

# Optimal Configuration of OSPF Aggregates

Rajeev Rastogi, Yuri Breitbart, Minos Garofalakis, *Associate Member, IEEE*, and Amit Kumar

**Abstract**—*Open Shortest Path First (OSPF)* is a popular protocol for routing within an autonomous system (AS) domain. In order to scale for large networks containing hundreds and thousands of subnets, OSPF supports a two-level hierarchical routing scheme through the use of *OSPF areas*. Subnet addresses within an area are aggregated, and this aggregation is a crucial requirement for scaling OSPF to large AS domains, as it results in significant reductions in routing table sizes, smaller link-state databases, and less network traffic to synchronize the router link-state databases. On the other hand, address aggregation also implies loss of information about the length of the shortest path to each subnet, which in turn, can lead to *suboptimal* routing.

In this paper, we address the important practical problem of configuring OSPF aggregates to minimize the error in OSPF shortest-path computations due to subnet aggregation. We first develop an optimal dynamic programming algorithm that, given an upper bound  $k$  on the number of aggregates to be advertised and a weight assignment function for the aggregates, computes the  $k$  aggregates that result in the minimum cumulative error in the shortest-path computations for all source–destination subnet pairs. Subsequently, we tackle the problem of assigning weights to OSPF aggregates such that the cumulative error in the computed shortest paths is minimized. We demonstrate that, while for certain special cases (e.g., unweighted cumulative error) efficient optimal algorithms for the weight assignment problem can be devised, the general problem itself is NP-hard. Consequently, we have to rely on search heuristics to solve the weight assignment problem. To the best of our knowledge, our work is the first to address the algorithmic issues underlying the configuration of OSPF aggregates and to propose efficient configuration algorithms that are *provably optimal* for many practical scenarios.

**Index Terms**—Area border routers, dynamic programming, IP address aggregation, optimal routing, Open Shortest Path First (OSPF), OSPF advertisements, OSPF weights.

## I. INTRODUCTION

**O**PEN Shortest Path First (OSPF) is a widely used protocol for routing within an autonomous system (AS) domain in today's Internet [1]–[3]. To scale for large AS networks, OSPF implements a two-level hierarchical routing scheme through the deployment of *OSPF areas*. Each OSPF area comprises a collection of subnets interconnected by routers. Detailed information about links and subnets within an OSPF area is flooded throughout the elements connected to the area. As a result, every router knows the exact topology of its enclosing OSPF area; this includes the subnets and the links connecting routers within the area. On the other hand, details of an area's topology are not

advertised beyond the area's borders and are, thus, hidden from other areas in the same AS. Instead, subnet addresses within each area are grouped into *aggregates* and only these aggregates are flooded into the rest of the network (thus, making an area's subnets reachable from the remainder of the AS network). This task of advertising aggregate information about subnets in an area is carried out by *area border routers (ABRs)*, that is, routers attached to two or more areas.

OSPF areas and address aggregation are crucial in enabling OSPF to scale for AS domains comprising hundreds or thousands of subnets; specifically, they play an important role in optimizing router and network resource consumption, as explained below.

- 1) **Router Memory:** For OSPF areas *not* directly connected to a router in the AS, the router's routing tables only need to contain entries corresponding to subnet aggregates rather than individual subnet addresses. In other words, a router stores individual subnet addresses in its routing table only for the OSPF areas that are directly connected to it. This obviously leads to smaller routing table sizes and, thus, lower memory requirements at routers.
- 2) **Router Processing Cycles:** The link-state database maintained at each router is much smaller, since it only needs to include summary information for subnets belonging to OSPF areas not directly connected to the router. Consequently, the computational cost of the shortest-path calculation decreases substantially.
- 3) **Network Bandwidth:** For subnets within each OSPF area, only aggregate address information (rather than individual subnet addresses) is flooded into the rest of the AS network. As a result, the volume of OSPF flooding traffic necessary to synchronize the link-state databases of the AS routers is significantly reduced.

Nevertheless, despite its obvious benefits, OSPF address aggregation involves important practical tradeoffs. This is because address aggregation typically results in loss of information which, in turn, can lead to *suboptimal* routing paths. To see this, we need to delve in more depth into how OSPF routing works in the presence of address aggregation. Briefly, each ABR attaches a *weight* to each aggregate that it advertises to the rest of the network. This weight is critical in determining the path used by a router external to the area to reach subnets covered by the aggregate. More specifically, among all the ABRs advertising the aggregate (with possibly different weights), an external router chooses the ABR (say,  $b$ ) that minimizes the sum of the following two quantities: 1) the length of the shortest path from the external router to the border router  $b$  and 2) the weight advertised by  $b$  for the aggregate. Once such an ABR  $b$  is chosen, IP packets from the external router to every subnet covered by the aggregate are

Manuscript received November 15, 2001; revised May 6, 2002; approved by IEEE/ACM TRANSACTIONS ON NETWORKING Editor J. Rexford.

R. Rastogi, M. Garofalakis, and A. Kumar are with Bell Laboratories, Murray Hill, NJ 07974 USA (e-mail: rastogi@lucent.com; minos@lucent.com; amitk@lucent.com).

Y. Breitbart is with the Department of Computer Science, Kent State University, Kent, OH 44240 USA (e-mail: yuri@cs.kent.edu).

Digital Object Identifier 10.1109/TNET.2003.810317

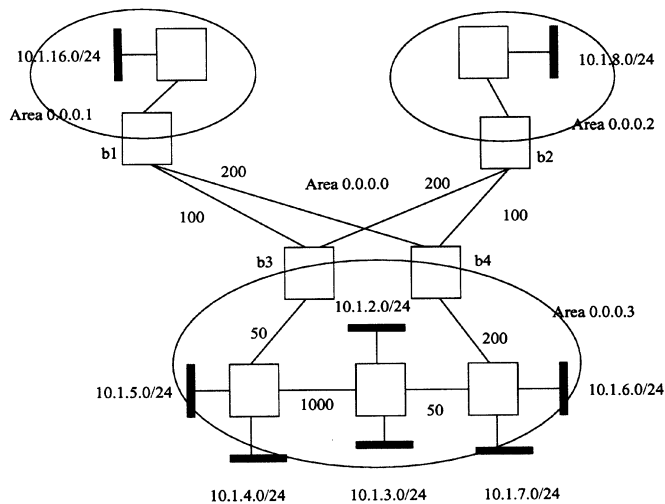


Fig. 1. Example of suboptimal routing due to address aggregation.

forwarded along the shortest path from the external router to  $b$  and, subsequently, along the shortest path from  $b$  to the subnet. However, for certain subnet(s) covered by the aggregate, this path may be significantly suboptimal, since there can be a much shorter path from the external router to the subnet through a different ABR. This is illustrated in the following example.

*Example 1:* Consider the AS network consisting of the four areas 0.0.0.0, 0.0.0.1, 0.0.0.2, and 0.0.0.3 shown in Fig. 1. (Area 0.0.0.0 corresponds to the ASs *backbone area* that interconnects the ABRs of the different OSPF areas in the AS.) The boxes in the figure are routers, while the thin black rectangles denote subnets. The figure also illustrates the various subnet addresses and the weight of each link connecting a pair of routers. ABR  $b_1$  belongs to area 0.0.0.1,  $b_2$  to area 0.0.0.2, and  $b_3$  and  $b_4$  to area 0.0.0.3. The subnet addresses in area 0.0.0.3 can be aggregated to different degrees. For instance, the aggregate 10.1.0.0/21 covers all the subnets in the area. In contrast, 10.1.4.0/22 covers subnets 10.1.4.0/24, 10.1.5.0/24, 10.1.6.0/24 and 10.1.7.0/24, while 10.1.2.0/23 covers subnets 10.1.2.0/24 and 10.1.3.0/24.

Suppose one of the aggregate addresses advertised by the ABRs of area 0.0.0.3 is 10.1.4.0/22. Suppose further that each ABR assigns a weight to the aggregate that equals the distance of the furthest component subnet in the aggregate from the router (as suggested, for example, by Moy [3]). Thus, router  $b_3$  advertises 10.1.4.0/22 with a weight of 1100 (distance of subnet 10.1.6.0/24 from  $b_3$ ), while router  $b_4$  advertises 10.1.4.0/22 with a weight of 1250 (distance of subnet 10.1.4.0/24 from  $b_4$ ). Thus, external router  $b_1$  belonging to area 0.0.0.1 forward all packets to subnets in 10.1.4.0/22 through border router  $b_3$ , since the shortest path to the aggregate through  $b_3$  has a length of  $100 + 1100 = 1200$ , while the shortest path through  $b_4$  has length  $200 + 1250 = 1450$ . Note, however, that the path from  $b_1$  to subnets 10.1.6.0/24 and 10.1.7.0/24 passing through router  $b_3$  has length 1200 and is suboptimal since the shortest path from  $b_1$  to both subnets is through border router  $b_4$  and its length is only 400. Thus, even in this simple scenario, address aggregation results in an error of  $1200 - 400 = 800$  in the optimal route (i.e., shortest-path) computation between  $b_1$  and each of the subnets 10.1.6.0/24 and 10.1.7.0/24.

Further, note that considering different weight assignments for the aggregates does not alleviate the problem. The root of the problem is that a single border router is selected by  $b_1$  for reaching all subnets in 10.1.4.0/22. If  $b_4$  is chosen instead of  $b_3$ , then the paths from  $b_1$  to subnets 10.1.4.0/24 and 10.1.5.0/24 through  $b_4$  become much longer (their length is  $200 + 1250 = 1450$ ) compared to the shortest paths to the subnets that pass through  $b_3$  (whose length is  $100 + 50 = 150$ ).  $\square$

The primary reason for suboptimal paths being selected when subnets are aggregated is that a single weight is used by each ABR for all the subnets covered by the aggregate; obviously, a single weight may be incapable of accurately capturing the distance of the border router to every covered subnet. This problem is exacerbated when the aggregated subnets are spread across the area with some subnets being in close proximity to distinct border routers in the area. This was precisely the case in Example 1, where aggregate 10.1.4.0/22 spans two subnets: 10.1.5.0/24 (close to  $b_3$ ) and 10.1.6.0/24 (close to  $b_4$ ), the distance between which is greater than 1000. The end result is that the single weight advertised for aggregate 10.1.4.0/22 is not representative of the true distance between any border router (either  $b_3$  or  $b_4$ ) and the two subnets in 10.1.4.0/22.

One possible way to reduce the error in suboptimal OSPF routing paths in the presence of aggregation is to avoid aggregating distant subnets that are close to multiple border routers. Thus, in Example 1, instead of advertising the single aggregate 10.1.4.0/22, one can choose to advertise two aggregates 10.1.4.0/23 (with weights 50 and 1250 at  $b_3$  and  $b_4$ , respectively) and 10.1.6.0/23 (with weights 1100 and 200 at  $b_3$  and  $b_4$ , respectively). This clearly reduces the error in the selected paths to zero, since the assigned weights capture the ABRs' distances to the aggregated subnets *exactly*. Thus, there is an important tradeoff between the number of aggregates advertised (and, consequently, the size of the routing tables) and the error in the selected shortest paths. This tradeoff is further complicated by the fact that the aggregates advertised by OSPF border routers do not have to be disjoint—it is entirely possible for one advertised aggregate to be *completely contained* in another. In such a scenario, the *longest match property* of IP routing causes the more specific aggregate to take precedence for route computation to subnets within the aggregate. Also, configuring the *weight assignments* used by ABRs for address aggregates is another important mechanism for controlling the quality of the OSPF routing paths in the presence of aggregation. In the following example, we illustrate how, by carefully selecting the address aggregates as well as the associated weights, all the subnets in area 0.0.0.3 (Fig. 1) can be advertised using only two (overlapping) aggregates while incurring zero error in shortest-path computations.

*Example 2:* Consider the AS depicted in Fig. 1. One way to ensure that the error in the selected paths to area 0.0.0.3's subnets is zero is to have the ABRs advertise the following three aggregates: 10.1.4.0/23, 10.1.2.0/23, and 10.1.6.0/23. The reason for this is that it is possible to choose weights for each aggregate at each ABR such that the weight equals the exact distance between the border router and every subnet covered by the aggregate. For instance, for aggregate 10.1.2.0/23, weights of 1050 and 250 at ABRs  $b_3$  and  $b_4$ , respectively, reflect the exact distances of the border routers to subnets in it.

