

Covert Communication In The Cloud With Lambdas

Anil Yelam
UC San Diego

Shibani Subbareddy*
Salesforce Inc.

Keerthana Ganesan*
Facebook Inc.

Stefan Savage
UC San Diego

Ariana Mirian
UC San Diego

ABSTRACT

“Serverless” cloud services, such as AWS lambdas, are one of the fastest growing segment of the cloud services market. These services are light-weight and provide more flexibility in scheduling and cost, which contributes to their popularity, however the security issues associated with serverless computing are not well understood. In this work, we explore the feasibility of constructing a practical covert channel from lambdas. We establish that a fast co-residence detection for lambdas is key to enabling such a covert channel, and proceed to develop a reliable, and scalable co-residence detector based on the memory bus hardware. Our technique enables dynamic neighbor discovery for co-resident lambdas and is incredibly fast, executing in a matter of seconds. We evaluate our approach for correctness and scalability, and use it to establish covert channels and perform data transfer on AWS lambdas. We show that we can establish tens of individual covert channels for every 1000 lambdas, and each of those channels can send data at a rate of 200bps, demonstrating the feasibility of covert communication.

CCS CONCEPTS

• Security and privacy → Virtualization and security.

KEYWORDS

cloud, cartography, serverless, coresidency, covert channels

ACM Reference Format:

Anil Yelam, Shibani Subbareddy, Keerthana Ganesan, Stefan Savage, and Ariana Mirian. 2021. Covert Communication In The Cloud With Lambdas. In *Proceedings of the Web Conference 2021 (WWW '21), April 19–23, 2021, Ljubljana, Slovenia*. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3442381.3450100>

1 INTRODUCTION

Over the last decade, organizations have increasingly offloaded their data processing and storage needs to third-party “cloud” platforms. However, the economics of cloud platforms is predicated on high levels of statistical multiplexing and thus *co-tenancy* – the

*Work done while at UC San Diego

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

WWW '21, April 19–23, 2021, Ljubljana, Slovenia

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-8312-7/21/04.

<https://doi.org/10.1145/3442381.3450100>

contemporaneous execution of computation from disparate customers on the same physical hardware – is the norm. The risks associated with this arrangement, both data leakage and interference, are well-appreciated and have generated both a vast research literature (starting with Ristenpart et al. [19]) as well a wide-array of technical isolation countermeasures employed by cloud platform providers. Most of this work has focused squarely on the risks of information channels between long-lived, heavy-weight virtual machines (“instances” in Amazon parlance) used to virtualize the traditional notion of dedicated network-connected servers.

However, over the last six years, most of the largest cloud providers have introduced a *new* “serverless” service modality that executes short-lived, lightweight computations on demand (e.g., Amazon’s Lambda [14], Google’s Cloud Functions [10] and Microsoft’s Azure functions [5]). These services, by design, use lighter-weight tenant isolation mechanisms (so-called “micro-VMs” or containers) as well as a fixed system environment to provide low-latency startup and a reduced memory footprint. In return, serverless systems can support even higher levels of statistical multiplexing and thus can offer significant cost savings to customers whose needs are able to match this model (e.g., event-driven computations with embedded state). However, the security issues associated with serverless computing are far less well understood than their heavier weight brethren. While the transient and dynamic nature of serverless computing pose inherent challenges for attackers, their low-cost and light-weight isolation potentially present offer new points of purchase as well.

In our work, we explore these issues through the lens of a singular question: can a practical covert channel be constructed entirely from existing “serverless” cloud services¹?

Covert channels provide a means of transmitting data that bypasses traditional monitoring or auditing – typically by encoding data into some resource access that is not normally deemed a communications medium but is externally visible. In virtualized environments, covert channels typically involve some shared resource (e.g. a cache) for which contention provides a means of signaling. In the serverless context, the threat model is that an adversary is able to launch, or inject code into, lambdas from inside a target organization and wishes to communicate information to parties outside the organization (i.e., to their own lambdas) without offering any clear evidence of such (e.g., opening network connections, etc.)

However, the serverless context presents a number of unique challenges for implementing covert channels. First, the physical location of a lambda is unknown, as the scheduling and placement of lambdas is managed by the cloud service provider. Thus, there is no way to arrange that a sending lambda and a receiving lambda will

¹We will use the term lambdas to stand for all such services going forward.

execute on the same physical hardware, *let alone at the same time*. Second, given this reality, any serverless covert communications protocol must repeatedly launch lambdas in the hope that at least two sending and receiving lambdas are co-resident on the same hardware at the same time. The extent to which this is practicable, on existing cloud platforms and reasonable cost, is unknown. Third, it is not enough to simply *achieve* co-residency, but any lambdas lucky enough to be co-resident must be able to quickly determine this fact, and then use the balance of their limited lifetimes to effect communications. Finally, since rendezvous in a serverless system is inherently statistical, any such protocol must anticipate the potential for interference (i.e., when multiple sending lambdas happen to be co-resident with one or more receiving lambdas).

In this paper we address each of these issues in turn and demonstrate the feasibility of covert communication entirely in the context of the Amazon’s serverless cloud platform. In particular, we make three key technical contributions:

- **Fast co-residence detection.** Leveraging the memory-bus contention work of Wu et. al [26], we develop and implement a lambda co-residence detector that is generic, reliable, scalable and, most importantly, fast, executing in a matter of seconds for thousands of concurrent lambdas.
- **Dynamic neighbor discovery.** We present a novel protocol for the co-resident lambdas to communicate their IDs using a hardware-based covert channel and discover one another. A key enabler of our co-residence detector, the protocol also helps a sender or receiver lambda to enumerate *all* its co-resident neighbors, a requirement to avoid unwanted communication interference while performing covert communication.
- **Covert channel demonstration** We use our co-residence detector to establish lambda covert channels on AWS. We then perform a study on the feasibility of covert communication by 1) estimating the capacity of each channel and 2) measuring coresidence density – that is, for a given number of short-lifetime lambdas launched at a point in time, how many would be expected to become co-resident? We conduct these measurements across a range of Amazon data centers to establish that there is ample co-residence of lambdas and covert communication is practicable.

Our implementation is publicly available at <https://github.com/anikyelam/columbus>.

2 BACKGROUND & RELATED WORK

We begin with a brief background on related topics.

2.1 Lambdas/Serverless Functions

We focus on serverless functions in this paper, as they are one of the fastest-growing cloud services and are less well-studied from a security standpoint. Offered as lambdas on AWS [14], and as cloud functions on GCP [10] and Azure [5], these functions are of interest because they do not require the developer to provision, maintain, or administer servers. In addition to this low maintenance, lambdas are much more cost-efficient than virtual machines (VMs) as they allow more efficient packing of functions on servers. Moreover, lambdas execute as much smaller units and are more ephemeral

than virtual machines. For example, on AWS, the memory of lambdas is capped at 3 GB, with a maximum execution limit of 15 minutes. As with other cloud services, the user has no control over the physical location of the server(s) on which their lambdas are spawned.

While lambdas are limited in the computations they can execute (typically written in high-level languages like Python, C#, etc), they are conversely incredibly lightweight and can be initiated and deleted in a very short amount of time. Cloud providers run lambdas in dedicated containers with limited resources (e.g., Firecracker [1]), which are usually cached and re-used for future lambdas to mitigate cold-start latencies [2]. The ephemeral nature of serverless functions and their limited flexibility increases the difficulty in detecting co-residency, as we will discuss later. While previous studies that profiled lambdas [24] focused on the performance aspects like cold start latencies, function instance lifetime, and CPU usage across various clouds, the security aspects remain relatively understudied.

2.2 Covert Channels in the Cloud

In our attempt to shed light on the security aspects of lambdas, we focus particularly on the feasibility of establishing a reliable covert channel in the cloud using lambdas. Covert channels enable a means of transmitting information between entities that bypasses traditional monitoring or auditing. Typically, this is achieved by communicating data across unintended channels such as signalling bits by causing contention on shared hardware media on the server [16, 18, 25–27]. Past work has demonstrated covert channels in virtualized environments like the clouds using various hardware such as caches [19, 27], memory bus [26], and even processor temperature [16].

Memory bus covert channel Of particular interest to this work is the covert channel based on memory bus hardware introduced by Wu et al. [26]. In x86 systems, atomic memory instructions designed to facilitate multi-processor synchronization are supported by cache coherence protocols as long as the operands remain within a cache line (generally the case as language compilers make sure that operands are aligned). However, if the operand is spread across two cache lines (referred to as “exotic” memory operations), x86 hardware achieves atomicity by locking the memory bus to prevent any other memory access operations until the current operation finishes. This results in significantly higher latencies for such locking operations compared to traditional memory accesses. As a result, a few consecutive locking operations could cause contention on the memory bus that could be exploited for covert communication. Wu et al. achieved a data rate of 700 bps on the memory bus channel in an ideal laboratory setup.

