# Secure IP Address Allocation at Cloud Scale

Eric Pauley*‡, Kyle Domico*, Blaine Hoak*, Ryan Sheatsley*, Quinn Burke*,
Yohan Beugin*, Engin Kirda†, Patrick McDaniel*

*University of Wisconsin–Madison
‡Email: epauley@cs.wisc.edu
†Northeastern University

*Abstract*—Public clouds necessitate dynamic resource allocation and sharing. However, the dynamic allocation of IP addresses can be abused by adversaries to source malicious traffic, bypass rate limiting systems, and even capture traffic intended for other cloud tenants. As a result, both the cloud provider and their customers are put at risk, and defending against these threats requires a rigorous analysis of tenant behavior, adversarial strategies, and cloud provider policies. In this paper, we develop a practical defense for IP address allocation through such an analysis. We first develop a statistical model of cloud tenant deployment behavior based on literature and measurement of deployed systems. Through this, we analyze IP allocation policies under existing and novel threat models. In response to our stronger proposed threat model, we design *IP scan segmentation*, an IP allocation policy that protects the address pool against adversarial scanning even when an adversary is not limited by number of cloud tenants. Through empirical evaluation on both synthetic and real-world allocation traces, we show that IP scan segmentation reduces adversaries' ability to rapidly allocate addresses, protecting both address space reputation and cloud tenant data. In this way, we show that principled analysis and implementation of cloud IP address allocation can lead to substantial security gains for tenants and their users.

## I. INTRODUCTION

Cloud providers allow near limitless scalability to tenants while reducing or eliminating upfront costs. One component that enables this architecture is the reuse of scarce IPv4 addresses across tenants as services scale. Though a practical necessity, this reuse–combined with the use of IP addresses as a *security principal*–enables malicious cloud tenants to abuse IP address reputation [33, 45, 52], pollute the address space for future tenants [8], and even collect sensitive information intended for previous tenants [34, 12, 39]. We observe that these seemingly disparate attack spaces share a common thread: the ability of adversaries to easily sample large numbers of IP addresses from provider pools.

While prior works have identified and confirmed the issue of IP address reuse, and proposed some preliminary mitigations [12, 39], the community still lacks a complete understanding of the security provided by these measures, especially against a more powerful or adaptive adversary. For instance, prior works that attempt to reassign addresses to the same tenant can be defeated by adversaries using many disconnected cloud accounts (a form of Sybil attack). Developing secure policies for IP address allocation necessitates a fine-grained analysis of tenant behaviors and adversarial strategies. Such an analysis, and the stronger defenses that analysis enables, are the key focus of this work.

Towards this goal, we propose a novel, comprehensive model for IP address allocation on public clouds. By considering realistic distributions of benign tenant behaviors, configuration management, and cloud provider allocation policies, our new model enables us to concretely evaluate the effectiveness of attacks against the address pool. Implemented in the Elastic IP Simulator (EIPSIM), tenant and adversarial behaviors enable the key goal of our work: developing new allocation strategies that reduce the ability of adversaries to allocate, measure, and exploit many IP addresses. Our model is validated via real-world data on cloud tenant allocations, as well as data collected on cloud configuration management practices and discussions with major cloud providers. In this way, our model enables the development of new defenses against a broad class of attacks against cloud services.

Our model enables us to characterize and defend against a stronger adversary than considered in prior work. This *adaptive adversary* performs a Sybil attack against the cloud provider, creating many accounts to continually allocate new IP addresses from the pool. Hence, this attacker effectively defeats the protections provided by prior works. We propose *IP scan segmentation*, a novel IP allocation policy that heuristically identifies adversarial behavior across many cloud tenants, and effectively segments the pool to prevent such adversaries from allocating many unique IPs and exploiting vulnerabilities.

We use EIPSIM to evaluate the security properties (i.e., adversarial ability to discover unique IPs and exploitable configurations) of our studied allocation policies and tenant/adversarial behaviors in real cloud settings. Our analysis, spanning over 250 years of simulated IP address allocation, highlights the marked impact of IP allocation policies on the exploitability of IP address reuse. Indeed, our analysis shows that IP scan segmentation reduces adversarial success by $83.8\%$ over the IP allocation policies deployed by cloud providers, and by $70.1\%$ compared to prior explored techniques. Because our model concretely parallels the actual behavior of cloud providers and tenants, the techniques studied in this work can be directly implemented by providers to protect their customers and network resources. We have shared our findings with providers and release our models and policies

as open source artifacts[1] to support practical security of IP address allocation.

IP address reuse poses a practical security concern, but principled study of new allocation techniques can lead to practical defenses, making this reuse less exploitable in practice. Our work provides such a defense, as well as a basis on which future research in IP allocation can be measured.

## II. BACKGROUND

Our work addresses security properties of IP address allocation for public clouds. As such, we briefly describe considerations in IP allocation generally, as well as contemporary work in cloud security related to IP address allocation.

### A. IP Address Allocation

Network hosts require an IP address for communication. This can be manually assigned or managed out of band, or it can be provisioned through some automation. In home and corporate networks, the standard solution to automatic IP allocation is DHCP [18]. Likewise, in public clouds such as Amazon Web Services [3], Microsoft Azure [35] or Google Cloud [25], servers are allocated a private (i.e., RFC1918 [36]) IP address via DHCP [18]. While the DHCP standard does not specify how addresses are assigned, they are generally drawn from a pool either sequentially or based on the physical (MAC) address of the requesting machine [18]. For workloads with only private or outbound communications, these addresses are sufficient, as outbound connections can be mapped to publicly-routable IPs via Network Address Translation (NAT) [30].

When services need to receive connections from the broader Internet, they require a public IP address (usually, at a minimum, an IPv4, though support is increasing for IPv6 [46]). These addresses could be configured directly in the machine or over DHCP. However, cloud providers generally opt to use NAT [30] to route public IP addresses to the private IPs of servers. This has multiple benefits, including flexibility (public IPs can be changed dynamically without host involvement), security (tenants cannot spoof IPs), and ease of management (centralized view of IP address allocations).

**Cloud Provider IP Allocation.** When a tenant requests an IP address, cloud providers have a choice to return any unused address they control, subject to their own internal policy. For instance, a recent work [39] showed that Amazon Web Services samples their pool of available addresses pseudo-randomly subject to a 30-minute delay between reusing any given address. Another study [1] found that IP reuse followed a random process, though the ranges of used IP addresses could be inferred from many samples of the pool. Other works have found that Microsoft Azure [12] and Google Cloud Platform [1] show allocation behavior consistent with random allocation. While this random allocation can have the positive effect of allowing for a moving-target defense [1], wherein tenants move around the IP address pool to evade attack, it can also lead to severe security weaknesses as discussed below.
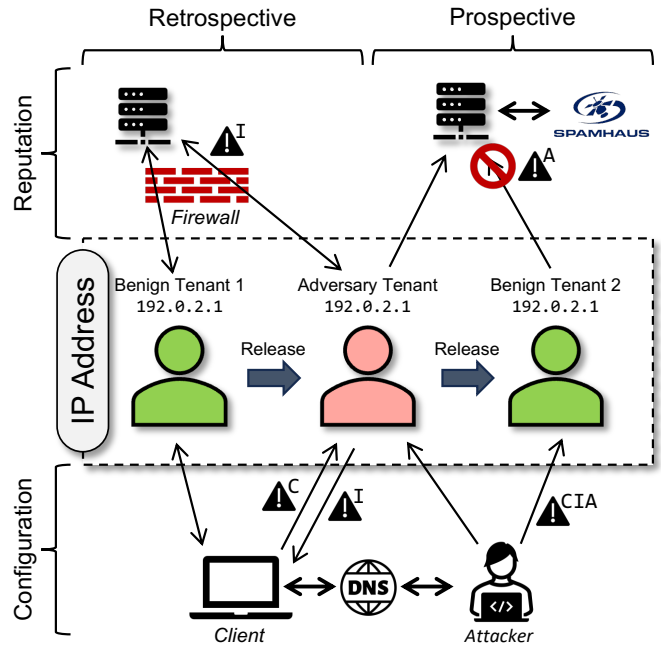
[1] https://github.com/MadSP-McDaniel/eipsim/



Fig. 1: Taxonomy of threats (⚠) to the (C)onfidentiality, (I)ntegrity, and (A)vailability of cloud-based network services from IP address reuse. Threats apply to previous tenants (retrospective), future tenants (prospective), and leverage the reputation of IP addresses or associated configuration.

**The Security Role of IP Addresses.** When viewed solely as a means to route traffic, IP addresses serve little security role. However, addresses have long been used in the capacity of *security principals*, i.e., control of an IP mediates access to resources, is associated with reputation, and can lead to the receipt of sensitive data. Firewall rules may filter access to specific IP addresses [33], servers may block messages from historical spam IPs [49, 8], and DNS can cause clients to send data to addresses [34, 12, 39].

### B. Exploiting IP Address Reuse

Due to the use of IP addresses as security principals, the (necessary) reuse of IPv4 addresses by cloud providers opens a set of vulnerabilities to attackers [33, 45, 52, 8, 34, 12, 39]. Depicted in Figure 1, these vulnerabilities allow adversaries to compromise the confidentiality, integrity, and availability guarantees of the network to other tenants in a variety of ways. We taxonomize such vulnerabilities into those that affect previous tenants (retrospective) and those that affect future tenants (prospective). Further, vulnerabilities may be related to the reputation of the IP address (and associated accessibility of other network services) or to configuration associated with that IP (and associated inbound traffic). Described below, these threat scenarios present different avenues for exploitation, though all rely on the ability for adversaries to acquire and route traffic over a sampling of cloud IP addresses.

2

**Reputation Attacks.** Source IP address is used to mediate access to a variety of resources on the public Internet. When an adversary uses an address to abuse other services (e.g., by sending spam email, malicious traffic, or large request volume) services may respond by blocking the address [33, 45, 52] and reporting to centralized reputation services (e.g., Spamhaus for email [49]). This poses a prospective threat to network availability for future tenants. When a future tenant attempts to access services, their address may be blocked because of the actions of previous tenants. Cloud providers pay careful attention to the reputation of their IP pools for services such as managed email sending [20], and routinely pay services to *clean* the reputation of their address space. The reputation of IP address ranges is also a key factor in the sale of address blocks [31]. Indeed, it is clear that prospective reputation threats to future tenants are a widespread and important issue, though one that to-date has seen little attention in terms of affecting address allocation.

Reputation can also pose risks retrospectively, though such attacks have not yet been observed in practice. Consider a service that mediates access via IP address allowlists [33]. A benign tenant may be granted access to restricted systems via their cloud-allocated IP, and then later release the IP address to the pool. An adversary can then allocate the IP address, and have access to the restricted service via the firewall allow rule. This attack is exceedingly difficult to perform, as the adversary must acquire the IP through random sampling then *also* determine additional services that may be accessible. However, if a service is known to authorize access to a variety of customers via their IP address, it may be possible to quickly enumerate a cloud provider's address space and discover authorized IPs.

**Configuration Attacks.** Tenants use IP addresses to refer to resources hosted on cloud providers, causing clients to connect to the resources and establishing trust relationships. Recent works have shown that, when tenants fail to remove the configurations referring to IP addresses they no longer control, these *latent* configurations can be exploited by future tenants [34, 12, 39]. Clients continue to send sensitive data, which is often unencrypted due to trust in the network isolation of the cloud provider. This *retrospective configuration* vulnerability is relatively easy for adversaries to exploit *en masse* on popular cloud providers, as the rapid and random reuse of IP addresses leaves little time for organizations to correct latent configurations. This leaves a long *window of vulnerability* during which adversaries could identify and exploit latent configuration. The community has proposed methods for correcting configurations such that they do not become latent, but changes to IP address allocation can also play a role when tenants fail to take action.

