On the Accuracy of Formal Verification of Selective Defenses for TDoS Attacks

Marcilio O. O. Lemos

Federal University of Paraíba, João Pessoa, Brazil.

Yuri Gil Dantas

Technische Universität Darmstadt, Darmstadt, Germany

Iguatemi E. Fonseca

Federal University of Paraíba, João Pessoa, Brazil.

Vivek Nigam Federal University of Paraíba, João Pessoa, Brazil.

Abstract

Telephony Denial of Service (TDoS) attacks target telephony services, such as Voice over IP (VoIP), not allowing legitimate users to make calls. There are few defenses that attempt to mitigate TDoS attacks, most of them using IP filtering, with limited applicability. In our previous work, we proposed to use selective strategies for mitigating HTTP Application-Layer DDoS Attacks demonstrating their effectiveness in mitigating different types of attacks. Developing such types of defenses is challenging as there are many design options, e.q., which dropping functions and selection algorithms to use. Our first contribution is to demonstrate both experimentally and by using formal verification that selective strategies are suitable for mitigating TDoS attacks. We used our formal model to help decide which selective strategies to use with much less effort than carrying out experiments. Our second contribution is a detailed comparison of the results obtained from our formal models and the results obtained by carrying out experiments. We demonstrate that formal methods is a powerful tool for specifying defenses for mitigating Distributed Denial of Service attacks allowing to increase our confidence on the proposed defense before actual implementation.

Email addresses: marciliolemos@ci.ufpb.br (Marcilio O. O. Lemos), dantas@mais.informatik.tu-darmstadt.de (Yuri Gil Dantas), iguatemi@ci.ufpb.br (Iguatemi E. Fonseca), vivek.nigam@gmail.com (Vivek Nigam)

Preprint submitted to Journal of Logical and Algebraic Methods in ProgrammingMay 16, 2017

1 1. Introduction

Telephony Denial of Service (TDoS) attacks is a type of Denial of Service (DoS) attack that target telephony services, such as Voice over IP (VoIP). With the increase in the popularity of VoIP services, we have witnessed an increase in TDoS attacks being used to target hospital line systems [1, 2] and systems for emergency lines (like the American 911 system) [3]. According to the FBI, 200 TDoS attacks have been identified only in 2013 [2].

This paper investigates the use of *selective defenses* [4] for mitigating one type of TDoS attack called Coordinated Call [5] attack. The Coordinated Call attack [5] exploits the fact that pairs of attackers, Alice and Bob, can collude to 10 exhaust the resources of the VoIP server. Assume that Alice and Bob are valid 11 registered users.¹ The attack goes by Alice simply calling Bob and trying to 12 stay in the call as long as she can. Since the server allocates resources for each 13 call, by using enough pairs of attackers, attackers can exhaust the resources of 14 the server denying service to honest participants. This is a simple, but ingenious 15 attack, as only a relatively low rate of incoming calls is needed generating a small 16 network traffic (when compared to SIP flooding attack for example). Thus it is 17 hard for the network administrator to detect and counter-measure such attack. 18 Formal methods and, in particular, rewriting logic can help developers to 19 design defenses for mitigating DDoS attacks. In our previous work [4] we used 20 selective strategies in the form of the tool SeVen for mitigating HTTP Low-21 Rate Application-Layer DDoS attacks targeting web-servers. We formalized 22 different attack scenarios in Maude [6] and since our strategies are constructed 23 over some probability functions, we used statistical model checking [7], namely 24 PVeStA [8], to validate our defense. Due to our reasonable preliminary results, 25 we implemented SeVen and carried out experiments over the network obtaining 26 similar results to the ones obtained using formal methods. It took us only 3 27 person months to obtain our results using formal methods, while it took us 24 28 *person months* to obtain our first experimental results. Although we strongly 29 believe that systems should also be validated by means of experiments, the 30 confidence acquired from our formal analysis was invaluable for the success of 31 this project.² 32

This paper provides more evidence supporting the claim that formal methods can help specifiers in designing selective defenses. We systematically consider a number of selective defenses used for mitigating TDoS attacks. We compare the results obtained using our formal specification and the results obtained implementing such defenses and carrying out experiments on the network. Our results show a high accuracy for most of the results, specially on availability,

¹This can be easily done for many VoIP services.

²Notice that although our experiments on the network were controlled experiments, they used off-the-shelf tools, such as Apache web-servers, which implement a number of optimizations not modeled in our formal specification. Moreover our experiments suffered from interference that cannot be controlled, such as network latency. The same is true for our results involving the VoIP server Asterisk used in our experimental results.

³⁹ and less accurate on results involving time measurements.

40 Our contributions are three-fold:

• We formalized in Maude the Coordinated Call attack and three selective defenses based on SeVen: the first using a uniform selection strategy, the second with roulette selection strategy [9], and the third with a tournament selection strategy [10]. We also considered two models for legitimate call duration: an exponential call duration which models traditional telephony [11] and lognormal call duration which models VoIP telephony [12].

We carried out a number of simulations using PVeStA to test the efficiency
of each version the defense used under the two different assumptions on call
duration. Our simulation results suggest that SeVen mitigates the Coordinated Call attack;

• We implemented the different selective defenses analysed using our formal models, and integrated them with the VoIP server Asterisk [13] using the SIP-protocol. We also implemented the Coordinated Call attack. Our experimental results demonstrate in practice that our selective defenses can mitigate the Coordinated Call attack;

• Finally, we compare the results obtained from our formal analysis with the results obtained from our experimental results to analyze the accuracy of the results obtained from our formal analysis. This comparison demonstrates that formal methods are of great value as they can be used early on to develop and evaluate new defense mechanisms for mitigating TDoS attacks with much less effort than implementing defenses and carrying out experiments on the network.

A small subset of experimental and simulation results appearing in this pa-63 per appeared in our previous work [14, 15] which only considered scenarios 64 where call duration followed a uniform probability and a single mechanism for 65 dropping calls, namely the roulette strategy. This paper extends our previous 66 work by considering different assumptions on call duration, namely lognormal 67 distribution, modeling usual VoIP calls, and exponential distribution, modeling 68 usual telephony, *i.e.*, non VoIP calls. Moreover, we consider here different mech-69 anisms for dropping calls, namely uniform, roulette and tournament dropping 70 strategies. In terms of total time of experiments, the results in this paper add 71 more than 40 hours of experimental results when compared to the results in our 72 previous work [14, 15]. 73

This paper is organized as follows. Section 2 we review the Session Initiation 74 Protocol (SIP) used for initiating a VoIP call and also explain the Coordinated 75 Call attack. Section 3 describes how SeVen works, while Section 4 details its 76 formalization in Maude. Sections 5 and 6 contain our simulation and exper-77 imental results including our main assumptions, results and discussion of the 78 results obtained and Section 7 discusses the accuracy of our simulation results. 79 We discuss in Section 8 related and future work. Finally, the implementation 80 used to carry out our simulations is available for download at [16]. 81



Figure 1: Exchange of messages between the server and two users (Alice and Bob) during a normal execution of the SIP protocol.

⁸² 2. VoIP Protocols and the Coordinated Call Attack

We now review the Session Initiation Protocol [17], which is one of the main protocols used to establish Voice over IP (VoIP) connections. Figure 1 shows the message exchanges performed to establish a connection between two registered users, Alice and Bob, where Alice tries to initiate a conversation with Bob. It also contains the messages exchanged to terminate the connection.

For initiating a call, Alice sends an INVITE message to the SIP server in-88 forming that she wants to call Bob. If Bob or Alice is not registered as valid 89 users, the server sends a reject message to Alice. Otherwise, the server sends 90 an INVITE message to Bob.³ At the same time, the server sends a TRYING 91 message to Alice informing her the server is waiting for Bob's response to Al-92 ice's invitation. The server waits for a RINGING message from Bob indicating 93 that Bobs telephone is ringing. Bob might reject the request, in which case the 94 server informs Alice (not shown in the Figure), or accept the call by sending 95 the message OK. Finally, the server sends the message OK to Alice who sends 96 an ACK message back to the server which forwards it to Bob. 97

At this point, the communication is established and Alice and Bob should be able to communicate as long as they need/want. (This is represented by the three ellipses in Figure 1.) The call is then terminated once one of the parties (Alice) sends a BYE message to the server. The server then sends a BYE message to the other party (Bob), which then answers with the message OK, which is forwarded to Alice, and the connection is terminated.

¹⁰⁴ Coordinated VoIP Attack [5]. A pair of colluding attackers, A_1 and A_2 , that ¹⁰⁵ are registered in the VoIP service,⁴ call each other and stay in the call for ¹⁰⁶ as much time as they can. Once the call is established, the attackers stay ¹⁰⁷ in the call for indefinite time. They might be disconnected by some Timeout ¹⁰⁸ mechanism establishing some time bounds on the amount of time that two users

 $^{^{3}}$ In fact, we omit some steps carried out by the server to find Bob in the network. This step can lead to DDoS amplification attacks [18] for which known solutions exists. Such amplification attacks are not, however, the main topic of this paper.

