

An Adaptive Management Approach to Resolving Policy Conflicts

SELMA YILMAZ

IBRAHIM MATTA

Computer Science Department
Boston University
Boston, MA 02215, USA

{selma,matta}@cs.bu.edu

Technical Report BUCS-TR 2006-008

Abstract

The Border Gateway Protocol (BGP) is the current inter-domain routing protocol used to exchange reachability information among Autonomous Systems (ASes) in the Internet. BGP supports policy-based routing which allows each AS to independently define a set of local policies regarding which routes it accepts and advertises from/to other networks, as well as which route it prefers when more than one route becomes available. However, independently chosen local policies may cause global conflicts, which result in protocol divergence. In this paper, we propose a new algorithm, called Adaptive Policy Management (APM), to resolve policy conflicts in a distributed manner. Akin to distributed feedback control systems, each AS independently classifies the state of the network as either conflict-free or potentially conflicting by observing its local history only (namely, route flaps). Based on the degree of measured conflicts, each AS dynamically adjusts its own path preferences—increasing its preference for observably stable paths over flapping paths. APM also includes a mechanism to distinguish route flaps due to topology changes, so as not to confuse them with those due to policy conflicts. The correctness and convergence analysis of APM derives from the sub-stability property of chosen paths. Implementation in the SSFNet simulator is performed, and simulation results for different performance metrics are presented. The metrics capture the dynamic performance (in terms of instantaneous throughput, delay, etc.) of APM and other competing solutions, thus exposing the often neglected aspects of performance.

Key Words: Inter-domain Routing; Border Gateway Protocol (BGP); Feedback Control; Convergence Analysis; Simulation.

I. INTRODUCTION

The Border Gateway Protocol (BGP) plays a major role in the performance of the Internet, and is known to have properties that are far from ideal. BGP allows policy-based routing; each AS independently defines a set of local policies regarding which routes to accept and advertise from/to other networks, as well as on which route it prefers when more than one route becomes available. However, independently defined local policies may lead to policy conflicts. Policy conflicts occur when neighboring ASes have opposite interests over routes. For example, assume AS u and AS v are neighbors, and AS v has two permitted paths $p1$ and $p2$, where $p1$ is preferred over $p2$. If extensions of $p1$ and $p2$, *i.e.* $(u, v)p1$ and $(u, v)p2$, are permitted at AS u , and $(u, v)p2$ is preferred over $(u, v)p1$, then there is a policy conflict. When AS v improves its best

path from $p2$ to $p1$, AS u will be forced to give up its more preferred path for the less preferred one. Any policy conflict can be resolved by changing the preference of the ASes over their paths, *i.e.* local policies.

Although not all policy conflicts are harmful, a group of ASes may define conflicting policies that cannot be satisfied simultaneously, causing BGP to *diverge*. Assume AS u , v , and z form such group. The scenario of divergence may take place as follows: When AS u improves its best path, it forces AS v to give up its best path for a less preferred path, which in turn gives AS z an opportunity to improve its best path, which forces AS u to give up its best path for a less preferred path, and so on. Each AS in such conflict repeatedly selects the same sequence of routes, never converging on any one set of routes. Therefore, route oscillations due to policy conflicts are *persistent*, and require some kind of intervention to stop.

As the commercial infrastructure of the Internet continues to grow, so does the potential for developing persistent route oscillation because of the growth of policies both in size and complexity [1], [2].

Several studies [3], [2], [4] have examined the dynamic behavior of inter-domain routing and highlighted the negative impacts of unstable routes. Instabilities taking place across ASes may negatively impact end-to-end network performance and efficiency of the Internet. A network that has not yet reached convergence may drop packets or deliver packets out of order. Routers may experience severe CPU load and memory problems: Because of repeated advertising and withdrawal of routes, routers need to rerun the BGP decision process to select the best paths, and update routing and forwarding tables. Frequent changes in the routes that are advertised by the other domains also make traffic engineering through an AS very difficult. BGP is crucial for a healthy and efficient global routing, and it is imperative to guarantee convergence of BGP independent of the locally selected policies.

Contribution of This Paper: There have been a number of studies (reviewed in Section II) on guaranteeing safety, *i.e.* convergence, of BGP independent of the locally selected policies [5], [6], [7], [8], [9], [10], [11], [12]. In our previous work [13], we introduced the idea of dynamically detecting and suppressing BGP oscillations through probabilistic change of path ranks (preferences). The algorithm is designed to detect policy conflicts by using local histories only. This paper extends and completes our preliminary idea [13] in many ways: (1) we augment the algorithm of path rank change so that an AS might choose a less preferred but *observably stable* path over a more preferred but oscillating path, thus it becomes natural for an AS to implicitly assign a higher cost (and hence

less preference value) to oscillating (flapping) paths; (2) with new additions, the algorithm enables the nodes to dynamically adapt to any state of the network. After the system stabilizes, we let the nodes attempt to restore some of the local preference values of their paths which they have modified so as to keep the overall path rank change minimal¹; (3) a new mechanism is added to distinguish route flaps due to topology changes, so as not to confuse them with those due to policy conflicts; (4) BGP extensions of the proposed algorithm are specified; (5) a correctness and convergence analysis of the proposed algorithm is developed based on the sub-stability property of chosen paths; (6) the proposed algorithm is implemented in the SSFNet simulator [15] is performed, and simulation results for different performance metrics are presented. The metrics capture the dynamic performance (in terms of instantaneous throughput, delay, routing load, etc.) of our algorithm as well as other competing solutions, thus exposing often neglected aspects of performance. Although our exposition is BGP-specific, the problem of inconsistent policies at independent distributed entities is more general.

The paper is organized as follows: Section II reviews background and related work. Section III describes our algorithm, and Section IV presents convergence and correctness analysis. Simulation results and conclusion are presented in Section V and Section VI, respectively.

II. BACKGROUND AND RELATED WORK

A. Border Gateway Protocol Abstraction

We use the abstraction of BGP proposed by Griffin *et al.* [7], which is called *Safe Path Vector Protocol (SPVP)*. SPVP is a distributed algorithm for solving the so-called *Stable Paths Problem (SPP)*. This model abstracts away low-level details of BGP and makes it easier to reason about convergence related issues.

Informally, SPP consists of an undirected graph with a single destination. Each node in the graph has a set of *permitted paths* to the destination, which are the routes learned from peers, and allowed by the local policy of the node. Each node also has a *ranking function* to impose an order of preference on the paths, such that more preferable paths have higher values assigned to them. A solution of an SPP is an assignment of permitted paths to the nodes that is consistent with the path chosen by its next-hop neighbor: Node u may choose the path $P = \langle u, v, w, \dots, \text{destination} \rangle$ only if the current path at node v is $\langle v, w, \dots, \text{destination} \rangle$ and path P is the current best path of node u .

The formal definition of SPP is as follows: A network is represented as a simple, undirected, connected graph $G = (V, E)$, where $V = \{0, 1, \dots, n\}$ is the set of nodes connected by edges from E . Nodes represent BGP routers and edges represent BGP sessions. For a node u , its set of *peers* is $\text{peers}(u) = \{w | \{u, w\} \in E\}$. Node 0 is the destination to which all other nodes are trying to find paths. A *path* P in G is a sequence of nodes $(v_k, v_{k-1}, \dots, v_1, v_0)$, such that $(v_i, v_{i-1}) \in E$, for all $i, 1 \leq i \leq k$.

An empty path, ϵ , indicates that a router cannot reach the destination. Nonempty paths $P = \langle v_1, v_2, \dots, v_k \rangle$ and $Q = \langle w_1, w_2, \dots, w_m \rangle$ can be *concatenated* as follows $PQ = \langle v_1, v_2, \dots, v_k, w_2, \dots, w_m \rangle$ if $v_k = w_1$. For every path P , concatenation with the empty path returns the path itself: $P\epsilon = \epsilon P = P$.

For every $v \in V - \{0\}$, the set \mathcal{P}^v denotes the permitted paths from v to the destination. Let $\mathcal{P} = \{\mathcal{P}^v | v \in V - \{0\}\}$ denotes the set of all permitted paths. For every $v \in V - \{0\}$, there is a ranking function $\lambda^v : \mathcal{P}^v \rightarrow \mathbf{N}$. $\lambda^v(P)$ denotes the degree of preference that node v gives to the path $P \in \mathcal{P}^v$. More preferable paths have higher values of λ^v . Let $\Lambda = \{\lambda^v | v \in V - \{0\}\}$ be the set of all ranking functions.

An instance of a *Stable Paths Problem (SPP)* $S = (G, \mathcal{P}, \Lambda)$, is a graph with the permitted paths and ranking function at each node if the following conditions are satisfied for every $v \in V - \{0\}$:

- (1) *Empty path is permitted*: $\epsilon \in \mathcal{P}^v$.
- (2) *Empty path is the lowest ranked path*: $\lambda^v(\epsilon) = 0$.
- (3) *Strictness*: If $\lambda^v(P_1) = \lambda^v(P_2)$, then $P_1 = P_2$ or $P_1 = (v, u)P'_1$ and $P_2 = (v, u)P'_2$ for some node u .
- (4) *Simplicity*: If path $P \in \mathcal{P}^v$, then P does not have repeated nodes, *i.e.* P is loop free.

Given a node u , and $W \subseteq \mathcal{P}^u$ with distinct next-hops, the *maximal path* in W , $\text{max}(u, W)$, is defined to be the highest ranked path in W . A *path assignment* is a function π that maps each node $u \in V$ to a permitted path $\pi(u) \in \mathcal{P}^u$. π defines the path chosen by each node to reach the destination. Given a path assignment π and a node u , the set of permitted paths that are one-hop extension of paths through neighbors is defined as

$$\text{choices}(u, \pi) = \{(u, v)\pi(v) | \{u, v\} \in E\} \cap \mathcal{P}^u.$$

The path assignment π is called *stable at node u* if $\pi(u) = \text{max}(u, \text{choices}(u, \pi))$. The path assignment π is called *stable* if it is stable at every node $u \in V$.

An SPP instance $S = (G, \mathcal{P}, \Lambda)$ is *solvable* if there exists a stable path assignment π for S . Every such assignment is called a *solution* for S and written as (P_1, P_2, \dots, P_n) , where $\pi(u) = P_u$. An instance of SPP may have no solution, or one or more solutions.

SPVP is an abstraction of BGP. Every node runs a copy of the SPVP process. With this abstraction, messages are simply paths. Each node maintains two data structures: $\text{rib}(u)$ is the current path that node u is using to reach the *destination*, and $\text{rib_in}(u \leftarrow w)$ denotes the path that has been most recently advertised by peer w and processed at node u . The set of paths available at node u is updated as

$$\text{choices}(u) = \{(u, w)\text{rib_in}(u \leftarrow w) | w \in \text{peers}(u)\} \cap \mathcal{P}^u$$

and the best path at u is

$$\text{best}(u) = \text{max}(u, \text{choices}(u)) \text{ and } \text{rib}(u) = \text{best}(u)$$

As long as node u receives advertisements from its peers, $\text{best}(u)$ is recomputed with the most recent $\text{choices}(u)$, and stored in $\text{rib}(u)$. Just as it is the case with BGP, when u changes its current path, it notifies its current peers about the change. This may cause the peers to send advertisements to

¹Akin to distributed recovery mechanisms, e.g. congestion avoidance of TCP [14].

their peers. The network reaches a *stable state* when there is no node which would change its current path to the destination. If such a state is reached, then the resulting state is the *solution* of the Stable Paths Problem (SPP). If SPP has no solution, then SPVP diverges. Figure 1 shows an example of a policy conflict leading to divergence.

