

Columbus: Fast, Reliable Co-residence Detection for Lambdas

ABSTRACT

Ariana ▶ Cloud computing has seen explosive growth in the past decade. While efficient sharing of infrastructure among tenants has contributed to this growth, the same principles also open avenues for covert-channels, or the capability to share information between processes that should be hypothetically isolated. Providers like AWS and Azure have traditionally relied on obfuscation to prevent this leakage of information, but recent works have repeatedly found detection techniques that break this encapsulation, prompting the providers to harden isolation on their platforms. In this work, we find yet another such covert-channel for lambdas based on the memory bus, which is more pervasive, reliable, and harder to fix. We show that this technique can be used to reliably perform co-operative co-residence detection for thousands of AWS lambdas within a few seconds, which could aid attackers in performing DDoS attacks or learn cloud's internal mechanisms. In this paper, we present the technique in detail, evaluate it, and use it to perform a measurement study on lambda activity across AWS regions. Through this work, we hope to motivate the need to address this covert channel in the cloud.◀

Cloud computing has seen explosive growth in the past decade. This is made possible by efficient sharing of infrastructure among tenants, which unfortunately also opens avenues for security attacks (e.g., side-channels). Providers, like AWS and Azure, have traditionally relied on hiding the co-residency information to prevent targeted attacks in their clouds. But recent works have repeatedly found co-residence detection techniques that break this encapsulation, prompting the providers to address them and harden isolation on their platforms. In this work, we find yet another such technique based on a memory bus covert channel, which is more pervasive, reliable and harder to fix. We show that we can use this technique to reliably perform co-operative co-residence detection for thousands of AWS lambdas within a few seconds, which could aid attackers in performing DDoS attacks or learn cloud's internal mechanisms. In this paper, we present the technique in detail, evaluate it and use it to perform a small study on lambda activity across AWS regions. Through this work, we hope to motivate the need to address this covert channel in the cloud.

CCS CONCEPTS

- Security and privacy → Side-channel analysis and countermeasures; Virtualization and security.

KEYWORDS

cloud, cartography, serverless, coresidency, covert channels

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.

The Web Conference 2021, April 19–23, 2021, Ljubljana, Slovenia

© 2020 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

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

ACM Reference Format:

. 2020. Columbus: Fast, Reliable Co-residence Detection for Lambdas. In *The Web Conference 2021, April 19–23, 2021, Ljubljana, Slovenia*. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Cloud computing is a fast growing technology that is being widely used around the globe ▶ *cite; can also throw in numbers to make this more convincing*◀. Major cloud providers like AWS [4], Microsoft Azure [5] and Google Cloud [10] provide a multitude of compute and storage services and have seen immense growth over the past decade ▶ *cite*◀. While there are many benefits of using cloud services, like increased scalability and decreased IT costs [2], cloud also brings new security concerns to the table as multiple tenants share the same underlying physical infrastructure. For example, by multiplexing applications from multiple tenants on the same server, cloud computing opens up new avenues and opportunities for side-channel attacks through shared server hardware like caches [15, 24]. While the virtualization platforms used to share the infrastructure have been constantly improving to harden the isolation between the tenants, it's a never ending process and new attack surfaces keep showing up ▶ *cite, HELP!*◀. Also, hardening isolation usually comes at the cost of performance which may offset the benefits of sharing the infrastructure that enabled cloud computing in the first place.

Traditionally, clouds have relied on hiding their placement mechanisms (i.e., how they pack tenants onto their servers) and coresidence information (i.e., which particular tenants are packed together in a server) as a first line of defence against targeted attacks. By reducing attacker's ability to get on the same server as the victim, attackers may have to go with trial-and-error or brute-force solutions that can be very expensive with low yield. However, the encapsulation mechanisms are not perfect. A plethora of previous work, kick-started by Ristenpart et al. [18] a decade ago, have exploited various covert channels to enable co-residency detection on clouds like AWS, and used them to demonstrate targeted attacks or shed light on cloud's internal placement and resource allocation mechanisms ▶ *cite*◀ that further weaken the encapsulation. Clouds have since promptly fixed many of these covert channels and came up with containers that provide better isolation, like AWS Firecracker[1] for example. In this work, we set out to prove that the job is not done yet.

We found yet another way to detect coresidency among cloud instances (i.e., whether they are running on the same physical server), one which is more pervasive, reliable and harder to fix than with previous approaches. We use a covert channel based on memory bus hardware, first introduced by Wu et al.[23], which we show to be omni-present in all clouds. Unlike software-based covert channels that can be fixed with new releases, this one is inherent to x86 hardware and is harder to fix. Finally, while previous approaches have made simplistic use of this covert channel[20] due to noise and synchronization issues, we overcome these to communicate bits

117 reliably over the channel that enabled us to achieve lightning-fast
 118 colocation detection for thousands of cloud instances.

119 In this work, we chose serverless functions as our cloud
 120 containers of choice. Serverless functions have seen increased interest
 121 in recent years with most clouds providing these services, such
 122 as lambdas on AWS [14] and cloud functions on GCP [9] (We use
 123 the terms lambdas, serverless functions and cloud instances inter-
 124 changeably hereafter). They provide a more challenging environ-
 125 ment for coresidence detection techniques based on covert channels
 126 as they usually have restricted runtimes that limit low-level code
 127 sometimes required to access such covert channels, and are more
 128 ephemeral (forcing the coresidence detection to be faster). They
 129 are also significantly cheaper, providing us a cost-effective way
 130 to demonstrate the technique. Note that whatever we are able to
 131 achieve with lambdas can be easily replicated with other traditional
 132 containers, as they are less, not more, restricted environments.

133 In this paper, we only focus on *cooperative* coresidence detec-
 134 tion, to start with. That is, all the lambdas are assumed to be in
 135 our/attacker's control. While this doesn't by itself help attacker
 136 target a victim, it can be used to perform targeted DDoS attacks
 137 **TODO** ► cite power attacks?◀ or learn cloud's internal mechanisms
 138 that may aid this. We propose a technique which, for given a set of
 139 lambdas deployed onto the cloud, can figure out which of them ran
 140 on the same server in the cloud. We show that we can do this with
 141 100% success rate for bigger lambdas and furthermore, we can do
 142 this within a minute for thousands of lambdas. We then perform a
 143 minor study on colocation patterns in various AWS regions with
 144 some insights on lambda activity. **TODO** ► any other takeaways?◀
 145 Through this, we hope to motivate the need to address the memory
 146 bus covert channel in all the three clouds.

147 The remainder of this paper is organized as follows. Section 2
 148 presents some background from the literature. Sections 3 and 4
 149 present our co-residence detection mechanism in detail. We eval-
 150 uate the mechanism in section 5 and conclude with a placement
 151 study of AWS lambdas using our technique in section 6.

