

## Journal Pre-proof

KASLR-MT: Kernel address space layout randomization for multi-tenant cloud systems

Fernando Vano-Garcia, Hector Marco-Gisbert



PII: S0743-7315(19)30409-5  
DOI: <https://doi.org/10.1016/j.jpdc.2019.11.008>  
Reference: YJPDC 4153

To appear in: *J. Parallel Distrib. Comput.*

Received date: 27 May 2019  
Revised date: 11 September 2019  
Accepted date: 8 November 2019

Please cite this article as: F. Vano-Garcia and H. Marco-Gisbert, KASLR-MT: Kernel address space layout randomization for multi-tenant cloud systems, *Journal of Parallel and Distributed Computing* (2019), doi: <https://doi.org/10.1016/j.jpdc.2019.11.008>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2019 Elsevier Inc. All rights reserved.

# KASLR-MT: Kernel Address Space Layout Randomization for Multi-Tenant Cloud Systems

Fernando Vano-Garcia, Hector Marco-Gisbert

School of Computing, Engineering and Physical Sciences, University of the West of Scotland, High St, Paisley PA1 2BE, UK

## ARTICLE INFO

### Keywords:

Cloud  
Virtualization  
Security  
KASLR  
Memory Deduplication  
Memory Management  
Operating Systems

## ABSTRACT

Cloud computing has completely changed our lives. This technology dramatically impacted on how we play, work and live. It has been widely adopted in many sectors mainly because it reduces the cost of performing tasks in a flexible, scalable and reliable way. To provide a secure cloud computing architecture, the highest possible level of protection must be applied. Unfortunately, the cloud computing paradigm introduces new scenarios where security protection techniques are weakened or disabled to obtain a better performance and resources exploitation. Kernel ASLR (KASLR) is a widely adopted protection technique present in all modern operating systems. KASLR is a very effective technique that thwarts unknown attacks but unfortunately its randomness have a significant impact on memory deduplication savings. Both techniques are very desired by the industry, the first one because of the high level of security that it provides and the latter to obtain better performance and resources exploitation. In this paper, we propose KASLR-MT, a new Linux kernel randomization approach compatible with memory deduplication. We identify why the most widely and effective technique used to mitigate attacks at kernel level, KASLR, fails to provide protection and shareability at the same time. We analyze the current Linux kernel randomization and how it affects to the shared memory of each kernel region. Then, based on the analysis, we propose KASLR-MT, the first effective and practical Kernel ASLR memory protection that maximizes the memory deduplication savings rate while providing a strong security. Our tests reveal that KASLR-MT is not intrusive, very scalable and provides strong protection without sacrificing the shareability.

## 1. Introduction

Cloud computing has become a significant aspect of our lives. It allows a provider to share pools of configurable resources (hardware/software) through virtualization, yielding new complex business models that were unpredictable some years ago. Cloud computing has been widely adopted in many sectors, mainly because it reduces the cost of performing tasks in a flexible, scalable and reliable way. From the user's point of view, they can benefit from vast computing power and storage without the need to possess the necessary hardware resources.

*Infrastructure as a Service* (IaaS) [28] is considered one of the fundamental building blocks for cloud services because at this level, a client is able to configure virtualized environments with high flexibility without having to concern about deploying large rooms of physical computers. Thence the service provider supplies the storage, networking and virtualization so that the client has full control over the system from OS layer upwards. Efficient resource management is fundamental to deal with a proper cloud infrastructure [43, 51, 20]. Hardware resources are a critical asset in the business, and it must be managed and utilized adequately. A cloud service provider will obtain more benefits if he/she is able to operate more virtual machines with the same resources [37, 48].

Given the significance that cloud computing has in people's lives, it is imperative to offer confidentiality, integrity, and availability in any cloud computing architecture. Furthermore, there is a stack of services relying on IaaS, for example *Platform as a Service* (PaaS), *Software as a Service* (SaaS) or *Function as a Service* (FaaS) to name a few. Any security issue affecting the base will affect the upper layers as well. For that reason, IaaS service providers must ensure to reach the highest possible level of security in order to guarantee a suitable quality to their clients.

Since the finding of the first computer bug back in 1947 [3], the sphere has been changing swiftly. Attackers are aware of the fact that computers are the building blocks of our society. For this reason, bugs with security implications are being abused in order to make profit of those vulnerabilities. Given the asymmetric nature of the confrontation between attackers and defenders, the former enjoy their tactical advantage, while the latter design and develop defense mechanisms to prevent the successful exploitation of the most complex attacks. In recent years, there has been a transition in the *battlefield* from userland to kernel, since the complexity of userland exploits has overtaken the kernel ones [36]. Furthermore, a successful exploitation of a kernel vulnerability is much more dainty for an attacker.

Although cloud service providers desire to yield as much security as possible in their infrastructure, it is not always possible. Current security protection mechanisms are far from perfect and, in some cases, they introduce prohibitive overheads. Kernel randomization is a widely adopted security mechanism that introduces a prohibitive overhead [44] to the memory savings of the memory deduplication. Since

✉ Fernando.Vano-Garcia@uws.ac.uk (F. Vano-Garcia);

Hector.Marco@uws.ac.uk (H. Marco-Gisbert)

www.fervagar.com (F. Vano-Garcia); www.hmarco.org (H.

Marco-Gisbert)

ORCID(s): 0000-0001-6158-8694 (F. Vano-Garcia);

0000-0001-6976-5763 (H. Marco-Gisbert)

it is not recommended to disable the kernel randomization mechanism, cloud providers will sustain a forfeit of memory resources.

The rest of the paper is organized as follows. Section 2 provides a detailed background on memory deduplication and kernel randomization. Section 3 briefly describes how kernel randomization is performed in the Linux kernel v5.0.6 for the x86\_64 architecture. Section 4 points out the reasons why the kernel randomization penalizes the memory deduplication. In section 5 the Linux kernel randomized zones are identified and a comprehensive analysis on the deduplication effect of randomizing them is presented. Based on this analysis, section 6 describes the proposed solution, while the implementation of the *KASLR-MT* technique in Linux is presented in section 7. Section 8 provides an evaluation of the proposed implementation. The paper finishes in section 9 with a discussion on Linux kernel modules and how position-independent code could alleviate the issue, followed by some conclusions and future work in section 10.

## 2. Background

In this section we explain the memory deduplication technique used to save physical memory. Then, we explain the kernel randomization concept and the two main approaches followed to randomize code. Finally the attacker model is presented.

### 2.1. Memory Deduplication

Memory deduplication is a physical memory saving technique. Given the noteworthy importance of efficient memory resource utilization on behalf of cloud computing providers, deduplication is a desired feature. It is able to reduce the memory footprint across virtual machines [1, 16], decreasing the total cost of managing and ownership.

Although in the first instance deduplication was designed to be used in hypervisors [40, 1, 29, 47], it was gently applied for memory contents of non-virtualized environments as well. Then, it was widely adopted by most of the operating systems. For example the Linux kernel included Kernel Samepage Merging (KSM) in the version 2.6.32 and Windows introduced Memory Page Combining in Windows 8 [41]. Furthermore, data deduplication is commonly used in other areas such as databases or web contents [8, 31, 46].

When used with virtualization technologies, memory deduplication is usually operating at the hypervisor layer, along with the memory manager of the physical host machine. In almost all modern operating systems, memory is organized in pages [6]. Typically, each memory page consumes 4096 bytes of physical memory. This strategy facilitates an efficient management of memory resources and enables virtual memory, which is a fundamental feature for virtualization. Thanks to this memory organization, memory deduplication can merge all pages with identical content into only one. Note that swapped pages can not be deduplicated but only the ones that reside in physical memory.

Different architectures support different page sizes [49]. A greater page size implies less page table entries and less

TLB faults, resulting in higher performance [30]. However, this reduces the chances of finding two pages with equal content, which reduces the memory saving rate [15]. Thanks to the memory management virtualization support, it is possible for hypervisors to implement a memory deduplication using pages of 4096 bytes independently of the page size guest view.

In virtualized environments, deduplication is commonly applied to the entire guest memory region corresponding to the virtual machine (often called guest physical memory) [2]. Hence, those pages belonging to that memory region will be candidates for being shared across all virtual machines. On the one hand, this increases the chances of finding matching pages to share. On the other hand, it enables different types of side-channels [17, 42, 23] that might compromise the confidentiality. Different solutions have been researched [45, 32, 21, 35, 21] and, depending on the adversary model of a cloud provider, memory deduplication can be used without the need to sacrifice security.

### 2.2. Address Space Layout Randomization

Address Space Layout Randomization (ASLR) [34] is a security technique that consists in placing the memory regions of a program in random locations. The objective of this technique is to hinder the successful exploitation of vulnerabilities that rely on the knowing program addresses [39, 24].