While of lower impact, configuration can also pose risks to services prospectively. Here, a tenant (denoted as adversarial although they may be relatively benign) may host services and create a configuration that causes large volumes of traffic to be sent to the address. This traffic could be sourced from legitimate services or from attackers targeting deployed software with exploits or denial of service attacks due to the services a tenant hosts. After releasing the IP, it is allocated to a new (benign) tenant, which then receives the malicious or high-volume traffic targeted at the previous tenant. At a minimum, such high-volume traffic can impose a cost on the new tenant, since cloud providers still charge for outbound bandwidth due to unwanted requests.

### C. Preventing exploitation of IP Reuse.

A commonality of all the above attacks is that they rely on the adversary allocating a vulnerable IP address. While the random nature of IP address allocation ostensibly makes the attacks untargeted, prior works have shown that adversaries can easily allocate thousands or even millions of addresses. Because allocation by major providers is currently pseudo-random [12, 39], vulnerabilities spanning many IP addresses become akin to the *birthday paradox*, wherein the probability of some adversary IP overlapping with some vulnerable IP quickly approaches 100%.

Changes to IP allocation policies have been shown to reduce the exploitability of IP Reuse. The goal here is to both (a) reduce the number of IPs that an adversary can allocate, (b) reduce the number of vulnerable tenants associated with those IPs, and (c) increase the window of time between reuse such that associated factors (configuration and reputation) have time to decay. While initial techniques towards achieving this have been proposed [12, 39], the community's understanding of the space of attacks and countermeasures here remains incomplete: that is, the ways in which an adversary might adapt to new techniques have not yet been modeled, and resulting further improvements to IP allocation strategies have yet to be explored. Hence, such important questions are the key focus of our work.

### III. MODELING THE IP ADDRESS POOL

Here, we present a comprehensive, novel framework for modeling secure IP address allocation. Towards this, we propose statistical models for tenant behavior (resource allocation and latent configuration), describe algorithms for allocation policies (including our proposed *IP Scan Segmentation* policy), and define threat models under which adversaries might exploit cloud resources. In each case, our methodology is informed by prior works, and validated based on real-world allocation and configuration datasets. *Note:* a reference of symbols used throughout the paper can be found in Appendix A.

### A. Tenant Behavior

Cloud providers lease resources (e.g., IPs) to tenants under two general paradigms: static and dynamic [15, 32, 14, 50, 11]. Static allocation allows tenants to acquire a specified amount of resources (perhaps for a fixed period of time); such resources are often used to handle workloads with known or predictable behavior. On the other hand, dynamic allocation allows tenants to acquire and release resources on-demand (to specified upper and lower limits); such resources are

typically backed by auto-scalers and other automation tools to handle less predictable workloads efficiently [42]. As such, we model the behavior of tenants within a spectrum of potential allocation strategies (defined in terms of the number of IPs currently allocated to the tenant) spanning static and dynamic resource allocation.

Benign tenants independently allocate IP addresses at some time $t_a$ from the pool and release those addresses at a later time $t_r > t_a$ (here, the IP is said to be allocated for $d_a = t_r - t_a$). Tenants also associate configuration with IP addresses, which is dissociated from the IP at $t_c$. Each tenant's overall behavior $B_i$ with respect to IP allocation can therefore be described as a set of timestamps:

$$B_i = \{(t_{a,0}, t_{r,0}, t_{c,0}), ..., (t_{a,n}, t_{r,n}, t_{c,n})\},$$

where $n$ is the total number of IPs allocated to the tenant. A single tenant's behavior then has a maximum limit of $S_{max}$ servers and minimum limit of $S_{min}$ servers; this can capture both static ($S_{max} = S_{min}$) and dynamic ($S_{max} > S_{min}$) resource allocation. For the purposes of our experiments, we focus primarily on dynamic allocations using auto-scalers, as we found this to be most representative of cloud tenant workloads [42, 51].

We next model each tenant's behavior as being independently sampled from a distribution of potential tenant behaviors: $B_i \sim \mathcal{B}$. We approximate $\mathcal{B}$ as a randomized $n$-term Fourier series with a base period of one day [51]. The intuition is that a given tenant's resource needs will likely vary throughout the day as demand peaks and subsides, but for a given tenant, this pattern will likely be similar from day to day. One work [51] suggests modeling with a period of 1 week for more precision. Our framework is flexible in this regard, but simulations are performed with 1-day periods. Recall that, by the Shannon-Nyquist sampling theorem [44], *any* daily-periodic function can be approximated by a Fourier series of sufficient terms. We compute the tenant's server utilization as a function of the current time $t$ ($0 \leq t \leq 1$), where 0 and 1 represent the beginning and end of the day, respectively. We then model the mean server usage of the tenant ($\bar{S} = \frac{S_{max} + S_{min}}{2}$) and the relative deviation from the mean server usage using the Fourier series:

$$\frac{S(t) - \bar{S}}{S_{max} - S_{min}} = \frac{\sum_{i=1}^{n} \frac{a_i}{i} \sin(2\pi i(t + \phi_i))}{\sum_{i=1}^{n} \frac{a_i}{i}},$$

where the Fourier amplitudes ($a_i$) and phases ($\phi_i$) are randomly sampled from the range $[0, 1]$. This series has an expected range of $[-0.5, 0.5]$, spanning from $S_{min}$ to $S_{max}$ throughout a simulated day. The tenant then allocates or releases IP addresses to respond to this change in compute needs [29]. In keeping with the behavior of a major cloud provider [4], IP addresses allocated under autoscale behavior are selected at random for release when a tenant scales down infrastructure.

Modeling autoscaling behavior as a Fourier series creates traces of tenant allocation that are sufficiently realistic to simulate allocation policies. However, on its own, it fails to account for the fact that IP allocations in a given cloud provider region would likely be correlated (due to the local geographies served by that region [28]). We account for this by biasing the sampling of the lowest-frequency phase of the Fourier series ($\phi_1$): enforcing that $\phi_1 < 0.5$, for instance, will roughly align peak loads to one half of the day. Moreover, tenants may have multiple workloads deployed under the same account that exhibit a hybrid of the above and other behaviors. While evaluation of these hybrid allocation behaviors is beyond the scope of this work, we note that EIPSIM can also be extended to support other models (or distributions) of tenant behavior, as well as real-world allocation traces. Analysis on real allocations (Section V-B) further support findings based on Fourier-distributed allocations, though effectiveness could vary on other workloads.

### B. Latent Configuration

As discussed above, tenants associate configuration with IP addresses when they are allocated. In most cases, this configuration is dissociated from the IP when or before the IP is released ($t_c \leq t_r$). In some cases, however, the configuration remains ($t_c > t_r$). If an adversary manages to allocate the IP address before $t_c$, we consider the adversary to have exploited the configuration. The time between IP release and latent configuration ($t_c - t_r$) is the *duration of vulnerability $d_v$* for a given tenant and IP.

Tenant behavior in dissociating configuration can be highly diverse. For feasibility, we model this configuration dissociation as a Poisson process. We assume that with some probability ($p_c$, a simulation parameter) the tenant leaves latent configuration. If latent configuration is left, it will be dissociated from the IP after some duration $d_v = t_c - t_r$. We model this as an exponential distribution

$$d_v \sim \text{Exponential}(1/d_a),$$

where the duration of vulnerability is distributed proportionally to the duration of allocation. Recall the probability density function of such a distribution:

$$f(d_v) = \begin{cases} \frac{1}{d_a} e^{-\frac{d_v}{d_a}} & d_v \geq 0 \\ 0 & d_v < 0 \end{cases}.$$

This distribution approximates the relationship between the duration of vulnerability and duration of allocation. It reflects empirical observations of cloud deployments [47], where relatively short-lived allocations are often orchestrated by automation tools and receive frequent configuration updates (and thus are less prone to having latent configurations), and relatively long-lived allocations are often configured manually and receive infrequent configuration updates (and thus are more prone to having latent configurations). In analysis of data on real-world latent configuration (Section V-A), we find additional evidence supporting this distribution.

## C. Adversarial Behavior

Within a public cloud, an adversary aims to obtain a large number of IP addresses with the goal of exploiting IP reputation, trust, or configurations by previous tenants. We proceed by describing the threat model and capabilities of such adversaries, followed by two modes of behavior: single-tenant (proposed by a prior work [39]) and multi-tenant (a new consideration of this work).

### 1) Threat Model

Our work considers an adversary attempting to scan a cloud provider's IP address pool to exploit both address reputation [33, 45, 52] and latent configuration left by other tenants (as demonstrated in [12, 34, 39]). This adversary has no privileged access to cloud resources, and bypasses no security controls in place. Instead, they can only provision resources using paid cloud accounts on a platform. In addition, the adversary can perform a Sybil attack, wherein they control a large number of cloud accounts that are indistinguishable from unique paid customers (e.g., by stealing credentials from other accounts, a common attack vector [22, 24, 37]). We parameterize adversaries by their compute budget (in unique IPs allocated simultaneously) and number of cloud accounts. While these may not be a direct financial cost to an adversary who steals accounts or payment details, they do still represent an opportunity cost, as these credentials could be used for other profitable purposes. The goal of this work is to decrease the effectiveness and increase the cost of such an attack as much as possible.

Within our scenario, the adversary has the capability to allocate IP addresses through public cloud offerings (e.g., Amazon EC2). Because we assume the cloud provider cannot soundly determine which tenants are controlled by the adversary, it must serve all tenant requests that are within policy.[2] For instance, allocating many instances and IP addresses is commonly used for autoscaling and short-lived tasks[51]. A cloud provider's actions must be a subset of those that would occur under existing offerings. For instance, while a cloud provider must allocate IPs to paying tenants, it may choose any free IP address to allocate. Based on this threat model, prior works [12, 39] have proposed a single-tenant adversary that allocates IPs under one tenant. This work considers a stronger adversary that has access to multiple tenants, defeating existing defenses through a Sybil attack.

### 2) Single-tenant Adversary

Discussed in prior works [34, 12, 39], a single-tenant adversary provisions IP addresses under a cloud account with the aim of allocating many unique addresses. In most cases, the most effective means by which to do this is to rent virtual servers with an associated IP address. A tenant allocates many of these servers simultaneously, runs them for the minimum time required to discover vulnerable configuration or leverage the address reputation, and then releases the IPs back to the

provider (or retains the server if there is interesting configuration associated). In this way, the tenant can easily sample from the IP address space unless the provider takes steps to prevent it. In line with cloud provider service quotas on concurrent allocations [6], our simulated single-tenant adversary allocates up to 60 IPs simultaneously for 10 minutes each, before releasing the IPs and allocating new ones.

### 3) Multi-tenant Adversary

The multi-tenant adversary adapts to protective allocation policies by leveraging multiple tenants for allocations. An adversary could create multiple tenants using Virtual Private Networks and private credit cards to evade detection[3]. Under this threat model, we also assume that a cloud provider must make allocation decisions based solely on tenant behavior, and cannot identify collusion between tenants otherwise. Further, the cloud provider must prioritize availability, and so must grant tenant allocation requests even if they believe the tenant to be malicious. Due to these factors, the multi-tenant adversary represents a stronger threat model that existing allocation policies may not protect against. In the worst case (and as simulated in Section IV-C), the adversary would continually use new accounts after allocating the maximum concurrent IPs on a single account.