Achieving zero error with only two aggregates is more challenging. Note that all the subnets in area 0.0.0.3 can be covered by the single aggregate 10.1.0.0/21; however, this aggregate by itself cannot result in zero error. Another possibility is to consider the two disjoint aggregates 10.1.4.0/22 and 10.1.2.0/23, which cover all the subnets. However, as we saw earlier, since subnets 10.1.4.0/23 and 10.1.6.0/23 covered by 10.1.4.0/22 are closer to different routers (and distant from each other), this cannot result in zero error, either. Thus, the key to optimizing the error is to bundle 10.1.4.0/23 into one aggregate, and 10.1.2.0/23 and 10.1.6.0/23 into the other. It turns out that this can be achieved by advertising the following two aggregates: 10.1.0.0/21 and 10.1.4.0/23. The longest match characteristic of IP routing causes the latter aggregate to be used for routing to subnets in 10.1.4.0/23 and the former to be used to route to subnets in 10.1.2.0/23 and 10.1.6.0/23.

One question still remains: what weights should be assigned to each aggregate? While this is straightforward for the aggregate 10.1.4.0/23 (since the two subnets 10.1.4.0/24 and 10.1.5.0/24 covered by it are at the same distance from any given border router), it is somewhat less obvious for the aggregate 10.1.0.0/21. Simply setting the weight equal to the distance of the ABR to the most distant covered subnet (see, e.g., [3]) may not result in the least error. To see this, suppose  $b_3$  and  $b_4$  advertise 10.1.0.0/21 with weights 1100 and 1250, respectively (i.e., the maximum distance to a subnet contained in the aggregate). This causes  $b_1$  to select  $b_3$  for subnets in 10.1.0.0/21, and the resulting cumulative error in the selected paths from  $b_1$  to all subnets in area 0.0.0.3 is  $2 * (150 - 150) + 2 * (1150 - 450) + 2 * (1200 - 400) = 3000$ . On the other hand, a lower cumulative error can be achieved if  $b_3$  and  $b_4$  advertise 10.1.0.0/21 with weights 730 and 570, respectively (i.e., the average distance to the subnets contained in the aggregate). In this case,  $b_1$  selects border router  $b_4$  to access subnets in 10.1.0.0/21, resulting in a lower cumulative error of  $2 * (1450 - 150) + 2 * (450 - 450) + 2 * (400 - 400) = 2600$  (assuming again that only 10.1.0.0/21 is advertised for area 0.0.0.3).

Further, configuring border routers  $b_3$  and  $b_4$  to advertise aggregates 10.1.0.0/21 and 10.1.4.0/23,  $b_3$  with weights 730 and 50, and  $b_4$  with weights 570 and 1250, causes the cumulative error for  $b_1$  to reduce to zero. This is because  $b_1$  selects ABR  $b_3$  for subnets in 10.1.4.0/23 and ABR  $b_4$  for the remaining subnets (that is, subnets in 10.1.0.0/21 but not contained in 10.1.4.0/23). Thus, the selected paths after aggregation are indeed the shortest paths from  $b_1$  to the subnets in area 0.0.0.3.  $\square$

#### A. Contributions of This Paper

In this paper, we address the important practical problem of configuring OSPF aggregates to minimize the error in OSPF shortest-path computations due to address aggregation. From our discussion above, we can see that OSPF aggregate configuration involves two key subproblems: 1) selecting the aggregates to advertise at each ABR and 2) assigning weights to each advertised aggregate at each ABR. We address each of these two problems separately. We first develop an *optimal* dynamic programming algorithm that, given an upper bound  $k$  on the

number of aggregates to be advertised by the ABRs and a weight assignment function for the aggregates, computes the  $k$  aggregates that result in the minimum cumulative/maximum error in the OSPF shortest-path computations for all source–destination subnet pairs. This problem is obviously relevant when there is a limit on the number of aggregates that can be advertised within an AS in order to bound the routing table sizes, number of entries in the link-state database, or the amount of network traffic due to OSPF advertisements. The objective then is to choose the  $k$  aggregates to advertise such that the selected paths are as close to the shortest paths as possible (where “closeness” is measured in terms of either the total or the maximum over all source–destination subnet pairs in the AS). Furthermore, our proposed algorithm can be easily extended to optimally solve the dual OSPF configuration problem, where the goal is to compute the minimum number of aggregates so that the (cumulative or maximum) error in selected paths is less than a certain user-specified threshold.

We then address our second subproblem of selecting weights for OSPF aggregates at each ABR such that the deviation of selected paths from the shortest paths is minimized. More specifically, we attack the following problem: Given an address aggregate  $x$ , determine an assignment of weights to  $x$  at each ABR in its area such that the (cumulative or maximum) error in the selected paths between source–destination subnet pairs is minimized. We demonstrate that, while for certain restricted (but interesting) cases the above problem can be solved in polynomial time, the general problem itself is NP-hard. Consequently, we have to rely on search heuristics to solve the weight assignment problem. We also propose a randomized search strategy for the general case of weighted cumulative error, and an optimal pseudopolynomial time algorithm for the maximum error case.

The second subproblem involving weight selection for OSPF aggregates is clearly important since, as shown in Example 2, assigning to each aggregate (at an ABR), a weight equal to the maximum distance from the ABR of subnets covered by the aggregate (as suggested, for example, by Moy [3]), may not minimize error. This is because the maximum distance may, at times, be a poor estimate of the distance between subnets of the aggregate and the ABR. A better alternative that has the ability to capture the distance between an aggregate and an ABR fairly accurately, is the average distance between the aggregate’s subnets and the ABR. In fact, in this paper, we show that choosing the average distance as the weight for an aggregate minimizes the cumulative error for a single advertised aggregate and is, thus, in many respects, more representative of the distance between an aggregate and an ABR than maximum distance. Consequently, it is important to consider alternative weight assignments for aggregates at ABRs, besides the maximum distance advocated by [3]—essentially, depending on the error metric we wish to optimize, one set of weights may be more representative than others of distances between aggregates and ABRs, and may ultimately lead to the selection of a better set of aggregates (solution to the first subproblem).

Combining our algorithmic results for the two subproblems provides us with an effective integrated solution for configuring (in many cases, optimally) OSPF aggregates in an AS domain. The idea is to first apply our weight assignment algorithms to

determine “good” weights for all candidate aggregates at each ABR, and then employ our dynamic programming algorithm to select the optimal subset of aggregates to advertise (given these weight assignments). To the best of our knowledge, our work is the first to address the algorithmic issues underlying the configuration of OSPF aggregates and to propose efficient configuration algorithms that are *provably optimal* for many practical scenarios. One point to note, however, is that we solve the aggregate-selection and weight assignment subproblems separately. Although they can be combined into a single problem, we believe the joint problem to be intractable. The only solution we know of has exponential complexity and is, thus, impractical.

Finally, we must point out that the problem of selecting aggregates is considerably simplified if networks and areas are designed carefully, and IP addresses are assigned to subnets within areas in a systematic manner. For example, one possible strategy is to segment an area into  $k$  smaller regions such that routers within each region are close to each other. Then, it is possible to achieve good aggregation by assigning consecutive IP addresses to subnets within a region, and defining a single aggregate per region. The above-mentioned approach, however, may not work in practice since networks are seldom static and tend to continuously evolve due to the addition of new network elements, subnets, and links. For instance, addition of new routers may require new IP aggregates to be allocated to a region. Similarly, deletion or failure of an existing link within a region, or addition of a new link between regions, may cause distances between routers to be substantially altered, thus, rendering the previously selected aggregates suboptimal. Our algorithms allow OSPF aggregates to be optimally configured in such continuously changing dynamic networks, and also to work on-line: emergence of spurious subnets (through external advertisements), link failures, and topology changes would trigger the computation of changes to the optimal set of advertised aggregates/weights. This computation can also be carried out incrementally and efficiently.

## II. SYSTEM MODEL AND PROBLEM FORMULATION

### A. System Model and Notation

We model the network as an undirected graph. Nodes in the graph correspond to either routers or subnets. Edges in the graph connect routers with other routers or subnets. A link exists between two routers if the two routers can directly exchange IP packets without going through an intermediate router (that is, the two routers are either connected to the same subnet or are connected by a point-to-point link). A link exists between a subnet and a router if the subnet is connected to the router. Each link has an associated weight, which is the OSPF weight assigned to the link (that is, the interface to which the link is connected to). For simplicity, we assume that the link is assigned the same weight at both ends—our algorithms, however, are applicable even if link weights are not symmetric. Table I describes the notation employed in this paper.

The set of subnets  $\mathcal{S}$  in the network are partitioned into disjoint areas. The set of areas is denoted by  $\mathcal{R}$  and the set of subnets in area  $R_i \in \mathcal{R}$  is denoted by  $S_i$ . A router is said to be attached

TABLE I  
NOTATION USED IN THE PAPER

Symbol	Description
$s, t$	Source, destination subnets
$\mathcal{R}$	Set of OSPF areas
$R_i$	Generic term for OSPF area
$\mathcal{S}$	Set of all subnets in autonomous system domain
$S_i$	Set of subnets in area $R_i$
$\mathcal{B}$	Set of all ABRs
$B_i$	Set of ABRs for area $R_i$
$b, c$	Generic letters for an ABR
$\mathcal{A}$	Set of aggregates eligible for advertising
$A_i$	Set of aggregates eligible for advertising in $R_i$
$x, y$	Generic advertised address aggregates
$X, Y$	Generic sets of advertised address aggregates
$D(s, t)$	Degree of importance of the source-destination subnet pair $(s, t)$
$W_X$	Weight assignment function for aggregates in $X$
$lsp(s, t)$	Length of shortest path between subnets $s$ and $t$
$lsp(s, t, X, W_X)$	Length of shortest path between subnets $s$ and $t$ when aggregates in $X$ are advertised with weights in $W_X$

to area  $R_i$  if it is directly connected to a subnet in  $S_i$ . A router that is attached to two or more areas is called an ABR. We denote by  $B_i$  the set of ABRs attached to area  $R_i$ . In addition to area  $R_i$  (and possibly other areas), every ABR in  $B_i$  is also attached to a special *backbone* area. The backbone area serves to connect the subnets in the various other areas. We denote by  $lsp(s, t)$  the length of the shortest path between  $s$  and  $t$ , where  $s$  and  $t$  can be subnets or routers within the AS. Note that, if  $s$  and  $t$  belong to the same area  $R_i$ , then the shortest path between  $s$  and  $t$  is defined to be over links in area  $R_i$ . If, instead,  $s$  and  $t$  belong to distinct areas (say,  $R_i$  and  $R_j$ , respectively), then the shortest path between  $s$  and  $t$  involves two ABRs  $b \in B_i$  and  $c \in B_j$  and consists of three path segments: the first is the shortest path between  $s$  and  $b$  involving links in  $R_i$ , the second is the shortest path between  $b$  and  $c$  over links in the backbone area, and the final segment is the shortest path between  $c$  and  $t$  all of whose links belong to area  $R_j$ . Note that  $lsp$  can be defined in a similar fashion if either of the subnets  $s$  and  $t$  above are replaced by routers.

In OSPF, information relating to links and subnets in an area are flooded throughout the area. Consequently, routers attached to area  $R_i$  have detailed knowledge of  $R_i$ 's topology. As a result, IP packets originating in any subnet  $s$  belonging to area  $R_i$  destined to a subnet  $t$  in the same area are forwarded along the shortest path between  $s$  and  $t$ . However, in order to ensure smaller routing table sizes and reduce network traffic overhead, detailed information about individual subnets within an area are typically not advertised beyond the area's borders. Instead, area  $R_i$ 's ABRs will typically be configured to advertise a set of aggregates  $X$  that cover subnets in  $S_i$  and a separate weight for each aggregate in  $X$ . We denote by  $W_X(x, b)$ , the weight assigned to an aggregate  $x \in X$  by ABR  $b \in B_i$ . Each ABR in  $B_i$  floods in the entire backbone area, every aggregate  $x \in X$  along with the weight assigned to  $x$  by it—this causes ABRs belonging to every other area to receive  $x$ . An ABR  $c \in B_j$ , in turn, floods  $x$  into area  $R_j$  with an adjusted weight equal to  $\min_{b \in B_i} \{lsp(b, c) + W_X(x, b)\}$ . Thus, a subnet  $s$  in  $S_j$ , in order to reach aggregate  $x$  that covers subnets in  $S_i$ , selects a path passing through ABR  $b \in B_i$  for which  $lsp(s, b) + W_X(x, b)$  is minimum.

