

Can SDN Accelerate BGP Convergence?

A Performance Analysis of Inter-domain Routing Centralization

Pavlos Sermpezis
FORTH, Greece
sermpezis@ics.forth.gr

Xenofontas Dimitropoulos
FORTH / University of Crete, Greece
fontas@ics.forth.gr

Abstract—The Internet is composed of Autonomous Systems (ASes) or domains, i.e., networks belonging to different administrative entities. Routing between domains/ASes is realised in a distributed way, over the Border Gateway Protocol (BGP). Despite its global adoption, BGP has several shortcomings, like slow convergence after routing changes, which can cause packet losses and interrupt communication even for several minutes. To accelerate convergence, inter-domain routing centralization approaches, based on Software Defined Networking (SDN), have been recently proposed. Initial studies show that these approaches can significantly improve performance and routing control over BGP. In this paper, we complement existing system-oriented works, by analytically studying the gains of inter-domain SDN. We propose a probabilistic framework to analyse the effects of centralization on the inter-domain routing performance. We derive bounds for the time needed to establish data plane connectivity between ASes after a routing change, as well as predictions for the control-plane convergence time. Our results provide useful insights (e.g., related to the penetration of SDN in the Internet) that can facilitate future research. We discuss applications of our results, and demonstrate the gains through simulations on the Internet AS-topology.

I. INTRODUCTION

The Border Gateway Protocol (BGP) is globally used, since the early days of the Internet, to route traffic between *Autonomous Systems* (ASes) or *domains*, i.e., networks belonging to different administrative entities. BGP is a distributed, shortest path vector protocol, over which ASes exchange routing information with their neighbors, and establish route paths.

Although BGP is known to suffer from a number of issues related to security [1], [2], or slow convergence [3], [4], [5], deployment of other protocols or modified versions of BGP is difficult, due to its widespread use, and the entailed political, technical, and economic challenges. Hence, any advances and proposed solutions, should be seamless to BGP.

Taking this into account, it has been proposed recently that Software Defined Networking (SDN) principles could be applied to improve BGP and inter-domain routing [6], [7], [8], [9], [10], [11]. The SDN paradigm has been successfully applied in enterprise (i.e., *intra-AS*) networks, like LANs, data centers, or WANs (e.g., Google). However, its application to inter-domain routing (i.e., between different ASes) has to overcome many challenges, like the potential unwillingness of some ASes to participate in the routing centralization. For instance, a small ISP might not have incentives (due to the high investment costs) to change its network configuration. This led previous works on inter-domain SDN to consider

(a) partial deployment, only by a fraction of ASes, and (b) interoperability with BGP.

The proposed solutions have demonstrated that bringing SDN to inter-domain routing can indeed improve the convergence performance of BGP [12], offer new routing capabilities [6], or lay the groundwork for new services and markets [13], [7]. However, most of previous works are system-oriented: they propose new systems or architectures, and focus on design or implementation aspects. Hence, despite some initial evaluations (e.g., experiments, emulations, simulations) we still lack a clear understanding about the interplay between inter-domain centralization and routing performance.

To this end, in this paper, we study *in an analytic way* the effects of centralization on the performance of inter-domain routing. We focus on the potential improvements on the (slow) BGP convergence, a long-standing issue that keeps on concerning industry and researchers [14]. Our goal is to complement previous (system-oriented) works, obtain an analytic understanding, and answer questions such as: “*To what extent can inter-domain centralization accelerate BGP convergence? How many ASes need to cooperate (partial deployment) for a significant performance improvement? Is the participation of certain ASes more crucial? Will all ASes experience equal performance gains?*” Specifically, our contributions are:

- We propose a model (Section II) and methodology (Sections III and IV) for the performance analysis of inter-domain routing centralization. To our best knowledge, we are the first to employ a probabilistic approach to study the performance of inter-domain SDN.
- We analyse the time that the network needs to establish connectivity after a routing change. In particular, we derive upper and lower bounds for the time needed to achieve data-plane connectivity between two ASes (Section III), and exact expressions and approximations for the time till control-plane convergence over the entire network (Section IV). Our results are given by closed-form expressions, as a function of network parameters, like network size, path lengths, and number of SDN nodes.
- Based on the theoretical expressions, as well as on extensive simulation results, we provide insights for potential gains of centralization, inter-domain SDN deployment strategies, network economics, etc.

We believe that our study can be useful in a number of directions. Research in inter-domain SDN can be accelerated

and facilitated, since a fast performance evaluation with our results can precede and limit the volume of required emulations/simulations. The probabilistic framework we propose can be used as the basis (and be extended and/or modified) to study other problems or aspects relating to inter-domain routing, e.g., BGP prefix hijacking, or anycast. Finally, the provided insights can be taken into account in the design of protocols, systems, architectures, pricing policies, etc.

II. MODEL

A. Network

We consider a network, e.g., the whole Internet or a part of it, that consists of N autonomous systems (ASes). We represent each AS as a *single node* that operates as a BGP router; this abstraction that is common in related literature [3], [12], allows to hide the details of the intra-AS structure and functionality, and focus on inter-domain routing. When two ASes are connected (transit, peering, etc., relation), we consider that a link exists between the corresponding routers, over which data traffic and BGP messages can be exchanged.

B. SDN Cluster

ASes can be ISPs, enterprises, CDNs, IXPs, etc., belong to different administrative entities, and span a wide range of topological, operational, economic, etc., characteristics. As a result, not all ASes should be expected to be willing to cooperate for and/or participate in an inter-domain centralization effort. Routing centralization is envisioned to begin from a group of a few ASes cooperating with each other, e.g., at an IXP location [6], [7]; then, more ASes could be attracted (performance or economics related incentives) to join the group, or form another group.

To this end, we assume that $k \in [1, N]$ ASes, i.e., a fraction of the entire network, cooperate in order to centralize their inter-domain routing. In the remainder, we refer to the set of these k ASes, as the *SDN cluster*¹. To avoid delving into system-specific issues of the centralization implementation, we assume the following setup, which captures main characteristics of several proposed solutions (e.g., [12], [9], [15]), and is generic enough to accommodate future solutions: ASes in the SDN cluster exchange routing information with a central entity, which we call *multi-domain SDN controller*. The multi-domain SDN controller might be an SDN controller that directly controls the BGP routers of the ASes (e.g., as in [12]), or a central server that only provides information or sends BGP messages to the ASes (e.g., similar to [15]).

C. BGP Updates

Each node has a routing table (Routing Information Base, RIB), in which each entry contains an IP Prefix, and the corresponding AS-path (i.e., sequence of ASes) through which this prefix can be reached. RIBs are built from the information received by the neighbor ASes: upon a routing change, the “source” AS (e.g., the AS that originates a prefix) sends BGP

updates to its neighbors to notify them about the change; when an AS receives a BGP update, it calculates the needed updates (if any) for its RIB, and sends BGP updates to its own neighbors. In this way, BGP updates propagate over the entire network, and paths to prefixes are built in a distributed way.

Let us assume that an AS receives a BGP update at time t_1 and forwards it to a neighbor AS at time t_2 . We call *BGP update time*, and denote T_{bgp} , the time between the reception of a BGP update in an AS and its forwarding to a neighbor AS, i.e., $T_{bgp} = t_2 - t_1$. The BGP update times may vary a lot among different ASes and/or connections, since they depend on a number of parameters: routers’ hardware/software (e.g., time to process BGP data and update RIB) and/or configuration (e.g., MRAI timers), intra-AS network structure (e.g., number of routers, topology) and/or operation (e.g., iBGP configuration, intra-AS SDN), etc.

Knowing all these parameters for every AS is not possible, and using (upper) bounds for T_{bgp} would not lead to practical conclusions [3]. Thus, to be able to perform a useful analysis, we follow a probabilistic approach, and model the BGP update times as follows.

Assumption 1 (BGP updates - renewal process). *The BGP update times T_{bgp} are independent and identically distributed random variables, drawn from an arbitrary distribution $f_{bgp}(t)$, with $E[T_{bgp}] = \mu_{bgp}$.*

Under Assumption 1, BGP update times are given by a renewal process. The model is very generic, since it allows to use any valid function $f_{bgp}(t)$, and thus describe a wide range of scenarios with different parameters. Real measurements can be used to make a realistic selection for the distribution $f_{bgp}(t)$, as we show in Appendix A; however, a detailed study for fitting the $f_{bgp}(t)$ is beyond the scope of this paper.

D. Inter-domain SDN Routing

Routing information in the SDN cluster propagates in a centralized way, through the multi-domain SDN controller. When an AS in the SDN cluster receives a BGP update from a neighbor AS (not in the SDN cluster), it forwards the update to the SDN controller. The SDN controller, which is aware of the topology in the SDN cluster and the connections/paths to external ASes, informs every AS in the SDN cluster about the needed changes in the routing paths. The ASes that receive the updated routes from the controller, notify their non-SDN neighbors using the standard BGP mechanism.

