

The attached DRAFT document (provided here for historical purposes) has been superseded by the following publication:

Publication Number: **NIST Internal Report (NISTIR) 8176**

Title: *Security Assurance Requirements for Linux Application
Container Deployments*

Publication Date: **October 2017**

- Final Publication: <https://doi.org/10.6028/NIST.IR.8176> (direct link: <http://nvlpubs.nist.gov/nistpubs/ir/2017/NIST.IR.8176.pdf>).
- Information on other NIST Computer Security Division publications and programs can be found at: <https://csrc.nist.gov/>

The following information was posted with the attached DRAFT document:

August 1, 2017

Comments Due: August 25, 2017

Email Comments to: NISTIR8176@nist.gov

Author(s)

Ramaswamy Chandramouli (NIST)

Announcement

NIST requests comments on the release of **Draft NISTIR 8176, *Security Assurance Requirements for Linux Application Container Deployments***.

Application Containers are slowly finding adoption in enterprise IT infrastructures. To address security concerns associated with deployment of application container platforms, NIST Special Publication 800-190 (2nd Draft), *Application Container Security Guide*, identified security threats to the components of the platform hosting the containers and related artifacts involved in building, storing and using container images. It has also proposed countermeasures for the following components: Hardware, Host OS, Container Runtime, Image, Registry and Orchestrator.

To implement the countermeasures one or more security solutions are needed. To assess the effectiveness of the security solutions implemented based on these recommendations, it is necessary to analyze them and outline the security assurance requirements they must satisfy to meet their intended objectives. This is the contribution of Draft NISTIR 8176. The focus is on application containers on Linux platforms.

The security solutions for which security assurance requirements have been derived cover the following areas:

1. Hardware-based root of trust providing integrity for boot process,
2. Configuration options using host OS kernel features and kernel loadable modules,
3. Protection measures for building and storing container images, and
4. Configuration options in Orchestrator tools used for rolling out a production infrastructure involving multiple containers and multiple hosts.

1
2
3
4
5
6
7
8
9
10
11
12
13
14

Security Assurance Requirements for Linux Application Container Deployments

Ramaswamy Chandramouli

Draft NISTIR 8176

Security Assurance Requirements for Linux Application Container Deployments

Ramaswamy Chandramouli
Computer Security Division
Information Technology Laboratory

August 2017



U.S. Department of Commerce
Wilbur L. Ross, Jr., Secretary

National Institute of Standards and Technology
Kent Rochford, Acting NIST Director and Under Secretary of Commerce for Standards and Technology

46
47

National Institute of Standards and Technology Internal Report 8176
36 pages (August 2017)

48

49 Certain commercial entities, equipment, or materials may be identified in this document in order to describe an
50 experimental procedure or concept adequately. Such identification is not intended to imply recommendation or
51 endorsement by NIST, nor is it intended to imply that the entities, materials, or equipment are necessarily the best
52 available for the purpose.

53 There may be references in this publication to other publications currently under development by NIST in accordance
54 with its assigned statutory responsibilities. The information in this publication, including concepts and methodologies,
55 may be used by federal agencies even before the completion of such companion publications. Thus, until each
56 publication is completed, current requirements, guidelines, and procedures, where they exist, remain operative. For
57 planning and transition purposes, federal agencies may wish to closely follow the development of these new
58 publications by NIST.

59 Organizations are encouraged to review all draft publications during public comment periods and provide feedback to
60 NIST. Many NIST cybersecurity publications, other than the ones noted above, are available at
61 <http://csrc.nist.gov/publications>.

62

63 **Public comment period: *August 01, 2017* through *August 25, 2017***

64 National Institute of Standards and Technology
65 Attn: Computer Security Division, Information Technology Laboratory
66 100 Bureau Drive (Mail Stop 8930) Gaithersburg, MD 20899-8930
67 Email: NISTIR8176@nist.gov

68 All comments are subject to release under the Freedom of Information Act (FOIA).

69

70

Reports on Computer Systems Technology

71 The Information Technology Laboratory (ITL) at the National Institute of Standards and
72 Technology (NIST) promotes the U.S. economy and public welfare by providing technical
73 leadership for the Nation's measurement and standards infrastructure. ITL develops tests, test
74 methods, reference data, proof-of-concept implementations, and technical analyses to advance the
75 development and productive use of information technology. ITL's responsibilities include the
76 development of management, administrative, technical, and physical standards and guidelines for
77 the cost-effective security and privacy of other than national security-related information in federal
78 information systems.

79

80

Abstract

81 Application Containers are slowly finding adoption in enterprise IT infrastructures. Security
82 guidelines and countermeasures have been proposed to address security concerns associated with
83 the deployment of application container platforms. To assess the effectiveness of the security
84 solutions implemented based on these recommendations, it is necessary to analyze them and
85 outline the security assurance requirements they must satisfy to meet their intended objectives.
86 This is the contribution of this document. The focus is on application containers on a Linux
87 platform.

88

89

90

Keywords

91 application container; capabilities; Cgroups; container image; container registry; kernel loadable
92 module; Linux kernel; namespace; Trusted Platform Module.

93

94

Acknowledgements

95

96

97

98

99

100

101

102

Audience

103 The target audience for this document includes system architects and system administrators for
104 container stacks in enterprise infrastructures or in infrastructures used for offering container
105 service as part of an overall cloud service.

106

107

Trademark Information

108 All registered trademarks or trademarks belong to their respective organizations.

109

110

111 **Executive Summary**

112 Application containers are now slowly finding adoption in production environments due to the
113 following advantages: short development and deployment cycle, resource efficiency through
114 lightweight virtualization, and availability of tools for automating the processes involved. At the
115 same time, security concerns are dictating the pace of adoption. To address these concerns,
116 security guidelines and countermeasures have been proposed by NIST through the Application
117 Container Security Guide (NIST Special Publication 800-190).

118 The Application Security Guide identified security threats to the components of the platform
119 hosting the containers and related artifacts involved in building containers and storing them prior
120 to launch. Taking into consideration the overall security implications for the entire ecosystem
121 involving containers, the document also provided security countermeasures for and through six
122 entities including Hardware, Host OS, Container Runtime, Image, Registry and Orchestrator.

123 To carry out these recommendations in the form of countermeasures, one or more security
124 solution are needed. In order for these security solutions to effectively meet their security
125 objectives, it is necessary to analyze those security solutions and detail the metrics they must
126 satisfy in the form of security assurance requirements. This is the objective and contribution of
127 this document.

128 Linux and its various distributions form the predominant host OS component of the deployed
129 container platforms. Since they are open-source products, sufficient security related information
130 is available to analyze the security solutions that can be configured using features provided by
131 Linux. Hence the focus of this document is on security assurance requirements for security
132 solutions for application containers hosted on Linux. The target audience includes system
133 security architects and administrators who are responsible for the actual design and deployment
134 of security solutions in enterprise infrastructures hosting containerized hosts.

135

136

137 **Table of Contents**

138 **Executive Summary iv**

139 **1 Introduction 1**

140 1.1 Scope of the Document 1

141 1.2 Document Structure 3

142 **2 Security Solutions for Linux Application Container Stack..... 5**

143 2.1 Linux Kernel Feature – Namespaces..... 5

144 2.2 Linux Kernel Feature – Cgroups 5

145 2.3 Linux Kernel Feature – Capabilities 6

146 2.4 Kernel Loadable Modules (or Linux Security Module or LSM) 6

147 2.5 Application Container Security Configuration Process..... 6

148 **3 Hardware-based Security Solutions for Containers 7**

149 3.1 vTPM in the host OS Kernel – Security Assurance Requirements 7

150 3.2 vTPM in a dedicated Container – Security Assurance Requirements..... 8

151 3.3 Leveraging Trusted Execution Support of Hardware 9

152 **4 Assurance Requirements for Host OS Protection 10**

153 4.1 Requirements for Generic Host OS Protection 10

154 4.2 Assurance Requirements for Host OS Protection for Container Escape 10

155 **5 Assurance Requirements for Container Runtime Configuration..... 12**

156 5.1 Requirements for Secure Connection 12

157 5.2 Requirements for Isolation-based Configurations 12

158 5.2.1 Process Isolation for Containers..... 12

159 5.2.2 Filesystem Isolation for Containers..... 13

160 5.2.3 IPC Isolation for Containers..... 14

161 5.2.4 Network Isolation for Containers 14

162 5.2.5 User and Group-level Isolation for Containers..... 16

163 5.3 Requirements for Resource Limiting Solutions 16

164 5.4 Requirements for Least Privilege Configuration for Containers 17

165 5.5 Requirements for Device Isolation solutions 17

166 5.6 Requirements for Container Launching Options 19

167 **6 Assurance Requirements for Image Integrity Solutions 22**

168 **7 Assurance Requirements for Image Registry Protection..... 23**

169 **8 Assurance Requirements for Orchestration Functions..... 24**

170 **9 Adverse Side effect of some Security Solutions..... 25**

171 9.1 Resource Limiting using Cgroups 25

172 9.2 Syscall filters using Seccomp 25

173 **10 Summary and Conclusions..... 26**

174

175

176

List of Appendices

177 **Appendix A— Acronyms 277**

178 **Appendix B— References 288**

179

180

List of Figures

181 Figure 1 – Container Technology Stack 2

182 Figure 2 – vTPM Implemented in a Kernel Module 8

183 Figure 3 – vTPM located in a dedicated Container 9

184

185

List of Tables

186 Table 1– Linux Resource Control using Cgroups 17

187 Table 2 – Prohibited Options for Container Launching..... 19

188

189 **1 Introduction**

190 Application containers are now slowly finding adoption in production environments due to the
191 following advantages: short development and deployment cycle, resource efficiency through
192 lightweight virtualization, and availability of tools for automating the processes involved. To
193 address the security concerns in these environments, the Application Container Security Guide
194 (NIST Special Publication 800-190) [1] (referred to in the rest of this document as *Container*
195 *Security Guide*) identified security threats to the components of the platform hosting the
196 containers as well as related artifacts involved in building containers and storing them prior to
197 launch. Taking into consideration the overall security implications for the entire ecosystem
198 involving containers, the document also provided security countermeasures for and through six
199 entities including Hardware, Host OS, Container Runtime, Image, Registry and Orchestrator.