Griffin *et al.* [16] present the structure called *dispute*

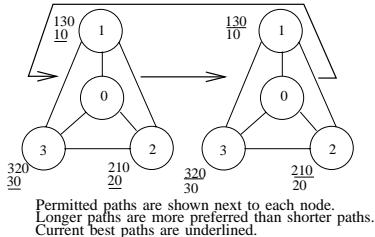


Fig. 1. An example of a divergence. Due to the cyclic conflict, this group of nodes cannot reach a stable state and keep oscillating between the shown states.

wheel for the purpose of checking the existence of a solution, and show that the lack of a dispute wheel is a sufficient condition which guarantees that SPP has a unique solution. A dispute wheel of size k is a structure that consists of nodes, u_1, u_2, \dots, u_k , and the set of paths Q_1, Q_2, \dots, Q_k , and R_1, R_2, \dots, R_k . For each $1 \leq i \leq k$, the following conditions are true: (1) R_i is a path from u_i to u_{i+1} ($u_1 = u_{k+1}$); (2) Q_i is a permitted path at u_i ; (3) $R_i Q_{i+1}$ is a permitted path at u_i ($Q_1 = Q_{k+1}$); (4) Q_i is less preferred than $R_i Q_{i+1}$ at node u_i . The dispute wheel of size k is shown in Figure 2. Q_i s

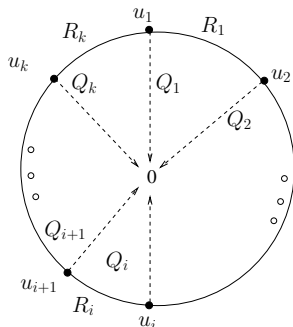


Fig. 2. Dispute wheel of size k

are called *spokes* of the dispute wheel and each spoke must be a simple (loop-free) path, *i.e.* no repeated nodes. None of the paths Q_i , R_i or Q_{i+1} can include node u_i . The paths R_i s are called the *rims* of the wheel, and each rim is also a simple (loop-free) path. The nodes at the ends of the paths R_i s are called *active nodes*. The active nodes are the nodes at which route preferences cause the dispute wheel.

Note that the presence of a dispute wheel does not imply that the system will diverge. However, if the system diverges, there exists a dispute wheel, and the oscillation must be either because of multiple solutions or lack of a solution as we demonstrate later.

B. Related Work

The possibility of BGP divergence due to policy conflicts is first shown by Varadhan *et al.* [1]. Since then, many studies

proposed approaches to guarantee the safety, *i.e.* convergence, of BGP independent of the locally selected policies [5], [6], [7], [8], [9], [13], [10], [11], [12]. These approaches can be broadly classified into static and dynamic solutions. Static solutions are centralized and require analyzing routing policies to verify that they are conflict free and cannot lead to protocol divergence, whereas dynamic solutions are distributed and require some mechanism to detect and resolve policy conflicts that are leading to divergence at run time.

Static solutions: Govindan *et al.* [5] propose a static solution which involves keeping policies in a repository called Internet Route Registry and verifying that they do not contain policy conflicts that could lead to protocol divergence. However, Griffin *et al.* [17] show that such kind of verification is computationally very expensive, and hard to achieve due to the private nature of the policies.

To avoid global coordination required in [5], Gao *et al.* [6], [18] propose another static solution which restricts the routing policies to the hierarchical structure that arises from commercial relationships between ASes, which may either be *provider-customer* or *peer-peer* relationship. Gao *et al.* give policy configuration guidelines in which each AS prefers routes heard from customers to the routes heard from providers and peers. Then the system is guaranteed to converge. This solution requires a database to keep relationships between ASes. Static periodic checks are required and performed by a global authority to verify conformance with these guidelines. Gao *et al.* algorithm may lead to unnecessary disabling of many routes from the start to guarantee the stability of the system, which restricts the flexibility in the choice of routing policies. Feamster *et al.* [19] argue that there may be legitimate reasons to deviate from the guidelines proposed in [6].

Other static solutions [12], [11], [10] suggest different constraints that also prevent policy based oscillations in advance.

Dynamic Solutions: Although *route flap damping* [20] can suppress temporary instabilities very well, it cannot detect or eliminate policy conflicts leading to persistent oscillation. Therefore, when there is a policy conflict leading to persistent oscillation, using route flap damping only makes the oscillations run in slow motion.

Griffin *et al.* [7] suggest extending BGP to carry additional information called *history* with each routing update message. A possible trace of SPVP for the system shown in Figure 1 is shown in Figure 3(a). *History* allows each router to describe the exact sequence of events that led to the selection of a path as the best path. An event $(+P)$ indicates that the node has chosen path P as its best path, and P is more preferred than its previous best path. Similarly, an event $(-P)$ indicates that the node has updated its best path, and the current best path is less preferred path than its previous best path P . A *history* containing loops is an indication of a potential protocol divergence. At step 4 of Figure 3(a), all 3 nodes have a cycle in the histories of their current best paths. SPVP assumes that such paths are problematic, and therefore eliminates them. For the assumed timing of events, with SPVP the system converges to unreachable destination for all nodes.

Since a cycle in the *history* is a necessary but not sufficient

step	node	best path	history
0	1	(10)	◊
	2	(20)	◊
	3	(30)	◊
1	1	(130)	(+130)
	2	(210)	(+210)
	3	(320)	(+320)
2	1	(10)	(-130)(+320)
	2	(20)	(-210)(+130)
	3	(30)	(-320)(+210)
3	1	(130)	(+130)(-320)(+210)
	2	(210)	(+210)(-130)(+320)
	3	(320)	(+320)(-210)(+130)
4	1	(10)	(-130) (+320)(-210) (+130)
	2	(20)	(-210) (+130)(-320) (+210)
	3	(30)	(-320) (+210)(-130) (+320)
5	1	ε	(-10)
	2	ε	(-20)
	3	ε	(-30)

(a)

step	node	count of node	best path
0	1	0	(10)
	2	0	(20)
	3	0	(30)
1	1	0	(130)
	2	0	(210)
	3	0	(320)
2	1	1	(10)
	2	1	(20)
	3	1	(30)
3	1	1	(130)
	2	1	(210)
	3	1	(320)
4	1	2	(10)
	2	2	(20)
	3	2	(30)
5	1	won't use (130) since count(3)≥2 will stabilize on (10)	
	2	won't use (210) since count(1)≥2 will stabilize on (20)	
	3	won't use (320) since count(2)≥2 will stabilize on (30)	

(b)

Fig. 3. (a) A possible trace of SPVP for the system shown in Figure 1; (b) A possible trace of the Cobb and Musunuri algorithm for the system shown in Figure 1 assuming threshold value for *count* is 2. All nodes stabilize on their lowest preferred paths.

condition for divergence, there may be false positives. Carrying *history* with each update creates communication overhead, and may also reveal private information about the preferences of ASes over the routes. APM uses the idea of keeping track of history of path changes, but does it only locally. By keeping histories as local information and avoiding exchanging such information helps overcome related privacy concerns and communication overhead.

Recently, Cobb *et al.* propose two dynamic algorithms [8], [9]. The first work [8] proposes a mechanism to enforce monotonic path orderings. With this algorithm, convergence is guaranteed by preventing the selection of a path with higher order (preference) in one node, if doing so would cause a conflict with other nodes along the *routing tree*, *i.e.* a node would be forced to use another less preferred path. Preventing the violation of monotonic path ordering property is realized via diffusing computation along the routing tree whenever a node wants to update its path.

In their second work, Cobb and Musunuri [9] associate an integer cost with each node and exchange this cost value with each update message. The cost increases monotonically if the system diverges. Therefore, discarding advertisements from nodes whose cost is greater than a threshold is suggested. Assuming threshold value of 2, Figure 3(b) shows a possible trace of the Cobb and Musunuri algorithm for the system. Since the cost of the nodes involved in the same conflict grows in tandem, all of the nodes simultaneously give up their most preferred paths and stabilize on their lowest preferred paths.

A weakness of this algorithm is keeping per node cost, which causes aggregation of the paths through the same node. One flapping path may cause all the alternative paths (through the same node) to be eliminated. With APM, we extend the idea of using *count* to keep per-path state at each node instead of per-node state, which prevents aggregation of the paths through the same node. Empowered with this extra information together with probabilistic update of path ranks, APM can pinpoint the paths causing problems, and lead to

fewer path elimination. Another limitation of the approach is that Cobb and Musunuri suggest resetting the cost of all nodes once a week or once a month via a distributed reset protocol [21], which is based on a diffusing computation over a min-hop spanning tree. Resetting costs allows nodes to follow one more time their routing policies, where possibly this time no conflicts exist. However, to limit the number of path eliminations, resetting costs should be performed as soon as the state of the system changes.

Other Work: Another line of work concentrated on solving policy conflicts by treating the routing problem as a game in which the ASes are strategic agents [22]. For the case where AS policies are restricted to the hierarchical relationships, Feigenbaum *et al.* [22] present a *strategy-proof mechanism* that can be computed in polynomial time in a centralized computational model. However, it is shown that this mechanism is incompatible with BGP: If this mechanism is computed by a BGP-like distributed algorithm with similar data structures and communication patterns, it may cause BGP to converge very slowly, and/or can trigger extensive amount of update messages. [22] also show that if AS policies are not restricted, then it is NP-hard to compute a routing tree that maximizes the overall utility of the ASes.

III. ADAPTIVE POLICY MANAGEMENT (APM)

A. Overview

We propose a new algorithm to dynamically detect and eliminate policy conflicts leading to BGP divergence. The idea is to locally detect the paths involved in a conflict, and eliminate the conflict by changing the relative preference of such paths. Note that such adaptation is limited to the node's set of permitted paths, any of which the AS is willing to use albeit at different preference level.

Each node involved in a particular conflict observes *route flaps*: Constantly chooses a path as its best path and later gives it up for another path. For example, in Figure 1, node 1 constantly upgrades its current best path to (130), but later it is forced to give up (130) for its less preferred path (10) as a result of its neighbors' response to this upgrade. The nodes observing constant route flaps can stop such behavior by sticking to their less preferred but more stable path, even when a better alternative is advertised. This can be achieved by changing the local preference of the paths. When the node stops advertising the paths alternately, the cyclic effect of the global conflict will be broken. In Figure 1, for example, if node 1 changes its local preferences to prefer (10) over (130), the system stabilizes on the following path assignment: (10)(210)(30).

To be able to locally detect route flaps and the paths whose preference cause divergence, each node needs to keep some form of *local history*. We suggest keeping track of the paths that have been recently selected as best path, and their *counts* indicating how many times the path has been chosen as best path and later given up. Figure 4 shows how *counts* keep increasing during divergence of the system shown in Figure 1. Nodes involved in the conflict can detect divergence by comparing *counts* against a threshold called *min.threshold*.

step	node	best path	local history (path, count)	path preference
0	1	(10)	((10),1)	(130)>(10)
	2	(20)	((20),1)	(210)>(20)
	3	(30)	((30),1)	(320)>(30)
1	1	(130)	((130),1), ((10),1)	(130)>(10)
	2	(210)	((210),1), ((20),1)	(210)>(20)
	3	(320)	((320),1), ((30),1)	(320)>(30)
2	1	(10)	((130),1), ((10),2)	(130)>(10)
	2	(20)	((210),1), ((20),2)	(210)>(20)
	3	(30)	((320),1), ((30),2)	(320)>(30)
3	1	(130)	((130),2), ((10),2)	(130)>(10)
	2	(210)	((210),2), ((20),2)	(210)>(20)
	3	(320)	((320),2), ((30),2)	(320)>(30)
4	1	(10)	((130),2), ((10),3)	count(10) > $min_threshold$, change rank with probability 1/2 assume this takes place: (10)>(130)
	2	(20)	((210),2), ((20),3)	count(20) > $min_threshold$, change rank with probability 1/2 assume this does not take place: (210)>(20)
	3	(30)	((320),2), ((30),3)	count(30) > $min_threshold$, change rank with probability 1/2 assume this does not take place: (320)>(30)
5	1	(10)	stabilizes on lower preferred and available path	(10)>(130)
	2	(210)	stabilizes on most preferred and available path	(210)>(20)
	3	(30)	stabilizes on most preferred and available path	(320)>(30)

Fig. 4. A possible trace of APM for the system shown in Figure 1. $min_threshold = 2$.