Let t_1 be the time that the first AS belonging to the SDN cluster receives a BGP update from a non-SDN neighbor, and t_2 the time till *all* ASes in the SDN cluster have been informed (by the controller) for the BGP updates. We denote as T_{sdn} the time needed for all the SDN cluster to be informed after a member has received a BGP update, i.e., $T_{sdn} = t_2 - t_1$. The times T_{sdn} would depend on the system implementation. However, it was shown that system designs can achieve $T_{sdn} \ll T_{bgp}$ [16]. Hence, in the remainder -for the sake of presentation- we assume that $T_{sdn} \rightarrow 0$. Nevertheless, our results can be easily modified for arbitrary T_{sdn} (even for

¹Although we use the term *SDN*, our framework does not require necessarily that routing centralization is implemented on an SDN architecture.

TABLE I: Important Notation

N	network size (total # of nodes)	
k	SDN cluster size (total # of SDN-nodes)	
T_{bgp}	BGP update time	
$f_{bgp}(t)$	distribution of BGP update times	Assumption 1
d	path length	
k'	# of SDN-nodes on a path	
T_{SD}	data-plane connectivity time in a SD-path	Theorem 1
T_c	BGP convergence time	Theorem 2
T_ℓ	ℓ -partial BGP convergence time	Corollary 1

cases with $E[T_{sdn}] > E[T_{bgp}]$, without this affecting the main conclusions of the study.

E. Preliminaries and Problem Statement

In our analysis, we consider the following setup:

Every node in the network knows at least one (BGP) path to every other node.

A node initiates a routing change that affects the inter-domain routing (e.g., node n_0 in Fig. 1). This could be an announcement or withdrawal of an IP prefix, an interruption of an AS connection (e.g., a link is down), etc. Here, we consider that a node, which we call the ‘‘source node’’, announces a new IP prefix; this routing change affects the entire network, every node will install a path for this prefix in its RIB upon the reception of the BGP update.

Nodes in the SDN cluster, receive route information from the SDN controller, and add an entry in their RIB for the prefix to the source node; even if the path is not established in the node preceding in this path (e.g., in Fig. 1 node n_j might receive the update before node n_{j-1}). In this case only the node in the SDN cluster knows how to route traffic to the new prefix, therefore, if the SDN node sends traffic to the new prefix, this would not necessarily reach the source-node. The connectivity will be established when every AS in the path has been informed about the BGP update.

BGP updates do not propagate backwards in the path; this would create loops or longer paths, which are discarded or not preferred by BGP.

We call ‘‘SD-path’’ the final path, i.e., the shortest conforming to the routing policies, between the source node (‘‘S’’) and another node (‘‘destination’’, or ‘‘D’’).

In the remainder of the paper we investigate the effects of routing centralization on: (a) the data-plane connectivity between the source node (‘‘S’’) and any node (‘‘D’’) in the network, i.e., the time needed till all nodes in an SD-path have installed the updated BGP paths after a routing change (Section III); and (b) the control-plane convergence, i.e. the time needed till the entire network has established the final paths corresponding to the routing change (Section IV).

For ease of reference, we summarize the notation in Table I.

III. DATA-PLANE CONNECTIVITY

A. Analysis

A source node ‘‘S’’ announces a new IP prefix, and SD-path is the final path from S to a ‘‘destination’’ node D; see, e.g.,

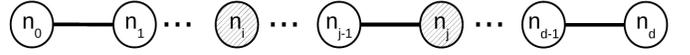


Fig. 1: SD path of size d . The node n_0 initiates the routing change; nodes n_i and n_j belong to the SDN cluster.

Fig. 1. Theorem 1 bounds the expectation of the time T_{SD} needed to establish data-plane connectivity in the path.

Theorem 1. *The expectation of the time T_{SD} in a path of length d with $k' \in [0, d + 1]$ nodes in the SDN cluster, is bounded as follows*

$$LB(d, k') \cdot E[T_{bgp}] \leq E[T_{SD} | d, k'] \leq UB(d, k') \cdot E[T_{bgp}] \quad (1)$$

where

$$LB(d, k') = \begin{cases} 0 & , d \leq k' \leq d + 1 \\ \frac{d}{k'+1} & , 0 \leq k' < d \end{cases} \quad (2)$$

and

$$UB(d, k') = \begin{cases} d - k' + 1 & , 2 \leq k' \leq d + 1 \\ d & , 0 \leq k' < 2 \end{cases} \quad (3)$$

Proof. The proof is given in Appendix B. \square

We provide the intuition behind the proof of Lemma 1 in relation to Fig. 1. When a node in the SD-path that belongs to the SDN cluster receives the BGP update (e.g., node n_i), then every other node in the SDN cluster (e.g., node n_j) is informed about the update, sometimes even before its preceding node(s) (e.g., n_{j-1}). Hence, the BGP update can propagate on *different sections* of the SD-path *simultaneously* (e.g., from n_i up to n_{j-1} , and -at the same time- from n_j to n_d). The length of these SD-path sections (which determine the BGP update propagation time) depend on the positions of the SDN nodes on the path. The bounds are derived based on the ‘‘best’’ and ‘‘worst’’ possible positions of the SDN nodes on the SD-path.

B. Network Topology and Routing Centralization

Based on Theorem 1 we can calculate the average time $E[T_{SD} | d]$ over all paths of the same size d (or, equivalently, for an average path of size d), using the property of conditional expectation:

$$E[T_{SD} | d] = \sum_{i=0}^{d+1} E[T_{SD} | d, k' = i] \cdot P\{k' = i | d\} \quad (4)$$

where $P\{k' = i | d\}$ denotes the probability that i nodes (out of the total $d + 1$ nodes on the path) belong to the SDN cluster.

Topology-independent SDN cluster. If the SDN cluster is formed independently of the network topology, the quantity k' follows an *hypergeometric distribution* with parameters N (population size), k (number of successes in the population), and $d + 1$ (number of draws), and probability mass function

$$P\{k' = i | d\} = \frac{\binom{k}{i} \cdot \binom{N-k}{d+1-i}}{\binom{N}{d+1}} \quad (5)$$

Topology-related SDN cluster. On the other hand, if the participation of ASes in the SDN cluster is related to the topology, e.g., because ASes are explicitly selected based on topological characteristics (e.g., centrality), or the incentives of cooperation are inherently related to their connectivity (e.g., SDN deployment on tier-1 ISPs, or IXPs [6], [7]), then k' might not be captured accurately by Eq. (5). Therefore, the actual distribution $P\{d, k'\}$ needs to be calculated; however, this might be a difficult (or infeasible) task.

Alternatively, in certain cases, the distribution $P\{k' = i|d\}$ could be approximated with variations of the standard hypergeometric distribution that are able to take into account the fact that different nodes appear in shortest paths with different probabilities. For instance the *Fisher's noncentral hypergeometric distribution* can be used to consider biased selection of ASes for the SDN cluster: let ω_i be the betweenness centrality [17] of a node n_i , and ω_{sdn} and ω_{bgp} the averages among the nodes in the respective sets, i.e.,

$$\omega_{sdn} = \frac{\sum_{n_i \in SDN} \omega_i}{|\{n_i : n_i \in SDN\}|}, \quad \omega_{bgp} = \frac{\sum_{n_i \notin SDN} \omega_i}{|\{n_i : n_i \notin SDN\}|}$$

Denoting $\omega = \frac{\omega_{sdn}}{\omega_{bgp}}$, the probability $P\{k' = i\}$ is approximately given by

$$P\{k' = i|d\} = \frac{\binom{k}{i} \cdot \binom{N-k}{d+1-i} \cdot \omega^i}{\sum_{j=0}^{d+1} \binom{k}{j} \cdot \binom{N-k}{d+1-j} \cdot \omega^j} \quad (6)$$

In the above distribution, the higher the betweenness centrality of the ASes in the SDN cluster, the more skewed towards the higher values of k' the distribution $P\{k'|d\}$ is, and, thus, the lower the delay T_{SD} .

Internet AS-topology vs. SDN cluster. We now focus on the Internet topology, which is of higher interest, and apply our -generic- theoretical results to investigate the effects of routing centralization.

We first build the Internet AS graph from a large experimentally collected dataset [18] (consisting of $N = 55567$ ASes), and infer routing policies over existing links based on the Gao-Rexford conditions [19] (this is the most common approach in related literature; more details can be found in Appendix F). We consider about 10^6 different SD-paths, from which we calculate the path length distribution $P\{d\}$ (see Fig. 3), and the betweenness centrality for each node.

We consider different scenarios with variable SDN cluster size $k = 1, \dots, N$, where the set of nodes in the SDN cluster are selected (a) *randomly*, or (b) based on their *betweenness centrality* (i.e., the top k nodes with the highest betweenness centrality values). From Theorem 1, we calculate the lower and upper bounds for the average T_{SD} time over all path lengths, i.e., $E[T_{SD}] = \sum_d E[T_{SD}|d] \cdot P\{d\}$, where $E[T_{SD}|d]$ is given by Eq. (4), and $P\{k'|d\}$ from Eq. (5) or Eq. (6) for the aforementioned cases (a) and (b), respectively.