Achieving such ideal performance, however, is generally not possible in cloud environments. Cloud platforms employ virtualization to enable statistical multiplexing, and communication on the covert channel may be affected by 1) scheduler interruptions as sender or the receiver may only get intermittent access to the channel and 2) interference from other non-participating workloads. This may result in errors in the transmitted data and require

additional mechanisms like error correction [26] to ensure reliable communication.

2.3 Co-residence Detection

In the cloud context, enabling communication over traditional covert channels comes with an additional challenge of placing sender and receiver on the same machine. However, such co-residency information is hidden, even if the entities belong to the same tenant. Past research has used various strategies to achieve co-residency in order to demonstrate various covert channel attacks in the cloud. Typically, the attacker launches a large number of cloud instances (VMs, Lambdas, etc.), following a certain launch pattern, and employs a co-residence detection mechanism for detecting if any pair of those instances are running on the same machine. Traditionally, such detection has been based on software runtime information that two instances running on the same server might share, like public/internal IP addresses [19], files in *procfs* or other environment variables [24, 26], and other such logical side-channels [21, 29].

As virtualization platforms moved towards stronger isolation between instances (e.g. AWS’ Firecracker VM [1]), these logical covert-channels have become less effective or infeasible. Furthermore, some of these channels were only effective on container-based platforms that shared the underlying OS image and were thus less suitable for hypervisor-based platforms. This prompted a move towards using hardware-based covert channels, such as the ones discussed in section 2.2, which can bypass software isolation and are usually harder to fix. For example, Varadarajan et al. [22] use the memory bus covert channel to detect co-residency for EC2 instances. However, as we will show, their approach does not extend well to lambdas as it is slow and not scalable.

3 MOTIVATION

We discuss the covert channel attack scenario that we are targeting, the challenges lambdas would pose in enabling such an attack and motivate the need for a co-residence detector for lambdas which is going to be our main focus in the rest of the paper.

Threat Model Covert channel attacks require an “insider” to send the data over a covert medium for exfiltration. We assume that the attacker uses social engineering techniques or some other means (beyond the scope of this work) to introduce such insiders in the victim system. In case of lambdas, this insider code could be in the lambda itself or in a system that controls lambda deployments for an organization; and already possesses the sensitive data that needs to be exfiltrated. We further assume that the attacker has the knowledge of the cloud region where the victim is operating and can deploy lambdas under its own account.

In a typical attack, the attacker launches a set of lambdas (receivers) in the cloud region where the victim lambdas (senders) are expected to operate. The attacker and (compromised) victim lambda(s) can then work together² to exchange the data over a covert channel like the memory bus hardware discussed in earlier section. However, as mentioned earlier, there are few unique challenges before we can use a traditional covert channel in the cloud.

²We do not explicitly differentiate attack and victim lambdas hereafter as they are all assumed to be in attacker’s control

We need: 1) to colocate the sender and receiver on the same server which requires co-residence detection and 2) to handle the interruptions on such channel introduced by noisy neighbors and inconsistent access to the channel due to scheduling.

While these challenges have been handled for other cloud platforms like VMs (see related work), lambdas are inherently different in that they have very short lifetimes. A covert channel between two co-resident lambdas will not last very long. However, while lambdas are not persistent, it is trivial and cheap to launch lambdas in large numbers at once and establish multiple rendezvous points to allow for more covert communication. Additionally, lambdas are also densely packed than VMs exacerbating the noisy neighbor problem.

The ephemeral, numerous and dense nature of lambdas require a fast, scalable, and reliable co-residence detector. Such a detector will allow the attacker to identify all the servers with two or more co-resident lambdas and establish a rendezvous on each such server. Moreover, the detector should precisely identify which lambdas are co-resident allowing the attacker to pick any two lambdas on a given machine, and use the covert channel without interference from neighbors.

4 CO-RESIDENCE DETECTOR FOR LAMBDA

In this section, we present a novel co-residence detector for lambdas, previous solutions to this problem, and the unique challenges we faced with lambdas and how we address each of them.

4.1 Specification

Given a set of cloud instances (VMs, Containers, Functions, etc) deployed to a public cloud, a co-residence detection mechanism should identify, for each pair of instances in the set, whether the pair is running on the same physical server at some point. Paraphrasing Varadarajan et al. [22], for any such mechanism to be useful across a wide range of launch strategies, it should have the following properties:

- **Generic** The technique should be applicable across a wide range of server architectures and software runtimes. In practice, the technique would work across most third-party cloud platforms and even among different platforms within a cloud.
- **Reliable** The technique should have a reasonable detection success with minimal false negatives (co-resident instances not detected) and even less false positives (non co-resident instances categorized as co-resident).
- **Scalable** A launch strategy may require hundreds or even thousands of instances to be deployed, and must scale such that the technique will detect all co-resident pairs at a reasonable cost.

We add another property to that list which is relevant to lambdas:

- **Fast** The technique should be fast, preferably finishing in the order of seconds. As lambdas are ephemeral (with some clouds restricting their execution times to as low as a minute), the technique should leave ample time for other activities that make use of the resulting co-resident information.

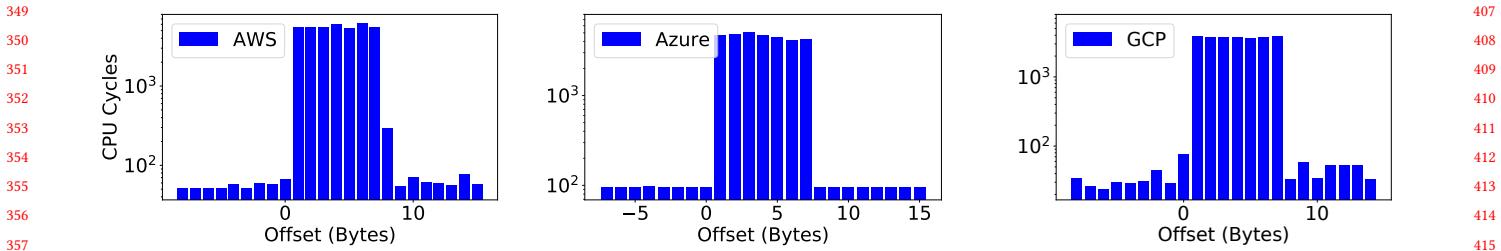


Figure 1: The plots show the latencies of atomic memory operations performed on an 8B memory region as we slide from one cache line across the boundary into another on AWS, Azure, and Google (GCP) respectively. The latencies are orders of magnitude higher when the 8B region falls across the two cache lines (offsets 0–7B), demonstrating the presence of the memory bus covert channel on all these cloud providers.

4.2 Design

4.2.1 Using Covert Channels for Generality. As mentioned in section 2.3, co-residence detection was previously done using unique software identifiers that revealed underlying server to the tenants. Such software identifiers present the fastest and perhaps most reliable(Anil ►reliable as in trustworthy information◀) way for co-residence detection as all it takes is for each lambda to quickly "read" the identifier and communicate this information with all other lambdas (through the network or shared remote storage) to identify the neighbors. However, such information can be easily obfuscated by platform providers; for example, currently there is no such identifier available for AWS lambdas. So we turn to hardware-based covert channels that are, by definition, also accessible to all the tenants. They are also generally more difficult to remove, and more pervasive too, given that hardware is more homogenous across computing platforms than software.

4.2.2 Covert Channels Cannot Support Multiple Parties. TODO ►think about this paragraph◀ As covert channels can send information, naturally we can expect to use the channel to communicate identity information (like unique ids) between co-resident lambdas all of which have access to the channel. However, covert channels generally presume that the sender and the receiver are already known and that there will be only two parties performing communication. But when multiple parties share a channel, we would need a channel access control, such as MAC protocols in Ethernet networks, to arbitrate access to the channel and handle collisions as they are in the same collision domain. However, covert channels are often binary channels (i.e., parties can either only send or receive a bit but not both) with no capability for collision detection and have very limited bandwidth (often only tens of bits per second) that such channel arbitration can be infeasible or complex enough to present significant overhead on channel bandwidth.