As a consequence, an attacker must obtain an addresses where code or data is located in memory in order to trigger a malicious payload [7, 12, 25]. This in part is because the *Data Execution Prevention / No-Execute* (DEP/NX) prevents to execute injected code.

ASLR embraces some requirements in order to be applied correctly. It needs a high-quality entropy source for generating cryptographically secure random numbers to determine the addresses where the program will be loaded [18, 26]. If this requirement is not fulfilled, the predictability of the generated addresses will be high and then it will be easier for an attacker to guess correctly a valid address [10, 11].

Kernel ASLR (KASLR) or *kernel randomization* is the application of this technique to the kernel [4]. Locations of kernel memory regions are determined at boot time and they are not changed until next shutdown/reboot. Each implementation has its particularities, but code and data regions are commonly randomized. It is currently being used by all the main operating systems: Apple introduced it into Mac OS X Mountain Lion (Mac OS X 10.8) [50], Microsoft introduced it in Windows Vista and Linux in the kernel version 3.14.

Since the randomization is only applied at boot time, any *information leak* revealing a kernel address will derandomize the kernel, and the bypass will be effective until next reboot [13, 19].

### 2.3. Randomizing Code

In order to randomize code, the loader must be able to choose an arbitrary address and then load the code at that address. To achieve this, the code being randomized must

be compiled and linked to allow loaders to put them at random memory locations [38]. There are mainly two different approaches to enable code to be randomized: Position-Independent Code (PIC) and relocations. Both approaches can load and run code at random virtual addresses, but they have implementation differences that have an impact in the performance and also in the shared memory.

Position-independent code is a piece of machine code that can be executed regardless of where it is loaded in memory without any code modification. This feature is crucial for shared libraries, to keep the code segment as non-writable and to allow memory sharing by several processes using the Copy-On-Write (COW) mechanism [9]. Instead of referencing symbols with absolute addresses, PIC uses indexed addressing using relative offsets. For example, a data variable can be referenced relatively using the program's instruction pointer. For efficiently accessing symbols and data beyond the addressable limits of an architecture, userland programs use a look-up table called Global Offset Table (GOT). This table is placed in a data segment, private for each process, and contains absolute addresses of global exported symbols. Hence, PIC allows to randomize code memory regions without the need of altering its contents.

On the other hand, relocatable code is a piece of machine code that also can be executed regardless of where it is loaded in memory, but the executable segment needs to be patched when it is loaded [22, 33]. Relocation information is consulted by the loader in order to adjust symbol references through different parts of the program with final absolute addresses in-place. Although the final program with the absolute references can be more efficient than position-independent code, the load time work of relocatable code is considerably heavier, as every reference in code must be fixed-up. In contrast to PIC, relocatable code is not suitable for shared libraries. The relocations in code pages trigger the COW mechanism, generating private copies of these pages for each process using the same library loaded at different virtual addresses. The result is a decrement of memory sharing among processes, which is the main purpose of shared libraries.

#### 2.4. Attacker Model

Disabling the kernel randomization protection in favour of having the highest possible memory deduplication saving rate introduces serious weaknesses. Vulnerabilities that rely on knowing where the kernel has been mapped will have 100% of success, since those addresses are not longer a secret that attackers need to obtain.

In a cloud environment where virtual machines can interact with each other, this is even more risky. A kernel vulnerability could compromise the entire physical cloud and all its virtual machines, even if they are running in different tenants and physical machines. Local, remote, inter-VM, intra-VM, inter-Tenant and intra-Tenant attackers do not need to perform a prior kernel attack to know where the kernel resides in memory, they already know, and their attacks will always be successful.

Our goal is to provide kernel randomization while introducing a negligible impact in the memory savings, to enable cloud providers to use the KASLR protection along with the memory deduplication benefits. We assume that the kernel contains a software vulnerability that requires to bypass the kernel randomization protection to be successfully exploited. That is, our goal is that attackers exploiting a kernel vulnerability from local, remote, inter-VM, intra-VM, inter-Tenant and intra-Tenant must face the full kernel randomization protection with almost no effect in the memory deduplication.

### 3. Linux Kernel Randomization

This section describes how kernel randomization is performed in the Linux kernel v5.0.6 for the x86\_64 architecture. We describe how the bootloader and Linux kernel randomizes its memory.

In Section 2.3, we described two different approaches to randomize code: position-independent code and relocations. Current Linux kernel implementation is not PIC compliant. It uses text relocations, patching dynamically all the position-dependent references after the final address of the code memory region is randomly calculated.

At the early stages of the boot process, the Linux kernel is decompressed in memory by the bootloader. One of its purposes, along with other system initialization operations, is to place the regular kernel into a random location.

The random numbers used by KASLR are requested at boot time, and at this stage there is not much entropy available. Having quality random numbers is key to have an effective KASLR, otherwise the addresses will be predictable and the protection useless. In order to obtain random numbers, Linux uses different sources of entropy such as `rand`, `rdtsc`, system timers, etc. The outcome of all the available sources is combined by using the XOR operation, and the result value is diffused by using circular multiplication.

The decompressor uses two random numbers to randomize both the physical address and the virtual address where the kernel will be loaded and mapped respectively. The physical address determines where the regular kernel image is going to be decompressed in memory. This image is an Executable and Linkable Format (ELF) file with a relocation table appended to it. Relocation information is generated by the static linker and it is needed for patching the virtual addresses of every position-dependent reference. Each virtual address that needs to be updated has its corresponding entry in this relocation table.

Once the kernel is decompressed into the randomly chosen physical address, the kernel image is parsed to extract the information about the segments that compose the kernel virtual memory. All the loadable segments are placed into their corresponding location. Then, relocations are processed. The kernel needs to be patched taking into account the new base address that differs from the base address that it was linked to. Listing 1 shows the pseudo-code of the loop processing and fixing relocations.



---

```

for each relocation do
    reloc_addr = <read relocation entry>
    // Calculate its current physical location
    ptr = (reloc_addr + phys_map)
    check_memory_bounds(ptr)

    // Update virtual address
    *ptr += kaslr_virt_offset

```

---

**Listing 1:** Pseudo-code of relocations update

After the relocations have been applied by the boot-loader, the execution control is transferred to the kernel. At this point, some system initialization operations are performed, including the randomization of the remaining four memory regions. First, the physical direct mapping (physmap), the dynamic memory region (vmalloc) and the virtual memory map (vmemmap) are randomized. Later, when the first module is loaded, the modules base address is randomly calculated (modules).

Although each region is randomized independently, there is a position order. The physical direct mapping will always be placed at a lower virtual address than the vmalloc region, which also will be placed at a lower virtual address than the vmemmap region.

The randomization algorithm tries to use the available entropy in the most efficient way possible. The first region is bounded to the lower third of the entire virtual space available for these three regions. A random address within the bounds is selected to map the memory region. The remaining space, subtracting a padding to avoid region overlapping, is then divided by two in order to place the second region, following the same approach. Similarly, the last region is placed at a random address within the remaining space after subtracting the padding corresponding to the second one.

The final memory region to be randomized is the zone where modules are loaded. It uses a different logic from the previous memory regions, since it is not calculated until the first module is loaded in the system. When the first module is loaded, a random number of pages between 1 and 1024 is determined and it is added to a static base address, establishing the virtual base address where the allocation of all the loadable modules starts. From that point, subsequent modules are sequentially allocated in the order as they are loaded. Once a module is loaded into its final location, relocations need to be done to fix-up their references. These relocations include absolute addresses and relative offsets. We have found modules with references to different parts of the memory, including the kernel code and data, the physmap region, their own memory, and even the memory of other modules. It is worth noting that even relative offsets depend on address randomization if they refer to memory regions being also randomized. In that case, the relative offset may differ, since the distance between referee and referrer is not constant. In addition, the module load order is not determin-

istic so the final order, and therefore the relocations applied to modules, can be different in different kernel boots. This issue is discussed in more detail in Section 9.

To summarize, the Linux kernel randomizes a total of six different addresses:

- Kernel Base Physical Address: Physical address where the kernel is decompressed.
- Kernel Base Virtual Address: Virtual address of the kernel text mapping, containing the code and data segments.
- Physmap: Direct mapping of all physical memory. Also used to dynamically allocate physically contiguous memory.
- Vmalloc: Memory region to dynamically allocate virtually contiguous memory.
- Vmemmap: Kernel virtual memory map, containing metadata of physical page frames.
- Modules: Virtual base address where modules are loaded.