In Fig. 2, we present the lower (LB) and upper (UB) bounds for $E[T_{SD}]$ for different SDN cluster sizes k , normalized over the case without routing centralization ($k = 0$). When the

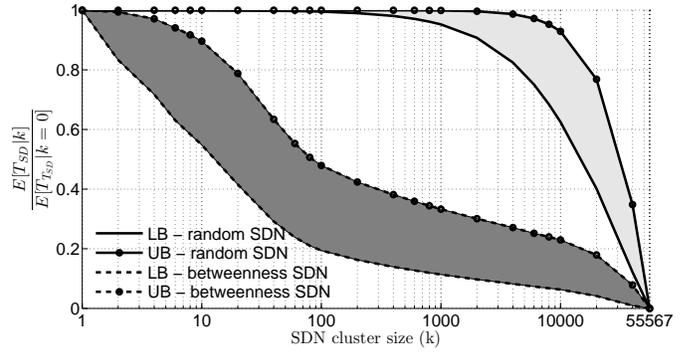


Fig. 2: Bounds for the average data-plane connectivity time, normalized over the no SDN scenario, i.e., $\frac{E[T_{SD}|k]}{E[T_{SD}|k=0]}$, in the Internet AS-graph. Upper (UB) and lower (LB) bounds enclose the colored areas: nodes in the SDN cluster are selected (i) *randomly* (light grey area) and (ii) with decreasing *betweenness centrality* (dark grey area).

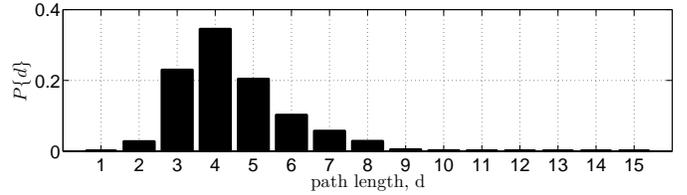


Fig. 3: Path length distribution on the Internet AS-topology.

nodes in the SDN cluster are selected *randomly*, i.e., independently of the topology, a significant decrease in the average connectivity time can be achieved only when at least 20% (around $k = 10000$) of the nodes participate in the SDN cluster (note the log scale of the x-axis). This observation, which is in accordance with previous findings [12], is a rather grim message for the efficiency of (a randomly deployed) inter-domain routing centralization, since even if a few hundreds or thousands of ASes were willing to cooperate, the gains would be marginal.

On the contrary, as is shown in Fig. 2, when the SDN cluster consists of ASes with high *betweenness centrality*, with only a few tens of nodes the average delay can decrease up to 50%. This new insight (compared to previous understanding of the effects of routing centralization) brings optimism for the feasibility of inter-domain centralization: even if deployed incrementally, e.g., starting from a few tier-1 ISPs², the Internet can immediately see significant performance improvements.

C. Simulation Results and Implications

To validate our theoretical results, we conduct simulations on scenarios with varying (a) *network topologies*: synthetic graphs such as full-mesh, Poisson graph, Barabasi-Albert (power law graph), Newman-Watts-Strogatz (small world graph), as well as, the real Internet AS-graph; (b) *SDN cluster sizes*: $k = 0, \dots, N$; and (c) *distributions* $f_{bgp}(t)$: exponential

²Large ISPs are central in the Internet topology, with high betweenness centrality. For example, the top-10 ASes with the highest betweenness centrality values belong to the list of the top-50 ASes with the largest number of ASes in customer cone [20]

with rate $\lambda = 1$ and uniform in $[0, 2]$, both with $\mu_{bgp} = 1$. In the following we present a subset of representative results, and discuss some important observations.

The average values of T_{SD} in the simulations, are *always* within the bounds of Theorem 1 for all pairs $\{d, k'\}$ in every scenario we tested.

In Fig. 4 we compare the simulation results for $E[T_{SD}|d]$ (average over all k') against the theoretical bounds, which are calculated from Eq. (4) by using the expressions of Eq. (5) (topology-independent SDN cluster) and Theorem 1. For both cases of $f_{bgp}(t)$, the bounds are very tight for $k = 50$, when only a small fraction (5%) of the nodes belong to the SDN cluster (top plots). For larger SDN cluster sizes ($k = 200$, or 20%; bottom plots), the bounds are still very tight for small path lengths (e.g., $d < 4$), while the range $[lower\ bound, upper\ bound]$ increases with d . In summary, the accuracy of the bounds increases for smaller k or d .

For $k = 200$ and $d = 7$ (rightmost points in bottom plots), while the simulated value lies in the middle of the two bounds in the exponential $f_{bgp}(t)$ case (Fig. 4(a)), it is closer to the upper bound in the uniform $f_{bgp}(t)$ case (Fig. 4(b)). Among all the scenarios we tested, we did not observe any tendency of the values to be closer either to the upper or lower bound. This is an indication that there is probably a limit on how much tighter bounds can be derived.

In Table II, we show how the times T_{SD} change for increasing SDN cluster size k . Comparing the two cases, $d = 2$ and $d = 5$, we can see that the effect of the routing centralization is higher for longer paths. The simulated data-plane connectivity times decrease more and faster for $d = 5$, and this is captured also by the relative changes of the theoretical bounds.

Similar behavior is observed also in Fig. 5, in simulation scenarios on the Internet topology where the SDN cluster comprises nodes with high betweenness centrality. For $k = 10$, paths of length $d = 3$, $d = 6$, and $d = 9$, see a relative decrease on the average connectivity time of about 10%, 20%, and 40%, respectively. The corresponding values for $k = 50$ are about 25%, 40%, and 60% (i.e., almost double than $k = 10$), while for larger SDN cluster sizes ($k > 50$) the extra gain is small.

These findings (Table II and Fig. 5) demonstrate that ASes which have (on average) longer paths to other ASes, e.g., stub networks or small ISPs at the edge of the Internet, would see a higher benefit from routing centralization than central ASes (e.g., tier-1 ISPs) or well connected ASes such as CDNs [21]. Hence, the *node closeness centrality* [17] can be used as a metric to evaluate (or rank) the improvement in the performance of ASes: the lower the closeness centrality, the higher the benefit from routing centralization.

The above observation sheds light on an interesting trade-off related to which nodes participate to the SDN cluster and which nodes benefit from routing centralization. As shown in Section III-B, nodes with high betweenness centrality improve more the performance if they participate in the SDN cluster (see, e.g., Fig. 2). However, their own gain is smaller since they are central nodes in the network (betweenness and closeness centrality are positively correlated measures). As a

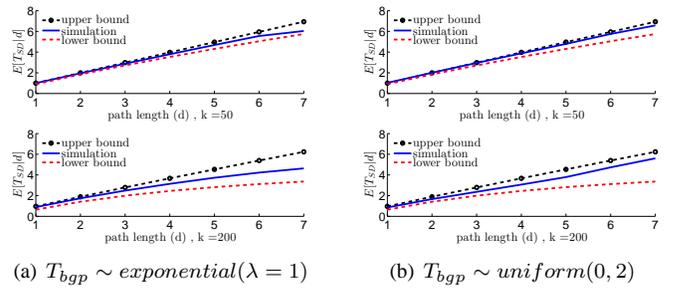


Fig. 4: Data-plane connectivity time $E[T_{SD}|d]$ (y-axis), vs. size of network cluster k (x-axis). Simulation scenarios: Poisson graph network topology of size $N = 1000$ and $p = 0.005$, with (a) $T_{bgp} \sim exponential(\lambda = 1)$ and (b) $T_{bgp} \sim uniform(0, 2)$.

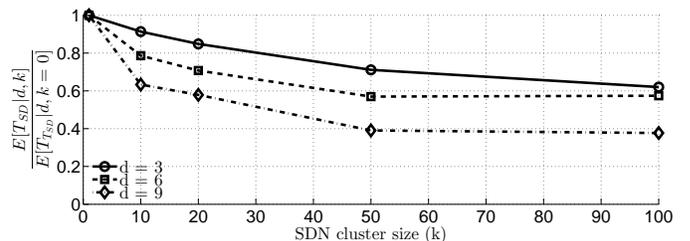


Fig. 5: Data-plane connectivity time normalized over the no SDN scenario $\frac{E[T_{SD}|d, k]}{E[T_{SD}|d, k=0]}$ (y-axis) vs. size of network cluster k (x-axis). Curves correspond to the averages for three different path lengths (i) $d = 3$, (ii) $d = 6$, and (iii) $d = 9$, in simulation scenarios over the Internet AS-graph, with $T_{bgp} \sim exponential(\lambda = 1)$, and nodes in the SDN cluster selected with decreasing *betweenness centrality*.

result, incentives -other than performance- might be needed for attracting central ASes to cooperate for routing centralization. For instance, tier-1 ISPs could deploy inter-domain centralization in order to offer new services (related to the improved BGP convergence performance) to their customers.