⁴Or alternatively two honest users that have been infected to be zombies by some attacker.

might call. During the time that A_1 and A_2 are communicating, they are using 109 resources of the server. Many VoIP servers have an upper-bound on the number 110 of simultaneous calls they can handle. If enough pairs of attackers collude, 111 then the resources of the server can be quickly exhausted. This attack is hard 112 to detect using network analyzers because the traffic generated by attackers 113 is similar to the traffic generated by legitimate clients. The attackers follow 114 correctly the SIP protocol and, moreover, there is no need to generate a large 115 burst of calls, but rather place calls in a moderate pace. Eventually, the server's 116 capacity will be exhausted. 117

There are many reasons why VoIP devices participate in a Coordinated Call attack. Pairs of legitimate users may be unsatisfied with the VoIP provider and participate in such attacks. Attackers may also use botnets with some infected malware. There has been evidence of the use of botnets in 2007 [19]. Tools that can place and receive calls (SIPp [20]) and tutorials on the Internet help automate the steps for carrying out the Coordinated Call attack.

Indeed, we have done so and as we demonstrate in Section 6, the Coordinated
 Call attack can reduce considerably availability to levels around 5% without
 generating large amount of traffic.

127 3. Selective Strategies

We proposed in our previous work [4] a new defense mechanism, called 128 SeVen, for mitigating Application-Layer DDoS attacks (ADDoS) using selective 129 strategies [21]. An application using selective strategies does not immediately 130 process incoming messages, but waits for a period of time, T_S , called a round. 131 During a round, SeVen accumulates messages received in an internal buffer. As-132 sume that k is the maximum capacity of the service being protected. For VoIP 133 servers k is the number of calls that the application can handle simultaneously. 134 If the number of messages accumulated reaches k and a new incoming request 135 R arrives, SeVen behaves as follows: 136

1. SeVen decides whether to process R or not based on a probability P_1 . P_1 is defined using the counter PMod following [22]:

$$\frac{k}{k + \texttt{PMod}}$$

At the beginning of the round, we set PMod = 0. PMod is incremented whenever the application's capacity is exhausted and a new incoming request arrives reducing thus the probability of new incoming request being selected by SeVen during a round. The intuition is that P_1 reduces with the increase of incoming traffic thus reducing the impact of high number of request to the application;

2. If SeVen decides to process R, then as the application is overloaded, it should decide which request currently being processed should be dropped. This decision is governed by P_2 , a distribution probability which might depend on the state of the existing request; 3. Otherwise, SeVen simply drops the request R without affecting the requests currently being processed and sends a message to the requesting user informing that the service is temporarily unavailable.

At the end of the round, SeVen processes the requests that are in its internal buffer (surviving the selective strategy) and sends them to the application.

SeVen mitigates attacks only when the maximum capacity of the VoIP server 152 reached. When this happens, SeVen has two mechanisms for dropping reis153 quests. The first one is by using probability P_1 and the other by using P_2 . The 154 main goal of the former is to mitigate the impact of volumetric attacks [21]. 155 This is because whenever the defense receives high volume attack the value 156 of PMod increases rapidly increasing rapidly the chance of dropping a request. 157 This is also reflected on the round time T_S which is typically of some hundred 158 milliseconds in order to avoid PMod from reaching too high values even under 159 normal traffic. We used $T_S = 100$ ms. 160

As Coordinated Calls do not generate a very large number of requests, the 161 mechanism using P_1 is not the main mitigation mechanism used by SeVen for 162 this attack, but rather the mechanism using P_2 . There is, however, much space 163 for specifying the probability distribution P_2 governing SeVen. In [4], we showed 164 that by using simple *uniform distributions* for dropping existing requests, SeVen 165 can be used to mitigate a number of ADDoS attacks using the HTTP protocol. 166 such as the Slowloris and HTTP POST attacks even in the presence of a large 167 number of attackers. 168

For mitigating the Coordinated Call attack described in Section 2, we set the probability P_2 , governing which call to be dropped from the internal buffer, to depend on (1) the status of the call and (2) on the duration of a call. We consider two types of call status:

• WAITING: A call is WAITING if it has already sent an INVITE message, but it is still waiting for the responder to join the call, that is, it has not completed the initiation part of the SIP protocol;

INCALL: A call is INCALL if the messages of initiation part of SIP have been completed and the initiator and the responder are already communicating (or simply in a call).

Thus, any incoming INVITE requests assume the status of WAITING, and these can change its status to INCALL once the initiation part of SIP is completed.

We assume here that it is preferable to a VoIP server, when overloaded, to 181 drop WAITING requests than INCALL requests that are communicating not for a 182 *very long duration.* In many cases, it is true that interrupting an existing call is 183 considered to be more damaging to server's reputation than not allowing a user 184 to start a new call. This could also be modeled by configuring the probability 185 distributions of SeVen accordingly. To determine whether a call is taking too 186 long, we assume that the server knows what is the average duration, t_M , of 187 calls.⁵ 188

⁵The value of t_M can be obtained by the history of a VoIP provider's usage.



Figure 2: Graph (not in scale) illustrating the behavior of SeVen according to the status of a call and its duration. p_{WAIT} is the probability of dropping a WAITING call, while p_{IN} the probability of dropping a INCALL call.

The dropping factor of an INCALL request increases exponentially once the call duration is greater than t_M . Figure 2 depicts roughly the dropping factor used to drop requests. The actual function d (for drop factor) is of the form, where t is the call duration where α is a parameter:

$$d(t) = \begin{cases} p_{\text{WAIT}} & \text{if } t = 0\\ p_{\text{IN}} & \text{if } 0 < t \le t_M \\ p_{\text{WAIT}} + e^{\alpha t/t_M} & \text{if } t > t_M \end{cases}$$
(1)

Given this dropping factor, we consider in our analysis three ways for selecting which call to drop. Assume the server has capacity of k simultaneous calls. Moreover, assume c_1, \ldots, c_k are the calls currently being processed by the server and they have dropping factors of, respectively, d_1, \ldots, d_k . We considered three different selection strategies described in the literature, namely, uniform [4], roulette [9] and *n*-tournament [10]:

- Uniform: In this strategy, the dropping factor of a call is not considered. We select using uniform probability which call is going to be dropped. Thus any call independent on its duration and status can be selected to be dropped by SeVen;
- Roulette: In the roulette strategy [9], we select randomly a call c_i to be dropped where the probability of being dropped is proportional to its dropping factor. Thus in the roulette strategy a call c_i has twice the chance of being dropped than a call c_i if $d_i = 2 \times d_i$.
- For instance, consider k = 4 and that the server is serving the calls c_1, c_2, c_3, c_4 with dropping factors 2, 3, 1, 6 respectively. We select using uniform distribution a number r between 0 and 2 + 3 + 1 + 6 = 12. If $0 \le r < 2$, then the call c_1 is selected, if $2 \le r < 5$, then the call c_2 is selected, if $5 \le r < 6$, then the call c_3 is selected, and otherwise if $6 \le r < 12$ then the call c_4 is selected. In this way, the call c_4 has 6 times more chance to be selected than the call c_3 for example;
- *n*-Tournament: In the *n*-tournament strategy [10], we first select *n* calls randomly using uniform probability to be part of the *tournament*. Then, the call to be dropped will be the call with the greatest dropping factor among the *n* selected calls. In case there is more than one possible call with the greatest dropping factor, we select one of them at random.

For instance, in the example above, if n = 2, then we would select randomly two out of the four calls c_1, c_2, c_3, c_4 to participate in the tournament. Say the calls c_2 and c_3 are chosen to be part of the tournament. In this case, the call c_2 is selected to be dropped as it has the greatest drop factor.

Notice that if n is chosen to be too low when compared to k, the n-223 tournament behaves closer to the uniform dropping strategy. In fact, if 224 n = 1, then the *n*-tournament can be shown to be equivalent to the uni-225 form dropping strategy. On the other hand, if n is chosen to be too high, 226 then the *n*-tournament behaves closer to a deterministic dropping strategy 227 that selects the call with the greatest dropping factor. Indeed, if k = n, 228 then the *n*-tournament strategy is deterministic. In our experiments, we 229 used n = k/2, that is, a strategy between the uniform and a deterministic 230 strategies. 231

²³² While the attackers attempt to stay in a call for very long periods of time, ²³³ legitimate clients do not behave so. The literature models legitimate call du-²³⁴ ration using the following distributions, where the parameters λ , σ and μ are ²³⁵ computed accordingly to the mean time assumed t_M (see [23] for more details):

• Exponential: Typical telephony models [11], *i.e.*, not VoIP, assume that the call duration of legitimate clients follows an exponential density distribution:

$$f(x,\lambda) = \begin{cases} \lambda e^{-\lambda x} & x \ge 0\\ 0 & x < 0 \end{cases}$$
(2)

Since the coordinated call attack can also be carried out in standard tele phony systems, we considered call duration following this distribution.

