

Towards Stress Testing the Internet Inter-Domain Routing System ‘in Silico’ with Domino

Elham Ehsani Moghadam
ETH Zurich

Fabián E. Bustamante
Northwestern University

Adrian Perrig
ETH Zurich

Walter Willinger
NIKSUN, Inc.

Abstract—In just a few decades, the Internet has evolved from a research prototype to a cyber-physical infrastructure of critical importance for modern society and the global economy. Surprisingly, despite its new role, the survivability of the Internet—its ability to fulfill its mission in the presence of large-scale failures—has received limited attention. We introduce Domino, our initial design and implementation of a testbench tool for stress testing the Internet’s routing system, a key element of the critical Internet infrastructure. The simulation-based testbench consists of a comprehensive and flexible framework that allows for the incorporation of diverse survivability metrics, provides a platform for specifying, evaluating, and comparing different topologies of the underlying Internet infrastructure, and can account for modifications to networking protocols and architectural components. By demonstrating the utility of the proposed testbench with a number of illustrative examples, we make a case for stress testing as a viable approach to evaluating the Internet’s survivability in the face of evolving challenges.

Index Terms—Internet Survivability, Large-Scale Failures, Survivability Testbench, Inter-domain Routing

I. INTRODUCTION

In a few decades, the Internet has transformed from a research prototype into a critical infrastructure that underpins modern society and the economy [1]–[3]. This development has been largely organic, driven by continuous expansion in physical networks, a diverse range of stakeholders, and the steady introduction of new applications and services. Other critical infrastructures—including transportation, banking, and water and wastewater management—now rely on the Internet to function effectively. As our awareness of the consequences of disconnection grows, particularly for extended periods, we also recognize vulnerabilities in the network that may not have been foreseen during its early years.

Although the Internet has generally shown resilience over the years, its designation as critical infrastructure prompts the question: is it able to fulfill this role? Specifically, can it endure catastrophic events like large-scale earthquakes or prolonged power outages, withstand cyber attacks orchestrated by nation-states targeting its critical infrastructure systems, and navigate extreme space weather phenomena such as Coronal Mass Ejection (CME) events [4]?

During crisis situations, the Internet becomes even more essential, acting as a key facilitator for effective crisis management. It serves as a lifeline for efficient communication and coordination among first responders and facilitates the timely dissemination of critical information, from providing

individuals with life-saving updates to weather forecasts and evacuation routes.

Internet resilience involves multiple layers. One key aspect is physical layer diversity, which provides redundant connectivity to help withstand failures. However, this redundancy is only effective if routing protocols can detect and use the available paths. Inter-domain and intra-domain routing protocols are expected to maintain connectivity by identifying and utilizing functioning routes, even when parts of the network are disrupted.

While we are interested in stress testing Internet critical infrastructure in general, our focus in this paper is on inter-domain routing, particularly the Border Gateway Protocol (BGP), because of its crucial role in maintaining global connectivity.

In this paper, we present our testbench tool Domino for evaluating inter-domain networking survivability. Specifically, our work makes the following main contributions:

- We present the capabilities of Domino, a simulation/emulation-based testbench that enables comparisons between different topologies, architectures, configurations, and BGP modifications by stress testing the inter-domain routing system. Beyond including diverse input environments and settings, a key novelty of Domino is an algorithm that estimates disconnection time for each router-prefix pair using only BGP update messages—independent of data-plane protocols and transparent to underlying transport behavior.
- We demonstrate the capabilities of Domino with illustrative examples that employ different router topologies and mimic various failure scenarios and their dynamics. We report on observed trends that are revealed by our stress tests and quantified by a selection of metrics.
- We make Domino available to the research community; the testbench is built with an extensible plug-and-play design that invites third parties to use their own simulation/emulation engines, consider survivability metrics of their own choosing, and leverage their preferred router topologies to perform comparative evaluation studies of the inter-domain routing system. In this sense, it can also serve as a basis for stress testing other critical Internet infrastructures such as DNS or CAs.

Ethical Considerations: This work does not raise any ethical issues.

II. MOTIVATION AND BACKGROUND

To motivate our work, we examine the type and dynamics of the failure scenarios on which the Internet has built its survivability reputation and also review the existing literature on mathematical network robustness studies based on abstract network-as-a-graph models.

A. Inter-Domain Routing

Inter-domain routing (BGP as default) plays a critical role in maintaining Internet connectivity and enabling recovery from failures. However, BGP has well-known limitations in resilience, and many works have proposed enhancements through protocol changes or configuration adjustments [5]–[7]. Despite these efforts, there is a lack of comprehensive tools to evaluate how changes to the logical topology of BGP sessions, protocol behaviors, and configuration defaults affect network resilience under diverse failure scenarios, operational scales, and dynamic conditions.

To measure and improve inter-domain routing, evaluation should focus on forwarding tables rather than the data plane, which involves factors like congestion control and transport behavior. There already exist several simulators and emulators that can capture BGP network dynamics and produce BGP update messages as output [8]–[12]. However, a survivability measurement tool must provide the simulator with realistic network state and failure scenarios, and process the resulting update messages to enable meaningful analysis of routing survivability.

B. Failure Scenario

While simulating a link failure is straightforward in most network simulators, meaningful impact assessment requires more than just matching the failure size—it also demands applying consistent temporal, geographical, or topological patterns. For instance, showing that a BGP modification improves resilience to random failures does not imply it will perform similarly under cascaded or correlated failure scenarios.

BGP-related failures extend beyond link/router failures and include complex events such as BGP hijacking and route leaks. However, our focus in this paper is on measuring the impact of structural/temporal failure types.

Event (date)	Scale (# of prefixes)	Duration (hours)
Code Red II (8/2001) [13]	NA	10^2
9/11 (9/2001) [14]	10^3	10^1
Northeast US blackout (8/2003) [15]	10^3	10^1
Italian blackout (9/2003) [15]	10^2	10^1
Hurricane Katrina (8/2005) [16]	10^2	10^2
Taiwan earthquake (12/2006) [17]	10^3	10^1
Egyptian Internet Shutdown (1/2011) [18]	10^3	10^2
Japan earthquake (4/2011) [19]	10^2	10^1
Hurricane Sandy (10/2012) [20]	10^3	10^1
War in Ukraine (2022) [21]	NA	10^4
Hurricane Ian (9/2022) [22]	NA	10^2

Table I: A subset of past events impacting the Internet.

C. Beyond Localized Failures

Over time, the Internet has faced multiple tests of its resilience, ranging from natural disasters like earthquakes, wildfires, and hurricanes to cyber-attacks initiated by individuals or organized entities, including nation-states. Table I lists a sampling of such reported incidents and shows that most of the listed incidents have been highly localized in both geographical space and duration, typically affecting only a small number of prefixes, mostly within the affected region.

As the Internet and its operating landscape continually evolve, it is crucial not only to consider past failure scenarios but also to anticipate new events. The changing physical environment of the Internet, such as climate change, and shifts in its operational dynamics, like the emergence of smart grids as power sources, suggest that these events may occur at unprecedented geographic scales, persist for extended durations, or present entirely novel disaster scenarios. For instance, the growing interdependence between the Internet and power grids [23]–[25], especially with the transition to future smart grids from the current highly centralized designs, introduces new possibilities but also brings forth novel and unknown vulnerabilities, necessitating careful attention.

Additional concerns stem from the ongoing evolution of the Internet, driven by economic, political, or societal factors. Examples encompass challenges arising from trends towards consolidation [26], an expanding digital divide, the rise of closed ecosystems (e.g., walled gardens operated by major cloud providers), and heightened risks of balkanization (e.g., censorship and the establishment of digital borders).

Existing studies on the survivability of the Internet often overlook extreme events—termed as low-probability and high-impact—that intermittently capture public attention. While absent from conventional research, these catastrophic events, including massive Coronal Mass Ejections (CME) [4], high-altitude electromagnetic pulse attacks (EMP) [27], and state-sponsored cyber warfare [28], possess the potential to disrupt vast areas, ranging from entire nations to continents, for extended periods, spanning days to weeks or even months. For instance, a massive CME event could disrupt satellite communications, parts of the global Internet infrastructure, and critical power grid functionalities simultaneously, leading to cascading failures and extensive outages with global repercussions.

Figure 1 provides a contextual overview of historical events, placing them based on their network impact scale (i.e., the number of impacted prefixes) and duration. It contrasts these events with hypothetical instances of new failure scenarios discussed earlier. The failure scenarios we analyzed in this work are drawn from the reddish-shaded region in Figure 1. Given the Internet’s growing significance as the “nervous system” of modern society and the global economy, the figure is also a reminder of the urgent need for a more comprehensive understanding of its actual robustness in the face of increasingly complex and severe realistic failure scenarios. This understanding necessitates the development of new tools and frameworks.