Since the algorithm we are proposing is distributed and based on using only local information, there may be many nodes synchronously detecting the same conflict and lowering the preferences of their higher preferred paths. If we assume $min_threshold=2$ for each node in Figure 4, at step 4, all 3 nodes simultaneously change their local preferences to prefer their shorter paths, which are more stable in the sense that they are always available. Note that the conflict can be broken even if only one of the nodes performs the path rank change. To prevent this kind of simultaneous and unnecessary path preference changes, we suggest changing relative preferences with probability 1/2.

Because of the probabilistic adjustment of path preferences, even though the effect of a particular conflict is observed several times, it is possible that the conflict remains unresolved. $max_threshold$ is introduced to handle such cases: When the *count* associated with a particular path exceeds $max_threshold$, then the path is removed from the set of permitted paths, and added to the set of *bad paths*. The

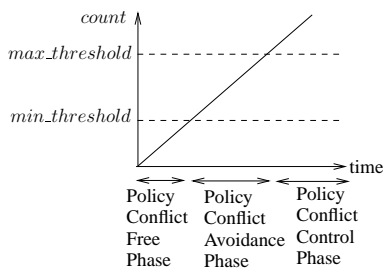


Fig. 5. Three phases of Adaptive Policy Management (APM)

bad paths set is a data structure that keeps the list of paths which the node believes that their adoption leads to a conflict. Therefore, they are excluded from further consideration in the best path selection process (until they are restored as the algorithm adapts to a conflict-free state), even if they are advertised by peers and permitted by original local policies. Setting $max_threshold$ to higher values helps reduce the number of paths placed in *bad paths*. However, smaller values

may reduce convergence time. *Count* values kept in the local history of a node are compared against $min_threshold$, and $max_threshold$ to detect and handle divergence as follows: (a) **Policy conflict-free phase**: When the *counts* are smaller than $min_threshold$, then the node assumes that there is no persistent oscillation. Therefore, setting $min_threshold$ to higher values helps prevent path preference changes when the oscillation is transient; (b) **Policy conflict-avoidance phase**: If any *count* exceeds $min_threshold$, but stays lower than $max_threshold$, then the node assumes that there is a policy conflict leading to persistent oscillation, which can be avoided by changing the relative preference (rank) of the paths; (c) **Policy conflict-control phase**: If any *count* exceeds $max_threshold$, then the path associated with this *count* is added to a set of *bad paths*, and excluded from further consideration in the best path selection process. Figure 5 shows these three different phases of our algorithm.

B. Details of APM

There may be different instantiations of the algorithm depending on the exact nature of the path information kept in the local history, and the way *count* values are associated with the paths. In this section, we describe the instantiation that we have chosen. Throughout this section, we assume that there is a single destination. Node u has *adopted path* p means that node u has chosen path p as its best path, and currently using it to reach the destination. Node u has *abandoned path* p means that node u was using path p to reach destination, and since node u 's best path has changed, node u is no longer using p to reach the destination.

1) *Data Structures*: In addition to data structures required for BGP, APM requires the usage of the following data structures:

- Each node u keeps a *local history* in the form of $(path, count)$ tuples, where *path* indicates a path that has been recently adopted by node u , and *count* indicates how many times the *path* has been adopted and later abandoned. When there is divergence, some *count* values keep increasing because of constant flapping.
- *peerStability* is an integer associated with each peer of node u . When $w \in peers(u)$ sends an update for a particular destination that advertises a path p that is different from the one that has been advertised previously, i.e. $rib_in(u \leftarrow w) \neq p$, the *peerStability* value corresponding to peer w is increased. The purpose of this counter is to differentiate the peers, and hence the paths advertised by those peers, that are stable. Therefore, if node u is observing a route flap, the flap can be stopped by adopting a path advertised by a *stable peer*: The smaller value of *peerStability* indicates a more stable peer. If *peerStability* is equal to one for peer w , it means the path advertised by w never changed later. We refer to such paths as *safe paths* and the peers advertising these paths as *stable peers*. To make stability last, after adopting a stable path, node u also changes its preference of the paths to reflect this choice. Node u updates the local preference of the stable path so that it will be the most

preferred path, *i.e.* $rank(\text{safe path}) = 1$, where $rank(p)$ is the index of path p among the current alternative paths in the order of decreasing local preference value. If there are more than one safe path, the one that is originally preferred more is chosen.

Note that *count* values associated with *paths* in *local history* cannot be used to measure stability of peers. A path p advertised by w may have a high *count* value associated with it, even if w never changes this advertisement. If the advertisement sent by peer w does not change, *peerStability* will always be equal to 1 for peer w even if the node receiving this update may constantly adopt p and later abandon it for a higher preferred path as a result of a policy conflict. In such case, *count* associated with p keeps increasing as the node constantly adopt p and later abandon it for a higher preferred path.

- B indicates the *bad path* set, which keeps the paths whose *count* exceeds the *max.threshold* value. Such paths are eliminated from the permissible set of paths at node u and not considered in the best path selection process even though they may be advertised by a peer.
- *keepaliveCount* is used to count the number of times a KEEPALIVE message is received from a peer w . If the value of *keepaliveCount* exceeds a threshold, *ka.threshold*, for all peers of node u , then node u concludes that the system has stabilized and there are no more policy conflicts. After this point, node u probabilistically (and more conservatively) resets some of the local preferences back to their original values². Although there is a possibility of introducing instability back into the system, our algorithm adapts to the changes dynamically and stabilizes at some state of path preferences eventually.

The state of the system is defined by the values kept in these data structures, as well as the path orderings at each node, which correspond to different policies.

2) *Update Handling*: Figure 6 shows the pseudo-code of the Adaptive Policy Management (APM) scheme for handling routing updates. The process runs at each node u in response to a received update. When node u adopts a path $p \in \mathcal{P}^u$, it informs each of its peers by sending an update message. $rib(u)$ indicates the current best path to the destination selected at node u . $rib_in(u \leftarrow w)$ indicates the most recent path sent from $w \in peers(u)$, and processed at node u . The set of path choices available at node u that are considered for best path selection, excluding the bad paths in $B(u)$, is defined as

$$choices_B(u) = \{(u, w)rib_in(u \leftarrow w) - B(u) | w \in peers(u)\} \cap \mathcal{P}^u$$

and the best path as

$$best_B(u) = \max(u, choices_B(u)).$$

As long as node u receives advertisements from its peers, $best_B(u)$ is recomputed with the most recent $choices_B(u)$. When $rib(u)$ changes, node u notifies its peers by sending an

update message.

When the process goes into the policy conflict-avoidance phase, after successfully choosing a safe path and changing its rank to be most preferred, the state of the system changes. The state of the system also changes when the process places a path in the *bad path* set, *i.e.* in the policy conflict-control phase. In either case, this new state corresponds to a different Stable Paths Problem (SPP), possibly a stable one. Therefore, counters are reset to give opportunity for a fresh start and to see if the change is enough to reach stability. When the process goes into the policy conflict-avoidance phase, if there is no *safe path* for node u , *i.e.* $peerStability(w) \neq 1$ for any $w \in peers(u)$, then node u does not do anything to stop the oscillation. However, if there is a safe path P_{safe} , with probability 1/2, the path ordering at node u is changed such that $rank(P_{safe}) = 1$. If there are more than one safe path, then the most preferred one is chosen as the highest ranked path.

3) *Restoring Local Preferences*: Each node keeps track of the number of KEEPALIVE messages received from its peers, and compares this value against a threshold, denoted by *ka.threshold*, to test the stability of the system. Figure 7 shows how node u probabilistically restores some rank changes for its paths after the system has stabilized. Since policies are placed for a purpose by each node, such as traffic engineering or security, it is important for ASes not to change them unless they are conflicting with the policies of other nodes and absolutely necessary to eliminate route oscillations. Although it is safe to restore rank changes that do not compromise the current stability, there is no way for node u to know which changes are safe to restore. Therefore, node u uses a probabilistic (albeit more conservative) approach, and risks introducing instability back into the system. Contrary to update handling, node u increases the local preference of a path with a much smaller probability, 1/4. We allow for bringing paths out of B with probability 1/4 as well, since the above argument is also true for the paths currently in the *bad path* set. If node u performs a rank change and/or remove (restore) a path from B , counters kept in the local history are reset because this new state corresponds to a different SPP.

Note that although re-using suppressed paths and/or restoring original ranks once the system reaches a stable state may introduce instability back into the system, this is the nature of all adaptive feedback control systems. By using a much smaller probability for reset, *i.e.* 1/4, we provide a conservative way of probing the network state. This is akin to the congestion avoidance mechanism of TCP, during which the current state of the network is probed at a slower rate. Cobb and Musunuri algorithm [9] suggests periodically resetting the cost of all nodes via a distributed reset protocol [21] for this purpose. However, the algorithm lacks any mechanism to prevent introducing the same conflict back into the system when the costs are reset. Therefore, this approach causes the system to oscillate between the stable and the unstable state unless policy conflicts disappear in the meantime due to a topology/policy change. SPVP [7] resets *history* when a cycle is detected and a path is eliminated. However, to be able to truly adapt to the system dynamics, SPVP should also reset

²This is akin to increase/decrease adaptation rules employed in many adaptive feedback-control systems.

```

process APMS_Update_Handling[u] //Update Handling
receive Update m from peer w →
  keepaliveCount(w)=0
  if rib.in(u ← w) ≠ m then
    peerStability(w)++
  rib.in(u ← w) = m
  if rib(u) ≠ bestB(u) then
    Pold=rib(u)
    Pnew = bestB(u)
    if (Pnew ≠ ε) then
      count(Pnew)++
    if count(Pnew) > max_threshold then
      B(u) = B(u) ∪ {Pnew} //Policy Conflict–Control Phase
      Pnew = bestB(u)
      count(Q)=0 for each path Q ∈ localHistory
      peerStability(v)=0 for each v ∈ peers(u)
    else if count(Pnew) > min_threshold then
      do with probability= 1/2
        find the most preferred safe path, Psafe //Policy Conflict–Avoidance Phase
        rank(Psafe)=1
        Pnew = Psafe
        count(Q)=0 for each path Q ∈ localHistory
        peerStability(v)=0 for each v ∈ peers(u)
  if Pnew ≠ Pold then
    rib(u) = Pnew
    for each v ∈ peers(u) do
      send rib(u) to v

```

Note: The code to the right of the → is assumed to be executed in one atomic step

Fig. 6. Adaptive Policy Management (APM): Update Handling

```

process APMS_Keepalive_Handling[u] //Keepalive Handling
receive keepalive from w →
  keepaliveCount(w)++
  if keepaliveCount(v) ≥ ka_threshold for every v ∈ peers(u)
    for each v ∈ peers(u)
      r=rib.in(u ← v)
      if (localpref(r) ≠ originallocalpref(r)) || (r ∈ B(u)) then
        do with probability= 1/4
          if r ∈ B(u) then
            remove r from B(u)
          localpref(r)= originallocalpref(r)
          count(Q)=0 for each path Q ∈ localHistory
          peerStability(v)=0 for each v ∈ peers(u)
          keepaliveCount(v)=0 for each v ∈ peers(u)
          Pnew=bestB(u)
          if Pnew ≠ rib(u) then
            rib(u) = Pnew
    for each v ∈ peers(u) do
      send rib(u) to v

```

Note: The code to the right of the → assumed to be executed one atomic step

Fig. 7. Adaptive Policy Management (APM): Restoring Local Preferences once Stability is Reached

the state by re-introducing all suppressed paths. Otherwise, the algorithm prevents usage of the suppressed paths even if the policy conflict disappears due to a topology/policy change. Similar to the Cobb and Musunuri algorithm, such reset may cause oscillation between the original unstable state and the stable state. APM resets the state of the system conservatively using the mechanism shown in Figure 7 so as to adapt to every state of the network automatically while oscillating in the vicinity of a stable state at a very slow rate. With APM, the backoff and recovery of path preferences could also be guided by explicit input from local AS administrators. We will investigate this approach in a future report.

