

# Columbus: Fast, Reliable Co-residence Detection for Lambdas

Anonymous Author(s)

## ABSTRACT

Cloud computing has seen explosive growth and popularity, especially “serverless” service modalities, such as AWS lambdas. These services are lighter-weight and provide more flexibility in scheduling and cost, which contributes to their popularity, however the security issues associated with serverless computing is not well understood. In this work, we construct a covert channel, a means of transmitting data outside of traditional channels, that is practical for AWS lambdas. We extend the work in Wu et. al to utilize the memory-bus to develop a co-residence detector that is generic, reliable, and scalable. This technique performs dynamic neighbor discovery and is incredibly fast, executing in a matter of seconds. In addition to this protocol for lambda co-residence detection, we also perform and evaluate a series of measurements that exemplify the practicality of this approach in a cloud setting. Through this work, we show that efforts to secure co-residency detection in cloud providers is not yet complete.

## CCS CONCEPTS

- Security and privacy → Virtualization and security.

## KEYWORDS

cloud, cartography, serverless, coresidency, covert channels

## ACM Reference Format:

Anonymous Author(s). 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, 12 pages. <https://doi.org/10.1145/nnnnnnnnnnnnnn>

## 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 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 service 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 services?

Covert channels, as a general matter, 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 such a channel. First, 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 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

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](mailto: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/nnnnnnnnnnnnnn>

117 the balance of their limited lifetimes to effect communications. Finally, since rendezvous in a serverless system is inherently statistical,  
 118 any such protocol must anticipate the potential for interference (i.e., when multiple sending lambdas happen to be co-resident  
 119 with a single receiving lambda).

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

- 122 • **Fast co-residence detection.** Leveraging the memory-bus  
 123 contention work of Wu et. al [26], we develop and implement  
 124 a lambda co-residence detector that is generic, reliable,  
 125 scalable and, most importantly, fast, executing in a matter of  
 126 seconds for thousands of concurrent lambdas.
- 127 • **Dynamic neighbor discovery.** We extend our co-residence  
 128 detector to enumerate the set of co-resident neighbors, a require-  
 129 ment to avoid unwanted communication interference  
 130 between neighbors.
- 131 • **Serverless density measurement.** We empirically estab-  
 132 lish measures of serverless function density – that is, for a  
 133 given number of short-lifetime lambdas launched at a point  
 134 in time, how many would be expected to become co-resident?  
 135 We conduct these measurements across a range of Amazon  
 136 data centers to establish that such an approach is practica-  
 137 ble.

138 The remainder of this paper is organized as follows. In section 2,  
 139 we present some background and motivate the need for co-residence  
 140 detection for lambdas. Sections 3 and 4 present our co-residence de-  
 141 tector mechanism in detail. We evaluate the mechanism in section  
 142 5 and perform a study of AWS lambda placement density using our  
 143 technique to demonstrate the practicality of lambda covert chan-  
 144 nels in section 6. Finally, we conclude with a discussion on TODO  
 145 ►what◀ in section 7.

## 146 2 BACKGROUND & MOTIVATION

147 We begin with a brief background on relevant topics as we mo-  
 148 tivate the need for fast and scalable co-residence detection in en-  
 149 abling covert channels with lambdas.

### 150 2.1 Lambdas/Serverless Functions

151 We focus on serverless functions in this paper, as they are one of  
 152 the fastest-growing cloud services and are less well-studied from  
 153 a security standpoint. Offered as lambdas on AWS [14], and cloud  
 154 functions on GCP [10] and Azure [5], these functions are of interest  
 155 because they do not require the developer to provision, maintain,  
 156 or administer servers. In addition to this low overhead, lambdas  
 157 are much more cost-efficient than virtual machines (VMs) as they  
 158 allow more efficient packing of functions on servers. Lambdas ex-  
 159 ecute as much smaller units than containers and virtual machines,  
 160 and are more ephemeral. For example, on AWS, the memory of  
 161 lambdas is capped at 3 GB, with a maximum execution limit of 15  
 162 minutes. As with other cloud services, the user has no control over  
 163 the physical location of the server(s) on which their lambdas are  
 164 spawned.

165 While lambdas are limited in the computations they can execute  
 166 (typically written in high-level languages like Python, C#, etc),

167 they are conversely incredibly lightweight and can be initiated  
 168 and deleted in a very short amount of time. Cloud providers run  
 169 lambdas in dedicated containers with limited resources (e.g., Fire-  
 170 cracker [1]), which are usually cached and re-used for future lamb-  
 171 das to mitigate cold-start latencies [2]. The ephemeral nature of  
 172 serverless functions and their limited flexibility increases the diffi-  
 173 culty in detecting co-residency, as we will discuss later. Previous  
 174 studies that profiled lambdas [24] focused on the performance as-  
 175 pects like cold start latencies, function instance lifetime, CPU us-  
 176 age, etc. across various clouds, while the security aspects remain  
 177 relatively understudied.

## 178 2.2 Covert Channels in the Cloud

179 In our attempt to shed light on the security aspects of lambdas,  
 180 we focus particularly on the feasibility of establishing a reliable  
 181 covert channel in the cloud using lambdas. Covert channels enable  
 182 a means of transmitting information between entities that bypasses  
 183 traditional monitoring or auditing. Typically, this is achieved by  
 184 communicating data across unintended channels such as signalling  
 185 bits by causing contention on shared hardware media on the server [16,  
 186 18, 25–27]. Past work has demonstrated covert channels in virtualized  
 187 environments like the clouds using various hardware such as caches [19, 27], memory bus [26], and even processor tempera-  
 188 ture [16].

189 **Memory bus covert channel** Of particular interest to this work  
 190 is the covert channel based on memory bus hardware introduced  
 191 by Wu et al. [26]. In x86 systems, atomic memory instructions de-  
 192 signed to facilitate multi-processor synchronization are supported  
 193 by cache coherence protocols as long as the operands stay within a  
 194 cache line (generally the case as language compilers make sure that  
 195 operands are aligned). However, if the operand is spread across  
 196 two cache lines (referred to as "exotic" memory operations), x86  
 197 hardware achieves atomicity by locking the memory bus to pre-  
 198 vent any other memory access operations until the current opera-  
 199 tion finishes. This results in significantly higher latencies for such  
 200 locking operations compared to traditional memory accesses. As  
 201 a result, a few consecutive locking operations could cause con-  
 202 tention on the memory bus that could be exploited for covert com-  
 203 munication. Wu et al. achieved a data rate of 700 bps on the mem-  
 204 ory bus channel in an ideal laboratory setup.