Due to the longest match property of IP routing, the most specific aggregate covering a subnet determines the path to the subnet. We say that an aggregate  $x$  is more specific than an aggregate  $y$  if  $x$  is contained in  $y$ , which we denote by  $x \in y$ . Thus, for a subnet  $t$  in  $S_i$ , if  $x$  is the most specific aggregate in  $X$  that covers  $t$ , then a subnet  $s$  in  $S_j$ , in order to reach  $t$ , selects the path comprising of the shortest path from  $s$  to  $b$  and then from  $b$  to  $t$ , where  $b \in B_i$  is the ABR for which  $\text{lsp}(s, b) + W_X(x, b)$  is minimum. We denote the length of this selected path from  $s$  to  $t$  for the set of advertised aggregates  $X$  and weight assignment  $W_X$  by  $\text{lsp}(s, t, X, W_X)$ . Thus,  $\text{lsp}(s, t, X, W_X) = \text{lsp}(s, b) + \text{lsp}(b, t)$ , and the error in the selected path is simply  $\text{lsp}(s, t, X, W_X) - \text{lsp}(s, t)$ . When  $s$  and  $t$  belong to the same area, we define  $\text{lsp}(s, t, X, W_X)$  to be equal to  $\text{lsp}(s, t)$ . Note that  $\text{lsp}(s, t, X, W_X) = \infty$  if  $X$  does not contain an aggregate that covers  $t$  (the implication here is that  $t$  is unreachable from  $s$ ).

### B. Problem Statement

We address the problem of computing the set of aggregates  $X$  advertised across all the areas in an AS and the weight assignment function  $W_X$  such that the error in the selected paths is minimized. Clearly, we need to impose certain restrictions on  $X$  and  $W_X$  in order to ensure the reachability of remote subnets in a different area. First, we require that  $X$  be *complete*, that is, every subnet in  $\mathcal{S}$  be covered by some aggregate in  $X$ . The next two restrictions serve to ensure that an ABR cannot advertise an aggregate covering a subnet in  $S_i$  unless it belongs to  $B_i$ . We say that an aggregate  $x$  is *eligible* if all the subnets in  $\mathcal{S}$  covered by it belong to a single area. Thus, in the network of Fig. 1, aggregate 10.1.0.0/21 is eligible, since it only covers subnets in Area 0.0.0.3; however, aggregate 10.1.0.0/20 is not, since it covers subnets 10.1.8.0/24 and 10.1.4.0/24, which belong to different areas. Let  $\mathcal{A}$  denote the set of all eligible aggregates such that every aggregate in  $\mathcal{A}$  covers at least one subnet in  $\mathcal{S}$ . Note that  $\mathcal{S} \subseteq \mathcal{A}$ . Further, let  $A_i \subseteq \mathcal{A}$  denote the set of eligible aggregates that cover subnets in  $S_i$ . We require that the set of advertised aggregates  $X \subseteq \mathcal{A}$ . We also require that only ABRs in  $B_i$  advertise aggregates in  $A_i$ . One way to model this is by requiring that  $W_X(x, b) = \infty$  if  $x \in A_i$  and  $b \notin B_i$ .

We are now in a position to state the two basic problems addressed in the remainder of this paper.

- 1) **Aggregate Selection Problem:** Given a  $k$  and a weight assignment function  $W_{\mathcal{A}}$ , compute a complete set  $X \subseteq \mathcal{A}$  containing at most  $k$  aggregates such that  $\sum_{s, t \in \mathcal{S}} (\text{lsp}(s, t, X, W_{\mathcal{A}}) - \text{lsp}(s, t))$  is minimum.
- 2) **Weight Selection Problem:** For an aggregate  $x \in A_i$ , compute a weight assignment function  $W_{\{x\}}$  such that  $\sum_{s, t \in \mathcal{S}, t \in x} (\text{lsp}(s, t, \{x\}, W_{\{x\}}) - \text{lsp}(s, t))$  is minimum.

The rationale for our objective function which aims to minimize the cumulative routing-path lengths is as follows. A link's weight typically is a measure of its desirability for routing (Cisco recommends setting weights inversely proportional to link capacities). Thus, paths with small weights are more desirable, and our problem formulation minimizes the sum of all routing-path weights.

Note that in both problem statements above, every source–destination subnet pair is assigned the same degree of importance. In other words, in the final error, the error in the selected path between every subnet pair is treated equally, that is, given an equal degree of importance. However, this is somewhat restrictive since minimizing the error in the selected paths for certain source–destination subnets may be more important. This may happen, for instance, for subnet pairs between which there is a disproportionately high volume of traffic, or subnet pairs carrying high-priority or delay-sensitive traffic such as voice. Thus, we can consider a *degree of importance* function  $D$  which for a pair of subnets  $s, t$  returns a real value  $D(s, t)$  that reflects the importance of minimizing the error in the selected path between subnets  $s$  and  $t$ . Note that  $D(s, t)$  can be any arbitrary function of the volume/priority of traffic flowing between subnets  $s$  and  $t$ . Subnet pairs for which the error in the selected path does not matter (either due to very low traffic volume or due to low-priority data traffic) can be assigned low values for  $D(s, t)$  or even zero. The generalized aggregate and weight selection problems incorporating the degrees of importance are then as follows.

- 1) **Generalized Aggregate Selection Problem:** Given a  $k$  and a weight assignment function  $W_{\mathcal{A}}$ , compute a complete set  $X \subseteq \mathcal{A}$  containing at most  $k$  aggregates such that  $\sum_{s, t \in \mathcal{S}} D(s, t) * (\text{lsp}(s, t, X, W_{\mathcal{A}}) - \text{lsp}(s, t))$  is minimum.
- 2) **Generalized Weight Selection Problem:** For an aggregate  $x \in A_i$ , compute a weight assignment function  $W_{\{x\}}$  such that  $\sum_{s, t \in \mathcal{S}, t \in x} D(s, t) * (\text{lsp}(s, t, \{x\}, W_{\{x\}}) - \text{lsp}(s, t))$  is minimum.

In all the problems outlined above, our goal is to minimize the (weighted) cumulative error across all source–destination subnet pairs. An alternative statement of the OSPF configuration problem aims to minimize the *maximum* error across all the source–destination subnets. The corresponding aggregate and weight selection problems can be formulated in a similar fashion, except that instead of minimizing  $\sum_{s, t \in \mathcal{S}}()$ , the objective is to minimize  $\max_{s, t \in \mathcal{S}}()$ .

### III. GENERALIZED AGGREGATE SELECTION PROBLEM

In this section, we present a dynamic programming algorithm for the generalized aggregate selection problem. Our algorithm exploits the fact that the containment structure of aggregates in  $\mathcal{A}$  is a set of trees (termed *aggregate trees*). We define the notion of error for each aggregate tree when certain aggregates in it are selected, and demonstrate that the cumulative error in the shortest-path computation when a subset  $X$  of aggregates is advertised is equal to the sum of the corresponding aggregate-tree errors when  $X$  is selected. We present our dynamic programming algorithm for selecting  $k$  aggregates in a single aggregate tree that minimize the tree's error in Section III-B. Section III-C then presents the algorithm for combining the results for the collection of aggregate trees to derive the final  $k$  aggregates that yield the minimum overall error. Finally, in Section III-E, we show how our algorithms can be extended for minimizing the maximum error in the OSPF routing paths.

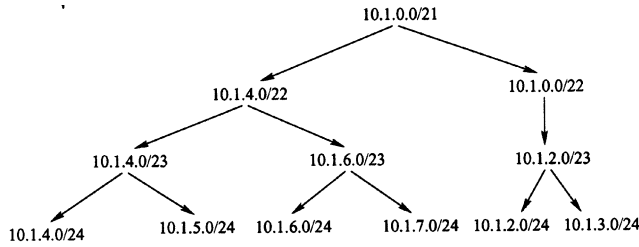


Fig. 2. Aggregate tree for eligible aggregates covering subnets in area 0.0.0.3.

### A. Aggregate Trees

The containment relationship among the eligible aggregates in  $\mathcal{A}$  naturally form a set of trees. For aggregates  $x, y \in \mathcal{A}$ , an edge from  $x$  to  $y$  exists if  $x$  covers  $y$  and every other aggregate  $z \in \mathcal{A}$  that covers  $y$  also covers  $x$ . Fig. 2 illustrates the aggregate tree for the eligible aggregates that cover subnets in Area 0.0.0.3 (from Fig. 1). Observe that since  $\mathcal{A}$  contains aggregates 10.1.4.0/22, 10.1.6.0/23, and 10.1.7.0/24, there is an edge from 10.1.6.0/23 to 10.1.7.0/24; however, there is no edge from 10.1.4.0/22 to 10.1.7.0/24 in the tree. Also, the internal nodes in the aggregate tree have either one or two children, but no more than two children. For instance, in Fig. 2, 10.1.0.0/22 has only one child since 10.1.0.0/23 does not exist in the network (subnet 10.1.1.0/24 does not exist in the network). Note that each leaf of an aggregate tree is a subnet in  $\mathcal{S}$ . Further, the root  $r(T)$  of tree  $T$  is basically an aggregate that is not covered by any other eligible aggregate.

We next define the error of a tree  $T$  when a set of aggregates in it have been selected. Suppose  $x$  is an aggregate in the tree  $T$ ,  $y$  is the most specific selected aggregate covering  $x$  in the tree, and  $X$  is the set of selected aggregates. Then the error  $E(x, y, X, W_{\mathcal{A}})$  of the subtree rooted at  $x$  is recursively defined as given in the equation at the bottom of the page. The error for an entire tree  $T$  with the set of selected aggregates  $X$  is then simply  $E(r(T), \epsilon, X, W_{\mathcal{A}})$  ( $\epsilon$  denotes the empty aggregate that does not cover any other aggregate). Note that each recursive invocation of  $E$  on  $x$ 's children propagates  $x$  as the most specific selected aggregate if  $x \in X$ . Consequently, whenever  $E$  is invoked for the subtree rooted at an aggregate  $x$ ,  $y$  is always the most specific selected aggregate covering  $x$ . As a result, the error of a tree is simply the sum of the errors of all the leaf subnets in it, where the error of a subnet  $t$  is  $\sum_{s \in \mathcal{S}} D(s, t) * (\text{lsp}(s, t, \{y\}, W_{\mathcal{A}}) - \text{lsp}(s, t))$ , where  $y$  is the most specific aggregate in  $X$  that covers  $t$ . Thus, since every subnet in  $\mathcal{S}$  is contained in one of the trees, the sum of errors of all the trees is essentially the cumulative error in the selected paths, which is the metric we are interested in minimizing.

Thus, we have reduced the aggregate selection problem to the problem of computing a set  $X$  containing at most  $k$  aggre-

gates such that the sum of the errors of all the trees is minimum. We break this into two subproblems which we address in Sections III-B and III-C. First, we present our dynamic programming algorithm to compute the optimal subset of aggregates for a single aggregate tree. Then, we show how to select a combination of aggregates from different trees (i.e., OSPF areas) that minimize the overall error for the entire collection of aggregate trees.

### B. Computing Optimal Aggregates for a Single Tree

We begin by presenting below a set of recursive equations for computing a tight lower bound on the error of a tree assuming that at most  $k$  arbitrary aggregates in the tree can be selected. The equations form the basis of our dynamic programming algorithm and can be used to compute the  $k$  best aggregates to select in order to minimize the error of the tree. Suppose that  $x$  is an aggregate in the tree  $T$  and  $y$  is the most specific aggregate in the tree covering  $x$  that has already been selected. Then, the minimum error  $\min E(x, y, k, W_{\mathcal{A}})$  of the subtree rooted at  $x$ , if we are allowed to choose at most  $k$  aggregates in the subtree, is as follows.

- If  $k = 0$ :  $\min E = E(x, y, \emptyset, W_{\mathcal{A}})$ .
- If  $k > 0$  and  $x$  has a single child  $u$ :

$$\min E = \min\{\min E(u, y, k, W_{\mathcal{A}}), \min E(u, x, k - 1, W_{\mathcal{A}})\}.$$

- If  $k > 0$  and  $x$  has children  $u, v$ :  $\min E$  is the minimum of

$$\begin{aligned} & \min_{0 \leq i \leq k} \{\min E(u, y, i, W_{\mathcal{A}}) + \min E(v, y, k - i, W_{\mathcal{A}})\} \\ & \min_{0 \leq i \leq k} \{\min E(u, y, i, W_{\mathcal{A}}) + \min E(v, x, k - i, W_{\mathcal{A}})\}. \end{aligned}$$

- If  $k > 0$  and  $x$  is a leaf:  $\min E = E(x, y, \{x\}, W_{\mathcal{A}})$ .