IV. CONTROL-PLANE CONVERGENCE

In this section we derive results for the control plane convergence time, i.e., the time needed after a routing change till *every* AS in the network has updated and established the final (i.e., shortest, conforming to routing policies) paths.

The control-plane convergence time is equal to the maximum of the T_{SD} times over all the SD-paths. Due to the involved order statistics, proceeding similarly to Section III, would lead to complex computations and loose bounds. Hence, in this section, we proceed to an approximate analysis that allows us to provide useful insights for the effects of routing centralization on the BGP convergence time.

Specifically, we first narrow the Assumption 1, by assuming that the renewal process for the BGP update times T_{bgp} is a Poisson process; this allows to study the problem using a Markovian framework. Our experiments and measurements in the real Internet (Appendix A), support the selection of the Poisson assumption for the times T_{bgp} .

TABLE II: Data-plane connectivity time normalized over the no SDN scenario, $\frac{E[T_{SD}|k]}{E[T_{SD}|k=0]}$.

Upper bound / Simulation / Lower bound	$k = 20$	$k = 50$	$k = 100$	$k = 200$
$d = 2$	99.9% / 99.2% / 97.0%	99.6% / 97.7% / 92.5%	98.6% / 92.9% / 85.1%	94.4% / 85.1% / 70.4%
$d = 5$	99.9% / 97.8% / 94.2%	99.3% / 93.9% / 86.2%	97.4% / 86.4% / 74.5%	90.1% / 75.6% / 56.4%

Assumption 2 (BGP updates - Poisson process). *The times T_{bgp} are iid random variables, drawn from an exponential distribution with rate $\lambda = \frac{1}{\mu_{bgp}}$, and mean value $E[T_{bgp}] = \mu_{bgp}$.*

Under Assumption 2, we can build a *transient* Markov Chain to model the propagation of BGP updates, where each state denotes the set of nodes that have updated the paths in their RIBs. However, analysing such a Markov chain is still very complex, since the state space contains $2^N - 1$ states, and the transition rates depend on the topology of the network, which cannot be known exactly in most practical cases.

To this end, we first consider the case of a full-mesh network (a common approach in related literature [12], [3], [22]), which can be described by a much simpler Markov chain, and compute the control-plane convergence time as a function of the network size N , SDN cluster size k , and rate λ (Section IV-A). Then, we generalize the results, and derive approximations for sparse topologies, which are of higher practical interest (Section IV-B). Simulation results show that the insights stemming from our analysis are valid also for the (much more complex) Internet AS-graph (Section IV-C).

A. Analysis: Full-Mesh Topology

In a full-mesh network, every pair of nodes has a direct connection, and, thus, the shortest path (i.e., BGP path) to each node is the direct path of size $d = 1$. Hence, every node receives the BGP update from the source node. Moreover, since all nodes in the SDN cluster are informed the time any of them receives the BGP update ($T_{sdn} \ll T_{bgp}$, or $T_{sdn} \rightarrow 0$), the SDN cluster can be considered as a single node.

As a result, a Markov Chain as this in Fig. 6 can be used to model the propagation of BGP updates. Each time a node (a single AS or the SDN cluster) receives the BGP update, the Markov chain moves to the next state. We start from the moment/state (time $t = 0$ / state 0) just before the routing change takes place. Control-plane convergence is achieved at state C , when all nodes have the updated paths in their RIBs.

To calculate the transition rates λ'_i , we first define the following quantities.

Definition 1 (bgp-eligible nodes & bgp-degree).

– A *bgp-eligible node* is a node the (a) has not received the BGP update, and (b) lies on a BGP (shortest) path where the previous node has the updated route in its RIB.

– The *bgp-degree* at step i , $D(i)$, is the number of nodes that are bgp-eligible nodes.

Under the above definition, the time to move from a step/state i to the next step/state, is the time needed till the first of the bgp-eligible nodes receives the update. Under Assumption 2, it follows that this time is the minimum of $D(i)$ iid random variables exponentially distributed with rate λ .

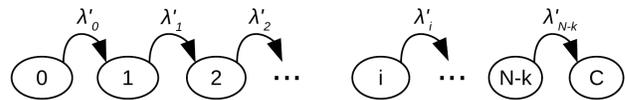


Fig. 6: Markov Chain for the BGP update dissemination process.

Therefore the transition time is also exponentially distributed with rate (i.e., the transition rate)

$$\lambda'_i = \lambda \cdot D(i) \quad (7)$$

Now, in a full-mesh network, bgp-eligible nodes are all the nodes that have not received the BGP update (since all nodes are directly connected to the source node). We denote as $n(i)$ the number of nodes that have received the BGP update at step i . From the above discussion it follows $n(i)$ depends on which step the SDN cluster received the BGP update. Denoting as x the state/step that the first node in the SDN cluster receives the BGP update, we can write

$$n(i|x) = \begin{cases} i & , i \leq x \\ i + k - 1 & , i > x \end{cases} \quad (8)$$

and the bgp-degree is easily shown to be given by Lemma 1.

Lemma 1. *The bgp-degree $D(i|x)$, $i \in [1, N - k]$, $x \in [0, N - k]$, in a full-mesh network topology is given by*

$$D(i|x) = N - n(i|x) \quad (9)$$

Up to this point, we have calculated the transition rates of the Markov chain of Fig. 6 conditionally on x (see, Eq. (7) and Lemma 1). To compute the control-plane convergence time, we need also the probabilities $P_{sdn}(x)$ that the SDN cluster receives the BGP update at step x . In the following lemma, we derive the expression for the probabilities $P_{sdn}(x)$.

Lemma 2. *The probability that the SDN cluster receives the update at step x is given by*

$$P_{sdn}(x) = \frac{k}{N - x} \cdot \prod_{j=0}^{x-1} \left(1 - \frac{k}{N - j}\right) \quad (10)$$

Proof. The proof is given in Appendix C. \square

Now, using Lemmas 1 and 2, we proceed and derive the following result for the distribution of the control-plane convergence time T_c . Specifically, Lemma 3 gives a closed form expression for the moment generating function (MGF)³ of the time T_c .

³ We remind that the MGF of a random variable X is defined as $M_X(\theta) = E[e^{\theta \cdot X}]$, $\theta \in \mathbb{R}$, and completely characterizes a random variable (equivalently to its distribution), since all the moments of X can be calculated from its MGF.

Lemma 3. *The moment generating function (MGF) $M_{T_c}(\theta)$ of the BGP convergence time T_c is given by*

$$M_{T_c}(\theta) = \sum_{x=0}^{N-k} \prod_{i=1}^{N-k} \left(1 - \frac{\theta}{\lambda \cdot D(i|x)}\right)^{-1} \cdot P_{sdn}(x) \quad (11)$$

Proof. The proof is given in Appendix D. \square

Using the above lemma, and applying the property

$$E[X^n] = \frac{d^n M_X(\theta)}{(d\theta)^n} \Big|_{\theta=0} \quad (12)$$

we can calculate the moments of T_c . The following theorem gives the mean value (first moment) of T_c as a function of $D(i|x)$ (Lemma 1) and $P_{sdn}(x)$ (Lemma 2), or, equivalently, as a function of the parameters N , k , and λ .

Theorem 2. *The expectation of the BGP convergence time T_c is*

$$E[T_c] = \frac{1}{\lambda} \cdot \sum_{x=0}^{N-k} \sum_{i=1}^{N-k} \frac{1}{D(i|x)} \cdot P_{sdn}(x) \quad (13)$$

The methodology in the proof of Lemma 3 can be applied to derive useful expressions for other quantities that are of practical interest, and allow us to obtain a better understanding of the effects of routing centralization on control-plane convergence. For example, the following corollary quantifies the speed of the control-plane convergence process.

Corollary 1. *The expectation of the ℓ -Partial BGP Convergence Time, T_ℓ , i.e., the time needed till ℓ ($\ell \leq N$) nodes have the final BGP updates, is given by*

$$E[T_\ell] = \frac{1}{\lambda} \cdot \sum_{x=0}^{N-k} \sum_{i=1}^{M(\ell,x)} \frac{1}{D(i|x)} \cdot P_{sdn}(x) \quad (14)$$

where

$$M(\ell, x) = \begin{cases} \ell - 1 & , \quad 0 < \ell \leq x + 1 \\ x & , \quad x + 1 < \ell \leq x + k \\ \ell - k & , \quad x + k < \ell \leq N \end{cases} \quad (15)$$

B. Analysis: Sparse Topologies

As mentioned earlier, computing the control-plane convergence for an arbitrary topology is very complex. For instance, applying the methodology of Section IV-A, the set of bgp-eligible nodes at a step i depends on the exact paths \mathcal{P} that the BGP updates have been propagated. Hence, we need to consider all $S \in \mathcal{P}$ (with $|\mathcal{P}| \sim O(2^N)$), and we need to keep track of all $D(i|x, S \in \mathcal{P})$ and $P_{sdn}(x|S \in \mathcal{P})$. However, approximating sparse topologies with a Poisson (or, Erdos-Renyi) random graph $G(N, p)$, we derive expressions for the BGP convergence time in the following result. As we show in the validation Section IV-C, our approximations describe well effects of routing centralization also in more generic/realistic topologies, like power-law graphs or the Internet AS-graph.