4) *Handling Transient Oscillations due to Topology Changes*: If there is a topology change, path updates experi-

enced as a result of a change in topology may interfere with diagnosing policy conflicts. For example, link or node failure or recovery may create a route flap, and increase the chance of false positives. More importantly, topology changes affect policy dynamics. Even if the original topology had policy conflicts, the new topology may be conflict-free, or vice versa. Therefore, during the process of resolving policy conflicts, it is important for the nodes to be aware of link/node failure and recovery events and distinguish them from route flaps due to policy conflicts.

Assume that node u 's next-hop along the path to the

```

process APMS_TopologyChange_Handling[u] //Handling topology changes
//when node u learns that its next-hop link along
//the path to the destination has failed or restored
//or the next-hop node failed or restored
if link(u, v) has failed or node v has failed
  r = rib.in(u ← v)
  if r is permissible
    feasible(r)=false
  if link(u, v) is restored or node v is restored
    receive Update m from v
    rib.in(u ← v) = path(m)
//reset local states
B(u)={ }
peerStability(w)= 0 for every w ∈ peers(u)
keepaliveCount(w)=0 for every w ∈ peers(u)
count(Q)=0 for each path Q ∈ localHistory
for every w ∈ peers(u)
  p ∈ rib.in(u ← w)
  if p is permissible
    if localpref(p) ≠ originallocalpref(p)
      localpref(p)=originallocalpref(p)
sequenceNumber++ //increase the sequence number of u
rib(u)=bestB(u)
for every w ∈ peers(u)
  send update message m
  where path(m)=rib(u), originator(m)=u, sequenceNumber(m)=sequenceNumber

```

Fig. 8. Adaptive Policy Management (APM): Handling Failures and Recovery

```

process APMS_TopologyChange_Handling[u] //Handling topology changes
//when node u gets an update message m from node v with originator(m) ≠ null
if u == originator(m)
  //process this update as shown in APMS_Update_Handling[u]
  if there is a new update to be sent
    set originator = null before sending update to the peers
else if (sequenceNumber(m) > sequenceNumber of originator(m) at node u)
  B(u)={ }
  peerStability(w)= 0 for each w ∈ peers(u)
  keepaliveCount(w)=0 for every w ∈ peers(u)
  count(Q)=0 for each path Q ∈ localHistory
  for every w ∈ peers(u)
    p ∈ rib.in(u ← w)
    if p is permissible
      if localpref(p) ≠ originallocalpref(p)
        localpref(p)=originallocalpref(p)
  //update sequenceNumber of originator(m) at node u
  sequenceNumber of originator(m) at u=sequenceNumber(m)
  rib(u)=bestB(u)
  for every w ∈ peers(u)
    send update message n
    where path(n)=rib(u), originator(n)=originator(m)
    sequenceNumber(n)=sequenceNumber(m)

```

Fig. 9. Adaptive Policy Management (APM): Handling Failures and Recovery

destination is node v , i.e. $P = \langle u, v, \dots, destination \rangle$. When the link between u and v goes down, as soon as node u detects the change, it discards the route learned from v , and recomputes its best path to the destination. At this point, node u knows that its best path has changed because of a failure. We suggest that while sending the resulting update message, node u includes some information about the failure so that the other nodes which are not in the neighborhood of the failed link can deduce that the update they are receiving is triggered by a failure. Since from node u 's perspective, the effects of failure of the peer v is similar, it can be handled the same way. Furthermore, link/peer restorations can also be handled in a similar way by observing OPEN messages exchanged when a peering TCP session is (re-)established.

We suggest adding two new fields to the update messages of BGP for this purpose: *originator* and *sequenceNumber*. While *originator* carries the ID of the node that has detected the topology change, *sequenceNumber* indicates the number of times the node detected a topology change. Each node also keeps track of the *sequenceNumber* for the other nodes to distinguish the most recent topology update when the *originator* is the same. For example, if the same link constantly going up and down quickly, the resulting updates will have the same *originator*. When these updates propagate at different speed, it is important to be able to distinguish which update reflects the latest topology change. With this addition, the messages exchanged between peers are triples $(P, originator, sequenceNumber)$. In the absence of topology changes, the *originator* field is set to *null*. The pseudocode of the topology handling algorithm is shown in Figures 8 and 9.

When node u detects a topology change, it first resets the local state, and increases its *sequenceNumber* as shown in Figure 8. The resulting update is sent with *originator*= u , and the current *sequenceNumber* of node u . As the pseudocode in Figure 9 shows, when node u receives an update message m with *originator*= v and the *sequenceNumber* indicating that this is a new topology change, node u temporarily turns off the policy conflict detection process and resets the local state. After updating the *sequenceNumber* of v for its records, node u sends the resulting update message with the same *originator* and *sequenceNumber* values to further propagate the change. When node u receives an update message with itself as the *originator*, it restarts the policy conflict detection process by assigning *null* to *originator*. Policy changes can be handled the same way as topology changes. Since a policy change is performed by a local authority, the AS performing policy change can use the above mechanism to provide temporary suspension of processing the resulting updates for policy conflict detection.

While this mechanism handles transient oscillations due to topology/policy changes, *min_threshold* still helps not to react too soon to transient oscillations arising from other causes.

IV. CONVERGENCE ANALYSIS OF APM

Different path orderings at each node correspond to different policies and define different states of the network. Among

these states, there are some stable configurations. Our goal is to show that starting with an arbitrary state of the system, the Adaptive Policy Management (APM) converges to a stable state within a finite number of steps. To that end, we list some formal definitions for terms we use henceforth. The state of the system changes over time and it is a function of the current policies and the update messages (advertisements and withdrawals) exchanged among the nodes. The following definitions assume that there is a single policy conflict in the system and the classification of the nodes is done with respect to this particular conflict.

Definition 4.1: Conflict-free node is a node which is not involved in the policy conflict, i.e. not an active node of a dispute wheel, and stabilized on its best path.

Definition 4.2: Non-flapping path, or stable path $P = \langle u, \dots, v, \dots, destination \rangle$ is the best path of a conflict-free node u , and this path does not change over time.

Definition 4.3: Observable safe path $P = \langle u, v, \dots, destination \rangle$ of a node u is a permitted path at u and advertised by peer v of node u . The path $\langle v, \dots, destination \rangle$ is a non-flapping path and none of the nodes along this path experiences route flaps due to the conflict.

Definition 4.4: Conflicting safe-alternative node is a node which is involved in the policy conflict, and observes a safe path which is not its most preferred path.

Definition 4.5: Conflicting node is a node which is involved in the policy conflict, and does not observe any safe path.

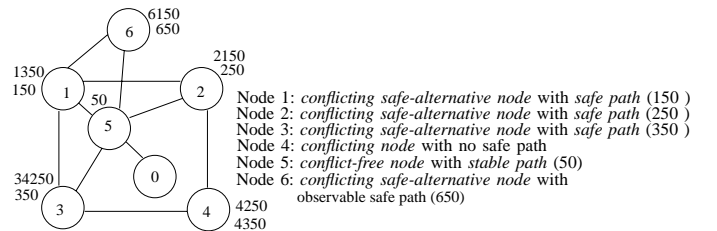


Fig. 10. An example showing different types of nodes when there is a single conflict in the system.

Figure 10 shows an example of a group of nodes, $\{1,2,3,4\}$, involved in a cyclic conflict. Node 5 is a *conflict-free* node with stable path (50). Although node 6 is not an active node on the dispute wheel for this conflict, it is not a *conflict-free* node yet according to our definition since it has not stabilized on its best path. At this point, node 6 is a *conflicting safe-alternative* node with single observable safe path (650). (6150) is not an *observable safe path* due to node 1's involvement in the conflict. The nodes actively involved in the conflict, $\{1,2,3,4\}$, are either *conflicting safe-alternative* or *conflicting*. *Conflicting safe-alternative nodes* can break the conflict by holding onto their safe paths, i.e. changing path rankings such that the safe path becomes their most preferred path. If a node does not have any safe paths, i.e. *conflicting node*, then it cannot stop oscillation through a rank change of its paths. However, as soon as one of its next-hop neighbor stabilizes, it will start to observe stable path advertisements coming from this newly

stabilized neighbor. At this point, the node becomes either a *conflicting safe-alternative node*, or a *conflict-free node*. For the example shown in Figure 10, as soon as node 2 changes its path preference to prefer (250) over (2150): Node 2 becomes a *conflict-free node*; path (250) becomes a *stable path* of node 2; path (4250) becomes a *safe path* at node 4 and node 4 stabilizes on it since it is its most preferred path; node 3 then chooses its most preferred path (34250) through node 4; node 1 has no choice but to stabilize on path (150) which is one of its permitted paths and available from its peer 5; and the conflict is resolved. Then node 6 will start to observe both of its safe paths and stabilize on its most preferred path (6150).

Assuming that there is a single cyclic conflict in the system, let N denote the set of nodes that are in this conflict, where $|N| \geq 2$. Since *conflicting safe-alternative nodes* are playing the key role in breaking conflicts, we would like to show that there must be some *conflicting safe-alternative nodes* in N for such a cyclic conflict to occur. Let M denote the set of such nodes. Obviously, the nodes in M have paths which they prefer over their safe path, thus causing a cyclic conflict. Throughout the conflict, the more preferred paths are constantly advertised and withdrawn. If any node in M changes its preference to pick its safe path over its more preferred but oscillating path, then it can break this cyclic conflict. Independent of the size of M , only one *conflicting safe-alternative node* suffices to perform path rank change for breaking the conflict. If we assume the nodes in M are u_1, \dots, u_k , they will form a dispute wheel as shown in Figure 2. Node u_i s, where $1 \leq i \leq k$, are the nodes at which route preferences cause the shown dispute wheel. Node u_i constantly changes its path from Q_i to its more preferred path $R_i Q_{i+1}$, which in return causes u_{i-1} to give up $R_{i-1} Q_i$ and use Q_{i-1} . Q_i is a safe path for u_i since we assumed that this is the only conflict in the system. If at least one node, u_i , changes the preference of its safe path so that Q_i is preferred over $R_i Q_{i+1}$, the dispute wheel cannot form.

Before we start the convergence proof, we would like to show that there must be at least 2 *conflicting safe-alternative nodes* for the existence of a persistent oscillation. Since it is obvious that if none of the nodes is *conflicting safe-alternative node* there cannot be persistent oscillation, it suffices to show that having only one *conflicting safe-alternative node* is also not enough to create a policy conflict leading to divergence.

Lemma 4.1: Under BGP, policies of a group of nodes cannot lead to any type of persistent oscillation if only one of the nodes in this group is a *conflicting safe-alternative node*.