The intuition underlying the above set of equations is that if  $k = 0$ , then since no aggregates in the subtree can be selected, the minimum error is simply the error of the subtree when no aggregates in it are chosen. In case  $k > 0$  and  $x$  has children, and if  $X$  is the set of aggregates in the subtree rooted at  $x$  that if selected result in the minimum error, then the following hold for  $X$ : 1) either  $x \in X$  or  $x \notin X$  and 2) of the remaining aggregates in  $X$ ,  $i$  are in the subtree rooted at its left child and the remaining  $k - i$  or  $k - i - 1$  aggregates (depending on whether  $x \in X$ ) are in the subtree rooted at its right child. Thus, since the error of the subtree with  $x$  as root is simply the sum of the errors of its left and right subtrees, we can compute the minimum error for (the subtree rooted at)  $x$  by first computing the minimum error of its left and right subtrees for  $0 \leq i \leq k$  selected aggregates and for the cases when  $x$  is either selected or not selected, and then choosing the combination with the smallest error. Finally, if  $k > 0$  and  $x$  is a leaf, then there are only two possible alternatives for selecting aggregates in  $x$ 's subtree: either to select  $x$  or not

$$E = \begin{cases} \sum_{(x,u) \in T} E(u, y, X, W_{\mathcal{A}}), & \text{if } x \text{ has children and } x \notin X \\ \sum_{(x,u) \in T} E(u, x, X, W_{\mathcal{A}}), & \text{if } x \text{ has children and } x \in X \\ \sum_{s \in \mathcal{S}} D(s, x) * (\text{lsp}(s, x, \{y\}, W_{\mathcal{A}}) - \text{lsp}(s, x)), & \text{if } x \text{ is a leaf and } x \notin X \\ \sum_{s \in \mathcal{S}} D(s, x) * (\text{lsp}(s, x, \{x\}, W_{\mathcal{A}}) - \text{lsp}(s, x)), & \text{if } x \text{ is a leaf and } x \in X \end{cases}$$

to select  $x$ . The minimum error for these two cases then yields the desired minimum error. In the following, we formally prove that  $\min E$  is indeed the minimum error of the subtree rooted at  $x$  if at most  $k$  aggregates can be selected in it.

*Theorem 1:*  $\min E(x, y, k, w_A)$  is equal to the minimum of  $E(x, y, X, W_A)$ , where  $x$  is any arbitrary set containing at most  $k$  aggregates in the subtree rooted at  $x$ .

*Proof:* The proof is by induction on the number of aggregates in the subtree rooted at  $x$ .

*Basis:* Suppose that there is only one aggregate in the subtree rooted at  $x$ , that is,  $x$  is a leaf. In case  $k = 0$ , then the minimum error of the subtree is simply  $E(x, y, \emptyset, W_A)$  since none of the aggregates in the subtree can be selected. On the other hand, if  $k > 0$ , then the two possibilities are that either  $x \in X$  or  $x \notin X$  and the minimum error is the minimum of the error for these two cases.

*Induction Step:* Suppose that  $X$  is the set containing at most  $k$  aggregates in the subtree rooted at  $x$  that minimizes the error of the subtree. We show that  $\min E(x, y, k, W_A)$  is equal to  $E(x, y, X, W_A)$  for one case (other cases can be handled in a similar fashion). The case we consider is when  $x \in X$  and has two children  $u$  and  $v$  with  $U$  and  $V$  denoting the aggregates selected in the subtrees rooted at  $u$  and  $v$ , respectively. Note that if  $|U| = i$ , then  $|V| \leq k - i - 1$ . From the definition of error,  $E(x, y, X, W_A) = E(u, x, U, W_A) + E(v, x, V, W_A)$ . Thus, since  $X$  minimizes the error of the subtree rooted at  $x$ ,  $U$  and  $V$  must be the sets containing at most  $i$  and  $k - i - 1$  aggregates, respectively, and that minimizes the error of subtrees rooted at  $U$  and  $V$ , respectively. Due to the induction hypothesis,  $\min E(u, x, i, W_A) = E(u, x, U, W_A)$  and  $\min E(v, x, k - i - 1, W_A) = E(v, x, V, W_A)$ . Thus, since  $\min E(x, y, k, W_A) \leq \min E(u, x, i, W_A) + \min E(v, x, k - i - 1, W_A)$ , we can conclude that  $\min E(x, y, k, W_A) \leq E(x, y, X, W_A)$ . Note that  $\min E(x, y, k, W_A)$  cannot be less than  $E(x, y, X, W_A)$ , since this would lead to a contradiction because  $X$  would not minimize the error of the subtree rooted at  $x$ .  $\square$

From Theorem 1, it follows that  $\min E(r(T), \epsilon, k, W_A)$  returns the minimum possible error for a tree  $T$  when at most  $k$  aggregates in the tree can be selected. Procedure COMPUTEMINE in Fig. 3 uses dynamic programming to compute the  $k$  aggregates that result in the minimum possible error for the subtree rooted at  $x$  and  $y$  is the most specific aggregate covering  $x$  that has already been selected. The procedure is invoked with arguments that include the root aggregate of the tree  $r(T)$ ,  $\epsilon$  and  $k$ . The key ideas are similar to those described earlier for the computation of  $\min E$ , the minimum error for the tree. For instance, if an aggregate  $x$  has children, then procedure COMPUTEMINE recursively invokes itself for each of its children for the cases when  $x$  is selected and when  $x$  is not selected. Furthermore, in the case that  $x$  has two children, the procedure is invoked for each child for all the possibilities for the number of aggregates in each child subtree.

The only difference is that in addition to the minimum error, the procedure also computes the aggregates that are responsible for minimizing the tree error. Thus, every invocation of procedure COMPUTEMINE, in addition to returning the minimum error for the subtree rooted at  $x$ , also returns the set of aggregates in the subtree that cause the error to be minimum. This set

```

procedure COMPUTEMINE( $x, y, l$ )
1. if  $sT[x, y, l].C = \text{true}$ 
2.   return  $[sT[x, y, l].E, sT[x, y, l].A]$ 
3.  $\min E := \min E1 := \min E2 := \infty$ 
4. if  $x$  is a leaf {
5.    $\min E1 := \sum_{s \in S} D(s, x) * (lsp(s, x, \{y\}, W_A) - lsp(s, x))$ 
6.   if  $l > 0$ 
7.      $\min E2 := \sum_{s \in S} D(s, x) * (lsp(s, x, \{x\}, W_A) - lsp(s, x))$ 
8.   if  $\min E1 \leq \min E2$ 
9.      $[sT[x, y, l].E, sT[x, y, l].A] := [\min E1, \emptyset]$ 
10.  else
11.     $[sT[x, y, l].E, sT[x, y, l].A] := [\min E2, \{x\}]$ 
12.  }
13. if  $x$  has a single child  $u$  {
14.   $[\min E1, A1] := \text{COMPUTEMINE}(u, y, l)$ 
15.  if  $l > 0$ 
16.     $[\min E2, A2] := \text{COMPUTEMINE}(u, x, l - 1)$ 
17.  if  $\min E1 \leq \min E2$ 
18.     $[sT[x, y, l].E, sT[x, y, l].A] := [\min E1, A1]$ 
19.  else
20.     $[sT[x, y, l].E, sT[x, y, l].A] := [\min E2, A2 \cup \{x\}]$ 
21.  }
22. if  $x$  has children  $u$  and  $v$  {
23.  for  $i := 0$  to  $l$  {
24.     $[\min E1, A1] := \text{COMPUTEMINE}(u, y, i)$ 
25.     $[\min E2, A2] := \text{COMPUTEMINE}(v, y, l - i)$ 
26.    if  $\min E1 + \min E2 < \min E$ 
27.       $\min E := \min E1 + \min E2$ 
28.       $A := A1 \cup A2$ 
29.  }
30.  for  $i := 0$  to  $l - 1$  {
31.     $[\min E1, A1] := \text{COMPUTEMINE}(u, x, i)$ 
32.     $[\min E2, A2] := \text{COMPUTEMINE}(v, x, l - i - 1)$ 
33.    if  $\min E1 + \min E2 < \min E$ 
34.       $\min E := \min E1 + \min E2$ 
35.       $A := A1 \cup A2 \cup \{x\}$ 
36.  }
37.   $[sT[x, y, l].E, sT[x, y, l].A] := [\min E, A]$ 
38. }
39.  $sT[x, y, l].C := \text{true}$ 
40. return  $[sT[x, y, l].E, sT[x, y, l].A]$ 

```

Fig. 3. Dynamic programming algorithm for computing the aggregates that minimize tree error.

is derived by taking the union of the optimal aggregates for the subtrees rooted at its children, and adding  $\{x\}$  to it if selecting  $x$  is required for minimizing the error (Steps 11, 20, and 35). Note also that in order to improve computational efficiency, the optimal aggregates and the minimum error for the subtree rooted at  $x$  with  $y$  as the most specific aggregate and at most  $l$  selected aggregates are stored in  $sT[x, y, l].A$  and  $sT[x, y, l].E$ , respectively. The first invocation of COMPUTEMINE( $x, y, l$ ) causes the body of the procedure to be executed, but subsequent invocations simply return the previously computed and stored values.

### C. Combining the Aggregates for Set of Trees

Suppose there are  $m$  aggregate trees  $T_1, T_2, \dots, T_m$ . Further, let  $T_i[j].E, A$  denote the minimum error and the set of at most  $j$  aggregates in  $T_i$  responsible for minimizing  $T_i$ 's error. Then,  $X_i[j].E, A$ , the minimum error for the set of trees  $T_1, \dots, T_i$  and the  $j$  aggregates that minimize their cumulative error can be computed using the result of the following theorem.

*Theorem 2:* For the set of trees  $T_1, \dots, T_m$

$$\begin{aligned}
 & X_i[j].E \\
 &= \begin{cases} T_i[j].E, & \text{if } i = 1 \\ \min_{0 \leq l \leq j} \{X_{i-1}[l].E + T_i[j-l].E\}, & \text{otherwise.} \end{cases}
 \end{aligned}$$

```

procedure COMBINEMINE()
1. for  $i = 1$  to  $m$ 
2.   for  $j = 0$  to  $k$  {
3.      $T_i[j].E, A := \text{COMPUTEMINE}(r(T_i), \epsilon, j)$ 
4.      $X_i[j].E, A := [\infty, \emptyset]$ 
5.   }
6. for  $j = 0$  to  $k$ 
7.    $X_1[j].E, A := T_1[j].E, A$ 
8. for  $i = 2$  to  $m$ 
9.   for  $j = 0$  to  $k$ 
10.    for  $l = 0$  to  $j$ 
11.      if  $(X_{i-1}[l].E + T_i[j-l].E < X_i[j].E)$  {
12.         $X_i[j].E = X_{i-1}[l].E + T_i[j-l].E$ 
13.         $X_i[j].A = X_{i-1}[l].A \cup T_i[j-l].A$ 
14.      }

```

Fig. 4. Combining aggregates for a set of aggregate trees.

*Proof:* The proof is by induction on  $i$ .

*Basis:* When  $i = 1$ , the statement of the theorem clearly holds since  $X_1[j].E = T_1[j].E$ .

*Induction Step:* Suppose for the trees  $T_1, \dots, T_i$ , the cumulative error is minimum for  $j$  aggregates when  $l$  aggregates are selected from  $T_1, \dots, T_{i-1}$  and the remaining  $j-l$  aggregates are selected from  $T_i$ . Note that  $X_i[j].E \leq X_{i-1}[l].E + T_i[j-l].E$ . Thus, since due to the induction hypothesis,  $X_{i-1}[l].E$  is equal to the minimum error for the first  $i-1$  trees when  $l$  aggregates are selected from them, it follows that  $X_i[j].E$  is less than or equal to the minimum error for the first  $i$  trees when  $j$  aggregates are selected. In fact, since  $X_i[j].E$  is set to the minimum of  $X_{i-1}[l].E + T_i[j-l].E$  for  $0 \leq l \leq j$ , due to the induction hypothesis, it must be equal to the minimum error for the first  $i$  trees when  $j$  aggregates are selected.  $\square$

Procedure COMBINEMINE in Fig. 4 computes in  $X_m[k].E, A$  the minimum cumulative error and the  $k$  aggregates that minimize the error for the trees  $T_1, \dots, T_m$ . After computing the error and aggregates for each individual tree in Steps 1–5, in each iteration of Steps 8–14, the  $X_i[j]$ s are computed for increasing values of  $i$  based on the individual tree errors and the  $X_{i-1}[j]$ s computed in the previous iteration (as stated in Theorem 2). For each  $X_i[j]$ , the aggregates are computed by taking the union of the aggregates for the  $X_{i-1}[l]$  and  $T_i[j-l]$  that result in the minimum error for  $X_i[j]$ .