153 2 BACKGROUND

154 We begin with a brief background on relevant topics.

156 2.1 Lambdas/Serverless Functions

158 **TODO** ► needs sprucing◀ One of the fast-growing cloud services in
 159 recent years are serverless functions such as lambdas on AWS [14]
 160 and cloud functions on GCP [9]. This interest stems from the fact
 161 that serverless architecture does not require the developer to worry
 162 about provisioning, maintaining, and administering servers. Addi-
 163 tionally, lambdas are much more cost-efficient than VMs as they
 164 allow more efficient packing of the servers. However, these func-
 165 tions are more ephemeral than containers and VMs, in many cases
 166 only few minutes. While this attribute provides more flexibility
 167 in cost and functionality, the nature of serverless functions also
 168 increases the difficulty in detecting co-residency and launching
 169 successful attacks.

170 We focus on AWS Lambdas in this paper but we show that our
 171 study is applicable to other clouds as well. Lambdas execute as
 172 much smaller units than containers and virtual machines. Lambda
 173 sizes range from 128 MB to 3 GB, and their maximum timeout

175 value is 15 minutes. While lambdas are limited in the computations
 176 they can execute, they are conversely incredibly lightweight and
 177 can be initiated and deleted in a very short amount of time. Since
 178 lambdas are short-lived and lightweight, the user has no control
 179 over the physical location of the server(s) on which their lambdas
 180 are spawned.

182 2.2 Co-residence Detection

183 To perform side-channel attacks against other tenants in a cloud
 184 setting, attackers need to co-locate their applications on the same
 185 servers as their victims. Past research has used various strategies
 186 to achieve co-residency for demonstrating such attacks. Typically,
 187 achieving co-residency includes a (VM/Container) launch strategy
 188 (varying number of instances, time of the day, etc) combined with a
 189 co-residence detection mechanism for detecting if two instances are
 190 running on the same machine. Traditionally, this was done based on
 191 software runtime information like public/internal IP addresses[18],
 192 files in *procfs* or other environment variables[22] and other such
 193 logical side-channels[19, 26] that two instances running on a same
 194 server might share.

195 As virtualization platforms move towards stronger isolation be-
 196 tween instances (e.g. AWS' Firecracker VM [1]), these logical side-
 197 channels have become less effective or infeasible. Furthermore,
 198 some of these side-channels were only effective on container-based
 199 platforms that share the underlying OS image and were less suitable
 200 for hypervisor-based platforms. This prompted a move towards
 201 hardware-based covert channels which can bypass software isolat-
 202 ion and are usually harder to fix. Typically, these covert channels
 203 involve sending/receiving information by causing contention on
 204 a shared hardware that results in observable performance fluctua-
 205 tions across applications. A number of such side-/covert channels
 206 based on shared hardware like last-level caches **TODO** ► references◀,
 207 memory bus [20, 23, 29] and storage devices **TODO** ► references◀
 208 have been explored in the past, some of which have already been
 209 addressed and are no longer even feasible in most clouds (e.g., last-
 210 level cache-based channels[25]).

212 2.3 Memory Bus Covert Channel

214 One shared hardware that we examine in our study is the memory
 215 bus. The memory bus is the piece of hardware that connects the
 216 memory controller to main memory. Memory bus contention can
 217 be caused by initiating repeated memory accesses to saturate the
 218 bus bandwidth and cause observable latency spikes. However, this
 219 turns out to be very challenging given the multiple levels of caches
 220 on today's servers, which prevent repeated memory accesses, and
 221 the high memory bandwidth that cannot be saturated by few CPU
 222 cores.

223 In x86 systems, atomic memory instructions designed to facil-
 224 itate multi-processor synchronization are supported by cache co-
 225 herence protocols as long as the operands stay within a cache line
 226 (which is generally the case as language compilers make sure that
 227 operands are aligned). However, if the operand is spread across
 228 two cache lines (referred to as "exotic" memory operations), x86
 229 hardware achieves atomicity by locking the memory bus to prevent
 230 any other memory access operations until the current operation
 231 finishes. This results in significantly higher latencies for the other

operations which cannot use the ample memory bandwidth due to the lock[23]. Furthermore, this behaviour persists even in the presence of multiple processor sockets, making the locking effects visible to all the cores on the machine. We exploit this property of x86 hardware to cause contention on the memory bus and use the resulting observable variations in performance as a covert channel for detecting co-residency.

3 METHODOLOGY

Ariana ► make sure that all these terms are defined ahead of time: co-location, co-operative, co-residence, serverless, lambdas◀

Our goal is to determine a co-operative co-residence detection mechanism for serverless functions. In other words, given a series of spawned lambdas in a given region on a cloud service, how can we determine the lambdas that are co-located on the same machines? In this section, we discuss the details of such a mechanism, previous solutions to this problem, and the unique challenges we faced with lambda co-residence.

Given a set of cloud instances (VMs, Containers, Functions, etc) deployed to a public cloud, a co-residence detection mechanism would identify, for each pair of instances in the set, whether the pair was running on the same physical server at some point. Paraphrasing Varadarajan et al.[20], for a co-detection mechanism to be useful across a wide range of launch strategies, we observe that it should have the following desirable 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 clouds 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 be fast and scalable such that the technique will take less time to detect all co-resided pairs at a reasonable cost.

Given these properties, we decide to investigate hardware-based covert channels. Hardware-based covert-channels are more difficult to remove and obfuscate than software-based covert channels, and are also more ubiquitous, given that hardware is more homogenous in nature than software.

3.0.1 RNG Hardware. **Ariana** ► This feels a bit awkward. We might want to move this to discussion◀ We first examined covert channels based on Random Number Generator (RNG) hardware[7]. Modern processors support this shared hardware module to generate true random numbers. To produce true random numbers, information from this module is routed from the host machine to the /dev/random file in the guest virtual machine for cryptographic operations. Since the hardware is shared, if one guest consumes these random bits within an infinite **Ariana** ► does it need to be infinite?◀ loop, another user could notice a spike in random operations, indicating contention on the machine. **Ariana** ► do we have an old graph to back up this point of being too noisy? to just drive the point home that RNG is not great for these purposes◀ However, our experiments on AWS using RNG hardware indicated that this channel is unreliable for lambda co-detection.

We hypothesize that perhaps, because causing contention is easy, the channel gets too noisy as a result to accurately use for our purposes.

3.0.2 Memory bus channel. We next explored (and ultimately used) the memory bus covert channel described in section 2.3 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 [19?], and we found that they are still vulnerable today. **Ariana** ► I think there should be another sentence or two describing how we got to this figure right here. It sort of feels like we just jump into results◀ Moreover, we were able to demonstrate this behavior through serverless function instances, whose runtimes are mostly restricted to high-level languages that prevent the pointer arithmetic required to perform these exotic operations. To demonstrate this attack, we used the unsafe environments (C++ on AWS, Unsafe Go on GCP, Unsafe C# On Azure) that these clouds allowed. Figure 1 shows that all three major cloud providers still exhibit significant difference in latencies for the "exotic" memory locking operations when compared to regular memory access latencies. This figure shows the applicability of using memory bus across different kinds of cloud instances as well.