• **Lognormal:** While standard telephony calls are paid per duration, in VoIP calls have fixed rates or are even for free. This difference impacts legitimate call duration which in VoIP follows a lognormal density distribution [12]:

$$f(x,\mu,\sigma) = \frac{1}{x\sigma\sqrt{2\pi}} \exp\left[-\frac{(\ln(x)-\mu)^2}{2\sigma^2}\right].$$
 (3)

238 3.1. Sample Execution

Consider the following buffer, \mathcal{B}_i , at the beginning of a round and assume that k = 3, PMod = 0, the current time is 9 and the average call duration is $t_M = 5$ time units:

 $\mathcal{B}_1 = [\langle id_1, \mathsf{WAITING}, \mathsf{undef} \rangle, \langle id_2, \mathsf{INCALL}, 0.5 \rangle]$

²³⁹ $\langle id, st, tm \rangle$ specifies that the call *id* has status *st* and the call started at time ²⁴⁰ *tm* where *tm* is **undef** whenever st = WAITING. This buffer specifies that id_1 ²⁴¹ is waiting the responding party to answer (with a OK message) his invitation ²⁴² request and that id_2 is currently in a call. This means that id_2 is calling already ²⁴³ for way more than the expected average.

Assume that a message $\langle id_1, OK \rangle$ at time 9.5 arrives specifying that the responder of the request id_1 answered the call. The buffer is updated to the following:

$$\mathcal{B}_2 = [\langle id_1, \mathsf{INCALL}, 9.5 \rangle, \langle id_2, \mathsf{INCALL}, 0.5 \rangle]$$

Then the message $\langle id_3, \text{INVITE} \rangle$ arrives. Since the buffer is not yet full, a new request is inserted in the buffer and the message TRYING is sent to the requesting user. Notice that the RINGING message is not yet sent to the responding user. The buffer changes to:

 $\mathcal{B}_3 = [\langle id_1, \mathsf{INCALL}, 9.5 \rangle, \langle id_2, \mathsf{INCALL}, 0.5 \rangle, \langle id_3, \mathsf{WAITING}, \mathsf{undef} \rangle]$

Suppose now that another message $m_1 = \langle id_4, \text{INVITE} \rangle$ arrives at time 10.5. As the buffer is now full, it sets PMod to 1 and the application has to decide whether it will keep m_1 . SeVen generates a random number in the interval [0,1] using uniform distribution. Say that this number is less than (3/3+1), which means that it will select to process m_1 . However, it has to drop some existing request. The current requests id_1, id_2, id_3 have dropping factors following Figure 2:

• id_1 has dropping factor p_{IN} to be dropped because it is calling for a duration less than t_M : 10.5 - 9.5 < 5;

• id_2 has a much higher dropping factor because it is calling for twice t_M : ²⁵³ $10.5 - 0.5 = 2 \times 5;$

• id_3 has dropping factor of p_{WAIT} as it has WAITING status.

The application decides which one to drop either using *uniform probability*, in which case the dropping factor of requests is not considered, or the *roullete strategy*, in which case id_2 has a greater probability of being dropped, or the *n*-tournament strategy in which case it would depend on *n*.

Suppose that the application decides to drop id_2 , which means that the call is interrupted by the application. The resulting buffer is:

 $\mathcal{B}_4 = [\langle id_1, \mathsf{INCALL}, 9.5 \rangle, \langle id_4, \mathsf{WAITING}, \mathsf{undef} \rangle, \langle id_3, \mathsf{WAITING}, \mathsf{undef} \rangle]$

Assume that now the round time is elapsed. The application sends a RINGING message to the responder of the requests id_3 and id_4 .

²⁶² 4. Formal Specification

Our specification follows [4, 24, 25] by specifying test scenarios using actors where attackers, clients, and the server send and receive messages. These messages are stored in a scheduler that maintains a queue of messages. The attackers do not take control over the channel. Instead they share a channel with clients.

We formalize all actors in Maude [6] and carry out simulations by using the 268 statistical model checker PVeStA [8]. For simplicity, we considered the server 269 and SeVen as one actor, which means that SeVen is also able to operate as 270 a normal SIP Server, e.g., processing and establishing call connections. Such 271 decision does not affect the analysis of our results as in practice SeVen and the 272 VoIP server are in the same machine and thus share a quick communication 273 channel. In the following, we describe our Maude specification. The complete 274 model can be found in [16]. 275

We refer to [6, 26] for more detailed description of Maude and its underlying foundations on Rewriting Logic.

```
278 4.1. Key Sorts and Functions
```

```
279 Actor. The elements of the sort Actor is constructed by the operator
```

```
280 op <name:_|_> : Address AttributeSet -> Actor .
```

which takes an Address, which can be a string, and a set of attributes, AttributeSet.
In our formalization, we use the following attributes:

```
283 op req-cnt:_ : Float -> Attribute .
284 op b-set:_ : NBuffer -> Attribute .
285 op server:_ : Address -> Attribute .
286 op status:_ : Status -> Attribute
```

The attribute req-cnt stores the value of PMod, b-set stores the internal buffer of the server of sort NBuffer and server stores the address of the server. Finally, the attribute status specifies the status of the call which may be any of the one of the following state constants:

- none a call that has not yet been placed;
- invite a call that has been placed and is waiting for the responder to answer;
- incall a call where the participants are currently communicating;
- complete a call that has been completed, *i.e.*, the parties have communicated for the expected time;
- incomplete a call that was interrupted while communicating by SeVen.

The SeVen buffer has sort NBuffer and is constructed by pairing a number and a buffer:

300 op [_|_] : Nat Buffer -> NBuffer .

The number specifies the elements in the buffer which is a list of elements of sort EleBuf. These elements are constructed using a 3-tuple of the form:

```
303 op <___> : Address State Float -> EleBuf .
```

The first position stores the address of the call. The second position denotes the state of the call. The third position stores the time of the first request of the call and is used to compute the call duration.

Messages. Actors process messages of sort Msg which are constructed using the
 following operator

309 op _<-_ : Address Contents -> Msg .

The first parameter of sort Address specifies to which actor the message is
directed and the second parameter of sort Contents is payload of the message.
An ActiveMsg is a timestamped Msg constructed using the following operator:

313 op {_,_} : Float Msg -> ActiveMsg .

The first parameter specifies the time when the paired message is to be processed.

We assume that an active message {gt1,msg1} is always going to be processed before an active message {gt2,msg2} whenever gt1 < gt2. This is accomplished by using a Scheduler as in [4, 24, 25]. A scheduler has sort Scheduler and is constructed by the following operator

```
320 op [_|_] : Float ActiveMsgList -> Scheduler .
```

The first parameter is the global time while the second parameter contains the list of active messages ordered by their delivery time. The following function returns the scheduler obtained by inserting at the correct position, *i.e.*, according to the messages timestamp, a list of active messages into a given scheduler.

```
325 op insert : Scheduler ActiveMsgList -> Scheduler .
```

```
326 Configuration. A configuration of sort Config is a collection of Actors and
327 Scheduler:
```

```
subsort Actors < Config .</p>
subsort Scheduler < Config .</p>
op __ : Config Config -> Config [assoc comm id: none]
```

```
<sup>331</sup> For example, the initial configuration is defined by the equation:
```

```
332 eq initState =
```

It specifies a configuration with three actors: an attacker generator, a client generation and a server. The global time is 0.0 and there are three active messages in the scheduler, the first two directed to the actors attacker-generate and client-generate, respectively. Here we follow the Shared Channel Model [21] where clients and attackers share the same channel. Thus the application does not distinguish malicious traffic from legitimate one as it is usually the case, in particular, for the Coordinated Call attack.

Intuitively, the user specifies the rate at which new client and attacker actors are created. Then, for instance, when the message client-generate is processed a new client actor is created and new message client-generate is scheduled so that a new client actor is created according to the rate given. Similarly, when processing attacker-generate which creates a new attacker actor and a new attacker-generate is scheduled to be processed according to the attacker generating rate. These rewrite rules are omitted.

The third message scheduled at the time **Ts** is directed to the **server** which is implementing the SeVen strategy establishing when a SeVen round ends. The following function extracts the first scheduled active message in a scheduler and returns a new scheduler with the global time advanced to its delivery time.

358 op mytick : Scheduler -> Config .

³⁵⁹ For instance, let msg be a message and SL an ActiveMsgList. Then

360 mytick([0.0 | {1.0,msg}; SL])

returns msg [1.0 | SL] containing the message msg and the scheduler [1.0
 | SL]. Intuitively, the message msg is going to be processed next. Message
 processing is formalized by rewriting rules.

Selective Strategies. Finally, we specified the three selective strategies described
 in Section 3, namely, Uniform, Roulette and Tournament. The function

366 op select : Float Buffer -> ActorBuffer .

implements one of the selective strategies. Given the global time and a buffer,
this function returns a pair of sort ActorBuffer with the actor name of the
selected element that has been selected to be removed and the resulting buffer
obtained by removing this element from the given buffer.

