectively. Systematic placement of ACLs ensures the design correctly conforms to the operator's security objectives. In contrast, today's ad-hoc design processes can result in inconsistencies such as those we pointed in our analysis. Finally, our approach can be customized to optimize for operator-preferred design strategies, and can produce designs tailored to network parameters such as traffic patterns and router resource constraints.

For future work, we hope to gain experience with our approach on a wider range of enterprise networks, and apply the systematic approach to other enterprise design tasks.

### IX. ACKNOWLEDGMENTS

We thank our colleagues in the Information Technology Department at Purdue (ITaP), for providing access to the data, and for being generous with their time. Particular thanks are due to Duane Kyburz and Brad Devine, who were our primary point of contact at ITaP, and enthusiastically met us on several occasions. We also thank David Collins and Sunil Dath Krothapalli for their help in evaluating our systematic design. This work was supported by NSF awards CNS-0721488, CNS-0520210, and CNS-0721574.

### REFERENCES

- [1] D. Maltz, G. Xie, J. Zhan, H. Zhang, G. Hjalmtysson, and A. Greenberg, "Routing design in operational networks: A look from the inside," in *Proc. ACM SIGCOMM*, 2004.
- [2] F. Le, G. G. Xie, D. Pei, J. Wang, and H. Zhang, "Shedding light on the glue logic of the internet routing architecture," in *Proc. ACM SIGCOMM*, 2008.
- [3] P. Garimella, Y.-W. E. Sung, N. Zhang, and S. Rao, "Characterizing vlan usage in an operational network," in *Proc. of ACM SIGCOMM INM workshop*, 2007.
- [4] G. Xie, J. Zhan, D. A. Maltz, H. Zhang, A. Greenberg, G. Hjalmtysson, and J. Rexford, "On static reachability analysis of IP networks," in *Proc. IEEE INFOCOM*, 2005.
- [5] Z. Kerravala, "Configuration management delivers business resiliency," The Yankee Group, Nov. 2002.
- [6] S. Narain, "Network configuration management via model finding," in *Proc. Large Installations Systems Administration (LISA) Conference*, 2005.
- [7] C. Alaettinoglu, C. Villamizar, E. Gerich, D. Kessensand, D. Meyer, T. Bates, D. Karenberg, and M. Terpstra, *Routing Policy Specification Language (RPSL)*, Internet Engineering Task Force, June 1999, rFC 2622.
- [8] T. G. Griffin and J. L. Sobrinho, "Metarouting," in *Proc. ACM SIGCOMM*, 2005.
- [9] J. Gottlieb, A. Greenberg, J. Rexford, and J. Wang, "Automated provisioning of BGP customers," in *IEEE Network Magazine*, Dec. 2003.
- [10] M. R. Garey and D. S. Johnson, *Computers and Intractability: A Guide to the Theory of NP-Completeness*. New York: W.H. Freeman, 1979.
- [11] R. E. Burkard, E. Çela, P. M. Pardalos, and L. S. Pitsoulis, "Quadratic assignment problems," *European Journal of Operational Research*, vol. 15, no. 3, pp. 283–289, March 1984.
- [12] S. Voss, "Heuristics for nonlinear assignment problems," *Nonlinear Assignment Problems: Algorithms and Applications*, pp. 175–215, 2000.
- [13] P. E. Dunne, "An annotated list of selected np-complete problems," [http://www.csc.liv.ac.uk/~ped/teachadmin/COMP202/annotated\\_np.html](http://www.csc.liv.ac.uk/~ped/teachadmin/COMP202/annotated_np.html), 2008.
- [14] A. X. Liu, E. Torng, and C. Meiners, "Firewall compressor: An algorithm for minimizing firewall policies," in *Proc. IEEE INFOCOM*, April 2008.
- [15] S. Acharya, J. Wang, Z. Ge, T. Znati, and A. Greenberg, "Simulation study of firewalls to aid improved performance," in *Proc. ANSS*, Washington, DC, USA, 2006, pp. 18–26.
- [16] T. H. Cormen, C. Stein, R. L. Rivest, and C. E. Leiserson, *Introduction to Algorithms*. New York: McGraw-Hill, 2001.
- [17] R. Pang, M. Allman, M. Bennett, J. Lee, V. Paxson, and B. Tierney, "A first look at modern enterprise traffic," in *Proc. ACM SIGCOMM IMC*, 2005.
- [18] H. Boehm, A. Feldmann, O. Maennel, C. Reiser, and R. Volk, "Network-wide inter-domain routing policies: Design and realization," Apr. 2005, draft.
- [19] R. Rastogi, Y. Breitbart, M. Garofalakis, and A. Kumar, "Optimal configuration of ospf aggregates," *IEEE/ACM Transaction on Networking*, 2003.
- [20] M. Caesar, D. Caldwell, N. Feamster, J. Rexford, A. Shaikh, and Jacobus van der Merwe, "Design and implementation of a Routing Control Platform," in *Proc. NSDI*, 2005.
- [21] N. Feamster, H. Balakrishnan, J. Rexford, A. Shaikh, and J. van der Merwe, "The case for separating routing from routers," in *Proc. ACM SIGCOMM Workshop on Future Directions in Network Architecture*, 2004.
- [22] M. Casado, T. Garfinkel, A. Akella, M. Freedman, D. Boneh, N. McKeown, and S. Shenker, "SANE: A protection architecture for enterprise networks," in *Proc. USENIX Security*, 2006.
- [23] A. Greenberg, G. Hjalmtysson, D. A. Maltz, A. Myers, J. Rexford, G. Xie, H. Yan, J. Zhan, and H. Zhang, "A clean slate 4D approach to network control and management," *ACM Computer Communication Review*, October 2005.
- [24] J. Rexford, A. Greenberg, G. Hjalmtysson, D. A. Maltz, A. Myers, G. Xie, J. Zhan, and H. Zhang, "Network-wide decision making: Toward a wafer-thin control plane," in *Proc. ACM SIGCOMM HotNets Workshop*, November 2004.
- [25] H. Ballani and P. Francis, "Conman: a step towards network manageability," in *Proc. ACM SIGCOMM*, 2007.
- [26] "Cisco IP solution center," <http://www.cisco.com/en/US/products/sw/netmgtsw/ps4748/index.html>.
- [27] "Intelliden," <http://www.intelliden.com/>.
- [28] "Opsware," <http://www.opsware.com/>.
- [29] "Voyence," <http://www.voyence.com/>.
- [30] W. Enck, P. McDaniel, S. Sen, P. Sebos, S. Spoerel, A. Greenberg, S. Rao, and W. Aiello, "Configuration management at massive scale: System design and experience," in *Proc. USENIX*, 2007.
- [31] J. Case, M. Fedor, M. Schoffstall, and J. Davin, "A simple network management protocol (SNMP)," RFC 1157, May 1990.
- [32] Distributed Management Task Force, Inc., <http://www.dmtf.org>.
- [33] "DSL forum TR-069," <http://www.broadband-forum.org/technical/download/TR-069.pdf>.