205 Cloud environments, however, pose challenges to both estab-  
 206 lishing such covert channels and using them at their ideal perfor-  
 207 mance. First, in the cloud setting, the sender and receiver entities  
 208 are not known to be on the same server as this (coresidency) in-  
 209 formation is hidden from the tenants (even if the entities belong  
 210 to the same tenant). Second, clouds present a virtualized platform  
 211 that affects the communication on the covert channel as it intro-  
 212 duces scheduling uncertainties between the participating entities  
 213 and interference from other non-participating workloads. While  
 214 studies generally deal with the latter problem by employing tradi-  
 215 tional error correction techniques [26] to overcome errors, the  
 216 former challenge of establishing the covert channel (by detecting  
 217 coresidency) in the first place requires a co-residence detection  
 218 mechanism which the attacker needs to employ to get the sender  
 219 and receiver on the same machine.

### 233 2.3 Co-residence Detection

234 Past research has used various strategies to achieve co-residency  
 235 for demonstrating various covert channel attacks in the cloud. Typically,  
 236 the attacker launches a large number of cloud instances  
 237 (VMs, Containers, etc.) following a certain launch strategy and  
 238 employs a co-residence detection mechanism for detecting if any  
 239 pair of those instances are running on the same machine. Traditionally,  
 240 the co-residence detection was done based on software  
 241 runtime information like public/internal IP addresses [19], files in  
 242 *procfs* or other environment variables [24, 26] and other such logical  
 243 side-channels [21, 29] that two instances running on a same  
 244 server might share.

245 As virtualization platforms move towards stronger isolation between  
 246 instances (e.g. AWS’ Firecracker VM [1]), these logical covert-  
 247 channels have become less effective or infeasible. Furthermore, some  
 248 of these channels were only effective on container-based platforms  
 249 that shared the underlying OS image and were thus less suitable for  
 250 hypervisor-based platforms. This prompted a move towards using  
 251 hardware-based covert channels, such as the ones discussed in the  
 252 earlier section, which can bypass software isolation and are usually  
 253 harder to fix. For example, Varadarajan et al. [22] uses the memory  
 254 bus covert channel to detect co-residency for EC2 instances. How-  
 255 ever, as we will show, traditional approaches do not extend well  
 256 to lambdas as they are too slow or are not scalable. Thus, we set  
 257 forth to determine a fast and scalable co-residence detection for  
 258 lambdas.

## 260 3 METHODOLOGY

262 Our goal is to determine a (co-operative) co-residence detection  
 263 mechanism for lambdas. In other words, given a series of spawned  
 264 lambdas in a certain region of a cloud service, how can we deter-  
 265 mine the lambdas that are co-located on the same machines? In  
 266 this section, we discuss the details of such a mechanism, previous  
 267 solutions to this problem, and the unique challenges we faced with  
 268 lambda co-residence.

269 Given a set of cloud instances (VMs, Containers, Functions, etc)  
 270 deployed to a public cloud, a co-residence detection mechanism  
 271 would identify, for each pair of instances in the set, whether the  
 272 pair was running on the same physical server at some point. Para-  
 273 phrasing Varadarajan et al.[22], for any such mechanism to be useful  
 274 across a wide range of launch strategies, we observe that it  
 275 should have the following desirable properties:

- 276 • **Generic** The technique should be applicable across a wide  
 277 range of server architectures and software runtimes. In prac-  
 278 tice, the technique would work across most third-party clouds  
 279 and even among different platforms within a cloud.
- 280 • **Reliable** The technique should have a reasonable detection  
 281 success with minimal false negatives (co-resident instances  
 282 not detected) and even less false positives (non co-resident  
 283 instances categorized as co-resident).
- 284 • **Scalable** A launch strategy may require hundreds or even  
 285 thousands of instances to be deployed, and must scalable  
 286 such that the technique will take less time to detect all co-  
 287 resided pairs at a reasonable cost.

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

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

294 Given these properties, we decide to investigate hardware-based  
 295 covert channels. Hardware-based covert-channels are more diffi-  
 296 cult to remove and obfuscate than software-based covert channels,  
 297 and are also ubiquitous, given that hardware is more homogenous  
 298 across computing platforms than software.

299 *3.0.1 Memory bus channel.* We chose the memory bus covert chan-  
 300 nel described in section 2.2 as it exploits a fundamental hardware  
 301 vulnerability that is present across all generations of x86 hardware.  
 302 Historically, multiple public cloud services have been vulnerable  
 303 to this channel [21, 32], and we found that they are still vul-  
 304 nerable today. To demonstrate the presence of the vulnerability, we  
 305 measure the latency of atomic operations on a 4B memory region  
 306 as we slide the region from one cacheline into another across the  
 307 cacheline boundary. We perform this experiment on three major  
 308 clouds (AWS, Google and Microsoft Azure) and show the latencies  
 309 observed in Figure 1. From the figure, we can see that all three  
 310 clouds still exhibit significant difference in latencies for the “exotic”  
 311 memory locking operations (where the memory region falls  
 312 across cacheline boundary) when compared to regular memory ac-  
 313 cesses, demonstrating the presence of this covert channel on all of  
 314 them. Moreover, we were able to run these experiments on server-  
 315 less function instances. Since lambdas have runtimes that are gen-  
 316 erally restricted to high-level languages (that prevent the pointer  
 317 arithmetic required to perform these exotic operations), we used  
 318 the unsafe environments on these clouds – C++ on AWS, Unsafe  
 319 Go on GCP, Unsafe C# On Azure. This shows the applicability of  
 320 using the covert channel across different kinds of cloud instances  
 321 as well.

322 *3.0.2 Previous approaches.* Previous works that used the memory  
 323 bus for co-residence detection divide the deployed instances into  
 324 sender and receiver roles, and attempt to colocate the sender role  
 325 with a receiver role. The sender instances continually lock the mem-  
 326 ory bus (locking process) for a certain duration (~10 seconds) while  
 327 the receiver samples the memory for any spike in access latencies  
 328 (probing process). If all the deployed instances try the detection i.e.,  
 329 locking and probing at once, (some of) the receivers may see lock-  
 330 ing effects, but there would be no way of knowing which or how  
 331 many senders co-resided with a particular receiver and caused the  
 332 locking. This provides no information about the number of phys-  
 333 ical servers that ran these instances or the amount of co-location.  
 334 The only information we can deduce is that receivers were prob-  
 335 ably co-located with just a single sender.