## D. IP Allocation Policies

When tenants request an IP address from the cloud provider, the provider can choose which IP to assign to the tenant. Here, we assume (and prior works have shown [34]) that the cloud provider can freely choose to assign any free IP address within some zone to a tenant, and that there is no technical restriction on when IPs get reused. As noted (Section II-A), cloud providers use NAT to route public IP addresses, so assignment of these addresses can happen instantaneously and without any restriction from the underlying network topology.

Within this framework, the policy is a stateful set of functions that ALLOCATE, RELEASE, and INIT IP addresses:

**ALLOCATE**$(T, \theta) \longrightarrow (ip, \theta')$: Accepts a tenant id $T$ and an opaque state $\theta$ (for tracking IP allocation parameters) and returns a new, usable IP for the tenant, as well as an updated opaque state $\theta'$.

**RELEASE**$(ip, \theta) \longrightarrow (\theta')$: Accepts an allocated $ip$ (previously allocated by some tenant id $T$) and an opaque state $\theta$ and releases the IP back to the pool, returning an updated opaque state $\theta'$.

**INIT**$(ip, \theta) \longrightarrow (\theta')$: Accepts a new $ip$ into the pool that was never previously allocated, and returns an updated state $\theta'$.

All calls to ALLOCATE and RELEASE are paired in order, such that IP addresses are in use by at most one tenant at a time.

We next describe different allocation policies considered in our evaluation, providing their implementation in natural

---

[2]Under a more relaxed threat model, the cloud provider may refuse to service allocations from accounts that are likely malicious. Such actions are fully compatible with and complementary to our approach.

[3]Note that our threat model assumes the adversary is still cost-limited, either directly or in ability to acquire usable stolen credit card numbers.

language and pseudocode. Note that if the RELEASE and INIT interfaces are not provided, it is assumed that the default implementations presented in Algorithm 1 are used. Of these policies, the RANDOM policy is implemented in practice by cloud providers [39, 12], and the LRU and TAGGED policies were proposed by a prior work [39]. In addition to these policies that encompass the current state of the art, we propose and evaluate a new policy, *IP scan segmentation*.

---

**1 Function** RELEASE $(ip, \theta)$ **:**
**2**    $ip.t_r \leftarrow$ currentTime();
**3**    $\theta' \leftarrow$ setIpNotAllocated$(\theta, ip)$;
**4**    **return** $\theta'$
**5 end**
**6 Function** INIT $(ip, \theta)$ **:**
**7**    $\theta' \leftarrow$ createIp$(\theta, ip)$;
**8**    **return** $\theta'$
**9 end**

**Algorithm 1:** Default RELEASE and INIT interfaces

---

*1) Pseudorandom (RANDOM, Algorithm 2).*

The most basic IP allocation policy (and that used by major cloud providers [1, 39]) is pseudorandom allocation. Here, IPs are sampled randomly from the pool of available addresses, with the only restriction being that IPs cannot be used within $d_{reuse}$ (30 min as observed by prior work [39]). It has benefits for ease of use and understanding, as minimal information needs to be associated with the address. Further, the pool could be managed in a distributed fashion (such as within separate datacenters).

---

**1 Function** RANDOM.ALLOCATE $(T, \theta)$ **:**
**2**    $ip \leftarrow$ randomSample$(\mathcal{I} \setminus \mathcal{I}_{A_t})$;
**3**    **while** currentTime$() - ip.t_r < d_{reuse}$ **do**
**4**      $ip \leftarrow$ randomSample$(\mathcal{I} \setminus \mathcal{I}_{A_t})$;
**5**    **end**
**6**    $\theta' \leftarrow$ setIpAllocated$(\theta, ip)$;
**7**    **return** $ip, \theta'$
**8 end**

**Algorithm 2:** (RANDOM) IP Allocation

---

*2) Least Recently Used (LRU, Algorithm 3).*

The LRU policy seeks to maximize the median time between reuse of IP addresses. It does this by always allocating the IP address that has been in the pool the longest. Such an algorithm can either be implemented deterministically (e.g., using a FIFO queue), or stochastically (e.g., by sampling a subset of the IPs in the pool and returning the oldest of that batch). Such stochastic approaches have been shown to achieve acceptable performance in practice for caches [43].

*3) IP Tagging (TAGGED, Algorithm 4).*

Recent work [39] presented the first IP allocation policy specifically intended to prevent adversaries from scanning the IP pool. Referred to as *IP Tagging*, the authors describe that,

---

**1 Function** LRU.ALLOCATE $(T, \theta)$ **:**
**2**    $ip \leftarrow \underset{ip \in \mathcal{I} \setminus \mathcal{I}_{A_t}}{\arg\min} (ip.t_r)$;
**3**    $\theta' \leftarrow$ setIpAllocated$(\theta, ip)$;
**4**    **return** $ip, \theta'$
**5 end**

**Algorithm 3:** LRU IP Allocation

---

intuitively, released IP addresses are tagged with the tenant ID that released them. When allocating an IP, tenants first preference the IP addresses that they are tagged to, followed by addresses tagged to any other tenant using LRU allocation. Our implementation additionally stipulates that tagged IP addresses are selected in an LRU fashion, though other variants such as selecting the most-recently-used tagged IP may also be valid approaches. In any case, selecting a tagged IP address inherently exposes no additional IP address or tenant configuration to an adversary. Our evaluation also further characterizes IP Tagging beyond the metrics performed in prior work to assess the generality of the technique to stronger adversaries.

---

**1 Function** TAGGED.ALLOCATE $(T, \theta)$ **:**
**2**    **if** $\exists ip \in \mathcal{I} \setminus \mathcal{I}_{A_t} \mid ip.ID = T$ **then**
**3**      $ip \leftarrow \underset{ip \in \mathcal{I} \setminus \mathcal{I}_{A_t}}{\arg\min} (ip.t_r \mid ip.ID = T)$;
**4**    **else**
**5**      $ip, \_ \leftarrow$ LRU.ALLOCATE$(T, \theta)$;
**6**    **end**
**7**    $ip.ID \leftarrow T$;
**8**    $\theta' \leftarrow$ setIpAllocated$(\theta, ip)$;
**9**    **return** $ip, \theta'$
**10 end**

**Algorithm 4:** TAGGED IP Allocation

---

*4) IP Scan Segmentation (Algorithm 5).*

While IP tagging provides protection against a single-tenant adversary, the technique could be susceptible if an adversary spreads allocations across many tenants, bypassing the tagging entirely. In response to this threat, and our more powerful characterization of pool scanning adversaries (Section III-C3), we propose a new IP allocation policy that aims to prevent IP scanning by adversaries even when the adversary has access to an arbitrary number of cloud tenants.

Our proposed policy, *IP scan segmentation* (shown in Figure 2), works by identifying tenant behavior that is indicative of (and necessary for) IP pool scanning. The pool tracks the *mean allocation time* ($\bar{d}_a$) for each tenant $T$: relatively long-lived resources will lead to high $\bar{d}_a$, and adversarial scanning (which inherently must allocate many IPs) would require a low $\bar{d}_a$ to be economically feasible. IP addresses are tagged with both (a) the ID of the most recent tenant, and (b) the duration the IP was allocated for (this decays over time, see *cooldown time*). If the IP was previously held for longer, this
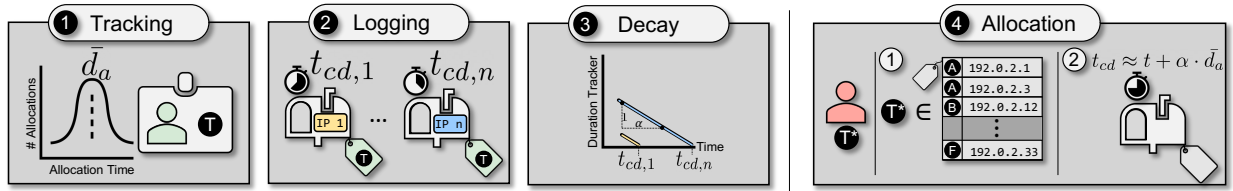
Fig. 2: IP Scan Segmentation - ❶ The mean IP allocation duration for tenant $T$ is tracked (i.e., $\bar{d}_a$), ❷ each released IP $n$ is first associated (i.e., tagged) with tenant $T$ & the allocated duration , ❸ the duration associated with the IP $n$ then decays linearly with rate $1/\alpha$ (stored as the cooldown time $t_{cd}$), and ❹ when an IP is allocated for tenant $T^*$, preference is first given to a $T^*$-tagged IP, then to an IP from the general pool whose $t_{cd}$ is closest to $t + \alpha \cdot \bar{d}_a$.

value does not change (so that a short allocation does not delete the protection from a previous longer allocation).

When a tenant allocates an IP address, preference is first given to an IP tagged to that tenant (as in IP Tagging), followed by an IP from the pool that was previously allocated for as close as possible to $\bar{d}_a$. In this way, adversary tenants that scan the IP space will in turn be allocated IP addresses that were previously allocated for short periods of time, either by another adversary tenant or by tenants deploying short-lived workloads (which are less likely to have associated latent configuration).

**Cooldown time.** As noted above, each IP is tagged with the longest duration it has been held for. Over time, this approach alone would cause more and more IPs to be tagged with long duration, leaving fewer and fewer with short durations and eventually allowing scanners to allocate the IPs that should be protected. Due to the scarcity of IP addresses, granting every IP address high protection means no IP receives protection.

To prevent this, the SEGMENTED policy applies a cooldown to the allocation duration over time with rate $1/\alpha$. The duration associated with an IP is therefore $d_a - (t - t_r)/\alpha$. Rather than continually update this duration in data structures, the SEGMENTED policy tracks the x-intercept of this function. This intercept, the *cooldown time* of the IP, is the time when the IP will no longer be provided any protection by the SEGMENTED policy. To select the IP with the most similar allocation duration for a given tenant, the policy minimizes $|(t_{cd}-t) - \alpha \cdot \bar{d}_a|$. In this way, tenants receive IP addresses that have exhibited similar allocation behavior to their past allocation behavior. Additionally, new tenants start with $\bar{d}_a = 0$, so they will receive IPs that have been segmented for allocation to scanners.

Since adversarial scanning would require a low $\bar{d}_a$ to be economical, an adversary tenant would then be matched with IPs that were either released a long time ago or were kept for a very short amount of time, mitigating some of the risk of an adversary acquiring an IP with latent configuration. Further, IPs recently released by the adversary would have a $t_{cd}$ consistent with their average allocation duration, increasing the likelihood that they receive the same IP back even under a different tenant.

```
1  Function SEGMENTED.ALLOCATE (T, θ) :
2      T.n_a ← T.n_a + 1;
3      if ∃ip ∈ I \ I_{A_t} | ip.ID = T then
4          ip ← arg min (ip.t_r | ip.ID = T);
              ip∈I\I_{A_t}
5      else
6          ip ← arg min (|ip.t_cd − currentTime() − α ·
              ip∈I\I_{A_t}
           T.d̄_a|);
7      end
8      ip.ID ← T;
9      ip.t_a ← currentTime();
10     θ' ← setIpAllocated(θ, ip);
11     return ip, θ'
12  end
13  Function SEGMENTED.RELEASE (ip, θ) :
14     ip.t_r ← currentTime();
15     ip.t_cd ← ip.t_r + α · (ip.t_r − ip.t_a);
16     T ← ip.ID;
17     T.d_a ← T.d_a + ip.t_r − ip.t_a;
18     θ' ← setIpNotAllocated(θ, ip);
19     return θ'
20  end
```

**Algorithm 5:** SEGMENTED IP Allocation

*E. Implementation*

To empirically study the distribution of IP allocation behaviors, adversarial techniques, and cloud provider defenses, we develop an IP pool simulator (EIPSIM). EIPSIM implements an extensible and configurable architecture for simulating interactions with IP address pools. Described in Figure 3 and Appendix C, EIPSIM allows us to evaluate the efficacy of allocation policies across a variety of scenarios over hundreds of years of simulated cloud provider allocation.