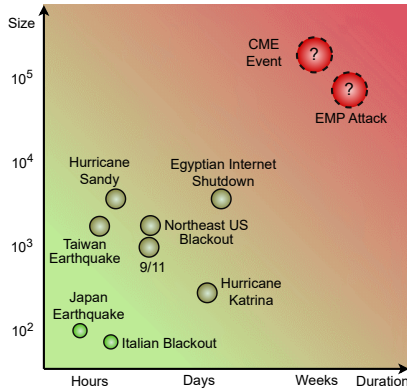


Figure 1: Historic events represented by impact (number of affected prefixes) and duration. Code Red II is not included due to a lack of reliable data

D. Beyond on/off Failure Modes

Previous research on Internet failures and their impact commonly adopts a binary “on/off” failure mode¹, where “off” indicates a failed network element or part of the network not functioning at all, and “on” signifies everything working as intended [5], [30]–[32]. This “failing-off” approach simplifies mathematical modeling, analysis, and simulation-based studies of network robustness but captures only one aspect of realistic failure scenarios, often not the most critical one. In contrast, a “failing-on” mode, where network elements or parts intermittently work and may function at a diminished capacity even when operational, introduces a dynamic element into the behavior of failure scenarios. Despite reflecting reality more accurately, this mode has been relatively understudied. This dynamic aspect becomes especially relevant in prolonged and large-scale events, where understanding changes during failures, considering external factors or internal effects, becomes crucial during, rather than just after, the failures occur [33], [34].

Intermittent network failures can arise in various forms. During prolonged power outages, colocation facilities and cell towers with backup systems may experience fail-on mode due to delays in fuel supply or battery replacement. In emerging massively distributed smart grids, isolated connectivity islands with varying power availability can emerge unpredictably within affected regions, depending on how smart grids balance power supply and demand [23], [25]. Similarly, worm propagation can generate excessive traffic, overloading routers and causing memory exhaustion, reduced capacity, or forced reboots, leading to intermittent or complete router failures [35]–[38].

E. Beyond Abstract Models

Traditional approaches to studying network robustness typically model networks as abstract graphs, where nodes and links are systematically removed based on predefined rules to simulate specific failure scenarios [39]–[44]. While these stylized models can provide valuable insights into abstract notions

¹Failure modes refers to the ways in which something might fail [29].

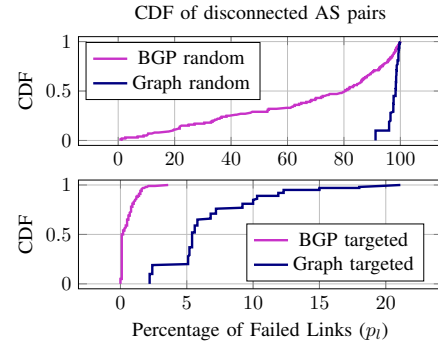


Figure 2: Topological vs. Protocol view of connectivity: The CDF for AS pairs disconnection by percentages of random (top) or targeted (bottom) link failures.

of network vulnerabilities, they lack the means to account for the complexities inherent in real-world Internet-like networks as they are subjected to realistic failure scenarios.

A key limitation of abstract graph models for Internet-like networks is their oversimplification of network dynamics. By treating networks as static graphs, they overlook the interaction between physical infrastructure—shaped by technical and economic factors—and adaptive routing protocols. These protocols determine how remaining connectivity is used following failures, guiding traffic flow when links or nodes go down [45].

To illustrate, Figure 2 shows how focusing only on the topology graph can be misleading. With 15 Tier 1 and 5 Tier 2 Autonomous Systems (ASes), the cumulative distribution functions (CDFs) display the number of disconnected AS pairs (y-axis) under different random and targeted link failure scenarios (x-axis). The random failure involves the random selection of failing links, whereas the targeted failure entails the earlier failure of links included in more routes. This simplified illustration highlights the disparity in AS pair disconnection between the topology with the inter-domain routing protocol (BGP) and without it.

Notably, the topology graph without BGP demands a substantial number of link disruptions to generate any disconnected AS pairs. In contrast, its BGP-based counterpart experiences disconnections with significantly fewer failures. This discrepancy arises because BGP, unlike the topology graph assumption, does not maintain all available paths between AS pairs and lacks support for AS-level multi-path routing.

III. TESTING FOR SURVIVABILITY

Before introducing the design of Domino, we provide a definition of survivability tailored to inter-domain routing and discuss key features that guide our design, including survivability metrics, failure scenarios, and modes, and the distinction between affected and unaffected regions.

a) *Defining Survivability*: Ellison et al. [46] define *survivability* as “the capability of a system [used in the broadest possible sense, including networks] to fulfill its mission, in a timely manner, in the presence of attacks, failures, or accidents.” As the core mission of the network layer is to ensure

connectivity, to “fulfill its mission” translates into “maintain connectivity” in our context.

b) Survivability Metrics: Following from our definition, we measure the disconnection time between each router-prefix pair and use the average disconnection time among all existing pairs as the key metric for measuring survivability. For each router-prefix pair, disconnection time is the period during which the router lacks a valid path to the prefix. We explain the details of calculating disconnection time in §IV-B.

While disconnection time serves as the main metric, to provide a more comprehensive perspective, our assessment of survivability includes additional metrics such as *convergence time* and the *number of update messages* but can also accommodate other metrics defined by a third-party. *Convergence time* captures the period during which routers update their forwarding tables. These updates can result in forwarding inconsistencies, leading to extended connectivity outages—quantified as *disconnection time* in our analysis. A large *number of route update* messages propagated in the network suggests that the network is experiencing a high level of instability, with routers frequently updating their routing tables and advertising new routes. High rates of route announcements can cause routers to experience substantial resource overheads, inducing delays in packet forwarding and processing. In extreme cases, this burden may escalate to the extent of routers crashing or malfunctioning. Moreover, monitoring the number of BGP update messages helps gauge how effectively the system manages and contains failures, preventing them from cascading and causing widespread disruption.

While disconnection time reflects network downtime, it does not fully capture Ellison’s timing aspects, such as task-specific guarantees. However, it provides a basis for future exploration of broader timing properties.

c) Different Failure Scenarios: Large-scale failures can be expected to exhibit highly diverse behaviors, varying greatly in their predictability, severity, impact, and geographical extent. Restricting stress testing to a single failure scenario may yield incomplete insights, as a system’s resilience can vary across different failure scenarios. In this paper, we consider four fundamental failure scenarios: random [31], [41], [47], [48], cascaded [49]–[55], regional [31], [56]–[59], and depeering [31], [60], [61]. The specific details and descriptions of these scenarios are given in §IV. In practical scenarios, a failure event may exhibit characteristics aligning with one or a combination of these scenarios.

d) Failing-off and Failing-on Modes: To accommodate the above-described dynamic failure modes, we do not restrict ourselves to simple fail-off scenarios, where failed network elements remain non-operational throughout the experiment. Specifically, we delve into modeling fail-on modes, where network elements undergo intermittent failures, capturing the dynamic nature of real-world network failures (§IV-D).

e) Affected and Unaffected Regions: In large-scale failure scenarios, the network can be roughly divided into two regions: parts directly impacted by the failure (affected region) and parts that are not (unaffected region). The affected region

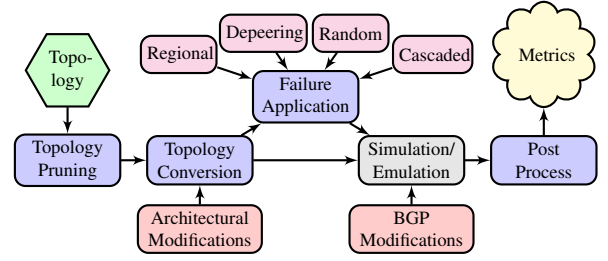


Figure 3: Key components of Domino illustrating the modularity of the tool and the flexibility in experimental evaluation of complex scenarios.

may be non-contiguous and can impact the unaffected region with ripple effects like the churn of routing updates and transient disconnections. In the unaffected region, the network should use available connectivity to maintain the appearance of a fully functioning Internet. In the affected region, connectivity has to be treated as a critical resource—detecting and preserving any remaining connectivity is crucial. Distinct survivability criteria in these regions necessitate separate and tailored evaluations for each.

IV. A TESTBENCH FOR STRESS-TESTING

In this section, we present Domino’s design and implementation by describing the components, metrics, failure scenarios, and modes of the current version of the testbench.

A. Testbench components

Figure 3 illustrates the different components of the Domino testbench, including topology conversion and pruning, simulation/emulation, failure application, and postprocessing. Domino accepts a network topology as input and allows the experimenter to select among different failure scenarios and modes for evaluation.

a) Topology: Domino accepts a network topology as its primary input. This topology consists of ASes and their interconnecting links (e.g., CAIDA AS-rel dataset [62]), including redundant links. In this case, individual links can be mapped to specific routers within each AS. Additionally, the input topology can be augmented with router location and intra-AS topology information. The input format for the topology is a text file, where each line follows the format: *from_AS|to_AS|relationship|from_rt_loc|to_rt_loc*, where *from_rt_loc* and *to_rt_loc* are the space-separated longitude and latitude of the respective routers, if available.