## 4. The Problem: KASLR vs Deduplication

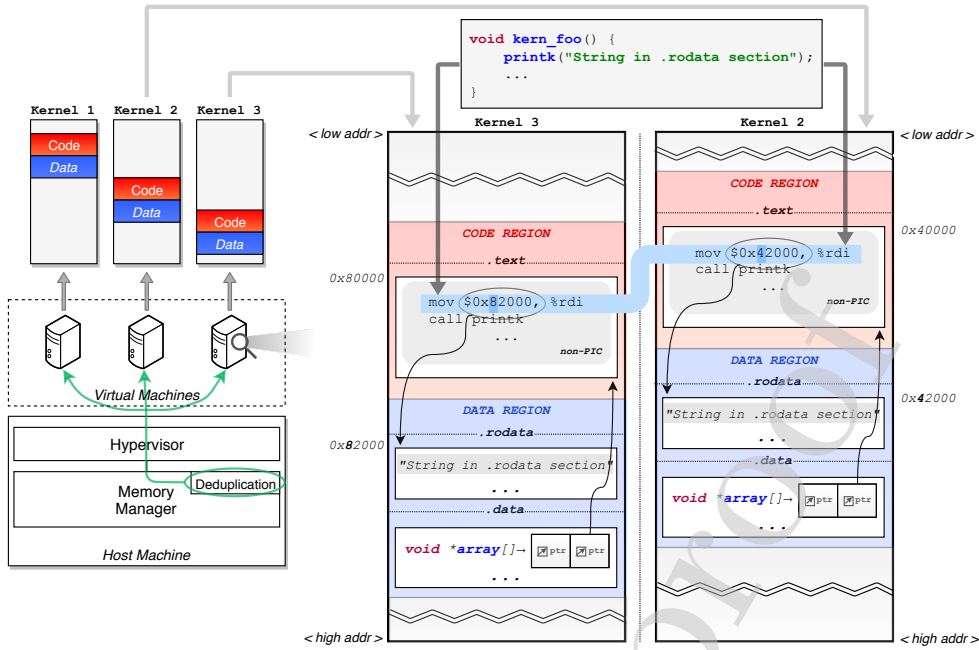
In this section, we present how kernel randomization reduces the effectiveness of memory sharing by deduplication in virtualized systems.

As detailed in section 2.1, the memory deduplication mechanism merges pages with equal content. Regardless of whether they comprise code or data and independently of their virtual address, two or more pages will be merged if their contents are identical. On the other hand, as described in section 2.3, kernel randomization could affect the host's memory sharing effectiveness if the relocation randomization approach is used. Unfortunately, as we present in section 3, we found that current Linux versions uses the relocations approach for randomizing the kernel.

The main issue is that current Linux kernel randomization is following a relocation approach that modifies code at boot time depending on random addresses, but memory deduplication requires to have equal content to merge pages. This conflict results on either having KASLR enabled but reducing the shared memory or disabling KASLR to have the maximum deduplicated memory.

### 4.1. Breaking Shared Code

Memory deduplication fails to merge kernel code because of the randomization. Figure 1 exemplifies this concept. It shows the same hypothetical `printk()` kernel call from 2 kernels that were mapped at different base addresses. Kernel 2 base address is `0x40000` and kernel 3 base address is `0x80000`. The resulting `mov $addr, %rdi` is patched by the kernel loader accordingly to the random kernel base address at boot time.



**Figure 1:** Overview scenario, with an hypervisor using memory deduplication and running three virtual machines. Kernel memories of two virtual machines are detailed, showing their load address along with some contents of the code and data regions.

In more detail, the instruction is placing the virtual address of the string into a register and passing it to the `printk()` function, following the *System V AMD64 ABI* convention [27]. The absolute reference in the `mov` instruction differs because the data region was loaded at different addresses. These absolute addresses need to be fixed-up at load time, before even running the kernel, altering the page content. As a consequence, those pages cannot be merged by deduplication because their content after the relocation will differ. This is just an example to illustrate the issue that prevents memory deduplication from merging Linux code. The Linux kernel is a modern and advanced operating system containing many coding tricks to improve both the readability and efficiency of the code, and this issue can appear in different ways.

Therefore, a full analysis is necessary, presented in section 5 to determine, on the one hand, to what extent the code is being modified because of the randomization and, on the other hand, which randomization zones produce most kernel and modules code modifications.

#### 4.2. Breaking Shared Data

Relocations do not only affect to code but also to data variables whose contents depend on the virtual address of the memory location being referenced. For example, a data pointer containing the address of a dynamically allocated object. Therefore, the problem occurs when memory contents depend on the position of the memory location being referenced and this location is randomized. A very known example in userland applications is the usage of a GOT/PLT to jump to libraries. The GOT contains pointers to where library functions reside in memory. Therefore, the fact of randomizing libraries prevents the merging of the page holding the GOT.

Similarly, the Linux kernel contains structures that hold pointers to functions. Those pointers depend on the kernel base address. As a consequence, all pages containing one single pointer referencing to a randomized address will not be merged by memory deduplication. This can be extended to any memory in Linux that holds data, such as `vmalloc`, `vmemmap` and modules data.

Returning to figure 1, we can observe an array of pointers in the data section of both kernels, where the second element points to an address of its corresponding code section. Similar to what occurs with the relocation in the `mov` instruction, given that the base address of these two kernels differ, the pointer will differ as well. As a consequence, the host cannot share these memory pages.

### 5. KASLR Influence on Deduplication

Randomizing the memory layout of guest kernels reduces the effectiveness of memory deduplication. In this section, in order to design a new compatible KASLR, we present a comprehensive analysis of the impact of each randomized area to the memory deduplication.

Our goal is to precisely measure to what extent a particular memory region (Linux code, Linux data, modules code, modules data, `vmalloc` and `vmemmap`) differs when randomizing memory base addresses (kernel physical address, kernel virtual address, physical mapping, `vmalloc`, `vmemmap` and modules). Since there are six memory base addresses to be randomized, we have tested  $2^6 = 64$  combinations per each memory region. For example, enabling randomization for all memory base addresses except one, all except two, etc. All tests were executed using Ubuntu 19.04 with Linux v5.0.6 carrying out 8 probes per each combination, resulting in a

[illegible]

**Table 1**  
Analysis results of the Linux code memory region.

total of  $64 * 8 = 512$  virtual machine executions. By comparing all the probes of each combination, we can obtain the percentage of equal pages. By doing this, we can precisely achieve the following:

1. **The KASLR influence:** of each combination to iden-

[illegible]

**Table 2**  
Analysis results of the Linux data memory region.

tify which randomized memory base address have more impact on memory deduplication.

2. **The best combination:** for our purpose. In our case, we aim to fully randomize as many memory base addresses as possible. To achieve that, we need to obtain

Randomized Base Addresses						
modules	vmemmap	vmalloc/ ioremap	physical mapping	kernel vaddr.	kernel paddr.	% equal pages
<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	30.4
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	29.4
<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	27.2
<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	26.8
<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	25.9
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	25.6
<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	25.5
<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	24.9
<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	24.7
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	24.5
<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	24.4
<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	24.3
<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	23.8
<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	23.7
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	23.7
<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	23.4
<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	12.2
<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	12.1
<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	12.1
<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	12.1
<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	12.1
<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	12.0
<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	12.0
<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	12.0
<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	12.0
<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	12.0
<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	12.0
<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	12.0
<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	12.0
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	12.0
<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	12.0
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	12.0
<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	12.0
<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	11.9
<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	11.9
<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	11.9
<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	11.9
<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	11.9
<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	11.9
<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	11.9
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	11.9
<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	11.9
<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	11.8
<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	11.8
<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	11.8
<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	11.8
<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	11.8
<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	11.8
<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	11.8
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	11.8
<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	11.8
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	11.8
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	11.8
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	11.8

**Table 3**  
Analysis results of the modules code memory region.

those that have low or zero impact on memory deduplication. As detailed in the solution section 6, those memory base addresses will be randomized as usual and we are not modifying the kernel code that randomizes them.

[illegible]

**Table 4**  
Analysis results of the modules data memory region.

3. **Not biased and Independent Results:** of the number of virtual machines. Instead of calculating the memory deduplication differences before and after randomization, we are calculating how different or equal a memory region is after applying randomiza-



[illegible]

**Table 5**  
Analysis results of the `vmalloc` memory region.

tion. The first approach could be tuned to obtain almost any number by choosing a high number of virtual machines. For example, if half of the memory of a virtual machine can be deduplicated, then if we have two identical machines, the memory deduplica-

[illegible]

**Table 6**  
Analysis results of the `vmemmap` memory region.

tion will save 50% of the memory but if we have three, the deduplication will report 66% of memory saved. Hence, although this is reflecting the actual savings, it is not appropriate to know the real impact of KASLR.