3.0.3 Previous approaches using Memory bus. Previous works that used the memory bus for co-residence detection divide the deployed instances into attack and (co-operative) victim roles, and attempt to co-locate the attacker instances with a victim instance. The attack roles continually lock the memory bus (locking process) for a certain duration (~10 seconds) while the victims sample the memory for any spike in access latencies (probing process). If all the deployed instances try the detection i.e., locking and probing at once, (some of) the victims may see locking effects, but there would be no way of knowing which or how many attack roles co-resided with a particular victim and caused the locking. This provides no information about the number of physical servers that ran these instances or the amount of co-location. The only information we can deduce is that victims were probably co-located with just a single attacker.

An alternative method is to try pair-wise detection where only one attack instance locks and one victim instance probes at a time revealing co-residence of this pair, and repeating this serially for each pair. However, this technique is too slow and scales quadratically with the number of instances e.g., a hundred instances take more than 10 hours assuming 10 secs for each pair. **TODO** ► get a cost estimate on ec2◀. Varadarajan et al.[19] speeds this process significantly by performing detection for mutually-exclusive subsets in parallel, allowing for false-positives and later eliminating the false-positives sequentially. **Ariana** ► might want to elaborate on this; on its own may be a bit confusing◀ This would still only scale linearly in the best case **TODO** ► (or not even that?)◀, which is still expensive - with a thousand instances, for example, the whole detection process takes well over 2 hours to finish, which is infeasible for lambdas that are, by nature, ephemeral. Thus, one challenge in this work is creating a faster neighbor detection algorithm.

3.0.4 The Path to Scalability. One method to quicken the co-location process is by cutting down on the time it takes for single attack-victim pair to determine co-residence i.e., improving upon probing

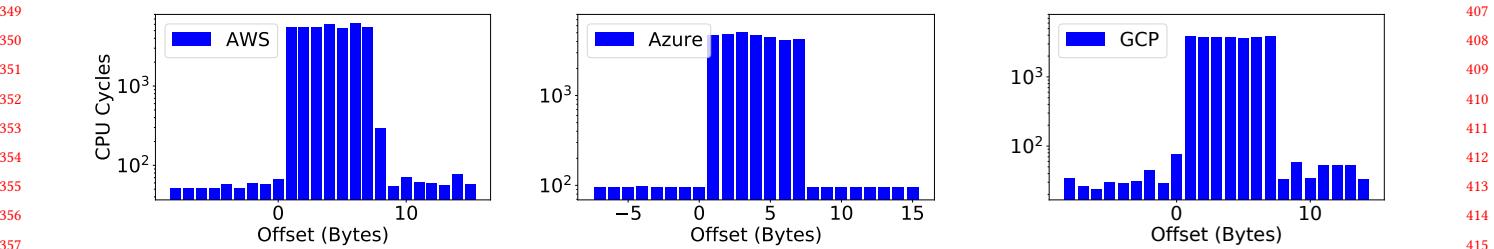


Figure 1: From left to right, the plots show the latencies of atomic memory operations performed on an 8B memory region as we slide it from one cache line across the boundary into another, on AWS, GCP and Azure clouds respectively. TODO ► Get plot for GCP◀. The latencies are orders of magnitude higher 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 clouds.

time and accuracy of the victim. However, this only affects total time by a constant factor. To improve scalability, we need to be able to run detection for two Ariana ► why two?◀ attack-victim pairs in parallel without sacrificing the certainty when they are run serially. For example, when both pairs see co-residence, we must be certain that each victim experienced co-residence because of its own attacker pair, which is not possible if co-residence is ascertained based a simple yes/no signal from the attack instance.

The memory bus convert channel was used to exchange more complex information like keys in previous work[23], and at first sight, can be used to exchange information such as IDs between the co-resided instances to solve our problem. However, the original work assumes that there is only one sender and one receiver who know exactly how and when to communicate. As we will see in the next section, this model is not sustainable when there exist many parties that have no knowledge of each other but try to communicate on the same channel.

To solve some of the challenges mentioned previously, we propose a protocol in which we use the memory busy covert channel to exchange IDs between instances, and the co-resided instances reliably exchange their IDs with each other to discover their neighbors. Ariana ► this sentence is a bit confusing◀ The protocol takes time on the order of number of instances involved, which is limited by the maximum number of co-located instances on a single server (tens) - something that is orders of magnitude less than total number of instances deployed to the cloud (hundreds to thousands). This lets us scale our co-residence detection significantly to run in less than a minute. Ariana ► seconds? minutes?◀.

4 NEIGHBOR DISCOVERY PROTOCOL

As noted earlier, co-residence detection scales well when co-residing instances on each server communicate among themselves and discover each other. Assuming that the instances have unique IDs, this requires the co-resided instances to exchange these (integer) IDs with each other using the memory bus channel as a transmission medium. In this section, we present a communication protocol that the co-resided instances can use to achieve communication in a fast and reliable way. We first discuss the challenges we faced in making the channel reliable before examining the protocol itself.

Algorithm 1 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

```

4.1 Reliable Transmission

Senders and receivers can reliably communicate 0-bits and 1-bits by causing contention on the memory busy. Consider the simple scenario where there is one sender and one receiver instance on a machine, and the sender has a set of bits that it needs to communicate with the receiver on the memory bus covert channel. To communicate a 1-bit, the sender instance causes contention on the memory busy by locking it using the special memory locking operations (discussed in section 2.3). The receiver would then 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). Pseudo-code for the sender instance is shown in Algorithm 1.

4.1.1 Sensing contention. There are two ways in which the receivers could detect memory bus 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 [20]) and observe a spike in their latencies to detect contention. 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.

Thought two methods exist for detecting memory bus contention, we decide to use the memory locking receiver for our experiments. Previous studies[20, 23] have established that both memory probing and memory locking receivers experience significant latency

Algorithm 2 Reading a bit in the receiver