We use two datasets to construct the network topology and obtain geographical information. First, we leverage the AS-rel-geo dataset [63]. This dataset is an annotated version of the CAIDA AS relationships dataset [62] and includes the links between ASes, along with *best-effort* estimates of the geographic location of these links. Additionally, we use the Macroscopic Internet Topology Data Kit (ITDK) [64], which offers the assignment of routers to ASes and includes the geolocation of network routers. The ITDK dataset, in comparison with the AS-rel-geo dataset, includes router geolocations but lacks Multi-Lateral Peering (MLP) links.

b) *Topology Pruning*: This component is tasked with selecting a subset of ASes from the full network topology when the employed simulator experiences scalability limitations. The objective is to ensure that the resulting topology reflects the structure of the real-scale topology as closely as possible. One effective approach is to prune the topology while preserving the hierarchical structure of the Internet topology. The default pruning strategy begins with 15 Tier 1 ASes and uses the CAIDA provider-peer customer cones dataset [62] to identify subsequent hierarchies of ASes. The topology’s depth and width are customizable, enabling manual configuration.

c) *Topology Conversion*: The initial task of this component involves converting the input topology format to the Free Range Routing (FRR) configuration format. Alongside this conversion, the component incorporates the intra-AS topology, defaulting to a star topology with a reflector positioned at the center. Furthermore, the component can configure a fixed number of routers per AS (default value is set to 10). This parameter is designed to be sufficiently larger than one to accommodate intra-AS complexities while avoiding excessive values that could compromise scalability. Then, links are assigned to these routers based on geographical clustering, utilizing geo-information from the topology if available. Another task of this component is to modify the topology according to specified architectural adjustments. This capability is particularly useful when a study demands modifications to the network architecture, such as grouping a set of ASes into a core with distinct policies. Lastly, the component configures the simulation parameters, including the Minimum Route Advertisement Interval (MRAI) timer.

d) *Failure Application*: The Failure Event Application component introduces failure events into the topology. It includes the injection of specific failure patterns or scenarios to simulate realistic failure occurrences. This step involves determining the location and timing of failures within the network and applying them to the appropriate nodes or links.

e) *Simulation/Emulation*: Domino incorporates a BGP simulator/emulator to explore the behavior of the routers under various failure events. The input of this component is a topology in FRR format, a set of announced prefixes, and a set of failed links. The output includes the updates/withdrawals in the forwarding table of network routers.

f) *Post Process*: The process component handles the analysis and processing of the simulator’s output data. It involves extracting relevant information, calculating the survivability metrics, and aggregating the results.

B. Metrics

We next describe how we calculate each of the survivability metrics mentioned in §III.

a) *Convergence time*: The metric is computed individually for each prefix, representing the time elapsed between the observation of the first and last update messages for that prefix. Essentially, it quantifies the duration required for all BGP routers in the network to establish valid and stable paths to the specified prefix.

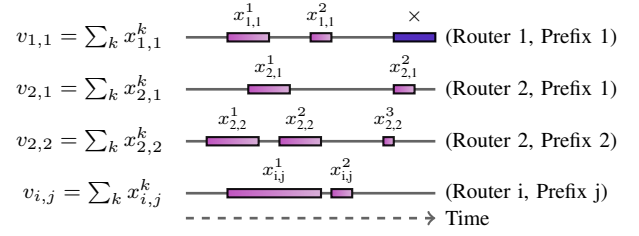


Figure 4: Disconnection time for router-prefix pairs $(v_{i,j})$: the cumulative duration of disconnection periods (in purple) along time. Persistent disconnection periods are represented in blue.

b) *Disconnection time*: We calculate the disconnection time by analyzing control-plane updates. To measure the transient disconnection periods, for each update, we examine its impact on the connectivity between the source router and the corresponding prefix. This examination extends to all routers whose route to the prefix, constructed iteratively through each router’s next hop for the respective prefix, involves the examined source router. For example, if a router eliminates a prefix from its forwarding table, it results in the loss of packets forwarded from other routers to this specific router for reaching the said prefix, creating what is referred to as a “black hole.” As depicted in Figure 4, disconnection intervals (in purple) may occur along the time axis for a (router, prefix) pair. The cumulative duration of these intervals along the time axis represents the disconnection time for each (router, prefix) pair. The blue period (denoted by \times) signifies a disconnection due to graph disconnectivity or BGP policy compliance and is excluded from the disconnection time metric.

Algorithms 1 to 3 demonstrate how we derive this metric from control-plane messages. In these algorithms, the variables *next_hop* and *pre_routers* are implemented as key-value dictionaries. The *next_hop* dictionary associates each (router, prefix) pair with its corresponding next hop, while the *pre_routers* dictionary maps each (router, prefix) pair to the list of routers that have this router as their next hop for reaching that specific prefix. The set of routers whose traffic must pass through a particular router to reach a given prefix is defined as the descendants of that router for that prefix.

For every update, if it is a withdrawal, the algorithm logs the (router, prefix) pair into the set of disconnected pairs, denoted as *dis_set*, and records the corresponding timestamp in *dis_t*. Subsequently, this process is extended to all descendants, identified by calling the *find_desc(router, prefix)* function (refer to line 3 of Algorithm 1). Upon receiving an update with a valid path, if the router is in *dis_set*, indicating having no path or an invalid path to the prefix (lines 3 and 16 in Algorithm 3), the router is removed from *dis_set*, and the time-lapse between disconnection and reconnection is recorded as a disconnection period for both the router and all its descendants. In the case of an update with an invalid path, if the router previously had a valid path (as per line 21 of Algorithm 3), the (router, prefix) pair is added to *dis_set*, and the corresponding timestamp is recorded in *dis_t*. This

Algorithm 1 Disconnection Time Calculation

```
1: Input: BGP updates (updates)
2: OUTPUT: discon. time for each router-prefix pair (v)
3: Initialize v, dis_set, dis_t, next_hop, and pre_routers
4: for u  $\in$  updates do
5:   t  $\leftarrow$  u.timestamp; s  $\leftarrow$  u.router; p  $\leftarrow$  u.prefix
6:   if u.type is withdrawal then
7:     HandleWUpdate(s, p)
8:   else if u.type is announcement then
9:     HandleAUpdate(s, p)
```

Algorithm 2 HandleWUpdate(*s*, *p*)

```
1: if (s, p)  $\in$  next_hop then
2:   nh_previous  $\leftarrow$  next_hop[(s, p)]
3:   Remove (s, p) from next_hop
4:   Remove s from pre_routers[nh_previous, p]
5:   Add (s, p) to dis_set; dis_t[(s, p)]  $\leftarrow$  t
6:   dsct  $\leftarrow$  find_dsct(pre_routers, s)
7:   for d  $\in$  dsct do
8:     Add (d, p) to dis_set; dis_t[(d, p)]  $\leftarrow$  t
```

process is repeated for all descendants.

Algorithm 3 HandleAUpdate(*s*, *p*)

```
1: if (s, p)  $\notin$  next_hop then
2:   Update next_hop and pre_routers accordingly
3:   if (nh, p)  $\notin$  dis_set then
4:     Remove (s, p) from dis_set
5:     v[(s, p)]  $+=$  t - dis_t[(s, p)]
6:     dsct  $\leftarrow$  find_dsct(pre_routers, s)
7:     for each d  $\in$  dsct do
8:       Remove (d, p) from dis_set
9:       v[(d, p)]  $+=$  t - dis_t[(d, p)]
10:  else
11:    nh_previous  $\leftarrow$  next_hop[(s, p)]
12:    if nh  $\neq$  nh_previous then
13:      next_hop[(s, p)]  $\leftarrow$  nh
14:      add s to pre_routers[(nh, p)]
15:      remove s from pre_routers[nh_previous, p]
16:      if (s, p)  $\notin$  dis_set and (nh, p)  $\in$  dis_set then
17:        add (s, p) to dis_set; dis_t[(s, p)]  $\leftarrow$  t
18:        dsct  $\leftarrow$  find_dsct(pre_routers, s)
19:        for d  $\in$  dsct do
20:          add (d, p) to dis_set; dis_t[(d, p)]  $\leftarrow$  t
21:        if (s, p)  $\in$  dis_set and (nh, p)  $\notin$  dis_set then
22:          remove (s, p) from dis_set
23:          v[(s, p)]  $+=$  t - dis_t[(s, p)]
24:          dsct  $\leftarrow$  find_dsct(pre_routers, s)
25:          for d  $\in$  dsct do
26:            remove (d, p) from dis_set
27:            v[(d, p)]  $+=$  t - dis_t[(d, p)]
```

c) Number of update messages: We measure this metric per router as the count of BGP update messages that a router exchanges during the network's recovery from a failure.

C. Failure Scenarios

Domino allows the experimenter to select among different failure scenarios and modes. Below, we explain our approach to modeling each of the main categories of failure scenarios and failure modes.

a) Regional: A regional failure can cause network elements in a specific geographic area to fail; however, the impact can ripple across the broader network, triggering subsequent failures. To model these failures, we simultaneously fail multiple routers within one specific geographical region and study the impact on the whole network. The way the testbench accomplishes this task depends on whether or not the input topology includes information about the geolocation of routers. With this information, the topology conversion module schedules failure events for the routers in the user-specified region. Otherwise, the conversion module schedules the failure events based on a distribution function. In this function, each link has a probability of failure determined by its end ASes. Users can provide this distribution as input. Otherwise, the topology conversion module derives these probabilities from the ITDK dataset for the user-specified region. The details of how the conversion module achieves this are explained in Appendix A.