336 An alternative method is to try pair-wise detection where only  
 337 one sender instance locks and one receiver instance probes at a  
 338 time revealing co-residence of this pair, and repeating this seri-  
 339 ally for each pair. However, this technique is too slow and scales  
 340 quadratically with the number of instances e.g., a hundred instances  
 341 take more than 10 hours assuming 10 secs for each pair. Varadar-  
 342 jan et al.[21] speeds up this process significantly by performing  
 343 detection for mutually-exclusive subsets in parallel, allowing for  
 344

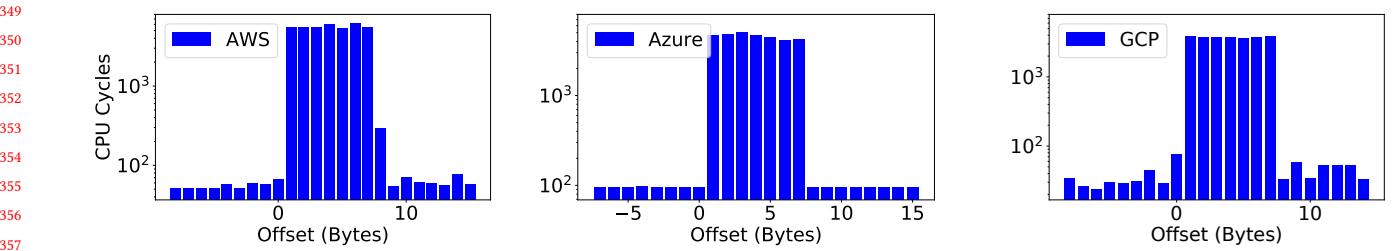


Figure 1: From left to right, the plots show the latencies of atomic memory operations performed on an 8B memory region as we slide it from one cache line across the boundary into another, on AWS, Azure and Google (GCP) clouds respectively. The latencies are orders of magnitude higher when the 8B region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

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 only scale linearly in the best case, 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.

**3.0.3 The Path to Scalability.** One method to quicken the co-location process is by decreasing the time a single sender-receiver pair requires to determine co-residence i.e., improving upon probing time and accuracy of the receiver. However, this method only affects total time by a constant factor. To improve scalability, we need to be able to run the detection for different sender-receiver pairs in parallel without sacrificing the certainty of information we get when they are run serially. For example, when two pairs observe co-residence, we must be certain that each receiver experienced co-residence because of its own paired sender instance, which is not possible if co-residence is ascertained based a simple yes/no signal from the sender instances.

Previous work utilized the memory bus channel to exchange more complex information like keys [26]. At first sight, it appears that the memory bus can be used to exchange information, such as unique IDs, between the co-resided instances, which would allow us to ascertain receiver/sender pairs. However, the previous work assumes that there is only one sender and one receiver that 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 bus covert channel to exchange information between the instances that have access to the channel. Using the protocol, co-resided instances can reliably exchange their unique IDs with each other to discover their neighbors. 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), a number that is orders of magnitude less than total number of instances deployed to the cloud

