

# Modular Control Plane Verification via Temporal Invariants

Timothy Alberdingk Thijm  
*Princeton University*

Ryan Beckett  
*Microsoft Research*

Aarti Gupta  
*Princeton University*

David Walker  
*Princeton University*

## Abstract

Satisfiability Modulo Theory (SMT)-based tools for network control plane analysis make it possible to reason exhaustively about interactions with peer networks and to detect vulnerabilities such as accidental use of a network as transit or prefix hijacking. SMT-based reasoning also facilitates synthesis and repair. To scale SMT-based verification to large networks, we introduce TIMEPIECE, a new modular control plane verification system. While past verifiers like Minesweeper [8] were based on analysis of stable paths, we show that such models, when deployed naïvely in service of modular verification, are unsound. To rectify the situation, we adopt a routing model based around a logical notion of time and develop a sound, expressive, and scalable verification engine.

Our system requires that a user specifies interfaces between module components. We develop methods for defining these interfaces using predicates inspired by temporal logic, and show how to use those interfaces to verify a range of network-wide properties such as reachability, “no transit,” and “no hijacking.” Verifying a prefix-filtering policy using a non-modular verification engine times out on a 320-node fattree network after 4 hours. However, TIMEPIECE verifies a 4,500-node fattree in 6.5 minutes on a 96-core virtual machine. Modular verification of individual routers is embarrassingly parallel and completes in seconds, which allows verification to scale beyond non-modular engines, while still allowing the full power of SMT-based symbolic reasoning.

## 1 Introduction

Configuring the routing protocols of large-scale networks, such as the data centers and wide-area networks of major cloud providers, is an immensely difficult task. Individual routers may have tens or hundreds of thousands of lines of configuration and these networks consist of hundreds or thousands of routers [30]. Moreover, these configurations are not always stable. Updates occur non-stop and are often a cause of unintentional errors [44, 45, 54].

Over the past several years, advances in verification of control plane routing protocols have been impressive, with a wide variety of algorithms, data structures, and optimizations deployed to make verification more scalable [2, 8–10, 17, 21, 35, 43, 52, 53]. One branch of Satisfiability Modulo Theories (SMT)-based tools, beginning with Minesweeper [8], operate by describing the *stable states* of a routing protocol as a constraint system for an SMT solver [7] to solve. Doing so

avoids having to reason about the many transient states that occur leading up to those stable states. These SMT-based tools can reason *symbolically* about network environments, and hence can consider, for example, all possible routing announcements that an external neighbor may generate, and whether those external announcements would interfere with internal routing. SMT-based verification methods also provide a foundation for systems for configuration synthesis or repair, such as ConfigAssure [39], NetComplete [16], or AED [1].

Still, these SMT-based tools do not scale well, with analysis topping out at hundreds of nodes, and sometimes much less when policy is more complex [2, 8, 43]. To handle large-scale networks such as industrial data centers and WANs containing thousands of devices, the community turned to simulation-based tools (sometimes extended with limited symbolic reasoning), such as FastPlane [35], ShapeShifter [10], Plankton [43], and Hoyan [53] for the job. These systems can often scale to thousands of nodes, but they do so by executing protocols *concretely* rather than symbolically. These tools illustrated it was possible to analyze large networks, and even consider fault tolerance properties. However, it is not clear whether (and how) these extended simulation-based tools could reason about all possible external BGP announcements a neighbor might send (*e.g.*, what BGP communities or AS path [14] their routes may have) that may interfere with local routing policy—an enormous source of concern and much network downtime [24]. Hence the question remains: how to achieve high-performance, scalable verification methods with full symbolic reasoning about unknown routes or policy.

One path forward is to break a network up into smaller components and to tackle verification of each component in isolation. This approach has proven successful in the context of software verification [6, 19, 22, 27, 29] and in network data plane verification [30]. In such a *modular* verification system, one might declare an *interface* for each component that describes the set of routing messages the component can generate in a stable state. Then, for each component  $c$ , one might check that, if all neighbors of  $c$  adhere to their interfaces, then  $c$  does as well. If one verified every component  $c$  in this way, and the interfaces imply some useful property, such as reachability or access control, then one might conclude the system as a whole also satisfies that property. If each network component is small enough, such an SMT-based verification method will be fast and general, and give access to powerful symbolic reasoning facilities.

**A temporal model.** Unfortunately, a naïve combination of modular reasoning and Minesweeper-style analysis of stable states *is unsound*. To our surprise, the best way to recover soundness, while maintaining the system’s generality, is to drop the idea of analyzing stable states alone and move to a *temporal model of network execution* that reasons over messages produced at all times, not just once the system is stable.

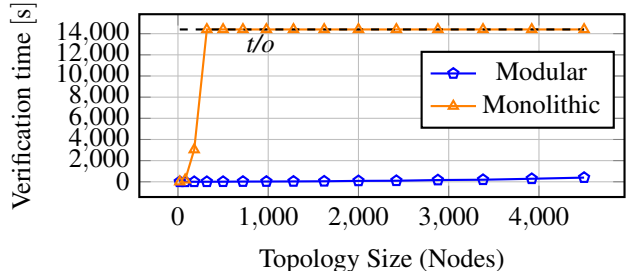
This temporal model appears to ask the verification engine to do a lot more work: the system must verify that all the messages produced *at all times* are consistent with a user-supplied interface for each network component. Nevertheless, because reasoning is modular, ensuring individual problems are small, the system scales with the size of the largest *component* rather than the size of the network. This modular reasoning is general and any symbolic method (*e.g.*, symbolic simulation, model checking) could use it to verify individual components. We use an SMT-based method in this work. As a preview of modularity’s benefits, Figure 1 shows the time it takes our modular verification system TIMEPIECE to verify connectivity for variable-sized fattree topologies with external WAN route announcements using the eBGP routing protocol, compared with a Minesweeper-style network-wide stable paths encoding.

TIMEPIECE does require more work of users than monolithic, non-modular systems: users must supply interfaces that characterize the routes each network component may generate at each time. Still, the presence of these interfaces, once constructed, will have the typical benefits of interfaces in any software engineering context. First, they localize exactly where an error occurs: if a component is not consistent with its interface, then one must only search that component for the mistake, and a counterexample from the SMT solver can help pinpoint it. Second, router configurations change rapidly, and these changes are often the source of network-wide problems [54]. Well-defined interfaces will be stable over time. As users update their configurations, they may easily recheck them against the stable, local interface for problems.

Inspired by temporal logic [42], we developed a simple language to help users specify their interfaces. Through this language, users may state that they expect to see certain sets of messages *always*, *eventually* (by some specified time  $t$ , to be more precise), or *until* (some approximate specified time). Moreover, the interface language allows users to generate abstract specifications that need not characterize irrelevant features of protocol messages, and instead provide only what is necessary to prove a desired safety property. For instance, a user might specify a reachability property simply by stating a node must “eventually receive some route,” without saying which route it must receive.

To summarize, the key contributions of this paper are:

- We demonstrate in depth why a natural, but naïve modular control plane analysis based on an analysis of stable states is unsound.



**Figure 1: Verification time for connectivity in a data center with a prefix filtering policy to the wide-area network.**

- We develop a new theory for modular control plane analysis based on time, and prove it sound and complete with respect to the network semantics. This theory is general, and can verify individual components using other verification methods [10, 43].
- We design and implement a new, modular control plane verification tool, TIMEPIECE, based directly on this theory, which uses an SMT-based backend to reason symbolically about all possible routes at all times.
- We evaluate the tool and check a variety of policies at individual nodes in hundreds of milliseconds. Thanks to its embarrassingly parallel modular procedure, TIMEPIECE scales to networks with thousands of nodes.

This work raises no ethical concerns.

## 2 Key Ideas

This section introduces the stable routing model of network control planes, which serves as a foundation for many past network verification tools [8–10, 20]. It illustrates in depth why naïve adoption of this model for modular verification is unsound. It then introduces a new temporal model for control plane verification and provides the intuition for why the revised model is superior. This section is long but contains a substantial payoff: the essence of why a sound and general modular control plane analysis should be based off a temporal model of control plane behavior.

### 2.1 Background

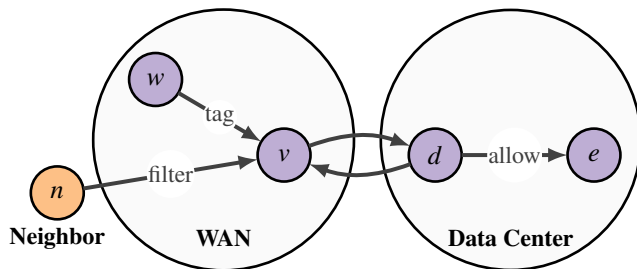
To determine how to deliver traffic between two endpoints, routers (also called nodes) run distributed routing protocols such as BGP [36], OSPF [38], RIP [28], or ISIS [40]. Each node participating in a protocol receives *messages* (also called *routes*) from its neighboring nodes. After receiving routes from its neighbors, a node will select its “best” route—the route it will use to forward traffic. Different protocols use different metrics to compare routes and select the best among those received. For instance, RIP compares hop count; OSPF

uses the shortest weighted-length path; and BGP uses a complex, user-configurable combination of metrics. Finally, each router sends its chosen route to its neighbors, possibly modifying the route along the way (for instance, by prepending its identifier to the path represented by the route).

**Routing algebras.** Routing algebras [26, 50] are abstract models that capture the similarities between different distributed routing protocols. Prior work on control plane verification [8, 23, 25] uses similar abstract models to formalize route computation. We adopt this standard abstract model of routing protocols, which specifies the following components.

- A directed graph  $G$  that defines the network topology, including its nodes ( $V$ ) and edges ( $E$ ). We typically use lower case letters ( $u, v, w, \text{etc.}$ ) for nodes and pairs ( $uv$ ) to indicate directed edges between those nodes.
- A set  $S$  of routes that communicate routing information between nodes.
- A initialization function  $I_v \in S$  that provides an initial route for each router  $v$  in the network.
- A set of edge transfer functions  $F$ . Each transfer function  $f_e \in F$  transforms routes as they traverse the edge  $e$ .
- A binary function  $\oplus$  (a.k.a. merge or the selection function) selects the best route between two options.