All selective strategies we formalized use the following function:

```
372 op sampleUniWithInt : Nat -> Nat
```

which for a given input n returns a random natural number between 0 and n. The following function uses this function to select at random an element from a given buffer and thus to implement the uniform selection strategy.

```
376 op selectRandom : Buffer -> EleBuf .
```

For the roulette strategy, we compute using the dropping factor

378 op roulette : Float Buffer -> EleBuf .

which creates a roulette by assigning weights to the elements of the buffer according to the dropping factor and then randomly selects one.

A tournament for the n-tournament selection strategy is created by the following function:

383 op creatingTour : NBuffer Nat Buffer -> Buffer .

It takes an NBuffer and a natural number, specifying the size of the tournament,
and accumulates the selected elements to the tournament in the third argument.
Once the tournament is created, we use the following function to select the one
with the greatest dropping factor:

388 op selectGreatest : Buffer -> EleBuf .

which takes a buffer with the tournament traversing it to find the element withthe greatest dropping factor.

```
391 4.2. Rewrite Rules
```

The rewrite rules modify elements from **Conf** and specify the operational semantics of a system. We describe next the main rewrite rules used in our formalization.

The first action we described is when an actor receives a pool message indicating that it should start a call at time gt + delay.

The following rewrite rule specifies the behavior of a client upon receiving a RINGING message from the server. It changes the client's state from *invite* to *connected* and generates a message BYE, scheduled to be sent after some time. This means that all legitimate clients do not overpass the average time of the duration of calls using one of the call duration models, exponential or lognormal, described in Section 3. This means that the client called c(i) is expected to end its call at time gt + callDur(tMedio).

⁴¹⁷ SeVen may, however, drop a call before the call is finished. We classify such a call
⁴¹⁸ as an incomplete call, *i.e.* the dropped client's status is changed to incomplete.
⁴¹⁹ We omit this rule.

The rewrite rules for the attackers are similar to the client rules. The only
difference is that no BYE message is generated, thus, specifying the Coordinated
Call attack where attackers attempt to stay in the call for indefinite time. We
elide these rules.

Figure 3 depicts the rules implementing SeVen's strategy. For each INVITE 424 message received by some actor Actor, the rule SeVen-RECEIVE-INVITE checks 425 whether the buffer of the server reached its maximum. If not, then the incoming 426 request is added to the server's buffer (ConfAcc2) and a message TRYING to the 427 corresponding actor is created. Otherwise, SeVen throws a coin (p1) to decide 428 whether the incoming request will be processed using pmod. If SeVen decides 429 to process the incoming request, then some request being processed is selected 430 to be dropped using the function select. It returns the name of the actor 431 ActorDr and the resulting buffer nBuf. The incoming request is added to nBuf 432 and pmod gets incremented, resulting in the configuration ConfAcc. Moreover, 433

```
crl [SeVen-RECEIVE-INVITE] :
    <name: Ser | req-cnt: pmod , b-set: [lenB | B], AS >
    {Ser <- INVITE(Actor)} [gt | SL]</pre>
 => if (lenB >= lenBufSeVen) then
          if p1 then ConfAcc myTick(SchAcc)
                else ConfRej myTick(SchRej)
          fi
    else ConfAcc2 mytick(SchAcc2)
    fi
if p1 := sampleBerWithP(accept-prob(pmod))
 /\ { ActorDr, bufDr } := select(gt,B)
 /\ nBuf := add([lenB + (- 1) | bufDr], < Actor invite gt >)
 /\ ConfAcc := <name: Ser | req-cnt: (pmod + 1.0), b-set: nBuf , AS >
 /\ SchAcc := insert([gt | SL],
                      {gt, Actor <- TRYING} ; {gt, ActorDr <- poll})</pre>
 /\ ConfRej := <name: Ser | req-cnt: (pmod + 1.0), b-set: [lenB | B], AS >
 /\ SchRej := insert([gt | SL],
                      {gt + delay , Actor <- poll})</pre>
 /\ b-setNu := add( [lenB | B], < Actor invite gt > )
 /\ ConfAcc2 := <name: Ser | req-cnt: pmod , b-set: b-setNu, AS >
 /\ SchAcc2 := insert([gt | SL], {gt + delay, Actor <- TRYING}) .</pre>
rl [SeVen-APP-ROUND] :
     <name: Ser | req-cnt: pmod , b-set: [lenB | B], AS >
     {Ser <- ROUND} [gt | SL]
 =>
     <name: Ser | req-cnt: 0.0, b-set: [lenB | B], AS >
     mytick(insert([gt | SL],
               {gt, reply(Ser, B, gt)} {gt + Ts, Ser <- ROUND})) .</pre>
```

Figure 3: Rewrite rules specifying SeVen's selective strategy.

a poll message to ActorDr and a TRYING to Actor are created. Otherwise,
the incoming request is rejected and pmod is incremented without affecting the
server's buffer resulting in the configuration ConfRej. A poll message to Actor
is also created.

The rule SeVen-APP-ROUND specifies that when the round finishes, all surviving requests in the server's buffer are answered, where a new round starts and pmod is re-set.

441 5. Simulations

We detail our simulation results obtained from our formal specification using
the statistical model checker PVeStA [8]. Our simulations are parametric in the
following values:

- Average time of a call $-t_M$: This is the assumed average time of the calls of honest users. For the simulations, we assumed $t_M = 5$ time units;
- **Dropping Factor** We assume the following values for the dropping factor function (Equation 1):
- 445 runetion (12 449 $- p_{\rm IN} = 2;$
 - $p_{\rm IN} = 2$,

450

451

- $-p_{\mathsf{WAIT}} = 8;$
- $-\alpha = 1.89.$
- These values were chosen so that the dropping factor increases in a reasonable fashion for calls with duration greater than t_M . Sample values are shown below, recalling that $t_M = 5$ seconds:

	Call Duration (mins)	Dropping Factor			
455	6	12.37			
	8	17.31			
	10	27.84			

That is, dropping factor of a call with duration of 10 minutes, *i.e.*, $2 \times t_M$, is approximately 3 times greater than the dropping factor of a call whose status is WAITING (27.84/8). This is a reasonable ratio. However, according to the specific application other values can be set for p_{IN} , p_{WAIT} and α . Finally, the choice of setting $p_{\text{WAIT}} = 4 \times p_{\text{IN}}$ was selected so that the calls with duration less than t_M have much less chance of being dropped than the calls that are still waiting for the responder.

- Size of Buffer -k: This is the upper-bound on the size of the server's buffer denoting the processing capacity of the application. k = 24;
- Rate of Calls (R): We fixed the total rate of calls to be R which is the result of summing the rate of legitimate calls, R_L , and the rate of attacker calls, R_A . That is, $R_L + R_A = R$.

The value of R is computed using standard techniques⁶ so that if $R_L = R$, *i.e.*, the server is not under attack, then the server will not be overloaded. With R fixed, we set R_L and R_A to be $R_L + R_A = R$, but considered scenarios with different proportions for R_A and R_L . This reflects the fact that Coordinated Call attack uses low traffic and therefore, can bypass usual defenses based on network traffic analysis which normally monitor the total rate of calls R.

We considered 5 different proportions for R_L and R_A expressed in the percentage of the total number of calls R:

⁶using the Erlang model which computes R by taking into account k and t_M

Legitimate Calls (R_L)	Attacker Calls (R_A)
83%	17%
67%	33%
50%	50%
33%	67%
17%	83%

• Total time of the simulation - *total*: This is the total time of the simulation using PVeStA. We used in our simulations *total* equal to 40 time units, similar to the time used in [24];

• Delay of the Network: We also assumed a delay of 0.1 time units of message in the network;

Degree of confidence for the simulation: Our simulations were carried out with a degree of confidence of 99% (see [27, 7] for more details on probabilistic model checking).

Quality Measures. In our simulation, we use quality measures specific for VoIP 486 services. These are specified by expressions of the QuaTEx quantitative, prob-487 abilistic temporal logic defined in [27]. We perform statistical model checking 488 of our defense in the sense of [7]: once a QuaTEx formula and desired degree of 489 confidence are specified, a sufficiently large number of Monte Carlo simulations 490 are carried out allowing for the verification of the QuaTEx formula. The Monte 491 Carlo simulations are carried out by the computational tool Maude [6] and the 492 statistical model checking is carried out by PVeStA. 493

The QuaTEx formulas, *i.e.*, the quality measures, that we use in our simulations are defined below. The operator \bigcirc is a temporal modality that specifies the advancement of the global time to the time of the next event (see [27] for more details).

• Complete: How many honest calls were able to stay in the INCALL status for the expected duration.

$$complete(total) =$$
 if $time > total$ then $\frac{countComplete}{countHonest}$
else $\bigcirc complete(total)$

where *countComplete* is a counter that is incremented whenever an honest call is completed.

• Incomplete: How many honest calls were able to have the INCALL status but were dropped before completing the call, *i.e.* not staying in INCALL status for the expected duration;

$$incomplete(total) = if time > total then \frac{count momplete}{count Honest}$$

$$else \bigcirc incomplete(total)$$

where countIncomplete = countIncall - countComplete and countIncall is

a counter that is incremented whenever an honest calls changes from status
 WAITING to INCALL.