#### D. Time and Space Complexity

Suppose that  $d$  is the maximum depth of an aggregate tree, the number of aggregates in  $\mathcal{A}$  is  $N$  and the number of subnets in  $\mathcal{S}$  is  $n$ . Note that for 32-bit IP addresses,  $d \leq 32$ . Then the time complexity of the procedure COMPUTEMINE can be shown to be  $O(n^3 + Ndk^2)$ . The reason for this is that  $\text{lsp}(s, t)$ , the shortest path between subnets  $s$  and  $t$ , needs to be computed for all subnet pairs. The time complexity of this step is  $O(n^3)$ . Also, for each subnet  $x$  and every aggregate  $y$  covering it, one can precompute and store  $\sum_{s \in \mathcal{S}} D(s, x) * (\text{lsp}(s, x, \{y\}, W_{\mathcal{A}}) - \text{lsp}(s, x))$ , thus, enabling this information to be accessed in constant time.

Further, the body of COMPUTEMINE is executed at most once for each combination of  $x, y$ , and  $l$ . For a specific  $x$ , there are at most  $dk$  different possibilities for  $y$  and  $l$  for which the body of the procedure is executed. This is because  $y$  has to be an ancestor of  $x$  in the tree and  $l \leq k$ . Each execution of the body of COMPUTEMINE makes at most  $2l + 1$  recursive calls, and thus, since there are  $N$  possible aggregates, the total number of times COMPUTEMINE is invoked is  $O(Ndk^2)$ . As a result, the time complexity of procedure COMPUTEMINE is  $O(n^3 + Ndk^2)$ . Further, the space complexity of the procedure is  $O(n^2 + Ndk)$ ,  $O(n^2)$  to store the shortest path and error information for subnets, and  $O(Ndk)$  to store the error and aggregate values for each of the  $Ndk$  possible combinations of values for  $x, y$  and  $l$ .

It is fairly straightforward to observe that the three **for** loops spanning Steps 8–14 of procedure COMBINEMINE execute  $O(mk^2)$  steps. Thus, the overall time complexity of the procedure is  $O(n^3 + Ndk^2 + mk^2)$ , where the first two terms are the time complexity of computing the aggregates that minimize the error for the  $m$  trees. Note that even though COMBINEMINE makes independent successive invocations to  $\text{COMPUTEMINE}(r(T_i), \epsilon, j)$  for  $j = 0, \dots, k$ , the results computed in subTree during an invocation are shared between the invocations. The space complexity of procedure COMBINEMINE is simply  $O(mk)$  to store the  $X_i$  and  $T_i$  arrays.

#### E. Minimizing Maximum Error

Note that instead of minimizing the cumulative error over source destination subnet pairs, our algorithms can be adapted to minimize the maximum error over source destination pairs. In order to do this, we simply need to redefine the error of a tree to be the maximum error of the leaf subnets in it (instead of the sum of errors). Thus, the recursive definition of the error of the subtree rooted at  $x$  given that  $y$  is the most specific selected aggregate covering  $x$  in the tree and  $X$  is the set of selected aggregates, is as given in the equation shown at the bottom of the page.

Further, the minimum error of the subtree rooted at  $x$  if at most  $k$  aggregates in the subtree can be chosen (given that  $y$  is the most specific aggregate in the tree covering  $x$  and that has already been selected), is as follows.

- If  $k = 0$ :  $\min E = E(x, y, \emptyset, W_{\mathcal{A}})$ .
- If  $k > 0$  and  $x$  has a single child  $u$ :

$$\min E = \{\min E(u, y, k, W_{\mathcal{A}}), \min E(u, x, k-1, W_{\mathcal{A}})\}$$

- If  $k > 0$  and  $x$  has children  $u, v$ :  $\min E$  is the minimum of

$$\min_{0 \leq i \leq k} \{\max\{\min E(u, y, i, W_{\mathcal{A}}), \min E(v, y, k-i, W_{\mathcal{A}})\}\}$$

$$\min_{0 \leq i \leq k-1} \{\max\{\min E(u, x, i, W_{\mathcal{A}}), \min E(v, x, k-1-i, W_{\mathcal{A}})\}\}.$$

- If  $k > 0$  and  $x$  is a leaf:  $\min E = E(x, y, \{x\}, W_{\mathcal{A}})$ .

$$E = \begin{cases} \max_{(x,u) \in T} E(u, y, X, W_{\mathcal{A}}), & \text{if } x \text{ has children and } x \notin X \\ \max_{(x,u) \in T} E(u, x, X, W_{\mathcal{A}}), & \text{if } x \text{ has children and } x \in X \\ \max_{s \in \mathcal{S}} D(s, x) * (\text{lsp}(s, x, \{y\}, W_{\mathcal{A}}) - \text{lsp}(s, x)), & \text{if } x \text{ is a leaf and } x \notin X \\ \max_{s \in \mathcal{S}} D(s, x) * (\text{lsp}(s, x, \{x\}, W_{\mathcal{A}}) - \text{lsp}(s, x)), & \text{if } x \text{ is a leaf and } x \in X \end{cases}$$



Note that unlike the cumulative error case, where we were interested in the distributing the aggregates among the subtrees of  $x$  rooted at children  $u$  and  $v$  so that the sum of the errors of the subtrees was minimized, for the maximum error case, we are interested in minimizing the maximum of the errors of the two subtrees (since the error of the subtree rooted at  $x$  is the maximum of the errors of its two child subtrees). Thus, the following modifications need to be made to procedure COMPUTEMINE to compute the  $k$  aggregates that minimize the (maximum) error for the tree.

- 1) Replace  $\sum_{s \in \mathcal{S}} D(s, x) * (\text{lsp}(s, x, \{y\}, W_A) - \text{lsp}(s, x))$  in Steps 5 and 7 with  $\max_{s \in \mathcal{S}} D(s, x) * (\text{lsp}(s, x, \{y\}, W_A) - \text{lsp}(s, x))$ .
- 2) Replace  $\min E1 + \min E2$  in Steps 26, 27, 33, and 34 with  $\max\{\min E1, \min E2\}$ .

Similarly, the following simple modification to procedure COMBINEMINE enables it to compute the minimum error of a set of trees for the maximum error case: replace  $X_{i-1}[l].E + T_i[j - l].E$  in Steps 11 and 12 with  $\max\{X_{i-1}[l].E, T_i[j - l].E\}$ .

#### IV. WEIGHT SELECTION PROBLEM

In Section III, for a given weight assignment function  $W_A$ , we proposed algorithms for computing the optimal set of aggregates  $X$  for which the error in the selected paths is minimized. However, the final error and set of optimal aggregates  $X$  are very sensitive to the weight that a border router advertises for each aggregate. Thus, the weight assignment problem is important for ensuring that selected paths are of high quality, and is the subject of this section.

Recall that the weight assignment problem is to compute a weight assignment function  $W_{\{x\}}$  for a single aggregate  $x \in A_i$  such that the error in the selected paths from all subnets to destination subnets covered by  $x$  is minimized. The weight assignment function  $W_{\{x\}}$  assigns a weight to  $x$  at each ABR  $b \in B_i$ . Note that we are interested in computing the optimal weights for  $x$  under the assumption that no other aggregates covering subnets in  $x$  are concurrently being advertised. Also, since the aggregate  $x$  for whom we wish to compute weights is fixed, we drop the subscript  $\{x\}$  for  $W$ —thus, we will use  $W(b)$  to denote the weight assigned to  $x$  by ABR  $b \in B_i$ .

Intuitively, since  $W(b)$  is supposed to represent the distance between  $b$  and subnets covered by  $x$ , two possible logical choices for  $W(b)$  are the following:

- 1)  $\max_{t \in x} \{\text{lsp}(b, t)\}$ ;
- 2)  $(1/|x|) \sum_{t \in x} \text{lsp}(b, t)$ .

The first choice, recommended in [3], is simply the maximum distance of a subnet in aggregate  $x$  from the border router  $b$ , while the latter is the average distance of subnets in  $x$  from  $b$ . Note that since both choices are oblivious of the source subnets (not covered by  $x$ ) and the error to be minimized, as illustrated in the examples below, for most cases, neither choice optimizes our objective error function. In the following two examples, we show that choosing  $W(b)$  to be  $\max_{t \in x} \{\text{lsp}(b, t)\}$  minimizes neither the cumulative error nor the maximum error.

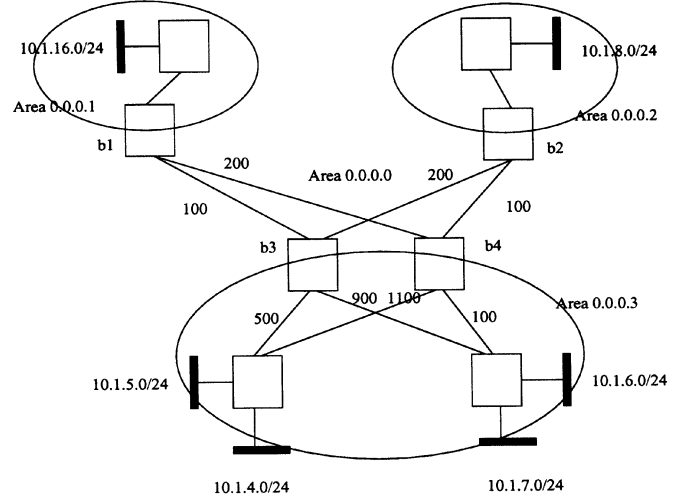


Fig. 5. Example of an AS where choosing maximum-distance weights does not minimize the maximum error.

*Example 3:* Consider the network in Fig. 1. Suppose we are interested in computing weights for the aggregate 10.1.0.0/21 that covers all the subnets in Area 0.0.0.3. If each border router chose the maximum distance to a subnet in 10.1.0.0/21 as the weight for it,  $b_3$  would assign 10.1.0.0/21 a weight of 1100 (distance of  $b_3$  from 10.1.6.0/24) and  $b_4$  would assign to 10.1.0.0/21 a weight of 1250 (distance between  $b_4$  and 10.1.6.0/24). Consequently, both subnets 10.1.16.0/24 and 10.1.8.0/24 select the path through ABR  $b_3$  to access the subnets in 10.1.0.0/21 which has a cumulative error of  $2 * 0 + 2 * 700 + 2 * 800 = 3000$  (for 10.1.16.0/24) and  $2 * 0 + 2 * 900 + 2 * 1000 = 4800$  (for 10.1.8.0/24). In contrast, assigning weights 1000 and 500 to 10.1.0.0/21 at ABRs  $b_3$  and  $b_4$ , respectively, causes the selected paths to be through  $b_4$  which results in much smaller cumulative errors of  $2 * 1300 + 2 * 0 + 2 * 0 = 2600$  (for 10.1.16.0/24) and  $2 * 1100 + 2 * 0 + 2 * 0 = 2200$  (for 10.1.8.0/24).  $\square$

*Example 4:* Consider the network in Fig. 5. Suppose we are interested in computing weights for the aggregate 10.1.0.0/21 that covers all the subnets in Area 0.0.0.3. If each border router chose the maximum distance to a subnet in 10.1.0.0/21 as the weight for it,  $b_3$  would assign 10.1.0.0/21 a weight of 900 (distance of  $b_3$  from 10.1.6.0/24) and  $b_4$  would assign to 10.1.0.0/21 a weight of 1100 (distance between  $b_4$  and 10.1.4.0/24). Consequently, both subnets 10.1.16.0/24 and 10.1.8.0/24 select the path through ABR  $b_3$  to access the subnets in 10.1.0.0/21 which has a maximum error of  $\max\{0, 700\} = 700$  (for 10.1.16.0/24) and  $\max\{0, 900\} = 900$  (for 10.1.8.0/24). In contrast, assigning weights 1000 and 500 to 10.1.0.0/21 at ABRs  $b_3$  and  $b_4$ , respectively, causes the selected paths to be through  $b_4$  which results in lower values for maximum error— $\max\{700, 0\} = 700$  (for 10.1.16.0/24) and  $\max\{500, 0\} = 500$  (for 10.1.8.0/24).