**An idealized example.** Many large cloud providers deploy data center networks to scale up their compute capacity. They connect those data centers to each other and the rest of the internet via a wide-area network (WAN). To illustrate the challenges of modular network verification, we will explore verification of an idealized cloud provider network with WAN and data center components. The picture below presents a highly abstracted view of the topology of our network.



The data center network contains routers  $d$  and  $e$  where  $d$  connects to the corporate WAN and  $e$  connects to data center servers. The WAN consists of routers  $w$  and  $v$ . Router  $v$  connects to the data center as well as to a neighboring network  $n$ , which is not controlled by our cloud provider.<sup>1</sup>

<sup>1</sup>Any of the edges could be bi-directional, allowing routes to pass in both directions, but for pedagogic reasons we strip down the example to the barest minimum, retaining edges that flow from left-to-right except for at  $v$  and  $d$  where routes may flow back and forth.

The default routing policy uses shortest-paths. However, in addition, the network administrators want  $e$  to be reachable from all cloud-provider-owned devices (i.e.,  $w, v, d$ ), but not to be reachable from outsiders (i.e.,  $n$ ). They intend to enforce this property by tagging all routing messages originating from their network ( $w$ ) as “internal” (e.g., using BGP community tags [14]) and allowing those routes to traverse the  $de$  edge. Doing so should allow  $e$  to communicate with internal machines but not external machines. Furthermore, to protect routes from outside interference, the cloud provider applies route filters to external peers to drop erroneous advertised routes that may “hijack” [18] internal routes.

**Modelling the example.** To model our example network, we define the network topology as the graph  $G$  pictured earlier. We assume all routers participate in an idealized variant of eBGP [36], which is commonly used in both wide-area networks and data centers [3]. The set of routes  $S$  used in this protocol are records with 3 fields: (i) an integer “local preference” that lets users overwrite default preferences, (ii) an integer path length, and (iii) a boolean tag field that is set to true if a route comes from an internal source and false otherwise. Lastly,  $S$  includes  $\infty$ , a message that indicates absence of a route.

Let’s consider what happens when starting with a specific route at WAN node  $w$ ,  $\langle 100, 0, \text{false} \rangle$  (local preference 100, path length of 0, not tagged “internal”). The  $I$  function would assign  $w$  that route, and assign the  $\infty$  route to all other nodes.

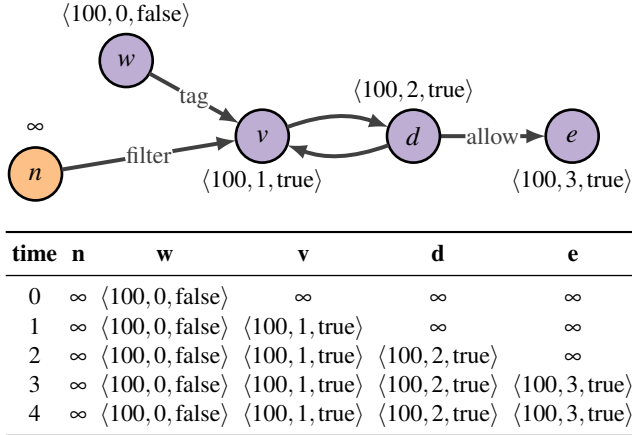
The transfer function  $f_e$  will increment the length field of every route by one across every edge  $e$ . In addition, edge  $wv$  sets the internal tag field to true and edge  $nv$  drops all routes (transforms them into  $\infty$ ). Finally, edge  $de$  drops all routes not tagged internal/true.

The merge function  $\oplus$  always prefers some route over the  $\infty$  route. In addition,  $\oplus$  prefers routes with higher local preference over lower local preference. If the local preference is the same, it chooses a route with a shorter path length.  $\oplus$  ignores the tag field. For example, merge operates as follows:

$$\begin{aligned} \langle 100, 2, \text{false} \rangle \oplus \infty &= \langle 100, 2, \text{false} \rangle \\ \langle 100, 2, \text{false} \rangle \oplus \langle 200, 5, \text{true} \rangle &= \langle 200, 5, \text{true} \rangle \\ \langle 200, 2, \text{false} \rangle \oplus \langle 200, 5, \text{true} \rangle &= \langle 200, 2, \text{false} \rangle \end{aligned}$$

**Network simulation.** A *state* of a network is a mapping from nodes to the “best routes” they have computed so far. One may carry out a simulation by starting in the initial state and repeatedly computing new states (i.e., new “best routes” for particular nodes). Eventually, well-behaved networks converge to *stable states* where no node can compute any better routes, given the routes provided by its neighbors.

To compute a new best route at a particular node, say  $v$ , we apply the  $f$  function to each best route computed so far at its neighbors  $w, n$ , and  $d$ , and then select the best route among the results and the initial value at  $v$ , using the merge



**Figure 2: Simulation of the example network for a fixed set of initial routes. Node  $e$  learns a route to  $w$  since the route is tagged as “internal”, and the network stabilizes at time 3.**

$\oplus$  function. More precisely:

$$v_{new} = f_{wv}(w_{old}) \oplus f_{nv}(n_{old}) \oplus f_{dv}(d_{old}) \oplus I_v$$

The table in Figure 2 presents an example simulation. At each time step, all nodes compute their best route using the equation above, given the routes supplied by their neighbors at the previous time step. After time step 3, no node ever computes any new route—the system has reached a *stable state*. The picture in Figure 2 annotates each node in the diagram with the stable route it computes.<sup>2</sup>

**Network verification.** Since the edge from  $d$  to  $e$  only allows routes tagged internal,  $w$ ’s route would not reach  $e$  if  $v$  were to receive a better route from  $n$  (e.g., if the route filter from  $n$  was implemented incorrectly). In other words, the simulation demonstrates that the network correctly operates when  $n$  sends no route ( $\infty$ ). But what about other routes? Will all routes from  $n$  get filtered correctly? SMT-based tools like Minesweeper [8] and Bagpipe [52] can answer such questions by translating the routing problem into a set of constraints for a Satisfiability Modulo Theory (SMT) solver to solve. In doing so, one may represent the set of all possible external route announcements symbolically and reason simultaneously about any external route announcement (something other verifiers such as Tiramisu [2], Plankton [43], Shapeshifter [10], Hoyan [53], or DNA [54] either do not do, or do partially). On the other hand, while the SMT-based approach is effective for small networks, it does not scale well in practice [2, 8, 43, 52].

<sup>2</sup>For simplicity, our model assumes a synchronous time semantics to simplify our examples and theory — however, we believe we could extend our approach to asynchronous time in the style of [15].

## 2.2 The Challenge of Modular Verification

A system for modular verification will partition a network into components and verify each component separately, possibly in parallel. However, since routes computed at a node in one component depend on the routes sent by nodes in neighboring components, each component must make some assumptions about the routes produced by its neighbors.

**Interfaces.** In our case, for simplicity (though this is not necessary), we place every node in its own component and define for it an *interface* that attempts to *overapproximate* the set of routes that a node might produce in a stable state. The interface for the network as a whole is a function  $A$  from nodes to sets of routes where  $A(x)$  is the interface for node  $x$ .

The person attempting to verify the network will supply these interfaces. Of course, interfaces may be *wrong*—that is, they might not include some route computed by a simulation (and hence might not be a proper overapproximation). Indeed, when there are bugs in the network, the interfaces supplied by the user are likely to be wrong! The user *expects* the network to behave one way, producing a certain set of routes, but, in fact, the network behaves differently due to an error in its configuration. A sound modular verification procedure must detect such errors. On the other hand, a useful modular verification procedure should allow interfaces to overapproximate the routes produced, when users find it convenient. Overapproximations are sound for verifying safety properties, and often simplify reasoning, allowing users to think more abstractly.

In our examples (and our tool), we often use predicates  $\phi$  to define interfaces, where  $\phi$  stands in for the set of routes  $\{s \mid s \in \mathcal{S}, \phi(s)\}$ . Returning to our running example, one might define the interface for neighbor  $w$  using the predicate  $s.lp = 100 \wedge s.len = 0 \wedge \neg s.tag$ . Such an interface would include exactly the one route generated by  $w$  in our example:  $\langle 100, 0, false \rangle$ . However, path length is unimportant in the current context; to avoid having to think about it, a user could instead provide a weaker interface representing infinitely many possible routes, such as  $s.lp = 100 \wedge \neg s.tag$ . This interface relieves the user of having to figure out the exact path length (not so hard in this simple example, but potentially challenging in an arbitrary wide-area network), and instead specifies only the local preference and the tag. In general, admitting overapproximations make it possible for users to ignore any features of a routing algorithm that are not actually relevant for analyzing the properties of interest.

**The Strawperson Verification Procedure.** For a given node  $x$ , the *component centered at  $x$*  is the subgraph of the network that includes node  $x$  and all edges that end at  $x$ . Given a network interface  $A$ , our strawperson verification procedure (SV) will consider the component centered at each node  $x$  independently. Suppose a node  $x$  has neighbors  $n_1, \dots, n_k$ .



For that node  $x$ ,  $SV$  checks that  $\forall s_1 \in A(n_1), \dots, \forall s_k \in A(n_k)$ ,

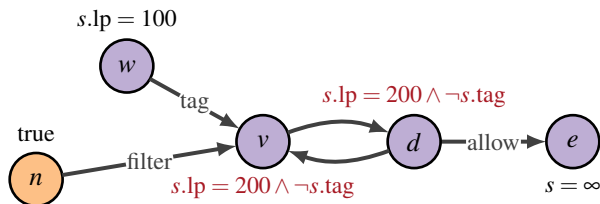
$$f_{n_1 x}(s_1) \oplus \dots \oplus f_{n_k x}(s_k) \oplus I_x \in A(x) \quad (1)$$

This check is akin to performing one local step of simulation, checking that all possible inputs from neighbors give rise to an output route that conforms to the interface. One might *hope* that by performing such a check on *all* components independently, one would be able to guarantee that all nodes converge to stable states described by their interfaces. If that were the case, then one could verify properties by:

1. Checking that all components guarantee their interfaces, under the assumption their neighbors do as well; and
2. Checking that the interfaces imply the network property of interest (*e.g.*, reachability, access control, no transit).