200 To implement these countermeasures, one or more security solutions are needed. This document
201 discusses potential security solutions that provide the functionality necessary in countermeasures
202 and the kind of security assurance requirements each should satisfy. These security solutions can
203 be broadly classified as:

- 204 (a) Hardware-based root of trust providing integrity for boot process
- 205 (b) Configuration options using host OS kernel features and kernel loadable modules
- 206 (c) Protection measures for building and storing container images
- 207 (d) Configuration options in Orchestrator tools used for rolling out a production
208 infrastructure that involves multiple containers and multiple hosts

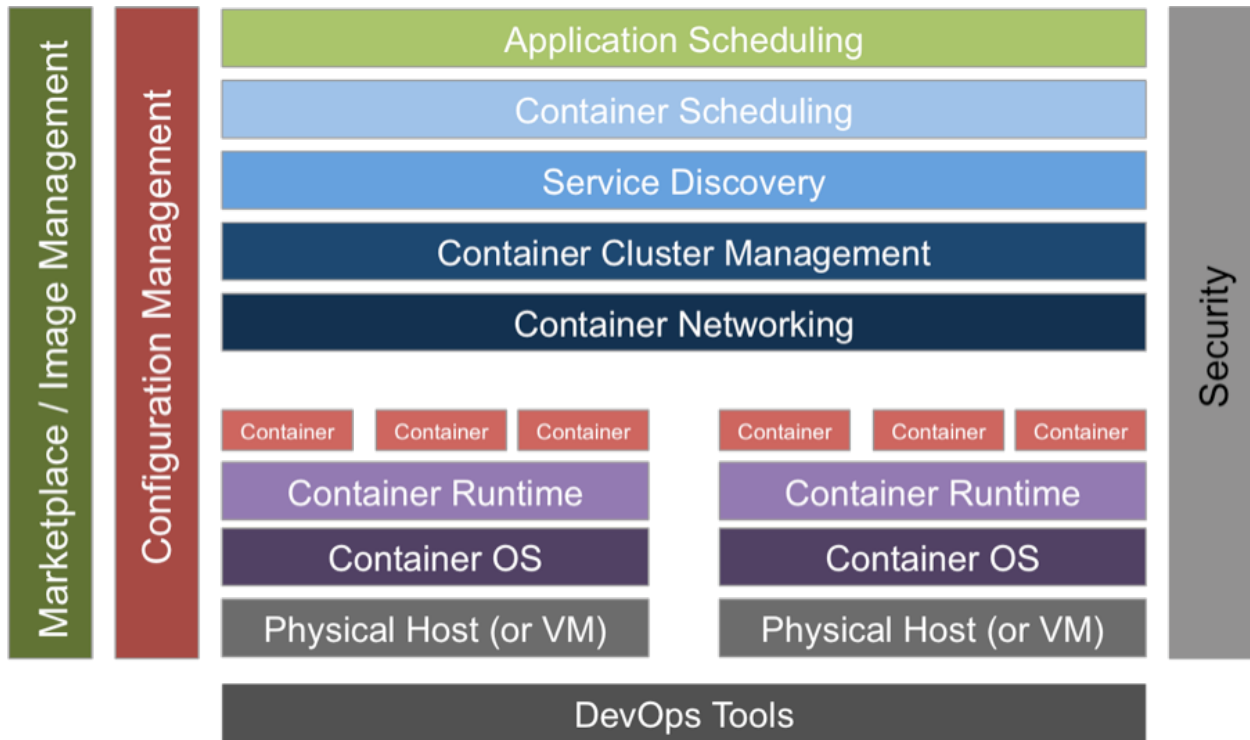
209
210 The purpose of this document is to examine each of the security solutions in the context of the
211 security objectives they are designed to meet and to develop assurance requirements that they
212 should satisfy in order to be effective. The host OS considered is Linux due to the following:

- 213 (a) Ubiquitous adoption in container stacks
- 214 (b) Linux distributions are open-source and allow for sufficient security related information
215 to be made publicly available

216 **1.1 Scope of the Document**

218 The functional architecture diagram of a container technology stack is shown in figure 1. In this
219 diagram, the stack is comprised of the Physical Host (or VM), Container OS (which we will refer
220 to as Host OS in this document), Container Runtime, and the multiple containers. Additionally,
221 tasks such as creating a virtual network linking containers within and across container hosts
222 (Container Networking), creating clusters of container hosts (Container Cluster Management),
223 creating pathway programs to identify and discover a specific container providing a particular
224 service (Service Discovery), scheduling of containers across a cluster (Container Scheduling),
225 and scheduling of specific business applications within various containers (Application
226 Scheduling) that are all performed by multiple tools are incorporated under the umbrella of an
227 Orchestrator software. Before actually launching them as containers on various container hosts,

228 templates of components that constitute a container called Container Image are created using
 229 various DevOps Tools. These container images are stored in a container registry (Image
 230 Management) and are then pulled into container hosts and launched as containers using
 231 Container Runtime tools. The container runtime also provides the interfaces for configuring host
 232 OS parameters and settings associated with kernel-loadable modules to enable secure
 233 deployment of various containers.



234

235

Figure 1 – Container Technology Stack

236

237 • As depicted in figure 1, the security functional layer spans all functional layers of the
 238 container technology stack. The security solutions covering these layers, however, must
 239 be implemented through the following components:

- 240 (a) Physical Host (i.e., hardware, since container hosting on VMs is out of scope for
 241 this document)
- 242 (b) Container OS (Host OS) interfaces
- 243 (c) Container Runtime interfaces
- 244 (d) Image Management and Registry Interfaces
- 245 (e) Orchestrator Interfaces

246 The containers running in the container stack can either be system containers or application
 247 containers. A container that behaves like a full OS and runs programs such as sshd (secure
 248 session establishment) and syslogd (logging capability) is called a system container, while one
 249 that runs only an application is called an application container [2]. This document focuses on
 250 application containers. Before analyzing the security solutions and identifying the assurance

251 requirements they should satisfy, it is necessary to state the execution model of the application
252 containers and the assumed attack model. First, the application is run within a container as a
253 single operating system process. The container has a copy of the application code itself as well as
254 the software stack (consisting of binaries and libraries) [3]. In most cases, this stack can be
255 assembled using some type of library system, avoiding the need for the developer to build and
256 configure the stack from scratch. These quickly assembled stacks are given different names in
257 different container product offerings (e.g., buildpacks, cartridges, etc.). There are stacks for
258 many of the popular programming language runtimes such as Java, PHP, Node.js, and Ruby. For
259 specialized applications, developers can create their own customized stack. The deployment
260 model in a container architecture may involve running copies of the same application in parallel
261 with separate containers, even those spread across different container hosts. In this scenario, the
262 infrastructure may have a mechanism to distribute incoming requests across all instances of the
263 same application using some form of load balancer.

264 The attack model assumed here is that the vulnerability in the application code of the container
265 or its faulty configuration (e.g., the container is configured to run in privileged mode) has been
266 exploited by an attacker to take control of and compromise the privilege code in container
267 runtime and host OS kernel where the latter is trusted by the application code in the container to
268 provide some protection guarantees such as process isolation [4]. An example of such an attack
269 is the replaying, recording, modifying, and dropping of a network packet or a file system access.
270 The security solutions discussed in this document are intended to protect the container runtime
271 and host OS against these types of attacks. Solutions to address the inherent insecure
272 characteristics of the application code itself, such as programming bugs, design flaws or
273 execution models, are beyond the scope of this document.