```

465   1: now  $\leftarrow$  time.now()
466   2: end  $\leftarrow$  now + sampling_duration
467   3: sampling_rate  $\leftarrow$  num_samples/sampling_duration
468   4: address  $\leftarrow$  cache_line_boundary - 2
469   5: samples  $\leftarrow$  {}
470   6: while now < end do
471     7:   before  $\leftarrow$  RDTSC()
472     8:   _ATOMIC_FETCH_ADD(address)
473     9:   after  $\leftarrow$  RDTSC()
474    10:  samples  $\leftarrow$  samples  $\cup$  {(after - before)}
475    11:  wait until NEXT_POISSON(sampling_rate)
476    12:  now  $\leftarrow$  time.now()
477    13: end while
478    14: ks_val  $\leftarrow$  KOLMOGOROV_SMIRINOV(samples, baseline)
479    15: return ks_val < ksvalue_threshold
480

```

overhead during memory bus contention, making them both avenues for sensing the covert-channel. Memory probing involves regular (un-cached) memory accesses, which is universal, unlike the locking operations which are rarely used, if at all, by standard applications. This makes memory probing the only viable option for **non-cooperative** co-residence detection, where victims are not under attacker's control and cannot be assumed to perform locking operations. Furthermore, memory probing can be done on multiple receivers constantly 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 TODO ► cite some papers ◀. 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. Thus, we decide to use the memory locking receiver.

4.1.2 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 [20], 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-resided receivers and sampling takes a few microseconds each time.

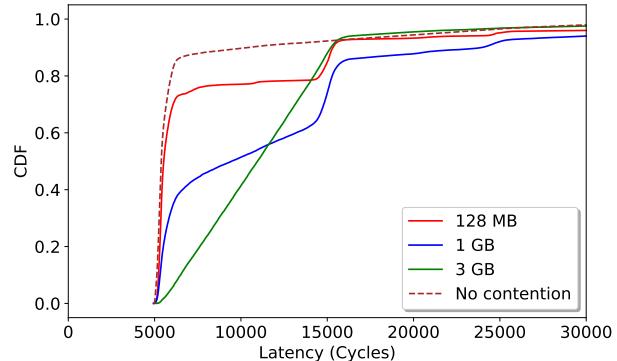


Figure 2: Shows CDF of latencies observed by 128 MB, 1 GB and 3 GB Lambdas during contention. The 128 MB lambda pair sees less contention due to more context switching, whereas the 1 GB and 3 GB lambdas see progressively more contention compared to baseline which we attribute to their relative stability on the underlying physical cores.

4.1.3 Sampling duration. 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, in practice, the time-sharing of processors produces difficulties. The sender is not continually causing contention, and neither is the receiver sensing it, as they are context-switched by the scheduler to run 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 the Lambdas[14] of various sizes on AWS, and compare the distribution of latencies seen by the receiver during the contention in each case. Figure 2 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 a lot of other instances may need to pause for more time and collect more samples to make up for lost communication due to scheduling.

4.1.4 Overcoming noise. Along with context switching and sensing noise, there are other imperfections in measurement apparatus that cause (minor) noise. For example, we use the difference in Ariana ► is this the first time we mention RDTSC? do we need to define it for a community like this? ◀ RDTSC timer readings before and after the locking operation to measure its Ariana ► what is its in this case? ◀ latency in cycles. If the receiver process is context-switched in between the timer readings (e.g., at line 8 in Algorithm 2), 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 2. To overcome missed samples and noise, we record hundreds of samples and compare it to the baseline

581 distribution of latencies sampled without contention. We then need
 582 to compare and differentiate the observed sample of latencies from
 583 the baseline to establish contention. Since we effectively need to
 584 compare the empirical CDFs of the baseline measurements to the
 585 contention measurements, we decide to use a variant of the two-
 586 sample Kolmogorov-Smirnov (KS) test. By taking the mean of the
 587 absolute differences between the empirical CDFs, we can categorize
 588 a KS-value above a certain threshold (KS-threshold) as a 1-bit or a
 589 0-bit.

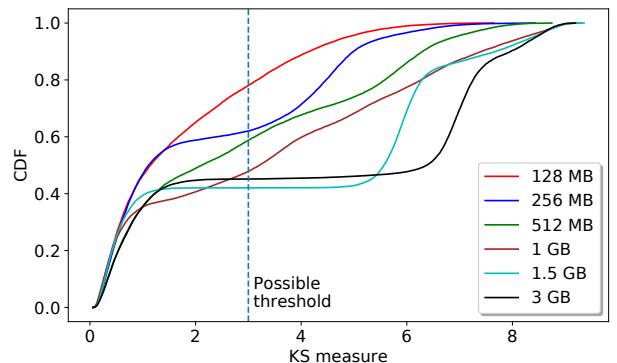
590 To determine the threshold, we deploy a large number of lambdas
 591 across AWS regions. Some of these lambdas cause contention (aka
 592 senders) while others observe contention by collecting samples
 593 of latencies (aka receivers). Each of the samples may or may not
 594 have observed contention depending on whether the receiver was
 595 co-located with a sender lambda (an unknown at this point). We
 596 then calculate the KS-value for each sample against the baseline and
 597 plot a CDF of these values for lambdas of different sizes in Figure 3.
 598 Ideally, we expect a bi-modal distribution (stepped CDF) with the
 599 lower and upper peaks corresponding to samples that have not and
 600 have seen contention respectively, and a big gap between the two
 601 (long step). Fortunately, we examine this differentiation with larger
 602 lambda sizes (which allows us to choose a clear threshold), but we
 603 do not examine a clear differentiation with smaller lambdas, where
 604 scheduling instability causes lossy communication (discussed in
 605 4.1.3). This trend also reflects in the reliability of our technique
 606 across various lambda sizes, as we will show in our evaluation.
 607 Based on the plot, we picked a KS-threshold at 3.0 which seems to
 608 be constant across AWS regions, suggesting that this is a platform
 609 constant.

610 We present the pseudo-code of a receiver lambda in Algorithm 2,
 611 which includes all the limitations discussed thus far. Ariana ► I'm
 612 not a huge fan of this sentence...need to find a way to rework◀

613 4.1.5 *Clock synchronization*. Since communicating each bit of in-
 614 formation takes time (i.e., receiver sampling duration), our algo-
 615 rithm requires synchronizing sender and receiver at the start of each
 616 bit. In traditional analog channels, this is achieved either using a sepa-
 617 rate clock signal or a self-clocking signal encoding. For example,
 618 [?] uses differential Manchester encoding for clock synchroniza-
 619 tion for the memory bus covert channel. Self-clocking encodings
 620 become much trickier (Ques ► why?◀) when there are multiple
 621 senders and receivers. In this work, we use the system clock for
 622 synchronizing communication. All the instances involved in the
 623 communication would be running on the same physical server and
 624 so they share the server's clock. The system clock on AWS lambdas,
 625 for example, is precise up to nanoseconds with a sub-microsecond
 626 drift between different lambdas running on the same server, which
 627 is accurate enough as we only work in the millisecond regime due
 628 to sampling noise constraints Ariana ► is there some sort of citation for
 629 this?◀

630 4.2 Protocol

631 In the preceding section, we discussed a communication chan-
 632 nel with synchronized time slots. In each time slot, an instance
 633 can reliably send (broadcast) or receive (listen) a bit by causing
 634 or sensing for contention. Given that there are multiple instances
 635 that may want to broadcast information on the channel, we next



633 **Figure 3: Shows CDF of KS values observed for various**
 634 **lambda sizes. A bimodal distribution with longer step lets**
 635 **us pick a KS-threshold that enables our technique to dif-**
 636 **ferentiate between 0-bit and 1-bit with high confidence.**

637 must determine which instance broadcasts first, to avoid collisions.
 638 Traditional channels like Ethernet or Wireless detect and avoid
 639 collisions by employing a random Ariana ► colloquial?◀ "back-off"
 640 mechanism. This type of mechanism will be challenging to im-
 641 plement in methodology for two reasons. First, lambda instances
 642 do not have the capability of sensing the channel while sending
 643 a bit, which is required for detecting collisions; instances can ei-
 644 ther cause contention or sense it, but not both. Note that senders
 645 do experience a higher latency for locking operations when other
 646 senders are simultaneously causing contention. However, reliably
 647 judging this higher latency requires each sender to already have
 648 calculated a baseline of latencies without collisions Ariana ► chicken
 649 and egg is a bit too colloquial...need to figure out how to communicate this
 650 thought better◀. Second, even implementing a random "back-off"
 651 mechanism will introduce significant overhead before any mean-
 652 ingful communication occurs. This overhead will also increase as
 653 the number of instances involved increases. Since each time slot
 654 takes up to 1 second, the additional overhead can be detrimental to
 655 the efficacy of the communication channel. Anil ► get information
 656 theoretic bounds on the capacity of our channel?◀

657 Ariana ► I think in general we need a bit more sign-posting to guide the
 658 reader. I will think about this a bit more◀. However, a communication
 659 channel for lambda co-residence detection need not be general and
 660 expressive Ariana ► why?◀. Thus, we assume that each instance
 661 involved has a unique fixed-length (say n) bit-string correspond-
 662 ing to its ID that must be communicated. As such, we propose
 663 a communication protocol that exchanges these bit-strings while
 664 allowing for collisions. We divide the running time of the protocol
 665 into phases, with each phase executing for an interval of n bit-slots.
 666 Each phase has a set of participating instances, which in the first
 667 phase would be all of the co-located instances. In each bit-slot k of
 668 n slots in a phase, every participating instance broadcasts a bit if
 669 the k^{th} bit of its bit-string (ID) is 1, otherwise it listens for a 0 or 1.
 670 If an instance senses a 1 while listening, it stops participating, and
 671 listens for the rest of the phase. Thus, only the instances with the
 672 highest ID among the initial set of participating lambdas continues

Algorithm 3 ID exchange protocol TODO ▶ Improve pseudo-code◀