**The problem: Execution interference.** It turns out this simple and natural verification procedure is unsound: users can supply interfaces that, when analyzed in isolation, satisfy equation (1) above, but wind up excluding the stable states computed by simulation. Hence, the second verification step is pointless: a destination that appears reachable according to an interface may not be; conversely, a route that appears blocked may not be.

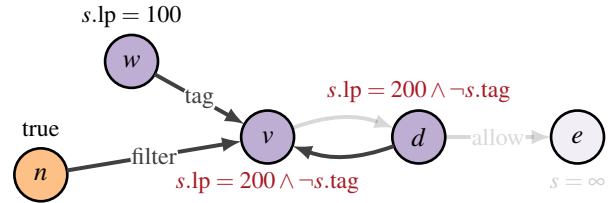
Let us reconsider the running example, where we assign  $w$  an initial route with local preference 100, and let us assume that the external neighbor  $n$  can send us any route (true). A user could provide the interfaces shown below in order to falsely conclude that  $e$  will not receive a route from  $w$ .



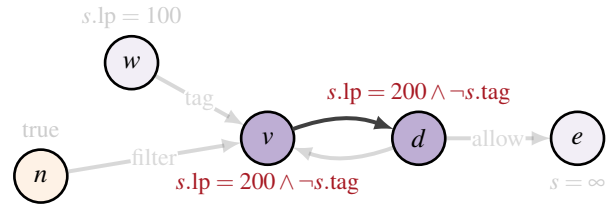
Here, it is easy to check that nodes  $n$  and  $w$  satisfy equation (1). Node  $n$ 's interface is simply any route (and it receives no routes from any other node). Node  $w$ 's route can be any route with a local preference of 100.<sup>3</sup>

The surprise comes at node  $v$  where its interface *only includes* routes that satisfy  $\neg s.tag$ , *i.e.*, routes not tagged as internal. Those routes have  $s.lp = 200$  and may have any path length. But the route from  $w$  is tagged *true* along the edge  $wv$  — why is such a route erroneously excluded from  $v$ 's interface? We show the component centered at  $v$  below. When computing its stable state,  $v$  will compare the routes it receives from  $w$  and  $d$ : because all routes from  $w$  have a local preference set to 100 by  $f_{wv}$ , whereas all routes from  $d$  have a better local preference of 200,  $v$  will always wind up selecting the route from  $d$  over the route from  $w$ .

<sup>3</sup>It could be any route (true), as the edge  $wv$  applies the default preference of 100, but for clarity we label the routes at  $w$  with preference 100.



But how then did  $d$  acquire these preferential routes tagged false? Such routes came in turn from  $v$ 's interface. Here is the component centered at  $d$ .



What has happened is that  $v$  transmits its spurious routes to  $d$ , enabling  $d$  to justify its own spurious routes.  $d$  transmits these back again to  $v$ , where  $d$ 's routes interact with the legitimate routes from  $w$ . Since  $w$ 's routes have lower local preference,  $v$  discards them during computation of stable states. In a nutshell, the routes from an imaginary simulation proposed by our interface overrule legitimate routes from the true simulation, crowding out the true routes at  $v$  and  $d$ . How might we prevent this *execution interference*?

**Other approaches.** There are a few ways to modify the suggested verification procedure to make it sound, but such modifications typically limit the power of the verification procedure or the expressiveness of the properties it can prove.

One approach is to limit every interface to exactly one route. Doing so avoids introducing any imaginary executions in the first place. Kirigami [5] takes this approach, but the cost is that a network engineer analyzing their network must know *exactly* which routes appear at which locations. Computing those routes exactly can be difficult in practice, and would seem unnecessary if all one cares about is a high-level property such as reachability. Moreover, it makes the interfaces brittle in the face of change—any change in network configuration likely necessitates a change in interface. A superior system would allow operators to define *durable* and *abstract* interfaces that imply key properties, and to check router configuration updates against those interfaces.

Another approach is to limit the set of properties that the system can check to only those that say what *does not happen* in the network rather than what does happen. This is the approach Lightyear [51] takes. For instance, Lightyear allows one to check that node A will *not* be able to reach node B, but cannot prove that A and B will have connectivity — a common requirement in networks.

A final approach is to try to impose an ordering on the components, and to verify each component according to this

ordering, using no information from the not-yet-verified components. This approach avoids the circular reasoning between  $d$  and  $e$  above, but is still unnecessarily conservative. The running example is overly simple as it shows routes propagated through a network in a single direction from left to right. In realistic networks, multiple destinations may broadcast routes in multiple directions at once. In such situations, there may be no way to order the components, and verification may not be possible. Moreover, we found that being unable to make assumptions from some neighbors made verifying certain properties, such as reachability, impossible in most cases.

### 2.3 The Solution: A Temporal Model

Our key insight is to change the model: rather than focus exclusively on the final stable states of a system, as a Minesweeper-style verifier would, we ensure that the model preserves the *entirety* of every step-by-step execution. To make this work, we need to add information to the model: a notion of *logical time*. By associating every route with the time at which a node computes it, we can (i) ensure that *all* routes at a particular time are properly considered, and their executions extended a time step, and (ii) ensure that we avoid collisions between routes computed at different times.

To verify such routing systems modularly, we once again must specify interfaces, but this time the interface for each node will specify the set of routes that may appear *at any time*. We write this now as an interface  $A$  that takes both a node *and a time* and returns an overapproximation of the set of routes that may appear at the node *at that time*. For example,  $A(x)(t)$  now gives the set of routes for node  $x$  and time  $t$ . To check the interfaces, we use a verification procedure structured inductively with respect to time, as follows:

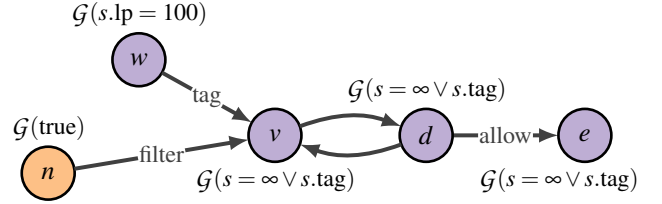
- At every node  $x$ , check  $I_x$  is included in  $A(x)(0)$
- Consider each node  $x$  with neighbors  $n_1, \dots, n_k$ . At time  $t + 1$ , check that merging any combination of routes  $s_1 \in A(n_1)(t), \dots, s_k \in A(n_k)(t)$  from neighbors' interfaces at time  $t$  produces a route in  $A(x)(t + 1)$ :

$$f_{n_1 x}(s_1) \oplus \dots \oplus f_{n_k x}(s_k) \oplus I_x \in A(x)(t + 1) \quad (2)$$

Because this procedure is structured inductively, we can prove, by induction on time, that all states at all times are included in their respective interfaces—the procedure is *sound*.

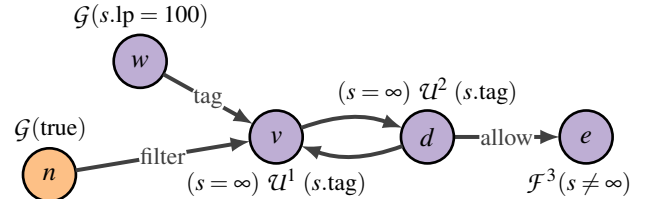
To reason over times, we borrow notation from temporal logic to specify interfaces. We write  $\mathcal{G}(P)$  (“globally  $P$ ”) when a node’s interface includes the routes that satisfy predicate  $P$  for all times  $t$ . We write  $P_1 \mathcal{U}^t P_2$  (“ $P_1$  until  $P_2$ ”) when a node may have routes satisfying  $P_1$  until time  $t - 1$  and  $P_2$  afterwards. Finally, we write  $\mathcal{F}^t(P)$  (“finally  $P$ ”) to mean that eventually at time  $t$  routes start satisfying  $P$ .

**Verifying correct interfaces.** The picture below presents an interface we may verify with this model.



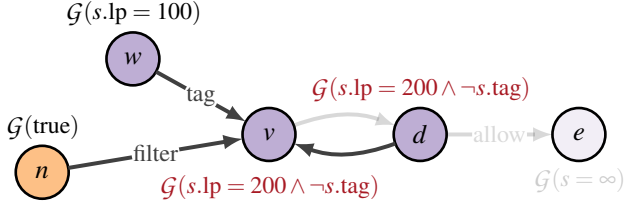
We once again assume that node  $n$  can send any route at any time, denoted by the interface:  $\mathcal{G}(\text{true})$ . Similarly, we assume  $w$  has some route with default local preference:  $\mathcal{G}(s.\text{lp} = 100)$ . The interesting part is at nodes  $v$  and  $d$  where the interfaces state that there is always either no route (e.g., at the beginning of time), or a tagged route:  $\mathcal{G}(s = \infty \vee s.\text{tag})$ . Node  $e$  is now able to prove a weak property: *if it receives a route*, then the route will be tagged internal. Node  $v$  is able to prove its interface since routes are always tagged on import from node  $w$ , routes from  $n$  are correctly dropped, and any routes from  $d$  must also have a tag per its interface. In fact, all the nodes can prove their interface given their neighbors’ interfaces.

**Proving reachability.** The previous interfaces were not strong enough to prove that  $w$  will be able to reach  $e$ . The problem is that we were trying to reason about *all* time, and yet  $e$  will not always have a route to  $w$  (i.e., from time 0 onward). Instead, we know that  $e$  will *eventually* be able to reach  $w$ . Consider now the stronger interfaces shown below:



As before, we allow  $n$  and  $w$  to send any route. However, now nodes  $v$  and  $d$  declare that they will not have a route *until* a specified (logical) time, at which point they receive a tagged route. We give precise witness times for  $v$  and  $d$ ’s interfaces, as otherwise  $v$  could give  $d$  a non-null route (or vice-versa) that would violate the interface before its witness time.  $e$ ’s interface simply requires that  $e$  receives some route at the witness time (allowing arbitrary routes before the witness time). These interfaces are sufficient to prove that  $e$  will eventually receive a route to  $w$ , since  $d$  will eventually have a route tagged as internal, and hence  $e$  will allow it.