---

#### 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

```

---

(hundreds to thousands). Removing this dependency on total number of instances deployed lets us scale our co-residence detection significantly, as we will see in the the next section.

## 4 NEIGHBOR DISCOVERY PROTOCOL

Co-residence detection for lambdas must be quick, as lambdas are ephemeral by nature. As noted earlier, co-residence detection is faster when co-residing instances on each server communicate among themselves and discover each other. Assuming that the instances have unique IDs, co-resided instances can exchange integer IDs with each other using the memory bus channel as a transmission medium, which allows us to scale co-residency detection. 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.

### 4.1 Reliable Transmission

First, we must determine how to reliably transmit information between the sender and receiver. Senders and receivers can accurately communicate 0-bits and 1-bits by causing contention on the memory bus. 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 via the memory bus covert channel. To communicate a 1-bit, the sender instance

465 causes contention on the memory bus by locking it using the spe-  
 466 cial memory locking operations (discussed in section ??). The re-  
 467 ceiver would then sample the memory bus for contention, infer-  
 468 ring whether the communication is a 1-bit (when contention is ob-  
 469 served) or a 0-bit (when contention is not observed). Pseudo-code  
 470 for the sender instance is shown in Algorithm 1.

471 *4.1.1 Sensing contention.* Next, we determine how best for the re-  
 472 ceivers to sense contention on the memory bus. There are two  
 473 ways to achieve this. When the memory bus is locked, any non-  
 474 cached memory accesses will queue and therefore see higher laten-  
 475 cies. The receiver can then continually make un-cached memory  
 476 accesses (referred to as the *memory probing* receiver in previous  
 477 literature [22]) and observe a spike in their latencies to detect con-  
 478 tention. On the other hand, the receiver can also detect memory  
 479 bus contention by using the same memory locking operations as  
 480 the sender (referred to as *memory locking* receiver) to probe the  
 481 memory bus. Since only one processor core can lock the memory  
 482 bus at a given time, any other concurrent locking operation will  
 483 see higher latency.

484 Of these two methods, we decide to use the memory locking  
 485 receiver for our experiments. Previous studies [22, 26] have estab-  
 486 lished that both memory probing and memory locking receivers  
 487 experience significant latency overhead during memory bus con-  
 488 tention, making them both viable avenues for sensing the covert-  
 489 channel. Memory probing involves regular (un-cached) memory  
 490 accesses, which is universal, unlike the locking operations which  
 491 are rarely used, if at all, by standard applications. This makes mem-  
 492 ory probing the only viable option for **non-cooperative** co-residence  
 493 detection, where receivers are not under the senders’s control and  
 494 cannot be assumed to perform locking operations. Furthermore,  
 495 memory probing can be done on multiple receivers constantly with-  
 496 out affecting each other (due to the high memory bandwidth), which  
 497 prevents noise in measurements. This is an important attribute,  
 498 as memory locking receivers must contend with this noise. How-  
 499 ever, bypassing multi-levels of caches in today’s servers to per-  
 500 form memory accesses with reliable consistency is a challenging  
 501 task. Even with a reliable cache-bypassing technique, the variety  
 502 of cache architectures and sizes that we encounter on different  
 503 clouds would make tuning the technique to suit these architectures  
 504 an arduous task while reducing the applicability of our overall co-  
 505 residence detection mechanism.

506 *4.1.2 Sampling frequency.* Another challenge for our protocol is  
 507 determining an adequate sampling frequency. Ideally, a memory  
 508 locking receiver would loop locking operations and determine con-  
 509 tention in real-time by identifying a decrease in the moving aver-  
 510 age of the number of operations. Note that, in this case, there is es-  
 511 sentially no difference between the sender and receiver (i.e., both  
 512 continually issue locking operations) except that the receiver is tak-  
 513 ing measurements. This is adequate when there is a single sender  
 514 and receiver [22], but when there are multiple receivers, the mere  
 515 act of sensing the channel by one receiver causes contention and  
 516 other receivers cannot differentiate between a silent (0-bit) and a  
 517 locking (1-bit) sender. To avoid this, we space the sampling of mem-  
 518 ory bus such that no two receivers would sample the bus at the  
 519 same time, with high probability. We achieve this by using large  
 520 intervals between successive samples and a poisson-sampling to

**Algorithm 2** Reading a bit in the receiver

---

1: <i>now</i> $\leftarrow$ <i>time.now()</i>	523
2: <i>end</i> $\leftarrow$ <i>now</i> + <i>sampling_duration</i>	524
3: <i>sampling_rate</i> $\leftarrow$ <i>num_samples</i> / <i>sampling_duration</i>	525
4: <i>address</i> $\leftarrow$ <i>cache_line_boundary</i> - 2	526
5: <i>samples</i> $\leftarrow$ {}	527
6: <b>while</b> <i>now</i> < <i>end</i> <b>do</b>	528
7: <i>before</i> $\leftarrow$ <i>RDTSC()</i>	529
8: $\text{__ATOMIC\_FETCH\_ADD}(\text{address})$	530
9: <i>after</i> $\leftarrow$ <i>RDTSC()</i>	531
10: <i>samples</i> $\leftarrow$ <i>samples</i> $\cup \{(after - before)\}$	532
11: <b>wait until</b> <i>NEXT_POISSON(sampling_rate)</i>	533
12: <i>now</i> $\leftarrow$ <i>time.now()</i>	534
13: <b>end while</b>	535
14: <i>ks_val</i> $\leftarrow$ <i>KOLMOGOROV_SMIRINOV(samples, baseline)</i>	536
15: <b>return</b> <i>ks_val</i> < <i>ksvalue_threshold</i>	537

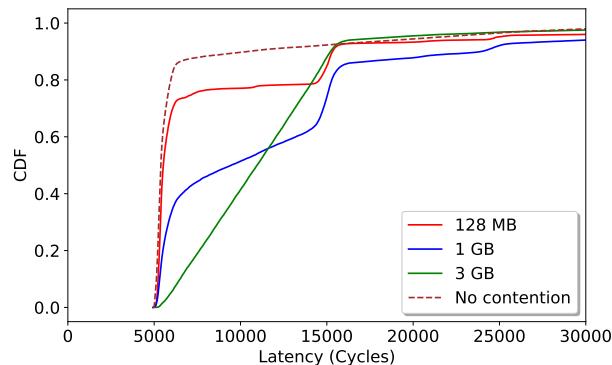
---