4.2.3 A fast ID broadcast designed for covert channels. For co-residence detection, lambdas only need to communicate their IDs with one another and as such, we don't need the channel to be very general and expressive. Thus, we assume that each lambda involved has a unique fixed-length (say N) bit-string corresponding to its ID that must be communicated. In regards to the covert

Algorithm 1 Neighbor discovery protocol

```

1: sync_point ← Start time for all instances
2: ID ← Instance ID
3: N ← Number of bits in ID
4: advertising ← TRUE
5: instances ← {}
6: WAIT_TILL(sync_point)
7: while id_read do
8:   slots ← 0
9:   id_read ← 0
10:  participating ← advertising
11:  while slots < N do
12:    bit ← slotsth most significant bit of ID
13:    if participating and bit then
14:      WRITE_BIT() (Alg. 2)
15:      bit_read ← 1
16:    else
17:      bit_read ← READ_BIT() (Alg. 3)
18:      if bit_read then
19:        participating ← FALSE
20:      end if
21:    end if
22:    id_read ← 2 * id_read + bit_read
23:    slots ← slots + 1
24:  end while
25:  if id_read = ID then
26:    advertising ← FALSE
27:  end if
28:  instances ← instances ∪ {id_read}
29: end while
30: return instances

```

channel, we only assume that any lambda with access to the channel can choose to either send or listen for a bit, and if at least one lambda chooses to send a 1-bit, all the listeners would record a 1-bit. We also assume that lambdas can synchronize themselves between sending or listening a bit in each time slot. We show both these assumptions to be reasonable in section 4.3. As such, we propose a communication protocol that efficiently broadcasts just these bit-strings on channels with such restricted properties.

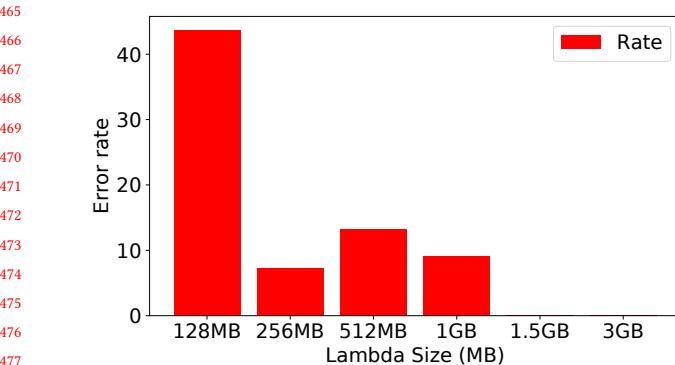


Figure 2: This figure presents the error rate (as a fraction of 1000 lambdas deployed) for different lambda sizes in the AWS Middle-East region. **TODO** ► make it print friendly ◀

We divide the total running time of the protocol into phases, with each phase executing for an interval of N bit-slots. Each phase has a set of *participating* lambdas, which in the first phase would be all of the lambdas. In each bit-slot K of N slots in a phase, every participating lambda broadcasts a bit if the K^{th} bit of its bit-string (ID) is 1, otherwise it listens for a 0 or 1. If a lambda senses a 1 while listening, it stops participating, and listens for the rest of the phase. Thus, only the lambdas with the highest ID among the initial set of participating lambdas continues broadcasting until the end of the protocol, effectively advertising its full ID to the rest of the (now listening) lambdas. In the next phase, the lambda with the previously highest ID now only listens, allowing the next highest lambda to advertise its ID, and so on. If the IDs are unique, there will always be only one lambda that broadcasts in every phase. The protocol ends after x phases (where x is number of co-resident lambdas), when none of the lambdas broadcast for N consecutive bit-slots. The pseudo-code of the protocol is provided in Algorithm 1.

4.2.4 Time Complexity. Assuming N total deployed lambdas to the cloud, the bit-string needs to be $\log_2 N$ bits to uniquely identify each lambda. If a maximum K of those lambdas are launched on the same server, the protocol executes for K phases of $\log_2 N$ bit-slots each, taking $(K + 1) * \log_2 N$ bit-slots for the whole operation. In fact, it is not necessary to run the protocol for all K phases. After the first phase, all the co-resident lambdas would know one of their neighbors (as each phase reveals the ID of the biggest participating lambda to others). If we use IDs that are globally unique, all the co-resident lambdas will see the same ID. The lambdas can then exchange these IDs offline (e.g., through the network) to infer the rest of their neighbors. This simplification removes the dependency on number of co-resident lambdas (K) and decreases the complexity to $O(\log_2 N)$.

4.3 Implementation

We implemented the above protocol using a covert channel based on memory bus hardware that is available for lambdas. We show how we used the hardware to reliably send and listen for bits in order to meet the requirements for the protocol.

Algorithm 2 Writing 1-bit from the sender

```

now ← time.now()
end ← now + sampling_duration
address ← cache_line_boundary - 2
while now < end do
    _ATOMIC_FETCH_ADD(address)
    now ← time.now()
end while

```

4.3.1 Introducing Memory Bus Covert Channel. We utilize the memory bus covert channel described in section 2.2 as it exploits a fundamental hardware vulnerability that is present across all generations of x86 hardware. Historically, multiple public cloud services have been vulnerable to this channel [21, 32], and we find that they are still vulnerable today. To demonstrate the presence of the vulnerability, we measure the latency of atomic operations on a 8B memory region as we slide the region from one cacheline into another across the cacheline boundary. We perform this experiment on three major cloud platforms (AWS, Google and Microsoft Azure) and show the latencies observed in Figure 1. From the figure, we can see that all three cloud platforms still exhibit a significant difference in latencies for the "exotic" memory locking operations where the memory region falls across cacheline boundary. When compared to regular memory accesses, it demonstrates the presence of this covert channel on all of them. Moreover, we were able to execute these experiments on serverless function instances. Since lambdas have runtimes that are generally restricted to high-level languages (that prevent the pointer arithmetic required to perform these exotic operations), we used the unsafe environments on these clouds — C++ on AWS, Unsafe Go on GCP, Unsafe C# On Azure. This shows the applicability of using the covert channel across different kinds of cloud instances as well.

4.3.2 Sending a bit. Senders and receivers can accurately communicate 0-bits and 1-bits by causing contention on the memory bus. To communicate a 1-bit, the sender instance causes contention on the memory bus by locking it using the special memory locking operations (discussed in section 2.2). Pseudo-code for the sender instance is shown in Algorithm 2.

4.3.3 Listening for a bit. The receiver can simply sample the memory bus for contention, inferring whether the communication is a 1-bit (when contention is observed) or a 0-bit (when contention is not observed). There are two ways to listen for contention. When the memory bus is locked, any non-cached memory accesses will queue and therefore see higher latencies. The receiver can then continually make un-cached memory accesses (referred to as the *memory probing* receiver in previous literature [22]) and observe a spike in their latencies to detect contention. On the other hand, the receiver can also detect memory bus contention by using the same memory locking operations as the sender (referred to as *memory locking* receiver) to probe the memory bus. Since only one processor core can lock the memory bus at a given time, any other concurrent locking operation will see higher latency.

Mechanism Of these two methods, we decide to use the memory locking receiver for our experiments. Since memory probing involves regular (un-cached) memory accesses, it can be done on

Algorithm 3 Reading a bit in the receiver

```

581 1: now ← time.now()
582 2: end ← now + sampling_duration
583 3: sampling_rate ← num_samples/sampling_duration
584 4: address ← cache_line_boundary - 2
585 5: samples ← {}
586 6: while now < end do
587 7:   before ← RDTSC()
588 8:   __ATOMIC_FETCH_ADD(address)
589 9:   after ← RDTSC()
590 10:  samples ← samples ∪ {(after - before)}
591 11:  wait until NEXT_POISSON(sampling_rate)
592 12:  now ← time.now()
593 13: end while
594 14: ks_val ← KOLMOGOROV_SMIRINOV(samples, baseline)
595 15: return ks_val < ksvalue_threshold
596
597
598
599

```

multiple receivers concurrently without affecting each other (due to the high memory bandwidth), which prevents noise in measurements. This is an important attribute, as memory locking receivers must contend with this noise. However, bypassing multi-levels of caches in today's servers to perform memory accesses with reliable consistency is a challenging task. Even with a reliable cache-bypassing technique, the variety of cache architectures and sizes that we encounter on different clouds would make tuning the technique to suit these architectures an arduous task while reducing the applicability of our overall co-residence detection mechanism.