Therefore, our analysis focus on obtaining the percent-

age of pages that are equal for each case, independently of the region's size. Note that our analysis obtains separately the randomization effect produced to Linux code and Linux data regions, but it is not possible to randomize both separately. On the other hand, the fact of randomizing the physical mapping can influence one or more memory regions, but this is a special case; it is a virtual mapping to all physical memory and not actually a memory region, so its contents must be ignored in the analysis.

Even though the Linux kernel guests use pages of 2 MiB, the memory deduplication implementation in Linux works with 4 KiB pages. Therefore, our analysis calculating how equal a memory region is after applying all randomization detects differences at page level, that is, a single bit difference in a page is reported as a full 4 KiB page mismatch.

Tables 1 and 2 show the analysis of the Linux code and data respectively. Tables 3 and 4 do the same for the code and data of the modules, table 5 shows the analysis of `vmalloc` and finally table 6 presents the results of the `vmemmap` memory region. The column *% equal pages* on the tables indicates the percentage of similarity between a particular memory region across all memory base address randomization combinations.

### 5.1. Linux Code Impact

The results of the randomization effects on the Linux code shown in table 1 point out that the Linux code memory region is only and drastically affected by the randomization of the kernel virtual address. When the randomization is applied to the kernel virtual address, the percentage of equal pages is reduced from 100% to 2.6%, regardless of whether other areas are randomized or not. Although we discussed in section 3 that Linux is using relocations to implement the kernel randomization, the results indicate that almost all kernel code pages are altered in this relocation process, resulting in a non-shared Linux code region.

This result was a bit surprising, not the fact that randomizing the Linux code will affect to the memory deduplication of its code pages but obtaining that 97.4% are patched due to the randomization process, which points out that the relocations are spread throughout the Linux code memory area.

### 5.2. Linux Data Impact

The results for the Linux data memory region shown in table 2 show that the percentage of equal pages when all areas are randomized is 64% in average, and by only disabling the randomization of the kernel virtual address, this value is increased to 80%. Hence, only by randomizing the kernel virtual address causes deduplication to decrease around a 20%.

The Linux data contains absolute addresses referencing to parts of the kernel itself (to either code and data regions). Considering that a single absolute reference affects the entire page, approximately 15% of pages in this region contain this kind of references. Randomizing other memory areas has negligible effect on the percentage of equal pages in the Linux data.

### 5.3. Modules Code Impact

Regarding the Linux modules code, table 3 points out that when both kernel virtual address and the modules base address are not randomized, then around 25-30% percent have equal content. Any other combination will reduce by around 52% the sharing rate reaching 11-12% of equal pages. Therefore, randomizing either the kernel virtual address or modules is enough to reduce modules code sharing rate by around 52%.

The latter is not surprising because when modules are loaded at a random position, all the relocations with absolute addresses alter the contents of the module's code region (e.g., accessing to a variable of its data region). However, why randomizing only the kernel virtual address has the same effect than randomizing the modules base address could be less evident. The reason of such behaviour is that any reference to kernel symbols (e.g., `kmalloc()` function) must be fixed up. Then, if the kernel virtual address is randomized, these relocations will differ. In the case of kernel references, the issue occurs with any type of relocation (absolute or relative) [27].

A result that might be seen as unexpected is not having the highest percentage of equal pages for the modules code when the randomization of all memory areas is disabled. This occurs because modules load order is not preserved across different kernel boots, even when the kernel randomization is fully disabled. This issue is further discussed in Section 9.

### 5.4. Modules Data Impact

The minimum percentage of equal pages in the data section of modules is achieved when the modules base address is randomized. The obtained value, as table 4 shows, is 30% and enabling or disabling the rest of memory areas does not affect to this value. On average, enabling the modules base address randomization will reduce around 12% the number of equal pages of the modules data.

When the modules base address is not randomized, the percentage of equal pages of module data is between 38-48% (42% average). From this range, the lowest sharing rate is obtained when the kernel virtual address is randomized. This confirms that the modules data region contains pointers to the kernel. However, it can be less obvious to realize that module data could contain pointers to its own data memory region, which breaks the shareability because of the unpredictable load module order even when the modules base address is not randomized. Around 64% of pages in modules data contain at least one absolute reference to its own module memory (code and data), and 23% to the kernel region.

### 5.5. Vmalloc Space Impact

The results of the `vmalloc` memory region shown in table 5 reveals a really low percentage of equal pages in all cases. Therefore the Linux kernel randomization is not affecting the contents of the `vmalloc` memory region, having in all cases a 3% of equal pages.

This is not surprising since `vmalloc` is mainly used to dynamically and temporally store data by the kernel, and ran-

Randomized Addresses	Impact on Linux Kernel Memory Regions					
	Linux Code	Linux Data	Modules Code	Modules Data	vmalloc	vmemmap
Kernel Base (Physical)	none	none	none	none	none	low
Kernel Base (Virtual)	high	med	high	low	none	none
Physical Direct Mapping	none	none	none	none	none	none
Vmalloc/ioremap	none	none	none	none	none	none
Vmemmap	none	none	none	none	none	high
Modules	none	none	high	med	none	none

**Table 7**

Impact of randomizing Linux base addresses on the different memory regions.

none &lt; 5% | low 5-15% | med 15-30% | high &gt; 30%

domizing memory areas does not affect to that content. This memory is always random in our proposal since it has not effect on memory deduplication.

### 5.6. Virtual Memory Map Impact

The results for the `vmemmap` memory region shown in table 6 indicates that only by randomizing the `vmemmap` address, the percentage of equal pages for this memory region drops to zero compared to the 38% obtained on average when the `vmemmap` memory region is not randomized.

The `vmemmap` memory region contains lists of objects with references to previous and next objects that reside in the `vmemmap` memory region itself. Therefore, by randomizing this area all those pointers will be different and the pages will not be mergeable.

### 5.7. Analysis summary

In order to summarize the impact of the kernel randomization on each memory region, we classified the relative sharing loss obtained from tables 1-6 into none, low, medium and high impact. This measures the relative percentage of equal pages that can not be merged when the kernel randomization is enabled. For example, the `modules` data memory region loses 28.5% of relative equal pages (from 42% to 30%) when the `modules` base address is randomized.

Table 7 summarizes the kernel randomization impact on the different memory regions. Three memory regions are highly impacted because of the kernel randomization (`Linux` code, `modules` code and `vmemmap`) and two have a medium impact (`Linux` data and `modules` data). As detailed in section 6, KASLR-MT will handle those five memory regions in order to restore back the memory sharing rate.

## 6. KASLR-MT: A Multi-Tenancy KASLR

In this section, we present KASLR-MT, a kernel randomization design for multi-tenant cloud systems which is compatible with memory deduplication while providing a strong level of security. As summarized in section 5.7, there are memory regions with none/low impact and others with medium/high impact on memory deduplication. Based on that, we have designed KASLR-MT to keep a strong security while maximizing memory deduplication sharing rate,

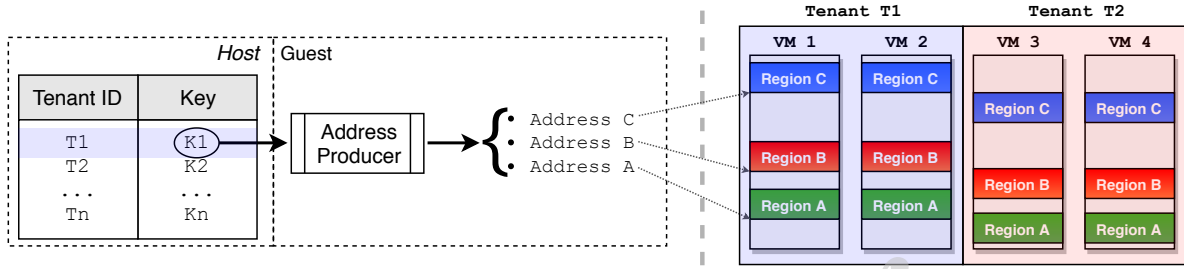
as analyzed in section 8.

As stated in section 4 and studied in section 5, kernel randomization produces alterations in memory contents that breaks memory sharing. To remedy this situation and yield memory sharing by deduplication in the host, kernel memory layouts of guest virtual machines should be as similar as possible.

A quick workaround to achieve this, although it is not a solution by itself, would be to disable kernel randomization in the guests. If kernel randomization is completely disabled, guest kernels are deterministically placed in a default address settled at compile time. Comparing guest instances of the same kernel, all of them produce the same kernel memory layout. As a consequence, relocations that are different when the kernel is randomized are equal when it is not. Even though this approach can be a suitable method in certain environments, for example private clouds properly secured for external attacks, we do not recommend disabling any security mechanism. As stated in section 2.4, the fact of disabling KASLR to obtain better memory deduplication savings introduces serious weaknesses.