Choosing  $W(b) = (1/|x|) \sum_{t \in x} \text{lsp}(b, t)$  yields somewhat better results because intuitively this is more representative of the distance between  $b$  and subnets in  $x$  than  $\max_{t \in x} \{\text{lsp}(b, t)\}$ . As a matter of fact, setting  $W(b)$  to be the average distance of  $b$  to subnets in  $x$  can be shown to minimize the cumulative error for the weight selection problem. However, it does not minimize the maximum error, as illustrated by the example below.

*Example 5:* Consider the network in Fig. 1. Suppose we are interested in computing weights for the aggregate 10.1.0.0/21 that covers all the subnets in Area 0.0.0.3. If each border router chose the average distance to a subnet in 10.1.0.0/21 as the weight for it,  $b_3$  would assign 10.1.0.0/21 a weight of 730 and  $b_4$  would assign to 10.1.0.0/21 a weight of 570. Consequently, both subnets 10.1.16.0/24 and 10.1.8.0/24 select the path through ABR  $b_4$  to access the subnets in 10.1.0.0/21 which has a maximum error of 1300 for 10.1.16.0/24 and 1100 for 10.1.8.0/24. In contrast, assigning weights 500 and 1000 to 10.1.0.0/21 at ABRs  $b_3$  and  $b_4$ , respectively, causes the selected paths to be through  $b_3$  which results in lower values for maximum error—700 for 10.1.16.0/24 and 1000 for 10.1.8.0/24.

In Sections IV-A–E, we first show that selecting  $W(b) = (1/|x|) \sum_{t \in x} \text{lsp}(b, t)$  results in the minimum cumulative error and is a solution to the weight selection problem. However, the generalized weight selection problem that involves minimizing the product of the cumulative error of selected paths and their degrees of importance is an NP-hard problem [5]. Consequently, we present search-based heuristics to solve the generalized weight selection problem and a pseudopolynomial time algorithm to solve the weight selection problem when the objective is to minimize the maximum error.

#### A. Problem Formulation

In this subsection, we simplify some of the notation and introduce some new terminology that we need in order to address the weight selection problem, which is: For an aggregate  $x \in A_i$ , compute a weight assignment function  $W$  such that  $\sum_{s \in \mathcal{S}, t \in x} (\text{lsp}(s, t, \{x\}, W) - \text{lsp}(s, t))$  is minimum. For each source  $s$ , the selected paths to subnets covered by  $x$  is through the ABR  $b \in B_i$  for which  $\text{lsp}(s, b) + W(b)$  is minimum (among all the ABRs). We denote the ABR selected for source  $s$  by  $B(s, W)$ . Note that for  $t \in x$ ,  $\text{lsp}(s, t, \{x\}, W) = \text{lsp}(s, B(s, W)) + \text{lsp}(B(s, W), t)$ . Further, suppose  $e(s, b)$  denotes the error in the selected paths to subnets in  $x$  if ABR  $b$  is selected for source  $s$ . Thus,  $e(s, b) = \sum_{t \in x} \text{lsp}(s, b) + \text{lsp}(b, t) - \text{lsp}(s, t)$ . Then,  $e(s, B(s, W)) = \sum_{t \in x} \text{lsp}(s, t, \{x\}, W) - \text{lsp}(s, t)$ , and, thus, the weight selection problem becomes that of computing a weight assignment  $W$  such that  $\sum_{s \in \mathcal{S}} e(s, B(s, W))$  is minimum.

The above problem formulation is for minimizing the cumulative error. If we wish to minimize the maximum error, then  $e(s, b) = \max_{t \in x} \{\text{lsp}(s, b) + \text{lsp}(b, t) - \text{lsp}(s, t)\}$  and the weight assignment  $W$  must be such that  $\max_{s \in \mathcal{S}} e(s, B(s, W))$  is minimum.

#### B. Weight Selection Problem (Cumulative Error)

For the cumulative error case, it can be shown that choosing  $W(b)$  to be the average distance of  $b$  to subnets in  $x$  minimizes the cumulative error in the selected paths between sources and destination subnets in  $x$ .

*Theorem 3:* The weight assignment function  $W$  which assigns a weight  $W(b) = (1/|x|) \sum_{t \in x} \text{lsp}(b, t)$  to ABR  $b$  results in the minimum value for  $\sum_{s \in \mathcal{S}} e(s, B(s, W))$ .

*Proof:* Suppose the weight assignment function  $W$  assigns a weight  $W(b) = (1/|x|) \sum_{t \in x} \text{lsp}(b, t)$  to ABR  $b$ . For each source  $s$ ,  $\text{lsp}(s, B(s, W)) + W(B(s, W))$  is less than or equal to  $\text{lsp}(s, b) + W(b)$  for all ABRs  $b$ . Thus, first expanding  $W$  and then multiplying by constant  $|x|$  and subtracting constant  $\sum_{t \in x} \text{lsp}(s, t)$  from both sides, we get  $\sum_{t \in x} \text{lsp}(s, B(s, W)) + \text{lsp}(B(s, W), t) - \text{lsp}(s, t) \leq \sum_{t \in x} \text{lsp}(s, b) + \text{lsp}(b, t) - \text{lsp}(s, t)$  for all ABRs  $b$ . As a result, for every source  $s$ ,  $e(s, B(s, W)) \leq e(s, b)$  for all ABRs  $b$  and, thus,  $\sum_{s \in \mathcal{S}} e(s, B(s, W))$  is minimum for the weight assignment  $W$ .  $\square$

#### C. Generalized Weight Selection Problem (Cumulative Error)

For the cumulative error case,  $e(s, b) = \sum_{t \in x} \text{lsp}(s, b) + \text{lsp}(b, t) - \text{lsp}(s, t)$  is closely related to the criterion for selecting an ABR  $b$  for  $s$  which is that  $\text{lsp}(s, b) + W(b)$  is minimum (note that  $\sum_{t \in x} \text{lsp}(s, t)$  is a constant). However, for the generalized cumulative error case,  $e(s, b) = \sum_{t \in x} D(s, t) * (\text{lsp}(s, b) + \text{lsp}(b, t) - \text{lsp}(s, t))$  and, thus,  $e(s, b)$  can be any arbitrary value based on the value of  $D(s, t)$ . This fact that  $e(s, b)$  can be any arbitrary value makes the problem of computing a weight assignment function  $W$  that minimizes  $\sum_{s \in \mathcal{S}} e(s, B(s, W))$  intractable, as the following theorem demonstrates. The proof of the theorem involves a rather complex reduction from 3-SAT and, in the interest of space, has been omitted from the paper. Details of the proof, however, can be found in [4].

*Theorem 4:* For arbitrary values of  $e(s, b)$  and constant  $E$ , determining if there exists a weight assignment function  $W$  for which  $\sum_{s \in \mathcal{S}} e(s, B(s, W)) \leq E$  is NP-hard.

A simple iterative greedy search heuristic can be used to compute a weight assignment  $W$  that results in a low value for the cumulative error. The basic idea is to start with a set  $\mathcal{W}$  of random weight assignments. Then, in each subsequent iteration, for each  $W \in \mathcal{W}$ , from a number of candidate modifications, the one that minimizes the cumulative error is greedily chosen and applied to  $W$ . For a  $W \in \mathcal{W}$ , each candidate modification consists of adjusting the weight  $W(b)$  for a single ABR  $b \in B_i$ . Thus, for each weight assignment  $W$ , we are interested in computing the ABR  $b$  and a weight  $w$  such that setting  $W(b) = w$  (and leaving the weights for other ABRs unchanged) results in the smallest value for the cumulative error  $\sum_s e(s, B(s, W))$ . For  $n$  sources and  $m$  ABRs, this can be computed in  $O(mn(\log n + m))$  time as follows. First, we compute for each ABR  $b$  the weight  $w$  such that setting  $W(b) = w$  minimizes the error. Then, we choose from all the (ABR, weight) pairs  $(b, w)$  the one that results in the minimum error.

In order to compute the optimal weight  $w$  for an ABR  $b$ , we first compute for every source  $s$ , the ABR  $c$  in  $B_i - \{b\}$  for which  $\text{lsp}(s, c) + W(c)$  is minimum (this can be achieved in  $O(mn)$  steps). For the source  $s$ , let  $v(s) = \text{lsp}(s, c) + W(c) - \text{lsp}(s, b)$  and  $e(s) = e(s, c)$ . Suppose that  $v(s_1), \dots, v(s_n)$  are the values for the sources in sorted (increasing) order (sorting the  $n$  values takes  $O(n \log n)$  steps). Also, let  $v(s_{n+1}) = \infty$ . Then, choosing a value for  $W(b)$ , such that  $v(s_{j-1}) < W(b) \leq v(s_j)$  causes the cumulative error to be  $\sum_{l=1}^{j-1} e(s_l) + \sum_{l=j}^n e(s_l, b)$  (since ABR  $b$  is selected for sources  $s_j, \dots, s_n$  when  $W(b) \leq v(s_j)$ ). Thus, in a single pass over the sequence  $v(s_1), \dots, v(s_n)$ , the optimal

```

procedure COMPUTEWEIGHTSMAX( $Q$ )
1. for each  $b \in B_i$ , set  $W_{old}(b) := 0$ 
2. while  $(\sum_{b \in B_i} W_{old}(b) \leq (\frac{|B_i| * (|B_i| - 1)}{2}) * lsp_{max})$  {
3.   Let  $Q'$  be a new set of inequalities that result when the value  $W_{old}(b)$  is
   substituted for each variable  $W(b)$  only on the LHS of each inequality in  $Q$ 
4.   for each  $b \in B_i$ , set  $W_{new}(b)$  to the smallest possible value such that each
   inequality in  $Q'$  is satisfied when  $W_{new}(b)$  is substituted for variable  $W(b)$  in  $Q'$ 
5.   if  $W_{new} = W_{old}$ 
6.     return  $W_{new}$ 
7.   else
8.      $W_{old} := W_{new}$ 
9.   }
10. return “there does not exist a weight assignment  $W$ ”
    
```

Fig. 6. Algorithm for computing weights for two ABRs.

weight  $v(s_j)$  for  $W(b)$  which minimizes the cumulative error  $\sum_{l=1}^{j-1} e(s_l) + \sum_{l=j}^n e(s_l, b)$  can be computed. Thus, the optimal weight for each ABR can be computed in  $O(n(\log n + m))$  steps and for all ABRs, in  $O(mn(\log n + m))$  time.

To recap, the greedy heuristic, in each iteration, modifies each weight assignment  $W \in \mathcal{W}$  by setting  $W(b) = w$ , where ABR  $b$  and weight  $w$  result in the minimum error for  $W$  and are computed as described above. It terminates the search computation either after a fixed number of iterations or if the improvement in cumulative error during an iteration due to modifying every  $W \in \mathcal{W}$  drops below an error threshold  $\epsilon$ .

#### D. Generalized Weight Selection Problem (Maximum Error)

Recall that if we are interested in minimizing the maximum error, then  $e(s, b) = \max_{t \in x} \{lsp(s, b) + lsp(b, t) - lsp(s, t)\}$  and the weight assignment  $W$  must be such that  $\max_{s \in \mathcal{S}} e(s, B(s, W))$  is minimum. Thus, we can employ an algorithm similar to the greedy search heuristic described earlier to compute a weight assignment function that minimizes the maximum error (instead of the cumulative error).

However, if we assume that weight assignments  $W$  and shortest path distances  $lsp$  to be nonnegative integers, then we can devise a more efficient pseudopolynomial time algorithm for computing the weight assignment that minimizes the maximum error. Suppose we could devise a procedure  $P$  that computes a weight assignment  $W$  (if there exists one) such that  $\max_{s \in \mathcal{S}} e(s, B(s, W)) \leq E$  for some constant  $E$ . Then, a simple procedure for computing the weight assignment that minimizes the maximum error is as follows:

- 1) sort the errors  $e(s, b)$  between (source, ABR) pairs—let  $E_1, \dots, E_r$  be the errors in order of increasing value;
- 2) repeatedly invoke the procedure  $P$  for increasing values of  $i$ , until  $P$  returns a weight assignment  $W$  for which  $\max_{s \in \mathcal{S}} e(s, B(s, W)) \leq E_i$ .

Thus,  $E_i$  is the smallest value for which a weight assignment exists and represents the minimum possible value for the maximum error. Further,  $W$  is the weight assignment that minimizes the maximum error. Note that instead of considering each  $E_i$  sequentially, one can also use a binary search procedure to compute the minimum value for the maximum error more efficiently.