**Debugging erroneous interfaces.** Let us revisit the example from before where the user provided unsound interfaces by introducing spurious routes with local preference 200. The equivalent interfaces are now shown below.



Unlike before, the verification procedure detects an error: the interfaces at nodes  $v$  and  $d$  do not include the initial route  $\infty$  at time 0. As a result, the user will receive a counterexample for time  $t = 0$  when verifying  $v$  or  $d$ . Suppose our imaginative user tries to circumvent this issue by also including the initial route in the interfaces for  $v$  and  $d$  with the interface:

$$\mathcal{G}((s.lp = 200 \wedge \neg s.tag) \vee (s = \infty))$$

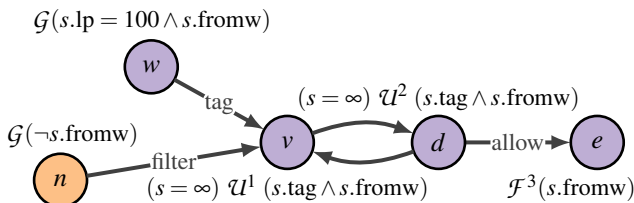
However, doing so merely pushes the problem “one step forward in time”—there is no way to circumvent our temporal analysis. If  $d$ ’s route may be  $\infty$ ,  $v$ ’s interface must also consider what routes it selects when that is the case, including tagged routes such as  $\langle 100, 1, \text{true} \rangle$ . The user might receive a counterexample at time  $t = 1$  where  $v$ ’s route is the following:

$$f_{wv}(\langle 100, 0, \text{false} \rangle) \oplus f_{mv}(\infty) \oplus f_{dv}(\infty) = \langle 100, 1, \text{true} \rangle$$

where  $A(v)(1)$  does not contain the result  $\langle 100, 1, \text{true} \rangle$ . This counterexample should reveal the fact that there is an error in either the specification (as in this case) or the configuration (e.g., if a buggy configuration tagged routes from  $w$  false rather than true as expected).

**Properties and ghost state.** Since modular properties can only reference the route at any single node, an observant reader may wonder if this limits the kinds of properties that we can verify. In the examples so far, we checked to see whether or not node  $e$  receives *some* route, but did not guarantee that the route originated at  $w$ . While at first glance this may seem limiting, there is a simple fix.

Enter ghost state, a technique used in a wide variety of settings (see Dafny [33], for instance). Users may model routes as containing additional “ghost” fields that play no role in a protocol’s routing behavior, yet can capture end-to-end properties. For instance, we could add an additional “ghost” field to determine if a route initiated at node  $w$ —let us call that field “fromw.” We assume this field is initially true at  $w$ , false at all other nodes, and that transfer functions preserve this field. With this addition, we can now check that  $e$  receives a route from  $w$  and no other node:



Property	Added ghost state
reachability to $d$ [20]	1 bit to mark routes from $d$
isolation [8]	1 bit per isolation domain
ordered waypoint [31].	$\log_2(k)$ bits for $k$ waypoints
unordered waypoint [8]	$k$ bits for $k$ waypoints
no-transit [12]	mark with {peer, prov, cust}
fault tolerance [8]	up to $ E $ bits
bounded path length [34]	integer length field

Figure 3: Ghost state for selected example properties.

Ghost state allows us to specify and check many network properties; Figure 3 presents a variety of other possibilities. That said, while ghost state is general and flexible, it can only capture information about the history of a route at a *single* node and is thus not a panacea. For instance, properties involving the routes at more than one node, such as a formulation of local equivalence [8], where  $\sigma(u)(t) = \sigma(v)(t)$  for some arbitrary  $u, v$  and  $t$ , is inexpressible using our verifier.

### 3 Formal Model with Temporal Invariants

Figure 4 provides a summary of the key definitions and notation needed to formalize our verification procedure. Much of the notation follows from the previous section. For instance, a network instance  $N = (G, S, I, F, \oplus)$  contains the key components introduced earlier. To refer to the message computed by a network simulator at node  $v$  at time  $t$ , we use the notation  $\sigma(v)(t)$  (defined as before—see Figure 4). Figure 4 also presents our language of temporal operators with lifted versions of some common set operations.

As before, we use  $A$  to denote network interfaces. A valid interface is an *inductive invariant* [22]. Such interfaces satisfy the initial and inductive conditions specified in Figure 4. Valid interfaces may be used to prove safety properties, again, as specified in Figure 4 (see *safety condition*).

The most important property of our system is *soundness*: the simulation states are included in any interface  $A$  that satisfies the initial and inductive conditions.

**Theorem 3.1 (Soundness).** *Let  $A$  satisfy initial and inductive conditions. Then  $\forall t \in \mathbb{N}, \forall v \in V, \sigma(v)(t) \in A(v)(t)$ .*

*Proof.* By induction on  $t$ . See A.1.  $\square$

Since initial and inductive conditions suffice to prove that simulation states are included within interfaces, it is safe in turn to use interfaces to check safety properties.

**Corollary 3.2 (Safety).** *Let  $A$  satisfy initial and inductive conditions. Let  $P$  satisfy the safety condition with respect to  $A$ . Then  $\forall t \in \mathbb{N}, \forall v \in V, \sigma(v)(t) \in P(v)(t)$ .*

*Proof.* From definitions. See A.2.  $\square$

**Network instances**  $N = (G, S, I, F, \oplus)$

$G = (V, E)$	network topology
$V$	topology nodes
$E \subseteq V \times V$	topology edges
$S$	set of network routes
$s \in S$	a route
$I : V \rightarrow S$	node initialization functions
$I_v \in S$	initial route at node $v$
$F : E \rightarrow (S \rightarrow S)$	edge transfer functions
$f_e : S \rightarrow S$	transfer function for edge $e$
$\oplus : S \times S \rightarrow S$	merge function

**Network semantics**  $\sigma : V \rightarrow \mathbb{N} \rightarrow S$

$\sigma(v)(t) \in S$	state at node $v$ at time $t$
$\text{preds}(v) = \{u \mid u \in V, uv \in E\}$	in-neighbors of $v$

$$\sigma(v)(0) = I_v \quad (3)$$

$$\sigma(v)(t+1) = I_v \oplus \bigoplus_{u \in \text{preds}(v)} f_{uv}(\sigma(u)(t)) \quad (4)$$

### Interfaces and Properties

$A : V \rightarrow \mathbb{N} \rightarrow 2^S$	node interfaces/invariants
$P : V \rightarrow \mathbb{N} \rightarrow 2^S$	safety properties
$\varphi : 2^S$	sets of states

**Temporal operators**  $Q : \mathbb{N} \rightarrow 2^S$

$\mathcal{G}\varphi$	$= \lambda t. \varphi$	globally
$\varphi_1 \mathcal{U}^r \varphi_2$	$= \lambda t. \text{if } t < \tau$ $\text{then } \varphi_1 \text{ else } \varphi_2$	until
$\mathcal{F}^r \varphi$	$= \lambda t. S \mathcal{U}^r \varphi$	finally
$Q_1 \sqcap Q_2$	$= \lambda t. Q_1(t) \cap Q_2(t)$	intersection (lifted)
$Q_1 \sqcup Q_2$	$= \lambda t. Q_1(t) \cup Q_2(t)$	union (lifted)
$\sim Q$	$= \lambda t. S \setminus Q(t)$	negation (lifted)

### Verification Conditions

*Initial condition for A:*

$$\forall v \in V, I_v \in A(v)(0) \quad (5)$$

*Inductive condition for A:*

$$\begin{aligned} &\forall v \in V, u_1, u_2, \dots, u_n \in \text{preds}(v), \\ &\forall t \in \mathbb{N}, \\ &\forall s_1 \in A(u_1)(t), \\ &\forall s_2 \in A(u_2)(t), \\ &\dots, \\ &\forall s_n \in A(u_n)(t), \end{aligned} \quad (6)$$

$$\left( I_v \oplus \bigoplus_{i \in \{1, \dots, n\}} f_{u_i v}(s_i) \right) \in A(v)(t+1)$$

*Safety condition for P with respect to A:*

$$\forall v \in V, \forall t \in \mathbb{N}, A(v)(t) \subseteq P(v)(t) \quad (7)$$

**Figure 4: Summary of our formal routing model, notation, and presentation of interfaces.**

Our verification procedure is also *complete* in the sense that for any network, there exists an interface that characterizes its behavior exactly. One of the consequences of completeness is that our modular verification procedure is powerful enough to prove any safety property that we could prove via simulation.

**Theorem 3.3** (Completeness). *Let  $\sigma$  be the state of the network. Then  $A(v)(t) = \{\sigma(v)(t)\}$  for all  $v \in V$  and all  $t \in \mathbb{N}$  satisfies the initial and inductive conditions.*

*Proof.* By construction of the interface. See A.3.  $\square$

## 4 SMT Algorithms for Verification

Each of the verification conditions (initial, inductive, safety) universally quantify over the (finitely many) nodes in the network. Consequently, a modular verifier can enumerate the network's nodes and independently check each VC for a concrete choice of node. To check an instance of a VC, an off-the-shelf SMT solver will attempt to prove the condition

is valid (*i.e.*, true for all choices of  $t$ ). If the instance is *not valid*, the solver can provide us with a *counterexample*—the state of a node at a particular time that the VC does not hold. Counterexamples to initial or inductive conditions indicate that the interface does not approximate the network's behavior, while a counterexample to the safety condition indicates that the interface is not strong enough to prove the property. The latter case may occur because the property is simply not true (indicating a bug), or alternatively because we must strengthen the given interface in order to prove the property.

To clarify how our modular procedure differs from prior work, we present a Minesweeper-style [8] monolithic checking procedure in Algorithm 1. This algorithm does not refer to time, but instead encodes the *stable states* of the network as a single formula  $\psi$  (as presented in §2.1). Given a safety property  $P$  over the stable states (note that here,  $P$  also ignores time), it checks if  $P$  always holds given these states by calling ISVALID to ask the solver if  $\psi \rightarrow P$  is always true.