274 **1.2 Document Structure**

275 The remainder of this document is organized into the following sections and appendices:

- 276 ▪ Section 2 provides an overview of the functions of various Linux kernel features
277 (Namespaces, Cgroups, Capabilities) and kernel loadable modules in providing security for
278 the containerized stack;
- 279 ▪ Section 3 discusses hardware-based security solutions for container environments;
- 280 ▪ Section 4 outlines host OS protection measures and their associated assurance requirements;
- 281 ▪ Section 5 presents, in detail, several container runtime configuration solutions that guarantee
282 container isolation for artifacts such as processes, filesystems, IPC, and networks. It also
283 presents solutions for limiting resources and ensuring least privilege. All solutions are
284 analyzed, and a set of assurance requirements that must be satisfied are presented;
- 285 ▪ Section 6 defines assurance requirements for building and maintaining container images;
- 286 ▪ Section 7 briefly discusses assurance requirements for container registry protection;
- 287 ▪ Section 8 outlines basic security assurance requirements for Orchestration tools;
- 288 ▪ Section 9 identifies some undesirable side effects of some security solutions and the need to
289 exercise caution in the use of such solutions;

- 290 ▪ Section 10 summarizes the various security solution areas that were covered in the document;
- 291 ▪ Appendix A provides the definition for acronyms used in the document; and
- 292 ▪ Appendix B contains a list of references.
- 293

2 Security Solutions for Linux Application Container Stack

In section 1.1, the host OS (in this context, Linux) interfaces were listed as mechanisms for implementing security solutions for a container stack. There are two types of interfaces: Linux kernel interfaces and Kernel Loadable Module (or Linux Security Module or LSM) interfaces. The Linux kernel features associated with the former type of interfaces are: Namespaces, Cgroups, and Capabilities. Out of these, the Namespaces and Cgroups kernel features provide isolation of processes running on top of the host OS and can be the driving features for development of the concept of containers. The salient functions of Linux kernel features and kernel-loadable module features are briefly described in the following sections to provide context for the security configurations and solutions analyzed in the subsequent sections.

2.1 Linux Kernel Feature – Namespaces

Namespaces divide the identifier tables and other structures associated with kernel global resources into separate instances. These partition filesystems, processes, users, network stacks, Inter-process communication (IPC) objects, host names, and other components into separate pieces. For example, each filesystem namespace has its own root directory and mount table [2]. These distinct namespaces can then be bundled in any frequency or combination to provide a unique view of resources for each container and subsequent accessibility to them. The restricted view of resources for a process within a container can be extended to a child process. Configuration capabilities, such as remapped root file systems and virtual network devices, are some of the security solutions that can be enabled using the Namespaces feature. The assurance of a security solution based on namespaces depends on the methods used to enforce namespace isolation, which in turn depends on the kind of metadata associated with each namespace that implements the appropriate access control.

The namespace concept has expanded into a general framework for isolating a range of kernel global resources, the former scope of which was system-wide. Thus, the associated API has also grown to include several system calls. However, there are still some resources that are not namespace-aware (e.g., devices).

2.2 Linux Kernel Feature – Cgroups

Control Groups (Cgroups) are a kernel mechanism for specifying and enforcing hardware resource limits and access controls to a process or a group of processes. Their goal is to prevent a process from hogging all available resources and starving other processes and containers on the host. Thus, Cgroups isolate and limit a given resource over a group of processes to control performance or security. Controlled resources include CPU shares, RAM, network bandwidth, and disk I/O [5]. It can also be used for task control.

The security protection provided by Cgroups are:

- (a) Preventing Denial-of-Service Attacks: It can provide protection against denial-of-service attacks preventing situations such as runaway containers by using features such as task freezing via SIGSTOP, setting limits on process ID (PID) using PID Cgroup to restrict

332 the maximum number of processes per user, and specifying network control parameters
333 such as buffer limits and traffic priority levels (enforced by iptables).

334 (b) Device Integrity Protection: It can restrict access to devices using mandatory access
335 control (MAC) or using a feature that allows the specification a device whitelist.

336 The configuration of Cgroups is enabled by mounting a special Cgroup virtual filesystem
337 (pseudo-filesystem) similar to /proc or /sys that allows viewing of the state of namespaces and
338 controls. The vulnerability of this mechanism is that attacks, such as unmounting or mounting-
339 over, can invalidate the resource limits set by Cgroups configurations. Cgroups can be
340 configured and managed outside of the container management frameworks since it is a
341 configuration feature purely associated with the kernel of the host OS.

342 **2.3 Linux Kernel Feature – Capabilities**

343 The *Capabilities* feature in Linux kernel helps to partition the extensive set of privileges
344 available to root so that processes (in our context, containers) can be allocated just the privileges
345 needed to perform a specific function. Prior to the introduction of the Capabilities feature, a
346 process that needs to open network sockets must be run as a root to perform this single function.
347 This meant that a bug in the corresponding binary, such as /bin/ping, could allow attackers to
348 gain all privileges for the root on the system [6]. By enabling the capability CAP_NET_RAW, a
349 version of ping can be created that has only the privileges enabled by this capability rather than
350 full root privileges. The security consequence of this is that the potential attackers would gain
351 significantly fewer privileges from exploiting the ping utility.

352

353 **2.4 Kernel Loadable Modules (or Linux Security Module or LSM)**

354 Kernel Loadable Modules, as the name implies, are modules loaded into the Linux kernel and
355 provide security functions. Examples include SELinux, AppArmor, and Seccomp. SELinux and
356 AppArmor enable specification and enforcement of mandatory access control (MAC) on
357 processes and objects. Seccomp enables specification of system call restrictions, and thus
358 reduces the Linux kernel attack surface.

359

360 **2.5 Application Container Security Configuration Process**

361 The Linux host OS kernel features—such as namespaces, Cgroups, and Capabilities—can be
362 leveraged to create a secure configuration for each container. Many container runtime products
363 offer APIs to create secure configurations for containers within a host. A typical container
364 runtime, generally accessed through a client, contains a library that directly makes the syscalls
365 and performs work on behalf of its client such as creating the required kernel namespaces,
366 Cgroups, and management of capabilities. Other administrative functions that may have security
367 implications (e.g., lack of availability due to uneven workloads) such as distribution of
368 containers across hosts and the creation of host clusters are managed by a set of tools called
369 Orchestrators.

3 Hardware-based Security Solutions for Containers

371 The *Container Security Guide*, under the topic of Hardware Countermeasures, recommends a
372 trusted computing model that starts with the measured/secured boot, provides a verified system
373 platform, and builds a chain of trust rooted in hardware. This chain of trust then extends to
374 bootloaders, the OS kernel, and the OS components to enable cryptographic verification of boot
375 mechanisms, system images, container runtimes, and container images. The technical solutions
376 for implementing a trusted computing module (TPM) for a containerized host are outlined in [7].
377 Two such approaches are discussed in this document as well as the security assurance required
378 for each solution.

379 Both approaches involve a combination of hardware-based, or physical, TPM and a software-
380 based vTPM (virtual TPM). The difference between the two approaches is in the location where
381 vTPM is placed in the container stack. The security solution where vTPM is placed in the Linux
382 kernel is discussed in section 3.1, and the solution where vTPM is placed in a dedicated
383 container is the topic of section 3.2.

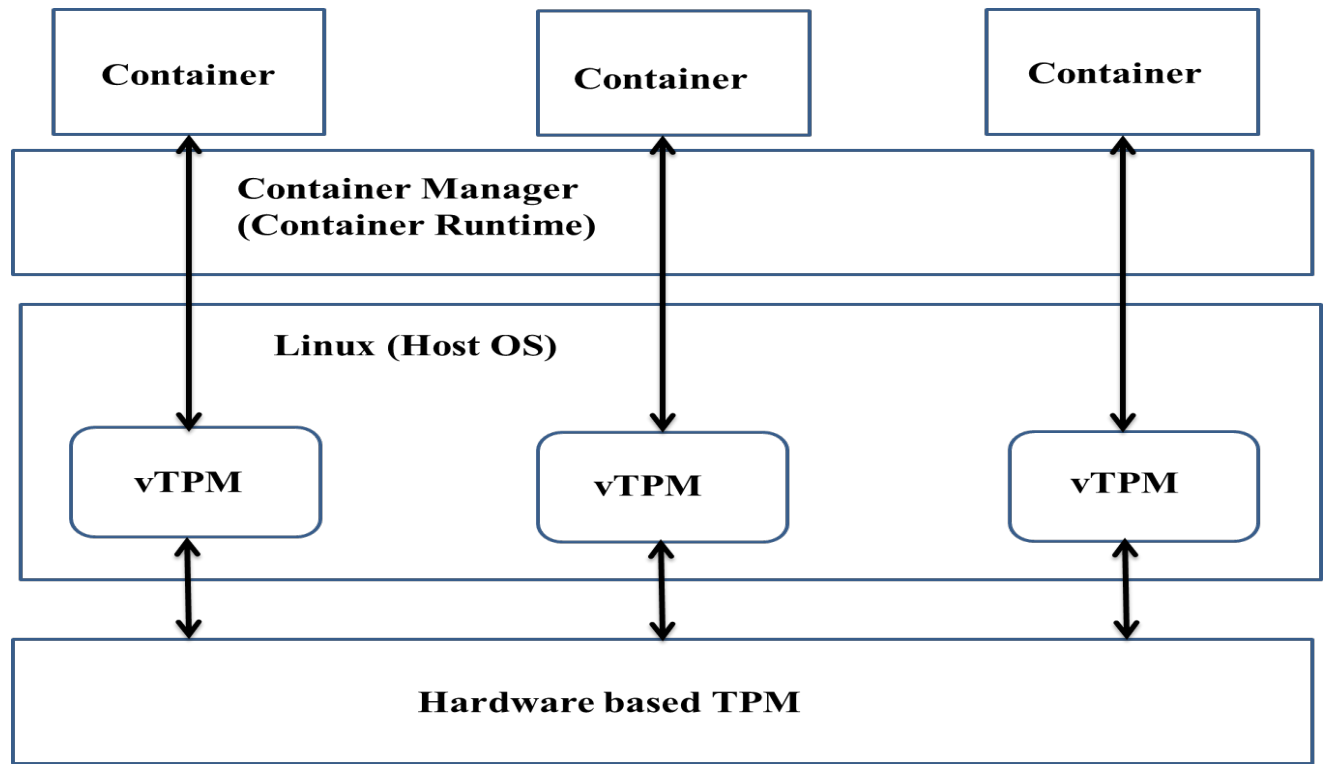
384 Building a TPM architecture is not the only type of approach for providing trust rooted in
385 hardware for the container stack. Another type of approach that has been proposed is to leverage
386 the trusted execution support of some CPU architectures to protect processes running in a
387 container against attacks from sources inside the same container stack. This includes privileged
388 software in the same stack such as the container runtime and host OS kernel [8]. A mechanism or
389 security solution based on this type of approach is discussed and analyzed in section 3.3.