Sampling frequency Another challenge for our protocol is determining an adequate sampling frequency. Ideally, a memory locking receiver would loop locking operations and determine contention in real-time by identifying a decrease in the moving average of the number of operations. Note that, in this case, there is essentially no difference between the sender and receiver (i.e., both continually issue locking operations) except that the receiver is taking measurements. This is adequate when there is a single sender and receiver [22], but when there are multiple receivers, the mere act of sensing the channel by one receiver causes contention and other receivers cannot differentiate between a silent (0-bit) and a locking (1-bit) sender. To avoid this, we space the sampling of memory bus such that no two receivers would sample the bus at the same time, with high probability. We achieve this by using large intervals between successive samples and a poisson-sampling to prevent time-locking of receivers. We determined that a millisecond poisson gap between samples is reasonable to minimize noise due to collisions in receiver sampling 1, assuming ten co-resident receivers and a few microsecond sampling time. TODO | ▶ R1 needed clarification◀

Sample Size In addition to adequate sampling frequency, we must also determine sample size. A receiver can confirm contention with high confidence with only a few samples, assuming that the sender is actively causing contention on the memory bus and the receiver is constantly sampling the memory bus throughout the sampling duration. However the time-sharing of processors produces difficulties. The sender is not continually causing contention, and neither is the receiver sensing it, as they are

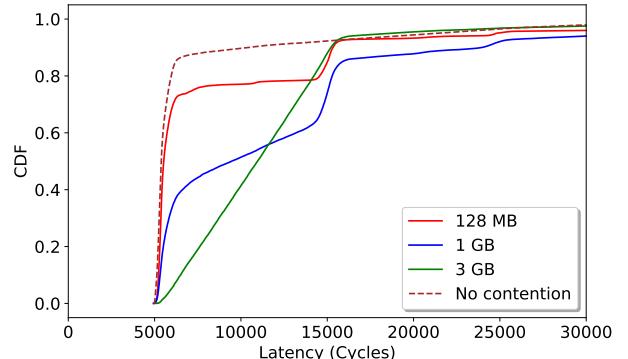


Figure 3: We present a CDF of latencies observed by 128 MB, 1 GB and 3 GB Lambdas during contention. The 128 MB lambda pair sees less contention due to more context switching, whereas the 1 GB and 3 GB lambdas see progressively more contention compared to the baseline, which we attribute to their relative stability on the underlying physical cores. TODO | ▶ make it print friendly◀

context-switched by the scheduler, which runs other processes. Assuming that the sender and receiver are running on different cores, the amount of time they are actively communicating depends on the proportion of time they are allocated on each core and how they are scheduled.

To illustrate such behavior, we run a sender-receiver pair using lambdas [14] of various sizes on AWS, and compare the distribution of latencies seen by the receiver during the contention in each case. Figure 3 shows that the much smaller 128 MB lambdas (which probably share a CPU core and are thus context-switched) exhibit less active communication than the bigger 3 GB lambdas (which may run on dedicated cores). This means that smaller instances that tend to share processor cores with many other instances may need to pause for more time and collect more samples to make up for lost communication due to scheduling.

Overcoming noise Along with context switching and sensing noise, there are other imperfections in the measurement apparatus that may cause noise. For example, we use the difference in readings from the timestamp counter of the processor (RDTSC) before and after the locking operation to measure the latency of the operation in cycles. If the receiver process is context-switched in between the timer readings (e.g., at line eight in Algorithm 3), the latency measured from their difference will be orders of magnitude higher as it includes the waiting time of the receiver process in the scheduler queue - which we believe is what contributes to the long tail in Figure 3. To overcome missed samples and noise, we record hundreds of samples and compare it to the baseline distribution of latencies sampled without contention. We then need to compare and differentiate the observed sample of latencies from the baseline to establish contention. To do this, we use a variant of the two-sample Kolomogorov-Smirinov (KS) test, which typically

697 compares the maximum of the absolute difference between empirical
 698 CDFs of samples (in our variant, we take the *mean* of the absolute
 699 difference instead of the maximum to reduce sensitivity to outliers). Using this measure, we can categorize a KS-value above
 700 a certain threshold as a 1-bit (contention) and a value below the
 701 threshold as 0-bit (baseline).

702 To determine the KS-threshold, we deploy a large number of
 703 lambdas across AWS regions. Some of these lambdas cause con-
 704 tention (aka senders) while others observe contention by collect-
 705 ing samples of latencies (aka receivers). Each of the samples may
 706 or may not have observed contention depending on whether the
 707 receiver was co-resident with a sender lambda (an unknown at this
 708 point). We then calculate the KS-value for each sample against the
 709 baseline and plot a CDF of these values for lambdas of different
 710 sizes in Figure 4. Ideally, we expect a bimodal distribution (stepped
 711 CDF) with the upper and lower peaks corresponding to samples
 712 that have and have not seen contention, and a big gap between
 713 the two (long step). Fortunately, we observe this differentiation
 714 with larger lambda sizes (which allows us to choose a clear thresh-
 715 old), but we do not observe a clear differentiation with smaller
 716 lambdas, where scheduling instability causes lossy communication
 717 (discussed in 4.3.3). This trend also reflects in the reliability of our
 718 technique across various lambda sizes, as we will show in our eval-
 719 uation. Based on the plot, we picked a KS-threshold at 3.0 which
 720 seems to be consistent across AWS regions, suggesting that this
 721 threshold is a platform constant.

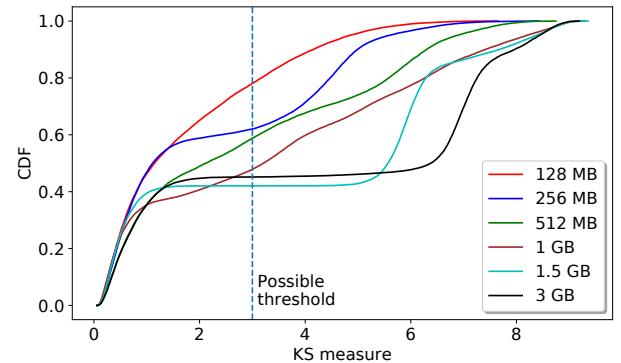
722 We present the pseudo-code of a receiver lambda in Algorithm 3,
 723 which includes all the challenges and subsequent solutions dis-
 724 cussed thus far.

725 **4.3.4 Synchronization.** A major enabler of our protocol in section
 726 4.2.3 is the ability to synchronize all the co-resident lambdas when
 727 sending and receiving bits. As all these lambdas are running on the
 728 same physical server, they share the server’s clock. On AWS, for
 729 example, we observe that the system clock on lambdas is precise up
 730 to nanoseconds. Assuming that clocks between different lambdas
 731 only exhibit a drift in the order of microseconds, sampling at a
 732 millisecond scale should provide us a margin for synchronization
 733 mismatch. Since we do not observe any synchronization-related
 734 noise in our results, we believe that this is a reasonable assumption.

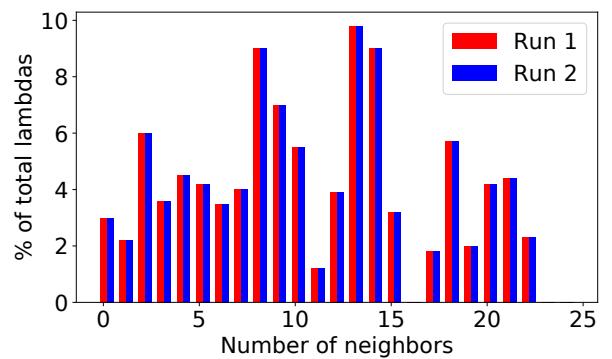
735 4.4 Evaluation

736 We examine our co-residence detector with respect to reliability
 737 and scalability, the desirable detection properties mentioned in
 738 section 4. We run all of our experiments with AWS lambdas [3].
 739 Though we decide to focus on only one of the cloud providers as a
 740 case study, we have previously shown in section 4 that this covert
 741 channel exists on the other clouds, and thus these experiments can
 742 be replicated on their serverless functions as well. We use the C++
 743 runtime in AWS lambdas as it allows pointer arithmetic that is re-
 744 quired to access the covert channel. **TODO** ► add note about this ful-
 745 filling fast requirement? ◀

746 **4.4.1 Setup.** For each experiment, we deploy a series of instances
 747 from an AWS lambda account. Once deployed, each instance par-
 748 ticipates in the first phase of the protocol as noted in section 4.2.4,
 749 thereby learning the largest ID of their neighbors. As bit-flip errors



750 **Figure 4:** We present a CDF of KS values observed for
 751 various lambda sizes. A bimodal distribution with a longer step
 752 allows us to pick a KS-threshold that enables our technique
 753 to differentiate between 0-bit and 1-bit with high confidence.