Thus, the crucial task for us is to develop the procedure  $P$  that computes a weight assignment  $W$  (if there exists one) such that  $\max_{s \in \mathcal{S}} e(s, B(s, W)) \leq E$  for some constant  $E$ . We show that the problem of computing a  $W$  such that the maximum error is

at most  $E$  is equivalent to solving a set of inequalities involving the  $W(b)$ s as variables. For a source  $s$ , let  $R(s)$  denote the set of ABRs  $b \in B_i$  for which  $e(s, b) \leq E$ —thus, for the remaining ABRs  $b \in B_i - R(s)$ ,  $e(s, b) > E$ . Consequently, since the error for each source can be at most  $E$ , the computed  $W$  must be such that one of the ABRs in  $R(s)$  is selected for  $s$ . For this, we require  $W$  to satisfy the following set of inequalities for all  $c \in B_i - R(s)$ :

$$\min_{b \in R(s)} \{W(b) + lsp(s, b)\} \leq W(c) + lsp(s, c).$$

Thus, for each source  $s$ , we obtain the set of inequalities described above.<sup>1</sup> Note that the  $W$ 's in the equations are variables and the  $lsp$ 's are constants. Also, for each ABR  $b$ ,  $W$  also needs to satisfy the constraint  $0 \leq W(b)$ . Suppose  $Q$  denotes the set of inequalities over all the sources and ABRs. It is straightforward to observe that for a  $W$  the maximum error is  $E$  if and only if  $W$  is a solution for the set of equations  $Q$ . Thus, we simply need to focus on computing a  $W$  that satisfies the inequalities in  $Q$ . Observe that if for a source  $s$ , the set  $R(s)$  is empty, then there does not exist a  $W$  for which the set of inequalities  $Q$  is satisfiable. The reason for this is that for the source, we obtain inequalities of the form  $\min\{\} \leq W(c) + lsp(s, c)$  which cannot be satisfied since  $\min\{\} = \infty$ . Also, no equations are generated for a source  $s$  if  $R(s) = B_i$  (that is, the error for the source  $s$  is at most  $E$  irrespective of the chosen ABR).

Procedure COMPUTEWEIGHTSMAX in Fig. 6 is an iterative pseudopolynomial time algorithm for computing a  $W$  that satisfies  $Q$ . In the procedure,  $W_{new}$  and  $W_{old}$  store the weight assignment values prior to and after each iteration. In each iteration, a new weight assignment  $W_{new}$  is computed after substituting the previous weight assignment  $W_{old}$  for the  $W$ 's only on the LHS of each inequality in  $Q$  (Steps 3–4). (We use LHS and RHS to denote the left- and right-hand side of an inequality, respectively.) Note that each inequality in  $Q'$  has the form  $C \leq W(c) + lsp(s, c)$ , where  $C$  and  $lsp(s, c)$  are constants and  $W(c)$  is a variable. Also,  $lsp_{max}$  is the maximum value for  $lsp(s, b)$  for a (source, ABR) pair.

<sup>1</sup>We assume that if for two ABRs  $b$  and  $c$ ,  $W(b) + lsp(s, b) = W(c) + lsp(s, c)$ , then the ABR with a smaller error is selected for  $s$ . In case this assumption does not hold, the following stronger inequality can be employed for all  $c \in B_i - R(s)$ :

$$\min_{b \in R(s)} \{W(b) + lsp(s, b)\} + 1 \leq W(c) + lsp(s, c).$$

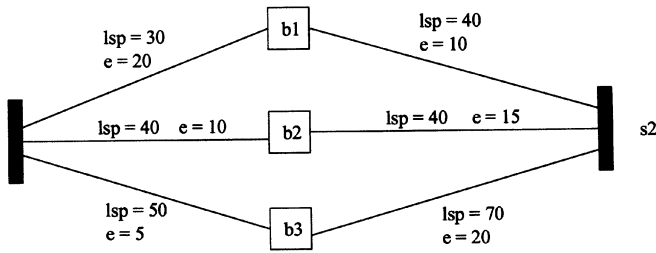


Fig. 7. Example network for tracing execution of procedure COMPUTEWEIGHTSMAX.

*Example 6:* Consider sources  $s_1$  and  $s_2$  and ABRs  $b_1$ ,  $b_2$ , and  $b_3$ , shown in Fig. 7. Let the functions  $\text{lsp}$  and  $e$  for the (source, ABR) pairs be as depicted in the figure. For  $E = 10$ ,  $R(s_1) = \{b_2, b_3\}$  and  $R(s_2) = \{b_1\}$ . Thus, for source  $s_1$ ,  $Q$  contains the following inequality:

$$\min\{W(b_2) + 40, W(b_3) + 50\} \leq W(b_1) + 30. \quad (1)$$

And for source  $s_2$ ,  $Q$  contains the following two inequalities:

$$W(b_1) + 40 \leq W(b_2) + 40 \quad (2)$$

$$W(b_1) + 40 \leq W(b_3) + 70. \quad (3)$$

Note that since weight assignments cannot be negative, even though we do not explicitly state this,  $Q$  and  $Q'$  always contain the following three constraints:  $0 \leq W(b_1)$ ,  $0 \leq W(b_2)$ , and  $0 \leq W(b_3)$ .

We now trace the execution of COMPUTEWEIGHTSMAX for the above set of inequalities in  $Q$ . Initially,  $W_{\text{old}}(b_i)$  is set to 0 for all three ABRs. In the first iteration, substituting 0 for all the  $W(b_i)$  variables on the LHS of inequalities (1)–(3) results in the following set of equations  $Q'$  (Step 3 of the procedure):  $\min\{10, 20\} \leq W(b_1)$  [due to inequality (1)],  $0 \leq W(b_2)$  [due to inequality (2)], and  $-30 \leq W(b_3)$  [due to inequality (3)]. As a result,  $W_{\text{new}}(b_1)$  is set to 10, while  $W_{\text{new}}$  for  $b_2$  and  $b_3$  continue to be 0. At the beginning of the second iteration, thus,  $W_{\text{old}}(b_1) = 10$ . Consequently, after substitution of  $W_{\text{old}}(b_i)$  for the  $W(b_i)$  variables on the LHS of inequalities in  $Q$ ,  $Q'$  contains equations  $\min\{10, 20\} \leq W(b_1)$ ,  $10 \leq W(b_2)$ , and  $-20 \leq W(b_3)$ . This causes  $W_{\text{old}}(b_1)$  and  $W_{\text{old}}(b_2)$  to be set to 10 and 10, respectively, at the end of the iteration (in Step 8). Similarly, it can be shown that at the end of the third iteration,  $W_{\text{old}}(b_1)$  and  $W_{\text{old}}(b_2)$  are set to 20 and 10, respectively, and during the fourth iteration,  $W_{\text{old}}(b_1)$  and  $W_{\text{old}}(b_2)$  are set to 20 and 20, respectively. In the fifth and final iteration, equations in  $Q'$  after substitution are  $20 \leq W(b_1)$ ,  $20 \leq W(b_2)$ , and  $-10 \leq W(b_3)$ , causing  $W_{\text{new}}$  to be equal to  $W_{\text{old}}$ . Thus, in the final weight assignment returned by COMPUTEWEIGHTSMAX,  $W(b_1) = 20$ ,  $W(b_2) = 20$ , and  $W(b_3) = 0$ .

In the following, we show that COMPUTEWEIGHTSMAX returns a  $W$  that is a solution to  $Q$  if and only if  $Q$  is satisfiable. In order to show this, in the following lemmas, we show that: 1) for any  $W$  that satisfies  $Q$ ,  $W_{\text{old}}(b) \leq W(b)$  and 2)

if  $Q$  is satisfiable, then there exists a  $W$  that is a solution to  $Q$  for which  $\sum_{b \in B_i} W(b) \leq (|B_i| * (|B_i| - 1)/2) * \text{lsp}_{\text{max}}$ . Thus, since  $W_{\text{old}}$  does not decrease between successive iterations and the procedure terminates only when a  $W$  is found or  $\sum_{b \in B_i} W_{\text{old}}(b)$  becomes greater than  $(|B_i| * (|B_i| - 1)/2) * \text{lsp}_{\text{max}}$ , COMPUTEWEIGHTSMAX computes  $W$  correctly. In the proofs,  $W_{\text{old}}^j$  denote the value of  $W_{\text{old}}$  at the end of the  $j$ th iteration. Also,  $Q'_j$  denotes the set of inequalities in  $Q$  when  $W_{\text{old}}^j$  is substituted for all variables on the LHS of each inequality in  $Q$ .

*Lemma 1:* For every ABR  $b \in B_i$ ,  $W_{\text{old}}^{j+1}(b) \geq W_{\text{old}}^j(b)$ .

*Proof:* The proof is by induction on  $j$ .

*Basis:* Clearly, the lemma holds for  $j = 0$ . At the end of the zeroth iteration (that is, initially),  $W_{\text{old}}^0(b) = 0$  for every ABR  $b \in B_i$ . Since  $Q$  contains the inequality  $0 \leq W(b)$  for every ABR  $b$ ,  $W_{\text{new}}(b) \geq 0$  and, thus,  $W_{\text{old}}^1(b) \geq 0$ .

*Induction Step:* Consider a  $j > 0$ . Since due to the induction hypothesis,  $W_{\text{old}}^j(b) \geq W_{\text{old}}^{j-1}(b)$ , it follows that the LHS of each inequality in  $Q'_{j-1}$  is less than or equal to the LHS of the corresponding inequality in  $Q'_j$ . Thus, since  $W_{\text{old}}^j(b)/W_{\text{old}}^{j+1}(b)$  is the smallest possible value for  $W(b)$  for which the inequalities in  $Q'_{j-1}/Q'_j$  with  $W(b)$  on the RHS are satisfied, it follows that  $W_{\text{old}}^{j+1}(b) \geq W_{\text{old}}^j(b)$ .  $\square$

*Lemma 2:* For every weight assignment  $W'$  that is a solution to  $Q$ ,  $W'(b) \geq W_{\text{old}}(b)$ , for every ABR  $b \in B_i$ .

*Proof:* We use induction to show that the lemma holds at the end of each iteration  $j$ .

*Basis:* Clearly, the lemma holds for  $j = 0$ . At the end of the zeroth iteration (that is, initially),  $W_{\text{old}}^0(b) = 0$  for every ABR  $b \in B_i$ . Thus, since  $W'(b) \geq 0$ , the lemma holds at the end of the 0 iteration.

*Induction Step:* Consider a  $j > 0$ . At the start of the  $j$  iteration, due to the induction hypothesis, we have that  $W'(b) \geq W_{\text{old}}^{j-1}(b)$ . Suppose that  $Q'$  denotes the set of inequalities in  $Q$  when all variables on the LHS of every inequality in  $Q$  are substituted with  $W'$ . Due to the induction hypothesis, it follows that the LHS of each inequality in  $Q'_{j-1}$  is less than or equal to the LHS of the corresponding inequality in  $Q'$ . Thus, since  $W_{\text{old}}^j(b)/W'(b)$  is the smallest possible value for  $W(b)$  for which the inequalities in  $Q'_{j-1}/Q'$  with  $W(b)$  on the RHS are satisfied, it follows that  $W'(b) \geq W_{\text{old}}^j(b)$ .  $\square$

*Lemma 3:* If  $Q$  is satisfiable, then there exists a weight assignment  $W'$  that is a solution to  $Q$  and for which  $\sum_{b \in B_i} W'(b) \leq (|B_i| * (|B_i| - 1)/2) * \text{lsp}_{\text{max}}$ .

*Proof:* Let  $W'$  be a weight assignment that is a solution for  $Q$  with the smallest value for  $\sum_{b \in B_i} W'(b)$  and further suppose that  $\sum_{b \in B_i} W'(b) > (|B_i| * (|B_i| - 1)/2) * \text{lsp}_{\text{max}}$ . Suppose  $v$  is the smallest value for  $W'(b)$  among all the ABRs. Then,  $v$  must be zero since the weight assignment  $W''$  where  $W''(b) = W'(b) - v$  is also a solution to  $Q$  and  $\sum_{b \in B_i} W''(b) < \sum_{b \in B_i} W'(b)$ —which leads to a contradiction. Without loss of generality, let  $W'(b_1), \dots, W'(b_{|B_i|})$  be the weights for ABRs sorted in increasing order. Thus,  $W'(b_1) = 0$ . Since  $\sum_{b \in B_i} W'(b) > (|B_i| * (|B_i| - 1)/2) * \text{lsp}_{\text{max}}$ , it must be the case that for a pair of consecutive ABRs  $b_l$  and  $b_{l+1}$ ,  $W'(b_{l+1}) - W'(b_l) > \text{lsp}_{\text{max}}$ . We show that  $W''$  is a solution for  $Q$ , where  $W''(b_j) = W'(b_j)$ , for  $1 \leq j \leq l$ , and

$W''(b_j) = W'(b_j) - 1$ , for  $l + 1 \leq j \leq |B_i|$ . However, this leads to a contradiction since  $\sum_{b \in B_i} W''(b) < \sum_{b \in B_i} W'(b)$ .