prevent time-locking of receivers. We determined that a millisecond-  
 542 one poisson gap between samples is reasonable to minimize noise  
 543 due to collisions in receiver sampling 1, assuming ten co-resided  
 544 receivers and a few microsecond sampling time.

545 *4.1.3 Sample Size.* In addition to adequate sampling frequency,  
 546 we must also determine sample size. A receiver can confirm con-  
 547 tention with high confidence with only a few samples, assuming  
 548 that the sender is actively causing contention on the memory bus  
 549 and the receiver is constantly sampling the memory bus through-  
 550 out the sampling duration. However, in practice, the time-sharing  
 551 of processors produces difficulties. The sender is not continually  
 552 causing contention, and neither is the receiver sensing it, as they  
 553 are context-switched by the scheduler to run other processes. As-  
 554 suming that the sender and receiver are running on different cores,  
 555 the amount of time they are actively communicating depends on  
 556 the proportion of time they are allocated on each core and how  
 557 they are scheduled.

558 To illustrate such behavior, we run a sender-receiver pair using  
 559 lambdas[14] of various sizes on AWS, and compare the distribu-  
 560 tion of latencies seen by the receiver during the contention in each  
 561 case. Figure 2 shows that the much smaller 128 MB lambdas (which  
 562 probably share a CPU core and are thus context-switched) exhibit  
 563 less active communication than the bigger 3 GB lambdas (which  
 564 may run on dedicated cores). This means that smaller instances  
 565 that tend to share processor cores with many other instances may  
 566 need to pause for more time and collect more samples to make up  
 567 for lost communication due to scheduling.

568 *4.1.4 Overcoming noise.* Along with context switching and sens-  
 569 ing noise, there are other imperfections in the measurement appa-  
 570 ratus that may cause noise. For example, we use the difference in  
 571 readings from the timestamp counter of the processor (RDTSC) be-  
 572 fore and after the locking operation to measure the latency of the  
 573 operation in cycles. If the receiver process is context-switched in  
 574 between the timer readings (e.g., at line 8 in Algorithm 2), the lat-  
 575 ency measured from their difference will be orders of magnitude  
 576 higher as it includes the waiting time of the receiver process in the  
 577 scheduler queue - which we believe is what contributes to the long  
 578

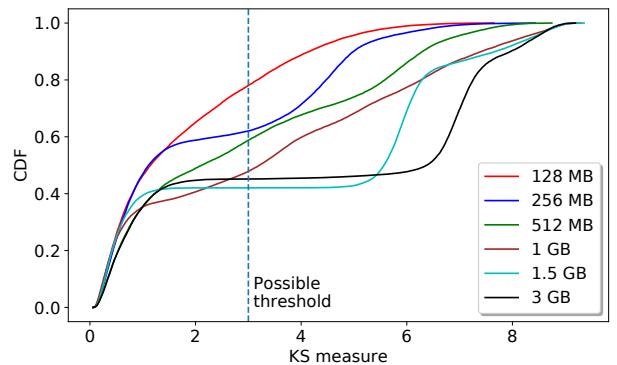


**Figure 2:** Shows CDF of latencies observed by 128 MB, 1 GB and 3 GB Lambdas during contention. The 128 MB lambda pair sees less contention due to more context switching, whereas the 1 GB and 3 GB lambdas see progressively more contention compared to baseline which we attribute to their relative stability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

tail in Figure 2. 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 Kolmogorov-Smirnov (KS) test that takes the maximum of absolute difference between empirical CDFs of the samples as measure of comparison (In our variant, we take the *mean* of absolute difference instead of the maximum). Using this measure, we can categorize a KS-value above a certain threshold as a 1-bit (contention) and a value below the threshold as 0-bit (baseline).

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

We present the pseudo-code of a receiver lambda in Algorithm 2, which includes all the challenges and subsequent solutions discussed thus far.



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

**4.1.5 Clock synchronization.** Since communicating each bit of information takes time (i.e., receiver sampling duration), our algorithm requires synchronizing sender and receiver at the start of each bit. In traditional analog channels, this is achieved either using a separate clock signal or a self-clocking signal encoding. For example, prior work [?] uses differential Manchester encoding for clock synchronization for the memory bus covert channel. Using self-clocking encodings becomes more challenging when there are multiple senders and receivers. In this work, we use the system clock for synchronizing communication. All the instances involved in the communication are executing on the same physical server and share the server’s clock. On AWS, for example, we observe that the system clock on lambdas is precise up to nanoseconds. While we have no evidence of it, we assume that the clocks between different lambdas drift only in the order of microseconds. This gives us enough space for synchronization mismatch as we take measurements only on millisecond scales due to noise constraints. We believe this to be a reasonable assumption as we did not observe any synchronization-related noise in our results.

**4.1.6 Handling Collisions.** In the preceding section, we discussed using a communication channel with synchronized time slots. In each time slot, an instance can reliably send (broadcast) or receive (listen) a bit by causing or sensing for contention. Given that there are multiple instances that may want to broadcast information on the channel, we next must determine which instance broadcasts first, to avoid collisions. Traditional channels like Ethernet or Wireless detect and avoid collisions by employing a random exponential backoff mechanism. Such a mechanism will be challenging to implement in case of our channel for two reasons. First, lambda instances do not have the capability of sensing the channel while sending a bit, which is required for detecting collisions; instances can either cause contention or sense it, but not both. Note that senders do experience a higher latency for locking operations when other senders are simultaneously causing contention. However, reliably judging this higher latency requires each sender to already have calculated a baseline of latencies without collisions, which

**Algorithm 3** ID exchange protocol TODO ► Improve pseudo-code ◀