390 3.1 vTPM in the host OS Kernel – Security Assurance Requirements

391 In an architectural approach suggested in [7], a software-based module called vTPM (virtual
392 TPM) is placed into the OS kernel. To make this module available to several containers, it needs
393 to be virtualized. This is accomplished using a kernel module that provides an arbitrary number
394 of software-based vTPMs, which are exposed to containers through the usual mechanisms and
395 present a character device type interface to the container userspace. This functionality can be
396 implemented by having the container runtime (or container manager) ask the host OS kernel to
397 create a new vTPM and assign the virtual device to a container. The vTPMs are linked to the
398 TPM implemented in the hardware platform (referred to as “physical TPM”) that hosts the
399 container stack. The schematic diagram of this architectural approach is illustrated in figure 2.

400 The security assurance requirements for the above discussed architectural approach can be
401 looked at for the following scenarios:

402 The host OS is completely trusted: The trust-in-host OS can be established by extending the root
403 of trust from the hardware using the hardware-based, or physical, TPM. Since the host OS is
404 trusted to prevent unauthorized access by containers and processes, it can also be trusted to
405 prevent unauthorized access to the in-kernel vTPM. Moreover, there is the assurance that
406 containers cannot modify the host kernel by loading new modules or by exploiting vulnerabilities
407 in the kernel. Containers can therefore reliably attest to their own state by using the hash extend
408 feature of the vTPM.



409

410

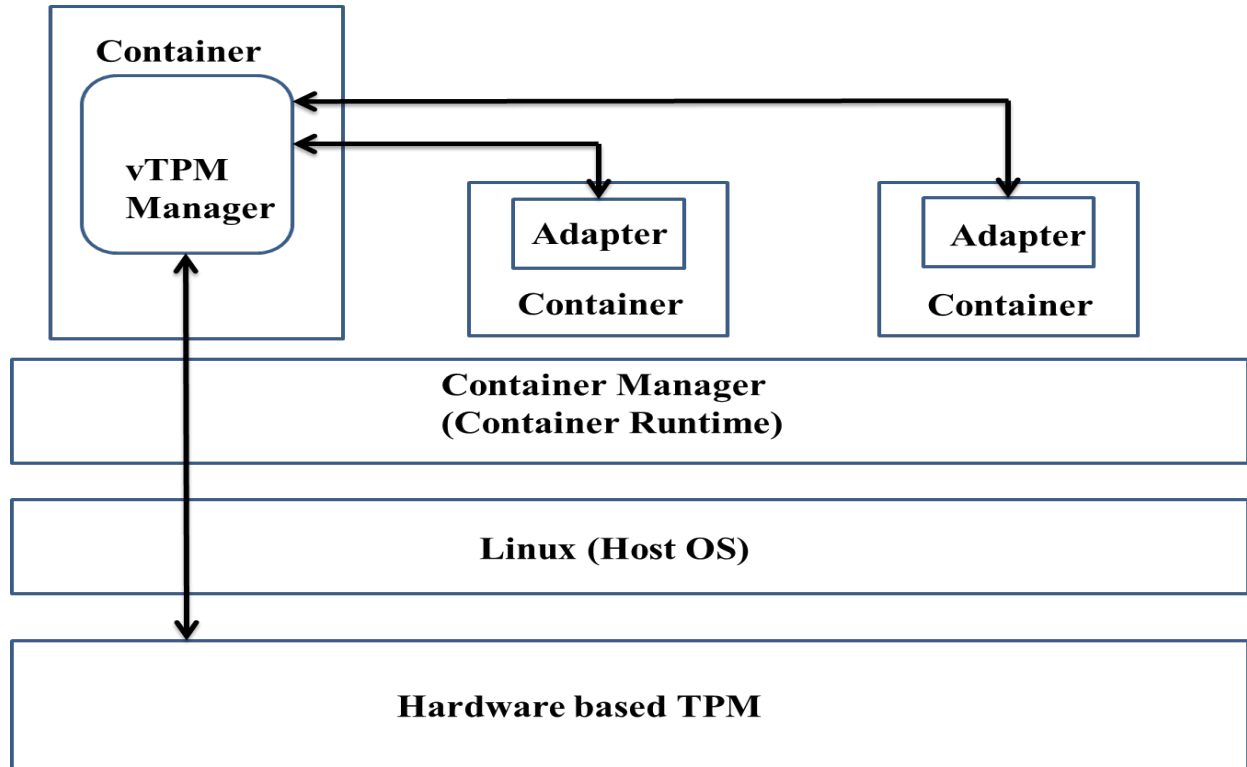
Figure 2 – vTPM Implemented in a Kernel Module

411 The host OS is not completely trusted, and independent trust is needed on vTPM: To implement
 412 trust on vTPM, a scheme using the same mechanism used for establishing hardware TPM
 413 (physical TPM) trust has been referred to in [7]. In the physical TPM, the hardware platform
 414 provider signs an endorsement key (EK) stating that the TPM is trustworthy. This is then
 415 extended by giving each vTPM instance its own endorsement key and deploying protocols for
 416 signing the endorsement keys of vTPMs using of the hardware-based TPM.

417 3.2 vTPM in a dedicated Container – Security Assurance Requirements

418 The software-based vTPM with the same functionality described in section 3.1 is built and
 419 hosted in a dedicated container (referred to as vTPM management container). The schematic
 420 diagram of this architectural approach is given in figure 3. This vTPM has two primary features:

- 421 (a) Access to hardware-based (physical) TPM
- 422 (b) Exposes the vTPM interface to other containers through a communication channel, which
 423 can be a local UNIX domain socket or another IPC mechanism. If the IPC mechanism is
 424 employed, the container using the vTPM service requires an additional piece of software
 425 (denoted as “adapter” in figure 3) that presents the IPC interface as a standard character
 426 device. In the container that is hosting the vTPM, a daemon will process requests from
 427 other containers instead of a kernel module as it was in the previous case.



428

429

Figure 3 – vTPM located in a dedicated Container

430 The security assurance provided by this architectural approach is the same as the one provided
 431 by the host OS in the container stack. A host OS, such as Linux, provides isolation between
 432 processes belonging to different containers through the Namespaces feature. If this functionality
 433 works correctly, no process belonging to a different container can access the state of the vTPM
 434 deployed in a dedicated container. In other words, the security of this implementation is
 435 jeopardized only in the event of a container escape attack. Still, this approach provides less
 436 protection than the approach in section 3.1 (vTPM in the host Kernel) since the kernel is more
 437 reliable in limiting the kind of access it exposes to the userspace.

438 3.3 Leveraging Trusted Execution Support of Hardware

439 In 2015, Intel released the Software Guard eXtensions (SGX) [8] for their CPUs, which provided
 440 the hardware mechanism for protecting user-level software from privileged system software
 441 using the concept of secure enclaves. An enclave page cache (EPC) is a region of protected
 442 physical memory where application code and data reside and are protected by CPU access
 443 controls. When code and data in EPC pages are moved to DRAM, they are instantaneously
 444 encrypted using an on-chip memory encryption engine (MEE) and then decrypted when they are
 445 transferred from DRAM to EPC pages. The integrity of the enclave memory itself is also
 446 protected by mechanisms that detect memory modifications and rollbacks. Thus, enclaves are
 447 trusted execution environments provided by SGX to applications residing in the container.