754 **Figure 5:** This figure shows the fraction of lambdas by
 755 the number of neighbors they identify for two inde-
 756 pendent runs that use same set of underlying AWS containers.
 757 The perfect correlation shows that both runs depict the co-
 758 residence status of those containers regardless of the lamb-
 759 das that ran on them, providing evidence for the correctness
 760 of our approach. **TODO** ► make it print friendly ◀

761 are possible, we repeat the same phase for two more (independent)
 762 "rounds" and take the majority result to record the ID seen by this
 763 instance. If all three rounds result in different IDs, we classify this
 764 instance as erroneous and report it in the error rate. We group all
 765 the instances that saw the same ID as successful and neighbors.
 766 We repeat the experiments for different lambda sizes and in vari-
 767 ous cloud regions. **TODO** ► add cost here ◀

768 **4.4.2 Reliability.** We consider the results of the technique reliable
 769 when 1) most of the deployed instances successfully see the same
 770 result in majority of the independent rounds (indicating lesser bit-
 771 flip errors) and 2) the resulting co-resident groups we see match the
 772 ground truth. For goal #1, we ran an experiment with 1000 AWS
 773 regions.

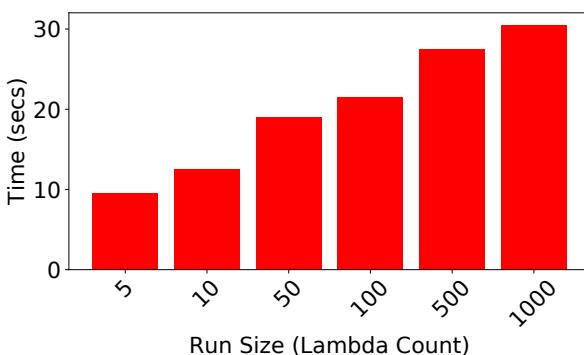


Figure 6: We present the average execution time of the lambdas for co-resident runs with a varying number of lambdas. The execution time increases logarithmically with the number of lambdas demonstrating the scalability of co-residence detection with our technique.

lambdas and compared the error rate across different lambda sizes (the error rate indicates the fraction of these 1000 lambdas that did not have a majority result). Figure 2 indicates that smaller lambdas exhibit many more errors. This is expected because, as discussed in section 4.3.3, these lambdas experience lossy communication making it harder for our technique to sense contention. Lambdas that are 1.5 GB and larger, though, exhibit a 100% success rate. **TODO**

►*rebuttal: clarify false positives and negatives*◀

Correctness **TODO** ►*promised clarity here in rebuttal*◀ To determine correctness, we require ground truth on which instances are co-resident with one another. While such information is not available, we are able to ascertain correctness of our approach by utilizing an AWS caching mechanism. On AWS, each lambdas runs in a dedicated container (sandbox). After execution, AWS caches these containers in order to reuse them [2] for repeat lambdas and mitigate “cold start” latencies. For C++ lambdas, we found that the data structures declared in the global namespace are tied to containers and are not cleared on each lambda invocation, so we can use a global array to record all the lambdas that were ever executed in a particular container. This indicates that, for a given lambda, we can precisely note all the lambdas that previously ran in the same container (aka predecessors). Using this, we are able to validate that identical experiments repeated within minutes of one another will use the same set of underlying containers for running the deployed lambdas. Since lambda co-residence is essentially co-residence of their containers, and given that containers persist across experiments that are executed within minutes of one another, lambda co-residence results must agree with the co-residence of their underlying containers for true correctness.

To demonstrate the correctness of our technique using this insight, we run an experiment with 1000 1.5GB cold-started lambdas (ID'ed 1 to 1000) in one of densest AWS regions (AWS MiddleEast), which resulted in many co-resident groups. We repeat the experiment within a few seconds, thereby ensuring that all 1000 lambdas are warm-started on the second trial (i.e., they use the same set

of containers from the previous experiment). For each co-resident group of lambdas in the latter experiment, we observed that their predecessor lambdas (that used the same set of containers) in the former experiment formed a co-resident group as well. That is, while the lambdas to the underlying container mapping is different across both experiments, the results of the experiments agree perfectly on the container colocation. Figure 5 shows that both experiments saw the same number of co-residing groups of different sizes. This proves the correctness of the results of our mechanism.

4.4.3 Scalability. One of the key properties of this technique is its execution speed. Since communicating each binary bit of the ID takes one second, we are able to scale the technique logarithmically with the number of lambdas involved. Figure 6 shows this result with experiments involving different number of lambdas. For example, in an experiment with 1000 lambdas, each lambda can find its neighbors within a minute of its invocation, leaving ample time for the attacker to then establish the covert channel and use it to send information. The logarithmic scale of our method also indicates that the cost per lambda scales logarithmically, making neighbor detection cost-effective. **TODO** ►*note here about sublinear time?*◀

5 PRACTICALITY OF COVERT COMMUNICATION

In this section, we perform a study to demonstrate the practicality of data transfer using the lambda covert channels discovered with our co-residence detector. The amount of information that can be transferred depends on two factors: 1) the capacity of each channel and 2) the number of co-resident clusters of lambdas, or rendezvous points, that materialize during the attack. We first produce an estimate on the capacity of the covert channels established, and then examine the co-residence density in various AWS regions to understand the number of rendezvous points and factors that affect it.

5.1 Covert Channel Capacity

Once co-residence between any two lambdas is established, the attacker can then use the same memory bus hardware to perform covert communication. Wu et al. [26], who first introduced covert channel based on this hardware channel, also presented an efficient and error-free communication protocol targeting cloud-based platforms like VMs. While such a protocol should theoretically work for lambdas, extending it is beyond the scope of this work. We do, however, use a much simpler (albeit more inefficient) protocol to report a conservative estimate of the capacity of each covert channel.

Our protocol for data transfer uses the same mechanisms used for co-residence detection in section 4.3 to send and receive bits and perform clock synchronization. However, since we can use our co-residence detector to identify lambdas on a machine and target the two that we wish to label as the sender and receiver, we are not concerned about noise from multiple receivers, and as such can allow the receiver to sample continuously (section 4.3.3) and sample for extremely small duration (milliseconds instead of seconds). While we want the sampling duration to be as small as possible (in

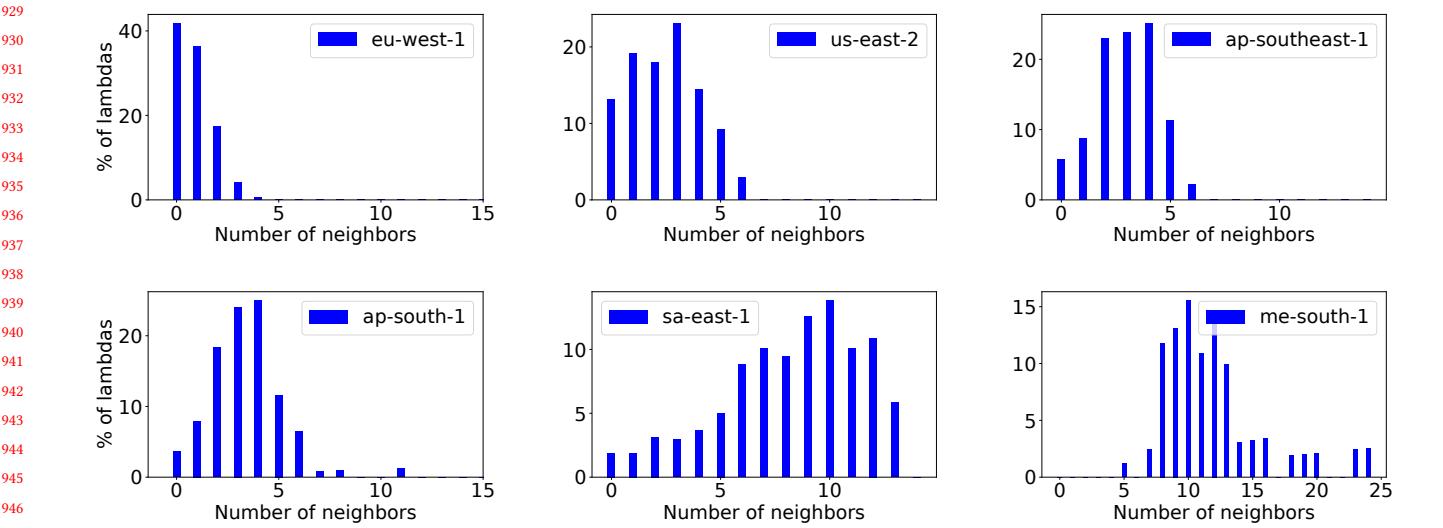


Figure 7: We present co-residence results from a 1000-lambda run in different AWS regions. Each bar shows the fraction of those 1000 lambdas (in %) that discovered a certain number of neighbors. The total amount and density of co-residence vary widely across regions, perhaps based on the size of those regions and the lambda activity within them.