```

697
698 1: sync_point ← Start time for all instances
699 2: ID ← Instance ID
700 3: N ← Number of bits in ID
701 4: advertising ← TRUE
702 5: instances ← {}
703 6: WAIT_TILL(sync_point)
704 7: while id_read do
705 8:   slots ← 0
706 9:   id_read ← 0
707 10:  participating ← advertising
708 11:  while slots < N do
709 12:    bit ← slotsth most significant bit of ID
710 13:    if participating and bit then
711 14:      WRITE_BIT() (Alg. 1)
712 15:      bit_read ← 1
713 16:    else
714 17:      bit_read ← READ_BIT() (Alg. 2)
715 18:      if bit_read then
716 19:        participating ← FALSE
717 20:      end if
718 21:    end if
719 22:    id_read ← 2 * id_read + bit_read
720 23:    slots ← slots + 1
721 24:  end while
722 25:  if id_read = ID then
723 26:    advertising ← FALSE
724 27:  end if
725 28:  instances ← instances ∪ {id_read}
726 29: end while
727 30: return instances
728
729

```

broadcasting until the end of the protocol, effectively advertising its full ID to the rest of the (now listening) lambdas). In the next cycle, the lambda with the previously highest ID now only listens, allowing the next highest instance to advertise its ID, and so on. Since the IDs are unique, there will always be only one instance that broadcasts in every phase. The protocol ends after x phases (where x is number of co-located instances), when none of the instances broadcast for n consecutive bit-slots. A pseudo-code of the protocol is provided in Algorithm 3. Note that the protocol itself is channel-agnostic and can be extended for other (future) covert channels with similar channel properties.

4.2.1 Complexity. Assuming N total deployed instances to the cloud, the bit-string needs to be $\log_2 N$ bits to uniquely identify each instance. If a maximum K of those instances 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 thing. For example, assuming 10,000 deployed lambdas and a maximum of 10 co-located instances on each server, the entire co-residence detection requires around 4 minutes to fully execute (with 1-second time slots). In fact, it is not necessary to run the protocol for all K phases. After the first round, all the co-located instances would know one of their neighbors. Ariana ▶ confused about this sentence◀ Instances can exchange the globally unique IDs offline (through

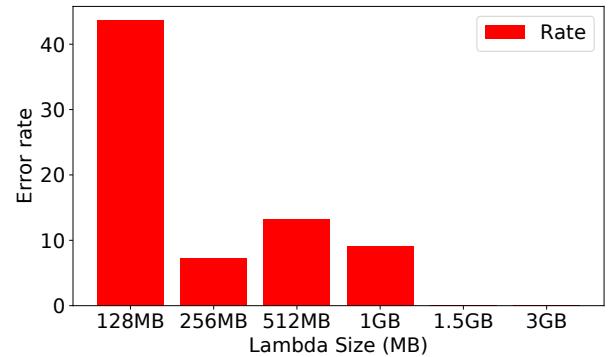


Figure 4: Shows the error rate (as a fraction of 1000 lambdas deployed) for different lambda sizes in AWS Middle-East region.

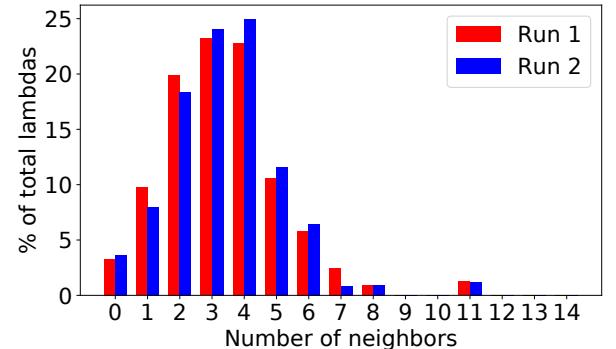


Figure 5: Shows the fraction of lambdas by the number of neighbors they saw for two independent runs that use same set of underlying AWS containers. The perfect correlation shows that both runs depict the colocation status of those containers regardless of the lambdas that ran on them, providing an evidence for the correctness of our technique.

TODO ▶ this doesn't show perfect correlation because only 99% of lambdas warm started. Get another perfect run?◀

the network) to determine the rest of their neighbors. This simplification removes the dependency on number of co-located instances (K) and decreases the complexity to $O(\log_2 N)$, allowing the entire protocol to finish within a minute instead of four.

5 EVALUATION

In this section, we evaluate the effectiveness of our co-residence detection technique with respect to the desirable properties mentioned in section 3 i.e., reliability and scalability.

5.1 Setup

We run all our experiments with AWS[4] lambdas. (We have showed that this covert channel exists on other clouds and so can be easily

replicated with their serverless functions). Activating the memory bus covert channel requires address manipulation using pointers and is not supported by most serverless platform runtimes that only allow high-level languages like Python. However, we were able to find at least one exception to this on all clouds: AWS allows C++ programs while Azure and GCP allow unsafe versions of C# and Golang respectively. Once deployed, each instance participates in the first phase of the protocol as noted in section 4.2.1, thereby learning the ID of their biggest neighbor. As bit-flip errors are possible, we repeat the same phase for two more (independent) rounds and take the majority result to record the ID seen by this instance. If all three rounds resulted in different IDs, we classify this instance as erroneous and report it in the error rate. We group all the successful instances that saw the same ID as neighbors. We repeat the experiments for different lambda sizes and in various cloud regions.

5.2 Reliability

We consider the results of the technique reliable when 1) most of the deployed instances successfully see the same result in majority of the independent rounds (indicating lesser bit-flip errors) and 2) the resulting co-located groups we see match the ground truth. For 1, we ran an experiment with 1000 AWS 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). From figure 4, we can see that smaller lambdas see lot more errors. This is expected because, as discussed in section 4.1.4, these lambdas experience lossy communication making it harder for our technique to sense contention. The lambdas above 1.5 GB though see a 100% success rate.

Correctness Obtaining the ground truth on which instances were "actually" co-located is not possible, considering that that is the purpose of our technique. However, we found a way to validate our results with high confidence. AWS caches the containers used to run lambdas for a while (**TODO** ▶ *how long?*◀) to reuse them[3] for later lambdas and mitigate cold start latencies. For C++ lambdas, we found that the data structures declared in 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 run in a particular container. This means, for a given lambda, we can precisely tell all the lambdas that previously ran in the same container (aka predecessors). Using this, we validated that identical experiments repeated within a few minutes of each other will use the same set of underlying containers for running the deployed lambdas. Since the lambda co-location is essentially co-location of their containers and given that these containers persist across experiments (run within few minutes of each other), co-location results from such experiments must agree on the co-location of their underlying containers.