We present our modular checking procedure in Algorithm 2. The outer CHECKMOD function iterates over each node of



---

**Algorithm 1** Minesweeper-style checking algorithm.

---

```
proc CHECKMONO(network ( $G, S, I, F, \oplus$ ), property  $P$ )  
   $\psi \leftarrow \text{ENCODESTABLESTATES}(G)$   
  return ISVALID( $\psi \rightarrow P$ )
```

---

the network and encodes the underlying formula (for the current node) of our three verification conditions by calling ENCODEINITCOND, ENCODEINDCOND and ENCODESAFECOND (respectively, (5), (6) and (7)). As in Algorithm 1, we then ask the solver if every encoded formula is valid using ISVALID. If ISVALID returns false for any check, we then can ask for the relevant counterexample model  $m$ .

Importantly, unlike Algorithm 1, by using temporal invariants we are able to separate the task of checking  $P$  on the network instance into three independent verification tasks, with the encoding of initial and inductive conditions roughly proportional in size to the complexity of the policy at the given node (which in turn is related to the in-degree of the node—denser networks that include nodes with higher in-degree are more expensive to check). The encoding of the safety condition is proportional to the size of the formulae describing the interface and safety property (and is generally tiny). In addition to reducing the size of each SMT formula, the factoring of the problem into independent conditions makes it possible to use parallelism to check conditions on nodes simultaneously. We will discuss the performance implications of our procedure further in §6.

---

**Algorithm 2** The modular checking algorithm.

---

```
proc CHECKMOD(network ( $G, S, I, F, \oplus$ ), interface  $A$ ,  
  property  $P$ )  
  for all  $v \in V$  do in parallel  
     $\psi_1 \leftarrow \text{ENCODEINITCOND}(v)$  ▷ (5)  
     $\psi_2 \leftarrow \text{ENCODEINDCOND}(v)$  ▷ (6)  
     $\psi_3 \leftarrow \text{ENCODESAFECOND}(v)$  ▷ (7)  
    if  $\bigvee_{1 \leq i \leq 3} \neg \text{ISVALID}(\psi_i)$  then return false  
  return true
```

---

If we get back a counterexample  $m$ , we can inspect it to understand why our checks failed. This can provide insight into how to strengthen the invariants, or pinpoint a bug in our policy. Anecdotally, this feature was critical to helping us design interface functions for our own experiments, and for discovering bugs in our modelling of network policies.

## 5 Implementation

We implemented our modular verification procedure as part of a library called TIMEPIECE, written in C#. The library allows users to construct models of networks and then to verify those models modularly. Like the modelling framework NV [23], TIMEPIECE allows users to customize their network models

by choosing the kinds of routes (which may involve integers, strings, booleans, bitvectors, records, optional data, lists, or sets) and the way initialization, transfer and merge functions process them. This modelling language makes it easy to add ghost state to routes, as described earlier. It is also possible to declare and use symbolic values in the model. Hence, one may reason about all possible prefixes or more generally about all possible external routing announcements.

For example, to model a network running eBGP, we would adopt many of the modelling choices made in Minesweeper [8]. For instance, we use integers and the SMT theory of integers to model path length. We use bitvectors to model local preference and MED. We use the theory of arrays to model sets of community tags.

Under the hood, TIMEPIECE uses Microsoft’s Zen verification library [11], which in turn serves as an interface to the Z3 SMT solver. Hence, the only practical limits to a user’s network model are those that arise from the features and theories supported by Z3.

TIMEPIECE uses multi-threading to run modular checks in parallel.<sup>4</sup> As each check is independent, the time to set up additional threads is the only overhead for parallelization.

## 6 Evaluation

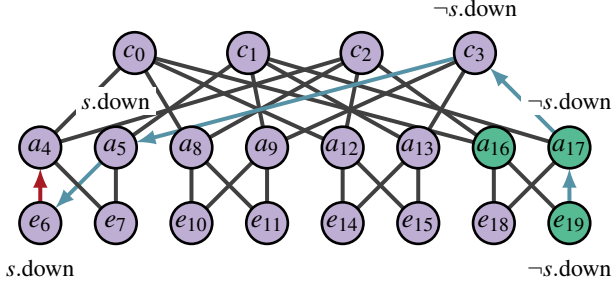
To evaluate TIMEPIECE and illustrate its scaling trends, we generated a series of synthetic fattree [4] data center networks of up to 4,500 nodes and 216,000 edges. We verified several properties and discuss the process we used to craft interfaces capable of proving those properties.

To compare our implementation against a baseline, we implemented a monolithic, network-wide Minesweeper-style procedure Ms from Algorithm 1 and compared its performance against TIMEPIECE. Ms analyzes stable states, which are independent of time. To compare Ms with TIMEPIECE, we first crafted properties for TIMEPIECE, which employs timed invariants. We erased the temporal components of these invariants to generate properties that Ms could manage. For instance, when TIMEPIECE would verify properties of the form  $G\phi$ ,  $F^t\phi$ , or  $\phi_2 \mathcal{U}^t \phi$ , Ms would instead verify that the network’s stable states satisfy  $\phi$ .

We ran all our benchmarks on a Microsoft Azure D96s v5 virtual machine with 96 vCPUs and 384GB of RAM. We took advantage of the machine’s multi-core processor to run all the modular checks in parallel, while monolithic checks necessarily ran on a single thread. We timed out any benchmark that did not complete in 4 hours. We reported four times for each of our benchmarks: (i) the total time until all TIMEPIECE threads finished (Tp); (ii) the median node check time; (iii) the 99<sup>th</sup> percentile node check time (99% of checks completed in less than this much time); and (iv) the total time taken by Ms.

---

<sup>4</sup>We use C#’s Parallel LINQ library [37], which can run up to 512 concurrent threads.



**Figure 5: An example fattree network for  $k = 4$ , showing how Vf sets  $s.down$  at the routes between the destination node  $e_{19}$  and  $e_6$ .  $a_4$  will drop the route from  $e_6$ , as accepting it would form a valley (as indicated by  $s.down$ ).**

**Experiments.** We parameterize our fattree networks by their number of pods  $k$ : a  $k$ -fattree has  $1.25k^2$  nodes and  $k^3$  edges: Figure 5 shows an example 4-fattree as part of our Vf policy. We considered multiples of 4 for  $4 \leq k \leq 60$  to assess TIMEPIECE’s scalability: whereas we expected a monolithic verifier to time out on larger topologies, we hypothesized that TIMEPIECE would scale to these networks. We present how verification time grows with respect to the number of nodes in each fattree in Figure 6. Note that the figure shows verification time on a logarithmic scale.

We considered five different properties: single-edge reachability (SpReach), all-edge reachability (ApReach), bounded path length (SpLen), valley freedom (Vf) and route filtering (Hijack). Our routes modelled the eBGP protocol in these networks. Figure 7 summarizes the eBGP fields represented and how we modelled them in SMT. We model the major common elements of eBGP routing: a route destination as a 32-bit integer (representing an IPv4 prefix); administrative distance, local preference, multi-exit discriminators as 32-bit integers (encoded as bitvectors); eBGP origin type as a ternary value; the AS path length as an (unbounded) integer; and BGP communities as a set of strings.

**SpReach.** SpReach demonstrates the simplest possible routing behavior and serves as a useful baseline. The policy simply increments the path length of a route on transfer. We initialized one destination edge node with a route to itself, and all other nodes with no route ( $\infty$ ). Our goal is to prove every node eventually has a route to the destination (*i.e.*, its route is not  $\infty$ ). More precisely, because our network has diameter 4, each node  $v$  should acquire a route in 4 time steps.

$$P_{\text{SpReach}}(v) \equiv \mathcal{F}^4(s \neq \infty)$$

Interfaces for these benchmarks mirror the simplicity of the policy and property. In order to establish that a node  $v$  eventually has a route  $s \neq \infty$  at time  $t$ , we must show that one of  $v$ ’s neighbors has a route at an earlier time  $t - 1$ . The destination node will have a route at time 0, its adjacent aggregation nodes will have routes at time 1, and so on for nodes

reachable from them at time 2. For brevity, we introduce a function  $dist$  to specify these times for each node, where  $dist(v)$  returns the length of the shortest path between  $v$  and the node originating the route. The interface will then be:

$$A_{\text{SpReach}}(v) \equiv \mathcal{F}^{dist(v)}(s \neq \infty)$$

SpReach’s policy and property are so simple that Tp is actually slightly *slower* than Ms, as shown in Figure 6a. We conjecture the Ms encoding reduces to a particularly easy SAT instance. That said, we can already see that individual checks in TIMEPIECE take only a fraction of the time that Ms takes, with 99% of node checks completing in at most 1.5 seconds, even for our largest benchmarks.

**ApReach.** For our first variation on SpReach, we modelled all-pairs reachability by using a symbolic value  $dest$  to consider any choice of edge node as the destination node. A node  $v$  will receive a route at different times depending on which pod  $dest$  is in relative to  $v$ . We define a new function  $apdist(v, dest)$  as a nested chain of if statements, returning a different time depending on whether  $v$  is (i) an aggregation node in  $dest$ ’s pod (1 hop,  $t = 1$ ); (ii) a core node, or an edge node in  $dest$ ’s pod (2 hops,  $t = 2$ ); (iii) an aggregation node in another pod ( $t = 3$ ); or (iv) an edge node in another pod ( $t = 4$ ). Our interface is then:

$$A_{\text{ApReach}}(v) \equiv \mathcal{F}^{apdist(v, dest)}(s \neq \infty)$$

Figure 6b shows that, perhaps from the burden of modelling the symbolic  $dest$ , Ms times out at  $k = 24$ . Tp verifies  $k = 60$  in under 15 minutes, with 99% of individual node checks taking under 15 seconds.

**SpLen.** Our next benchmark uses the same policy as SpReach, but considers a stronger property: every node eventually has a route of at most 4 hops to the destination.

$$P_{\text{SpLen}}(v) \equiv \mathcal{F}^4(s.len \leq 4)$$

To prove this property, our interfaces specify that path lengths in routes should not exceed the distance to the destination:  $s.len \leq dist(v)$ . In addition, because local preference influences routing, we fix the local preference to the default for all routes when present:  $s.lp = 100$ . The final interface is:

$$A_{\text{SpLen}}(v) \equiv \underbrace{\mathcal{G}(s = \infty \vee s.lp = 100)}_{\text{no better routes appear}} \sqcap \underbrace{\mathcal{F}^{dist(v)}(s.len \leq dist(v))}_{\text{eventually the route appears}}$$

Figure 6c shows that, because SpLen forces the solver to reason over path lengths, monolithic verification times out at  $k = 12$ . By contrast, modular verification is able to solve our largest benchmark ( $k = 60$  with 4,500 nodes) in under 7 minutes, with 99% of nodes verified in under 7 seconds.