## IV. EVALUATION

We proceed by evaluating the security properties of IP allocation policies, first in synthetic settings (this section) followed by evaluations of realism and practicality (Section V). We begin by defining simulation parameters and defender objectives, followed by comparison of policies in both benign
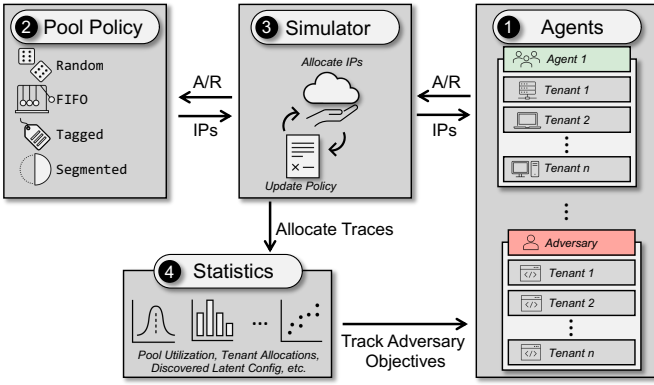
Fig. 3: Overview of our analysis - ❶ We first define agents who (A)llocate and (R)elease IP addresses in varying modalities (including adversarial behaviors), ❷ we then evaluate a suite of IP pool allocation policies that govern IPs associated with tenants, ❸ we then simulate interactions between agents and policies, and ❹ collect various statistics concerning pool utilization, adversarial goals, etc.

and adversarial settings.[4] Our analysis highlights the positive impact of new allocation policies: in both synthetic and real-world traces, IP Scan Segmentation reduces exploitable latent configuration and IP address yield compared to existing techniques, even when a strong adversary can use many cloud tenants.

### A. Simulation Parameters and Objectives

Our evaluation aims to quantify the impact of environmental, policy, and adversarial conditions on security properties. As such, we perform sweeps of multiple parameters that affect allocation policy performance. Within this setting, a cloud provider (defender) aims to reduce the ability of adversaries to exploit address reputation and previous tenants, objectives that are further defined here.

**IP Count and Utilization.** The size of the overall IP pool ($|\mathcal{I}|$), and the number of IPs allocated at any given time $|\mathcal{I}_{A_t}|$) has a substantial impact on allocation performance. If the majority of IP addresses are assigned, for instance, the pool policy has fewer choices when a tenant requests a new IP, and strategies that age, tag, or segment the addresses will therefore be less effective. Here, we can study performance by varying the max pool *allocation ratio* ($AR_{max} = \max_t \frac{|\mathcal{I}_{A_t}|}{|\mathcal{I}|}$) between simulations. Our evaluated simulation scenarios have $\max_t |\mathcal{I}_{A_t}| \approx 680\,\text{k}$, and compute $|\mathcal{I}|$ using $AR_{max}$ (set in each experiment).

**Allocation Duration.** Benign and adversarial tenants allocate IP addresses and hold them for some period of time. Study of the duration for which tenants and adversaries allocate IPs can yield insights on countermeasures. Our simulated adversary holds IPs for 10 minutes.

---

[4]Our main analysis considers a multi-tenant adversary. For completeness, we also include evaluations on single-tenant adversaries (those considered in prior works) in Appendix B.

**Free Duration.** Pools hold free addresses available for allocation, and holding an address for longer decreases the likelihood of associated latent configuration. As such, understanding the distribution of how long pools keep IPs free can suggest measures towards reducing latent configuration.

**Latent Configuration Probability.** In all simulations, we use a fixed probability of a given tenant leaving latent configuration, $p_c = 0.1$. In separate evaluations, we found that results varied roughly linearly with this parameter, making it less interesting for extensive study. However, future works could use more complex models for latent configuration where this constant plays a greater role.

**Defender Objectives.** Recall that the goal of the defender is to protect benign tenants against both *retrospective* (i.e., against previous tenants) and *prospective* (i.e., against future tenants) attacks by adversaries:

- *Retrospective attacks* are prevented by reducing the amount of *latent configuration* that an adversary can detect per IP allocated (proportional to total cost). We measure this quantity as *latent configuration yield*, the fraction of IP allocations which yield a (1) unique IP address with (2) some associated latent configuration.
- *Prospective attacks* are mitigated by shielding future tenants from the negative effects of adversaries. This can be achieved by reducing the number of unique IPs that adversaries can obtain. We therefore also measure *unique IP yield*, which is the fraction of IP allocations which yield a new unique IP address. Reducing unique IP yield likewise reduces the number of IPs whose reputation can be harmed by adversaries. Prospective attacks can also be evaluated by the rate of future tenant allocations that have been *poisoned* by an adversary. In ancillary evaluations, we found this metric to mirror that of unique IP yield, so we focus on the latter for extensive study.
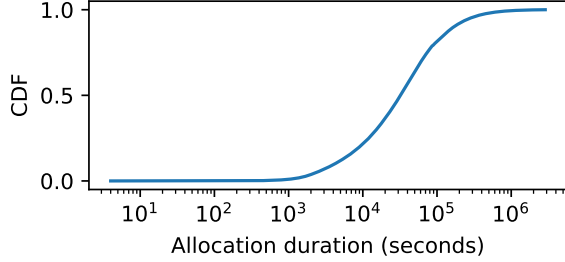
### B. Non-adversarial Scenario

To understand the aggregate performance of the various IP allocation policies, we first perform a simulation of the pool with no adversary. Here, agents allocate and release IP addresses on behalf of simulated tenants, and we study the effect of these policies on the configurations associated with allocated addresses. Experiments run for 180 simulated days, with 64 total experiments across allocation ratios and policy parameters (32 years of total simulated allocation).
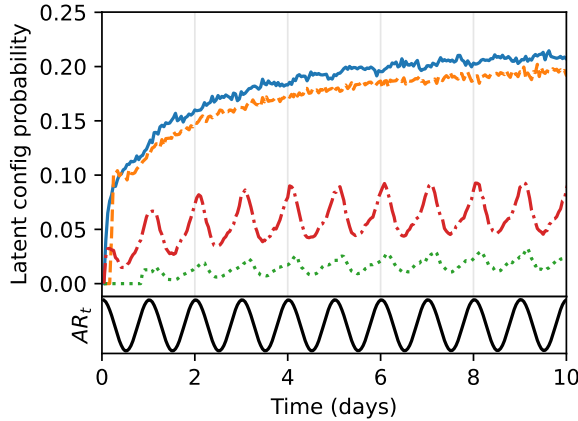
Results are shown in Figure 4. From these simulation results, we can distill several conclusions about the efficacy of our model and simulator, strength of existing and new allocation policies, and insights towards development of new policies.

**Tenant Behavior.** We first analyze the distribution of tenant allocation durations (Figure 4a). Here, we see that simulated allocations span several orders of magnitude in duration, representing a diverse distribution of behavior. Furthermore, to allocate within the distribution of other tenants an adversary would need to hold IPs and associated servers for an extended
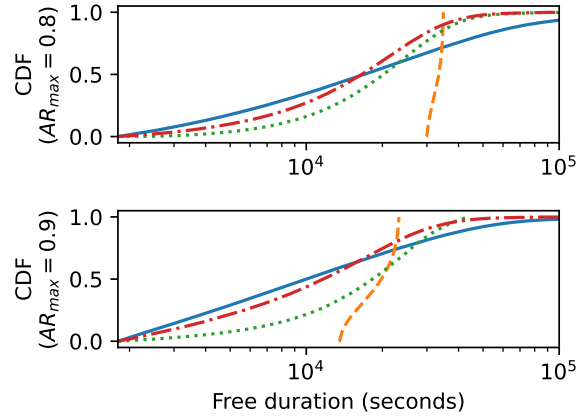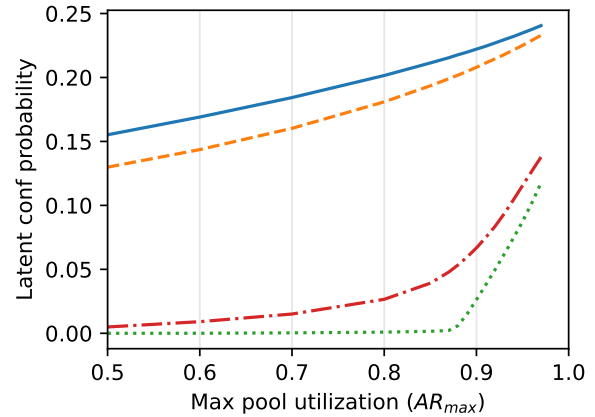
(a) Distribution of tenant allocation durations under the scenario described in Section IV-B. The broad range of allocation durations demonstrates the generality of the simulation parameters for evaluating pool policies.



(b) Latent configuration prevalence over time ($AR_{max} \approx 0.97$). Here, the lower plot shows the instantaneous allocation ratio of the pool over time ($AR_t$). As the pool reaches max allocation, strategies tend to allocated IPs with more latent configuration as addresses must be reused more quickly.



(c) Distribution of time between reuse across policies for two values of $RA_{max}$. The LRU pool sees the most impact from having more IPs available, as addresses can be aged for longer.



(d) Overall latent configuration varying max allocation ratio $AR_{max}$. IP Tagging and Segmentation policies further reduce prevalence even at high allocation ratios.

RANDOM ——  LRU -- -- --  TAGGED ·········  SEGMENTED —·—·—

Fig. 4: Modeling tenant allocations ($p_c = 0.1$).

period of time, reducing yield for a given cost. This provides hope that adversarial behavior in the pool could be identified and segmented from legitimate users.

**Time Between Reuse.** Next, we can see differences in how long policies keep IP addresses between reuse (Figure 4c). Results are shown for two allocation ratios ($AR_{max} = 0.8$ and $AR_{max} = 0.9$). These represent low- and high-contention scenarios for the pool, respectively. Beyond $AR_{max} = 0.97$, the policies cannot consistently age IPs for at least 30 minutes before reuse. In both cases, allocation schemes other than LRU perform similarly, reusing IP addresses in as little as 30 minutes, whereas LRU consistently maximizes the minimum time between reuse, by design. While this figure implies that LRU may be superior for preventing latent configuration, other policies that specifically target adversarial allocations may perform better in practice due to other factors.

**Pool Behavior Over Time.** Looking at prevalence of latent configuration over time in Figure 4b, we initially see lower prevalence as the pool has unused IP addresses to allocate. Beyond that, prevalence for RANDOM and LRU allocation approaches $p_c$ (note that prevalence can exceed $p_c$ as multiple tenants have the opportunity to associate configuration with a given IP address). LRU unsurprisingly outperforms RANDOM slightly, due to the higher time between reuse of IP addresses. While higher time between reuse most clearly reduces aggregate exposure of latent configuration under our posited exponential distribution, cloud providers could also use EIPSIM with other models of latent configuration to validate against their unique scenarios. We expect similar results from any monotonic distribution of $d_v$.

**Effect of Pool Utilization.** IP addresses are a scarce resource, so cloud providers should aim to achieve the best secu-

rity against latent configurations while incurring minimal pool size overhead. In Figure 4d, we see that the studied allocation policies have differing behavior as pool size changes. While both TAGGED and SEGMENTED outperform the RANDOM and LRU policies, TAGGED performs slightly better. This is because benign tenants preferentially receive IPs from other tenants exhibiting benign behavior, making other IPs available for segmentation to heuristically malicious tenants. Our experiments demonstrate that allocation policies can have marked impact on overall latent configuration exposure even for high IP allocation ratios.