To demonstrate that this is the case, we run an experiment with 1000 1.5GB cold-started lambdas (ID'ed 1 to 1000) in the one of densest AWS regions (AWS MiddleEast), which resulted in many co-located groups. We repeat the experiment within few seconds, making sure that all 1000 lambdas are warm-started this time (i.e., they use the same set of containers from the previous experiment). For each co-located group of lambdas in the latter experiment,

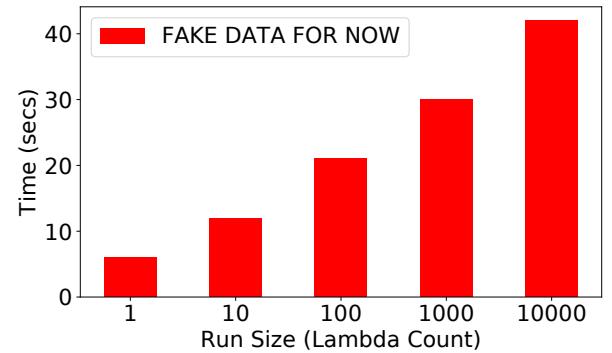


Figure 6: Shows the average runtime of a lambda for co-location runs of different sizes. The run time increases logarithmically with the number of lambdas as it is proportional to the number of bits required to uniquely identify all the lambdas. **TODO** ▶ *Get real data.*◀

we checked whether their predecessor lambdas in the former one formed a co-located group as well. We note that while different lambdas used the containers across the experiments, their co-located groups are perfectly correlated. Figure 5 shows that both experiments saw the same number of groups of different sizes. This proves the correctness of our co-location results.

5.3 Scalability

One of the key properties of this technique is that it's really fast. It takes only a second to communicate each binary bit of the ID, enabling it to scale logarithmically with the number of lambdas involved. Figure 6 shows this for experiments involving different number of lambdas. For a run with 10000 lambdas for example, each lambda can find its neighbors within a minute of its invocation, leaving ample time to perform other things using this information. This also means the cost per lambda also scales logarithmically, making it very cost-effective.

6 PLACEMENT STUDY

Our goal in developing this technique is to demonstrate how easy it is for attackers to exploit the pervasive memory bus covert channel and obtain co-residency information, thereby motivating the need to address it. This information can be used by attackers in aiding a lot of attack scenarios **TODO** ▶ *for example?*◀ or simply learn the internal mechanisms of a cloud. In the section, we explore some ways in which the tool can be used to gain some insights into lambda activity in some AWS regions. Unless specified otherwise, all the experiments are performed with 1.5 GB lambdas and ran successfully with zero error rate. We find that **TODO** ▶ *summary of takeaways from the study.*◀

6.1 Co-residence across AWS regions

We ran co-residence detection in different AWS regions with 1000 1.5GB Lambdas. Figure 8 shows multiple plots showing the co-located groups, one per region, with each bar in the plot showing the

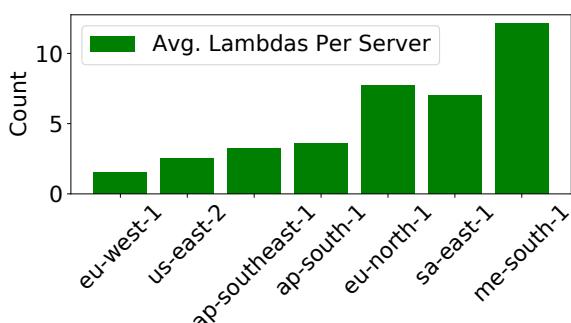


Figure 7: Shows the average number of lambdas per server i.e., colocation density seen in various AWS regions for a 1000-lambda run.

fraction of lambdas that saw a certain number of neighbors (i.e., that belong to a colocated group of certain size). Plots that are right-heavy (towards bottom-right) indicate higher colocation density compared to the left-heavy (towards top-left) ones (which is also illustrated in figure 7). We can see that most regions have almost all lambdas see at least one neighbor (smaller or non-existent bar at 0). Assuming that the cloud placement scheduler tries to efficiently bin pack lambdas, we hypothesize that the colocation is dependent on the total number of servers and the lambda activity in the region, both of which can be assumed to be lower in newer AWS regions. We note that most right-heavy plots correspond to the relatively newer AWS regions. The maximum size of a co-located group we ever saw was 25 (1.5G) lambdas on a single machine. TODO ►any other insights I missed?◀

6.2 Weekly & Daily Patterns

TODO ►Run experiments in different times of the day or using different deployment strategies that may affect co-location.◀

6.3 Different accounts

TODO ►Run experiments with lambdas from different user accounts and see how the colocation is affected between single vs different accounts.◀

7 DISCUSSION

TODO ►KG ►dssd What are the limitations of the technique? What can an attacker do with this colocation information? How can AWS prevent this? Something about firecracker? Helps performance isolation studies such as [22]. This technique can be used with other containers like VMs. While the co-operative coresidence detection itself does not directly help attackers locate their victims, it can help attackers perform highly sophisticated DDOS attacks once they did find the victim.

8 RELATED WORK

TODO ►needs polishing◀

Covert Channels & Cloud Attacks Co-residency is possible because of covert channels, so we begin our related work with an