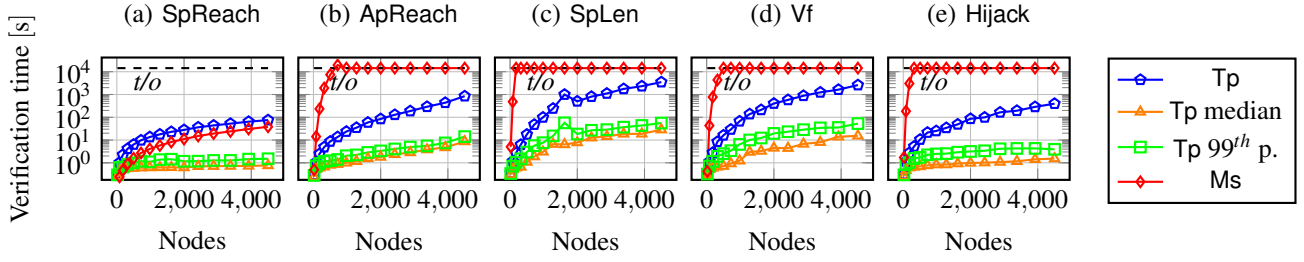


Figure 6: Ms vs. Tp verification times for fattree benchmarks with the SpReach, ApReach, SpLen, Vf, and Hijack policies.

Route field	Modelled type in SMT
Route destination	bitvector [47]
Administrative distance	bitvector [47]
eBGP local preference	bitvector [47]
eBGP multi-exit discriminator	bitvector [47]
eBGP origin type	{egp, igp, unknown} [47]
eBGP AS path length	integer [48]
eBGP communities	set<string> [46, 49]

Figure 7: eBGP route fields modelled by TIMEPIECE using SMT solvers for benchmark networks.

**Vf.** Vf extends SpReach with policy to prevent up-down-up (valley) routing [12, 13, 41], where routes transit an intermediate pod. To implement this policy, we add a BGP community  $D$  along “down” edges in the topology (*i.e.*, from a core node or from an aggregation node to an edge node), and drop routes with  $D$  on “up” edges (see the example network in Figure 5). For brevity, we write “ $s.down$ ” to mean “ $D \in s.tags$ ”. We test the same reachability property as SpReach.

The legitimate routes in the fattree all start as routes travelling up from the destination node’s pod, *e.g.*, the nodes in green ( $a_{16}, a_{17}, e_{19}$ ) in Figure 5. We refer to these nodes as “adjacent nodes” with a shorthand  $adj(v)$ : they provide the core nodes (and through them, the rest of the network) with routes via their up edges: these edges will drop the routes if  $s.down$ , so we require that  $adj(v) \rightarrow \neg s.down$ . To ensure this, we add conjuncts to our interfaces requiring that nodes’ final routes are no better than the shortest path’s route:  $s.lp = 100 \wedge s.len = dist(v)$ . This ensures our inductive condition holds after every node has a route: otherwise, a core node (for instance) could offer a spurious route with  $s.len < 1 \wedge s.down$  to an adjacent node. The final interface is:

$$A_{Vf}(v) \equiv \underbrace{(s = \infty)}_{\text{no better routes appear}} \underbrace{\mathcal{Q}^{dist(v)} \left( (s.lp=100 \wedge s.len=dist(v)) \sqcap (adj(v) \rightarrow \neg s.down) \right)}_{\text{adjacent nodes will share routes}}$$

Figure 6d shows that Ms verifies up to  $k = 16$  before timing out. As with SpLen, Tp time grows gradually in proportion to the number of nodes, topping out at 43 minutes for  $k = 60$ ,

with all node checks completing in under a minute.

**Hijack.** Hijack models a fattree with an additional “hijacker” node  $h$  connected to the core nodes.  $h$  represents a connection to the internet from outside the network, which may advertise any route. We add a boolean ghost state tag to  $S$  for this policy to mark routes as external (from  $h$ ) or internal. The destination node will advertise a route with  $s.prefix = p$ , where  $p$  is a symbolic value representing an announcement to an internal address (*i.e.*, one with a prefix length  $\geq 24$ ): the core nodes will then drop any routes from  $h$  which falsely claim  $s.prefix = p$ , but allow other routes through. Apart from this filtering, routing functions as in the SpReach benchmark. For this network, we verified that every node eventually has a route for prefix  $p$  which is not via the hijacker ( $\neg s.tag$ ), assuming nothing about the hijacker’s route ( $A_{Hijack}(h) \equiv P_{Hijack}(h) \equiv \mathcal{G}(\text{true})$ ). The property is thus:

$$P_{Hijack}(v) \equiv \mathcal{F}^4(s.prefix=p \wedge \neg s.tag)$$

The interface for this policy is straightforward. As our property does not depend on path length, local preference or communities, we must simply re-affirm that nodes with internal prefixes never have external routes:  $s.prefix = p \rightarrow \neg s.tag$ . Once nodes have received a route from the destination at time  $dist(v)$ , they should keep that route forever, and hence their route will have both  $s.prefix = p$  and  $\neg s.tag$ .

$$A_{Hijack}(v) \equiv \underbrace{\mathcal{F}^{dist(v)}(s.prefix=p \wedge \neg s.tag)}_{\text{route will be internally reachable}} \sqcap \underbrace{\mathcal{G}(s.prefix=p \rightarrow \neg s.tag)}_{\text{no hijack route is ever used}}$$

The monolithic and modular verification times follow now-familiar patterns. Monolithic verification times out again at  $k = 16$ , whereas modular verification times grow far more slowly — after initially tripling or doubling for each successive  $k$ , total modular time gradually slows to growing by a factor of 1.5. 99% of nodes complete their checks in under 4 seconds, with our longest check taking 19.5 seconds for an aggregation node in  $k = 60$ . As with our other benchmarks, this trend is due to how the in-degree of each node — which determines the size of the SMT encoding of our inductive condition — grows linearly for successive values of  $k$ .

## 7 Related Work

Our work is most closely related to other efforts in control plane verification [2, 5, 8–10, 20, 21, 23, 35, 43, 51–53]. The following paragraphs recommend different classes of tools depending on the specifics of one’s verification problem.

If your network is small (in the tens of nodes), then use an SMT-based tool such as Minesweeper [8] or Bagpipe [52] for their ease-of-use, generality, and symbolic reasoning. Minesweeper supports a broad range of properties including reachability, waypointing, no blackholes and loops, and device equivalence.

If your network is larger, and neither incremental recomputation after device update nor fully symbolic reasoning is important, use a simulation-based tool [10, 20, 23, 35, 43, 53]. Some of these tools also employ symbolic reasoning in limited ways to provide useful capabilities. For example, inspired by effective work on data plane analysis [32], Plankton [43] first analyzes configurations to identify IP prefix equivalence classes. Once the equivalence classes are known, they may be treated symbolically in the rest of the computation. Plankton might be able to reason symbolically about other attributes, but doing so would require additional custom engineering to find the appropriate sort of equivalence class ahead of time (*e.g.*, for BGP AS paths or communities). With TIMEPIECE, any component may be treated symbolically and the solver effort is passed off to the underlying SMT engine.

If your network is large and symbolic reasoning is important, then there are fewer options to consider. Bonsai exploits symmetry to derive smaller abstractions of a network [9], but does not work if the network has topology or policy asymmetries or one considers failures (which break topological symmetries). Exploiting the modularity in network designs is another way to achieve scalability in verification. Among recent efforts, Kirigami [5] applies assume-guarantee reasoning for control plane verification, but bases its underlying model on stable states and requires specifying exact routes as interfaces from the user. Another effort called Lightyear [51] focuses on safety properties of BGP configurations, where users provide local invariants on BGP policies on individual nodes and edges. Like TIMEPIECE, it performs local checks to verify the invariants and to ensure that they imply an end-to-end safety property. TIMEPIECE supports a more general network model that can support a wide range of protocols including BGP. On the other hand, Lightyear supports more modular checks (all neighbors considered at once in TIMEPIECE *vs.* one neighbor at a time in Lightyear) and thus is more scalable. In exchange, TIMEPIECE’s theory allows for the specification and verification of temporal invariants that are necessary to establish the existence of routes. Reachability falls into this class, a property of keen interest.

Daggitt *et al.* also use a timed model [15]. However, our verification method and target properties differ—they focus on convergence properties of routing protocols, whereas we

analyze properties that depend upon a network’s topology and configuration such as reachability.

Other inspirations for our work are SecGuru and RCDC [30] in data plane verification. Unlike our work, they use non-temporal invariants, which they extract from the network topology and assume as ground truth for policies on individual devices. Whereas our experiments likewise used the topology to define local invariants, our verification procedure checks that these invariants are in fact guaranteed by the other devices in the network.

**Modular verification techniques.** Our work is inspired by the success of modular verification based on *assume-guarantee reasoning* [22]. This has been widely used in verification of software, hardware and reactive systems [6, 19, 27, 29]. Many such applications use temporal logic for specifying assumptions and guarantees at component interfaces, along with correctness properties [42]. When the proof rules for modular reasoning are *circular*, *i.e.*, there is a circular chain of dependencies between components via their assumptions and guarantees, then one has to provide some well-founded ordering to carefully break the cycle for sound verification. TIMEPIECE uses time to provide this well-founded ordering.

## 8 Conclusion

Ensuring correct routing is critical to the operation of reliable networks that serve as a foundation for online interactions, communications, and services. Network verification of control plane routing protocols has long been a noble aspiration, yet existing tools and techniques have been shackled by a tradeoff between exhaustiveness associated with symbolic reasoning and scalability. In this paper we proposed TIMEPIECE, a radically different approach to verifying network routing based on splitting up the network into small modular components that we verify efficiently, in parallel. We believe that TIMEPIECE unleashes the full potential of control plane verification by achieving high scalability without sacrificing expressiveness. In return, it demands that users provide local interfaces to carry out its procedure. We proved that TIMEPIECE is sound and complete with respect to the network semantics, and argue that its temporal foundation is an excellent choice for a modular verification procedure.