Our non-adversarial experiments show that EIPSIM and its associated models are a compelling means by which to study the behavior of IP address pools, spanning a broad range of resulting tenant allocations. The parameters of our initial simulation prove interesting for further study, as the variety of tenant behaviors leads to differentiated performance across allocation policies.

### C. Multi-tenant Adversary

Next, we evaluate how pool implementations defend against a sophisticated adversary who can use many cloud tenants to bypass the protections of existing allocation policies. We also evaluate how an adversary's success varies with the number of tenants they can create. Simulations are run in the benign setting for 180 days, followed by 30 days of adversarial exploitation. 280 total simulations are performed across allocation ratio, policy, and adversarial tenant count parameters (160 years of total simulated allocation).

**Unique IPs.** Figure 5b shows the number of unique IPs discoverable by an unlimited-tenant adversary as pool utilization varies. While the non-tenant-aware policies RANDOM and LRU show no difference from the single-tenant adversary, tenant-aware policies show surprising results. In both cases, unique IPs reduce as utilization increases to some critical point, then increases again. The increased unique IP yield at low utilization results from the large number of free IPs, resulting in increased availability for both benign and malicious tenants. As a result, higher unique IP yield does not result in discovery of exploitable latent configurations. Above the critical point, both TAGGED and SEGMENTED *must* allocate potentially-dangerous IPs to tenants, but SEGMENTED successfully identifies behavior patterns across adversary tenants and reduces the number of unique IPs seen. In this way, SEGMENTED successfully protects a larger portion of the IP space.

**Latent Configuration.** Figure 5c shows how the multi-tenant adversary's yield of latent configuration varies with allocation ratio. Here, we see the complete effect of tenant-aware allocation policies: below the policy's critical point, allocated IPs have minimal associated latent configuration, so a high unique IP yield does not allow exploitation. Above this, strategies offer only mild protection (i.e., approaching that of non-tenant-aware policies). Most importantly, however, this plot emphasizes the advantages of IP scan segmentation:

SEGMENTED reduces latent configuration yield by $83.8\%$ compared to RANDOM, whereas TAGGED only reduces yield by $46.0\%$. In other words, SEGMENTED achieves an additional reduction of $70.1\%$. When considering an adversary with the ability to use multiple cloud tenants, SEGMENTED offers superior protection to prior works and currently-deployed policies.

Looking at a time-series plot of allocations (Figure 5a), we see that TAGGED performance improves over time, but fails to protect the IP pool against exploration for high $AR_{max}$. However, TAGGED still provides some protection even at these high allocation ratios, likely because it reduces the number of unique IPs with which tenants associate latent configuration. In contrast, SEGMENTED has an initial spike in exposed configurations, which then quickly stabilizes at a lower yield.

**Effect of Tenant Count.** A realistic adversary may not have access to create an unlimited number of tenants in the public cloud, due to billing and other compliance measures taken by the provider. As such, it is helpful to understand how adversarial capability scales with number of tenants under various allocation policies. In Figure 5d, we see the marked effect of scaling tenant counts on effectiveness against TAGGED. An adversary begins to increase latent configuration yield above 20 tenants, with peak yields reached at 60 tenants. In contrast, SEGMENTED provides only a slight increase in yields even with no limit on tenants[5].
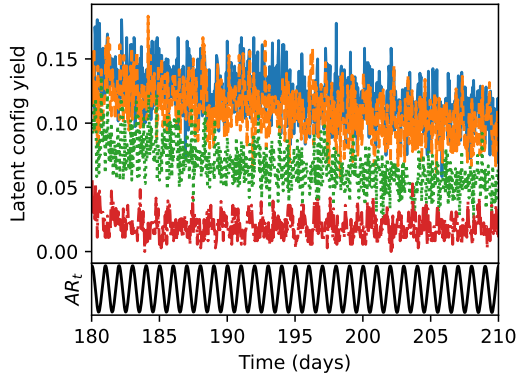
Our analysis of the multi-tenant adversary demonstrates the limitations of existing allocation policies, as an adversary using many tenants can still discover latent configuration. In contrast, IP scan segmentation's heuristics more effectively segment pool scanning based on the characteristics of allocations, rather than just tenant identifiers, and so are resistant to these attacks. Further, the SEGMENTED pool achieves improved performance even at very high pool contention, approaching the practical limit while maintaining existing minimum reuse durations.
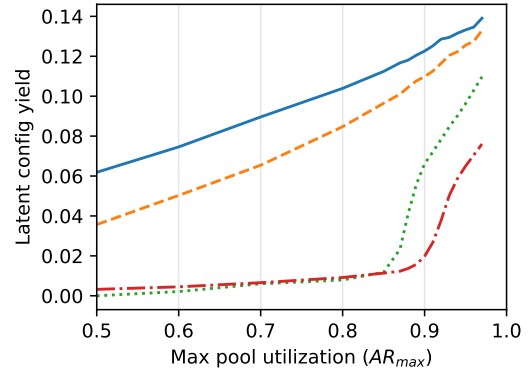
### D. Tuning Segmentation

Because the SEGMENTED policy is parameterized by some $\alpha$, it is important to tune the parameter for optimal performance. Here, we seek to understand how varying $\alpha$ affects adversarial yields under our simulation, and also how cloud providers might model policies on their own traces. To do this, we perform simulations of an unlimited-tenant adversary against a SEGMENTED pool. We vary $\alpha$ and study the yields of unique IPs and latent configuration. Simulation timelines are the same as in the multi-tenant evaluation (180 benign + 30 adversarial days), with 101 total experiments performed (58 total years of simulated allocation).

Our results (Figure 6) demonstrate the effect of varying $\alpha$. Varying $\alpha$ can make up to a $25\%$ variation in unique IP yields, and up to a $40\%$ variation in latent configuration yield. The relationship between $\alpha$ and adversarial objectives
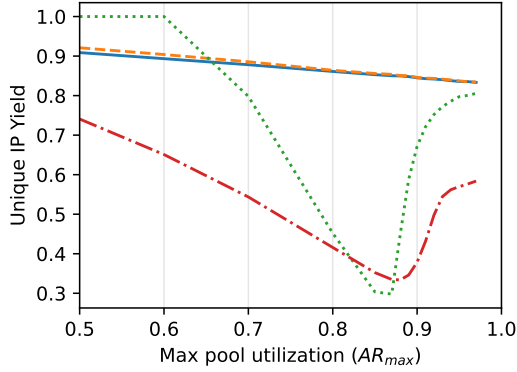
---

[5]A limit of $10^4$ in this scenario allows the adversary to never reuse a tenant, so the tenant count is effectively unlimited.
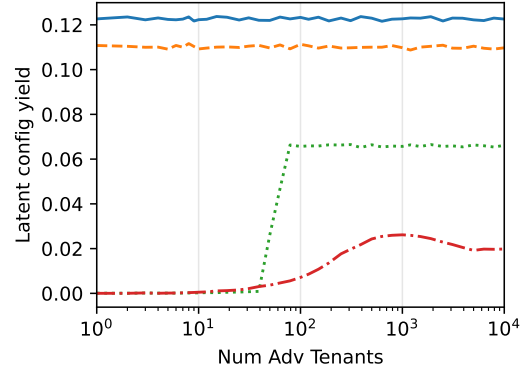
(a) Yield of new latent configurations over time ($AR_{max} = 0.9$). Data points are bucketed by hour.



(c) Effect of pool utilization on discovered latent configs (unlimited adversary tenants)



(b) Effect of pool utilization on discovered unique IPs (unlimited adversary tenants)



(d) Effect of total adversary tenants on latent configuration yield ($AR_{max} = 0.9$)

RANDOM    LRU    TAGGED    SEGMENTED
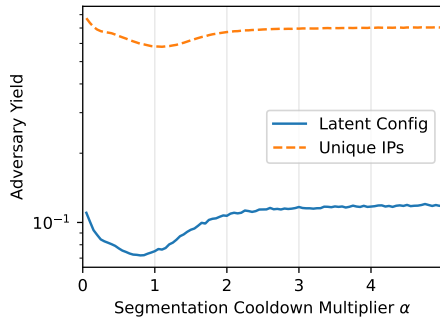
Fig. 5: Modeling the multi-tenant adversary.



Fig. 6: Effect of varying Segmentation parameter $\alpha$ on adversarial yield. ($AR_{max} = 0.9$)

is convex, leading to a clear global optimum for configuration of a deployed system.

In addition to demonstrating an optimal value of $\alpha$ in our simulation setting, our results also suggest that modeling configurations of the SEGMENTED policy could be performed without making as strict of assumptions about latent configuration. Recall that our analysis assumes exponentially-distributed latent configuration durations. While a cloud provider could substitute real IP allocation traces, collecting data on concrete configurations is far more difficult. In our results, however, we show that latent configuration and IP yield are highly correlated, so a cloud provider could model IP address yields on concrete data and be confident in applicability to latent configuration yields, as well.

## V. MODEL REALISM

We next demonstrate the realism of our model and policies by comparing our statistical models against data observed in real-world cloud settings (Section V-A), evaluating on real-world allocations (Section V-B), and demonstrating the scalability of allocation policies to major cloud providers (Section V-C). When evaluated in a realistic setting, IP Scan Segmentation performs even more favorably than under synthetic benchmarks. Further, our performance figures demonstrate that all evaluated policies can scale to the size and performance requirements of major providers.
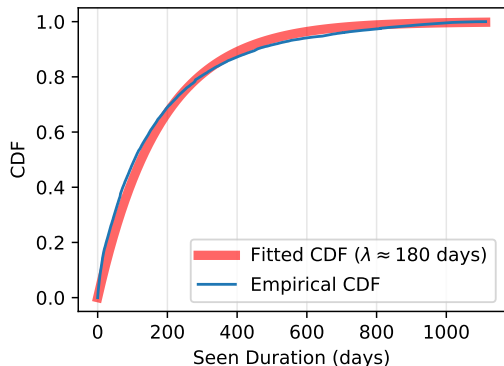
Fig. 7: Distributions of latent configuration durations collected on real-world cloud traffic. Latent configuration durations fit to the hypothesized exponential distribution.

### A. Validating Latent Configuration

Finally, we evaluate the realism of our model of latent configuration by analyzing the distribution of latent configuration in deployed systems. To do this, we analyze real-world latent configurations as visible from a cloud-deployed Internet telescope [38]. We use the DScope HTTP(S) request dataset spanning three years from March 2021-March 2024. Because many of the IP addresses studied are received many times, prevalence of latent configuration can be analyzed over time. For configurations seen multiple times in the study, we compute the maximum duration each configuration was seen, and use maximum likelihood estimation to fit a corresponding exponential distribution.

The resulting distribution (Figure 7) is consistent with our hypothesized exponential distribution of latent configuration. While characterizing the distribution of latent configuration with respect to the underlying IP address allocation would require data from cloud providers, these empirical results support our statistical models and the effectiveness of our studied policies.

### B. Evaluating on Real Allocation

We next seek to understand whether our model of IP allocation generalizes to real-world scenarios. To this end, we perform an evaluation of allocation policies against server allocation traces from Google's `clusterdata-2019` dataset [48]. This dataset contains real-world server allocations and usage traces across 31 days in eight independent clusters in a hybrid cloud setting.

To extract corresponding IP address allocation traces from allocations in `clusterdata-2019`, we take all Job (groups of processes running as a single collection) allocations across all eight clusters, remove malformed jobs or those running beyond the scope of 31 days, and extract the User of these jobs as a tenant ID. Each Job is assumed to have a public IP address allocated, and latent configuration is modeled over these jobs as previously discussed. The resulting traces contain $24\,\mathrm{M}$ allocations across $21\,\mathrm{k}$ tenants, with $\max_t |\mathcal{I}_{A_t}| \approx 119\,\mathrm{k}$.