4 Assurance Requirements for Host OS Protection

4.1 Requirements for Generic Host OS Protection

Installing a container-specific OS (as opposed to a generic OS distribution), keeping OS versions up-to-date and patched, logging features that can track anomalous accesses to the OS, and executing privileged operations form the crux of Host OS countermeasures in the *Container Security Guide*. In addition to the above countermeasures, it is also a good OS security practice to disable all unused interfaces (Serial or Proprietary) on the host and minimize the user and administrative accounts and groups. In addition to these, there are Linux-specific patches, such as grsecurity [9] and PaX [10], that are available for Linux distributions. All measures combined should provide the following security assurance for the host OS:

- (a) Prevent manipulation of program execution by modifying memory (e.g., buffer overflow attacks)
- (b) Prevent attempts to reroute code to existing procedures (e.g., system calls in common libraries)

4.2 Assurance Requirements for Host OS Protection for Container Escape

The host OS should be protected to mitigate threats that result from container escape or breakout, and all containers should be protected from other containers on the host. There are many solutions available in Linux environments that enable these protections, but the three solutions analyzed in this document are SELinux, AppArmor, and Seccomp, all of which utilize kernel-loadable modules (referred to using the acronym LKM, or Linux Kernel Module). SELinux, or Security Enhanced Linux, can be used to assign labels (e.g., type) and categories to processes and objects (e.g., files, sockets) and specify access restrictions (Mandatory Access Control or MAC) between resources belonging to certain combinations of labels and categories. For example, a specific SELinux label can be applied to a container to enforce a security policy (e.g., a container hosting a Webserver can only open ports 80 or 443) [6]. AppArmor is another LKM product that helps enforce mandatory access control policies by applying profiles to processes that enable restriction of privileges they have at the level of Linux capabilities and file access. The controls are thus data-centric and are at a coarser level of granularity compared to SELinux. SECure COMPUting (Seccomp) is a module that can define and enforce an access control method that enables specification of the number of system calls available for an application within a container to interface with the kernel. Limiting system calls provides a restricted execution environment and thus reduces the kernel attack surface. The allowed list (i.e., whitelist) and prohibited list (i.e., blacklist) of system calls for a process are set up using the syscall filter [11].

The overall goal of the kernel-loadable modules, or LKMs, described above is to provide another level of security checks on the access rights of processes and users beyond that provided by the standard file-level access control (discretionary access control, or DAC) in Linux [6]. This goal then drives the following security assurance requirements that need to be satisfied:

- 487 (a) A user authorized to run applications in the container should not be allowed access to the
488 above described kernel-loadable modules
- 489 (b) If using SELinux, the chcon utility used to label the files and parent folders should be
490 used at the correct levels in the file system hierarchy such that it results in least privileges
- 491 (c) If using Seccomp, both a syscall whitelist (a list of allowable calls) and a syscall blacklist
492 (a list of prohibited calls) should be generated. The choice of syscalls in the whitelist for
493 a container should be based on type of application(s) hosted in the container, deployment
494 situation, and container size. The syscalls included in the blacklist are for high risk,
495 possibly vulnerable, known dangerous, and explicitly disallowed ones [11]. Some
496 examples in this category include syscalls that allow for loading kernel modules,
497 rebooting, triggering mount operations, and other administrative calls.
- 498 (d) If using Seccomp, the sandboxes created by seccomp filters must not allow the use of the
499 ptrace command. If ptrace is allowed, the tracer can modify the process's system call to
500 bypass the filter and therefore call blocked or restricted system calls.
- 501 (e) It should be possible to create container-specific profiles using a combination of
502 configuration options provided by these security modules
- 503 (f) A minimal configuration feature that should be available is one that allows for the
504 partitioning of containers in the host to different security domains
- 505 (g) It should prevent containers' ability to mount/remount sensitive directories and/or
506 specific system directories critical to security enforcement (Cgroups, procfs, sysfs)
- 507 (h) It should be possible to create a security profile for the administrators of container
508 runtime using a combination of the above features

5 Assurance Requirements for Container Runtime Configuration

As already described in section 2.5, all security configuration parameters for containers, except for those dealing with cluster management and scheduling, are set using APIs provided by container runtime. Although most of them involve Linux kernel features (Namespaces, Cgroups, Capabilities) and Linux kernel modules, these tasks have been included under this section since they are performed by the container runtime making syscalls to Linux host OS interfaces. The overall organization of this section is as follows:

- (a) Section 5.2 discusses configurations involving Linux's Namespace feature, which provides isolation for various resources
- (b) Section 5.3 discusses configurations using the Cgroups feature, which is primarily utilized for setting resource limits and thus preventing denial of service attacks
- (c) Section 5.4 discusses configurations using the Capabilities feature, which enables the allocation of least privileges
- (d) Section 5.5 discusses the configurations for device isolation, which can be enabled using a combination of Cgroups and kernel-loadable MAC enforcement modules
- (e) Section 5.6 discusses configuration parameters that can be set at the time of launching the containers rather than being pre-configured using the functions discussed above

Before analyzing these functions, the need for a configuration feature for the container runtime itself is outlined in section 5.1.

5.1 Requirements for Secure Connection

Container runtime module is implemented with a daemon that listens through a Unix socket and thus enables remote administration of the runtime. It is possible under certain circumstances for members in the administrative group to change the Unix socket to a TCP socket [10]. Any connection to this TCP socket can allow attackers to pull and run any container in privileged mode, thereby giving them root access to the host. The security assurance requirement for the TLS connection involves the encryption and authentication of both sides (container runtime module as well as the client tool used for remote administration) of the connection before establishing the TLS session.

5.2 Requirements for Isolation-based Configurations

5.2.1 Process Isolation for Containers

Process Isolation is a core security requirement for containers to ensure the integrity of various applications running in different containers as well as in the host. A process isolation mechanism in a container environment should meet the following requirements [4]:

- (a) Ability to distinguish processes running in different containers from each other and from those running on the host
- (b) Limit cross-container process visibility

- 545 (c) Prevent certain type of attacks such as:
- 546 (i.) A process running in one container influencing a process running in another
547 container using interfaces provided by the OS for process management (e.g.,
548 signals and interrupts)
- 549 (ii.) A process running in one container and directly accessing the memory of a
550 process running in another container by using special system calls (e.g., the
551 ptrace() allows a debugger process to attach and monitor the memory of a
552 debugged process)

553 To provide process isolation, a Linux kernel feature called process id (PID) namespace is used.
554 A PID namespace is a mechanism that groups processes and controls their ability to see (e.g., via
555 proc pseudo-filesystem) and interact (e.g., sending signals) with one another. A PID namespace
556 is created using clone() or unshare() system call and is associated with one or more containers.
557 The first process carries the id PID1, and the identifiers for subsequent processes increases
558 sequentially. Thus, the PID namespaces feature also provides PID virtualization. Two processes
559 in different PID namespaces can have the same PID.

560 5.2.2 Filesystem Isolation for Containers

561 The goal of filesystem isolation is to prevent illegitimate access to filesystem objects from one
562 container to another and from any container to the host. The filesystem is an OS interface that
563 allows processes to store and share data as well as interact with one another. Access to data for a
564 container application is determined by its access to file systems through the filesystem mount
565 points. Therefore, access to data can be restricted by making the list of filesystem mount points
566 visible and accessible to a container application. This is accomplished through the mount
567 namespace. First, a named mount namespace is created along with a set of file system mount
568 points. This mount namespace is then associated with a process that can only see and issue
569 system calls such as mount () or unmount () on those mount points. It also operates on files that
570 are within that mount namespace and accessible through those mount points. The following are
571 the security solutions for filesystem isolation and their limitations:

- 572 (a) All Linux-based OS virtualization solutions utilize a *mount namespace* that allows for the
573 separation of mounts between the containers and the host. This is intended to facilitate
574 customization of the environment visible to users and processes. This feature does not
575 guarantee data isolation between the containers since containers inherit the view of
576 filesystem mounts from their parent and can access all parts of the filesystem even though
577 each container is created within a new mount namespace.
- 578 (b) The typical solution for process filesystem access containment is by using the chroot ()
579 system call, which binds a process to a subtree of the filesystem hierarchy. This allows a
580 container to share resources with the host by mounting them within the subtree visible
581 inside the container. However, this feature cannot provide the requisite protection in the
582 presence of privileged processes (i.e., processes with the CAP_SYS_CHROOT
583 privilege), which can escape the chroot jail due to the fact that the chroot () system call
584 only affects the pathname resolution.
- 585 (c) A better protection for filesystem objects is provided by modifying the root filesystem for
586 processes in a container as opposed to just modifying the root directory (which the chroot

587 () system call enables) [4]. This is enabled by the `pivot_root ()` call, which moves the
588 mountpoint of the old root filesystem to a directory under the new root filesystem and
589 puts the new root filesystem in its place. This provides filesystem level protection since
590 the old root filesystem can be unmounted when it is carried out inside the mount
591 namespace of the container, thus rendering the host root filesystem inaccessible for
592 processes inside the container.