```

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

```

takes us back to the original problem of avoiding collisions. Second, even implementing a random or exponential backoff mechanism will introduce significant overhead before any meaningful communication occurs. This overhead will also increase as the number of instances involved increases. Since each time slot could take up to 1 second, this additional overhead can make the whole communication very slow. We address this by designing a specialized communication protocol that allows for collisions but pays a price in terms of expressiveness, as we will see in the next section.

## 4.2 Protocol

Ariana: ► is it worth “recapping” and listing out all the challenges we faced, to really drive the point home to the reader?◀

When considering a communication channel for lambda co-residence detection, we note that the channel need not be general and expressive as lambdas only need to communicate their IDs with one another. Thus, we assume that each instance involved has a unique fixed-length (say  $n$ ) bit-string corresponding to its ID that must be communicated. As such, we propose a communication protocol that exchanges these bit-strings while allowing for collisions. We divide the 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 instances, which in the first phase would be all of

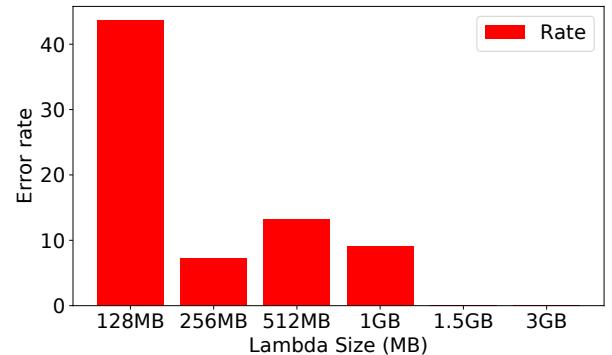


Figure 4: Shows the error rate (as a fraction of 1000 lambdas deployed) for different lambda sizes in AWS Middle-East region.

the co-resided instances. In each bit-slot  $k$  of  $n$  slots in a phase, every participating instance broadcasts a bit if the  $k^{th}$  bit of its bit-string (ID) is 1, otherwise it listens for a 0 or 1. If an instance senses a 1 while listening, it stops participating, and listens for the rest of the phase. Thus, only the instances 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 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-resided 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 Time 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 operation. For example, assuming 10,000 deployed lambdas and a maximum of 10 co-resided instances on each server, the entire co-residence detection requires around four 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 phase, all the co-resided instances would know one of their neighbors (as each phase reveals the ID of the biggest participating instance to others). If we use IDs that are globally unique, all (and only) the co-resided instances see the same ID. The instances 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-resided instances ( $K$ ) and decreases the complexity to  $O(\log_2 N)$ , allowing the entire protocol to finish within a minute instead of four.

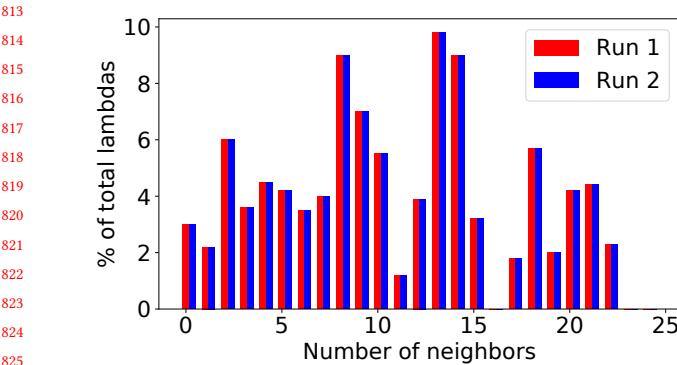


Figure 5: Shows the fraction of lambdas by the number of neighbors they saw for two independent runs that use same set of underlying AWS containers. The perfect correlation shows that both runs depict the colocation status of those containers regardless of the lambdas that ran on them, providing an evidence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

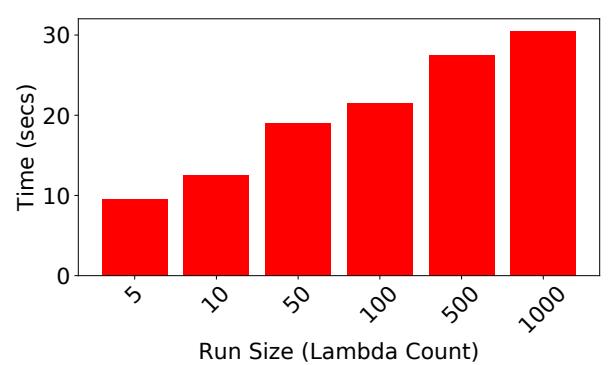


Figure 6: Shows the average execution time of the lambdas for co-location runs with varying number of lambdas. The execution time increases logarithmically with the number of lambdas demonstrating the scalability of co-residence detection with our technique.

## 5 EVALUATION

In this section, we evaluate the effectiveness of our co-residence detection technique with respect to reliability and scalability, the desirable detection properties mentioned in section 3. We run all of our experiments with AWS lambdas [3]. Though we decide to focus on one of the cloud providers, we have previously shown in section 3 that this covert channel exists on the other clouds, and thus these experiments can be replicated on their serverless functions as well. We use C++ runtime in AWS lambdas as it allows pointer arithmetic that is required to access the covert channel.

### 5.1 Setup

For each experiment, we deploy a series of instances from an AWS lambda account. Once deployed, each instance participates in the first phase of the protocol as noted in section 4.2.1, thereby learning the largest ID of their neighbors. 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 result in different IDs, we classify this instance as erroneous and report it in the error rate. We group all the instances that saw the same ID as successful and neighbors. We repeat the experiments for different lambda sizes and in various cloud regions.

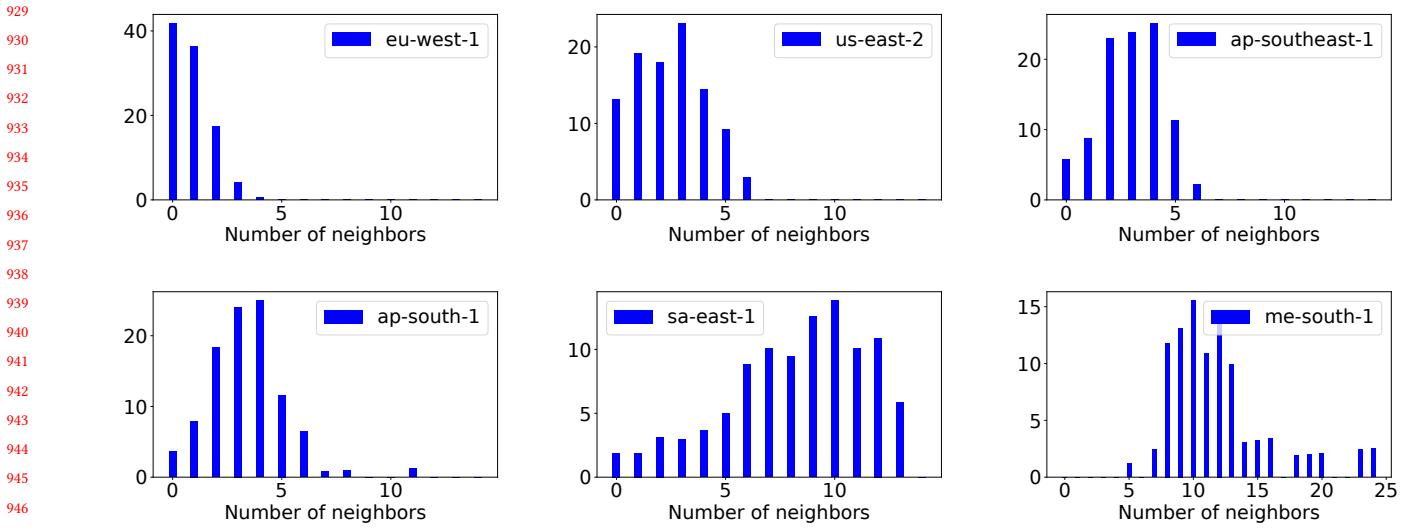
### 5.2 Reliability

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** To determine correctness, we require ground truth on which instances are co-located 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 lambda runs in a dedicated container (sandbox). After execution, AWS caches these containers 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 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 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-location is essentially co-location of their containers, and given that containers persist across experiments that are executed within minutes of one another, lambda co-location results must agree with the co-location 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-located 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-located 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-located 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



**Figure 7:** Shows co-residence results from 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-residence vary widely across regions, perhaps based on the size of those regions and the lambda activity within them.

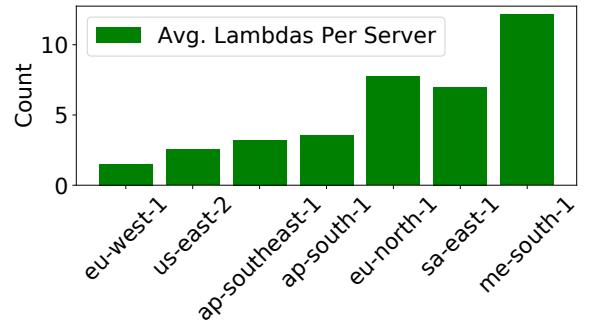
both experiments saw the same number of co-residing groups of different sizes. This proves the correctness of the results of our mechanism.

### 5.3 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.

## 6 MEASUREMENT STUDY

In this section, we present a variety of measurements on serverless function density on AWS using our co-residence detector, and the factors that may affect this density. As we discussed earlier, the key challenge in enabling traditional covert channels on the cloud is getting the sender and receiver on the same machine. So we attempt to answer the following question: Assuming that the attacker launches a number of the (sender and receiver) lambdas at a point in time in the expectation that a set of them would end up co-resident, what is the expected number of such co-resident pairs that she might see? We deploy a large number of lambdas on various AWS regions and observe the co-residence density – that is, the average number of lambdas that end up co-resident on each server. The higher the co-residence density, the easier it is



**Figure 8:** 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.

for the attacker to ultimately establish covert channels with lambdas. Unless specified otherwise, all the experiments are performed with 1.5 GB lambdas and executed successfully with **zero error** in co-residence detection.