In order to properly tackle the randomization-vs-sharing problem, we propose Kernel Address Space Layout Randomization Multi-Tenancy (KASLR-MT), a paravirtualization based solution that enables the hypervisor to provide the kernel memory layout to the guest virtual machines in multi-tenant cloud environments. It gives the host machine the ability to decide the locations of different kernel memory regions of guest virtual machines. There are two possibilities for each kernel memory region:

1. **None/Low impact memory regions.** If a memory region has low or negligible impact when all memory base address are randomized, then this memory region can be randomized without sharing its base address with any other kernel. Note that the base address of this memory region will be different among virtual machines in both, same and different tenants.
2. **Medium/High impact memory regions:** If a memory region has medium or high impact when any memory base address is randomized, then this memory region must be randomized in a per-tenant basis, *sharing* its base address with other kernels belonging to the



**Figure 2:** Design of KASLR-MT. A tenant key is transferred to the guest to deterministically produce the final addresses of kernel memory regions.

same tenant. Therefore, the base address of this memory region will be the same among virtual machines belonging to the same tenant and different among virtual machines of different tenants.

This solution is designed for multi-tenant cloud systems, where different tenants are owners of different groups of virtual machines. Thus, virtual machines belonging to a same tenant can attain a common kernel memory layout, allowing the host machine to share more memory. Security is kept, since the kernel memory layout is still unpredictable from the attacker perspective. This is valid even if virtual machines are not exactly identical, because the granularity of memory deduplication is 4096 bytes, which enables it to deduplicate parts of the virtual machines. KASLR-MT will provide the maximum possible benefit compared with the memory deduplication when the kernel randomization is enabled.

The cloud infrastructure maintains a table with one-to-one correspondence, linking Tenant-ID and a unique random key. Each unique key serves to deterministically produce base addresses of guest kernel memory regions. The table relation must be bijective so that there are no duplicate keys in the cloud infrastructure, and keys must be unpredictable. Otherwise, since the algorithm to produce addresses from a given key must be highly deterministic, an attacker who can predict the key of a victim tenant will be able to guess the kernel address space layout of the victim's virtual machines.

Figure 2 shows the design of KASLR-MT. The table relating each tenant with its key is stored in the host machine, so that the hypervisor transfers the corresponding key to a virtual machine when it is turned on. Guest kernels will use the key as input to the Address Producer to get memory addresses of regions that want to randomize. On the right part of the figure, the kernel memory of four different virtual machines are depicted. Tenant T1 owns two of them (VM-1 and VM-2) while the other virtual machines (VM-3 and VM-4) belong to Tenant T2. Both VM-1 and VM-2 obtain the same memory layout because both kernels are using the same key (K1) as input to the Address Producer algorithm. On the other hand, VM-3 and VM-4 use the key of Tenant T2 (K2), obtaining a different memory layout.

The table information needs to be distributed to all the host machines forming the cloud infrastructure. Coherence of this distributed state in the infrastructure is important to

support virtual machine migration. If, otherwise, the keys were only kept locally in host machines, a migrated virtual machine whose kernel memory layout was generated with a different key would introduce layout diversity within a same group of virtual machines.

The lifetime of a tenant's key goes from when its first virtual machine starts to when its last active virtual machine turns off, independently of the state of other virtual machines belonging to different tenants. A key cannot be changed if any virtual machine is running, since it would introduce layout diversity as well; i.e., kernel memory layout of guests started after a key change would likely differ from those started before. When the last virtual machine of a given tenant turns off, the key can be safely purged. In fact, we strongly recommend to purge it, in order to force the creation of a new key for subsequent guest kernels.

With KASLR-MT, kernel memory base addresses can be randomized in two different ways:

1. **Per-Tenant:** For a given memory region, a random key is generated by the host and associated to a tenant. The key is shared across all virtual machines belonging to the same tenant to enable them to generate the random base address. Guests will have the same memory base address (random) for a particular memory region. KASLR-MT uses a Per-Tenant approach for memory regions where the impact is high or medium.
2. **Per-VM:** For a given memory region, a random key is generated by the host and associated to a virtual machine. Every time a guest reboots, the region will have a different memory base address and it will not be shared across virtual machines belonging to same or other tenants. KASLR-MT uses a Per-VM approach for memory regions where the impact is low or none.

KASLR-MT combines the deduplication effectiveness of disabling kernel randomization and the statistical defense provided by the kernel randomization protection. The approach followed by KASLR-MT is similar to the one commonly used by some major operating systems (e.g., Windows and Mac OS) to randomize the virtual address of userspace libraries, using a per-boot ASLR scheme: a random system-wide value is computed once at system startup, and it is used to calculate the virtual address where libraries



are loaded at. This random value is not changed until the next system reboot. Similarly, our design uses a random key, which determines the guest kernel memory layout, for all the virtual machines of a tenant until the moment when the last one shuts down. However, with per-boot userspace ASLR, a local attacker already knows the address space layout of the target application. This is not the case for KASLR-MT, which offers the same protection as KASLR. Further details can be found in section 8.

## 7. KASLR-MT Linux Implementation

In this section, we present the KASLR-MT implementation in the Linux kernel based on the results of section 5.

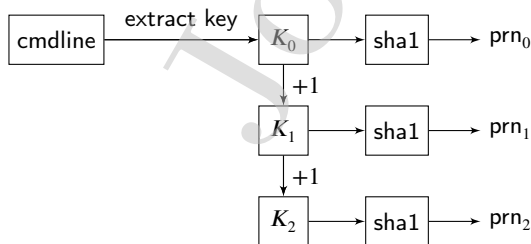
Table 8 show the randomization approach followed by KASLR-MT to maximize the deduplication sharing rate in the Linux kernel. Following this approach, it is expected to obtain a memory sharing rate similar to that obtained when kernel randomization is disabled. Although the kernel base physical address influences on the vmemmap memory region, we consider that it is not worth to randomize this address Per-Tenant because of its low impact.

Kernel Base Addresses	Randomization
Kernel Physical Address	Per-VM
Kernel Virtual Address	Per-Tenant
Physical Direct Mapping	Per-VM
Vmalloc/ioremap	Per-VM
Vmemmap	Per-Tenant
Modules	Per-Tenant

**Table 8**  
KASLR-MT randomization approach for the different Linux kernel memory base addresses.

To communicate the random key to the guest Linux kernel, we use the kernel's command-line parameters, passing directly its corresponding key value. Since Linux normally prints the contents of the cmdline to the kernel ring buffer at the beginning of its boot process, we need to retrieve the key and clean-up the cmdline buffer before it is printed out, to avoid leaking this information. If the key is leaked, every address being randomized from that key could be calculated.

Once the key is obtained, three addresses need to be produced from it. Actually, for our proof of concept, we need to get a pseudo-random number per address, which will be



**Figure 3:** Block diagram describing our proof-of-concept implementation of the Address Producer component.

```

776 static unsigned long find_random_virt_addr(...)
778 {
779     ...
791     slots = (KERNEL_IMAGE_SIZE - minimum - image_size) /
792             CONFIG_PHYSICAL_ALIGN + 1;
793
794     random_addr = kaslr_get_random_long("Virtual") % slots;
795     if (kaslrmt_enabled)
796         random_addr = kaslrmt_prn[0];
797     else
798         random_addr = kaslr_get_random_long("Virtual") % slots;
799
800     return random_addr * CONFIG_PHYSICAL_ALIGN + minimum;
801 }

```

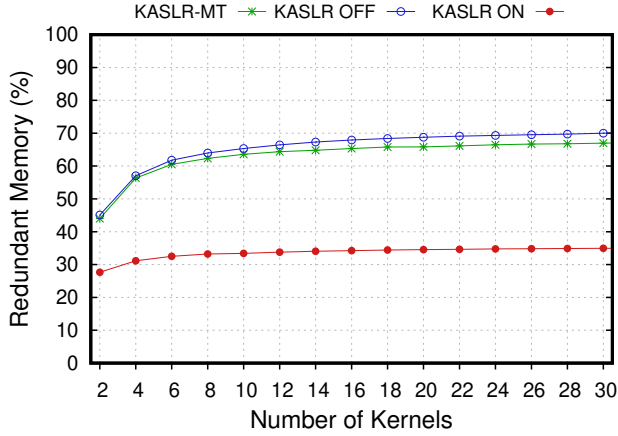
**Listing 2:** Changes in arch/x86/boot/compressed/kaslr.c file of the Linux source code.

used to calculate deterministically the final virtual address. To achieve this in a secure and robust way, we could use a cryptographically secure pseudorandom number generator (CSPRNG), seeding it with the transferred key to produce pseudo-random numbers, or a proper key derivation function (KDF) chain, as used in some encryption algorithms (e.g., AES-GCM-SIV [14]). However, since our purpose is to develop a proof of concept to evaluate our design, we have followed a simple approach to simplify the comprehension, inspired in these kind of algorithms. The general idea is depicted in figure 3. After extracting the key from the kernel command-line, a digest of the key itself is calculated by using the sha1 algorithm, and the result is treated as a pseudo-random number to get the final kernel base virtual address. The same procedure is done two more times, to extract the pseudo-random numbers for obtaining the vmemmap address and the modules base, incrementing the key by one before the digest calculation.