593 (d) Another filesystem-level protection strategy is to disallow mounting and unmounting of
594 filesystems for processes running inside a jail by default and enforce granular control of
595 this privilege using options in the `allow_mount*` command.

596 (e) Another mechanism to strengthen filesystem isolation is to designate a separate user
597 namespace per container, which maps the user and group ids to a lesser privileged range
598 of host UIDs and groups.

599 Because of the limitation of each of the above security solutions, the assurance requirements for
600 total filesystem-level protection involves a combination of configurations including mount
601 namespace, `chroot`, `pivot_root`, and user namespace needed for:

- 602 • Isolating mount points by mount namespace
- 603 • Changing the root directory for each process using `chroot`
- 604 • Changing the root filesystem visible to each process (container) using `pivot_root`
- 605 • Restricting user access scope using user namespace

606

607 **5.2.3 IPC Isolation for Containers**

608 Inter-process communication (IPC) isolation for containers means that processes in a container
609 must be restricted to communicate via certain IPC primitives only within that same container. An
610 IPC object (or associated mechanism) can be either a filesystem-based IPC object or non-
611 filesystem-based. Filesystem-based IPC objects, such as domain sockets and named pipes, can be
612 isolated using a combination of mount namespace and `pivot_root` features (section 5.2.2 above)
613 since they prevent processes from accessing filesystem paths outside of their own container.

614 However, there are other IPC objects such as System V IPC objects, semaphore sets (arrays),
615 shared memory segments, and message queues. These IPC objects can be isolated in Linux with
616 the help of IPC namespaces that allow the creation of a completely disjointed set of IPC objects.
617 Each IPC namespace has its own set of System V IPC identifiers and its own POSIX message
618 queue filesystem. Objects created in an IPC namespace are visible to all other processes that are
619 members of that namespace but are not visible to processes in other IPC namespaces. IPC objects
620 accessible for a process can be listed using the `ipcs` command and removed using the `ipcrm`
621 command.

622 **5.2.4 Network Isolation for Containers**

623 Network level isolation for containers is provided through the network namespace feature. For
624 each network namespace that is created, a set of network devices, IP addresses, IP routing
625 tables, `/proc/net` directory, and port numbers can be associated with it. Each container can have
626 its own virtual network device and applications that bind to the per-namespace port number

627 space. Suitable routing rules in the host system can direct network packets to the network device
628 associated with a specific container. It is therefore possible to have, for example, multiple
629 containerized web servers on the same host system with each server bound to port 80 in its (per-
630 container) network namespace.

631 Network connectivity is a core requirement for all production grade applications running on
632 containers such as web apps and multi-tier apps. The containers can be connected using a logical
633 IP network called the overlay network. The typical network configuration on a container
634 platform (consisting of containers, container runtime, host OS and the physical host) involves
635 creating a network bridge on the container host. Each container on a host is connected to that
636 bridge. A router captures Ethernet packets from its bridge-connected interface in promiscuous
637 mode, and captured packets are forwarded over the user datagram protocol (UDP) to router peers
638 running on other container hosts. These UDP "connections" are duplex, can traverse firewalls,
639 and can be encrypted [12]. Each container is connected to the bridge using a layer 2 (link layer)
640 virtualized network interface (VNI) with a valid Link Layer address or a Network Address
641 Translation (NAT) for layer 3 connectivity. The Linux Layer 2 network isolation is based on the
642 concept of Network Namespace, which allows for the creation of several networking stacks that
643 provide a view of being completely independent of the containers [4].

644 The simplest configuration for network isolation using layer 2 VNI involves defining a pair of
645 virtually linked Ethernet (veth) interfaces. One of the interfaces is assigned to the same network
646 namespace as the container and the other to the host namespace. A virtual link is then established
647 between the two interfaces, thus connecting the container to physical networks. There are two
648 options for enabling this link [4]:

- 649 (a) Network Bridge Device: The veth interface and the host physical interface are connected
650 using a virtual network bridge device. In this option, all container and host interfaces are
651 attached to the same link layer bridge and thus receive all link layer traffic on the bridge.
- 652 (b) Routing Tables: Another option is to utilize routing tables to forward the traffic between
653 virtual network interface (to which the container is connected) and physical network
654 interfaces (resident at the host). In this option, containers can communicate with each
655 other only when a network route is explicitly provided.

656 Security Analysis: The network isolation functionality provided by these two options forces a
657 container process to use a designated virtual network segment or a designated network route
658 (e.g., over a VPN connection). Between the two options, the routing table use presents a slightly
659 higher security assurance than the network bridge device solution since the latter allows a
660 container address to be visible to all containers connected to the bridge.

661 Another approach to provide network connectivity for containers is to use the MACVLAN
662 interface [13], which also allows each container to have its own separate link layer address. The
663 Virtual Ethernet Port Aggregator (VEPA) is the most widely used mode for configuring this
664 option for isolating the containers. However, complete assurance of network isolation can be
665 provided at the process level in containers only if the namespace-based approaches are
666 augmented with mandatory access controls (MAC) and the isolation of the process from other
667 global namespaces.

668 **5.2.5 User and Group-level Isolation for Containers**

669 Some processes may need some subset of root privileges. The user namespaces feature can be
 670 used to restrict the privileges of some user IDs to that needed subset. The user namespace
 671 isolates the user and group ID number spaces. In other words, a process's user and group IDs can
 672 be different inside and outside of a user namespace. The most interesting case here is that a
 673 process can have a normal unprivileged user ID outside of a user namespace while at the same
 674 time having a user ID 0 inside of the namespace. This means that the process has full root
 675 privileges for operations inside the user namespace, but is unprivileged for operations outside the
 676 namespace.

677 Starting in Linux 3.8, unprivileged processes can create user namespaces, which opens a raft of
 678 interesting new possibilities for applications. Since an otherwise unprivileged process can hold
 679 root privileges inside the user namespace, unprivileged applications now have access to
 680 functionality that was formerly limited to root [4].

681 **5.3 Requirements for Resource Limiting Solutions**

682 The primary protection mechanism for denial-of-service attacks in Linux container environments
 683 is the Cgroups feature that enables setting limits for various resources. The “limits” specification
 684 feature is restricted not only to hardware artifacts such as CPU, memory, and storage, but also to
 685 processes and tasks. In addition to the limits feature, Cgroups enables the designation of a
 686 collection of potential “resource hogging tasks” that can be frozen by sending a SIGSTOP signal.
 687 It can later be unfrozen by sending a SIGCONT signal [11].

688 In addition to its main role of preventing against denial-of-service attacks, the Cgroups feature
 689 also provides marginal network-level protection with a method (using network classifier Cgroup)
 690 that tags network packets with a “classid” value. This can then be used as a parameter for
 691 filtering certain packets. (The classid value can also be used for priority handling based on
 692 Quality of Service (QoS) requirements, though that feature falls under performance enhancement
 693 and not strictly security.)

694 The following table provides the list of hardware resources for which the Cgroups feature either
 695 enables setting up of resource limits or access control.

696 **Table 1– Linux Resource Control using Cgroups**

Resource	“Limit” Feature or Access Control
CPU	Specific number of CPUs or amount of “CPU Shares” for a group of processes
Memory	“Hard” and “Soft” memory allocation units for a group of processes
BLKIO	Set disk read or write speeds, operations per second, queue controls, and wait times on block devices designated by major and minor numbers; provides more granular access control compared to filesystem specific controls

Devices	Create a whitelist for devices based on either: (a) Type (character vs block) or (b) Major and Minor numbers
---------	--

697

698 Cgroups configuration should provide the following assurances:

699 (a) It should not expose container host information, such as the kernel ring buffer via `dmesg`,
700 which can assist in kernel exploitation or information leaks

701 (b) It should not allow local disk access, even within user namespaces and mount restricted
702 namespaces via `raw disk`, `device`, or `mknod` access [11].

703

704 5.4 Requirements for Least Privilege Configuration for Containers

705 As already mentioned, the Capabilities feature in Linux can be used to partition the set of root
706 privileges. All container runtime products, such as LXC, Docker, and CoreOS Rkt, come with a
707 default capability profile where some capabilities for containers are enabled and some are
708 disabled [11]. Due to the privilege needs of the application running in the container, some of the
709 defaults have be modified (i.e., some capabilities that have been enabled by default need to be
710 disabled, and some capabilities disabled by default need to be enabled). However, for most
711 applications hosted in containers, the following assurance requirements must be satisfied while
712 configuring the Capabilities feature in Linux:

713 (a) Capabilities that provide the privilege to manipulate a non-name spaced kernel parameter
714 (e.g., `Sys Time`) will have the effect of that parameter modified not only for the container
715 but also for the host and for all other containers. Hence such capabilities (e.g.,
716 `CAP_SYS_TIME`) should not be enabled.