### 6.1 Across AWS regions

We ran our co-residence detector with 1000 1.5 GB Lambdas in various AWS regions. Figure 7 comprises of multiple plots depicting the co-resided groups per region, with each bar indicating the fraction of lambdas that detected a certain number of neighbors (i.e., that belong to a co-resided group of a certain size). Plots that

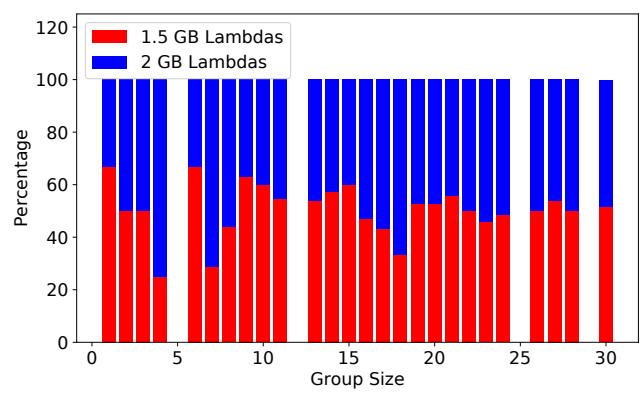
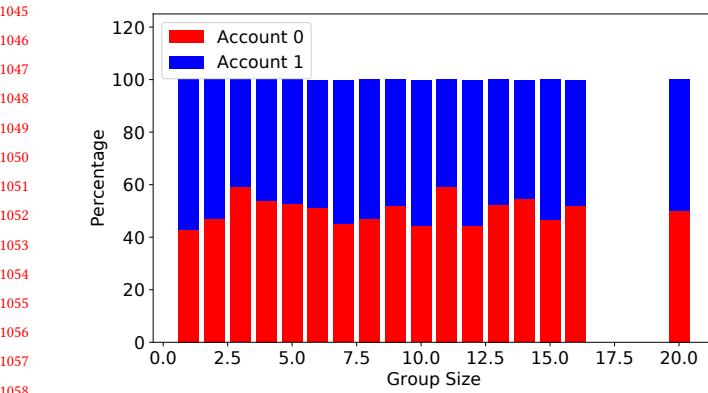


Figure 9: The left plot shows the breakdown of co-resided groups (of varying sizes) of lambdas by two different accounts in an experiment of 1000 lambdas where 500 lambdas are launched from each account. The uniformity of the split shows that lambda scheduler might be invariant to the account the lambdas are launched from. Similar results are shown for different lambda sizes in the right plot.

skew to the right indicate a higher co-residence density when compared to the plots skewed to the left (also illustrated in Figure 8). We note that, in most regions, almost all lambdas recognize at least one neighbor (indicated by smaller or non-existent first bar in each plot). We hypothesize that the co-residence density is (inversely) 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, hence the higher co-residence density in those regions as we can see in Figure 8. The ample co-residence in general across all the regions shows that lambdas provide a fertile ground for covert channel attacks.

## 6.2 Other factors

We examine how the co-residence is affected by various launch strategies that the attacker may use like deploying lambdas from multiple AWS accounts and different lambda sizes. For example, does it matter if the attacker uses one account for sender lambdas and another for the receivers (which is normally the case with covert channels)? Or one lambda size for senders and another for receivers? Do they still get colocated on the same server? To answer these questions, we run an experiment with 1000 lambdas of which we launch 500 lambdas from one account (senders) and 500 from other deployed in a random order. The co-residence observed was comparable to the case where all the lambdas were launched from one account. In the left subfigure of Figure 9, we show the breakdown of co-resided group of lambdas of each size among the two accounts. We can see that among the co-resided groups of all sizes, roughly half of lambdas came from either account. This show that lambda scheduler is agnostic to the accounts the lambdas were launched from. We see similar results for different lambda sizes, as shown in the right subfigure of figure 9

From our limited experiments, we also observe that co-residence density in a region barely changes during course of the day or the week (data not shown in any figure). This gives the attacker the freedom to attempt the attack at any time and expect similar results.

## 7 DISCUSSION

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

- While our co-residence detector does not directly help attackers locate their victims in the cloud, it can aid them in performing devastating DDOS attacks once by concentrating a number of attack instances on the victim machine. Also, the attacker could try to gain a wider surface area for targeted attacks in a cost-effective way by turning on/off her co-resided instances as necessary.
- Previous studies on performance aspects of lambdas (like performance isolation) [24] generally need a way to co-locate some lambdas. As software-level logical channels begin to disappear, our tool might provide a reliable alternative?
- Burst parallel frameworks like gg [8] that orchestrate lambdas can use our co-residence detector as a locality indicator to take advantage of server locality.

**Mitigation** In previous section, we showed that our co-residence detector makes the covert channels practical with lambdas, so it is important that clouds address this issue. One way to disable our co-residence detector is to fix the underlying memory bus channel that it employs. However, this only works for newer generation of servers and is not practical for existing infrastructure. An easier solution, one that is only practical with lambdas, is to disable the lambda support for low-level languages (or unsafe versions of high-level languages) by the cloud providers. This will prevent pointer arithmetic that is required to activate this channel. If that is not an option, clouds may look at more expensive solutions like BusMonitor [20] that isolate memory bus usage for different tenants. We leave such exploration to future work.

## 1161 8 RELATED WORK

1162 **Cloud Attacks** Co-residency is possible because of covert channels, so we begin our related work with an investigation into cloud  
 1163 attacks. Initial papers in co-residency detection utilized host information and network addresses arising due to imperfect virtualization [19]. However, these channels are now obsolete, as cloud provides have strengthened virtualization and introduced Virtual Private  
 1164 Clouds [4]. Later work used cache-based channels in various levels of the cache [12, 15, 28, 35] and hardware based channels like thermal covert channels [17], RNG module [7] and memory bus [26] have also been explored in the recent past. Moreover, studies have found that VM performance can be significantly degraded using memory DDoS attacks [31], while containers are susceptible to power attacks from adjacent containers [9].

1165 Our work focuses on using the memory bus as a covert channel  
 1166 for determining cooperative co-residency. Covert channels using  
 1167 memory bus were first introduced by Wu et. al [26], and subsequently has been used for co-residency detection on VMs and  
 1168 containers [21? ]. Wu et. al [26] introduced a new technique to lock  
 1169 the memory bus by using atomic memory operations on addresses  
 1170 that fall on multiple cache lines, a technique we rely on in our own  
 1171 work.

1172 **Co-residency** One of the first pieces of literature in detecting VM  
 1173 co-residency was introduced by Ristenpart et al., who demonstrated  
 1174 that VM co-residency detection was possible and that these techniques  
 1175 could be used to gather information about the victim machine (such as keystrokes and network usage) [19]. This initial  
 1176 work was further expanded in subsequent years to examine co-  
 1177 residency using memory bus locking [30] and active traffic analysis  
 1178 [6], as well as determining placement vulnerabilities in multi-  
 1179 tenant Public Platform-as-a-Service systems [22, 33]. Finally, Zhang  
 1180 et al. demonstrated a technique to detect VM co-residency detec-  
 1181 tion via side-channel analyses [34]. Our work expands on these  
 1182 previous works by investigating co-residency for lambdas.

1183 **Lambdas** While lambdas are a much newer technology than VMs,  
 1184 there still exists literature on the subject. Recent studies examined  
 1185 cost comparisons of running web applications in the cloud on lambdas  
 1186 versus other architectures [23], and also examined the lambdas  
 1187 have been studied in the context of cost-effectiveness of batching  
 1188 and data processing with lambdas [13]. Further research has shown  