Each pseudo-random number is used instead of a per-region random number. Listing 2 exemplifies how the virtual address of the kernel base is obtained, using the KASLR-MT pseudo-random number instead of calling to kaslr\_get\_random\_long(). The remaining addresses are calculated similarly.

Regarding the key length, given that we are using the kernel command-line as input method, we consider that keys should use common ascii alphanumeric characters. However, keys should have enough entropy to resist brute force attacks against a leaked kernel address. If a kernel address is leaked, an attacker could be able to recover the key by brute forcing the hash algorithm. Based on that the maximum entropy that can be obtained through get\_random\_bytes() is 64 bits, we consider this value as a proper entropy for a KASLR-MT key value. In addition, we want to use common printable ascii characters, avoiding the space and DEL characters. Therefore, the valid ascii characters are from value 0x21 ('!') until value 0x7e (~). This is a total of 94 characters in our alphabet of valid characters for a key. Considering a desired





**Figure 4:** Redundant memory curve for different number of simultaneous kernels.

entropy of 64 bits, the key size must be at least 10 characters:

$$x * \log_2(94) = 64 \rightarrow x \approx 9.76$$

## 8. Evaluation

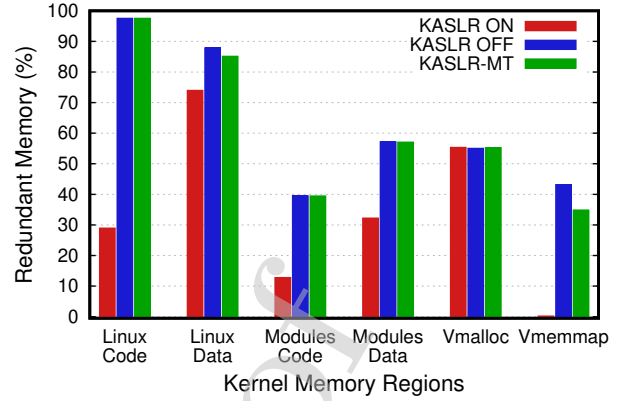
In this section, we evaluate the memory deduplication saving effectiveness of KASLR-MT as well as the security considerations.

### 8.1. Memory Deduplication Savings

In order to measure how much kernel memory is being saved by deduplication with KASLR-MT compared with Linux KASLR, we have launched another experiment, running from 2 to 30 simultaneous virtual machines. Each virtual machine configuration is the same as used in section 5, a generic GNU/Linux Ubuntu 19.04 with Linux v5.0.6, with Gnome desktop and full networking (NAT mode provided by Qemu). The physical machine used to run the experiments has an Intel Xeon W-2155 processor (Skylake server microarchitecture) and 32 GiB of SDRAM memory. The hypervisor used is KVM (Linux kernel 4.19-ARCH) along with Qemu VMM version 4.0.0.

Userspace activities modify the kernel state (changes in data structures as existing processes, open files, etc.). Inevitably, even though these alterations are not related with kernel randomization, they are present in the measurements. The possible combinations of userspace workloads are endless. For this reason, virtual machines run the GNU/Linux Ubuntu distribution in the experiment with the aim of obtaining a real environment with a representative userland workload. However, since the scope of this paper is to study the effects that kernel randomization has in memory contents, memory pages belonging to userland should be discarded.

Figure 4 shows the resulting percentage of redundant memory for each case. With only two kernels, if kernel randomization is enabled there is a 27.65% of redundant memory, which increases to 45.09% by disabling it. With our solution, the value obtained is close to the latter, a 44.02%. The



**Figure 5:** Percentage of redundant memory per kernel region, for 30 kernels.

percentage increases logarithmically as more kernel memories are added, until they approach a theoretical limit, where the curve grows slower. However, when kernel randomization is enabled, the curve has less slope; this is caused by a greater number of pages with matchless contents. This limit is approximately 70% when kernel randomization is not enabled, 67% in our solution, and 35% when kernel randomization is enabled.

Taking a stabilized case, for example with 30 kernels, we have split the memories of all the kernels by region, as shown in figure 5. From it, we can confirm which are the kernel regions benefited from our solution. The Linux code region gets exactly the same sharing as when kernel randomization is disabled, while the Linux data region (including .data, .bss and other data sections) is close, with only 2.73% of sharing loss. The modules code and data regions are also benefited despite the fact that the loading order factor, which is independent of kernel randomization, is affecting their contents similarity (elucidated in section 9). Similarly, vmemmap can share a 34.89% of its contents with our solution, an 8.29% less than when kernel randomization is disabled; but still a good portion, considering that its contents are almost entirely matchless (0.23%) when kernel randomization is enabled. With regard to vmalloc, we can see that it has similar redundant memory, independently of kernel randomization. Consequently, figure 5 shows that our solution is close to the best case scenario (i.e., KASLR OFF) in terms of redundant memory.

### 8.2. Security Considerations

Regarding security implications of KASLR-MT, it extends the memory saving benefits keeping the protection provided by KASLR.

Our proposed solution is intended to be beneficial in multi-tenant cloud systems, where several tenants are owners of one or more groups of virtual machines, and all of them share the resources of a single physical machine by the use of virtualization technologies. With KASLR-MT, the memory deduplication savings are restored back without disabling the kernel randomization, enabling Linux systems

Boot 1		Boot 2	
Module	Virtual Address	Module	Virtual Address
floppy	0xffffffffc0002000	floppy	0xffffffffc0002000
pata_acpi	0xffffffffc0017000	pata_acpi	0xffffffffc0017000
i2c_piix4	0xffffffffc001c000	i2c_piix4	0xffffffffc001c000
e1000	0xffffffffc002f000 →	← e1000	0xffffffffc0028000
ip_tables	0xffffffffc0052000	autofs4	0xffffffffc004b000
psmouse	0xffffffffc005b000	psmouse	0xffffffffc005b000
autofs4	0xffffffffc0081000 →	← x_tables	0xffffffffc0086000
x_tables	0xffffffffc008d000	ip_tables	0xffffffffc0097000
sch_fq_codel	0xffffffffc00b8000	sch_fq_codel	0xffffffffc00b7000
irqbypass	0xffffffffc0178000	qemu_fw_cfg	0xffffffffc0172000
joydev	0xffffffffc017f000 →	← irqbypass	0xffffffffc0178000
mac_hid	0xffffffffc019b000	mac_hid	0xffffffffc01a3000
serio_raw	0xffffffffc01c3000	serio_raw	0xffffffffc01b8000
input_leds	0xffffffffc01e0000	input_leds	0xffffffffc01c1000
crct10dif_pclmul	0xffffffffc0200000	crct10dif_pclmul	0xffffffffc01c6000
nfit	0xffffffffc0205000	crc32_pclmul	0xffffffffc01fe000
qemu_fw_cfg	0xffffffffc0237000	joydev	0xffffffffc020e000
crc32_pclmul	0xffffffffc0266000	kvm	0xffffffffc024b000
kvm	0xffffffffc026b000	nfit	0xffffffffc0276000
kvm_intel	0xffffffffc029d000	kvm_intel	0xffffffffc028e2000

**Figure 6:** Name and final virtual address of a list of modules being auto-loaded into the Linux kernel at boot time, for two different kernel boots, without enabling kernel randomization. The loading order is not preserved across boots.

to have both, high ratios of memory savings and security.

The trade-off for more memory sharing is to share akin kernel address space layouts among virtual machines belonging to the same group. Effectively, this solution maintains an equivalent protection offered by kernel randomization. On the one hand, it protects against external attackers including network applications interacting with the kernel and unknown tenants running virtual machines in the same host machine. On the other hand, attacks from those virtual machines that share the same kernel address space layout are equivalent to local attacks (userspace applications attacking its own kernel). In all of these cases, the kernel memory layout is unpredictable from the attacker perspective.

There is a case where KASLR-MT has a drawback: an attacker successfully bypassing the kernel randomization of a particular machine will be able to use that leak to bypass the kernel randomization of another machine belonging to the same tenant. However, the effort required to bypass KASLR-MT is the same as KASLR.

## 9. Discussion

In this section, we discuss some affairs that have not been thoroughly detailed for being out of scope.