In the fail-on regional scenario, failure events happen with the same pattern but with a delay within a certain margin.

b) De-peering: In a de-peering scenario, two ASes de-peer as a result of a failure event, such as a targeted attack aiming to sever their connectivity by failing all the links between them. Understanding the survivability of the network under such failure scenarios is important for identifying or pinpointing network vulnerabilities in the form of weak or "soft" spots and can also be beneficial for studying the impact of intentional attacks on a network's infrastructure. To model these scenarios, the testbench selects a set of neighboring AS pairs to de-peer—by probabilistically selecting each pair with a probability of p_s , i.e., de-peering probability—and schedules failure events for the selected AS pairs. At each failure event associated with a pair of ASes, the simulator cuts all the links between them. In a fail-on version, the testbench schedules such failure events to re-appear with a specific pattern.

c) Cascading: In cascading failures, secondary failures emerge as a consequence of preceding primary failures that propagate through interconnected elements, exacerbating the initial disruption and expanding its scope. Such a chain reaction can result in a significant disruption of network communication and connectivity as the failure spreads through interconnected routers and their associated links. To model such scenarios, the testbench first selects a random link and schedules its failure as a primary failure. Then, it schedules consequent failures accordingly: for each router connected to a primarily failed link, it iterates over all links connected to the router and decides whether each link should fail with a probability p_p , i.e., failure propagation probability. In the fail-on version, the testbench repeats the pattern of each failure for its consequent failure.

d) Random: Random failures, where a random subset of links fails, do not necessarily represent real-world scenarios but serve as valuable tools for modeling certain components or aspects of real-world failures. Randomly failing links tests a network's response to (unpredictable) failures that are uniformly distributed across the different parts of the network. To

model this scenario, the testbench uses the commonly studied case in which it randomly selects a fraction of network links to fail, where the percentage of the links experiencing failure (off or on) is denoted by p_l . In the fail-on version of a random scenario, each failing link has an independent failure pattern from other links.

D. Failure Modes

As discussed in §II-D, a significant portion of prior research on Internet failures has used a binary on/off model. However, to fully capture the complexities of large-scale failures, it is essential to consider “fail-on” dynamics, where individual network elements operate intermittently.

In our approach to modeling the fail-on dynamics, the failure pattern of a failing link is characterized by intervals of link status changes between up and down. We opt for the Pareto distribution for modeling these intervals due to its heavy-tailed nature, which allows for a high degree of variability in the lengths of these intervals, with most of them being short in duration and a few lasting for extremely long periods of time. This feature makes the Pareto distribution a widely accepted choice for modeling losses from catastrophic events [65], [66]. Specifically, we utilize the Pareto Type II (Lomax) distribution: $f(x; \lambda, \alpha) = \frac{\alpha}{\lambda} \left(1 + \frac{x}{\lambda}\right)^{-(\alpha+1)}$, where α is the shape parameter, λ is the scale parameter, and the mean failure interval is $2 * \frac{\lambda}{1-\alpha}$.

E. Flexibility

Domino is designed with a high degree of flexibility, ensuring adaptability to diverse testing scenarios. The modularity of Domino allows for the replacement of individual components. Notably, the employed simulator can be substituted with an alternative simulator or an emulator, provided that the output format of update messages remains consistent.

F. Simulation Scale

Given that our objective is to explore unprecedented large-scale failure scenarios, we need to rely on simulation/emulation, as empirical data is unavailable for such instances. The choice of the experimental setting for such a simulation-based study—both the scale and level of detail of the simulated network—is crucial for ensuring that Domino helps draw insightful conclusions and guide Internet design.

Our focus on prolonged large-scale failures, where understanding the dynamic behavior during the event is required, necessitates a simulator capable of capturing the transient states of network elements (e.g., the changes in a router’s forwarding table). This level of detail requires an event-driven simulation. Note that in most other studies where only the before and after states of the BGP routers are of interest, the simulator does not necessarily need to be event-driven, allowing for improved scalability [30], [31], [67], [68].

While event-driven simulators like SimBGP [9] offer the depth of detail we need, they have scalability limitations [8], [10], [30], [69]–[71]. Two approaches to address this are: (1) reducing the network topology size until the simulations

become tractable, and (2) using a distributed architecture of emulators like Kathara [11] or Seed [12] across multiple machines. We opt to follow approach (1) in this paper because this improves the ease-of-use of Domino and the reproducibility of the results.

Our future work involves replacing the current simulator (SimBGP) with the Kathara emulator and running it at scale using Megalos [72], which is a scalable architecture for the virtualization of large network scenarios. This will enable us to apply the same failure scenarios to larger-scale topologies.

G. Implementation and Simulation Setup

We implemented Domino with ~ 3000 lines of code (excluding the simulator component) in Python. We run Domino on a high-performance cluster with a share size of 128 cores and 512 GB of RAM. All code and data we have used will be made available as open-source. We use SimBGP [9], an event-driven BGP simulator in Python that provides the detailed control plane updates our experiments require, and has been validated in prior work [73]. Due to the scalability limitations of the simulation, we restrict our explorations to moderate-size topologies. For the majority of the experiments, we use a topology with over 2,700 routers and 20,000 links across 250 ASes pruned from the CAIDA AS-rel-geo topology, preserving the hierarchical structure of the Internet. For comparisons between affected and unaffected regions, we employ a topology pruned from the ITDK topology, featuring 2,200 routers and over 18,000 links spanning 200 ASes (refer to §IV-A for details on the distinctions between these topologies). These topologies are comparable to or larger than those used in previous related studies [8], [10], [30].

We set the MRAI timer to the default value of 30s for BGP connections and 5s for IBGP connections [74]. Processing delay and link delay are uniformly distributed between $[1ms, 10ms]$ and $[10ms, 100ms]$, respectively, inherited from the version of SimBGP used in prior work [73].

H. Limitations

Conducting real-world stress tests on the Internet is infeasible, rendering direct verification of results through real-world repetitions impractical. It is also not feasible to reproduce past events; even if route collectors have data from the time of an event (such as the 2006 Taiwan earthquake), we lack detailed information on which specific routers in our employed topology failed and when. However, Domino focuses on relative comparisons, instead of reproducing the ground truth, acknowledging that simulators/emulators cannot capture all the real-world effects.

The reliability of simulated BGP behavior depends on the chosen simulator/emulator. This necessitates using a widely-adopted and validated simulator/emulator, such as SimBGP [9].

V. EVALUATION RESULTS

In this section, we begin by evaluating the post-processing component of Domino with real-world historical data to verify the relevance of the metrics calculated in this component.

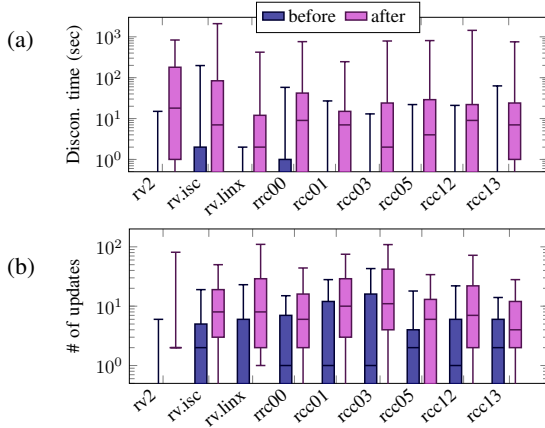


Figure 5: a) Disconnection time and b) count of update messages per prefix observed by different collectors, spanning a 1-hour period *before* and *after* the 2006 Taiwan earthquake.

Subsequently, we demonstrate the capabilities of Domino, with illustrative examples and report on observed trends that are revealed by our stress tests.

A. Real-world Data: The Taiwan 2006 Earthquake

To validate the survivability metrics considered by Domino, We employ the post-processing component of Domino to evaluate the disconnection time observed by nine BGP route collectors from RouteViews and RIPE RIS, spanning one hour before and one hour after the 2006 Taiwan earthquake (refer to Table I). Figure 5 shows the results. Given that these collectors represent only a subset of nodes within the Internet network, we apply a tailored disconnection time calculation approach. The disconnection period in this approach is the interval between withdrawing a prefix and identifying a route for that prefix.

Our observation in Figure 5a reveals a notable contrast in disconnection time pre- and post-failure event, affirming the efficacy of our metric in capturing network failover dynamics. It is crucial to highlight that our disconnection time measurements exclusively pertain to prefixes undergoing transient disconnection, not those entirely offline. Importantly, these affected prefixes may not necessarily originate from the affected region. In Figure 5b, we depict the count of BGP updates received by each collector. Once again, notable differences become evident before and after the failure, underscoring the sensitivity of these metrics to disruptive events.

B. Connectivity in Topology versus Protocol

Figure 6a provides insights into the spread and median convergence times of the prefixes after a cold boot (i.e., starting the entire network with empty routing tables) without any failures. The differing variability across sizes stems from pruning topologies from the CAIDA AS-level graph rather than generating them synthetically, resulting in distinct connectivity structures. Figure 6b illustrates the normalized reachability of routers over time. Here, “reachability” refers to the count of stably connected pairs of (router, prefix) identified through time. We observe that in a fully connected topology

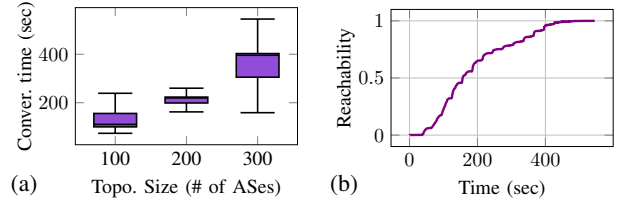


Figure 6: a) Cold-boot BGP convergence time in three topology sizes. b) Reachability over time for the 300 ASes topology.

and without any failures, routers exhibit notable convergence times. Specifically, the median convergence time for prefixes within the 300 ASes topology is ~ 400 s, while achieving 50% reachability across the network takes more than 150s. This highlights the prolonged time-to-connectivity in BGP, as the inter-domain routing protocol, emphasizing the need to move beyond abstract consideration of connectivity solely as topological connectivity (refer to §II-E).