• Unsuccessful: How many honest calls were not even able to reach the INCALL status. That is, how many calls were not even able to start talking

477

500



Figure 4: Client Success Ratio: Simulation Results when not using SeVen.

512	between each other.
	$unsuccessful(total) = $ if $time > total$ then $\frac{countUnsuccessful}{countHonest}$
513	$else \bigcirc unsuccessful(total)$
514	where $countUnsuccessful = countHonest - countIncall$.
515	• The average of client incomplete calls: We also measure the average propor-
516	tion of time legitimate clients were able to talk in an incomplete call before
517	they were dropped.
	$avgInCall(total) = $ if $time > total$ then $\frac{totalTimeInCall}{totalIncompleteCall}$
518	else $\bigcirc avgInCall(total)$

where *totalTimeInCall* is the sum of how much percent of time clients 519 were able to talk before being interrupted and the totalIncompleteCall is 520 the total of clients the were not able to finish their call. 521

We carried out simulations with the three different types of dropping strate-522 gies described in Section 3, namely uniform, roulette and $\frac{k}{2}$ -tournament. We 523 also carried out simulations with a scenario without SeVen. 524

5.1. No Defense 525

Our simulations results are depicted in Figure 4. They suggest that the 526 Coordinated Call attack is indeed effective in reducing the availability of a 527 VoIP service when assuming both an exponential and a lognormal call dura-528 tion. Increasing the proportion of attackers rapidly increases the proportion of 529 Unsuccessful calls, *i.e.*, calls that did not even start a conversation, while the 530 proportion of Complete calls falls. As expected there are no Incomplete calls as 531 the VoIP server does not interrupt calls. 532

5.2. Uniform Dropping Strategy 533



Figure 5: Client Success Ratio: Simulation Results when a uniform dropping strategy.

Figure 5 depicts the availability results when using SeVen with a uniform dropping strategy. It suggest that SeVen can indeed mitigate the Coordinated Call attack. The proportion of Complete calls remains at high levels when assuming both an exponential, above 60% of legitimate calls, and a lognormal call duration, above 80%.



Figure 6: Average Time of Incomplete Calls: Simulation Results when using SeVen with a uniform dropping strategy.

As SeVen may interrupt calls in the middle of a conversation, there are In-539 complete calls, *i.e.*, calls where the parties have started to communicate, but 540 were interrupted before communicating for the expected time. For exponential 541 call duration, around 30% of legitimate calls were interrupted, while for lognor-542 mal call duration around 10% of legitimate calls were interrupted. However, 543 the average time of incomplete calls depicted in Figure 6 suggests that although 544 these calls are interrupted, they still are able to communicate for long periods 545 of time, above 70% of the expected time for exponential call duration and above 546 89% of the expected time for lognormal call duration. 547

548 5.3. Roulette Dropping Strategy

Figures 7 and 8 depict our simulation results when using a roulette dropping strategy. The results are similar to the results obtained with the uniform dropping strategy. For exponential call duration, above 60% of legitimate calls



Figure 7: Client Success Ratio: Simulation Results when a roulette dropping strategy.

were completed and around 30% were interrupted by SeVen. For lognormal call duration, above 80% of legitimate calls were completed and around 10% were interrupted by SeVen. Moreover, as depicted in Figure 8, the interrupted calls stayed communicating in average above 70% of the expected call duration for exponential call duration and above 88% of the expected call duration for lognormal call duration.



Figure 8: Average Time of Incomplete Calls: Simulation Results when using SeVen with a roulette dropping strategy.

558 5.4. $\frac{k}{2}$ -tournament

Figures 9 and 10 depict the simulation results obtained by using a $\frac{k}{2}$ -tournament 559 dropping strategy. They suggest that this strategy is better than the roulette 560 and uniform strategies. The proportion of complete calls is above 62% for expo-561 nential call duration and above 86% for lognormal call duration. In other words, 562 around 30% of legitimate calls were interrupted by SeVen when assuming ex-563 ponential call duration and around 10% of legitimate calls were interrupted by 564 SeVen when assuming lognormal call duration. Moreover, the average time of 565 incomplete calls is always greater than 72% for exponential call duration and 566 above 84% for lognormal call duration. 567



Figure 9: Client Success Ratio: Simulation Results when a roulette dropping strategy.



Figure 10: Average Time of Incomplete Calls: Simulation Results when using SeVen with a tournament dropping strategy.

568 6. Experiments

In our experiments, we used Asterisk version 13.6.0 which is a SIP server 569 widely used by small and mid size companies for implementing their VoIP ser-570 vices. We assume there are honest users and malicious attackers which try to 571 make the VoIP unavailable. Both the traffic of the honest users and the attack-572 ers are emulated using the tool SIPp [20] version 3.4.1. SIPp generates calls 573 which may be configured as the caller or the callee. Thus, in our experiments, 574 we used pairs of SIPp, one pair for generating the honest user calls and the 575 other pair for generating the attacker calls. Finally, we developed the SeVen 576 proxy in C++ which implements the selective strategy described in Section 3. 577



578

The figure above illustrates the topology of the experiments we carried out. 579 To make a call, the pairs of SIPp send messages to the SeVen proxy which on 580 the other hand forwards them to Asterisk. Similarly, any message generated 581 by Asterisk is forward to the SeVen proxy which then forwards them to the 582 corresponding users. Therefore, SeVen is acting as an *Outbound Proxy* for both 583 Asterisk and the pairs of SIPp. For our experiments, it is enough to use a single 584 machine. We used a machine with configuration Intel(R) Core(TM) i7-4510U 585 CPU @ 2.00GHz and 8 GB of RAM. 586

⁵⁸⁷ *Parameters.* We use the following parameters to configure our experiments:

- Average Call Duration (t_M) We assume known what is the average duration of calls. This can be determined in practice by analyzing the history of calls. We assume in our experiments that $t_M = 160$ seconds (approximately 2.6 minutes);
- SIP Sever Capacity (k) This is the number of simultaneous calls the SIP server can handle. We set k = 200 which is a realistic capacity for a small company allowing 400 users (2×200) to use the service at the same time.
- Experiment Total Time (T) Each one of our experiments had a duration of 60 minutes which corresponds to 3600 calls in each experiment. With this duration, it was already possible to witness the damage caused by the Coordinated Call Attack as well as the efficiency of our solution for mitigating this attack.
- Traffic Rate (R) Using a server with capacity of k, we calculated using standard techniques [23] what would be a typical traffic of such a server. It is R = 60 calls per minute. This value is computed using traditional techniques [23] (Erlang model) taking into account k = 200 and $t_M = 160$. Thus the service can handle R legitimate calls per second. However, as the attacker does not behave as legitimate placing calls with much greater durations, the server can be subject to this attack.
- In our experiments, we split this rate among clients and attackers. This is because we want to emulate the fact that Coordinated Call attack can deny

service using low traffic thus bypassing usual defenses [28, 29, 30, 31, 32, 33] 610

based on network traffic analysis or monitoring the number of incoming 611

calls. This means that the total traffic (attacker + client) is always less or 612 equal than R.

613

The following graph illustrates client usage in normal conditions, *i.e.*, with-614 out suffering an attack, using the parameters as specified above. 615



616

It shows that the server is indeed well dimensioned for this rate of (legitimate) 617 calls. Clients occupy in average approximately 85% of Asterisk's capacity thus 618 not overloading the server, but still being heavily used. 619

Finally, we set the duration of the calls generated by SIPp as follows: 620

• Total Call Duration of Clients: As described in Section 3, we used two 621 models for the call duration of legitimate client calls, namely the exponential 622 model suitable for non-VoIP calls and the lognormal model suitable for 623 VoIP calls. The parameters, λ, μ, σ , in Equations 2 and 3 were computed 624 as described in [11, 12] using the average call duration t_M . Whenever we 625 generate a new call, we generate randomly according to the used model 626 (exponential or lognormal) the call duration. SIPp ends the call when its corresponding call duration is reached. 628

Call Duration of Attackers: Following the Coordinated Call Attack, we 629 do not limit the call duration of an attacker call. His call communicate for 630 indefinite time. 631

Quality Measures. For our experiments, we used the following three quality 632 measures for our calls: 633

- Complete Call: A call is complete whenever its status changed from 634 WAITING to INCALL and it is able to stay in status INCALL for its cor-635 responding call duration. That is, the caller was able to communicate with 636 the responder for all the prescribed duration; 637
- Incomplete Call: A call is incomplete whenever its status changed from 638 WAITING to INCALL, but it was not able to stay in status INCALL for its 639 corresponding call duration. That is, the caller was interrupted before completing the call; 641
- **Unsuccessful Call:** A call is unsuccessful if it did not even change its 642 status from WAITING to INCALL. That is, the caller did not even have the 643 chance to speak with the responder. 644



Figure 11: Client Success Ratio: Experimental Results when not using SeVen.

Intuitively, complete calls are better than incomplete calls which are better than unsuccessful calls. In order to support this claim, we also computed the average duration call of the incomplete calls, that is, the time that users in average were able to stay communicating before they were interrupted by SeVen.