Result 1. *Lemma 3, Theorem 2, and Corollary 1, with $E[D(i|x)]$ (instead of $D(i|x)$), approximate the control-plane convergence time in a Poisson graph network topology; where*

$E[D(i|x)]$ is the expectation of the bgp-degree $D(i|x)$, $i \in [1, N - k]$, $x \in [0, N - k]$, in a Poisson graph

$$E[D(i|x)] = (N - n(i|x)) \cdot \left(1 - (1 - p)^{n(i|x)}\right) \quad (16)$$

Proof. The proof is given in Appendix E. \square

C. Simulation Results and Implications

We evaluate the accuracy of our theoretical results in various simulation scenarios, including also scenarios where the assumptions for (i) exponential $f_{bgp}(t)$, and (ii) full-mesh or Poisson graph networks, do not hold.

In scenarios of full-mesh networks, where the times T_{bgp} are exponentially distributed, our theoretical expressions of Section IV-A predict the simulation results for the expected convergence time $E[T_c]$ with very high accuracy.

For the validation of the theoretical expressions in sparse networks (Section IV-B), we simulate various sparse topologies, like Poisson, Barabasi-Albert (power law), and Newman-Watts-Strogatz (small world) graphs. Although the theoretical results are derived under the Poisson graph assumption, our simulations show that they can predict the performance with similar accuracy in the all the topologies we tested.

In Fig. 7 we present a representative subset of our results that demonstrate how the routing centralization can decrease the BGP convergence time. We plot the partial convergence time, normalized over the scenario without centralization, i.e., $\frac{E[T_\ell|k]}{E[T_\ell|k=0]}$. We consider three cases, $\ell = 100$ (or, $0.1 \cdot N$) in Fig. 7(a), $\ell = 500$ (or, $0.5 \cdot N$) in Fig. 7(b), and $\ell = N = 1000$ that corresponds to the control-plane convergence in Fig. 7(c).

A first observation is that our results can capture well the relative changes⁴ in the (partial) convergence time, not only for scenarios with exponential $f_{bgp}(t)$ (as we assume in our analysis), but also for scenarios with uniform $f_{bgp}(t)$.

In Fig. 7(c), we can see that the control-plane convergence time does not significantly improve as the SDN cluster size k increases. For instance, even for $k = 500$ (i.e., 50% of the nodes belong to the SDN cluster), the decrease in the convergence time is less than 30%. This comes to verify the results of [12], which showed that significant gains can be achieved only for high values ($> 50\%$) of SDN penetration.

However, when it comes to the partial control-plane convergence (Figs. 7(a) and 7(b)), the effects of routing centralization are higher. The time needed till 10% of the nodes ($\ell = 100$ - Fig. 7(a)) to receive the updated routing information, decreases quickly; e.g., to 0.5 of its no-SDN ($k = 0$) value, only with $k = 100$ nodes (10%) participating in the SDN cluster.

This reveals an important aspect, relating to the effects of routing centralization, which has not been shown in previous works (e.g., [12]): although the control-plane convergence can significantly improve only if a high percentage ($> 50\%$) of nodes cooperate, we can have very large gains in the *partial convergence* even with small sizes of SDN clusters.

⁴The accuracy of the theoretical results (approximations), when we consider the *actual* -not normalized- values, is lower.

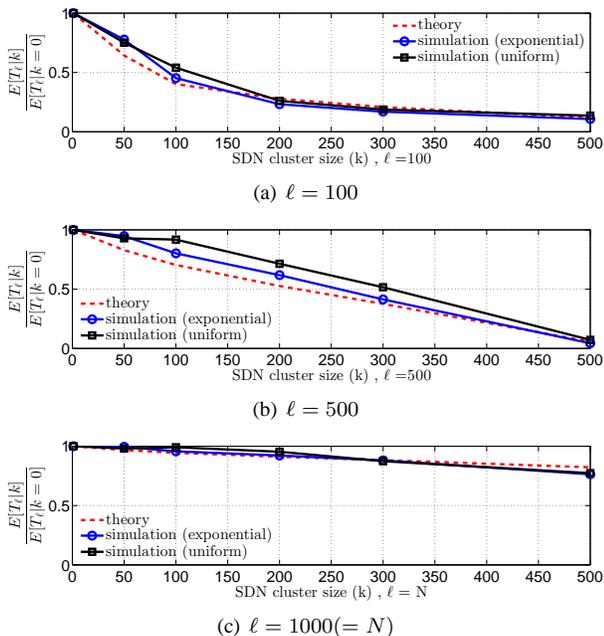


Fig. 7: Partial convergence time, normalized over the no SDN scenario, $\frac{E[T_\ell|k]}{E[T_\ell|k=0]}$ (y-axis), vs. size of SDN cluster k (x-axis). Simulation scenarios: Barabasi-Albert topology with $N = 1000$ and average node degree 10; $T_{bgp} \sim exponential(\lambda = 1)$ (black line - squares) and $T_{bgp} \sim uniform(0, 2)$ (blue line - circles).

In Fig. 8 we present simulation results on the Internet AS-graph⁵, where the top betweenness centrality nodes form the SDN cluster. Despite the fact that the simulated scenario deviates from our assumptions, our main theoretical findings are still valid: centralization can significantly accelerate the connectivity time with a large percentage of ASes (i.e., ℓ -partial convergence, see, e.g., curves for $\ell = 0.1 \cdot N$ and $\ell = 0.5 \cdot N$), while the time needed till every AS has received the updated routes (i.e., total convergence $E[T_c]$) improves more slowly with the SDN cluster size k . Moreover, we can see that the efficiency of inter-domain centralization is quite impressive; with only $k = 50$ central nodes in the SDN cluster, the time needed to establish updated paths with half of the Internet nodes ($\ell = 0.5 \cdot N$) is 50% less than in the case without centralization.

V. RELATED WORK

Inter-domain SDN is a new research area that attracts increasing attention [6], [7], [11], [8], [9], [10], [12]. In [6] authors propose and implement SDX, a software-defined component for IXPs, which increases the capabilities on routing control. Another IXP-based system that enables novel services for establishing QoS route paths is described in [7]. In [11] a solution for incremental deployment of inter-domain SDN, which is seamless to traditional IP networks, is proposed, and [8] contributes in this direction by proposing an SDN-

⁵For scalability issues, we did not consider here stub ASes and ASes with less than 3 neighbors, resulting in a reduced Internet graph with $N = 11527$.

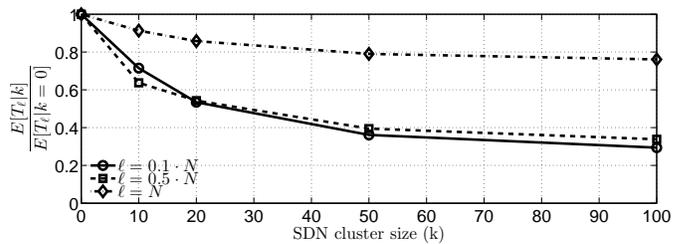


Fig. 8: Partial convergence time, normalized over the no SDN scenario, $\frac{E[T_\ell|k]}{E[T_\ell|k=0]}$ (y-axis), vs. size of SDN cluster k (x-axis). Simulation scenarios on the Internet AS-graph. Nodes in the SDN cluster are selected with decreasing *betweenness centrality*.

based methodology for decoupling BGP policy control from routing. [9] proposes an SDN-based architecture to enhance inter-domain routing, and [10] proposes an component to enable inter-domain SDN. Finally, authors in [12] build a realistic emulator, and use it to investigate the effects of routing centralization on BGP convergence time.

The slow convergence of BGP has been extensively studied through measurements in [3], [4], [5]. It has been shown that BGP can take several minutes to converge after a routing change, and this can cause severe packet losses [3] and performance degradation [4].

Finally, analytic approaches for the BGP convergence can be found in [22], [3], [23]. In [22], a probabilistic model and automata theory is used to study the BGP convergence (probability of convergence, and convergence time). [3] studies analytically the BGP convergence with respect to the number of exchanged messages, while [23] performs a worst-case analysis of BGP convergence

VI. CONCLUSION

In this paper, we analytically studied the effects of inter-domain SDN on the time needed for establishing connectivity and convergence after a routing change. We proposed a probabilistic model, and derived results for the expected data-plane connectivity time (lower/upper bounds) and control-plane convergence time (exact predictions and approximations).