C. The Impact of Failure Scale and Dynamics

Figures 7 and 8 depict the results of the stress tests for the different failure scenarios. Note that Figures 7 to 12 display normalized values to emphasize trends rather than absolute values.

a) *Failure Scale*: We observe an exponential rise in the disconnection time as the size of the failure increases. Additionally, there is a substantial rise in the number of update messages. However, we observe that convergence time does not necessarily mirror the trends observed in disconnection time, indicating that convergence does not imply disconnection.

b) *Failure Dynamics*: In fail-on mode, where links are intermittently active during the event duration, in contrast to fail-off mode where links remain off after failing, we still observe a higher transient disconnection time and number of updates. This underscores the significance of taking into account the dynamics of failures.

D. Affected and Unaffected regions

Figure 9 illustrates variations in the average disconnection time across affected/unaffected regions, categorized into four groups. The A-A category represents the disconnection time between routers and prefixes within the affected region, averaged by the product of the number of routers and prefixes in that region. The A-U category pertains to disconnection times between routers within the affected region and prefixes within the unaffected region. Conversely, the U-A category focuses on disconnection times for unaffected routers to affected routers, while the U-U category examines disconnection times for unaffected routers to unaffected prefixes. Figure 9a presents the results for failure sizes of 50% and 20% in the US. Figure 9b compares a 50% large failure in the US with a similar failure in Germany.

Domino enables selecting specific routers and prefixes to investigate their disconnection during failure events. For instance, it allows an examination of the disconnection time

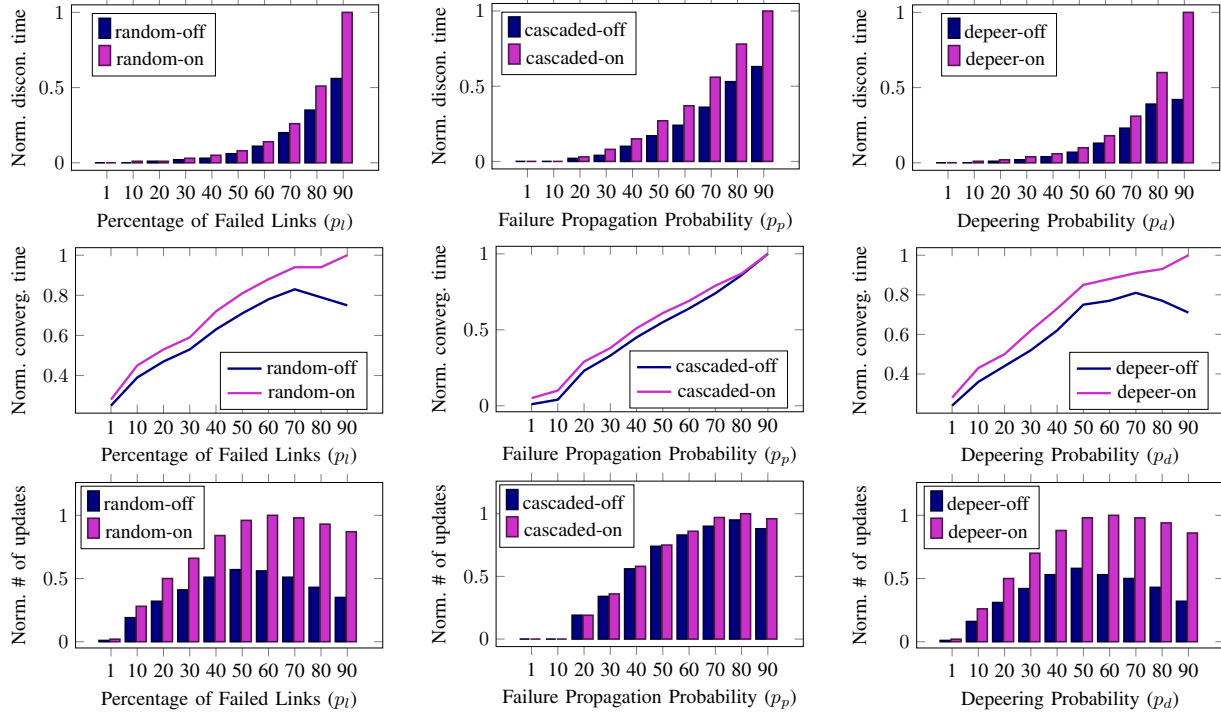


Figure 7: The impact of failure scale on disconnection time, convergence time, and the number of update messages highlighting the dynamics of fail-on/fail-off for the three event scenarios: random, cascading, and de-peering.

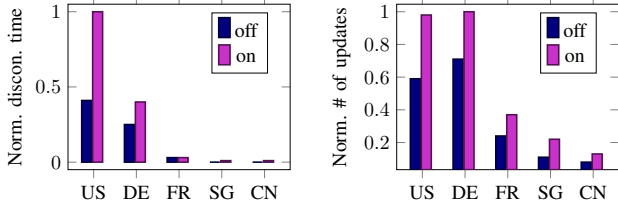


Figure 8: The impact of fail-on/fail-off dynamics on regional failures for different affected countries: the United States (US), Germany (DE), France (FR), Singapore (SG), and China (CN).

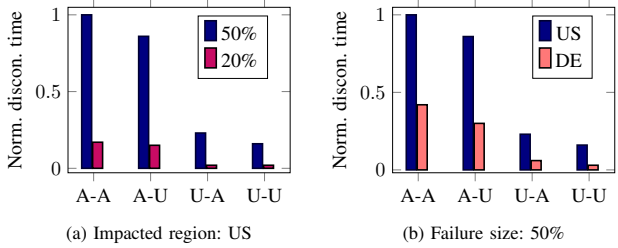


Figure 9: Disconnection time for Affected-to-Affected (A-A), Affected-to-Unaffected (A-U), Unaffected-to-Affected (U-A), and Unaffected-to-Unaffected (U-U) router-prefix pairs for two failure sizes: 50% and 20% (left) and two affected regions: US and Germany (right).

between Region A's routers and Region B's prefixes during a failure event with size $x\%$ in Region C.

E. The Impact of Fail-On Pattern Parameters

Domino provides the flexibility to choose diverse patterns for fail-on mode. In Figure 10, we exemplify the distinction

between two sets of parameters (refer to §IV-D): one with $\alpha = 1.5$ and $\lambda = 10$ for both the failure and revival periods, resulting in a mean failure/revival period of 20s (average of 40s between failures), and the other with $\alpha = 1.5$ and $\lambda = 5$, yielding an average interval of 20s between consecutive failures. While fail-on mode exhibits a higher transient disconnection time than fail-off mode (observed in Figure 7), for more frequent failures, the impact is more effectively regulated by the MRAI timer, improving disconnection time.

F. BGP Modifications

Domino enables the evaluation of various BGP modifications, ranging from adjustments to input configuration parameters to the incorporation of an entirely new simulator, as long as the input/output format remains consistent. Figure 11 shows the disconnection time and the number of updates for MRAI timer values other than the default values, i.e., 12s rather than 30s for eBGP and 2s rather than 5s for iBGP. We observe that the MRAI timer proves effective in mitigating transient disconnections and reducing the overall number of messages by delaying the propagation of update messages.

G. The impact of Failure Scenario

Figure 12 illustrates the disconnection time and number of updates versus the percentage of failed links in different failure scenarios. We observe that the cascaded scenario exhibits the most severe impact even with a smaller number of failed links. The reason for the decline in the largest failure sizes is that the topology has fragmented into smaller islands that

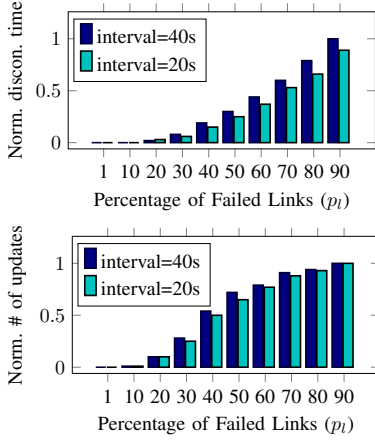


Figure 10: The impact of the fail-on pattern parameters on the cascaded scenario for two mean failure intervals: 40s and 20s.

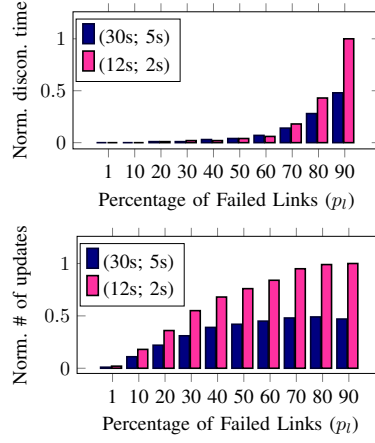


Figure 11: The impact of the MRAI timers, for two timer pairs (eBGP; iBGP) in seconds, in the random scenario.

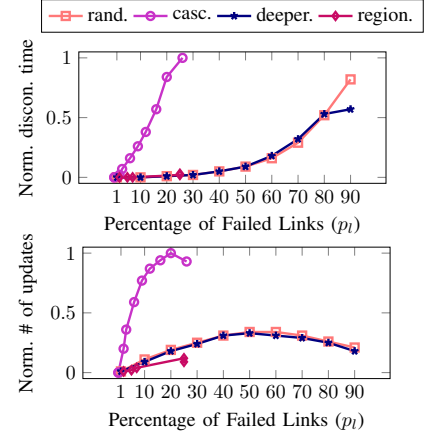


Figure 12: The impact of the failure scenario on the disconnection time and number of updates versus the percentage of failed links.

converge faster within themselves, but overall connectivity in the topology has deteriorated.