Proof: We prove this by contradiction: Assume that $u_1 \dots u_n$ are the group of nodes who are involved in the persistent oscillation, and only one of the nodes in this group has a safe path. Let u_1 denote this only *conflicting safe-alternative node*, and P_1 denote its safe path. The other nodes in the group, $u_2 \dots u_n$, are *conflicting nodes*, and their permitted paths must be traversing at least one node involved in the conflict because otherwise their paths would be safe. Let u_i be a *conflicting node*, where $2 \leq i \leq n$. Permitted paths of u_i must be in the form of $P_i = \langle \dots, u_j, \dots, destination \rangle$, where u_j is a node in the conflict, *i.e.* $1 \leq j \leq n$ but $j \neq i$. (Since BGP requires the paths to be *simple*, *i.e.* no

repeated nodes, the path P_i cannot traverse u_i .) Without losing generality assume the following: If there are multiple nodes involved in the conflict along path P_i , node u_j is the *closest* one to the destination. There are two cases: (1) $j \neq 1$; (2) $j = 1$. Case (1) cannot be true, because it implies that the subpath of P_i following u_j to the destination does not traverse any node in the conflict, and therefore this path is a safe path for node u_j . However, this makes node u_j a *conflicting safe-alternative node*, which contradicts our assumption that u_1 is the only *conflicting safe-alternative node* in the conflict. Therefore, case (2) holds for every node u_i where $2 \leq i \leq n$: The *conflicting nodes* reach the destination through u_1 , which contradicts our assumption that they do not observe any safe paths. ■

Lemma 4.1 implies that there must be 2 or more *conflicting safe-alternative nodes* in a persistent oscillation.

For the above discussion, we have assumed that there is only one conflict involving a group of nodes. If there are multiple conflicts in the system, nodes may get involved in many conflicts simultaneously. The safe paths of the nodes that are in a particular conflict may not be *observable* due to other conflicts. The following definitions relax the single policy conflict assumption that we have made earlier.

Definition 4.6: Conflict-free node is a node which is not involved in *any* policy conflict, *i.e.* not an active node of a dispute wheel, and stabilized on its best path.

Definition 4.7: Observable safe path $P = \langle u, v, \dots, destination \rangle$ of a node u is a permitted path of u which is advertised by peer v . The path $\langle v, \dots, destination \rangle$ is a non-flapping path and none of the nodes along this path experiences route flaps due to *any* conflict.

Definition 4.8: Innermost conflict along the path $P = \langle u_k, u_{k-1}, \dots, u_2, u_1, destination \rangle$ is the conflict that involves node u_i , where u_i is the *closest* node to the destination and involved in a conflict. In this case the **innermost safe path** along the path P is $\langle u_i, u_{i-1}, \dots, destination \rangle$.

Definition 4.9: An inactive node with observable safe path is a node that is not an active node on a dispute wheel, and its most preferred path is an *observable safe path*.

Definition 4.10: Conflicting safe-alternative node is a node which is involved in a policy conflict, and observes a safe path which is not its most preferred path.

Definition 4.11: Conflicting node is a node which is involved in a policy conflict, and does not observe any safe path.

Figure 11 shows an example of groups of nodes, $\{1,2,3\}$ and $\{4,5,6\}$, that are in 2 different conflicts, *conflict*₁ and *conflict*₂, respectively. Node 3 is a *conflicting safe-alternative node* with safe path (350) with respect to *conflict*₁. However, due to node 5's involvement in *conflict*₂, (350) is not an observable safe path. The innermost conflict along (350) is *conflict*₂.

Once the *innermost conflicts* along the safe paths are broken, the nodes of the outer conflicts start to observe their safe paths, and have a chance to break their conflict under APM by sticking to their safe (albeit less preferred) path. For each conflict, as we have shown earlier, we have at least 2

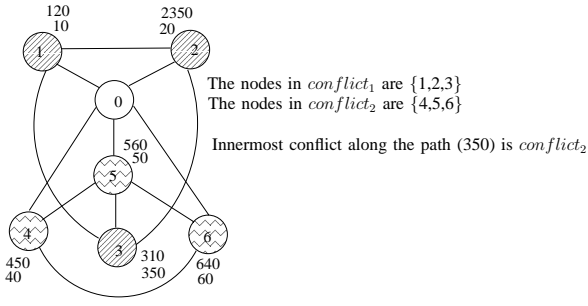


Fig. 11. An example showing different types of nodes when there are multiple conflicts in the system.

conflicting safe-alternative nodes. Starting from the innermost conflicts, by having at least one of these *conflicting safe-alternative nodes* perform path rank change at each step, we break the cyclic conflicts, and increase the number of *conflict-free* nodes that are stabilized on their paths. This is the idea behind the convergence proof of our Adaptive Policy Management (APM).

Denote by S the set of nodes that are *conflict-free* and stabilized on their paths during the execution of the algorithm. During the k^{th} step of APM, the nodes that are currently in S advertise their paths to their peers. The peers process these advertisements, and enter into S after stabilizing on their paths. However, the k^{th} step does not end until the resulting new peers that are *inactive nodes with observable safe path* are also processed iteratively (thus becoming *conflict-free* and entering S) until no more *inactive nodes with observable safe path* can be reached.

A node changes its state over time as it runs APM. APM helps the node stabilize by increasing the preference of its safe path. A *conflicting safe-alternative node* becomes *conflict-free node* when it performs rank change to make its observable safe path the most preferred path and stabilizes on it. A *conflicting node* u with flapping paths $(u, v)P_v$ and $(u, w)P_w$, becomes a *conflicting safe-alternative node* when the node advertising path P_v enters S stabilizing on path P_v .

Lemma 4.2: During the execution of the APM algorithm, the size of the set of nodes that are *conflict-free* (set S) increases monotonically if we perform path rank changes whenever conflict is detected.

Proof: The nodes in S form a routing tree of the paths on which they are stabilized. This routing tree is rooted at the destination and grows as the nodes outside of S adopt and stabilize on the extension of the paths advertised by the nodes in S .

To show that S grows monotonically, we need to show that at each step of the algorithm, at least one conflict is resolved and at least one node is added to the set until the system converges. We use induction based on the number of nodes in S .

Basis: At the beginning, while S was empty, the destination is added. Hence it holds for the base case.

Hypothesis: At step k of the execution, assume the size of the set S is n and up to this point the set S grew monotonically.

Induction Step: We need to show that at step $k + 1$, the size

of S will be greater than n . When node u is about to be a member of S , one of the following is true:

Case 1) Node u has just received a path advertisement P_v from node v in S , where $(u, v)P_v$ is a permitted path at node u , but node u is not stabilizing on this path. Since $(u, v)P_v$ is an observable safe path of node u , node u must be a *conflicting safe-alternative node*. At this point, node u performs rank change and sticks to the path $(u, v)P_v$, becomes a *conflict-free node*, and is added to S . As we have shown earlier, if there is a conflict, there will be at least 2 *conflicting safe-alternative nodes* and this is the step where they break the conflict. Consequent *inactive nodes with observable safe paths* are also handled iteratively, and entered into S until hitting the nodes that are not *inactive nodes with observable safe paths*.

Case 2) Node u was a *conflicting node* at previous step (step k), and observing path flap in the form of $P1_u, P2_u, P1_u, P2_u, \dots$. Assume that $P1_u$ is preferred over $P2_u$ at node u , and $P1_u = (u, v)P_v$, where P_v is advertised by node v and $P2_u = (u, w)P_w$, where P_w is advertised by node w . At step k , neither v nor w was in S , since otherwise node u wouldn't be a *conflicting node*. At the end of step k due to a resolved conflict, node v and/or node w might have entered S . The possible cases at step $(k + 1)$ due to a resolved conflict at step k are: (i) Node v is in S and stabilized on path P_v , and node w is not in S ; (ii) Node w is in S and stabilized on path P_w , and node v is not in S ; (iii) Both node v and w are in S and stabilized on paths P_v and P_w , respectively. If (i) happens, node u stabilizes on $P1_u$ and becomes a *conflict-free node*. If (ii) happens, node u becomes a *conflicting safe-alternative node* with observable safe path $P2_u$. Note that $P1_u$ is not an *observable safe path* since node v has not stabilized on P_v yet, *i.e.* node v not in S . Then case 1 applies (*i.e.* node u makes path $P2_u$ its most preferred path, and stabilizes on it). If (iii) happens, node u stabilizes on its $P1_u$ and becomes a *conflict-free node*.

Case 3) If none of the paths advertised by nodes in S are permitted at receiving node u , then nothing will happen and the size of S will stay the same. However, node u will become a *conflict-free node* converging to ϵ , *i.e.* unreachable destination. This also implies that not only immediate neighbors of nodes in S , but all the nodes outside of the set S at that point will converge to ϵ . Then the APM algorithm returns with a stable routing tree. ■

Note that once a node is in S , during the execution of APM, it does not get out of S . This is because as node u enters in S , it updates its path preferences so that its observable safe path is the most preferred path. Therefore, node u will not give up this path for any other path that node u may learn in the future. The only time when node u enters S and does not update its path preferences is when node u is an *inactive node with observable safe path*. However, in such a case, as node u enters S , it stabilizes on its best path.

Theorem 4.1: Starting from an arbitrary state of the system, the Adaptive Policy Management (APM) converges to a stable state within a finite number of steps with a reasonable probability. In the worst case, the number of steps is $(|V| - 1)$, where $|V|$ is the number of nodes in the topology, and the probability is $(1/2)^{k+l+\dots+((|V|-1)/2)}$, where k, l, \dots, i are

positive integers.

Proof: By Lemma 4.2, we can show that APM runs in a finite number of steps. Since Lemma 4.2 shows that the size of the nodes in the conflict-free set, S , monotonically increases, at each step of the algorithm, the size of the set of nodes yet to be explored, *i.e.* $\{V - S\}$, must be decreasing monotonically too. Since there are only $|V|$ nodes, after a finite number of steps, the algorithm converges with all nodes moving to the set S . However, this is true only if APM performs path rank changes (*i.e.* *conflicting safe-alternative nodes* stick to their lower ranked but safe path) whenever a conflict is detected.

If the rank change is done probabilistically, then reaching a stable state in a finite number of steps will take longer. The worst case happens when each conflict in the system has only two *conflicting safe-alternative nodes*, and the conflicts in the system are nested and independent from each other as shown in Figure 12. Since the nodes use probability $1/2$ for rank change, the worst-case probability of breaking a cycle with only 2 *conflicting safe-alternative nodes* happens when only one of such nodes performs rank change and this probability is $1/2$. Once the rank change is performed and conflict is resolved, both nodes will be *conflict-free*. Since the expected number of tries for successfully realizing this rank change is 2, in the worst case, the average number of steps is $(|V| - 1)$.

Since the probability of breaking the current innermost conflict after j attempts is $(1/2)^j(1/2)$, in the worst case, the probability of reaching a conflict-free state is $((1/2)^k(1/2))((1/2)^l(1/2)) \dots ((1/2)^i(1/2)) = (1/2)^{k+l+\dots+i+(|V|-1)/2}$, where k, l, \dots, i are positive integers.

However, such worst-case scenario is not very realistic. In more practical cases, breaking a conflict may simultaneously resolve more than one conflict, or several independent conflicts may break simultaneously, which shortens the convergence time. ■

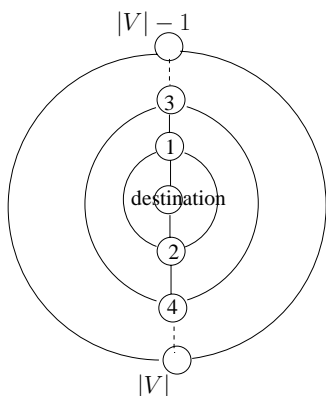


Fig. 12. Worst-case scenario for APM when there is only 2 conflicting safe-alternative nodes in each conflict and the conflicts are nested and independent from each other. The depth is $(|V| - 1)/2$.

The analysis presented in this section does not consider adaptation and *recovery* of path rank changes. The next section (Section V) serves this purpose.

V. SIMULATION RESULTS

We have simulated the algorithms in the SSFNet simulator [15]. We present two sets of results for two different topologies, which are presented in subsection V-A and subsection V-C, respectively. We first define our performance metrics:

- The *average of the percentage of paths that are eliminated per node at time t among permitted paths* to provide stability. The smaller value of this metric indicates better performance, since eliminating permitted paths (*i.e.* moving them to the *bad paths* set) may strain reachability, or force the router to choose a less preferred path to reach a destination.