As advanced in previous sections, the loading order of Linux loadable modules is not deterministic. We are focusing on the kernel's automatic module loader at boot time, without explicit user/administrator action once the system is running [5]. Linux can perform the module auto-loading in two different ways: by using hardware-driven mechanisms sending events to userland daemons (e.g., udev) when hardware is discovered, or by using an older mechanism that uses the `request_module()` kernel function to load the module from userspace. In all cases, this job is done when a certain

module is needed and has not yet been loaded. This is common for device drivers compiled as modules.

Most of the cases, modules auto-loading is done asynchronously, exploiting the advantages of parallelism to increase performance. As a consequence, even if kernel randomization is disabled, a similar effect occurs in the code and data memory regions of the modules, caused by the intrinsic randomization derived from the unpredictability of their loading order. Furthermore, taking into account that the list and number of modules being loaded into a kernel may vary depending on the task in which it is involved, even if a pre-established loading order is determined for a certain group of Linux modules, the boundless variety of different possible modules makes it a complex issue.

Figure 6 compares a list of modules being auto-loaded into the Linux kernel at boot time for two different kernel boots (*boot 1* and *boot 2*) when the kernel randomization is fully disabled. There are modules loaded in different positions and virtual addresses. For example, the `ip_tables` is loaded in position 5 in *boot 1* and in position 8 in *boot 2*, resulting in different virtual addresses where they have been loaded. A situation that could seem unexpected is when a module is loaded in the same position in both boots but in different virtual address. An example is the module `e1000`, where the base address in *boot 1* is `0xffffffffc002f000` and `0xffffffffc0028000` in *boot 2*.

This behaviour alters modules memory contents in a similar way as kernel randomization does. As a consequence, all their absolute references (in code and data) will not be equal when different kernel memories are compared.

Therefore, even fully disabling the kernel randomization, the load order module and addresses where they are allocated are not the same across different kernel reboots. The asynchronous userspace module loading requests provide a

flexible and time saving boot time but, unfortunately, it also reduces the memory sharing rate. As a future work, having modules compiled and linked as position-independent code (PIC) could alleviate this issue. Unfortunately, current Linux versions do not support PIC and an analysis to obtain the real benefit would be necessary since, as we have discussed, modules also contain references to other kernel parts that can break the memory sharing.

## 10. Conclusions

In this paper we presented KASLR-MT, a new Linux kernel randomization design for multi-tenant cloud systems, compatible with memory deduplication, which maximizes memory savings rate while providing a strong security.

We identified why the most widely and effective technique used to mitigate attacks at kernel level, KASLR, fails to provide protection and shareability at the same time. To design KASLR-MT, we performed a deep analysis on how kernel randomization affects to the shared memory. Then, we propose KASLR-MT, the first effective and practical Kernel ASLR memory protection that maximizes the memory deduplication savings rate while providing a strong security.

We have implemented and tested KASLR-MT in the Linux kernel and our results showed that KASLR-MT is not intrusive, highly scalable and maximizes the memory savings rate while providing a strong security.

## References

- [1] Arcangeli, A., Eidus, I., Wright, C., 2009. Increasing memory density by using ksm, in: Proceedings of the linux symposium, Citeseer. pp. 19–28.
- [2] Binu, A., Kumar, G.S., 2011. Virtualization techniques: a methodical review of xen and kvm, in: International Conference on Advances in Computing and Communications, Springer. pp. 399–410.
- [3] Computer History Museum, . What happened on september 9th. URL: <http://www.computerhistory.org/t dih/September/9/>. [Retrieved: Sep, 2018].
- [4] Cook, K., 2013. Kernel address space layout randomization. URL: <https://outflux.net/slides/2013/lss/kaslr.pdf>. [Retrieved: Sep, 2019].
- [5] Corbet, J., 2017. Restricting automatic kernel-module loading. URL: <https://lwn.net/Articles/740455/>. [Retrieved: Sep, 2019].
- [6] Daley, R.C., Dennis, J.B., 1967. Virtual memory, processes, and sharing in multics, in: Proceedings of the first ACM symposium on Operating System Principles, ACM. pp. 12–1.
- [7] Evtushkin, D., Ponomarev, D., Abu-Ghazaleh, N., 2016. Jump over aslr: Attacking branch predictors to bypass aslr, in: The 49th Annual IEEE/ACM International Symposium on Microarchitecture, IEEE Press. p. 40.
- [8] Filipe, R., Barreto, J., 2011. End-to-end data deduplication for the mobile web, in: 2011 IEEE 10th International Symposium on Network Computing and Applications, IEEE. pp. 334–337.
- [9] Fábrega, F.J.T., Javier, F., Guttman, J.D., 1995. Copy on write.
- [10] Ganz, J., Peisert, S., 2017. Aslr: How robust is the randomness?, in: 2017 IEEE Cybersecurity Development (SecDev), IEEE. pp. 34–41.
- [11] Gisbert, H.M., Ripoll, I., 2014. On the effectiveness of nx, ssp, renews, and aslr against stack buffer overflows, in: 2014 IEEE 13th International Symposium on Network Computing and Applications, IEEE. pp. 145–152.
- [12] Gras, B., Razavi, K., Bosman, E., Bos, H., Giuffrida, C., 2017. Aslr on the line: Practical cache attacks on the mmu., in: NDSS, p. 26.
- [13] Gruss, D., Maurice, C., Fogh, A., Lipp, M., Mangard, S., 2016. Prefetch side-channel attacks: Bypassing smap and kernel aslr, in: Proceedings of the 2016 ACM SIGSAC conference on computer and communications security, ACM. pp. 368–379.
- [14] Gueron, S., Langley, A., Lindell, Y., 2017. Aes-gcm-siv: Specification and analysis. IACR Cryptology ePrint Archive 2017, 168.
- [15] Guo, F., Li, Y., Xu, Y., Jiang, S., Lui, J.C., 2017. Smartmd: A high performance deduplication engine with mixed pages, in: 2017 {USENIX} Annual Technical Conference ({USENIX}{ATC} 17), pp. 733–744.
- [16] Gupta, D., Lee, S., Vrabie, M., Savage, S., Snoeren, A.C., Varghese, G., Voelker, G.M., Vahdat, A., 2010. Difference engine: Harnessing memory redundancy in virtual machines. Communications of the ACM 53, 85–93.
- [17] Harnik, D., Pinkas, B., Shulman-Peleg, A., 2010. Side channels in cloud services: Deduplication in cloud storage. IEEE Security & Privacy 8, 40–47.
- [18] Herlends, W., Hobson, T., Donovan, P.J., 2014. Effective entropy: Security-centric metric for memory randomization techniques, in: 7th Workshop on Cyber Security Experimentation and Test ({CSET} 14).
- [19] Jang, Y., Lee, S., Kim, T., 2016. Breaking kernel address space layout randomization with intel tsx, in: Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, ACM. pp. 380–392.
- [20] Kaur, G., Bala, A., Chana, I., 2019. An intelligent regressive ensemble approach for predicting resource usage in cloud computing. Journal of Parallel and Distributed Computing 123, 1 – 12. URL: <http://www.sciencedirect.com/science/article/pii/S0743731518306063>, doi:<https://doi.org/10.1016/j.jpdc.2018.08.008>.
- [21] Kim, S., Kim, H., Lee, J., Jeong, J., 2014. Group-based memory oversubscription for virtualized clouds. Journal of Parallel and Distributed Computing 74, 2241 – 2256. URL: <http://www.sciencedirect.com/science/article/pii/S0743731514000033>, doi:<https://doi.org/10.1016/j.jpdc.2014.01.001>.
- [22] Koo, H., Chen, Y., Lu, L., Kemerlis, V.P., Polychronakis, M., 2018. Compiler-assisted code randomization, in: 2018 IEEE Symposium on Security and Privacy (SP), IEEE. pp. 461–477.
- [23] Lindemann, J., Fischer, M., 2018. A memory-deduplication side-channel attack to detect applications in co-resident virtual machines, in: Proceedings of the 33rd Annual ACM Symposium on Applied Computing, ACM, New York, NY, USA. pp. 183–192. URL: <http://doi.acm.org/10.1145/3167132.3167151>, doi:10.1145/3167132.3167151.
- [24] Luo, B., Yang, Y., Zhang, C., Wang, Y., Zhang, B., 2018. A survey of code reuse attack and defense, in: International Conference on Intelligent and Interactive Systems and Applications, Springer. pp. 782–788.
- [25] Marco-Gisbert, H., Ripoll, I., 2018. return-to-csu: a new method to bypass 64-bit linux aslr. URL: <https://www.blackhat.com/asia-18/black-hat-asia-2018/Black-Hat-Conference-date-20-03-2018-Through-23-03-2018>.
- [26] Marco-Gisbert, H., Ripoll, I., 2019. Address space layout randomization next generation. Applied Sciences 9, 2928.
- [27] Matz, M., Hubicka, J., Jaeger, A., Mitchell, M., 2013. System v application binary interface. AMD64 Architecture Processor Supplement, Draft v0 99.
- [28] Mell, P., Grance, T., et al., 2011. The nist definition of cloud computing.
- [29] Mihós, G., Murray, D.G., Hand, S., Fetterman, M.A., 2009. Satori: Enlightened page sharing, in: Proceedings of the 2009 conference on USENIX Annual technical conference, pp. 1–1.
- [30] Mittal, S., 2017. A survey of techniques for architecting tlbs. Concurrency and Computation: Practice and Experience 29, e4061.
- [31] Neves, P., Ferreira, P., Barreto, J., 2013. Leveraging web prefetching systems with data deduplication, in: 2013 IEEE 12th International Symposium on Network Computing and Applications, IEEE. pp. 259–262.
- [32] Oliverio, M., Razavi, K., Bos, H., Giuffrida, C., 2017. Secure Page Fusion with VUion, in: Proceedings of the 26th Sympo-