 1189 how lambdas perform with scalability and hardware isolation, indi-  
 1190 cating some flaws in the lambda architecture [24]. From a security  
 1191 perspective, Izhikevich et. al examined lambda co-residency using  
 1192 RNG and memory bus techniques (similar to techniques utilized in  
 1193 VM co-residency) [11]. However, our work differs from this study  
 1194 in that our technique informs the user of which lambdas are on the  
 1195 same machine, not only that the lambdas experience co-residency.

## 1211 9 ETHICAL CONSIDERATIONS

1212 As with any large scale measurement project, there are ethical con-  
 1213 siderations to take into account. First, there are security and pri-  
 1214 vacy concerns of using this technique to uncover other consumer's  
 1215 lambdas. However, since we focus on co-operative co-residence de-  
 1216 tection, we only determine co-residence for the lambdas we launched,

1217 and do not gain insight into other consumer's lambdas. Second,  
 1218 there is concern that our experiments may cause performance is-  
 1219 sues with other lambdas, as we may block their access to the mem-  
 1220 ory bus. We believe this concern is small, for a number of rea-  
 1221 sons. Memory accesses are infrequent due to the multiple levels  
 1222 of caches; we would only be affecting a small number of opera-  
 1223 tions. Memory accesses and locking operations are FIFO, which  
 1224 prevents starvation of any one of the lambdas sharing a machine.  
 1225 Moreover, lambdas are generally not recommended for latency-  
 1226 sensitive workloads, due to their cold-start latencies. Thus, the small  
 1227 amount of lambdas that we might affect should not, in practice, be  
 1228 affected in their longterm computational goals.

## 1229 10 CONCLUSION & FUTURE WORK

1230 **TODO** ▶ *Copied abstract*◀ Cloud computing has seen explosive growth  
 1231 in the past decade. This is made possible by efficient sharing of  
 1232 infrastructure among tenants, which unfortunately also raises se-  
 1233 curity challenges like preventing side-channel attacks. Providers,  
 1234 like AWS and Azure, have traditionally relied on hiding the co-  
 1235 residency information to prevent targeted attacks in their clouds.  
 1236 But recent works have repeatedly found co-residence detection  
 1237 techniques that break this encapsulation, prompting the providers  
 1238 to address them and harden isolation on their platforms. In this  
 1239 work, we find yet another such technique based on a memory bus  
 1240 covert channel that is more pervasive, reliable and harder to fix.  
 1241 We show that we can use this technique to reliably perform co-  
 1242 operative co-residence detection for thousands of AWS lambdas  
 1243 within a few seconds, which opens a way for attackers to perform  
 1244 DDoS attacks or learn cloud's internal mechanisms. We present  
 1245 this technique in detail, evaluate it and use it to perform a small  
 1246 study on lambda activity across a few AWS regions. Through this  
 1247 work, we hope to motivate the need to address this covert channel  
 1248 in the cloud.

## 1249 REFERENCES

- [1] Alexandru Agache, Marc Brooker, Alexandra Iordache, Anthony Liguori, 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>
- [7] Dmitry Etyushkin 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] Sajjad Fouladi, Francisco Romero, Dan Iter, Qian Li, Shuvo Chatterjee, Christos Kozyrakis, Matei Zaharia, and Keith Winstein. 2019. From Laptop to Lambda: Outsourcing Everyday Jobs to Thousands of Transient Functional Containers. In *2019 USENIX Annual Technical Conference (USENIX ATC 19)*. USENIX Association, Renton, WA, 475–488. <https://www.usenix.org/conference/atc19/presentation/fouladi>
- [9] 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*

1254  
 1255  
 1256  
 1257  
 1258  
 1259  
 1260  
 1261  
 1262  
 1263  
 1264  
 1265  
 1266  
 1267  
 1268  
 1269  
 1270  
 1271  
 1272  
 1273  
 1274  
 1275  
 1276

- 1277      Annual IEEE/IFIP International Conference on. IEEE, IEEE, 237–248.  
 1278 [10] gcp 2019. Google Cloud Functions. <https://cloud.google.com/functions/>.  
 1279 [11] Elizabeth Izhikevich. 2018. *Building and Breaking Burst-Parallel Systems*. Master's thesis. University of California, San Diego.  
 1280 [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>  
 1281 [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.  
 1282 [14] Lambda 2019. AWS Lambda. <https://aws.amazon.com/lambda/>.  
 1283 [15] 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>  
 1284 [16] Ramya Jayaram Masti, Devendra Rai, Aanjan Ranganathan, Christian Müller, Lothar Thiele, and Srdjan Capkun. 2015. Thermal Covert Channels on Multi-core 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>  
 1285 [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>  
 1286 [18] Clementine Maurice, Manuel Weber, Michael Schwarz, Lukas Giner, Daniel Gruss, Carlo Alberto Boano, Stefan Mangard, and Kay Römer. 2017. Hello from the Other Side: SSH over Robust Cache Covert Channels in the Cloud. <https://doi.org/10.14722/ndss.2017.23294>  
 1287 [19] 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>  
 1288 [20] Brendan Saltaformaggio, D. Xu, and X. Zhang. 2013. BusMonitor : A Hypervisor-Based Solution for Memory Bus Covert Channels.  
 1289 [21] 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>  
 1290 [22] 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>  
 1291 [23] Mario Villamizar, Oscar Garcés, 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>  
 1292 [24] 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>  
 1293 [25] Zhenghong Wang and Ruby B. Lee. 2006. Covert and Side Channels Due to Processor Architecture. In *Proceedings of the 22nd Annual Computer Security Applications Conference (ACSAC '06)*. IEEE Computer Society, USA, 473–482. <https://doi.org/10.1109/ACSAC.2006.20>  
 1294 [26] Zhenyu Wu, Zhang Xu, and Haining Wang. 2012. Whispers in the Hyper-space: 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>  
 1295 [27] 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>  
 1296 [28] 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>  
 1297 [29] 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>  
 1298 [30] 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>  
 1299 [31] Tianwei Zhang, Yinqian Zhang, and Ruby Lee. 2016. Memory DoS Attacks in Multi-tenant Clouds: Severity and Mitigation.  
 1300 [32] 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.  
 1301 [33] 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](https://doi.org/10.1007/978-3-319-50011-9_28)  
 1302 [34] 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>  
 1303 [35] 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>

1335  
 1336  
 1337  
 1338  
 1339  
 1340  
 1341  
 1342  
 1343  
 1344  
 1345  
 1346  
 1347  
 1348  
 1349  
 1350  
 1351  
 1352  
 1353  
 1354  
 1355  
 1356  
 1357  
 1358  
 1359  
 1360  
 1361  
 1362  
 1363  
 1364  
 1365  
 1366  
 1367  
 1368  
 1369  
 1370  
 1371  
 1372  
 1373  
 1374  
 1375  
 1376  
 1377  
 1378  
 1379  
 1380  
 1381  
 1382  
 1383  
 1384  
 1385  
 1386  
 1387  
 1388  
 1389  
 1390  
 1391  
 1392