717 (b) Capabilities that provide the broad set of privileges almost equal to that of the root should
718 not be enabled (e.g., `CAP_SYS_ADMIN`).

719 (c) There is no need to enable the capability `CAP_SYS_MODULE`, which allows for the
720 loading and unloading of kernel modules as this will lead to insecure privilege escalation.

721 (d) The Capabilities feature should always be used in conjunction with user namespace as
722 any privilege escalation to the process due to enabling some Capabilities by error will be
723 limited to the namespace.

724

725 5.5 Requirements for Device Isolation solutions

726 In Linux, access to devices is enabled by device nodes, which are special files that provide an
727 interface to the host device drivers. Device nodes are separated from the rest of the filesystem,
728 and their nodes are placed in the `/dev` directory. These nodes are not namespace-aware. The
729 creation of device nodes is performed by the `udev` daemon process issuing the `mknod` system
730 call. The permission for a process to create device nodes (for accessing block or character
731 devices) is provided by the `CAP_SYS_MKNOD` capability. Containers are given access to device
732 nodes if the corresponding devices are to be shared among containers or between different

733 containers and the host. However, device nodes are security-sensitive since they provide
734 interfaces to device drivers. These drivers present significant attack vectors because they expose
735 interfaces (particularly the storage interface) to code running in the kernel space, which may be
736 abused to gain illegitimate data access, escalate privileges, or mount other attacks.

737 One possible solution for providing device-level isolation between containers is the use of
738 “device namespace,” provided the referenced input/output (physical) devices are namespace-
739 aware. Unfortunately, many Linux kernel distributions do not support the device namespace
740 feature. Where available, this feature can be used to create virtual devices for each container,
741 which can be multiplexed for access to a physical host device. Further, when Linux device
742 drivers controlling physical devices are not namespace-aware and the devices assume only one
743 controlling master host, access privileges for them are hard to securely grant for unprivileged
744 containers unless the device is used exclusively by a single container.

745 In the absence of the device namespace feature, two features are utilized for controlling access to
746 devices for containers. They are: (a) control groups, or Cgroups; and (b) Mandatory Access
747 Control (MAC) enforcement. The Cgroups subsystem for devices is used to create a whitelist,
748 formatted for devices based on type (i.e., character vs block) and device major and minor
749 numbers. The wild card “all” applies to all device types and major and minor numbers, and it is
750 typically used as a default deny before whitelisting explicit devices [11].

751 There are two MAC enforcement methods available in Linux environments: Security-Enhanced
752 Linux (SELinux) and Apparmor. In SELinux, Multi-Level Security (MLS) labels or
753 classification labels are applied on processes and data/devices, and the system applies fine-
754 grained policy and type enforcement across the different labels. AppArmor is another MAC
755 system that offers a pathname-based access control (as opposed to filesystem nodes within
756 SELinux). The restrictions can be aggregated to define a profile for a specific application,
757 process, or container. A common weakness for all MAC systems is that the controls it provides
758 can be subverted through direct execution of system calls.

759 The assurance requirements for device isolation solutions therefore are:

- 760 (a) All containers must be prevented from creating new device nodes, and the
761 CAP_SYS_MKNOD capability should not be enabled for them
- 762 (b) All mountpoints inside containers should have the nodev flag set to prevent them from
763 being used to create files to access device drivers
- 764 (c) All containers should only be allowed to access the following set of devices since they
765 are characterized as safe [4] due to observations given below:
- 766 • *Purely virtual devices* – such as pseudo-terminals and virtual network interfaces; the
767 security guarantee comes from the fact that these devices are explicitly created for
768 each container and not shared
 - 769 • *Stateless devices* – such as random, null, and others; sharing these devices among all
770 containers and the host is safe because they are stateless
 - 771 • *User namespace-aware devices* – if the device (through the device driver code)
772 supports verifying capabilities of the process in the corresponding user namespace,

773 then such a device can be safely exposed to a container since the specified restrictions
 774 will be enforced

775 (d) When Cgroups and MAC enforcement systems are both used for controlling access to
 776 devices, care should be taken to ensure that their respective rules do not create conflict.

777

778 **5.6 Requirements for Container Launching Options**

779 Every container runtime product has a command to launch containers that carry many options.
 780 The assurance requirements associated with the secure use of this command are stated as a set of
 781 options that should be avoided [4]. As a best security practice, containers should not use options
 782 that will enable sharing any namespaces associated with the container host when launched [11].
 783 If this is not the case, it may not only enable the container to view the resources/objects
 784 associated with that namespace but also manipulate those resources/objects by subverting the
 785 isolation provided by static configuration of namespaces for the container. The following table
 786 provides the list of namespaces for which sharing the corresponding host counterpart should not
 787 be used in the container launch options.

788

Table 2 – Prohibited Options for Container Launching

Namespace/ Example Resource-Object	Brief Description	Security Threat
UTS	All containers are assigned their own UTS namespace and thus have no need to know the UTS namespace of the host	Processes within the container can see and manipulate the hostname and domain of the host
IPC/ Shared Memory Segment	Shared Memory segments for inter-process communication between application modules are set up for faster communication as they are faster than REST API calls	Processes within the container can see and manipulate host IPC object
Filesystem	Host-sensitive directories should not be mounted in read-write mode as container volumes	Gives containers the ability to modify the files in those directories with a potential to jeopardize host security
Setting net=host in the container launching command	The networking mode for the container should not be set equal to host	This will give privileges to a container that only a host should have (e.g., shutting itself down) or access to networking services that only

Namespace/ Example Resource-Object	Brief Description	Security Threat
		the host needs
Publishing container ports to the host	This is done for setting up communication to and from that container	The default option of publishing to all interfaces should not be used; by specifying the interface that the port should bind to explicitly, traffic into and from the container is restricted to the given interface
Inter-container communication	If it exists, the option to enable blanket inter-container communication must not be enabled; instead, explicit communication channels must be set up between two containers that need to communicate.	Any compromised container can attack any other container on the host

789

790 In addition to container launch options that involve objects shared with the host, there are some
 791 parameters exclusively applicable to the container that should be set when launching containers.

792 (a) Containers should always be launched with a specific memory limit to prevent denial-of-
 793 service attacks or certain applications leaking memory that may eventually consume all
 794 the memory on the host

795 (b) Containers should always be launched by specifying the number of CPU shares. The
 796 default value (Total CPU/number of containers) may not be sufficient for some
 797 containers, resulting in denial of service. The number of CPU shares assigned to a
 798 container should be such that no container can starve others with default settings. Further,
 799 if there exists a group of containers that dominate others in CPU usage, then a lower
 800 default value should be assigned to containers in that group to ensure fair distribution of
 801 CPU shares.

802 (c) If the host OS Linux distribution supports a MAC system (e.g., SELinux), a policy
 803 template should be set up, the container engine should be started with an option to
 804 recognize the template, and the container launching API should have an option to
 805 recognize the policy template parameter and include it as part of the launch parameter.

806 (d) Containers should be launched only with “required” capabilities by initially dropping all
807 capabilities and then adding only the required ones. The following capabilities in general
808 should not be present (i.e., NET_ADMIN, SYS_ADMIN, SYS_MODULE).

809 **6 Assurance Requirements for Image Integrity Solutions**

810 The integrity of the container images is of paramount importance since they are converted to
811 running instances, some of which may host mission-critical applications. The image
812 countermeasures covered in the *Container Security Guide* include recommendations for
813 monitoring images for malware and other vulnerabilities, proper image configuration, separating
814 secrets from image files, and ensuring trust in images through cryptographic signatures and
815 regular updates. The security solutions needed for carrying out these recommendations should
816 include the following assurance requirements:

- 817 (a) There should exist a means to create metadata linking each image to its base image
- 818 (b) There should exist a feature to rebuild the image automatically if the linked base image
819 changes [6]
- 820 (c) When any changes are made to the base image or dependent image (e.g., patching a
821 vulnerability), changes should not be made to the running containers. Instead, the
822 corresponding image should be recreated and the container re-launched using the
823 modified image. Thus, a single master, or golden image, is to be maintained for any
824 service.
- 825 (d) When employing “image signing” solutions for digitally signing and uniquely identifying
826 each image, the following requirements should be met [6]:
 - 827 1. There should be a robust key management to minimize the possibility of key
828 compromise. One approach is to have a PKI system that issues a certificate to each
829 developer exclusively for signing the image. The private key associated with this
830 certificate will then be the “signing key” that is used to sign all container images in a
831 repository.
 - 832 2. Replay attacks must be mitigated by embedding expiration timestamps in signed
833 container images. Alternatively, a special key can be used to sign the metadata for
834 the repository, ensuring that the images in the repository do not contain stale
835 versions of the image with valid signatures.
- 836 (e) In addition to creating a unique identifier for an image using digital signatures, the
837 integrity of individual components of the image can be ensured by using labels such as
838 key/value pairs for each component.
- 839 (f) Images should be built such that the application(s) in them are not used for any privilege
840 escalation attacks. This can be achieved by disabling the `chmod a-s` command, which
841 removes the `suid` bit, or removing `setuid` and `setgid` binaries in them [6].