- sium on Operating Systems Principles, ACM, New York, NY, USA. pp. 531–545. URL: <http://doi.acm.org/10.1145/3132747.3132781>, doi:10.1145/3132747.3132781.
- [33] Pappas, V., Polychronakis, M., Keromytis, A.D., 2014. Dynamic reconstruction of relocation information for stripped binaries, in: International Workshop on Recent Advances in Intrusion Detection, Springer. pp. 68–87.
- [34] PaX, 2003. Pax address space layout randomization (aslr). URL: <https://pax.grsecurity.net/docs/aslr.txt>. [Retrieved: Sep, 2018].
- [35] Payer, M., 2016. HexPADS: A Platform to Detect “Stealth” Attacks, in: Caballero, J., Bodden, E., Athanasopoulos, E. (Eds.), Engineering Secure Software and Systems, Springer International Publishing, Cham. pp. 138–154.
- [36] Perla, E., Oldani, M., 2010. A Guide to Kernel Exploitation: Attacking the Core. Syngress Publishing.
- [37] Ranjbari, M., Torkestani, J.A., 2018. A learning automata-based algorithm for energy and sla efficient consolidation of virtual machines in cloud data centers. Journal of Parallel and Distributed Computing 113, 55 – 62. URL: <http://www.sciencedirect.com/science/article/pii/S074373151730285X>, doi:<https://doi.org/10.1016/j.jpdc.2017.10.009>.
- [38] Schwarz, B., Debray, S., Andrews, G., 2002. Disassembly of executable code revisited, in: Ninth Working Conference on Reverse Engineering, 2002. Proceedings., IEEE. pp. 45–54.
- [39] Shacham, H., et al., 2007. The geometry of innocent flesh on the bone: return-into-libc without function calls (on the x86), in: ACM conference on Computer and communications security, New York., pp. 552–561.
- [40] Sharma, P., Kulkarni, P., 2012. Singleton: system-wide page deduplication in virtual environments, in: Proceedings of the 21st international symposium on High-Performance Parallel and Distributed Computing, ACM. pp. 15–26.
- [41] Steven, S., 2011. Reducing runtime memory in windows 8. URL: <https://blogs.msdn.microsoft.com/b8/2011/10/07/reducing-runtime-memory-in-windows-8/>. [Retrieved: Sep, 2018].
- [42] Suzaki, K., Iijima, K., Yagi, T., Artho, C., 2011. Memory deduplication as a threat to the guest os, in: Proceedings of the Fourth European Workshop on System Security, ACM, New York, NY, USA. pp. 1:1–1:6. URL: <http://doi.acm.org/10.1145/1972551.1972552>, doi:10.1145/1972551.1972552.
- [43] Tziritas, N., Khan, S.U., Xu, C.Z., Loukopoulos, T., Lalis, S., 2013. On minimizing the resource consumption of cloud applications using process migrations. Journal of Parallel and Distributed Computing 73, 1690 – 1704. URL: <http://www.sciencedirect.com/science/article/pii/S0743731513001585>, doi:<https://doi.org/10.1016/j.jpdc.2013.07.020>. heterogeneity in Parallel and Distributed Computing.
- [44] Vaño, F., Marco, H., 2018a. How Kernel Randomization is Canceling Memory Deduplication in Cloud Computing Systems, in: 2018 IEEE 17th International Symposium on Network Computing and Applications (NCA), IEEE. pp. 373–376. URL: <https://ieeexplore.ieee.org/abstract/document/8548338>.
- [45] Vaño, F., Marco, H., 2018b. Slicedup: A Tenant-Aware Memory Deduplication for Cloud Computing, in: The Twelfth International Conference on Mobile Ubiquitous Computing, Systems, Services and Technologies (UBICOMM 2018), pp. 15–20. URL: [https://www.thinkmind.org/download.php?articleid=ubicomm\\_2018\\_1\\_30\\_10065](https://www.thinkmind.org/download.php?articleid=ubicomm_2018_1_30_10065).
- [46] VMware Docs, 2018. VMware vsan: Using deduplication and compression. URL: <https://docs.vmware.com/en/VMware-vSphere/6.5/com.vmware.vsphere.virtualsan.doc/GUID-3D2D80CC-444E-454E-9B8B-25C3F620EFED.html>. [Retrieved: Sep, 2018].
- [47] Waldspurger, C.A., 2002. Memory resource management in vmware esx server. ACM SIGOPS Operating Systems Review 36, 181–194.
- [48] Wang, Z., Sun, D., Xue, G., Qian, S., Li, G., Li, M., 2018. Ada-things: An adaptive virtual machine monitoring and migration strategy for internet of things applications. Journal of Parallel and Distributed Computing URL: <http://www.sciencedirect.com/science/article/pii/S0743731518304404>, doi:<https://doi.org/10.1016/j.jpdc.2018.06.009>.
- [49] Wienand, I., 2006. A survey of large-page support. University of New South Wales, 1–52.
- [50] Xu, Z., Liu, G., Wang, T., Xu, H., 2017. Exploitations of uninitialized uses on macos sierra, in: 11th {USENIX} Workshop on Offensive Technologies ({WOOT} 17).
- [51] Zhou, H., Li, Q., Choo, K.K.R., Zhu, H., 2018. Dadta: A novel adaptive strategy for energy and performance efficient virtual machine consolidation. Journal of Parallel and Distributed Computing 121, 15 – 26. URL: <http://www.sciencedirect.com/science/article/pii/S0743731518304520>, doi:<https://doi.org/10.1016/j.jpdc.2018.06.011>.



**Fernando Vano-Garcia** is a PhD researcher at the University of the West of Scotland, United Kingdom. His main research interests include cybersecurity, memory management in cloud computing critical infrastructures and virtualization technologies, among others. He is also technical program committee member of international scientific conferences. He completed his BSc Computer Engineering degree at Universitat Politècnica de Valencia, and his MSc in Cybersecurity at Universidad Carlos III de Madrid, Spain.



**Hector Marco** is an associate professor and cybersecurity researcher at the University of the West of Scotland, UK. He holds a PhD in Computer Science, Cybersecurity, from Universitat Politècnica de Valencia, Spain. Hector is senior member of the Institute of Electrical and Electronics (IEEE), and member of the Engineering and Physical Sciences Research Council (EPSRC) in UK. Previously, he was research associate at the Universitat Politècnica de Valencia where he co-founded the “cybersecurity research group”. Hector was part of the team developing the multi-processor version of the XtratuM hypervisor to be used by the European Space Agency in its space crafts. He participated in multiple research projects as Principal Investigator and Co-Investigator. Hector is author of many papers of computer security and cloud computing. He has been invited multiple times to reputed cybersecurity conferences such as Black Hat and DeepSec. Hector has published more than 10 Common Vulnerabilities and Exposures (CVE) affecting important software such as the Linux kernel. He has received honors and awards from Google, Packet Storm Security and IBM for his security contributions to the design and implementation of the Linux ASLR. Hector’s professional interests include low level cybersecurity, kernel and userland security, virtualization security and applied cryptography.

**Highlights:**

- We discuss the effects on memory caused by the kernel randomization protection technique in virtualized environments.
- We perform a comprehensive analysis of the influences of Linux kernel randomization on different kernel memory regions.
- We identify why KASLR fails to provide protection and shareability at the same time.
- We provide a new Linux kernel randomization design for multi-tenant cloud systems, compatible with memory deduplication.



The authors declare that there is no conflict of interest regarding the publication of this paper.

Journal Pre-proof