Our results can be used to quickly evaluate the effects of different network parameters, like network size, topology, path lengths, or number of SDN nodes, on the routing performance. Hence, they can complement previous system-oriented studies and facilitate future research. Moreover, our methodology and results can be a useful tool for studying important problems relating to routing changes in the Internet. Finally, they can be applied in practical design problems, like selecting the nodes to participate in the SDN cluster based on performance criteria (i.e., which node can have the highest impact), or for network economics purposes (e.g., detecting the potential incentives for an AS to participate in inter-domain routing centralization).

REFERENCES

- [1] S. Kent, C. Lynn, and K. Seo, "Secure border gateway protocol (s-bgp)," *IEEE Journal on Selected Areas in Communications*, vol. 18, no. 4, pp. 582–592, 2000.

- [2] L. Subramanian, V. Roth, I. Stoica, S. Shenker, and R. Katz, "Listen and whisper: Security mechanisms for bgp," in *Proc. NSDI*, 2004.
- [3] C. Labovitz, A. Ahuja, A. Bose, and F. Jahanian, "Delayed Internet routing convergence," *ACM SIGCOMM Computer Communication Review*, vol. 30, no. 4, pp. 175–187, 2000.
- [4] N. Kushman, S. Kandula, and D. Katabi, "Can you hear me now?: it must be BGP," *ACM SIGCOMM Computer Communication Review*, vol. 37, no. 2, pp. 75–84, 2007.
- [5] R. Oliveira, B. Zhang, D. Pei, and L. Zhang, "Quantifying path exploration in the Internet," *Networking, IEEE/ACM Transactions on*, vol. 17, no. 2, pp. 445–458, 2009.
- [6] A. Gupta, L. Vanbever, M. Shahbaz, S. Donovan, B. Schlinker, N. Feamster, J. Rexford, S. Shenker, R. Clark, and E. Katz-Bassett, "SDX: A software defined internet exchange," *ACM SIGCOMM Comput. Commun. Rev.*, vol. 44, no. 4, pp. 551–562, 2014.
- [7] Kotronis, V. and Kloeti, R. and Rost, M. and Georgopoulos, P. and Ager, B. and Schmidt, S. and Dimitropoulos, X., "Stitching inter-domain paths over IXPs," in *ACM SOSP*, 2016.
- [8] P. W. Thai and J. C. de Oliveira, "Decoupling bgp policy from routing with programmable reactive policy control," in *Proc. ACM CoNEXT Student*, 2012.
- [9] C. Rothenberg, M. Nascimento, M. Salvador, C. Corrêa, S. C. de Lucena, and R. Raszuk, "Revisiting routing control platforms with the eyes and muscles of software-defined networking," in *Proc. ACM HotSDN*, 2012.
- [10] R. Benesby, E. Mota, P. Fonseca, and A. Passito, "Innovating on interdomain routing with an inter-sdn component," in *Proc. IEEE AINA*, 2014.
- [11] P. Lin, J. Hart, U. Krishnaswamy, T. Murakami, M. Kobayashi, A. Al-Shabibi, K.-C. Wang, and J. Bi, "Seamless interworking of sdn and ip," in *Proc. ACM SIGCOMM 2013 (demo)*, 2013.
- [12] V. Kotronis, A. Gämperli, and X. Dimitropoulos, "Routing centralization across domains via sdn: A model and emulation framework for BGP evolution," *Computer Networks*, pp. –, 2015.
- [13] G. Gibb, H. Zeng, and N. McKeown, "Outsourcing network functionality," in *Proc. ACM HotSDN*, 2012.
- [14] NANOG mailing list archives, "A survey on BGP convergence." <http://seclists.org/nanog/2017/Jan/55>, Jan. 2017.
- [15] S. Vissicchio, O. Tilmans, L. Vanbever, and J. Rexford, "Central control over distributed routing," *ACM SIGCOMM Computer Communication Review*, vol. 45, no. 4, pp. 43–56, 2015.
- [16] M. Alan Chang, T. Holterbach, M. Happe, and L. Vanbever, "Supercharge me: Boost router convergence with sdn," *ACM SIGCOMM Comput. Commun. Rev.*, vol. 45, no. 4, pp. 341–342, 2015.
- [17] M. Newman, *Networks: An Introduction*. Oxford University Press, Inc., 2010.
- [18] "The CAIDA AS relationships dataset." <http://data.caida.org/datasets/as-relationships/>, 01 Nov. 2016.
- [19] L. Gao and J. Rexford, "Stable internet routing without global coordination," *IEEE/ACM Transactions on Networking (TON)*, vol. 9, no. 6, pp. 681–692, 2001.
- [20] "AS-rank, CAIDA." <http://as-rank.caida.org>.
- [21] Y.-C. Chiu, B. Schlinker, A. B. Radhakrishnan, E. Katz-Bassett, and R. Govindan, "Are we one hop away from a better internet?," in *Proc. ACM IMC*, 2015.
- [22] R. Viswanathan, K. K. Sabnani, R. J. Holt, and A. N. Netravali, "Expected convergence properties of bgp," *Computer Networks*, vol. 55, no. 8, pp. 1957–1981, 2011.
- [23] R. Sami, M. Schapira, and A. Zohar, "Searching for stability in inter-domain routing," in *Proc. IEEE INFOCOM*, 2009.
- [24] B. Schlinker, K. Z. I., Cunha, N. Feamster, and E. Katz-Bassett, "Peering: An as for us," in *Proc. ACM HotNets 2014*, <https://peering.usc.edu>.
- [25] G. Chavarias, P. Gigis, P. Sermpezis, and X. Dimitropoulos, "ARTEMIS: Real-time detection and automatic mitigation for bgp prefix hijacking," in *Proc. ACM SIGCOMM Demo*, 2016.
- [26] "BGPmon." <http://www.bgpmon.io>.
- [27] "RIPE RIS." <http://ris.ripe.net/>.
- [28] <http://www.caida.org/tools/utilities/looking-glass-api/>.
- [29] G. W. Oehlert, "A note on the delta method," *The American Statistician*, vol. 46, no. 1, pp. 27–29, 1992.
- [30] P. Gill, M. Schapira, and S. Goldberg, "Let the market drive deployment: A strategy for transitioning to bgp security," in *ACM SIGCOMM Computer Communication Review*, vol. 41, pp. 14–25, 2011.
- [31] S. Goldberg, M. Schapira, P. Hummon, and J. Rexford, "How secure are secure interdomain routing protocols?," vol. 70, pp. 260–287, Elsevier, 2014.
- [32] A. Cohen, Y. Gilad, A. Herzberg, and M. Schapira, "Jumpstarting bgp security with path-end validation," in *Proc. ACM SIGCOMM 2016*, 2016.
- [33] Y. Gilad, A. Cohen, A. Herzberg, M. Schapira, and H. Shulman, "Are we there yet? on RPKI's deployment and security." NDSS 2017, to appear, <http://eprint.iacr.org/2016/1010.pdf>, 2016.
- [34] M. Luckie, B. Huffaker, A. Dhamdhere, V. Giotsas, *et al.*, "As relationships, customer cones, and validation," in *Proc. ACM IMC*, 2013.
- [35] V. Giotsas, S. Zhou, M. Luckie, and k. claffy, "Inferring multilateral peering," in *Proc. ACM CoNEXT*, 2013.

APPENDIX A

DISTRIBUTION OF BGP UPDATE TIMES T_{bgp}

To investigate if and how well our modeling assumptions can describe the BGP update times in the Internet, we compare them against real measurement data.

We conducted experiments in the Internet using the PEERING testbed [24], which owns IP prefixes and ASNs, peers with networks in different locations around the world, and allows users to make real BGP announcements. In our experiments/measurements, we follow a similar methodology as in [25]: we (i) announce a /24 prefix from a site of the PEERING testbed, and (ii) use publicly available control-plane monitoring services (route collectors and looking glass servers) [26], [27], [28] to measure the time needed till different ASes receive our announcements.

We collected BGP updates, as seen from the monitors, from $M = 40$ ASes. We repeated the experiments 14 times; each time making a BGP announcement either from the PEERING site at an IXP at Amsterdam (NL), or at an ISP at Los Angeles (US). From each received BGP update i , we consider (a) $T_{SD}(i)$, the time needed till the BGP update i received by the monitor (i.e., timestamp of the BGP update i minus the timestamp of our BGP announcement), and (b) $d(i)$, the length of the AS-path included in the BGP update i .

We group the times $T_{SD}(i)$ by the respective path lengths $d(i)$, and plot the distribution (CCDF) of the measured times T_{SD} in Fig. 9 for two example cases with $d = 2$ and $d = 5$.

Then, we fit the real data with a distribution $f_{bgp}(t)$ (cf. Section II), where we select $f_{bgp}(t) \sim exponential(\lambda)$ in order to test the validity of (the stronger) Assumption 2. We estimate the average BGP update time from the measured data as $\hat{E}[T_{bgp}] = \frac{\sum_i T_{SD}(i)}{\sum_i d(i)}$ and set the rate $\lambda = \frac{1}{\hat{E}[T_{bgp}]}$.