649 6.1. Experimental Results

We carried out the corresponding experiments to the scenarios used in Section 5. That is, we tested the efficiency of the Coordinated Call attack when the server is not running any defense. We also carried out experiments with scenarios using an exponential and lognormal call duration with the uniform, roulette and 100-tournament dropping strategies.

655 6.1.1. No Defense

Figures 11 and 12 illustrate our main results when assuming exponential 656 and lognormal call duration and not using any defense mechanism. Our results 657 demonstrate the efficiency of the Coordinated Call attack. We observed that 658 the VoIP availability decreases considerably when increasing the proportion of 659 attacker calls in the rate R (Figure 11). In particular, the number of unsuccessful 660 call increases to level near 100%, while the number of completed calls falls to 661 close to 0%. Moreover, there are no incomplete calls, which is expected since 662 no calls are interrupted. 663

We also measured the number of attacker calls that the server serves during the experiment (Figure 12). As expected from the profile of the Coordinated Call attack, the attacker is able to deny service by slowly (after 10 minutes) occupying all the available calls in the server and therefore deny its service to legitimate clients.



Figure 12: Attacker Call Occupancy in Buffer: Experimental Results when not using SeVen.

669 6.1.2. Uniform Defense



Figure 13: Client Success Ratio: Experimental Results when SeVen and a uniform dropping strategy.

We carried out experiments to test the efficiency of SeVen when using a uniform dropping strategy. Our results are summarized in Figures 13, 14 and 15.

The graphs depicted in Figure 13 show that SeVen when using a uniform 673 dropping strategy can mitigate the Coordinated Call attack. The results assum-674 ing a lognormal call duration is slightly better than the results when assuming 675 an exponential call duration. In both cases, the proportion of completed calls 676 stayed above 50% levels even when the attacker call rate is 5 times more than 677 the client call rate. The proportion of incomplete calls was more affected by 678 the call duration model. Under the exponential call duration assumption, the 679 proportion of incomplete call was around 35% of all legitimate calls, while under 680 the lognormal call duration assumption, the proportion of incomplete calls was 681 of around 28%. The proportion of unsuccessful calls stays below 10% under 682 both assumptions of call duration. 683

Figure 15 illustrates how the attacker is able to occupy the resources of the server. While when not running SeVen the attacker was able to occupy all the server's resources (Figure 12), when using SeVen with a uniform dropping strategy, the attacker is only able to occupy around 70% of the server's resources.



Figure 14: Average Time of Incomplete Calls: Experimental Results when using SeVen with a uniform dropping strategy.



Figure 15: Attacker Call Occupancy in Buffer: Experimental Results when using SeVen with a uniform dropping strategy.

This may seem to be high value, but the graph hides the fact that attackers are 688 dropped by the defense strategy and thus the high levels of availability obtained. 689 Finally, we also measured the average time of incomplete calls, that is, the 690 proportion of time that incomplete calls were able to stay in a call before they 691 were dropped by the SeVen defense strategy. The values are depicted in Fig-692 ure 14. When assuming both an exponential call duration and a lognormal call 693 duration, the incomplete call were in average around 40% of the expected call 694 time. 695

696 6.1.3. Roulette Defense

Figures 16, 17 and 18 depict our experimental results when using SeVen with the roulette dropping strategy. As with the uniform dropping strategy, the defense was able to mitigate the Coordinated Call attack. Furthermore, the roulette strategy performed better than the uniform strategy.

The availability depicted in Figure 16 remained at high levels under both assumptions on call duration (exponential and lognormal). The proportion of completed calls was above 70% percent when assuming an exponential call duration and above 75% when assuming a lognormal call duration. In both cases, the proportion of incomplete calls was less than 6%.



Figure 16: Client Success Ratio: Experimental Results when using SeVen with a roulette dropping strategy.



Figure 17: Average Time of Incomplete Calls: Experimental Results when using SeVen with a roulette dropping strategy.

Despite the availability results using the roulette strategy being better than 706 the availability results obtained using the uniform strategy, the attacker was 707 still able to occupy a similar proportion of the server's resource as depicted 708 in Figure 18. It occupied at most 70% of the server's resources. Intuitively, 709 the difference in the availability between the uniform and roulette strategies is 710 because the roulette strategy tends to drop calls with greater duration. This 711 means that the attacker calls are more likely to be selected leaving more chance 712 for a legitimate call to access the service. 713

Finally, we also measured the average time of incomplete calls. These results are depicted in Figure 17. They show that these calls were able to communicate for more than 50% of the expected time. These results are better than the results obtained using a uniform dropping strategy.



Figure 18: Attacker Call Occupancy in Buffer: Experimental Results when using SeVen with a roulette dropping strategy.

718 6.1.4. 100-Tournament Defense



Figure 19: Client Success Ratio: Experimental Results when using SeVen with a 100-tournament dropping strategy.

Our last set of experiments evaluated the efficiency of SeVen with a 100tournament dropping strategy. Figures 19, 20 and 21 depict our main results. The 100-tournament dropping strategy resulted in the best results when compared with the uniform and roulette strategies.

The attacker was still able to use roughly the same amount of resources of the 723 server as when using the uniform and roulette strategy, namely around 70% of its 724 resources as depicted in Figure 18. Moreover, the availability results depicted 725 in Figure 19 were slightly better than the results obtained with the roulette 726 strategy (Figure 16). In both assumptions of call duration (following an expo-727 nential and a lognormal distributions). SeVen was able to maintain high levels 728 of availability. More than 70% (respectively, 80%) of calls were completed when 729 assuming call duration following an exponential distribution (respectively, log-730 normal distribution). The proportion of incomplete calls reached levels around 731 25% of all calls when assuming an exponential cal duration and reached 16%732 of all calls when assuming lognormal call duration. Thus, more than 95% of all 733 legitimate calls were able to reach the incall status, which means that they were 734 able to communicate. 735



Figure 20: Average Time of Incomplete Calls: Experimental Results when using SeVen with a 100-tournament dropping strategy.



Figure 21: Attacker Call Occupancy in Buffer: Experimental Results when using SeVen with a 100-tournament dropping strategy.

Finally, as depicted in Figure 20, incomplete calls were interrupted by SeVen
 after communicating more than 60% of the expected time when assuming expo nential call duration and more than 50% when assuming lognormal call duration.

739 6.2. Impact of Using SeVen When not Suffering an Attack

We investigate additionally the impact of using SeVen as a proxy and implementing the selective strategy on the performance of Asterisk. Asterisk modules
provide many statistics on the performance of the system. The following tables
contains the number of request according to the time intervals for the time to
respond (TTR) when using and not using SeVen during normal situation, that
is, when only receiving legitimate calls.

	Not Using SeVen					
746	TTR (ms)	[0,1]	[3, 4]	$[4,\!5]$	[8, 9]	[10,20]
	Num Requests	1837	1	960	2	2

	Using SeVen											
748	TTR (ms)	[0,1]	$[3,\!4]$	[4, 5]	$[7,\!8]$	[10, 20]	[20, 30]	[30, 40]	[40, 50]	[50, 100]	[100, 150]	
	Num Requests	19	2	404	341	657	434	148	96	70	8	

As one can observe, there is an impact to TTR when using SeVen. While without SeVen most of the requests are responded within 5ms, with SeVen, most of the requests are responded within 100ms.

⁷⁵²Such a delay does not greatly impact user experience as 100 ms is negligible ⁷⁵³with respect to the time users wait until establishing a call, *e.g.*, waiting until ⁷⁵⁴Bob accepts the call which normally takes some seconds to happen. Moreover, ⁷⁵⁵as SeVen only acts on SIP messages (Initiation and Termination phases of Fig-⁷⁵⁶ure 1), SeVen does not affect user experience during when the parties are in a ⁷⁵⁷call (Communication Phase in Figure 1).

Finally, it seems possible to improve SeVen's performance by improving our implementation, *e.g.*, implementing it as a module instead of a proxy. This is left, however, to future work.

761 7. Comparison between Simulation and Experimental Results

762 7.1. Differences between Simulations and Experiments

There are some important differences between the formal model and the experimental set-up. For a starter, the formal specification abstracts several aspects present in the experimental set-up. For example:

- We did not model how Asterisk actually manages its workers/threads. As
 Asterisk has a number of modules that among other things, maintain call
 statistics, convert calls encoded some codex to another, etc;
- In our experiments, there are other applications running in parallel with
 Asterisk that have to compete for resources (CPU and network interface for
 example). Our formal model does not incorporate this;

• We use a simplified model for network latency with constant latency time;

While the parameters, e.g., k, incoming client and attacker traffic, used in the simulations were proportional to the parameters used in the experiments, they were much lower to the ones used in the experimental results. If we used the actual values for these parameters, simulations would have taken much longer;

• In our experiments, SeVen is used as a proxy, while in our formal model the defense was incorporated into the application.