- The *average of the percentage of the paths whose rank has been changed per node at time t* . Since changing the rank of the paths means changing locally configured policies that have been carefully placed for specific purposes, an algorithm causing a lot of rank changes would be undesirable.

- The *average of the percentage of the preference loss per node at time t among permitted paths*. Preference loss of a path is the difference between its original local preference and current preference value. If a path is placed in the bad path set, its preference loss is equal to its original local preference value. This metric helps us quantify the total effect of both path elimination and rank change.

- The *number of updates exchanged between routers* is an indication of stability. When the system is not stable, the routers constantly exchange update messages. Therefore, smaller number of exchanged update messages reflects the efficiency of the protocol dealing with conflicts. To compute this metric, we have measured the average number of updates carried over the last 2000 seconds.

- The *number of bytes carried by update messages* is used to evaluate the overhead of the algorithms. Longer update messages takes longer to process and transmit. This overhead may negatively affect the overall performance of the system. We computed this metric by measuring the average number of bytes carried over the last 2000 seconds.

- The *average extra storage used at time t (in bytes)* is another metric for evaluating the overhead of the algorithms. For BGP4 [23], the value of this metric is always zero. For SPVP, *history* is the newly added path attribute, and the main contributor of extra storage in routing tables, *i.e.* *rib*, *rib_in* and *rib_out*. The other source of extra storage is due to the bad path set since BGP4 does not have such set. The extra storage for the Cobb and Musunuri algorithm is the per peer *cost* kept at each node, which indicates the number of times a node has observed a route flap. For APM, the size of *local history*, and *bad path set* are the main contributors of extra storage. We also added per peer *peerStability* and *keepaliveCount* as well as the list of *sequenceNumbers* that a node has learned during the execution of the algorithm, which are just integers. Depending on the occurrence of topology changes, the contribution of *sequenceNumbers* may be zero.

- *Power* is used to measure the ratio of throughput (average total number of packets delivered over the last 50 seconds) and delay (average delay of delivered packets over the last 50 seconds). The power metric captures the desire of achieving as high throughput as possible while keeping delay as small as possible.

- Percentage of nodes that cannot reach the destination at time t is also measured.

The performance plots presented next show 90% confidence intervals for these metrics.

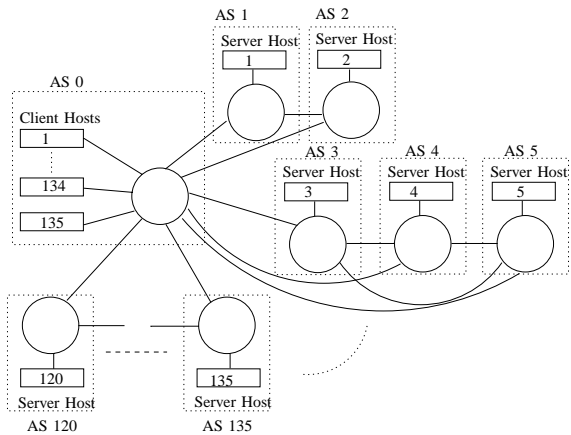


Fig. 13. Topology of Simulation Set I

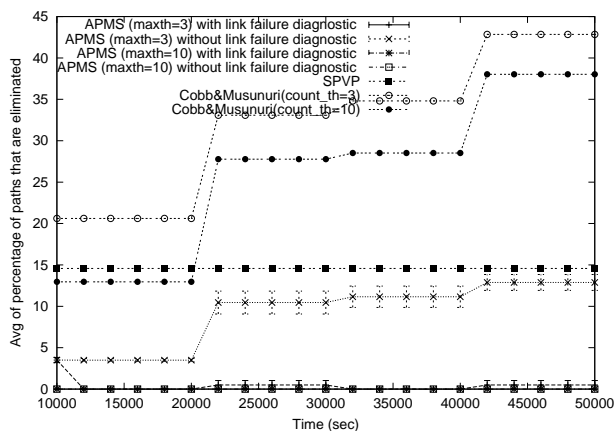


Fig. 14. Average Percentage of Paths Eliminated per Node

A. Simulation Set I

The topology is shown in Figure 13, and consists of 15 independent dispute wheels, where each AS in a dispute wheel has a direct connection to the destination AS 0. The destination AS has 135 client hosts, to whom there is constant data flow from the servers located in the other ASes. Each router within each AS has 3 permitted paths: The path through its clockwise neighbor, the direct path, and the path through its counter-clockwise neighbor. The policies are set to create policy conflicts, *i.e.* each AS prefers going through its clockwise neighbor rather than its direct path, which is preferred over going through the counter-clockwise neighbor—Local preference values assigned at each node are 100, 80, and 40, respectively.

Simulation is run for 50000 seconds, and data flow from servers to clients continues for the whole duration. We also introduced periodic link failures, during which all ASes lose their connection to the destination AS 0. After the system stabilizes at 10000 seconds, link failures are introduced.

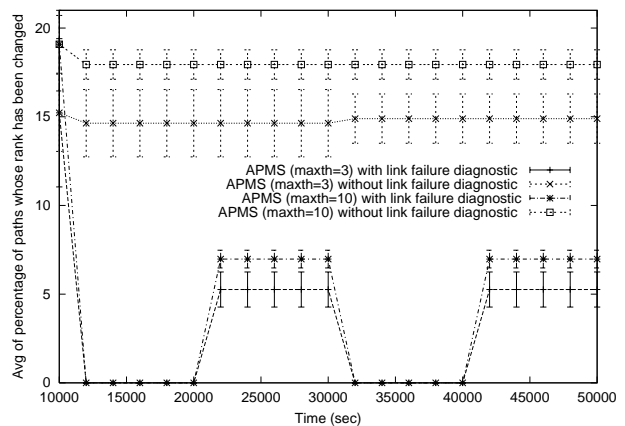


Fig. 15. Average Percentage of Paths Whose Rank Changed per Node

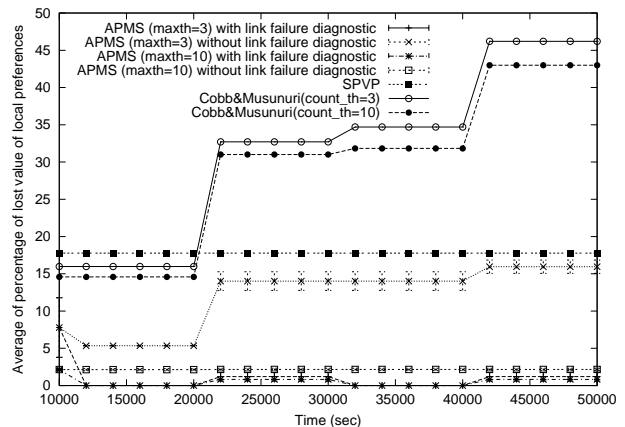


Fig. 16. Average Percentage of Lost Preference Value per Node

Recoveries and failures are then scheduled alternately every 10000 seconds.

The variations of APM include using different values for $max_threshold$ of 3 and 10, and whether route flaps due to topology change are distinguished. We set $min_threshold$ to 2, and $ka_threshold$ to 6. We have compared APM against the SPVP [7], the Cobb and Musunuri algorithm [9], and BGP4 [23], where the details of these algorithms can be found in Section II-B. We have two versions of the Cobb and Musunuri algorithm, where the threshold for cost is set to either 3 or 10 to be consistent with the values assigned to $max_threshold$ for the variations of APM. With SPVP, since there is no built-in mechanism to differentiate between transient oscillations and persistent oscillations due to policy conflicts, Griffin *et al.* [7] suggest suppressing routes only after they are seen to contain some number of dispute cycles, or after the length of the history has exceeded some limit. In our simulations, we have used the first approach, and suppressed the routes only after seeing the same policy cycle twice. This is consistent with the value we have chosen for APM, for which $min_threshold$ is assigned a value of 2.

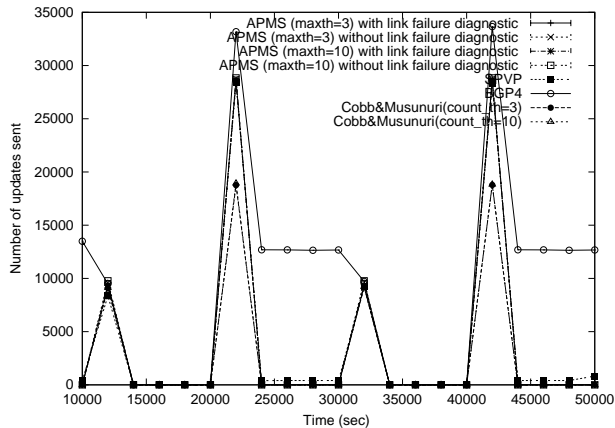


Fig. 17. Number of Updates Messages Sent

B. Results

The results for the average of the percentage of paths that are eliminated per node at time t is shown in Figure 14. BGP4 does not resolve conflicts, hence does not eliminate any paths. Therefore, the value of this metric is zero for BGP4 and not shown in the figure. In general, APM resolves conflicts by means of path rank change instead of path elimination. APM waits to see the same route flap $max_threshold$ times before eliminating the path involved in the route flap. When $max_threshold$ is larger, the system stays in the policy conflict-avoidance phase longer and tries harder to resolve conflicts through path rank change. When $max_threshold$ is smaller, the system enters the policy conflict-control phase sooner, which causes elimination of more paths. With SPVP however, eliminating some of the flapping paths is its only way to deal with policy conflicts. Therefore, the performance of SPVP is worse than APM for both small and larger values of the $max_threshold$. With SPVP, on average each node eliminates 14.4% of the available paths. Cobb and Musunuri algorithm eliminates paths advertised by the peer whose cost has reached the threshold, *i.e.* 3 or 10 for this experiment. Both versions of the Cobb and Musunuri algorithm eliminate most of the paths, upto 42.84% (for threshold=3) and 38.02% (for threshold=10), because of both aggregation of the paths through the high cost nodes and absence of any mechanism to distinguish route flaps due to topology change (false positives).

The version of APM lacking any topology change diagnosis and using a smaller value of $max_threshold$, 3, keeps eliminating the available paths since there is no mechanism to differentiate route flaps due to failure and/or recovery of the links from those triggered by policy conflicts. In this case, paths seem to be flapping more often, and *count* values increase faster. For the same value of $max_threshold$, adding topology change diagnosis to the algorithm provides big improvement: On average each node eliminates only 0.48% of the available paths to reach stability during non-fail periods. When $max_threshold$ is larger, resolving conflicts by means of path rank change dominates that by path elimination. As we can see, this leads to minimal path elimination for both versions of APM; with and without topology change diagnosis. However, this result remains valid for APM without topology

change diagnosis only when the $max_threshold$ value is larger than the number of route flaps resulting from topology change. Otherwise, there will be false positives, and more path elimination.

For the previous metric, we have seen that for large enough value of $max_threshold$, APM can avoid eliminating paths performing path rank changes to resolve policy conflicts. However for APM, we would like to see whether the mechanism is causing a lot of path rank changes, and hence significantly altering the policies that have been carefully placed for specific purposes. Figure 15 shows the results for the average of the percentage of the paths whose rank has been changed per node. Since APM is the only algorithm that changes the policies to stabilize the system, for SPVP, Cobb and Musunuri algorithm and BGP4, the value of this metric is 0 and therefore not shown in the figure. For larger values of $max_threshold$, we observe higher number of path rank changes due to longer policy conflict-avoidance phase. Using topology change diagnosis improves performance in the presence of failure and recovery of the links, and drops the metric value from 18% to 7.0% for non-fail periods. This corresponds to changes in the rank of only about 2 paths out of about 20 available paths per dispute wheel on average (plots not shown for lack of space), without eliminating a single path. Smaller value of $max_threshold$ shows lower percentage of path rank change due to shorter policy conflict-avoidance phase. However, as we have seen in Figure 14, this version of the algorithm eliminates more paths to deal with conflicts.