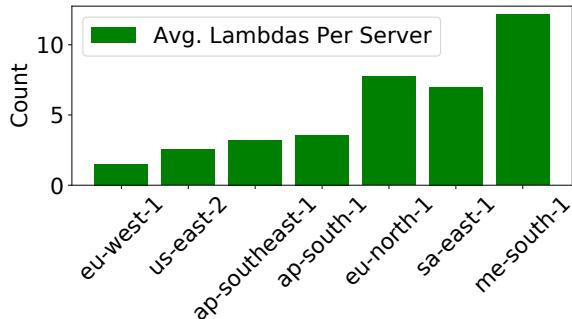


Figure 8: This figure shows the average number of lambdas per server i.e., the co-residence density seen in various AWS regions for the runs shown in Figure 7. The ample co-residence across regions demonstrates the practicality of establishing covert channels with lambdas in these regions.

order to increase the rate of bits transferred), the chances of erasures or errors also increases as the sender and receiver may get descheduled during this time.

To demonstrate this, we launched hundreds of 3 GB lambdas on AWS and use our co-residence detector to establish tens of covert channels. We then send data over these channels at various bitrates and record the error ratio (for byte-sized data segments). Figure 9 shows the mean error ratio at 50% and 95% one-sided confidence intervals, both of which increase with the bitrate.

To correct these errors, we use Reed-Solomon coding, block-based error correction codes that is suitable for burst-errors caused

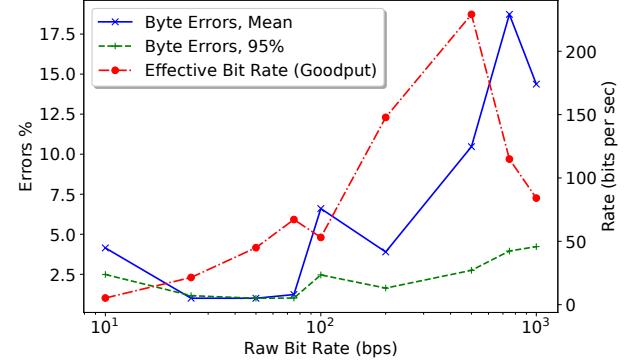


Figure 9: TODO ▶ increase font◀

by descheduling, with byte-sized symbols [26]. However, error correction comes with an overhead; Reed-Solomon requires twice as many extra symbols. Ariana ▶ why dont we just use the term bytes here instead of symbols? - that would be ◀ as there are errors to correct. So, for each bitrate, we must compute effective bitrate by subtracting the overhead of error correction symbols. From Figure 9, we can see that effective bitrate rises to a maximum of over 200 bits per second (bps) (at 500bps raw rate) before falling again due to high error rate. We confirmed this by sending Reed-Solomon encoded data over the covert channels at this rate and observed near-zero data corruption. Thus, we conclude that, by a conservative estimate, we can safely send data across each of these covert channels at a rate of 200 bps.

1045 5.2 Rendezvous point density

1046 **TODO** ► *Changes due to above subsection*◀ Next, we present measurements
 1047 on covert channel density on AWS using our co-residence
 1048 detector, and discuss the factors that may affect this density. As
 1049 we discussed in section 3, each co-resident group of lambdas rep-
 1050 resent an individual server in the cloud and hence can enable an
 1051 independent covert channel wherever the group has more than
 1052 two lambdas. So we attempt to answer the following question: as-
 1053 suming that the user launches a number of (sender and receiver)
 1054 lambdas at a specific point in time, what is the expected number
 1055 of such co-resident groups (with two or more lambdas) that they
 1056 might see? We deploy a large number of lambdas on various AWS
 1057 regions and report the co-residence density, that is, the average
 1058 number of such co-residence groups. The higher the co-residence
 1059 density, the easier it is for the user to ultimately establish covert
 1060 channels with lambdas, and the more information they can send.
 1061 Unless specified otherwise, all the experiments discussed in this
 1062 section are performed with 1.5 GB lambdas and executed success-
 1063 fully with **zero error** in co-residence detection.

1064
 1065 5.2.1 *Across AWS regions.* **TODO** ► *clarify the attacker model here, per*
 1066 *R1 comments*◀ We execute our co-residence detector with 1000 1.5
 1067 GB Lambdas in various AWS regions. Figure 7 comprises multiple
 1068 plots depicting the co-resident groups per region, with each bar in-
 1069 dicating the fraction of lambdas that detected a certain number of
 1070 neighbors (i.e., that belong to a co-resident group of a certain size).
 1071 Plots that skew to the right indicate a higher co-residence density
 1072 when compared to the plots skewed to the left (also illustrated in
 1073 Figure 8). We note that, in most regions, almost all lambdas rec-
 1074 ognize at least one neighbor (indicated by smaller or non-existent
 1075 first bar in each plot). We hypothesize that the co-residence den-
 1076 sity is (inversely) dependent on the total number of servers and
 1077 the lambda activity in the region, both of which can be assumed
 1078 to be lower in newer AWS regions, hence the higher co-residence
 1079 density in those regions as we can see in Figure 8. The ample co-
 1080 residence in general across all the regions shows that lambdas pro-
 1081 vide a fertile ground for covert channels.

1082
 1083 5.2.2 *Other factors.* We also examine how co-residence is affected
 1084 by various launch strategies that the user may use, like deploying
 1085 lambdas from multiple AWS accounts and different lambda sizes.
 1086 In particular, we wish to determine if our mechanism exhibits
 1087 different results when: 1) the user deploys sender lambdas and re-
 1088 ceiver lambdas on two separate accounts (normally the case with
 1089 covert channels) and 2) the senders and receivers are created with
 1090 different lambdas sizes. To answer these questions, we run an ex-
 1091 periment with 1000 lambdas of which we launch 500 lambdas from
 1092 one account (senders) and 500 from other deployed in a random or-
 1093 der. The co-residence observed was comparable to the case where
 1094 all the lambdas were launched from one account. In the left sub-
 1095 figure of Figure 10, we show the breakdown of co-resident group
 1096 of lambdas of each size among the two accounts. We can see that
 1097 among the co-resident groups of all sizes, roughly half of lamb-
 1098 das came from either account. This shows that lambda scheduler
 1099 is agnostic to the accounts the lambdas were launched from. We
 1100 see similar results for different lambda sizes, as shown in the right
 1101 subfigure of Figure 10.

1103 6 DISCUSSION

1104 **Alternate use cases** Our main motivation behind proposing a co-
 1105 residence detector for lambdas is demonstrating the feasibility of
 1106 a covert channel. However, there are other scenarios where such
 1107 tool can be (ab)used, of which we provide some examples.

- 1108 • While our co-residence detector does not directly help at-
 1109 tackers locate their victims in the cloud, it can aid them
 1110 in performing devastating DDOS attacks once by concen-
 1111 trating a number of attack instances on the victim machine.
 1112 Also, the attacker could try to gain a wider surface area for
 1113 targeted attacks in a cost-effective way by turning on/off
 1114 her co-resident instances as necessary.
- 1115 • Previous studies on performance aspects of lambdas (like
 1116 performance isolation) [24] generally need a way to find co-
 1117 resident lambdas. As software-level logical channels begin
 1118 to disappear, our tool might provide a reliable alternative.
- 1119 • Burst parallel frameworks [8] that orchestrate lambdas can
 1120 use our co-residence detector as a locality indicator to take
 1121 advantage of server locality.

1122 **Mitigation** In previous section, we showed that our co-
 1123 residence detector makes the covert channels practical with lamb-
 1124 das, so it is important that clouds address this issue. One way to
 1125 disable our co-residence detector is to fix the underlying memory
 1126 bus channel that it employs. However, this only works for newer
 1127 generation of servers and is not practical for existing infrastruc-
 1128 ture. An easier solution, one that is only practical with lambdas, is
 1129 to disable the lambda support for low-level languages (or unsafe
 1130 versions of high-level languages) by the cloud providers. This will
 1131 prevent pointer arithmetic that is required to activate this channel.
 1132 If that is not an option, cloud providers may look at more expen-
 1133 sive solutions like BusMonitor [20] that isolate memory bus usage
 1134 for different tenants by trapping the atomic operations to the hy-
 1135 pervisor. We leave such exploration to future work. **TODO** ► *small*
 1136 *note on limitations and future work?*◀

1137 7 RELATED WORK

1138 **Cloud Attacks** Co-residency is possible because of covert chan-
 1139 nels, so we begin our related work with an investigation into cloud
 1140 attacks. Initial papers in co-residency detection utilized host infor-
 1141 mation and network addresses arising due to imperfect virtualization
 1142 [19]. However, these channels are now obsolete, as cloud provi-
 1143 ders have strengthened virtualization and introduced Virtual Pri-
 1144 vate Clouds [4]. Later work used cache-based channels in various
 1145 levels of the cache [12, 15, 28, 35] and hardware based channels
 1146 like thermal covert channels [17], RNG module [7] and memory
 1147 bus [26] have also been explored in the recent past. Moreover, stud-
 1148 ies have found that VM performance can be significantly degraded
 1149 using memory DDoS attacks [31], while containers are susceptible
 1150 to power attacks from adjacent containers [9].