investigation into covert channels and cloud attacks. Initial papers in co-residency detection utilized host information and network addresses arising due to imperfect virtualization [18]. However, these covert channels are now obsolete, as cloud provides have strengthened virtualization and introduced Virtual Private Clouds TODO ►cite◀. Later work used cache-based channels in various levels of the cache [12, 16, 24, 32] and hardware based covert channels like thermal covert channels [17], RNG module [7] and memory bus [23] have also been explored in the recent past. Moreover, studies have found that VM performance can be significantly degraded using memory DDoS attacks [28], while containers are susceptible to power attacks from adjacent containers [8]. Our work focuses on using the memory bus as a covert channel for determining co-operative co-residency. Ariana ►work this in?◀ Covert channels using memory bus were first introduced by Wu et al. [?], and subsequently has been used for co-residency detection on VMs and Containers [19?]. Wu et. al [?] introduced a new technique to lock the memory bus by using atomic memory operations on addresses that fall on multiple cache lines. .

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) [18]. This initial work was further expanded in subsequent years to examine co-residency using memory bus locking [27] and active traffic analysis [6], as well as determining placement vulnerabilities in multi-tenant Public Platform-as-a-Service systems [20, 30]. Finally, Zhang et al. demonstrated a technique to detect VM co-residency detection via side-channel analyses [31]. Our work expands on these previous works by investigating co-residency for lambdas.

Lambdas TODO ►pretty rough paragraph...needs a rewrite◀ While lambdas are a newer technology than VMs, there still exists a variety of literature. For example, recent studies examine cost comparisons of running web applications in the cloud on lambdas versus other architectures [21]. Moreover, lambdas have been studied in the context of cost-effective batching and data processing [13]. Further research has shown how lambdas perform with scalability and hardware isolation, indicating some flaws in the lambda architecture [22]. From a security perspective, Izhikevich et. al examined lambda co-residency using RNG and memory bus techniques (similar to techniques used looking at 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.

9 ETHICAL CONSIDERATIONS

As with any large scale measurement project, there are ethical considerations to take into account. 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-location for the lambdas we launched, and do not gain insight into other consumer's lambdas. Secondly, there is concern that our experiments may cause performance issues

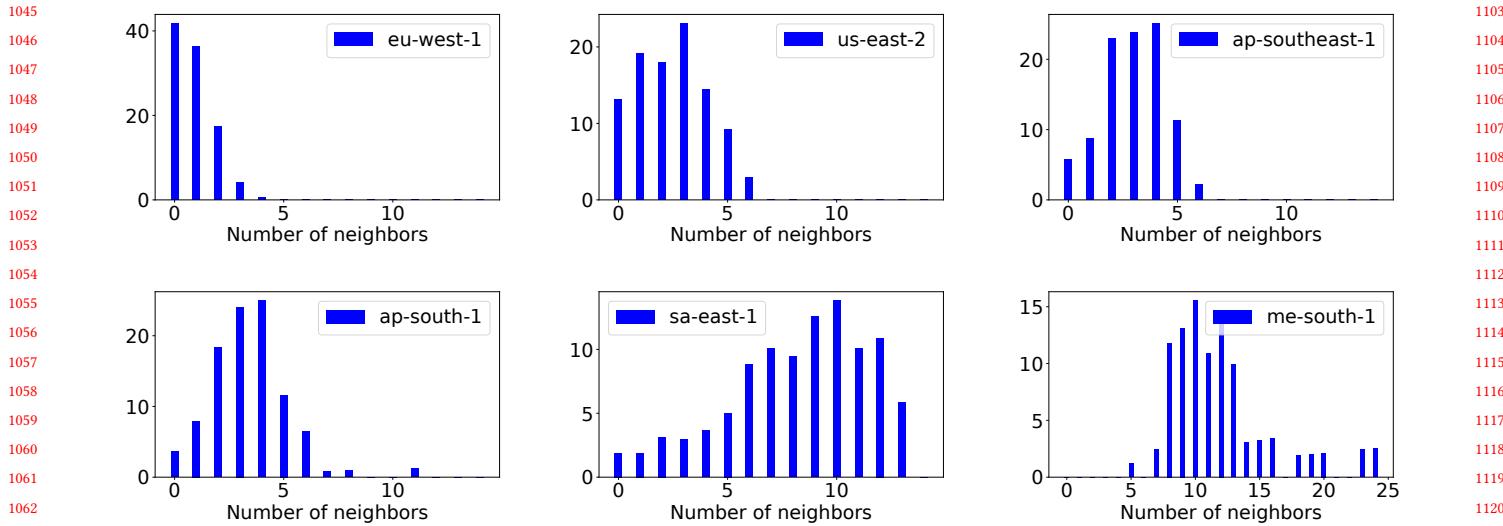


Figure 8: Shows co-location results for a 1000-lambda run in different AWS regions. Each bar shows the fraction of those 1000 lambdas (in %) that saw a certain number of neighbors. The total amount and density of co-location vary widely across regions, perhaps based on the size and lambda activity within those regions.

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. Moreover, memory accesses and locking operations are FIFO, which prevents starvation to the multiple processes 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, by definition, be affected in their longterm computational goals. **Ariana** ► is lambda cost by second? is it possible that by causing a longer lambda runtime, we are costing someone more money?◀

10 CONCLUSION & FUTURE WORK

TODO ► *Copied abstract* ◀ Cloud computing has seen explosive growth in the past decade. This is made possible by efficient sharing of infrastructure among tenants, which unfortunately also raises security challenges like preventing side-channel attacks. Providers, like AWS and Azure, have traditionally relied on hiding the co-residency information to prevent targeted attacks in their clouds. But recent works have repeatedly found co-residence detection techniques that break this encapsulation, prompting the providers to address them and harden isolation on their platforms. In this work, we find yet another such technique based on a memory bus covert channel that is more pervasive, reliable and harder to fix. We show that we can use this technique to reliably perform co-operative co-residence detection for thousands of AWS lambdas within a few seconds, which opens a way for attackers to perform DDoS attacks or learn cloud's internal mechanisms. We present this technique in detail, evaluate it and use it to perform a small study on lambda activity across a few AWS regions. Through this work, we hope to motivate the need to address this covert channel in the cloud.

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] Michael Armbrust, Armando Fox, Rean Griffith, Anthony D. Joseph, Randy Katz, Andy Konwinski, Gunho Lee, David Patterson, Ariel Rabkin, Ion Stoica, and Matei Zaharia. 2010. A View of Cloud Computing. *Commun. ACM* 53, 4 (April 2010), 50–58. <https://doi.org/10.1145/1721654.1721672>
- [3] aws 2018. Understanding Container Reuse in AWS Lambda. <https://aws.amazon.com/blogs/compute/container-reuse-in-lambda/>.
- [4] aws 2019. Amazon EC2. <https://aws.amazon.com/ec2/>.
- [5] Azure 2019. Azure VMs. <https://azure.microsoft.com/en-us/services/virtual-machines/>.
- [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>
- [7] Dmitry Evtyushkin and Dmitry Ponomarev. 2016. Covert Channels Through Random Number Generator: Mechanisms, Capacity Estimation and Mitigations. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security (Vienna, Austria) (CCS '16)*. ACM, New York, NY, USA, 843–857. <https://doi.org/10.1145/2976749.2978374>
- [8] Xing Gao, Zhongshu Gu, Mehmet Kayaalp, Dimitrios Pendarakis, and Haining Wang. 2017. ContainerLeaks: Emerging Security Threats of Information Leaks in Container Clouds. In *Dependable Systems and Networks (DSN), 2017 47th Annual IEEE/IFIP International Conference on*. IEEE, IEEE, 237–248.
- [9] gcp 2019. Google Cloud Functions. <https://cloud.google.com/functions/>.
- [10] GoogleCloud 2019. Google Compute. <https://cloud.google.com/products/compute/>.
- [11] Elizabeth Izhevich. 2018. *Building and Breaking Burst-Parallel Systems*. Master’s thesis. University of California, San Diego.