H. Main Takeaways

The following points summarize the key insights from our simulation results:

- Fail-on scenarios are up to 60% more disruptive than fail-off scenarios.
- Convergence time is similar across fail-on and fail-off modes, but disconnection time differs significantly, indicating that convergence time alone does not adequately capture downtime.
- A cascading failure affecting just 3% of the links can disrupt connectivity almost as much as a 40% random failure, highlighting the outsized impact of structurally targeted failures.
- While a higher MRAI-timer delays updates, the default 30s setting turns out to be more efficient than a lower setting (e.g., 12s).
- Disconnection time is asymmetric: routers in affected regions experience longer delays accessing unaffected prefixes as compared to routers in unaffected regions accessing affected prefixes.
- The impact of regional failures is not strictly proportional to the size of the geographic area.

VI. RELATED WORK

There exists an extensive body of literature concerned with assessing or enhancing the resiliency of the Internet to attacks and natural or man-made disasters. For example, the physics and complex networks literature typically considers highly stylized models of networks that are amenable to mathematical analysis or simulation-based studies (e.g., [39], [41], [75]), but this level of abstraction is of little relevance for real-world outage scenarios in networks such as the Internet with

its protocols, some explicitly designed to ensure continued operation in the presence of failures [45], [76], [77].

In the networking literature, one category of studies deals with resilience at the physical rather than the network layer [78]–[83], topology-level resilience [31], or assessing intra-domain reliability [84]–[86]. While largely complementary to our work, most of these efforts focus on traditional failure scenarios of the fail-off type and are not directly relevant for highly volatile scenarios that our tool can capture and that typically arise in the context of large-scale disasters where network elements intermittently failing on is the rule rather than the exception. A few works have undertaken empirical analyses of past events to examine the impact of failures, referred to as large-scale in these studies, albeit on a smaller scale compared to catastrophic yet plausible events [87]–[89].

Many other research studies have employed various evaluation methods in conjunction with their proposed network-layer resilience enhancement approaches. These approaches range from modifications to BGP [5], [6], [90]–[98], such as the addition of fast fail-over mechanisms [99], [100], to the proposal of new internet architectures [30], [101], [102]. However, their employed evaluation schemes have various limitations, particularly in terms of incorporating diverse scales, dynamicity, and the range of failure scenarios that can be represented, which are hindering a comprehensive understanding of the resilience capabilities of the proposed schemes and their applicability to real-world network scenarios.

These solutions need to be thoroughly evaluated against large-scale and dynamic failure scenarios, such as those caused by solar superstorms or massive coronal mass ejection [4] and that differ from commonly-studied events in size (e.g., geographic extent, duration, intensity) as well as in their dynamic nature (i.e., highly intermittent outage patterns in both time and space). This has been a major motivation for our work.

VII. CONCLUSION

As an essential component of any critical infrastructure, the implementation of a continuous assessment and improvement cycle is key for identifying and potentially mitigating issues related to availability. Large-scale failure events, such as earthquakes or region-wide blackouts, can introduce unforeseen complexities that challenge system resilience, including intricate dependencies leading to cascading impacts and convergence problems. How can we detect and mitigate such effects and increase the resilience of the Internet infrastructure?

We propose a first step toward achieving this vision with Domino, a system for comparative evaluation of the survivability of the inter-domain routing infrastructure in the face of large-scale failures. We make Domino available as open-source software and publish all datasets we have used to enable independent validation of all results in this paper. As a next step, we anticipate expanding Domino to include additional critical infrastructure such as the DNS and web infrastructure.

REFERENCES

- [1] T. E. Commission, "Critical infrastructure," https://home-affairs.ec.europa.eu/pages/page/critical-infrastructure_en, 2013, accessed: 2023-2-7.
- [2] U. Cybersecurity and I. security agency, "Critical infrastructure sectors," <https://www.cisa.gov/critical-infrastructure-sectors>, 2020, accessed: 2023-2-7.
- [3] T. O. of American States, "Declaration protection of critical infrastructure from emerging threats," <https://www.oas.org/es/sms/cicte/ciberseguridad/publicaciones/CICTE%20DOC%201%20DECLARATION%20CICTE00955E04.pdf>, accessed: 2023-2-7.
- [4] S. A. Jyothi, "Solar superstorms: planning for an internet apocalypse," in *Proceedings of the 2021 ACM SIGCOMM 2021 Conference*, 2021, pp. 692–704.
- [5] D. Pei, X. Zhao, L. Wang, D. Massey, A. Mankin, S. F. Wu, and L. Zhang, "Improving bgp convergence through consistency assertions," in *Proceedings. Twenty-First Annual Joint Conference of the IEEE Computer and Communications Societies*, vol. 2. IEEE, 2002, pp. 902–911.
- [6] A. Bremner-Barr, Y. Afek, and S. Schwarz, "Improved bgp convergence via ghost flushing," in *IEEE INFOCOM 2003. Twenty-second Annual Joint Conference of the IEEE Computer and Communications Societies (IEEE Cat. No. 03CH37428)*, vol. 2. IEEE, 2003, pp. 927–937.
- [7] S. Deshpande and B. Sikdar, "On the impact of route processing and mrai timers on bgp convergence times," in *IEEE Global Telecommunications Conference, 2004. GLOBECOM'04.*, vol. 2. IEEE, 2004, pp. 1147–1151.
- [8] X. A. Dimitropoulos and G. F. Riley, "Efficient large-scale bgp simulations," *Computer Networks*, vol. 50, no. 12, pp. 2013–2027, 2006.
- [9] J. L. Sobrinho, L. Vanbever, and F. Le, "Distributed route aggregation (dragon) simulator based on simbgp," https://github.com/network-aggregation/dragon_simulator, 2014.
- [10] J. H. Cowie, D. M. Nicol, and A. T. Ogielski, "Modeling the global internet," *Computing in Science & Engineering*, vol. 1, no. 1, pp. 42–50, 1999.
- [11] M. Scazzariello, L. Ariemma, and T. Caiazzi, "Kathará: A lightweight network emulation system," in *NOMS 2020-2020 IEEE/IFIP Network Operations and Management Symposium*. IEEE, 2020, pp. 1–2.
- [12] W. Du, H. Zeng, and K. Won, "Seed emulator: An internet emulator for research and education," in *Proceedings of the 21st ACM Workshop on Hot Topics in Networks*, 2022, pp. 101–107.
- [13] J. Cowie, A. Ogielski, B. Premore, and Y. Yuan, "Global routing instabilities triggered by code red ii and nimda," 2001.
- [14] N. R. Council *et al.*, *The Internet under crisis conditions: learning from September 11*. National Academies Press, 2003.
- [15] J. H. Cowie, A. T. Ogielski, B. Premore, E. A. Smith, and T. Underwood, "Impact of the 2003 blackouts on internet communications," *Preliminary Report, Renesys Corporation (updated March 1, 2004)*, 2003.
- [16] J. Cowie, A. Popescu, and T. Underwood, "Impact of hurricane katrina on internet infrastructure," *Report, Renesys*, 2005.
- [17] A. Popescu, T. Underwood, and E. Zmijewski, "Quaking tables: The Taiwan earthquakes and the Internet routing table," in *Proc. of NANOG*, vol. 39, 2007.
- [18] A. Dainotti, C. Squarcella, E. Aben, K. C. Claffy, M. Chiesa, M. Russo, and A. Pescapé, "Analysis of country-wide Internet outages caused by censorship," in *tons*, 2014.
- [19] K. Cho, C. Pelsser, R. Bush, and Y. Won, "The japan earthquake: the impact on traffic and routing observed by a local isp," in *Proceedings of the Special Workshop on Internet and Disasters*, 2011, pp. 1–8.
- [20] J. Heidemann, L. Quan, and Y. Pradkin, "A preliminary analysis of network outages during hurricane Sandy," USC/Information Sciences Institute, Tech. Rep. ISI-TR-2008-685b, 2012.
- [21] M. Tripathy, "How is starlink changing connectivity?" <https://www.smithsonianmag.com/science-nature/how-is-starlink-changing-connectivity-180980735/>, 10 2022.
- [22] F. C. Commission, "Communications status report for areas impacted by hurricane ian," September 2022.
- [23] J. D. Taft and A. S. Becker-Dippmann, "The emerging interdependence of the electric power grid & information and communication technology," Pacific Northwest National Lab.(PNNL), Richland, WA (United States), Tech. Rep., 2015.
- [24] L. Martins, R. Girao-Silva, L. Jorge, A. Gomes, F. Musumeci, and J. Rak, "Interdependence between power grids and communication networks: A resilience perspective," in *DRCN 2017-Design of Reliable Communication Networks; 13th International Conference*. VDE, 2017, pp. 1–9.
- [25] Z. Wang, G. Chen, L. Liu, and D. J. Hill, "Cascading risk assessment in power-communication interdependent networks," *Physica A: Statistical Mechanics and its Applications*, vol. 540, p. 120496, 2020.
- [26] Internet Society, "Consolidation in the Internet economy," 2019 Global Internet Report, Feb 2020. [Online]. Available: <https://future.internetsociety.org/2019/>
- [27] N. Bajema, "One atmospheric nuclear explosion could take out the power grid," *IEEE Spectrum*, September 2021.
- [28] N. R. Council, *Terrorism and the Electric Power Delivery System*. The National Academies Press, 2012. [Online]. Available: <https://www.nap.edu/catalog/12050/terrorism-and-the-electric-power-delivery-system>
- [29] A. S. for Quality (ASQ), "Failure mode and effect analysis (fmea)," <https://asq.org/quality-resources/fmea>, 2019.
- [30] P. B. Godfrey, I. Ganichev, S. Shenker, and I. Stoica, "Pathlet routing," *ACM SIGCOMM Computer Communication Review*, vol. 39, no. 4, pp. 111–122, 2009.
- [31] J. Wu, Y. Zhang, Z. M. Mao, and K. G. Shin, "Internet routing resilience to failures: analysis and implications," in *Proceedings of the 2007 ACM CoNEXT conference*, 2007, pp. 1–12.
- [32] F. Wang, Z. M. Mao, J. Wang, L. Gao, and R. Bush, "A measurement study on the impact of routing events on end-to-end internet path performance," *ACM SIGCOMM Computer Communication Review*, vol. 36, no. 4, pp. 375–386, 2006.
- [33] L. Wang, X. Zhao, D. Pei, R. Bush, D. Massey, A. Mankin, S. F. Wu, and L. Zhang, "Observation and analysis of bgp behavior under stress," in *Proceedings of the 2nd ACM SIGCOMM Workshop on Internet measurement*, 2002, pp. 183–195.
- [34] D.-F. Chang, R. Govindan, and J. Heidemann, "An empirical study of router response to large bgp routing table load," in *Proceedings of the 2nd ACM SIGCOMM Workshop on Internet measurement*, 2002, pp. 203–208.
- [35] M. Roughan, J. Li, R. Bush, Z. Mao, and T. Griffin, "Is bgp update storm a sign of trouble: Observing the internet control and data planes during internet worms," 2006.
- [36] M. Lad, X. Zhao, B. Zhang, D. Massey, and L. Zhang, "Analysis of bgp update surge during slammer worm attack," in *IWDC*, vol. 2918, 2003, pp. 66–79.
- [37] J. Cowie, "Global routing instabilities during code red 11 and nimda worm propagation," http://www.renesys.com/projects/bgp_instability/, 2001.