1151 Our work focuses on using the memory bus as a covert chan-
 1152 nel for determining cooperative co-residency. Covert channels
 1153 using memory bus were first introduced by Wu et. al [26], and
 1154 subsequently has been used for co-residency detection on VMs
 1155 and containers [21]. Wu et. al [26] introduced a new technique
 1156 to lock the memory bus by using atomic memory operations on
 1157 addresses that fall on multiple cache lines, a technique we rely on
 1158

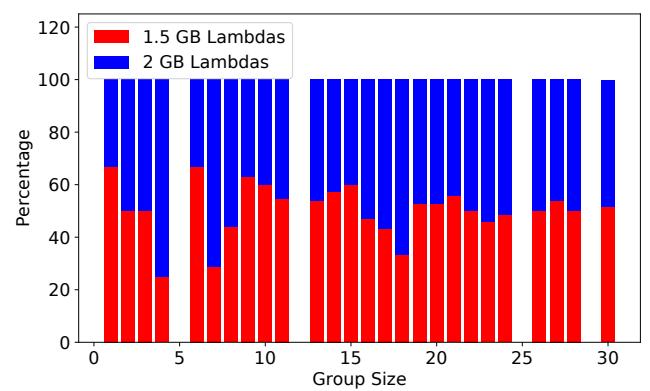
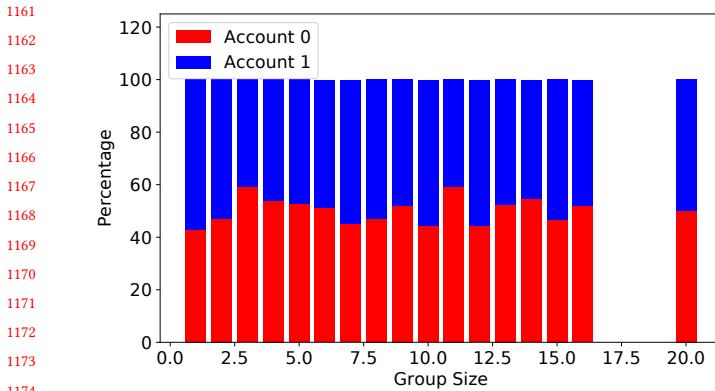


Figure 10: The left plot shows the breakdown of co-resident groups (of varying sizes) of lambdas by two different accounts in an experiment of 1000 lambdas, where 500 lambdas are launched from each account. The uniformity of the split indicates that the lambda scheduler might be invariant to the account the lambdas are launched from. Similar results are shown for different lambda sizes in the right plot.

in our own work.

Co-residency One of the first pieces of literature in detecting VM co-residency was introduced by Ristenpart et al., who demonstrated that VM co-residency detection was possible and that these techniques could be used to gather information about the victim machine (such as keystrokes and network usage) [19]. This initial work was further expanded in subsequent years to examine co-residency using memory bus locking [30] and active traffic analysis [6], as well as determining placement vulnerabilities in multi-tenant Public Platform-as-a-Service systems [22, 33]. Finally, Zhang et al. demonstrated a technique to detect VM co-residency detection via side-channel analyses [34]. Our work expands on these previous works by investigating co-residency for lambdas.

Lambdas While lambdas are a much newer technology than VMs, there still exists literature on the subject. Recent studies examined cost comparisons of running web applications in the cloud on lambdas versus other architectures [23], and also examined the lambdas have been studied in the context of cost-effectiveness of batching and data processing with lambdas [13]. Further research has shown how lambdas perform with scalability and hardware isolation, indicating some flaws in the lambda architecture [24]. From a security perspective, Izhikevich et. al examined lambda co-residency using RNG and memory bus techniques (similar to techniques utilized in VM co-residency) [11]. However, our work differs from this study in that our technique informs the user of which lambdas are on the same machine, not only that the lambdas experience co-residency.

8 ETHICAL CONSIDERATIONS

As with any large scale measurement project, we discuss the ethical considerations. First, there are security and privacy concerns of using this technique to uncover other consumer’s lambdas. However, since we focus on co-operative co-residence detection, we only determine co-residence for the lambdas we launched, and do

not gain insight into other consumer’s lambdas. Second, there is concern that our experiments may cause performance issues with other lambdas, as we may block their access to the memory bus. We believe this concern is small, for a number of reasons. Memory accesses are infrequent due to the multiple levels of caches; we would only be affecting a small number of operations. Memory accesses and locking operations are FIFO, which prevents starvation of any one of the lambdas sharing a machine. Moreover, lambdas are generally not recommended for latency-sensitive workloads, due to their cold-start latencies. Thus, the small amount of lambdas that we might affect should not, in practice, be affected in their longterm computational goals.

9 CONCLUSION

In this paper, we have demonstrated techniques to build robust, scalable and efficient covert channels entirely using serverless cloud functions such as AWS lambdas. To achieve this goal, we developed a fast and reliable co-residence detector for lambdas, and evaluated it for correctness and scalability. Finally, we have empirically demonstrated the practicality of such covert channels by studying the co-residence density of lambdas on various AWS regions. TODO ► add acknowledgements ◀

REFERENCES

- [1] Alexandru Agache, Marc Brooker, Alexandra Iordache, Anthony Ligouri, Rolf Neugebauer, Phil Piwonka, and Diana-Maria Popa. 2020. Firecracker: Lightweight Virtualization for Serverless Applications. In *17th USENIX Symposium on Networked Systems Design and Implementation (NSDI 20)*. USENIX Association, Santa Clara, CA, 419–434. <https://www.usenix.org/conference/nsdi20/presentation/agache>
- [2] aws 2018. Understanding Container Reuse in AWS Lambda. <https://aws.amazon.com/blogs/compute/container-reuse-in-lambda/>.
- [3] aws 2019. Amazon EC2. <https://aws.amazon.com/ec2/>.
- [4] aws 2019. Amazon VPC. <https://aws.amazon.com/vpc/>.
- [5] azure 2019. Azure Functions. <https://azure.microsoft.com/en-us/services/functions/>.
- [6] Adam Bates, Benjamin Mood, Joe Pletcher, Hannah Pruse, Masoud Valafar, and Kevin Butler. 2012. Detecting co-residency with active traffic analysis techniques. <https://doi.org/10.1145/2381913.2381915>