Results (Figure 8) largely confirm the effectiveness of new IP allocation policies. Here, we see that SEGMENTED prevents discovery of latent configurations by an adversary with unlimited tenants, even at high pool utilization. Notably, `clusterdata-2019`'s composition of short-lived allocations for batch jobs represents a worst-case scenario, with many of these allocations seemingly indistinguishable from those used by an adversary. Yet, the SEGMENTED policy reduces the sharing of long-lived IP allocations with these short-lived tenants, preventing the adversary from discovering IPs with associated latent configuration.

One interesting phenomenon visible on these real-world allocation traces is the non-monotonic effect of tenant count on attack effectiveness. Here, we see that an adversary achieves increasing latent configuration yield with more tenants, then reduced effectiveness once tenants are no longer reused. This is a result of the default reputation of tenants: a new tenant has a $d_a/n_a$ of 0, which is then increased by allocating and releasing IPs. Reusing tenants with this (minimal) increase in reputation affords greater yield, especially when legitimate tenants have similar IP allocation behavior. While this worst-case scenario emphasizes a weakness of the SEGMENTED policy, it is unlikely that similar behavior would be seen in a public cloud, where job-based products such as AWS Lambda and Batch do not assign public IPs to short-lived instances.

### C. Performance & Scalability

TABLE I: Performance scaling of IP Scan Segmentation with pool size. Speedup is the amount of simulated time (100 days) divided by time to simulate. IP Scan Segmentation scales to pools with millions of IPs and hundreds of millions of allocations.

| # IPs | Runtime | Speedup | Allocations | Allocs/s |
|-------|---------|---------|-------------|----------|
| 100   | 500 ms  | 17 M    | 4.2 k       | 8.3 k    |
| 1 k   | 530 ms  | 16 M    | 26 k        | 48 k     |
| 10 k  | 2 s     | 4.3 M   | 220 k       | 110 k    |
| 100 k | 14 s    | 630 k   | 2.2 M       | 160 k    |
| 1 M   | 187 s   | 46 k    | 22 M        | 120 k    |
| 10 M  | 2.3 ks  | 3.8 k   | 220 M       | 97 k     |

TABLE II: Largest major cloud compute regions.

| Provider   | Largest Region | # Zones | # IPs  |
|------------|----------------|---------|--------|
| GCP [23]   | `us-central-1` | 4       | 2.8 M  |
| Azure [10] | `eastus`       | 3       | 3.3 M  |
| AWS [9]    | `us-east-1`    | 5       | 16 M   |

Our work aims to provide practical security improvements through IP allocation policies, and it is therefore important that such policies are realizable. To this end, we evaluate the performance of IP Scan Segmentation on various IP address pool sizes. Notably, because EIPSIM simulates each concrete IP address allocation, the compute requirements required to simulate allocations are similar to that of a production allocation pool. We simulate non-adversarial scenarios on a
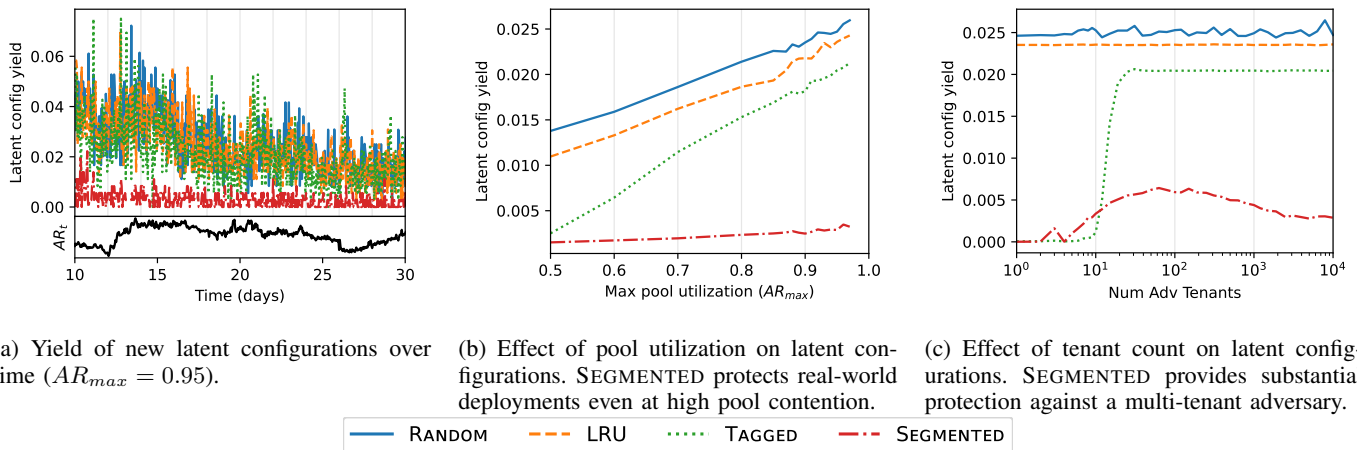
(a) Yield of new latent configurations over time ($AR_{max} = 0.95$).

(b) Effect of pool utilization on latent configurations. SEGMENTED protects real-world deployments even at high pool contention.

(c) Effect of tenant count on latent configurations. SEGMENTED provides substantial protection against a multi-tenant adversary.

——— RANDOM      - - - LRU      ······· TAGGED      —·—· SEGMENTED

Fig. 8: Evaluating allocation policies on real-world traces from `clusterdata-2019`.

commodity x64 server with 64vCPU and 192GB of RAM, though simulations use only one CPU thread. In each case, $|\mathcal{I}|/10$ tenants were used[6] with a max concurrent allocation of 10 per tenant. Simulations run for 100 (virtual) days. Results (Table I) show runtime and allocation rates with respect to pool size, demonstrating that EIPSIM scales with pool size to millions of allocations. The largest cloud compute regions (Table II) can reallocate at most a few thousand addresses per second, well within the performance of SEGMENTED on a single CPU core.

Studied performance numbers align with the sizes of the largest cloud compute regions (see Table II), demonstrating the proposed techniques scale to the size of major providers.

We further demonstrate the achievability of new policies by evaluating the real-world behavior of an existing provider and how those map to the information storage requirements of our proposed SEGMENTED policy. In the case of AWS, while allocation is random, AWS also already tags IP addresses with their previous tenant, and allows tenants to reuse released IPs if they have not been allocated to another tenant [19]. This currently-stored data is sufficient to perform the tenant tagging used by SEGMENTED and TAGGED policies. The remainder of the SEGMENTED policy requires associating an additional timestamp with each IP. Candidate IPs are then randomly sampled (as under current policies) and a best-fit IP is selected based on the heuristic. In this way, the SEGMENTED policy can be achieved using the existing data structures implemented by a major provider.

## VI. LIMITATIONS

IP allocation policies are a heuristic mitigation, rather than a sound solution, for abuse of a cloud provider's IP address pool. Under our threat model, providers *must* allocate some address to a tenant on request. Further, the provider will always

---

[6]The policies discussed store $O(1)$ data per tenant, so per-tenant compute overhead in EIPSIM is largely caused by simulating tenant agents, rather than the allocation policies themselves.

face limited information, as a Sybil attack is not soundly distinguishable from benign new tenants.

**Effects on Benign Tenants** While IP Scan Segmentation reduces the ability of adversaries to observe latent configuration and harm future customers *in expectation*, it may also affect benign tenants, especially new customers. These new customers would be treated the same as new adversarial tenants, and may therefore receive disproportionately more low-quality IP addresses that were previously controlled by an adversary. Note that this applies only to prospective threats (i.e., a new benign tenant receiving an address that has been polluted by an adversary), as new tenants that then hold IPs for long periods will receive protection from retrospective threats.

To mitigate harms to new tenants, providers can add additional signals to the allocation process. The multi-tenant adversary requires access to many payment credentials that are likely of low quality (e.g., stolen credit cards), so countermeasures can privilege behaviors likely not associated with these. For instance, customers that purchase high-margin non-compute products (i.e., those not useful for adversarial IP allocation) or those with commercial contracts and vetted relationships, such as new accounts under existing billing arrangements. The effective price of leased addresses can also be considered (Section VII-A).

**Provider Feedback and Implementation** Our proposed policies have not yet been deployed by major cloud providers. While our evaluation demonstrates that allocation policies can be effectively implemented from a technical perspective, other factors may prevent their adoption, such as associated reputational risks as providers take responsibility for client configurations. Cloud providers operate under a *shared responsibility model* [2], wherein they take responsibility for infrastructure security and customers are responsible for their workloads. Defending against retrospective threats to IP allocation blurs this boundary, and potentially exposes providers to increased risk or scrutiny when protections fail. Similar actions have been taken in other shared domains, such as cloud storage

security [16], providing hope that evaluations of IP allocation effectiveness may lead to practical improvements in security.

Ultimately, we recommend that providers continue to embrace the shared responsibility model while protecting tenants when possible. In the case of IP address allocation, this would entail soft adoption of proposed allocation policies. Notably, all discussed policies are *fully compatible* with the existing documented behavior of major providers. Non-documented adoption of new allocation policies would provide heuristic protection for customers while minimizing associated responsibility for the security of customer configurations.

## VII. Discussion & Related Work

### A. Allocation Pricing Signals

IP Scan Segmentation aims to increase the cost associated with IP scanning by tracking the amount spent per IP based on allocation time. In real cloud providers, this could motivate further extending the policy by incorporating other pricing signals from the cloud provider. For example, a tenant allocating powerful servers for short periods of batch processing is indistinguishable from scanning using just allocation traces, but the cloud provider could measure the total cost associated with these allocations and distinguish the activity as legitimate. The IP pool is a scarce resource, and so reducing the number of scanner-segmented IPs allocated to these resources will leave more available for scanners, improving policy effectiveness. We anticipate that cloud providers can extend the EIPSIM framework to incorporate these pricing signals and further improve practical security.

### B. Configuration Management

While the security of IP address allocation has had limited study, the way that organizations *configure* services has seen extensive related research. This is especially true in the cloud setting, which introduces new challenges in managing service configurations. Prior work has demonstrated that configuration complexity may increase substantially with the scale of the service [13, 47] and from the added tasks associated with making services cloud native (i.e., using advanced features such as auto-scaling [21, 17]). Automated configuration management tools (such as Puppet [41], Chef [40], and Ansible [7]) have eased this complexity to some extent. Further, infrastructure-as-code (IaC) [26] tools (such as AWS CloudFormation [5] or Terraform [27]) have made configuration management almost entirely non-interactive. However, while automation tools can eliminate most human-errors at runtime, a large proportion of configuration errors have been attributed to subtle bugs in the configuration files themselves (or ambiguities in the code generating them) [47] and other improper lifecycle management practices [13] (e.g., failing to remove configurations pointing to released IPs [34, 12, 39]).

## VIII. Conclusion

The way in which cloud IP addresses are allocated has a substantial impact on the security of hosted applications. Our work proposes new defenses for cloud IP allocation, and evaluates these defenses through a comprehensive model of tenant and adversarial behavior. Our proposed new policy, IP scan segmentation, successfully reduces an adversary's ability to scan the IP pool even if they can create new cloud tenants without limit. We anticipate that new IP allocation policies, such as IP scan segmentation, will prove useful to providers in protecting their customers and address pools. To that end, we release both our models and policies as open source for use by providers. We are also hopeful that modeling of IP allocation, such as that implemented in EIPSIM, will enable further improvements in the security of cloud provider offerings.