- [38] M. Lad, X. Zhao, B. Zhang, D. Massey, and L. Zhang, "Analysis of BGP update surge during Slammer worm attack," in *Proc. of the International Workshop on Distributed Computing (IWDC)*, 2003.
- [39] A. Baumann and B. Fabian, "How robust is the Internet? – insights from graph analysis," in *Proc. of CRIStIS: Risks and Security of Internet and Systems*, 2015.
- [40] S. Carmi, S. Havlin, S. Kirkpatrick, Y. Shavitt, and E. Shir, "A model of internet topology using k-shell decomposition," *Proceedings of the National Academy of Sciences*, vol. 104, no. 27, June 2007.
- [41] R. Cohen, K. Erez, D. Ben-Avraham, and S. Havlin, "Resilience of the internet to random breakdowns," *Physical review letters*, vol. 85, no. 21, p. 4626, 2000.
- [42] N. Berger, C. Borgs, J. Chayes, and A. Saberi, "On the spread of viruses on the internet," in *Proceedings of the 16th ACM-SIAM Symposium on Discrete Algorithm (SODA)*, 2005, pp. 301–310.
- [43] A.-L. Barabási and R. Albert, "Emergence of scaling in random networks," *Science*, vol. 286, no. 5439, pp. 509–512, 1999.
- [44] W. Ren, J. Wu, X. Zhang, R. Lai, and L. Chen, "A stochastic model of cascading failure dynamics in communication networks," *IEEE Transactions on Circuits and Systems II: Express Briefs*, vol. 65, no. 5, pp. 632–636, 2018.
- [45] W. Willinger, D. Alderson, and J. C. Doyle, "Mathematics and the Internet: A source of enormous confusion and great potential," *Notice of the AMS*, vol. 56, no. 5, pp. 586–599, May 2009.
- [46] R. J. Ellison, D. A. Fisher, R. C. Linger, H. F. Lipson, and T. Longstaff, "Survivable network systems: An emerging discipline," Carnegie-mellon Univ Pittsburgh PA Software Engineering Inst, Tech. Rep., 1997.
- [47] D. Dolev, S. Jamin, O. O. Mokryn, and Y. Shavitt, "Internet resiliency to attacks and failures under bgp policy routing," *Computer Networks*, vol. 50, no. 16, pp. 3183–3196, 2006.
- [48] B. Wang, H. Tang, C. Guo, and Z. Xiu, "Entropy optimization of scale-free networks' robustness to random failures," *Physica A: Statistical Mechanics and its Applications*, vol. 363, no. 2, pp. 591–596, 2006.
- [49] E. G. Coffman Jr, Z. Ge, V. Misra, and D. Towsley, "Network resilience: exploring cascading failures within bgp," in *Proc. 40th Annual Allerton Conference on Communications, Computing and Control*, 2002.
- [50] Y. Liu, W. Peng, J. Su, and Z. Wang, "Assessing the impact of cascading failures on the interdomain routing system of the internet," *New Generation Computing*, vol. 32, pp. 237–255, 2014.
- [51] J. Wang, C. Jiang, and J. Qian, "Robustness of internet under targeted attack: a cascading failure perspective," *Journal of Network and Computer Applications*, vol. 40, pp. 97–104, 2014.
- [52] X. Wu, R. Gu, Y. Ji, and H. E. Stanley, "Dynamic behavior analysis of an internet flow interaction model under cascading failures," *Physical Review E*, vol. 100, no. 2, p. 022309, 2019.
- [53] J. Wang, Y.-h. Liu, J.-q. Zhu, and Y. Jiao, "Model for cascading failures in congested internet," *Journal of Zhejiang University-SCIENCE A*, vol. 9, no. 10, pp. 1331–1335, 2008.
- [54] K. M. Lhakmana, Y. Murakami, and T. Ishida, "Analysis of large-scale service network tolerance to cascading failure," *IEEE Internet of Things Journal*, vol. 3, no. 6, pp. 1159–1170, 2016.
- [55] S. V. Buldyrev, R. Parshani, G. Paul, H. E. Stanley, and S. Havlin, "Catastrophic cascade of failures in interdependent networks," *Nature*, vol. 464, no. 7291, pp. 1025–1028, 2010.
- [56] B. Bassiri and S. S. Heydari, "Network survivability in large-scale regional failure scenarios," in *Proceedings of the 2nd Canadian Conference on Computer Science and Software Engineering*, 2009, pp. 83–87.
- [57] W. Peng, Z. Li, J. Su, and M. Dong, "Evaluation of topological vulnerability of the internet under regional failures," in *Availability, Reliability and Security for Business, Enterprise and Health Information Systems: IFIP WG 8.4/8.9 International Cross Domain Conference and Workshop, ARES 2011, Vienna, Austria, August 22-26, 2011. Proceedings 6*. Springer, 2011, pp. 164–175.
- [58] T. Leighton, "Improving performance on the internet," *Communications of the ACM*, vol. 52, no. 2, pp. 44–51, 2009.
- [59] C. Hu, K. Chen, Y. Chen, and B. Liu, "Evaluating potential routing diversity for internet failure recovery," in *2010 Proceedings IEEE INFOCOM*. IEEE, 2010, pp. 1–5.
- [60] C. Hu, K. Chen, Y. Chen, B. Liu, and A. V. Vasilakos, "A measurement study on potential inter-domain routing diversity," *IEEE Transactions on Network and Service Management*, vol. 9, no. 3, pp. 268–278, 2012.
- [61] W. Deng, P. Zhu, N. Xiong, Y. Xiao, and X. Hu, "How resilient are individual ases against as-level link failures?" in *2011 IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS)*. IEEE, 2011, pp. 982–987.
- [62] C. UCSD/SDSC, "The caida as relationships dataset," <https://www.caida.org/catalog/datasets/as-relationships/>, 2022.
- [63] —, "The CAIDA UCSD, AS relationships – with geographic annotations," <https://publicdata.caida.org/datasets/as-relationships-geo>, 2022.
- [64] —, "The caida ucsd internet topology data kit (itdk)," <https://www.caida.org/catalog/datasets/internet-topology-data-kit>, 2022.
- [65] Y. Chen, P. Embrechts, and R. Wang, "An unexpected stochastic dominance: Pareto distributions, catastrophes, and risk exchange," *arXiv preprint arXiv:2208.08471*, 2022.
- [66] P. Embrechts, S. I. Resnick, and G. Samorodnitsky, "Extreme value theory as a risk management tool," *North American Actuarial Journal*, vol. 3, no. 2, pp. 30–41, 1999.
- [67] P. Gill, M. Schapira, and S. Goldberg, "Modeling on quicksand: Dealing with the scarcity of ground truth in interdomain routing data," *ACM SIGCOMM Computer Communication Review*, vol. 42, no. 1, pp. 40–46, 2012.
- [68] J. Furuness, C. Morris, R. Morillo, A. Herzberg, and B. Wang, "Bgpy: The bgp python security simulator," in *Proceedings of the 16th Cyber Security Experimentation and Test Workshop*, 2023, pp. 41–56.
- [69] J. Nykvist and L. Carr-Motyko, "Simulating convergence properties of bgp," in *Proceedings. Eleventh International Conference on Computer Communications and Networks*. IEEE, 2002, pp. 124–129.
- [70] B. K. Szymanski, Y. Liu, and R. Gupta, "Parallel network simulation under distributed genesis," in *Seventeenth Workshop on Parallel and Distributed Simulation, 2003.(PADS 2003). Proceedings*. IEEE, 2003, pp. 61–68.
- [71] J. Pan, S. Paul, and R. Jain, "A survey of the research on future internet architectures," *IEEE Communications Magazine*, vol. 49, no. 7, pp. 26–36, 2011.
- [72] M. Scazzariello, L. Ariemma, G. Di Battista, and M. Patrignani, "Megalos: A scalable architecture for the virtualization of large network scenarios," *Future internet*, vol. 13, no. 9, p. 227, 2021.
- [73] J. L. Sobrinho, L. Vanbever, F. Le, and J. Rexford, "Distributed route aggregation on the global network," in *Proceedings of the 10th ACM International on Conference on emerging Networking Experiments and Technologies*, 2014, pp. 161–172.
- [74] Y. Rekhter and T. Li, "Rfc1771: A border gateway protocol 4 (bgp-4)," 1995.
- [75] R. A. Barabási, H. Jeong, and Albert-László, "Error and attack tolerance of complex networks," *Nature*, no. 406, pp. 378–382, 2000.
- [76] J. C. Doyle, D. L. Alderson, L. Li, S. Low, M. Roughan, S. Shalunov, R. Tanaka, and W. Willinger, "The 'robust yet fragile' nature of the internet," *Proceedings of the National Academy of Sciences*, vol. 102, no. 41, pp. 14497–14502, 2005.
- [77] D. L. Alderson, J. C. Doyle, and W. Willinger, "Lessons from" a first-principles approach to understanding the internet's router-level topology," *ACM SIGCOMM Computer Communication Review*, vol. 49, no. 5, pp. 96–103, 2019.
- [78] J. Rak and D. Hutchison, *Guide to disaster-resilient communication networks*. Springer Nature, 2020.
- [79] S. Neumayer and E. Modiano, "Network reliability under geographically correlated line and disk failure models," *Computer Networks*, vol. 94, pp. 14–28, 2016.
- [80] E. Ayanoglu, I. Chih-Lin, R. D. Gitlin, and J. E. Mazo, "Diversity coding: Using error control for self-healing in communication networks," in *Proceedings. IEEE INFOCOM'90: Ninth Annual Joint Conference of the IEEE Computer and Communications Societies@ m_The Multiple Facets of Integration*. IEEE, 1990, pp. 95–104.
- [81] J. Strand, A. L. Chiu, and R. Tkach, "Issues for routing in the optical layer," *IEEE Communications magazine*, vol. 39, no. 2, pp. 81–87, 2001.
- [82] A. Miu, H. Balakrishnan, and C. E. Koksal, "Improving loss resilience with multi-radio diversity in wireless networks," in *Proceedings of the 11th annual international conference on Mobile computing and networking*, 2005, pp. 16–30.
- [83] E. E. Moghaddam, H. Beyranvand, and J. A. Salehi, "Crosstalk-aware resource allocation in survivable space-division-multiplexed elastic optical networks supporting hybrid dedicated and shared path protection," *Journal of Lightwave Technology*, vol. 38, no. 6, pp. 1095–1102, 2019.