Despite these important differences/abstractions, as we compare in more detail next, the simulation results corresponded to many of the experimental results. For example in terms of availability, *i.e.*, proportion of completed, incomplete and unsuccesful calls, the simulation results correctly indicated the power of the attack and the efficiency of SeVen mitigating this attack. They also correctly indicated which dropping strategy is better and how availability changes with the increase on the attack rate. The simulation results were less

747

accurate in predicting the time of incomplete calls reaching a difference of 30%.
It is not completely clear the reasons for this discrepancy, but we believe it has
to do with the abstractions mentioned above. We leave this investigation to
future work

791 7.2. Detailed Comparisons

Typically each simulation takes about 30 seconds to be completed. In contrast, each experiment carried out on the network took 60 minutes. This means that specifiers can quickly test different selective strategies using formal verification before implementing the necessary machinery and carrying out experiments. Once a selective strategy is shown by formal verification to have reasonable results, experiments can be carried out to validate the chosen defense.

In this section, we compare the results obtained through formal verification detailed in Section 5 and the results obtained by carrying out experiments detailed in Section 6. In general, availability results obtained through formal verification indeed corresponded to our experimental results showing a high degree of accuracy for our simulation results.

Efficiency of the Coordinated Call Attack:. Our simulation results with the scenario without defense showed that the Coordinated Call is effective if the server does not have any defense mechanism and in both assumptions of duration calls (exponential and lognormal). This result was also observed in our experimental results. Figure 22 on the service availability illustrate this correspondence.



Figure 22: Client Success Ratio: Comparison between simulation and experimental results for lognormal call duration.

Efficiency of SeVen: All our simulations results using scenarios with SeVen indicated that SeVen is indeed a good defense for mitigating the Coordinated Call attack. The greater proportion of calls were completed calls, while a smaller proportion of the calls were incomplete and a minority of calls were unsuccessful. The same behavior was observed by our experiments. This is illustrated by Figure 23.



Figure 23: Client Success Ratio: Comparison between simulation and experimental results for lognormal call duration when using SeVen with a tournament dropping strategy.

Dropping Strategies Evaluation:. Our simulations were able to predict that the
tournament dropping strategy would perform best. However, it was not able
to forecast that the roulette strategy would perform better than the uniform
strategy. It is not clear to us why this was the case.

Better Performance for Lognormal Call Duration:. Our simulations results also predicted that selective strategies would perform better in scenarios where call duration of legitimate clients follows a lognormal distribution, such as in VoIP communication, than scenarios with exponential call duration, such as in traditional telephony. This was the case independent on the dropping strategy used (uniform, roulette or $\frac{k}{2}$ -tournament). The same behavior was observed in our experimental results.

Average time of Incomplete Calls:. Our simulations results also indicated that the average time that incomplete calls stayed communicating before being dropped by SeVen was relatively high: more than 80% of the expected time when using the tournament strategy under a lognormal call duration. This results diverged from the experimental results which observed an average time of incomplete calls of around 60%. This is illustrated by Figure 24.



Figure 24: Average Time of Incomplete Calls: Comparison between simulation and experimental results when using SeVen with the $\frac{k}{2}$ -tournament dropping strategy.

However, for other scenarios, the difference between simulation and experi-831 mental results was greater reaching 30% of difference. It is not clear to us what 832 are the reasons for this difference. We suspect that the modeling of the time 833 delays should be improved. In any case, our results suggest that while simula-834 tions provide quite accurate results on availability, it is less accurate on specific 835 timing analysis. We observed a similar behavior during our experiments and 836 simulations when modeling our defense for mitigating Application-Layer DDoS 837 attacks [4]. 838

839 8. Related and Future Work

This paper formalized a new selective defense, called SeVen, for mitigating 840 Coordinated Call attacks. We have shown that using state-dependent proba-841 bility distributions for selecting which calls are to be processed results in high 842 levels of availability. We proposed three defenses based on the dropping strat-843 egy method (uniform, roulette and $\frac{k}{2}$ -tournament). We carried out simulations 844 and experiments assuming traditional telephony and VoIP communications. In 845 both cases, we observed that our SeVen was able to mitigate the Coordinated 846 Call attack. Finally, we compared the results obtained using our formal analysis 847 with the results obtained by experimentation obtaining a high accuracy. This 848 means further supports the value of formal analysis during the development of 849 selective defenses for mitigating DoS attacks. 850

Most of the existing work [28, 29, 30, 31, 32, 33] on mitigating DoS attacks 851 on VoIP services focuses on flooding attacks, such as SIP-Flooding attack. They 852 analyze the network traffic and whenever they observe an abrupt increase in the 853 traffic load, they activate their defenses. The network traffic is usually modeled 854 using some statistical approach, such as correlating the number of INVITE 855 requests and the number of requests that completed the SIP initiation phase [28] 856 or using more complicated metrics such as Helling distance to monitor traffic 857 probability distributions [29, 30, 31]. Other solutions place a lower priority on 858 INVITE messages, which are only processed when there are no other types of 859 request to be processed [32, 33]. 860

As the Coordinated Call Attack emulates legitimate client traffic not causing 861 an unexpected sudden increase in traffic, all these defenses are not effective in 862 mitigating the Coordinated Call Attack. The few solutions we found in the 863 literature for this type of attack are commercial tools that act as a firewall 864 which monitors all the call traffic and the signaling [34, 35] or analyze audio 865 samples [36] in order to differentiate the fraudulent calls from the legitimate 866 ones. Less sophisticated mechanisms [37] monitors all the incoming requests 867 and rejects those whose IPs do not belong to a list of trusted IPs. Clearly 868 such approaches does not work well when the attackers are malicious users 869 870 whose IPs are in the trusted list and are not using automation to make the calls. In addition, these commercial tools can be expensive for small businesses 871 to purchase and maintain, and they do require technical expertise for proper 872 installation. 873

One main advantage of our proposed solution is that it is not tailored using many specific assumptions on type of service. The only assumption used is a previous knowledge of the average call duration, which can be easily inferred from the service call history. Moreover, our solution can be easily integrated with other mechanisms such as the IP filtering approach used in [37].

[38] proposes a filtering mechanism for SIP flooding attacks. It is not clear whether such mechanisms will be enough for mitigating the Coordinated VoIP attack, as the number of messages needed to carry out such attack is much less. Wu *et al.* [39] have proposed mechanism to identify intruders using SIP by analyzing the traffic data. Although we do not tackle the identification of intruders problem, we find it an interesting future direction.

The formalization of DDoS attacks and their defenses has been subject of 885 other papers. For example, Meadows proposed a cost based model in [40], while 886 others use branching temporal logics [41]. This paper takes the approach used in 887 [42, 25, 24], where one formalizes the system in Maude and uses the Statistical 888 Model Checker PVeStA to carry out analyses. While [42, 25, 24] modeled tra-889 ditional DDoS attacks exploiting stateless protocols on the transport/network 890 layers, we are modeling stateful Application Layer DDoS attacks. Moreover, 891 the quality measures used for VoIP services under TDoS attacks, described in 892 Section 3, are different to the quality measures considered in the previous work. 893 More recently [4], we proposed SeVen showing that it can be used to mitigate 894 ADDoS attacks that exploit the HTTP protocol. This paper shows that SeVen 895 can also be used to mitigate DDoS attacks in VoIP protocols, but in order to 896 do so one needs state-dependent probabilistic distributions. This is because of 897 the quality requirements that we need in VoIP communications. We would like 898 to give a priority to the types of call that should be given more chances to keep 899 using resources of the server. In particular, we give preference to calls that do 900 not take more than the average duration time. Such quality measures are not 901 present in HTTP protocols that we analyzed in [4]. 902

For future work, we are going to investigate the mitigating of other types 903 of attacks, such as volumetric and amplification attacks. We are also thinking 904 on intrusion detection mechanisms. We are also interested in building defenses 905 for mitigating amplification attacks [18]. We have also been using SeVen for 906 mitigating High-Rate ADDoS attacks using Software Defined Networks [43]. 907 We are also investigating ways to improve simulation accuracy by improving 908 the modeling of timing aspects of the system, such as processing and network 909 delay. 910

Acknowledgments. This work has been funded by the DFG as part of the project
 Secure Refinement of Cryptographic Algorithms (E3) within the CRC 1119
 CROSSING, by RNP project GT-ACTIONS, by Capes Science without Borders
 grant 88881.030357/2013-01 and CNPq.

915 **References**

916 [1] Cyber Threat Bulletin: Boston Hospital TDoS Attack, http://voipsecurityblog.typepad.