- 1277 [7] Dmitry Evtyushkin and Dmitry Ponomarev. 2016. Covert Channels Through
1278 Random Number Generator: Mechanisms, Capacity Estimation and Mitigations.
1279 In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security* (Vienna, Austria) (CCS '16). ACM, New York, NY, USA, 843–857.
1280 <https://doi.org/10.1145/2976749.2978374>
- 1281 [8] Sajjad Fouladi, Francisco Romero, Dan Iter, Qian Li, Shuvo Chatterjee, Christos
1282 Kozyrakis, Matei Zaharia, and Keith Winstein. 2019. From Laptop to Lambda:
1283 Outsourcing Everyday Jobs to Thousands of Transient Functional Containers.
1284 In *2019 USENIX Annual Technical Conference (USENIX ATC '19)*. USENIX Association, Renton, WA, 475–488. <https://www.usenix.org/conference/atc19/presentation/fouladi>
- 1285 [9] Xing Gao, Zhongshu Gu, Mehmet Kayaalp, Dimitrios Pendarakis, and Haining
1286 Wang. 2017. ContainerLeaks: Emerging Security Threats of Information Leaks
1287 in Container Clouds. In *Dependable Systems and Networks (DSN), 2017 47th Annual IEEE/IFIP International Conference on*. IEEE, IEEE, 237–248.
- 1288 [10] gcp 2019. Google Cloud Functions. <https://cloud.google.com/functions/>.
- 1289 [11] Elizabeth Izhikevich. 2018. *Building and Breaking Burst-Parallel Systems*. Master's thesis. University of California, San Diego.
- 1290 [12] Mehmet Kayaalp, Nael Abu-Ghazaleh, Dmitry Ponomarev, and Aamer Jaleel.
1291 2016. A high-resolution side-channel attack on last-level cache. 1–6. <https://doi.org/10.1145/2897937.2897962>
- 1292 [13] Mariam Kiran, Peter Murphy, Inder Monga, Jon Dugan, and Sartaj Singh Baveja.
1293 2015. Lambda architecture for cost-effective batch and speed big data processing.
1294 In *2015 IEEE International Conference on Big Data (Big Data) (2015)*, 2785–2792.
- 1295 [14] Lambda 2019. AWS Lambda. <https://aws.amazon.com/lambda/>.
- 1296 [15] Fangfei Liu, Yuval Yarom, Qian ge, Gernot Heiser, and Ruby Lee. 2015. Last-
1297 Level Cache Side-Channel Attacks are Practical, Vol. 2015. 605–622. <https://doi.org/10.1109/SP.2015.43>
- 1298 [16] Ramya Jayaram Masti, Devendra Rai, Aanjan Ranganathan, Christian Müller,
1299 Lothar Thiele, and Srdjan Capkun. 2015. Thermal Covert Channels on Multi-core
1300 Platforms. In *24th USENIX Security Symposium (USENIX Security 15)*. USENIX Association, Washington, D.C., 865–880. <https://www.usenix.org/conference/usenixsecurity15/technical-sessions/presentation/masti>
- 1301 [17] Ramya Jayaram Masti, Devendra Rai, Aanjan Ranganathan, Christian Müller,
1302 Lothar Thiele, and Srdjan Capkun. 2015. Thermal Covert Channels on Multi-core
1303 Platforms. In *Proceedings of the 24th USENIX Conference on Security Symposium* (Washington, D.C.) (SEC'15). USENIX Association, Berkeley, CA, USA, 865–880.
1304 <http://dl.acm.org/citation.cfm?id=2831143.2831198>
- 1305 [18] Clementine Maurice, Manuel Weber, Michael Schwarz, Lukas Giner, Daniel
1306 Gruss, Carlo Alberto Boano, Stefan Mangard, and Kay Römer. 2017. Hello
1307 from the Other Side: SSH over Robust Cache Covert Channels in the Cloud.
1308 <https://doi.org/10.14722/ndss.2017.23294>
- 1309 [19] Thomas Ristenpart, Eran Tromer, Hovav Shacham, and Stefan Savage. 2009. Hey,
1310 You, Get off of My Cloud: Exploring Information Leakage in Third-party Compute
1311 Clouds. In *Proceedings of the 16th ACM Conference on Computer and Communications Security* (Chicago, Illinois, USA) (CCS '09). ACM, New York, NY, USA,
1312 199–212. <https://doi.org/10.1145/1653662.1653687>
- 1313 [20] Brendan Saltaformaggio, D. Xu, and X. Zhang. 2013. BusMonitor : A Hypervisor-
1314 Based Solution for Memory Bus Covert Channels.
- 1315 [21] Venkatanathan Varadarajan, Yinqian Zhang, Thomas Ristenpart, and Michael
1316 Swift. 2015. A Placement Vulnerability Study in Multi-Tenant Public Clouds.
1317 In *24th USENIX Security Symposium (USENIX Security 15)*. USENIX Association, Washington, D.C., 913–928. <https://www.usenix.org/conference/usenixsecurity15/technical-sessions/presentation/varadarajan>
- 1318 [22] Venkatanathan Varadarajan, Yinqian Zhang, Thomas Ristenpart, and Michael M.
1319 Swift. 2015. A Placement Vulnerability Study in Multi-tenant Public Clouds.
1320 *CoRR* abs/1507.03114. arXiv:1507.03114 <http://arxiv.org/abs/1507.03114>
- 1321 [23] Mario Villamizar, Oscar Garcés, Lina Ochoa, Harold Castro, Lorena Salamanca,
1322 Mauricio Verano Merino, Rubby Casallas, Santiago Gil, Carlos Valencia, Angee
1323 Zambrano, and Mery Lang. 2016. Infrastructure Cost Comparison of Running
1324 Web Applications in the Cloud Using AWS Lambda and Monolithic and Microservice
1325 Architectures. <https://doi.org/10.1109/CCGrid.2016.37>
- 1326 [24] Liang Wang, Mengyuan Li, Yinqian Zhang, Thomas Ristenpart, and Michael
1327 Swift. 2018. Peeking Behind the Curtains of Serverless Platforms. In *2018 USENIX
1328 Annual Technical Conference (USENIX ATC 18)*. USENIX Association, Boston, MA,
1329 133–146. <https://www.usenix.org/conference/atc18/presentation/wang-liang>
- 1330 [25] Zhenghong Wang and Ruby B. Lee. 2006. Covert and Side Channels Due to
1331 Processor Architecture. In *Proceedings of the 22nd Annual Computer Security
1332 Applications Conference (ACSAC '06)*. IEEE Computer Society, USA, 473–482.
1333 <https://doi.org/10.1109/ACSAC.2006.20>
- 1334 [26] Zhenyu Wu, Zhang Xu, and Haining Wang. 2012. Whispers in the Hyper-space:
1335 High-speed Covert Channel Attacks in the Cloud. In *Presented as part of the 21st
1336 USENIX Security Symposium (USENIX Security 12)*. USENIX, Bellevue, WA, 159–
1337 173. <https://www.usenix.org/conference/usenixsecurity12/technical-sessions/presentation/wu>
- 1338 [27] Yunjing Xu, Michael Bailey, Farnam Jahanian, Kaustubh Joshi, Matti Hiltunen,
1339 and Richard Schlichting. 2011. An Exploration of L2 Cache Covert Channels
1340 in Virtualized Environments. In *Proceedings of the 3rd ACM Workshop on Cloud
Computing Security Workshop* (Chicago, Illinois, USA) (CCSW '11). Association
1341 for Computing Machinery, New York, NY, USA, 29–40. <https://doi.org/10.1145/2046660.2046670>
- 1342 [28] Yunjing Xu, Michael Bailey, Farnam Jahanian, Kaustubh Joshi, Matti Hiltunen,
1343 and Richard Schlichting. 2011. An Exploration of L2 Cache Covert Channels
1344 in Virtualized Environments. In *Proceedings of the 3rd ACM Workshop on Cloud
Computing Security Workshop* (Chicago, Illinois, USA) (CCSW '11). ACM, New
1345 York, NY, USA, 29–40. <https://doi.org/10.1145/2046660.2046670>
- 1346 [29] Zhang Xu, Haining Wang, and Zhenyu Wu. 2015. A Measurement Study on Co-
1347 residence Threat inside the Cloud. In *24th USENIX Security Symposium (USENIX
1348 Security 15)*. USENIX Association, Washington, D.C., 929–944. <https://www.usenix.org/conference/usenixsecurity15/technical-sessions/presentation/xu>
- 1349 [30] Zhang Xu, Haining Wang, and Zhenyu Wu. 2015. A Measurement Study on Co-
1350 residence Threat inside the Cloud. In *24th USENIX Security Symposium (USENIX
1351 Security 15)*. USENIX Association, Washington, D.C., 929–944. <https://www.usenix.org/conference/usenixsecurity15/technical-sessions/presentation/xu>
- 1352 [31] Tianwei Zhang, Yinqian Zhang, and Ruby Lee. 2016. Memory DoS Attacks in
1353 Multi-tenant Clouds: Severity and Mitigation.
- 1354 [32] Weijuan Zhang, Xiaoqi Jia, Chang Wang, Shengzhi Zhang, Qingjia Huang, Ming-
1355 sheng Wang, and Peng Liu. 2016. A Comprehensive Study of Co-residence
1356 Threat in Multi-tenant Public PaaS Clouds. In *Information and Communications
1357 Security*, Kwok-Yan Lam, Chi-Hung Chi, and Sihan Qing (Eds.). Springer International
1358 Publishing, Cham, 361–375.
- 1359 [33] Weijuan Zhang, Xiaoqi Jia, Chang Wang, Shengzhi Zhang, Qingjia Huang, Ming-
1360 sheng Wang, and Peng Liu. 2016. A Comprehensive Study of Co-residence
1361 Threat in Multi-tenant Public PaaS Clouds, Vol. 9977. 361–375. https://doi.org/10.1007/978-3-319-50011-9_28
- 1362 [34] Yinqian Zhang, Ari Juels, Alina Oprea, and Michael Reiter. 2011. HomeAlone: Co-
1363 Residency Detection in the Cloud via Side-Channel Analysis. 313 – 328. <https://doi.org/10.1109/SP.2011.31>
- 1364 [35] Yinqian Zhang, Ari Juels, Michael K. Reiter, and Thomas Ristenpart. 2014. Cross-
1365 Tenant Side-Channel Attacks in PaaS Clouds. In *Proceedings of the 2014 ACM
1366 SIGSAC Conference on Computer and Communications Security* (Scottsdale, Arizona,
1367 USA) (CCS '14). ACM, New York, NY, USA, 990–1003. <https://doi.org/10.1145/2660267.2660356>