- [12] Mehmet Kayaalp, Nael Abu-Ghazaleh, Dmitry Ponomarev, and Aamer Jaleel. 2016. A high-resolution side-channel attack on last-level cache. 1–6. <https://doi.org/10.1145/2897937.2897962>
- [13] Mariam Kiran, Peter Murphy, Inder Monga, Jon Dugan, and Sartaj Singh Baveja. 2015. Lambda architecture for cost-effective batch and speed big data processing. *2015 IEEE International Conference on Big Data (Big Data) (2015)*, 2785–2792.
- [14] Lambda 2019. AWS Lambda. <https://aws.amazon.com/lambda/>.
- [15] Moritz Lipp, Michael Schwarz, Daniel Gruss, Thomas Prescher, Werner Haas, Anders Fogh, Jann Horn, Stefan Mangard, Paul Kocher, Daniel Genkin, Yuval Yarom, and Mike Hamburg. 2018. Meltdown: Reading Kernel Memory from User Space. In *27th USENIX Security Symposium (USENIX Security 18)*.

1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160

- [16] Fangfei Liu, Yuval Yarom, Qian ge, Gernot Heiser, and Ruby Lee. 2015. Last-Level Cache Side-Channel Attacks are Practical, Vol. 2015. 605–622. <https://doi.org/10.1109/SP.2015.43>

[17] Ramya Jayaram Masti, Devendra Rai, Aanjan Ranganathan, Christian Müller, Lothar Thiele, and Srdjan Capkun. 2015. Thermal Covert Channels on Multi-core Platforms. In *Proceedings of the 24th USENIX Conference on Security Symposium* (Washington, D.C.) (SEC'15). USENIX Association, Berkeley, CA, USA, 865–880. <http://dl.acm.org/citation.cfm?id=2831143.2831198>

[18] Thomas Ristenpart, Eran Tromer, Hovav Shacham, and Stefan Savage. 2009. Hey, You, Get off of My Cloud: Exploring Information Leakage in Third-party Compute Clouds. In *Proceedings of the 16th ACM Conference on Computer and Communications Security* (Chicago, Illinois, USA) (CCS '09). ACM, New York, NY, USA, 199–212. <https://doi.org/10.1145/1653662.1653687>

[19] Venkatanathan Varadarajan, Yinqian Zhang, Thomas Ristenpart, and Michael Swift. 2015. A Placement Vulnerability Study in Multi-Tenant Public Clouds. 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>

[20] Venkatanathan Varadarajan, Yinqian Zhang, Thomas Ristenpart, and Michael M. Swift. 2015. A Placement Vulnerability Study in Multi-tenant Public Clouds. *CoRR* abs/1507.03114. arXiv:1507.03114 <http://arxiv.org/abs/1507.03114>

[21] Mario Villamizar, Oscar GarcÃAs, Lina Ochoa, Harold Castro, Lorena Salamanca, Mauricio Verano Merino, Rubby Casallas, Santiago Gil, Carlos Valencia, Angee Zambrano, and Mery Lang. 2016. Infrastructure Cost Comparison of Running Web Applications in the Cloud Using AWS Lambda and Monolithic and Microservice Architectures. <https://doi.org/10.1109/CCGrid.2016.37>

[22] Liang Wang, Mengyuan Li, Yinqian Zhang, Thomas Ristenpart, and Michael Swift. 2018. Peeking Behind the Curtains of Serverless Platforms. In *2018 USENIX Annual Technical Conference (USENIX ATC 18)*. USENIX Association, Boston, MA, 133–146. <https://www.usenix.org/conference/atc18/presentation/wang-liang>

[23] Zhenyu Wu, Zhang Xu, and Haining Wang. 2012. Whispers in the Hyperspace: High-speed Covert Channel Attacks in the Cloud. In *Presented as part of the 21st USENIX Security Symposium (USENIX Security 12)*. USENIX, Bellevue, WA, 159–173. <https://www.usenix.org/conference/usenixsecurity12/technical-sessions/presentation/wu>

[24] Yunjing Xu, Michael Bailey, Farnam Jahanian, Kaustubh Joshi, Matti Hiltunen, and Richard Schlichting. 2011. An Exploration of L2 Cache Covert Channels in Virtualized Environments. In *Proceedings of the 3rd ACM Workshop on Cloud Computing Security Workshop* (Chicago, Illinois, USA) (CCSW '11). ACM, New York, NY, USA, 29–40. <https://doi.org/10.1145/2046660.2046670>

[25] Yunjing Xu, Michael Bailey, Farnam Jahanian, Kaustubh Joshi, Matti Hiltunen, and Richard Schlichting. 2011. An Exploration of L2 Cache Covert Channels in Virtualized Environments. In *Proceedings of the 3rd ACM Workshop on Cloud Computing Security Workshop* (Chicago, Illinois, USA) (CCSW '11). Association for Computing Machinery, New York, NY, USA, 29–40. <https://doi.org/10.1145/2046660.2046670>

[26] Zhang Xu, Haining Wang, and Zhenyu Wu. 2015. A Measurement Study on Co-residence Threat inside the Cloud. In *24th USENIX Security Symposium (USENIX Security 15)*. USENIX Association, Washington, D.C., 929–944. <https://www.usenix.org/conference/usenixsecurity15/technical-sessions/presentation/xu>

[27] Zhang Xu, Haining Wang, and Zhenyu Wu. 2015. A Measurement Study on Co-residence Threat inside the Cloud. In *24th USENIX Security Symposium (USENIX Security 15)*. USENIX Association, Washington, D.C., 929–944. <https://www.usenix.org/conference/usenixsecurity15/technical-sessions/presentation/xu>

[28] Tianwei Zhang, Yinqian Zhang, and Ruby Lee. 2016. Memory DoS Attacks in Multi-tenant Clouds: Severity and Mitigation.

[29] Weijuan Zhang, Xiaoqi Jia, Chang Wang, Shengzhi Zhang, Qingjia Huang, Ming-sheng Wang, and Peng Liu. 2016. A Comprehensive Study of Co-residence Threat in Multi-tenant Public PaaS Clouds. In *Information and Communications Security*, Kwok-Yan Lam, Chi-Hung Chi, and Sihan Qing (Eds.). Springer International Publishing, Cham, 361–375.

[30] Weijuan Zhang, Xiaoqi Jia, Chang Wang, Shengzhi Zhang, Qingjia Huang, Ming-sheng Wang, and Peng Liu. 2016. A Comprehensive Study of Co-residence Threat in Multi-tenant Public PaaS Clouds, Vol. 9977. 361–375. https://doi.org/10.1007/978-3-319-50011-9_28

[31] Yinqian Zhang, Ari Juels, Alina Oprea, and Michael Reiter. 2011. HomeAlone: Co-Residency Detection in the Cloud via Side-Channel Analysis. 313 – 328. <https://doi.org/10.1109/SP.2011.31>

[32] Yinqian Zhang, Ari Juels, Michael K. Reiter, and Thomas Ristenpart. 2014. Cross-Tenant Side-Channel Attacks in PaaS Clouds. In *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security* (Scottsdale, Arizona, USA) (CCS '14). ACM, New York, NY, USA, 990–1003. <https://doi.org/10.1145/2660267.2660356>