The results for the average of the percentage of the preference loss per node is shown in Figure 16. Since BGP4 has no mechanism to deal with conflicts, there is no loss in terms of preference value, and hence BGP4 is not shown in the figure. SPVP causes loss of around 18%, which is contributed only by eliminated paths. Cobb and Musunuri algorithm causes loss of 46.2% for cost threshold 3, and 42.0% for cost threshold 10, which are contributed only by eliminated paths. With APM, the performance is always better than SPVP and the Cobb and Musunuri algorithm. Using larger values of $max_threshold$, 10, and topology change diagnosis significantly improves performance to less than 1% loss in path preference.

Figure 17 presents the results for the number of update messages sent for each 2000 interval of time. Topology changes in the form of link failure or restoration cause a burst of updates, and the burst is smaller in the case of link failures due to limited reachability. When each node loses its reach to the destination, Hold Timer expires, and all of the paths to the destination become infeasible and are withdrawn. When the links are restored, BGP sessions are re-established, and the whole routing tables are exchanged. The biggest routing table is of node 0, since it has a BGP session with every other node, and this table contributes a very big portion of the high peaks in the graph right after the link restorations (at times 20000, 40000, ...). For topology changes in the form of link failures, we observe a smaller burst of updates due to restricted reachability: After the first 2000 seconds of each fail period, the number of updates sent is zero for all algorithms. This is because the new topology does not have any policy conflicts,

and therefore it stabilizes independent of the algorithm used. For non-fail periods, since with BGP4 the system does not stabilize, the number of updates sent under BGP4 does not get close to zero. SPVP and APM show very close performance regarding this metric, and the number of updates sent is much smaller than BGP4. Since reachability is the most restricted with both versions of the Cobb and Musunuri algorithm (see Figure 14), they have the lowest value for this metric.

To evaluate of the overhead of the algorithms, *the number*

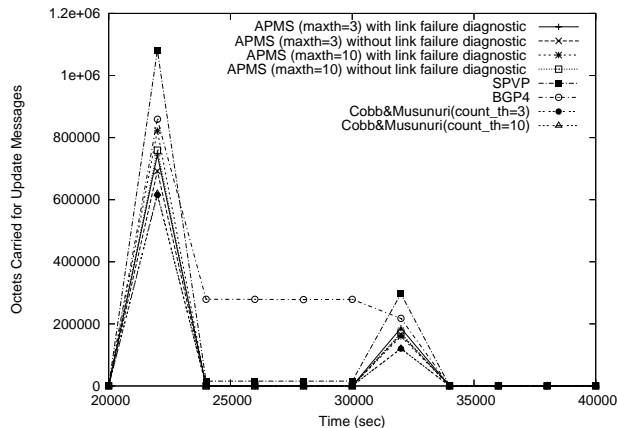


Fig. 18. Number of Bytes Carried by Update Messages

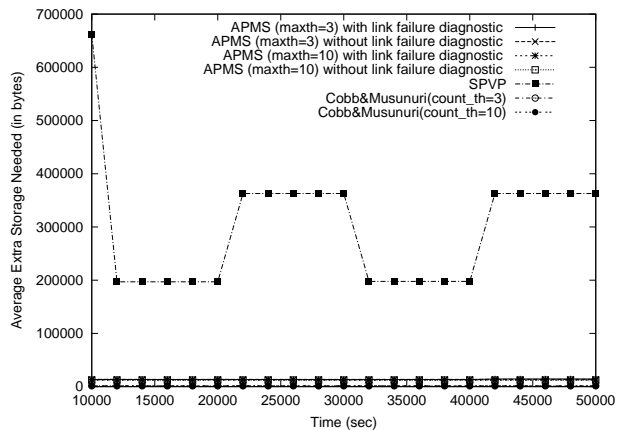


Fig. 19. Average Extra Storage Needed at Each Node by the Algorithms (in bytes)

of bytes carried by update messages is shown in Figure 18. In Figure 17, we have seen that with SPVP the total number of updates exchanged was much smaller than BGP4, and very close to APM. However, the divergence detection mechanism of SPVP requires carrying the sequence of path change events in each update, *i.e.* *history*. Thus, SPVP has the highest number of bytes carried by its update messages. All versions of APM show very close performance, and all are better than both BGP4 and SPVP. When topology change diagnosis is deployed in APM, the update messages carry some extra information in the *originator* and *sequenceNumber* fields (6 bytes in total). However, as we can observe, the mechanism used by APM to distinguish temporary oscillations (due to topology/policy changes) is much more efficient than the SPVP mechanism of observing repeated cycles in the history. Although the only

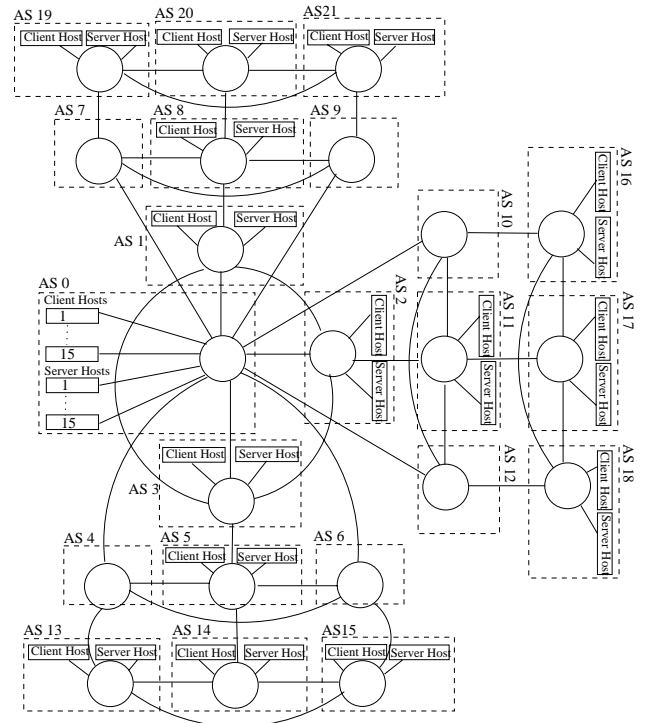


Fig. 20. Topology Used for Simulation Set II

Node	Permitted Paths	Local Preference	Node	Permitted Paths	Local Preference
1	(1 2 0)	100	13	(13 15 6 0)	100
	(1 0)	80		(13 15 6 5 3 0)	100
	paths learned from 8	1		(13 15 6 5 3 1 0)	100
				(13 4 0)	80
2	(2 3 0)	100	14	(14 13 4 0)	100
	(2 0)	80		(14 13 4 6 0)	100
	paths learned from 11	1		(14 5 3 0)	80
				(14 5 3 1 0)	80
3	(3 1 0)	100	15	(15 14 5 3 0)	100
	(3 0)	80		(15 14 5 3 1 0)	100
	paths learned from 5	1		(15 14 5 4 0)	100
				(15 6 0)	80
4	(4 6 0)	100	5	(15 6 5 3 0)	80
	(4 0)	80		(15 6 5 3 1 0)	80
	paths learned from 13	1			
5	(5 4 0)	100	6	(6 5 3 0)	100
	(5 3 0)	80		(6 5 3 1 0)	100
	(5 3 1 0)	80		(6 0)	80
	paths learned from 14	1		paths learned from 15	1
6	(6 5 3 0)	100	7	(7 8 1 2 0)	100
	(6 5 3 1 0)	100		(7 8 1 0)	100
	(6 0)	80		(7 0)	80
	paths learned from 15	1		paths learned from 19	1
7	(7 8 1 2 0)	100	8	(8 9 0)	100
	(7 8 1 0)	100		(8 1 0)	80
	(7 0)	80		(8 1 2 0)	80
	paths learned from 19	1		paths learned from 20	1
8	(8 9 0)	100	9	(9 7 0)	100
	(8 1 0)	80		(9 0)	80
	(8 1 2 0)	80		paths learned from 21	1
	paths learned from 20	1			
9	(9 7 0)	100	10	(10 11 2 0)	100
	(9 0)	80		(10 11 2 3 0)	100
	paths learned from 21	1		(10 0)	80
				paths learned from 16	1
10	(10 11 2 0)	100	11	(11 12 0)	100
	(10 11 2 3 0)	100		(11 2 0)	80
	(10 0)	80		(11 2 3 0)	80
	paths learned from 16	1		paths learned from 17	1
11	(11 12 0)	100	12	(12 10 0)	100
	(11 2 0)	80		(12 0)	80
	(11 2 3 0)	80		paths learned from 18	1
	paths learned from 17	1			
12	(12 10 0)	100	17	(17 18 12 0)	100
	(12 0)	80		(17 18 12 10 0)	100
	paths learned from 18	1		(17 11 2 0)	80
				(17 11 2 3 0)	80
13	(13 15 6 0)	100	18	(18 16 10 0)	100
	(13 15 6 5 3 0)	100		(18 16 10 11 2 0)	100
	(13 15 6 5 3 1 0)	100		(18 16 11 2 3 0)	100
	(13 4 0)	80		(18 12 10 0)	80
14	(14 13 4 0)	100	19	(18 12 0)	80
	(14 13 4 6 0)	100		(19 20 8 9 0)	100
	(14 5 3 0)	80		(19 20 8 1 0)	100
	(14 5 3 1 0)	80		(19 20 8 1 2 0)	100
15	(14 5 4 0)	80	20	(19 7 0)	80
	(15 14 5 3 0)	80		(19 7 8 1 2 0)	80
	(15 14 5 3 1 0)	80		(19 7 8 1 0)	80
	(15 6 0)	80		(20 21 9 0)	100
16	(15 6 5 3 0)	80	21	(20 21 9 7 0)	100
	(15 6 5 3 1 0)	80		(20 8 9 0)	80
	(16 17 11 12 0)	100		(20 8 1 0)	80
	(16 17 11 2 0)	100		(20 8 1 2 0)	80
17	(16 17 11 2 0)	100	19	(21 19 7 0)	100
	(16 17 11 2 3 0)	100		(21 19 7 8 1 2 0)	100
	(16 10 0)	80		(21 19 7 8 1 0)	100
	(16 10 11 2 0)	80		(21 9 0)	80
18	(16 10 11 2 0)	100	20	(21 9 7 0)	80
	(16 11 2 3 0)	100			
	(18 12 10 0)	80			
	(18 12 0)	80			
19	(18 12 0)	80	21	(21 19 7 0)	100
	(19 20 8 9 0)	100		(21 19 7 8 1 2 0)	100
	(19 20 8 1 0)	100		(21 19 7 8 1 0)	100
	(19 20 8 1 2 0)	100		(21 9 0)	80
20	(19 7 0)	80	21	(21 9 7 0)	80
	(19 7 8 1 2 0)	80			
	(19 7 8 1 0)	80			
	(20 8 1 0)	80			
21	(20 8 1 2 0)	80	21	(21 19 7 0)	100
	(20 8 1 0)	80		(21 19 7 8 1 2 0)	100
	(20 8 1 2 0)	80		(21 19 7 8 1 0)	100
	(20 8 1 0)	80		(21 9 0)	80

Fig. 21. Path Rankings for Topology Used in Simulation Set II

extra information carried with update messages in the Cobb and Musunuri algorithm is *cost* associated with each node, this is not the only reason for its best performance. Cobb and Musunuri algorithm has the lowest value because of the least number of exchanged update messages due to limited reachability (see Figure 17 and Figure 14), which is a result of both aggregation (and elimination) of all paths through high cost nodes and absence of any mechanism to distinguish route flaps due to topology change (false positives).

Figure 19 shows the results for *the average extra storage used at time t (in bytes)*. Due to the size of *history*, the storage required by SPVP is much larger than that required by APM and Cobb and Musunuri algorithm. For SPVP, the value of the metric is higher for non-fail periods than fail periods due to better reachability during the non-fail periods. For this experiment, while APM requires around 10KB extra storage, SPVP requires 200KB-360KB extra storage. APM requires a little bit more extra storage than the Cobb and Musunuri algorithm due to bigger local state kept at each node.