We generate from $f_{bgp}(t)$ a large number of times T_{SD} for paths of length $d = 2$ and $d = 5$, calculate their CCDFs, and compare them against the real data in Fig. 9. As we can observe, there is a good match between the generated and real data. This indicates that Assumption 2 is a realistic and reasonable assumption, and, thus, emphasizes the practicality of our theoretical and simulation findings in real settings

APPENDIX B

PROOF OF THEOREM 1

Proof. Let us assume a SD-path of length d and denote the ASes/nodes in the path as n_0, n_1, \dots, n_d , where $n_0 \equiv S$ and $n_d \equiv D$. The total number of ASes on the SD-path is $d + 1$

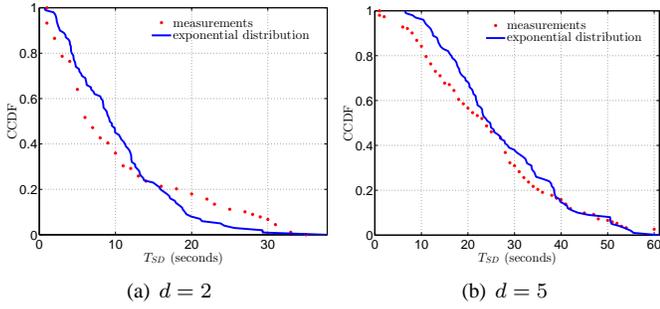


Fig. 9: CCDF of the times T_{SD} for SD-paths of length (a) $d = 2$ and (b) $d = 5$. Comparison of times T_{SD} from *measurements* in the Internet (where we found $E[T_{bgp}] = 6.27$), and times T_{SD} generated from our model with $f_{bgp}(t) \sim \text{exponential distribution}(\lambda = \frac{1}{E[T_{bgp}]})$.

(including nodes S and D). Let us denote as k' , $0 \leq k' \leq d+1$ the number of these nodes that belong also to the SDN cluster.

If none of the nodes comprising the SD-path belong to the SDN cluster (i.e., $k' = 0$), the BGP updates propagate from $n_0 \equiv S$ to n_1 , then from n_1 to n_2 , etc., till they reach the destination node $n_d \equiv D$. Therefore, the time T_{SD} is equal to

$$T_{SD} = T_{n_0, n_1} + T_{n_1, n_2} + \dots + T_{n_{d-1}, n_d} = \sum_{i=0}^{d-1} T_{n_i, n_{i+1}} \quad (17)$$

and since the times $T_{n_i, n_{i+1}}$ are iid random variables, i.e., $T_{n_i, n_{i+1}} \sim f_{bgp}(t)$, (Assumption 1), the expectation of T_{SD} is

$$E[T_{SD}|d, k' = 0] = \sum_{i=0}^{d-1} E[T_{n_i, n_{i+1}}] = d \cdot E[T_{bgp}] \quad (18)$$

Now assume that the node n_j , $j = 1, \dots, d$, is the only node on the SD-path that belongs in the SDN cluster (i.e., $k' = 1$). Let $T_1 = \sum_{i=1}^{j-1} T_{n_i, n_{i+1}}$ be the time needed for the update to propagate from $n_0 \equiv S$ to n_j , and $T_2 = \sum_{i=j}^{d-1} T_{n_i, n_{i+1}}$ the time needed for the update to propagate from n_j to the destination $n_d \equiv D$.

The node n_j is first informed about the BGP update at time $T' \leq T_1$: either from the previous node in the path ($T' = T_1$), or at an earlier time ($T' < T_1$) from the SDN cluster, if the SDN cluster has received (through another path) the BGP update earlier.

Therefore, the total time needed for all the nodes in the SD-path to receive the BGP update can be expressed as

$$T_{SD} = \max\{T_1, T' + T_2\} \quad (19)$$

Lower Bound:

To derive the lower bound of the expectation of T_{SD} , we take the expectations on Eq. (19) and proceed as follows.

$$E[T_{SD}|d, k' = 1] = E[\max\{T_1, T' + T_2\}] \quad (20)$$

$$\geq E[\max\{T_1, T_2\}] \quad (21)$$

$$\geq \max\{E[T_1], E[T_2]\} \quad (22)$$

$$= \max\{E[T_{bgp}] \cdot d_1, E[T_{bgp}] \cdot d_2\} \quad (23)$$

$$= E[T_{bgp}] \cdot \max\{d_1, d_2\} \quad (24)$$

$$\geq E[T_{bgp}] \cdot \min_{d_1, d_2} \{\max\{d_1, d_2\}\} \quad (25)$$

which gives

$$E[T_{SD}|d, k' = 1] \geq \begin{cases} 0 & , d = 1 \\ E[T_{bgp}] \cdot \frac{d}{2} & , d > 1 \end{cases} \quad (26)$$

where

- Eq. (21) follows since $T' \geq 0$ ($T' = 0$ denotes the event that the SDN cluster receives the BGP update immediately after the routing change takes place).
- The inequality of Eq. (22) follows since the times T_1 and T_2 are independent random variables, and thus it holds

$$\begin{aligned} P\{\max\{T_1, T_2\} \leq t\} &= P\{T_1 \leq t\} \cdot P\{T_2 \leq t\} \Rightarrow \\ P\{\max\{T_1, T_2\} \leq t\} &\leq P\{T_i \leq t\} \Rightarrow \\ P\{\max\{T_1, T_2\} > t\} &\geq P\{T_i > t\}, \quad \forall i = \{1, 2\} \end{aligned}$$

and for a positive r.v. X it also holds that $E[X] = \int_0^\infty P\{X > x\} dx$, and thus taking the integral in the above inequality it follows

$$\begin{aligned} \int_0^\infty P\{\max\{T_1, T_2\} > t\} dt &\geq \int_0^\infty P\{T_i > t\} dt \Rightarrow \\ E[\max\{T_1, T_2\}] &\geq E[T_i], \quad \forall i = \{1, 2\} \end{aligned}$$

or, equivalently, $E[\max\{T_1, T_2\}] \geq \max\{E[T_i]\}$.

- The expectations $E[T_i]$, $i = \{1, 2\}$ are substituted in Eq. (23) with $E[T_{bgp}] \cdot d_i$ since T_i is the sum of d_i iid r.v. with expected value $E[T_{bgp}]$.
- In Eq. (25) we consider all the possible combinations of d_1 and d_2 (under the condition $d_1 + d_2 = d$), whose max value is minimized when $d_1 = d_2 = \frac{d}{2}$ (Eq. (26)).

Now, if there are k' nodes in the SD-path that belong to the SDN cluster, proceeding similarly to the above case $k' = 1$ leads to the following generic inequality

$$E[T_{SD}|d, k'] \geq \begin{cases} 0 & , d \leq k' \\ E[T_{bgp}] \cdot \frac{d}{k'+1} & , d > k' \end{cases} \quad (27)$$

which gives the lower bound of Theorem 1.

Upper Bound:

For $k' = 0$, the expectation of T_{SD} is given by Eq. (18). For $k' = 1$, since $T' \leq T_1$, we can use Eq. (19) and write

$$E[T_{SD}|d, k' = 1] = E[\max\{T_1, T' + T_2\}] \quad (28)$$

$$\leq E[\max\{T_1, T_1 + T_2\}] \quad (29)$$

$$= d \cdot E[T_{bgp}] \quad (30)$$

where the last equality follows from Eq. (18).

In the case of $k' > 1$, it is probable that, after the SDN cluster is informed about the routing change, the BGP update propagates simultaneously on more than one sections on the SD-path. For example, in Fig. 1, after the SDN cluster is informed (n_i and n_j receive the update at the same time), the BGP update will propagate *simultaneously* in the sub-paths $n_i \rightarrow \dots \rightarrow n_{j-1}$ and $n_j \rightarrow \dots \rightarrow n_d$. This, accelerates the propagation process, and, thus, decreases the time T_{SD} .