In order to show that  $Q$  is satisfiable for  $W''$ , we need to consider the following two cases for each inequality.

- 1) Variable  $W(b_j)$  is on the RHS of the inequality ( $1 \leq j \leq l$ ). In this case, the value of the RHS of the inequality is identical for both  $W'$  and  $W''$ , while the value of the LHS for  $W''$  is less than or equal to the LHS for  $W'$ . Thus, if the inequality was satisfiable for  $W'$ , it is also satisfiable for  $W''$ .
- 2) Variable  $W(b_j)$  is on the RHS of the inequality ( $l+1 \leq j \leq |B_i|$ ). In this case, the value of the RHS of the inequality decreases by one for  $W''$  compared to  $W'$ . If the value of the LHS of the inequality also decreases by one for  $W''$  compared to  $W'$ , then since the inequality is satisfiable for  $W'$ , it is also satisfiable for  $W''$ .

The other case we need to consider is when the value of the LHS of the inequality is the same for  $W'$  and  $W''$ . This corresponds to the case when the minimum value of the LHS (for both  $W'$  and  $W''$ ) is due to the term  $W(b_k) + \text{lsp}(s, b_k)$  for some  $1 \leq k \leq l$ . Since the inequality holds for  $W'$ ,  $W'(b_k) + \text{lsp}(s, b_k) \leq W'(b_j) + \text{lsp}(s, b_j)$ . Further, since  $k \leq l < j$ ,  $W'(b_j) - W'(b_k) > \text{lsp}_{\max}$ . Rearranging terms,  $W'(b_j) > W'(b_k) + \text{lsp}_{\max}$ , or alternately,  $W'(b_j) > W'(b_k) + \text{lsp}(s, b_k)$  (since  $\text{lsp}(s, b_k) \leq \text{lsp}_{\max}$ ). Thus,  $W'(b_j) + \text{lsp}(s, b_j) > W'(b_k) + \text{lsp}(s, b_k)$  and  $W'(b_j) + \text{lsp}(s, b_j) - 1 \geq W'(b_k) + \text{lsp}(s, b_k)$ . However, since  $W''(b_j) = W'(b_j) - 1$  and  $W''(b_k) = W'(b_k)$ , the inequality holds for  $W''$ .

Thus,  $W''$  is a solution for  $Q$ . However, this leads to a contradiction since  $\sum_{b \in B_i} W''(b) < \sum_{b \in B_i} W'(b)$ , and, thus, it must be the case that  $\sum_{b \in B_i} W'(b) \leq (|B_i| * (|B_i| - 1)/2) * \text{lsp}_{\max}$ .  $\square$

Clearly, if a weight assignment  $W$  is returned by procedure COMPUTEWEIGHTSMAX, then this is a solution to  $Q$ . The reason for this is that for the returned  $W$ ,  $W = W_{\text{new}} = W_{\text{old}}$ —thus, when the value for  $W(b)$  returned by the procedure is substituted for the occurrence of variable  $W(b)$  in each inequality, every inequality is satisfied. However, we also need to show that our procedure finds a solution for  $Q$  if one exists.

*Theorem 5:* If  $Q$  is satisfiable, then procedure COMPUTEWEIGHTSMAX returns a weight assignment  $W$  that is a solution to  $Q$ .

*Proof:* Suppose  $Q$  is satisfiable. Due to Lemma 3, it follows that there exists a  $W$  that is a solution to  $Q$  such that  $\sum_{b \in B_i} W(b) \leq (|B_i| * (|B_i| - 1)/2) * \text{lsp}_{\max}$ . Further, due to Lemma 2,  $W_{\text{old}}(b) \leq W(b)$  and so for  $W_{\text{old}}$ ,  $\sum_{b \in B_i} W_{\text{old}}(b) \leq (|B_i| * (|B_i| - 1)/2) * \text{lsp}_{\max}$ . Also, due to Lemma 1,  $\sum_{b \in B_i} W_{\text{old}}(b)$  at the end of an iteration is greater than or equal to its value at the end of the previous iteration. Thus, since  $\sum_{b \in B_i} W_{\text{old}}(b)$  cannot exceed  $(|B_i| * (|B_i| - 1)/2) * \text{lsp}_{\max}$ , at some point during the execution of the procedure, the value of  $W_{\text{old}}$  between two consecutive iterations does not change. This weight assignment  $W$  is returned by the procedure and is a solution to  $Q$ . The reason for this is that for the returned  $W$ ,  $W = W_{\text{new}} = W_{\text{old}}$ —thus, when the value for  $W(b)$  returned by the procedure is substituted for the occurrence of variable  $W(b)$  in each inequality, every inequality is satisfied.  $\square$

The worst case time complexity of the overall procedure to compute a  $W$  that minimizes the maximum error (by repeatedly invoking COMPUTEWEIGHTSMAX for error values in  $E_1, \dots, E_r$ ,  $r = mn$ ) can be shown to be  $O(m^4 \log(mn) \text{lsp}_{\max})$ , where  $m = |B_i|$  is the number of ABRs and  $n$  is the number of sources. The  $\log(mn)$  term is due to the binary search over all the errors to determine the minimum error for which a weight assignment  $W$  can be computed by COMPUTEWEIGHTSMAX. In the worst case, COMPUTEWEIGHTSMAX performs  $O(m^2 \text{lsp}_{\max})$  iterations since the procedure terminates when the sum of the weights become  $(|B_i| * (|B_i| - 1)/2) * \text{lsp}_{\max}$  and the sum increases by at least one in each iteration. Finally, the time complexity for computing  $W_{\text{new}}$  in each iteration is  $O(m^2)$  since in the worst case, the number of inequalities in  $Q$  is  $m^2$  (one inequality for every pair of  $W$  variables, one on the LHS and the other on the RHS of the inequality).

## V. CONCLUDING REMARKS

Address aggregation within OSPF areas is critical for scalability since it can result in significant reductions in routing table sizes, smaller link-state databases, and less network traffic to synchronize the router link-state databases. However, address aggregation can also lead to the selection of *suboptimal* OSPF routing paths between source–destination subnet pairs that span different areas. In this paper, we addressed the important practical problem of configuring OSPF aggregates at ABRs to minimize the error in OSPF shortest-path computations due to subnet aggregation. We first developed an optimal dynamic programming algorithm that, given an upper bound  $k$  on the number of aggregates to be advertised by the ABRs and a weight assignment function for the aggregates, computes the  $k$  aggregates that result in the minimum cumulative/maximum error in the shortest-path computations for all source–destination subnet pairs. Subsequently, we tackled the problem of assigning weights to OSPF aggregates such that the cumulative/maximum error in the computed shortest paths is minimized. We showed that, while for certain special cases (e.g., unweighted cumulative error) efficient optimal algorithms for the weight assignment problem can be devised, the general problem itself is NP-hard. We proposed a randomized search strategy for the general case of weighted cumulative error. To the best of our knowledge, our work is the first to carry out a systematic study of the algorithmic issues underlying the configuration of OSPF aggregates and to propose efficient configuration algorithms that are *provably optimal* for many practical scenarios.

While the problem of selecting OSPF aggregates is simplified if IP addresses are systematically assigned to carefully designed networks and areas, such an approach may not work in practice since networks are seldom static and tend to evolve continuously due to the addition of new network elements, subnets, and links. Our proposed algorithms allow aggregates to be optimally configured in such continuously changing dynamic networks, and can also work on-line: emergence of spurious subnets (through external advertisements), link failures, and topology changes

would trigger the computation of changes to the optimal set of advertised aggregates/weights. This computation can also be carried out incrementally and efficiently.

#### ACKNOWLEDGMENT

The authors would like to thank M. Yannakakis for pointing out the optimality of the weight assignment that assigns to each aggregate the average of the weights in subnets covered by it.

#### REFERENCES

- [1] B. Fortz and M. Thorup, "Internet traffic engineering by optimizing OSPF weights," in *Proc. IEEE INFOCOM*, Tel Aviv, Israel, Mar. 2000, pp. 519–528.
- [2] C. Huitema, *Routing in the Internet*, 2nd ed. Englewood Cliffs, NJ: Prentice-Hall, 2000.
- [3] J. T. Moy, *OSPF Anatomy of an Internet Routing Protocol*. Reading, MA: Addison-Wesley, 1998.
- [4] Y. Breitbart, M. Garofalakis, A. Kumar, and R. Rastogi, "Optimal Configuration of OSPF Aggregates," Bell Labs., Tech. Memo., Aug. 2000.
- [5] M. R. Garey and D. S. Johnson, *Computers and Intractability: A Guide to the Theory of NP-Completeness*. San Francisco, CA: Freeman, 1979.



**Rajeev Rastogi** received the B.Tech. degree in computer science from the Indian Institute of Technology, Bombay, India, in 1988, and the M.S. and Ph.D. degrees in computer science from the University of Texas, Austin, in 1990 and 1993, respectively.

He is currently the Director of the Internet Management Research Department, Bell Laboratories, Lucent Technologies, Murray Hill, NJ. He joined Bell Laboratories in 1993 and became a Distinguished Member of Technical Staff in 1998.

He is active in the field of databases and has served as a Program Committee Member for several conferences in that area. His work has appeared in many professional conferences and journals. His research interests include database systems, network management, storage systems, and knowledge discovery. His most recent research has focused on the areas of network topology discovery, monitoring, configuration and provisioning, data mining, and high-performance transaction systems.

Dr. Rastogi currently serves on the Editorial Board of IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING.



**Yuri Breitbart** received the D.Sc. degree in computer science from the Technion–Israel Institute of Technology, Haifa, Israel.

He is currently the Ohio Board of Regents Distinguished Professor of Computer Science with the Department of Computer Science, Kent State University, Kent, OH. He joined the Department in January 2002. From 1996 to 2002, he was a Member of Technical Staff and Distinguished Member of Technical Staff at Bell Laboratories. From 1986 to 1996, he was a Professor with the Department of Computer Science, University of Kentucky, Lexington, where he was the Department Chair for seven years. Prior to 1986, he led the Database Research group, first at ITT Research Center (1979–1981) and then at Amoco Production Company Research Center (1981–1986). From 1975 to 1979, he was an Assistant and then Associate Professor in the Departments of Computer Science at the State University of New York, Albany, and the University of Wisconsin, Milwaukee. He has served on numerous program committees and National Science Foundation panels. His research focuses on distributed information system, network management systems, and bio-informatics.

Dr. Breitbart is a Fellow of the Association for Computing Machinery (ACM) and a member of the IEEE Computer Society and the ACM SIGMOD.



**Minos Garofalakis** (S'94–A'98) received the B.Sc. degree in computer engineering and informatics from the University of Patras (UOPCEID), Patras, Greece, in 1992, and the M.Sc. and Ph.D. degrees in computer science from the University of Wisconsin, Madison, in 1994 and 1998, respectively.

In 1998, he joined Bell Laboratories, Lucent Technologies, Murray Hill, NJ, where he is currently a Member of Technical Staff with the Internet Management Research Department. His current research interests are in the areas of network management, data streaming, approximate query processing, data mining, and XML databases.

Dr. Garofalakis is a member of the Association for Computing Machinery (ACM), and has served as a Program Committee Member for several conferences in the database area, including ACM SIGMOD, VLDB, ACM SIGKDD, and IEEE ICDE.



**Amit Kumar** received the B.Tech. degree in computer science from the Indian Institute of Technology, Kanpur, India, in 1997 and the Ph.D. degree in computer science from Cornell University, Ithaca, NY, in 2002.

He is currently a Member of Technical Staff with the Internet Management Research Department, Bell Laboratories, Lucent Technologies, Murray Hill, NJ. His research interests are in theoretical computer science, approximation algorithms, and designing efficient algorithms for network management.