## References

[1] Hussain M. J. Almohri, Layne T. Watson, and David Evans. "Predictability of IP Address Allocations for Cloud Computing Platforms". In: *IEEE Transactions on Information Forensics and Security* 15 (2020). Conference Name: IEEE Transactions on Information Forensics and Security, pp. 500–511. ISSN: 1556-6021. DOI: 10.1109/TIFS.2019.2924555.

[2] Gui Alvarenga. *What is the Shared Responsibility Model? - CrowdStrike*. en. Nov. 2022. URL: https://www.crowdstrike.com/cybersecurity-101/cloud-security/shared-responsibility-model/ (visited on 07/14/2024).

[3] *Cloud Services - Amazon Web Services (AWS)*. en-US. URL: https://aws.amazon.com/ (visited on 08/19/2021).

[4] *Work with Amazon EC2 Auto Scaling termination policies - Amazon EC2 Auto Scaling*. URL: https://docs.aws.amazon.com/autoscaling/ec2/userguide/ec2-auto-scaling-termination-policies.html (visited on 06/09/2022).

[5] *AWS CloudFormation*. URL: https://aws.amazon.com/cloudformation/.

[6] *AWS service quotas - AWS General Reference*. URL: https://docs.aws.amazon.com/general/latest/gr/aws_service_limits.html (visited on 06/16/2022).

[7] *Ansible*. URL: https://www.ansible.com/.

[8] Manos Antonakakis et al. "Building a dynamic reputation system for DNS". In: *19th USENIX Security Symposium (USENIX Security 10)*. 2010. URL: https://www.usenix.org/event/sec10/tech/full_papers/Antonakakis.pdf (visited on 04/05/2024).

[9] *AWS IP address ranges - AWS General Reference*. URL: https://docs.aws.amazon.com/general/latest/gr/aws-ip-ranges.html (visited on 08/16/2021).

[10] *Azure IP address ranges*. URL: https://www.microsoft.com/en-us/download/details.aspx?id=56519.

[11] BH Bhavani and HS Guruprasad. "Resource provisioning techniques in cloud computing environment: a survey". In: *International Journal of Research in Computer and Communication Technology* 3.3 (2014), pp. 395–401.

[12] Kevin Borgolte et al. "Cloud Strife: Mitigating the Security Risks of Domain-Validated Certificates". en. In: *Proceedings 2018 Network and Distributed System Security Symposium*. San Diego, CA: Internet Society, 2018. ISBN: 978-1-891562-49-5. DOI: 10.14722/ndss.2018.23327. URL: https://www.ndss-symposium.org/wp-content/uploads/2018/02/ndss2018_06A-4_Borgolte_paper.pdf (visited on 08/04/2021).

[13] Aaron B Brown and Joseph L Hellerstein. "An approach to benchmarking configuration complexity". In: *Proceedings of the 11th workshop on ACM SIGOPS European workshop*. 2004, 18–es.

[14] Rodrigo N Calheiros, Rajiv Ranjan, and Rajkumar Buyya. "Virtual machine provisioning based on analytical performance and QoS in cloud computing environments". In: *2011 International Conference on Parallel Processing*. IEEE. 2011, pp. 295–304.

[15] Sivadon Chaisiri et al. "Cost minimization for provisioning virtual servers in amazon elastic compute cloud". In: *2011 IEEE 19th Annual International Symposium on Modelling, Analysis, and Simulation of Computer and Telecommunication Systems*. IEEE. 2011, pp. 85–95.

[16] Thomas Claburn. *AWS simplifies Simple Storage Service to prevent data leaks*. en. Dec. 2022. URL: https://www.theregister.com/2022/12/14/aws_simple_storage_service_simplified/ (visited on 07/14/2024).

[17] Tharam Dillon, Chen Wu, and Elizabeth Chang. "Cloud computing: issues and challenges". In: *2010 24th IEEE international conference on advanced information networking and applications*. Ieee. 2010, pp. 27–33.

[18] Ralph Droms. *Dynamic Host Configuration Protocol*. Request for Comments RFC 2131. Num Pages: 45. Internet Engineering Task Force, Mar. 1997. DOI: 10.17487/RFC2131. URL: https://datatracker.ietf.org/doc/rfc2131 (visited on 06/10/2022).

[19] *Elastic IP addresses - Amazon Elastic Compute Cloud*. URL: https://docs.aws.amazon.com/AWSEC2/latest/UserGuide/elastic-ip-addresses-eip.html#using-eip-recovering (visited on 02/13/2023).

[20] Holly Esquivel, Aditya Akella, and Tatsuya Mori. "On the effectiveness of IP reputation for spam filtering". In: *2010 Second International Conference on COMmunication Systems and NETworks (COMSNETS 2010)*. ISSN: 2155-2509. Jan. 2010, pp. 1–10. DOI: 10.1109/COMSNETS.2010.5431981. URL: https://ieeexplore.ieee.org/abstract/document/5431981 (visited on 04/05/2024).

[21] Dennis Gannon, Roger Barga, and Neel Sundaresan. "Cloud-native applications". In: *IEEE Cloud Computing* 4.5 (2017), pp. 16–21.

[22] Thomas Gardner and Cody Betsworth. *How adversaries infiltrate AWS cloud accounts*. en. 2023. URL: https://redcanary.com/blog/aws-sts/ (visited on 04/16/2024).

[23] *GCP IP address ranges*. URL: https://www.gstatic.com/ipranges/cloud.json.

[24] Dan Goodin. *Developers can't seem to stop exposing credentials in publicly accessible code*. en-us. Nov. 2023. URL: https://arstechnica.com/security/2023/11/developers-cant-seem-to-stop-exposing-credentials-in-publicly-accessible-code/ (visited on 04/16/2024).

[25] *Cloud Computing Services — Google Cloud*. URL: https://cloud.google.com/ (visited on 08/19/2021).

[26] Michele Guerriero et al. "Adoption, support, and challenges of infrastructure-as-code: Insights from industry". In: *2019 IEEE International Conference on Software Maintenance and Evolution (ICSME)*. IEEE. 2019, pp. 580–589.

[27] *Terraform by HashiCorp*. URL: https://www.terraform.io/.

[28] Keqiang He et al. "Next stop, the cloud: Understanding modern web service deployment in ec2 and azure". In: *Proceedings of the 2013 conference on Internet measurement conference*. 2013, pp. 177–190.

[29] Nikolas Roman Herbst, Samuel Kounev, and Ralf Reussner. "Elasticity in cloud computing: What it is, and what it is not". In: *10th International Conference on Autonomic Computing (ICAC 13)*. 2013, pp. 23–27.

[30] Matt Holdrege and Pyda Srisuresh. *IP Network Address Translator (NAT) Terminology and Considerations*. Request for Comments RFC 2663. Num Pages: 30. Internet Engineering Task Force, Aug. 1999. DOI: 10.17487/RFC2663. URL: https://datatracker.ietf.org/doc/rfc2663 (visited on 06/10/2022).

[31] Lee Howard. *Networking Education: IP Blocklist & Removal*. en. June 2020. URL: https://ipv4.global/blog/ip-blocklist-blacklist/ (visited on 04/05/2024).

[32] Ren-Hung Hwang et al. "Cost optimization of elasticity cloud resource subscription policy". In: *IEEE Transactions on Services Computing* 7.4 (2013), pp. 561–574.

[33] Sotiris Ioannidis et al. "Implementing a distributed firewall". en. In: *Proceedings of the 7th ACM conference on Computer and Communications Security*. Athens Greece: ACM, Nov. 2000, pp. 190–199. ISBN: 978-1-58113-203-8. DOI: 10.1145/352600.353052. URL: https://dl.acm.org/doi/10.1145/352600.353052 (visited on 03/12/2024).

[34] Daiping Liu, Shuai Hao, and Haining Wang. "All Your DNS Records Point to Us: Understanding the Security Threats of Dangling DNS Records". en. In: *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*. Vienna Austria: ACM, Oct.
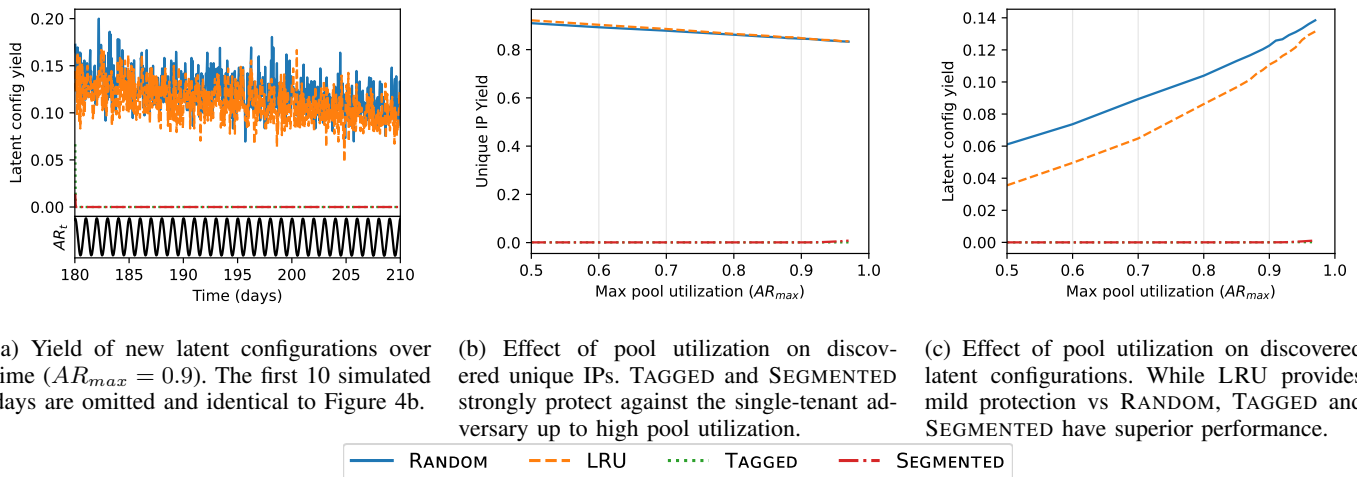
2016, pp. 1414–1425. ISBN: 978-1-4503-4139-4. DOI: 10.1145/2976749.2978387. URL: https://dl.acm.org/doi/10.1145/2976749.2978387 (visited on 09/14/2020).

[35] *Cloud Computing Services — Microsoft Azure*. en. URL: https : / / azure . microsoft . com / en - us/ (visited on 08/19/2021).

[36] Robert Moskowitz et al. *Address Allocation for Private Internets*. Request for Comments RFC 1918. Num Pages: 9. Internet Engineering Task Force, Feb. 1996. DOI: 10.17487/RFC1918. URL: https://datatracker.ietf.org/doc/rfc1918 (visited on 06/07/2022).

[37] Zack Whittaker and Carly Page. *Microsoft employees exposed internal passwords in security lapse*. en-US. Apr. 2024. URL: https : / / techcrunch . com / 2024 / 04 / 09/microsoft-employees-exposed-internal-passwords-security-lapse/ (visited on 04/16/2024).

[38] Eric Pauley, Paul Barford, and Patrick McDaniel. "DScope: A Cloud-Native Internet Telescope". In: *Proceedings of the 32nd USENIX Security Symposium (USENIX Security 2023, to appear)*. Anaheim, CA: USENIX Association, Aug. 2023.

[39] Eric Pauley et al. "Measuring and Mitigating the Risk of IP Reuse on Public Clouds". English. In: *2022 IEEE Symposium on Security and Privacy (SP)*. ISSN: 2375-1207. IEEE Computer Society, Apr. 2022, pp. 1523–1523. ISBN: 978-1-66541-316-9. DOI: 10 . 1109 / SP46214 . 2022 . 00094. URL: https : / / www . computer . org / csdl / proceedings - article / sp / 2022 / 131600b523/1CIO7rpcgSs (visited on 05/28/2022).