## Acknowledgements

This work was supported in part by the National Science Foundation awards NeTS 1703493 and FMitF 1837030.

## References

- [1] Anubhavnidhi Abhashkumar, Aaron Gember-Jacobson, and Aditya Akella. *AED: Incrementally Synthesizing*



- Policy-Compliant and Manageable Configurations*, page 482–495. Association for Computing Machinery, New York, NY, USA, 2020. <https://doi.org/10.1145/3386367.3431304>.
- [2] Anubhavnidhi Abhashkumar, Aaron Gember-Jacobson, and Aditya Akella. Tiramisu: Fast multilayer network verification. In *17th USENIX Symposium on Networked Systems Design and Implementation (NSDI 20)*, pages 201–219, 2020. <https://www.usenix.org/system/files/nsdi20-paper-abhashkumar.pdf>.
- [3] Anubhavnidhi Abhashkumar, Kausik Subramanian, Alexey Andreyev, Hyojeong Kim, Nanda Kishore Salem, Jingyi Yang, Petr Lapukhov, Aditya Akella, and Hongyi Zeng. Running BGP in data centers at scale. In *18th USENIX Symposium on Networked Systems Design and Implementation (NSDI 21)*, pages 65–81. USENIX Association, April 2021. <https://www.usenix.org/conference/nsdi21/presentation/abhashkumar>.
- [4] Mohammad Al-Fares, Alexander Loukissas, and Amin Vahdat. A scalable, commodity data center network architecture. In *SIGCOMM*, 2008. <https://dl.acm.org/doi/10.1145/1402946.1402967>.
- [5] Tim Alberdingk Thijm, Ryan Beckett, Aarti Gupta, and David Walker. Kirigami, the verifiable art of network cutting, 2022. <https://arxiv.org/abs/2202.06098>.
- [6] Rajeev Alur and Thomas A Henzinger. Reactive modules. *Formal methods in system design*, 15(1):7–48, 1999. <https://doi.org/10.1023/A:1008739929481>.
- [7] Clark Barrett and Cesare Tinelli. Satisfiability modulo theories. In *Handbook of model checking*, pages 305–343. Springer, 2018. [https://doi.org/10.1007/978-3-319-10575-8\\_11](https://doi.org/10.1007/978-3-319-10575-8_11).
- [8] Ryan Beckett, Aarti Gupta, Ratul Mahajan, and David Walker. A general approach to network configuration verification. In *SIGCOMM*, August 2017. <https://doi.org/10.1145/3098822.3098834>.
- [9] Ryan Beckett, Aarti Gupta, Ratul Mahajan, and David Walker. Control plane compression. In *Proceedings of the 2018 Conference of the ACM Special Interest Group on Data Communication, SIGCOMM '18*, pages 476–489, New York, NY, USA, 2018. ACM. <https://doi.org/10.1145/3230543.3230583>.
- [10] Ryan Beckett, Aarti Gupta, Ratul Mahajan, and David Walker. Abstract interpretation of distributed network control planes. *Proc. ACM Program. Lang.*, 4(POPL), dec 2019. <https://doi.org/10.1145/3371110>.
- [11] Ryan Beckett and Ratul Mahajan. A general framework for compositional network modeling. In *Proceedings of the 19th ACM Workshop on Hot Topics in Networks*, pages 8–15, 2020. <https://doi.org/10.1145/3422604.3425930>.
- [12] Ryan Beckett, Ratul Mahajan, Todd Millstein, Jitendra Padhye, and David Walker. Don’t mind the gap: Bridging network-wide objectives and device-level configurations. In *SIGCOMM*, 2016. <https://doi.org/10.1145/2934872.2934909>.
- [13] Ryan Beckett, Ratul Mahajan, Todd Millstein, Jitu Padhye, and David Walker. Network configuration synthesis with abstract topologies. In *PLDI*, June 2017. <https://doi.org/10.1145/3062341.3062367>.
- [14] CISCO. Using bgp community values to control routing policy in upstream provider network, August 2005. <https://www.cisco.com/c/en/us/support/docs/ip/border-gateway-protocol-bgp/28784-bgp-community.html>.
- [15] Matthew L Daggitt, Alexander JT Gurney, and Timothy G Griffin. Asynchronous convergence of policy-rich distributed Bellman-Ford routing protocols. In *Proceedings of the 2018 Conference of the ACM Special Interest Group on Data Communication*, pages 103–116. ACM, 2018. <https://doi.org/10.1145/3230543.3230561>.
- [16] Ahmed El-Hassany, Petar Tsankov, Laurent Vanbever, and Martin Vechev. Netcomplete: Practical network-wide configuration synthesis with autocompletion. In *Proceedings of the 15th USENIX Conference on Networked Systems Design and Implementation, NSDI’18*, page 579–594, USA, 2018. USENIX Association. <https://www.usenix.org/conference/nsdi18/presentation/el-hassany>.
- [17] Seyed K. Fayaz, Tushar Sharma, Ari Fogel, Ratul Mahajan, Todd Millstein, Vyas Sekar, and George Varghese. Efficient network reachability analysis using a succinct control plane representation. In *OSDI*, 2016. <https://www.usenix.org/system/files/conference/osdi16/osdi16-fayaz.pdf>.
- [18] Nick Feamster and Hari Balakrishnan. Detecting BGP configuration faults with static analysis. In *NSDI*, May 2005. <https://www.usenix.org/legacy/events/nsdi05/tech/feamster/feamster.pdf>.
- [19] Cormac Flanagan and Shaz Qadeer. Thread-modular model checking. In *International SPIN Workshop on Model Checking of Software*, pages 213–224. Springer, 2003. [https://doi.org/10.1007/3-540-44829-2\\_14](https://doi.org/10.1007/3-540-44829-2_14).

- [20] Ari Fogel, Stanley Fung, Luis Pedrosa, Meg Walraed-Sullivan, Ramesh Govindan, Ratul Mahajan, and Todd Millstein. A general approach to network configuration analysis. In *NSDI*, October 2015. <https://www.usenix.org/system/files/conference/nsdi15/nsdi15-paper-fogel.pdf>.
- [21] Aaron Gember-Jacobson, Raajay Viswanathan, Aditya Akella, and Ratul Mahajan. Fast control plane analysis using an abstract representation. In *SIGCOMM*, August 2016. <https://doi.org/10.1145/2934872.2934876>.
- [22] Dimitra Giannakopoulou, Kedar S Namjoshi, and Corina S Păsăreanu. Compositional reasoning. In *Handbook of Model Checking*, pages 345–383. Springer, 2018. [https://doi.org/10.1007/978-3-319-10575-8\\_12](https://doi.org/10.1007/978-3-319-10575-8_12).
- [23] Nick Giannarakis, Devon Loehr, Ryan Beckett, and David Walker. NV: An intermediate language for verification of network control planes. In *PLDI*, page 958–973, New York, NY, USA, 2020. Association for Computing Machinery. <https://doi.org/10.1145/3385412.3386019>.
- [24] Sharon Goldberg. Why is it taking so long to secure internet routing? *Communications of the ACM*, 57(10):56–63, 2014. <https://queue.acm.org/detail.cfm?id=2668966>.
- [25] Timothy G. Griffin, F. Bruce Shepherd, and Gordon Wilfong. The stable paths problem and interdomain routing. *IEEE/ACM Trans. Networking*, 10(2), 2002. <https://ieeexplore.ieee.org/abstract/document/993304>.
- [26] Timothy G. Griffin and João Luís Sobrinho. Metarouting. In *SIGCOMM*, pages 1–12, August 2005. [10.1145/1080091.1080094](https://doi.org/10.1145/1145/1080091.1080094).
- [27] Orna Grumberg and David E Long. Model checking and modular verification. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 16(3):843–871, 1994. <https://doi.org/10.1145/177492.177725>.
- [28] C. Hedrick. Routing Information Protocol. Internet Request for Comments, October 1988. <https://datatracker.ietf.org/doc/html/rfc1058>.
- [29] Thomas A Henzinger, Shaz Qadeer, and Sriram K Rajamani. You assume, we guarantee: Methodology and case studies. In *International Conference on Computer Aided Verification*, pages 440–451. Springer, 1998. <https://doi.org/10.1007/BFb0028765>.
- [30] Karthick Jayaraman, Nikolaj Bjørner, Jitu Padhye, Amar Agrawal, Ashish Bhargava, Paul-Andre C Bissonnette, Shane Foster, Andrew Helwer, Mark Kasten, Ivan Lee, Anup Namdhari, Haseeb Niaz, Aniruddha Parkhi, Hanukumar Pinnamraju, Adrian Power, Neha Milind Raje, and Parag Sharma. Validating datacenters at scale. In *Proceedings of the ACM Special Interest Group on Data Communication, SIGCOMM '19*, page 200–213, New York, NY, USA, 2019. Association for Computing Machinery. <https://doi.org/10.1145/3341302.3342094>.
- [31] Peyman Kazemian, Michael Chang, Hongyi Zeng, George Varghese, Nick McKeown, and Scott Whyte. Real time network policy checking using header space analysis. In *NSDI*, pages 99–112, April 2013. <https://www.usenix.org/system/files/conference/nsdi13/nsdi13-final8.pdf>.
- [32] Ahmed Khurshid, Xuan Zou, Wenxuan Zhou, Matthew Caesar, and P. Brighten Godfrey. Veriflow: Verifying network-wide invariants in real time. In *NSDI*, April 2013. <https://www.usenix.org/system/files/conference/nsdi13/nsdi13-final100.pdf>.
- [33] K Rustan M Leino. Dafny: An automatic program verifier for functional correctness. In *International Conference on Logic for Programming Artificial Intelligence and Reasoning*, pages 348–370. Springer, 2010. [https://doi.org/10.1007/978-3-642-17511-4\\_20](https://doi.org/10.1007/978-3-642-17511-4_20).
- [34] Nuno P. Lopes, Nikolaj Bjørner, Patrice Godefroid, Karthick Jayaraman, and George Varghese. Checking beliefs in dynamic networks. In *NSDI*, 2015. <https://www.usenix.org/system/files/conference/nsdi15/nsdi15-paper-lopes.pdf>.
- [35] Nuno P Lopes and Andrey Rybalchenko. Fast BGP simulation of large datacenters. In *International Conference on Verification, Model Checking, and Abstract Interpretation*, pages 386–408. Springer, 2019. <https://web.ist.utl.pt/nuno.lopes/pubs/fastplane-vmcai19.pdf>.
- [36] K. Lougheed and Y. Rekhter. A Border Gateway Protocol 3 (BGP-3). Internet Request for Comments, October 1991. <https://datatracker.ietf.org/doc/html/rfc1267>.
- [37] Microsoft. Introduction to PLINQ. <https://docs.microsoft.com/en-us/dotnet/standard/parallel-programming/introduction-to-plinq>, 2021.
- [38] J. Moy. Open Shortest Path First Protocol Version 2. Internet Request for Comments, April 1998. <https://datatracker.ietf.org/doc/html/rfc2328>.