917 918		com/files/cyber-threat-bulletin-13-06-boston-hospital-telephony-denial-of-service-attack.pdf, accessed: 2015-27-09.
919 920 921	[2]	TDoS – Extortionists Jam Phone Lines of Public Services Including Hospitals, https://nakedsecurity.sophos.com/pt/2014/01/22/ tdos-extortionists-jam-phone-lines-of-public-services-including-hospitals/, accessed: 2015-27-09.
922 923 924	[3]	Situational Advisory: Recent Telephony Denial of Services (TDoS) Attacks, http://voipsecurityblog.typepad.com/files/ky-fusion_tdos_3-29-13-2.pdf/, accessed: 2015-27-09.
925 926	[4]	Y. G. Dantas, V. Nigam, I. E. Fonseca, A Selective Defense for Application Layer DDoS Attacks, in: JISIC 2014, 2014, pp. 75–82.
927 928 929	[5]	The Surging Threat of Telephony Denial of Service Attacks, http://voipsecurityblog.typepad.com/files/tdos_paper_4-11-13.pdf (accessed: 2015-09-28).
930 931 932	[6]	M. Clavel, F. Durán, S. Eker, P. Lincoln, N. Martí-Oliet, J. Meseguer, C. Talcott, All About Maude: A High-Performance Logical Framework, Vol. 4350 of LNCS, Springer, 2007.
933 934	[7]	K. Sen, M. Viswanathan, G. Agha, On Statistical Model Checking of Stochastic Systems, in: CAV, 2005, pp. 266–280.
935 936	[8]	M. AlTurki, J. Meseguer, PVeStA: A Parallel Statistical Model Checking and Quantitative Analysis Tool, in: CALCO, 2011, pp. 386–392.
937 938	[9]	A. Lipowski, D. Lipowska, Roulette-wheel selection via stochastic acceptance, CoRR abs/1109.3627.
939 940 941	[10]	T. Blickle, L. Thiele, A Mathematical Analysis of Tournament Selection, in: Proceedings of the 6th International Conference on Genetic Algorithms, San Francisco, CA, USA, 1995, pp. 9–16.
942 943 944	[11]	L. Brown, N. Gans, A. Mandelbaum, A. Sakov, H. Shen, S. Zeltyn, L. Zhao, Statistical Analysis of a Telephone Call Center: A Queueing-Science Per- spective, Journal of the American Statistical Association 100 (2005) 36–50.
945 946 947	[12]	P. Galiotos, T. Dagiuklas, S. Kotsopoulos, Call-Level VoIP Traffic Modelling Based on Data from a Real-Life VoIP Service Provider, in: 2015 IEEE Globecom Workshops (GC Wkshps), 2015, pp. 1–7.
948 949	[13]	I. Digium, "SAsterisk Private Branch eXchange", http://www.asterisk. org/, accessed: 2015-09-28 (2015).
950 951 952 953	[14]	M. O. O. Lemos, Y. G. Dantas, I. Fonseca, V. Nigam, G. Sampaio, A Selective Defense for Mitigating Coordinated Call Attacks, in: 34th Brazilian Symposium on Computer Networks and Distributed Systems (SBRC), 2016.

34

- Y. G. Dantas, M. O. O. Lemos, I. Fonseca, V. Nigam, Formal Specification and Verification of a Selective Defense for TDoS Attacks, in: 11th International Workshop on Rewriting Logic and its Applications (WRLA), 2016.
- ⁹⁵⁸ [16] SeVen https://github.com/ygdantas/SeVen.git (2016).
- ⁹⁵⁹ [17] Session Initiation Protocol, http://www.ietf.org/rfc/rfc3261.txt.
- [18] R. Shankesi, M. AlTurki, R. Sasse, C. A. Gunter, J. Meseguer, Model Checking DoS Amplification for VoIP Session Initiation, in: ESORICS,
 2009, pp. 390–405.
- [19] M. Hines, Attackers Get Chatty on VOIP, http://www.pcworld.com/
 article/132389/article.html, accessed: 2015-27-09 (2007).
- R. Gayraud, O. Jacques, SIPp SIP Traffic Generator, http://sipp.
 sourceforge.net, accessed: 2015-09-28 (2014).
- ⁹⁶⁷ [21] S. Khanna, S. Venkatesh, O. Fatemieh, F. Khan, C. Gunter, Adaptive
 ⁹⁶⁸ Selective Verification: An Efficient Adaptive Countermeasure to Thwart
 ⁹⁶⁹ DoS Attacks, Networking, IEEE/ACM Transactions on 20 (3) (2012) 715–
 ⁹⁷⁰ 728.
- [22] S. Khanna, S. S. Venkatesh, O. Fatemieh, F. Khan, C. A. Gunter, Adaptive
 Selective Verification, in: INFOCOM, 2008, pp. 529–537.
- J. Jewett, J. Shrago, B. Yomtov, Designing Optimal Voice Networks for
 Businesses, Government, and Telephone Companies., 1980.
- J. Eckhardt, T. Mühlbauer, M. AlTurki, J. Meseguer, M. Wirsing, Stable
 Availability under Denial of Service Attacks through Formal Patterns, in:
 FASE, 2012, pp. 78–93.
- J. Eckhardt, T. Mühlbauer, J. Meseguer, M. Wirsing, Statistical Model
 Checking for Composite Actor Systems, in: WADT, 2012, pp. 143–160.
- ⁹⁸⁰ [26] J. Meseguer, Twenty years of rewriting logic, J. Log. Algebr. Program.
 ⁹⁸¹ 81 (7-8) (2012) 721-781.
- ⁹⁸² [27] G. Agha, J. Meseguer, K. Sen, PMaude: Rewrite-based Specification Lan⁹⁸³ guage for Probabilistic Object Systems, Electron. Notes Theor. Comput.
 ⁹⁸⁴ Sci. 153 (2) (2006) 213–239.
- [28] D.-Y. Ha, H.-K. Kim, K.-H. Ko, C.-Y. Lee, J.-W. Kim, H.-C. Jeong, Design and Implementation of SIP-aware DDoS Attack Detection System, in: ICIS
 '09, ACM, New York, NY, USA, 2009, pp. 1167–1171.
- J. Tang, Y. Cheng, C. Zhou, Sketch-Based SIP Flooding Detection Us ing Hellinger Distance, in: Global Telecommunications Conference, 2009.
 GLOBECOM 2009. IEEE, 2009, pp. 1–6.

- [30] J. Tang, Y. Cheng, Y. Hao, Detection and prevention of SIP flooding attacks in voice over IP networks, in: INFOCOM, 2012 Proceedings IEEE, 2012, pp. 1161–1169.
- J. Tang, Y. Cheng, Y. Hao, W. Song, SIP Flooding Attack Detection with
 a Multi-Dimensional Sketch Design, Dependable and Secure Computing,
 IEEE Transactions on 11 (6) (2014) 582–595.
- [32] X.-Y. Wan, Z. Li, Z.-F. Fan, A SIP DoS flooding attack defense mechanism
 based on priority class queue, in: WCNIS 2010, 2010, pp. 428–431.
- [33] F. Zi-Fu, Y. Jun-Rong, W. Xiao-Yu, A SIP DoS Flooding Attack Defense
 Mechanism Based on Custom Weighted Fair Queue Scheduling, in: ICMT
 2010, 2010, pp. 1–4.
- [34] SecureLogix: Telephony Denial of Service (TDoS) Solutions, http://www.
 securelogix.com/solutions/telephony-denial-of-service-TDoS.
 html, accessed: 2015-27-09.
- 1005[35] TransNexusNexOSS,http://transnexus.com/1006telephony-denial-service-attacks/, accessed: 2015-27-09.
- [36] PINDROP: Protecting your Call Centers Against Phone Fraud & Social Engineering, https://www.pindrop.com/wp-content/uploads/2016/ 01/pindrop_overview_whitepaper_fi_20141121_v2.pdf, accessed: 2015-27-09.
- 1011 [37] TDoS Attack Mitigation, http://www.cisco.com/c/en/us/td/
 1012 docs/ios-xml/ios/voice/cube_proto/configuration/15-mt/
 1013 cube-proto-15-mt-book/voi-cube-tdos-attack-mitigation.pdf,
 1014 accessed: 2015-27-09.
- [38] F. Huici, S. Niccolini, N. D'Heureuse, Protecting SIP Against Very Large
 Flooding DoS Attacks, in: GLOBECOM'09', IEEE Press, Piscataway, NJ,
 USA, 2009, pp. 1369–1374.
- [39] Y.-S. Wu, S. Bagchi, S. Garg, N. Singh, T. Tsai, SCIDIVE: A Stateful and Cross Protocol Intrusion Detection Architecture for Voice-over-IP Environments, in: DSN'04, IEEE Computer Society, Washington, DC, USA, 2004, pp. 433-.
- [40] C. Meadows, A Formal Framework and Evaluation Method for Network
 Denial of Service, in: CSFW, 1999, pp. 4–13.
- [41] A. Mahimkar, V. Shmatikov, Game-Based Analysis of Denial-of-Service
 Prevention Protocols, in: CSFW, 2005, pp. 287–301.
- [42] M. AlTurki, J. Meseguer, C. A. Gunter, Probabilistic Modeling and Analysis of DoS Protection for the ASV Protocol, Electr. Notes Theor. Comput.
 Sci. 234 (2009) 3–18.

- ¹⁰²⁹ [43] J. Henrique, I. E. Fonseca, V. Nigam, SHADE: Uma Estratégia Seletiva
- para Mitigar Ataques DDoS na Camada de Aplicação em Redes Definidas
 por Software, XXXIV Simpósio Brasileiro de Telecomunicações e Proces-
- 1032 samento de Sinais (SBrT 2016) (2016).