[40] *Progress Chef*. URL: https://www.chef.io/.

[41] *Puppet*. URL: https://puppet.com/.

[42] Chenhao Qu, Rodrigo N Calheiros, and Rajkumar Buyya. "Auto-scaling web applications in clouds: A taxonomy and survey". In: *ACM Computing Surveys (CSUR)* 51.4 (2018), pp. 1–33.

[43] Salvatore Sanfilippo. *Random notes on improving the Redis LRU algorithm*. URL: http://antirez.com/news/109 (visited on 06/14/2022).

[44] Claude E Shannon. "Communication in the presence of noise". In: *Proceedings of the IRE* 37.1 (1949), pp. 10–21.

[45] Sushant Sinha, Michael Bailey, and Farnam Jahanian. "Shades of grey: On the effectiveness of reputation-based blacklists". en. In: *2008 3rd International Conference on Malicious and Unwanted Software (MALWARE)*. Alexandria, VA, USA: IEEE, Oct. 2008, pp. 57–64. ISBN: 978-1-4244-3288-2. DOI: 10 . 1109 / MALWARE.2008.4690858. URL: http://ieeexplore.ieee.org/document/4690858/ (visited on 04/05/2024).

[46] Stephen Strowes. "IPv6 Adoption in 2021". en-US. In: *RIPE Labs* (2021). URL: https://labs.ripe.net/author/stephen_strowes/ipv6 - adoption - in - 2021/ (visited on 12/01/2021).

[47] Chunqiang Tang et al. "Holistic configuration management at facebook". In: *Proceedings of the 25th Symposium on Operating Systems Principles*. 2015, pp. 328–343.

[48] Abhishek Verma et al. "Large-scale cluster management at Google with Borg". In: *Proceedings of the European Conference on Computer Systems (EuroSys)*. Bordeaux, France, 2015.

[49] *Who is Spamhaus - the leader in IP & domain reputation data*. en. URL: https://www.spamhaus.org/who-is-spamhaus/ (visited on 03/12/2024).

[50] Andreas Wolke et al. "More than bin packing: Dynamic resource allocation strategies in cloud data centers". In: *Information Systems* 52 (2015), pp. 83–95.

[51] Danny Yuan et al. *Scryer: Netflix's Predictive Auto Scaling Engine — Part 2*. en. Apr. 2017. URL: https: / / netflixtechblog . com / scryer - netflixs - predictive - auto - scaling - engine - part - 2 - bb9c4f9b9385 (visited on 06/15/2022).

[52] Jian Zhang, Phillip A. Porras, and Johannes Ullrich. "Highly predictive blacklisting." In: *USENIX security symposium*. 2008, pp. 107–122. URL: https : / / www . usenix.org/legacy/events/sec08/tech/full_papers/zhang/zhang.pdf (visited on 04/05/2024).

APPENDIX A
SYMBOLOGY

| Symbol | Meaning |
|---|---|
| $a_i$ | Fourier amplitude |
| $AR_{max}$ | maximum pool allocation ratio |
| $AR_t$ | pool allocation ratio |
| $\mathcal{B}$ | distribution of tenant behaviors |
| $B_i$ | behavior of tenant $i$ |
| $d$ | duration |
| $d_a$ | duration of allocation |
| $d_{reuse}$ | minimum time duration before an IP can be reused |
| $d_v$ | duration of vulnerability |
| $\mathcal{I}$ | set of all IP addresses |
| $\mathcal{I}_{A_t}$ | set of IP addresses currently allocated |
| $p$ | probability |
| $p_c$ | probability of latent configuration |
| $S_{max}$ | maximum number of servers |
| $S_{min}$ | minimum number of servers |
| $t$ | time |
| $t_a$ | allocation time |
| $t_c$ | time of configuration dissociation |
| $t_{cd}$ | cooldown time |
| $t_r$ | release time |
| $T$ | tenant ID |
| $\alpha$ | segmentation cooldown multiplier |
| $\phi_i$ | Fourier phase |
| $\theta$ | opaque state |

APPENDIX B
SINGLE-TENANT ADVERSARY

For completeness, we simulate an adversary that is attempting to explore the IP space and discover latent configurations using only a single account. To do this, we model each simulation as in Section IV-C, but with a tenant count of one. For each simulated adversary, we seek to answer two questions:

(a) Yield of new latent configurations over time ($AR_{max} = 0.9$). The first 10 simulated days are omitted and identical to Figure 4b.

(b) Effect of pool utilization on discovered unique IPs. TAGGED and SEGMENTED strongly protect against the single-tenant adversary up to high pool utilization.

(c) Effect of pool utilization on discovered latent configurations. While LRU provides mild protection vs RANDOM, TAGGED and SEGMENTED have superior performance.

—— RANDOM　　- - - LRU　　······ TAGGED　　-·-·- SEGMENTED

Fig. 9: Modeling the single-tenant adversary.

1) How many unique IPs can the adversary discover based on their allocation scheme?
2) How many new latent configurations does the adversary discover associated with those IP addresses?

**Unique IPs.** Figure 9b displays the adversary's ability to discover new IPs across policies and allocation ratio. The RANDOM and LRU policies exhibit roughly identical behavior: IP yield is reduced as the pool gets smaller ($AR_{max}$ gets higher) because the adversary is more likely to receive the same IPs back. Likewise, TAGGED and SEGMENTED both almost completely eliminate the single tenant adversary's ability to discover new IPs. This is unsurprising, as both strategies tag IPs to the most recent tenant and reallocate those IPs back to the tenant. SEGMENTED exhibits a slight increase in adversarial IP yield at very high allocation ratios, as other tenant allocations interfere with the IPs tagged to the adversary–this does not occur in TAGGED because the LRU backup queue prevents the tenant's tagged IPs from being taken.

**Latent Configuration.** While an adversary might directly seek to observe a high number of IPs, the end goal is to discover IPs that actually have associated configuration. Our results (Figure 9c) demonstrate a marked difference here as well, with both TAGGED and SEGMENTED performing equivalently well against the single-tenant adversary. As seen in the non-adversarial scenario, LRU also slightly outperforms RANDOM as IP addresses are held in the pool longer before reuse, though this effect is diminished as the allocation ratio increases since the policies are best effort and must allocate some available IP to the adversary.

Our tool also allows us to model adversarial objectives over time (Figure 9a). Here, we see that the bulk of latent configuration discovered under TAGGED and SEGMENTED occurs early in the experiment. Beyond this, the pool returns the same IP addresses to the adversary and no latent configurations are discovered.

## APPENDIX C
## ADDITIONAL IMPLEMENTATION DETAILS

To empirically study the interaction of IP allocation behaviors, adversarial techniques, and deployable defenses, we develop the Elastic IP Simulator (EIPSIM). While the scale of cloud computing is astronomical, the allocation of IP addresses occurs and can be simulated independently. Although the space of these addresses is still quite large ($\approx 16\,\mathrm{M}$) for the largest AWS cloud region, this is still within the realm of exact simulation. To this end, EIPSIM simulates concrete tenant and cloud provider behavior at IP- and second-level granularity.

Within this architecture, *agents* perform the behavior of tenants, either by simulating tenant behavior or by replaying allocation traces from a previous run of the simulator or from actual tenants. The simulator fulfills IP allocation requests from these agents by referring to implementations of an *IP pool policy*. Each agent has the ability to allocate IPs under multiple tenant IDs, and the simulator treats these allocations as though they come from different tenants. While processing allocations, the simulator records statistics on the lifecycle of addresses, associated latent configuration, and adversarial objectives. Importantly, while these results are aggregated across addresses and tenants, they are a product of granular simulation of each tenant IP allocation.

### A. Tenant Agents

EIPSIM relies on tenant agents to perform the allocations of tenants. At each time-step ($1\,\mathrm{s}$) the simulator allows each agent to perform actions. Benign behaviors can be simulated by one of two agents:

- The *benign tenant agent* simulates the allocation behavior of tenants scaling cloud resources (Section III-A). For each tenant managed by the agent, and at each time-step, the agent checks if the tenant should allocate or release IP addresses, and passes these actions back to the simulator.
- The *file agent* allows loading of tenant behaviors from a time-series file. This file contains the timestamps and

tenant IDs of each IP allocation and release from either a previous run of EIPSIM or recorded from a live cloud environment (such as Google's clusters in Section V-B).

The *adversarial agent* is a specialized agent designed to simulate and analyze the behavior of a single- or multi-tenant adversary. The adversarial agent performs allocations exactly as it would on a real system (except that allocation requests are passed to the simulator instead of a cloud provider), and proceeds in several steps:

1) The agent requests IP addresses from the provider up to some quota (the maximum number of IPs it will hold at once). It records previous tenants and latent configuration associated with these for analysis (in reality an adversary would listen for network traffic or search DNS databases to identify these [39, 34, 12]).
2) The agent holds these IPs for a fixed duration. In EIPSIM, this is accomplished by performing no action when called by the simulator during this time.
3) The agent releases IPs that have been held for the specified duration back to the pool.
4) In the case of a multi-tenant adversary, new IPs are allocated under new tenant IDs. After a maximum tenant ID is reached, the adversary loops back to the initial tenant, simulating an adversary with access to only a fixed number of tenants.

While the techniques employed by the adversary could be performed by any cloud customer, the adversarial agent has access to the internal data structures of the simulator to be able to record time-series data on the functioning of the pool. For instance, when the agent allocates IP addresses it can access the list of previous tenants associated with that address (as this is used by our analysis).

### B. Allocation Policies

When the simulator receives a request for an IP address from a tenant, it forwards it to an allocation policy for servicing. While the simulator tracks what IP addresses are in use at any time, it is ultimate up to the policy to determine which free IP address is allocated to a given tenant. The policy receives the tenant ID associated with each allocation, but is not told the agent performing the request, or if the tenant is adversarial. The policy must also service all requests, though it may return any free IP for a given request.

The policy contains data structures that can track the history of a given IP address. For instance, the SEGMENTED policy tracks the most recent tenant ID for each IP, the cooldown time, and the average allocation durations of tenants. When a tenant requests an IP address, it heuristically samples available IPs that best conform to the policy based on this data. Considerations for deploying policies in practice are discussed in Section V-C.

### C. Extending the EIPSIM Framework

EIPSIM supports expansion to new policies, behaviors, and adversaries as academics and practitioners continue to study cloud IP allocation. EIPSIM defines `interfaces` between components, and new components can be added either as part of the EIPSIM package, or within a separate program that uses EIPSIM as a library. EIPSIM provides convenience functions to ease in the development of new components: for example, our studied allocation policies were implemented in an average of 71 lines of code, and new parameter sweep tests can be built on top of EIPSIM in around 70 lines of code. We expect that, by encouraging the development of new components on top of our framework, the community can reach a unified means to compare threat models and defenses. EIPSIM also supports allocation traces collected by cloud providers through custom agents. Practitioners can directly read allocations as tuples of $(T, t_a, t_r)$ and use EIPSIM to simulate adversarial and pool behavior.

## APPENDIX D
### DATA AVAILABILITY

EIPSIM is proposed as a theoretical and applied framework for studying the security of IP address allocation. As such, code and evaluation artifacts are being made available to reproduce results and support further study by practitioners. Code for EIPSIM, the evaluation, and figures is released as open source. The `clusterdata-2019` dataset (used in Section V-B) is freely available, and data preparation code is provided. Data used in Section V-A is provided by the authors of the referenced prior work, and cannot be openly shared for privacy reasons (as the underlying data refers to concrete instances of vulnerable configurations).