- [39] Sanjai Narain, Gary Levin, Sharad Malik, and Vikram Kaul. Declarative infrastructure configuration synthesis and debugging. *Journal of Network Systems Management*, 16(3), 2008.
- [40] D. Oran. OSI IS-IS Intra-domain Routing Protocol. Internet Request for Comments, February 1990. <https://datatracker.ietf.org/doc/html/rfc1142>.
- [41] Ivan Pepelnjak. Valley-free routing in data center fabrics. <https://blog.ipSPACE.net/2018/09/valley-free-routing-in-data-center.html>, 2018.
- [42] Amir Pnueli. In transition from global to modular temporal reasoning about programs. In *Logics and Models of Concurrent Systems - Conference proceedings*, volume 13 of *NATO ASI Series*, pages 123–144. Springer, 1984. [https://doi.org/10.1007/978-3-642-82453-1\\_5](https://doi.org/10.1007/978-3-642-82453-1_5).
- [43] Santhosh Prabhu, Ali Kheradmand, Brighten Godfrey, and Matthew Caesar. Predicting network futures with Plankton. In *Proceedings of the First Asia-Pacific Workshop on Networking*, APNet’17, pages 92–98, August 2017. <https://dl.acm.org/doi/10.1145/3106989.3106991>.
- [44] Simon Sharwood. IBM blames ‘external’ network provider, incorrect routing, traffic flood for its two-hour cloud outage. [https://www.theregister.com/2020/06/11/ibm\\_cloud\\_outage\\_report/](https://www.theregister.com/2020/06/11/ibm_cloud_outage_report/), 2020.
- [45] Simon Sharwood. Facebook rendered spineless by buggy audit code that missed catastrophic network config error. [https://www.theregister.com/2021/10/06/facebook\\_outage\\_explained\\_in\\_detail/](https://www.theregister.com/2021/10/06/facebook_outage_explained_in_detail/), 2021.
- [46] SMT-LIB. ArraysEx. <https://smtlib.cs.uiowa.edu/theories-ArraysEx.shtml>, 2010.
- [47] SMT-LIB. FixedSizeBitVectors. <https://smtlib.cs.uiowa.edu/theories-FixedSizeBitVectors.shtml>, 2010.
- [48] SMT-LIB. Ints. <https://smtlib.cs.uiowa.edu/theories-Ints.shtml>, 2010.
- [49] SMT-LIB. Unicode strings. <http://smtlib.cs.uiowa.edu/theories-UnicodeStrings.shtml>, 2020.
- [50] João Luís Sobrinho. An algebraic theory of dynamic network routing. *IEEE/ACM Trans. Netw.*, 13(5):1160–1173, October 2005. <https://ieeexplore.ieee.org/abstract/document/1528502>.
- [51] Alan Tang, Ryan Beckett, Karthick Jayaraman, Todd Millstein, and George Varghese. LIGHTYEAR: Using modularity to scale BGP control plane verification, 2022. <https://arxiv.org/abs/2204.09635>.
- [52] Konstantin Weitz, Doug Woos, Emina Torlak, Michael D. Ernst, Arvind Krishnamurthy, and Zachary Tatlock. Formal semantics and automated verification for the border gateway protocol. In *NetPL*, March 2016. <https://www.dougwoos.com/papers/bagpipe-netpl16.pdf>.
- [53] Fangdan Ye, Da Yu, Ennan Zhai, Hongqiang Harry Liu, Bingchuan Tian, Qiaobo Ye, Chunsheng Wang, Xin Wu, Tianchen Guo, Cheng Jin, Duncheng She, Qing Ma, Biao Cheng, Hui Xu, Ming Zhang, Zhiliang Wang, and Rodrigo Fonseca. Accuracy, scalability, coverage: A practical configuration verifier on a global WAN. In *Proceedings of the Annual Conference of the ACM Special Interest Group on Data Communication on the Applications, Technologies, Architectures, and Protocols for Computer Communication*, SIGCOMM ’20, page 599–614, New York, NY, USA, 2020. Association for Computing Machinery. <https://doi.org/10.1145/3387514.3406217>.
- [54] Peng Zhang, Aaron Gember-Jacobson, Yueshang Zuo, Yuhao Huang, Xu Liu, and Hao Li. Differential network analysis. In *19th USENIX Symposium on Networked Systems Design and Implementation (NSDI 22)*, pages 601–615, Renton, WA, April 2022. USENIX Association. <https://www.usenix.org/conference/nsdi22/presentation/zhang-peng>.

## A Proofs

We present the full proofs of Theorem 3.1, Corollary 3.2 and Theorem 3.3 below.

**Theorem A.1 (Soundness).** *Let  $A$  satisfy initial and inductive conditions. Then  $\forall t \in \mathbb{N}, \forall v \in V, \sigma(v)(t) \in A(v)(t)$ .*

*Proof.* Straightforward by induction over time. Let  $v$  be an arbitrary node in  $V$ .

At time 0, we have  $\sigma(v)(0) = I_v$  by the definition of  $\sigma$  and  $I_v \in A(v)(0)$  by the definition of an inductive invariant, so by substitution we have  $\sigma(v)(0) \in A(v)(0)$ .

For the inductive case, we must show that

$$\begin{aligned} (\forall v \in V, \sigma(v)(t) \in A(v)(t)) \\ \Downarrow \\ (\forall v \in V, \sigma(v)(t+1) \in A(v)(t+1)) \end{aligned}$$

We assume the antecedent (the inductive hypothesis). Consider the neighbors  $u_1, \dots, u_k$  of  $v$ . By the definition of an inductive invariant (6), we have:

$$\begin{aligned} \forall s_1 \in A(u_1)(t), \\ \forall s_2 \in A(u_2)(t), \\ \dots, \\ \forall s_n \in A(u_n)(t), \\ \left( I_v \oplus \bigoplus_{i \in \{1, \dots, n\}} f_{u_i v}(s_i) \right) \in A(v)(t+1) \end{aligned}$$

Then by our inductive hypothesis, we have that  $\sigma(u_i)(t) \in A(u_i)(t)$  for all  $u_i$ , so we can instantiate the universal quantifiers with these routes and substitute, which gives us:

$$I_v \oplus \bigoplus_{i \in \{1, \dots, n\}} f_{u_i v}(\sigma(u_i)(t)) \in A(v)(t+1)$$

The left-hand side is equal to the definition of  $\sigma(v)(t+1)$  in (4), so we have what we wanted to show:

$$\sigma(v)(t+1) \in A(v)(t+1)$$

Then  $\forall t \in \mathbb{N}, \forall v \in V, \sigma(v)(t) \in A(v)(t)$ .  $\square$

**Corollary A.2 (Safety).** *Let  $A$  satisfy initial and inductive conditions. Let  $P$  satisfy the safety condition with respect to  $A$ . Then  $\forall t \in \mathbb{N}, \forall v \in V, \sigma(v)(t) \in P(v)(t)$ .*

*Proof.* Let  $v$  be a node and  $t$  be a time. By Theorem 3.1,  $\sigma(v)(t) \in A(v)(t)$ . By the safety condition (7),  $A(v)(t) \subseteq P(v)(t)$ . Then  $\sigma(v)(t) \in P(v)(t)$ , i.e.,  $P$  holds for  $v$  at  $t$ .  $\square$

**Theorem A.3 (Completeness).** *Let  $\sigma$  be the state of the network. Then  $A(v)(t) = \{\sigma(v)(t)\}$  for all  $v \in V$  and all  $t \in \mathbb{N}$  satisfies the initial and inductive conditions.*

*Proof.* Let  $A$  be a function from nodes and times to singleton sets of routes such that  $A(v)(t) = \{\sigma(v)(t)\}$  for all  $v \in V$  and all  $t \in \mathbb{N}$ . Let  $v$  be an arbitrary node in  $V$  and let  $t$  be an arbitrary time. We want to show that  $\{\sigma(v)(t)\}$  is an inductive invariant.

Starting with the initial condition case, at time 0, we have  $\sigma(v)(0) = I_v$  by the definition of  $\sigma$ . Since  $\sigma(v)(0) \in \{\sigma(v)(0)\}$ , the initial condition holds.

For the inductive condition case, we want to show that:

$$\begin{aligned} \forall s_1 \in \{\sigma(u_1)(t)\}, \\ \forall s_2 \in \{\sigma(u_2)(t)\}, \\ \dots, \\ \forall s_n \in \{\sigma(u_n)(t)\}, \\ \left( I_v \oplus \bigoplus_{i \in \{1, \dots, n\}} f_{u_i v}(s_i) \right) \in \{\sigma(v)(t+1)\} \end{aligned}$$

Checking set inclusion of  $s \in \{\sigma(v)(t)\}$  is equivalent to checking that  $s = \sigma(v)(t)$ , so we can simplify the expression further by substituting  $\sigma(u_i)(t)$  for  $s_i$ :

$$\left( I_v \oplus \bigoplus_{i \in \{1, \dots, n\}} f_{u_i v}(\sigma(u_i)(t)) \right) = \sigma(v)(t+1)$$

This is now simply (4) flipped, so the inductive condition holds.

Then  $A$  is an inductive invariant.  $\square$