It is easy to see, that the smaller decrease (on average) on T_{SD} , will take place when the k' nodes that belong to the SDN cluster are located consecutively on the SD-path. Without loss of generality, let assume that the first k' nodes $n_0, \dots, n_{k'-1}$ are the nodes that belong to the SDN cluster, and denote the time T_{SD} for this (worst) case as T_{SD}^{max} . Then, the time T_{SD}^{max} is given by

$$T_{SD}^{max} = \sum_{i=0}^{k'-2} T_{n_i, n_{i+1}} + \sum_{i=k'-1}^{d-1} T_{n_i, n_{i+1}} = \sum_{i=k'-1}^{d-1} T_{n_i, n_{i+1}} \quad (31)$$

since $\sum_{i=0}^{k'-2} T_{n_i, n_{i+1}} = T_{sdn} \equiv 0$. The expectation of T_{SD}^{max} is derived similarly to Eq. (17) and Eq. (18), i.e.,

$$E[T_{SD}^{max}] = E \left[\sum_{i=k'-1}^{d-1} T_{n_i, n_{i+1}} \right] = (d - (k' - 1)) \cdot E[T_{bgp}] \quad (32)$$

Combining Eq. (18), Eq. (30), and Eq. (32), gives the upper bound of Theorem 1. \square

APPENDIX C PROOF OF LEMMA 2

Proof. Considering all the cases for which node initiates the routing change, the probability that the source node belongs to the SDN cluster (and thus $x = 0$) is

$$P_{sdn}(0) \equiv P_{sdn}(x = 0) = \frac{k}{N} \quad (33)$$

If the source node does not belong to the SDN cluster, then at step 1 there are $N - 1$ bgp-eligible nodes, of which k belong to the SDN cluster. This gives

$$P_{sdn}(1 | x > 0) = \frac{k}{N - 1} \quad (34)$$

and, consequently,

$$P_{sdn}(1) = P_{sdn}(1 | x > 0) \cdot P_{sdn}(x > 0) = \frac{k}{N-1} \cdot \left(1 - \frac{k}{N}\right)$$

Proceeding recursively, we derive Eq. (10) that gives the probability $P_{sdn}(x)$. \square

APPENDIX D PROOF OF THEOREM 3

Proof. The convergence time is T_c is calculated by the sum of the transition times of the Markov Chain of Fig. 6, i.e.,

$$T_c = T_{1,2} + T_{2,3} + \dots + T_{N-k,C} = \sum_{i=1}^{N-k} T_{i,i+1} \quad (35)$$

where we denote $T_{N-k, N-k+1} \equiv T_{N-k, C}$. Hence, the MGF of T_c is expressed as

$$M_{T_c}(\theta) = E \left[e^{\theta \cdot \sum_{i=1}^{N-k} T_{i,i+1}} \right] \quad (36)$$

$$= E \left[\prod_{i=1}^{N-k} e^{\theta \cdot T_{i,i+1}} \right] \quad (37)$$

$$= \sum_{x=0}^{N-k} E \left[\prod_{i=1}^{N-k} e^{\theta \cdot T_{i,i+1}} \middle| x \right] \cdot P_{sdn}(x) \quad (38)$$

$$= \sum_{x=0}^{N-k} \prod_{i=1}^{N-k} E \left[e^{\theta \cdot T_{i,i+1}} \middle| x \right] \cdot P_{sdn}(x) \quad (39)$$

$$= \sum_{x=0}^{N-k} \prod_{i=1}^{N-k} \left(1 - \frac{\theta}{\lambda \cdot D(i|x)} \right)^{-1} \cdot P_{sdn}(x) \quad (40)$$

where

- In Eq. (38) we consider the conditional expectation, given that the SDN cluster receives the update at step x .
- Eq. (39) follows from the fact that the times $T_{i,i+1}$ are independent under a given x ; due to Assumption 2, they depend only on the number of infected nodes, which is determined by the step i and the value of x .
- We derive Eq. (40), since $T_{i,i+1}$ is an exponential random variable with rate $\lambda'_{i,i+1} = \lambda \cdot D(i|x)$, and the MGF of an exponential r.v. with rate μ is given by $(1 - \theta/\mu)^{-1}$. \square

APPENDIX E PROOF OF RESULT 1

Proof. To derive the MGF of T_c we apply the methodology in the proof of Lemma 3; here, we highlight only the key points and differences from the full-mesh case.

$$M_{T_c}(\theta) = E \left[\prod_{i=1}^{N-k} e^{\theta \cdot T_{i,i+1}} \right] \quad (41)$$

$$= \sum_{S \in \mathcal{P}} \sum_{x=0}^{N-k} E \left[\prod_{i=1}^{N-k} e^{\theta \cdot T_{i,i+1}} \middle| x, S \right] \cdot P\{x, S\} \quad (42)$$

$$= \sum_{x=0}^{N-k} \sum_{S \in \mathcal{P}} E \left[\prod_{i=1}^{N-k} e^{\theta \cdot T_{i,i+1}} \middle| x, S \right] \cdot P\{S\} \cdot P_{sdn}(x) \quad (43)$$

$$= \sum_{x=0}^{N-k} \sum_{S \in \mathcal{P}} \prod_{i=1}^{N-k} \left(1 - \frac{\theta}{\lambda \cdot D(i|x, S)} \right)^{-1} \cdot P\{S\} \cdot P_{sdn}(x) \quad (44)$$

$$= \sum_{x=0}^{N-k} E \left[\prod_{i=1}^{N-k} \left(1 - \frac{\theta}{\lambda \cdot D(i|x, S)} \right)^{-1} \right] \cdot P_{sdn}(x) \quad (45)$$

$$\approx \sum_{x=0}^{N-k} \prod_{i=1}^{N-k} \left(1 - \frac{\theta}{\lambda \cdot E_{\mathcal{P}}[D(i|x)]} \right)^{-1} \cdot P_{sdn}(x) \quad (46)$$

where

- After expressing the MGF in Eq. (41), we apply the conditional expectation property to write Eq. (42), where x is the step that the SDN cluster received the BGP update,

S is the set of nodes that have the BGP update, and with $P\{x, S\}$ we denote the respective joint probability.

- Since we assume the SDN cluster to be formed independently of the topology, it holds (for any topology) that the variables x and S are independent. Hence, $P\{x, S\} = P\{x\} \cdot P\{S\}$, where $P\{x\} \equiv P_{sdn}(x)$ and its value is given by Theorem 2. Also, we can reorder the summations over x and S , which gives Eq. (43).
- Eq. (44) follows by making similar arguments as in the proof of Lemma 3, and can be written as Eq. (45), where the expectation is taken over the set $S \in \mathcal{P}$.
- Since the expectation in Eq. (45) is difficult to compute (see above discussion), we approximate it with the *Delta method* [29]. In the Delta method the expectation of a function (i.e., the product in Eq. (45)) of a random variable (i.e., $D(i|x, S)$) is approximated by the function of the expectation of the random variable (i.e., $E_{\mathcal{P}}[D(i|x)]$).

From Eq. (46), it can be seen that the approximation of $M_{T_c}(\theta)$ is given by an expression as in Lemma 3, where $D(i|x)$ is replaced by $E_{\mathcal{P}}[D(i|x)]$. Moreover, it is easy to see that all the consequent results for the full-mesh network can be similarly modified for the Poisson graph case.

Now, we need only to calculate the expected bgp-degree $E_{\mathcal{P}}[D(i|x)]$: Let assume that we are at step i , and $n(i)$ nodes (see Eq. (8)) have received the BGP updates; we denote the set of these nodes as S_i . A node $m \notin S_i$ is connected with a node $j \in S_i$ with probability $P(m, j) = p$ (by the definition of a Poisson graph). Hence, the probability that m is a bgp-eligible node (i.e., is connected with *any* of the nodes $j \in S_i$, where $|S_i| = n(i)$), is given by

$$P(m, S_i) = 1 - (1 - p)^{n(i)} \quad (47)$$

Finally, we note that there are $N - n(i)$ nodes without the update, with each of them being a bgp-eligible node with any of the nodes $j \in S_i$ with (equal) probability $P(m, S_i)$. As a result, the total number of bgp-eligible nodes (or, as defined in Def. 1, the *bgp-degree* $D(i)$) is a binomially distributed random variable, whose expectation is given by

$$E[D(i)] = (N - n(i)) \cdot (1 - (1 - p)^{n(i)}) \quad (48)$$

□

APPENDIX F

INTERNET TOPOLOGY AND ROUTING POLICIES

To approximate the routing system of the Internet, we use a methodology similar to many previous works [30], [31], [32], [33]. We first build the Internet topology graph from a large experimentally collected dataset [18], and infer routing policies over existing links based on the Gao-Rexford conditions [19].

A. Building the Internet Topology

We build the Internet topology graph from the AS-relationship dataset of CAIDA [18], which is collected based on the methodology of [34] and enriched with many extra peering (p2p) links [35]. The dataset contains a list of AS

pairs with a peering link, which is annotated based on their relationship as *c2p* (*customer to provider*) or *p2p* (*peer to peer*).

B. Selecting Routing Policies

When an AS learns a new route for a prefix (or, announces a new prefix), it updates its routing table and, if required, sends BGP updates to its AS neighbors. The update and export processes are defined by its routing policies. Similarly to previous works [30], [31], [32], [33], we select the routing policies based on the Gao-Rexford conditions that guarantee BGP convergence and stability [19]:

- C.1 Paths learned from customers are preferred to paths learned from peers or providers. Paths learned from peers are preferred to paths learned from providers.
- C.2 Between paths that are equivalent with respect to **C.1**, shorter paths (in number of AS-hops) are preferred.
- C.3 Between paths that are equivalent with respect to **C.1** and **C.2**, the path learned from the AS neighbor with the highest *local preference* is preferred.
- C.4 Paths learned from customers, are advertised to all AS neighbors. Paths learned from peers or providers, are advertised only to customers.

In practice, the local preferences (see, **C.3**) are selected by an AS based on factors related to its intra-domain topology, business agreements, etc. Since it is not possible to know and emulate the real policies for every AS, we assign randomly the local preferences.