7 Assurance Requirements for Image Registry Protection

843 The suggested registry countermeasures in *Container Security Guide* include developing secure
844 connections to registries and ensuring that they do not contain out-of-date vulnerable images by
845 pruning them out through an automated process or controlling their accidental deployment
846 through use of discrete version numbers. Some assurance requirements unrelated to these
847 countermeasures yet still critical to processes involving creating, posting, and removing images
848 into and from registries are:

- 849 (a) The number of accounts accessing the registry must be limited since the common threat
850 in some environments is account hijacking when a diverse set of clients has access to a
851 container registry. One such environment is the registry maintained by cloud service
852 providers who offer container services.
- 853 (b) The permission to create container image registries and add or remove content to
854 registries must be cryptographically protected.

8 Assurance Requirements for Orchestration Functions

856 The use of an Orchestration platform (consisting of a suite of tools) in a containerized
857 infrastructure is intended to perform the following functions:

- 858 • Enable the definition of a cluster (a named group of container hosts that can be managed as a
859 single entity) and schedule containers into the cluster. The cluster configuration should
860 support specification of parameters such as the amount of CPU/Memory to reserve, the
861 number of replicas (i.e., duplicate copies of same container to be run), and the circumstances
862 under which a container should continue to run or be taken offline.
- 863 • Enable automated deployment of containers in various clusters/hosts (container scheduling).
864 This is achieved by integrating various automation tools to execute automation scripts as part
865 of an orchestrated workflow and to obtain feedback and status results for those automation
866 tasks. This kind of integration depends on the interfaces that the automation tools provide
867 and the type of formats (open or closed) that they follow [14].
- 868 • Provisioning, or defining new container hosts and attaching them to existing clusters

869 The suggested orchestration countermeasures in the *Container Security Guide* include granular
870 access control of administrative actions based on hosts, containers and images as parameters, use
871 of enterprise-grade authentication services using strong credentials and directories, and isolating
872 containers to separate hosts based on the sensitivity level of the applications running in them. In
873 addition to these countermeasures, the orchestration artifacts should satisfy the following
874 security assurance requirements:

- 875 (a) Clusters should have capabilities for logging and monitoring the resource consumption
876 patterns of individual containers to avoid unanticipated spikes in resource usage leading
877 to non-availability of critical resources
- 878 (b) The Orchestration platform must be usable on containerized infrastructures with more
879 than one host OS. In other words, the orchestration tools used must be container-host OS-
880 neutral. Using different tools for different container host OS platforms increases the
881 probability of denial-of-service attacks in those environments since the enterprise is not
882 able to obtain a global picture of resource usage for all running containers in the entire
883 containerized infrastructure of the enterprise.

884 **9 Adverse Side effect of some Security Solutions**

885 While discussing a security solution (e.g., using mount namespace) in the context of a security
886 objective (i.e., filesystem isolation), certain augmenting solutions are recommended since the
887 solution under discussion cannot meet the objective by itself. However, there are some security
888 solutions that, irrespective of any augmenting controls, impose certain limitations on the
889 functionality and performance of certain container functions. Despite their direct impact
890 affecting only functional and performance aspects, they may have an indirect impact on certain
891 security parameters. These are discussed below.

892 **9.1 Resource Limiting using Cgroups**

893 The use of Cgroups to limit resource access for processes/containers is included as a security
894 solution because of its potential to mitigate the chances of denial-of-service situations. The Linux
895 control groups (Cgroups) subsystem is used to group processes and manage their aggregate
896 resource consumption. It is commonly used to limit the memory and CPU consumption of
897 containers. A container can be resized by simply changing the limits of its corresponding
898 Cgroup. However, processes running inside a container are not aware of their resource limits [2].
899 For example, a process can see all the CPUs in the system even if it is only allowed to run on a
900 subset of them; the same applies to memory. If an application attempts to automatically tune
901 itself by allocating resources based on the total system resources available, it may over-allocate
902 when running in a resource-constrained container, thus resulting in denial-of-service to other
903 applications within the same container [2].

904 **9.2 Syscall filters using Seccomp**

905 Setting up system call filters (with whitelist and blacklist) using Seccomp is used as a security
906 solution since system calls are not namespace-aware (ruling out the use of the namespaces
907 feature), though in the presence of malicious processes, this can introduce accidental leakage
908 between containers. However, the choice of system calls to be allowed is based on a current set
909 of applications in the container, and this security solution has the potential to introduce
910 application incompatibility since applications can be migrated between containers for load-
911 balancing reasons.

912 **10 Summary and Conclusions**

913 The security solutions analyzed in this document can be summarized as follows:

- 914 (a) Providing authenticity and attestation of integrity for software components of a container
915 stack such as Linux (Host OS), container runtime, and the containers using hardware-
916 based root-of-trust solutions such as TPM and vTPM
- 917 (b) Hardware-based protection for shielding one container from another as well as containers
918 from higher privileged software, such as Linux kernel, using the safe execution model
919 provided by hardware architecture (e.g., Intel SGX)
- 920 (c) Linux kernel features (Namespaces, Cgroups, Capabilities) and loadable kernel module
921 (LKM) features for protection of the Linux kernel itself and for protecting one container
922 from another
- 923 (d) Protection measures for container runtime, container images, container registry, and
924 container orchestration tools.

925 The conclusion from the analysis is that every security solution must satisfy some security
926 assurance requirements to effectively provide necessary and sufficient security guarantees.

927 **Appendix A—Acronyms**

928 Selected acronyms and abbreviations used in this paper are defined below.

EPC	Enclave Page Cache
IPC	Inter-process Communication
MAC	Mandatory Access Control
MEE	Memory Encryption Engine
NAT	Network Address Translation
PID	Process ID
PKI	Public Key Infrastructure
SGX	Software Guard eXtensions
TPM	Trusted Platform Module
UDP	User Datagram Protocol
UTS	UNIX Timesharing System
VM	Virtual Machine
VNI	Virtualized Network Interface

929

Appendix B—References

- 930
- 931 [1] NIST Special Publication (SP) 800-190, Application Container Security Guide (Draft),
932 National Institute of Standards and Technology, Gaithersburg, Maryland, April 2017.
933
- 934 [2] W.Felter, A.Ferreira, R.Rajamony, J.Rubio., An Updated Performance Comparison of
935 Virtual Machines and Linux Containers, IBM Research Report, RC25482 (AUS1407-
936 001), July 21, 2014.
937
- 938 [3] Practical Guide to Platform-as-a-Service, Version 1.0, Cloud Standards Customer
939 Council, <http://www.cloud-council.org/deliverables/CSCC-Practical-Guide-to-PaaS.pdf>,
940 September 2015.
941
- 942 [4] E. Reshetova et al. Security of OS-level virtualization technologies, Cornell University
943 Library, <https://arxiv.org/abs/1407.4245>
944
- 945 [5] T. Combe, A.Martin, R. Pietro., To Docker or Not to Docker: A Security Perspective,
946 IEEE Computer, September/October 2016, pp, 54-62.
947
- 948 [6] A. Mouat, Docker Security, O'Reilly Media, 2015.
949
- 950 [7] S. Hosseinzadeh, S. Lauren, and V. Leppanen, Security in container-based Virtualization
951 through vTPM. In Proceedings of IEEE/ACM 9th International Conference on Utility and
952 Cloud Computing, Shanghai, China, Pages 214-219, December 2016.
- 953 [8] S. Arnautov et al., SCONE: Secure Linux Containers with Intel SGX, Proceedings of the
954 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI
955 '16). November 2–4, 2016 • Savannah, GA, USA
956 (<https://www.usenix.org/system/files/conference/osdi16/osdi16-arnautov.pdf>).
- 957 [9] grsecurity, <https://grsecurity.net/features.php>
- 958 [10] Home Page of The PaX Team, <https://pax.grsecurity.net/>
- 959 [11] A. Grattafiori, Understanding and Hardening Linux Containers – Version 1.1, NCC
960 Group Whitepaper, June 2016.
- 961 [12] N. Kratzke, About Microservices, Containers and their Underestimated Impact on
962 Network Performance, CLOUD COMPUTING 2015: The Sixth International Conference
963 on Cloud Computing, GRIDs, and Virtualization, pp. 165-169, 2015.
- 964 [13] LxC project: <http://linuxcontainers.org>
- 965 [14] B.Kirsch, *What to choose from the top orchestration software on the market*,
966 [http://searchitoperations.techtarget.com/feature/What-to-choose-from-the-top-orchestration-](http://searchitoperations.techtarget.com/feature/What-to-choose-from-the-top-orchestration-software-on-the-market)
967 [software-on-the-market](http://searchitoperations.techtarget.com/feature/What-to-choose-from-the-top-orchestration-software-on-the-market)
968