C. Simulation Set II

To be able to observe *throughput* and *delay* better, in the second setup we have used a smaller topology shown in Figure 20. The topology consists of 7 dispute wheels, where each AS in a dispute wheel prefers the path through its clockwise neighbor. The wheels consist of the following group of ASes: {AS 1, AS 2, AS 3}, {AS 4, AS 5, AS 6}, {AS 7, AS 8, AS 9}, {AS 10, AS 11, AS 12}, {AS 13, AS 14, AS 15}, {AS 16, AS 17, AS 18}, {AS 19, AS 20, AS 21}. Permitted paths and path preferences are shown in Figure 21. Simulation is run for 350

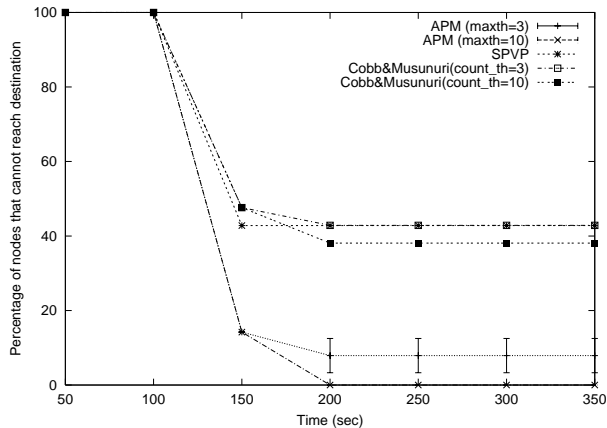


Fig. 22. Percentage of nodes that cannot reach the destination

seconds, and data flow from servers to clients continues for the whole duration. For this experiment set, there is no topology change. Buffer size is 50000 bytes, and routing packets are given priority over data packets when there is congestion at the buffers. Our findings are presented in Section V-D.

D. Results

Figure 22 shows the percentage of the nodes that cannot reach the destination AS 0. SPVP and Cobb and Musunuri algorithms eliminate a high number of paths while enforcing

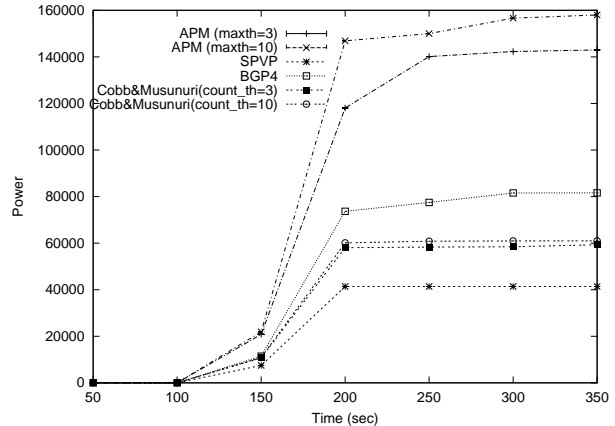


Fig. 23. Power

stability, and therefore leave a higher number of nodes with unreachable destination. Different versions of APM perform much better than SPVP and the Cobb and Musunuri algorithms. With APM, higher *max_threshold* value helps forcing conflict resolution through changing path preferences, and therefore minimizing the number of path eliminations. For *max_threshold=10*, the system stabilizes to a state where each node has a way to reach the destination. Higher threshold value also helps the Cobb and Musunuri algorithm achieve better performance, but the improvement is not much because of the simultaneous elimination of the paths through the same high cost node.

Figure 23 shows the results for *power*. Different versions of APM have higher value than SPVP and the Cobb and Musunuri algorithms because APM maximizes throughput. The different performance for throughput stems from both unreachable destinations, and/or competition for the limited buffer size. SPVP and the Cobb and Musunuri algorithms leave a higher number of nodes with unreachable destination (Figure 22). SPVP has the longest update messages, which take longer to process and require more memory to be stored in the buffers. Although BGP4 causes constant exchange of updates due to divergence, its performance is better than SPVP and the Cobb and Musunuri algorithms! This is because BGP4 does not cause permanent path elimination, even though some packets may not reach the destination temporarily due to instability.

VI. SUMMARY AND FUTURE WORK

Unlike static solutions (e.g. Gao *et al.* [6] algorithm) which may lead to unnecessary disallowance of the usage of many routes from the start to guarantee the stability of the system, APM is a dynamic algorithm, and allows ASes to adapt to the current state of the network, either conflict-free or potentially conflicting. APM makes minimal changes to local policies as the primary means of removing policy conflicts (and associated routing oscillations). Path elimination with APM happens only if probabilistic rank change of paths does not resolve the conflict³. APM attempts to keep as many paths as possible to have better connectivity, and more flexibility in

³Since the approach is probabilistic, there is a non-zero probability that none of the nodes change the rank of their paths.

path selection for the stabilized system. APM is distributed, and does not require a global authority or global database.

Compared to other dynamic algorithms (e.g. SPVP [7] and Cobb and Musunuri [9]), APM has several advantages: (1) APM eliminates the need to carry possibly large amount of information like *history* in the update messages. Instead, APM uses local state information. Therefore, with APM there is no communication overhead, nor any concerns about revealing private information about the preferences of ASes over their routes; (2) APM minimizes path elimination by path rank changes and by adapting to a conflict-free state by (conservatively) restoring some of the original preferences as well as eliminated paths; (3) APM has a more effective mechanism for clearly distinguishing route flaps due to topology change, which helps minimize false positives and communication overhead; (4) APM automatically adapts to the dynamics of the system by observing either path changes or keepalive messages without requiring an expensive protocol, such as diffusing computation [9]; and (5) APM deals with the problem at the path level instead of node level, which prevents aggregation of paths through the same node, and hence elimination of many paths whose preference does not cause conflicts.

Route flap damping [20] is not an alternative solution to APM, and only effective in suppressing temporary instabilities. APM is still needed on top of route flap damping to detect and resolve policy conflicts leading to divergence and/or adapt to dynamically changing network state. Usage of route flap damping may delay diagnosing persistent oscillations via APM due to suppressed updates. If there is a topology change, APM needs to be aware of this change so it handles the resulting route flaps differently from the flaps resulting from policy conflicts until the system adapts to the new network state. However, there may be no need to react and adapt to every topology change. For example, an unstable link may go down and up constantly for a period of time. In such cases, it is better that APM does not react to such short-term changing link status. APM may take route flap damping into consideration to hide such shorter term topology changes from other ASes.

If only some of the nodes in the network upgraded to deploy APM, APM still can catch and resolve policy conflicts since the algorithm is based on only local information kept at each node. However, in such heterogenous settings, the nodes deploying APM will be the only ones which may give up their preferred paths for the sake of the network stability without knowing whether or not the other nodes are working for the same purpose. As future work, to be able to prevent selfish behavior, and improve cooperation among ASes to deploy APM, incentives should be proposed.

To increase transparency of APM, we plan to investigate allowing local AS administrators to explicitly guide the backoff and recovery probabilities for lowering and restoring the rank of paths. With such a human impact, it may be possible to resolve policy conflicts in a more efficient way albeit at a longer time scale.

APM is explained using an abstract model of BGP (in Section III). Although the implementation of APM in SSFNet

relaxes some of these abstract model assumptions (Section V), there are still many details to be addressed such as IGP/BGP interactions, duplicate updates, session resets etc. These details are left for future work.

REFERENCES

- [1] K. Varadhan, R. Govindan, and D. Estrin, "Persistent Route Oscillations in Inter-Domain Routing," *Computer Networks*, vol. 32:1-16, 2000.
- [2] C. Labovitz, G. Malan, and F. Jahanian, "Internet Routing Instability," *IEEE/ACM Transactions on Networking*, vol. 6, no. 5, pp. 515-528, 1997.
- [3] R. Govindan and A. Reddy, "An Analysis of Interdomain Routing Topology and Route Stability," in *Proceedings of the Conference on Computer Communications (IEEE Infocom)*, Kobe Japan, April 1997.
- [4] V. Paxson, "End-to-end Routing Behavior in the Internet," *IEEE/ACM Transactions on Networking*, Vol.9, No.4, pp. 392-403, August 2001.
- [5] R. Govindan, C. Alaettinoglu, G. Eddy, D. Kessens, S. Kumar, and W. Lee, "An Architecture for Stable, Analyzable Internet Routing," *IEEE Network*, vol. 13(1):29-35, 1999.
- [6] L. Gao and J. Rexford, "Stable Internet Routing without Global Coordination," in *Proceedings of ACM SIGMETRICS*, Santa Clara CA, June 2000.
- [7] T. Griffin and G. Wilfong, "A Safe Path Vector Protocol," in *Proceedings of IEEE INFOCOM*, Tel Aviv Israel, March 2000.
- [8] J. A. Cobb, M. G. Gouda, and R. Musunuri, "A Stabilizing Solution to the Stable Paths Problem," *Symposium on Self-Stabilizing Systems, Springer Verlag Lecture Notes in Computer Science*, vol. 2704, pp. 169-183, 2003.
- [9] J. A. Cobb and R. Musunuri, "Enforcing Convergence in Inter-Domain Routing," in *Proceedings of IEEE Global Communications (GLOBECOM) Conference*, Dallas TX, December 2004.
- [10] T. Griffin, A. Jaggard, and V. Ramachandran, "Design Principles of Policy Languages for Path Vector Protocols," in *Proceedings of ACM SIGCOMM*, August 2003.
- [11] A. Jaggard and V. Ramachandran, "Robustness of Class-based Path-Vector Systems," in *Proceedings of IEEE International Conference on Network Protocols (ICNP)*, March 2004.
- [12] J. Sobrihho, "Network Routing with Path Vector Protocols: Theory and Applications," in *Proceedings of ACM SIGCOMM*, Karlsruhe Germany, August 2003.
- [13] S. Yilmaz and I. Matta, "A Randomized Solution to BGP Divergence," in *Proceedings of the 2nd IASTED International Conference on Communication and Computer Networks (CCN'04)*, Cambridge MA, November 2004.
- [14] V. Jacobson, "Congestion Avoidance and Control," in *ACM SIGCOMM '88*, Stanford CA, August 1988, pp. 314-329.
- [15] SSFNet: Scalable Simulation Framework, "<http://www.ssfnet.org>."
- [16] T. Griffin, F. Shepherd, and G. Wilfong, "Policy Disputes in Path-Vector Protocols," in *Proceedings of IEEE International Conference on Network Protocols (ICNP)*, Toronto Canada, November 1999.
- [17] T. Griffin and G. Wilfong, "An Analysis of BGP Convergence Properties," in *Proceedings of ACM SIGCOMM*, Cambridge MA, September 1999.
- [18] L. Gao, T. Griffin, and J. Rexford, "Inherently Safe Backup Routing with BGP," in *Proceedings of IEEE INFOCOM*, Anchorage AK, April 2001.
- [19] N. Feamster, H. Balakrishnan, and J. Rexford, "Some Foundational Problems in Interdomain Routing," in *ACM SIGCOMM Workshop on Hot Topics in Networking (HOTNets-III)*, San Diego CA, November 2004.
- [20] C. Villamizar, R. Chandra, and R. Govindan, "BGP Route Flap Damping," RFC 2439, 1998.
- [21] A. Arora and M. Gouda, "Distributed Reset," *IEEE Transactions on Computers*, vol. 43, pp. 1026-1038, 1994.
- [22] J. Feigenbaum, R. Sami, and S. Shenker, "Mechanism Design for Policy Routing," in *Proceedings of the Symposium on Principles of Distributed Computing (PODC)*, Newfoundland Canada, July 2004.
- [23] Y. Rekhter and T. Li, "A Border Gateway Protocol," RFC 1771, 1995.