For each neighboring pair of ASes in ITDK, we determine the probability of a link between these ASes failing by calculating the ratio of the failed links to all their redundant links. For relationships in the main topology (e.g., CAIDA as-rel) that do not exist in ITDK, we adopt the following approach: For each AS, we find the ratio of border routers outside the failing region to all the border routers and consider it as an estimation of the probability that a border router of the AS survives the failure (denoted as p_i^s). Then for a link between two ASes, the probability that the link survives is the probability that both the router at the ends of the link survive. Considering the location of the end routers is independent, then the probability of a link between AS A and AS B surviving the failure is the multiplication of the estimated probability of each end router surviving the failure:

$$\begin{aligned} &\text{Probability}((i, j) \text{ fails} | (i, j) \text{ geospatially uncorrelated}) \\ &= 1 - p_i^s \times p_j^s \end{aligned} \quad (1)$$

- [84] J. P. Sterbenz, D. Hutchison, E. K. Çetinkaya, A. Jabbar, J. P. Rohrer, M. Schöller, and P. Smith, "Resilience and survivability in communication networks: Strategies, principles, and survey of disciplines," *Computer networks*, vol. 54, no. 8, pp. 1245–1265, 2010.
- [85] K.-W. Kwong, L. Gao, R. Guérin, and Z.-L. Zhang, "On the feasibility and efficacy of protection routing in ip networks," *IEEE/ACM Transactions on Networking*, vol. 19, no. 5, pp. 1543–1556, 2011.
- [86] B. Yang, J. Liu, S. Shenker, J. Li, and K. Zheng, "Keep forwarding: Towards k-link failure resilient routing," in *IEEE INFOCOM 2014-IEEE Conference on Computer Communications*. IEEE, 2014, pp. 1617–1625.
- [87] F. Palmieri, U. Fiore, A. Castiglione, F.-Y. Leu, and A. De Santis, "Analyzing the internet stability in presence of disasters," in *Security Engineering and Intelligence Informatics: CD-ARES 2013 Workshops: MoCrySEn and SeCIHD, Regensburg, Germany, September 2-6, 2013. Proceedings 8*. Springer, 2013, pp. 253–268.
- [88] T. Gomes, J. Tapolcai, C. Esposito, D. Hutchison, F. Kuipers, J. Rak, A. De Sousa, A. Iossifides, R. Travanca, J. André *et al.*, "A survey of strategies for communication networks to protect against large-scale natural disasters," pp. 11–22, 2016.
- [89] J. Li, Z. Wu, and E. Purpus, "Cam04-5: Toward understanding the behavior of bgp during large-scale power outages," in *IEEE Globecom 2006*. IEEE, 2006, pp. 1–5.
- [90] J. Luo, J. Xie, R. Hao, and X. Li, "An approach to accelerate convergence for path vector protocol," in *Global Telecommunications Conference, 2002. GLOBECOM'02. IEEE*, vol. 3. IEEE, 2002, pp. 2390–2394.
- [91] N. Gvozdiev, B. Karp, and M. Handley, "Loup: The principles and practice of intra-domain route dissemination," in *Proceedings of the 10th USENIX Conference on Networked Systems Design and Implementation*, ser. nsdi'13. USA: USENIX Association, 2013, p. 413–426.
- [92] L. Subramanian, M. Caesar, C. T. Ee, M. Handley, M. Mao, S. Shenker, and I. Stoica, "HLP: A next generation inter-domain routing protocol," *ACM SIGCOMM Computer Communication Review*, vol. 35, no. 4, pp. 13–24, 2005.
- [93] D. Walton, "Advertisement of Multiple Paths in BGP (first draft)," Nov. 2002, library Catalog: tools.ietf.org. [Online]. Available: <https://tools.ietf.org/html/draft-walton-bgp-add-paths-00>
- [94] V. Van den Schrieck, P. Francois, and O. Bonaventure, "BGP Add-Paths: The Scaling/Performance Tradeoffs," *IEEE Journal on Selected Areas in Communications*, vol. 28, no. 8, pp. 1299–1307, Oct. 2010, conference Name: IEEE Journal on Selected Areas in Communications.
- [95] A. Retana, "Advertisement of multiple paths in BGP: Implementation report," Working Draft, IETF Secretariat, Internet-Draft draft-ietf-idr-add-paths-implementation-00, February 2015. [Online]. Available: <http://www.ietf.org/internet-drafts/draft-ietf-idr-add-paths-implementation-00.txt>
- [96] V. Van den Schrieck, "Analysis of paths selection modes for Add-Paths," Jul. 2009. [Online]. Available: <https://tools.ietf.org/id/draft-vvds-add-paths-analysis-00.html>
- [97] N. Kushman, S. Kandula, D. Katabi, and B. M. Maggs, "R-BGP: staying connected in a connected world," in *4th Symposium on Networked Systems Design and Implementation (NSDI 2007), April 11-13, 2007, Cambridge, Massachusetts, USA, Proceedings*, H. Balakrishnan and P. Druschel, Eds. USENIX, 2007. [Online]. Available: <http://www.usenix.org/events/nsdi07/tech/kushman.html>
- [98] I. Van Beijnum, J. Crowcroft, F. Valera, and M. Bagnulo, "Loop-freeness in multipath bgp through propagating the longest path," in *2009 IEEE International Conference on Communications Workshops*. IEEE, 2009, pp. 1–6.
- [99] T. Holterbach, S. Vissicchio, A. Dainotti, and L. Vanbever, "Swift: Predictive fast reroute," in *Proceedings of the Conference of the ACM Special Interest Group on Data Communication*. ACM, 2017, pp. 460–473.
- [100] T. Holterbach, E. C. Molero, M. Apostolaki, A. Dainotti, S. Vissicchio, and L. Vanbever, "Blink: Fast connectivity recovery entirely in the data plane," in *nsdi*, Feb. 2019.
- [101] X. Zhang, H.-C. Hsiao, G. Hasker, H. Chan, A. Perrig, and D. G. Andersen, "Scion: Scalability, control, and isolation on next-generation networks," in *2011 IEEE Symposium on Security and Privacy*. IEEE, 2011, pp. 212–227.
- [102] X. Yang, "Nira: A new internet routing architecture," *ACM SIGCOMM Computer Communication Review*, vol. 33, no. 4, pp. 301–312, 2003.